



# Inteligencia artificial

## NetApp Solutions

NetApp  
April 25, 2024

# Tabla de contenidos

- Soluciones de inteligencia artificial de NetApp ..... 1
  - Infraestructuras convergentes de IA ..... 1
  - Canalizaciones de datos, lagos de datos y gestión ..... 31
  - Casos de uso ..... 179

# Soluciones de inteligencia artificial de NetApp

## Infraestructuras convergentes de IA

### AFF A400 de NetApp con Lenovo ThinkSystem SR670 V2 para formación de modelos de IA y ML

#### TR-4810: NetApp AFF A400 con Lenovo ThinkSystem SR670 V2 para formación de modelos AI y ML

Sathish Thyagarajan, David Arnette, NetApp Mircea Troaca, Lenovo

Esta solución presenta una arquitectura de clúster de gama media con almacenamiento de NetApp y servidores Lenovo optimizados para cargas de trabajo de inteligencia artificial (IA). Está destinado a empresas pequeñas y medianas para las que la mayoría de tareas informáticas sean de un único nodo (una única o varias GPU) o distribuidas en varios nodos informáticos. Esta solución se alinea con la mayoría de los trabajos de formación de IA diarios para muchas empresas.

En este documento se tratan las pruebas y la validación de una configuración de computación y almacenamiento compuesta por ocho servidores Lenovo SR670V2 de la GPU, un sistema de almacenamiento AFF A400 de NetApp de gama media y un switch de interconexión de 100 GbE. Para medir el rendimiento, utilizamos ResNet50 con el conjunto de datos ImageNET, un tamaño de lote de 408, media precisión, CUDA y cuDNN. Esta arquitectura proporciona una solución eficiente y rentable para organizaciones pequeñas y medianas con el inicio de iniciativas de IA que requieren las funcionalidades de clase empresarial del almacenamiento de datos conectado al cloud ONTAP de NetApp.

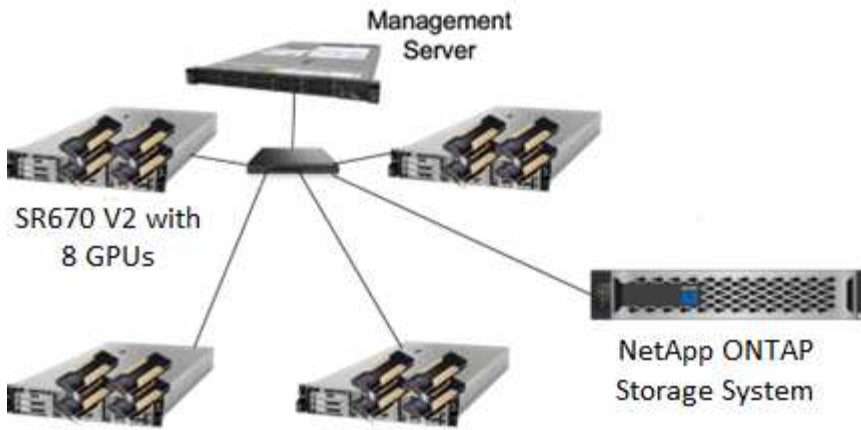
#### Público objetivo

Este documento está dirigido a los siguientes destinatarios:

- Científicos de datos, ingenieros de datos, administradores de datos y desarrolladores de sistemas de IA
- Arquitectos empresariales que diseñan soluciones para el desarrollo de modelos de IA
- Científicos e ingenieros de datos que buscan formas eficientes de alcanzar objetivos de desarrollo DE aprendizaje profundo (DL) y aprendizaje automático (ML)
- Líderes empresariales y responsables DE la toma de decisiones de OT/TECNOLOGÍA que desean conseguir el plazo de comercialización más rápido posible para iniciativas de IA

#### Arquitectura de la solución

Esta solución con los servidores Lenovo ThinkSystem y el sistema ONTAP de NetApp con almacenamiento AFF ha sido diseñado para gestionar el entrenamiento de IA en grandes conjuntos de datos utilizando la potencia de procesamiento de GPU junto con CPU tradicionales. Esta validación demuestra un alto rendimiento y una gestión de datos óptima con una arquitectura de escalado horizontal que utiliza uno, dos o cuatro servidores Lenovo SR670 V2 junto con un único sistema de almacenamiento AFF A400 de NetApp. En la siguiente figura, se proporciona una descripción general de la arquitectura.



Esta solución de NetApp y Lenovo ofrece las siguientes ventajas clave:

- Rendimiento muy eficaz y rentable al ejecutar varias tareas de formación en paralelo
- Rendimiento escalable basado en diferentes números de servidores Lenovo y diferentes modelos de controladoras de almacenamiento de NetApp
- Protección de datos sólida para cumplir con los objetivos de punto de recuperación (RPO) y los objetivos de tiempo de recuperación (RTO) bajos sin pérdida de datos
- Gestión de datos optimizada con copias Snapshot y clones para optimizar los flujos de trabajo de desarrollo

### Información general de la tecnología

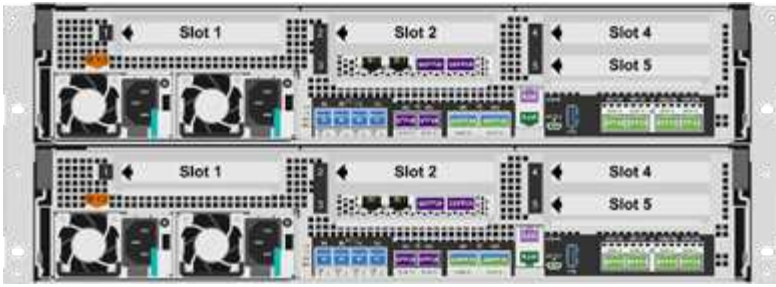
En esta sección se presentan con mayor detalle los principales componentes de esta solución.

#### Sistemas AFF de NetApp

Los sistemas de almacenamiento AFF de NetApp permiten a las empresas cumplir sus requisitos de almacenamiento con un rendimiento líder del sector, una flexibilidad superior, integración con el cloud y la mejor gestión de datos. Los sistemas AFF han sido diseñados específicamente para flash y ayudan a acelerar, gestionar y proteger los datos esenciales para la empresa.







NetApp AFF A400 es un sistema de almacenamiento flash NVMe de gama media que incluye las siguientes funciones:

- Capacidad efectiva máxima: ~20PB
- Escalado horizontal máximo: 2-24 nodos (12 parejas de alta disponibilidad)
- Compatibilidad con host de FC de 25 GbE y 16 GB
- Conectividad RDMA over Converged Ethernet (roce) de 100 GbE con bandejas de almacenamiento de ampliación NVMe
- Los puertos roce de 100 GbE se pueden utilizar para la conexión adjunta de red de host si no se conectan las bandejas NVMe
- Bandejas de ampliación de almacenamiento con conectividad SAS completa de 12 Gbps
- Disponible en dos configuraciones:
  - Ethernet: 4 puertos Ethernet de 25 GB (SFP28)
  - Fibre Channel: 4 puertos FC de 16 GB (SFP+)
- 100 % lectura aleatoria de 8 KB @,4 ms 400 k IOPS

Las funciones AFF A250 de NetApp para puestas en marcha DE IA/ML de gama básica incluyen lo siguiente:

- Capacidad efectiva máxima: 35 PB
- Escalado horizontal máximo: 2-24 nodos (12 parejas de alta disponibilidad)
- 440 000 IOPS lecturas aleatorias a 1 ms
- Basado en la última versión de ONTAP de NetApp, ONTAP 9.8 o posterior
- Dos puertos Ethernet de 25 GB para alta disponibilidad e interconexión de clúster

NetApp también ofrece otros sistemas de almacenamiento, como A800 y AFF A700 de AFF, que proporcionan un mayor rendimiento y escalabilidad para puestas en marcha de IA/ML a gran escala.

### **ONTAP de NetApp**

ONTAP 9, la última generación del software de gestión del almacenamiento de NetApp, permite a las empresas modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. ONTAP ofrece las mejores capacidades de gestión de datos y permite la gestión y protección de los datos con un solo conjunto de herramientas, sin importar dónde residan. Los datos también pueden trasladarse libremente allí donde sea necesario: El perímetro, el núcleo o el cloud. ONTAP 9 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen los datos esenciales y preparan su infraestructura para el futuro con arquitecturas de cloud híbrido.

## Simplificar la gestión de los datos

La gestión de los datos es crucial para las operaciones TECNOLÓGICAS empresariales, de modo que se utilicen recursos apropiados para las aplicaciones y conjuntos de datos. ONTAP incluye las siguientes funciones para facilitar y simplificar las operaciones, y reducir el coste total de las operaciones:

- **Compactación de datos inline y deduplicación expandida.** la compactación de datos reduce el espacio perdido dentro de los bloques de almacenamiento, y la deduplicación aumenta significativamente la capacidad efectiva. Esto es aplicable a los datos almacenados localmente y a los datos organizados en niveles en el cloud.
- **Calidad de servicio (QoS) mínima, máxima y adaptativa.** los controles granulares de QoS ayudan a mantener los niveles de rendimiento para aplicaciones críticas en entornos altamente compartidos.
- **ONTAP FabricPool.** Esta función organiza automáticamente en niveles los datos fríos en opciones de almacenamiento en cloud público y privado, incluido el almacenamiento de objetos Amazon Web Services (AWS), Azure y StorageGRID de NetApp.

## Acelere y proteja sus datos

ONTAP no solo ofrece niveles de rendimiento y protección de datos superiores, sino que amplía estas capacidades de las siguientes maneras:

- **Rendimiento y menor latencia.** ONTAP ofrece el rendimiento más alto posible con la menor latencia posible.
- **Protección de datos.** ONTAP ofrece capacidades integradas de protección de datos con administración común en todas las plataformas.
- **Cifrado de volumen de NetApp.** ONTAP ofrece cifrado nativo a nivel de volumen con compatibilidad para la gestión de claves incorporada y externa.

## Infraestructura preparada para futuros retos

ONTAP 9 ayuda a satisfacer las exigentes y siempre cambiantes necesidades de su empresa:

- **Escalado sencillo y operaciones no disruptivas.** ONTAP admite la adición sin interrupciones de capacidad a las controladoras existentes y a los clústeres de escalado horizontal. Los clientes pueden empezar a utilizar tecnologías punteras como NVMe y FC 32 GB, sin necesidad de realizar costosas migraciones de datos y sin cortes.
- **Conexión en cloud.** ONTAP es el software de gestión del almacenamiento con mejor conexión en cloud, con opciones de almacenamiento definido por software (ONTAP Select) e instancias nativas del cloud (Cloud Volumes Service de NetApp) en todos los clouds públicos.
- **Integración con aplicaciones emergentes.** ONTAP ofrece servicios de datos de clase empresarial para plataformas y aplicaciones de próxima generación como OpenStack, Hadoop y MongoDB usando la misma infraestructura que soporta las aplicaciones empresariales existentes.

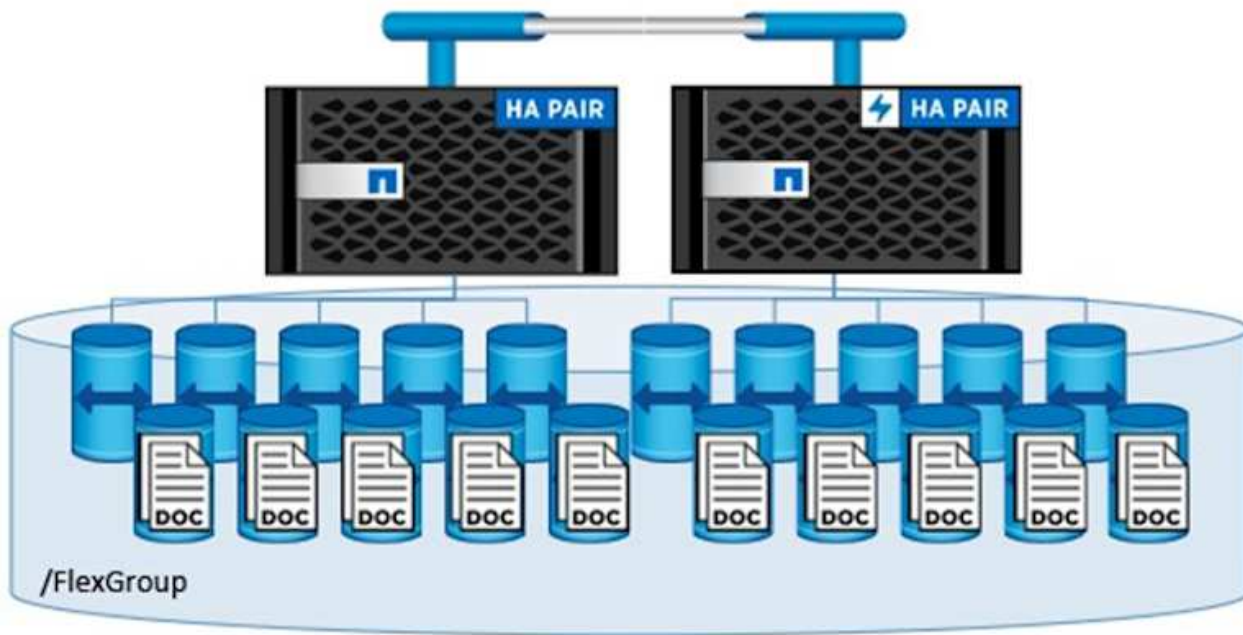
## FlexGroup Volumes de NetApp

Los conjuntos de datos de entrenamiento suelen ser una colección de potencialmente miles de millones de archivos. Pueden ser archivos de texto, de audio, de vídeo o cualquier otra forma de datos no estructurados que deban almacenarse y procesarse para su lectura en paralelo. El sistema de almacenamiento debe almacenar muchos archivos pequeños y debe leerlos en paralelo, con una entrada y salida secuencial o aleatoria

Un volumen FlexGroup (figura a continuación) es un espacio de nombres único compuesto por múltiples

volúmenes constituyentes que se gestiona y actúa como un volumen FlexVol de NetApp entre los administradores de almacenamiento. Los archivos de un volumen de FlexGroup se asignan a volúmenes miembro individuales y no están repartidos en volúmenes o nodos. Ofrecen las siguientes capacidades:

- Hasta 20 petabytes de capacidad y baja latencia predecible para cargas de trabajo con una gran cantidad de metadatos
- Hasta 400 000 millones de archivos en un mismo espacio de nombres
- Operaciones en paralelo para cargas de trabajo NAS entre varias CPU, nodos, agregados y volúmenes FlexVol constituyentes



"

### Gama Lenovo ThinkSystem

Los servidores Lenovo ThinkSystem incluyen hardware, software y servicios innovadores que resuelven los desafíos actuales de los clientes y ofrecen un enfoque de diseño modular, evolutivo y adecuado para su propósito para afrontar los desafíos del futuro. Estos servidores se capitalizan en las mejores tecnologías estándar del sector, junto con innovaciones diferenciadas de Lenovo, para proporcionar la mayor flexibilidad posible en servidores x86.

Entre las ventajas clave de la implementación de servidores Lenovo ThinkSystem se incluyen las siguientes:

- Diseños modulares y de gran escalabilidad que crecen con el negocio
- Resiliencia líder en el sector para ahorrar horas de costosos tiempos de inactividad no programados
- Tecnologías flash rápidas para reducir las latencias, acelerar los tiempos de respuesta y gestionar los datos de forma más inteligente en tiempo real

En el ámbito de la IA, Lenovo está adoptando un enfoque práctico para ayudar a las empresas a comprender y adoptar las ventajas DEL APRENDIZAJE AUTOMÁTICO y la IA para sus cargas de trabajo. Los clientes de Lenovo pueden explorar y evaluar las ofertas de IA de Lenovo en los centros de innovación de IA de Lenovo para comprender por completo el valor de su caso de uso en particular. Con el fin de mejorar la rentabilidad de la inversión, este enfoque centrado en el cliente ofrece a los clientes pruebas de concepto para plataformas de desarrollo de soluciones listas para usar y optimizadas para la IA.

## Lenovo SR670 V2

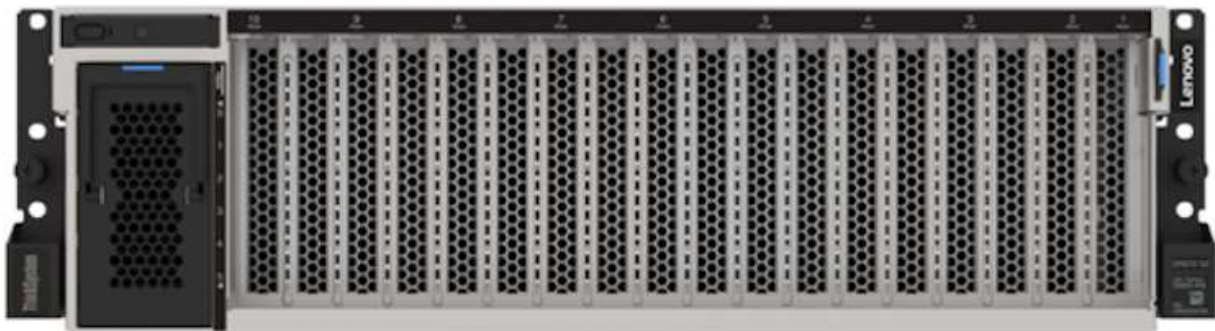
El servidor en rack Lenovo ThinkSystem SR670 V2 ofrece un rendimiento óptimo para una IA acelerada y una informática de alto rendimiento (HPC). Con soporte para hasta ocho GPU, la SR670 V2 es adecuada para los requisitos de cargas de trabajo informáticas intensivas DE ML, DL e inferencia.



4x SXM GPUs with 8x 2.5-inch HS drives and 2x PCIe I/O slots



4x double-wide or 8x single-wide GPU slots and 2x PCIe I/O slots  
with 8x 2.5-inch or 4x 3.5-inch HS drives



8x double-wide GPU slots with 6x EDSFF HS drives and 2x PCIe I/O slots

Con las CPU Intel Xeon más recientes y escalables que admiten GPU de gama alta (incluida la GPU 8x PCIe NVIDIA A100 de 80 GB), ThinkSystem SR670 V2 ofrece un rendimiento optimizado y acelerado para cargas de trabajo de IA y computación de alto rendimiento.

Puesto que en más cargas de trabajo se utiliza el rendimiento de los aceleradores, ha aumentado la demanda de densidad de GPU. Sectores como el comercio minorista, los servicios financieros, la energía y la sanidad utilizan GPU para obtener una mayor información e impulsar la innovación con APRENDIZAJE AUTOMÁTICO, DL y técnicas de inferencia.

ThinkSystem SR670 V2 es una solución empresarial optimizada para poner en marcha cargas de trabajo aceleradas de HPC e IA en la producción, maximizando el rendimiento del sistema a la vez que mantiene la densidad del centro de datos para los clústeres de supercomputación con plataformas de última generación.

Entre otras funciones se incluyen las siguientes:

- Compatibilidad con I/O RDMA directa de GPU en la que los adaptadores de red de alta velocidad están conectados directamente a las GPU para maximizar el rendimiento de I/O.
- Compatibilidad con almacenamiento directo de GPU en el que las unidades NVMe están conectadas directamente a las GPU para maximizar el rendimiento del almacenamiento.

**Rendim. MLPerf**

MLPerf es el conjunto de pruebas de rendimiento líder del sector para evaluar el rendimiento de la IA. En esta validación, utilizamos su punto de referencia de clasificación de imágenes con MXNet, uno de los marcos de IA más populares. El script de formación MXNet\_Benchmark se utilizó para impulsar la formación de IA. El script contiene implementaciones de varios modelos convencionales populares y está diseñado para ser lo más rápido posible. Puede ejecutarse en una sola máquina o ejecutarse en modo distribuido entre varios hosts.

**Plan de pruebas**

En esta validación, realizamos el entrenamiento de reconocimiento de imágenes según lo especificado por MLPerf v2.0. Específicamente, entrenamos el modelo ResNet v2.0 con el conjunto de datos ImageNET hasta que alcanza una precisión del 76.1 %. La métrica principal es el tiempo para alcanzar la precisión deseada. También reportamos formación de ancho de banda en imágenes por segundo para valorar mejor la eficiencia del escalado horizontal.

El caso de prueba principal evaluó la ejecución simultánea de varios procesos de entrenamiento independientes (uno por nodo). Simula el caso práctico principal, un sistema compartido utilizado por varios científicos de datos. En el segundo caso de prueba se evaluó la eficiencia del escalado horizontal.

**Resultados de la prueba**

La tabla siguiente resume los resultados de todas las pruebas realizadas para esta solución.

Descripción de la prueba	Resumen de resultados
Formación de reconocimiento de imágenes: Varios trabajos simultáneos	Rendimiento muy eficiente. Todos los trabajos se ejecutaron a toda velocidad incluso cuando se utilizaba al máximo el clúster. Los sistemas de almacenamiento de NetApp proporcionaron un rendimiento de formación comparable al almacenamiento SSD local a la vez que posibilitan un uso compartido sencillo de los datos entre servidores.

Descripción de la prueba	Resumen de resultados
Formación para el reconocimiento de imágenes: Escalado horizontal	Altamente eficiente para un máximo de cuatro nodos. En ese momento, la escalabilidad horizontal era menos eficiente pero todavía factible. El uso de una red computacional de mayor velocidad mejora la escalabilidad. El sistema de almacenamiento de NetApp proporcionó un rendimiento de formación comparable al almacenamiento en unidades de estado sólido locales a la vez que permitió compartir fácilmente los datos entre servidores.

## Configuración de prueba

En este apartado se describen las configuraciones probadas, la infraestructura de red, el servidor SR670 V2 y los detalles sobre el aprovisionamiento de almacenamiento NetApp.

## Arquitectura de la solución

Hemos utilizado los componentes de la solución enumerados en la siguiente tabla para esta validación.

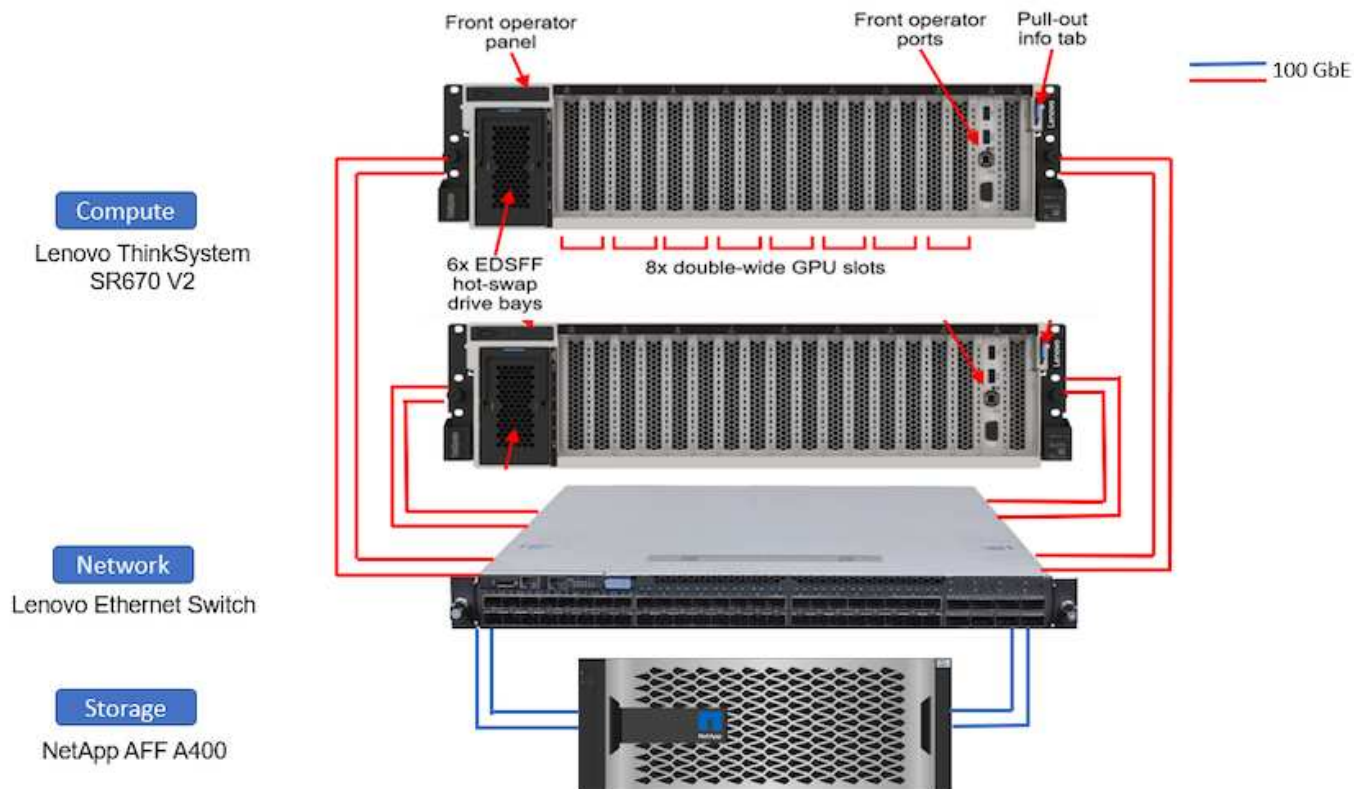
Componentes de la solución	Detalles
Servidores Lenovo ThinkSystem	<ul style="list-style-type: none"> <li>• Dos servidores SR670 V2 cada uno con ocho tarjetas GPU NVIDIA A100 de 80 GB</li> <li>• Cada servidor contiene 2 CPU Intel Xeon Platinum 8360Y (28 núcleos físicos) y 1 TB de RAM</li> </ul>
Linux (Ubuntu – 20.04 con CUDA 11.8)	
Sistema de almacenamiento AFF de NetApp (pareja de alta disponibilidad)	<ul style="list-style-type: none"> <li>• Software ONTAP 9.10.1 de NetApp</li> <li>• 24 unidades SSD de 960 GB</li> <li>• Protocolo NFS</li> <li>• 1 grupo de interfaces (ifgrp) por controladora, con cuatro direcciones IP lógicas para los puntos de montaje</li> </ul>

En esta validación, utilizamos ResNet v2.0 con la base ImageNET establecida como lo especifica MLPerf v2.0. El conjunto de datos se almacena en un sistema de almacenamiento AFF de NetApp con el protocolo NFS. Los SR670 se conectaron al sistema de almacenamiento AFF A400 de NetApp a través de un switch de 100 GbE.

ImageNET es un conjunto de datos de imágenes utilizado con frecuencia. Contiene casi 1.3 millones de imágenes por un tamaño de 144 GB. El tamaño medio de la imagen es de 108 KB.

En la siguiente figura, se muestra la topología de red de la configuración probada.





### Controladora de almacenamiento

La siguiente tabla enumera la configuración de almacenamiento.

Controladora	Agregado	Volumen FlexGroup	Tamaño del agregado	Tamaño del volumen	Punto de montaje del sistema operativo
Control1	Agr1	/a400-100g	9,9 TB	19 TB	/a400-100g
Control2	Agr2	/a400-100g	9,9 TB		/a400-100g



La carpeta /a400-100g contiene el conjunto de datos utilizado para la validación ResNet.

### Procedimiento de prueba y resultados detallados

En esta sección se describen los resultados detallados del procedimiento de prueba.

#### El entrenamiento de reconocimiento de imágenes mediante ResNet en ONTAP

Ejecutamos la prueba de rendimiento ResNet50 con uno y dos servidores SR670 V2. Esta prueba utilizó el contenedor MXNet 22.04-py3 NGC para realizar el entrenamiento.

Hemos utilizado el siguiente procedimiento de prueba en esta validación:

1. Se borró la caché del host antes de ejecutar el script para garantizar que los datos no se almacenaban en caché:

```
sync ; sudo /sbin/sysctl vm.drop_caches=3
```

2. Ejecutamos el script de pruebas de rendimiento con el conjunto de datos ImageNET en el almacenamiento de servidores (almacenamiento SSD local) y en el sistema de almacenamiento AFF de NetApp.
3. Validamos el rendimiento del almacenamiento local y en red con el dd comando.
4. Para la ejecución de un solo nodo, utilizamos el siguiente comando:

```
python train_imagenet.py --gpus 0,1,2,3,4,5,6,7 --batch-size 408 --kv
-store horovod --lr 10.5 --mom 0.9 --lr-step-epochs pow2 --lars-eta
0.001 --label-smoothing 0.1 --wd 5.0e-05 --warmup-epochs 2 --eval-period
4 --eval-offset 2 --optimizer sgdwfastlars --network resnet-v1b-stats-fl
--num-layers 50 --num-epochs 37 --accuracy-threshold 0.759 --seed 27081
--dtype float16 --disp-batches 20 --image-shape 4,224,224 --fuse-bn-relu
1 --fuse-bn-add-relu 1 --bn-group 1 --min-random-area 0.05 --max-random
-area 1.0 --conv-algo 1 --force-tensor-core 1 --input-layout NHWC --conv
-layout NHWC --batchnorm-layout NHWC --pooling-layout NHWC --batchnorm
-mom 0.9 --batchnorm-eps 1e-5 --data-train /data/train.rec --data-train
-idx /data/train.idx --data-val /data/val.rec --data-val-idx
/data/val.idx --dali-dont-use-mmap 0 --dali-hw-decoder-load 0 --dali
-prefetch-queue 5 --dali-nvjpeg-memory-padding 256 --input-batch
-multiplier 1 --dali-threads 6 --dali-cache-size 0 --dali-roi-decode 1
--dali-preallocate-width 5980 --dali-preallocate-height 6430 --dali-tmp
-buffer-hint 355568328 --dali-decoder-buffer-hint 1315942 --dali-crop
-buffer-hint 165581 --dali-normalize-buffer-hint 441549 --profile 0
--e2e-cuda-graphs 0 --use-dali
```

5. Para las ejecuciones distribuidas, utilizamos el modelo de paralelización del servidor de parámetros. Se utilizaron dos servidores de parámetros por nodo y se establece el número de épocas igual que para la ejecución de un único nodo. Lo hicimos porque el entrenamiento distribuido a menudo toma más épocas debido a la sincronización imperfecta entre procesos. El número diferente de épocas puede sesgar comparaciones entre casos de un solo nodo y distribuidos.

#### Velocidad de lectura de los datos: Local frente a almacenamiento en red

La velocidad de lectura se probó con el dd Comando en uno de los archivos para el conjunto de datos ImageNET. Específicamente, ejecutamos los siguientes comandos para datos locales y de red:

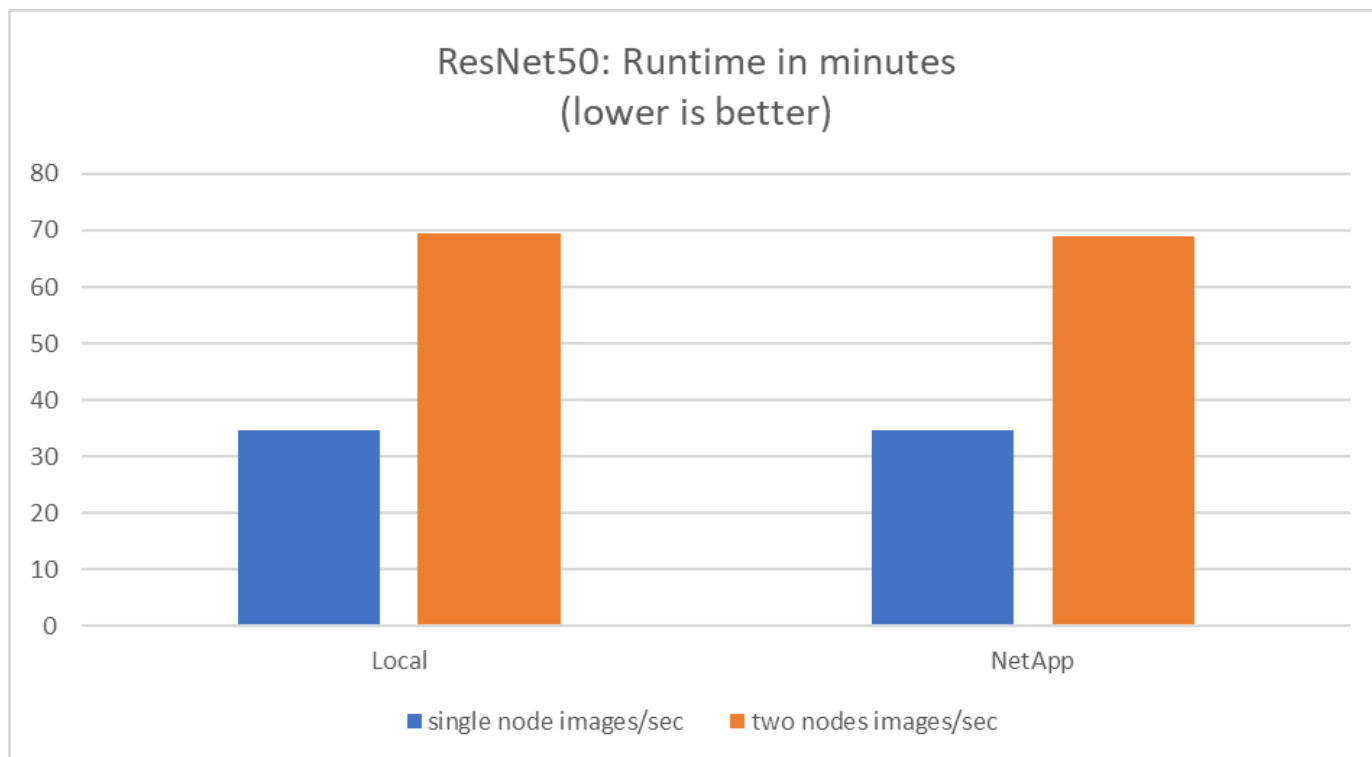
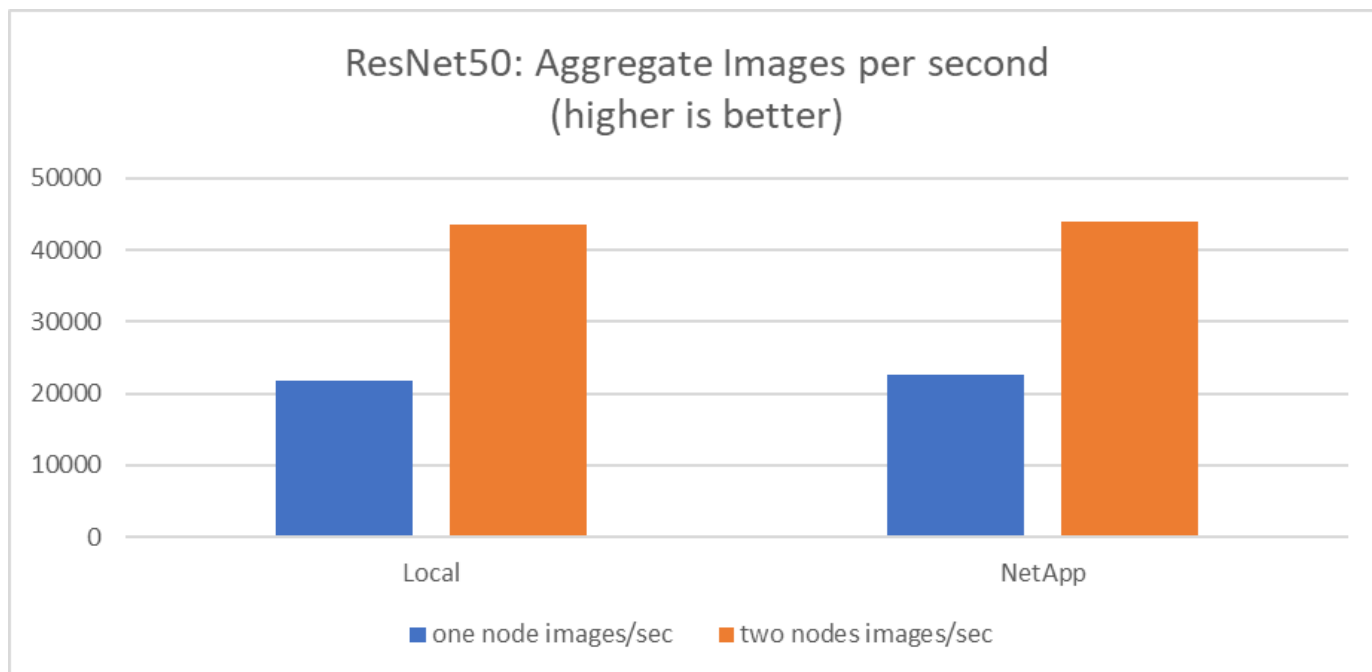
```
sync ; sudo /sbin/sysctl vm.drop_caches=3dd if=/a400-100g/netapp-
ra/resnet/data/preprocessed_data/train.rec of=/dev/null bs=512k
count=2048Results (average of 5 runs):
Local storage: 1.7 GB/s Network storage: 1.5 GB/s.
```

Ambos valores son similares, lo que demuestra que el almacenamiento en red puede ofrecer datos a un ritmo similar al del almacenamiento local.



### Caso de uso compartido: Trabajos múltiples, independientes y simultáneos

Esta prueba simuló el caso de uso esperado para esta solución: Formación de IA para varios trabajos y varios usuarios. Cada nodo ejecutó su propia formación mientras utiliza el almacenamiento de red compartida. Los resultados se muestran en la siguiente figura, que muestra que el caso de la solución proporcionó un rendimiento excelente con todos los trabajos que se ejecutan básicamente a la misma velocidad que los trabajos individuales. El rendimiento total se escaló de forma lineal con el número de nodos.



Estos gráficos presentan el tiempo de ejecución en minutos y las imágenes agregadas por segundo para los nodos de computación que utilizaron ocho GPU de cada servidor en una red de cliente de 100 GbE, combinando tanto el modelo de entrenamiento simultáneo como el único modelo de entrenamiento. El tiempo

de ejecución medio del modelo de entrenamiento fue de 35 minutos y 9 segundos. Los tiempos de ejecución individuales fueron de 34 minutos y 32 segundos, 36 minutos y 21 segundos, 34 minutos y 37 segundos, 35 minutos y 25 segundos, y 34 minutos y 31 segundos. El promedio de imágenes por segundo para el modelo de entrenamiento fue de 22,573, y las imágenes individuales por segundo fueron de 21,764; 23,438; 22,556; 22,564; y 22,547.

Según nuestra validación, un modelo de entrenamiento independiente con tiempo de ejecución de datos de NetApp fue de 34 minutos y 54 segundos con 22,231 imágenes/s. Un modelo de entrenamiento independiente con un tiempo de ejecución de datos locales (DAS) era de 34 minutos y 21 segundos con 22,102 imágenes/s. Durante estas operaciones, el uso medio de la GPU era del 96 %, como se observó en nvidia-smi. Tenga en cuenta que esta media incluye la fase de prueba, durante la cual no se utilizaron las GPU, mientras que el uso de CPU fue del 40 % medido por mpstat. Esto demuestra que la tasa de entrega de datos es suficiente en cada caso.

## **Ajustes de arquitectura**

La configuración utilizada para esta validación se puede ajustar para ajustarse a otros casos prácticos.

### **Ajustes de CPU**

Utilizamos un procesador Intel Xeon Platinum 8360Y Skylake para esta validación, según lo recomendado por Lenovo. Esperamos que la CPU equivalente de Cascade Lake, un procesador Intel Xeon Gold 6330, proporcione un rendimiento similar porque esta carga de trabajo no está vinculada a la CPU.

### **Aumente la capacidad de almacenamiento**

En función de sus necesidades de capacidad de almacenamiento, puede aumentar el espacio de almacenamiento compartido (volumen NFS) bajo demanda siempre que disponga de las bandejas de discos y los modelos de controladora adicionales. Puede hacerlo desde la CLI o desde la interfaz web de NetApp del controlador de almacenamiento como usuario administrador.

## **Conclusión**

La solución de NetApp y Lenovo validada aquí es una arquitectura flexible de escalado horizontal ideal para IA empresarial de gama media.

El almacenamiento de NetApp proporciona el mismo rendimiento o mejor que el almacenamiento SSD local y ofrece las siguientes ventajas a los científicos de datos, ingenieros de datos y responsables de la toma DE decisiones TECNOLÓGICAS:

- Uso compartido de datos sin esfuerzo entre sistemas de IA, análisis y otros sistemas de negocio cruciales. Este uso compartido de datos reduce la sobrecarga de la infraestructura, mejora el rendimiento y optimiza la gestión de datos en toda la empresa.
- Almacenamiento y cálculo escalables independientemente para minimizar los costes y mejorar la utilización de recursos.
- Flujos de trabajo de desarrollo y puesta en marcha optimizados que utilizan copias Snapshot y clones integrados para espacios de trabajo de usuario instantáneos con gestión eficiente del espacio, control de versiones integrado y puesta en marcha automatizada.
- Protección de datos de clase empresarial para recuperación ante desastres y continuidad empresarial.

## Reconocimientos

- Karthikeyan Nagalingam, Ingeniero Técnico de Marketing de NetApp
- Jarrett Upton, Admin, AI Lab Systems, Lenovo

## Dónde encontrar información adicional

Para obtener más información sobre la información descrita en este documento, consulte los siguientes documentos y/o sitios web:

- Página de producto de cabinas all-flash de NetApp  
["https://www.netapp.com/us/products/storage-systems/all-flash-array/aff-a-series.aspx"](https://www.netapp.com/us/products/storage-systems/all-flash-array/aff-a-series.aspx)
- Página AFF A400 de NetApp  
["https://docs.netapp.com/us-en/ontap-systems/a400/index.html"](https://docs.netapp.com/us-en/ontap-systems/a400/index.html)
- Página de producto del software de gestión de datos ONTAP de NetApp  
["http://www.netapp.com/us/products/data-management-software/ontap.aspx"](http://www.netapp.com/us/products/data-management-software/ontap.aspx)
- Rendim. MLPerf  
["https://mlperf.org"](https://mlperf.org)
- Prueba de rendimiento de TensorFlow  
["https://github.com/tensorflow/benchmarks"](https://github.com/tensorflow/benchmarks)
- SMI de NVIDIA (nvidia-smi)  
["https://developer.nvidia.com/nvidia-system-management-interface"](https://developer.nvidia.com/nvidia-system-management-interface)

## IA de NetApp con NVIDIA

Descripción general de las soluciones de infraestructura convergente de ONTAP AI de NetApp y NVIDIA.

### AI Pod de NetApp con sistemas NVIDIA DGX

- ["NetApp AI Pod con sistemas NVIDIA DGX"](#)

### ONTAP AI de NetApp con sistemas NVIDIA DGX A100

- ["Guía de diseño"](#)
- ["Guía de puesta en marcha"](#)

### ONTAP AI de NetApp con sistemas NVIDIA DGX A100 y switches Mellanox Spectrum Ethernet

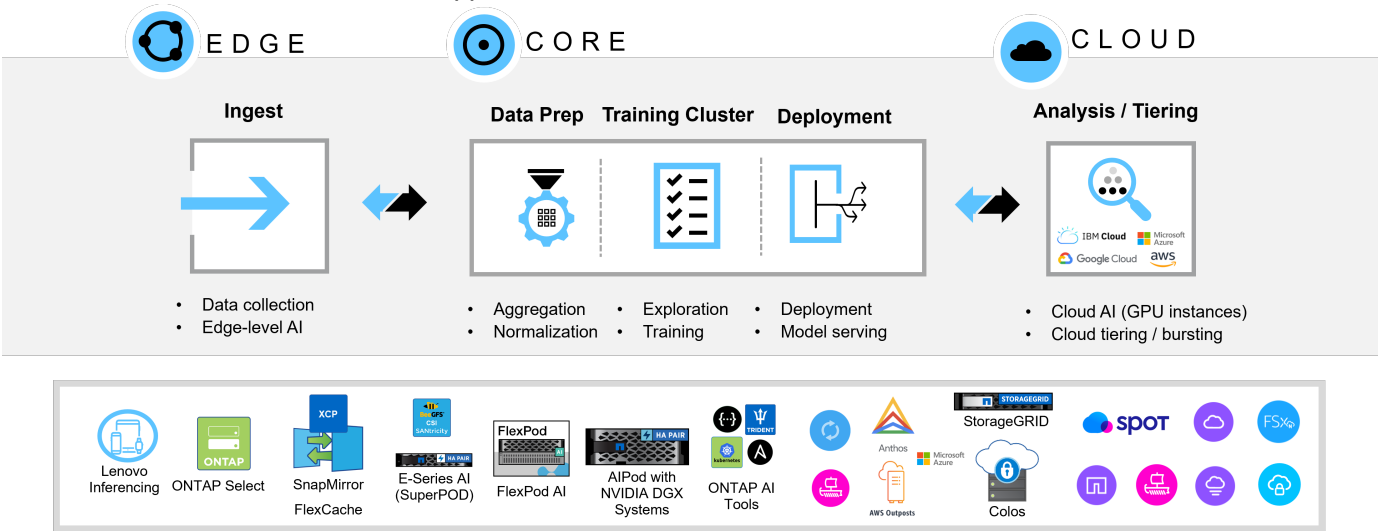
- ["Guía de diseño"](#)
- ["Guía de puesta en marcha"](#)

NetApp AIPod con sistemas NVIDIA DGX: Introducción

Ingeniería de soluciones de NetApp

NetApp™ AIPod con NVIDIA DGX™ los sistemas y los sistemas de almacenamiento conectados al cloud de NetApp simplifican las puestas en marcha de infraestructuras para cargas de trabajo de aprendizaje automático (ML) e inteligencia artificial (IA) al eliminar complejidades y conjeturas en la fase de diseño. Basándose en el diseño de NVIDIA DGX BasePOD para proporcionar un rendimiento informático excepcional para cargas de trabajo de nueva generación, AIPod con sistemas NVIDIA DGX añade sistemas de almacenamiento NetApp AFF que permiten a los clientes empezar poco a poco y crecer sin interrupciones, a la vez que gestionan de forma inteligente los datos desde el perímetro al núcleo y al cloud, y viceversa. NetApp AIPod forma parte de la cartera más amplia de soluciones de IA de NetApp, muestra la figura a continuación:

Cartera de soluciones de IA de NetApp



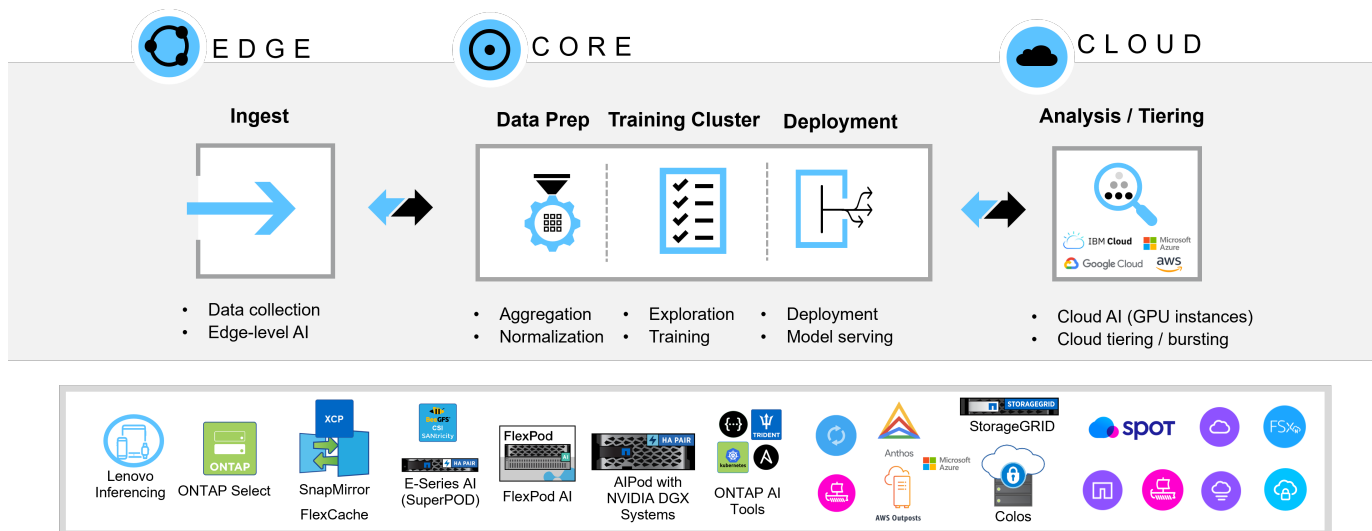
En este documento se describen los componentes clave de la arquitectura de referencia AIPod, la información sobre conectividad del sistema y la guía para dimensionar la solución. Este documento está destinado a ingenieros de soluciones de NetApp y partners, y para aquellos clientes interesados en poner en marcha una infraestructura de alto rendimiento para cargas de trabajo de análisis y ML/DL.

NetApp AIPod con sistemas NVIDIA DGX: Introducción

Ingeniería de soluciones de NetApp

NetApp™ AIPod con NVIDIA DGX™ los sistemas y los sistemas de almacenamiento conectados al cloud de NetApp simplifican las puestas en marcha de infraestructuras para cargas de trabajo de aprendizaje automático (ML) e inteligencia artificial (IA) al eliminar complejidades y conjeturas en la fase de diseño. Basándose en el diseño de NVIDIA DGX BasePOD para proporcionar un rendimiento informático excepcional para cargas de trabajo de nueva generación, AIPod con sistemas NVIDIA DGX añade sistemas de almacenamiento NetApp AFF que permiten a los clientes empezar poco a poco y crecer sin interrupciones, a la vez que gestionan de forma inteligente los datos desde el perímetro al núcleo y al cloud, y viceversa. NetApp AIPod forma parte de la cartera más amplia de soluciones de IA de NetApp, muestra la figura a continuación:

Cartera de soluciones de IA de NetApp



En este documento se describen los componentes clave de la arquitectura de referencia AI Pod, la información sobre conectividad del sistema y la guía para dimensionar la solución. Este documento está destinado a ingenieros de soluciones de NetApp y partners, y para aquellos clientes interesados en poner en marcha una infraestructura de alto rendimiento para cargas de trabajo de análisis y ML/DL.

## AI Pod de NetApp con sistemas NVIDIA DGX: Componentes de hardware

### Sistemas de almacenamiento de NetApp AFF

Con los vanguardistas sistemas de almacenamiento de NetApp AFF, los departamentos TECNOLÓGICOS pueden satisfacer los requisitos de almacenamiento del negocio con un rendimiento de primera clase, una flexibilidad superior, integración con el cloud y la mejor gestión de datos. Los sistemas AFF han sido diseñados específicamente para flash y ayudan a acelerar, gestionar y proteger los datos esenciales para la empresa.

### Sistemas de almacenamiento A900 de AFF

La serie NetApp AFF A900, impulsada por el software para la gestión de datos NetApp ONTAP, ofrece protección de datos integrada, funcionalidades opcionales anti-ransomware y el alto rendimiento y la resiliencia necesarios para dar soporte a las cargas de trabajo empresariales más importantes. Acaba con las interrupciones de las operaciones esenciales para la empresa, reduce al máximo el ajuste del rendimiento y protege sus datos de los ataques de ransomware. Proporciona:

- Rendimiento líder en el sector
- Seguridad de los datos sin concesiones
- Actualizaciones simplificadas y sin interrupciones

### Sistema de almacenamiento A900 de NetApp AFF



### **Rendimiento líder del sector**

El AFF A900 gestiona con facilidad las cargas de trabajo de última generación, como el aprendizaje profundo, la IA y los análisis de alta velocidad, así como las bases de datos empresariales tradicionales como Oracle, SAP HANA, Microsoft SQL Server y aplicaciones virtualizadas. Mantiene las aplicaciones vitales para el negocio ejecutándose a máxima velocidad con hasta 2,4m 000 IOPS por pareja de alta disponibilidad y una latencia de hasta 100µs, y aumenta el rendimiento hasta en un 50 % con respecto a los modelos de NetApp anteriores. Con NFS a través de RDMA, pNFS y Session Trunking, los clientes pueden lograr el alto nivel de rendimiento de red necesario para las aplicaciones de nueva generación usando la infraestructura de red existente del centro de datos.

Además, los clientes pueden escalar y crecer con la compatibilidad multiprotocolo unificada para el almacenamiento SAN, NAS y de objetos y ofrecer la máxima flexibilidad con un único software de gestión de datos ONTAP unificado, tanto para los datos locales como en el cloud. Además, el estado del sistema se puede optimizar con los análisis predictivos basados en IA que ofrecen Active IQ y Cloud Insights.

### **Seguridad de los datos garantizada**

Los sistemas A900 de AFF contienen una suite completa del software de NetApp para la protección de datos consistente con las aplicaciones e integrado. Ofrece una protección de datos integrada y soluciones innovadoras anti-ransomware que facilitan la recuperación previa y posterior al ataque. Los archivos maliciosos pueden bloquearse para que no se escriban en el disco, y las anomalías de almacenamiento se supervisan fácilmente para obtener información.

### **Actualizaciones simplificadas y fluidas**

La AFF A900 está disponible como una actualización fluida e integrada para los clientes actuales de A700.

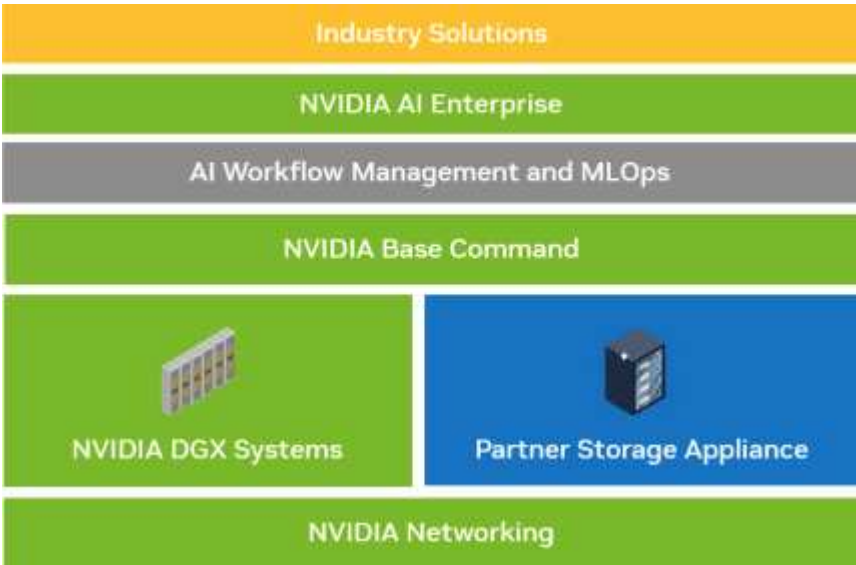
NetApp facilita la actualización y la eliminación de las interrupciones de las operaciones vitales para la empresa gracias a nuestras funciones avanzadas de fiabilidad, disponibilidad, aptitud y gestión (RASM). Además, NetApp aumenta aún más la eficiencia operativa y simplifica las actividades cotidianas de los equipos TECNOLÓGICOS, ya que el software ONTAP aplica automáticamente las actualizaciones de firmware de todos los componentes del sistema.

Para las mayores puestas en marcha, los sistemas AFF A900 ofrecen las opciones más elevadas de rendimiento y capacidad, mientras que otros sistemas de almacenamiento de NetApp, como AFF A800, AFF C800, AFF A400, AFF C400 y AFF A250 ofrecen opciones para puestas en marcha de menor tamaño a puntos de coste.

**DGX BasePOD de NVIDIA**

NVIDIA DGX BasePOD es una solución integrada que consta de componentes de hardware y software de NVIDIA, soluciones de MLOps y almacenamiento de terceros. Al aprovechar las mejores prácticas del diseño de sistemas de escalado horizontal con productos NVIDIA y soluciones de partners validadas, los clientes pueden implementar una plataforma eficiente y gestionable para el desarrollo de la IA. La figura 1 destaca los distintos componentes de NVIDIA DGX BasePOD.

*Solución NVIDIA DGX BasePOD*



**SISTEMAS DGX H100 DE NVIDIA**

El sistema NVIDIA DGX H100™ es el centro neurálgico de la IA que se acelera gracias al innovador rendimiento de la GPU Tensor Core de NVIDIA H100.

*SISTEMA NVIDIA DGX H100*



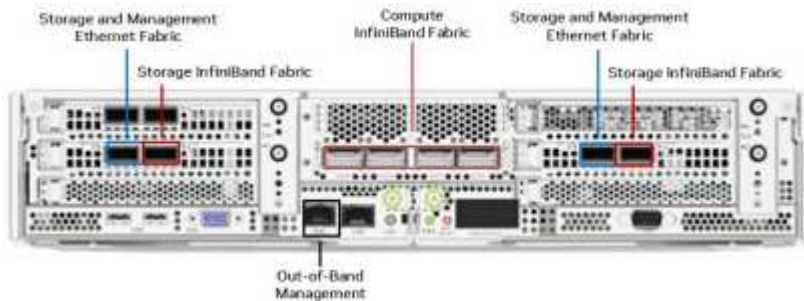
Las especificaciones clave del sistema DGX H100 son:



- Ocho GPU NVIDIA H100.
- 80 GB de memoria GPU por GPU, para un total de 640GB.
- Cuatro chips NVIDIA NVSwitch™.
- Procesadores Intel® Xeon® Platinum 8480 duales de 56 núcleos con soporte PCIe 5,0.
- 2 TB de memoria del sistema DDR5.
- Cuatro puertos OSFP que sirven a ocho adaptadores NVIDIA ConnectX-7 (InfiniBand/Ethernet) de un puerto y dos adaptadores NVIDIA ConnectX-7 (InfiniBand/Ethernet) de doble puerto.
- Dos unidades NVMe M,2 de 1,92 TB para el sistema operativo DGX y ocho unidades NVMe U,2 de 3,84 TB para almacenamiento/caché.
- 10,2 kW de potencia máxima.

A continuación, se muestran los puertos posteriores de la bandeja de CPU DGX H100. Cuatro de los puertos OSFP sirven a ocho adaptadores ConnectX-7 para la estructura de computación InfiniBand. Cada par de adaptadores ConnectX-7 de doble puerto proporciona vías paralelas para los tejidos de almacenamiento y gestión. El puerto fuera de banda se utiliza para el acceso BMC.

#### *Panel posterior NVIDIA DGX H100*



### **Conexión a redes NVIDIA**

#### **Switch NVIDIA Quantum-2 QM9700**

##### *Conmutador InfiniBand NVIDIA Quantum-2 QM9700*



Los switches NVIDIA Quantum-2 QM9700 con conectividad InfiniBand de 400GB Gb/s potencian la estructura de computación en las configuraciones InfiniBand BasePOD de NVIDIA Quantum-2. Se utilizan adaptadores de un solo puerto ConnectX-7 para la estructura de computación InfiniBand. Cada sistema NVIDIA DGX tiene conexiones dobles a cada switch de QM9700 GbE, lo que proporciona múltiples rutas de alto ancho de banda y baja latencia entre los sistemas.

#### **Conmutador NVIDIA Spectrum-3 SN4600**

##### *Conmutador NVIDIA Spectrum-3 SN4600*



Los switches NVIDIA Spectrum-3 SN4600 ofrecen 128 puertos en total (64 por switch) con el fin de proporcionar conectividad redundante para la gestión en banda del DGX BasePOD. El switch NVIDIA SN4600



puede proporcionar velocidades de entre 1 GbE y 200 GbE. Para los dispositivos de almacenamiento conectados con Ethernet, también se utilizan los switches NVIDIA SN4600. Los puertos de los adaptadores ConnectX-7 de doble puerto NVIDIA DGX se utilizan para la conectividad de almacenamiento y gestión en banda.

### **Conmutador NVIDIA Spectrum SN2201**

*Conmutador NVIDIA Spectrum SN2201*



Los switches NVIDIA Spectrum SN2201 ofrecen puertos 48 GbE para proporcionar conectividad para la gestión fuera de banda. La gestión fuera de banda proporciona conectividad de gestión consolidada para todos los componentes de DGX BasePOD.

### **Adaptador NVIDIA ConnectX-7**

*Adaptador NVIDIA ConnectX-7*



El adaptador NVIDIA ConnectX-7 puede proporcionar 25/50/100/200/400G de rendimiento. Los sistemas NVIDIA DGX utilizan los adaptadores ConnectX-7 de puerto único o doble para proporcionar flexibilidad en puestas en marcha de DGX BasePOD con InfiniBand de 400GB Gb/s y Ethernet 100/200GB Gb.

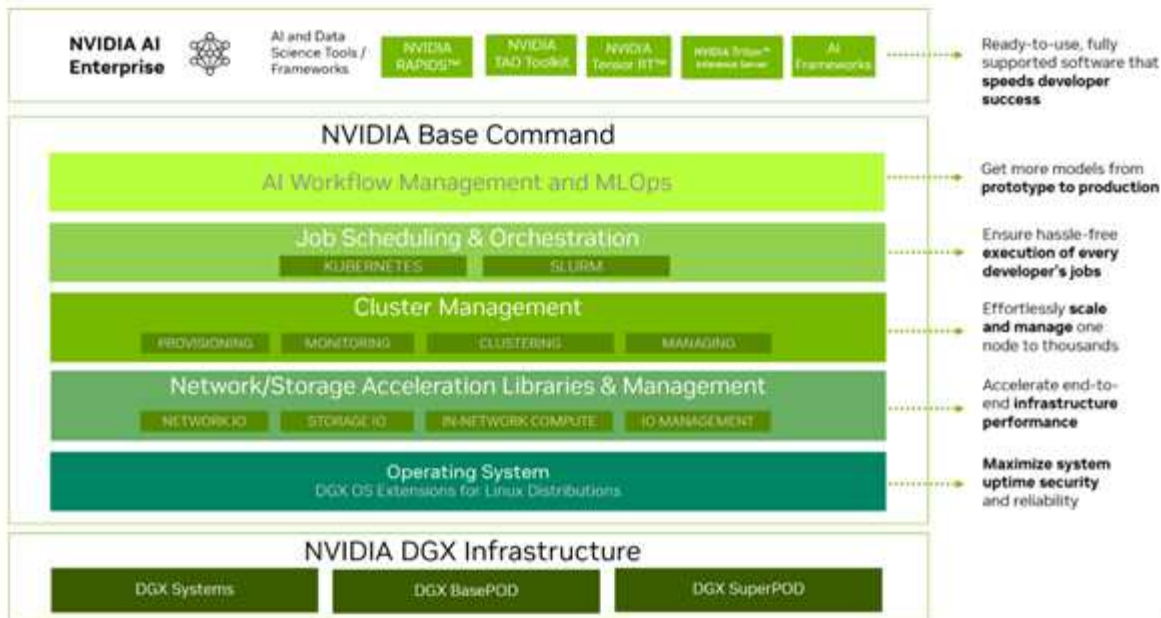
### **AIPOD de NetApp con sistemas NVIDIA DGX: Componentes de software**

#### **Software NVIDIA**

#### **Comando base NVIDIA**

NVIDIA Base Command™ impulsa todos los DGX BasePOD, lo que permite a las organizaciones aprovechar lo mejor de la innovación en software de NVIDIA. Las empresas pueden liberar todo el potencial de su inversión con una plataforma contrastada que incluye orquestación y gestión de clústeres de clase empresarial, bibliotecas que aceleran la computación, la infraestructura de redes y almacenamiento, y un sistema operativo (SO) optimizado para cargas de trabajo de IA.

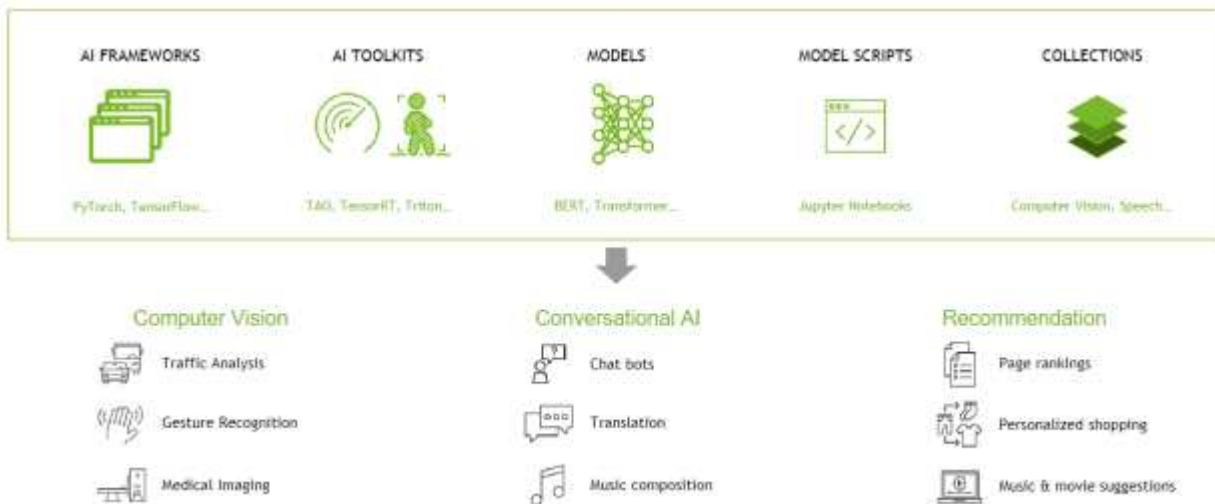
*Solución NVIDIA BaseCommand*



## GPU CLOUD DE NVIDIA (NGC)

NVIDIA NGC™ proporciona software para satisfacer las necesidades de los científicos de datos, desarrolladores e investigadores con varios niveles de experiencia en IA. El software alojado en NGC se somete a análisis contra un conjunto agregado de vulnerabilidades y exposiciones comunes (CVE), cripto y claves privadas. Está probado y diseñado para escalar a varias GPU y, en muchos casos, a varios nodos, garantizando que los usuarios maximicen la inversión en sistemas DGX.

### NVIDIA GPU Cloud



## IA Enterprise de NVIDIA

NVIDIA AI Enterprise es la plataforma de software integral que lleva la IA generativa al alcance de todas las empresas, lo que proporciona el tiempo de ejecución más rápido y eficiente para modelos de base de IA generativa optimizados para ejecutarse en la plataforma NVIDIA DGX. Gracias a la seguridad, estabilidad y capacidad de gestión del nivel de producción, se optimiza el desarrollo de soluciones de IA generativas. NVIDIA AI Enterprise se incluye en DGX BasePOD para que los desarrolladores empresariales accedan a modelos previamente entrenados, marcos optimizados, microservicios, bibliotecas aceleradas y soporte empresarial.

## Software NetApp

### ONTAP de NetApp

ONTAP 9, la última generación del software de gestión del almacenamiento de NetApp, permite a las empresas modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. ONTAP ofrece las mejores capacidades de gestión de datos y permite la gestión y protección de los datos con un solo conjunto de herramientas, sin importar dónde residan. También puede mover los datos libremente a donde sea necesario: El perímetro, el núcleo o el cloud. ONTAP 9 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen los datos esenciales y permiten disfrutar de funcionalidades de infraestructura de nueva generación en arquitecturas de cloud híbrido.

### Acelere y proteja sus datos

ONTAP no solo ofrece niveles de rendimiento y protección de datos superiores, sino que amplía estas capacidades de las siguientes maneras:

- Rendimiento y menor latencia. ONTAP ofrece la mayor salida posible con la menor latencia posible, incluida la compatibilidad con el almacenamiento GPUDirect (GDS) de NVIDIA mediante NFS a través de RDMA, Parallel NFS (pNFS) y conexión de enlaces de sesiones NFS.
- Protección de datos. ONTAP ofrece funcionalidades de protección de datos integradas y la mayor garantía antiransomware del sector, con una gestión común en todas las plataformas.
- Cifrado de volúmenes de NetApp (NVE). ONTAP ofrece cifrado nativo en el nivel de volumen y permite la gestión de claves incorporada o externa.
- Almacenamiento multi-tenancy y autenticación multifactor. ONTAP permite compartir recursos de infraestructura con los niveles más altos de seguridad.

### Simplificar la gestión de los datos

La gestión de los datos es crucial para las operaciones TECNOLÓGICAS empresariales y los científicos de datos, para que se utilicen recursos apropiados para las aplicaciones de IA y para entrenar conjuntos de datos de IA/ML. La siguiente información adicional sobre las tecnologías de NetApp no está disponible para esta validación, pero puede ser relevante en función de su puesta en marcha.

El software para la gestión de datos ONTAP incluye las siguientes funciones para mejorar y simplificar las operaciones, y reducir el coste total de funcionamiento:

- Las copias Snapshot y clones permiten la colaboración, la experimentación paralela y una gestión de datos mejorada para los flujos de trabajo DE ML/DL.
- SnapMirror posibilita un movimiento fluido de datos en entornos de cloud híbrido y multisitio, para permitir que los datos se muevan cuando y donde sea necesario.
- Compactación de datos inline y deduplicación expandida. La compactación de datos reduce el espacio perdido dentro de los bloques de almacenamiento, mientras que la deduplicación aumenta la capacidad efectiva de forma significativa. Esto es aplicable a los datos almacenados localmente y a los datos organizados en niveles en el cloud.
- Calidad de servicio (AQoS) mínima, máxima y adaptativa. Los controles granulares de calidad de servicio (QoS) ayudan a mantener los niveles de rendimiento para aplicaciones críticas en entornos altamente compartidos.
- Los FlexGroup de NetApp permiten distribuir datos a través de todos los nodos del clúster de almacenamiento, lo que proporciona una capacidad masiva y un mayor rendimiento para conjuntos de datos extremadamente grandes.

- FabricPool de NetApp. Proporciona la organización automática en niveles de datos fríos en opciones de almacenamiento en cloud privado como Amazon Web Services (AWS), Azure y la solución de almacenamiento StorageGRID de NetApp. Para obtener más información sobre FabricPool, consulte ["TR-4598: Prácticas recomendadas de FabricPool"](#).
- NetApp FlexCache Proporciona funcionalidades de almacenamiento en caché de volúmenes remotos que simplifican la distribución de archivos, reduce la latencia de redes WAN y reduce los costes de ancho de banda WAN. FlexCache permite el desarrollo de productos distribuidos en varios sitios, así como el acceso acelerado a los conjuntos de datos corporativos desde ubicaciones remotas.

## Infraestructura preparada para futuros retos

ONTAP ayuda a satisfacer las exigentes y siempre cambiantes necesidades de su empresa con las siguientes funciones:

- Escalado sencillo y funcionamiento sin interrupciones. ONTAP admite la adición en línea de capacidad a las controladoras existentes y a los clústeres de escalado horizontal. Los clientes pueden empezar a utilizar tecnologías punteras como NVMe y FC 32 GB, sin necesidad de realizar costosas migraciones de datos y sin cortes.
- Conexión de cloud. ONTAP es el software de gestión de almacenamiento con mejor conexión de cloud e incluye opciones de almacenamiento definido por software (ONTAP Select) e instancias nativas del cloud (NetApp Cloud Volumes Service) en todos los clouds públicos.
- Integración con aplicaciones emergentes. ONTAP ofrece servicios de datos de clase empresarial para plataformas y aplicaciones de última generación, como vehículos autónomos, ciudades inteligentes e Industria 4.0, utilizando la misma infraestructura que da soporte a las aplicaciones empresariales existentes.

## Kit de herramientas de operaciones de datos de NetApp

El kit de herramientas DataOps de NetApp es una herramienta basada en Python que simplifica la gestión de espacios de trabajo de desarrollo/formación y servidores de inferencia respaldados por un almacenamiento de NetApp de escalado horizontal y de alto rendimiento. DataOps Toolkit puede funcionar como una utilidad independiente y es incluso más eficaz en entornos de Kubernetes si se utiliza Astra Trident de NetApp para automatizar las operaciones de almacenamiento. Estas son algunas funcionalidades clave:

- Aprovisiona con rapidez nuevos espacios de trabajo de JupyterLab de alta capacidad respaldados por el almacenamiento de NetApp de escalado horizontal y de alto rendimiento.
- Aprovisiona rápidamente nuevas instancias del servidor de inferencia de NVIDIA Triton, respaldadas por un almacenamiento empresarial de NetApp.
- Clonación casi instantánea de espacios de trabajo JupyterLab de gran capacidad para permitir la experimentación o la iteración rápida.
- Copias Snapshot casi instantáneas de espacios de trabajo de gran capacidad JupyterLab para backup o trazabilidad/creación de bases de datos.
- Aprovisionamiento, clonado y copias Snapshot casi instantáneos de volúmenes de datos de gran capacidad y alto rendimiento.

## Astra Trident de NetApp

Astra Trident es un orquestador de almacenamiento de código abierto y totalmente compatible para contenedores y distribuciones de Kubernetes, incluido Anthos. Trident funciona con toda la cartera de almacenamiento de NetApp, incluida NetApp ONTAP, y también es compatible con conexiones NFS, NVMe/TCP e iSCSI. Trident acelera el flujo de trabajo de DevOps al permitir que los usuarios finales

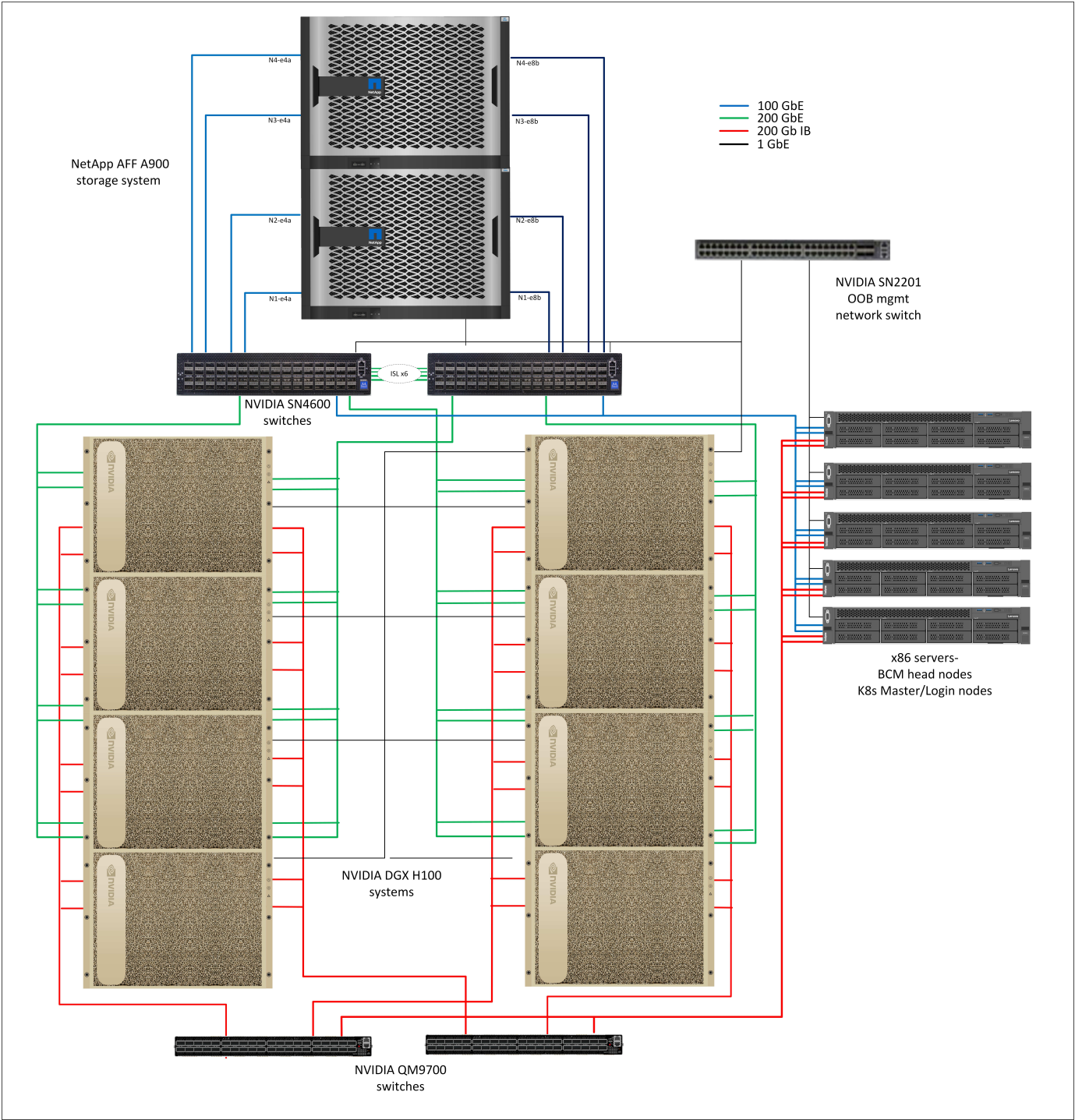
aprovisionen y gestionen el almacenamiento desde sus sistemas de almacenamiento de NetApp sin necesidad de intervención del administrador de almacenamiento.

**NetApp AI Pod con sistemas NVIDIA DGX: Arquitectura de la solución**

**AI Pod de NetApp con sistemas DGX H100**

Esta arquitectura de referencia aprovecha las estructuras separadas para la interconexión en clústeres de computación y el acceso al almacenamiento, con conectividad InfiniBand (IB) de 400GB Gb/s entre los nodos de computación. El siguiente dibujo muestra la topología general de la solución de NetApp AI Pod con sistemas DGX H100.

*Topología de la solución Alpod de NetApp*





## Configuración de red

En esta configuración, la estructura de clústeres de computación utiliza un par de switches IB de QM9700 400GB Gb/s, que se conectan entre sí para lograr una alta disponibilidad. Cada sistema DGX H100 está conectado a los switches mediante ocho conexiones, con puertos números pares conectados a un switch y puertos con número impar conectados al otro switch.

Para el acceso al sistema de almacenamiento, la gestión en banda y el acceso de clientes, se utiliza un par de switches Ethernet SN4600 Gb. Los switches están conectados con enlaces entre switches y se configuran con varias VLAN para aislar los distintos tipos de tráfico. Para implementaciones mayores, la red Ethernet se puede expandir a una configuración hoja-espina añadiendo pares de switches adicionales para los switches espina y hojas adicionales según sea necesario.

Además de la interconexión informática y las redes Ethernet de alta velocidad, todos los dispositivos físicos también están conectados a uno o más switches Ethernet de SN2201 Gb para la gestión fuera de banda. Para obtener más información sobre la conectividad del sistema DGX H100, consulte la ["Documentación de NVIDIA BasePOD"](#).

### Configuración de cliente para el acceso al almacenamiento

Cada sistema DGX H100 está aprovisionado con dos adaptadores ConnectX-7 de doble puerto para el tráfico de gestión y almacenamiento. Para esta solución, ambos puertos de cada tarjeta están conectados al mismo switch. Después, un puerto de cada tarjeta se configura en un vínculo LACP MLAG con un puerto conectado a cada switch, y las VLAN para la gestión en banda, el acceso de clientes y el acceso al almacenamiento en el nivel de usuario se alojan en este vínculo.

El otro puerto de cada tarjeta se utiliza para la conectividad con los sistemas de almacenamiento AFF A900 y se puede utilizar en varias configuraciones según los requisitos de la carga de trabajo. Para configuraciones que utilizan NFS a través de RDMA para admitir el almacenamiento Magnum IO GPUDirect de NVIDIA, los puertos se configuran en un vínculo activo-pasivo, ya que RDMA no es compatible con ningún otro tipo de vínculo. Para las puestas en marcha que no requieren RDMA, las interfaces de almacenamiento también se pueden configurar mediante enlaces LACP para ofrecer una alta disponibilidad y ancho de banda adicional. Con o sin RDMA, los clientes pueden montar el sistema de almacenamiento mediante pNFS v4,1 y Trunking de sesiones para permitir el acceso en paralelo a todos los nodos de almacenamiento del clúster.

### Configuración del sistema de almacenamiento

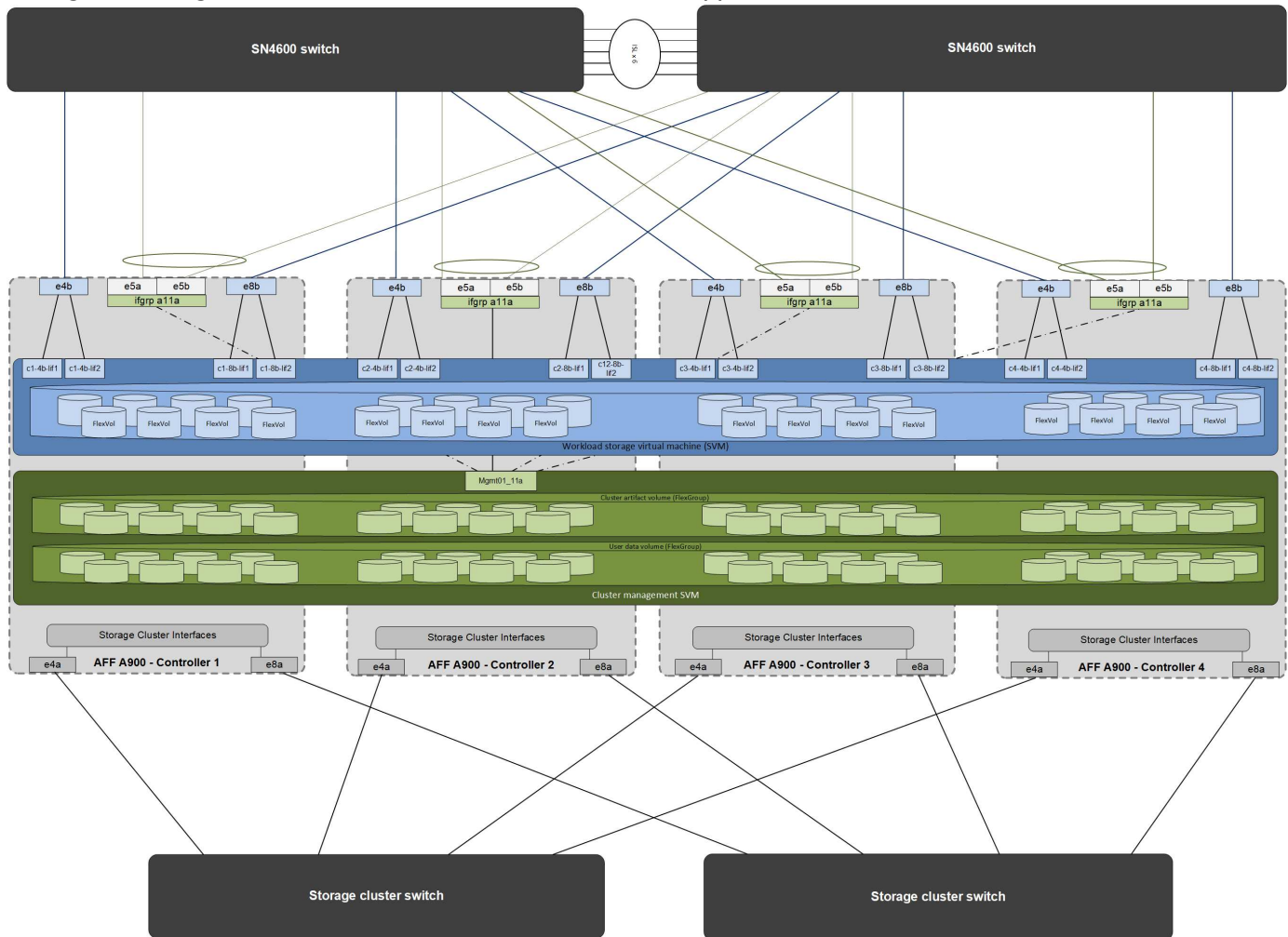
Cada sistema de almacenamiento A900 de AFF está conectado mediante cuatro puertos de 100 GbE desde cada controladora. Dos puertos de cada controladora se utilizan para el acceso a los datos de carga de trabajo desde los sistemas DGX y dos puertos de cada controladora están configurados como un grupo de interfaces LACP para admitir el acceso desde los servidores del plano de gestión para artefactos de gestión de clústeres y directorios iniciales de usuario. Todos los accesos a datos desde el sistema de almacenamiento se realizan mediante NFS, con una máquina virtual de almacenamiento (SVM) dedicada al acceso a las cargas de trabajo de IA y una SVM independiente dedicada a los usos de gestión del clúster.

La SVM de carga de trabajo está configurada con un total de ocho interfaces lógicas (LIF), con dos para cada puerto físico. Esta configuración proporciona el ancho de banda máximo, así como los medios para que cada LIF pueda conmutar por error a otro puerto de la misma controladora, de modo que ambas controladoras permanezcan activas en caso de un fallo de red. Esta configuración también admite NFS sobre RDMA para habilitar el acceso a almacenamiento GPUDirect. La capacidad de almacenamiento se aprovisiona en forma de un único volumen FlexGroup grande que abarca todas las controladoras de almacenamiento del clúster, con 16 volúmenes constituyentes en cada controladora. Puede accederse a esta FlexGroup desde cualquiera de las LIF en la SVM y al utilizar NFSv4,1 con pNFS y conexión de enlaces de sesiones, los clientes establecen conexiones a todas las LIF de la SVM, lo cual permite acceder a los datos locales de cada nodo de almacenamiento en paralelo con el fin de mejorar significativamente el rendimiento. La SVM de carga de

trabajo y cada LIF de datos también están configuradas para el acceso a protocolo RDMA. Para obtener más información sobre la configuración de RDMA para ONTAP, consulte la ["Documentación de ONTAP"](#).

La SVM de gestión solo requiere un solo LIF, que está alojado en los grupos de interfaz de 2 puertos configurados en cada controladora. Otros volúmenes FlexGroup se aprovisionan en la SVM de gestión con el fin de albergar artefactos de gestión del clúster, como imágenes de nodos de clúster, datos históricos de supervisión del sistema y directorios iniciales de usuarios finales. El siguiente dibujo muestra la configuración lógica del sistema de almacenamiento.

*Configuración lógica del clúster de almacenamiento de NetApp A900*



## Servidores del plano de gestión

Esta arquitectura de referencia también incluye cinco servidores basados en CPU para los usos del plano de gestión. Dos de estos sistemas se usan como nodos principales de NVIDIA Base Command Manager para la puesta en marcha y la gestión del clúster. Los otros tres sistemas se utilizan para proporcionar servicios de clúster adicionales, como los nodos maestros de Kubernetes o los nodos de inicio de sesión para las implementaciones que utilizan Slurm para la programación de tareas. Las puestas en marcha que utilizan Kubernetes pueden aprovechar el controlador CSI Astra Trident de NetApp para proporcionar aprovisionamiento automatizado y servicios de datos con almacenamiento persistente para cargas de trabajo de gestión y de IA en el sistema de almacenamiento A900 de AFF.

Cada servidor está conectado físicamente a los switches IB y Ethernet para permitir la puesta en marcha y gestión de clústeres. Además, está configurado con montajes NFS en el sistema de almacenamiento a través de la SVM de gestión para almacenamiento de artefactos de gestión de clústeres, tal como se ha descrito

anteriormente.

## **AIPOD de NetApp con sistemas NVIDIA DGX: Validación de la solución y guía de tamaño**

### **Validación de la solución**

La configuración del almacenamiento de esta solución se validó mediante una serie de cargas de trabajo sintéticas utilizando la herramienta de código abierto FIO. Estas pruebas incluyen patrones de I/O de lectura y escritura destinados a simular la carga de trabajo de almacenamiento generada por los sistemas DGX que realizan trabajos de entrenamiento de aprendizaje profundo. La configuración del almacenamiento se validó utilizando un clúster de servidores CPU de 2 sockets que ejecutaban las cargas de trabajo FIO de forma simultánea para simular un clúster de sistemas DGX. Cada cliente se configuró con la misma configuración de red descrita anteriormente, con la adición de los siguientes detalles.

Para esta validación se han utilizado las siguientes opciones de montaje:

- `vers=4,1` # Habilita pNFS para el acceso paralelo a varios nodos de almacenamiento
- `proto=rdma` # establece el protocolo de transferencia en RDMA en lugar del TCP predeterminado
- `puerto=20049` # Especifique el puerto correcto para el servicio NFS RDMA
- `max_connect=16` # Permite el trunking de sesión NFS para agregar ancho de banda de puerto de almacenamiento
- `write=eager` # mejora el rendimiento de escritura de escrituras en búfer
- `Rsize=262144,wsz=262144` # establece el tamaño de transferencia de I/O en 256K kb

Además, los clientes se configuraron con un valor de NFS `max_session_slots` de 1024. Dado que la solución se probó con NFS a través de RDMA, los puertos de redes de almacenamiento se configuraron con un vínculo activo-pasivo. Para esta validación se han utilizado los siguientes parámetros de enlace:

- `mode=backup activo` # establece el enlace al modo activo / pasivo
- `primario=<interface name>` # las interfaces primarias para todos los clientes se distribuyeron a través de los switches
- `mii-monitor-interval=100` # especifica el intervalo de monitoreo de 100ms
- `fail-over-mac-policy=active` # especifica que la dirección MAC del enlace activo es el MAC del enlace. Esto es necesario para el correcto funcionamiento de RDMA a través de la interfaz vinculada.

El sistema de almacenamiento se configuró como se describe con dos pares de alta disponibilidad A900 (4 controladoras) con dos bandejas de discos NS224 de 24 1,9TB unidades de disco NVMe añadidas a cada par de alta disponibilidad. Tal y como se indica en la sección de arquitectura, la capacidad de almacenamiento de todas las controladoras se combinó mediante el uso de un volumen FlexGroup y la información de todos los clientes se distribuyó entre todas las controladoras del clúster.

### **Directrices de tamaño del sistema de almacenamiento**

NetApp ha completado correctamente la certificación DGX BasePOD y los dos pares de alta disponibilidad A900 como pruebas pueden admitir fácilmente un clúster de ocho sistemas DGX H100. Para puestas en marcha de mayor tamaño con requisitos de rendimiento del almacenamiento superiores, es posible añadir sistemas AFF adicionales al clúster de NetApp ONTAP hasta 12 pares de alta disponibilidad (24 nodos) en un único clúster. Con la tecnología FlexGroup descrita en esta solución, un clúster de 24 nodos puede proporcionar más de 40 PB y hasta 300 Gbps de rendimiento en un solo espacio de nombres. Otros sistemas de almacenamiento de NetApp como AFF A400, A250 y C800 ofrecen un rendimiento menor y opciones de capacidad superior para puestas en marcha de menor tamaño a puntos de coste menores. Como ONTAP 9 admite clústeres de modelo mixto, los clientes pueden comenzar con una huella inicial pequeña e ir aumentando el sistema de almacenamiento a medida que crezcan los requisitos de capacidad y rendimiento. La siguiente tabla muestra una estimación aproximada del número de GPU A100 y H100 admitidas en cada modelo de AFF.

*Guía para ajustar el tamaño del sistema de almacenamiento de NetApp*



		Throughput <sup>2</sup>	Raw capacity (typical / max)	Connectivity	# NVIDIA A100 GPUs supported <sup>3</sup>	# NVIDIA H100 GPUs supported <sup>4</sup>
NetApp® AFF A900	1 HA pair <sup>1</sup>	28GB/s	182TB / 14.7PB	100 GbE	1 - 64	1-32
	12 HA pairs	336GB/s	2.1PB / 176.4PB		768	384
AFF A800	1 HA pair	25GB/s	368TB / 3.6PB	100 GbE	1 - 64	1-32
	12 HA pairs	300GB/s	4.4PB / 43.2PB		768	384
AFF C800	1 HA pair	21GB/s	368TB / 3.6PB	100 GbE	1-48	1-24
	12 HA pairs	252GB/s	4.4PB / 43.2PB		576	288
AFF A400	1 HA pair	11GB/s	182TB / 14.7PB	40/100 GbE	1 - 32	1-16
	12 HA pairs	132GB/s	2.1PB / 176.4PB		384	192
AFF C400	1 HA pair	8GB/s	182TB / 14.7PB	40/100 GbE	1 - 16	1-8
	12 HA pairs	128GB/s	2.1PB / 176.4PB		192	96
AFF A250	1 HA pair	7.4GB/s	91.2TB / 4.4PB	25 GbE 40/100GbE	1 - 16	1-8
	4 HA pairs	29.6GB/s	364.8TB / 17.6PB		64	32
AFF C250	1 HA pair	5 GB/s	91.2TB / 4.4PB	25 GbE 40/100GbE	1-8	1-4
	4 HA pairs	20 GB/s	364.8TB / 17.6PB		32	8

1 – 1 AFF = 1 HA pair = 2 Nodes. 12 HA pairs = 24 nodes  
2 – 100% sequential read

3 – Based on workload testing in NVA-1153  
4 – Based on BasePOD validation test results

## AIPod de NetApp con sistemas NVIDIA DGX: Conclusión e información adicional

### Conclusión

La arquitectura de DGX BasePOD es una plataforma de aprendizaje profundo de próxima generación que requiere funcionalidades de gestión de datos y almacenamiento igualmente avanzadas. Combinando DGX BasePOD con sistemas NetApp AFF, el AIPod de NetApp con la arquitectura de sistemas DGX se puede implementar a casi cualquier escala hasta 48 sistemas DGX H100 en un clúster AFF A900 de 24 nodos. Gracias a la superior integración con el cloud y las capacidades definidas por software de NetApp ONTAP, AFF ofrece una completa gama de canalizaciones de datos que abarcan el perímetro, el núcleo y el cloud y que le permitirán llevar a buen puerto sus proyectos de DL.

### Información adicional

Para obtener más información sobre la información descrita en este documento, consulte los siguientes documentos y/o sitios web:

- Software de gestión de datos ONTAP de NetApp: Biblioteca de información de ONTAP

["https://docs.netapp.com/us-en/ontap-family/"](https://docs.netapp.com/us-en/ontap-family/)

- Sistemas DE ALMACENAMIENTO A900 DE NetApp AFF-

["https://www.netapp.com/data-storage/aff-a-series/aff-a900/"](https://www.netapp.com/data-storage/aff-a-series/aff-a900/)

- Información de RDMA de NetApp ONTAP-

["https://docs.netapp.com/us-en/ontap/nfs-rdma/index.html"](https://docs.netapp.com/us-en/ontap/nfs-rdma/index.html)

- Kit de herramientas de operaciones de datos de NetApp

["https://github.com/NetApp/netapp-dataops-toolkit"](https://github.com/NetApp/netapp-dataops-toolkit)

- Astra Trident de NetApp

["https://docs.netapp.com/us-en/netapp-solutions/containers/rh-os-n\\_overview\\_trident.html"](https://docs.netapp.com/us-en/netapp-solutions/containers/rh-os-n_overview_trident.html)

- Blog sobre almacenamiento GPUDirect de NetApp-

["https://www.netapp.com/blog/ontap-reaches-171-gpudirect-storage/"](https://www.netapp.com/blog/ontap-reaches-171-gpudirect-storage/)

- DGX BasePOD de NVIDIA

["https://www.nvidia.com/en-us/data-center/dgx-basepod/"](https://www.nvidia.com/en-us/data-center/dgx-basepod/)

- SISTEMAS DGX H100 DE NVIDIA

["https://www.nvidia.com/en-us/data-center/dgx-h100/"](https://www.nvidia.com/en-us/data-center/dgx-h100/)

- Conexión a redes NVIDIA

["https://www.nvidia.com/en-us/networking/"](https://www.nvidia.com/en-us/networking/)

- Almacenamiento NVIDIA Magnum IO GPUDirect

["https://docs.nvidia.com/gpudirect-storage"](https://docs.nvidia.com/gpudirect-storage)

- Comando base NVIDIA

["https://www.nvidia.com/en-us/data-center/base-command/"](https://www.nvidia.com/en-us/data-center/base-command/)

- Administrador de comandos base de NVIDIA

["https://www.nvidia.com/en-us/data-center/base-command/manager"](https://www.nvidia.com/en-us/data-center/base-command/manager)

- IA Enterprise de NVIDIA

["https://www.nvidia.com/en-us/data-center/products/ai-enterprise/"](https://www.nvidia.com/en-us/data-center/products/ai-enterprise/)

## Reconocimientos

Este documento es el trabajo de los equipos de Ingeniería de Soluciones NetApp y ONTAP: David Arnette, Olga Kornievskaia, Dustin Fischer, Srikanth Kaligotla, Mohit Kumar y Rajeev Badrinath. Los autores también quieren agradecer a NVIDIA y al equipo de ingeniería de NVIDIA DGX BasePOD su apoyo continuo.

## Diseño de la arquitectura NVA-1151: Guía de diseño de sistemas ONTAP AI de NetApp con NVIDIA DGX A100

David Arnette y Sung-han Lin, NetApp

EL DISEÑO NVA-1151 describe una arquitectura verificada de NetApp para cargas de trabajo de aprendizaje automático y inteligencia artificial que utilizan sistemas de almacenamiento AFF A800 de NetApp, sistemas NVIDIA DGX A100 y switches de red NVIDIA Mellanox. También se incluyen los resultados de las pruebas de rendimiento de la arquitectura tal y como se han implementado.

## **NVA-1151-PUESTA en MARCHA: ONTAP AI de NetApp con sistemas NVIDIA DGX A100**

David Arnette, NetApp

NVA-1151-DEPLOY incluye instrucciones de puesta en marcha del sistema de almacenamiento para cargas de trabajo de arquitectura verificada de NetApp (NVA) para aprendizaje automático (ML) y inteligencia artificial (AI) con sistemas de almacenamiento AFF A800 de NetApp, sistemas NVIDIA DGX A100 y switches de red NVIDIA Mellanox. También incluye instrucciones para ejecutar pruebas de prueba de rendimiento de validación una vez completada la implementación.

["NVA-1151-PUESTA en MARCHA: ONTAP AI de NetApp con sistemas NVIDIA DGX A100"](#)

## **DISEÑO de Mellanox-1153: ONTAP AI de NetApp con sistemas NVIDIA DGX A100 y switches NVA Spectrum Ethernet**

David Arnette y Sung-han Lin, NetApp

EL DISEÑO NVA-1153 describe una arquitectura verificada de NetApp para cargas de trabajo de aprendizaje automático (ML) y de inteligencia artificial (IA) con los sistemas de almacenamiento AFF A800 de NetApp, los sistemas NVIDIA DGX A100 y los switches Ethernet de 200 GB de NVIDIA Mellanox Spectrum SN3700V. Este diseño presenta RDMA over Converged Ethernet (roce) para la estructura de interconexión de clústeres informáticos, que ofrece a los clientes una arquitectura completamente basada en ethernet para cargas de trabajo de alto rendimiento. Este documento también incluye los resultados de las pruebas de rendimiento para la arquitectura tal y como se ha implementado.

["DISEÑO de Mellanox-1153: ONTAP AI de NetApp con sistemas NVIDIA DGX A100 y switches NVA Spectrum Ethernet"](#)

## **Puesta EN MARCHA de Mellanox-1153: ONTAP AI de NetApp con sistemas NVIDIA DGX A100 y switches NVA Spectrum Ethernet**

David Arnette, NetApp

La puesta EN MARCHA de NVA-1153 incluye instrucciones para la puesta en marcha del sistema de almacenamiento para cargas de trabajo de NetApp Verified Architecture para el aprendizaje automático (ML) y de inteligencia artificial (IA) mediante los sistemas de almacenamiento AFF A800 de NetApp, los sistemas NVIDIA DGX A100 y los switches Ethernet de 200 GB de NVIDIA Mellanox Spectrum SN3700V. También incluye instrucciones para ejecutar pruebas de prueba de rendimiento de validación una vez completada la implementación.

["Puesta EN MARCHA de Mellanox-1153: ONTAP AI de NetApp con sistemas NVIDIA DGX A100 y switches NVA Spectrum Ethernet"](#)

## EF-Series AI de NetApp con NVIDIA

Descripción general de las soluciones de infraestructura convergente de EF-Series de NetApp y NVIDIA.

### IA de EF-Series con sistemas NVIDIA DGX A100 y BeeGFS

- ["Guía de diseño"](#)
- ["Guía de puesta en marcha"](#)
- ["Guía de implementación de BeeGFS"](#)

### DISEÑO NVA-1156: EF-Series AI de NetApp con sistemas NVIDIA DGX A100 y BeeGFS

Abdel Sadek, Tim Chau, Joe McCormick y David Arnette, NetApp

El DISEÑO de NVA-1156 describe una arquitectura verificada de NetApp para cargas de trabajo de aprendizaje automático (ML) y de inteligencia artificial (IA) mediante los sistemas de almacenamiento EF600 de NetApp, el sistema de archivos en paralelo BeeGFS, los sistemas NVIDIA DGX A100 y los switches IB de NVIDIA Quantum QM8700 de 200 Gbps. Este diseño incluye InfiniBand (IB) de 200 Gbps para la estructura de interconexión de clústeres de almacenamiento e informáticos con el fin de proporcionar a los clientes una arquitectura completamente basada en IB para las cargas de trabajo de alto rendimiento. Este documento también incluye los resultados de las pruebas de rendimiento para la arquitectura tal y como se ha implementado.

["DISEÑO NVA-1156: EF-Series AI de NetApp con sistemas NVIDIA DGX A100 y BeeGFS"](#)

### NVA-1156-PUESTA en MARCHA: EF-Series AI de NetApp con sistemas NVIDIA DGX A100 y BeeGFS

Abdel Sadek, Tim Chau, Joe McCormick y David Arnette, NetApp

En este documento se describe una arquitectura verificada de NetApp para cargas de trabajo de aprendizaje automático (ML) e inteligencia artificial (IA) mediante los sistemas de almacenamiento EF600 NVMe de NetApp, el sistema de archivos en paralelo ThinkParQ BeeGFS, los sistemas NVIDIA DGX A100 y los switches Mellanox Quantum QM8700 InfiniBand (IB) de 200 Gbps. También se incluyen instrucciones para ejecutar pruebas de referencia de validación una vez completada la implementación.

["NVA-1156-PUESTA en MARCHA: EF-Series AI de NetApp con sistemas NVIDIA DGX A100 y BeeGFS"](#)

### TR-4859: Puesta en marcha de escala de espectro de IBM con almacenamiento E-Series de NetApp: Instalación y validación

Chris Seirer, NetApp

TR-4859 describe el proceso de implementación de una solución de sistema de archivos en paralelo completa basada en la pila de software Spectrum Scale de IBM. TR-4859 está diseñado para proporcionar detalles sobre cómo instalar Spectrum Scale, validar la

infraestructura y gestionar la configuración.

["TR-4859: Puesta en marcha de escala de espectro de IBM con almacenamiento E-Series de NetApp: Instalación y validación"](#)

## **TR-4815: AFF A800 de NetApp y Fujitsu Server PRIMERGY GX2570 M5 para cargas de trabajo de entrenamiento de modelos AI y ML**

David Arnette, NetApp Takashi Oishi, Fujitsu

Esta solución se centra en una arquitectura de escalado horizontal para poner en marcha sistemas de inteligencia artificial con sistemas de almacenamiento de NetApp y servidores Fujitsu. La solución se validó con las pruebas de rendimiento del modelo MLperf v0.6 mediante servidores Fujitsu GX2570 y un sistema de almacenamiento AFF A800 de NetApp.

["TR-4815: AFF A800 de NetApp y Fujitsu Server PRIMERGY GX2570 M5 para cargas de trabajo de entrenamiento de modelos AI y ML"](#)

## **Canalizaciones de datos, lagos de datos y gestión**

### **AWS FSx para NetApp ONTAP (FSxN) para MLOps (en inglés)**

#### **Autor(es):**

Jian Jian (KEN), científico sénior de datos y aplicado, NetApp

En esta sección se profundiza en la aplicación práctica del desarrollo de la infraestructura de IA, ofreciendo un tutorial integral de la construcción de una canalización de MLOps con FSxN. Compuesto por tres ejemplos completos, te guía para satisfacer tus necesidades de MLOps a través de esta potente plataforma de gestión de datos.

Estos artículos se centran en:

1. ["1 parte: Integración de AWS FSx para NetApp ONTAP \(FSxN\) como bloque de S3 privado en AWS SageMaker"](#)
2. ["Parte 2: Aprovechamiento de AWS FSx para NetApp ONTAP \(FSxN\) como fuente de datos para el entrenamiento de modelos en SageMaker"](#)
3. ["Parte 3: Creación de Una canalización simplificada de MLOps \(CI/CT/CD\)"](#)

Al final de esta sección, habrá obtenido una sólida comprensión de cómo usar FSxN para agilizar los procesos de MLOps.

#### **1 parte: Integración de AWS FSx para NetApp ONTAP (FSxN) como bloque de S3 privado en AWS SageMaker**

#### **Autor(es):**

Jian Jian (KEN), científico sénior de datos y aplicado, NetApp

## Introducción

Utilizando SageMaker como ejemplo, esta página proporciona una guía para configurar FSxN como un bucket privado de S3.

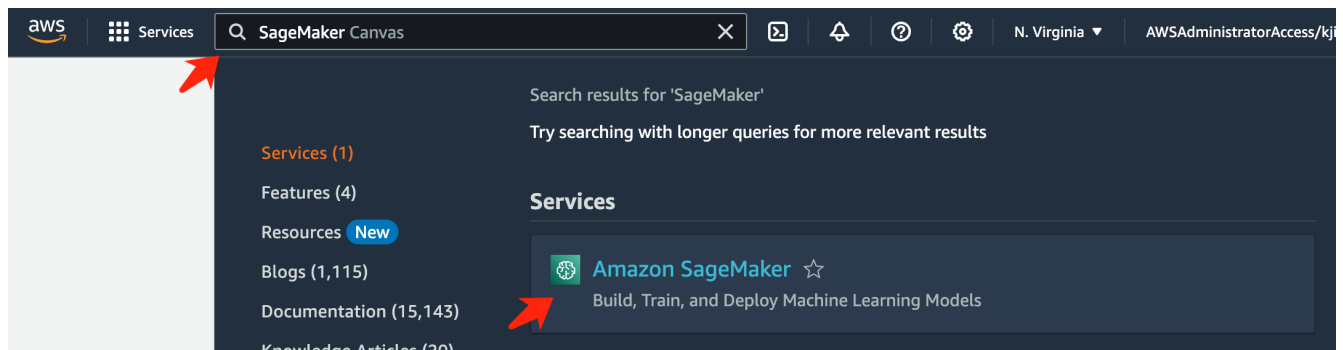
Para obtener más información sobre FSxN, por favor eche un vistazo a esta presentación ("[Enlace de vídeo](#)")

## Guía del usuario

### Creación de servidores

#### Cree una instancia de SageMaker Notebook

1. Abra la consola de AWS. En el panel de búsqueda, busque en SageMaker y haga clic en el servicio **Amazon SageMaker**.



2. Abra las instancias del bloc de notas \* en la pestaña del bloc de notas, haga clic en el botón naranja \* Crear instancia del bloc de notas \*.

The screenshot shows the Amazon SageMaker console. The left sidebar contains the navigation menu with the following items: Getting started, Studio, Studio Lab, Canvas, RStudio, TensorBoard, Profiler, Admin configurations, SageMaker dashboard, Search, JumpStart, Governance, Ground Truth, Notebook (expanded), Notebook instances (highlighted with a red arrow), Git repositories, and Processing. The main content area is titled 'Amazon SageMaker > Notebook instances'. It features a 'Notebook instances' header with an 'Info' link, a refresh button, an 'Actions' dropdown, and a 'Create notebook instance' button. Below this is a search bar labeled 'Search notebook instances' with a red arrow pointing to it. A table with columns 'Name', 'Instance', 'Creation time', 'Last updated', 'Status', 'Lifecycle config', and 'Actions' is shown, with a message 'There are currently no resources.' below it.

- En la página de creación,  
Introduzca el nombre de la instancia de **Notebook**  
Expanda el panel **Network**  
Deje otras entradas predeterminadas y seleccione un **VPC**, **Subnet** y **Grupo(s) de seguridad**. (Este **VPC** y **Subnet** se utilizarán para crear el sistema de archivos FSxN más adelante)  
Haga clic en el botón naranja **Crear instancia de bloc de notas** en la parte inferior derecha.

## Create notebook instance

Amazon SageMaker provides pre-built fully managed notebook instances that run Jupyter notebooks. The notebook instances include example code for common model training and hosting exercises. [Learn more](#)

### Notebook instance settings

Notebook instance name

fsxn-demo

Maximum of 63 alphanumeric characters. Can include hyphens (-), but not spaces. Must be unique within your account in an AWS Region.

Notebook instance type

ml.t3.medium

Elastic Inference [Learn more](#)

none

Platform identifier [Learn more](#)

Amazon Linux 2, Jupyter Lab 3

► Additional configuration

### Permissions and encryption

IAM role

Notebook instances require permissions to call other services including SageMaker and S3. Choose a role or let us create a role with the [AmazonSageMakerFullAccess](#) IAM policy attached.

AmazonSageMakerServiceCatalogProductsUseRole

Create role using the role creation wizard

Root access - optional

- ☒ Enable - Give users root access to the notebook
- ☐ Disable - Don't give users root access to the notebook  
Lifecycle configurations always have root access

Encryption key - optional

Encrypt your notebook data. Choose an existing KMS key or enter a key's ARN.

No Custom Encryption

### ▼ Network - optional

VPC - optional

Default vpc-0df3956ab1fca2ec9 (172.31.0.0/16)

Subnet

Choose a subnet in an availability zone supported by Amazon SageMaker.

subnet-00060df0d0f562672 (172.31.16.0/20) | us-east-1a

Security group(s)

sg-0a39b3985770e9256 (default) X

Direct internet access

- ☒ Enable — Access the internet directly through Amazon SageMaker
- ☐ Disable — Access the internet through a VPC  
To train or host models from a notebook, you need internet access. To enable internet access, make sure that your VPC has a NAT gateway and your security group allows outbound connections. [Learn more](#)

► Git repositories- optional

► Tags - optional

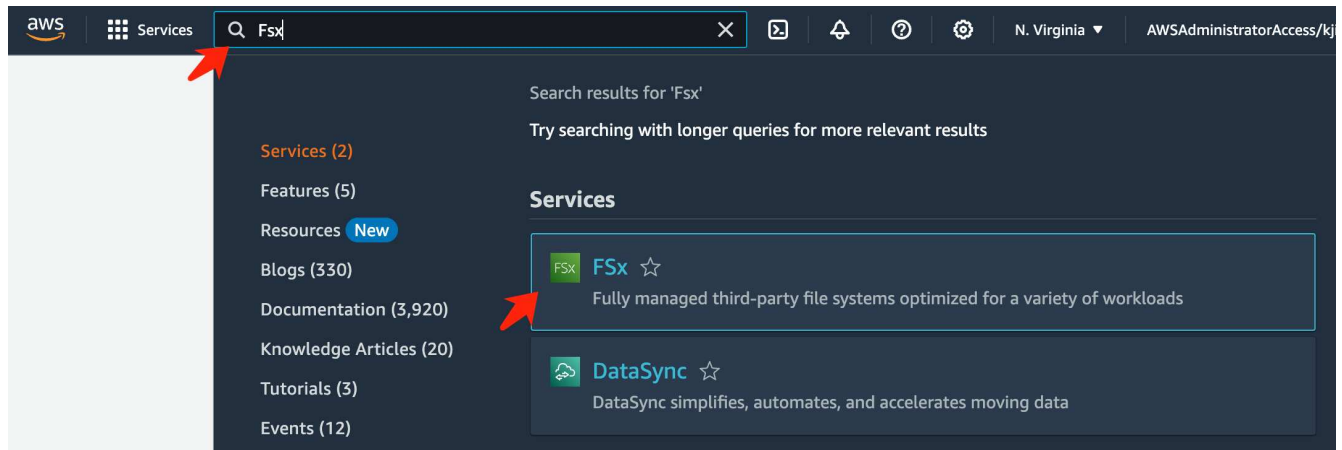
Cancel

Create notebook instance

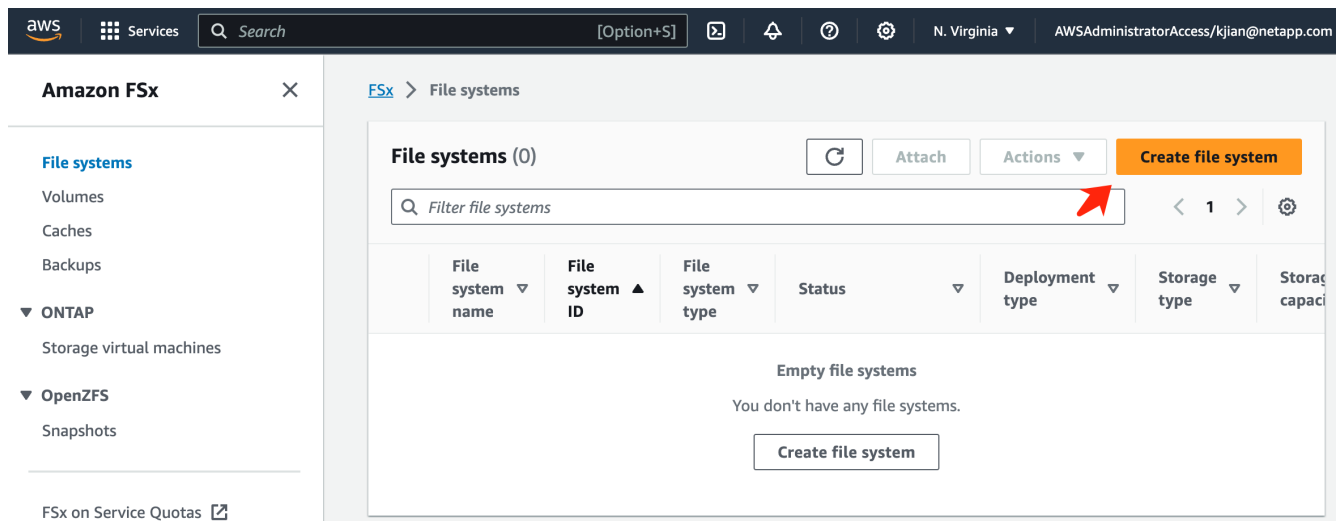


## Crear un sistema de archivos FSxN

1. Abra la consola de AWS. En el panel de búsqueda, busca FSX y haz clic en el servicio **FSX**.



2. Haga clic en **Crear sistema de archivos**.



3. Seleccione la primera tarjeta **FSx for NetApp ONTAP** y haga clic en **Siguiente**.

aws Services Search [Option+S] N. Virginia AWSAdministratorAccess/kjian@netapp

FSx > File systems > Create file system

Step 1  
**Select file system type**

Step 2  
Specify file system details

Step 3  
Review and create

### Select file system type

#### File system options

- ☒ Amazon FSx for NetApp ONTAP
- ☐ Amazon FSx for OpenZFS
- ☐ Amazon FSx for Windows File Server
- ☐ Amazon FSx for Lustre

#### Amazon FSx for NetApp ONTAP

Amazon FSx for NetApp ONTAP provides feature-rich, high-performance, and highly-reliable storage built on NetApp's popular ONTAP file system and fully managed by AWS.

- Broadly accessible from Linux, Windows, and macOS compute instances and containers (running on AWS or on-premises) via industry-standard NFS, SMB, and iSCSI protocols.
- Provides ONTAP's popular data management capabilities like Snapshots, SnapMirror (for data replication), FlexClone (for data cloning), and data compression / deduplication.
- Delivers hundreds of thousands of IOPS with consistent sub-millisecond latencies, and up to 3 GB/s of throughput.
- Offers highly-available and highly-durable single-AZ and multi-AZ deployment options, SSD storage with support for cross-region replication, and built-in, fully managed backups.
- Supports dynamic scaling of your file system to fit your storage capacity and throughput needs.
- Automatically tiers infrequently-accessed data to capacity pool storage, a fully elastic storage tier that can scale to petabytes in size and is cost-optimized for infrequently-accessed data.
- Integrates with Microsoft Active Directory (AD) to support Windows-based environments and enterprises.

Cancel Next

4. En la página de configuración de detalles.
- a. Seleccione la opción **Standard create**.

aws Services Search [Option+S] N. Virginia AWSAdministratorAccess/kjian@netapp

FSx > File systems > Create file system

Step 1  
[Select file system type](#)

Step 2  
**Specify file system details**

Step 3  
Review and create

### Specify file system details

#### Creation method

- ☐ Quick create  
Use recommended best-practice configurations. Most configuration options can be changed after the file system is created.
- ☒ Standard create  
You set all of the configuration options, including specifying performance, networking, security, backups, and maintenance.

- b. Introduzca el **Nombre del sistema de archivos** y la **Capacidad de almacenamiento SSD**.

## File system details

File system name - optional [Info](#)

fsxn-demo

Maximum of 256 Unicode letters, whitespace, and numbers, plus + - = . \_ : /

Deployment type [Info](#)

- ☒ Multi-AZ  
☐ Single-AZ

SSD storage capacity [Info](#)

1024

GiB

Minimum 1024 GiB; Maximum 192 TiB.

Provisioned SSD IOPS

Amazon FSx provides 3 IOPS per GiB of storage capacity. You can also provision additional SSD IOPS as needed.

- ☒ Automatic (3 IOPS per GiB of SSD storage)  
☐ User-provisioned

Throughput capacity [Info](#)

The sustained speed at which the file server hosting your file system can serve data. The file server can also burst to higher speeds for periods of time.

- ☒ Recommended throughput capacity  
128 MB/s  
☐ Specify throughput capacity

c. Asegúrese de usar **VPC** y **subnet** igual a la instancia **SageMaker Notebook**.

## Network & security

### Virtual Private Cloud (VPC) [Info](#)

Specify the VPC from which your file system is accessible.

vpc-0df3956ab1fca2ec9 (CIDR: 172.31.0.0/16) ▼

### VPC Security Groups [Info](#)

Specify VPC Security Groups to associate with your file system's network interfaces.

Choose VPC security group(s) ▼

sg-0a39b3985770e9256 (default) ✕

### Preferred subnet [Info](#)

Specify the preferred subnet for your file system.

subnet-00060df0d0f562672 (us-east-1a | use1-az4) ▼

### Standby subnet

subnet-02b029f24d03a4af2 (us-east-1b | use1-az6) ▼

### VPC route tables [Info](#)

Specify the VPC route tables to associate with your file system.

☒ VPC's main route table

☐ Select one or more VPC route tables

### Endpoint IP address range [Info](#)

Specify the IP address range in which the endpoints to access your file system will be created

☒ Unallocated IP address range from your VPC

Simplest option for access from other AWS services or peered / on-premises networks

☐ Floating IP address range outside your VPC

☐ Enter an IP address range

- d. Introduzca el nombre de la máquina virtual **Storage** y **especifique una contraseña** para su SVM (máquina virtual de almacenamiento).

### Default storage virtual machine configuration

Storage virtual machine name

Info

fsxn-svm-demo

SVM administrative password

Password for this SVM's "vsadmin" user, which you can use to access the ONTAP CLI or REST API. You can provide a password later if you don't provide one now.

☐ Don't specify a password

☒ Specify a password

Password

.....

Confirm password

.....

Volume security style

The security style of the volume determines whether preference is given to NTFS or UNIX ACLs for multi-protocol access. The MIXED mode is not required for multi-protocol access and is only recommended for advanced users.

Unix (Linux)

Active Directory

Joining an Active Directory enables access from Windows and MacOS clients over the SMB protocol.

☒ Do not join an Active Directory

☐ Join an Active Directory

e. Deja otras entradas predeterminadas y haz clic en el botón naranja **Siguiente** en la parte inferior derecha.

► Backup and maintenance - optional

► Tags - optional

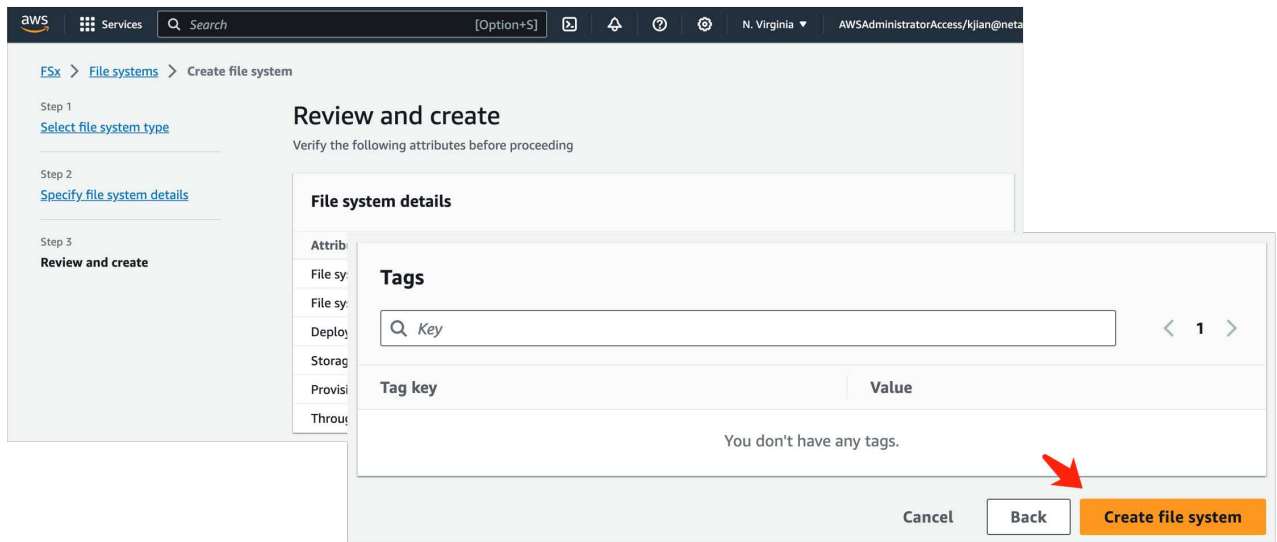
Cancel

Back

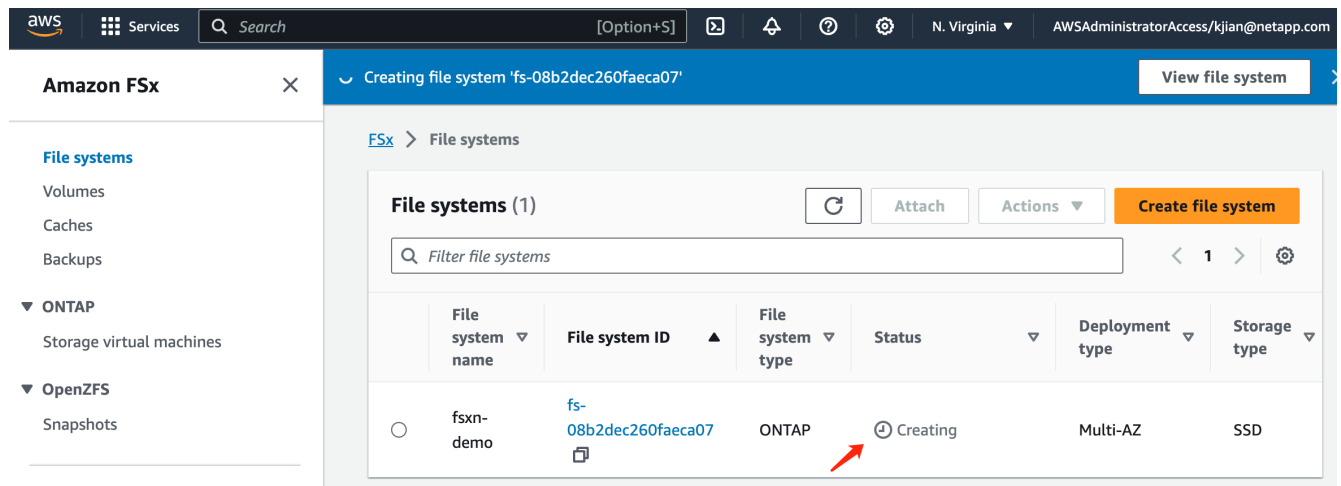
Next

f. Haga clic en el botón naranja **Crear sistema de archivos** en la parte inferior derecha de la página de revisión.

39



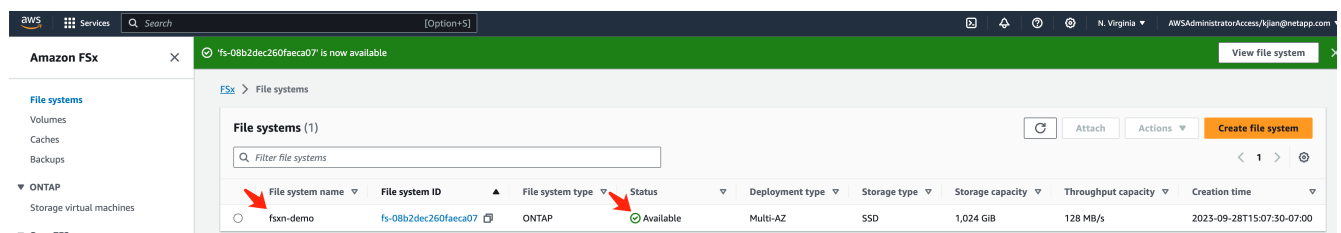
5. Puede tardar unos **20-40 minutos** en activar el sistema de archivos FSX.



## Configuración del servidor

### Configuración de ONTAP

1. Abra el sistema de archivos FSX creado. Por favor, asegúrese de que el estado es **disponible**.



2. Seleccione la pestaña **Administración** y mantenga el **Punto final de administración - dirección IP** y el **Nombre de usuario del administrador de ONTAP**.

**Amazon FSx**

File systems  
Volumes  
Caches  
Backups

▼ **ONTAP**  
Storage virtual machines

▼ **OpenZFS**  
Snapshots

FSx on Service Quotas

**fsxn-demo (fs-08b2dec260faeca07)**

**Summary**

File system ID fs-08b2dec260faeca07	SSD storage capacity 1024 GiB	Availability Zones us-east-1a (Preferred) us-east-1b (Standby)
Lifecycle state Creating	Throughput capacity 128 MB/s	Creation time 2023-09-28T14:41:50-07:00
File system type ONTAP	Provisioned IOPS 3072	
Deployment type Multi-AZ		

**ONTAP administration**

Management endpoint - DNS name management.fs-08b2dec260faeca07.fsx.us-east-1.amazonaws.com	Management endpoint - IP address 172.31.255.250	ONTAP administrator username fsxadmin
Inter-cluster endpoint - DNS name intercluster.fs-08b2dec260faeca07.fsx.us-east-1.amazonaws.com	Inter-cluster endpoint - IP address 172.31.31.157	ONTAP administrator password <a href="#">Update</a>

3. Abra la instancia creada de **SageMaker Notebook** y haga clic en **Abrir JupyterLab**.

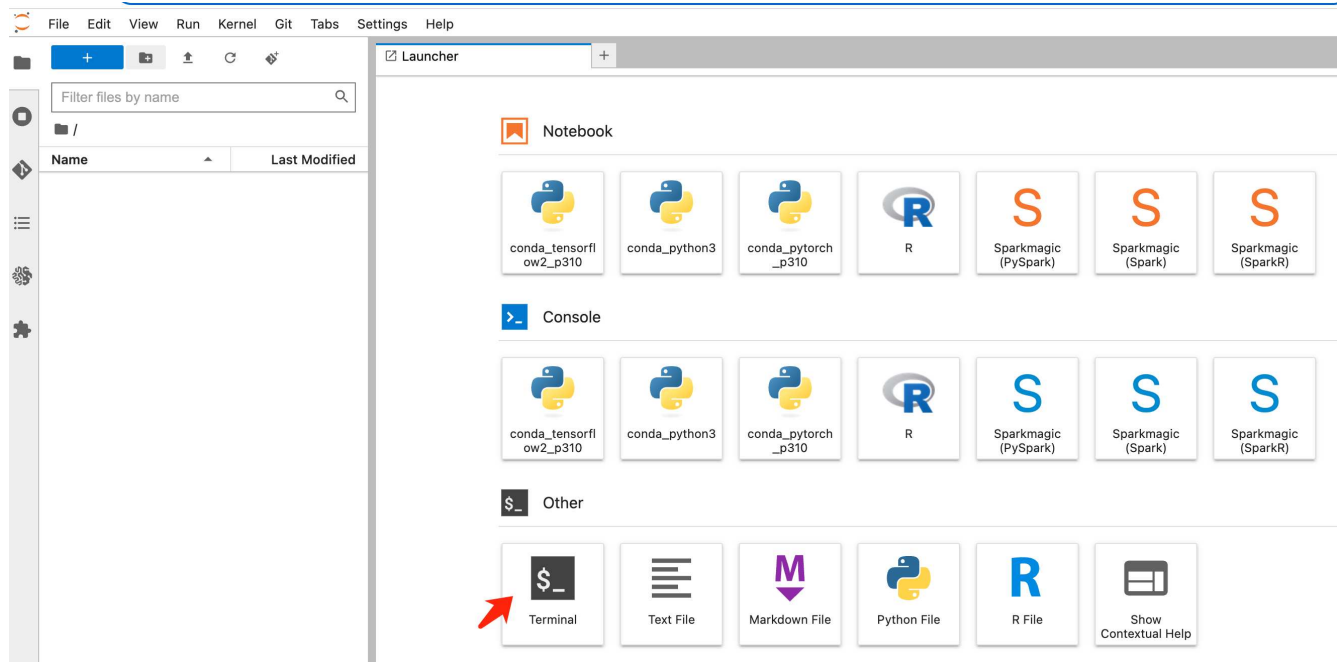
**Amazon SageMaker**

Getting started  
Studio  
Studio Lab  
Canvas  
RStudio  
TensorBoard

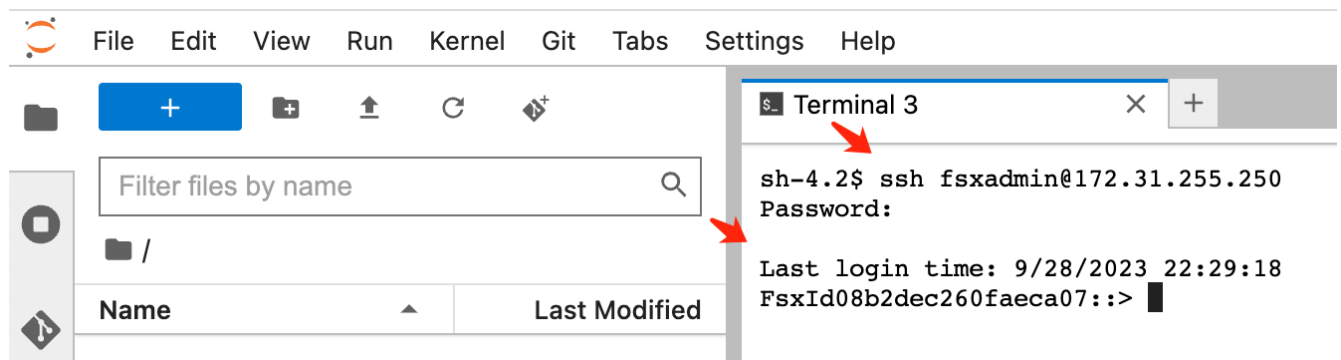
**Notebook instances**

Name	Instance	Creation time	Last updated	Status	Lifecycle config	Actions
fsxn-demo	ml.t3.medium	9/28/2023, 1:47:27 PM	9/28/2023, 1:50:28 PM	InService		<a href="#">Open Jupyter</a>   <a href="#">Open JupyterLab</a>

4. En la página Jupyter Lab, abre un nuevo **Terminal**.



- Introduzca el comando `ssh <nombre de usuario admin>@<IP del servidor de ONTAP>` para iniciar sesión en el sistema de archivos ONTAP FSxN. (El nombre de usuario y la dirección IP se recuperan del paso 2)  
Utilice la contraseña utilizada al crear la **Storage virtual machine**.



- Ejecute los comandos en el siguiente orden.  
Utilizamos **fsxn-ontap** como nombre para el **FSxN private S3 bucket name**.  
Utilice el **nombre de máquina virtual de almacenamiento** para el argumento **-Vserver**.



```

vserver object-store-server create -vserver fsxn-svm-demo -object-store
-server fsx_s3 -is-http-enabled true -is-https-enabled false

vserver object-store-server user create -vserver fsxn-svm-demo -user
s3user

vserver object-store-server group create -name s3group -users s3user
-policies FullAccess

vserver object-store-server bucket create fsxn-ontap -vserver fsxn-svm-
demo -type nas -nas-path /vol1

```



7. Ejecute los siguientes comandos para recuperar la IP de punto final y las credenciales para FSxN private S3.

```

network interface show -vserver fsxn-svm-demo -lif nfs_smb_management_1

set adv

vserver object-store-server user show

```

8. Conserve la IP del extremo y las credenciales para usarlo en el futuro.

Filter files by name

/

Name	Last Modified

```

Terminal 3
sh-4.2$ ssh fsxadmin@172.31.255.250
Password:
Last login time: 9/28/2023 22:32:42
FsxId08b2dec260faeca07:> network interface show -vserver fsxn-svm-demo -lif nfs_smb_management_1

Vserver Name: fsxn-svm-demo
Logical Interface Name: nfs_smb_management_1
Service Policy: default-data-files
Service List: data-core, data-nfs, data-cifs,
              management-ssh, management-https,
              data-s3-server, data-dns-server
(DEPRECATED)-Role: data
Data Protocol: nfs, cifs, s3
Network Address: Fsx IP Address
Netmask: 255.255.255.192
Bits in the Netmask: 26
Is VIP LIF: false
Subnet Name: -
Home Node: FsxId08b2dec260faeca07-01
Home Port: e0e
Current Node: FsxId08b2dec260faeca07-01
Current Port: e0e
Operational Status: up
Extended Status: -
Is Home: true
Administrative Status: up
Failover Policy: system-defined
(DEPRECATED)-Firewall Policy: data
Auto Revert: true
Fully Qualified DNS Zone Name: none
DNS Query Listen Enable: false
Failover Group Name: Fsx
FCP WWPN: -
Address family: ipv4
Comment: -
IPspace of LIF: Default
Is Dynamic DNS Update Enabled?: true
Probe-port for Cloud Load Balancer: -
Broadcast Domain: Fsx
Vserver Type: data
Required RDMA offload protocols: -

FsxId08b2dec260faeca07:> set adv
Warning: These advanced commands are potentially dangerous; use them only when directed to do so by NetApp personnel.
Do you want to continue? {y|n}: y

FsxId08b2dec260faeca07:> vserver object-store-server user show
Vserver  User      ID      Access Key      Secret Key
-----  -
fsxn-svm-demo
  Comment: Root User
fsxn-svm-demo
  s3user      1      AWS Access Key ID AWS Secret Access Key

2 entries were displayed.

FsxId08b2dec260faeca07:>

```

## Configuración del cliente

1. En la instancia de SageMaker Notebook, cree un nuevo cuaderno Jupyter.

File Edit View Run Kernel Git Tabs Settings Help

New

- New Launcher
- Open from Path...
- Open from URL...
- New View for
- New Console for Activity
- Close Tab
- Close and Shutdown
- Close All Tabs
- Save
- Save As...
- Save All
- Reload from Disk
- Revert to Checkpoint
- Rename...
- Download
- Save and Export Notebook As...
- Save Current Workspace As...
- Save Current Workspace
- Print...
- Log Out
- Shut Down

Console

Notebook

conda\_tensorflow2\_p310 conda\_python3 conda\_pytorch\_p310 R Sparkmagic (PySpark) Sparkmagic (Spark) Sparkmagic (SparkR)

Console

conda\_tensorflow2\_p310 conda\_python3 conda\_pytorch\_p310 R Sparkmagic (PySpark) Sparkmagic (Spark) Sparkmagic (SparkR)

Other

Terminal Text File Markdown File Python File R File Show Contextual Help

2. Utilice el siguiente código como solución alternativa para cargar archivos en el cubo privado de FSxN S3. Para obtener un ejemplo de código completo, consulte este cuaderno.

["fsxn\\_demo.ipynb"](#)

```
# Setup configurations
# ----- Manual configurations -----
seed: int = 77                                     # Random
seed
bucket_name: str = 'fsxn-ontap'                     # The bucket
name in ONTAP
aws_access_key_id = '<Your ONTAP bucket key id>'     # Please get
this credential from ONTAP
aws_secret_access_key = '<Your ONTAP bucket access key>' # Please get
this credential from ONTAP
fsx_endpoint_ip: str = '<Your FSxN IP address>'      # Please get
this IP address from FSxN
# ----- Manual configurations -----

# Workaround
## Permission patch
!mkdir -p voll
!sudo mount -t nfs $fsx_endpoint_ip:/voll /home/ec2-user/SageMaker/voll
!sudo chmod 777 /home/ec2-user/SageMaker/voll

## Authentication for FSxN as a Private S3 Bucket
!aws configure set aws_access_key_id $aws_access_key_id
!aws configure set aws_secret_access_key $aws_secret_access_key

## Upload file to the FSxN Private S3 Bucket
%%capture
local_file_path: str = <Your local file path>

!aws s3 cp --endpoint-url http://$fsx_endpoint_ip /home/ec2-user
/SageMaker/$local_file_path s3://$bucket_name/$local_file_path

# Read data from FSxN Private S3 bucket
## Initialize a s3 resource client
import boto3

# Get session info
region_name = boto3.session.Session().region_name

# Initialize FsxN S3 bucket object
# --- Start integrating SageMaker with FSxN ---
# This is the only code change we need to incorporate SageMaker with
FSxN
```

```
s3_client: boto3.client = boto3.resource(
    's3',
    region_name=region_name,
    aws_access_key_id=aws_access_key_id,
    aws_secret_access_key=aws_secret_access_key,
    use_ssl=False,
    endpoint_url=f'http://{fsx_endpoint_ip}',
    config=boto3.session.Config(
        signature_version='s3v4',
        s3={'addressing_style': 'path'}
    )
)

# --- End integrating SageMaker with FSxN ---

## Read file byte content
bucket = s3_client.Bucket(bucket_name)

binary_data = bucket.Object(data.filename).get()['Body']
```

Esto concluye la integración entre FSxN y la instancia de SageMaker.

#### Lista de comprobación de depuración útil

- Asegúrese de que la instancia de SageMaker Notebook y el sistema de archivos FSxN estén en la misma VPC.
- Recuerde ejecutar el comando **set dev** en ONTAP para establecer el nivel de privilegio en **dev**.

#### Preguntas frecuentes (a partir del 27 de septiembre de 2023)

P: ¿Por qué recibo el error **“Se ha producido un error (NotImplemented) al llamar a la operación CreateMultipartUpload: El comando S3 que solicitó no está implementado”** al cargar archivos a FSxN?

R: Como depósito privado de S3, FSxN admite la carga de archivos de hasta 100MB GB. Cuando se utiliza el protocolo S3, los archivos de más de 100MB MB se dividen en 100MB fragmentos y se llama a la función 'CreateMultipartUpload'. Sin embargo, la implementación actual de FSxN PRIVATE S3 no soporta esta función.

P: ¿Por qué recibo el error **“Se ha producido un error (ACCESSDENIED) al llamar a las operaciones PutObject: Acceso denegado”** al cargar archivos a FSxN?

R: Para acceder al bucket S3 privado FSxN desde una instancia de Notebook de SageMaker, cambie las credenciales de AWS a las credenciales FSxN. Sin embargo, otorgar permiso de escritura a la instancia requiere una solución provisional que implique montar el bucket y ejecutar el comando shell 'chmod' para cambiar los permisos.

P: ¿Cómo puedo integrar el cubo FSxN private S3 con otros servicios de SageMaker ML?

R: Desafortunadamente, el SDK de servicios de SageMaker no proporciona una forma de especificar el punto final para el cubo privado de S3. Como resultado, FSxN S3 no es compatible con los servicios de SageMaker tales como Sagemaker Data Wrangler, Sagemaker Clarify, Sagemaker Glue, Sagemaker Athena, Sagemaker

AutoML, Sagemaker AutoML, y otros.

## Parte 2: Aprovechamiento de AWS FSx para NetApp ONTAP (FSxN) como fuente de datos para el entrenamiento de modelos en SageMaker

### Autor(es):

Jian Jian (KEN), científico sénior de datos y aplicado, NetApp

### Introducción

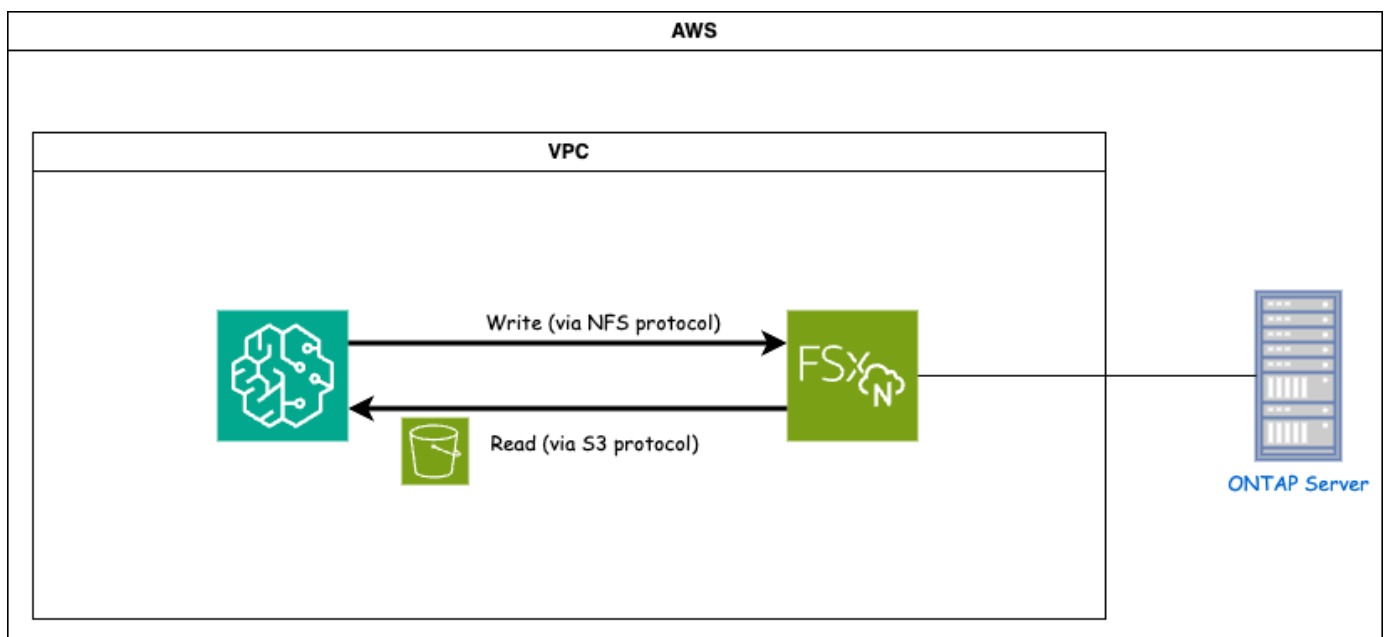
Este tutorial ofrece un ejemplo práctico de un proyecto de clasificación de visión por computadora, que proporciona experiencia práctica en la construcción de modelos ML que utilizan FSxN como fuente de datos dentro del entorno de SageMaker. El proyecto se centra en el uso de PyTorch, un marco de aprendizaje profundo, para clasificar la calidad de los neumáticos en función de las imágenes de los neumáticos. Hace hincapié en el desarrollo de modelos de aprendizaje automático utilizando FSxN como fuente de datos en Amazon SageMaker.

### Qué es FSxN

Amazon FSx para NetApp ONTAP es, de hecho, una solución de almacenamiento totalmente gestionada que ofrece AWS. Aprovecha el sistema de archivos ONTAP de NetApp para ofrecer un almacenamiento fiable y de alto rendimiento. Con su compatibilidad con protocolos como NFS, SMB e iSCSI, permite un acceso fluido desde diferentes instancias de computación y contenedores. El servicio está diseñado para ofrecer un rendimiento excepcional, lo que garantiza operaciones de datos rápidas y eficaces. También ofrece alta disponibilidad y durabilidad, lo que garantiza que sus datos permanezcan accesibles y protegidos. Además, la capacidad de almacenamiento de Amazon FSx para NetApp ONTAP es escalable, lo que te permite ajustarla fácilmente según tus necesidades.

### Requisito previo

### Entorno de red



FSxN (Amazon FSx para NetApp ONTAP) es un servicio de almacenamiento de AWS. Incluye un sistema de archivos que se ejecuta en el sistema NetApp ONTAP y una máquina virtual de sistema gestionado por AWS (SVM) que se conecta a él. En el diagrama proporcionado, el servidor NetApp ONTAP gestionado por AWS se

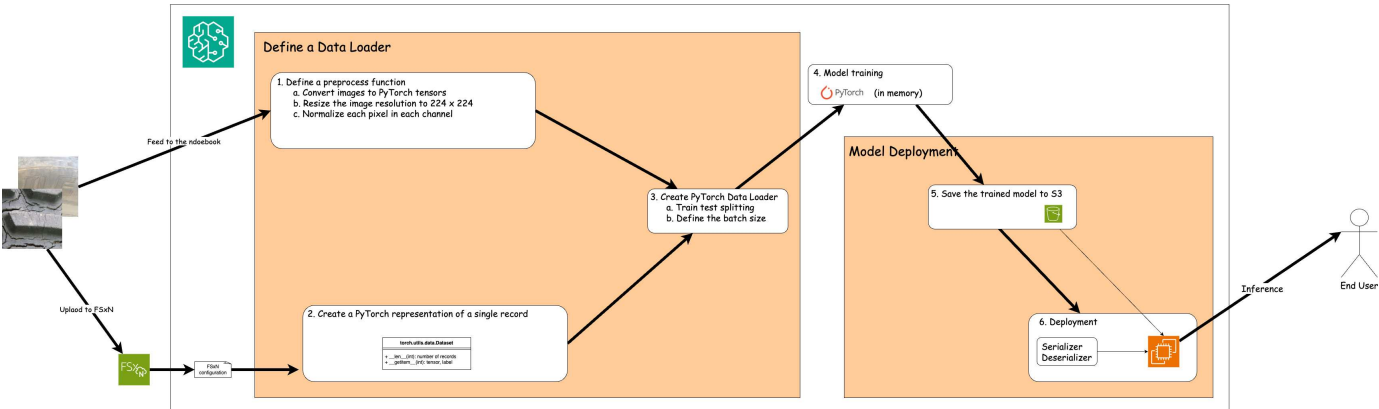
encuentra fuera del VPC. El SVM sirve como intermediario entre SageMaker y el sistema NetApp ONTAP, al recibir solicitudes de operaciones de SageMaker y reenviarlas al almacenamiento subyacente. Para acceder a FSxN, SageMaker debe colocarse dentro de la misma VPC que la implementación FSxN. Esta configuración garantiza la comunicación y el acceso a los datos entre SageMaker y FSxN.

Acceso a los datos

En escenarios del mundo real, los científicos de datos suelen utilizar los datos existentes almacenados en FSxN para crear sus modelos de aprendizaje automático. Sin embargo, a efectos de demostración, dado que el sistema de archivos FSxN está inicialmente vacío después de la creación, es necesario cargar manualmente los datos de entrenamiento. Esto se puede lograr mediante el montaje de FSxN como un volumen a SageMaker. Una vez que el sistema de archivos se ha montado correctamente, puede cargar su conjunto de datos en la ubicación montada, lo que lo hace accesible para el entrenamiento de sus modelos dentro del entorno de SageMaker. Este enfoque le permite aprovechar la capacidad de almacenamiento y las capacidades de FSxN mientras trabaja con SageMaker para el desarrollo y entrenamiento de modelos.

El proceso de lectura de datos implica la configuración de FSxN como un bucket S3 privado. Para obtener más información sobre las instrucciones de configuración detalladas, consulte "[1 parte: Integración de AWS FSx para NetApp ONTAP \(FSxN\) como bloque de S3 privado en AWS SageMaker](#)"

Visión General de la Integración



El flujo de trabajo del uso de datos de entrenamiento en FSxN para construir un modelo de aprendizaje profundo en SageMaker se puede resumir en tres pasos principales: Definición de cargador de datos, entrenamiento de modelos e implementación. En líneas generales, estos pasos forman la base de una canalización de MLOps. Sin embargo, cada paso implica varios subpasos detallados para una implementación integral. Estos subpasos abarcan diversas tareas, como el preprocesamiento de datos, la división de conjuntos de datos, la configuración del modelo, el ajuste de hiperparámetros, la evaluación de modelos, y la puesta en marcha de modelos. Estos pasos garantizan un proceso completo y eficaz para construir e implementar modelos de aprendizaje profundo utilizando datos de entrenamiento de FSxN dentro del entorno de SageMaker.

Integración paso a paso

Cargador de datos

Para entrenar una red de aprendizaje profundo de PyTorch con datos, se crea un cargador de datos para facilitar la alimentación de datos. El cargador de datos no sólo define el tamaño del lote, sino que también determina el procedimiento para leer y preprocesar cada registro del lote. Al configurar el cargador de datos, podemos manejar el procesamiento de datos en lotes, lo que permite el entrenamiento de la red de aprendizaje profundo.

El cargador de datos consta de 3 partes.

### Función de preprocesamiento

```
from torchvision import transforms

preprocess = transforms.Compose([
    transforms.ToTensor(),
    transforms.Resize((224, 224)),
    transforms.Normalize(
        mean=[0.485, 0.456, 0.406],
        std=[0.229, 0.224, 0.225]
    )
])
```

El fragmento de código anterior demuestra la definición de las transformaciones de preprocesamiento de imágenes utilizando el módulo **torchvision.transform**. En este tutorial, se crea el objeto de preproceso para aplicar una serie de transformaciones. En primer lugar, la transformación **ToTensor()** convierte la imagen en una representación tensora. Posteriormente, la transformación **Resize 224.224** cambia el tamaño de la imagen a un tamaño fijo de 224x224 píxeles. Finalmente, la transformación **Normalize()** normaliza los valores del tensor restando la media y dividiendo por la desviación estándar a lo largo de cada canal. Los valores de desviación media y estándar utilizados para la normalización se emplean comúnmente en modelos de redes neuronales pre-entrenados. En general, este código prepara los datos de la imagen para su posterior procesamiento o entrada en un modelo preentrenado convirtiéndolo en un tensor, ajustándolo y normalizando los valores de píxeles.

### La clase de conjunto de datos de PyTorch



```

import torch
from io import BytesIO
from PIL import Image

class FSxNImageDataset(torch.utils.data.Dataset):
    def __init__(self, bucket, prefix='', preprocess=None):
        self.image_keys = [
            s3_obj.key
            for s3_obj in list(bucket.objects.filter(Prefix=prefix).all())
        ]
        self.preprocess = preprocess

    def __len__(self):
        return len(self.image_keys)

    def __getitem__(self, index):
        key = self.image_keys[index]
        response = bucket.Object(key)

        label = 1 if key[13:].startswith('defective') else 0

        image_bytes = response.get()['Body'].read()
        image = Image.open(BytesIO(image_bytes))
        if image.mode == 'L':
            image = image.convert('RGB')

        if self.preprocess is not None:
            image = self.preprocess(image)
        return image, label

```

Esta clase proporciona funcionalidad para obtener el número total de registros en el conjunto de datos y define el método para leer datos para cada registro. Dentro de la función **getitem**, el código utiliza el objeto `bucket` boto3 S3 para recuperar los datos binarios de FSxN. El estilo de código para acceder a los datos de FSxN es similar a la lectura de datos de Amazon S3. La explicación subsiguiente profundiza en el proceso de creación del objeto privado S3 **bucket**.

## FSxN como repositorio S3 privado

```

seed = 77 # Random seed
bucket_name = '<Your ONTAP bucket name>' # The bucket
name in ONTAP
aws_access_key_id = '<Your ONTAP bucket key id>' # Please get
this credential from ONTAP
aws_secret_access_key = '<Your ONTAP bucket access key>' # Please get
this credential from ONTAP
fsx_endpoint_ip = '<Your FSxN IP address>' # Please get
this IP address from FSxN

```

```

import boto3

# Get session info
region_name = boto3.session.Session().region_name

# Initialize FsxN S3 bucket object
# --- Start integrating SageMaker with FSxN ---
# This is the only code change we need to incorporate SageMaker with FSxN
s3_client: boto3.client = boto3.resource(
    's3',
    region_name=region_name,
    aws_access_key_id=aws_access_key_id,
    aws_secret_access_key=aws_secret_access_key,
    use_ssl=False,
    endpoint_url=f'http://{fsx_endpoint_ip}',
    config=boto3.session.Config(
        signature_version='s3v4',
        s3={'addressing_style': 'path'}
    )
)
# s3_client = boto3.resource('s3')
bucket = s3_client.Bucket(bucket_name)
# --- End integrating SageMaker with FSxN ---

```

Para leer datos de FSxN en SageMaker, se crea un manejador que apunta al almacenamiento FSxN mediante el protocolo S3. Esto permite que FSxN se trate como un cubo privado de S3. La configuración del manejador incluye especificar la dirección IP de la SVM FSxN, el nombre del depósito y las credenciales necesarias. Para obtener una explicación completa sobre la obtención de estos elementos de configuración, consulte el documento en ["1 parte: Integración de AWS FSx para NetApp ONTAP \(FSxN\) como bloque de S3 privado en AWS SageMaker"](#).

En el ejemplo mencionado anteriormente, el objeto bucket se utiliza para instanciar el objeto de conjunto de datos PyTorch. El objeto del conjunto de datos se explicará con más detalle en la sección siguiente.

## El cargador de datos de PyTorch

```
from torch.utils.data import DataLoader
torch.manual_seed(seed)

# 1. Hyperparameters
batch_size = 64

# 2. Preparing for the dataset
dataset = FSxNImageDataset(bucket, 'dataset/tyre', preprocess=preprocess)

train, test = torch.utils.data.random_split(dataset, [1500, 356])

data_loader = DataLoader(dataset, batch_size=batch_size, shuffle=True)
```

En el ejemplo proporcionado, se especifica un tamaño de lote de 64, lo que indica que cada lote contendrá 64 registros. Al combinar la clase PyTorch **Dataset**, la función de preprocesamiento y el tamaño de lote de entrenamiento, obtenemos el cargador de datos para el entrenamiento. Este cargador de datos facilita el proceso de iteración por el conjunto de datos en lotes durante la fase de entrenamiento.

## Entrenamiento de modelos

```
from torch import nn

class TyreQualityClassifier(nn.Module):
    def __init__(self):
        super().__init__()
        self.model = nn.Sequential(
            nn.Conv2d(3, 32, (3, 3)),
            nn.ReLU(),
            nn.Conv2d(32, 32, (3, 3)),
            nn.ReLU(),
            nn.Conv2d(32, 64, (3, 3)),
            nn.ReLU(),
            nn.Flatten(),
            nn.Linear(64 * (224 - 6) * (224 - 6), 2)
        )
    def forward(self, x):
        return self.model(x)
```

```

import datetime

num_epochs = 2
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')

model = TyreQualityClassifier()
fn_loss = torch.nn.CrossEntropyLoss()
optimizer = torch.optim.Adam(model.parameters(), lr=1e-3)

model.to(device)
for epoch in range(num_epochs):
    for idx, (X, y) in enumerate(data_loader):
        X = X.to(device)
        y = y.to(device)

        y_hat = model(X)

        loss = fn_loss(y_hat, y)
        optimizer.zero_grad()
        loss.backward()
        optimizer.step()
        current_time = datetime.datetime.now().strftime("%Y-%m-%d
%H:%M:%S")
        print(f"Current Time: {current_time} - Epoch [{epoch+1}]/
{num_epochs}] - Batch [{idx + 1}] - Loss: {loss}", end='\r')

```

Este código implementa un proceso de entrenamiento estándar de PyTorch. Define un modelo de red neuronal llamado **TyreQualityClassifier** usando capas convolucionales y una capa lineal para clasificar la calidad de los neumáticos. El bucle de entrenamiento itera sobre los lotes de datos, calcula la pérdida y actualiza los parámetros del modelo mediante retropropagación y optimización. Además, imprime la hora actual, la época, el lote y la pérdida con fines de monitorización.

## Puesta en marcha de modelos

### Puesta en marcha

```

import io
import os
import tarfile
import sagemaker

# 1. Save the PyTorch model to memory
buffer_model = io.BytesIO()
traced_model = torch.jit.script(model)
torch.jit.save(traced_model, buffer_model)

# 2. Upload to AWS S3
sagemaker_session = sagemaker.Session()
bucket_name_default = sagemaker_session.default_bucket()
model_name = f'tyre_quality_classifier.pth'

# 2.1. Zip PyTorch model into tar.gz file
buffer_zip = io.BytesIO()
with tarfile.open(fileobj=buffer_zip, mode="w:gz") as tar:
    # Add PyTorch pt file
    file_name = os.path.basename(model_name)
    file_name_with_extension = os.path.splitext(file_name)[-1]
    tarinfo = tarfile.TarInfo(file_name_with_extension)
    tarinfo.size = len(buffer_model.getbuffer())
    buffer_model.seek(0)
    tar.addfile(tarinfo, buffer_model)

# 2.2. Upload the tar.gz file to S3 bucket
buffer_zip.seek(0)
boto3.resource('s3') \
    .Bucket(bucket_name_default) \
    .Object(f'pytorch/{model_name}.tar.gz') \
    .put(Body=buffer_zip.getvalue())

```

El código guarda el modelo de PyTorch en **Amazon S3** porque SageMaker requiere que el modelo se almacene en S3 para su implementación. Al subir el modelo a **Amazon S3**, se vuelve accesible para SageMaker, lo que permite la implementación e inferencia en el modelo desplegado.

```

import time
from sagemaker.pytorch import PyTorchModel
from sagemaker.predictor import Predictor
from sagemaker.serializers import IdentitySerializer
from sagemaker.deserializers import JSONDeserializer

class TyreQualitySerializer(IdentitySerializer):

```

```

CONTENT_TYPE = 'application/x-torch'

def serialize(self, data):
    transformed_image = preprocess(data)
    tensor_image = torch.Tensor(transformed_image)

    serialized_data = io.BytesIO()
    torch.save(tensor_image, serialized_data)
    serialized_data.seek(0)
    serialized_data = serialized_data.read()

    return serialized_data

class TyreQualityPredictor(Predictor):
    def __init__(self, endpoint_name, sagemaker_session):
        super().__init__(
            endpoint_name,
            sagemaker_session=sagemaker_session,
            serializer=TyreQualitySerializer(),
            deserializer=JSONDeserializer(),
        )

sagemaker_model = PyTorchModel(
    model_data=f's3://{bucket_name_default}/pytorch/{model_name}.tar.gz',
    role=sagemaker.get_execution_role(),
    framework_version='2.0.1',
    py_version='py310',
    predictor_cls=TyreQualityPredictor,
    entry_point='inference.py',
    source_dir='code',
)

timestamp = int(time.time())
pytorch_endpoint_name = '{}-{}-{}'.format('tyre-quality-classifier', 'pt',
timestamp)
sagemaker_predictor = sagemaker_model.deploy(
    initial_instance_count=1,
    instance_type='ml.p3.2xlarge',
    endpoint_name=pytorch_endpoint_name
)

```

Este código facilita el despliegue de un modelo PyTorch en SageMaker. Define un serializador personalizado, **TyreQualitySerializer**, que preprocesa y serializa los datos de entrada como un tensor PyTorch. La clase **TyreQualityPredictor** es un predictor personalizado que utiliza el serializador definido y un **JSONDeserializer**. El código también crea un objeto **PyTorchModel** para especificar la ubicación S3 del modelo, el rol IAM, la versión del marco y el punto de entrada para la inferencia. El código genera una marca

de tiempo y construye un nombre de punto final basado en el modelo y la marca de tiempo. Por último, el modelo se despliega mediante el método de despliegue, especificando el recuento de instancias, el tipo de instancia y el nombre de punto final generado. Esto permite que el modelo de PyTorch se despliegue y sea accesible para la inferencia en SageMaker.

## Inferencia

```
image_object = list(bucket.objects.filter('dataset/tyre'))[0].get()
image_bytes = image_object['Body'].read()

with Image.open(with Image.open(BytesIO(image_bytes)) as image:
    predicted_classes = sagemaker_predictor.predict(image)

print(predicted_classes)
```

Este es el ejemplo de utilizar el punto final desplegado para llevar a cabo la inferencia.

## Parte 3: Creación de Una canalización simplificada de MLOps (CI/CT/CD)

### Autor(es):

Jian Jian (KEN), científico sénior de datos y aplicado, NetApp

### Introducción

En este tutorial, aprenderá cómo aprovechar varios servicios de AWS para construir una canalización simple de MLOps que abarque la integración continua (CI), el entrenamiento continuo (CT) y la implementación continua (CD). A diferencia de las canalizaciones tradicionales de DevOps, MLOps requiere consideraciones adicionales para completar el ciclo operativo. Al seguir este tutorial, obtendrá información sobre la incorporación de CT en el bucle de MLOps, lo que permite el entrenamiento continuo de sus modelos y la implementación sin problemas para la inferencia. El tutorial le guiará a través del proceso de uso de los servicios de AWS para establecer este pipeline de MLOps de extremo a extremo.

### Manifiesto

Funcionalidad	Nombre	Comentar
Almacenamiento de datos	FSxN de AWS	Consulte " <a href="#">1 parte: Integración de AWS FSx para NetApp ONTAP (FSxN) como bloque de S3 privado en AWS SageMaker</a> ".
IDE de ciencia de datos	SageMaker de AWS	Este tutorial se basa en el cuaderno Jupyter que se presenta en la " <a href="#">Parte 2: Aprovechamiento de AWS FSx para NetApp ONTAP (FSxN) como fuente de datos para el entrenamiento de modelos en SageMaker</a> ".
Función para activar el pipeline de MLOps	Función AWS Lambda	-

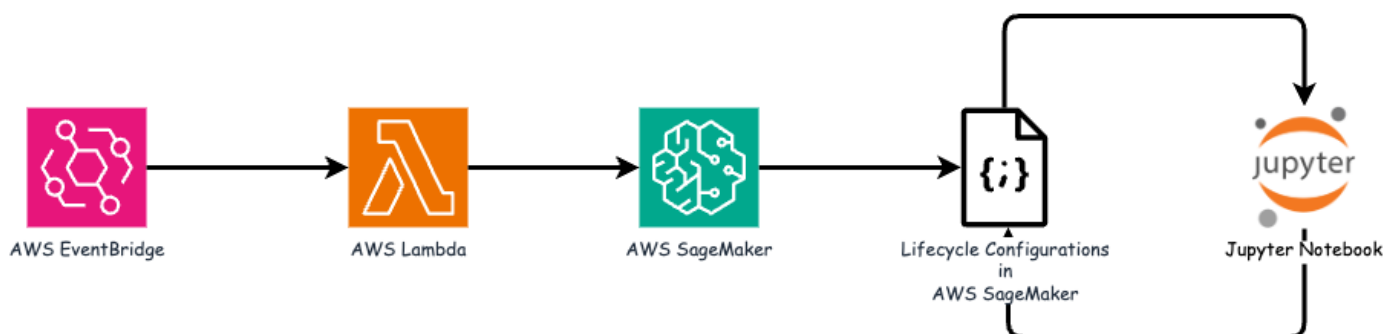


Funcionalidad	Nombre	Comentar
Disparador de trabajo CRON	EventBridge de AWS	-
Marco de aprendizaje profundo	PyTorch	-
SDK de AWS Python	boto3	-
Lenguaje de programación	Python	v3,10

### Requisito previo

- Un sistema de archivos FSxN preconfigurado. Este tutorial utiliza los datos almacenados en FSxN para el proceso de entrenamiento.
- Una instancia **SageMaker Notebook** que está configurada para compartir la misma VPC que el sistema de archivos FSxN mencionado anteriormente.
- Antes de activar la función **AWS Lambda**, asegúrese de que la instancia **SageMaker Notebook** esté en estado **Detenido**.
- El tipo de instancia **ML.g4dn.xlarge** es necesario para aprovechar la aceleración de GPU necesaria para los cálculos de redes neuronales profundas.

### Arquitectura



Esta canalización de MLOps es una implementación práctica que utiliza un trabajo cron para activar una función sin servidor, que a su vez ejecuta un servicio de AWS registrado con una función de devolución de llamada de ciclo de vida. El **AWS EventBridge** actúa como el trabajo cron. Invoca periódicamente una función **AWS Lambda** responsable de reciclar y reimplementar el modelo. Este proceso implica poner en marcha la instancia de **AWS SageMaker Notebook** para realizar las tareas necesarias.

### Configuración paso a paso

#### Configuraciones de ciclo de vida

Para configurar la función de devolución de llamada de ciclo de vida para la instancia de AWS SageMaker Notebook, utilizaría **Configuraciones de ciclo de vida**. Este servicio le permite definir las acciones necesarias que se deben realizar durante el giro de la instancia del bloc de notas. Específicamente, se puede implementar un script de shell dentro de las configuraciones de ciclo de vida \* para cerrar automáticamente la instancia de notebook una vez que se completen los procesos de entrenamiento e implementación. Esta es una configuración necesaria, ya que el coste es uno de los principales factores que hay que tener en cuenta en MLOps.

Es importante tener en cuenta que la configuración de **configuraciones de ciclo de vida** debe configurarse con antelación. Por lo tanto, se recomienda priorizar la configuración de este aspecto antes de continuar con

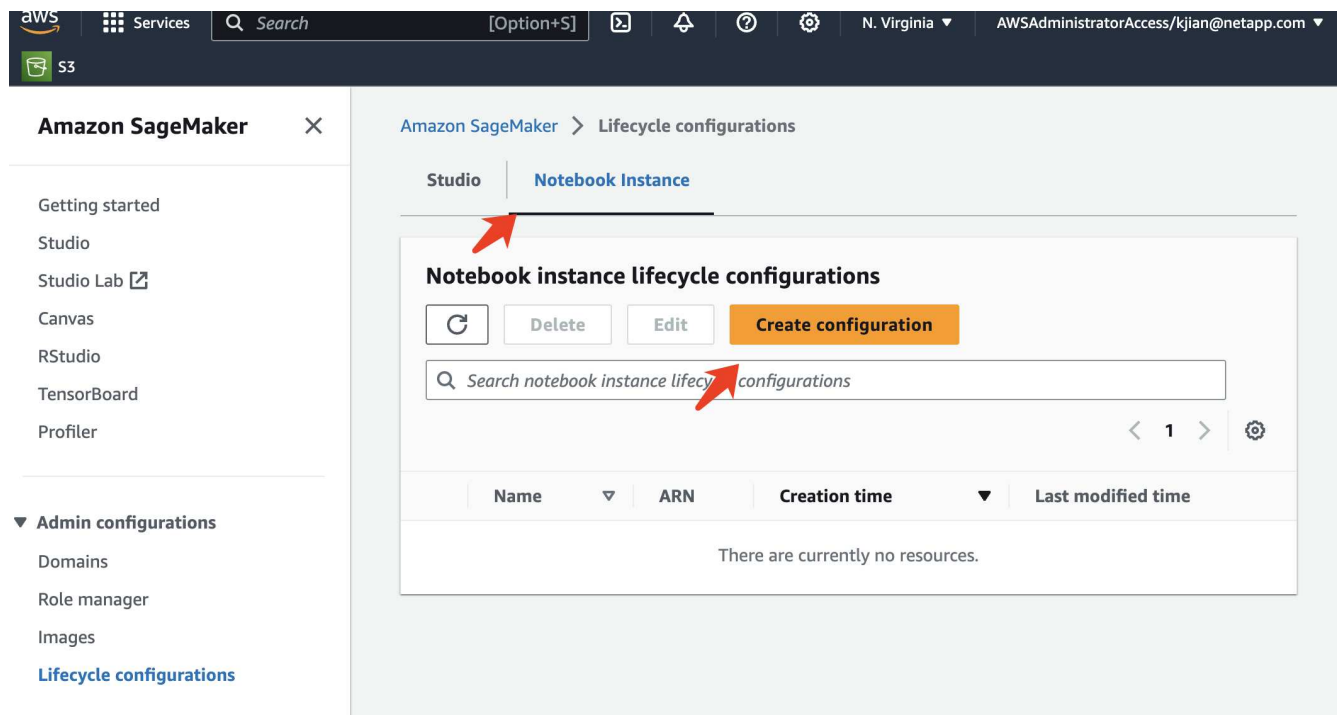
la otra configuración de pipeline de MLOps.

1. Para configurar una configuración de ciclo de vida, abra el panel **Sagemaker** y vaya a **Configuraciones de ciclo de vida** en la sección **Configuraciones de administración**.

The screenshot shows the Amazon SageMaker console interface. The top navigation bar includes the AWS logo, 'Services', and a search bar. Below this, there's a dark blue header with an 'S3' icon. The left sidebar is titled 'Amazon SageMaker' and contains a list of navigation items: 'Getting started', 'Studio', 'Studio Lab', 'Canvas', 'RStudio', 'TensorBoard', 'Profiler', 'Admin configurations' (expanded), 'Domains' (highlighted in blue), 'Role manager', 'Images', 'Lifecycle configurations' (indicated by a red arrow), 'SageMaker dashboard', 'Search', and 'JumpStart'. The main content area is titled 'Amazon SageMaker > Domains' and features a 'Domains' section with an 'Info' link. Below this, there's a 'Domain structure diagram' link and a 'Domains (4)' section with an 'Info' link. A search bar labeled 'Find domain name' is present. A table lists four domains: 'rdsml-east-1', 'rdsml-east-2', 'rdsml-east-3', and 'rdsml-east-4', each with a radio button for selection.

	Name
<input type="radio"/>	rdsml-east-1
<input type="radio"/>	rdsml-east-2
<input type="radio"/>	rdsml-east-3
<input type="radio"/>	rdsml-east-4

2. Seleccione la pestaña **Instancia de bloc de notas** y haga clic en el botón **Crear configuración**




3. Pegue el siguiente código en el área de entrada.


```
#!/bin/bash


set -e
sudo -u ec2-user -i <<'EOF'
# 1. Retraining and redeploying the model
NOTEBOOK_FILE=/home/ec2-
user/SageMaker/tyre_quality_classification_local_training.ipynb
echo "Activating conda env"
source /home/ec2-user/anaconda3/bin/activate pytorch_p310
nohup jupyter nbconvert "$NOTEBOOK_FILE"
--ExecutePreprocessor.kernel_name=python --execute --to notebook &
nbconvert_pid=$!
conda deactivate

# 2. Scheduling a job to shutdown the notebook to save the cost
PYTHON_DIR='/home/ec2-
user/anaconda3/envs/JupyterSystemEnv/bin/python3.10'
echo "Starting the autostop script in cron"
(crontab -l 2>/dev/null; echo "*/5 * * * * bash -c 'if ps -p
$nbconvert_pid > /dev/null; then echo \"Notebook is still running.\" >>
/var/log/jupyter.log; else echo \"Notebook execution completed.\" >>
/var/log/jupyter.log; $PYTHON_DIR -c \"import boto3;boto3.client(
\"sagemaker\").stop_notebook_instance(NotebookInstanceName=get_notebook_
name())\" >> /var/log/jupyter.log; fi'") | crontab -
EOF
```

4. Este script ejecuta el Jupyter Notebook, que se encarga del reciclaje y el redespliegue del modelo para la inferencia. Una vez finalizada la ejecución, el bloc de notas se apagará automáticamente en 5 minutos. Para obtener más información sobre la declaración del problema y la implementación del código, consulte ["Parte 2: Aprovechamiento de AWS FSx para NetApp ONTAP \(FSxN\) como fuente de datos para el entrenamiento de modelos en SageMaker"](#).

 Services  [Option+S]





Amazon SageMaker > Lifecycle configurations > Create lifecycle configuration

## Create lifecycle configuration

### Configuration setting

Name

Alphanumeric characters and "-", no spaces. Maximum 63 characters.

### Scripts

Start notebook


Create notebook

This script will be run each time an associated notebook instance is started, including during initial creation. If the associated notebook instance is already started, it will be run the next time it is stopped and started. [a curated list of sample scripts](#)

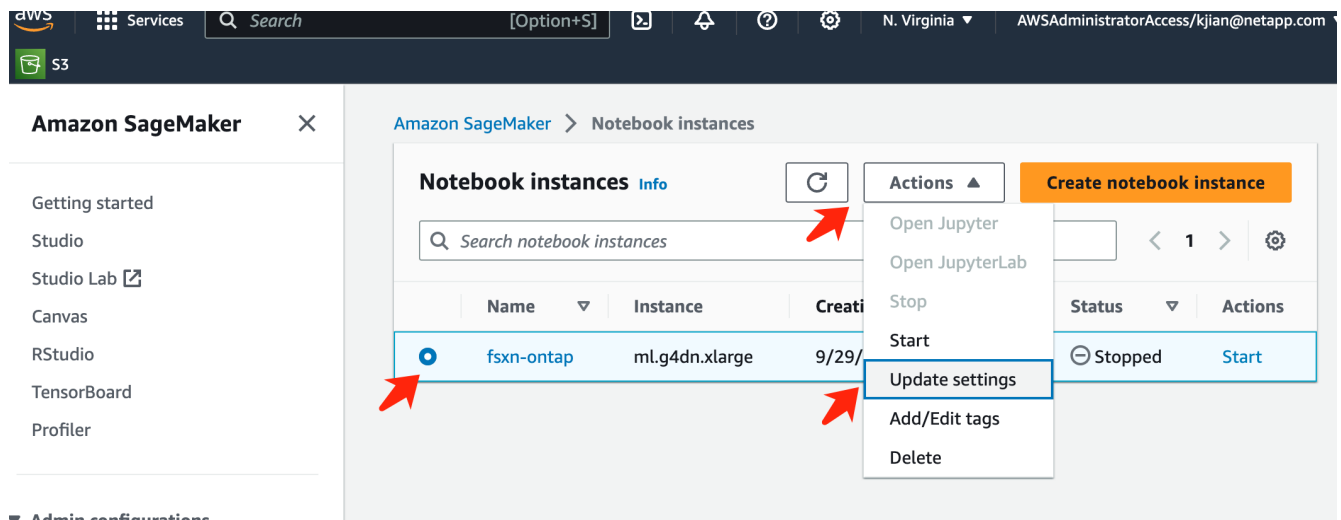
```
1 #!/bin/bash
2
3 set -e
4 sudo -u ec2-user -i <<'EOF'
5 # 1. Retraining and redeploying the model
6 NOTEBOOK_FILE=/home/ec2-user/SageMaker/tyre_quality_classification_local_training.ipynb
7 echo "Activating conda env"
8 source /home/ec2-user/anaconda3/bin/activate pytorch_p310
9 nohup jupyter nbconvert "$NOTEBOOK_FILE" --ExecutePreprocessor.kernel_name=python --execute --to n
10 nbconvert_pid=$!
11 conda deactivate
12
13 # 2. Scheduling a job to shutdown the notebook to save the cost
14 PYTHON_DIR='/home/ec2-user/anaconda3/envs/JupyterSystemEnv/bin/python3.10'
15 echo "Starting the autostop script in cron"
16 (crontab -l 2>/dev/null; echo "*/5 * * * * bash -c 'if ps -p $nbconvert_pid > /dev/null; then echo
17 EOF
```

Cancel

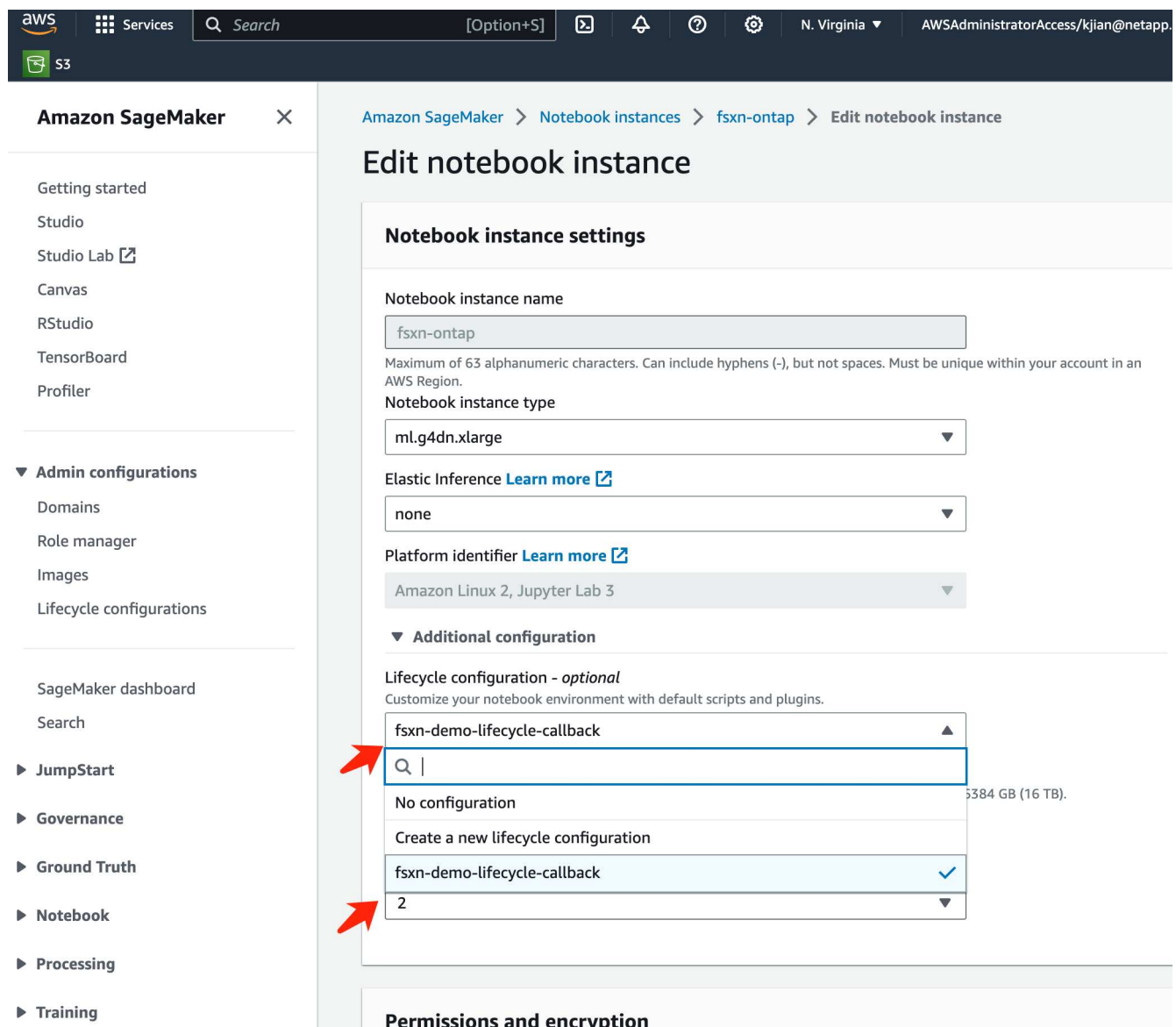
Create configuration

 CloudShell [Feedback](#)

5. Después de la creación, navegue a Instancias de bloc de notas, seleccione la instancia de destino y haga clic en **Actualizar configuración** en el menú desplegable Acciones.



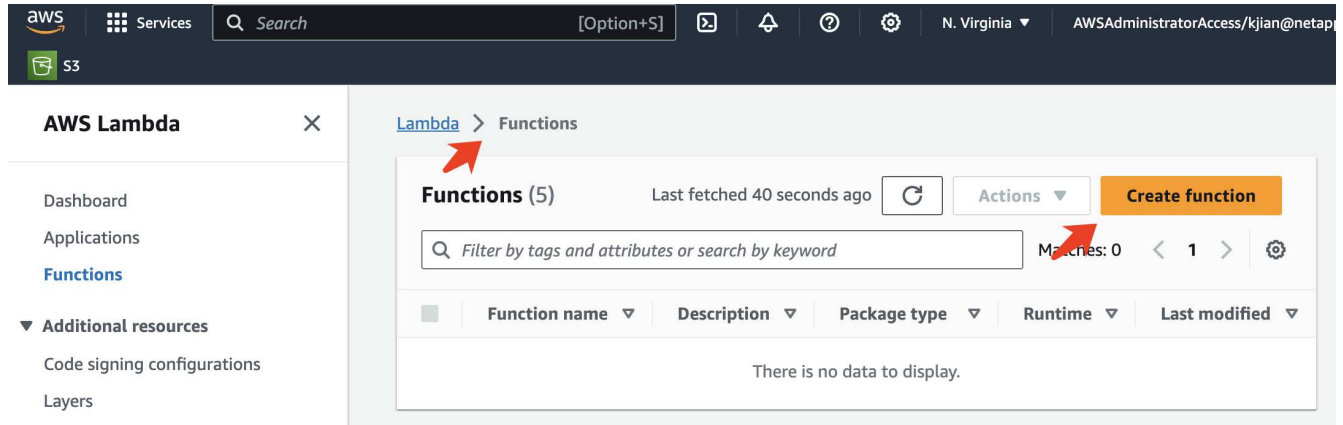
6. Seleccione la **Configuración de ciclo de vida** creada y haga clic en **Actualizar instancia de bloc de notas**.



## Función sin servidor de AWS Lambda

Como se mencionó anteriormente, la función **AWS Lambda** es responsable de poner en funcionamiento la instancia **AWS SageMaker Notebook**.

1. Para crear una función **AWS Lambda**, navegue hasta el panel correspondiente, cambie a la pestaña **Funciones** y haga clic en **Crear función**.



2. Por favor, archiva todas las entradas requeridas en la página y recuerda cambiar el tiempo de ejecución a **Python 3,10**.



aws Services Search [Option+S] N. Virgi AWSAdministratorAccess/kjian@

S3

Lambda > Functions > Create function

## Create function [Info](#)

AWS Serverless Application Repository applications have moved to [Create application](#).

☒ **Author from scratch**  
Start with a simple Hello World example.

☐ **Use a blueprint**  
Build a Lambda application from sample code and configuration presets for common use cases.

☐ **Container image**  
Select a container image to deploy for your function.

### Basic information

**Function name**  
Enter a name that describes the purpose of your function.

fsxn-demo-mlops

Use only letters, numbers, hyphens, or underscores with no spaces.

**Runtime** [Info](#)  
Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

Python 3.10

**Architecture** [Info](#)  
Choose the instruction set architecture you want for your function code.

☒ x86\_64

☐ arm64

**Permissions** [Info](#)  
By default, Lambda will create an execution role with permissions to upload logs to Amazon CloudWatch Logs. You can customize this default role later when adding triggers.

3. Verifique que el rol designado tiene el permiso requerido **AmazonSageMakerFullAccess** y haga clic en el botón **Crear función**.

aws Services Search [Option+S] N. Virgi AWSAdministratorAccess/kjian@

S3

Use only letters, numbers, hyphens, or underscores with no spaces.

**Runtime** [Info](#)  
Choose the language to use to write your function. Note that the console code editor supports only Node.js, Python, and Ruby.

Python 3.10

**Architecture** [Info](#)  
Choose the instruction set architecture you want for your function code.

☒ x86\_64  
☐ arm64

**Permissions** [Info](#)  
By default, Lambda will create an execution role with permissions to upload logs to Amazon CloudWatch Logs. You can customize this default role later when adding triggers.

▼ **Change default execution role**

**Execution role**  
Choose a role that defines the permissions of your function. To create a custom role, go to the [IAM console](#).

☐ Create a new role with basic Lambda permissions  
☒ Use an existing role  
☐ Create a new role from AWS policy templates

**Existing role**  
Choose an existing role that you've created to be used with this Lambda function. The role must have permission to upload logs to Amazon CloudWatch Logs.

service-role/fsxn-demo-mlops-role-585jzdny

[View the fsxn-demo-mlops-role-585jzdny role](#) on the IAM console.

► **Advanced settings**

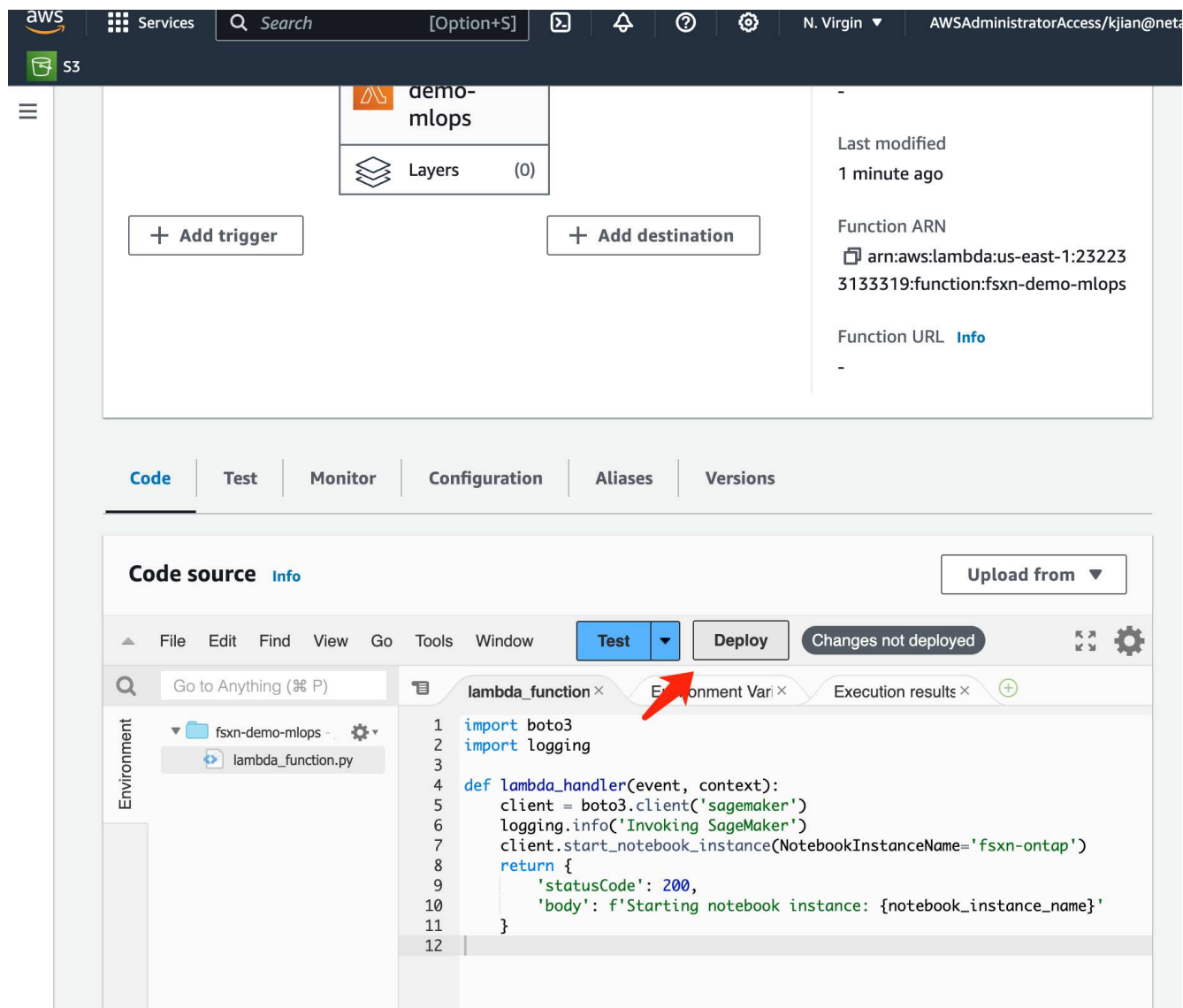
Cancel Create function

4. Seleccione la función Lambda creada. En la pestaña de código, copie y pegue el siguiente código en el área de texto. Este código inicia la instancia de notebook llamada **fsxn-ontap**.

```
import boto3
import logging

def lambda_handler(event, context):
    client = boto3.client('sagemaker')
    logging.info('Invoking SageMaker')
    client.start_notebook_instance(NotebookInstanceName='fsxn-ontap')
    return {
        'statusCode': 200,
        'body': f'Starting notebook instance: {notebook_instance_name}'
    }
```

5. Haga clic en el botón **Desplegar** para aplicar este cambio de código.



6. Para especificar cómo activar esta función de AWS Lambda, haga clic en el botón Agregar Disparador.

The screenshot shows the AWS Lambda console interface. At the top, the navigation bar includes the AWS logo, 'Services', a search bar, and the user's profile. The breadcrumb trail indicates the path: [Lambda](#) > [Functions](#) > fsxn-demo-mlops. The function name 'fsxn-demo-mlops' is prominently displayed. To the right of the name are buttons for 'Throttle', 'Copy ARN', and an 'Actions' dropdown menu. Below the function name, the 'Function overview' section is expanded, showing a card for the function with its icon and a 'Layers (0)' section. Two buttons, '+ Add trigger' and '+ Add destination', are visible. A red arrow points to the '+ Add trigger' button. On the right side of the overview, a details panel lists: 'Description' (empty), 'Last modified' (2 minutes ago), 'Function ARN' (arn:aws:lambda:us-east-1:232233133319:function:fsxn-demo-mlops), and 'Function URL' (empty with an 'Info' link).

7. Seleccione EventBridge en el menú desplegable y, a continuación, haga clic en el botón de opción con la etiqueta Crear una nueva regla. En el campo de expresión de programación, introduzca `rate(1 day)`, Y haga clic en el botón Agregar para crear y aplicar esta nueva regla de trabajo cron a la función AWS Lambda.

aws Services Search [Option+S] N. Virginia AWSAdministratorAccess

S3

[Lambda](#) > Add trigger

## Add trigger

**Trigger configuration** [Info](#)

**EventBridge (CloudWatch Events)**  
aws asynchronous schedule management-tools

**Rule**  
Pick an existing rule, or create a new one.

☒ Create a new rule  
☐ Existing rules

**Rule name**  
Enter a name to uniquely identify your rule.

mlops-retraining-trigger

**Rule description**  
Provide an optional description for your rule.

**Rule type**  
Trigger your target based on an event pattern, or based on an automated schedule.

☐ Event pattern  
☒ Schedule expression

**Schedule expression**  
Self-trigger your target on an automated schedule using [Cron or rate expressions](#). Cron expressions are in UTC.

rate(1 day)

e.g. rate(1 day), cron(0 17 ? \* MON-FRI \*)

Lambda will add the necessary permissions for Amazon EventBridge (CloudWatch Events) to invoke your Lambda function from this trigger. [Learn more](#) about the Lambda permissions model.

Cancel Add

Después de completar la configuración en dos pasos, diariamente, la función **AWS Lambda** iniciará el **SageMaker Notebook**, realizará el reciclaje del modelo utilizando los datos del repositorio **FSxN**, volverá a desplegar el modelo actualizado en el entorno de producción y cerrará automáticamente la instancia **SageMaker Notebook** para optimizar los costos. Esto garantiza que el modelo permanezca actualizado.

Esto concluye el tutorial para desarrollar un pipeline de MLOps.

## MLOps de multinube híbrida con Domino Data Lab y NetApp

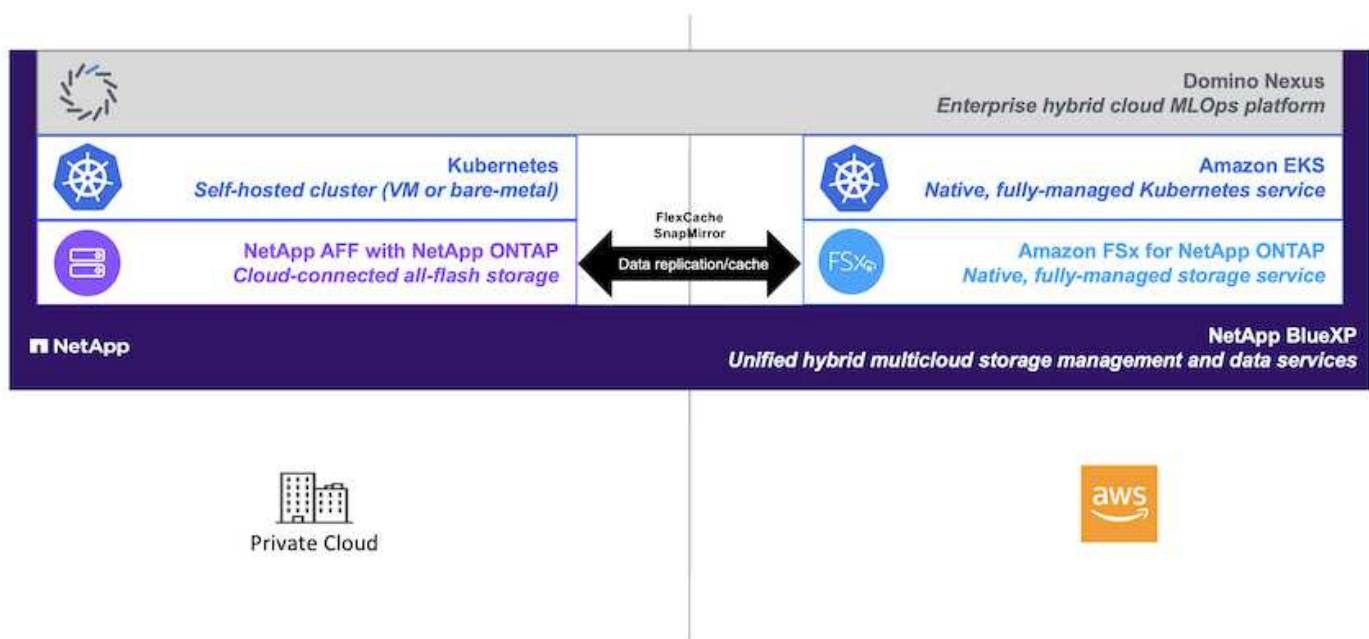
## MLOps de multinube híbrida con Domino Data Lab y NetApp

Mike Oglesby, NetApp

Organizaciones de todo el mundo están adoptando actualmente la IA para transformar sus empresas y sus procesos. Por ello, la infraestructura de computación preparada para la IA suele escasear. Las empresas están adoptando arquitecturas MLOps de multinube híbrida para aprovechar los entornos informáticos disponibles en diferentes regiones, centros de datos y nubes, logrando así un equilibrio entre costes, disponibilidad y rendimiento.

Domino Nexus, de Domino Data Lab, es un plano de control de MLOps unificado que te permite ejecutar cargas de trabajo de ciencia de datos y aprendizaje automático en cualquier clúster de computación, en cualquier nube, región o on-premises. Unifica los silos de ciencia de datos en toda la empresa, de modo que dispone de un lugar para crear, poner en marcha y supervisar modelos. Del mismo modo, las funcionalidades de gestión de datos en cloud híbrido de NetApp le permiten llevar sus datos a sus trabajos y espacios de trabajo, independientemente de dónde se ejecuten. Al emparejar Domino Nexus con NetApp, dispone de la flexibilidad necesaria para programar cargas de trabajo en diferentes entornos sin tener que preocuparse por la disponibilidad de los datos. En otras palabras, tiene la capacidad de enviar sus cargas de trabajo y sus datos al entorno de computación apropiado, lo que le permite acelerar sus puestas en marcha de IA a la vez que se rigen por las normativas en torno a la privacidad y soberanía de los datos.

Esta solución demuestra la puesta en marcha de un plano de control unificado de MLOps que incorpora un clúster de Kubernetes en las instalaciones y un clúster de Elastic Kubernetes Service (EKS) que se ejecuta en Amazon Web Services (AWS).



### Visión general de la tecnología

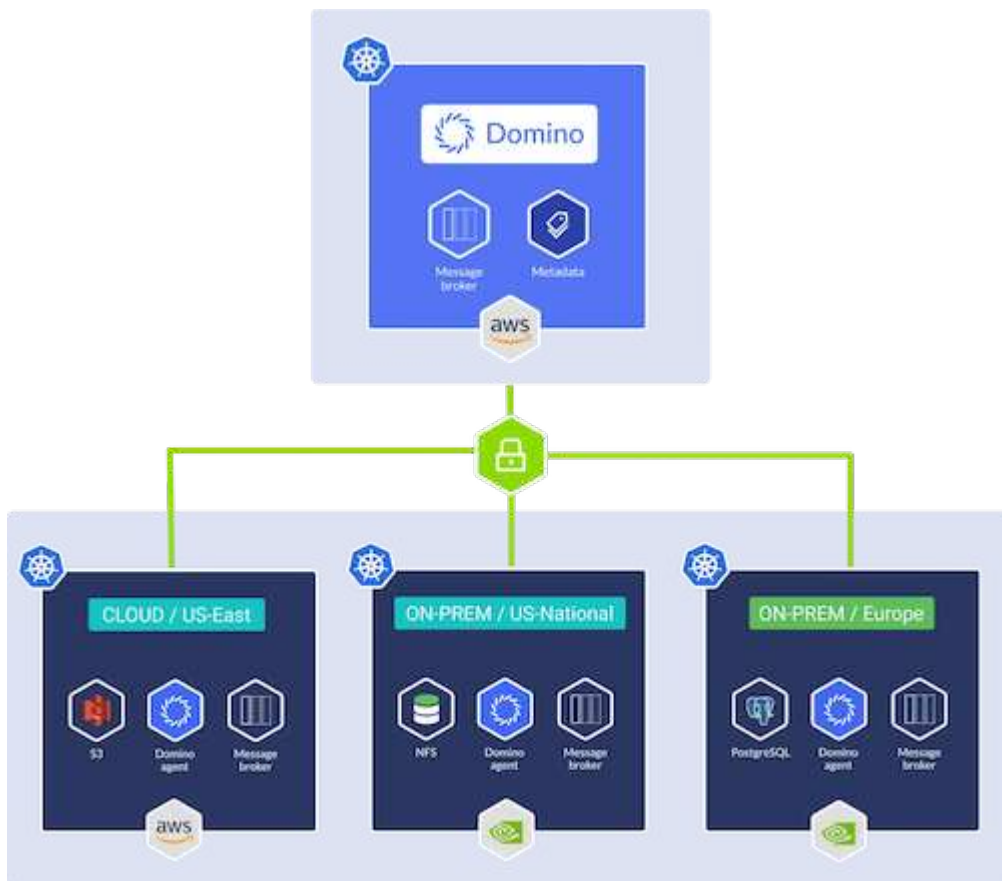
#### Domino Data Lab

Domino Data Lab impulsa empresas basadas en modelos con su plataforma líder de IA empresarial de confianza más del 20 % de las empresas Fortune 100. Domino acelera el desarrollo y la puesta en marcha del trabajo de ciencia de datos, a la vez que aumenta la colaboración y la gestión. Con Domino, las empresas de todo el mundo pueden desarrollar mejores medicamentos, cultivar cultivos más productivos, construir mejores

coches y mucho más. Fundada en 2013, Domino cuenta con el respaldo de Coatue Management, Great Hill Partners, Highland Capital, Sequoia Capital y otros inversores líderes.

Domino permite a las empresas y a sus científicos de datos crear, implementar y administrar IA en una plataforma unificada de extremo a extremo, de forma rápida, responsable y rentable. Los equipos pueden acceder a todos los datos, herramientas, computación, modelos y proyectos que necesitan en cualquier entorno, para que puedan colaborar, reutilizar el trabajo pasado, seguir los modelos en producción para mejorar la precisión, estandarizar con las mejores prácticas y hacer que la IA sea responsable y gobernada.

- **Abierto y flexible:** Acceda al ecosistema más amplio de herramientas comerciales y de código abierto, e infraestructura, para las mejores innovaciones y sin bloqueo de proveedores.
- **Sistema de registro:** Centro central de operaciones de IA y conocimiento en toda la empresa, lo que permite mejores prácticas, colaboración interfuncional, innovación más rápida y eficiencia.
- **Integrado:** Los flujos de trabajo integrados y la automatización — creados para procesos empresariales, controles y gobierno — satisfacen sus necesidades de cumplimiento y regulación.
- **Multinube híbrida:** Ejecuta cargas de trabajo de IA cerca de tus datos en cualquier lugar (on-premises, híbridas, en cualquier nube o en múltiples nubes) para un menor costo, un rendimiento óptimo y cumplimiento normativo.



### Nexus de Domino

Domino Nexus es un único panel que le permite ejecutar cargas de trabajo de ciencia de datos y aprendizaje automático en cualquier clúster de computación, en cualquier nube, región o on-premises. Unifica los silos de ciencia de datos en toda la empresa, de modo que dispone de un lugar para crear, poner en marcha y supervisar modelos.

## BlueXP de NetApp

NetApp BlueXP unifica todos los servicios de datos y almacenamiento de NetApp en una única herramienta que te permite crear, proteger y gobernar tus activos de datos multicloud híbridos. Proporciona una experiencia unificada para el almacenamiento y los servicios de datos en entornos en las instalaciones y en el cloud, y permite la sencillez operativa a través de la potencia de AIOps, con los parámetros de consumo flexible y la protección integrada que necesita el mundo actual dirigido por el cloud.

## ONTAP de NetApp

ONTAP 9, la última generación del software de gestión del almacenamiento de NetApp, permite a las empresas modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. ONTAP ofrece las mejores capacidades de gestión de datos y permite la gestión y protección de los datos con un solo conjunto de herramientas, sin importar dónde residan. También puede mover los datos libremente a donde sea necesario: El perímetro, el núcleo o el cloud. ONTAP 9 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen los datos esenciales y permiten disfrutar de funcionalidades de infraestructura de nueva generación en arquitecturas de cloud híbrido.

### Simplificar la gestión de los datos

La gestión de los datos es crucial para las operaciones TECNOLÓGICAS empresariales y los científicos de datos, para que se utilicen recursos apropiados para las aplicaciones de IA y para entrenar conjuntos de datos de IA/ML. La siguiente información adicional sobre las tecnologías de NetApp no está disponible para esta validación, pero puede ser relevante en función de su puesta en marcha.

El software para la gestión de datos ONTAP incluye las siguientes funciones para mejorar y simplificar las operaciones, y reducir el coste total de funcionamiento:

- Compactación de datos inline y deduplicación expandida. La compactación de datos reduce el espacio perdido dentro de los bloques de almacenamiento, mientras que la deduplicación aumenta la capacidad efectiva de forma significativa. Esto es aplicable a los datos almacenados localmente y a los datos organizados en niveles en el cloud.
- Calidad de servicio (AQoS) mínima, máxima y adaptativa. Los controles granulares de calidad de servicio (QoS) ayudan a mantener los niveles de rendimiento para aplicaciones críticas en entornos altamente compartidos.
- FabricPool de NetApp. Proporciona la organización automática en niveles de datos fríos en opciones de almacenamiento en cloud privado como Amazon Web Services (AWS), Azure y la solución de almacenamiento StorageGRID de NetApp. Para obtener más información sobre FabricPool, consulte "[TR-4598: Prácticas recomendadas de FabricPool](#)".

### Acelere y proteja sus datos

ONTAP no solo ofrece niveles de rendimiento y protección de datos superiores, sino que amplía estas capacidades de las siguientes maneras:

- Rendimiento y menor latencia. ONTAP ofrece la salida más alta posible con la menor latencia posible.
- Protección de datos. ONTAP ofrece capacidades integradas de protección de datos, con una administración común entre todas las plataformas.
- Cifrado de volúmenes de NetApp (NVE). ONTAP ofrece cifrado nativo en el nivel de volumen y permite la gestión de claves incorporada o externa.
- Multi-tenancy y autenticación multifactor. ONTAP permite compartir recursos de infraestructura con los niveles más altos de seguridad.



## **Infraestructura preparada para futuros retos**

ONTAP ayuda a satisfacer las exigentes y siempre cambiantes necesidades de su empresa con las siguientes funciones:

- Escalado sencillo y operaciones no disruptivas. ONTAP admite la adición no disruptiva de capacidad a las controladoras existentes y a clústeres de escalado horizontal. Los clientes pueden empezar a utilizar tecnologías punteras como NVMe y FC 32 GB, sin necesidad de realizar costosas migraciones de datos y sin cortes.
- Conexión de cloud. ONTAP es el software de gestión de almacenamiento con mejor conexión de cloud e incluye opciones de almacenamiento definido por software e instancias nativas del cloud en todos los clouds públicos.
- Integración con aplicaciones emergentes. ONTAP ofrece servicios de datos de clase empresarial para plataformas y aplicaciones de última generación, como vehículos autónomos, ciudades inteligentes e Industria 4.0, utilizando la misma infraestructura que da soporte a las aplicaciones empresariales existentes.

### **Amazon FSx para ONTAP de NetApp**

Amazon FSx para NetApp ONTAP es un servicio AWS de primer proveedor totalmente gestionado que proporciona un almacenamiento de archivos altamente fiable, escalable, de alto rendimiento y con muchas funciones integrado en el popular sistema de archivos ONTAP de NetApp. FSx para ONTAP combina las funciones, el rendimiento, las funcionalidades y las operaciones API de los sistemas de archivos de NetApp con la agilidad, la escalabilidad y la simplicidad de un servicio AWS totalmente gestionado.

### **Astra Trident de NetApp**

Astra Trident permite el consumo y la gestión de recursos de almacenamiento en todas las plataformas de almacenamiento de NetApp más conocidas, tanto en el cloud público como en las instalaciones, incluidos ONTAP (AFF, FAS, Select, Cloud, Amazon FSx para NetApp ONTAP), el software Element (NetApp HCI, SolidFire), el servicio Azure NetApp Files y Cloud Volumes Service en Google Cloud. Astra Trident es una interfaz de almacenamiento de contenedores (CSI) que ordena el almacenamiento dinámico conforme a la normativa que se integra de forma nativa con Kubernetes.

### **Kubernetes**

Kubernetes es una plataforma de orquestación de contenedores distribuida de código abierto que originalmente diseñada por Google y que ahora se mantiene mediante Cloud Native Computing Foundation (CNCF). Kubernetes permite la automatización de funciones de puesta en marcha, gestión y escalado para aplicaciones en contenedores, y es la plataforma de orquestación de contenedores dominante en entornos empresariales.

### **Amazon Elastic Kubernetes Service (EKS)**

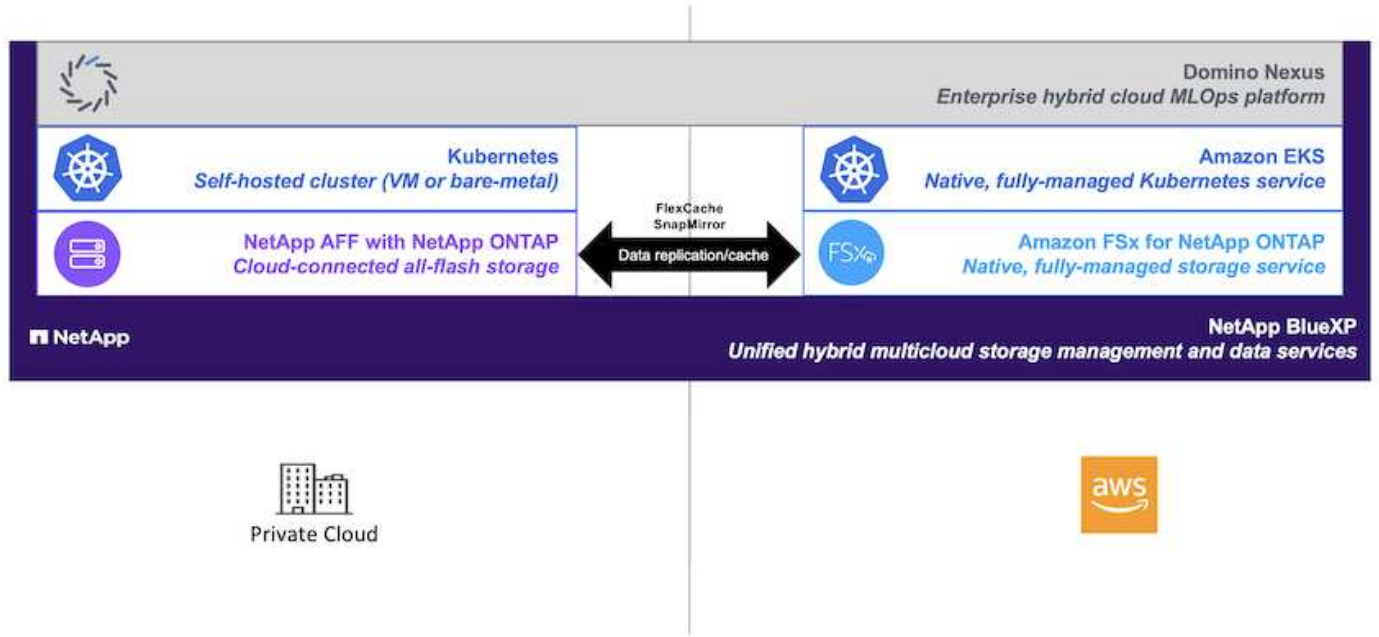
Amazon Elastic Kubernetes Service (Amazon EKS) es un servicio de Kubernetes gestionado en la nube de AWS. Amazon EKS gestiona automáticamente la disponibilidad y escalabilidad de los nodos del plano de control de Kubernetes responsables de programar contenedores, gestionar la disponibilidad de las aplicaciones, almacenar datos del clúster y otras tareas clave. Con Amazon EKS, puede aprovechar todo el rendimiento, la escalabilidad, la fiabilidad y la disponibilidad de la infraestructura de AWS, así como las integraciones con las redes y los servicios de seguridad de AWS.

## **Arquitectura**

Esta solución combina las capacidades de programación de cargas de trabajo de la

multinube híbrida de Domino Nexus con los servicios de datos de NetApp para crear una plataforma MLOps unificada de la nube híbrida. Consulte la siguiente tabla para obtener más información.

Componente	Nombre	Entorno Oracle
Plano de control de MLOps	"Domino Enterprise AI Platform con Domino Nexus"	AWS
Entornos informáticos de la plataforma MLOps	"Planos de datos Domino Nexus"	AWS, centro de datos en las instalaciones
Plataforma de informática en las instalaciones	"Kubernetes" con "Astra Trident de NetApp"	Centro de datos en las instalaciones
Plataforma de tecnología cloud	"Amazon Elastic Kubernetes Service (EKS)" con "Astra Trident de NetApp"	AWS
Plataforma de datos en las instalaciones	"Dispositivo de almacenamiento de NetApp" alimentado por "ONTAP de NetApp"	Centro de datos en las instalaciones
Plataforma de datos cloud	"Amazon FSx para ONTAP de NetApp"	AWS



### Configuración inicial

En esta sección se describen las tareas de configuración iniciales que se deben realizar para utilizar Domino Nexus con servicios de datos de NetApp en un entorno híbrido que incorpora un centro de datos local y AWS.

### Requisitos previos

Antes de realizar los pasos que se describen en esta sección, asumimos que ya ha realizado las siguientes

tareas:

- Ya ha puesto en marcha y configurado su plataforma de almacenamiento NetApp ONTAP en las instalaciones. Para obtener más información, consulte la ["Documentación de productos de NetApp"](#).
- Ya ha aprovisionado una instancia de Amazon FSx para NetApp ONTAP en AWS. Para obtener más información, consulte la ["Página de producto de Amazon FSx para NetApp ONTAP"](#).
- Ya ha aprovisionado un clúster de Kubernetes en su centro de datos en las instalaciones. Para obtener más información, consulte la ["Guía del administrador de Domino"](#).
- Ya ha aprovisionado un clúster de Amazon EKS en AWS. Para obtener más información, consulte la ["Guía del administrador de Domino"](#).
- Ha instalado Astra Trident de NetApp en su clúster de Kubernetes en las instalaciones. Además, ha configurado esta instancia de Trident para utilizar su plataforma de almacenamiento de NetApp ONTAP local al aprovisionar y gestionar recursos de almacenamiento. Para obtener más información, consulte la ["Documentación de NetApp Astra Trident"](#).
- Ha instalado Astra Trident de NetApp en su clúster de Amazon EKS. Además, ha configurado esta instancia de Trident de modo que utilice su instancia de Amazon FSx para NetApp ONTAP al aprovisionar y gestionar recursos de almacenamiento. Para obtener más información, consulte la ["Documentación de NetApp Astra Trident"](#).
- Debe disponer de conectividad de red bidireccional entre su centro de datos en las instalaciones y la nube privada virtual (VPC) en AWS. Para obtener más información sobre las distintas opciones para implementarlo, consulte la ["Documentación de Amazon Virtual Private Network \(VPN\)"](#).

### Instale Domino Enterprise AI Platform en AWS

Para instalar Domino Enterprise MLOps Platform en AWS, siga las instrucciones descritas en ["Guía del administrador de Domino"](#). Debe implementar Domino en el mismo clúster de Amazon EKS que había aprovisionado previamente. Además, Astra Trident de NetApp ya debe estar instalado y configurado en este clúster de EKS, y debe especificar una clase de almacenamiento gestionado por Trident como la clase de almacenamiento compartido de su archivo de configuración de instalación de domino.yml.



Consulte la ["Guía de referencia de configuración de instalación de Domino"](#) para obtener detalles sobre cómo especificar una clase de almacenamiento compartido en el archivo de configuración de instalación de domino.yml.



["Informe técnico TR-4952"](#) Explica la puesta en marcha de Domino en AWS con Amazon FSx para NetApp ONTAP y puede ser una referencia útil para solucionar los problemas que surjan.

### Habilite Domino Nexus

A continuación, debe habilitar Domino Nexus. Consulte la ["Guía del administrador de Domino"](#) para obtener más detalles.

### Ponga en marcha un plano de datos de Domino en su centro de datos local

A continuación, debe poner en marcha un plano de datos de Domino en su centro de datos local. Es necesario poner en marcha este plano de datos en el clúster de Kubernetes local que ya haya aprovisionado. Además, Astra Trident de NetApp ya debe estar instalado y configurado en este clúster de Kubernetes. Consulte la ["Guía del administrador de Domino"](#) para obtener más detalles.

## Expone los volúmenes NetApp existentes a Domino

Esta sección describe las tareas que se deben realizar para exponer los volúmenes NFS de NetApp ONTAP existentes a la plataforma MLOps de Domino. Estos mismos pasos se aplican tanto en las instalaciones como en AWS.

### ¿Por qué exponer volúmenes de NetApp ONTAP a Domino?

El uso de NetApp Volumes junto con Domino ofrece las siguientes ventajas:

- Puede ejecutar cargas de trabajo en conjuntos de datos extremadamente grandes aprovechando las funcionalidades de escalado horizontal de ONTAP de NetApp.
- Puede ejecutar cargas de trabajo en varios nodos de computación sin tener que copiar los datos en cada nodo.
- Puede aprovechar las funcionalidades de sincronización y movimiento de datos de multicloud híbrido de NetApp para acceder a sus datos en varios centros de datos o clouds.
- Quieres poder crear rápida y fácilmente una caché de tus datos en otro centro de datos o en otro cloud.

### Muestre los volúmenes de NFS existentes que Astra Trident no aprovisionó

Si tu volumen de NFS de NetApp ONTAP existente no se ha aprovisionado en Astra Trident, sigue los pasos descritos en esta subsección.

### Cree un VP y una RVP en Kubernetes



Para los volúmenes en las instalaciones, cree el VP y el PVC en su clúster de Kubernetes en las instalaciones. Para Amazon FSx para NetApp ONTAP Volumes, cree el VP y la RVP en Amazon EKS.

Primero, debe crear un volumen persistente (PV) y una reclamación de volumen persistente (PVC) en su clúster de Kubernetes. Para crear el VP y la RVP, utilice el ["EJEMPLO DE PV/PVC DE NFS"](#) De la guía de administración de Domino y actualice los valores para reflejar el entorno. Asegúrese de especificar los valores correctos para el namespace, `nfs.path`, y `nfs.server` campos. Además, recomendamos indicar a sus nombres VP y RVP únicos que representen la naturaleza de los datos que se almacenan en el volumen NFS de ONTAP correspondiente. Por ejemplo, si el volumen contiene imágenes de defectos de fabricación, puede nombrar el VP, `pv-mfg-defect-images`, Y el PVC, `pvc-mfg-defect-images`.

### Registrar volumen de datos externo en Domino

A continuación, debe registrar un volumen de datos externo en Domino. Para registrar un volumen de datos externo, consulte la ["instrucciones"](#) En la guía del administrador de Domino. Al registrar el volumen, asegúrese de seleccionar «NFS» en el menú desplegable «Tipo de volumen». Después de seleccionar «NFS», deberías ver tu RVP en la lista de volúmenes disponibles.

**Register an External Volume**

1 **Volume**  
NFS

2 **Configuration**  
Read-Only

3 **Access**  
Everyone

**Volume Type**

NFS

**Available Volumes**

☐ chatbot-data-cache

Cancel Next >

#### Exponer los volúmenes existentes que fueron aprovisionados por Astra Trident

Si el volumen existente se aprovisionó con Astra Trident, siga los pasos descritos en esta subsección.

#### Edite la RVP existente

Si tu volumen fue aprovisionado por Astra Trident, ya tienes una reclamación de volumen persistente (RVP) correspondiente al volumen. Para exponer este volumen a Domino, debe editar la PVC y agregar la siguiente etiqueta a la lista de etiquetas de la `metadata.labels` campo:

```
"dominodatalab.com/external-data-volume": "Generic"
```

#### Registrar volumen de datos externo en Domino

A continuación, debe registrar un volumen de datos externo en Domino. Para registrar un volumen de datos externo, consulte la ["instrucciones"](#) En la guía del administrador de Domino. Al registrar el volumen, asegúrese de seleccionar «Genérico» en el menú desplegable «Tipo de volumen». Después de seleccionar «Genérico», deberías ver tu PVC en la lista de «Volúmenes disponibles».

## Acceda a los mismos datos en diferentes entornos

Esta sección describe las tareas que se deben realizar para acceder a los mismos datos en diferentes entornos informáticos. En la plataforma MLOps de Domino, los entornos informáticos se denominan «planos de datos». Siga las tareas descritas en esta sección si sus datos residen en un volumen de NetApp en un plano de datos, pero necesita acceder a ellos en otro plano de datos. Este tipo de escenario se suele denominar «explosión» o, cuando el entorno de destino es el cloud, «explosión de cloud». Esta funcionalidad a menudo es necesaria cuando se utilizan recursos informáticos limitados o sobresuscritos. Por ejemplo, si su clúster de computación en las instalaciones está sobresuscrito, puede que desee programar cargas de trabajo en el cloud donde pueden iniciarse inmediatamente.

Hay dos opciones recomendadas para acceder a un volumen de NetApp que reside en un plano de datos diferente. Estas opciones se describen en las subsecciones siguientes. Elija una de estas opciones en función de sus requisitos específicos. En la siguiente tabla, se describen los beneficios e inconvenientes de las dos opciones.

Opción	Beneficios	Inconvenientes
Opción 1: Caché	<ul style="list-style-type: none"><li>- Flujo de trabajo más simple</li><li>- Capacidad de almacenar en caché un subconjunto de datos en función de las necesidades</li><li>- Capacidad para escribir datos de nuevo en la fuente</li><li>- No hay copia remota para administrar</li></ul>	<ul style="list-style-type: none"><li>- Aumento de la latencia en el acceso inicial a los datos a medida que la caché está hidratada.</li></ul>
Opción 2 - Espejo	<ul style="list-style-type: none"><li>- Copia completa del volumen fuente</li><li>- Sin aumento de la latencia debido a la hidratación de la caché (después de la operación del espejo está completa)</li></ul>	<ul style="list-style-type: none"><li>- Debe esperar a que se complete la operación del espejo antes de acceder a los datos</li><li>- Debe gestionar una copia remota</li><li>- No hay capacidad de escribir de nuevo a la fuente</li></ul>

### Opción 1: Crear una caché de un volumen que reside en un plano de datos diferente

Con "[Tecnología NetApp FlexCache](#)", Puede crear una caché de un volumen NetApp que reside en un plano de datos diferente. Por ejemplo, si tiene un volumen de NetApp en su plano de datos local y necesita acceder a dicho volumen en su plano de datos de AWS, puede crear una caché del volumen en AWS. En esta sección se describen las tareas que se deben realizar para crear una caché de un volumen NetApp que se encuentre en un plano de datos diferente.

### Crear volumen FlexCache en entorno de destino



Si el entorno de destino es su centro de datos en las instalaciones, creará el volumen FlexCache en su sistema ONTAP en las instalaciones. Si el entorno de destino es AWS, creará el volumen de FlexCache en tu instancia de Amazon FSx para NetApp ONTAP.

Primero, debe crear un volumen de FlexCache en el entorno de destino.

Recomendamos utilizar BlueXP para crear el volumen de FlexCache. Para crear un volumen de FlexCache con BlueXP, siga las instrucciones que se describen en la ["Documentación de almacenamiento en caché de volúmenes de BlueXP"](#).

Si prefiere no usar BlueXP, puede usar ONTAP System Manager o la interfaz de línea de comandos de ONTAP para crear el volumen de FlexCache. Para crear un volumen FlexCache con System Manager, consulte las instrucciones descritas en la ["Documentación de ONTAP"](#). Para crear un volumen de FlexCache con la CLI de ONTAP, consulte las instrucciones que se describen en ["Documentación de ONTAP"](#).

Si desea automatizar este proceso, puede utilizar el ["API de BlueXP"](#), la ["API REST de ONTAP"](#), o el ["Colección Ansible de ONTAP"](#).



System Manager no está disponible en Amazon FSx para NetApp ONTAP.

## **Exponga el volumen FlexCache a Domino**

A continuación, debe exponer el volumen FlexCache a la plataforma Domino MLOps. Para exponer el volumen FlexCache a Domino, siga las instrucciones que se describen en la subsección 'Exponga volúmenes NFS existentes que no se aprovisionaban con Astra Trident' de la ["Sección 'Exponga los volúmenes NetApp existentes a Domino'"](#) de esta solución.

Ahora, podrá montar el volumen FlexCache al iniciar trabajos y espacios de trabajo en el plano de datos de destino, como se muestra en las siguientes capturas de pantalla.

### **Antes de crear Volumen FlexCache**

Start a Job

✓

Execution

FILE: main.py

ENV: Domino Sta...

✓

Compute Cluster

(optional)

✓

Data

Data that will be mounted

NAME	DATA TYPE	DATA PLANE	KIND
quick-start	Dataset	Local	Project
image-data	EDV	rtp-aillab-kube02 ...	Nfs

Unavailable in selected Dataplane

Change your Hardware Tier to mount currently unavailable data.

NAME	DATA TYPE	DATA PLANE	KIND
chatbot-data	EDV	rtp-aillab-kube02	Nfs

Cancel

< Back

Start

Después de exponer el volumen de FlexCache a Domino

79



Start a Job

✓ Execution

FILE: model.py

ENV: Domino Sta...

✓ Compute Cluster

(optional)

3 Data

Data that will be mounted

NAME	DATA TYPE	DATA PLANE	KIND
quick-start	Dataset	Local	Project
image-data	EDV	rtp-aillab-kube02	Nfs
chatbot-data	EDV	rtp-aillab-kube02	Nfs

Unavailable in selected Dataplane

Change your Hardware Tier to mount currently unavailable data.

NAME	DATA TYPE	DATA PLANE	KIND
No data found			

Cancel

< Back

Start

## Opción 2: Replicar un volumen que reside en un plano de datos distinto

Con "[Tecnología de replicación de datos SnapMirror de NetApp](#)", Puede crear una copia de un volumen NetApp que reside en un plano de datos diferente. Por ejemplo, si tiene un volumen de NetApp en su plano de datos local y necesita acceder a dicho volumen en su plano de datos de AWS, puede crear una copia del volumen en AWS. En esta sección se describen las tareas que se deben realizar para crear una copia de un volumen NetApp que se encuentre en un plano de datos diferente.

### Crear una relación de SnapMirror

Primero, debe crear una relación de SnapMirror entre el volumen de origen y un nuevo volumen de destino en el entorno de destino. Tenga en cuenta que el volumen de destino se creará como parte del proceso de creación de la relación de SnapMirror.

Recomendamos usar BlueXP para crear la relación de SnapMirror. Para crear una relación de SnapMirror con BlueXP, siga las instrucciones que se describen en ["Documentación de replicación de BlueXP"](#).

Si prefiere no usar BlueXP, puede usar ONTAP System Manager o la interfaz de línea de comandos de ONTAP para crear la relación de SnapMirror. Para crear una relación de SnapMirror con System Manager, consulte las instrucciones descritas en ["Documentación de ONTAP"](#). Para crear una relación de SnapMirror con la CLI de ONTAP, consulte las instrucciones descritas en ["Documentación de ONTAP"](#).

Si desea automatizar este proceso, puede utilizar el ["API de BlueXP"](#), la ["API REST de ONTAP"](#), o el ["Colección Ansible de ONTAP"](#).



System Manager no está disponible en Amazon FSx para NetApp ONTAP.

## Interrumpir la relación de SnapMirror

Luego, debe interrumpir la relación de SnapMirror para activar el volumen de destino para el acceso a los datos. Espere hasta que se complete la replicación inicial antes de realizar este paso.



Puede determinar si la replicación está completa o no, compruebe el estado de reflejo en BlueXP, ONTAP System Manager o la CLI de ONTAP. Cuando se complete la replicación, el estado de reflejo será Snapmirror.

Recomendamos usar BlueXP para interrumpir la relación de SnapMirror. Para interrumpir una relación de SnapMirror con BlueXP, siga las instrucciones indicadas en ["Documentación de replicación de BlueXP"](#).

Si prefiere no usar BlueXP, puede usar ONTAP System Manager o la interfaz de línea de comandos de ONTAP para interrumpir la relación de SnapMirror. Para interrumpir una relación de SnapMirror con System Manager, consulte las instrucciones descritas en ["Documentación de ONTAP"](#). Para interrumpir una relación de SnapMirror con la CLI de ONTAP, consulte las instrucciones que se describen en ["Documentación de ONTAP"](#).

Si desea automatizar este proceso, puede utilizar el ["API de BlueXP"](#), la ["API REST de ONTAP"](#), o el ["Colección Ansible de ONTAP"](#).

## Exponga el volumen de destino a Domino

A continuación, debe exponer el volumen de destino a la plataforma Domino MLOps. Para exponer el volumen de destino a Domino, siga las instrucciones que se describen en la subsección 'Exponga volúmenes NFS existentes que no se aprovisionaban con Astra Trident' de la ["Sección 'Exponga los volúmenes NetApp existentes a Domino'"](#) de esta solución.

Ahora, podrá montar el volumen de destino al iniciar trabajos y espacios de trabajo en el plano de datos de destino como se muestra en las siguientes capturas de pantalla.

## Antes de crear una relación de SnapMirror

Start a Job

Execution

FILE: main.py

ENV: Domino Sta...

Compute Cluster

(optional)

Data

Data that will be mounted

NAME	DATA TYPE	DATA PLANE	KIND
quick-start	Dataset	Local	Project
image-data	EDV	rtp-aillab-kube02 ...	Nfs

Unavailable in selected Dataplane

Change your Hardware Tier to mount currently unavailable data.

NAME	DATA TYPE	DATA PLANE	KIND
chatbot-data	EDV	rtp-aillab-kube02	Nfs

Cancel

< Back

Start

Después de exponer el volumen de destino a Domino

82

Start a Job

✓ Execution  
FILE: model.py  
ENV: Domino Sta...

✓ Compute Cluster  
(optional)

3 Data

Data that will be mounted

NAME	DATA TYPE	DATA PLANE	KIND
quick-start	Dataset	Local	Project
image-data	EDV	rtp-aillab-kube02	Nfs
chatbot-data	EDV	rtp-aillab-kube02	Nfs

Unavailable in selected Dataplane

Change your Hardware Tier to mount currently unavailable data.

NAME	DATA TYPE	DATA PLANE	KIND
No data found			

Cancel

< Back

Start

### Dónde encontrar información adicional

Para obtener más información sobre la información descrita en este documento, consulte los siguientes documentos y/o sitios web:

- Domino Data Lab

["https://domino.ai"](https://domino.ai)

- Nexus de Domino

["https://domino.ai/platform/nexus"](https://domino.ai/platform/nexus)

83

- BlueXP de NetApp

["https://bluexp.netapp.com"](https://bluexp.netapp.com)

- Software de gestión de datos ONTAP de NetApp

["https://www.netapp.com/data-management/ontap-data-management-software/"](https://www.netapp.com/data-management/ontap-data-management-software/)

- Soluciones de IA de NetApp

["https://www.netapp.com/artificial-intelligence/"](https://www.netapp.com/artificial-intelligence/)

## Reconocimientos

- Josh Mineroff, director ejecutivo de SA de Tech Alliances en Domino Data Lab
- Nicholas Jablonski, director técnico de campo, Domino Data Lab
- Es/Prabu Arjunan, arquitecto de soluciones en NetApp
- Brian Young, director de alianza global, partners de alianza tecnológica de NetApp

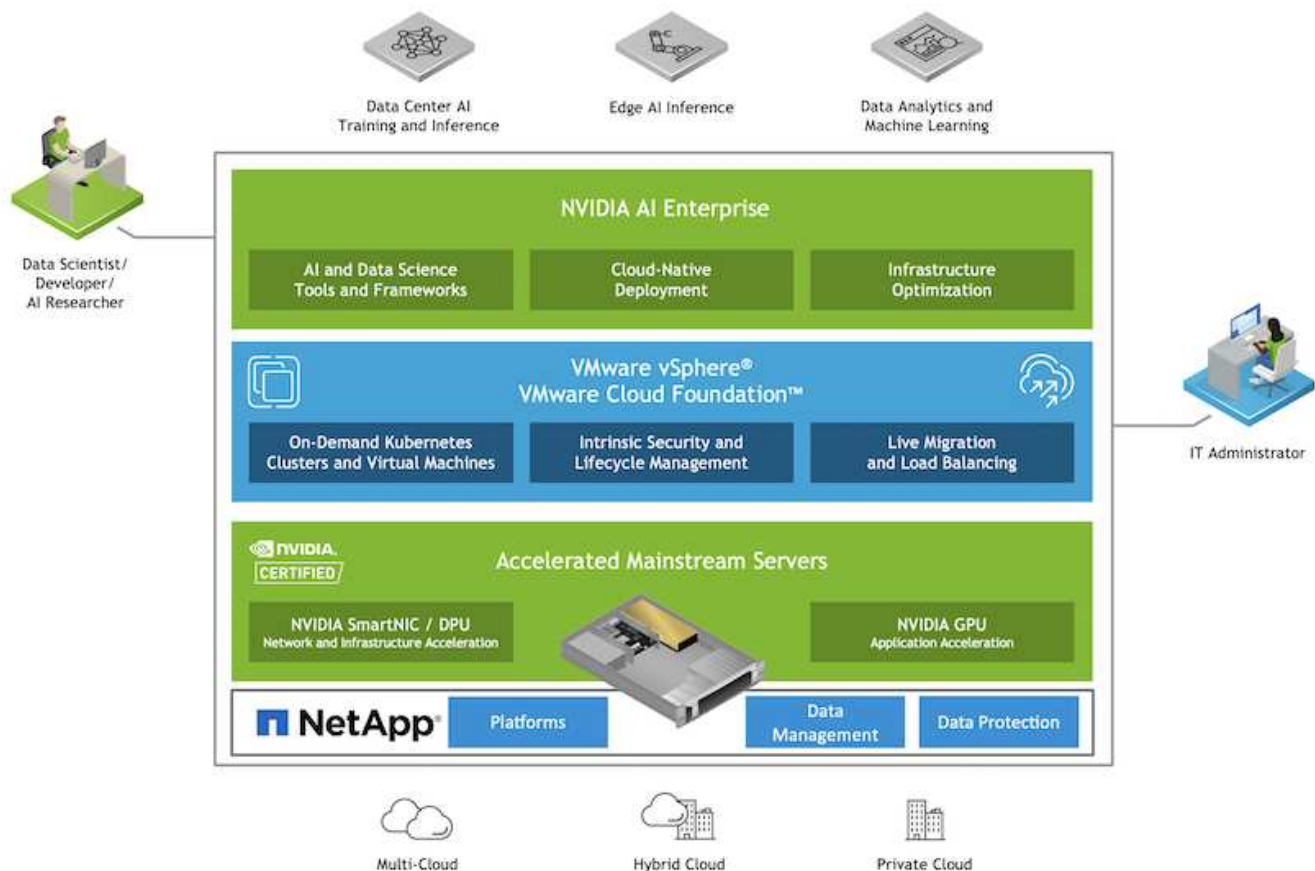
## NVIDIA AI Enterprise con NetApp y VMware

### NVIDIA AI Enterprise con NetApp y VMware

Mike Oglesby, NetApp

Para los arquitectos y administradores DE TECNOLOGÍA, las herramientas de IA pueden ser complicadas y no conocidas. Además, muchas plataformas de IA no están listas para las empresas. NVIDIA AI Enterprise, con la tecnología de NetApp y VMware, se ha creado para proporcionar una arquitectura de IA optimizada para la gran empresa.

NVIDIA AI Enterprise es una suite integral y nativa en el cloud de software de IA y análisis de datos optimizado, certificado y compatible con NVIDIA para ejecutarse en VMware vSphere con sistemas certificados por NVIDIA. Este software facilita la puesta en marcha, la gestión y el escalado simples y rápidos de las cargas de trabajo de IA en el entorno de cloud híbrido moderno. NVIDIA AI Enterprise, con la tecnología de NetApp y VMware, ofrece una excelente carga de trabajo de IA y gestión de datos en un paquete simplificado y conocido.



## Visión general de la tecnología

### IA Enterprise de NVIDIA

NVIDIA AI Enterprise es una suite integral y nativa en el cloud de software de IA y análisis de datos optimizado, certificado y compatible con NVIDIA para ejecutarse en VMware vSphere con sistemas certificados por NVIDIA. Este software facilita la puesta en marcha, la gestión y el escalado simples y rápidos de las cargas de trabajo de IA en el entorno de cloud híbrido moderno.

### GPU CLOUD DE NVIDIA (NGC)

NVIDIA NGC aloja un catálogo de software optimizado para GPU para que los profesionales de la IA puedan desarrollar sus soluciones de IA. También proporciona acceso a diversos servicios de IA, incluido NVIDIA base Command para el entrenamiento de modelos, NVIDIA Fleet Command para implementar y supervisar modelos, y el registro privado NGC para acceder y gestionar de forma segura el software de IA propio. Además, los clientes de NVIDIA AI Enterprise pueden solicitar soporte a través del portal NGC.

### VSphere de VMware

VMware vSphere es la plataforma de virtualización de VMware, que transforma los centros de datos en infraestructuras informáticas globales que incluyen recursos de CPU, almacenamiento y red. vSphere gestiona estas infraestructuras como un entorno operativo unificado y proporciona a los administradores las herramientas necesarias para gestionar los centros de datos que participan en ese entorno.

Los dos componentes principales de vSphere son ESXi y vCenter Server. ESXi es la plataforma de virtualización donde los administradores crean y ejecutan máquinas virtuales y dispositivos virtuales. vCenter Server es el servicio a través del cual los administradores gestionan varios hosts conectados en una red y un

pool de recursos de host.

## **ONTAP de NetApp**

ONTAP 9, la última generación del software de gestión del almacenamiento de NetApp, permite a las empresas modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. ONTAP ofrece las mejores capacidades de gestión de datos y permite la gestión y protección de los datos con un solo conjunto de herramientas, sin importar dónde residan. También puede mover los datos libremente a donde sea necesario: El perímetro, el núcleo o el cloud. ONTAP 9 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen los datos esenciales y permiten disfrutar de funcionalidades de infraestructura de nueva generación en arquitecturas de cloud híbrido.

### **Simplificar la gestión de los datos**

La gestión de los datos es crucial para las operaciones TECNOLÓGICAS empresariales y los científicos de datos, para que se utilicen recursos apropiados para las aplicaciones de IA y para entrenar conjuntos de datos de IA/ML. La siguiente información adicional sobre las tecnologías de NetApp no está disponible para esta validación, pero puede ser relevante en función de su puesta en marcha.

El software para la gestión de datos ONTAP incluye las siguientes funciones para mejorar y simplificar las operaciones, y reducir el coste total de funcionamiento:

- Compactación de datos inline y deduplicación expandida. La compactación de datos reduce el espacio perdido dentro de los bloques de almacenamiento, mientras que la deduplicación aumenta la capacidad efectiva de forma significativa. Esto es aplicable a los datos almacenados localmente y a los datos organizados en niveles en el cloud.
- Calidad de servicio (AQoS) mínima, máxima y adaptativa. Los controles granulares de calidad de servicio (QoS) ayudan a mantener los niveles de rendimiento para aplicaciones críticas en entornos altamente compartidos.
- FabricPool de NetApp. Proporciona la organización automática en niveles de datos fríos en opciones de almacenamiento en cloud privado como Amazon Web Services (AWS), Azure y la solución de almacenamiento StorageGRID de NetApp. Para obtener más información sobre FabricPool, consulte "[TR-4598: Prácticas recomendadas de FabricPool](#)".

### **Acelere y proteja sus datos**

ONTAP no solo ofrece niveles de rendimiento y protección de datos superiores, sino que amplía estas capacidades de las siguientes maneras:

- Rendimiento y menor latencia. ONTAP ofrece la salida más alta posible con la menor latencia posible.
- Protección de datos. ONTAP ofrece capacidades integradas de protección de datos, con una administración común entre todas las plataformas.
- Cifrado de volúmenes de NetApp (NVE). ONTAP ofrece cifrado nativo en el nivel de volumen y permite la gestión de claves incorporada o externa.
- Multi-tenancy y autenticación multifactor. ONTAP permite compartir recursos de infraestructura con los niveles más altos de seguridad.

### **Infraestructura preparada para futuros retos**

ONTAP ayuda a satisfacer las exigentes y siempre cambiantes necesidades de su empresa con las siguientes funciones:

- Escalado sencillo y operaciones no disruptivas. ONTAP admite la adición no disruptiva de capacidad a las controladoras existentes y a clústeres de escalado horizontal. Los clientes pueden empezar a utilizar tecnologías punteras como NVMe y FC 32 GB, sin necesidad de realizar costosas migraciones de datos y sin cortes.
- Conexión de cloud. ONTAP es el software de gestión de almacenamiento con mejor conexión de cloud e incluye opciones de almacenamiento definido por software (ONTAP Select) e instancias nativas del cloud (NetApp Cloud Volumes Service) en todos los clouds públicos.
- Integración con aplicaciones emergentes. ONTAP ofrece servicios de datos de clase empresarial para plataformas y aplicaciones de última generación, como vehículos autónomos, ciudades inteligentes e Industria 4.0, utilizando la misma infraestructura que da soporte a las aplicaciones empresariales existentes.

#### Kit de herramientas de operaciones de datos de NetApp

El kit de herramientas DataOps de NetApp es una herramienta basada en Python que simplifica la gestión de espacios de trabajo de desarrollo/formación y servidores de inferencia respaldados por un almacenamiento de NetApp de escalado horizontal y de alto rendimiento. Estas son algunas funcionalidades clave:

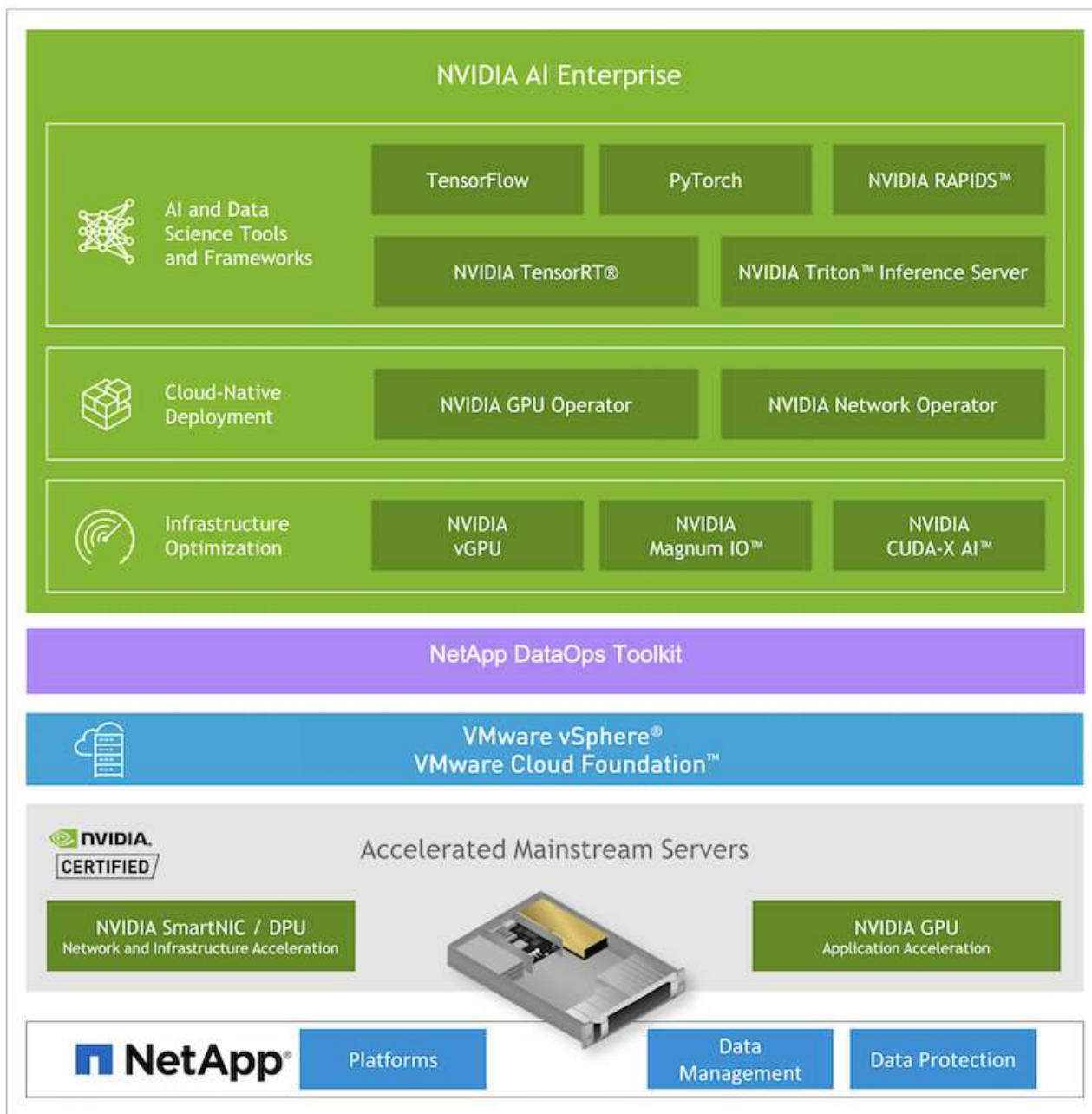
- Aprovechamiento con rapidez nuevos espacios de trabajo de JupyterLab de alta capacidad respaldados por el almacenamiento de NetApp de escalado horizontal y de alto rendimiento.
- Aprovechamiento rápidamente nuevas instancias del servidor de inferencia de NVIDIA Triton, respaldadas por un almacenamiento empresarial de NetApp.
- Clone por vía casi instantánea los espacios de trabajo de JupyterLab de gran capacidad para permitir la experimentación o la iteración rápida.
- Guarde copias Snapshot casi por vía instantánea de espacios de trabajo JupyterLab de gran capacidad para backup o seguimiento o creación de bases.
- Aprovechamiento, clonado y copias Snapshot de alta capacidad y alto rendimiento casi por vía casi instantánea.

#### Arquitectura

Esta solución se basa en una arquitectura probada y familiar que incluye NetApp, VMware y sistemas con certificación NVIDIA. Consulte la siguiente tabla para obtener más información.

Componente	Detalles
Software de IA y análisis de datos	<a href="#">"NVIDIA AI Enterprise para VMware"</a>
Plataforma de virtualización	<a href="#">"VSphere de VMware"</a>
Plataforma informática	<a href="#">"Sistemas con certificación NVIDIA"</a>
Plataforma de gestión de datos	<a href="#">"ONTAP de NetApp"</a>





## Configuración inicial

En esta sección se describen las tareas de configuración inicial que se deben realizar para utilizar NVIDIA AI Enterprise con NetApp y VMware.

### Requisitos previos

Antes de realizar los pasos descritos en esta sección, asumimos que ya ha implementado VMware vSphere y ONTAP de NetApp. Consulte la ["Matriz de compatibilidad de productos empresariales de NVIDIA AI"](#) Para obtener más detalles sobre las versiones de vSphere compatibles. Consulte la ["Documentación sobre las soluciones de NetApp y VMware"](#) Para obtener más información sobre cómo implementar VMware vSphere con ONTAP de NetApp.

## Instale el software de host de NVIDIA AI Enterprise

Para instalar el software host NVIDIA AI Enterprise, siga las instrucciones que se describen en las secciones 1-4 de la ["Guía de inicio rápido de NVIDIA AI Enterprise"](#).

## Utilice el software NVIDIA NGC

En esta sección, se describen las tareas que se deben realizar para utilizar el software empresarial NVIDIA NGC en un entorno de NVIDIA AI Enterprise.

### Configuración

En esta sección se describen las tareas de configuración inicial que se deben realizar para utilizar el software empresarial NVIDIA NGC en un entorno de NVIDIA AI Enterprise.

### Requisitos previos

Antes de realizar los pasos descritos en esta sección, asumimos que ya ha implementado el software host NVIDIA AI Enterprise siguiendo las instrucciones que se describen en la ["Configuración inicial"](#) página.

## Cree una máquina virtual Ubuntu Guest con vGPU

En primer lugar, debe crear una máquina virtual invitada Ubuntu 20.04 con vGPU. Para crear una VM huésped de Ubuntu 20.04 con vGPU, siga las instrucciones descritas en el ["Guía de puesta en marcha de NVIDIA AI Enterprise"](#).

## Descargue e instale el software invitado de NVIDIA

A continuación, debe instalar el software invitado NVIDIA requerido en la máquina virtual invitada que creó en el paso anterior. Para descargar e instalar el software invitado NVIDIA necesario en el equipo virtual invitado, siga las instrucciones que se describen en las secciones 5.1-5.4 de la ["Guía de inicio rápido de NVIDIA AI Enterprise"](#).



Al realizar las tareas de verificación descritas en la sección 5.4, es posible que necesite utilizar una etiqueta de versión de imagen de contenedor CUDA diferente ya que la imagen de contenedor CUDA se ha actualizado desde la escritura de la guía. En nuestra validación, utilizamos 'nvidia/cuda:11.0.3-base-ubuntu20.04'.

## Descargue contenedores de marcos de análisis/IA

A continuación, debe descargar las imágenes de contenedor de IA o marco de análisis necesarias de NVIDIA NGC para que estén disponibles en la máquina virtual invitada. Para descargar los contenedores de marco en el equipo virtual invitado, siga las instrucciones que se describen en ["Guía de puesta en marcha de NVIDIA AI Enterprise"](#).

## Instalación y configuración del kit de herramientas de operaciones de datos de NetApp

A continuación, debe instalar el kit de herramientas DataOPS de NetApp para entornos tradicionales en el equipo virtual invitado. El kit de herramientas DataOPS de NetApp se puede usar para gestionar los volúmenes de datos de escalado horizontal en su sistema ONTAP directamente desde el terminal dentro del equipo virtual «guest». Para instalar el kit de herramientas DataOPS de NetApp en el equipo virtual invitado, realice las siguientes tareas.

## 1. Instalar la tubería.

```
$ sudo apt update
$ sudo apt install python3-pip
$ python3 -m pip install netapp-dataops-traditional
```

## 2. Cierre la sesión en el terminal de la máquina virtual invitada y vuelva a iniciarla.

## 3. Configurar el kit de herramientas de operaciones de datos de NetApp. Para completar este paso, necesitará detalles de acceso a la API para su sistema ONTAP. Es posible que tenga que obtenerlos del administrador de almacenamiento.

```
$ netapp_dataops_cli.py config
```

```
Enter ONTAP management LIF hostname or IP address (Recommendation: Use
SVM management interface): 172.22.10.10
```

```
Enter SVM (Storage VM) name: NVAIE-client
```

```
Enter SVM NFS data LIF hostname or IP address: 172.22.13.151
```

```
Enter default volume type to use when creating new volumes
(flexgroup/flexvol) [flexgroup]:
```

```
Enter export policy to use by default when creating new volumes
[default]:
```

```
Enter snapshot policy to use by default when creating new volumes
[none]:
```

```
Enter unix filesystem user id (uid) to apply by default when creating
new volumes (ex. '0' for root user) [0]:
```

```
Enter unix filesystem group id (gid) to apply by default when creating
new volumes (ex. '0' for root group) [0]:
```

```
Enter unix filesystem permissions to apply by default when creating new
volumes (ex. '0777' for full read/write permissions for all users and
groups) [0777]:
```

```
Enter aggregate to use by default when creating new FlexVol volumes:
aff_a400_01_NVME_SSD_1
```

```
Enter ONTAP API username (Recommendation: Use SVM account): admin
```

```
Enter ONTAP API password (Recommendation: Use SVM account):
```

```
Verify SSL certificate when calling ONTAP API (true/false): false
```

```
Do you intend to use this toolkit to trigger BlueXP Copy and Sync
operations? (yes/no): no
```

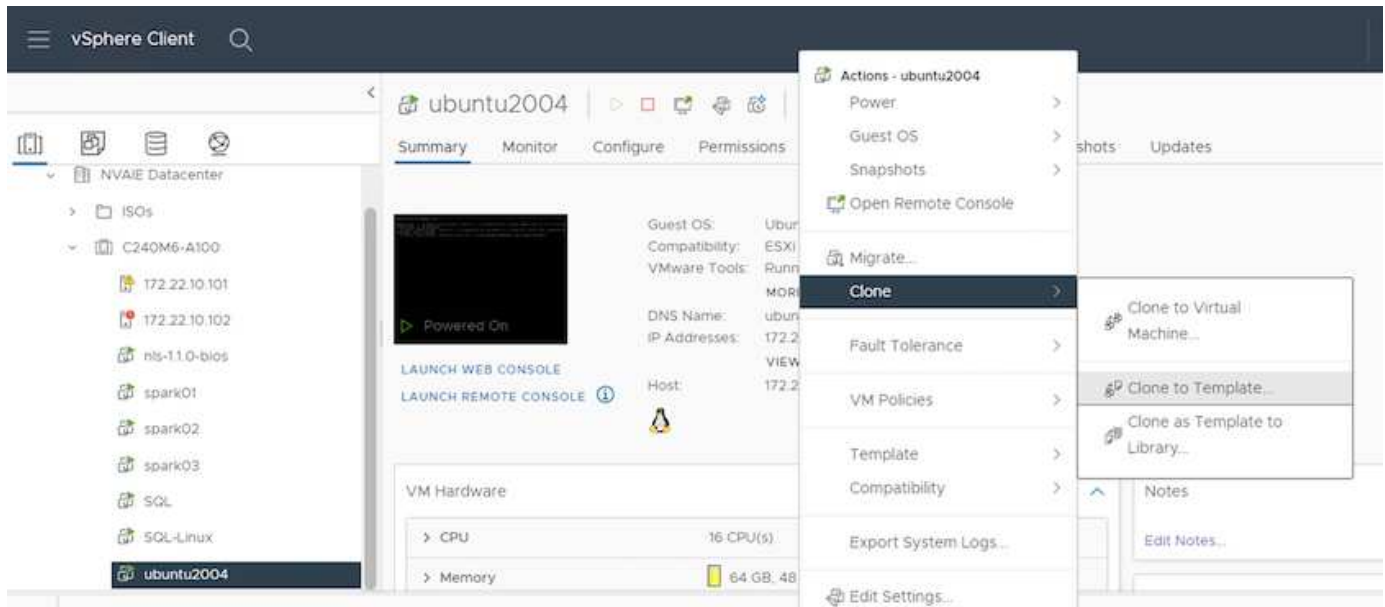
```
Do you intend to use this toolkit to push/pull from S3? (yes/no): no
```

```
Created config file: '/home/user/.netapp_dataops/config.json'.
```

## Cree una plantilla de máquina virtual invitada

Por último, debe crear una plantilla de equipo virtual basada en el equipo virtual «guest». Podrá utilizar esta plantilla para crear rápidamente equipos virtuales invitados con el software NVIDIA NGC.

Si desea crear una plantilla de máquina virtual basada en su equipo virtual invitado, inicie sesión en VMware vSphere, haga clic con el botón derecho en el nombre de la máquina virtual invitada, seleccione 'Clonar', elija 'Clonar en plantilla...' y, a continuación, siga el asistente.



#### Ejemplo de caso práctico: Trabajo de formación de TensorFlow

En esta sección se describen las tareas que se deben realizar para ejecutar una tarea de entrenamiento de TensorFlow dentro de un entorno NVIDIA AI Enterprise.

#### Requisitos previos

Antes de realizar los pasos descritos en esta sección, asumimos que ya ha creado una plantilla de equipo virtual invitado siguiendo las instrucciones que se describen en la ["Configuración"](#) página.

#### Crear máquina virtual invitada a partir de la plantilla

En primer lugar, debe crear un nuevo equipo virtual invitado a partir de la plantilla que ha creado en la sección anterior. Para crear una máquina virtual invitada nueva desde la plantilla, inicie sesión en VMware vSphere, haga clic en el nombre de la plantilla, seleccione "New VM from this Template..." y, a continuación, siga el asistente.

vSphere Client

<

vgpu-client-ubun

Summary

Monitor

Co

172.22.10.100

NVAIE Datacenter

Discovered virtual machine

vCLS

nls-1.1.0-bios

spark01

spark02

spark03

SQL

SQL-Linux

ubuntu2004

vgpu-client-ubuntu2

Guest OS:

Compatibility

VMware Tool

Actions - vgpu-client-ubuntu2004

New VM from This Template...

Convert to Virtual Machine...

Clone to Template...

Clone to Library...

Move to folder...

Rename...

Edit Notes...

Tags & Custom Attributes

Add Permission...

Alarms

Remove from Inventory

Delete from Disk

vSAN

Recent Tasks

Alarms

Task Name

Target

Delete virtual machine

Clone virtual machine

All

More Tasks

## Cree y monte el volumen de datos

A continuación, debe crear un nuevo volumen de datos en el que almacenar el conjunto de datos de entrenamiento. Puede crear rápidamente un nuevo volumen de datos con el kit de herramientas de operaciones de datos de NetApp. El siguiente comando de ejemplo muestra la creación de un volumen llamado 'imagenet' con una capacidad de 2 TB.

```
$ netapp_dataops_cli.py create vol -n imagenet -s 2TB
```

Antes de poder completar el volumen de datos con datos, debe montarlo en la máquina virtual invitada. Puede montar rápidamente un volumen de datos con el kit de herramientas Data OPS de NetApp. El comando de ejemplo siguiente muestra el bigote del volumen que se creó en el paso anterior.

```
$ sudo -E netapp_dataops_cli.py mount vol -n imagenet -m ~/imagenet
```

## Completar volumen de datos

Una vez aprovisionado y montado el volumen nuevo, el conjunto de datos de entrenamiento puede recuperarse de la ubicación de origen y colocarse en el volumen nuevo. Normalmente, esto implica extraer los datos de un lago de datos de S3 o Hadoop y a veces contará con la ayuda de un ingeniero de datos.

## Ejecutar la tarea de formación de TensorFlow

Ahora, estará listo para ejecutar su trabajo de formación de TensorFlow. Para ejecutar el trabajo de entrenamiento de TensorFlow, realice las siguientes tareas.

1. Tire de la imagen del contenedor NVIDIA NGC para TensorFlow empresarial.

```
$ sudo docker pull nvcr.io/nvaie/tensorflow-2-1:22.05-tfl-nvaie-2.1-py3
```

2. Inicie una instancia del contenedor NVIDIA NGC para TensorFlow empresarial. Utilice la opción '-v' para adjuntar el volumen de datos al contenedor.

```
$ sudo docker run --gpus all -v ~/imagenet:/imagenet -it --rm  
nvcr.io/nvaie/tensorflow-2-1:22.05-tfl-nvaie-2.1-py3
```

3. Ejecute su programa de formación TensorFlow en el contenedor. El comando de ejemplo siguiente muestra la ejecución de un programa de entrenamiento ResNet-50 de ejemplo que se incluye en la imagen contenedora.

```
$ python ./nvidia-examples/cnn/resnet.py --layers 50 -b 64 -i 200 -u  
batch --precision fp16 --data_dir /imagenet/data
```

## Dónde encontrar información adicional

Para obtener más información sobre la información descrita en este documento, consulte los siguientes documentos y/o sitios web:

- Software de gestión de datos ONTAP de NetApp: Biblioteca de información de ONTAP  
<http://mysupport.netapp.com/documentation/productlibrary/index.html?productID=62286>
- Kit de herramientas de operaciones de datos de NetApp  
["https://github.com/NetApp/netapp-dataops-toolkit"](https://github.com/NetApp/netapp-dataops-toolkit)
- NVIDIA AI Enterprise con VMware  
<https://www.nvidia.com/en-us/data-center/products/ai-enterprise/vmware/>^]

## Reconocimientos

- Bobby Oommen, Sr. Sénior, NetApp
- Ramesh Isaac, Administrador de sistemas, NetApp
- Roney Daniel, Ingeniero Técnico de Marketing de NetApp

## TR-4851: Lago de datos de StorageGRID de NetApp para cargas de trabajo de conducción autónoma - Diseño de la solución

David Arnette, NetApp

En el TR-4851 se muestra el uso del almacenamiento de objetos StorageGRID de NetApp como repositorio de datos y sistema de gestión para el desarrollo de software de aprendizaje automático (ML) y aprendizaje profundo (DL). Este documento describe el flujo de datos y los requisitos del desarrollo de software de vehículos autónomos y las funciones de StorageGRID que optimizan el ciclo de vida de los datos. Esta solución se aplica a cualquier flujo de trabajo de canalización de datos multietapa típico en los procesos de desarrollo DE APRENDIZAJE AUTOMÁTICO y aprendizaje profundo.

["TR-4851: Lago de datos de StorageGRID de NetApp para cargas de trabajo de conducción autónoma - Diseño de la solución"](#)

## Plano de control de IA de NetApp

### TR-4798: Plano de control de IA de NetApp

Mike Oglesby, NetApp

Empresas y organizaciones de todos los tamaños y sectores se están decantando por la inteligencia artificial (IA), el aprendizaje automático (ML) y el aprendizaje profundo (DL) para resolver problemas reales, ofrecer productos y servicios innovadores y obtener una ventaja en un mercado cada vez más competitivo. A medida que las organizaciones aumentan el uso de la IA, EL ML y el AP, deben hacer frente a numerosos retos, como la

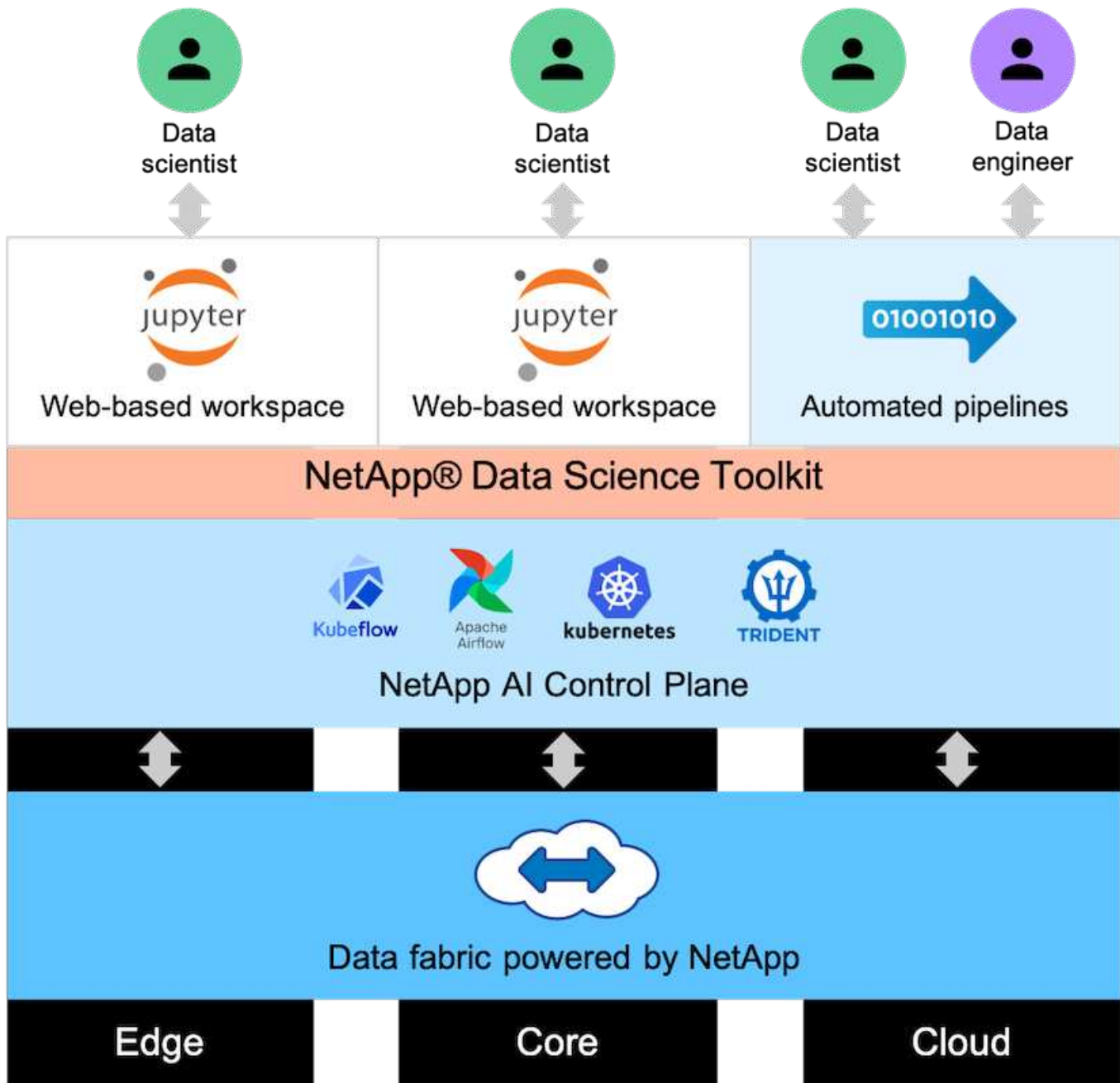
escalabilidad de la carga de trabajo y la disponibilidad de los datos. En este documento se demuestra cómo puede hacer frente a estos retos usando el plano de control de IA de NetApp, una solución que empareja las funcionalidades de gestión de datos de NetApp con marcos y herramientas de código abierto más populares.

Este informe le muestra cómo clonar rápidamente un espacio de nombres de datos. También le muestra cómo replicar datos sin problemas en sitios y regiones para crear una canalización de datos coherente y unificada de IA/ML/DL. Además, le guía por la definición e implementación de flujos de trabajo de formación de IA, ML y DL que incorporan la creación casi instantánea de datos y líneas de base de modelos para su trazabilidad y versionado. Con esta solución, puede rastrear cada entrenamiento de modelo que se ejecuta en el conjunto de datos exacto que se utilizó para entrenar o validar el modelo. Por último, este documento le muestra cómo aprovisionar rápidamente espacios de trabajo de Jupyter Notebook con acceso a conjuntos de datos masivos.

**Nota:** Para entrenamiento distribuido al estilo HPC a escala que implique un gran número de servidores GPU que requieran acceso compartido al mismo conjunto de datos o, si necesita o prefiere un sistema de archivos paralelo, consulte el documento "[TR-4890](#)". En este informe técnico se describe cómo incluir "[Solución de sistema de archivos en paralelo totalmente compatible de NetApp BeeGFS](#)". Como parte del plano de control de IA de NetApp. Esta solución está diseñada para escalar a partir de un puñado de sistemas NVIDIA DGX A100, hasta un SuperPOD de 140 nodos completo.

El plano de control de IA de NetApp se dirige a los científicos e ingenieros de datos y, por lo tanto, se necesita experiencia mínima en NetApp o en ONTAP® de NetApp. Con esta solución, las funciones de gestión de datos se pueden ejecutar utilizando interfaces y herramientas sencillas y conocidas. Si ya dispone de almacenamiento de NetApp en su entorno, puede probar el plano de control de IA de NetApp hoy mismo. Si desea probar la solución pero no dispone de almacenamiento de NetApp, visite "[cloud.netapp.com](#)". Y puede volver a tener una solución de almacenamiento de NetApp basada en cloud en cuestión de minutos. La siguiente figura ofrece una visualización de la solución.





## Conceptos y componentes

### Inteligencia artificial

La IA es una disciplina informática en la que las computadoras están entrenadas para imitar las funciones cognitivas de la mente humana. Los desarrolladores de IA entrenan computadoras para aprender y resolver problemas de una manera similar, o incluso superior a, humanos. El aprendizaje profundo y el aprendizaje automático son subcampos de la IA. Las organizaciones adoptan cada vez más IA, ML y DL para dar soporte a sus necesidades empresariales cruciales. Algunos ejemplos son los siguientes:

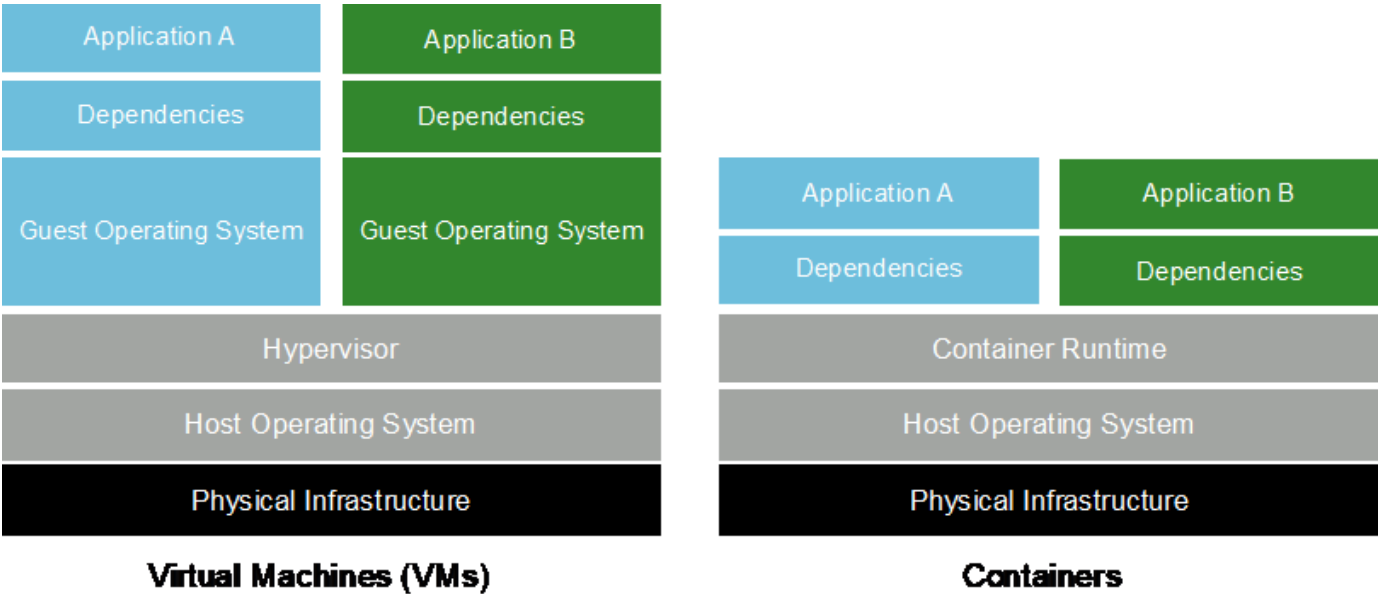
- Analizar grandes cantidades de datos para desconocer información empresarial anteriormente desconocida
- Interacción directa con los clientes mediante el procesamiento de lenguaje natural
- Automatización de diversos procesos y funciones empresariales

La formación de IA moderna y las cargas de trabajo de inferencia requieren de funcionalidades de computación en paralelo masivas. Por lo tanto, se están utilizando cada vez más GPU para ejecutar operaciones de IA, ya que las capacidades de procesamiento paralelo de las GPU son muy superiores a las de las CPU de uso general.

**Contenedores**

Los contenedores son instancias aisladas del espacio de usuario que se ejecutan sobre un kernel de sistema operativo host compartido. La adopción de contenedores aumenta rápidamente. Los contenedores ofrecen muchos de los mismos beneficios de uso de pruebas de espacio que las máquinas virtuales (VM). Sin embargo, debido a que se eliminan las capas de hipervisor y de sistema operativo «guest» de las que dependen las máquinas virtuales, los contenedores son mucho más ligeros. En la siguiente figura, se muestra una visualización de las máquinas virtuales en comparación con los contenedores.

Los contenedores también permiten el paquete eficiente de dependencias de aplicaciones, tiempos de ejecución, etc., directamente con una aplicación. El formato de embalaje de contenedor más utilizado es el contenedor Docker. Una aplicación que se haya contenedor en el formato de contenedor Docker se puede ejecutar en cualquier máquina que pueda ejecutar contenedores Docker. Esto es cierto incluso si las dependencias de la aplicación no están presentes en la máquina porque todas las dependencias están empaquetadas en el propio contenedor. Para obtener más información, visite la ["Sitio web de Docker"](#).



**Kubernetes**

Kubernetes es una plataforma de orquestación de contenedores distribuida de código abierto que originalmente diseñada por Google y que ahora se mantiene mediante Cloud Native Computing Foundation (CNCF). Kubernetes permite automatizar las funciones de puesta en marcha, gestión y escalado para aplicaciones en contenedores. En los últimos años, Kubernetes se ha convertido en la plataforma de orquestación de contenedores dominante. Aunque son compatibles con otros formatos de empaquetado y tiempos de ejecución, Kubernetes se utiliza con más frecuencia como un sistema de orquestación para contenedores Docker. Para obtener más información, visite la ["Sitio web de Kubernetes"](#).

**Trident de NetApp**

Trident es un orquestador de almacenamiento de código abierto desarrollado y mantenido por NetApp que simplifica en gran medida la creación, la gestión y el consumo de almacenamiento persistente para cargas de trabajo de Kubernetes. Trident, en sí misma una aplicación nativa de Kubernetes, se ejecuta directamente en

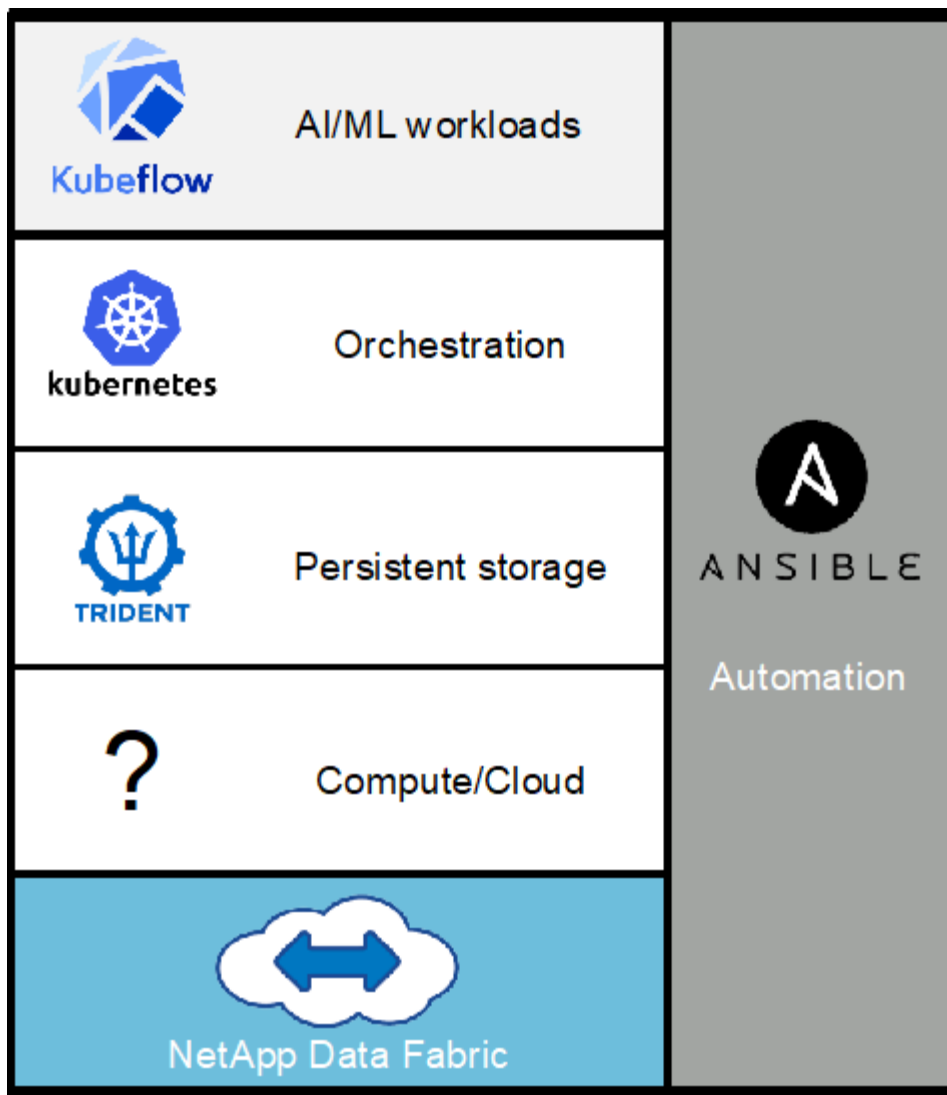
un clúster de Kubernetes. Con Trident, los usuarios de Kubernetes (desarrolladores, científicos de datos, administradores de Kubernetes, etc.) pueden crear, gestionar e interactuar con volúmenes de almacenamiento persistente en el formato Kubernetes estándar, con el que ya están familiarizados. Al mismo tiempo, pueden aprovechar las funciones avanzadas de gestión de datos de NetApp y un Data Fabric con tecnología de NetApp. Trident elimina las complejidades del almacenamiento persistente y facilita el consumo. Para obtener más información, visite la ["Sitio web de Trident"](#).

### **DeepOps de NVIDIA**

DeepOps es un proyecto de código abierto de NVIDIA que, con Ansible, automatiza la puesta en marcha de clústeres de servidores de GPU de acuerdo con las prácticas recomendadas. DeepOps es modular y se puede utilizar para realizar varias tareas de puesta en marcha. En este documento y en el ejercicio de validación descrito, DeepOps se utiliza para poner en marcha un clúster de Kubernetes que consta de nodos de trabajo de servidor GPU. Para obtener más información, visite la ["Sitio web DeepOps"](#).

### **Kubeflow**

Kubeflow es un kit de herramientas DE IA Y ML de código abierto para Kubernetes que fue desarrollado originalmente por Google. El proyecto Kubeflow hace que la puesta en marcha de flujos de trabajo de IA y ML en Kubernetes sea sencilla, portátil y escalable. —Kubeflow elimina las complejidades de Kubernetes, lo que permite a los científicos de datos centrarse en lo que saben más ciencia de datos. Consulte la siguiente figura para ver una visualización. Kubeflow ha ido ganando terreno a medida que los departamentos DE TI de las empresas se han estandarizado cada vez más en Kubernetes. Para obtener más información, visite la ["Sitio web de Kubeflow"](#).



### Tuberías de Kubeflow

Los oleoductos de Kubeflow son un componente clave de Kubeflow. Las canalizaciones de Kubeflow son una plataforma y un estándar para definir y poner en marcha flujos de trabajo DE IA Y ML escalables y portátiles. Para obtener más información, consulte "[Documentación oficial de Kubeflow](#)".

### Servidor de portátiles Jupyter

Un servidor Jupyter Notebook es una aplicación web de código abierto que permite a los científicos de datos crear documentos similares a wiki llamados portátiles Jupyter que contienen código en vivo así como pruebas descriptivas. Los portátiles Jupyter se utilizan ampliamente en la comunidad de AI Y ML como medio para documentar, almacenar y compartir proyectos de IA y ML. Kubeflow simplifica el aprovisionamiento y la puesta en marcha de servidores para portátiles Jupyter en Kubernetes. Para obtener más información sobre los Cuadernos Jupyter, visite "[Sitio Web de Jupyter](#)". Para obtener más información acerca de Jupyter Notebooks en el contexto de Kubeflow, consulte "[Documentación oficial de Kubeflow](#)".

### Flujo de aire Apache

Apache Airflow es una plataforma de gestión de flujos de trabajo de código abierto que permite la creación, programación y supervisión de programas para flujos de trabajo empresariales complejos. A menudo se utiliza para automatizar los flujos de trabajo de ETL y de canalización de datos, pero estos tipos de flujos de trabajo no se limitan a ellos. El proyecto de flujo de aire fue iniciado por Airbnb, pero desde entonces se ha vuelto

muy popular en la industria y ahora está bajo los auspicios de la Apache Software Foundation. El flujo de aire se escribe en Python, los flujos de trabajo del flujo de aire se crean a través de scripts Python y el flujo de aire está diseñado según el principio de "configuración como código". Muchos usuarios de flujo de aire empresarial ahora ejecutan el flujo de aire sobre Kubernetes.

### Gráficos de Acíclicos dirigidos (DAG)

En el flujo de aire, los flujos de trabajo se denominan gráficos Acíclicos dirigidos (DAG). Los DAG se componen de tareas que se ejecutan en secuencia, en paralelo o en una combinación de las dos, dependiendo de la definición DAG. El programador de flujo de aire ejecuta tareas individuales en una matriz de trabajadores y cumple con las dependencias a nivel de tarea especificadas en la definición DAG. Los DAG se definen y crean a través de scripts Python.

### ONTAP 9 de NetApp

ONTAP 9 de NetApp es la última generación del software de gestión del almacenamiento de NetApp y permite a empresas como la suya modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. Gracias a las capacidades de gestión de datos líderes del sector, ONTAP le permite gestionar y proteger sus datos con un solo conjunto de herramientas sin importar dónde residan. También puede mover los datos libremente a donde necesite: El extremo, el núcleo o el cloud. ONTAP 9 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen sus datos esenciales y preparan su infraestructura para el futuro con arquitecturas de cloud híbrido.

### Simplifique la gestión de los datos

La gestión de datos es vital para las operaciones TECNOLÓGICAS de su empresa, de modo que pueda usar los recursos adecuados para sus aplicaciones y conjuntos de datos. ONTAP incluye las siguientes funciones para facilitar y simplificar las operaciones, y reducir el coste total de funcionamiento:

- **Compactación de datos inline y deduplicación expandida.** la compactación de datos reduce el espacio perdido dentro de los bloques de almacenamiento, y la deduplicación aumenta significativamente la capacidad efectiva.
- **Calidad de servicio (QoS) mínima, máxima y adaptativa.** los controles granulares de QoS ayudan a mantener los niveles de rendimiento para aplicaciones críticas en entornos altamente compartidos.
- **ONTAP FabricPool.** esta función ofrece una organización automática en niveles de datos inactivos para opciones de almacenamiento en cloud público y privado, incluidos Amazon Web Services (AWS), Azure y almacenamiento basado en objetos StorageGRID de NetApp.

### Acelere y proteja sus datos

ONTAP no solo ofrece niveles de rendimiento y protección de datos superiores, sino que amplía estas funcionalidades con las siguientes funciones:

- **Alto rendimiento y baja latencia.** ONTAP ofrece el rendimiento más alto posible con la menor latencia posible.
- **La tecnología ONTAP FlexGroup de NetApp.** Un volumen FlexGroup es un contenedor de datos de alto rendimiento que se puede escalar linealmente hasta 20 PB y 400 000 millones de archivos, lo que proporciona un espacio de nombres único que simplifica la gestión de datos.
- **Protección de datos.** ONTAP ofrece capacidades integradas de protección de datos con administración común en todas las plataformas.
- **Cifrado de volumen de NetApp.** ONTAP ofrece cifrado nativo a nivel de volumen con compatibilidad para la gestión de claves incorporada y externa.

## Infraestructura preparada para futuros retos

ONTAP 9 le ayuda a satisfacer las exigentes y siempre cambiantes necesidades de su empresa:

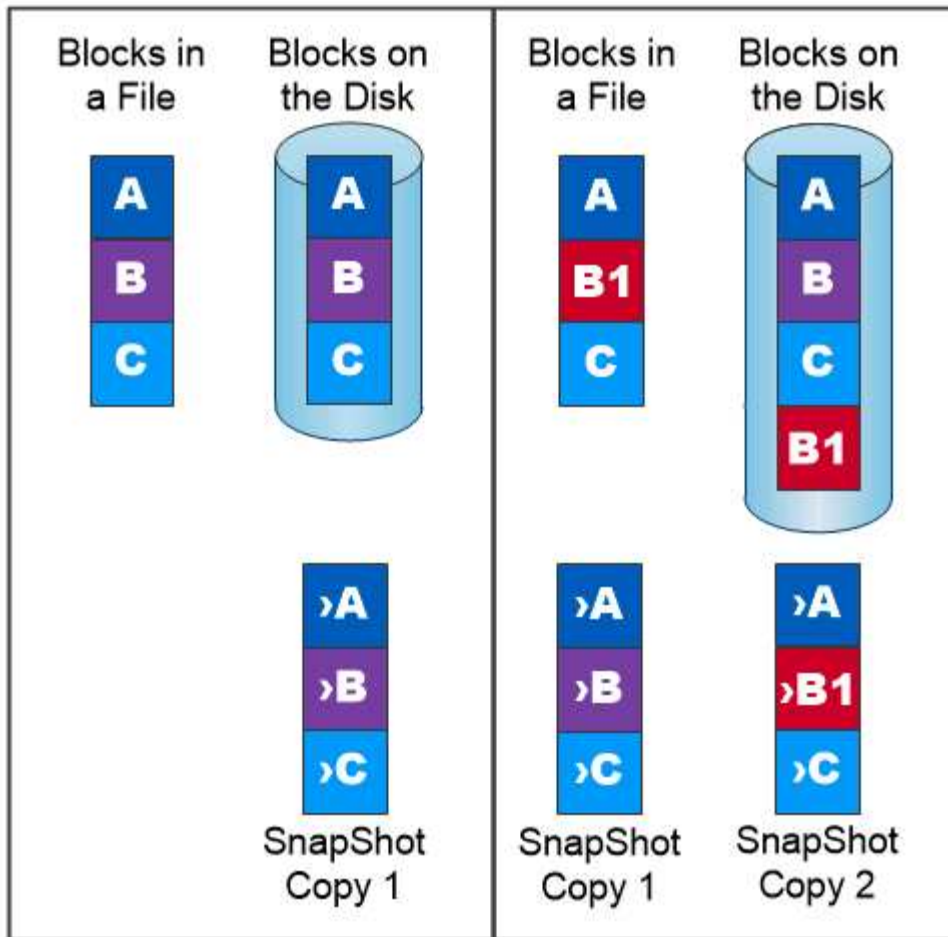
- **Escalado sencillo y operaciones no disruptivas.** ONTAP admite la adición sin interrupciones de capacidad a las controladoras existentes y a los clústeres de escalado horizontal. Puede empezar a utilizar tecnologías punteras como NVMe y FC 32 GB, sin necesidad de realizar costosas migraciones de datos y sin cortes.
- **Conexión al cloud.** ONTAP es uno de los programas de gestión del almacenamiento con mejor conexión al cloud e incluye opciones de almacenamiento definido por software (ONTAP Select) e instancias nativas del cloud (Cloud Volumes Service de NetApp) en todos los clouds públicos.
- **Integración con aplicaciones emergentes.** al utilizar la misma infraestructura que soporta las aplicaciones empresariales existentes, ONTAP ofrece servicios de datos de clase empresarial para plataformas y aplicaciones de próxima generación como OpenStack, Hadoop y MongoDB.

### Copias Snapshot de NetApp

Una copia Snapshot de NetApp es una imagen puntual de solo lectura de un volumen. La imagen consume un espacio de almacenamiento mínimo y tiene una sobrecarga del rendimiento mínima, ya que solo registra los cambios que se han realizado en los archivos creados desde que se realizó la última copia Snapshot, como se muestra en la siguiente figura.

Las copias Snapshot deben su eficiencia a la tecnología de virtualización del almacenamiento central de ONTAP, el sistema de archivos de escritura en cualquier lugar (WAFL). Al igual que una base de datos, WAFL utiliza metadatos para apuntar a los bloques de datos reales en el disco. Sin embargo, a diferencia de una base de datos, WAFL no sobrescribe los bloques existentes. Escribe los datos actualizados en un bloque nuevo y cambia los metadatos. Porque ONTAP hace referencia a los metadatos cuando crea una copia Snapshot, en lugar de copiar bloques de datos, es tan eficiente que las copias Snapshot. Al hacerlo, se elimina el tiempo de búsqueda que otros sistemas incurren en la localización de los bloques a copiar, así como el costo de hacer la copia misma.

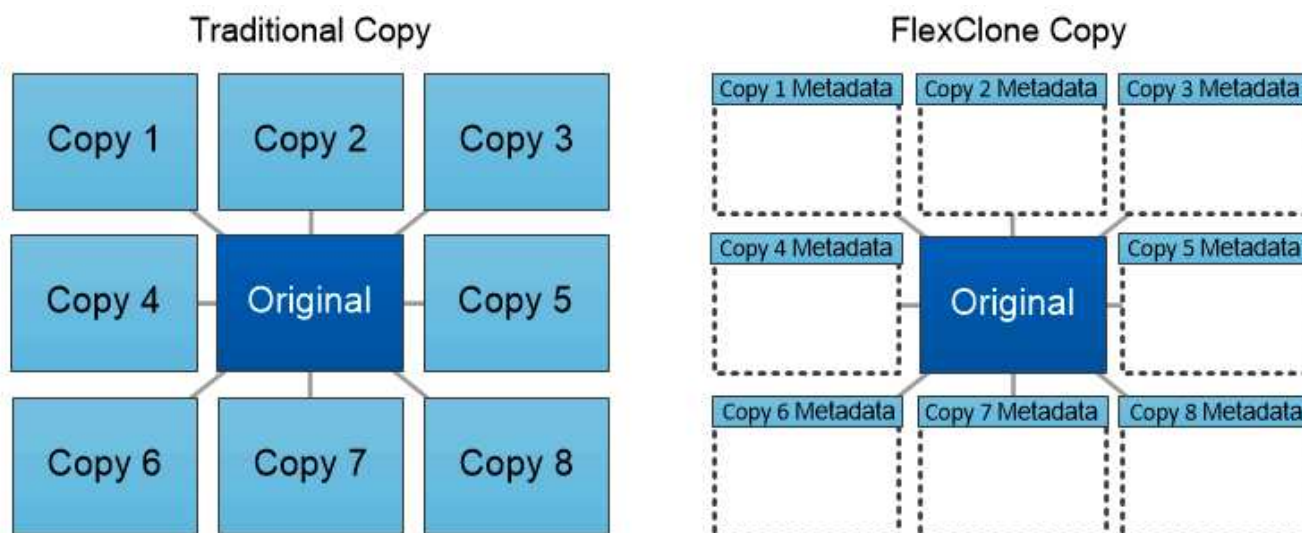
Puede utilizar una copia Snapshot para recuperar archivos o LUN individuales o para restaurar el contenido completo de un volumen. ONTAP compara la información de punteros de la copia Snapshot con los datos del disco para reconstruir el objeto faltante o dañado, sin tiempo de inactividad ni un coste de rendimiento significativo.



*A Snapshot copy records only changes to the active file system since the last Snapshot copy.*

#### Tecnología FlexClone de NetApp

La tecnología FlexClone de NetApp hace referencia a los metadatos de Snapshot para crear copias puntuales editables de un volumen. Las copias comparten bloques de datos con sus padres, sin consumir almacenamiento excepto lo que se necesita para los metadatos hasta que se escriben los cambios en la copia, como se muestra en la siguiente figura. Cuando se pueden crear copias tradicionales en minutos o incluso horas, el software FlexClone le permite copiar incluso los conjuntos de datos más grandes de forma casi instantánea. Esto lo convierte en la opción ideal para las situaciones en las que necesita varias copias de conjuntos de datos idénticos (un espacio de trabajo de desarrollo, por ejemplo) o copias temporales de un conjunto de datos (probar una aplicación contra un conjunto de datos de producción).



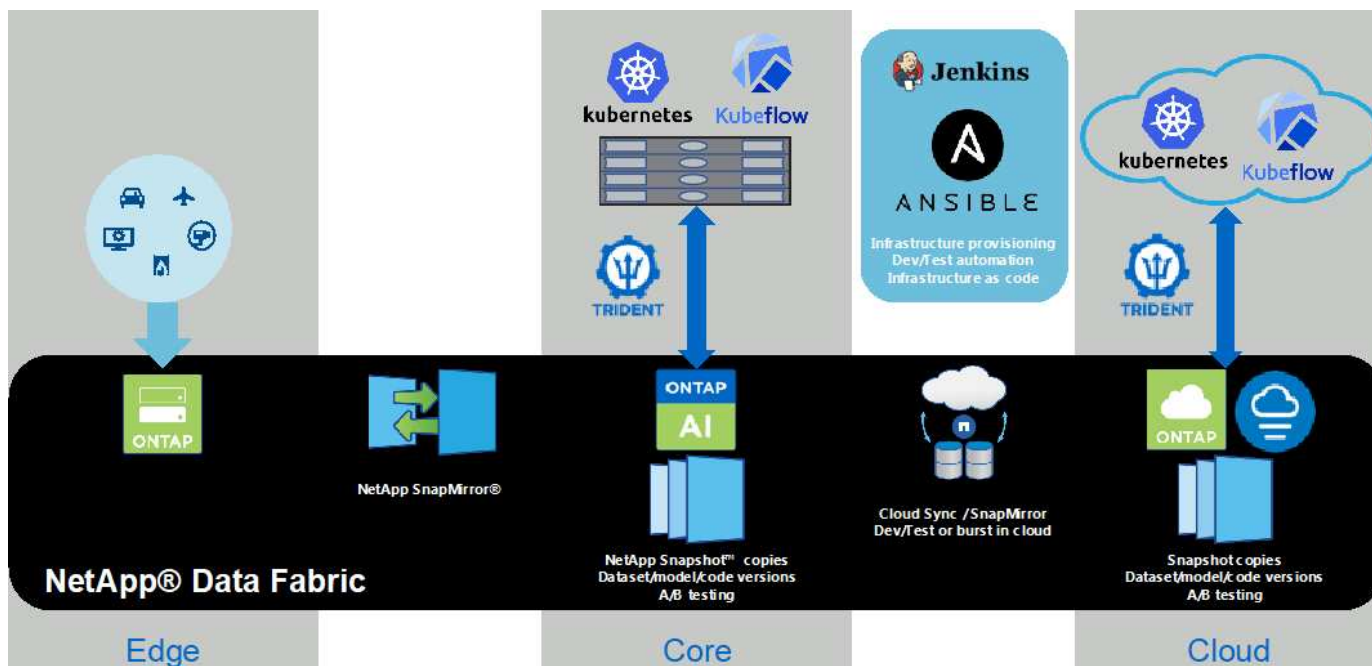
*FlexClone copies share data blocks with their parents, consuming no storage except what is required for metadata.*

#### Tecnología de replicación de datos de SnapMirror de NetApp

El software SnapMirror de NetApp es una solución de replicación unificada rentable y fácil de usar para todo Data Fabric. Replica datos a altas velocidades mediante LAN o WAN. Le proporciona una alta disponibilidad de datos y una rápida replicación de datos para todo tipo de aplicaciones, incluidas aplicaciones vitales para el negocio en entornos tanto virtuales como tradicionales. Al replicar datos en uno o varios sistemas de almacenamiento de NetApp y actualizar continuamente los datos secundarios, estos están siempre al día y disponibles cuando los necesite. No se requieren servidores de replicación externos. Consulte la figura siguiente para ver un ejemplo de una arquitectura que aprovecha la tecnología SnapMirror.

El software SnapMirror aprovecha las eficiencias del almacenamiento de ONTAP de NetApp y envía únicamente los bloques cambiados a través de la red. El software SnapMirror también usa la compresión de red incorporada para acelerar las transferencias de datos y reducir la utilización de ancho de banda hasta un 70 %. Con la tecnología SnapMirror, puede aprovechar un flujo de datos de thin replication para crear un único almacén que mantenga los reflejos activos y las copias de momentos específicos anteriores, lo que reduce el tráfico de red hasta un 50 %.





### Copia y sincronización de NetApp BlueXP

La copia y sincronización de BlueXP es un servicio de NetApp que ofrece una sincronización de datos rápida y segura. Ya tenga que transferir archivos entre recursos compartidos de archivos NFS o SMB en las instalaciones, NetApp StorageGRID, NetApp ONTAP S3, NetApp Cloud Volumes Service, Azure NetApp Files, AWS S3, AWS EFS, Azure Blob, Google Cloud Storage, o IBM Cloud Object Storage, BlueXP Copy and Sync mueve los archivos a donde los necesitas de forma rápida y segura.

Una vez transferidos los datos, estarán completamente disponibles para su uso tanto en origen como en destino. BlueXP Copy and Sync puede sincronizar los datos bajo demanda al activar una actualización o sincronizar continuamente los datos en función de una programación predefinida. Independientemente de ello, BlueXP Copy y Sync solo mueve los diferenciales, por lo que se reducen al mínimo el tiempo y el dinero que se invierten en la replicación de datos.

Copia y sincronización de BlueXP es una herramienta de software como servicio (SaaS) extremadamente sencilla de configurar y usar. Las transferencias de datos activadas por BlueXP Copy and Sync se llevan a cabo por agentes de datos. Los agentes de datos de BlueXP Copy y Sync se pueden poner en marcha en AWS, Azure, Google Cloud Platform o en las instalaciones.

### XCP de NetApp

XCP de NetApp es el software basado en cliente para migraciones de datos y análisis del sistema de archivos entre NetApp y NetApp. XCP se ha diseñado para escalar y lograr el máximo rendimiento utilizando todos los recursos del sistema disponibles para gestionar conjuntos de datos de gran volumen y migraciones de alto rendimiento. XCP le ayuda a obtener una visibilidad completa del sistema de archivos con la opción de generar informes.

XCP de NetApp está disponible en un único paquete compatible con los protocolos NFS y SMB. XCP incluye un binario de Linux para conjuntos de datos NFS y un ejecutable de Windows para conjuntos de datos SMB.

XCP File Analytics de NetApp es un software basado en host que detecta recursos compartidos de archivos, ejecuta análisis en el sistema de archivos y proporciona una consola para el análisis de archivos. XCP File Analytics es compatible con los sistemas NetApp y de otros proveedores, y se ejecuta en hosts Linux o Windows para proporcionar análisis en sistemas de archivos NFS y exportados SMB.

## ONTAP FlexGroup Volumes de NetApp

Un conjunto de datos de entrenamiento puede ser una colección con hasta miles de millones de archivos. Pueden ser archivos de texto, de audio, de vídeo o cualquier otra forma de datos no estructurados que deban almacenarse y procesarse para su lectura en paralelo. El sistema de almacenamiento debe almacenar un gran número de archivos pequeños y debe leerlos en paralelo, con una entrada y salida secuencial o aleatoria.

Un volumen FlexGroup es un espacio de nombres único que comprende varios volúmenes miembro constituyentes, tal y como se muestra en la siguiente figura. Desde el punto de vista de un administrador de almacenamiento, un volumen FlexGroup se gestiona y actúa como un volumen FlexVol de NetApp. Los archivos de un volumen de FlexGroup se asignan a volúmenes miembro individuales y no están repartidos en volúmenes o nodos. Ofrecen las siguientes capacidades:

- Los volúmenes FlexGroup proporcionan varios petabytes de capacidad y una baja latencia predecible para cargas de trabajo con una gran cantidad de metadatos.
- Permiten un máximo de 400 000 millones de archivos en un mismo espacio de nombres.
- Admiten operaciones en paralelo para cargas de trabajo NAS entre varias CPU, nodos, agregados y volúmenes FlexVol constituyentes.



## Requisitos de hardware y software

La solución de plano de control de IA de NetApp no depende de este hardware específico. La solución es compatible con cualquier dispositivo de almacenamiento físico, instancia definida por software o servicio cloud de NetApp compatible con Trident. Entre los ejemplos se incluyen un sistema de almacenamiento AFF de NetApp, Azure NetApp Files, Cloud Volumes Service de NetApp, una instancia de almacenamiento definido por software ONTAP Select de NetApp o una instancia de Cloud Volumes ONTAP de NetApp. Además, la solución puede implementarse en cualquier clúster de Kubernetes, siempre que la versión de Kubernetes utilizada sea compatible con Kubeflow y Trident de NetApp. Si desea ver una lista de las versiones de Kubernetes compatibles con Kubeflow, consulte la ["Documentación oficial de Kubeflow"](#). Si desea ver una lista de las versiones de Kubernetes compatibles con Trident, consulte ["Documentación de Trident"](#). Consulte

las siguientes tablas para obtener información detallada sobre el entorno que se utilizó para validar la solución.

Componente de la infraestructura	Cantidad	Detalles	Sistema operativo
Host de salto de implementación	1	MÁQUINA VIRTUAL	Sistema operativo Ubuntu 20.04.2 LTS
Nodos maestros de Kubernetes	1	MÁQUINA VIRTUAL	Sistema operativo Ubuntu 20.04.2 LTS
Nodos de trabajo de Kubernetes	2	MÁQUINA VIRTUAL	Sistema operativo Ubuntu 20.04.2 LTS
Nodos de trabajo de GPU de Kubernetes	2	NVIDIA DGX-1 (nativo)	NVIDIA DGX OS 4.0.5 (basado en Ubuntu 18.04.2 LTS)
Reducida	1 par de alta disponibilidad	AFF A220 de NetApp	ONTAP 9.7 P6 de NetApp

Componente de software	Versión
Flujo de aire Apache	2.0.1
Tabla de Helm para flujo de aire de Apache	8.0.8
Docker	19.03.12
Kubeflow	1.2
Kubernetes	1.18.9
Trident de NetApp	21.01.2
DeepOps de NVIDIA	Funcionalidad de implementación de Trident desde la rama maestra hasta el encargo "61898cdfda"; Todas las demás funciones de la versión 21.03

#### Soporte técnico

NetApp no ofrece compatibilidad empresarial con Apache Airflow, Docker, Kubeflow, Kubernetes u NVIDIA DeepOps. Si está interesado en una solución totalmente compatible con funcionalidades similares a la solución del plano de control de IA de NetApp, ["Póngase en contacto con NetApp"](#) Acerca de soluciones de IA/ML totalmente compatibles que NetApp ofrece en colaboración con partners.

#### Puesta en marcha de Kubernetes

En esta sección se describen las tareas que debe completar para poner en marcha un clúster de Kubernetes en el que desea implementar la solución de plano de control de IA de NetApp. Si ya tiene un clúster de Kubernetes, puede omitir esta sección siempre que ejecute una versión de Kubernetes compatible con Kubeflow y NetApp Trident. Si desea ver una lista de las versiones de Kubernetes compatibles con ubeflow, consulte la ["Documentación oficial de Kubeflow"](#). Si desea ver una lista de las versiones de Kubernetes compatibles con Trident, consulte ["Documentación de Trident"](#).

En el caso de puestas en marcha de Kubernetes en las instalaciones que incorporan nodos de reconstrucción completa con GPU de NVIDIA, NetApp recomienda usar la herramienta de puesta en marcha de Kubernetes DeepOps de NVIDIA. En esta sección, se describe la puesta en marcha de un clúster de Kubernetes con DeepOps.

### Requisitos previos

Antes de realizar el ejercicio de implementación descrito en esta sección, asumimos que ya ha realizado las siguientes tareas:

1. Ya ha configurado todos los nodos Kubernetes con configuración básica (por ejemplo, un sistema NVIDIA DGX que forma parte de un pod ONTAP AI), de acuerdo con instrucciones de configuración estándar.
2. Ha instalado un sistema operativo compatible en todos los nodos maestro y de trabajo de Kubernetes y en un host de salto de implementación. Para obtener una lista de los sistemas operativos compatibles con DeepOps, consulte ["Sitio DeepOps GitHub"](#).

### Utilice NVIDIA DeepOps para instalar y configurar Kubernetes

Para poner en marcha y configurar su clúster de Kubernetes con NVIDIA DeepOPS, realice las siguientes tareas desde un host de salto de implementación:

1. Descargue NVIDIA DeepOps siguiendo las instrucciones de ["Página de Inicio"](#) En el sitio de NVIDIA DeepOps GitHub.
2. Implemente Kubernetes en su clúster siguiendo las instrucciones que se indican en ["Página de la Guía de implementación de Kubernetes"](#) En el sitio de NVIDIA DeepOps GitHub.

### Implementación y configuración de Trident de NetApp

#### Implementación y configuración de Trident de NetApp

En esta sección se describen las tareas que debe completar para instalar y configurar NetApp Trident en su clúster Kubernetes.

### Requisitos previos

Antes de realizar el ejercicio de implementación descrito en esta sección, asumimos que ya ha realizado las siguientes tareas:

1. Ya existe un clúster de Kubernetes en funcionamiento y ejecuta una versión de Kubernetes compatible con Trident. Para obtener una lista de las versiones compatibles, consulte ["Documentación de Trident"](#).
2. Ya existe un dispositivo de almacenamiento de NetApp, una instancia definida por software o un servicio de almacenamiento en cloud que Trident admite.

### Instale Trident

Para instalar y configurar NetApp Trident en su clúster de Kubernetes, realice las siguientes tareas desde el host de salto de implementación:

1. Implemente Trident mediante uno de los siguientes métodos:
  - Si utilizó NVIDIA DeepOps para poner en marcha su clúster de Kubernetes, también puede usar NVIDIA DeepOps para implementar Trident en el clúster de Kubernetes. Para poner en marcha Trident con DeepOps, siga ["Instrucciones de puesta en marcha de Trident"](#) En el sitio de NVIDIA DeepOps

GitHub.

- Si no utilizó NVIDIA DeepOPS para poner en marcha su clúster de Kubernetes, o si simplemente prefiere poner en marcha Trident manualmente, puede implementar Trident siguiendo el ["instrucciones de puesta en funcionamiento"](#) En la documentación de Trident. Asegúrese de crear al menos un back-end de Trident y al menos un Kubernetes StorageClass para obtener más información acerca de cómo configurar ["Back-ends"](#) y.. ["Clases de almacenamiento"](#) Consulte los subapartados vinculados en los documentos de NetApp.



Si va a poner en marcha la solución del plano de control de IA de NetApp en un POD de IA de ONTAP, consulte ["Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI"](#) Para algunos ejemplos de diferentes Back-ends de Trident que es posible que desee crear y. ["Ejemplo de una historia Kubernetes para las puestas en marcha de IA de ONTAP"](#) Puede que desee crear algunos ejemplos de diferentes clases de almacenamiento de Kubernetes.

### Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI

Antes de poder usar Trident para aprovisionar recursos de almacenamiento de forma dinámica dentro de su clúster de Kubernetes, debe crear uno o varios Back-ends de Trident. Los siguientes ejemplos representan diferentes tipos de backends que puede crear si va a implementar la solución de plano de control de IA de NetApp en un pod ONTAP AI. Para obtener más información acerca de backends, consulte ["Documentación de Trident"](#).

1. NetApp recomienda crear un Back-end de Trident habilitado para FlexGroup para cada LIF de datos (interfaz de red lógica que proporciona acceso a los datos) que desee utilizar en su sistema AFF de NetApp. Esto le permitirá equilibrar los montajes de volumen entre LIF

Los siguientes comandos de ejemplo muestran la creación de dos Back-ends de Trident habilitados para FlexGroup para dos LIF de datos diferentes asociadas con la misma máquina virtual de almacenamiento (SVM) de ONTAP. Estas backends utilizan la `ontap-nas-flexgroup` controlador de almacenamiento. ONTAP admite dos tipos de volúmenes de datos principales: FlexVol y FlexGroup. Los volúmenes FlexVol tienen un tamaño limitado (a partir de la escritura, el tamaño máximo depende de la implementación específica). Por otro lado, los volúmenes FlexGroup se pueden escalar de forma lineal hasta 20 PB y 400 000 millones de archivos y, además, ofrecen un espacio de nombres único que simplifica enormemente la gestión de los datos. Por lo tanto, los volúmenes FlexGroup son óptimos para cargas de trabajo de IA y ML que dependen de grandes cantidades de datos.

Si está trabajando con una pequeña cantidad de datos y desea usar volúmenes de FlexVol en lugar de volúmenes de FlexGroup, puede crear Back-ends de Trident que utilizan `ontap-nas` controlador de almacenamiento en lugar de `ontap-nas-flexgroup` controlador de almacenamiento.

```
$ cat << EOF > ./trident-backend-ontap-ai-flexgroups-ifacel.json
{
  "version": 1,
  "storageDriverName": "ontap-nas-flexgroup",
  "backendName": "ontap-ai-flexgroups-ifacel",
  "managementLIF": "10.61.218.100",
  "dataLIF": "192.168.11.11",
  "svm": "ontapai_nfs",
```

```

    "username": "admin",
    "password": "ontapai"
}
EOF
$ tridentctl create backend -f ./trident-backend-ontap-ai-flexgroups-
iface1.json -n trident
+-----+-----+
+-----+-----+
|          NAME          | STORAGE DRIVER |
UUID                   | STATE  | VOLUMES |
+-----+-----+
+-----+-----+
| ontap-ai-flexgroups-iface1 | ontap-nas-flexgroup | b74cbddb-e0b8-40b7-
b263-b6da6dec0bdd | online |      0 |
+-----+-----+
+-----+-----+
$ cat << EOF > ./trident-backend-ontap-ai-flexgroups-iface2.json
{
    "version": 1,
    "storageDriverName": "ontap-nas-flexgroup",
    "backendName": "ontap-ai-flexgroups-iface2",
    "managementLIF": "10.61.218.100",
    "dataLIF": "192.168.12.12",
    "svm": "ontapai_nfs",
    "username": "admin",
    "password": "ontapai"
}
EOF
$ tridentctl create backend -f ./trident-backend-ontap-ai-flexgroups-
iface2.json -n trident
+-----+-----+
+-----+-----+
|          NAME          | STORAGE DRIVER |
UUID                   | STATE  | VOLUMES |
+-----+-----+
+-----+-----+
| ontap-ai-flexgroups-iface2 | ontap-nas-flexgroup | 61814d48-c770-436b-
9cb4-cf7ee661274d | online |      0 |
+-----+-----+
+-----+-----+
$ tridentctl get backend -n trident
+-----+-----+
+-----+-----+
|          NAME          | STORAGE DRIVER |
UUID                   | STATE  | VOLUMES |
+-----+-----+

```

```

+-----+-----+-----+
| ontap-ai-flexgroups-iface1 | ontap-nas-flexgroup | b74cbddb-e0b8-40b7-
b263-b6da6dec0bdd | online |          0 |
| ontap-ai-flexgroups-iface2 | ontap-nas-flexgroup | 61814d48-c770-436b-
9cb4-cf7ee661274d | online |          0 |
+-----+-----+-----+
+-----+-----+-----+

```

2. NetApp también recomienda crear uno o varios Back-ends de Trident habilitados para FlexVol. Si se utilizan FlexGroup Volumes para el almacenamiento de conjuntos de datos, es posible que se desee usar volúmenes de FlexVol para almacenar resultados, resultados, información de depuración, etc. Si se desean usar volúmenes de FlexVol, se deben crear uno o varios Back-ends de Trident habilitados para FlexVol. Los siguientes comandos de ejemplo muestran la creación de una única estructura de fondo Trident habilitada para FlexVol que utiliza una sola LIF de datos.

```
$ cat << EOF > ./trident-backend-ontap-ai-flexvols.json
{
  "version": 1,
  "storageDriverName": "ontap-nas",
  "backendName": "ontap-ai-flexvols",
  "managementLIF": "10.61.218.100",
  "dataLIF": "192.168.11.11",
  "svm": "ontapai_nfs",
  "username": "admin",
  "password": "ontapai"
}
EOF
$ tridentctl create backend -f ./trident-backend-ontap-ai-flexvols.json -n
trident
+-----+-----+-----+
+-----+-----+-----+
|          NAME          | STORAGE DRIVER |          UUID          |
| STATE | VOLUMES | |
+-----+-----+-----+
+-----+-----+-----+
| ontap-ai-flexvols      | ontap-nas      | 52bdb3b1-13a5-4513-
a9c1-52a69657fabe | online |      0 |
+-----+-----+-----+
+-----+-----+-----+
$ tridentctl get backend -n trident
+-----+-----+-----+
+-----+-----+-----+
|          NAME          | STORAGE DRIVER |          UUID          |
| STATE | VOLUMES | |
+-----+-----+-----+
+-----+-----+-----+
| ontap-ai-flexvols      | ontap-nas      | 52bdb3b1-13a5-4513-
a9c1-52a69657fabe | online |      0 |
| ontap-ai-flexgroups-iface1 | ontap-nas-flexgroup | b74cbddb-e0b8-40b7-
b263-b6da6dec0bdd | online |      0 |
| ontap-ai-flexgroups-iface2 | ontap-nas-flexgroup | 61814d48-c770-436b-
9cb4-cf7ee661274d | online |      0 |
+-----+-----+-----+
+-----+-----+-----+
```

#### Ejemplo de clases de almacenamiento Kubernetes para puestas en marcha de IA en ONTAP

Antes de poder usar Trident para aprovisionar recursos de almacenamiento de forma dinámica dentro del clúster de Kubernetes, debe crear una o varias clases de almacenamiento de Kubernetes. Los siguientes ejemplos representan diferentes tipos de



clases de almacenamiento que puede que desee crear si pone en marcha la solución de plano de control de IA de NetApp en un pool de IA de ONTAP. Si desea obtener más información sobre las clases de almacenamiento, consulte ["Documentación de Trident"](#).

1. NetApp recomienda crear un tipo de almacenamiento aparte para cada back-end de Trident habilitado para FlexGroup que haya creado en la sección ["Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI"](#), paso 1. Estas clases de almacenamiento granular le permiten añadir montajes NFS que corresponden a LIF específicas (los LIF especificados cuando se crearon las back-ends de Trident) como un back-end concreto especificado en el archivo SPEC de StorageClass. Los comandos de ejemplo siguientes muestran la creación de dos StorageClasses que corresponden a los dos backends de ejemplo que se crearon en la sección ["Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI"](#), paso 1. Si desea obtener más información sobre las clases de almacenamiento, consulte ["Documentación de Trident"](#).

Para que no se elimine un volumen persistente cuando se elimine la reclamación de volumen persistente (RVP) correspondiente, en el siguiente ejemplo se utiliza un `reclaimPolicy` valor de `Retain`. Para obtener más información acerca de `reclaimPolicy` consulte el funcionario ["Documentación de Kubernetes"](#).

```

$ cat << EOF > ./storage-class-ontap-ai-flexgroups-retain-iface1.yaml
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: ontap-ai-flexgroups-retain-iface1
provisioner: netapp.io/trident
parameters:
  backendType: "ontap-nas-flexgroup"
  storagePools: "ontap-ai-flexgroups-iface1:.*"
reclaimPolicy: Retain
EOF
$ kubectl create -f ./storage-class-ontap-ai-flexgroups-retain-
iface1.yaml
storageclass.storage.k8s.io/ontap-ai-flexgroups-retain-iface1 created
$ cat << EOF > ./storage-class-ontap-ai-flexgroups-retain-iface2.yaml
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: ontap-ai-flexgroups-retain-iface2
provisioner: netapp.io/trident
parameters:
  backendType: "ontap-nas-flexgroup"
  storagePools: "ontap-ai-flexgroups-iface2:.*"
reclaimPolicy: Retain
EOF
$ kubectl create -f ./storage-class-ontap-ai-flexgroups-retain-
iface2.yaml
storageclass.storage.k8s.io/ontap-ai-flexgroups-retain-iface2 created
$ kubectl get storageclass

```

NAME	PROVISIONER	AGE
ontap-ai-flexgroups-retain-iface1	netapp.io/trident	0m
ontap-ai-flexgroups-retain-iface2	netapp.io/trident	0m

2. NetApp también recomienda crear un StorageClass que se corresponda con el back-end Trident habilitado para FlexVol que ha creado en la sección ["Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI"](#), paso 2. Los comandos de ejemplo siguientes muestran la creación de un solo tipo de almacenamiento para volúmenes FlexVol.

En el siguiente ejemplo, no se ha especificado un back-end determinado en el archivo de definición StorageClass porque solo se creó un back-end Trident habilitado para FlexVol. Cuando se usa Kubernetes para administrar volúmenes que usan este clase de almacenamiento, Trident intenta usar cualquier back-end disponible que utilice ontap-nas controlador.

```

$ cat << EOF > ./storage-class-ontap-ai-flexvols-retain.yaml
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: ontap-ai-flexvols-retain
provisioner: netapp.io/trident
parameters:
  backendType: "ontap-nas"
reclaimPolicy: Retain
EOF
$ kubectl create -f ./storage-class-ontap-ai-flexvols-retain.yaml
storageclass.storage.k8s.io/ontap-ai-flexvols-retain created
$ kubectl get storageclass

```

NAME	PROVISIONER	AGE
ontap-ai-flexgroups-retain-iface1	netapp.io/trident	1m
ontap-ai-flexgroups-retain-iface2	netapp.io/trident	1m
ontap-ai-flexvols-retain	netapp.io/trident	0m

3. NetApp también recomienda crear un tipo de almacenamiento genérico para volúmenes FlexGroup. En los siguientes comandos de ejemplo, se muestra la creación de un solo tipo de almacenamiento genérico para volúmenes FlexGroup.

Tenga en cuenta que no se ha especificado un back-end determinado en el archivo de definición StorageClass. Por lo tanto, cuando se usa Kubernetes para administrar volúmenes que usan esta clase de almacenamiento, Trident intenta usar cualquier back-end disponible que utilice ontap-nas-flexgroup controlador.

```
$ cat << EOF > ./storage-class-ontap-ai-flexgroups-retain.yaml
apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: ontap-ai-flexgroups-retain
provisioner: netapp.io/trident
parameters:
  backendType: "ontap-nas-flexgroup"
reclaimPolicy: Retain
EOF
$ kubectl create -f ./storage-class-ontap-ai-flexgroups-retain.yaml
storageclass.storage.k8s.io/ontap-ai-flexgroups-retain created
$ kubectl get storageclass
```

NAME	PROVISIONER	AGE
ontap-ai-flexgroups-retain	netapp.io/trident	0m
ontap-ai-flexgroups-retain-iface1	netapp.io/trident	2m
ontap-ai-flexgroups-retain-iface2	netapp.io/trident	2m
ontap-ai-flexvols-retain	netapp.io/trident	1m

## Despliegue de Kubeflow

En esta sección se describen las tareas que debe completar para poner en marcha Kubeflow en su clúster de Kubernetes.

### Requisitos previos

Antes de realizar el ejercicio de implementación descrito en esta sección, asumimos que ya ha realizado las siguientes tareas:

1. Ya tiene un clúster de Kubernetes en funcionamiento y ejecuta una versión de Kubernetes que admite ubeflow. Para obtener una lista de las versiones compatibles, consulte ["Documentación oficial de Kubeflow"](#).
2. Ya ha instalado y configurado NetApp Trident en su clúster de Kubernetes, como se indica en ["Implementación y configuración de Trident"](#).

### Establezca el tipo de almacenamiento de Kubernetes predeterminado

Antes de poner en marcha Kubeflow, debe designar un clase de almacenamiento predeterminado dentro del clúster de Kubernetes. El proceso de implementación de Kubeflow intenta aprovisionar nuevos volúmenes persistentes mediante el tipo de almacenamiento predeterminado. Si no se designa StorageClass como clase de almacenamiento predeterminado, la implementación falla. Para designar un StorageClass predeterminado en el clúster, realice la siguiente tarea desde el host de salto de implementación. Si ya ha designado un tipo de almacenamiento predeterminado en el clúster, puede omitir este paso.

1. Designe una de las clases de almacenamiento existentes como clase de almacenamiento predeterminada. Los comandos de ejemplo siguientes muestran la designación de un StorageClass llamado `ontap-ai-flexvols-retain` Como el tipo de almacenamiento predeterminado.



La `ontap-nas-flexgroup` El tipo de backend de Trident tiene un tamaño de RVP mínimo que es bastante grande. De manera predeterminada, Kubeflow intenta suministrar EVs que son sólo unos pocos GBS en tamaño. Por lo tanto, no debe designar un StorageClass que utilice `ontap-nas-flexgroup` Tipo back-end como StorageClass predeterminado para la implementación de Kubeflow.

```
$ kubectl get sc
NAME                                     PROVISIONER             AGE
ontap-ai-flexgroups-retain              csi.trident.netapp.io   25h
ontap-ai-flexgroups-retain-iface1       csi.trident.netapp.io   25h
ontap-ai-flexgroups-retain-iface2       csi.trident.netapp.io   25h
ontap-ai-flexvols-retain                 csi.trident.netapp.io   3s
$ kubectl patch storageclass ontap-ai-flexvols-retain -p '{"metadata":
{"annotations":{"storageclass.kubernetes.io/is-default-class":"true"}}}'
storageclass.storage.k8s.io/ontap-ai-flexvols-retain patched
$ kubectl get sc
NAME                                     PROVISIONER             AGE
ontap-ai-flexgroups-retain              csi.trident.netapp.io   25h
ontap-ai-flexgroups-retain-iface1       csi.trident.netapp.io   25h
ontap-ai-flexgroups-retain-iface2       csi.trident.netapp.io   25h
ontap-ai-flexvols-retain (default)      csi.trident.netapp.io   54s
```

#### Utilice NVIDIA DeepOps para poner en marcha Kubeflow

NetApp recomienda usar la herramienta de puesta en marcha de Kubeflow que proporciona NVIDIA DeepOps. Para poner en marcha Kubeflow en su clúster de Kubernetes con la herramienta de puesta en marcha DeepOps, siga estas tareas desde el host de salto de implementación.



Como alternativa, puede implementar Kubeflow manualmente siguiendo la ["instrucciones de instalación"](#) En la documentación oficial de Kubeflow

1. Implemente Kubeflow en su clúster siguiendo el ["Instrucciones de despliegue de Kubeflow"](#) En el sitio de NVIDIA DeepOps GitHub.
2. Tenga en cuenta la URL del panel de Kubeflow que genera la herramienta de puesta en marcha de DeepOps Kubeflow.

```
$ ./scripts/k8s/deploy_kubeflow.sh -x
...
INFO[0007] Applied the configuration Successfully!
filename="cmd/apply.go:72"
Kubeflow app installed to: /home/ai/kubeflow
It may take several minutes for all services to start. Run 'kubectl get
pods -n kubeflow' to verify
To remove (excluding CRDs, istio, auth, and cert-manager), run:
./scripts/k8s_deploy_kubeflow.sh -d
To perform a full uninstall : ./scripts/k8s_deploy_kubeflow.sh -D
Kubeflow Dashboard (HTTP NodePort): http://10.61.188.111:31380
```

3. Confirmar que todos los POD implementados en el espacio de nombres Kubeflow muestran un STATUS de Running y confirmar que ningún componente puesto en marcha en el espacio de nombres se encuentra en un estado de error. El inicio de todos los pods puede tardar varios minutos.

```
$ kubectl get all -n kubeflow
```

NAME		READY
STATUS	RESTARTS	AGE
pod/admission-webhook-bootstrap-stateful-set-0		1/1
Running	0	95s
pod/admission-webhook-deployment-6b89c84c98-vrtbh		1/1
Running	0	91s
pod/application-controller-stateful-set-0		1/1
Running	0	98s
pod/argo-ui-5dcf5d8b4f-m2wn4		1/1
Running	0	97s
pod/centraldashboard-cf4874ddc-7hcr8		1/1
Running	0	97s
pod/jupyter-web-app-deployment-685b455447-gjhh7		1/1
Running	0	96s
pod/katib-controller-88c97d85c-kgq66		1/1
Running	1	95s
pod/katib-db-8598468fd8-5jw2c		1/1
Running	0	95s
pod/katib-manager-574c8c67f9-wtrf5		1/1
Running	1	95s
pod/katib-manager-rest-778857c989-fjbzn		1/1
Running	0	95s
pod/katib-suggestion-bayesianoptimization-65df4d7455-qthmw		1/1
Running	0	94s
pod/katib-suggestion-grid-56bf69f597-98vwn		1/1
Running	0	94s
pod/katib-suggestion-hyperband-7777b76cb9-9v6dq		1/1
Running	0	93s

pod/katib-suggestion-nasrl-77f6f9458c-2qzxq	1/1
Running 0 93s	
pod/katib-suggestion-random-77b88b5c79-164j9	1/1
Running 0 93s	
pod/katib-ui-7587c5b967-nd629	1/1
Running 0 95s	
pod/metacontroller-0	1/1
Running 0 96s	
pod/metadata-db-5dd459cc-swzkm	1/1
Running 0 94s	
pod/metadata-deployment-6cf77db994-69fk7	1/1
Running 3 93s	
pod/metadata-deployment-6cf77db994-mpbjt	1/1
Running 3 93s	
pod/metadata-deployment-6cf77db994-xg7tz	1/1
Running 3 94s	
pod/metadata-ui-78f5b59b56-qb6kr	1/1
Running 0 94s	
pod/minio-758b769d67-1lvdr	1/1
Running 0 91s	
pod/ml-pipeline-5875b9db95-g8t2k	1/1
Running 0 91s	
pod/ml-pipeline-persistenceagent-9b69ddd46-bt9r9	1/1
Running 0 90s	
pod/ml-pipeline-scheduledworkflow-7b8d756c76-7x56s	1/1
Running 0 90s	
pod/ml-pipeline-ui-79ffd9c76-fcwpd	1/1
Running 0 90s	
pod/ml-pipeline-viewer-controller-deployment-5fdc87f58-b2t9r	1/1
Running 0 90s	
pod/mysql-657f87857d-15k9z	1/1
Running 0 91s	
pod/notebook-controller-deployment-56b4f59bbf-8bvnr	1/1
Running 0 92s	
pod/profiles-deployment-6bc745947-mrdkh	2/2
Running 0 90s	
pod/pytorch-operator-77c97f4879-hmlrv	1/1
Running 0 92s	
pod/seldon-operator-controller-manager-0	1/1
Running 1 91s	
pod/spartakus-volunteer-5fdfd9db779-17qkm	1/1
Running 0 92s	
pod/tensorboard-6544748d94-nh8b2	1/1
Running 0 92s	
pod/tf-job-dashboard-56f79c59dd-6w59t	1/1
Running 0 92s	

```

pod/tf-job-operator-79cbfd6dbc-rb58c                                     1/1
Running    0          91s
pod/workflow-controller-db644d554-cwrnb                               1/1
Running    0          97s
NAME                                                                 TYPE
CLUSTER-IP      EXTERNAL-IP  PORT(S)          AGE
service/admission-webhook-service                                ClusterIP
10.233.51.169    <none>          443/TCP          97s
service/application-controller-service                          ClusterIP
10.233.4.54      <none>          443/TCP          98s
service/argo-ui                                              NodePort
10.233.47.191    <none>          80:31799/TCP     97s
service/centraldashboard                                    ClusterIP
10.233.8.36      <none>          80/TCP           97s
service/jupyter-web-app-service                                ClusterIP
10.233.1.42      <none>          80/TCP           97s
service/katib-controller                                    ClusterIP
10.233.25.226    <none>          443/TCP          96s
service/katib-db                                             ClusterIP
10.233.33.151    <none>          3306/TCP         97s
service/katib-manager                                       ClusterIP
10.233.46.239    <none>          6789/TCP         96s
service/katib-manager-rest                                   ClusterIP
10.233.55.32     <none>          80/TCP           96s
service/katib-suggestion-bayesianoptimization                ClusterIP
10.233.49.191    <none>          6789/TCP         95s
service/katib-suggestion-grid                                ClusterIP
10.233.9.105     <none>          6789/TCP         95s
service/katib-suggestion-hyperband                           ClusterIP
10.233.22.2      <none>          6789/TCP         95s
service/katib-suggestion-nasrl                               ClusterIP
10.233.63.73     <none>          6789/TCP         95s
service/katib-suggestion-random                              ClusterIP
10.233.57.210    <none>          6789/TCP         95s
service/katib-ui                                             ClusterIP
10.233.6.116     <none>          80/TCP           96s
service/metadata-db                                          ClusterIP
10.233.31.2      <none>          3306/TCP         96s
service/metadata-service                                    ClusterIP
10.233.27.104    <none>          8080/TCP         96s
service/metadata-ui                                          ClusterIP
10.233.57.177    <none>          80/TCP           96s
service/minio-service                                       ClusterIP
10.233.44.90     <none>          9000/TCP         94s
service/ml-pipeline                                          ClusterIP
10.233.41.201    <none>          8888/TCP,8887/TCP 94s

```



service/ml-pipeline-tensorboard-ui		ClusterIP		
10.233.36.207	<none>	80/TCP	93s	
service/ml-pipeline-ui		ClusterIP		
10.233.61.150	<none>	80/TCP	93s	
service/mysql		ClusterIP		
10.233.55.117	<none>	3306/TCP	94s	
service/notebook-controller-service		ClusterIP		
10.233.10.166	<none>	443/TCP	95s	
service/profiles-kfam		ClusterIP		
10.233.33.79	<none>	8081/TCP	92s	
service/pytorch-operator		ClusterIP		
10.233.37.112	<none>	8443/TCP	95s	
service/seldon-operator-controller-manager-service		ClusterIP		
10.233.30.178	<none>	443/TCP	92s	
service/tensorboard		ClusterIP		
10.233.58.151	<none>	9000/TCP	94s	
service/tf-job-dashboard		ClusterIP		
10.233.4.17	<none>	80/TCP	94s	
service/tf-job-operator		ClusterIP		
10.233.60.32	<none>	8443/TCP	94s	
service/webhook-server-service		ClusterIP		
10.233.32.167	<none>	443/TCP	87s	
NAME			READY	UP-
TO-DATE	AVAILABLE	AGE		
deployment.apps/admission-webhook-deployment			1/1	1
1	97s			
deployment.apps/argo-ui			1/1	1
1	97s			
deployment.apps/centraldashboard			1/1	1
1	97s			
deployment.apps/jupyter-web-app-deployment			1/1	1
1	97s			
deployment.apps/katib-controller			1/1	1
1	96s			
deployment.apps/katib-db			1/1	1
1	97s			
deployment.apps/katib-manager			1/1	1
1	96s			
deployment.apps/katib-manager-rest			1/1	1
1	96s			
deployment.apps/katib-suggestion-bayesianoptimization			1/1	1
1	95s			
deployment.apps/katib-suggestion-grid			1/1	1
1	95s			
deployment.apps/katib-suggestion-hyperband			1/1	1
1	95s			

deployment.apps/katib-suggestion-nasrl	1/1	1
1 95s		
deployment.apps/katib-suggestion-random	1/1	1
1 95s		
deployment.apps/katib-ui	1/1	1
1 96s		
deployment.apps/metadata-db	1/1	1
1 96s		
deployment.apps/metadata-deployment	3/3	3
3 96s		
deployment.apps/metadata-ui	1/1	1
1 96s		
deployment.apps/minio	1/1	1
1 94s		
deployment.apps/ml-pipeline	1/1	1
1 94s		
deployment.apps/ml-pipeline-persistenceagent	1/1	1
1 93s		
deployment.apps/ml-pipeline-scheduledworkflow	1/1	1
1 93s		
deployment.apps/ml-pipeline-ui	1/1	1
1 93s		
deployment.apps/ml-pipeline-viewer-controller-deployment	1/1	1
1 93s		
deployment.apps/mysql	1/1	1
1 94s		
deployment.apps/notebook-controller-deployment	1/1	1
1 95s		
deployment.apps/profiles-deployment	1/1	1
1 92s		
deployment.apps/pytorch-operator	1/1	1
1 95s		
deployment.apps/spartakus-volunteer	1/1	1
1 94s		
deployment.apps/tensorboard	1/1	1
1 94s		
deployment.apps/tf-job-dashboard	1/1	1
1 94s		
deployment.apps/tf-job-operator	1/1	1
1 94s		
deployment.apps/workflow-controller	1/1	1
1 97s		
NAME		
DESIRED CURRENT READY AGE		
replicaset.apps/admission-webhook-deployment-6b89c84c98		1
1 1 97s		

replicaset.apps/argo-ui-5dcf5d8b4f	1
1 1 97s	
replicaset.apps/centraldashboard-cf4874ddc	1
1 1 97s	
replicaset.apps/jupyter-web-app-deployment-685b455447	1
1 1 97s	
replicaset.apps/katib-controller-88c97d85c	1
1 1 96s	
replicaset.apps/katib-db-8598468fd8	1
1 1 97s	
replicaset.apps/katib-manager-574c8c67f9	1
1 1 96s	
replicaset.apps/katib-manager-rest-778857c989	1
1 1 96s	
replicaset.apps/katib-suggestion-bayesianoptimization-65df4d7455	1
1 1 95s	
replicaset.apps/katib-suggestion-grid-56bf69f597	1
1 1 95s	
replicaset.apps/katib-suggestion-hyperband-7777b76cb9	1
1 1 95s	
replicaset.apps/katib-suggestion-nasrl-77f6f9458c	1
1 1 95s	
replicaset.apps/katib-suggestion-random-77b88b5c79	1
1 1 95s	
replicaset.apps/katib-ui-7587c5b967	1
1 1 96s	
replicaset.apps/metadata-db-5dd459cc	1
1 1 96s	
replicaset.apps/metadata-deployment-6cf77db994	3
3 3 96s	
replicaset.apps/metadata-ui-78f5b59b56	1
1 1 96s	
replicaset.apps/minio-758b769d67	1
1 1 93s	
replicaset.apps/ml-pipeline-5875b9db95	1
1 1 93s	
replicaset.apps/ml-pipeline-persistenceagent-9b69ddd46	1
1 1 92s	
replicaset.apps/ml-pipeline-scheduledworkflow-7b8d756c76	1
1 1 91s	
replicaset.apps/ml-pipeline-ui-79ffd9c76	1
1 1 91s	
replicaset.apps/ml-pipeline-viewer-controller-deployment-5fdc87f58	1
1 1 91s	
replicaset.apps/mysql-657f87857d	1
1 1 92s	

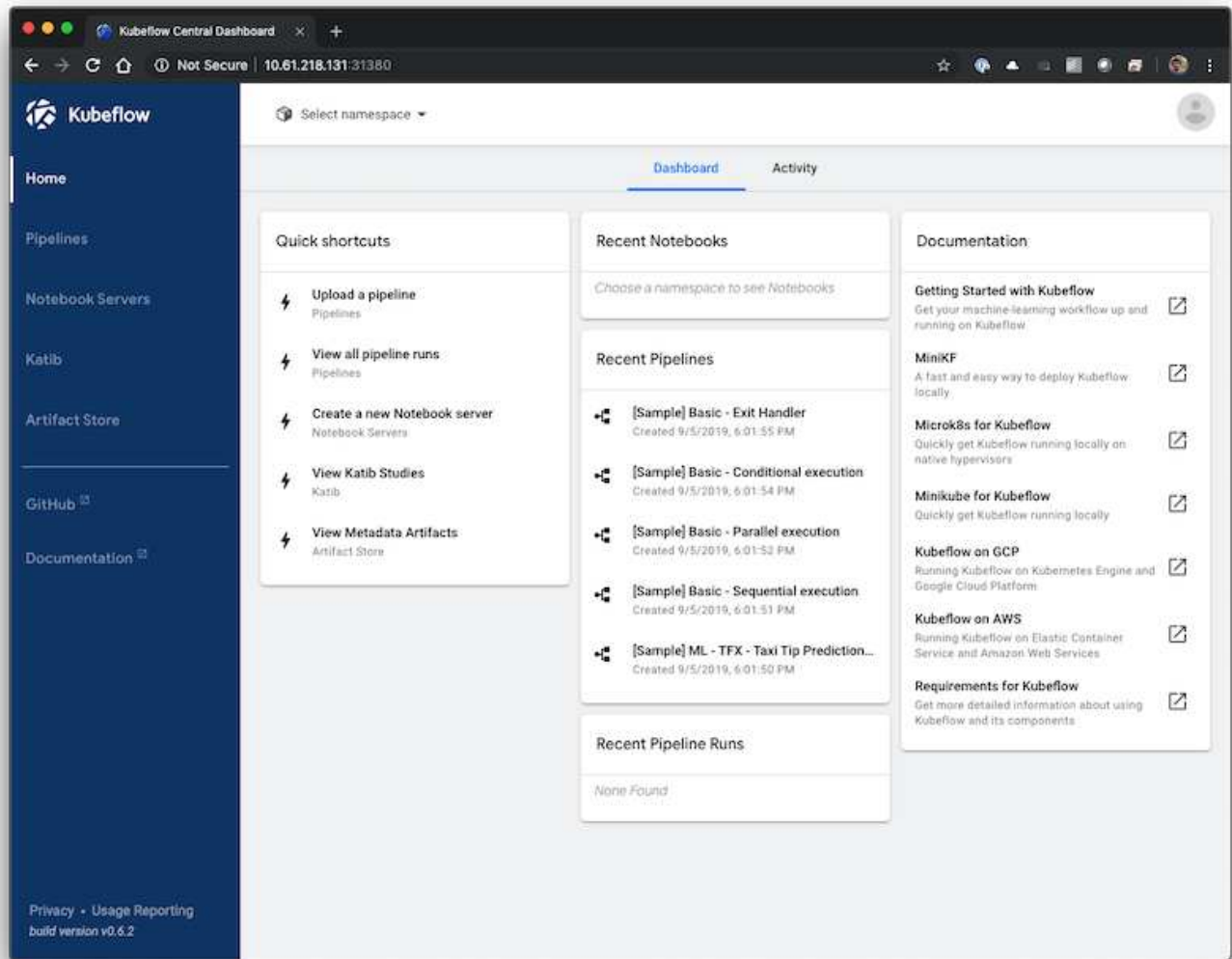
```

replicaset.apps/notebook-controller-deployment-56b4f59bbf      1
1          1          94s
replicaset.apps/profiles-deployment-6bc745947                  1
1          1          91s
replicaset.apps/pytorch-operator-77c97f4879                    1
1          1          94s
replicaset.apps/spartakus-volunteer-5fdfd6db779                1
1          1          94s
replicaset.apps/tensorboard-6544748d94                          1
1          1          93s
replicaset.apps/tf-job-dashboard-56f79c59dd                    1
1          1          93s
replicaset.apps/tf-job-operator-79cbfd6dbc                      1
1          1          93s
replicaset.apps/workflow-controller-db644d554                  1
1          1          97s
NAME                                                              READY   AGE
statefulset.apps/admission-webhook-bootstrap-stateful-set      1/1    97s
statefulset.apps/application-controller-stateful-set           1/1    98s
statefulset.apps/metacontroller                                1/1    98s
statefulset.apps/seldon-operator-controller-manager            1/1    92s
$ kubectl get pvc -n kubeflow
NAME          STATUS    VOLUME
CAPACITY      ACCESS MODES  STORAGECLASS          AGE
katib-mysql    Bound      pvc-b07f293e-d028-11e9-9b9d-00505681a82d
10Gi          RWO        ontap-ai-flexvols-retain  27m
metadata-mysql Bound      pvc-b0f3f032-d028-11e9-9b9d-00505681a82d
10Gi          RWO        ontap-ai-flexvols-retain  27m
minio-pv-claim Bound      pvc-b22727ee-d028-11e9-9b9d-00505681a82d
20Gi          RWO        ontap-ai-flexvols-retain  27m
mysql-pv-claim Bound      pvc-b2429afd-d028-11e9-9b9d-00505681a82d
20Gi          RWO        ontap-ai-flexvols-retain  27m

```

4. En su navegador web, acceda al panel central de Kubeflow navegando hasta la URL que anotó en el paso 2.

El nombre de usuario predeterminado es `admin@kubeflow.org`, y la contraseña predeterminada es `12341234`. Para crear usuarios adicionales, siga las instrucciones de ["Documentación oficial de Kubeflow"](#).



## Ejemplo de operaciones y tareas de Kubeflow

Esta sección incluye ejemplos de diversas operaciones y tareas que puede que desee realizar utilizando Kubeflow.

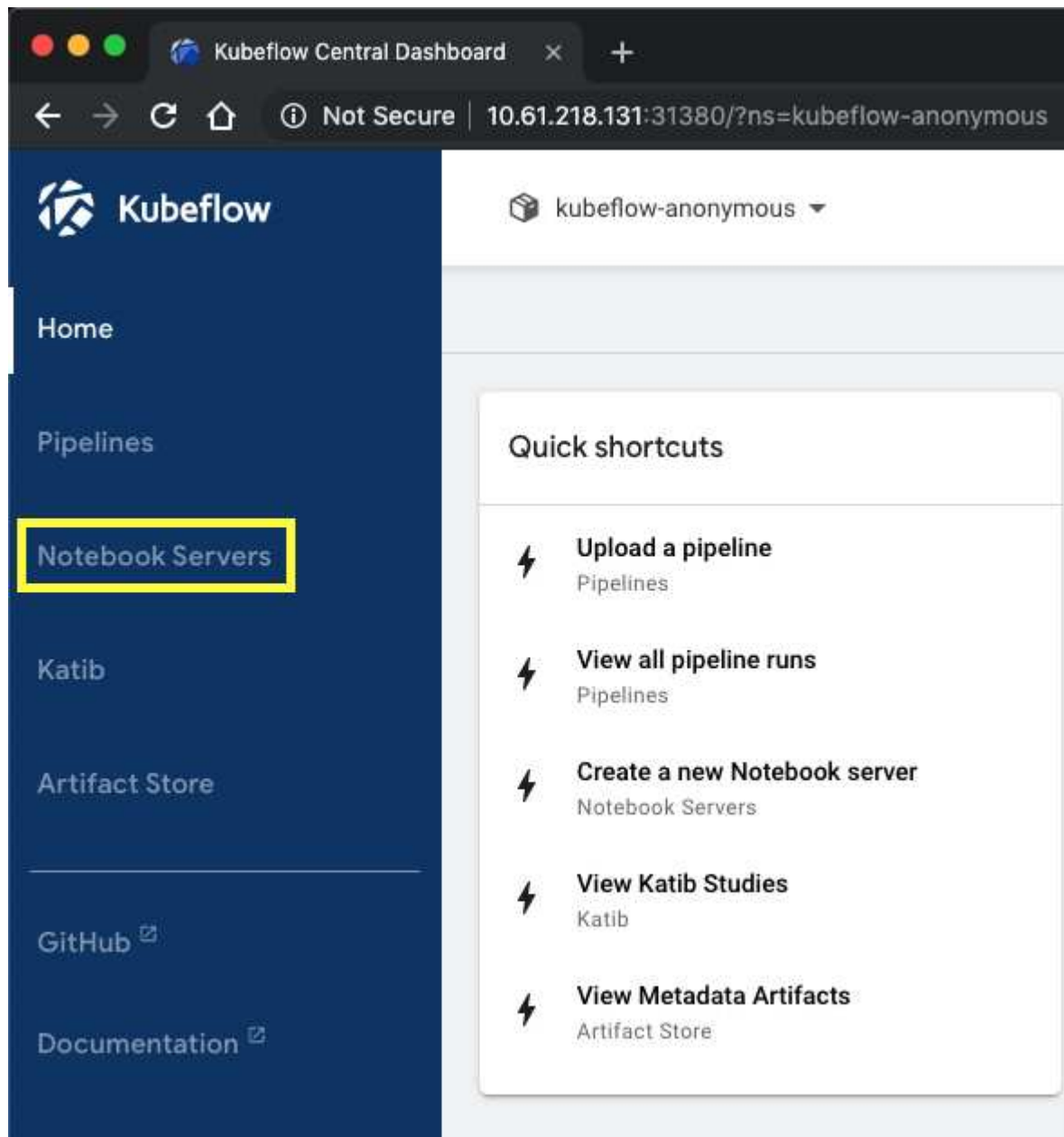
### Ejemplo de operaciones y tareas de Kubeflow

Esta sección incluye ejemplos de diversas operaciones y tareas que puede que desee realizar utilizando Kubeflow.

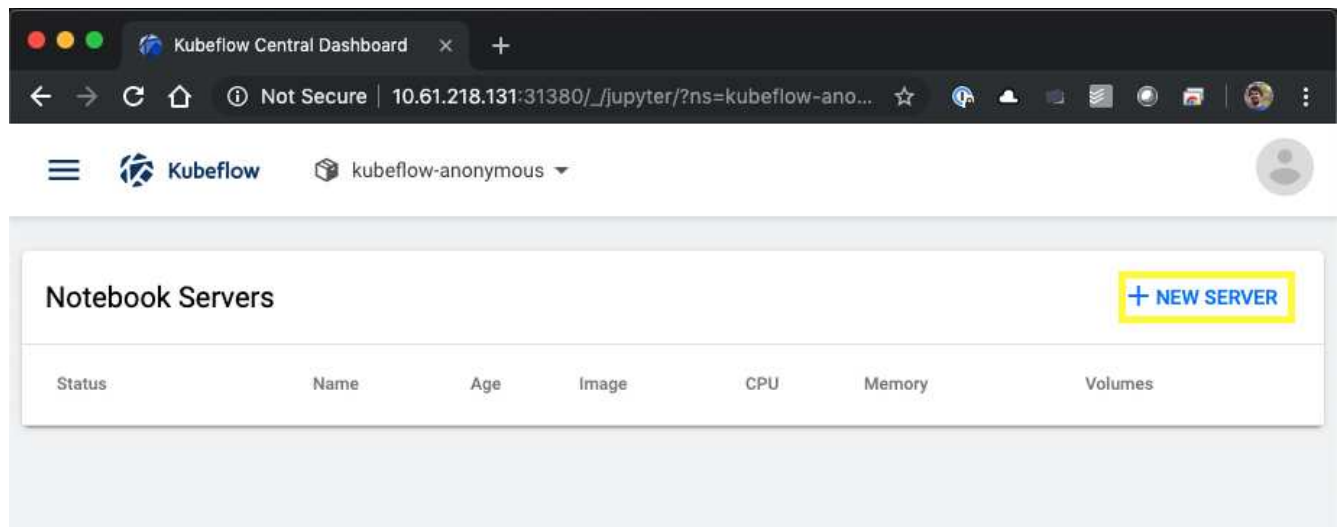
### Aprovisione un espacio de trabajo para portátiles Jupyter para uso científico de datos o desarrollador

Kubeflow es capaz de suministrar rápidamente nuevos servidores Jupyter Notebook para actuar como espacios de trabajo de científicos de datos. Para aprovisionar un nuevo servidor Jupyter Notebook con Kubeflow, realice las siguientes tareas. Para obtener más información acerca de Jupyter Notebooks dentro del contexto de Kubeflow, consulte ["Documentación oficial de Kubeflow"](#).

1. Desde el panel central de Kubeflow, haga clic en Notebook Servers en el menú principal para acceder a la página de administración del servidor Jupyter Notebook.



2. Haga clic en Nuevo servidor para aprovisionar un nuevo servidor Jupyter Notebook.



3. Asigne un nombre a su nuevo servidor, elija la imagen de Docker en la que desea que se base su servidor y especifique la cantidad de CPU y RAM que debe reservar su servidor. Si el campo espacio de nombres está en blanco, utilice el menú Seleccionar espacio de nombres en el encabezado de la página para elegir un espacio de nombres. El campo de espacio de nombres se rellena automáticamente con el espacio de nombres elegido.

En el siguiente ejemplo, la `kubeflow-anonymous` se elige espacio de nombres. Además, se aceptan los valores predeterminados de imagen Docker, CPU y RAM.

**Name**

Specify the name of the Notebook Server and the Namespace it will belong to.

Name:  Namespace:

**Image**

A starter Jupyter Docker Image with a baseline deployment and typical ML packages.

☐ Custom Image

Image:

**CPU / RAM**

Specify the total amount of CPU and RAM reserved by your Notebook Server. For CPU-intensive workloads, you can choose more than 1 CPU (e.g. 1.5).

CPU:  Memory:

- Especifique los detalles del volumen del área de trabajo. Si decide crear un nuevo volumen, ese volumen o RVP se aprovisionan con el tipo de almacenamiento predeterminado. Dado que se designó un clase de almacenamiento con Trident como el clase de almacenamiento predeterminado de la sección "[Despliegue de Kubeflow](#)", El volumen o PVC se aprovisiona con Trident. Este volumen se monta automáticamente como espacio de trabajo predeterminado dentro del contenedor servidor del portátil Jupyter. Los portátiles que cree un usuario en el servidor que no se guarden en un volumen de datos independiente se guardarán automáticamente en este volumen de área de trabajo. Por lo tanto, los portátiles se conservan entre reinicios.

**Workspace Volume**

Configure the Volume to be mounted as your personal Workspace.

☐ Don't use Persistent Storage for User's home

Type:  Name:  Size:  Mode:  Mount Point:

- Añadir volúmenes de datos. En el siguiente ejemplo, se especifica un RVP existente llamado 'pb-fg-All' y acepta el punto de montaje predeterminado.



**Data Volumes**

Configure the Volumes to be mounted as your Datasets.

[+ ADD VOLUME](#)

Type	Name	Size	Mode	Mount Point
Existing	pb-fg-all	10Gi	ReadWriteOnce	/home/jovyan/data-vol-1

6. **Opcional:** Solicite que se asigne el número deseado de GPU al servidor de su portátil. En el siguiente ejemplo, se solicita una GPU.

**Configurations**

Extra layers of configurations that will be applied to the new Notebook. (e.g. Insert credentials as Secrets, set Environment Variables.)

Configurations

**Extra Resources**

Specify extra resources that might be needed in the Notebook Server.

☒ **Enable Shared Memory**

Extra Resources \*

`{"nvidia.com/gpu": 1}`

Extra Resources available in the cluster (ex. NVIDIA GPUs)

[LAUNCH](#) [CANCEL](#)

7. Haga clic en Iniciar para aprovisionar su nuevo servidor de notebook.
8. Espere a que el servidor de su portátil esté completamente aprovisionado. Esto puede tardar varios minutos si nunca ha aprovisionado un servidor con la imagen Docker que ha especificado porque es necesario descargar la imagen. Una vez que el servidor ha sido completamente aprovisionado, verá una Marca de verificación verde en la columna Estado de la página de administración del servidor Jupyter Notebook.

**Notebook Servers** [+ NEW SERVER](#)

Status	Name	Age	Image	CPU	Memory	Volumes	
	mike	12 mins ago	tensorflow-1.13.1-notebook-cpu:v0.5.0	0.5	1.0Gi		<a href="#">CONNECT</a>

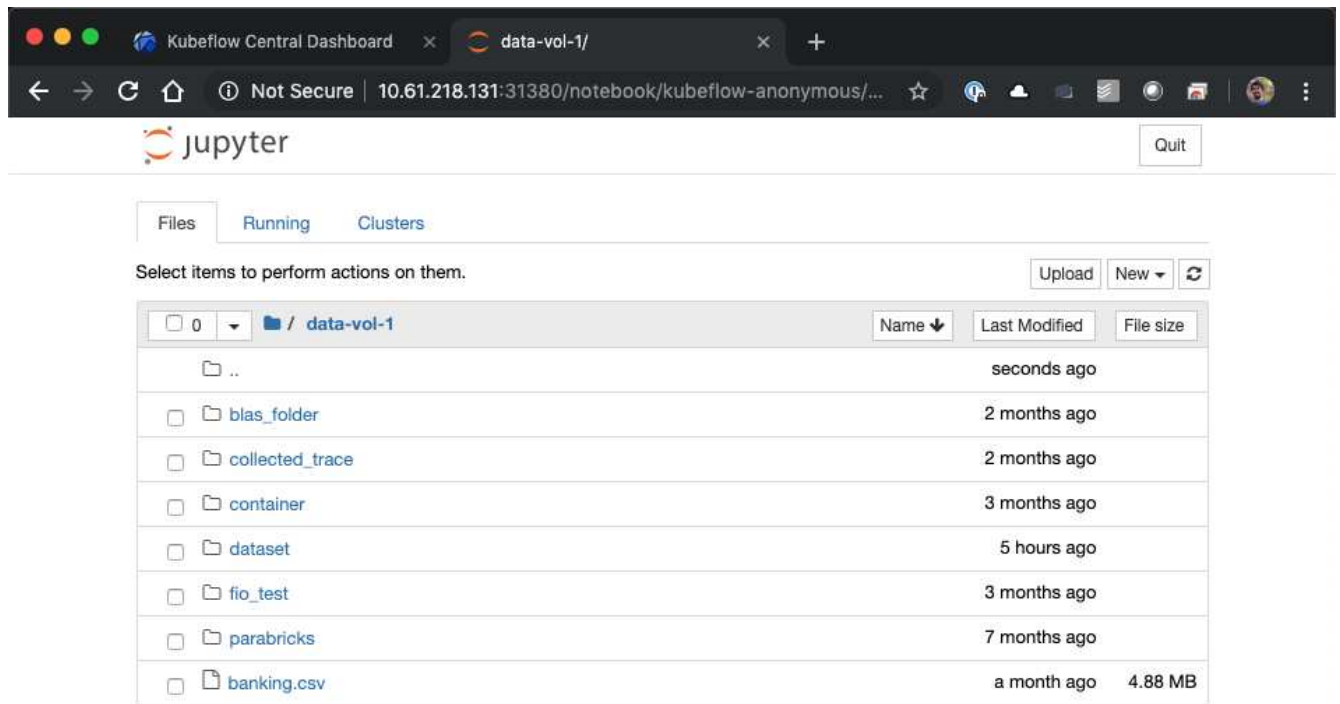
9. Haga clic en conectar para conectarse a la nueva interfaz web del servidor.
10. Confirme que el volumen del conjunto de datos especificado en el paso 6 está montado en el servidor. Tenga en cuenta que este volumen se monta en el espacio de trabajo predeterminado de forma predeterminada. Desde la perspectiva del usuario, ésta es sólo otra carpeta dentro del área de trabajo. El usuario, que es probable que un científico de datos y no sea un experto en infraestructura, no necesita tener experiencia en almacenamiento para utilizar este volumen.

**jupyter** [Quit](#)

Files Running Clusters

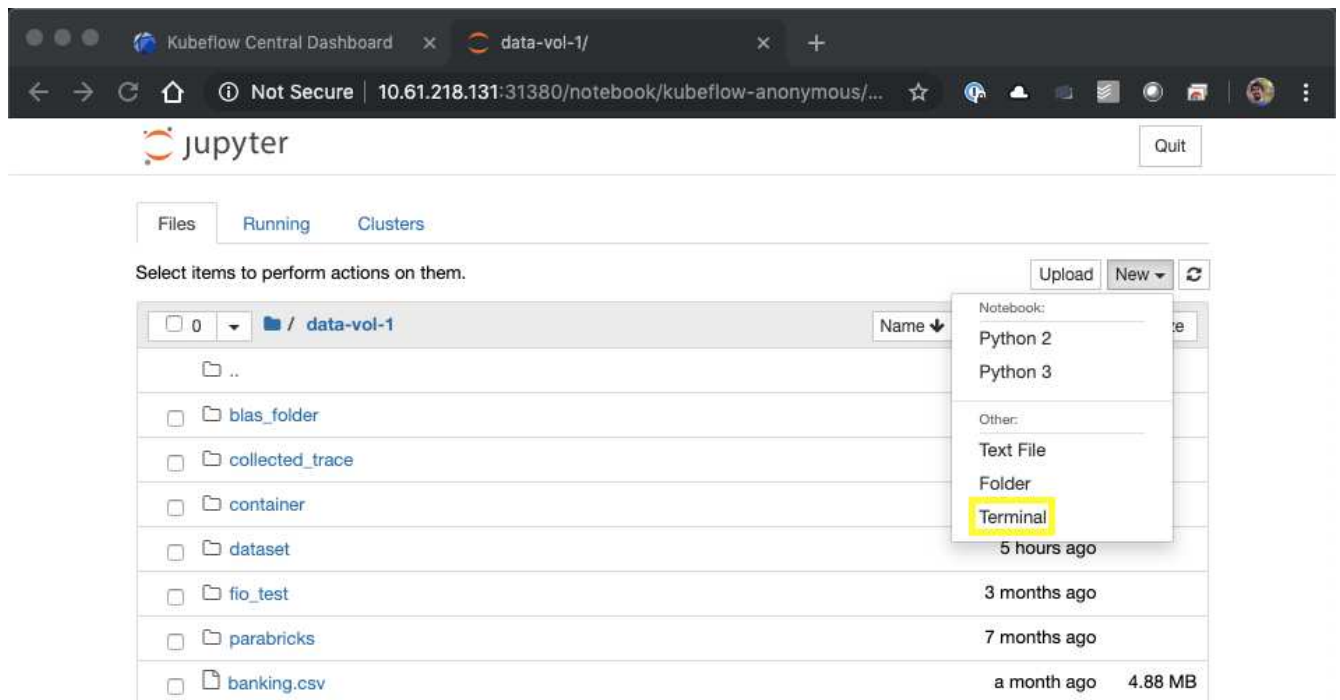
Select items to perform actions on them. [Upload](#) [New](#) [Refresh](#)

	Name	Last Modified	File size
<input type="checkbox"/>	/		
<input type="checkbox"/>	data-vol-1	a day ago	



11. Abra un terminal y, suponiendo que se haya solicitado un nuevo volumen en el paso 5, ejecute `df -h`. Para confirmar que un nuevo volumen persistente aprovisionado por Trident se monta como espacio de trabajo predeterminado.

El directorio de área de trabajo predeterminado es el directorio base con el que se presenta cuando se accede por primera vez a la interfaz web del servidor. Por lo tanto, cualquier artefacto que cree mediante la interfaz web se almacena en este volumen persistente aprovisionado por Trident.



```

$ df -h
Filesystem                                Size  Used Avail
Use% Mounted on
overlay                                  439G   34G  382G
9% /
tmpfs                                     64M    0   64M
0% /dev
tmpfs                                     252G    0  252G
0% /sys/fs/cgroup
/dev/sda2                                439G   34G  382G
9% /etc/hosts
192.168.11.11:/trident_pvc_3dcfe7e5_d5a9_11e9_9b9d_00505681a82d 10G 320K 10G
1% /home/jovyan
tmpfs                                     252G    0  252G
0% /dev/shm
192.168.11.11:/pb_fg_all                  10T   10T   47G
100% /home/jovyan/data-vol-1
tmpfs                                     252G   12K  252G
1% /run/secrets/kubernetes.io/serviceaccount
tmpfs                                     252G   12K  252G
1% /proc/driver/nvidia
tmpfs                                     51G   4.9M   51G
1% /run/nvidia-persistenced/socket
udev                                     252G    0  252G
0% /dev/nvidia5
tmpfs                                     252G    0  252G
0% /proc/acpi
tmpfs                                     252G    0  252G
0% /proc/scsi
tmpfs                                     252G    0  252G
0% /sys/firmware
$

```

- Con el terminal, ejecute `nvidia-smi` Para confirmar que se ha asignado el número correcto de GPU al servidor del portátil. En el siguiente ejemplo, se ha asignado una GPU al servidor de portátiles como se solicitó en el paso 7.

```

$ nvidia-smi
Fri Sep 13 13:52:15 2019
+-----+
| NVIDIA-SMI 410.104      Driver Version: 410.104      CUDA Version: N/A      |
+-----+
| GPU   Name               Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
+-----+-----+
|  0  Tesla V100-SXM2...    On         | 00000000:86:00:0 Off  |           0%         |
| N/A   38C    P0      46W / 300W | 0MiB / 32480MiB |           0%         |
+-----+-----+
+-----+
| Processes:                                                       GPU Memory |
|  GPU       PID    Type    Process name                        Usage    |
+-----+-----+
| No running processes found                                         |
+-----+
$

```

La ["Kit de herramientas para la ciencia de datos de NetApp para Kubernetes"](#) Se puede utilizar junto con Kubeflow. El uso del kit de herramientas para la ciencia de datos de NetApp con Kubeflow ofrece las siguientes ventajas:

- Los científicos de datos pueden llevar a cabo operaciones avanzadas de gestión de datos de NetApp directamente desde un portátil Jupyter.
- Las operaciones avanzadas de gestión de datos de NetApp pueden incorporarse en flujos de trabajo automatizados mediante el marco de canalizaciones de Kubeflow.

Consulte la ["Ejemplos de Kubeflow"](#) Sección dentro del repositorio de Data Science Toolkit de NetApp, GitHub para obtener información sobre el uso del kit de herramientas con Kubeflow.

### Distribución del flujo de aire de Apache

NetApp recomienda ejecutar el flujo de aire de Apache sobre Kubernetes. En esta sección se describen las tareas que debe completar para poner en marcha el flujo de aire en el clúster de Kubernetes.



Es posible poner en marcha un flujo de aire en plataformas distintas a Kubernetes. Esta solución no cubre la posibilidad de poner en marcha un flujo de aire en plataformas distintas a Kubernetes.

### Requisitos previos

Antes de realizar el ejercicio de implementación descrito en esta sección, asumimos que ya ha realizado las siguientes tareas:

1. Ya tiene un clúster de Kubernetes en funcionamiento.
2. Ya ha instalado y configurado NetApp Trident en su clúster de Kubernetes como se indica en la sección ["NetApp Trident Deployment and Configuration"](#).

### Instale el Helm

El flujo de aire se pone en marcha con Helm, un conocido administrador de paquetes para Kubernetes. Antes de implementar el flujo de aire, debe instalar Helm en el host de salto de la implementación. Para instalar Helm en el host de salto de despliegue, siga la ["instrucciones de instalación"](#) En la documentación oficial de Helm.

### Establezca el tipo de almacenamiento de Kubernetes predeterminado

Antes de implementar el flujo de aire, debe designar un tipo de almacenamiento predeterminado en el clúster de Kubernetes. El proceso de implementación de flujo de aire intenta aprovisionar nuevos volúmenes persistentes mediante el tipo de almacenamiento predeterminado. Si no se designa StorageClass como clase de almacenamiento predeterminado, la implementación falla. Para designar un tipo de almacenamiento predeterminado en el clúster, siga las instrucciones que se describen en la sección ["Despliegue de Kubeflow"](#). Si ya ha designado un tipo de almacenamiento predeterminado en el clúster, puede omitir este paso.

### Utilice Helm para desplegar el flujo de aire

Para poner en marcha el flujo de aire en su clúster de Kubernetes con Helm, realice las siguientes tareas

desde el host de salto de implementación:

1. Despliegue el flujo de aire con Helm siguiendo la ["instrucciones de puesta en funcionamiento"](#) Para el diagrama de flujo de aire oficial en el Hub de artefactos. Los comandos de ejemplo siguientes muestran la implementación del flujo de aire con Helm. Modifique, agregue o elimine valores en la `custom-values.yaml` fichero según sea necesario en función de su entorno y de la configuración deseada.

```
$ cat << EOF > custom-values.yaml
#####
# Airflow - Common Configs
#####
airflow:
  ## the airflow executor type to use
  ##
  executor: "CeleryExecutor"
  ## environment variables for the web/scheduler/worker Pods (for
  airflow configs)
  ##
  #
#####
# Airflow - WebUI Configs
#####
web:
  ## configs for the Service of the web Pods
  ##
  service:
    type: NodePort
#####
# Airflow - Logs Configs
#####
logs:
  persistence:
    enabled: true
#####
# Airflow - DAGs Configs
#####
dags:
  ## configs for the DAG git repository & sync container
  ##
  gitSync:
    enabled: true
    ## url of the git repository
    ##
    repo: "git@github.com:mboglesby/airflow-dev.git"
    ## the branch/tag/sha1 which we clone
    ##
    branch: master
```

```

revision: HEAD
## the name of a pre-created secret containing files for ~/.ssh/
##
## NOTE:
## - this is ONLY RELEVANT for SSH git repos
## - the secret commonly includes files: id_rsa, id_rsa.pub,
known_hosts
## - known_hosts is NOT NEEDED if `git.sshKeyscan` is true
##
sshSecret: "airflow-ssh-git-secret"
## the name of the private key file in your `git.secret`
##
## NOTE:
## - this is ONLY RELEVANT for PRIVATE SSH git repos
##
sshSecretKey: id_rsa
## the git sync interval in seconds
##
syncWait: 60
EOF
$ helm install airflow airflow-stable/airflow -n airflow --version 8.0.8
--values ./custom-values.yaml
...
Congratulations. You have just deployed Apache Airflow!
1. Get the Airflow Service URL by running these commands:
    export NODE_PORT=$(kubectl get --namespace airflow -o
jsonpath="{.spec.ports[0].nodePort}" services airflow-web)
    export NODE_IP=$(kubectl get nodes --namespace airflow -o
jsonpath="{.items[0].status.addresses[0].address}")
    echo http://$NODE_IP:$NODE_PORT/
2. Open Airflow in your web browser

```

2. Confirme que todos los pods de flujo de aire estén activos y en funcionamiento. Puede tardar varios minutos en comenzar todos los pods.

```

$ kubectl -n airflow get pod

```

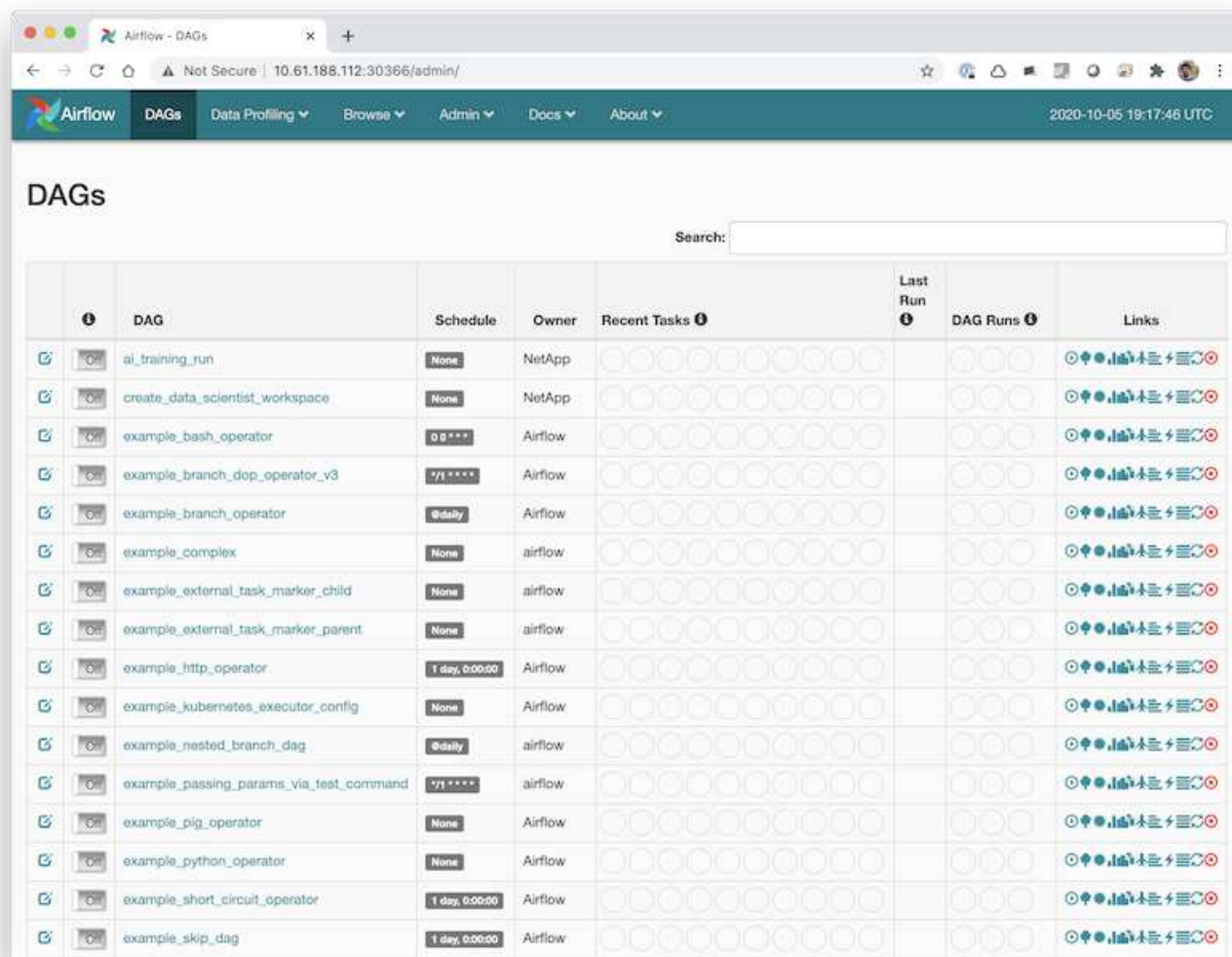
NAME	READY	STATUS	RESTARTS	AGE
airflow-flower-b5656d44f-h8qjk	1/1	Running	0	2h
airflow-postgresql-0	1/1	Running	0	2h
airflow-redis-master-0	1/1	Running	0	2h
airflow-scheduler-9d95fcd9-clf4b	2/2	Running	2	2h
airflow-web-59c94db9c5-z7rg4	1/1	Running	0	2h
airflow-worker-0	2/2	Running	2	2h

3. Obtenga la URL del servicio web de flujo de aire siguiendo las instrucciones que se imprimieron en la

consola cuando implementó el flujo de aire con Helm en el paso 1.

```
$ export NODE_PORT=$(kubectl get --namespace airflow -o
jsonpath="{.spec.ports[0].nodePort}" services airflow-web)
$ export NODE_IP=$(kubectl get nodes --namespace airflow -o
jsonpath="{.items[0].status.addresses[0].address}")
$ echo http://$NODE_IP:$NODE_PORT/
```

4. Confirme que puede acceder al servicio web de flujo de aire.



The screenshot shows the Airflow web interface in a browser. The URL is 10.61.188.112:30366/admin/. The page title is "Airflow - DAGs". The navigation bar includes "DAGs", "Data Profiling", "Browse", "Admin", "Docs", and "About". The main content area is titled "DAGs" and features a search bar. Below the search bar is a table listing various DAGs.

	ⓘ	DAG	Schedule	Owner	Recent Tasks ⓘ	Last Run ⓘ	DAG Runs ⓘ	Links
🔗	🔗	ai_training_run	None	NetApp	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	create_data_scientist_workspace	None	NetApp	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_bash_operator	0 0 * * *	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_branch_dop_operator_v3	* * * * *	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_branch_operator	@daily	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_complex	None	airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_external_task_marker_child	None	airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_external_task_marker_parent	None	airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_http_operator	1 day, 0:00:00	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_kubernetes_executor_config	None	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_nested_branch_dag	@daily	airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_passing_params_via_test_command	* * * * *	airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_pig_operator	None	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_python_operator	None	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_short_circuit_operator	1 day, 0:00:00	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗
🔗	🔗	example_skip_dag	1 day, 0:00:00	Airflow	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗		🔗🔗🔗🔗🔗🔗	🔗🔗🔗🔗🔗🔗🔗🔗🔗🔗

### Ejemplo de flujos de trabajo de flujo de aire de Apache

La ["Kit de herramientas para la ciencia de datos de NetApp para Kubernetes"](#) Se puede utilizar junto con el flujo de aire. El uso del kit de herramientas para la ciencia de datos de NetApp con flujo de aire le permite incorporar operaciones de gestión de datos de NetApp en flujos de trabajo automatizados orquestados por flujo de aire.

Consulte la ["Ejemplos de flujo de aire"](#) Sección dentro del repositorio de Data Science Toolkit de NetApp,



GitHub para obtener información sobre el uso del kit de herramientas con flujo de aire.

## Ejemplo de operaciones de Trident

En esta sección se incluyen ejemplos de diversas operaciones que quizás desee realizar con Trident.

### Importe un volumen existente

Si hay volúmenes existentes en su sistema/plataforma de almacenamiento de NetApp que desea montar en contenedores dentro de su clúster de Kubernetes, pero que no están ligados a las RVP en el clúster, debe importar estos volúmenes. Es posible usar la funcionalidad de importación de volúmenes de Trident para importar estos volúmenes.

Los comandos de ejemplo siguientes muestran la importación del mismo volumen, denominado `pb_fg_all`, Dos veces, una para cada backend de Trident que se creó en el ejemplo de la sección ["Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI"](#), paso 1. Importar el mismo volumen dos veces de esta manera le permite montar el volumen (un volumen FlexGroup existente) varias veces en diferentes LIF, como se describe en la sección ["Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI"](#), paso 1. Para obtener más información acerca de las EVs, consulte ["Documentación oficial sobre Kubernetes"](#). Para obtener más información sobre la funcionalidad de importación de volúmenes, consulte ["Documentación de Trident"](#).

An `accessModes` valor de `ReadOnlyMany` Se especifica en los archivos de especificaciones de PVC de ejemplo. Para obtener más información acerca de `accessMode` consulte ["Documentación oficial sobre Kubernetes"](#).



Los nombres de backend especificados en los siguientes comandos de importación de ejemplo corresponden a las backends que se crearon en el ejemplo de la sección ["Ejemplo de Back-ends de Trident para puestas en marcha de ONTAP AI"](#), paso 1. Los nombres de StorageClass especificados en el siguiente ejemplo de archivos de definición PVC corresponden a las clases de almacenamiento que se crearon en el ejemplo de la sección ["Ejemplo de clases de almacenamiento Kubernetes para puestas en marcha de IA en ONTAP"](#), paso 1.

```
$ cat << EOF > ./pvc-import-pb_fg_all-iface1.yaml
kind: PersistentVolumeClaim
apiVersion: v1
metadata:
  name: pb-fg-all-iface1
  namespace: default
spec:
  accessModes:
    - ReadOnlyMany
  storageClassName: ontap-ai-flexgroups-retain-iface1
EOF
$ tridentctl import volume ontap-ai-flexgroups-iface1 pb_fg_all -f ./pvc-
import-pb_fg_all-iface1.yaml -n trident
+-----+-----+
+-----+-----+
+-----+-----+-----+-----+
```

NAME	SIZE	STORAGE CLASS
default-pb-fg-all-iface1-7d9f1	10 TiB	ontap-ai-flexgroups-retain- iface1

```

$ cat << EOF > ./pvc-import-pb_fg_all-iface2.yaml
kind: PersistentVolumeClaim
apiVersion: v1
metadata:
  name: pb-fg-all-iface2
  namespace: default
spec:
  accessModes:
    - ReadOnlyMany
  storageClassName: ontap-ai-flexgroups-retain-iface2
EOF
$ tridentctl import volume ontap-ai-flexgroups-iface2 pb_fg_all -f ./pvc-
import-pb_fg_all-iface2.yaml -n trident

```

NAME	SIZE	STORAGE CLASS
default-pb-fg-all-iface2-85aee	10 TiB	ontap-ai-flexgroups-retain- iface2

```

$ tridentctl get volume -n trident

```

NAME	SIZE	STORAGE CLASS
default-pb-fg-all-iface2-85aee	10 TiB	ontap-ai-flexgroups-retain- iface2

```

+-----+-----+
+-----+-----+
+-----+-----+-----+
| default-pb-fg-all-iface1-7d9f1 | 10 TiB | ontap-ai-flexgroups-retain-
iface1 | file | b74cbddb-e0b8-40b7-b263-b6da6dec0bdd | online | true
|
| default-pb-fg-all-iface2-85aee | 10 TiB | ontap-ai-flexgroups-retain-
iface2 | file | 61814d48-c770-436b-9cb4-cf7ee661274d | online | true
|
+-----+-----+
+-----+-----+
+-----+-----+
$ kubectl get pvc
NAME                                STATUS    VOLUME                                CAPACITY
ACCESS MODES    STORAGECLASS    AGE
pb-fg-all-iface1    Bound    default-pb-fg-all-iface1-7d9f1
10995116277760    ROX      ontap-ai-flexgroups-retain-iface1    25h
pb-fg-all-iface2    Bound    default-pb-fg-all-iface2-85aee
10995116277760    ROX      ontap-ai-flexgroups-retain-iface2    25h

```

### Aprovisione un nuevo volumen

Puede usar Trident para aprovisionar un nuevo volumen en su plataforma o sistema de almacenamiento de NetApp. Los siguientes comandos de ejemplo muestran el aprovisionamiento de un volumen FlexVol nuevo. En este ejemplo, se aprovisiona el volumen con el tipo de almacenamiento que se creó en el ejemplo de la sección ["Ejemplo de clases de almacenamiento Kubernetes para puestas en marcha de IA en ONTAP"](#), paso 2.

An accessModes valor de ReadWriteMany Se especifica en el siguiente archivo de definición de PVC de ejemplo. Para obtener más información acerca de accessMode consulte ["Documentación oficial sobre Kubernetes"](#).

```

$ cat << EOF > ./pvc-tensorflow-results.yaml
kind: PersistentVolumeClaim
apiVersion: v1
metadata:
  name: tensorflow-results
spec:
  accessModes:
    - ReadWriteMany
  resources:
    requests:
      storage: 1Gi
  storageClassName: ontap-ai-flexvols-retain
EOF
$ kubectl create -f ./pvc-tensorflow-results.yaml
persistentvolumeclaim/tensorflow-results created
$ kubectl get pvc
NAME                                STATUS      VOLUME                                     CAPACITY   ACCESS MODES   STORAGECLASS          AGE
pb-fg-all-iface1                    Bound      default-pb-fg-all-iface1-7d9f1          10995116277760    ROX            ontap-ai-flexgroups-retain-iface1    26h
pb-fg-all-iface2                    Bound      default-pb-fg-all-iface2-85aee          10995116277760    ROX            ontap-ai-flexgroups-retain-iface2    26h
tensorflow-results                  Bound      default-tensorflow-results-2fd60        1073741824        RWX            ontap-ai-flexvols-retain            25h

```

### Ejemplo de tareas de alto rendimiento para puestas en marcha de ONTAP AI

En esta sección, se incluyen ejemplos de varias tareas de alto rendimiento que se pueden ejecutar cuando se pone en marcha Kubernetes en un pod ONTAP AI.

#### Ejemplo de tareas de alto rendimiento para puestas en marcha de ONTAP AI

En esta sección, se incluyen ejemplos de varias tareas de alto rendimiento que se pueden ejecutar cuando se pone en marcha Kubernetes en un pod ONTAP AI.

#### Ejecute una carga de trabajo de IA de un solo nodo

Para ejecutar una tarea DE IA y ML de un solo nodo en el clúster de Kubernetes, realice las siguientes tareas desde el host de puesta en marcha. Con Trident, puede crear de forma rápida y sencilla un volumen de datos, con potencialmente petabytes de datos al que se puede acceder una carga de trabajo de Kubernetes. Para que un volumen de datos de este tipo sea accesible desde un pod de Kubernetes, solo tiene que especificar una RVP en la definición del pod. Este paso es una operación nativa de Kubernetes, no se necesita experiencia en NetApp.



En esta sección se supone que ya ha realizado un contenedor (en el formato de contenedor de Docker) con la carga de trabajo específica DE IA y ML que intenta ejecutar en su clúster de Kubernetes.

1. Los siguientes comandos de ejemplo muestran la creación de un trabajo de Kubernetes para una carga de trabajo de prueba de ImageNET que utiliza el conjunto de datos de TensorFlow. Para obtener más información acerca del conjunto de datos ImageNET, consulte ["Sitio web de ImageNET"](#).

Este trabajo de ejemplo solicita ocho GPU y, por lo tanto, puede ejecutarse en un solo nodo de trabajo de GPU con ocho o más GPU. Este trabajo de ejemplo se puede enviar en un clúster para el que no hay un nodo de trabajo con ocho o más GPU o esté ocupado actualmente con otra carga de trabajo. Si es así, el trabajo permanece en estado pendiente hasta que dicho nodo de trabajo esté disponible.

Además, para maximizar el ancho de banda de almacenamiento, el volumen que contiene los datos de entrenamiento necesarios se monta dos veces en el pod que crea este trabajo. Otro volumen también se monta en el pod. Este segundo volumen se utilizará para almacenar resultados y métricas. Estos volúmenes se hacen referencia en la definición de trabajo utilizando los nombres de las RVP. Para obtener más información sobre los trabajos de Kubernetes, consulte ["Documentación oficial sobre Kubernetes"](#).

An `emptyDir` volumen con un `medium` valor de `Memory` está montado en `/dev/shm` en el pod que crea este trabajo de ejemplo. El tamaño predeterminado de `/dev/shm` El volumen virtual que se crea automáticamente mediante el tiempo de ejecución del contenedor Docker puede en ocasiones ser insuficiente para las necesidades de TensorFlow. Montaje de un `emptyDir` volumen como en el ejemplo siguiente proporciona un tamaño suficiente `/dev/shm` volumen virtual. Para obtener más información acerca de `emptyDir` volúmenes, consulte ["Documentación oficial sobre Kubernetes"](#).

El contenedor único que se especifica en esta definición de trabajo de ejemplo se proporciona un `securityContext > privileged` valor de `true`. Este valor significa que el contenedor tiene acceso raíz en el host de forma efectiva. Esta anotación se utiliza en este caso porque la carga de trabajo específica que se está ejecutando requiere acceso raíz. Específicamente, una operación de caché clara que ejecuta la carga de trabajo requiere acceso raíz. Si esto o no `privileged: true` la anotación es necesaria depende de los requisitos de la carga de trabajo específica que se esté ejecutando.

```
$ cat << EOF > ./netapp-tensorflow-single-imagenet.yaml
apiVersion: batch/v1
kind: Job
metadata:
  name: netapp-tensorflow-single-imagenet
spec:
  backoffLimit: 5
  template:
    spec:
      volumes:
      - name: dshm
        emptyDir:
          medium: Memory
      - name: testdata-iface1
        persistentVolumeClaim:
          claimName: pb-fg-all-iface1
      - name: testdata-iface2
```

```

    persistentVolumeClaim:
      claimName: pb-fg-all-iface2
- name: results
  persistentVolumeClaim:
    claimName: tensorflow-results
containers:
- name: netapp-tensorflow-py2
  image: netapp/tensorflow-py2:19.03.0
  command: ["python", "/netapp/scripts/run.py", "--
dataset_dir=/mnt/mount_0/dataset/imagenet", "--dgx_version=dgx1", "--
num_devices=8"]
  resources:
    limits:
      nvidia.com/gpu: 8
  volumeMounts:
- mountPath: /dev/shm
  name: dshm
- mountPath: /mnt/mount_0
  name: testdata-iface1
- mountPath: /mnt/mount_1
  name: testdata-iface2
- mountPath: /tmp
  name: results
  securityContext:
    privileged: true
  restartPolicy: Never
EOF
$ kubectl create -f ./netapp-tensorflow-single-imagenet.yaml
job.batch/netapp-tensorflow-single-imagenet created
$ kubectl get jobs
NAME                                COMPLETIONS   DURATION   AGE
netapp-tensorflow-single-imagenet   0/1            24s        24s

```

2. Confirme que el trabajo que ha creado en el paso 1 se está ejecutando correctamente. El siguiente comando de ejemplo confirma que se creó un solo pod para el trabajo, tal como se especifica en la definición de trabajos, y que este pod se ejecuta actualmente en uno de los nodos de trabajo de la GPU.

```

$ kubectl get pods -o wide
NAME                                READY   STATUS
RESTARTS   AGE
IP          NODE          NOMINATED NODE
netapp-tensorflow-single-imagenet-m7x92   1/1     Running   0
3m         10.233.68.61   10.61.218.154   <none>

```

3. Confirme que el trabajo que ha creado en el paso 1 se ha completado correctamente. Los siguientes

comandos de ejemplo confirman que el trabajo se ha completado correctamente.

```
$ kubectl get jobs
NAME                                COMPLETIONS  DURATION
AGE
netapp-tensorflow-single-imagenet  1/1           5m42s
10m
$ kubectl get pods
NAME                                READY  STATUS
RESTARTS  AGE
netapp-tensorflow-single-imagenet-m7x92  0/1    Completed
0        11m
$ kubectl logs netapp-tensorflow-single-imagenet-m7x92
[netapp-tensorflow-single-imagenet-m7x92:00008] PMIX ERROR: NO-
PERMISSIONS in file gds_dstore.c at line 702
[netapp-tensorflow-single-imagenet-m7x92:00008] PMIX ERROR: NO-
PERMISSIONS in file gds_dstore.c at line 711
Total images/sec = 6530.59125
===== Clean Cache !!! =====
mpirun -allow-run-as-root -np 1 -H localhost:1 bash -c 'sync; echo 1 >
/proc/sys/vm/drop_caches'
=====
mpirun -allow-run-as-root -np 8 -H localhost:8 -bind-to none -map-by
slot -x NCCL_DEBUG=INFO -x LD_LIBRARY_PATH -x PATH python
/netapp/tensorflow/benchmarks_190205/scripts/tf_cnn_benchmarks/tf_cnn_be
nchmarks.py --model=resnet50 --batch_size=256 --device=gpu
--force_gpu_compatible=True --num_intra_threads=1 --num_inter_threads=48
--variable_update=horovod --batch_group_size=20 --num_batches=500
--nodistortions --num_gpus=1 --data_format=NCHW --use_fp16=True
--use_tf_layers=False --data_name=imagenet --use_datasets=True
--data_dir=/mnt/mount_0/dataset/imagenet
--datasets_parallel_interleave_cycle_length=10
--datasets_sloppy_parallel_interleave=False --num_mounts=2
--mount_prefix=/mnt/mount_%d --datasets_prefetch_buffer_size=2000
--datasets_use_prefetch=True --datasets_num_private_threads=4
--horovod_device=gpu >
/tmp/20190814_105450_tensorflow_horovod_rdma_resnet50_gpu_8_256_b500_ima
genet_nodistort_fp16_r10_m2_nockpt.txt 2>&1
```

4. **Opcional:** limpiar artefactos de trabajo. Los siguientes comandos de ejemplo muestran la eliminación del objeto de trabajo creado en el paso 1.

Cuando se elimina el objeto de trabajo, Kubernetes elimina automáticamente todos los pods asociados.

```

$ kubectl get jobs
NAME                                     COMPLETIONS   DURATION
AGE
netapp-tensorflow-single-imagenet      1/1            5m42s
10m
$ kubectl get pods
NAME                                     READY   STATUS
RESTARTS   AGE
netapp-tensorflow-single-imagenet-m7x92 0/1     Completed
0         11m
$ kubectl delete job netapp-tensorflow-single-imagenet
job.batch "netapp-tensorflow-single-imagenet" deleted
$ kubectl get jobs
No resources found.
$ kubectl get pods
No resources found.

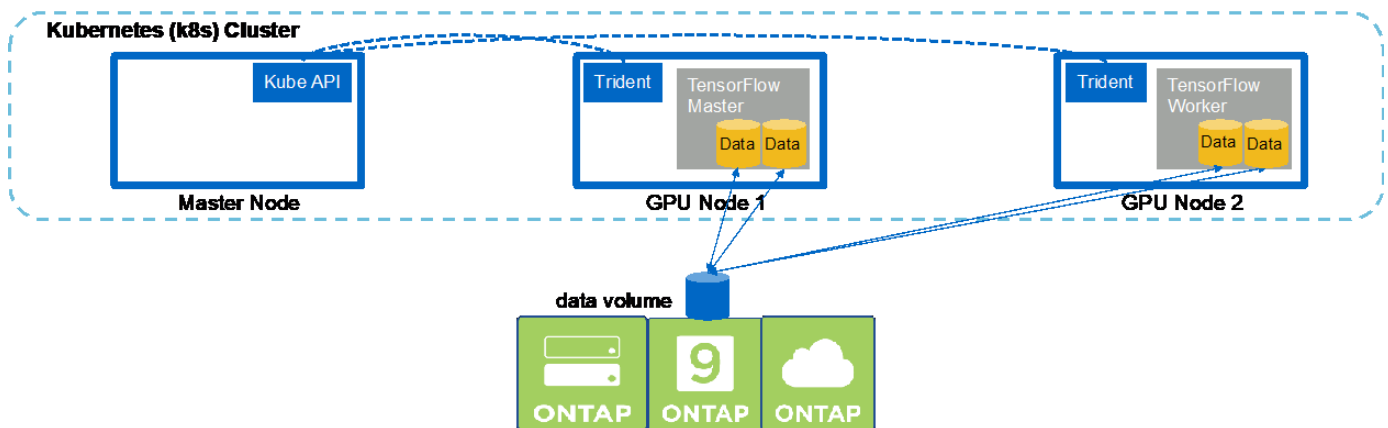
```

### Ejecute una carga de trabajo de IA distribuida síncrona

Para ejecutar un trabajo de IA y ML multinodo síncrono en un clúster de Kubernetes, lleve a cabo las siguientes tareas en el host de saltos de la puesta en marcha. Este proceso le permite aprovechar los datos almacenados en un volumen de NetApp y utilizar más GPU de las que puede proporcionar un único nodo de trabajo. Consulte la siguiente figura para obtener una descripción de un trabajo de IA distribuido síncrono.



Los trabajos distribuidos síncronos pueden ayudar a aumentar el rendimiento y la precisión de la formación en comparación con los trabajos distribuidos de manera asíncrona. Un análisis de los pros y los contras de los trabajos síncronos frente a los trabajos asíncronos está fuera del alcance de este documento.



1. En los siguientes comandos de ejemplo, se muestra la creación de un trabajador que participa en la ejecución síncrona y distribuida de la misma tarea de prueba de rendimiento TensorFlow que se ejecutó en un solo nodo en el ejemplo de la sección ["Ejecute una carga de trabajo de IA de un solo nodo"](#). En este ejemplo específico, sólo se implementa un único trabajador porque el trabajo se ejecuta en dos nodos de



trabajo.

Esta puesta en marcha de trabajo de ejemplo solicita ocho GPU y, por lo tanto, puede ejecutarse en un único nodo de trabajo de GPU con ocho o más GPU. Si los nodos de trabajo de la GPU tienen más de ocho GPU, para maximizar el rendimiento, es posible que desee aumentar este número para ser igual al número de GPU que disponen los nodos de trabajo. Para obtener más información sobre las puestas en marcha de Kubernetes, consulte ["Documentación oficial sobre Kubernetes"](#).

En este ejemplo se crea una puesta en marcha de Kubernetes, ya que este trabajador en contenedor específico nunca lo completaría por sí solo. Por lo tanto, no tiene sentido implementarlo usando la construcción de trabajos de Kubernetes. Si su trabajador está diseñado o escrito para completar por sí solo, entonces podría tener sentido utilizar la construcción del trabajo para desplegar a su trabajador.

El POD especificado en esta especificación de implementación de ejemplo recibe una `hostNetwork` valor de `true`. Este valor significa que el pod utiliza la pila de red del nodo de trabajo del host en lugar de la pila de red virtual que Kubernetes suele crear para cada pod. Esta anotación se utiliza en este caso porque la carga de trabajo específica depende de Open MPI, NCCL y Horovod para ejecutar la carga de trabajo de forma síncrona distribuida. Por lo tanto, requiere acceso a la pila de red del host. Un debate sobre Open MPI, NCCL y Horovod está fuera del alcance de este documento. Si esto o no `hostNetwork: true` la anotación es necesaria depende de los requisitos de la carga de trabajo específica que se esté ejecutando. Para obtener más información acerca de `hostNetwork` consulte ["Documentación oficial sobre Kubernetes"](#).

```
$ cat << EOF > ./netapp-tensorflow-multi-imagenet-worker.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: netapp-tensorflow-multi-imagenet-worker
spec:
  replicas: 1
  selector:
    matchLabels:
      app: netapp-tensorflow-multi-imagenet-worker
  template:
    metadata:
      labels:
        app: netapp-tensorflow-multi-imagenet-worker
    spec:
      hostNetwork: true
      volumes:
      - name: dshm
        emptyDir:
          medium: Memory
      - name: testdata-iface1
        persistentVolumeClaim:
          claimName: pb-fg-all-iface1
      - name: testdata-iface2
        persistentVolumeClaim:
          claimName: pb-fg-all-iface2
```

```

- name: results
  persistentVolumeClaim:
    claimName: tensorflow-results
containers:
- name: netapp-tensorflow-py2
  image: netapp/tensorflow-py2:19.03.0
  command: ["bash", "/netapp/scripts/start-slave-multi.sh",
"22122"]
  resources:
    limits:
      nvidia.com/gpu: 8
  volumeMounts:
  - mountPath: /dev/shm
    name: dshm
  - mountPath: /mnt/mount_0
    name: testdata-iface1
  - mountPath: /mnt/mount_1
    name: testdata-iface2
  - mountPath: /tmp
    name: results
  securityContext:
    privileged: true
EOF
$ kubectl create -f ./netapp-tensorflow-multi-imagenet-worker.yaml
deployment.apps/netapp-tensorflow-multi-imagenet-worker created
$ kubectl get deployments
NAME                                DESIRED   CURRENT   UP-TO-DATE
AVAILABLE   AGE
netapp-tensorflow-multi-imagenet-worker  1         1         1
1         4s

```

2. Confirme que el despliegue del trabajador que creó en el paso 1 se inició correctamente. Los siguientes comandos de ejemplo confirman que se creó un solo pod de trabajadores para la implementación, tal y como se indica en la definición de la puesta en marcha, y que este pod se ejecuta actualmente en uno de los nodos de trabajo de la GPU.

```

$ kubectl get pods -o wide
NAME                                READY
STATUS   RESTARTS   AGE   IP            NODE            NOMINATED NODE
netapp-tensorflow-multi-imagenet-worker-654fc7f486-v6725  1/1
Running   0          60s   10.61.218.154  10.61.218.154   <none>
$ kubectl logs netapp-tensorflow-multi-imagenet-worker-654fc7f486-v6725
22122

```

3. Cree un trabajo de Kubernetes en un maestro que se inicia, participe y realice un seguimiento de la ejecución de un trabajo de varios nodos síncronos. Los siguientes comandos de ejemplo crean un maestro que inicia sesión, participa en y realiza un seguimiento de la ejecución síncrona distribuida de la misma tarea de prueba de rendimiento TensorFlow que se ejecutó en un solo nodo del ejemplo de la sección ["Ejecute una carga de trabajo de IA de un solo nodo"](#).

Este trabajo maestro de ejemplo solicita ocho GPU y, por lo tanto, puede ejecutarse en un único nodo de trabajo de GPU con ocho o más GPU. Si los nodos de trabajo de la GPU tienen más de ocho GPU, para maximizar el rendimiento, es posible que desee aumentar este número para ser igual al número de GPU que disponen los nodos de trabajo.

El POD maestro especificado en esta definición de trabajo de ejemplo recibe una `hostNetwork` valor de `true`, así como se le dio a la cápsula de trabajo un `hostNetwork` valor de `true` en el paso 1. Consulte el paso 1 para obtener más información acerca de por qué es necesario este valor.

```
$ cat << EOF > ./netapp-tensorflow-multi-imagenet-master.yaml
apiVersion: batch/v1
kind: Job
metadata:
  name: netapp-tensorflow-multi-imagenet-master
spec:
  backoffLimit: 5
  template:
    spec:
      hostNetwork: true
      volumes:
      - name: dshm
        emptyDir:
          medium: Memory
      - name: testdata-iface1
        persistentVolumeClaim:
          claimName: pb-fg-all-iface1
      - name: testdata-iface2
        persistentVolumeClaim:
          claimName: pb-fg-all-iface2
      - name: results
        persistentVolumeClaim:
          claimName: tensorflow-results
      containers:
      - name: netapp-tensorflow-py2
        image: netapp/tensorflow-py2:19.03.0
        command: ["python", "/netapp/scripts/run.py", "--
dataset_dir=/mnt/mount_0/dataset/imagenet", "--port=22122", "--
num_devices=16", "--dgx_version=dgx1", "--
nodes=10.61.218.152,10.61.218.154"]
        resources:
          limits:
            nvidia.com/gpu: 8
```

```

    volumeMounts:
      - mountPath: /dev/shm
        name: dshm
      - mountPath: /mnt/mount_0
        name: testdata-iface1
      - mountPath: /mnt/mount_1
        name: testdata-iface2
      - mountPath: /tmp
        name: results
    securityContext:
      privileged: true
    restartPolicy: Never
EOF
$ kubectl create -f ./netapp-tensorflow-multi-imagenet-master.yaml
job.batch/netapp-tensorflow-multi-imagenet-master created
$ kubectl get jobs
NAME                                COMPLETIONS   DURATION   AGE
netapp-tensorflow-multi-imagenet-master  0/1            25s        25s

```

4. Confirme que el trabajo maestro que creó en el paso 3 se está ejecutando correctamente. El siguiente comando de ejemplo confirma que se creó un único pod maestro para el trabajo, tal como se indica en la definición de trabajos, y que este pod se ejecuta actualmente en uno de los nodos de trabajo de la GPU. También debe ver que el pod de trabajo que originalmente vio en el paso 1 sigue en ejecución y que los pods maestro y trabajador se ejecutan en diferentes nodos.

```

$ kubectl get pods -o wide
NAME                                READY
STATUS   RESTARTS   AGE   IP              NODE              NOMINATED NODE
netapp-tensorflow-multi-imagenet-master-ppwwj  1/1
Running   0           45s   10.61.218.152   10.61.218.152    <none>
netapp-tensorflow-multi-imagenet-worker-654fc7f486-v6725  1/1
Running   0           26m   10.61.218.154   10.61.218.154    <none>

```

5. Confirme que el trabajo maestro que ha creado en el paso 3 se ha completado correctamente. Los siguientes comandos de ejemplo confirman que el trabajo se ha completado correctamente.

```

$ kubectl get jobs
NAME                                COMPLETIONS   DURATION   AGE
netapp-tensorflow-multi-imagenet-master  1/1            5m50s      9m18s
$ kubectl get pods
NAME                                READY
STATUS   RESTARTS   AGE   IP              NODE              NOMINATED NODE
netapp-tensorflow-multi-imagenet-master-ppwwj  0/1
Completed 0           9m38s

```

```

netapp-tensorflow-multi-imagenet-worker-654fc7f486-v6725 1/1
Running 0 35m
$ kubectl logs netapp-tensorflow-multi-imagenet-master-ppwwj
[10.61.218.152:00008] WARNING: local probe returned unhandled
shell:unknown assuming bash
rm: cannot remove '/lib': Is a directory
[10.61.218.154:00033] PMIX ERROR: NO-PERMISSIONS in file gds_dstore.c at
line 702
[10.61.218.154:00033] PMIX ERROR: NO-PERMISSIONS in file gds_dstore.c at
line 711
[10.61.218.152:00008] PMIX ERROR: NO-PERMISSIONS in file gds_dstore.c at
line 702
[10.61.218.152:00008] PMIX ERROR: NO-PERMISSIONS in file gds_dstore.c at
line 711
Total images/sec = 12881.33875
===== Clean Cache !!! =====
mpirun -allow-run-as-root -np 2 -H 10.61.218.152:1,10.61.218.154:1 -mca
pml obl -mca btl ^openib -mca btl_tcp_if_include enpls0f0 -mca
plm_rsh_agent ssh -mca plm_rsh_args "-p 22122" bash -c 'sync; echo 1 >
/proc/sys/vm/drop_caches'
=====
mpirun -allow-run-as-root -np 16 -H 10.61.218.152:8,10.61.218.154:8
-bind-to none -map-by slot -x NCCL_DEBUG=INFO -x LD_LIBRARY_PATH -x PATH
-mca pml obl -mca btl ^openib -mca btl_tcp_if_include enpls0f0 -x
NCCL_IB_HCA=mlx5 -x NCCL_NET_GDR_READ=1 -x NCCL_IB_SL=3 -x
NCCL_IB_GID_INDEX=3 -x
NCCL_SOCKET_IFNAME=enp5s0.3091,enp12s0.3092,enp132s0.3093,enp139s0.3094
-x NCCL_IB_CUDA_SUPPORT=1 -mca orte_base_help_aggregate 0 -mca
plm_rsh_agent ssh -mca plm_rsh_args "-p 22122" python
/netapp/tensorflow/benchmarks_190205/scripts/tf_cnn_benchmarks/tf_cnn_be
nchmarks.py --model=resnet50 --batch_size=256 --device=gpu
--force_gpu_compatible=True --num_intra_threads=1 --num_inter_threads=48
--variable_update=horovod --batch_group_size=20 --num_batches=500
--nodistortions --num_gpus=1 --data_format=NCHW --use_fp16=True
--use_tf_layers=False --data_name=imagenet --use_datasets=True
--data_dir=/mnt/mount_0/dataset/imagenet
--datasets_parallel_interleave_cycle_length=10
--datasets_sloppy_parallel_interleave=False --num_mounts=2
--mount_prefix=/mnt/mount_%d --datasets_prefetch_buffer_size=2000 --
datasets_use_prefetch=True --datasets_num_private_threads=4
--horovod_device=gpu >
/tmp/20190814_161609_tensorflow_horovod_rdma_resnet50_gpu_16_256_b500_im
agenet_nodistort_fp16_r10_m2_nockpt.txt 2>&1

```

6. Elimine la implementación del trabajador cuando ya no la necesite. Los siguientes comandos de ejemplo muestran la eliminación del objeto de implementación de trabajo que se creó en el paso 1.

Cuando se elimina el objeto de implementación de trabajo, Kubernetes elimina automáticamente todos los pods de trabajador asociados.

```
$ kubectl get deployments
NAME                                DESIRED    CURRENT    UP-TO-DATE
AVAILABLE    AGE
netapp-tensorflow-multi-imagenet-worker  1          1          1
1          43m
$ kubectl get pods
NAME                                READY
STATUS    RESTARTS    AGE
netapp-tensorflow-multi-imagenet-master-ppwwj  0/1
Completed  0          17m
netapp-tensorflow-multi-imagenet-worker-654fc7f486-v6725  1/1
Running    0          43m
$ kubectl delete deployment netapp-tensorflow-multi-imagenet-worker
deployment.extensions "netapp-tensorflow-multi-imagenet-worker" deleted
$ kubectl get deployments
No resources found.
$ kubectl get pods
NAME                                READY    STATUS
RESTARTS    AGE
netapp-tensorflow-multi-imagenet-master-ppwwj  0/1      Completed  0
18m
```

7. **Opcional:** Limpie los artefactos del trabajo maestro. Los siguientes comandos de ejemplo muestran la eliminación del objeto de trabajo maestro que se creó en el paso 3.

Cuando se elimina el objeto de trabajo maestro, Kubernetes elimina automáticamente todos los pods maestros asociados.

```
$ kubectl get jobs
NAME                                COMPLETIONS    DURATION    AGE
netapp-tensorflow-multi-imagenet-master  1/1            5m50s      19m
$ kubectl get pods
NAME                                READY    STATUS
RESTARTS    AGE
netapp-tensorflow-multi-imagenet-master-ppwwj  0/1      Completed  0
19m
$ kubectl delete job netapp-tensorflow-multi-imagenet-master
job.batch "netapp-tensorflow-multi-imagenet-master" deleted
$ kubectl get jobs
No resources found.
$ kubectl get pods
No resources found.
```

## Pruebas de rendimiento

Hemos realizado una comparativa simple del rendimiento como parte de la creación de esta solución. Hemos ejecutado varias tareas estándar de prueba de rendimiento de IA de NetApp mediante Kubernetes, y comparamos los resultados de las pruebas de rendimiento con ejecuciones realizadas mediante un sencillo comando Docker run. No observamos diferencias notables en cuanto al rendimiento. Por lo tanto, llegamos a la conclusión de que el uso de Kubernetes para orquestar trabajos de entrenamiento de IA en contenedores no afecta negativamente al rendimiento. Consulte la siguiente tabla para ver los resultados de nuestra comparativa de rendimiento.

Prueba de rendimiento	Conjunto de datos	Docker Run (imágenes/s)	Kubernetes (imágenes/s)
TensorFlow de un solo nodo	Datos sintéticos	6,667.2475	6,661.93125
TensorFlow de un solo nodo	ImageNET	6,570.2025	6,530.59125
TensorFlow distribuido síncrono de dos nodos	Datos sintéticos	13,213.70625	13,218.288125
TensorFlow distribuido síncrono de dos nodos	ImageNET	12,941.69125	12,881.33875

## Conclusión

Empresas y organizaciones de todos los tamaños y sectores se están decantando por la inteligencia artificial (IA), el aprendizaje automático (ML) y el aprendizaje profundo (DL) para solucionar problemas del mundo real, ofrecer productos y servicios innovadores y obtener una ventaja en un mercado cada vez más competitivo. A medida que las organizaciones aumentan el uso de la IA, EL ML y el AP, deben hacer frente a numerosos retos, como la escalabilidad de la carga de trabajo y la disponibilidad de los datos. Estos desafíos se pueden abordar mediante el uso de la solución de plano de control de IA de NetApp.

Esta solución le permite clonar rápidamente un espacio de nombres de los datos. Además, le permite definir e implementar flujos de trabajo de entrenamiento de IA, ML y DL que incorporan la creación casi instantánea de bases de datos y modelos para su trazabilidad y versionado. Con esta solución, puede rastrear cada entrenamiento de modelo único que se ejecuta en los conjuntos de datos exactos con los que se ha entrenado el modelo y/o validado. Por último, esta solución le permite aprovisionar rápidamente espacios de trabajo de los portátiles Jupyter con acceso a conjuntos de datos masivos.

Dado que esta solución está dirigida a científicos e ingenieros de datos, se requiere una experiencia mínima en NetApp o ONTAP de NetApp. Con esta solución, las funciones de gestión de datos se pueden ejecutar utilizando interfaces y herramientas sencillas y conocidas. Además, esta solución utiliza componentes totalmente de código abierto y libres. Por lo tanto, si ya dispone de almacenamiento de NetApp en su entorno, puede implantar esta solución hoy mismo. Si desea probar esta solución pero no dispone de almacenamiento de NetApp, visite ["cloud.netapp.com"](https://cloud.netapp.com)Y puede tener una solución de almacenamiento de NetApp basada en cloud en un instante.

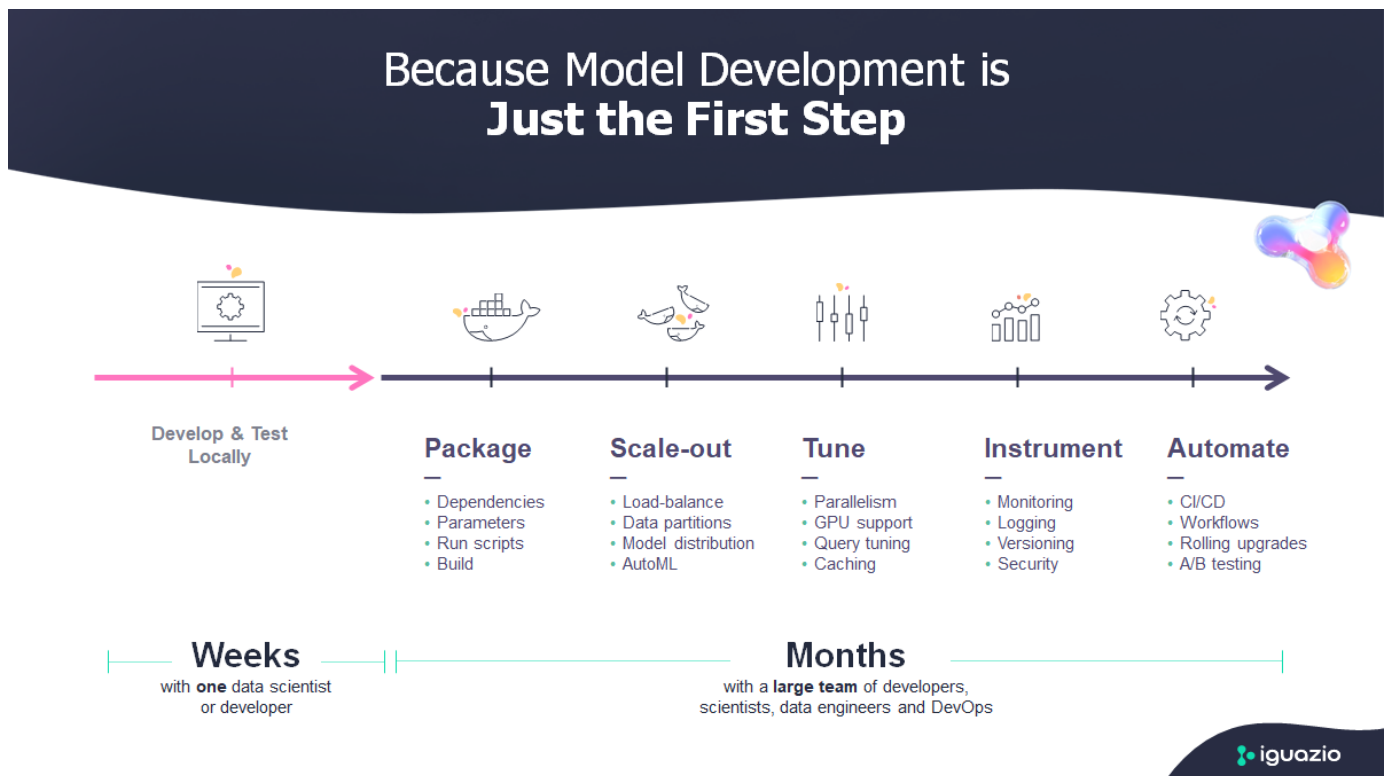
## MLRun Pipeline con Iguazio

### TR-4834: NetApp e Iguazio para el suministro de MLRun

Rick Huang, David Arnette, NetApp Marcelo Litovsky, Iguazio

Este documento contiene los detalles de la canalización MLRun usando NetApp ONTAP AI, el plano de control de IA de NetApp, el software Cloud Volumes de NetApp y la plataforma de ciencia de datos de Iguazio. Utilizamos la función Nuclio serverless, volúmenes persistentes de Kubernetes, volúmenes Cloud de NetApp, copias Snapshot de NetApp, consola Grafana, Y otros servicios en la plataforma Iguazio para crear una canalización de datos de extremo a extremo para la simulación de la detección de fallos de red. Hemos integrado las tecnologías de Iguazio y NetApp para permitir una rápida puesta en marcha de modelos, replicación de datos y supervisión de la producción tanto en las instalaciones como en el cloud.

El trabajo de un científico de datos debe centrarse en el entrenamiento y el ajuste de los modelos de aprendizaje automático y de inteligencia artificial (IA). Sin embargo, según una investigación realizada por Google, los científicos de datos gastan aproximadamente un 80 % de su tiempo en averiguar cómo hacer que sus modelos funcionen con aplicaciones empresariales y se ejecutan a escala, como se muestra en la siguiente imagen, en la que se ilustra el desarrollo de modelos en el flujo de trabajo de IA/ML.



Para gestionar proyectos de IA y ML integrales, se necesita una comprensión más amplia de los componentes empresariales. Aunque DevOps ha retomado la definición, la integración y la puesta en marcha de estos tipos de componentes, las operaciones de aprendizaje automático siguen un flujo similar que incluye proyectos de IA/ML. Para hacerse una idea de lo que puede tocar una canalización de IA/ML integral en la empresa, consulte la siguiente lista de componentes requeridos:

- Reducida



- Redes
- Oracle
- Sistemas de ficheros
- Contenedores
- Integración continua y canalización de puesta en marcha continua (CI/CD)
- Entorno de desarrollo integrado (IDE)
- Seguridad
- Políticas de acceso a los datos
- Hardware subyacente
- Cloud
- Virtualización
- Conjuntos de herramientas y bibliotecas de ciencia de datos

En este artículo mostramos cómo la asociación entre NetApp y Iguazio simplifica drásticamente el desarrollo de una canalización completa de IA/ML. Esta simplificación acelera el plazo de comercialización de todas sus aplicaciones de IA/ML.

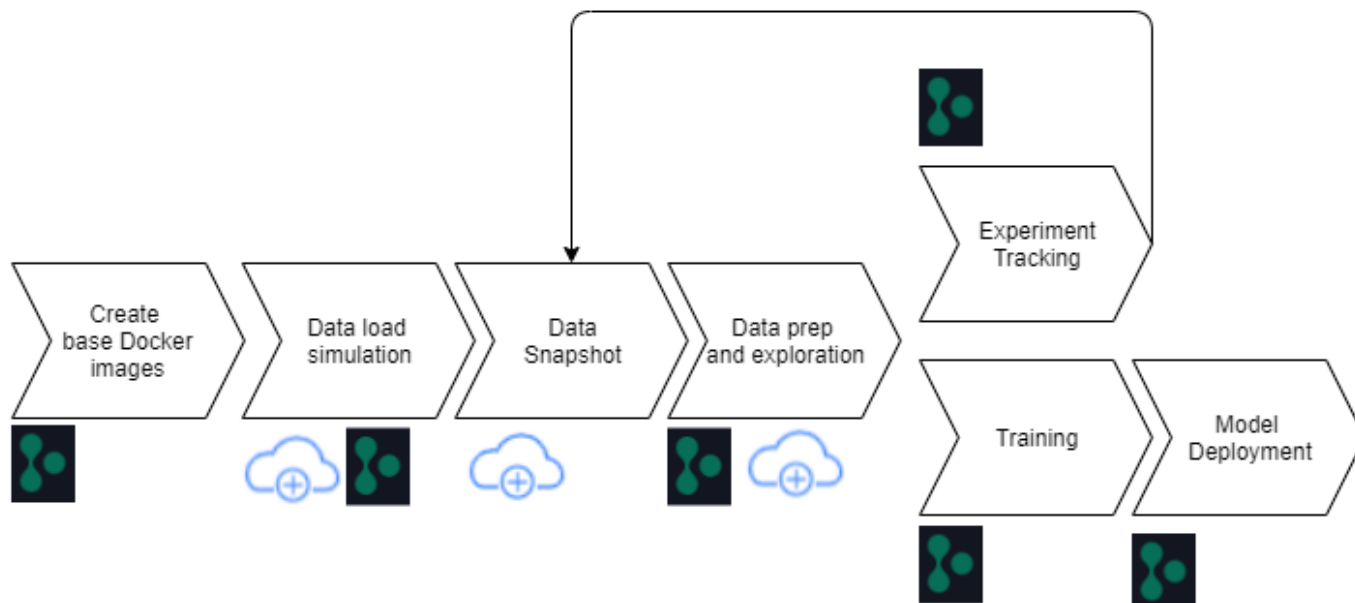
#### **Público objetivo**

El mundo de la ciencia de datos tiene múltiples disciplinas de tecnología de la información y negocio.

- El científico de datos necesita la flexibilidad para poder usar las herramientas y las bibliotecas de elección.
- El ingeniero de datos necesita saber cómo fluyen los datos y dónde residen.
- Un ingeniero de DevOps necesita herramientas para integrar nuevas aplicaciones de IA/ML en sus canalizaciones de CI/CD.
- Los usuarios empresariales quieren tener acceso a aplicaciones de IA/ML. Describimos cómo NetApp y Iguazio ayudan a cada una de estas funciones a aportar valor empresarial con nuestras plataformas.

#### **Descripción general de la solución**

Esta solución sigue el ciclo de vida de una aplicación de IA/ML. Empezamos con la labor de los científicos de datos para definir los diferentes pasos necesarios para preparar datos y entrenar y poner en marcha modelos. Seguimos con el trabajo necesario para crear una canalización completa con la capacidad de realizar un seguimiento de los artefactos, experimentar con la ejecución y desplegar en Kubeflow. Para completar el ciclo completo, integramos la canalización con NetApp Cloud Volumes para permitir el versionado de datos, tal y como se puede ver en la siguiente imagen.



## Visión general de la tecnología

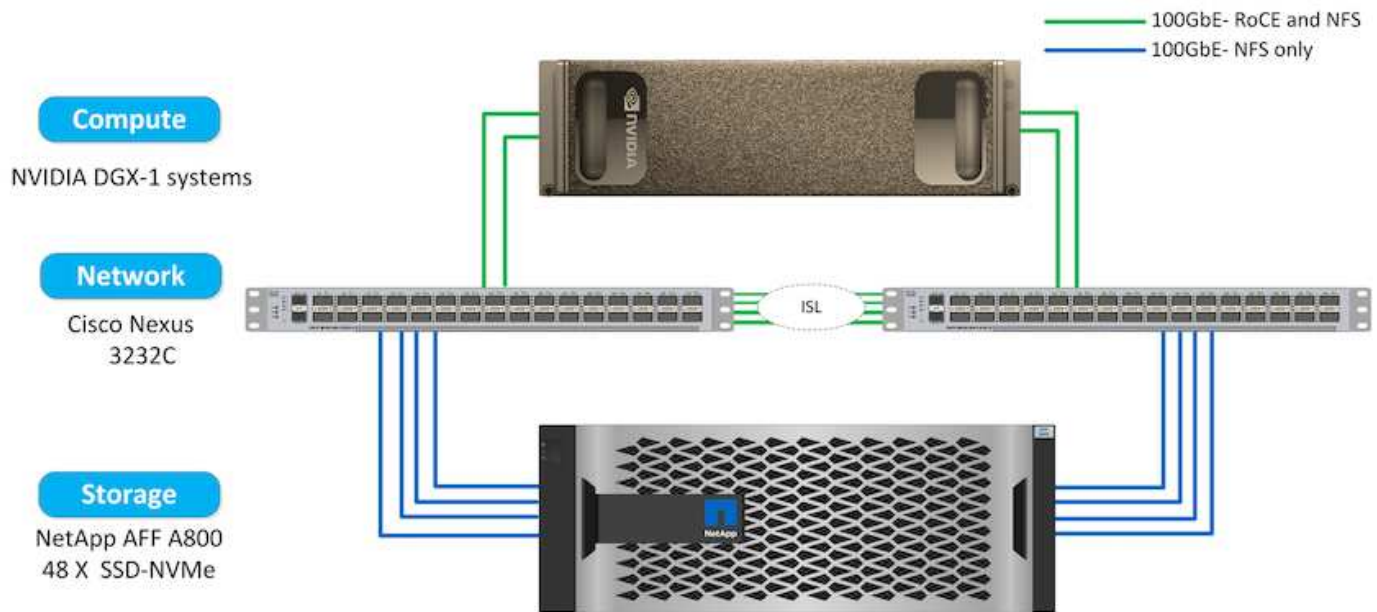
### Información general de NetApp

NetApp es un referente en materia de datos para el cloud híbrido. NetApp proporciona una gama completa de servicios de datos en el cloud híbrido que simplifican la gestión de aplicaciones y datos en entornos de cloud y en las instalaciones para acelerar la transformación digital. Junto con nuestros partners, NetApp permite a organizaciones globales aprovechar al máximo todo el potencial de sus datos para ampliar los puntos de contacto con los clientes, fomentar una mayor innovación y optimizar sus operaciones.

### ONTAP AI de NetApp

ONTAP AI de NetApp, impulsado por los sistemas NVIDIA DGX y el almacenamiento all-flash conectado al cloud de NetApp, optimiza el flujo de datos con total confianza y acelera el análisis, la formación y la inferencia con su Data Fabric, que abarca desde el perímetro al núcleo y al cloud. Proporciona a las organizaciones DE TI una arquitectura que ofrece las siguientes ventajas:

- Elimina las complejidades de diseño
- Permite un escalado independiente de las capacidades de computación y almacenamiento
- Permite a los clientes empezar con poco y escalar sin problemas
- Ofrece una gama de opciones de almacenamiento para diferentes tipos de rendimiento y coste. ONTAP AI ofrece pilas de infraestructuras convergentes que incorporan NVIDIA DGX-1, un sistema de IA a escala de petaflops y switches Ethernet de alto rendimiento NVIDIA Mellanox para unificar cargas de trabajo de IA, simplificar la puesta en marcha y acelerar el retorno de la inversión. Para este informe técnico, hemos aprovechado la inteligencia artificial de ONTAP con un sistema de almacenamiento DGX-1 y AFF A800 de NetApp. En la siguiente imagen, se muestra la topología de ONTAP AI con el sistema DGX-1 utilizado en esta validación.



### Plano de control de IA de NetApp

El plano de control de IA de NetApp permite aprovechar la IA Y EL APRENDIZAJE AUTOMÁTICO con una solución que ofrece una escalabilidad extrema, una puesta en marcha optimizada y una disponibilidad de datos ininterrumpida. La solución de plano de control de IA integra Kubernetes y Kubeflow con una estructura de datos habilitada por NetApp. Kubernetes, la plataforma de orquestación de contenedores estándar del sector para puestas en marcha nativas del cloud, permite la escalabilidad y la portabilidad de las cargas de trabajo. Kubeflow es una plataforma de aprendizaje automático de código abierto que simplifica la gestión y la puesta en marcha, lo que permite a los desarrolladores hacer más ciencia de datos en menos tiempo. Una estructura de datos habilitada por NetApp ofrece una disponibilidad y portabilidad de datos sin concesiones para garantizar que sus datos están accesibles en toda la canalización, desde el perímetro al núcleo y al cloud. En este informe técnico se utiliza el plano de control de IA de NetApp en una canalización MLRun. La siguiente imagen muestra la página de gestión de clústeres de Kubernetes, donde se pueden tener distintos extremos para cada clúster. Conectamos volúmenes persistentes de NFS al clúster de Kubernetes, y las siguientes imágenes muestran un volumen persistente conectado al clúster donde **"Trident de NetApp"** ofrece compatibilidad con almacenamiento persistente y capacidades de gestión de datos.

## 4 Kubernetes Clusters



kubernetes

<https://3.20.111.39:6443>  
Cluster Endpointv1.15.5  
Cluster Version19.07.1  
Trident Version0  
Working Environments

kubernetes

<https://172.31.14.31:6443>  
Cluster Endpointv1.15.5  
Cluster Version19.07.1  
Trident Version1  
Working Environments

## Persistent Volumes for Kubernetes

## Connected with Kubernetes Cluster

Cloud Volumes ONTAP is connected to 1 Kubernetes cluster. [View Cluster](#)

You can connect another Kubernetes cluster to this Cloud Volumes ONTAP system. If the Kubernetes cluster is in a different network than Cloud Volumes ONTAP, specify a custom export policy to provide access to clients.

## Kubernetes Cluster

Custom Export Policy *(Optional)*

Select Kubernetes Cluster

kubernetes

Custom Export Policy

172.31.0.0/16


☒ Set as default storage class☒ NFS ☐ iSCSI

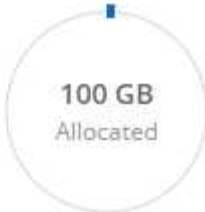
Connect

Cancel

## Volumes

4 Volumes | 300 GB Allocated | 1.43 GB Total Used

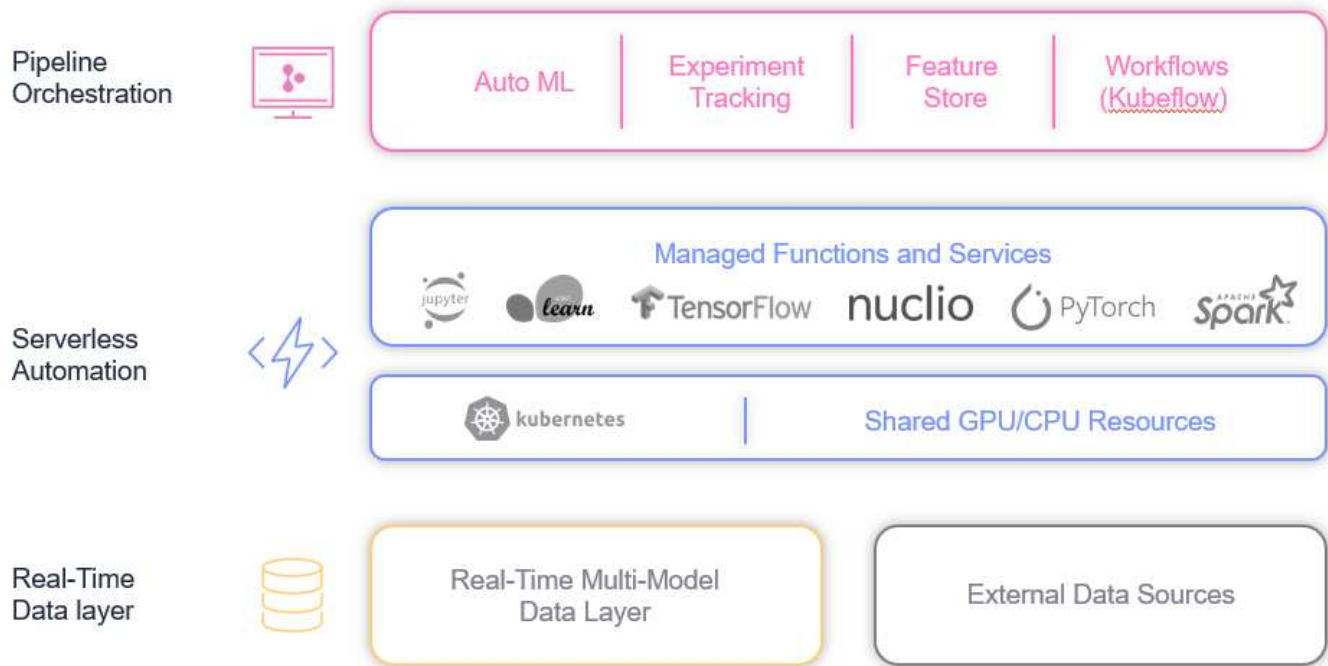

kubernetes\_trident\_pvc\_551720fa\_3758\_461...
ONLINE

INFO		CAPACITY	
Disk Type	GP2	 <div>1.25 GB EBS Used</div>	
Tiering Policy	None		
Backup	OFF		

### Visión General de Iguazio

La Plataforma de Ciencias de datos de Iguazio es una plataforma de ciencia de datos como servicio (PaaS) totalmente integrada y segura que simplifica el desarrollo, acelera el rendimiento, facilita la colaboración y aborda los retos operativos. Esta plataforma incorpora los siguientes componentes, y la Plataforma de Ciencias de datos de Iguazio se presenta en la siguiente imagen:

- Banco de trabajo de ciencia de datos que incluye portátiles Jupyter, motores de análisis integrados y paquetes Python
- Gestión de modelos con seguimiento de experimentos y funcionalidades de canalizaciones automatizadas
- Datos gestionados y servicios ML a través de un clúster de Kubernetes escalable
- Nuclio, un marco de funciones sin servidor en tiempo real
- Una capa de datos extremadamente rápida y segura compatible con SQL, NoSQL, bases de datos de series temporales, archivos (objetos sencillos) y streaming
- Integración con fuentes de datos de terceros como NetApp, Amazon S3, HDFS, bases de datos de SQL y protocolos de transmisión o mensajería
- Paneles en tiempo real basados en Grafana



## Requisitos de software y hardware

### Configuración de red

A continuación, se muestra el requisito de configuración de red para configurar en el cloud:

- El clúster de Iguazio y los volúmenes cloud de NetApp deben estar en el mismo cloud privado virtual.
- El administrador de la nube debe tener acceso al puerto 6443 en los nodos de la aplicación Iguazio.
- Utilizamos Amazon Web Services en este informe técnico. Sin embargo, los usuarios tienen la opción de poner en marcha la solución en cualquier proveedor de cloud. para las pruebas en las instalaciones en IA de ONTAP con NVIDIA DGX-1, hemos utilizado el servicio DNS alojado en Iguazio para mayor comodidad.

Los clientes deben poder acceder a los dominios DNS creados de forma dinámica. Los clientes pueden utilizar su propio DNS si lo desean.

### Requisitos de hardware

Puede instalar Iguazio en sus instalaciones en su propio clúster. Hemos verificado la solución en ONTAP AI de NetApp con un sistema NVIDIA DGX-1. La siguiente tabla enumera el hardware utilizado para probar esta solución.

Hardware subyacente	Cantidad
Sistemas DGX-1	1
Sistema AFF A800 de NetApp	1 pareja de alta disponibilidad (ha), incluye 2 controladoras y 48 SSD NVMe (3,8 TB o superior)
Switches de red Cisco Nexus 3232C	2

En la siguiente tabla se enumeran los componentes de software necesarios para las pruebas in situ:

De NetApp	Versión u otra información
Software de gestión de datos ONTAP de NetApp	9.7
Firmware de switch Cisco NX-OS	7.0(3)I6(1)
SO DGX DE NVIDIA	4.4 - Ubuntu 18.04 LTS
Plataforma contenedora Docker	19.03.5
Versión del contenedor	20.01-tf1-py2
Marco de aprendizaje automático	TensorFlow 1.15.0
Iguazio	La versión es 2.8 o posterior
Servidor ESX	6.5

Esta solución se ha probado totalmente con Iguazio versión 2.5 y NetApp Cloud Volumes ONTAP para AWS. El clúster de Iguazio y el software de NetApp se ejecutan en AWS.

De NetApp	Versión o tipo
Iguazio	La versión es 2.8 o posterior
Nodo de aplicación	M5.4xgrande
Nodo de datos	I3.4xgrande

## Resumen de casos de uso de predicción de fallos de dispositivos de red

Este caso de uso se basa en un cliente de Iguazio en el espacio de telecomunicaciones en Asia. Con 100 000 clientes empresariales y 125 000 eventos de interrupción de la red al año, era muy importante predecir y tomar medidas proactivas para evitar que los fallos de red afecten a los clientes. Esta solución ofreció las siguientes ventajas:

- Análisis predictivo para fallos de red
- Integración con un sistema de emisión de boletos
- Tomar medidas proactivas para evitar fallos de la red como resultado de esta implementación de Iguazio, el 60 % de los fallos se impidieron de forma proactiva.

## Descripción general de la configuración

Iguazio se puede instalar de forma local o en un proveedor de cloud.

### Instalación de Iguazio

El aprovisionamiento puede realizarse como servicio y gestionado por Iguazio o por el cliente. En ambos casos, Iguazio ofrece una aplicación de despliegue (Proventio) para implementar y gestionar clústeres.

Para la instalación en las instalaciones, consulte ["NVA-1121"](#) para la configuración de computación, red y almacenamiento. La implantación en las instalaciones de Iguazio es proporcionada por Iguazio sin costes adicionales para el cliente. Consulte ["esta página"](#) Para configuraciones de servidores DNS y SMTP. La página de instalación de Proventio se muestra a continuación.

New System (dev)

●

Installation Scenario

●

General

●

Clusters

●

Cloud

Bare metal / virtual machines

Installs the system on bare-metal or virtual-machine instances, pre-provisioned with prerequ...

AWS

Creates applicable compute/networking resources in AWS and installs the system on the in...

Azure

Creates applicable compute/networking resources in Azure and installs the system on the i...

AWS (pre-provisioned)

Installs the system on Amazon Web Services instances, manually provisioned beforehand

Azure (pre-provisioned)

Installs the system on Microsoft Azure instances, manually provisioned beforehand

Advanced

Show advanced options in the next steps

BACK

NEXT

## Configurar Kubernetes Cluster

Esta sección se divide en dos partes para la puesta en marcha en el cloud y en las instalaciones, respectivamente.

### Configuración de Kubernetes de puesta en marcha del cloud

A través de Cloud Manager de NetApp, puede definir la conexión al clúster de Kubernetes del Iguazio. Trident requiere acceso a varios recursos en el clúster para que el volumen esté disponible.

1. Para habilitar el acceso, obtenga el archivo de configuración de Kubernetes desde uno de los nodos de Iguazio. El archivo está ubicado en `/home/Iguazio/.kube/config`. Descargue este archivo en su escritorio.
2. Vaya a detectar clúster para configurar.



## 4 Kubernetes Clusters

 <b>kubernetes</b>			
 <a href="https://3.20.111.39:6443">https://3.20.111.39:6443</a> Cluster Endpoint	 v1.15.5 Cluster Version	 19.07.1 Trident Version	 0 Working Environments

 <b>kubernetes</b>			
 <a href="https://172.31.14.31:6443">https://172.31.14.31:6443</a> Cluster Endpoint	 v1.15.5 Cluster Version	 19.07.1 Trident Version	 1 Working Environments

3. Cargue el archivo de configuración de Kubernetes. Consulte la siguiente imagen.

## Upload Kubernetes Configuration File

Upload the Kubernetes configuration file (kubeconfig) so Cloud Manager can install Trident on the Kubernetes cluster.

Connecting Cloud Volumes ONTAP with a Kubernetes cluster enables users to request and manage persistent volumes using native Kubernetes interfaces and constructs. Users can take advantage of ONTAP's advanced data management features without having to know anything about it. Storage provisioning is enabled by using NetApp Trident.

Learn more about [Trident for Kubernetes](#).

Upload File

4. Ponga en marcha Trident y asocie un volumen con el clúster. Vea la siguiente imagen sobre la definición y asignación de un volumen persistente al clúster de Iguazio. Este proceso crea un volumen persistente (PV) en el clúster Kubernetes de Iguazio. Antes de poder usarlo, debe definir una solicitud de volumen persistente (PVC).

## Persistent Volumes for Kubernetes

### Connected with Kubernetes Cluster

Cloud Volumes ONTAP is connected to 1 Kubernetes cluster. [View Cluster](#) ⓘ

You can connect another Kubernetes cluster to this Cloud Volumes ONTAP system. If the Kubernetes cluster is in a different network than Cloud Volumes ONTAP, specify a custom export policy to provide access to clients.

#### Kubernetes Cluster

Select Kubernetes Cluster

kubernetes

#### Custom Export Policy *(Optional)* ⓘ

Custom Export Policy

172.31.0.0/16

☒ Set as default storage class

☒ NFS ☐ iSCSI

Connect

Cancel

## Configuración de Kubernetes de puesta en marcha en las instalaciones

Para ver la instalación local de Trident de NetApp, consulte ["TR-4798"](#) para obtener más detalles. Tras configurar su clúster de Kubernetes e instalar Trident de NetApp, puede conectar Trident al clúster Iguazio para habilitar las capacidades de gestión de datos de NetApp, como realizar copias Snapshot de sus datos y su modelo.

### Defina una solicitud de volumen persistente

1. Guarde el siguiente YAML en un archivo para crear un PVC del tipo Basic.

```
kind: PersistentVolumeClaim
apiVersion: v1
metadata:
  name: basic
spec:
  accessModes:
    - ReadWriteOnce
  resources:
    requests:
      storage: 100Gi
  storageClassName: netapp-file
```

## 2. Aplique el archivo YLMA al clúster de Kubernetes de Iguazio.

```
Kubectl -n default-tenant apply -f <your yaml file>
```

### Adjunte el volumen de NetApp al ordenador portátil Jupyter

Iguazio ofrece varios servicios gestionados para proporcionar a los científicos de datos una pila completa para el desarrollo y la puesta en marcha de aplicaciones de IA/ML. Puede obtener más información sobre estos componentes en la ["Visión General de Iguazio de Servicios y Herramientas de aplicación"](#).

Uno de los servicios gestionados es el portátil Jupyter. Cada desarrollador obtiene su propia implementación de un contenedor de portátiles con los recursos que necesitan para el desarrollo. Para darles acceso a Cloud Volume de NetApp, puede asignar el volumen a su contenedor y la asignación de recursos, la configuración de usuario en ejecución y la variable de entorno para las reclamaciones de volumen persistente se presentan en la siguiente imagen.

Para obtener información sobre una configuración en las instalaciones, puede consultar ["TR-4798"](#) En la configuración de Trident para habilitar las funcionalidades de gestión de datos de ONTAP de NetApp, como realizar copias Snapshot de sus datos o modelos para el control de versiones. Añada la siguiente línea en su archivo de configuración de back-end de Trident para que sean visibles los directorios de Snapshot:

```
{
  ...
  "defaults": {
    "snapshotDir": "true"
  }
}
```

Debe crear un archivo de configuración back-end de Trident en formato JSON y, a continuación, ejecutar lo siguiente ["Comando Trident"](#) como referencia:

```
tridentctl create backend -f <backend-file>
```

The screenshot shows the JupyterLab configuration interface. At the top, there is a toggle for 'Enabled' which is checked. Below it is an 'Inactivity window' slider with options 5m, 10m, 1h, 2h, and 4h. The 'Resources' section includes a link to 'Kubernetes documentation' and a note: 'The memory and CPU configurations are applied to each replica.' There are two rows of input fields for 'Request' and 'Limit' for 'Memory' and 'CPU'. The 'Running User' section has a text input field with 'admin' and a 'Username' dropdown menu.

The screenshot shows the JupyterLab configuration interface. At the top, there is a 'Flavor' dropdown menu with 'Full stack without GPU' selected. Below it is a 'Spark' dropdown menu with 'spark' selected and a 'Create new...' link. The 'Environment Variables' section has a 'Create a new environment variable' button. The 'Persistent Volume Claims (PVCs)' section has a table with columns 'Name' and 'Mount Path'. The table has one row with 'basic' in the 'Name' column and '/netapp' in the 'Mount Path' column. There is an 'Add PVC' button at the bottom.

## Implementación de la aplicación

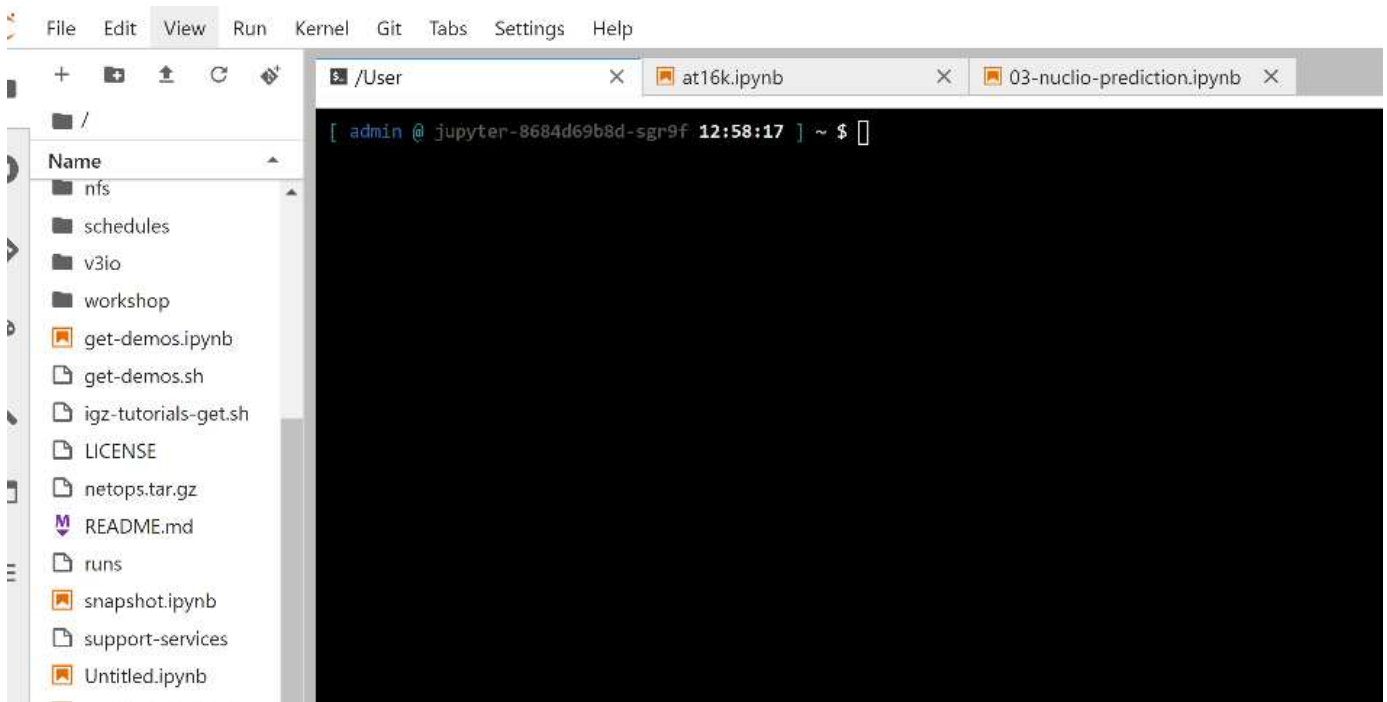
En las siguientes secciones se describe cómo instalar e implementar la aplicación.

### Obtenga el código de GitHub

Ahora que el volumen de cloud de NetApp o el volumen de Trident de NetApp están disponibles para el clúster de Iguazio y el entorno de desarrolladores, puede comenzar a revisar la aplicación.

Los usuarios tienen su propio espacio de trabajo (directorio). En cada notebook, la ruta al directorio de usuario es `/User`. La plataforma Iguazio administra el directorio. Si sigue las instrucciones anteriores, el volumen de cloud de NetApp está disponible en `/netapp` directorio.

Obtenga el código de GitHub con una terminal Jupyter.



En el prompt de la terminal de Jupyter, clone el proyecto.

```
cd /User
git clone .
```

Ahora debería ver la `netops- netapp` Carpeta en el árbol de archivos del espacio de trabajo Jupyter.

### Configurar entorno de trabajo

Copie el Notebook `set_env-Example.ipynb` como `set_env.ipynb`. Abra y edite `set_env.ipynb`. Este cuaderno establece variables para las credenciales, ubicaciones de archivos y controladores de ejecución.

Si sigue las instrucciones anteriores, los pasos siguientes son los únicos cambios que se deben realizar:

1. Obtenga este valor desde el panel de servicios de Iguazio: `docker_registry`

Ejemplo: `docker-registry.default-tenant.app.clusterq.iguazidev.com:80`

2. Cambiar admin A su nombre de usuario de Iguazio:

```
IGZ_CONTAINER_PATH = '/users/admin'
```

A continuación, se muestran los detalles de conexión del sistema ONTAP. Incluya el nombre del volumen que se generó al instalar Trident. La siguiente configuración es para un clúster ONTAP en las instalaciones:

```
ontapClusterMgmtHostname = '0.0.0.0'
ontapClusterAdminUsername = 'USER'
ontapClusterAdminPassword = 'PASSWORD'
sourceVolumeName = 'SOURCE VOLUME'
```

La siguiente configuración es para Cloud Volumes ONTAP:

```
MANAGER=ontapClusterMgmtHostname
svm='svm'
email='email'
password=ontapClusterAdminPassword
weid="weid"
volume=sourceVolumeName
```

## Cree imágenes Docker básicas

Todo lo que necesita para construir un ducto ML está incluido en la plataforma Iguazio. El desarrollador puede definir las especificaciones de las imágenes Docker necesarias para ejecutar la canalización y ejecutar la creación de imágenes desde el portátil Jupyter. Abra el portátil `create-images.ipynb` Y ejecute todas las celdas.

Este cuaderno crea dos imágenes que utilizamos en la canalización.

- `iguazio/netapp`. Se utiliza para manejar tareas ML.

### Create image for training pipeline

```
[4]: fn.build_config(image=docker_registry+'iguazio/netapp', commands=['pip install \
v3io_frames fsspec>=0.3.3 PyYAML==5.1.2 pyarrow==0.15.1 pandas==0.25.3 matplotlib seaborn yellowb
fn.deploy()
```

- `netapp/pipeline`. Contiene utilidades para gestionar las copias snapshot de NetApp.

### Create image for Ontap utilitites

```
[0]: fn.build_config(image=docker_registry + '/netapp/pipeline:latest', commands=['apt -y update', 'pip install v3io_frames netapp_ontap'
fn.deploy()
```

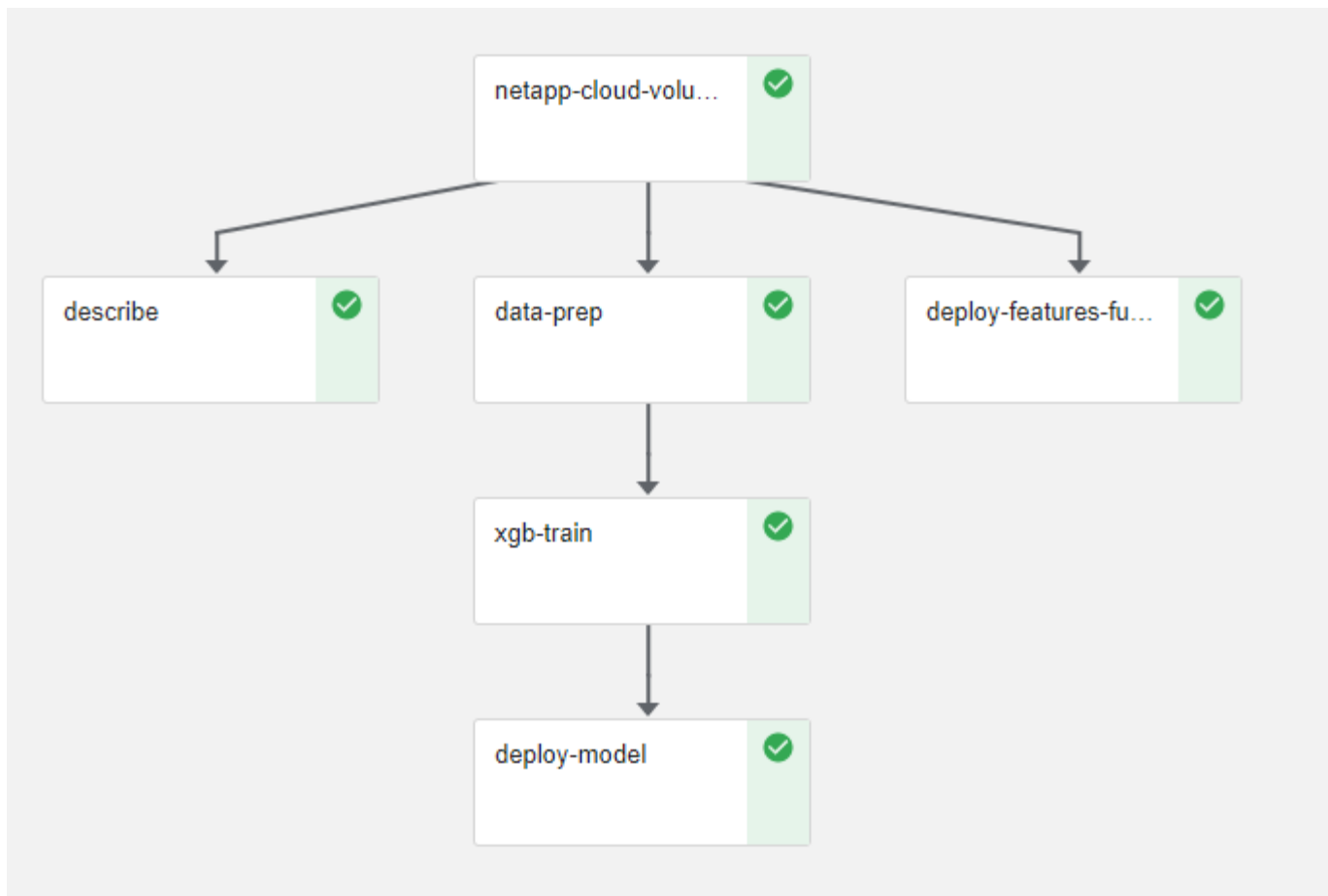
## Revise los cuadernos Jupyter individuales

En la siguiente tabla se enumeran las bibliotecas y los marcos que utilizamos para crear esta tarea. Todos estos componentes se han integrado plenamente con los controles de seguridad y acceso basados en funciones de Iguazio.

Bibliotecas/Marco	Descripción
MLRun	Un gestionado por Iguazio para permitir el montaje, ejecución y supervisión de una canalización ML/IA.
Nuclio	Un marco de funciones sin servidor integrado con Iguazio. También disponible como proyecto de código abierto gestionado por Iguazio.
Kubeflow	Un marco basado en Kubernetes para poner en marcha la canalización. Este es también un proyecto de código abierto al que colabora Iguazio. Se integra con Iguazio para una mayor seguridad e integración con el resto de la infraestructura.
Docker	Un registro Docker se ejecuta como servicio en la plataforma Iguazio. También puede cambiar esta opción para conectarse al registro.
Cloud Volumes de NetApp	Cloud Volumes que se ejecuta en AWS nos proporciona acceso a grandes cantidades de datos y la capacidad de realizar copias de Snapshot para versiones de los conjuntos de datos utilizados para el entrenamiento.
Trident	Trident es un proyecto de código abierto gestionado por NetApp. Facilita la integración con los recursos de almacenamiento y computación en Kubernetes.

Hemos utilizado varios portátiles para construir la canalización DE ML. Cada portátil puede probarse individualmente antes de ser reunido en la tubería. Cubrimos cada portátil individualmente tras el flujo de despliegue de esta aplicación de demostración.

El resultado deseado es una canalización que entrena un modelo basado en una copia Snapshot de los datos y pone en marcha el modelo para la inferencia. En la siguiente imagen se muestra un diagrama de bloque de una canalización MLRun completada.



### Ponga en marcha la función de generación de datos

En esta sección se describe cómo utilizamos las funciones de Nuclio serverless para generar datos de dispositivos de red. El caso de uso se adapta a partir de un cliente Iguazio que ha implementado la canalización y utiliza servicios Iguazio para supervisar y predecir los fallos del dispositivo de red.

Simulamos los datos procedentes de dispositivos de red. Ejecución del cuaderno Jupyter `data-generator.ipynb` Crea una función sin servidor que se ejecuta cada 10 minutos y genera un archivo Parquet con nuevos datos. Para implementar la función, ejecute todas las celdas de este portátil. Consulte ["Sitio Web de Nuclio"](#) para revisar cualquier componente desconocido en este cuaderno.

Al generar la función, se ignora una celda con el siguiente comentario. Se asume que todas las celdas del portátil forman parte de la función. Importe el módulo Nuclio para activarlo `%nuclio magic`.

```
# nuclio: ignore
import nuclio
```

En la especificación para la función, definimos el entorno en el que se ejecuta la función, cómo se activa y los recursos que consume.

```
spec = nuclio.ConfigSpec(config={"spec.triggers.inference.kind":"cron",
                                "spec.triggers.inference.attributes.interval" : "10m",
                                "spec.readinessTimeoutSeconds" : 60,
                                "spec.minReplicas" : 1},.....
```

La `init_context` La función es invocada por el marco Nuclio tras la inicialización de la función.

```
def init_context(context):
    ...
```

Cualquier código que no esté en una función se invoca cuando se inicializa la función. Cuando lo invoca, se ejecuta una función de controlador. Puede cambiar el nombre del controlador y especificarlo en la especificación de función.

```
def handler(context, event):
    ...
```

Puede probar la función desde el portátil antes de la implementación.

```
%%time
# nuclio: ignore
init_context(context)
event = nuclio.Event(body='')
output = handler(context, event)
output
```

La función puede desplegarse desde el portátil o puede desplegarse a partir de una canalización CI/CD (adaptando este código).

```
addr = nuclio.deploy_file(name='generator',project='netops',spec=spec,
tag='v1.1')
```

## Cuadernos de pipeline

Estos cuadernos no están diseñados para ejecutarse individualmente para esta configuración. Esto es sólo una revisión de cada bloc de notas. Los invocamos como parte del proyecto. Para ejecutarlas individualmente, revise la documentación MLRun para ejecutarlas como trabajos de Kubernetes.

### snap\_cv.ipynb

Este portátil gestiona las copias snapshot de Cloud Volume al principio de la canalización. Pasa el nombre del volumen al contexto de la canalización. Este cuaderno invoca un script de shell para manejar la copia



Snapshot. Mientras se ejecuta en la canalización, el contexto de ejecución contiene variables que ayudan a localizar todos los archivos necesarios para su ejecución. Mientras escribe este código, el desarrollador no tiene que preocuparse por la ubicación del archivo en el contenedor que lo ejecuta. Como se describe más tarde, esta aplicación se implementa con todas sus dependencias y es la definición de los parámetros de canalización que proporciona el contexto de ejecución.

```
command = os.path.join(context.get_param('APP_DIR'), "snap_cv.sh")
```

La ubicación de la copia Snapshot creada se coloca en el contexto de MLRun que consumirán los pasos de la canalización.

```
context.log_result('snapVolumeDetails', snap_path)
```

Los siguientes tres portátiles se ejecutan en paralelo.

### **data-prep.ipynb**

Las métricas sin formato deben convertirse en funciones para permitir el entrenamiento de modelos. Este cuaderno lee las métricas sin formato del directorio Snapshot y escribe las funciones de entrenamiento de modelos en el volumen de NetApp.

Cuando se ejecuta en el contexto de la canalización, la entrada `DATA_DIR` Contiene la ubicación de la copia Snapshot.

```
metrics_table = os.path.join(str(mlruncontext.get_input('DATA_DIR',
os.getenv('DATA_DIR', '/netpp'))),
                             mlruncontext.get_param('metrics_table',
os.getenv('metrics_table', 'netops_metrics_parquet')))
```

### **describa.ipynb**

Para visualizar las métricas entrantes, implementamos un paso de canalización que proporciona gráficos y gráficos disponibles a través de las interfaces de usuario de Kubeflow y MLRun. Cada ejecución tiene su propia versión de esta herramienta de visualización.

```
ax.set_title("features correlation")
plt.savefig(os.path.join(base_path, "plots/corr.png"))
context.log_artifact(PlotArtifact("correlation", body=plt.gcf()),
local_path="plots/corr.html")
```

### **deploy-feature-function.ipynb**

Supervisamos continuamente las métricas en busca de anomalías. Este bloc de notas crea una función sin servidor que genera las funciones que necesitan ejecutar la predicción en las métricas entrantes. Este cuaderno invoca la creación de la función. El código de función se encuentra en el portátil `data-prep.ipynb`. Observe que utilizamos el mismo bloc de notas como un paso en la tubería para este propósito.

## training.ipynb

Una vez que creamos las funciones, activamos la formación del modelo. El resultado de este paso es el modelo que se va a utilizar para la inferencia. También recopilamos estadísticas para realizar un seguimiento de cada ejecución (experimento).

Por ejemplo, el siguiente comando introduce la puntuación de precisión en el contexto de ese experimento. Este valor es visible en Kubeflow y MLRun.

```
context.log_result('accuracy', score)
```

## despliegue-inferencia-function.ipynb

El último paso de la canalización es poner en marcha el modelo como una función sin servidor para la inferencia continua. Este cuaderno invoca la creación de la función sin servidor definida en `nuclio-inference-function.ipynb`.

### Revisar y crear el canalización

La combinación de ejecutar todos los portátiles en una tubería permite que la ejecución continua de experimentos reevalúe la precisión del modelo con las nuevas métricas. En primer lugar, abra la `pipeline.ipynb` portátil. Le repasamos por los detalles que muestran cómo NetApp y Iguazio simplifican la puesta en marcha de esta canalización DE ML.

Utilizamos MLRun para proporcionar contexto y manejar la asignación de recursos a cada paso de la canalización. El servicio de API MLRun se ejecuta en la plataforma Iguazio y es el punto de interacción con los recursos de Kubernetes. Cada desarrollador no puede solicitar recursos directamente; la API gestiona las solicitudes y habilita los controles de acceso.

```
# MLRun API connection definition
mlconf.dbpath = 'http://mlrun-api:8080'
```

La canalización puede funcionar con Cloud Volumes de NetApp y los volúmenes en las instalaciones. Hemos creado esta demostración para usar Cloud Volumes, pero puede ver en el código la opción para ejecutarse en las instalaciones.

```
# Initialize the NetApp snap function once for all functions in a notebook
if [ NETAPP_CLOUD_VOLUME ]:
    snapfn =
code_to_function('snap',project='NetApp',kind='job',filename="snap_cv.ipyn
b").apply(mount_v3io())
    snap_params = {
        "metrics_table" : metrics_table,
        "NETAPP_MOUNT_PATH" : NETAPP_MOUNT_PATH,
        'MANAGER' : MANAGER,
        'svm' : svm,
        'email': email,
        'password': password ,
        'weid': weid,
        'volume': volume,
        "APP_DIR" : APP_DIR
    }
else:
    snapfn =
code_to_function('snap',project='NetApp',kind='job',filename="snapshot.ipyn
b").apply(mount_v3io())
...
snapfn.spec.image = docker_registry + '/netapp/pipeline:latest'
snapfn.spec.volume_mounts =
[snapfn.spec.volume_mounts[0],netapp_volume_mounts]
    snapfn.spec.volumes = [ snapfn.spec.volumes[0],netapp_volumes]
```

La primera acción necesaria para convertir un cuaderno Jupyter en un paso de Kubeflow es convertir el código en una función. Una función tiene todas las especificaciones necesarias para ejecutar ese portátil. A medida que se desplaza hacia abajo por el bloc de notas, puede ver que definimos una función para cada paso de la canalización.

Parte del portátil	Descripción
<code_to_function> (parte del módulo MLRun)	Nombre de la función: Nombre del proyecto. se utiliza para organizar todos los artefactos del proyecto. Esto es visible en la interfaz de usuario de MLRun. Amable. En este caso, un trabajo de Kubernetes. Esto podría ser DASK, mpi, sparkk8s, y más. Consulte la documentación de MLRun para obtener más detalles. Archivo. El nombre del portátil. También puede ser una ubicación en Git (HTTP).
imagen	El nombre de la imagen de Docker que estamos utilizando para este paso. Lo hemos creado anteriormente con el bloc de notas create-image.ipynb.
montajes_volúmenes y volúmenes	Detalles para montar el Cloud Volume de NetApp en tiempo de ejecución.

También definimos parámetros para los pasos.

```
params={  "FEATURES_TABLE":FEATURES_TABLE,
          "SAVE_TO" : SAVE_TO,
          "metrics_table" : metrics_table,
          'FROM_TSDB': 0,
          'PREDICTIONS_TABLE': PREDICTIONS_TABLE,
          'TRAIN_ON_LAST': '1d',
          'TRAIN_SIZE':0.7,
          'NUMBER_OF_SHARDS' : 4,
          'MODEL_FILENAME' : 'netops.v3.model.pickle',
          'APP_DIR' : APP_DIR,
          'FUNCTION_NAME' : 'netops-inference',
          'PROJECT_NAME' : 'netops',
          'NETAPP_SIM' : NETAPP_SIM,
          'NETAPP_MOUNT_PATH': NETAPP_MOUNT_PATH,
          'NETAPP_PVC_CLAIM' : NETAPP_PVC_CLAIM,
          'IGZ_CONTAINER_PATH' : IGZ_CONTAINER_PATH,
          'IGZ_MOUNT_PATH' : IGZ_MOUNT_PATH
        }
```

Después de tener la definición de función para todos los pasos, puede construir la canalización. Utilizamos la `kfp` módulo para realizar esta definición. La diferencia entre el uso de `MLRun` y la construcción por su cuenta es la simplificación y el acortamiento de la codificación.

Las funciones definidas se convierten en componentes de pasos mediante el `as_step` Función de `MLRun`.

### Definición de paso de instantánea

Inicie una función `Snapshot`, la salida y el montaje `v3io` como fuente:

```
snap = snapfn.as_step(NewTask(handler='handler',params=snap_params),
name='NetApp_Cloud_Volume_Snapshot',outputs=['snapVolumeDetails','training_
_parquet_file']).apply(mount_v3io())
```

Parámetros	Detalles
Nueva tarea	Newtask es la definición de la ejecución de la función.
(Módulo <code>MLRun</code> )	Manipulador. Nombre de la función Python que se va a invocar. Utilizamos el controlador de nombres en el portátil, pero no es necesario. parámetros. Los parámetros que pasamos a la ejecución. Dentro de nuestro código, utilizamos <code>context.get_param</code> («PARAMETER») para obtener los valores.

Parámetros	Detalles
paso_as	Nombre. Nombre del paso de la canalización de Kubeflow. salidas. Estos son los valores que el paso agrega al diccionario al terminar. Eche un vistazo al portátil SNAP_cv.ipynb. mount_v3io(). Esto configura el paso para montar /User para el usuario que ejecuta la canalización.

```

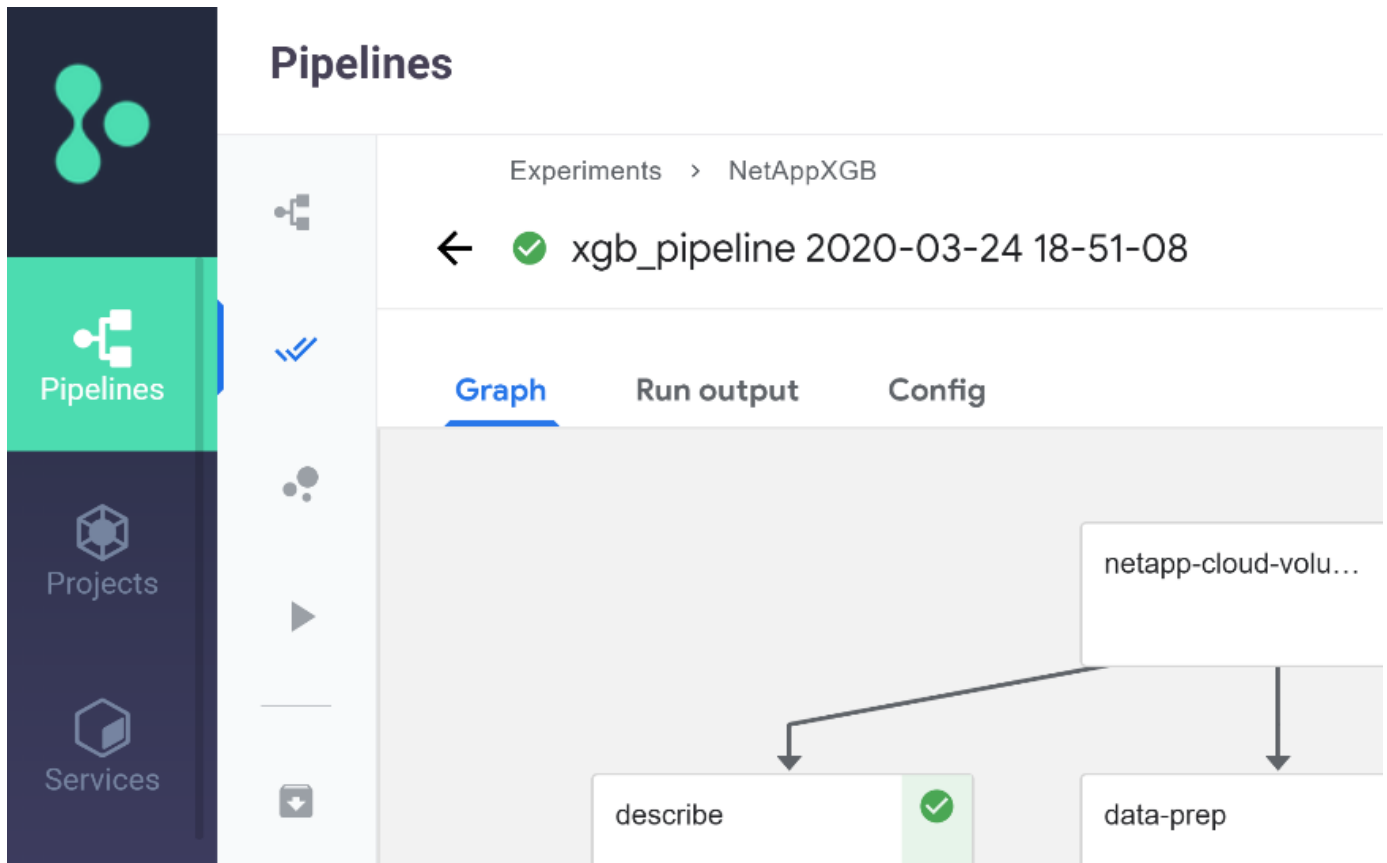
prep = data_prep.as_step(name='data-prep',
handler='handler',params=params,
                        inputs = {'DATA_DIR':
snap.outputs['snapVolumeDetails']} ,

out_path=artifacts_path).apply(mount_v3io()).after(snap)

```

Parámetros	Detalles
entradas	Puede pasar a un paso las salidas de un paso anterior. En este caso, snap.outputs[snapVolumeDetails] es el nombre de la copia Snapshot que creamos en el paso snap.
ruta de salida	Ubicación para colocar artefactos que generan utilizando el módulo MLRun log_Artifacts.

Puede ejecutar `pipeline.ipynb` de arriba a abajo. A continuación, puede ir a la pestaña tuberías desde el panel de control de Iguazio para supervisar el progreso tal y como se ve en la pestaña tuberías del panel de control de Iguazio.



Debido a que hemos registrado la precisión del paso de entrenamiento en cada carrera, tenemos un registro de precisión para cada experimento, como se ve en el registro de precisión de entrenamiento.

<input type="checkbox"/>	Run name	Status	Duration	Pipeline Version	Recurring ...	Start time	accuracy
<input type="checkbox"/>	xgb_pipeline 2020-03-24 18-51-...	✓	0:08:43	[View pipeline]	-	3/24/2020, 2:51:09 PM	0.985
<input type="checkbox"/>	xgb_pipeline 2020-03-19 13-31-...	✓	0:08:14	[View pipeline]	-	3/19/2020, 9:31:19 AM	0.980
<input type="checkbox"/>	xgb_pipeline 2020-03-18 12-56-...	✓	0:08:11	[View pipeline]	-	3/18/2020, 8:56:08 AM	0.990
<input type="checkbox"/>	xgb_pipeline 2020-03-17 19-49-...	✓	0:08:03	[View pipeline]	-	3/17/2020, 3:49:31 PM	0.985
<input type="checkbox"/>	xgb_pipeline 2020-03-17 18-34-...	✓	0:05:54	[View pipeline]	-	3/17/2020, 2:34:56 PM	0.980
<input type="checkbox"/>	xgb_pipeline 2020-03-17 17-34-...	✓	0:04:48	[View pipeline]	-	3/17/2020, 1:34:16 PM	0.982
<input type="checkbox"/>	xgb_pipeline 2020-03-17 17-01-...	✓	0:05:25	[View pipeline]	-	3/17/2020, 1:01:58 PM	0.987
<input type="checkbox"/>	xgb_pipeline 2020-03-16 16-47-...	✓	0:06:08	[View pipeline]	-	3/16/2020, 12:47:19 ...	0.983
<input type="checkbox"/>	xgb_pipeline 2020-03-16 13-57-...	✓	0:05:18	[View pipeline]	-	3/16/2020, 9:57:03 AM	0.980

Si selecciona el paso Snapshot, puede ver el nombre de la copia Snapshot que se utilizó para ejecutar este experimento.

netops-trainign-pipeline-with-netapp-volume-cloning-rtxdl-2910983943

Artifacts **Input/Output** Volumes Manifest Logs

Input artifacts

Output parameters

netapp-cloud-volume-snapshot-snapVolumeDetails	/netapp/.snapshot/kfp_20200324_185122
netapp-cloud-volume-snapshot-training_parquet_file	/netapp/.snapshot/kfp_20200324_18512...

Output artifacts

El paso descrito tiene artefactos visuales para explorar las métricas que utilizamos. Puede expandir para ver el trazado completo como se ve en la siguiente imagen.

netops-trainign-pipeline-with-netapp-volume-cloning-rtxdl-2

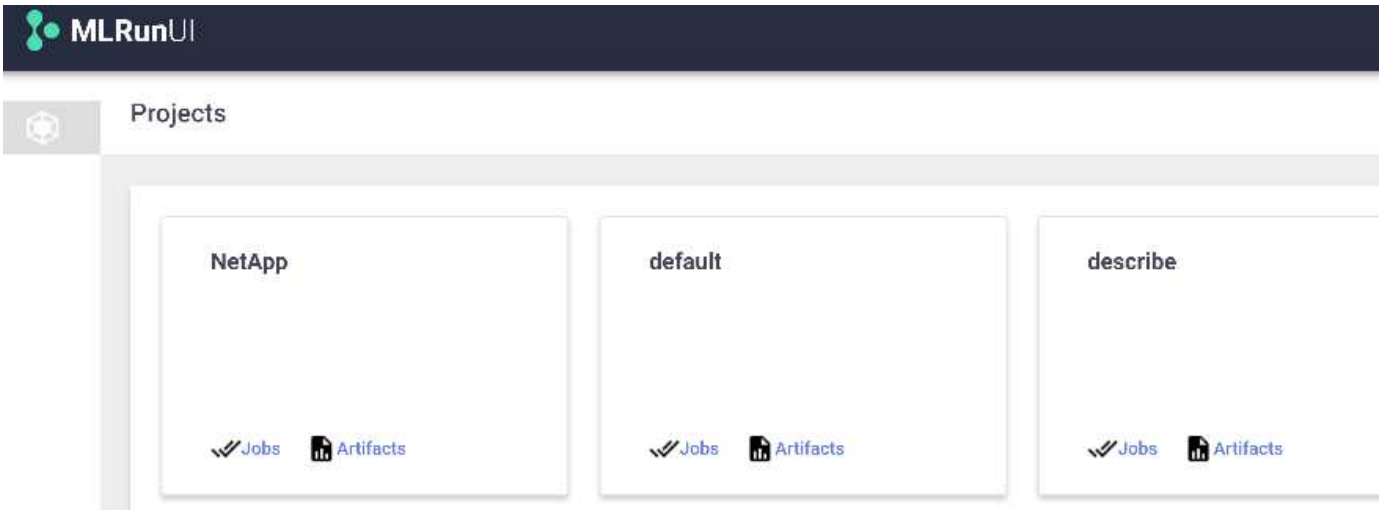
Artifacts **Input/Output** Volumes Manifest Logs

Static HTML

Class Balance for 48,008

40000

La base de datos de la API MLRun también realiza un seguimiento de las entradas, salidas y artefactos de cada ejecución organizada por el proyecto. En la siguiente imagen se puede ver un ejemplo de entradas, salidas y artefactos para cada secuencia.



Para cada trabajo, almacenamos detalles adicionales.

Name	
<a href="#">deploy-model</a> ● 24 Mar, 14:56:03 ...bcbe38e	
<a href="#">xgb_train</a> ● 24 Mar, 14:53:18 ...5c85949	
<a href="#">data-prep</a> ● 24 Mar, 14:52:46 ...126dc73	
<a href="#">describe</a> ● 24 Mar, 14:52:45 ...c2a460e	
<a href="#">deploy-features-function</a> ● 24 Mar, 14:52:43 ...50d8b83	
<a href="#">NetApp_Cloud_Volume_Sna</a> 24 Mar, 14:51:22 ...3108eb2	

### describe

24 Mar, 14:52:45 ●

InfoInputsArtifactsResultsLogs

UID

66ef22187efb4ad89e8da8433c2a460e

Start time

24 Mar, 14:52:45

Parameters

Completed ●

Results

class\_label... ▾

key: summary

label\_colu... ▾

Hay más información sobre MLRun que podemos cubrir en este documento. Los artefactos de al, incluida la definición de los pasos y las funciones, se pueden guardar en la base de datos de API, con versiones e invocados individualmente o como un proyecto completo. Los proyectos también se pueden guardar e insertar en Git para su uso posterior. Le animamos a obtener más información en la ["Sitio de MLRun GitHub"](#).

Implemente el panel de Grafana

Una vez que todo se pone en marcha, ejecutamos inferencias sobre nuevos datos. Los modelos predicen fallos en el equipo de dispositivo de red. Los resultados de la predicción se almacenan en una tabla de timbres de Iguazio. Puede visualizar los resultados con Grafana en la plataforma integrada con la política de acceso a datos y seguridad de Iguazio.








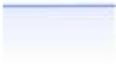











Puede implementar la consola importando el archivo JSON proporcionado en las interfaces de Grafana en el



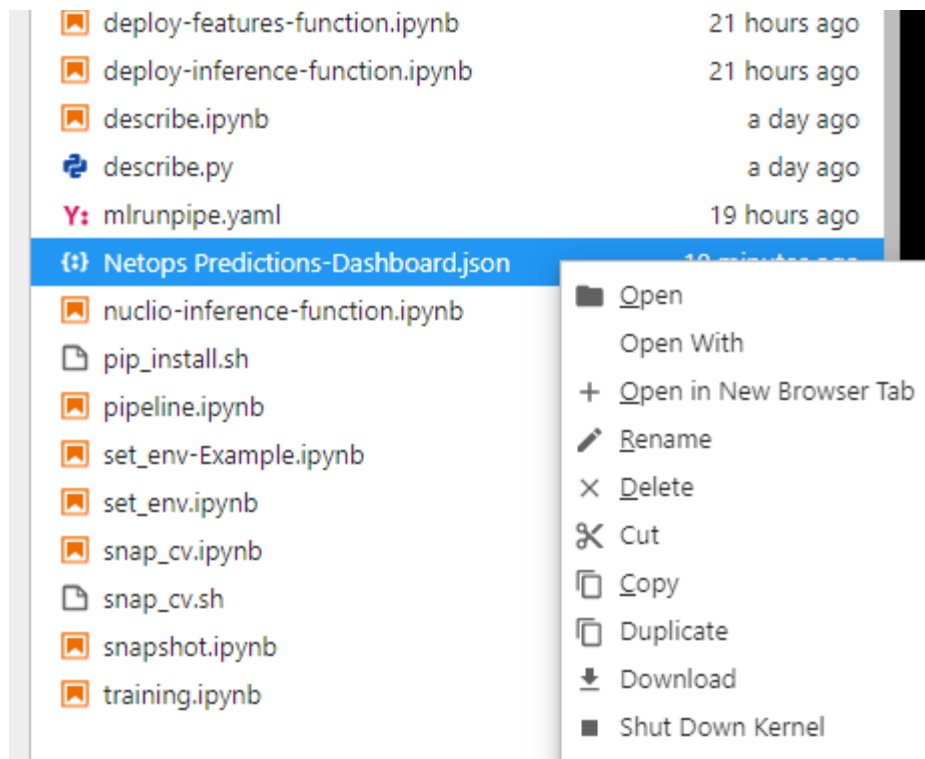
clúster.

1. Para verificar que el servicio Grafana se está ejecutando, busque Servicios.

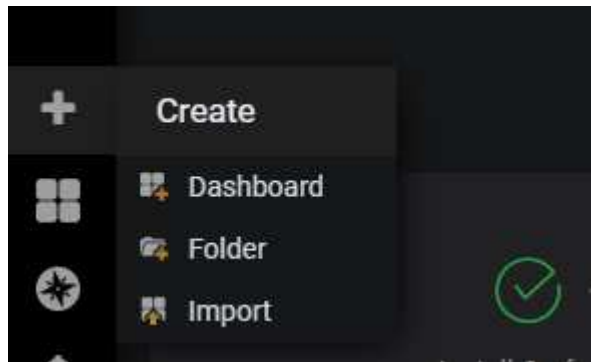
## Services

<input type="checkbox"/>	Name ↑	Running User	Version ✓	CPU (cores)	Memory	AF
<input type="checkbox"/>	 <b>docker-registry</b> Type: Docker Regi		2.7.1	96μ 	1.67 GB 	H
<input type="checkbox"/>	 <b>framesd</b> Type: V3IO Frame		0.6.10	369μ 	795.19 MB 	H
<input type="checkbox"/>	 <b>grafana</b> Type: Grafana		6.6.0	1m 	38.39 MB 	
<input type="checkbox"/>	 <b>jupyter</b> Type: Jupyter Note	admin	1.0.2	81m 	3.27 GB 	
<input type="checkbox"/>	 <b>log-forwarder</b> Type: Log forward		6.7.2	0 	0 bytes 	

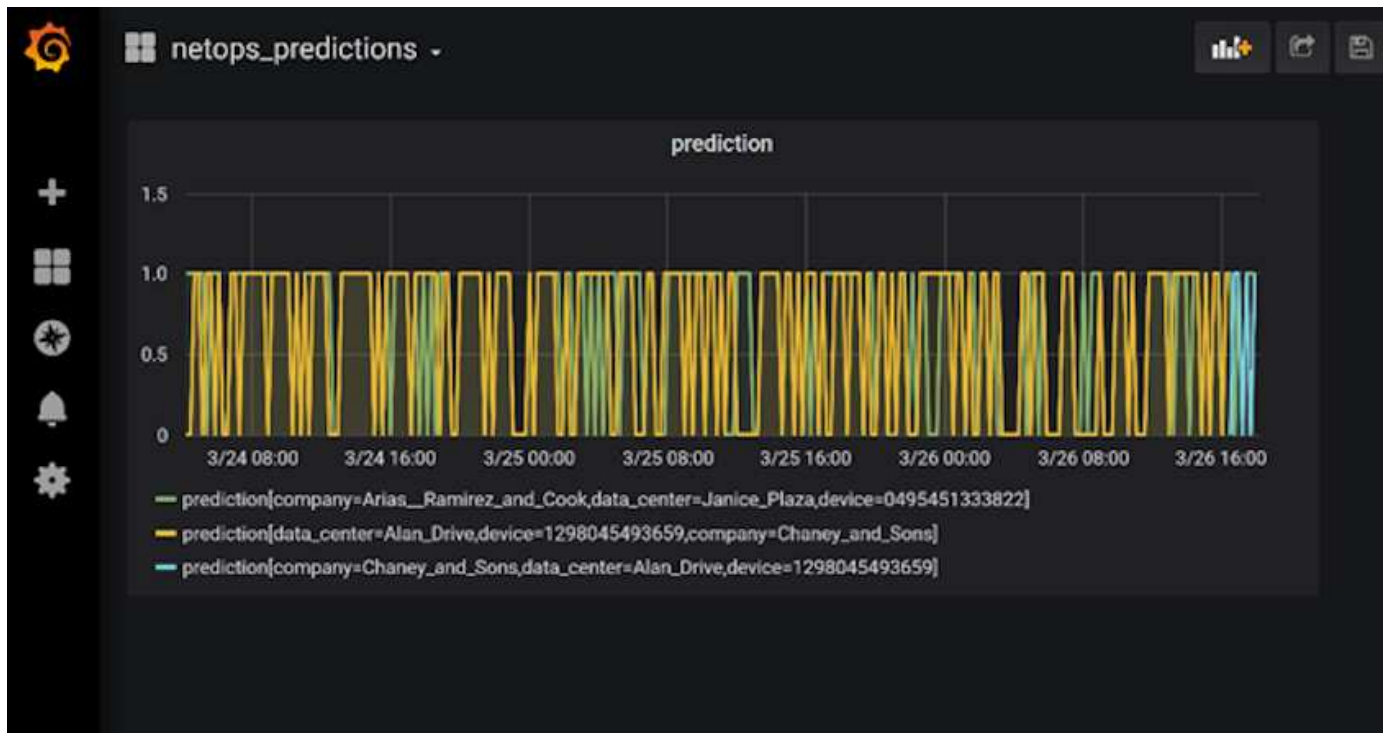
2. Si no está presente, implemente una instancia en la sección Servicios:
  - a. Haga clic en Nuevo servicio.
  - b. Seleccione Grafana de la lista.
  - c. Acepte los valores predeterminados.
  - d. Haga clic en Siguiente paso.
  - e. Introduzca su ID de usuario.
  - f. Haga clic en Guardar servicio.
  - g. Haga clic en aplicar cambios en la parte superior.
3. Para implementar el panel de control, descargue el archivo `NetopsPredictions-Dashboard.json` A través de la interfaz Jupyter.



4. Abra Grafana en la sección Servicios e importe el panel de control.



5. Haga clic en Upload \*.json File (Archivo) y seleccione el archivo que descargó anteriormente (NetopsPredictions-Dashboard.json). El panel se muestra una vez finalizada la carga.



## Despliegue la función de limpieza

Cuando genera una gran cantidad de datos, es importante mantener las cosas limpias y organizadas. Para ello, implemente la función de limpieza con `cleanup.ipynb` portátil.

## Beneficios

NetApp y Iguazio aceleran y simplifican la puesta en marcha de aplicaciones de IA y ML mediante la creación de marcos esenciales como Kubeflow, Apache Spark y TensorFlow, junto con herramientas de orquestación como Docker y Kubernetes. Al unificar la canalización de datos completa, NetApp y Iguazio reducen la latencia y la complejidad inherentes a muchas cargas de trabajo informáticas avanzadas, y esta brecha entre el desarrollo y las operaciones. Los científicos de datos pueden ejecutar consultas en grandes conjuntos de datos y compartir de forma segura datos y modelos algorítmicos con usuarios autorizados durante la fase de entrenamiento. Después de que los modelos en contenedores están listos para la producción, puede moverlos fácilmente desde entornos de desarrollo a entornos operativos.

## Conclusión

Cuando cree sus propias canalizaciones de IA/ML, configurar la integración, la gestión, la seguridad y la accesibilidad de los componentes en una arquitectura es una tarea ardua. Dar a los desarrolladores acceso y control de su entorno presenta otro conjunto de retos.

La combinación de NetApp e Iguazio aúna estas tecnologías como servicios gestionados para acelerar la adopción de tecnologías y mejorar los plazos de comercialización de las nuevas aplicaciones DE IA/ML.

## TR-4915: Movimiento de datos con E-Series y BeeGFS para flujos de trabajo de análisis e IA

Cody Harryman y Ryan Rodine, NetApp

En TR-4915 se describe cómo mover datos desde cualquier repositorio de datos a un

sistema de archivos BeeGFS respaldado por el almacenamiento SAN E-Series de NetApp. En el caso de aplicaciones de inteligencia artificial (IA) y aprendizaje automático (ML), es posible que los clientes deban mover de manera rutinaria grandes conjuntos de datos que superen muchos petabytes de datos en sus clústeres BeeGFS para el desarrollo de modelos. En este documento se explora cómo llevar a cabo este proceso mediante el XCP de NetApp y las herramientas de copia y sincronización de BlueXP de NetApp.

["TR-4915: Movimiento de datos con E-Series y BeeGFS para flujos de trabajo de análisis e IA"](#)

## Casos de uso

### IA responsable e inferencia confidencial: IA de NetApp con transformación de imagen Protopía

#### TR-4928: IA responsable e inferencia confidencial - NetApp AI con Protopía Image y transformación de datos

Sathish Thyagarajan, Michael Oglesby, NetApp Byung Hoon Ahn, Jennifer Cwaggenberg, Protopia

Las interpretaciones visuales se han convertido en parte integral de la comunicación con la aparición de la captura de imágenes y el procesamiento de imágenes. La inteligencia artificial (IA) en el procesamiento de imágenes digitales brinda nuevas oportunidades de negocio, como en el campo médico para la identificación del cáncer y otras enfermedades, en el análisis visual geoespacial para estudiar peligros ambientales, en el reconocimiento de patrones, en el procesamiento de vídeos para combatir la delincuencia, etc. Sin embargo, esta oportunidad también viene con responsabilidades extraordinarias.

Cuanto más decisiones las organizaciones se pongan en manos de la IA, más aceptan riesgos relacionados con la privacidad y la seguridad de los datos, así como con cuestiones legales, éticas y reguladoras. La IA responsable permite aplicar una práctica que permite a las empresas y organizaciones gubernamentales crear confianza y gobernanza, algo crucial para la IA a escala en grandes empresas. Este documento describe una solución de inferencia de IA validada por NetApp en tres situaciones diferentes utilizando las tecnologías de gestión de datos de NetApp con el software de ofuscación de datos de Protopía con el fin de privatizar datos confidenciales y reducir riesgos y preocupaciones éticas.

Cada día, los consumidores y entidades de negocio generan millones de imágenes con diversos dispositivos digitales. La consiguiente explosión masiva de datos y cargas de trabajo informáticas hace que las empresas recurran a plataformas de cloud computing para obtener escalado y eficiencia. Mientras tanto, la preocupación por la privacidad de la información confidencial que contienen los datos de imágenes surge como consecuencia de la transferencia a un cloud público. La falta de garantías de seguridad y privacidad se convierte en la principal barrera para el despliegue de sistemas de IA que procesan imágenes.

Además, está el ["derecho a borrar"](#) Según el RGPD, el derecho de una persona a solicitar que una organización borre todos sus datos personales. También está la ["Ley de Privacidad"](#), que establece un código de prácticas de información justas. Las imágenes digitales, como las fotografías, pueden constituir datos personales en virtud del RGPD, que rige la forma en que deben recopilarse, procesarse y borrarse los datos. El no hacerlo es un incumplimiento del RGPD, que puede conllevar multas cuantiosas por incumplimiento de las normativas que pueden resultar seriamente perjudiciales para las organizaciones. Los principios de

privacidad son uno de los pilares de la implementación de IA responsable que garantice la justicia en las predicciones de modelos de aprendizaje automático y aprendizaje profundo (DL) y reduce los riesgos asociados con la infracción de privacidad o el cumplimiento de normativas.

Este documento describe una solución de diseño validada en tres situaciones diferentes, con y sin la confusión de imágenes relevante para preservar la privacidad y poner en marcha una solución de IA responsable:

- **Escenario 1.** inferencia a petición dentro del cuaderno Jupyter.
- **Escenario 2.** inferencia por lotes en Kubernetes.
- **Escenario 3.** servidor de inferencia NVIDIA Triton.

Para esta solución, utilizamos el conjunto de datos y el punto de referencia de detección facial (Fddb), un conjunto de datos de regiones faciales diseñado para estudiar el problema de la detección facial sin restricciones, combinado con el marco de aprendizaje de máquinas PyTorch para la implementación de FaceBox. Este conjunto de datos contiene las anotaciones para 5171 caras en un conjunto de 2845 imágenes de varias resoluciones. Además, este informe técnico presenta algunas áreas de soluciones y casos de uso relevantes recopilados por clientes de NetApp e ingenieros de campo en situaciones en las que sea aplicable esta solución.

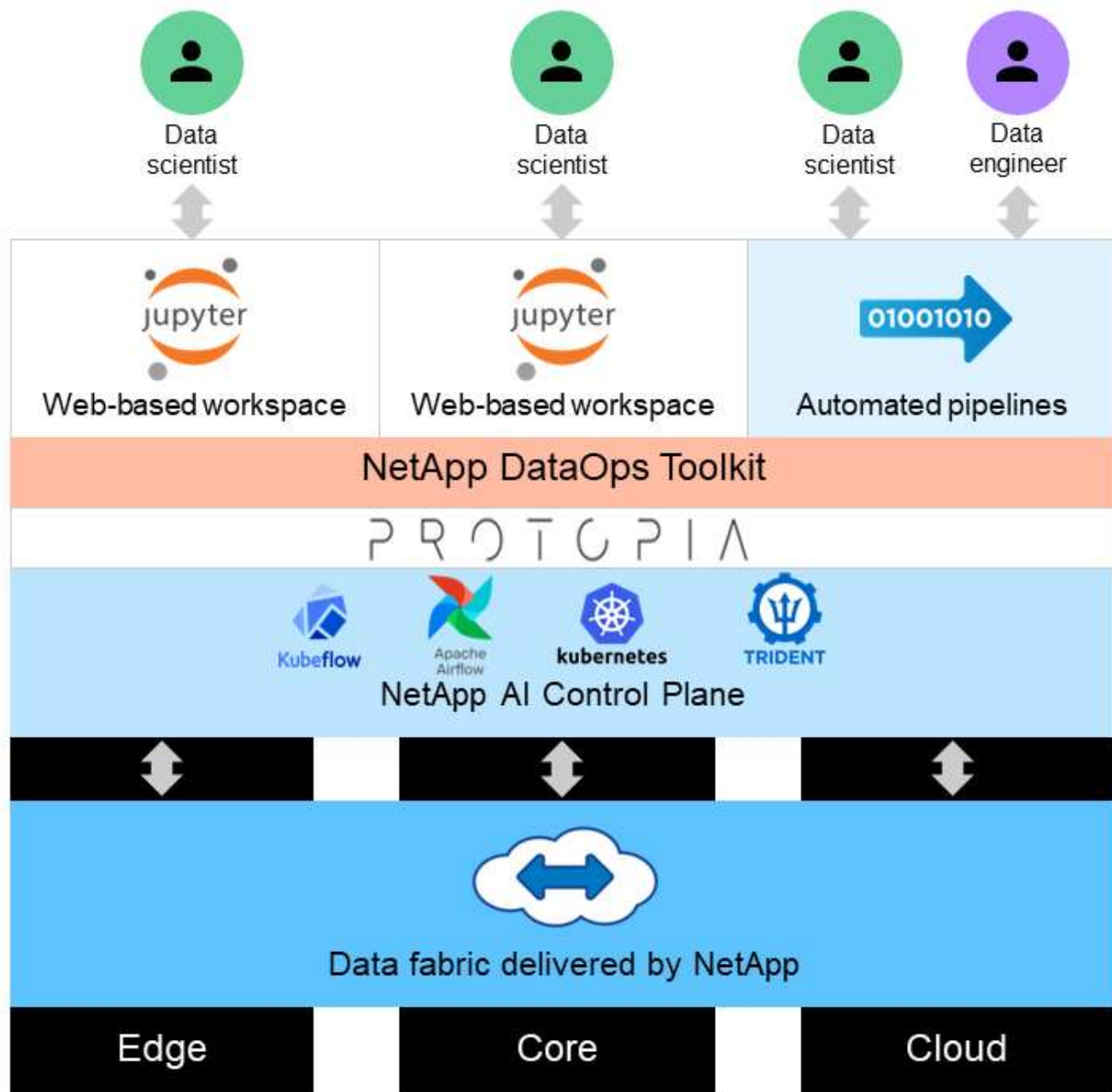
### **Público objetivo**

Este informe técnico está dirigido a los siguientes destinatarios:

- Líderes empresariales y arquitectos empresariales que desean diseñar y poner en marcha una IA responsable y abordar cuestiones de protección y privacidad de datos relacionadas con el procesamiento de imágenes faciales en espacios públicos.
- Científicos de datos, ingenieros de datos, investigadores DE IA/aprendizaje automático (ML) y desarrolladores de sistemas de IA/ML que pretenden proteger y preservar la privacidad.
- Arquitectos empresariales que diseñan soluciones de confusión de datos para modelos y aplicaciones de IA/ML que cumplen con las normativas tales como RGPD, CCPA o las organizaciones gubernamentales (Ley de privacidad del Departamento de Defensa) y la Ley de privacidad del Departamento de Defensa (DoD).
- Científicos e ingenieros de IA buscan formas eficientes de poner en marcha el aprendizaje profundo (DL) y modelos de inferencia de IA/ML/DL que protegen la información confidencial.
- Los administradores de dispositivos periféricos y los administradores de servidor perimetral son responsables de la puesta en marcha y la gestión de modelos de inferencia perimetrales.

### **Arquitectura de la solución**

Esta solución está diseñada para gestionar cargas de trabajo de IA de inferencia en lote y en tiempo real en grandes conjuntos de datos utilizando la potencia de procesamiento de las GPU junto con las CPU tradicionales. Esta validación demuestra la inferencia privacidad de ML y la gestión de datos óptima que necesitan las organizaciones que buscan una puesta en marcha de IA responsable. Esta solución proporciona una arquitectura adecuada para una plataforma Kubernetes de uno o varios nodos para el perímetro y el cloud computing interconectados con ONTAP AI de NetApp en el núcleo de las instalaciones, el kit de herramientas DataOPS de NetApp y el software de ofuscación por medio de las interfaces CLI y de Jupyter Lab. En la siguiente figura se muestra información general sobre la arquitectura lógica del Data Fabric con la tecnología de NetApp con el kit de herramientas de operaciones de datos y Protopia.



El software de ofuscación de Protopía se ejecuta sin problemas sobre el kit de herramientas DataOPS de NetApp y transforma los datos antes de salir del servidor de almacenamiento.

### Áreas de soluciones

El procesamiento de imágenes digitales ofrece muchas ventajas, lo que permite a muchas organizaciones aprovechar al máximo los datos asociados con las representaciones visuales. Esta solución de NetApp y Protopia ofrece un diseño único de inferencia de IA para proteger y privatizar datos de IA/ML en TODO el ciclo de vida DE ML/DL. Permite a los clientes conservar la propiedad de datos confidenciales, utilizar modelos de puesta en marcha de cloud público o híbrido para escalar y eficiencia al eliminar las preocupaciones relacionadas con la privacidad y poner en marcha la inferencia de IA en el perímetro.

## Inteligencia ambiental

Hay muchas maneras en que las industrias pueden aprovechar los análisis geoespaciales en las áreas de peligros ambientales. Los gobiernos y el departamento de obras públicas pueden obtener información práctica sobre la salud pública y las condiciones climáticas para asesorar mejor al público durante una pandemia o un desastre natural como los incendios forestales. Por ejemplo, puede identificar a un paciente con un COVID positivo en espacios públicos, como aeropuertos o hospitales, sin poner en peligro la privacidad del individuo afectado y alertar a las autoridades respectivas y al público cercano acerca de las medidas de seguridad necesarias.

## Dispositivos portátiles

En el ejército y en campos de batalla, se puede utilizar la inferencia de IA en el borde como dispositivos portátiles para rastrear la salud de los soldados, monitorear el comportamiento de los conductores y alertar a las autoridades sobre la seguridad y los riesgos asociados de acercarse a los vehículos militares mientras se preserva y protege la privacidad de los soldados. El futuro de los militares está yendo de alta tecnología con el Internet de Battlefield Things (IoBT) y el Internet de las cosas militares (IoMT) por llevar equipos de combate que ayudan a los soldados a identificar a los enemigos y a actuar mejor en la batalla mediante la computación de avanzada. Proteger y preservar los datos visuales recopilados de dispositivos periféricos como aviones teledirigidos y engranajes portátiles es crucial para mantener a raya a los hackers y al enemigo.

## Operaciones de evacuación no combatiente

Las operaciones de evacuación no combatiente (Neos) las lleva a cabo el Departamento de Defensa para ayudar a evacuar a ciudadanos y nacionales estadounidenses, personal civil del Departamento de Defensa y personas designadas (nación anfitriona (HN) y nacionales de terceros países (TCN) cuyas vidas están en peligro de un refugio seguro apropiado. Los controles administrativos vigentes utilizan en gran medida procesos manuales de detección de evacuados. Sin embargo, la precisión, la seguridad y la velocidad de la identificación del evacuado, el seguimiento del evacuado y la detección de amenazas podrían mejorarse utilizando herramientas de IA/ML altamente automatizadas combinadas con tecnologías de ofuscación por vídeo de IA/ML.

## Salud e investigación biomédica

El procesamiento de imágenes se utiliza para diagnosticar patologías para la planificación quirúrgica a partir de imágenes 3D obtenidas de tomografía computarizada (TC) o resonancia magnética (RM). Las reglas de privacidad de HIPAA rigen cómo deben recopilarse, procesar y borrar los datos para toda la información personal e imágenes digitales como fotografías. Para que los datos se puedan calificar para que se puedan compartir con la normativa HIPAA Safe Harbor, es necesario eliminar las imágenes fotográficas de cara completa y cualquier imagen comparable. Las técnicas automatizadas como la desidentificación o los algoritmos de decapado de cráneo utilizados para ocultar las características faciales de un individuo a partir de las imágenes estructurales de TC/RM se han convertido en una parte esencial del proceso de intercambio de datos para las instituciones de investigación biomédica.

## Migración al cloud de los análisis de IA/ML

Los clientes empresariales han formado y puesto en marcha modelos de IA/ML en sus instalaciones. Por razones de escalado y eficiencia, estos clientes amplían su capacidad para trasladar las funciones de IA/ML a puestas en marcha de cloud público, híbrido o multicloud. Sin embargo, están ligados a qué datos se pueden exponer a otras infraestructuras. Las soluciones de NetApp se plantean una amplia gama de amenazas de ciberseguridad para las que es necesario ["protección de datos"](#) La evaluación de la seguridad y, cuando se combina con la transformación de datos de Protopía, minimizan los riesgos asociados con la migración de cargas de trabajo DE IA/ML de procesamiento de imágenes al cloud.

Para ver casos de uso adicionales para la computación perimetral y la inferencia de IA en otros sectores,



consulte "[TR-4886 inferencia de IA en el perímetro](#)" Y el blog de IA de NetApp, "[Inteligencia frente a privacidad](#)".

## Información general de la tecnología

En esta sección se ofrece una descripción general de los distintos componentes técnicos necesarios para completar esta solución.

### Protopía

Protopia AI ofrece una solución discreta y exclusiva de software para la inferencia confidencial en el mercado actual. La solución Protopia ofrece una protección sin igual para servicios de inferencia al minimizar la exposición de información confidencial. La IA solo se alimenta de la información en el registro de datos que es realmente esencial para realizar la tarea a mano y nada más. La mayoría de las tareas de inferencia no utilizan toda la información que existe en cada registro de datos. Independientemente de si su IA consume imágenes, voz, vídeo o incluso datos tabulares estructurados, Protopia solo ofrece lo que el servicio de inferencia necesita. La tecnología de núcleo patentada utiliza el ruido matemáticamente curado para transformar estocamente los datos y engardar la información que no necesita un servicio DE ML dado. Esta solución no enmascara los datos; más bien, cambia la representación de datos mediante el uso de ruido aleatorio curado.

La solución Protopia formula el problema de cambiar la representación como un método de maximización de turbación basado en gradiente que aún conserva la información pertinente en el espacio de características de entrada con respecto a la funcionalidad del modelo. Este proceso de detección se ejecuta como un pase de ajuste preciso al final del entrenamiento del modelo ML. Después de que el paso genera automáticamente un conjunto de distribuciones de probabilidad, una transformación de datos de baja sobrecarga aplica muestras de ruido de estas distribuciones a los datos, ocultando dichos datos antes de pasarlos al modelo para la inferencia.

### ONTAP AI de NetApp

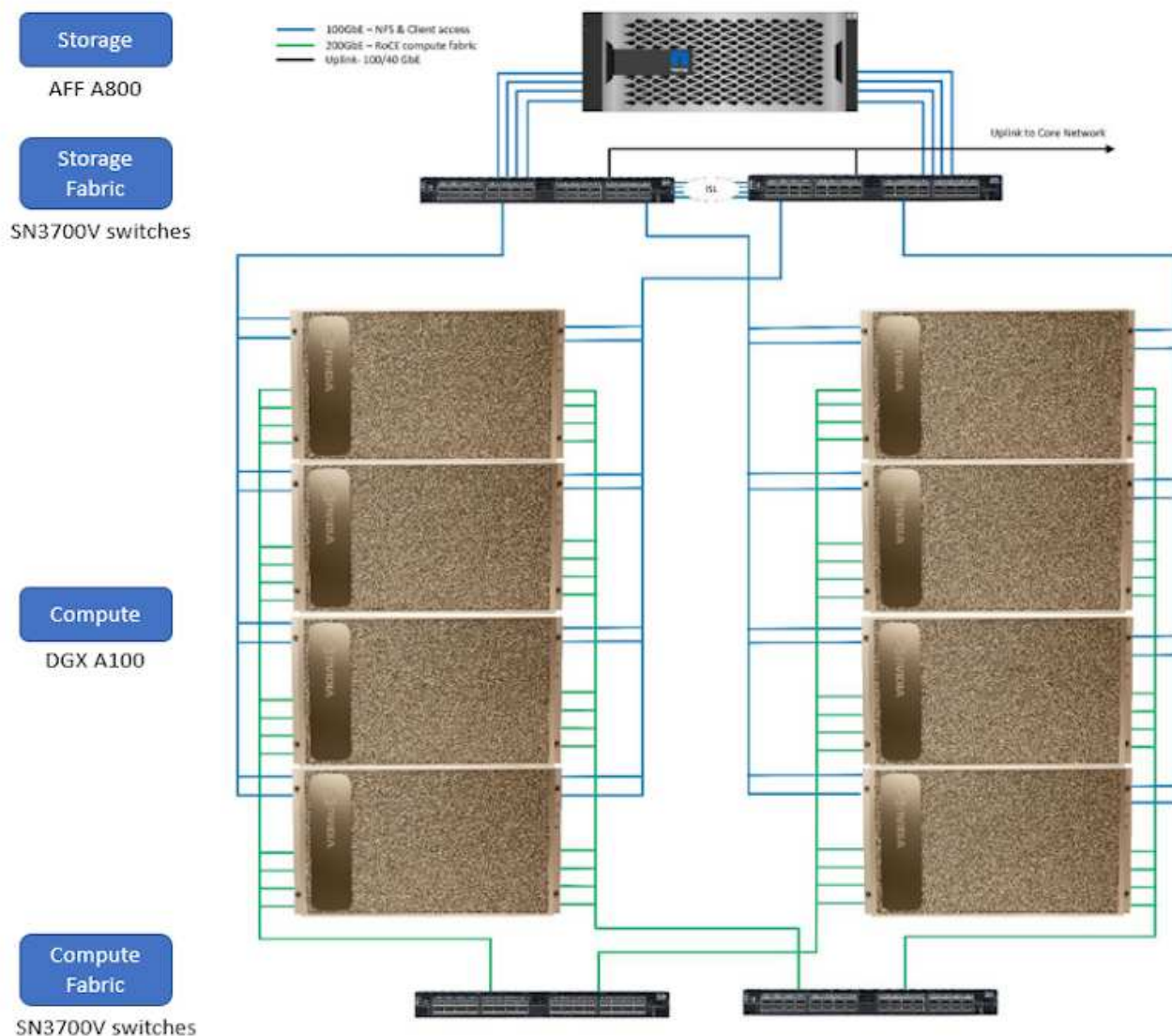
La arquitectura de referencia de ONTAP AI de NetApp, con sistemas DGX A100 y sistemas de almacenamiento conectados al cloud de NetApp, ha sido desarrollada y verificada por NetApp y NVIDIA. Proporciona a las organizaciones DE TI una arquitectura que ofrece las siguientes ventajas:

- Elimina las complejidades de diseño
- Permite un escalado independiente de las capacidades de computación y almacenamiento
- Permite a los clientes empezar con poco y escalar sin problemas
- Ofrece opciones de almacenamiento para distintos niveles de rendimiento y coste

ONTAP AI integra perfectamente los sistemas DGX A100 y los sistemas de almacenamiento AFF A800 de NetApp con una red de vanguardia. ONTAP AI simplifica las puestas en marcha de IA eliminando complejidades y conjeturas de diseño. Los clientes pueden empezar con poco e ir creciendo de forma no disruptiva a la vez que gestionan de forma inteligente los datos desde el perímetro hasta el núcleo, pasando por el cloud.

En la siguiente figura, se muestran varias variaciones de la familia de soluciones de IA de ONTAP con sistemas DGX A100. El rendimiento del sistema AFF A800 se verifica con hasta ocho sistemas DGX A100. Al añadir pares de controladoras de almacenamiento al clúster ONTAP, la arquitectura puede escalarse a varios racks y admitir muchos sistemas DGX A100 y petabytes de capacidad de almacenamiento con rendimiento lineal. Este enfoque permite alterar de forma independiente las tasas de computación a almacenamiento en función del tamaño de los modelos de AP utilizados y de las métricas de rendimiento necesarias.





Para obtener más información sobre ONTAP AI, consulte ["NVA-1153: ONTAP AI de NetApp con sistemas NVIDIA DGX A100 y switches Mellanox Spectrum Ethernet."](#)

## ONTAP de NetApp

ONTAP 9.11, la última generación del software de gestión del almacenamiento de NetApp, permite a las empresas modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. ONTAP ofrece las mejores capacidades de gestión de datos y permite la gestión y protección de los datos con un solo conjunto de herramientas, sin importar dónde residan. También puede mover los datos libremente a donde sea necesario: El perímetro, el núcleo o el cloud. ONTAP 9.11 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen los datos esenciales y permiten disfrutar de funcionalidades de infraestructura de nueva generación en arquitecturas de cloud híbrido.

## Kit de herramientas de operaciones de datos de NetApp

El kit de herramientas DataOps de NetApp es una biblioteca Python que facilita a los desarrolladores, científicos de datos, ingenieros de DevOps e ingenieros de datos la realización de varias tareas de gestión de datos, como el aprovisionamiento casi instantáneo de un nuevo volumen de datos o un espacio de trabajo JupyterLab, el clonado casi instantáneo de un volumen de datos o un espacio de trabajo JupyterLab. Y tomar instantáneas de un volumen de datos o espacio de trabajo JupyterLab para su trazabilidad o línea de

base. Esta biblioteca de Python puede funcionar como una utilidad de línea de comandos o una biblioteca de funciones que puede importar a cualquier programa de Python o a cualquier cuaderno de Jupyter.

### **Servidor de inferencia NVIDIA Triton**

El servidor de inferencia de NVIDIA Triton es un software de servicio de inferencia de código abierto que ayuda a estandarizar la puesta en marcha y ejecución de modelos para ofrecer IA rápida y escalable en producción. Triton Inference Server optimiza la inferencia de IA al permitir a los equipos poner en marcha, ejecutar y escalar modelos de IA entrenados desde cualquier marco en cualquier infraestructura basada en GPU o CPU. El servidor de inferencia de Triton admite los principales marcos de trabajo, como TensorRT, NVIDIA TensorRT, PyTorch, MXNet, OpenVINO, etc. Triton se integra con Kubernetes para la orquestación y el escalado, que puede utilizar en las principales plataformas de inteligencia artificial y Kubernetes. También está integrada con muchas soluciones de software de MLOPS.

### **PyTorch**

"PyTorch" Es un marco DE APRENDIZAJE AUTOMÁTICO de código abierto. Se trata de una biblioteca tensora optimizada para el aprendizaje profundo que utiliza GPU y CPU. El paquete PyTorch contiene estructuras de datos para tensores multidimensionales que proporcionan muchas utilidades para serializar eficazmente los tensores entre otras utilidades útiles. También tiene un par CUDA que permite ejecutar los cálculos tensores en una GPU de NVIDIA con capacidad de computación. En esta validación, utilizamos la biblioteca de OpenCV-Python (cv2) para validar nuestro modelo y aprovechar los conceptos de visión computarizada más intuitivos de Python.

### **Simplificar la gestión de los datos**

La gestión de los datos es crucial para las operaciones TECNOLÓGICAS empresariales y los científicos de datos, para que se utilicen recursos apropiados para las aplicaciones de IA y para entrenar conjuntos de datos de IA/ML. La siguiente información adicional sobre las tecnologías de NetApp no está disponible para esta validación, pero puede ser relevante en función de su puesta en marcha.

El software para la gestión de datos ONTAP incluye las siguientes funciones para mejorar y simplificar las operaciones, y reducir el coste total de funcionamiento:

- Compactación de datos inline y deduplicación expandida. La compactación de datos reduce el espacio perdido dentro de los bloques de almacenamiento, mientras que la deduplicación aumenta la capacidad efectiva de forma significativa. Esto es aplicable a los datos almacenados localmente y a los datos organizados en niveles en el cloud.
- Calidad de servicio (AQoS) mínima, máxima y adaptativa. Los controles granulares de calidad de servicio (QoS) ayudan a mantener los niveles de rendimiento para aplicaciones críticas en entornos altamente compartidos.
- FabricPool de NetApp. Proporciona la organización automática en niveles de datos fríos en opciones de almacenamiento en cloud privado como Amazon Web Services (AWS), Azure y la solución de almacenamiento StorageGRID de NetApp. Para obtener más información sobre FabricPool, consulte ["TR-4598: Prácticas recomendadas de FabricPool"](#).

### **Acelere y proteja sus datos**

ONTAP no solo ofrece niveles de rendimiento y protección de datos superiores, sino que amplía estas capacidades de las siguientes maneras:

- Rendimiento y menor latencia. ONTAP ofrece la salida más alta posible con la menor latencia posible.
- Protección de datos. ONTAP ofrece capacidades integradas de protección de datos, con una administración común entre todas las plataformas.

- Cifrado de volúmenes de NetApp (NVE). ONTAP ofrece cifrado nativo en el nivel de volumen y permite la gestión de claves incorporada o externa.
- Multi-tenancy y autenticación multifactor. ONTAP permite compartir recursos de infraestructura con los niveles más altos de seguridad.

## Infraestructura preparada para futuros retos

ONTAP ayuda a satisfacer las exigentes y siempre cambiantes necesidades de su empresa con las siguientes funciones:

- Escalado sencillo y operaciones no disruptivas. ONTAP admite la adición no disruptiva de capacidad a las controladoras existentes y a clústeres de escalado horizontal. Los clientes pueden empezar a utilizar tecnologías punteras como NVMe y FC 32 GB, sin necesidad de realizar costosas migraciones de datos y sin cortes.
- Conexión de cloud. ONTAP es el software de gestión de almacenamiento con mejor conexión de cloud e incluye opciones de almacenamiento definido por software (ONTAP Select) e instancias nativas del cloud (NetApp Cloud Volumes Service) en todos los clouds públicos.
- Integración con aplicaciones emergentes. ONTAP ofrece servicios de datos de clase empresarial para plataformas y aplicaciones de última generación, como vehículos autónomos, ciudades inteligentes e Industria 4.0, utilizando la misma infraestructura que da soporte a las aplicaciones empresariales existentes.

## Control Astra de NetApp

La familia de productos Astra de NetApp ofrece servicios de gestión de datos para aplicaciones y almacenamiento para aplicaciones de Kubernetes en las instalaciones y en el cloud público, con la tecnología de gestión de datos y almacenamiento de NetApp. Le permite realizar fácilmente backups de aplicaciones Kubernetes, migrar datos a un clúster diferente y crear, de forma instantánea, clones de aplicaciones de trabajo. Si necesita gestionar aplicaciones de Kubernetes que se ejecutan en un cloud público, consulte la documentación de "[Servicio de control Astra](#)". Astra Control Service es un servicio gestionado por NetApp que proporciona gestión de datos para aplicaciones de clústeres de Kubernetes en Google Kubernetes Engine (GKE) y Azure Kubernetes Service (AKS).

## Astra Trident de NetApp

Astra "[Trident](#)" De NetApp es un orquestador de almacenamiento dinámico de código abierto para Docker y Kubernetes que simplifica la creación, la gestión y el consumo de almacenamiento persistente. Trident, una aplicación nativa de Kubernetes, se ejecuta directamente dentro de un clúster de Kubernetes. Trident permite que los clientes implementen sin problemas imágenes de contenedores de DL en el almacenamiento de NetApp y proporciona una experiencia de clase empresarial para implementaciones de contenedores de IA. Los usuarios de Kubernetes (desarrolladores DE ML, científicos de datos, etc.) pueden crear, gestionar y automatizar la orquestación y el clonado para aprovechar las funcionalidades avanzadas de gestión de datos que se ofrecen con la tecnología de NetApp.

## Copia y sincronización de NetApp BlueXP

"[Copia y sincronización de BlueXP](#)" Es un servicio de NetApp que ofrece una sincronización de datos rápida y segura. Ya tenga que transferir archivos entre recursos compartidos de archivos NFS o SMB en las instalaciones, NetApp StorageGRID, NetApp ONTAP S3, NetApp Cloud Volumes Service, Azure NetApp Files, Amazon Simple Storage Service (Amazon S3), Amazon Elastic File System (Amazon EFS), Azure Blob, Google Cloud Storage, o IBM Cloud Object Storage, BlueXP Copy and Sync mueve los archivos a donde los necesites de forma rápida y segura. Una vez transferidos los datos, estarán completamente disponibles para su uso tanto en origen como en destino. BlueXP Copy y Sync sincronizan continuamente los datos en

función de tu programación predefinida, moviendo solo los deltas, por lo que se reducen al mínimo el tiempo y el dinero que se invierten en la replicación de datos. Copia y sincronización de BlueXP es una herramienta de software como servicio (SaaS) extremadamente sencilla de configurar y utilizar. Las transferencias de datos activadas por BlueXP Copy and Sync se llevan a cabo por agentes de datos. Puedes poner en marcha agentes de datos de BlueXP Copy y Sync en AWS, Azure, Google Cloud Platform o en las instalaciones.

**Clasificación de NetApp BlueXP**

Impulsado por potentes algoritmos de IA, "Clasificación de NetApp BlueXP" proporciona controles automatizados y control de datos en todos sus datos. Puede localizar con facilidad el ahorro de costes, identificar problemas relacionados con el cumplimiento de normativas y la privacidad, y buscar oportunidades de optimización. La consola de clasificación de BlueXP le ofrece la información necesaria para identificar los datos duplicados y eliminar la redundancia, asignar datos personales, no personales y confidenciales, así como activar alertas para datos confidenciales y anomalías.

**Plan de pruebas y validación**

Para este diseño de solución, se validaron los tres siguientes supuestos:

- Una tarea de inferencia, con y sin Protopía ofuscación, dentro de un espacio de trabajo JupyterLab que fue orquestada mediante el kit de herramientas DataOPS de NetApp para Kubernetes.
- Una tarea de inferencia por lotes, con y sin ofuscación de Protopía, en Kubernetes con un volumen de datos orquestado mediante el kit de herramientas DataOps de NetApp para Kubernetes.
- Una tarea de inferencia mediante una instancia del servidor de inferencia de NVIDIA Triton orquestada mediante el kit de herramientas NetApp DataOPS para Kubernetes. Aplicamos la confusión Protopia a la imagen antes de llamar a la API de inferencia Triton para simular el requisito común de que cualquier dato que se transmita a través de la red debe ser ocultado. Este flujo de trabajo es aplicable para casos en los que los datos se recopilan en una zona de confianza, pero debe pasarse fuera de esa zona de confianza para la inferencia. Sin la ocultación de Protopia, no es posible implementar este tipo de flujo de trabajo sin que los datos confidenciales salgan de la zona de confianza.

**Configuración de prueba**

La siguiente tabla describe el entorno de validación del diseño de la solución.

Componente	Versión
Kubernetes	1.21.6
Controlador Astra Trident CSI de NetApp	22.01.0
Kit de herramientas Data OPS de NetApp para Kubernetes	2.3.0
Servidor de inferencia NVIDIA Triton	21.11-py3

**Procedimiento de prueba**

En esta sección se describen las tareas necesarias para completar la validación.

**Requisitos previos**

Para ejecutar las tareas descritas en esta sección, debe tener acceso a un host Linux o MacOS con las siguientes herramientas instaladas y configuradas:

- Kubectl (se configura para acceder a un clúster de Kubernetes existente)
  - Se pueden encontrar instrucciones de instalación y configuración ["aquí"](#).
- Kit de herramientas Data OPS de NetApp para Kubernetes
  - Se pueden encontrar instrucciones de instalación ["aquí"](#).

### Escenario 1 – inferencia bajo demanda en JupyterLab

1. Cree un espacio de nombres de Kubernetes para las cargas de trabajo de inferencia de IA/ML.

```
$ kubectl create namespace inference
namespace/inference created
```

2. Utilice el kit de herramientas DataOPS de NetApp para aprovisionar un volumen persistente para almacenar los datos en los que realizará la inferencia.

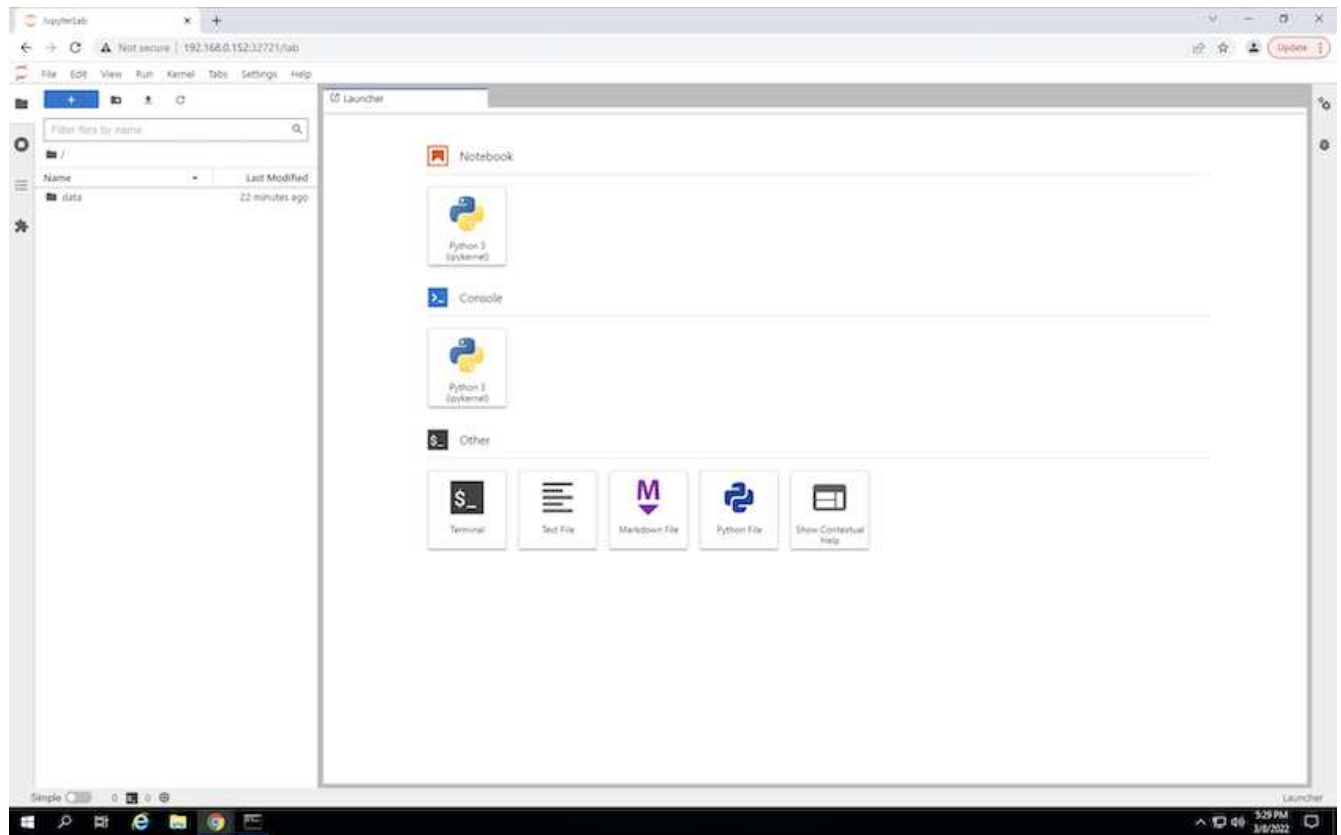
```
$ netapp_dataops_k8s_cli.py create volume --namespace=inference --pvc
-name=inference-data --size=50Gi
Creating PersistentVolumeClaim (PVC) 'inference-data' in namespace
'inference'.
PersistentVolumeClaim (PVC) 'inference-data' created. Waiting for
Kubernetes to bind volume to PVC.
Volume successfully created and bound to PersistentVolumeClaim (PVC)
'inference-data' in namespace 'inference'.
```

3. Use el kit de herramientas DataOPS de NetApp para crear un nuevo espacio de trabajo JupyterLab. Monte el volumen persistente que se creó en el paso anterior mediante el `--mount- pvc` opción. Asigne las GPU de NVIDIA al espacio de trabajo según sea necesario mediante el `-- nvidia-gpu` opción.

En el siguiente ejemplo, el volumen persistente `inference-data` Está montado en el contenedor de espacio de trabajo JJupyterLab en `/home/jovyan/data`. Cuando utilice las imágenes del contenedor de Jupyter del proyecto oficial, `/home/jovyan` Se presenta como el directorio de nivel superior dentro de la interfaz Web JupyterLab.

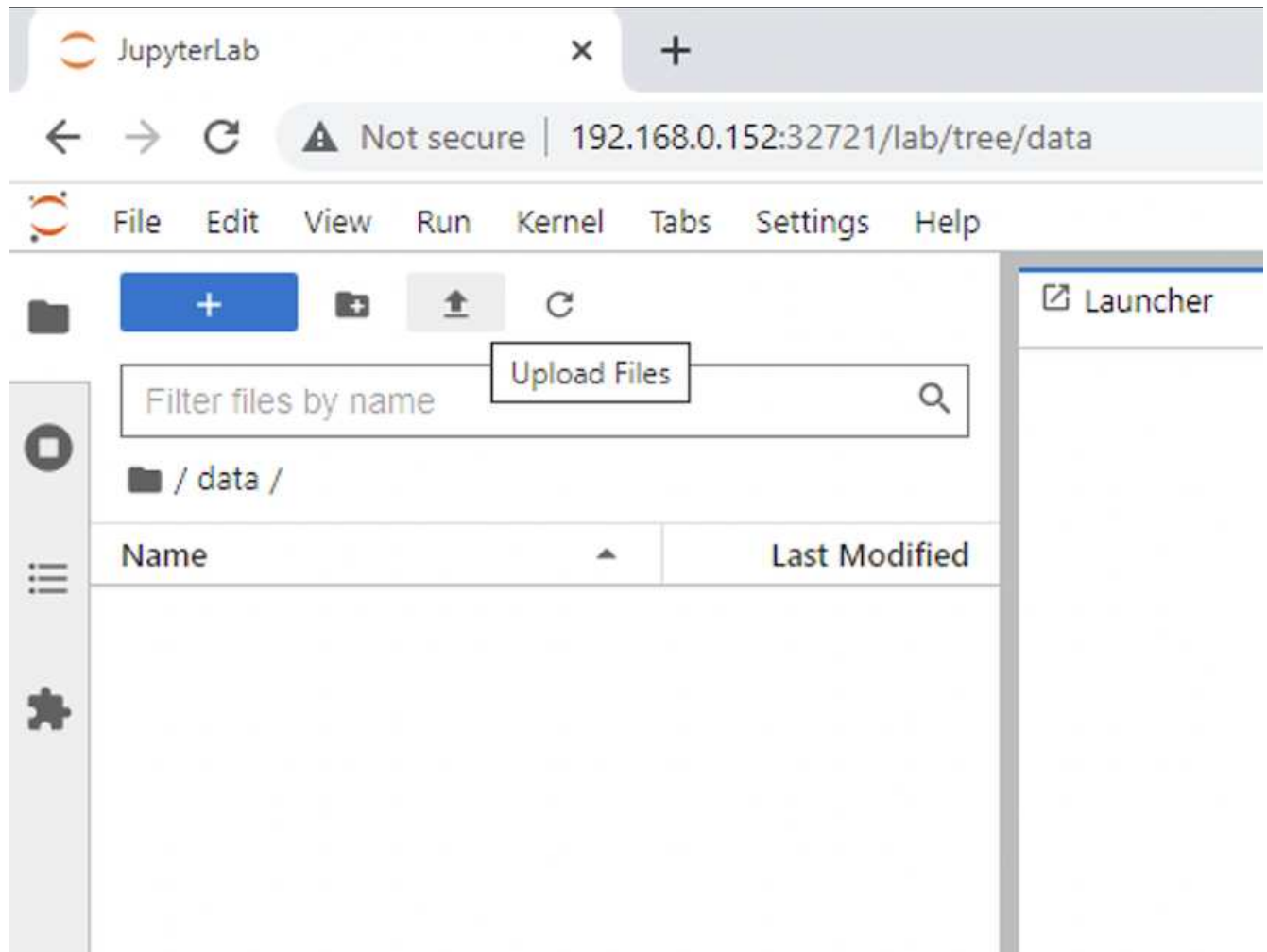
```
$ netapp_dataops_k8s_cli.py create jupyterlab --namespace=inference
--workspace-name=live-inference --size=50Gi --nvidia-gpu=2 --mount
-pvc=inference-data:/home/jovyan/data
Set workspace password (this password will be required in order to
access the workspace):
Re-enter password:
Creating persistent volume for workspace...
Creating PersistentVolumeClaim (PVC) 'ntap-dsutil-jupyterlab-live-
inference' in namespace 'inference'.
PersistentVolumeClaim (PVC) 'ntap-dsutil-jupyterlab-live-inference'
created. Waiting for Kubernetes to bind volume to PVC.
Volume successfully created and bound to PersistentVolumeClaim (PVC)
'ntap-dsutil-jupyterlab-live-inference' in namespace 'inference'.
Creating Service 'ntap-dsutil-jupyterlab-live-inference' in namespace
'inference'.
Service successfully created.
Attaching Additional PVC: 'inference-data' at mount_path:
'/home/jovyan/data'.
Creating Deployment 'ntap-dsutil-jupyterlab-live-inference' in namespace
'inference'.
Deployment 'ntap-dsutil-jupyterlab-live-inference' created.
Waiting for Deployment 'ntap-dsutil-jupyterlab-live-inference' to reach
Ready state.
Deployment successfully created.
Workspace successfully created.
To access workspace, navigate to http://192.168.0.152:32721
```

4. Acceda al espacio de trabajo JupyterLab utilizando la dirección URL especificada en la salida del `create jupyterlab` comando. El directorio de datos representa el volumen persistente que se montó en el espacio de trabajo.

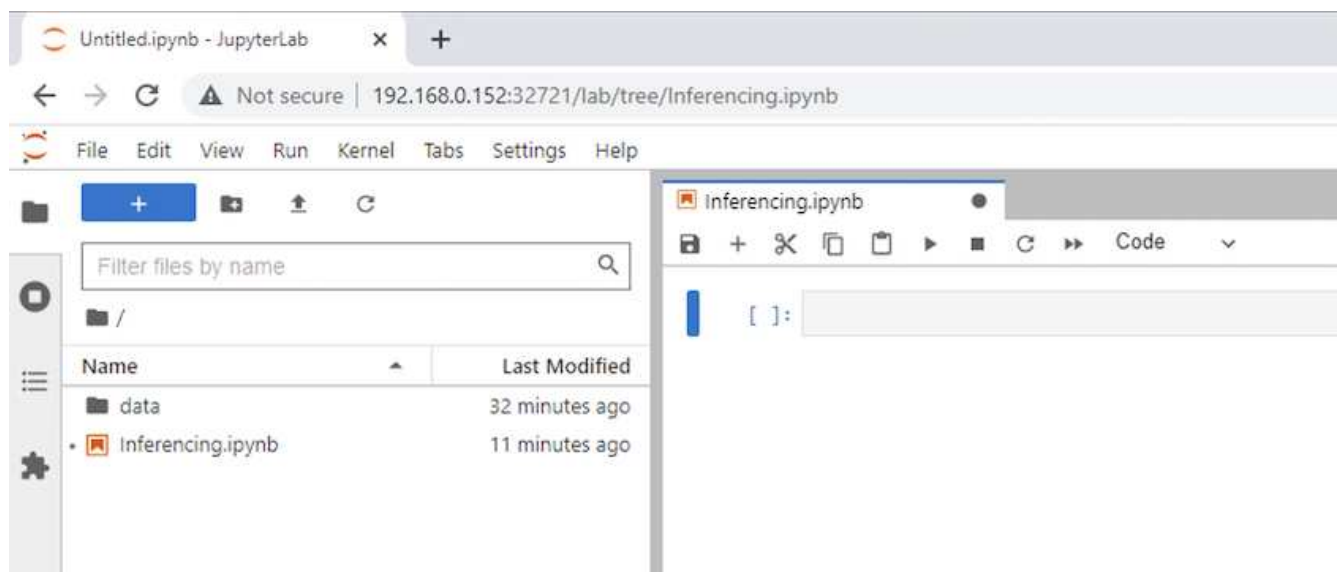


5. Abra el data directory y cargue los archivos en los que se va a realizar la inferencia. Cuando se cargan archivos en el directorio de datos, se almacenan automáticamente en el volumen persistente que se montó en el espacio de trabajo. Para cargar archivos, haga clic en el icono cargar archivos, como se muestra en la siguiente imagen.





6. Vuelva al directorio de nivel superior y cree un nuevo portátil.



7. Agregue el código de inferencia al cuaderno. En el siguiente ejemplo, se muestra el código de inferencia para un caso de uso de detección de imagen.



```
Launcher image-demo-pytorch.ipynb Python 3 (ipykernel)

STEP 3-1: Clean (Without obfuscation) detection

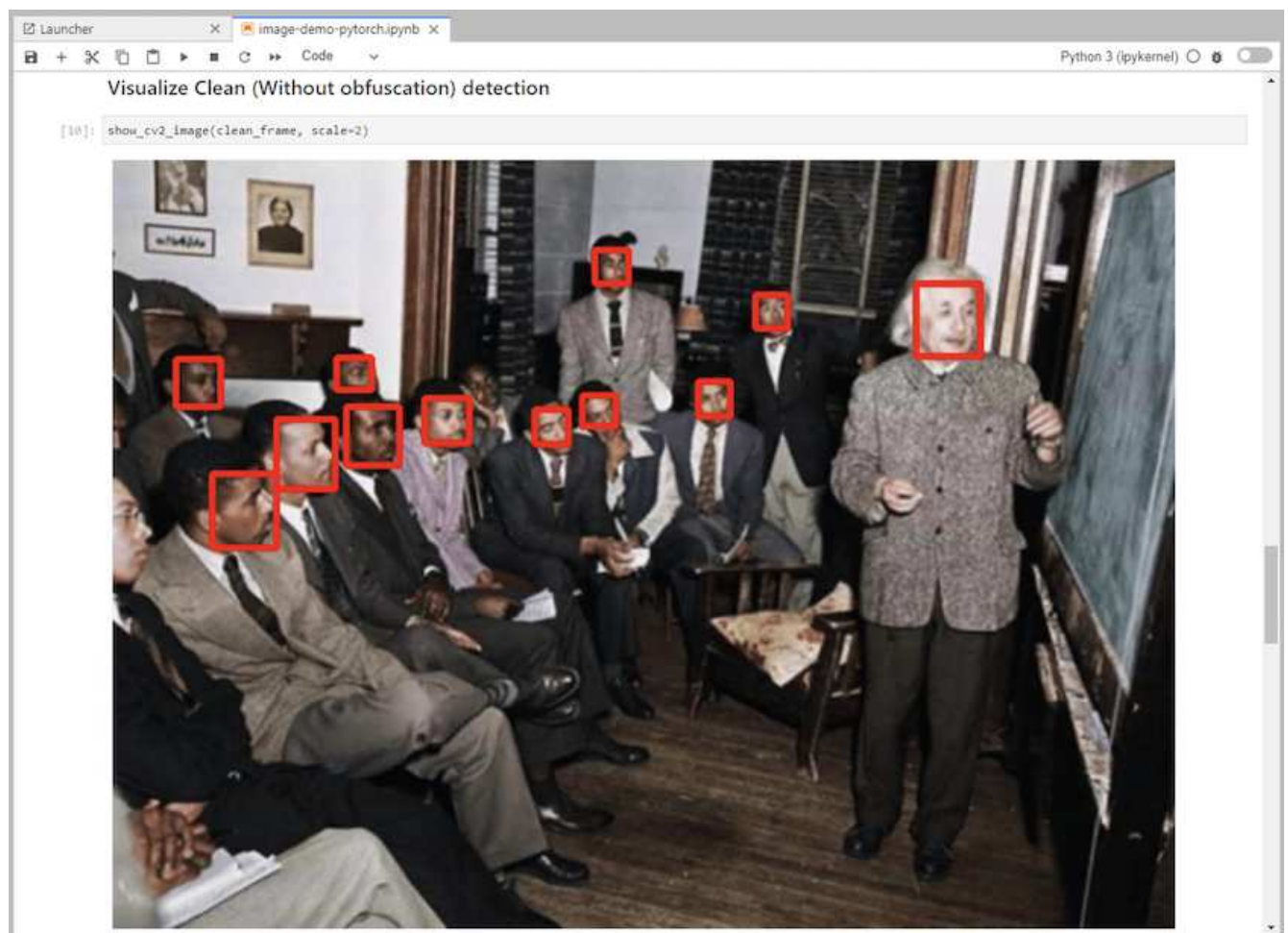
[9]: # get current frame
frame = input_image

# preprocess input
preprocessed_input = preprocess_input(frame)
preprocessed_input = torch.Tensor(preprocessed_input).to(device)

# run forward pass
clean_activation = clean_model.forward_head(preprocessed_input) # runs the first few layers
loc, pred = clean_model.forward_tail(clean_activation) # runs rest of the layers

# postprocess output
clean_pred = (loc.detach().cpu().numpy(), pred.detach().cpu().numpy())
clean_outputs = postprocess_outputs(
    clean_pred, [[input_image_width, input_image_height]], priors, THRESHOLD
)

# draw rectangles
clean_frame = copy.deepcopy(frame) # needs to be deep copy
for (x1, y1, x2, y2, s) in clean_outputs[0]:
    x1, y1 = int(x1), int(y1)
    x2, y2 = int(x2), int(y2)
    cv2.rectangle(clean_frame, (x1, y1), (x2, y2), (0, 0, 255), 4)
```



8. Agregue la ofuscación Protopia al código de inferencia. Protopia trabaja directamente con los clientes para proporcionar documentación específica para casos de uso y está fuera del alcance de este informe técnico. En el siguiente ejemplo se muestra el código de inferencia para un caso de uso de detección de imágenes con ofuscación Protopia agregada.

```
Launcher X image-demo-pytorch.ipynb X Python 3 (ipykernel)

STEP 3-2: Protopia AI (With obfuscation) detection

[11]: # get current frame
      frame = input_image

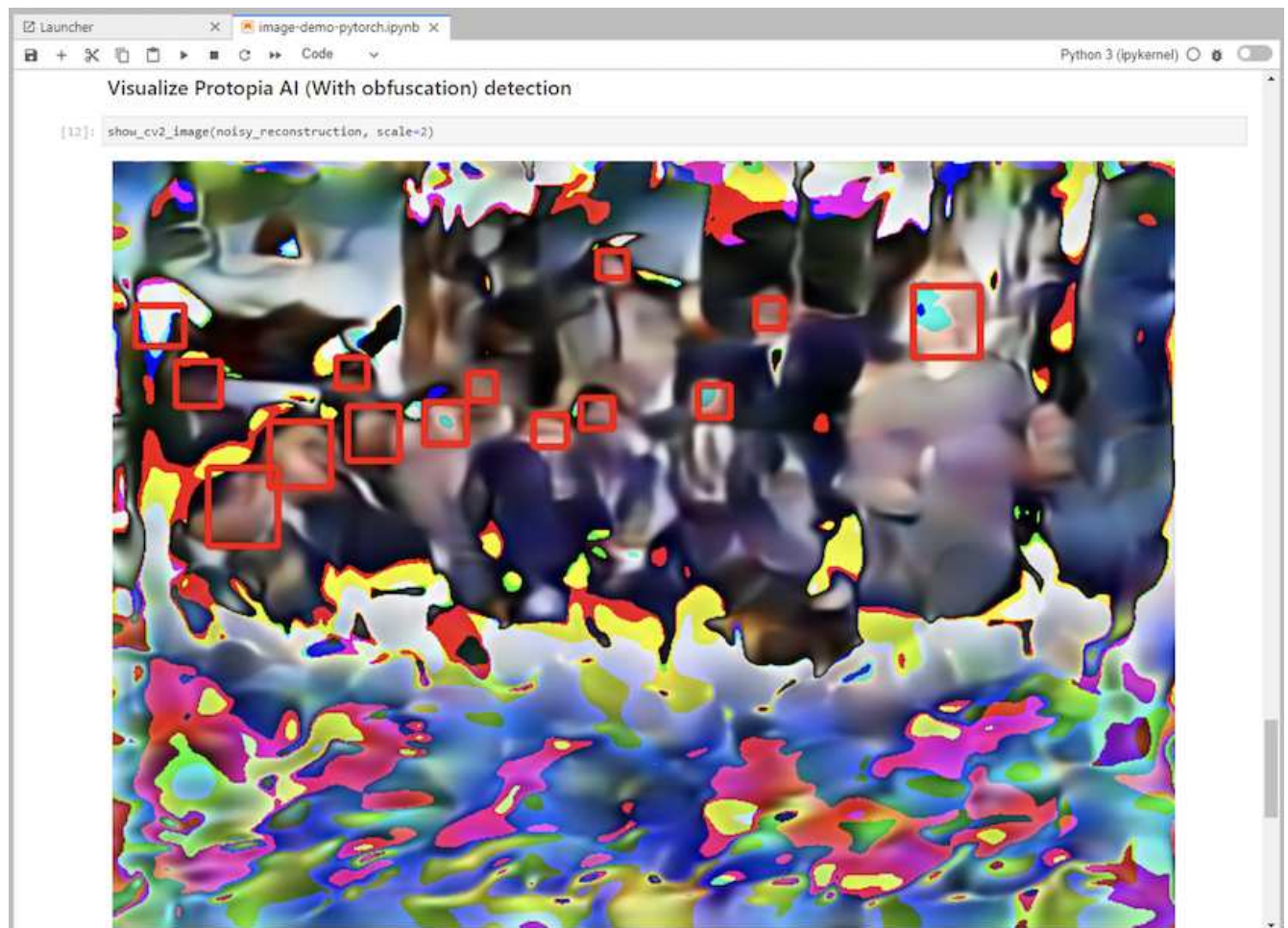
      # preprocess input
      preprocessed_input = preprocess_input(frame)
      preprocessed_input = torch.Tensor(preprocessed_input).to(device)

      # run forward pass
      not_noisy_activation = noisy_model.forward_head(preprocessed_input) # runs the first few layers
      #####
      # SINGLE ADDITIONAL LINE FOR PRIVATE INFERENCE #
      #####
      noisy_activation = noisy_model.forward_noise(not_noisy_activation)
      #####
      loc, pred = noisy_model.forward_tail(noisy_activation) # runs rest of the layers

      # postprocess output
      noisy_pred = (loc.detach().cpu().numpy(), pred.detach().cpu().numpy())
      noisy_outputs = postprocess_outputs(
          noisy_pred, [[input_image_width, input_image_height]], priors, THRESHOLD * 0.5
      )

      # get reconstruction of the noisy activation
      noisy_reconstruction = decoder_function(noisy_activation)
      noisy_reconstruction = noisy_reconstruction.detach().cpu().numpy()[0]
      noisy_reconstruction = unpreprocess_output(
          noisy_reconstruction, (input_image_width, input_image_height), True
      ).astype(np.uint8)

      # draw rectangles
      for (x1, y1, x2, y2, s) in noisy_outputs[0]:
          x1, y1 = int(x1), int(y1)
          x2, y2 = int(x2), int(y2)
          cv2.rectangle(noisy_reconstruction, (x1, y1), (x2, y2), (0, 0, 255), 4)
```



## Escenario 2: Inferencia por lotes en Kubernetes

1. Cree un espacio de nombres de Kubernetes para las cargas de trabajo de inferencia de IA/ML.

```
$ kubectl create namespace inference
namespace/inference created
```

2. Utilice el kit de herramientas DataOPS de NetApp para aprovisionar un volumen persistente para almacenar los datos en los que realizará la inferencia.

```
$ netapp_dataops_k8s_cli.py create volume --namespace=inference --pvc
-name=inference-data --size=50Gi
Creating PersistentVolumeClaim (PVC) 'inference-data' in namespace
'inference'.
PersistentVolumeClaim (PVC) 'inference-data' created. Waiting for
Kubernetes to bind volume to PVC.
Volume successfully created and bound to PersistentVolumeClaim (PVC)
'inference-data' in namespace 'inference'.
```

3. Rellene el nuevo volumen persistente con los datos en los que realizará la inferencia.

Existen varios métodos para cargar datos en un PVC. Si actualmente sus datos están almacenados en una plataforma de almacenamiento de objetos compatible con S3, como StorageGRID de NetApp o Amazon S3, podrá utilizar "[Funcionalidades de NetApp DataOPS Toolkit S3 Data mover](#)". Otro método simple es crear un espacio de trabajo JupyterLab y cargar archivos a continuación a través de la interfaz web JupyterLab, como se indica en los pasos 3 a 5 de la sección "[Escenario 1 – inferencia bajo demanda en JupyterLab](#)."

4. Cree un trabajo de Kubernetes para la tarea de inferencia de lotes. El siguiente ejemplo muestra un trabajo de inferencia en lote para un caso de uso de detección de imagen. Este trabajo realiza la inferencia en cada imagen de un conjunto de imágenes y escribe métricas de precisión de inferencia para su colocación.

```

$ vi inference-job-raw.yaml
apiVersion: batch/v1
kind: Job
metadata:
  name: netapp-inference-raw
  namespace: inference
spec:
  backoffLimit: 5
  template:
    spec:
      volumes:
      - name: data
        persistentVolumeClaim:
          claimName: inference-data
      - name: dshm
        emptyDir:
          medium: Memory
      containers:
      - name: inference
        image: netapp-protopia-inference:latest
        imagePullPolicy: IfNotPresent
        command: ["python3", "run-accuracy-measurement.py", "--dataset",
"/data/netapp-face-detection/Fddb"]
        resources:
          limits:
            nvidia.com/gpu: 2
        volumeMounts:
        - mountPath: /data
          name: data
        - mountPath: /dev/shm
          name: dshm
        restartPolicy: Never
$ kubectl create -f inference-job-raw.yaml
job.batch/netapp-inference-raw created

```

5. Confirme que el trabajo de inferencia se completó correctamente.

```

$ kubectl -n inference logs netapp-inference-raw-255sp
100%|██████████| 89/89 [00:52<00:00, 1.68it/s]
Reading Predictions : 100%|██████████| 10/10 [00:01<00:00, 6.23it/s]
Predicting ... : 100%|██████████| 10/10 [00:16<00:00, 1.64s/it]
===== Results =====
FDDb-fold-1 Val AP: 0.9491256561145955
FDDb-fold-2 Val AP: 0.9205024466101926
FDDb-fold-3 Val AP: 0.9253013871078468
FDDb-fold-4 Val AP: 0.9399781485863011
FDDb-fold-5 Val AP: 0.9504280149478732
FDDb-fold-6 Val AP: 0.9416473519339292
FDDb-fold-7 Val AP: 0.9241631566241117
FDDb-fold-8 Val AP: 0.9072663297546659
FDDb-fold-9 Val AP: 0.9339648715035469
FDDb-fold-10 Val AP: 0.9447707905560152
FDDb Dataset Average AP: 0.9337148153739079
=====
mAP: 0.9337148153739079

```

6. Agregue la ofuscación de Protopia a su trabajo de inferencia. Puede encontrar instrucciones específicas para casos de uso para agregar la ofuscación Protopia directamente desde Protopia, que está fuera del alcance de este informe técnico. El ejemplo siguiente muestra un trabajo de inferencia por lotes para un caso de uso de detección de cara con ofuscación Protopia agregada mediante un valor ALFA de 0.8. Este trabajo aplica la ofuscación Protopia antes de realizar la inferencia para cada imagen en un conjunto de imágenes y luego escribe las métricas de precisión de inferencia para el stdout.

Hemos repetido este paso para los valores ALFA 0.05, 0.1, 0.2, 0.4, 0.6, 0.8, 0.9 y 0.95. Puede ver los resultados en ["Comparación de precisión de inferencia."](#)

```

$ vi inference-job-protopia-0.8.yaml
apiVersion: batch/v1
kind: Job
metadata:
  name: netapp-inference-protopia-0.8
  namespace: inference
spec:
  backoffLimit: 5
  template:
    spec:
      volumes:
      - name: data
        persistentVolumeClaim:
          claimName: inference-data
      - name: dshm
        emptyDir:
          medium: Memory
      containers:
      - name: inference
        image: netapp-protopia-inference:latest
        imagePullPolicy: IfNotPresent
        env:
        - name: ALPHA
          value: "0.8"
        command: ["python3", "run-accuracy-measurement.py", "--dataset",
"/data/netapp-face-detection/Fddb", "--alpha", "$(ALPHA)", "--noisy"]
        resources:
          limits:
            nvidia.com/gpu: 2
        volumeMounts:
        - mountPath: /data
          name: data
        - mountPath: /dev/shm
          name: dshm
        restartPolicy: Never
$ kubectl create -f inference-job-protopia-0.8.yaml
job.batch/netapp-inference-protopia-0.8 created

```

7. Confirme que el trabajo de inferencia se completó correctamente.



```
$ kubectl -n inference logs netapp-inference-protopia-0.8-b4dkz
100%|██████████| 89/89 [01:05<00:00, 1.37it/s]
Reading Predictions : 100%|██████████| 10/10 [00:02<00:00, 3.67it/s]
Predicting ... : 100%|██████████| 10/10 [00:22<00:00, 2.24s/it]
===== Results =====
FDDb-fold-1 Val AP: 0.8953066115834589
FDDb-fold-2 Val AP: 0.8819580264029936
FDDb-fold-3 Val AP: 0.8781107458462862
FDDb-fold-4 Val AP: 0.9085731346308461
FDDb-fold-5 Val AP: 0.9166445508275378
FDDb-fold-6 Val AP: 0.9101178994188819
FDDb-fold-7 Val AP: 0.8383443678423771
FDDb-fold-8 Val AP: 0.8476311547659464
FDDb-fold-9 Val AP: 0.8739624502111121
FDDb-fold-10 Val AP: 0.8905468076424851
FDDb Dataset Average AP: 0.8841195749171925
=====
mAP: 0.8841195749171925
```

### Escenario 3: Servidor de inferencia NVIDIA Triton

1. Cree un espacio de nombres de Kubernetes para las cargas de trabajo de inferencia de IA/ML.

```
$ kubectl create namespace inference
namespace/inference created
```

2. Utilice el kit de herramientas DataOPS de NetApp para aprovisionar un volumen persistente y usarlo como repositorio de modelo para el servidor de inferencia NVIDIA Triton.

```
$ netapp_dataops_k8s_cli.py create volume --namespace=inference --pvc
-name=triton-model-repo --size=100Gi
Creating PersistentVolumeClaim (PVC) 'triton-model-repo' in namespace
'inference'.
PersistentVolumeClaim (PVC) 'triton-model-repo' created. Waiting for
Kubernetes to bind volume to PVC.
Volume successfully created and bound to PersistentVolumeClaim (PVC)
'triton-model-repo' in namespace 'inference'.
```

3. Almacene su modelo en el nuevo volumen persistente en un "formato" Reconocida por el servidor de inferencia NVIDIA Triton.

Existen varios métodos para cargar datos en un PVC. Un método simple es crear un espacio de trabajo JupyterLab y luego cargar archivos a través de la interfaz web JupyterLab, como se describe en los pasos 3 a 5 en ["Escenario 1 – inferencia bajo demanda en JupyterLab."](#)

4. Utilice el kit de herramientas DataOPS de NetApp para poner en marcha una nueva instancia del servidor de inferencia NVIDIA Triton.

```
$ netapp_dataops_k8s_cli.py create triton-server --namespace=inference
--server-name=netapp-inference --model-repo-pvc-name=triton-model-repo
Creating Service 'ntap-dsutil-triton-netapp-inference' in namespace
'inference'.
Service successfully created.
Creating Deployment 'ntap-dsutil-triton-netapp-inference' in namespace
'inference'.
Deployment 'ntap-dsutil-triton-netapp-inference' created.
Waiting for Deployment 'ntap-dsutil-triton-netapp-inference' to reach
Ready state.
Deployment successfully created.
Server successfully created.
Server endpoints:
http: 192.168.0.152: 31208
grpc: 192.168.0.152: 32736
metrics: 192.168.0.152: 30009/metrics
```

5. Utilice el SDK del cliente Triton para realizar una tarea de inferencia. El siguiente extracto de código de Python utiliza el SDK del cliente de Triton Python para realizar una tarea de inferencia para un caso de uso de detección facial. En este ejemplo se llama a la API de Triton y se pasa una imagen para la inferencia. A continuación, el servidor de inferencia Triton recibe la solicitud, invoca el modelo y devuelve la salida de inferencia como parte de los resultados de la API.

```
# get current frame
frame = input_image
# preprocess input
preprocessed_input = preprocess_input(frame)
preprocessed_input = torch.Tensor(preprocessed_input).to(device)
# run forward pass
clean_activation = clean_model_head(preprocessed_input) # runs the
first few layers
#####
#####
#           pass clean image to Triton Inference Server API for
inferencing           #
#####
#####
triton_client =
httpclient.InferenceServerClient(url="192.168.0.152:31208",
verbose=False)
model_name = "face_detection_base"
inputs = []
outputs = []
```



```

inputs.append(httpclient.InferInput("INPUT__0", [1, 128, 32, 32],
"FP32"))
inputs[0].set_data_from_numpy(clean_activation.detach().cpu().numpy(),
binary_data=False)
outputs.append(httpclient.InferRequestedOutput("OUTPUT__0",
binary_data=False))
outputs.append(httpclient.InferRequestedOutput("OUTPUT__1",
binary_data=False))
results = triton_client.infer(
    model_name,
    inputs,
    outputs=outputs,
    #query_params=query_params,
    headers=None,
    request_compression_algorithm=None,
    response_compression_algorithm=None)
#print(results.get_response())
statistics =
triton_client.get_inference_statistics(model_name=model_name,
headers=None)
print(statistics)
if len(statistics["model_stats"]) != 1:
    print("FAILED: Inference Statistics")
    sys.exit(1)

loc_numpy = results.as_numpy("OUTPUT__0")
pred_numpy = results.as_numpy("OUTPUT__1")
#####
#####
# postprocess output
clean_pred = (loc_numpy, pred_numpy)
clean_outputs = postprocess_outputs(
    clean_pred, [[input_image_width, input_image_height]], priors,
THRESHOLD
)
# draw rectangles
clean_frame = copy.deepcopy(frame) # needs to be deep copy
for (x1, y1, x2, y2, s) in clean_outputs[0]:
    x1, y1 = int(x1), int(y1)
    x2, y2 = int(x2), int(y2)
    cv2.rectangle(clean_frame, (x1, y1), (x2, y2), (0, 0, 255), 4)

```

6. Agregue la ofuscación Protopia al código de inferencia. Puede encontrar instrucciones específicas para casos de uso para agregar la ofuscación Protopia directamente desde Protopia; sin embargo, este proceso está fuera del alcance de este informe técnico. El ejemplo siguiente muestra el mismo código Python que se muestra en el paso anterior 5, pero con la ofuscación de Protopia agregada.

Tenga en cuenta que la confusión Protopia se aplica a la imagen antes de pasarla a la API de Triton. Así, la imagen no ofuscada nunca sale de la máquina local. Sólo la imagen oculta se pasa a través de la red. Este flujo de trabajo es aplicable para casos de uso en los que los datos se recopilan en una zona de confianza, pero luego debe pasarse fuera de esa zona de confianza para la inferencia. Sin la ocultación de Protopia, no es posible implementar este tipo de flujo de trabajo sin que haya datos confidenciales que salgan de la zona de confianza.

```
# get current frame
frame = input_image
# preprocess input
preprocessed_input = preprocess_input(frame)
preprocessed_input = torch.Tensor(preprocessed_input).to(device)
# run forward pass
not_noisy_activation = noisy_model_head(preprocessed_input) # runs the
first few layers
#####
#           obfuscate image locally prior to inferencing           #
#           SINGLE ADDITIONAL LINE FOR PRIVATE INFERENCE           #
#####
noisy_activation = noisy_model_noise(not_noisy_activation)
#####
#####
#####
#           pass obfuscated image to Triton Inference Server API for
inferencing           #
#####
#####
triton_client =
httpclient.InferenceServerClient(url="192.168.0.152:31208",
verbose=False)
model_name = "face_detection_noisy"
inputs = []
outputs = []
inputs.append(httpclient.InferInput("INPUT__0", [1, 128, 32, 32],
"FP32"))
inputs[0].set_data_from_numpy(noisy_activation.detach().cpu().numpy(),
binary_data=False)
outputs.append(httpclient.InferRequestedOutput("OUTPUT__0",
binary_data=False))
outputs.append(httpclient.InferRequestedOutput("OUTPUT__1",
binary_data=False))
results = triton_client.infer(
    model_name,
    inputs,
    outputs=outputs,
    #query_params=query_params,
    headers=None,
```

```

        request_compression_algorithm=None,
        response_compression_algorithm=None)
#print(results.get_response())
statistics =
triton_client.get_inference_statistics(model_name=model_name,
headers=None)
print(statistics)
if len(statistics["model_stats"]) != 1:
    print("FAILED: Inference Statistics")
    sys.exit(1)

loc_numpy = results.as_numpy("OUTPUT__0")
pred_numpy = results.as_numpy("OUTPUT__1")
#####

#####

# postprocess output
noisy_pred = (loc_numpy, pred_numpy)
noisy_outputs = postprocess_outputs(
    noisy_pred, [[input_image_width, input_image_height]], priors,
    THRESHOLD * 0.5
)
# get reconstruction of the noisy activation
noisy_reconstruction = decoder_function(noisy_activation)
noisy_reconstruction = noisy_reconstruction.detach().cpu().numpy()[0]
noisy_reconstruction = unpreprocess_output(
    noisy_reconstruction, (input_image_width, input_image_height), True
).astype(np.uint8)
# draw rectangles
for (x1, y1, x2, y2, s) in noisy_outputs[0]:
    x1, y1 = int(x1), int(y1)
    x2, y2 = int(x2), int(y2)
    cv2.rectangle(noisy_reconstruction, (x1, y1), (x2, y2), (0, 0, 255),
4)

```

### Comparación de precisión de inferencia

Para esta validación, realizamos la inferencia para un caso de uso de detección de imágenes mediante un conjunto de imágenes en bruto. A continuación, realizamos la misma tarea de inferencia en el mismo conjunto de imágenes con la ofuscación Protopía agregada antes de la inferencia. Repetimos la tarea usando diferentes valores DE ALFA para el componente de ofuscación de Protopia. En el contexto de la ofuscación de Protopia, el valor ALFA representa la cantidad de ofuscación que se aplica, con un valor ALFA más alto que representa un nivel más alto de ofuscación. A continuación, comparamos la precisión de la inferencia en estas carreras diferentes.

En las dos tablas siguientes se ofrecen detalles sobre nuestro caso de uso y se resumen los resultados.

Protopia trabaja directamente con los clientes para determinar el valor ALFA adecuado para un caso de uso específico.

Componente	Detalles
Modelo	FaceBoxes (PyTorch) -
Conjunto de datos	Conjunto de datos Fddb

Ofuscación Protopia	ALFA	Precisión
No	N.A.	0.9337148153739079
Sí	0.05	0.9028766627325002
Sí	0.1	0.9024301009661478
Sí	0.2	0.9081836283186224
Sí	0.4	0.9073066107482036
Sí	0.6	0.8847816568680239
Sí	0.8	0.8841195749171925
Sí	0.9	0.8455427675252052
Sí	0.95	0.8455427675252052

**Velocidad de ofuscación**

Para esta validación, aplicamos la ofuscación de Protopia a una imagen de 1920 x 1080 píxeles cinco veces y medimos la cantidad de tiempo que tardó en completar el paso de ofuscación cada vez.

Utilizamos PyTorch ejecutándose en una única GPU NVIDIA V100 para aplicar la ofuscación y borramos la memoria caché de la GPU entre ejecuciones. El paso de ofuscación tomó 5,47 ms, 5.27 ms, 4,54 ms, 5.24 ms y 4,84 ms respectivamente para completar las cinco carreras. La velocidad media era de 5,072 ms.

**Conclusión**

Los datos existen en tres estados: En reposo, en tránsito y en computación. Una parte importante de cualquier servicio de inferencia de IA debe ser la protección de datos contra amenazas durante todo el proceso. La protección de datos durante la inferencia es vital, ya que el proceso puede exponer información privada tanto sobre clientes externos como sobre la empresa que proporciona el servicio de inferencia. Protopia AI es una solución no obstrusiva de solo software para la inferencia de IA confidencial en el mercado actual. Con Protopia, la IA solo recibe la información transformada de los registros de datos que es esencial para llevar a cabo la tarea de IA/ML disponible y nada más. Esta transformación estocástica no es una forma de enmascarar y se basa en cambiar matemáticamente la representación de los datos utilizando el ruido curado.

Los sistemas de almacenamiento de NetApp con funcionalidades de ONTAP ofrecen el mismo rendimiento o

incluso mejor que el almacenamiento SSD local y, combinados con el kit de herramientas DataOPS de NetApp, ofrecen las siguientes ventajas a científicos de datos, ingenieros de datos, desarrolladores DE IA/ML y responsables DE la toma DE decisiones TECNOLÓGICAS empresariales o de negocio:

- Uso compartido de datos sin esfuerzo entre sistemas de IA, análisis y otros sistemas de negocio cruciales. Este uso compartido de datos reduce la sobrecarga de la infraestructura, mejora el rendimiento y optimiza la gestión de datos en toda la empresa.
- Almacenamiento y cálculo escalables de forma independiente para minimizar los costes y mejorar el uso de recursos.
- Flujos de trabajo de desarrollo y puesta en marcha optimizados que utilizan copias Snapshot integradas y clones para espacios de trabajo de usuario instantáneos con gestión eficiente del espacio, control de versiones integrado y una puesta en marcha automatizada.
- Protección de datos de clase empresarial y regulación de datos para la recuperación ante desastres, la continuidad del negocio y los requisitos normativos.
- Invocación simplificada de las operaciones de gestión de datos; realice rápidamente copias Snapshot de los espacios de trabajo de científicos de datos para realizar backups y trazabilidad desde el kit de herramientas de DataOps de NetApp en los portátiles Jupyter.

La solución de NetApp y Protopía proporciona una arquitectura flexible de escalado horizontal ideal para puestas en marcha de inferencia de IA de clase empresarial. Permite la protección de datos y proporciona privacidad para información confidencial en la que los requisitos confidenciales de inferencia de IA pueden satisfacerse con prácticas responsables de IA tanto en puestas en marcha en el entorno local como en cloud híbrido.

### **Dónde encontrar información adicional y reconocimientos**

Para obtener más información sobre la información descrita en este documento, consulte los siguientes documentos y/o sitios web:

- Software de gestión de datos ONTAP de NetApp: Biblioteca de información de ONTAP

<http://mysupport.netapp.com/documentation/productlibrary/index.html?productID=62286>

- Almacenamiento persistente de NetApp para contenedores: Trident de NetApp

["https://netapp.io/persistent-storage-provisioner-for-kubernetes/"](https://netapp.io/persistent-storage-provisioner-for-kubernetes/)

- Kit de herramientas de operaciones de datos de NetApp

["https://github.com/NetApp/netapp-dataops-toolkit"](https://github.com/NetApp/netapp-dataops-toolkit)

- Almacenamiento persistente de NetApp para contenedores: Astra Trident de NetApp

["https://netapp.io/persistent-storage-provisioner-for-kubernetes/"](https://netapp.io/persistent-storage-provisioner-for-kubernetes/)

- Protopía AI: Inferencia confidencial

["https://protopia.ai/blog/protopia-ai-takes-on-the-missing-link-in-ai-privacy-confidential-inference/"](https://protopia.ai/blog/protopia-ai-takes-on-the-missing-link-in-ai-privacy-confidential-inference/)

- Copia y sincronización de NetApp BlueXP

["https://docs.netapp.com/us-en/occm/concept\\_cloud\\_sync.html#how-cloud-sync-works"](https://docs.netapp.com/us-en/occm/concept_cloud_sync.html#how-cloud-sync-works)

- Servidor de inferencia NVIDIA Triton

["https://developer.nvidia.com/nvidia-triton-inference-server"](https://developer.nvidia.com/nvidia-triton-inference-server)

- Documentación de NVIDIA Triton inferencias Server

["https://docs.nvidia.com/deeplearning/triton-inference-server/index.html"](https://docs.nvidia.com/deeplearning/triton-inference-server/index.html)

- FaceBoxes en PyTorch

["https://github.com/zisianw/FaceBoxes.PyTorch"](https://github.com/zisianw/FaceBoxes.PyTorch)

## Reconocimientos

- Mark Cates, Director de producto, NetApp
- Sufian Ahmad, ingeniero técnico de marketing, NetApp
- Hadi Esmaeilzadeh, director técnico y profesor de Protopía AI

## Análisis de la confianza con IA de NetApp

### TR-4910: Análisis de sentimiento en las comunicaciones de los clientes con IA de NetApp

Rick Huang, Sathish Thyagarajan, y David Arnette, NetApp Diego Sosa-Coba, SFL Scientific

Este informe técnico proporciona directrices de diseño para que los clientes realicen análisis de opinión en un centro de soporte global de nivel empresarial usando las tecnologías de gestión de datos de NetApp con un marco de software de NVIDIA mediante la transferencia de aprendizaje e IA conversacional. Esta solución es aplicable a cualquier sector que desee obtener información de los clientes a partir de archivos de voz o texto grabados que representen registros de chat, correos electrónicos y otras comunicaciones de texto o audio. Hemos implementado una canalización integral para demostrar el reconocimiento automático de voz, el análisis de sensibilidad en tiempo real y el modelo de procesamiento de lenguaje natural y aprendizaje profundo:

Funcionalidades de nueva formación en un clúster informático acelerado por GPU con almacenamiento all-flash conectado al cloud de NetApp. Se pueden formar y optimizar modelos lingüísticos masivos de vanguardia para realizar inferencia rápidamente con el centro de soporte global y crear una experiencia de cliente excepcional y evaluaciones de rendimiento de empleados objetivas y a largo plazo.

El análisis del sentimiento es un campo de estudio dentro del procesamiento del lenguaje natural (NLP) por el cual se extraen sentimientos positivos, negativos o neutrales del texto. Los sistemas de IA conversacionales han aumentado hasta alcanzar un nivel de integración casi global a medida que cada vez más personas acuden para interactuar con ellos. El análisis de confianza tiene una variedad de casos de uso, desde determinar el rendimiento de los empleados del centro de soporte en conversaciones con los llamantes y proporcionar respuestas automatizadas adecuadas al bot conversacional para predecir el precio de las acciones de una empresa basándose en las interacciones entre representantes de la empresa y la audiencia en llamadas trimestrales a ganancias. Además, el análisis de confianza puede utilizarse para determinar la visión del cliente sobre los productos, servicios o asistencia proporcionados por la Marca.

Esta solución integral utiliza modelos NLP para realizar un análisis de confianza de alto nivel que permita

marcos analíticos del centro de soporte. Las grabaciones de audio se procesan en texto escrito y el sentimiento se extrae de cada frase de la conversación. Los resultados, agregados en un panel, se pueden elaborar para analizar los sentimientos de la conversación, tanto históricos como en tiempo real. Esta solución se puede generalizar a otras soluciones con modalidades de datos y necesidades de producción similares. Con los datos adecuados se pueden llevar a cabo otros casos de uso. Por ejemplo, las llamadas de beneficios de la empresa se pueden analizar para determinar la opinión utilizando la misma canalización de extremo a extremo. Otras formas de análisis de NLP, como el modelado de temas y el reconocimiento de entidades con nombre (TNM), también son posibles debido a la naturaleza flexible del ducto.

Estas implementaciones de IA fueron posibles gracias a NVIDIA RIVA, el kit de herramientas TAO de NVIDIA y el kit de herramientas DataOPS de NetApp trabajando juntos. Las herramientas de NVIDIA se utilizan para poner en marcha rápidamente soluciones de IA de gran rendimiento utilizando modelos y canalizaciones prediseñados. El kit de herramientas DataOPS de NetApp simplifica diversas tareas de gestión de datos para acelerar el desarrollo.

### **Valor para el cliente**

Las empresas ven el valor de una herramienta de evaluación de empleados y de reacción del cliente para la conversación de texto, audio y vídeo para el análisis de opiniones. Los gerentes se benefician de la información presentada en el panel, permitiendo una evaluación de los empleados y la satisfacción del cliente basada en ambos lados de la conversación.

Además, el kit de herramientas DataOps de NetApp gestiona el versionado y la asignación de datos dentro de la infraestructura del cliente. Esto lleva a frecuentes actualizaciones de los análisis presentados en el panel de control sin crear costes rígidos de almacenamiento de datos.

### **Casos de uso**

Debido al número de llamadas que procesa estos centros de soporte, la evaluación del rendimiento de la llamada puede llevar un tiempo considerable si se realiza manualmente. Los métodos tradicionales, como el conteo de bolsas de palabras y otros métodos, pueden lograr cierta automatización, pero estos métodos no capturan aspectos más matizados y un contexto semántico del lenguaje dinámico. Las técnicas de modelado de IA se pueden usar para realizar algunos de estos análisis más matices de forma automatizada. Además, con lo último en tecnología y herramientas de modelado preformadas publicadas por NVIDIA, AWS, Google y otras, puede ponerse en marcha una canalización integral con modelos complejos con una facilidad relativa personalizada.

Una canalización de extremo a extremo para el análisis de confianza del centro de soporte procesa archivos de audio en tiempo real a medida que los empleados converse con las personas que llaman. A continuación, estos archivos de audio se procesan para su uso en el componente voz a texto que los convierte en un formato de texto. Cada frase de la conversación recibe una etiqueta que indica el sentimiento (positivo, negativo o neutro).

El análisis de confianza puede proporcionar un aspecto esencial de las conversaciones para evaluar el rendimiento de la llamada. Estos sentimientos añaden un nivel adicional de profundidad a las interacciones entre empleados y personas que llaman. El panel de control de confianza asistido por IA proporciona a los administradores un seguimiento en tiempo real de la opinión dentro de una conversación, junto con un análisis retrospectivo de las llamadas pasadas del empleado.

Existen herramientas prediseñadas que se pueden combinar de manera potente para crear rápidamente una canalización de IA completa para solucionar este problema. En este caso, la biblioteca NVIDIA RIVA se puede

utilizar para realizar las dos tareas en serie: Transcripción de audio y análisis de sentimiento. El primero es un algoritmo de procesamiento de señales de aprendizaje supervisado y el segundo es un algoritmo de clasificación NLP de aprendizaje supervisado. Estos algoritmos listos para usar pueden ajustarse para cualquier caso de uso relevante con datos relevantes del negocio mediante el kit de herramientas TAO de NVIDIA. Esto lleva a que se estén creando soluciones más precisas y potentes por una fracción del coste y los recursos. Los clientes pueden incorporar la ["Maxine de NVIDIA"](#) Marco de trabajo para aplicaciones de videoconferencia aceleradas por GPU en el diseño del centro de soporte.

Los siguientes casos de uso son el núcleo de esta solución. En ambos casos de uso, se utiliza TAO Toolkit para realizar un ajuste preciso del modelo y RIVA para su implementación.

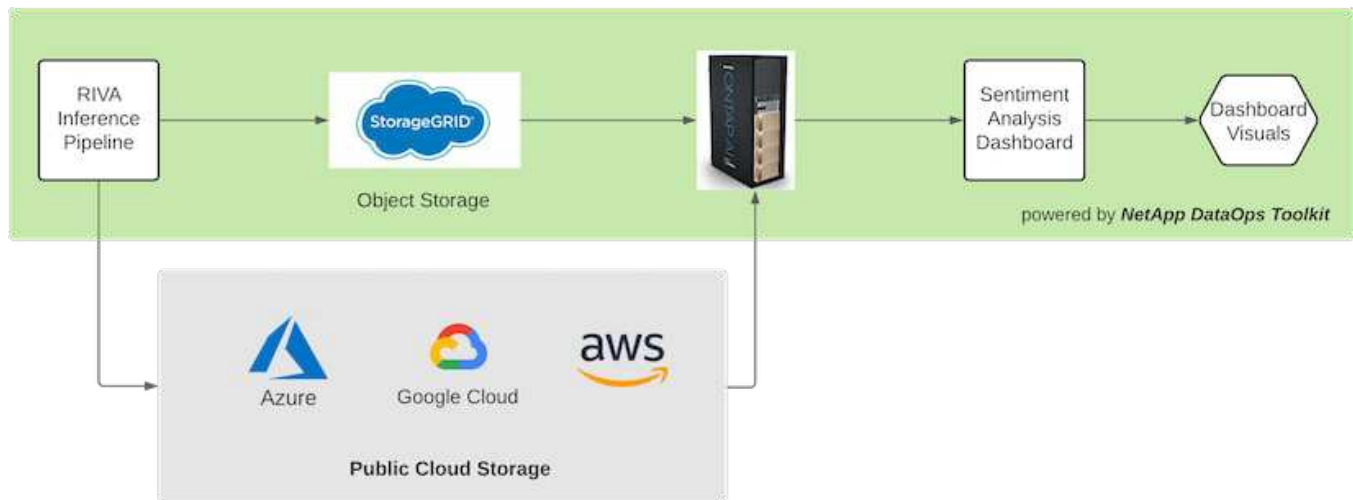
- Voz a texto
- Análisis de la confianza

Para analizar las interacciones del centro de soporte entre empleados y clientes, cada conversación del cliente en forma de llamadas de audio se puede realizar a través de la canalización para extraer sentimientos a nivel de frase. Esos sentimientos pueden ser verificados por un ser humano para justificar los sentimientos o ajustarlos según sea necesario. A continuación, los datos etiquetados se pasan al paso de ajuste fino para mejorar las predicciones de sentimientos. Si ya existen datos de confianza etiquetados, se puede acelerar el ajuste preciso del modelo. En cualquier caso, el ducto es generalizable a otras soluciones que requieren la ingestión de audio y la clasificación de frases.



Los resultados de opiniones de la IA se cargan en una base de datos cloud externa o en un sistema de almacenamiento gestionado por la empresa. Los resultados de la confianza se transfieren desde esta base de datos más grande al almacenamiento local para su uso dentro del panel de control que muestra el análisis de opinión de los administradores. La funcionalidad principal del panel es la de interactuar con el empleado del servicio de atención al cliente en tiempo real. Los gerentes pueden evaluar y proporcionar comentarios a los empleados durante sus llamadas con actualizaciones en vivo de la opinión de cada frase, así como una revisión histórica del desempeño pasado del empleado o de las reacciones del cliente.





La "[Kit de herramientas de operaciones de datos de NetApp](#)" Puede seguir gestionando sistemas de almacenamiento de datos incluso después de que la canalización de inferencia RIVA genere etiquetas de sentimiento. Estos resultados de la IA pueden cargarse en un sistema de almacenamiento de datos gestionado por el kit de herramientas DataOPS de NetApp. Los sistemas de almacenamiento de datos deben ser capaces de gestionar cientos de inserciones y seleccionar cada minuto. El sistema de almacenamiento del dispositivo local consulta el mayor almacenamiento de datos en tiempo real para su extracción. También es posible consultar a la instancia de almacenamiento de datos de mayor tamaño los datos históricos para mejorar aún más la experiencia de la consola. El kit de herramientas DataOPS de NetApp facilita ambos usos mediante el clonado rápido de los datos y la distribución entre todas las consolas que los utilizan.

### Público objetivo

El público objetivo de la solución incluye los siguientes grupos:

- Gerentes de empleados
- Ingenieros/científicos de datos
- Administradores DE TECNOLOGÍA (en las instalaciones, en cloud o híbrida)

El seguimiento de los sentimientos a lo largo de las conversaciones es una herramienta valiosa para evaluar el rendimiento de los empleados. Mediante el panel de instrumentos de inteligencia artificial, los administradores pueden ver cómo los empleados y las personas que llaman cambian sus sentimientos en tiempo real, lo que permite evaluaciones en directo y sesiones de orientación. Además, las empresas pueden obtener información valiosa de los clientes que participan en conversaciones vocales, bots de chat de texto y videoconferencias. Dichos análisis de clientes utilizan las funcionalidades de procesamiento multimodal a escala con modernos modelos de IA y flujos de trabajo de vanguardia.

En cuanto a los datos, el centro de soporte procesa diariamente un gran número de archivos de audio. El kit de herramientas DataOPS de NetApp facilita esta tarea de gestión de datos tanto para la sintonización periódica de modelos como para consolas de análisis de opiniones.

Los administradores DE TECNOLOGÍA también se benefician del kit de herramientas de NetApp DataOps ya que permite mover datos con rapidez entre entornos de implementación y de producción. Los entornos y servidores de NVIDIA también deben gestionarse y distribuirse para permitir la inferencia en tiempo real.

### Arquitectura

La arquitectura de esta solución de centro de soporte se basa en las herramientas

prediseñadas de NVIDIA y el kit de herramientas Data OPS de NetApp. Las herramientas de NVIDIA se utilizan para poner en marcha rápidamente soluciones de IA de alto rendimiento utilizando modelos y canalizaciones prediseñados. El kit de herramientas DataOPS de NetApp simplifica diversas tareas de gestión de datos para acelerar el desarrollo.

### Tecnología de soluciones

**"RIVA DE NVIDIA"** Es un SDK acelerado por GPU para crear aplicaciones de IA conversacionales multimodales que ofrecen rendimiento en tiempo real en las GPU. El kit de herramientas NVIDIA Train, Adapt y Optimize (TAO) proporciona un método más rápido y sencillo para acelerar la formación y crear rápidamente modelos de IA específicos para dominios con un alto rendimiento y precisión.

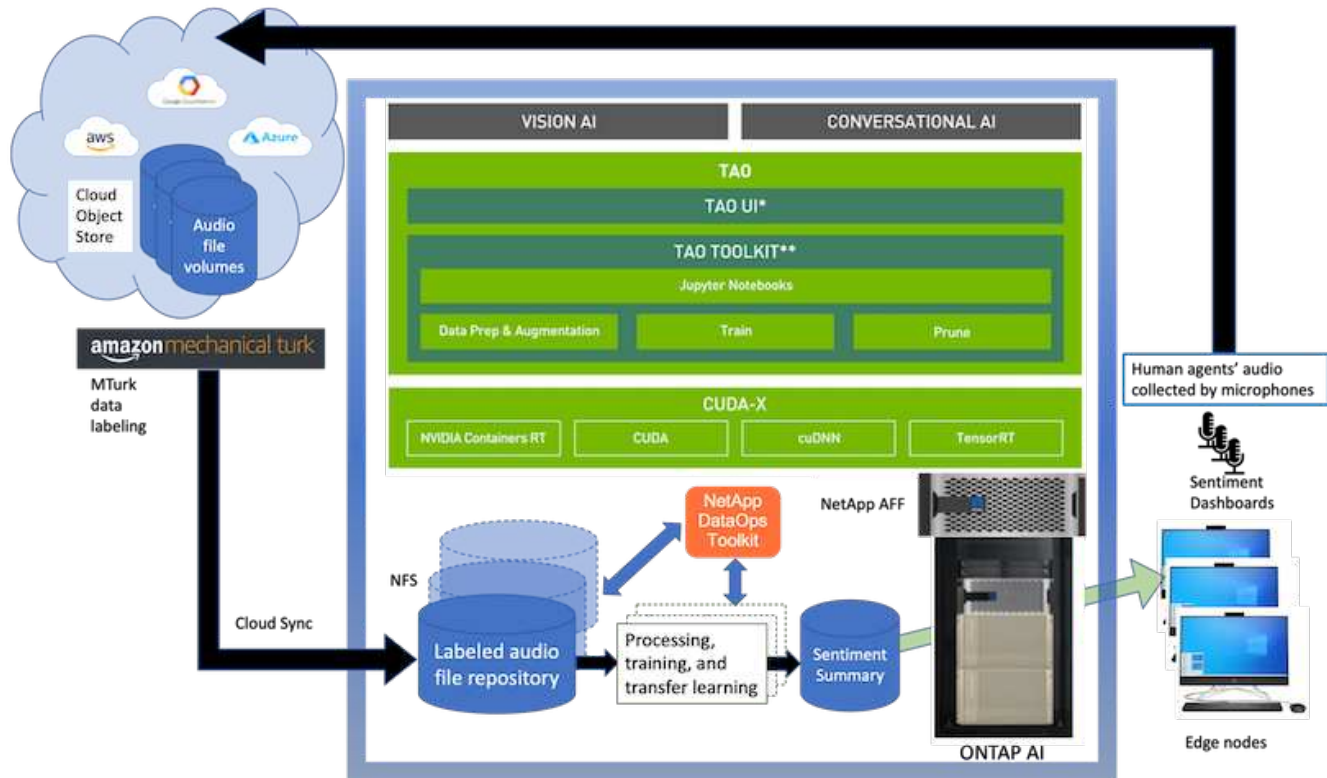
El kit de herramientas NetApp DataOPS es una biblioteca Python que facilita a los desarrolladores, científicos de datos, ingenieros de DevOps e ingenieros de datos la tarea de realizar varias tareas de gestión de datos. Esto incluye el aprovisionamiento casi instantáneo de un nuevo volumen de datos o espacio de trabajo JupyterLab, el clonado casi instantáneo de un volumen de datos o un espacio de trabajo JupyterLab y la creación casi instantánea de copias Snapshot de un volumen de datos o de un espacio de trabajo JupyterLab para su seguimiento y línea de base.

### Diagrama arquitectónico

El siguiente diagrama muestra la arquitectura de la solución. Existen tres categorías de entorno principales: El cloud, el núcleo y el perímetro. Cada una de las categorías puede estar geográficamente dispersa. Por ejemplo, la nube contiene almacenes de objetos con archivos de audio en bloques de diferentes regiones, mientras que el núcleo puede contener centros de datos vinculados a través de una red de alta velocidad o Copia y sincronización de NetApp BlueXP. Los nodos EDGE denotan las plataformas de trabajo diario del agente humano individual, en las que hay disponibles micrófonos y herramientas de panel interactivas para visualizar la opinión y recopilar datos de audio de conversaciones con los clientes.

En los centros de datos acelerados por GPU, las empresas pueden utilizar NVIDIA **"RIVA"** Marco para crear aplicaciones de IA conversacionales, con las que el **"Kit de herramientas Tao"** Se conecta para la finetunización de modelos y el reciclaje mediante técnicas de aprendizaje en L de transferencia. Estas aplicaciones y flujos de trabajo de computación funcionan con la **"Kit de herramientas de operaciones de datos de NetApp"**, Activación de las mejores capacidades de administración de datos que ONTAP tiene para ofrecer. El kit de herramientas permite a los equipos de datos corporativos realizar rápidamente prototipos de sus modelos con datos estructurados y no estructurados asociados a través de copias Snapshot y clones para llevar a cabo seguimientos, crear versiones, realizar pruebas A/B, proporcionando así seguridad, gobernabilidad, y cumplimiento de las normativas. Consulte la sección **"Diseño del almacenamiento"** para obtener más detalles.

Esta solución muestra los pasos detallados del procesamiento de archivos de audio, el entrenamiento de modelos NLP, el aprendizaje de transferencias y la administración de datos. La canalización integral resultante genera un resumen de sentimiento que se muestra en tiempo real en los paneles de control de los agentes de apoyo humano.



## Requisitos de hardware

En la siguiente tabla se enumeran los componentes de hardware necesarios para implementar la solución. Los componentes de hardware que se usan en cualquier implementación particular de la solución pueden variar en función de las necesidades del cliente.

Pruebas de latencia de respuesta	Tiempo (milisegundos)
Procesamiento de datos	10
Inferencia	10

Estas pruebas de tiempo de respuesta se realizaron en más de 50,000 archivos de audio en 560 conversaciones. Cada archivo de audio tenía un tamaño de ~100 KB como MP3 y de ~1 MB cuando se convirtió a WAV. El paso de procesamiento de datos convierte MP3s en archivos WAV. Los pasos de inferencia convierten los archivos de audio en texto y extraen un sentimiento del texto. Estos pasos son todos independientes entre sí y pueden ser paralelizados para acelerar el proceso.

Teniendo en cuenta la latencia de transferencia de datos entre almacenes, los administradores deben poder ver actualizaciones del análisis de opinión en tiempo real en un segundo del final de la frase.

## Hardware NVIDIA RIVA

Hardware subyacente	Requisitos
SO	Linux x86_64
Memoria GPU (ASR)	Modelos de streaming: Aprox. 5600 MB modelos sin secuencias: Aprox. 3100 MB
Memoria GPU (NLP)	~500 MB por modelo BERT

## Hardware del kit de herramientas TAO de NVIDIA

Hardware subyacente	Requisitos
RAM del sistema	32 GB
RAM DE GPU	32 GB
CPU	8 núcleos
GPU	NVIDIA (A100, V100 y RTX 30x0)
SSD	100 GB

## Sistema de almacenamiento flash

### ONTAP 9 de NetApp

ONTAP 9.9, la última generación del software de gestión del almacenamiento de NetApp, permite a las empresas modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. ONTAP ofrece las mejores capacidades de gestión de datos y permite la gestión y protección de los datos con un solo conjunto de herramientas, sin importar dónde residan. También puede mover los datos libremente a donde sea necesario: El perímetro, el núcleo o el cloud. ONTAP 9.9 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen los datos esenciales y permiten disfrutar de funcionalidades de infraestructura de nueva generación en arquitecturas de cloud híbrido.

### Copia y sincronización de NetApp BlueXP

"[Copia y sincronización de BlueXP](#)" Es un servicio de NetApp que ofrece una sincronización de datos rápida y segura que permite transferir archivos entre recursos compartidos de archivos NFS o SMB en las instalaciones a cualquiera de los siguientes destinos:

- StorageGRID de NetApp
- ONTAP S3 de NetApp
- Cloud Volumes Service de NetApp
- Azure NetApp Files
- Simple Storage Service (Amazon S3)
- Sistema de archivos Elastic de Amazon (Amazon EFS)
- Azure Blob
- Google Cloud Storage
- Almacenamiento de objetos en cloud de IBM

BlueXP Copy and Sync mueve los archivos a donde los necesitas de una forma rápida y segura. Una vez transferidos los datos, estarán completamente disponibles para su uso tanto en el origen como en el destino. BlueXP Copy and Sync sincroniza los datos de forma continua en función de tu programación predefinida, moviendo solo los deltas, por lo que se reducen al mínimo el tiempo y el dinero que se invierten en la replicación de datos. BlueXP Copy and Sync es una herramienta de software como servicio (SaaS) fácil de configurar y usar. Las transferencias de datos activadas por BlueXP Copy and Sync se llevan a cabo por agentes de datos. Puedes poner en marcha agentes de datos de BlueXP Copy y Sync en AWS, Azure, Google Cloud Platform o en las instalaciones.

## StorageGRID de NetApp

La suite de almacenamiento de objetos definida por software StorageGRID admite una amplia gama de casos de uso en entornos multicloud públicos, privados e híbridos sin problemas. Con innovaciones líderes del sector, StorageGRID de NetApp almacena, protege y preserva datos no estructurados para usos múltiples, incluida la gestión automatizada del ciclo de vida durante largos periodos de tiempo. Para obtener más información, consulte ["StorageGRID de NetApp" sitio](#).

### Requisitos de software

En la siguiente tabla se enumeran los componentes de software necesarios para implementar esta solución. Los componentes que se usan en cualquier implementación particular de la solución pueden variar en función de las necesidades del cliente.

Máquina host	Requisitos
RIVA (anteriormente JARVIS)	1.4.0
TAO Toolkit (antes Transfer Learning Toolkit)	3.0
ONTAP	9.9.1
SO DGX	5.1
DOTK	2.0.0

### Software NVIDIA RIVA

De NetApp	Requisitos
Docker	>19.02 (con nvidia-docker instalado)>=19.03 si no se utiliza DGX
Controlador NVIDIA	Más de 465.19.01 418.40+, 440.33+, 450.51+ y 460.27+ para las GPU del centro de datos
So del contenedor	Ubuntu 20.04
CUDA	11.3.0
CuBLAS	11.5.1.101
CuDNN	8.2.0.41
NCCL	2.9.6
TensorRT	7.2.3.4
Servidor de inferencia Triton	2.9.0

### Software NVIDIA TAO Toolkit

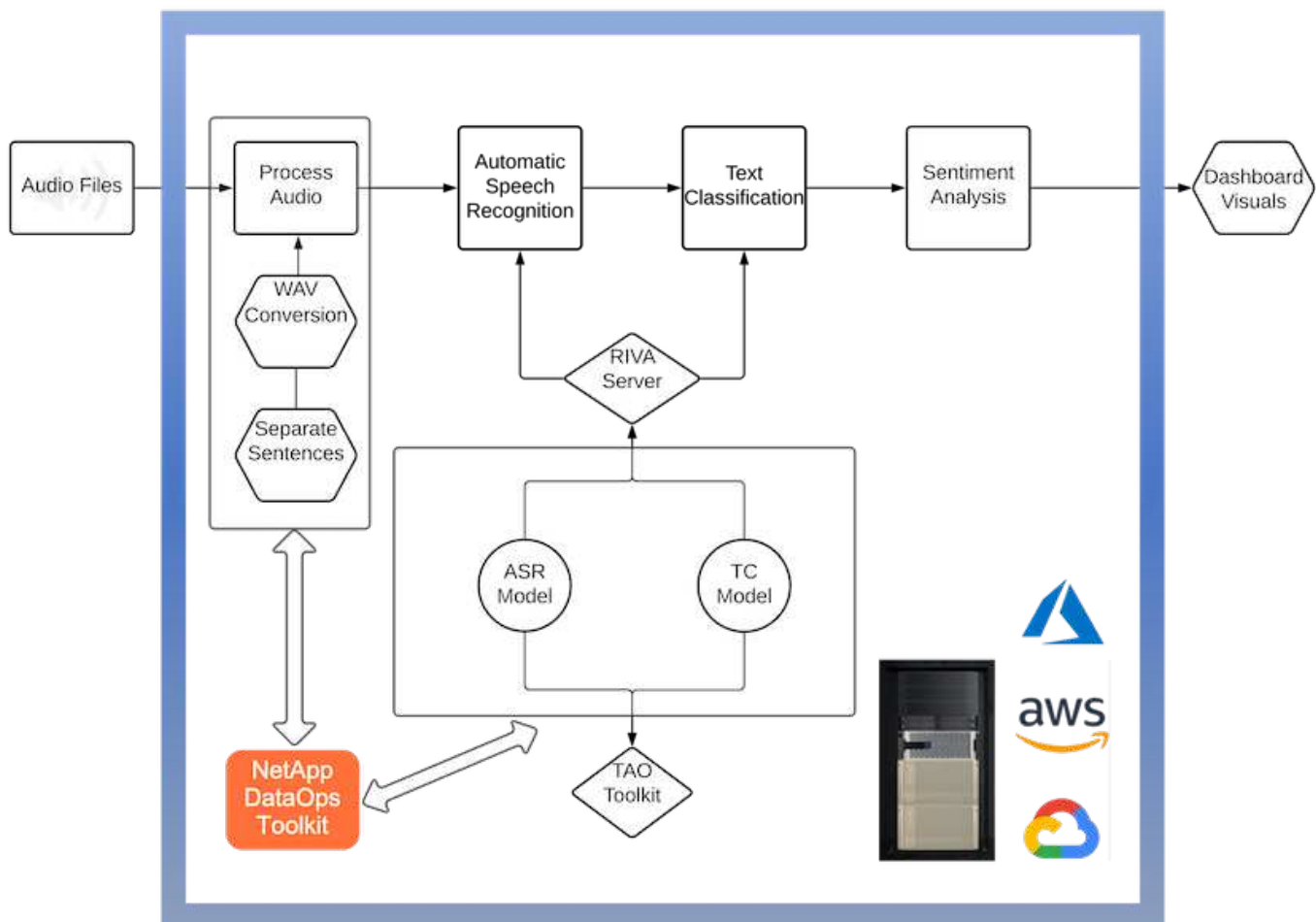
De NetApp	Requisitos
Sistema operativo Ubuntu 18.04 LTS	18.04
python	>=3.6.9
docker-ce	>19.03.5

De NetApp	Requisitos
docker-API	1.40
kit de herramientas de nvidia-container	>1.3.0-1
nvidia-container-runtime	3.4.0-1
nvidia-docker2	2.5.0-1
controlador nvidia	>455
python-pip	>21.06
nvidia-pyindex	Última versión

## Detalles de casos de uso

Esta solución se aplica a los siguientes casos de uso:

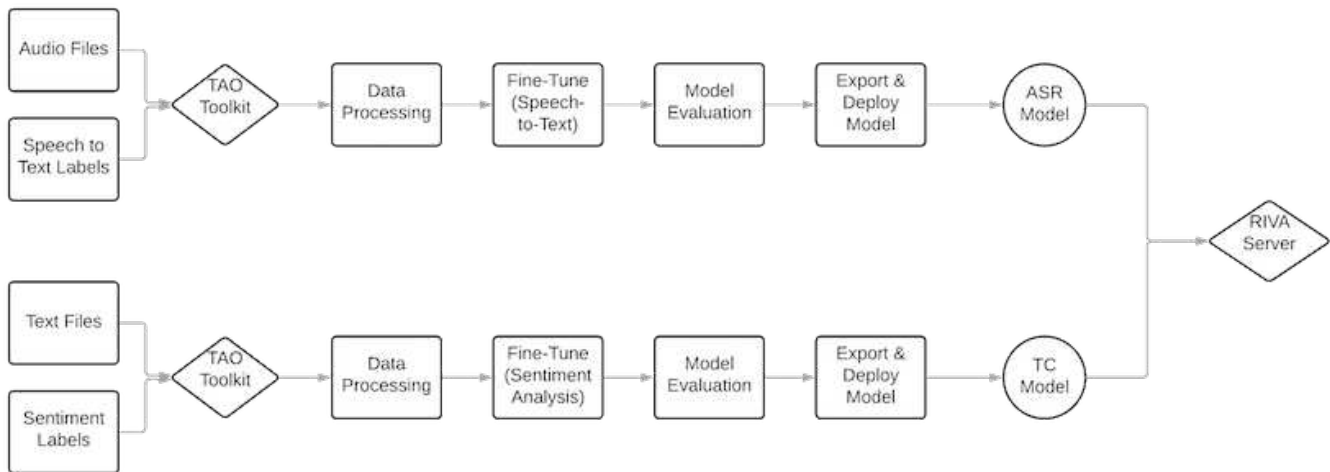
- Voz a texto
- Análisis de la confianza



El caso práctico de voz a texto comienza ingiriendo archivos de audio para los centros de soporte. Este audio se procesa entonces para ajustarse a la estructura requerida por RIVA. Si los archivos de audio aún no se han dividido en sus unidades de análisis, esto se debe hacer antes de pasar el audio a RIVA. Una vez procesado el archivo de audio, se pasa al servidor RIVA como una llamada API. El servidor emplea uno de los muchos

modelos que aloja y devuelve una respuesta. Este mensaje de voz a texto (parte del reconocimiento automático de voz) devuelve una representación de texto del audio. A partir de ahí, la canalización pasa a la parte del análisis de confianza.

Para el análisis de confianza, la salida de texto del reconocimiento automático de voz sirve como entrada a la clasificación de texto. Text Classification es el componente NVIDIA para clasificar el texto en cualquier número de categorías. Las categorías de sentimiento varían de positivo a negativo para las conversaciones del centro de apoyo. El rendimiento de los modelos se puede evaluar utilizando un conjunto de holdout para determinar el éxito del paso de ajuste fino.



Una canalización similar se utiliza tanto para el análisis de voz a texto como para el análisis de sentimiento dentro del TAO Toolkit. La principal diferencia es el uso de etiquetas que se requieren para la afinación de los modelos. La canalización del kit de herramientas TAO comienza con el procesamiento de los archivos de datos. A continuación, los modelos preformados (procedentes de "[Catálogo de NVIDIA NGC](#)") se ajustan con precisión mediante los datos del centro de soporte. Los modelos ajustados con precisión se evalúan en función de sus mediciones de rendimiento correspondientes y, si tienen un rendimiento superior al de los modelos preformados, se implementan en EL servidor RIVA.

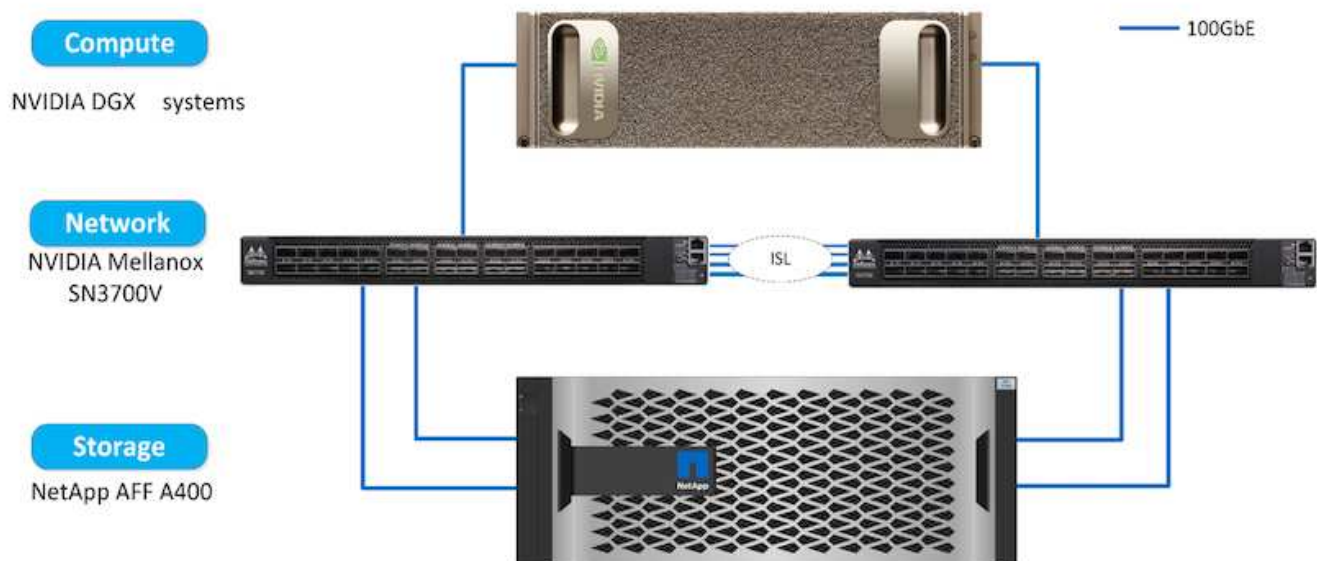
## Consideraciones de diseño

En esta sección se describen las consideraciones de diseño para los distintos componentes de esta solución.

### Diseño informático y de red

En función de las restricciones en la seguridad de los datos, todos los datos deben permanecer dentro de la infraestructura del cliente o en un entorno seguro.





### Diseño del almacenamiento

El kit de herramientas Data OPS de NetApp sirve como servicio principal para gestionar los sistemas de almacenamiento. El kit de herramientas DataOps es una biblioteca Python que facilita a los desarrolladores, científicos de datos, ingenieros de DevOps e ingenieros de datos la tarea de gestión de datos realizar varias tareas de gestión de datos, como el aprovisionamiento casi instantáneo de un nuevo volumen de datos o un espacio de trabajo JupyterLab, el clonado casi instantáneo de un volumen de datos o el espacio de trabajo JupyterLab. Y una copia de Snapshot casi instantánea de un volumen de datos o un espacio de trabajo JupyterLab para poder seguir su seguimiento o crear una línea de base. Esta biblioteca de Python puede funcionar como una utilidad de línea de comandos o una biblioteca de funciones que se pueden importar a cualquier programa de Python o a cualquier ordenador portátil Jupyter.

### Mejores prácticas de RIVA

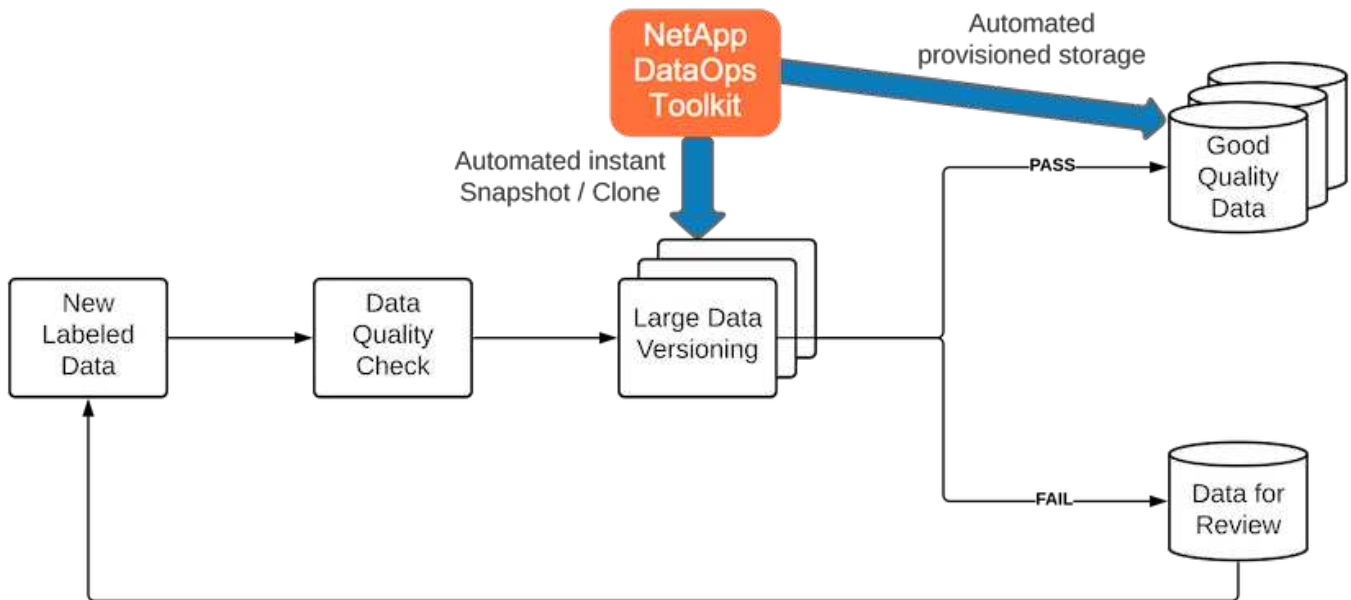
NVIDIA proporciona varias funciones generales "[mejores prácticas de datos](#)" Para utilizar RIVA:

- **Utilice formatos de audio sin pérdidas si es posible.** el uso de códecs con pérdida como MP3 puede reducir la calidad.
- **Aumentar los datos de entrenamiento.** la adición de ruido de fondo a los datos de entrenamiento de audio puede inicialmente disminuir la precisión y, sin embargo, aumentar la solidez.
- **Limite el tamaño del vocabulario si utiliza texto raspado.** muchas fuentes en línea contienen tipopos o pronombres auxiliares y palabras poco comunes. La eliminación de estos elementos puede mejorar el modelo de idioma.
- **Utilice una frecuencia de muestreo mínima de 16 kHz si es posible.** sin embargo, trate de no resampling, ya que al hacerlo se reduce la calidad de audio.

Además de estas mejores prácticas, los clientes deben priorizar la recopilación de un conjunto de datos de ejemplo representativo con etiquetas precisas para cada paso de la canalización. En otras palabras, el conjunto de datos de ejemplo debería reflejar proporcionalmente las características especificadas ejemplificadas en un conjunto de datos de destino. De forma similar, los anotadores del conjunto de datos tienen la responsabilidad de equilibrar la precisión y la velocidad del etiquetado, de modo que se maximice la calidad y la cantidad de los datos. Por ejemplo, esta solución de centro de soporte requiere archivos de audio, etiquetas de texto etiquetadas y etiquetas de sentimiento. La naturaleza secuencial de esta solución significa que los errores desde el principio de la canalización se propagan hasta el final Si los archivos de audio son de mala calidad, las transcripciones de texto y las etiquetas de sentimiento también serán.



Esta propagación de errores también se aplica a los modelos entrenados en estos datos. Si las predicciones de sentimiento son 100% exactas pero el modelo de voz a texto tiene un rendimiento deficiente, entonces la canalización final está limitada por las transcripciones iniciales de audio a texto. Es esencial que los desarrolladores consideren el rendimiento de cada modelo individualmente y como un componente de una canalización mayor. En este caso en particular, el objetivo final es desarrollar una canalización que pueda predecir con precisión el sentimiento. Por lo tanto, la métrica general sobre la cual evaluar la canalización es la precisión de los sentimientos, que afecta directamente la transcripción del habla al texto.



El kit de herramientas DataOPS de NetApp complementa la canalización de comprobación de la calidad de los datos mediante su tecnología de clonado de datos casi instantánea. Cada archivo etiquetado debe evaluarse y compararse con los archivos etiquetados existentes. La distribución de estas comprobaciones de calidad a través de distintos sistemas de almacenamiento de datos garantiza una ejecución rápida y eficaz de estas comprobaciones.

### Implementar el análisis de confianza del centro de soporte

La implementación de la solución consta de los siguientes componentes:

1. Kit de herramientas de operaciones de datos de NetApp
2. Configuración de NGC
3. Servidor NVIDIA RIVA
4. Kit de herramientas TAO de NVIDIA
5. Exportar modelos TAO a RIVA

Para realizar la implementación, lleve a cabo los siguientes pasos:

#### Kit de herramientas Data OPS de NetApp: Compatibilidad con el análisis de confianza del centro de soporte

Para utilizar la "[Kit de herramientas de operaciones de datos de NetApp](#)", lleve a cabo los siguientes pasos:

1. Instalación del kit de herramientas de PIP.

```
python3 -m pip install netapp-dataops-traditional
```

## 2. Configurar la gestión de datos

```
netapp_dataops_cli.py config
```

### Configuración de NGC: Respaldo el análisis de confianza del centro

Para configurar "NVIDIA NGC", lleve a cabo los siguientes pasos:

#### 1. Descargar el NGC.

```
wget -O ngccli_linux.zip  
https://ngc.nvidia.com/downloads/ngccli_linux.zip && unzip -o  
ngccli_linux.zip && chmod u+x ngc
```

#### 2. Agregue su directorio actual a la ruta de acceso.

```
echo "export PATH=\"\$PATH:$(pwd)\"" >> ~/.bash_profile && source  
~/.bash_profile
```

#### 3. Debe configurar la CLI de NGC para su uso con el fin de poder ejecutar los comandos. Introduzca el siguiente comando, incluida su clave de API cuando se le solicite.

```
ngc config set
```

Para sistemas operativos que no están basados en Linux, visite [aquí](#).

### NVIDIA RIVA Server: Análisis de opinión del centro de soporte

Para configurar "RIVA DE NVIDIA", lleve a cabo los siguientes pasos:

#### 1. Descargar los archivos RIVA de NGC.

```
ngc registry resource download-version  
nvidia/riva/riva_quickstart:1.4.0-beta
```

#### 2. Inicialice LA configuración DE RIVA (riva\_init.sh).

#### 3. Inicie EL servidor RIVA (riva\_start.sh).

#### 4. Inicie EL cliente RIVA (riva\_start\_client.sh).

5. En EL cliente RIVA, instale la biblioteca de procesamiento de audio ( ["FFMPEG"](#) )

```
apt-get install ffmpeg
```

6. Inicie el ["Jupyter"](#) servidor.

7. Ejecute el portátil de canalización de inferencia DE RIVA.

#### **Kit de herramientas TAO de NVIDIA: Análisis de opinión del centro de soporte**

Para configurar NVIDIA TAO Toolkit, lleve a cabo los siguientes pasos:

1. Prepare y active una ["entorno virtual"](#) Para TAO Toolkit.
2. Instale el ["paquetes requeridos"](#).
3. Tire manualmente de la imagen utilizada durante el entrenamiento y ajuste preciso.

```
docker pull nvcr.io/nvidia/tao/tao-toolkit-pyt:v3.21.08-py3
```

4. Inicie el ["Jupyter"](#) servidor.

5. Ejecute el cuaderno TAO de afinación fina.

#### **Exportar modelos TAO a RIVA: Apoyar el análisis de confianza del centro**

Para usar ["Modelos TAO Toolkit EN RIVA"](#), lleve a cabo los siguientes pasos:

1. Guarde los modelos en el cuaderno TAO de sintonización fina.
2. Copie los modelos TAO entrenados en el directorio del modelo RIVA.
3. Inicie EL servidor RIVA (`riva_start.sh`).

#### **Obstáculos para la implementación**

Estas son algunas cosas que debe tener en cuenta a medida que desarrolla su propia solución:

- El kit de herramientas Data OPS de NetApp se instala primero para garantizar que el sistema de almacenamiento de datos se ejecute de forma óptima.
- NVIDIA NGC debe instalarse antes de cualquier otra cosa porque autentica la descarga de imágenes y modelos.
- RIVA se debe instalar antes que TAO Toolkit. LA instalación DE RIVA configura el demonio docker para extraer imágenes según sea necesario.
- El DGX y el docker deben tener acceso a Internet para descargar los modelos.

#### **Resultados de validación**

Tal como se ha mencionado en la sección anterior, los errores se propagan por la canalización siempre que haya dos o más modelos de aprendizaje automático ejecutándose de forma secuencial. Para esta solución, el sentimiento de la frase es el factor más importante a la hora de medir el nivel de riesgo de las acciones de la

empresa. El modelo de voz a texto, aunque esencial para la canalización, sirve como unidad de preprocesamiento antes de que se puedan predecir los sentimientos. Lo que realmente importa es la diferencia en el sentimiento entre las frases de la verdad en el terreno y las oraciones predichas. Esto sirve como proxy para la velocidad de error de la palabra (WER). La precisión de voz a texto es importante, pero el WER no se utiliza directamente en la métrica de canalización final.

```
PIPELINE_SENTIMENT_METRIC = MEAN(DIFF(GT_sentiment, ASR_sentiment))
```

Estas métricas de sentimiento se pueden calcular para la puntuación F1, la recuperación y la precisión de cada frase. Los resultados se pueden agregar y mostrar dentro de una matriz de confusión, junto con los intervalos de confianza de cada métrica.

La ventaja de usar el aprendizaje de transferencia es un aumento en el rendimiento de los modelos por una fracción de los requisitos de datos, el tiempo de entrenamiento y el coste. Los modelos ajustados también deben compararse con sus versiones de referencia para garantizar que el aprendizaje de la transferencia mejore el rendimiento en lugar de deteriorarlo. En otras palabras, el modelo ajustado debería tener un mejor rendimiento en los datos del centro de soporte que el modelo preentrenado.

Evaluación de la canalización

Caso de prueba	Detalles
Número de prueba	Métrica de sentimiento de canalización
Probar los requisitos previos	Modelos ajustados para modelos de análisis de voz a texto y de sentimiento
Resultado esperado	La métrica de sentimiento del modelo ajustado funciona mejor que el modelo preentrenado original.

Métrica de sentimiento de canalización

1. Calcule la métrica de sentimiento para el modelo de referencia.
2. Calcule la métrica de sentimiento para el modelo ajustado.
3. Calcular la diferencia entre estas métricas.
4. Calcule la media de las diferencias entre todas las frases.

Vídeos y demostraciones

Existen dos portátiles que contienen la canalización de análisis de confianza: ["Support-Center-Model-Transfer-Learning-and-Fine-Tuning.ipynb"](#) y.. ["Support-Center-sentimiento-Analysis-Pipeline.ipynb"](#). Juntos, estos portátiles muestran cómo desarrollar una canalización para la incorporación de datos de centro de soporte y extraer sentimientos de cada frase mediante modelos de aprendizaje profundo de última generación ajustados a los datos del usuario.

## Centro de apoyo - Análisis de confianza Pipeline.ipynb

Este cuaderno contiene la canalización RIVA de inferencia para la incorporación de audio, la conversión a texto y la extracción de sentimientos para su uso en un panel externo. El conjunto de datos se descarga y procesa automáticamente si aún no se ha hecho. La primera sección del bloc de notas es la voz a texto que controla la conversión de archivos de audio a texto. A continuación se muestra la sección Análisis de sentimientos que extrae sentimientos para cada frase de texto y muestra dichos resultados en un formato similar al panel propuesto.



Este portátil debe ejecutarse antes del entrenamiento del modelo y del ajuste preciso porque el conjunto de datos MP3 debe descargarse y convertirse al formato correcto.

## Call Center - Sentiment Analysis Pipeline

This notebook demonstrates how to build a pipeline for sentiment analysis of call center conversations. The goal of this pipeline is to develop sentiment analysis for use within an external dashboard.

This tutorial will guide you through the use of [NVIDIA's RIVA](#) for automatic speech recognition and text classification. This tutorial uses NetApp cloud storage for data storage and a pre-trained RIVA model.

### Channels

These are the channels on which RIVA is hosting models.

- speech: 51051
- voice: 61051

These channels **must** be aligned with `riva_speech_api_port` and `riva_vision_api_port` within `config.sh`

```
In [4]: speech_channel = "localhost:51051"
voice_channel = "localhost:61051"
```

## Speech-To-Text

Automatic Speech Recognition (ASR) takes as input an audio stream or audio buffer and returns one or more text transcripts, along with additional optional metadata. ASR represents a full speech recognition pipeline that is GPU accelerated with optimized performance and accuracy. ASR supports synchronous and streaming recognition modes.

For more information on NVIDIA RIVA's Automatic Speech Recognition, visit [here](#).

## Constants

Use these constants to affect different aspects of this pipeline:

- `DATA_DIR` : base folder where data is stored
- `DATASET_NAME` : name of the call center dataset
- `COMPANY_DATE` : folder name identifying the particular call center conversation

## Centro de soporte - formación de modelos y afinación.ipynb

El entorno virtual del kit de herramientas de TAO debe configurarse antes de ejecutar el cuaderno (consulte la sección del kit de herramientas de TAO en la descripción general de comandos para obtener instrucciones de instalación).

Este cuaderno confía en el kit de herramientas TAO para ajustar los modelos de aprendizaje profundo en los datos de los clientes. Al igual que en el cuaderno anterior, este se separa en dos secciones para los

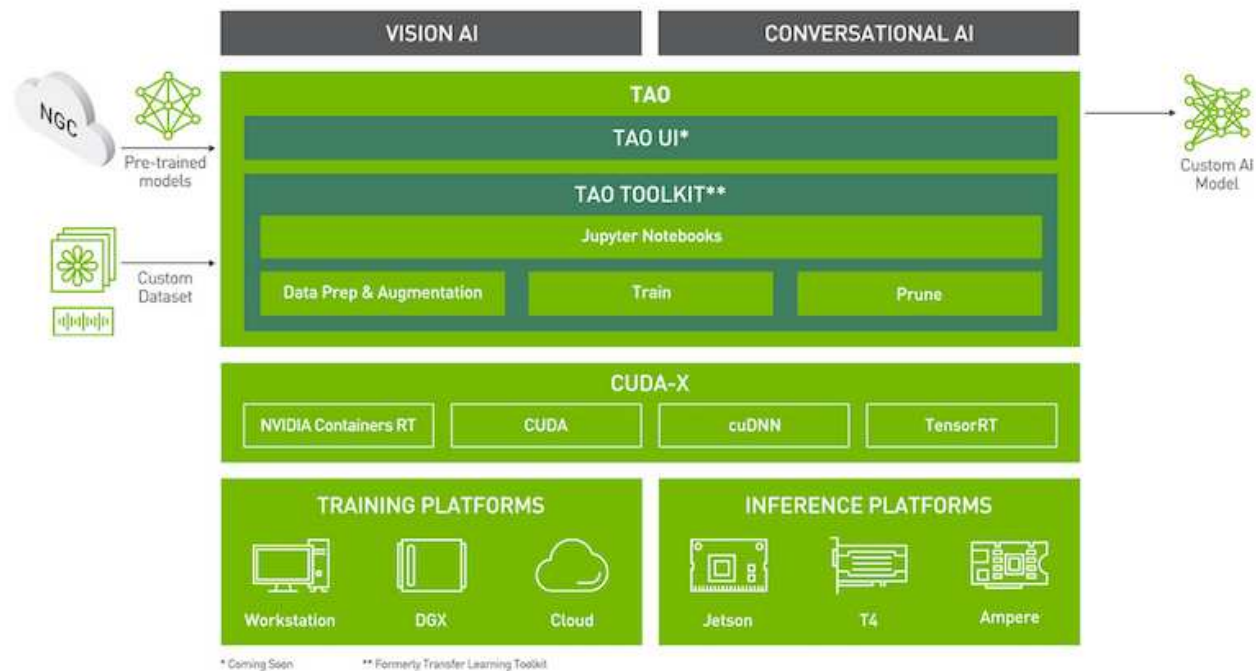
componentes de análisis de opinión y voz. Cada sección abarca el procesamiento de datos, la formación de modelos y el ajuste preciso, la evaluación de los resultados y la exportación de modelos. Por último, existe una sección final para implementar tanto sus modelos ajustados para su uso en RIVA.

## Call Center - Model Transfer Learning and Fine-Tuning

TAO Toolkit is a python based AI toolkit for taking purpose-built pre-trained AI models and customizing them with your own data. Transfer learning extracts learned features from an existing neural network to a new one. Transfer learning is often used when creating a large training dataset is not feasible in order to enhance the base performance of state-of-the-art models.

For this call center solution, the speech-to-text and sentiment analysis models are fine-tuned on call center data to augment the model performance on business specific terminology.

For more information on the TAO Toolkit, please visit [here](#).



### Installing necessary dependencies

For ease of use, please install TAO Toolkit inside a python virtual environment. We recommend performing this step first and then launching the notebook from the virtual environment. Please refer to the README for these instructions.

### Conclusión

A medida que se considera cada vez más la experiencia de los clientes como un campo de batalla competitivo clave, un centro de soporte global aumentado por IA se convierte en un componente fundamental que las empresas de casi todos los sectores no pueden permitirse el lujo de descuidar. La solución propuesta en este informe técnico ha quedado demostrada para respaldar la prestación de estas experiencias excepcionales a los clientes, y el reto consiste ahora en garantizar que las empresas tomen acciones para modernizar su infraestructura de IA y sus flujos de trabajo.

Las mejores implementaciones de IA en servicio al cliente no consisten en sustituir a los agentes humanos. Más bien, la IA puede empoderarlos para crear experiencias de cliente excepcionales mediante análisis de sensibilidad en tiempo real, escalado de disputas y computación afectiva multimodal para detectar señales verbales, no verbales y faciales con las que los modelos de IA integrales pueden hacer recomendaciones a escala y complementar lo que podría carecer un agente humano individual. Asimismo, la IA puede proporcionar una mejor unión entre un cliente concreto con los agentes disponibles actualmente. Gracias a la IA, las empresas pueden extraer un valioso sentimiento de los clientes en cuanto a sus pensamientos e impresiones sobre los productos, servicios y la imagen de la Marca del proveedor.

La solución también se puede usar para construir datos de series temporales para que los agentes de soporte sirvan como métrica de evaluación del rendimiento objetivo. Las encuestas convencionales de satisfacción de clientes a menudo no tienen respuestas suficientes. Al recopilar la opinión de los empleados y los clientes a largo plazo, los empleadores pueden tomar decisiones informadas con respecto al desempeño de los agentes de apoyo.

La combinación de NetApp, SFL Scientific, marcos de orquestación de código abierto y NVIDIA reúne las tecnologías más recientes como servicios gestionados con una gran flexibilidad para acelerar la adopción de tecnología y mejorar el plazo de comercialización de las nuevas aplicaciones de IA/ML. Estos servicios avanzados se ofrecen en las instalaciones que se pueden transportar fácilmente para entornos nativos del cloud así como arquitecturas de puesta en marcha híbrida.

### **Dónde encontrar información adicional**

Si quiere más información sobre el contenido de este documento, consulte los siguientes documentos o sitios web:

- Demostraciones interactivas en 3D

["www.netapp.com/ai"](http://www.netapp.com/ai)

- Conecte directamente con un especialista en IA de NetApp

["https://www.netapp.com/artificial-intelligence/"](https://www.netapp.com/artificial-intelligence/)

- Descripción de la solución NVIDIA base Command Platform con NetApp

<https://www.netapp.com/pdf.html?item=/media/32792-DS-4145-NVIDIA-Base-Command-Platform-with-NetApp.pdf>

- NetApp para IA 10: Infografía buenas razones

["https://www.netapp.com/us/media/netapp-ai-10-good-reasons.pdf"](https://www.netapp.com/us/media/netapp-ai-10-good-reasons.pdf)

- IA en la sanidad: Aprendizaje profundo para identificar las lesiones del COVID-19 en exploraciones pulmonares TC

<https://www.netapp.com/pdf.html?item=/media/31240-WP-7342.pdf>

- AI en la atención sanitaria: Monitorización del uso de la máscara facial en entornos sanitarios documento técnico

<https://www.netapp.com/pdf.html?item=/media/37490-NA-611-Monitoring-face-mask-usage-in-healthcare-settings.pdf>

- IA en la sanidad: Informe técnico de imágenes de diagnóstico



<https://www.netapp.com/pdf.html?item=/media/7395-tr4811.pdf>

- IA para minoristas: Inteligencia artificial de NetApp conversacional con NVIDIA RIVA

["https://docs.netapp.com/us-en/netapp-solutions/ai/cainvidia\\_executive\\_summary.html"](https://docs.netapp.com/us-en/netapp-solutions/ai/cainvidia_executive_summary.html)

- Resumen de la solución ONTAP AI de NetApp

<https://www.netapp.com/pdf.html?item=/media/6736-sb-3939.pdf>

- Resumen de la solución del kit de herramientas de operaciones de datos de NetApp

<https://www.netapp.com/pdf.html?item=/media/21480-SB-4111-1220-NA-Data-Science-Toolkit.pdf>

- Resumen de la solución del plano de control de IA de NetApp

<https://www.netapp.com/pdf.html?item=/media/6737-sb-4055.pdf>

- Libro electrónico "transformando el sector con los datos impulsan la IA"

["https://www.netapp.com/us/media/na-337.pdf"](https://www.netapp.com/us/media/na-337.pdf)

- Resumen de la solución de IA EF-Series de NetApp

<https://www.netapp.com/pdf.html?item=/media/26708-SB-4136-NetApp-AI-E-Series.pdf>

- Resumen de la solución de NetApp AI y Lenovo ThinkSystem para la inferencia de IA

<https://www.netapp.com/pdf.html?item=/media/25316-SB-4129.pdf>

- Resumen de la solución NetApp AI y Lenovo ThinkSystem para IA empresarial Y ML

<https://www.netapp.com/pdf.html?item=/media/25317-SB-4128.pdf>

- NetApp y NVIDIA: Redefiniendo lo que es posible con el vídeo de IA

<https://www.youtube.com/watch?v=38xw65SteUc>

## Formación distribuida en Azure: Predicción de frecuencias mediante clic

### TR-4904: Formación distribuida en Azure - predicción de frecuencias mediante clic

Rick Huang, Verron Martina, Muneer Ahmad, NetApp

El trabajo de un científico de datos debe centrarse en el entrenamiento y el ajuste de los modelos de aprendizaje automático y de inteligencia artificial (IA). Sin embargo, según una investigación de Google, los científicos de datos dedican aproximadamente el 80 % de su tiempo a averiguar cómo hacer que sus modelos funcionen con aplicaciones empresariales y se ejecuten a escala.

Para gestionar proyectos de IA y ML integrales, se necesita una comprensión más amplia de los componentes empresariales. Aunque DevOps ha retomado la definición, la integración y la puesta en marcha, estos tipos de componentes, LAS operaciones ML siguen un flujo similar que incluye proyectos de IA/ML. Para hacerse una idea de lo que puede tocar una canalización de IA/ML integral en la empresa, consulte la siguiente lista de



componentes requeridos:

- Reducida
- Redes
- Oracle
- Sistemas de ficheros
- Contenedores
- Integración continua y canalización de puesta en marcha continua (CI/CD)
- Entorno de desarrollo integrado (IDE)
- Seguridad
- Políticas de acceso a los datos
- Hardware subyacente
- Cloud
- Virtualización
- Conjuntos de herramientas y bibliotecas de ciencia de datos

### **Público objetivo**

El mundo de la ciencia de datos tiene múltiples disciplinas DE TECNOLOGÍA y negocio:

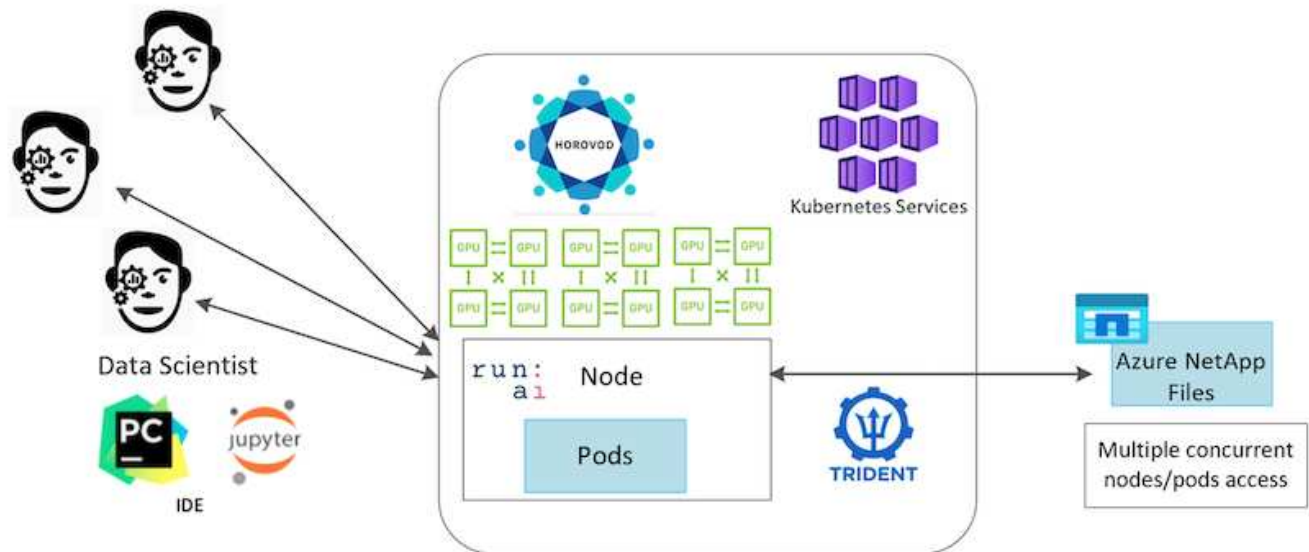
- El científico de datos necesita la flexibilidad para poder usar las herramientas y las bibliotecas de elección.
- El ingeniero de datos necesita saber cómo fluyen los datos y dónde residen.
- Un ingeniero de DevOps necesita herramientas para integrar nuevas aplicaciones de IA/ML en sus canalizaciones de CI/CD.
- Los administradores de cloud y arquitectos tienen que poder configurar y gestionar recursos de Azure.
- Los usuarios empresariales quieren tener acceso a aplicaciones de IA/ML.

En este informe técnico describimos cómo Azure NetApp Files, RAPIDS AI, Dink y Azure ayudan a cada uno de estos roles a aportar valor empresarial.

### **Descripción general de la solución**

Esta solución sigue el ciclo de vida de una aplicación de IA/ML. Empezamos con la labor de los científicos de datos para definir los diferentes pasos necesarios para preparar datos y entrenar modelos. Al aprovechar RAPIDS on Dink, realizamos formación distribuida en el clúster de Azure Kubernetes Service (AKS) para reducir drásticamente el tiempo de entrenamiento en comparación con el método convencional de aprendizaje del kit de ciencias de Python. Para completar el ciclo completo, integramos la canalización con Azure NetApp Files.

Azure NetApp Files ofrece varios niveles de rendimiento. Los clientes pueden comenzar con un nivel estándar y escalar horizontalmente y verticalmente a un nivel de alto rendimiento de forma no disruptiva y sin necesidad de mover datos. Esta funcionalidad permite a los científicos de datos entrenar modelos a escala sin problemas de rendimiento y evitar silos de datos en el clúster, como se muestra en la siguiente figura.



## Información general de la tecnología

Esta página proporciona una descripción general de la tecnología utilizada en esta solución.

### Microsoft y NetApp

Desde mayo de 2019, Microsoft ha ofrecido un servicio de portales nativo de Azure para servicios de archivos NFS y SMB empresariales basados en la tecnología ONTAP de NetApp. Este desarrollo está impulsado por una asociación estratégica entre Microsoft y NetApp, y amplía aún más el alcance de los servicios de datos de ONTAP de primera calidad para Azure.

### Azure NetApp Files

El servicio de Azure NetApp Files es un servicio de almacenamiento de ficheros de alto rendimiento y medida para empresas. Azure NetApp Files es compatible con cualquier tipo de carga de trabajo y está altamente disponible de manera predeterminada. Puede seleccionar los niveles de servicio y rendimiento, y configurar copias Snapshot a través del servicio. Azure NetApp Files es un servicio de primera parte de Azure para migrar y ejecutar las cargas de trabajo de archivo empresarial más exigentes en el cloud, incluidas bases de datos, SAP y aplicaciones de computación de alto rendimiento sin necesidad de modificar el código.

Esta arquitectura de referencia proporciona a las organizaciones DE TI las siguientes ventajas:

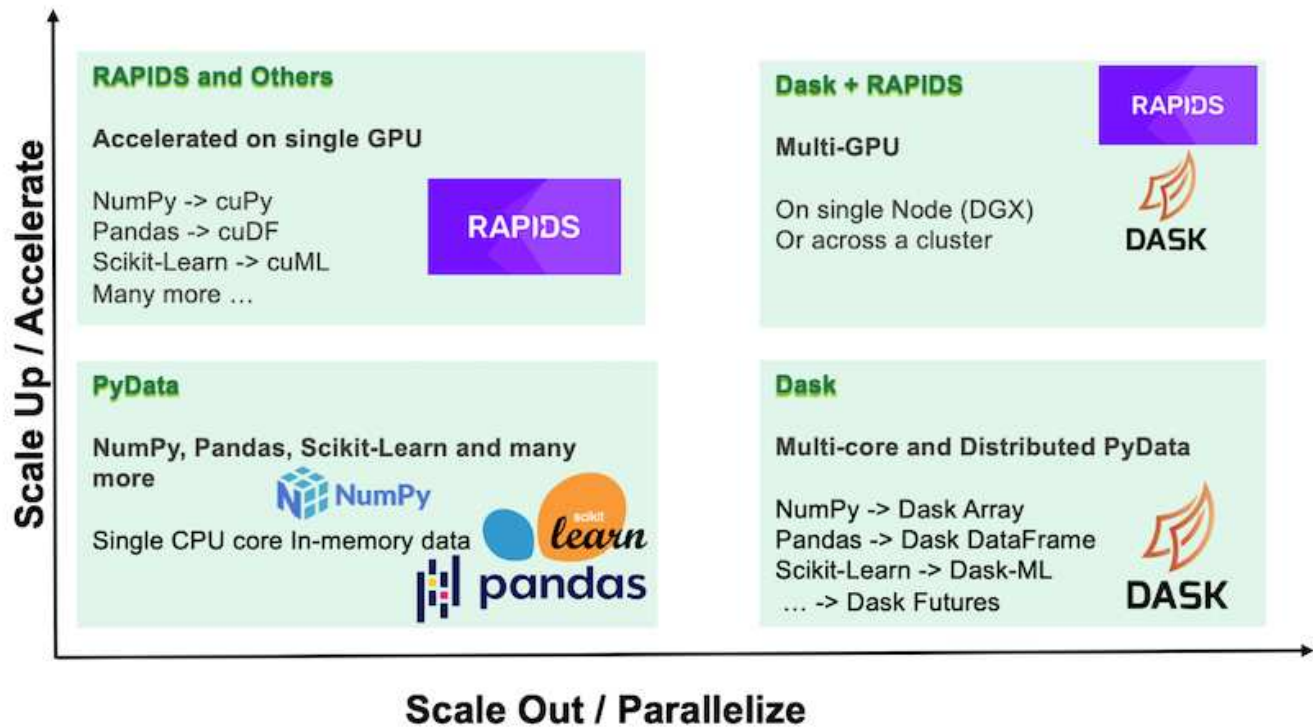
- Elimina las complejidades de diseño
- Permite un escalado independiente de las capacidades de computación y almacenamiento
- Permite a los clientes empezar con poco y escalar sin problemas
- Ofrece una gama de niveles de almacenamiento para distintos niveles de rendimiento y coste

### Información general sobre DASK y NVIDIA RAPIDS

DASK es una herramienta de computación paralela de código abierto que escala bibliotecas Python en varias máquinas y permite un procesamiento más rápido de grandes cantidades de datos. Proporciona una API similar a las bibliotecas de Python convencionales de un solo subproceso, como Pandas, numpy y scikit-Learn. Como resultado, los usuarios nativos de Python no se ven obligados a cambiar mucho su código existente para utilizar recursos en el clúster.

NVIDIA RAPIDS es un paquete de bibliotecas de código abierto que permite ejecutar flujos de trabajo completos de APRENDIZAJE AUTOMÁTICO y análisis de datos en GPU. Junto con DASK, puede escalar con facilidad desde una estación de trabajo con GPU (escalado vertical) a clústeres de varios nodos y varias GPU (escalado horizontal).

Para poner en marcha DASK en un clúster, podría utilizar Kubernetes para la orquestación de recursos. También podría escalar verticalmente o reducir los nodos de trabajo según los requisitos del proceso, lo cual, a su vez, puede ayudarle a optimizar el consumo de recursos de clúster, como se muestra en la siguiente figura.



**Requisitos de software**

En la siguiente tabla se enumeran los requisitos de software necesarios para esta solución.

De NetApp	Versión
Servicio Azure Kubernetes	1.18.14
Imagen de contenedor DE RAPIDS y DASK	Repositorio: "Rapidsai/rapidsai" etiqueta: 0.17-cuda11.0-Runtime-ubuntu18.04
Trident de NetApp	20.01.1
Timón	3.0.0

**Requisitos de recursos cloud**

Esta página describe la configuración de los recursos de cloud para Azure NetApp Files.

## Configure Azure NetApp Files

Configure Azure NetApp Files como se describe en ["Inicio rápido: Configure Azure NetApp Files y cree un volumen NFS"](#).

Puede pasar más allá de la sección "Crear volumen NFS para Azure NetApp Files" porque va a crear volúmenes a través de Trident. Antes de continuar, realice los siguientes pasos:

1. Regístrese para Azure NetApp Files y el proveedor de recursos de NetApp (a través de la shell de Azure) (["enlace"](#)).
2. Crear una cuenta en Azure NetApp Files (["enlace"](#)).
3. Configurar un pool de capacidad (un estándar o Premium de 4 TB como mínimo, según sus necesidades) (["enlace"](#)).la siguiente tabla enumera los requisitos de configuración de red para configurar en la nube. El clúster de DASK y Azure NetApp Files deben estar en la misma red virtual de Azure (vnet) o en una vnet con conexión entre iguales.

Recursos	Tipo/versión
Servicio Azure Kubernetes	1.18.14
Nodo de agente	3x Standard_DS2_v2
Nodo GPU	3x Standard_NC63_v3
Azure NetApp Files	Pool de capacidad estándar
Capacidad en TB	4

## Resumen de casos de uso de predicción de velocidad mediante clic

Este caso de uso se basa en el público disponible ["Terabyte haga clic en registros"](#) conjunto de datos de ["Laboratorio Criteo AI"](#). Con los recientes avances en las plataformas Y aplicaciones DE ML, ahora se presta mucha atención al aprendizaje a escala. La tasa de clics (CTR) se define como el número medio de clics-throughs por cien impresiones de anuncios en línea (expresado como porcentaje). Se ha adoptado ampliamente como métrica clave en diversos mercados verticales del sector y casos de uso, incluidos el marketing digital, el comercio minorista, el comercio electrónico y los proveedores de servicios. Algunos ejemplos de uso de CTR como una métrica importante para el tráfico potencial de clientes son los siguientes:

- **Marketing digital:** in ["Google Analytics"](#), CTR se puede usar para medir cómo se están realizando las palabras clave, los anuncios y los listados libres de un anunciante o comerciante. Un CTR alto es una buena indicación de que los usuarios encuentran sus anuncios y listados útiles y relevantes. CTR también contribuye a la CTR esperada de su palabra clave, que es un componente de ["Clasificación de anuncios"](#).
- **Comercio electrónico:** además de apalancar ["Google Analytics"](#), hay por lo menos algunas estadísticas de visitantes en un fondo de comercio electrónico. Aunque estas estadísticas pueden no parecer útiles a primera vista, suelen ser fáciles de leer y pueden ser más precisas que otras informaciones. Los conjuntos de datos de primera parte compuestos por estas estadísticas son de propiedad y, por lo tanto, los más relevantes para los vendedores, compradores y plataformas de comercio electrónico. Estos conjuntos de datos se pueden utilizar para establecer pruebas de rendimiento, comparar los resultados con el año pasado y el pasado construyendo una serie temporal para un análisis más profundo.
- **Retail:** los minoristas tradicionales pueden correlacionar el número de visitantes y el número de clientes con el CTR. El número de clientes se puede ver desde su historial de puntos de venta. El CTR de los sitios

web de los minoristas o del tráfico de anuncios puede resultar en las ventas mencionadas. Los programas de fidelidad son otro caso de uso, ya que los clientes redirigidos de anuncios en línea u otros sitios web podrían unirse para obtener recompensas. Los minoristas pueden conseguir clientes a través de programas de fidelidad y comportamientos récord de historias de ventas para crear un sistema de recomendaciones que no solo predice las conductas de compra de los consumidores en diferentes categorías, sino que también personaliza las cupones y reduce la pérdida de clientes.

- **Proveedores de servicios:** las empresas de telecomunicaciones y los proveedores de servicios de Internet tienen una abundancia de datos de telemetría de usuarios de primera parte para casos de uso de IA, ML y analítica perspicaz. Por ejemplo, las telecomunicaciones pueden aprovechar los registros diarios de historial de dominios de navegación web de sus suscriptores móviles para ajustar los modelos existentes y producir una segmentación de público actualizada, predecir el comportamiento de los clientes y colaborar con los anunciantes para colocar anuncios en tiempo real para una mejor experiencia en línea. En este flujo de trabajo de marketing basado en datos, CTR es una métrica importante para reflejar las conversiones.

En el contexto del marketing digital, "[Criteo Terabyte haga clic en registros](#)" Son ahora el conjunto de datos de referencia a la hora de evaluar la escalabilidad de las plataformas Y algoritmos ML. Al predecir la tarifa de clic, un anunciante puede seleccionar a los visitantes que tienen más probabilidades de responder a los anuncios, analizar su historial de navegación y mostrar los anuncios más relevantes basados en los intereses del usuario.

La solución proporcionada en este informe técnico destaca las siguientes ventajas:

- Las ventajas de Azure NetApp Files en formación distribuida o a gran escala
- RÁPIDO procesamiento de datos habilitado para CUDA (cuDF, cúpula, etc.) y algoritmos ML (cuML)
- El marco informático paralelo de DASK para la formación distribuida

Un flujo de trabajo integral basado en RAPIDS AI y Azure NetApp Files demuestra la drástica mejora del tiempo de entrenamiento de los modelos de bosques aleatorios en dos órdenes de magnitud. Esta mejora es significativa en comparación con el enfoque convencional de pandas al tratar los registros de clic del mundo real con 45GB de datos tabulares estructurados (en promedio) cada día. Esto equivale a un DataFrame que contiene aproximadamente veinte mil millones de filas. Demostraremos la configuración del entorno de clúster, la instalación de marcos y bibliotecas, la carga y el procesamiento de datos, la formación convencional frente a la distribuida, la visualización y supervisión, y compararemos los resultados fundamentales de tiempo de ejecución completo en este informe técnico.

## Configuración

### Instalar y configurar el clúster AKS

Para instalar y configurar el clúster AKS, consulte la página web "[Cree un clúster de AKS](#)" y, a continuación, realice los siguientes pasos:

1. Al seleccionar el tipo de nodo (nodos System [CPU] o worker [GPU]), seleccione lo siguiente:
  - a. Los nodos del sistema principal deben ser DS2v2 estándar (agentpool tres nodos predeterminados).
  - b. A continuación, agregue el pool Standard\_NC66s\_v3 del nodo de trabajo (tres nodos como mínimo) para el grupo de usuarios (para nodos GPU) denominado gpupool.

+ Add node pool		Delete		
Name	Mode	OS type	Node count	Node size
<input type="checkbox"/> agentpool	System	Linux	3	Standard_DS2_v2
<input type="checkbox"/> gpupool	User	Linux	3	Standard_NC6s_v

- La puesta en marcha tarda entre 5 y 10 minutos. Una vez finalizado, haga clic en Connect to Cluster.
- Para conectarse al clúster AKS recién creado, instale lo siguiente desde su entorno local (portátil/pc):
  - La herramienta de línea de comandos de Kubernetes que utiliza ["Instrucciones proporcionadas para su SO específico"](#)
  - El CLI de Azure tal como se describe en el documento, ["Instale la CLI de Azure"](#)
- Para acceder al clúster AKS desde el terminal, introduzca `az login` e introduzca las credenciales.
- Ejecute los dos comandos siguientes:

```
az account set --subscription xxxxxxxx-xxxx-xxxx-xxxx-xxxxxxxxxxxxxxxx
aks get-credentials --resource-group resourcegroup --name aksclustername
```

- Introduzca Azure CLI: `kubectl get nodes`.
- Si los seis nodos están en funcionamiento, como se muestra en el siguiente ejemplo, su clúster AKS estará listo y conectado a su entorno local

```
verronmartina@verron-mac-0 ~ % kubectl get nodes
NAME                                STATUS    ROLES    AGE   VERSION
aks-agentpool-34613062-vmss000000  Ready    agent    22m   v1.18.14
aks-agentpool-34613062-vmss000001  Ready    agent    22m   v1.18.14
aks-agentpool-34613062-vmss000002  Ready    agent    22m   v1.18.14
aks-gpupool-34613062-vmss000000     Ready    agent    20m   v1.18.14
aks-gpupool-34613062-vmss000001     Ready    agent    20m   v1.18.14
aks-gpupool-34613062-vmss000002     Ready    agent    20m   v1.18.14
verronmartina@verron-mac-0 ~ %
```

### Cree una subred delegada para Azure NetApp Files

Para crear una subred delegada para Azure NetApp Files, lleve a cabo los siguientes pasos:

- Acceda a Virtual Networks dentro del portal de Azure. Busque la red virtual que acaba de crear. Debe tener un prefijo como `aks-vnet`.
- Haga clic en el nombre de la vnet.

Microsoft Azure

Search resources, services, and docs (G+ /)

Dashboard >


## Virtual networks

seanlucelive (Default Directory)

+ Add Manage view Refresh Export to CSV Open query Assign tags Feedback

Filter by name... Subscription == AzureSub01 Resource group == all Location == all Add filter

Showing 1 to 5 of 5 records. No grouping List view

Name ↑↓	Resource group ↑↓	Location ↑↓	Subscription ↑↓
 aks-vnet-22885919	MC_sluce.rg_TridentDemo_eastus2	East US 2	AzureSub01

3. Haga clic en subredes y, a continuación, en +Subnet en la barra de herramientas superior.

Microsoft Azure

Search resources, services, and docs (G+ /)

Dashboard > Virtual networks > aks-vnet-22885919

## aks-vnet-22885919 | Subnets

Virtual network

Search (Ctrl+ /) + Subnet Gateway subnet Refresh Manage users Delete

Overview Activity log Access control (IAM) Tags Diagnose and solve problems Settings Address space Connected devices Subnets

Search subnets

Name ↑↓	IPv4 ↑↓	IPv6 (many availab... ↑↓	Delegated to ↑↓	Security group ↑↓
aks-subnet	10.240.0.0/16 (65530 av...	-	-	aks-agentpool-2288591...

4. Proporcione la subred con un nombre como ANF.sn Y, en el encabezado Delegación de subred, seleccione Microsoft.Netapp/volumes. No cambie nada más. Haga clic en Aceptar.



## Add subnet



Name \*

ANF.sn



Subnet address range \* ⓘ

10.0.0.0/24

10.0.0.0 - 10.0.0.255 (251 + 5 Azure reserved addresses)

☐

Add IPv6 address space ⓘ

NAT gateway ⓘ

None



Network security group

None



Route table

None



### SERVICE ENDPOINTS

Create service endpoint policies to allow traffic to specific azure resources from your virtual network over service endpoints. [Learn more](#)

Services ⓘ

0 selected



### SUBNET DELEGATION

Delegate subnet to a service ⓘ

Microsoft.Netapp/volumes



OK

Cancel

Los volúmenes Azure NetApp Files se asignan al clúster de aplicaciones y se consumen como reclamaciones de volúmenes persistentes (RVP) en Kubernetes. A su vez, este proceso le proporciona la flexibilidad para asignarlos a diferentes servicios, tales como portátiles Jupyter, funciones sin servidor, etc.

Los usuarios de servicios pueden consumir almacenamiento desde la plataforma de muchas maneras. Como este informe técnico trata sobre NFSS, los principales beneficios de Azure NetApp Files son:

- Brindar a los usuarios la capacidad de usar copias de Snapshot.
- Permitir a los usuarios almacenar grandes cantidades de datos en volúmenes de Azure NetApp Files.
- Utilización de las ventajas en el rendimiento de los volúmenes Azure NetApp Files cuando se ejecutan sus modelos en conjuntos de archivos de gran tamaño.



Para conectar el AKS vnet al vnet de Azure NetApp Files, lleve a cabo los siguientes pasos:

1. Introduzca redes virtuales en el campo de búsqueda.
2. Seleccione `vnet aks-vnet-name`. Haga clic en él e introduzca los peerings en el campo de búsqueda.
3. Haga clic en +Agregar.
4. Introduzca los siguientes descriptores:
  - a. El nombre del enlace de relación de paridad es `aks-vnet-name_to_anf`.
  - b. SubscriptionId y Azure NetApp Files vnet como partner de vnet peering.
  - c. Deje todas las secciones que no sean asteriscos con los valores predeterminados.
5. Haga clic en Añadir.

Para obtener más información, consulte ["Crear, cambiar o eliminar una conexión de red virtual entre iguales"](#).

### Instale Trident

Para instalar Trident con Helm, lleve a cabo los siguientes pasos:

1. Instale Helm (para obtener instrucciones de instalación, visite ["origen"](#)).
2. Descargue y extraiga el instalador de Trident 20.01.1.

```
$wget  
$tar -xf trident-installer-21.01.1.tar.gz
```

3. Cambie el directorio a `trident-installer`.

```
$cd trident-installer
```

4. Copiar `tridentctl` a un directorio del sistema `$PATH`.

```
$sudo cp ./tridentctl /usr/local/bin
```

5. Instale Trident en el clúster Kubernetes (K8s) con Helm ( ["origen"](#)):

- a. Cambie el directorio a `helm` directorio.

```
$cd helm
```

- b. Instale Trident.

```
$helm install trident trident-operator-21.01.1.tgz --namespace  
trident --create-namespace
```

c. Comprobar el estado de los pods de Trident.

```
$kubectl -n trident get pods
```

Si todos los pods están ya en funcionamiento, se instala Trident y se puede avanzar.

6. Configurar el back-end de Azure NetApp Files y la clase de almacenamiento para AKS.

a. Cree un principio de Azure Service.

El principal del servicio es cómo Trident se comunica con Azure para manipular sus recursos de Azure NetApp Files.

```
$az ad sp create-for-rbac --name ""
```

El resultado debería ser como el ejemplo siguiente:

```
{  
  "appId": "xxxxx-xxxx-xxxx-xxxx-xxxxxxxxxxxx",  
  "displayName": "netapptrident",  
  "name": "",  
  "password": "xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx",  
  "tenant": "xxxxxxxx-xxxx-xxxx-xxxx-xxxxxxxxxxxx"  
}
```

7. Cree un archivo del back-end json de Trident, nombre de ejemplo `anf-backend.json`.

8. Con el editor de texto preferido, complete los siguientes campos dentro del `anf-backend.json` archivo:

```
{
  "version": 1,
  "storageDriverName": "azure-netapp-files",
  "subscriptionID": "fakec765-4774-fake-ae98-a721add4fake",
  "tenantID": "fakef836-edc1-fake-bff9-b2d865eefake",
  "clientID": "fake0f63-bf8e-fake-8076-8de91e57fake",
  "clientSecret": "SECRET",
  "location": "westeurope",
  "serviceLevel": "Standard",
  "virtualNetwork": "anf-vnet",
  "subnet": "default",
  "nfsMountOptions": "vers=3,proto=tcp",
  "limitVolumeSize": "500Gi",
  "defaults": {
    "exportRule": "0.0.0.0/0",
    "size": "200Gi"
  }
}
```

9. Sustituya los siguientes campos:

- `subscriptionID`. Su ID de suscripción de Azure.
- `tenantID`. Su ID de inquilino de Azure de la salida de `az ad sp` en el paso anterior.
- `clientID`. Su AppID desde la salida de `az ad sp` en el paso anterior.
- `clientSecret`. Su contraseña de la salida de `az ad sp` en el paso anterior.

10. Indique a Trident que cree el back-end de Azure NetApp Files en la `trident` espacio de nombres con `anf-backend.json` como archivo de configuración:

```
$tridentctl create backend -f anf-backend.json -n trident
```

NAME	STORAGE DRIVER	UUID	STATE	VOLUMES
azurenetafiles_86181	azure-netapp-files	2ca85462-59ac-4946-be05-c03f5575a2ad	online	0

11. Cree una clase de almacenamiento. Los usuarios de Kubernetes aprovisionan volúmenes con RVP que especifican una clase de almacenamiento por nombre. Indique a K8S que cree una clase de almacenamiento `azurenetafiles` Que hace referencia al back-end de Trident creado en el paso anterior.
12. Cree una AYLMA (`anf-storage-class.yaml`) archivo para la clase de almacenamiento y copiar.

```

apiVersion: storage.k8s.io/v1
kind: StorageClass
metadata:
  name: azurenetappfiles
provisioner: netapp.io/trident
parameters:
  backendType: "azure-netapp-files"
$kubectl create -f anf-storage-class.yaml

```

13. Compruebe que la clase de almacenamiento se ha creado.

```
kubectl get sc azurenetappfiles
```

NAME	PROVISIONER	RECLAIMPOLICY	VOLUMEBINDINGMODE	ALLOWVOLUMEEXPANSION	AGE
azurenetappfiles	csi.trident.netapp.io	Delete	Immediate	false	98s

### Configurar Dask con LA implementación DE RAPIDS en AKS con Helm

Para configurar el despliegue de Dask con RAPIDS en AKS con Helm, lleve a cabo los siguientes pasos:

1. Cree un espacio de nombres para instalar DASK con RAPIDS.

```
kubectl create namespace rapids-dask
```

2. Crear una RVP para almacenar el conjunto de datos con velocidad de clic:

a. Guarde el siguiente contenido de YAML en un archivo para crear un PVC.

```

kind: PersistentVolumeClaim
apiVersion: v1
metadata:
  name: pvc-criteo-data
spec:
  accessModes:
    - ReadWriteMany
  resources:
    requests:
      storage: 1000Gi
  storageClassName: azurenetappfiles

```

b. Aplique el archivo YAML al clúster de Kubernetes.

```
kubectl -n rapids-dask apply -f <your yaml file>
```

3. Clone el rapidsai git repositorio ( "<https://github.com/rapidsai/helm-chart>").

```
git clone https://github.com/rapidsai/helm-chart helm-chart
```

4. Modificar `values.yaml` E incluya el PVC creado anteriormente para los trabajadores y el espacio de trabajo Jupyter.

- a. Vaya a la rapidsai directorio del repositorio.

```
cd helm-chart/rapidsai
```

- b. Actualice el `values.yaml` Coloque el archivo y monte el volumen con la RVP.

```
dask:
  ...
worker:
  name: worker
  ...
mounts:
  volumes:
    - name: data
      persistentVolumeClaim:
        claimName: pvc-criteo-data
  volumeMounts:
    - name: data
      mountPath: /data
  ...
jupyter:
  name: jupyter
  ...
mounts:
  volumes:
    - name: data
      persistentVolumeClaim:
        claimName: pvc-criteo-data
  volumeMounts:
    - name: data
      mountPath: /data
  ...
```

5. Vaya al directorio principal del repositorio e implemente Dask con tres nodos de trabajo en AKS utilizando

Helm.

```
cd ..
helm dep update rapidsai
helm install rapids-dask --namespace rapids-dask rapidsai
```

Niveles de rendimiento de Azure NetApp Files

Es posible cambiar el nivel de servicio de un volumen existente si se mueve el volumen a otro pool de capacidad que utiliza el nivel de servicio que se desea para el volumen. Esta solución permite a los clientes comenzar con un conjunto de datos pequeño y un número reducido de GPU en el nivel estándar y escalar horizontalmente o verticalmente hasta el nivel Premium a medida que aumenta la cantidad de datos y las GPU. El nivel Premium ofrece cuatro veces el rendimiento por terabyte como nivel estándar, así como la escalabilidad vertical se realiza sin tener que mover datos para cambiar el nivel de servicio de un volumen.

Cambie dinámicamente el nivel de servicio de un volumen

Para cambiar de forma dinámica el nivel de servicio de un volumen, complete los pasos siguientes:

- 1. En la página Volumes, haga clic con el botón derecho en el volumen cuyo nivel de servicio desea cambiar. Seleccione Cambiar pool.

NFSv3	10.28.254.4:/norootfor	Standard	pool0	...
NFSv4.1	NAS-735a.docs.lab:/fo	Premium		...
NFSv4.1	NAS-735a.docs.lab:/krt	Premium		...
NFSv3	10.28.254.4:/moveme0	Premium		...
NFSv3	10.28.254.4:/placeholder	Premium		...

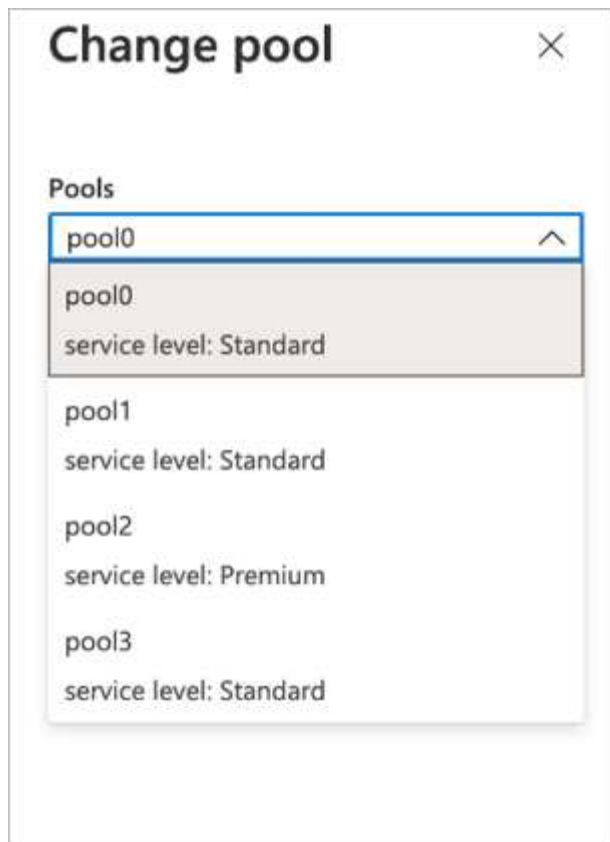
Resize

Edit

Change pool

Delete

- 2. En la ventana Cambiar pool, seleccione el pool de capacidad al que desea mover el volumen.



3. Haga clic en Aceptar.

### Automatice el cambio en el nivel de rendimiento

Existen las siguientes opciones para automatizar los cambios en el nivel de rendimiento:

- El cambio de nivel de servicio dinámico sigue en la vista previa pública en este momento y no está activado de forma predeterminada. Para habilitar esta función en la suscripción a Azure, consulte esta [documentación sobre cómo "Cambie dinámicamente el nivel de servicio de un volumen"](#).
- Se proporcionan comandos de cambio de pool de volúmenes de la interfaz de línea de comandos de Azure en ["documentación de cambio de pool de volúmenes"](#) y en el ejemplo siguiente:

```
az netappfiles volume pool-change -g mygroup --account-name myaccname  
--pool-name mypoolname --name myvolname --new-pool-resource-id  
mynewresourceid
```

- PowerShell: El ["Set-AzNetAppFilesVolumePool"](#) Cambia el pool de un volumen Azure NetApp Files y se muestra en el ejemplo siguiente:

```
Set-AzNetAppFilesVolumePool
-ResourceGroupName "MyRG"
-AccountName "MyAnfAccount"
-PoolName "MyAnfPool"
-Name "MyAnfVolume"
-NewPoolResourceId 7d6e4069-6c78-6c61-7bf6-c60968e45fbf
```

## Haga clic en el procesamiento de datos de predicción de velocidad y el entrenamiento de modelos

### Bibliotecas para el procesamiento de datos y el entrenamiento de modelos

En la tabla siguiente se enumeran las bibliotecas y los marcos que se utilizaron para generar esta tarea. Todos estos componentes se han integrado completamente con los controles de seguridad y acceso basados en roles de Azure.

Bibliotecas/marco de trabajo	Descripción
CuML DASK	Para QUE EL ML funcione en la GPU, el <a href="#">"Biblioteca de cuML"</a> Ofrece acceso al paquete cuML DE RAPIDS con DASK. RAPIDS cuML implementa algoritmos DE ML más conocidos, como los métodos de clustering, reducción de dimensiones y regresión, con implementaciones basadas en GPU de alto rendimiento que ofrecen una velocidad de hasta 100 veces superior a los métodos basados en CPU.
DASK cuDF	CuDF incluye varias otras funciones que admiten la extracción, transformación y carga (ETL) acelerada por GPU, como la subconfiguración de datos, transformaciones, codificación en caliente, etc. El equipo DE RAPIDS mantiene un <a href="#">"biblioteca dask-cudf"</a> Eso incluye métodos auxiliares para usar DASK y cuDF.
Formación en Scikit	Scikit-Learn proporciona docenas de algoritmos y modelos de aprendizaje automático integrados, llamados estimadores. Cada uno <a href="#">"estimator"</a> se puede ajustar a algunos datos mediante su <a href="#">"encajar"</a> método.

Hemos utilizado dos cuadernos para construir los gasoductos ML para su comparación; uno es el método convencional de curscikit-aprender de pandas, y el otro es el entrenamiento distribuido con RAPIDS y Dink. Cada portátil se puede probar individualmente para ver el rendimiento en términos de tiempo y escala. Cubrimos cada bloc de notas individualmente para demostrar las ventajas de la formación distribuida con RAPIDS y Dink.

**Cargue Criteo haga clic en el día 15 de los registros en pandas y entrena un cikit-aprende el modelo de bosque aleatorio**

En esta sección se describe cómo utilizamos Pandas y DASK DataFrames para cargar datos Click Logs del conjunto de datos Criteo Terabyte. El caso de uso es relevante en la publicidad digital para intercambios de anuncios para crear perfiles de usuarios al



predecir si se hará clic en anuncios o si el intercambio no está utilizando un modelo exacto en una canalización automatizada.

Se cargaron los datos del día 15 desde el conjunto de datos Click Logs, sumando 45 GB. Ejecutar la siguiente celda en el portátil Jupyter CTR-PandasRF-collated.ipynb Crea un DataFrame de pandas que contiene los primeros 50 millones de filas y genera un modelo de bosque aleatorio cikit-aprender.

```
%%time
import pandas as pd
import numpy as np
header = ['col'+str(i) for i in range (1,41)] #note that according to
criteo, the first column in the dataset is Click Through (CT). Consist of
40 columns
first_row_taken = 50_000_000 # use this in pd.read_csv() if your compute
resource is limited.
# total number of rows in day15 is 20B
# take 50M rows
"""
Read data & display the following metrics:
1. Total number of rows per day
2. df loading time in the cluster
3. Train a random forest model
"""
df = pd.read_csv(file, nrows=first_row_taken, delimiter='\t',
names=header)
# take numerical columns
df_sliced = df.iloc[:, 0:14]
# split data into training and Y
Y = df_sliced.pop('col1') # first column is binary (click or not)
# change df_sliced data types & fillna
df_sliced = df_sliced.astype(np.float32).fillna(0)
from sklearn.ensemble import RandomForestClassifier
# Random Forest building parameters
# n_streams = 8 # optimization
max_depth = 10
n_bins = 16
n_trees = 10
rf_model = RandomForestClassifier(max_depth=max_depth,
n_estimators=n_trees)
rf_model.fit(df_sliced, Y)
```

Para realizar la predicción utilizando un modelo de bosque aleatorio entrenado, ejecute el siguiente párrafo en este cuaderno. Tomamos las últimas filas de un millón del día 15 como conjunto de pruebas para evitar cualquier duplicación. La celda también calcula la precisión de la predicción, definida como el porcentaje de ocurrencias que el modelo predice con precisión si un usuario hace clic o no en un anuncio. Para revisar cualquier componente desconocido en este cuaderno, consulte ["documentación oficial de scikit-aprender"](#).

```
# testing data, last 1M rows in day15
test_file = '/data/day_15_test'
with open(test_file) as g:
    print(g.readline())

# dataframe processing for test data
test_df = pd.read_csv(test_file, delimiter='\t', names=header)
test_df_sliced = test_df.iloc[:, 0:14]
test_Y = test_df_sliced.pop('coll')
test_df_sliced = test_df_sliced.astype(np.float32).fillna(0)
# prediction & calculating error
pred_df = rf_model.predict(test_df_sliced)
from sklearn import metrics
# Model Accuracy
print("Accuracy:", metrics.accuracy_score(test_Y, pred_df))
```

### Cargar día 15 en DASK y entrenar un modelo de bosque aleatorio DASK cuML

De una manera similar a la sección anterior, cargue Criteo Click Logs Day 15 en Pandas y entrena un cikit-aprende el modelo de bosque aleatorio. En este ejemplo, realizamos la carga de DataFrame con DASK cuDF y entrenamos un modelo de bosque aleatorio en DASK cuML. Hemos comparado las diferencias en el tiempo de formación y el escalado en la sección ["Comparación del tiempo de formación"](#).

#### criteo\_dask\_RF.ipynb

Este portátil importa numpy, cuml, y lo necesario dask bibliotecas, como se muestra en el siguiente ejemplo:

```
import cuml
from dask.distributed import Client, progress, wait
import dask_cudf
import numpy as np
import cudf
from cuml.dask.ensemble import RandomForestClassifier as cumlDaskRF
from cuml.dask.common import utils as dask_utils
```

Inicie cliente DASK().

```
client = Client()
```

Si su clúster está configurado correctamente, puede ver el estado de los nodos de trabajo.

```

client
workers = client.has_what().keys()
n_workers = len(workers)
n_streams = 8 # Performance optimization

```

En nuestro clúster AKS, se muestra el siguiente estado:

Client	Cluster
<b>Scheduler:</b> tcp://rapidsai-scheduler:8786	<b>Workers:</b> 3
<b>Dashboard:</b> <a href="/proxy/rapidsai-scheduler:8787/status">/proxy/rapidsai-scheduler:8787/status</a>	<b>Cores:</b> 3
	<b>Memory:</b> 354.55 GB

Tenga en cuenta que DASK emplea el paradigma de ejecución lenta: En lugar de ejecutar el código de procesamiento al instante, DASK crea en su lugar un gráfico cíclico dirigido (DAG) de ejecución. DAG contiene un conjunto de tareas y sus interacciones que cada trabajador necesita ejecutar. Este diseño significa que las tareas no se ejecutan hasta que el usuario le indique a DASK que las ejecute de una forma u otra. Con DASK tiene tres opciones principales:

- **Call `comput()` en un `DataFrame`.** esta llamada procesa todas las particiones y, a continuación, devuelve los resultados al planificador para la agregación final y conversión a `cuDF DataFrame`. Esta opción debe usarse con moderación y sólo en resultados muy reducidos a menos que el nodo del programador se quede sin memoria.
- **Call `persistent()` en un `DataFrame`.** esta llamada ejecuta el gráfico, pero, en lugar de devolver los resultados al nodo del planificador, los mantiene en la memoria a través del clúster para que el usuario pueda reutilizar estos resultados intermedios en la canalización sin necesidad de volver a ejecutar el mismo procesamiento.
- **Call `head()` en un `DataFrame`.** al igual que con `cuDF`, esta llamada devuelve 10 registros al nodo del planificador. Esta opción se puede utilizar para comprobar rápidamente si el `DataFrame` contiene el formato de salida deseado o si los propios registros tienen sentido, en función del procesamiento y cálculo.

Por lo tanto, a menos que el usuario llama a cualquiera de estas acciones, los trabajadores se sientan inactivos esperando que el programador inicie el procesamiento. Este paradigma de ejecución perezosa es común en marcos informáticos modernos en paralelo y distribuidos como Apache Spark.

En el siguiente párrafo se entrena un modelo de bosque aleatorio mediante el uso de DASK `cuML` para computación acelerada por GPU distribuida y se calcula la precisión de predicción del modelo.

```

Adsf
# Random Forest building parameters
n_streams = 8 # optimization
max_depth = 10
n_bins = 16
n_trees = 10
cuml_model = cumlDaskRF(max_depth=max_depth, n_estimators=n_trees,
n_bins=n_bins, n_streams=n_streams, verbose=True, client=client)
cuml_model.fit(gdf_sliced_small, Y)
# Model prediction
pred_df = cuml_model.predict(gdf_test)
# calculate accuracy
cu_score = cuml.metrics.accuracy_score( test_y, pred_df )

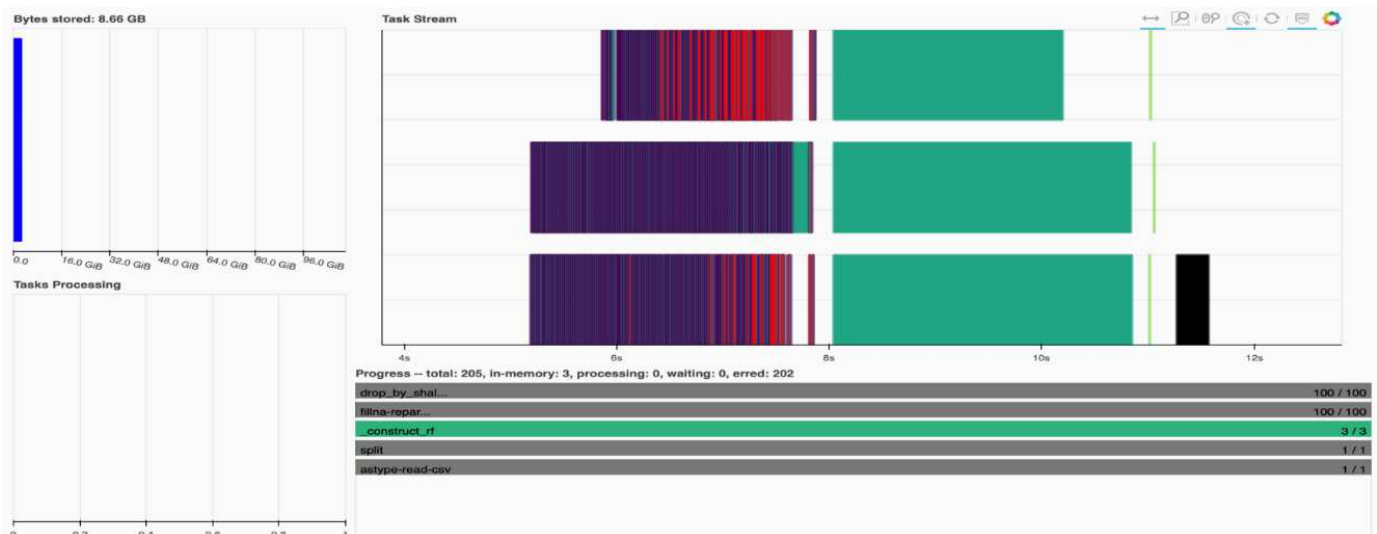
```

### Supervisión de tarea mediante el panel de control de flujos de tareas nativo

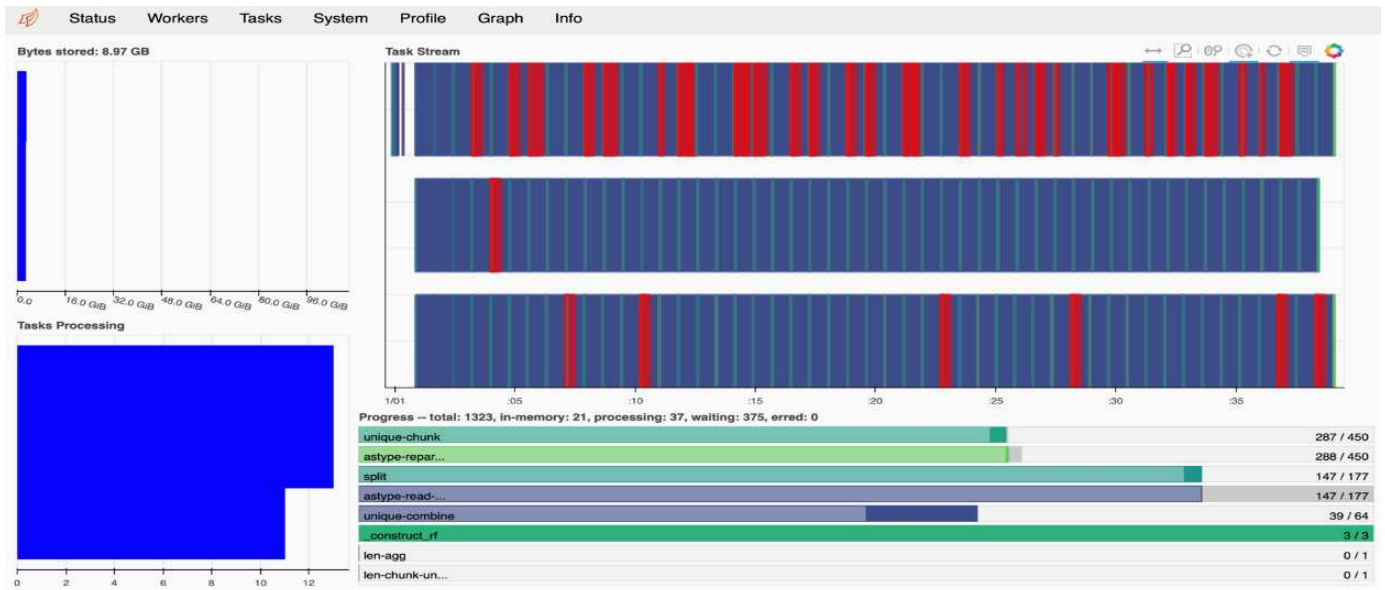
La "**Planificador distribuido DASK**" proporciona comentarios en directo de dos formas:

- Un panel interactivo que contiene muchos trazados y tablas con información en directo
- Una barra de progreso adecuada para uso interactivo en consolas o portátiles

En nuestro caso, la siguiente figura muestra cómo puede supervisar el progreso de la tarea, incluidos los bytes almacenados, el flujo de tareas con un desglose detallado del número de flujos y el progreso por los nombres de tareas con las funciones asociadas ejecutadas. En nuestro caso, debido a que tenemos tres nodos de trabajo, hay tres partes principales del flujo y los códigos de color denotan diferentes tareas dentro de cada flujo.



Tiene la opción de analizar tareas individuales y examinar el tiempo de ejecución en milisegundos o identificar cualquier obstáculo o impedimento. Por ejemplo, la siguiente figura muestra los flujos de tareas para la etapa de ajuste del modelo de bosque aleatorio. Se están ejecutando muchas más funciones, incluido el fragmento único para el procesamiento de DataFrame, `_construct_rf` para ajustar el bosque aleatorio, etc. La mayor parte del tiempo se ha empleado en operaciones DataFrame debido al gran tamaño (45GB) de un día de datos de los registros de clic de Criteo.



### Comparación del tiempo de entrenamiento

Esta sección compara el tiempo de entrenamiento del modelo utilizando pandas convencionales en comparación con el DASK. Para Pandas, cargamos una cantidad menor de datos debido a la naturaleza del tiempo de procesamiento más lento, para evitar que se desbordara la memoria. Por lo tanto, interpolamos los resultados para ofrecer una comparación justa.

La siguiente tabla muestra la comparación del tiempo de entrenamiento bruto cuando hay significativamente menos datos utilizados para el modelo de bosque aleatorio de pandas (50 millones de filas de 20 mil millones por día 15 del conjunto de datos). Esta muestra sólo utiliza menos del 0.25% de todos los datos disponibles. Mientras que para DASK-cuML entrenamos el modelo de bosque aleatorio en las 20 mil millones de filas disponibles. Los dos enfoques dieron lugar a un tiempo de capacitación comparable.

Enfoque	Tiempo de entrenamiento
Scikit-Learn: Usando sólo 50 m de filas en el día 15 como datos de entrenamiento	47 minutos y 21 segundos
RAPIDS-Dask: Utilizando todas las filas 20B del día 15 como datos de entrenamiento	1 hora, 12 minutos y 11 segundos

Si interpolamos los resultados del tiempo de entrenamiento linealmente, como se muestra en la siguiente tabla, hay una ventaja significativa a utilizar el entrenamiento distribuido con DASK. Tomaría el enfoque convencional de Pandas scikit-Learn 13 días para procesar y entrenar 45GB de datos para un solo día de registros tecleo, mientras que EL enfoque RAPIDS-DASK procesa la misma cantidad de datos 262.39 veces más rápido.

Enfoque	Tiempo de entrenamiento
Scikit-Learn: Usando todas las filas 20B en el día15 como datos de entrenamiento	13 días, 3 horas, 40 minutos y 11 segundos
RAPIDS-Dask: Utilizando todas las filas 20B del día 15 como datos de entrenamiento	1 hora, 12 minutos y 11 segundos

En la tabla anterior, puede ver que usando RAPIDS con Dink para distribuir el procesamiento de datos y el entrenamiento de modelos en varias instancias de GPU, el tiempo de ejecución es significativamente más corto en comparación con el procesamiento convencional de Pandas DataFrame con el entrenamiento de modelos scikit-Learn. Este marco permite un escalado vertical y horizontal en el cloud, así como en las instalaciones, en un clúster multinodo con varias GPU.

#### Monitorizar Dink y RAPIDS con Prometheus y Grafana

Una vez que todo se pone en marcha, ejecute inferencias sobre nuevos datos. Los modelos predicen si un usuario hace clic en un anuncio basado en actividades de navegación. Los resultados de la predicción se almacenan en un cuDF de DASK. Puede supervisar los resultados con Prometheus y visualizar en paneles Grafana.

Para obtener más información, consulte este tema ["RAPIDS AI Media Post"](#).

#### Creación de versiones de conjuntos de datos y modelos con el kit de herramientas de operaciones de datos de NetApp

El kit de herramientas DataOPS de NetApp para Kubernetes abstrae los recursos de almacenamiento y las cargas de trabajo de Kubernetes hasta el nivel de espacio de trabajo de ciencia de datos. Estas funciones se presentan en una interfaz sencilla y fácil de usar diseñada para científicos e ingenieros de datos. Utilizando la forma familiar de un programa de Python, el kit de herramientas permite a científicos e ingenieros de datos aprovisionar y destruir espacios de trabajo de JupyterLab en cuestión de segundos. Estas áreas de trabajo pueden contener terabytes o incluso petabytes de capacidad de almacenamiento, lo que permite a los científicos de datos almacenar todos sus conjuntos de datos de entrenamiento directamente en sus espacios de trabajo de proyectos. Han pasado los días de gestionar los espacios de trabajo y los volúmenes de datos por separado.

Para obtener más información, visite el Kit de herramientas ["Repositorio de GitHub"](#).

#### Portátiles Jupyter para referencias

Existen dos cuadernos Jupyter asociados a este informe técnico:

- ["CTR-PandasRF-collated.ipynb."](#) Este cuaderno carga el día 15 desde el conjunto de datos de registros Criteo Terabyte Click, procesa y formatea datos en un DataFrame de Pandas, entrena un modelo de bosque aleatorio Scikit-Learn, realiza predicción y calcula la precisión.
- ["criteo\\_dask\\_RF.ipynb."](#) Este cuaderno carga el día 15 desde el conjunto de datos de registros Criteo Terabyte Click, procesa y formatea datos en un CuDF DASK, entrena un modelo de bosque aleatorio DASK cuML, realiza predicción y calcula la precisión. Al aprovechar varios nodos de trabajo con GPU, este método de procesamiento y entrenamiento de datos distribuidos y modelos es altamente eficiente. Cuantos más datos procese, mayor será el ahorro de tiempo que se consigue con el método DE ML convencional. Puede implementar este portátil en el cloud, en las instalaciones o en un entorno híbrido en el que el clúster de Kubernetes contenga recursos informáticos y de almacenamiento en diferentes ubicaciones, siempre y cuando su configuración de red permita el movimiento libre de datos y la distribución de modelos.

## Conclusión

Azure NetApp Files, RAPIDS y Dink aceleran y simplifican la puesta en marcha del procesamiento y formación DE ML a gran escala gracias a la integración con herramientas de orquestación como Docker y Kubernetes. Al unificar la canalización de datos completa, esta solución reduce la latencia y la complejidad inherentes de muchas cargas de trabajo informáticas avanzadas, y permite salvar en la práctica la brecha entre el desarrollo y las operaciones. Los científicos de datos pueden ejecutar consultas en grandes conjuntos de datos y compartir de forma segura datos y modelos algorítmicos con otros usuarios durante la fase de entrenamiento.

Cuando cree sus propias canalizaciones de IA/ML, configurar la integración, la gestión, la seguridad y la accesibilidad de los componentes en una arquitectura es una tarea ardua. Dar a los desarrolladores acceso y control de su entorno presenta otro conjunto de retos.

Al crear un modelo de entrenamiento distribuido completo y una canalización de datos en el cloud, demostramos dos órdenes de mejora de magnitud en el tiempo de finalización total de los flujos de trabajo en comparación con un enfoque convencional de código abierto que no utilizaba el procesamiento de datos acelerado por GPU y los marcos informáticos.

La combinación de NetApp, Microsoft, marcos de orquestación de código abierto y NVIDIA reúne las tecnologías más recientes como servicios gestionados con una gran flexibilidad para acelerar la adopción de tecnología y mejorar el plazo de comercialización de las nuevas aplicaciones de IA/ML. Estos servicios avanzados se ofrecen en un entorno nativo del cloud que se puede realizar fácilmente en las instalaciones y en arquitecturas de puesta en marcha híbrida.

## Dónde encontrar información adicional

Si quiere obtener más información sobre el contenido de este documento, consulte los siguientes recursos:

- Azure NetApp Files.:
  - Página de arquitectura de soluciones para Azure NetApp Files  
["https://docs.microsoft.com/azure/azure-netapp-files/azure-netapp-files-solution-architectures"](https://docs.microsoft.com/azure/azure-netapp-files/azure-netapp-files-solution-architectures)
- Almacenamiento persistente de Trident para contenedores:
  - Azure NetApp Files y Trident  
["https://netapptrident.readthedocs.io/en/stablev20.07/kubernetes/operations/tasks/backends/anf.html"](https://netapptrident.readthedocs.io/en/stablev20.07/kubernetes/operations/tasks/backends/anf.html)
- DASK y RAPIDS:
  - DASK  
["https://docs.dask.org/en/latest/"](https://docs.dask.org/en/latest/)
  - Instalar el DASK  
["https://docs.dask.org/en/latest/install.html"](https://docs.dask.org/en/latest/install.html)
  - API DASK

["https://docs.dask.org/en/latest/api.html"](https://docs.dask.org/en/latest/api.html)

- Aprendizaje automático DASK

["https://examples.dask.org/machine-learning.html"](https://examples.dask.org/machine-learning.html)

- DASK Distributed Diagnostics

["https://docs.dask.org/en/latest/diagnostics-distributed.html"](https://docs.dask.org/en/latest/diagnostics-distributed.html)

- Marco Y herramientas DE ML:

- TensorFlow: Un marco de aprendizaje automático de código abierto para todos

["https://www.tensorflow.org/"](https://www.tensorflow.org/)

- Docker

["https://docs.docker.com"](https://docs.docker.com)

- Kubernetes

["https://kubernetes.io/docs/home/"](https://kubernetes.io/docs/home/)

- Kubeflow

["http://www.kubeflow.org/"](http://www.kubeflow.org/)

- Servidor de portátiles Jupyter

["http://www.jupyter.org/"](http://www.jupyter.org/)

## **TR-4896: Formación distribuida en Azure: Detección de carriles - diseño de soluciones**

Muneer Ahmad y Verron Martina, NetApp Ronen dar, RUN:AI

Desde mayo de 2019, Microsoft proporciona un servicio de portal nativo de Azure para servicios de archivos NFS y SMB empresariales basados en la tecnología ONTAP de NetApp. Este desarrollo está impulsado por una asociación estratégica entre Microsoft y NetApp, y amplía aún más el alcance de los servicios de datos de ONTAP de primera calidad para Azure.

NetApp, un proveedor líder de servicios de datos en el cloud, se ha Unido para EJECUTARSE: IA, una empresa que virtualiza la infraestructura de IA, para permitir una experimentación de IA más rápida con un uso completo de la GPU. Esta alianza permite a los equipos acelerar la IA ejecutando muchos experimentos en paralelo, con un acceso rápido a los datos y aprovechando recursos informáticos ilimitados. EJECUTAR: La IA permite el uso completo de la GPU al automatizar la asignación de recursos y la arquitectura contrastada de Azure NetApp Files permite que cada experimento se ejecute a la máxima velocidad al eliminar las obstrucciones de la canalización de datos.

NetApp y EJECUTE: La IA ha Unido sus fuerzas para ofrecer a sus clientes una plataforma preparada para el futuro en su viaje a la IA en Azure. Desde el análisis y la computación de alto rendimiento (HPC) hasta las decisiones autónomas (en las que los clientes pueden optimizar sus inversiones EN TECNOLOGÍA pagando



solo por lo que necesitan y cuando lo necesitan), la alianza entre NetApp Y RUN: La IA ofrece una única experiencia unificada en el cloud de Azure.

## **Descripción general de la solución**

En esta arquitectura, lo primordial es la parte de cálculo más intensivo del proceso de formación distribuida de detección de carriles POR inteligencia artificial o aprendizaje automático (ML). La detección de carriles es una de las tareas más importantes en la conducción autónoma, que ayuda a guiar los vehículos por la localización de las marcas de carril. Los componentes estáticos, como las marcas de carril, guían al vehículo a conducir por la carretera de forma interactiva y segura.

Los enfoques basados en la red neuronal convolucional (CNN) han empujado la comprensión y segmentación de la escena a un nuevo nivel. Aunque no funciona bien para objetos con estructuras largas y regiones que podrían ser ocluidas (por ejemplo, postes, sombra en el carril, etc.). La Red neuronal convolucional Espacial (SCNN) generaliza la CNN a un rico nivel espacial. Permite la propagación de información entre neuronas en la misma capa, lo que lo hace más adecuado para objetos estructurados como carriles, polos o camiones con oclusiones. Esta compatibilidad se debe a que la información espacial se puede reforzar y conserva la suavidad y la continuidad.

Es necesario inyectar miles de imágenes de escenas en el sistema para permitir que el modelo aprenda y distinga los diversos componentes del conjunto de datos. Estas imágenes incluyen el tiempo, día o noche, carreteras de varios carriles y otras condiciones de tráfico.

Para la formación, es necesario disponer de buena calidad y cantidad de datos. Una única GPU o varias GPU pueden tardar entre días y semanas para completar el entrenamiento. El entrenamiento con distribución de datos puede acelerar el proceso mediante varias GPU de varios nodos. Horovod es uno de esos marcos que concede entrenamiento distribuido pero la lectura de datos en clústeres de GPU puede ser un obstáculo. Azure NetApp Files proporciona un rendimiento elevado y una latencia baja constante ultrarrápida que proporciona funcionalidades de escalado horizontal y escalado vertical, de tal modo que se utilicen las GPU en la mejor capacidad computacional. Nuestros experimentos verificaron que todas las GPU del clúster se utilizan de media más del 96 % para entrenar la detección de carriles con SCNN.

## **Público objetivo**

La ciencia de datos incorpora varias disciplinas EN TECNOLOGÍA y negocio, por lo que múltiples personas forman parte de nuestro público objetivo:

- Los científicos de datos necesitan la flexibilidad necesaria para utilizar las herramientas y las bibliotecas que prefieran.
- Los ingenieros de datos necesitan saber cómo fluyen los datos y dónde residen.
- Expertos en casos de uso de conducción autónoma.
- Administradores de cloud y arquitectos para configurar y gestionar recursos de cloud (Azure).
- Un ingeniero de DevOps necesita herramientas para integrar nuevas aplicaciones de IA/ML en sus canalizaciones de integración continua y de puesta en marcha continua (CI/CD).
- Los usuarios empresariales quieren tener acceso a aplicaciones de IA/ML.

En este documento, describimos cómo Azure NetApp Files, EJECUTA: IA y Microsoft Azure ayudan a cada uno de estos roles a aportar valor empresarial.

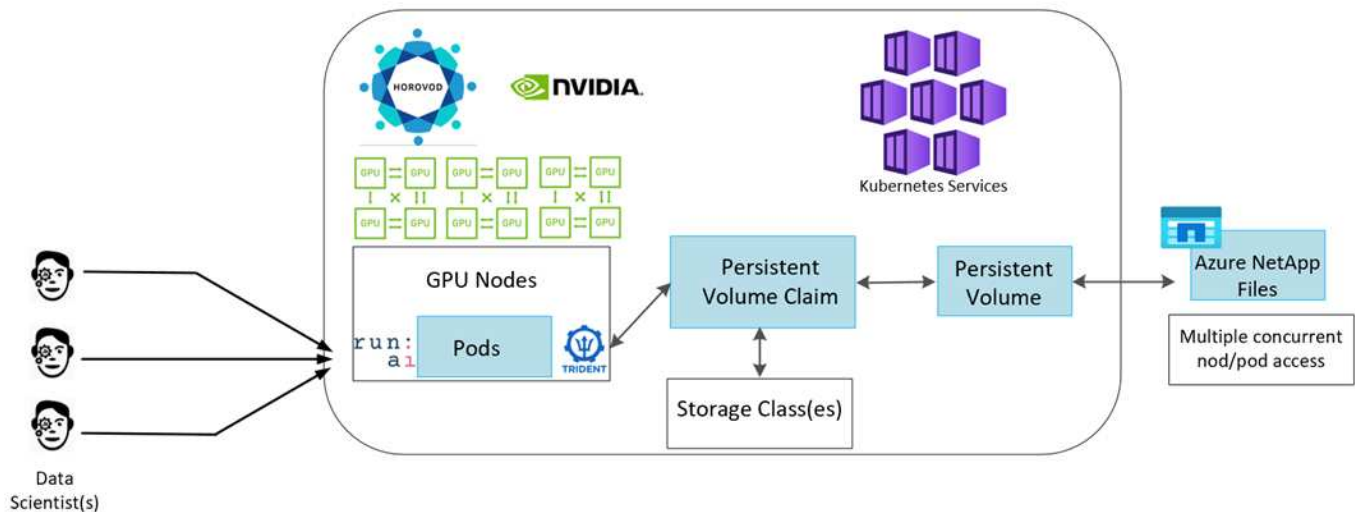
Tecnología de soluciones

Esta sección abarca los requisitos tecnológicos del caso práctico de detección de carriles al implementar una solución de formación distribuida a escala que se ejecuta completamente en el cloud de Azure. La siguiente figura muestra información general sobre la arquitectura de la solución.

Los elementos utilizados en esta solución son:

- Azure Kubernetes Service (AKS)
- SKU de Azure Compute con GPU de NVIDIA
- Azure NetApp Files
- EJECUCIÓN: IA
- Trident de NetApp

Los vínculos a todos los elementos mencionados aquí se enumeran en el "Información adicional" sección.



Requisitos de servicios y recursos cloud

En la siguiente tabla se enumeran los componentes de hardware necesarios para implementar la solución. Los componentes cloud que se usan en cualquier implementación de la solución pueden variar en función de las necesidades del cliente.

Cloud	Cantidad
AKS	Un mínimo de tres nodos de sistema y tres nodos de trabajo de GPU
Nodos del sistema de SKU de máquinas virtuales (VM)	Tres Standard_DS2_v2
Nodos de trabajo de GPU de VM SKU	Tres Standard_NC63_v3
Azure NetApp Files	Nivel estándar de 4 TB

Requisitos de software

En la siguiente tabla se enumeran los componentes de software necesarios para implementar la solución. Los componentes de software que se usan en cualquier implementación de la solución pueden variar en función

de las necesidades del cliente.

De NetApp	Versión u otra información
AKS - versión Kubernetes	1.18.14
EJECUCIÓN: CLI DE IA	v2.2.25
EJECUTAR: versión del operador de Kubernetes de orquestación de IA	1.0.109
Horovod	0.21.2
Trident de NetApp	20.01.1
Timón	3.0.0

### Detección de carriles: Formación distribuida con RUN:AI

En esta sección se ofrecen detalles sobre la configuración de la plataforma para la realización de formación distribuida de detección de carriles a escala mediante EL orquestador de IA. Se trata de la instalación de todos los elementos de la solución y de la ejecución del trabajo de formación distribuido en dicha plataforma. EL versionado DE ML se completa utilizando Snapshot™ de NetApp vinculado A RUN: Experimentos de IA para conseguir la reproducibilidad de los datos y los modelos. LA creación de versiones DE ML desempeña un papel fundamental en el seguimiento de modelos, el intercambio de trabajo entre miembros del equipo, la reproducibilidad de resultados, la implementación de nuevas versiones de modelos a la producción y la procedencia de los datos. El control de versiones DE NetApp ML (Snapshot) puede capturar versiones puntuales de los datos, modelos entrenados y registros asociados con cada experimento. Cuenta con un amplio soporte en la API que facilita la integración con la plataforma DE IA RUN: Solo tiene que activar un evento basado en el estado de entrenamiento. También hay que capturar el estado de todo el experimento sin cambiar nada en el código o los contenedores que se ejecutan sobre Kubernetes (K8s).

Por último, este informe técnico se resume en la evaluación del rendimiento en varios nodos habilitados para GPU en toda la serie AKS.

### Formación distribuida para el caso de uso de detección de carriles mediante el conjunto de datos TuSimple

En este informe técnico, la formación distribuida se realiza en el conjunto de datos TuSimple para la detección de carriles. Horovod se utiliza en el código de entrenamiento para realizar el entrenamiento con datos distribuidos en varios nodos de GPU simultáneamente en el clúster de Kubernetes a través de AKS. El código se presenta como imágenes de contenedor para la descarga y el procesamiento de datos de TuSimple. Los datos procesados se almacenan en volúmenes persistentes asignados por el complemento Trident de NetApp. Para el entrenamiento, se crea una imagen de contenedor más y utiliza los datos almacenados en volúmenes persistentes creados durante la descarga de los datos.

Para enviar el trabajo de datos y entrenamiento, use RUN: AI para orquestar la asignación y la gestión de recursos. RUN: AI le permite realizar operaciones de interfaz de paso de mensajes (MPI) que son necesarias para Horovod. Este diseño permite que varios nodos de GPU se comuniquen entre sí para actualizar los pesos de entrenamiento después de cada lote de entrenamiento. También permite supervisar el entrenamiento a través de la interfaz de usuario y la CLI, para facilitar la supervisión del progreso de los

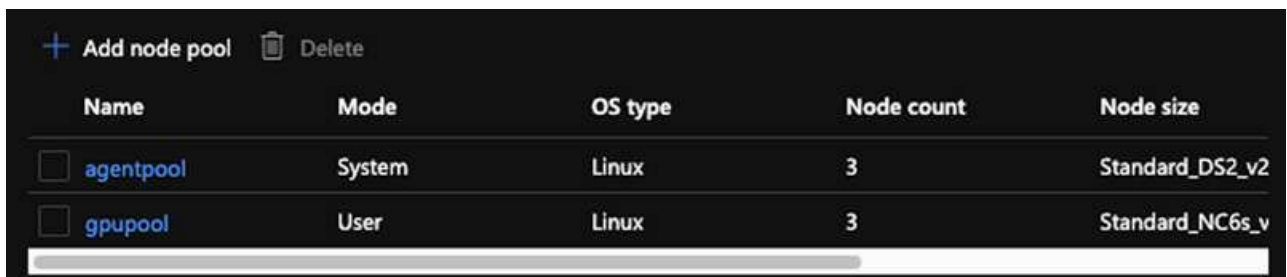
experimentos.

NetApp Snapshot se integra dentro del código de entrenamiento y captura el estado de los datos y el modelo entrenado para cada experimento. Esta funcionalidad le permite realizar un seguimiento de la versión de los datos y el código que se usan y del modelo entrenado asociado generado.

### Instalación e instalación de AKS

Para la configuración e instalación del clúster de AKS, vaya a ["Cree un clúster de AKS"](#). A continuación, siga estos pasos:

1. Al seleccionar el tipo de nodos (ya sean nodos del sistema (CPU) o de trabajadores (GPU)), seleccione lo siguiente:
  - a. Añada el nodo del sistema principal llamado `agentpool` en la `Standard_DS2_v2` tamaño. Utilice los tres nodos predeterminados.
  - b. Agregar nodo de trabajo `gpupool` con the `Standard_NC6s_v3` el tamaño del pool. Utilice un mínimo de tres nodos para los nodos GPU.



Name	Mode	OS type	Node count	Node size
<input type="checkbox"/> agentpool	System	Linux	3	Standard_DS2_v2
<input type="checkbox"/> gpupool	User	Linux	3	Standard_NC6s_v



La implementación tarda entre 5 y 10 minutos.

2. Una vez completada la implementación, haga clic en [Connect to Cluster](#). Para conectarse al clúster AKS recién creado, instale la herramienta de línea de comandos Kubernetes desde su entorno local (portátil/PC). Visite ["Instalar herramientas"](#) Para instalarlo según el sistema operativo.
3. ["Instale el CLI de Azure en su entorno local"](#).
4. Para acceder al clúster AKS desde el terminal, primero introduzca `az login` y coloque las credenciales.
5. Ejecute los dos comandos siguientes:

```
az account set --subscription xxxxxxxx-xxxx-xxxx-xxxx-xxxxxxxxxxxxxxxx
aks get-credentials --resource-group resourcegroup --name aksclustername
```

6. Introduzca este comando en la CLI de Azure:

```
kubectl get nodes
```



Si los seis nodos están en funcionamiento como se ve aquí, su clúster de AKS estará listo y conectado a su entorno local.

```
verronmartina@verron-mac-0 ~ % kubectl get nodes
NAME                                STATUS    ROLES    AGE   VERSION
aks-agentpool-34613062-vmss000000  Ready    agent    22m   v1.18.14
aks-agentpool-34613062-vmss000001  Ready    agent    22m   v1.18.14
aks-agentpool-34613062-vmss000002  Ready    agent    22m   v1.18.14
aks-gpupool-34613062-vmss000000     Ready    agent    20m   v1.18.14
aks-gpupool-34613062-vmss000001     Ready    agent    20m   v1.18.14
aks-gpupool-34613062-vmss000002     Ready    agent    20m   v1.18.14
verronmartina@verron-mac-0 ~ %
```

## Cree una subred delegada para Azure NetApp Files

Para crear una subred delegada para Azure NetApp Files, siga esta serie de pasos:

1. Acceda a redes virtuales en el portal de Azure. Busque la red virtual que acaba de crear. Debería tener un prefijo como el de los roles-vnet, como se ve aquí. Haga clic en el nombre de la red virtual.

Microsoft Azure

Search resources, services, and docs (G+/I)

Dashboard > Virtual networks

seanlucelive (Default Directory)

+ Add Manage view Refresh Export to CSV Open query Assign tags Feedback

Filter by name... Subscription == AzureSub01 Resource group == all Location == all Add filter

Showing 1 to 5 of 5 records. No grouping List view

Name	Resource group	Location	Subscription
aks-vnet-22885919	MC_sluce_rg_TridentDemo_eastus2	East US 2	AzureSub01

2. Haga clic en subredes y seleccione +Subnet en la barra de herramientas superior.

Microsoft Azure

Search resources, services, and docs (G+/I)

Dashboard > Virtual networks > aks-vnet-22885919

aks-vnet-22885919 | Subnets

Virtual network

Search (Ctrl+/) + Subnet Gateway subnet Refresh Manage users Delete

Search subnets

Name	IPv4	IPv6 (many availab...	Delegated to	Security group
aks-subnet	10.240.0.0/16 (65530 av...			aks-agentpool-22885919...

Overview Activity log Access control (IAM) Tags Diagnose and solve problems Settings Address space Connected devices Subnets

3. Proporcione la subred con un nombre como `ANF.sn` Y en el encabezado Delegación de subred, seleccione `Microsoft.NetApp/volumes`. No cambie nada más. Haga clic en Aceptar.

## Add subnet

Name \*

ANF.sn

Subnet address range \* ⓘ

10.0.0.0/24

10.0.0.0 - 10.0.0.255 (251 + 5 Azure reserved addresses)

☐ Add IPv6 address space ⓘ

NAT gateway ⓘ

None

Network security group

None

Route table

None

### SERVICE ENDPOINTS

Create service endpoint policies to allow traffic to specific azure resources from your virtual network over service endpoints. [Learn more](#)

Services ⓘ

0 selected

### SUBNET DELEGATION

Delegate subnet to a service ⓘ

Microsoft.Netapp/volumes

OK

Cancel

Los volúmenes Azure NetApp Files se asignan al clúster de aplicaciones y se consumen como reclamaciones de volúmenes persistentes (RVP) en Kubernetes. A su vez, esta asignación nos proporciona la flexibilidad para asignar volúmenes a diferentes servicios, ya sea ordenadores portátiles Jupyter, funciones sin servidor, etc.

Los usuarios de servicios pueden consumir almacenamiento desde la plataforma de muchas maneras. Las principales ventajas de Azure NetApp Files son:

- Ofrece a los usuarios la capacidad de usar copias Snapshot.
- Permite a los usuarios almacenar grandes cantidades de datos en volúmenes de Azure NetApp Files.
- Obtenga las ventajas en el rendimiento de los volúmenes de Azure NetApp Files cuando ejecute sus

modelos en conjuntos de archivos de gran tamaño.

### Configuración de Azure NetApp Files

Para completar la configuración de Azure NetApp Files, primero debe configurarlo como se describe en ["Inicio rápido: Configure Azure NetApp Files y cree un volumen NFS"](#).

Sin embargo, puede omitir los pasos para crear un volumen NFS para Azure NetApp Files a medida que creará volúmenes a través de Trident. Antes de continuar, asegúrese de que dispone de:

1. ["Registrado para Azure NetApp Files y el proveedor de recursos de NetApp \(a través de Azure Cloud Shell\)"](#).
2. ["Se creó una cuenta en Azure NetApp Files"](#).
3. ["Configure un pool de capacidad"](#) (Como mínimo, 4 TIB Standard o Premium, según sus necesidades).

### Agrupación de la red virtual de AKS y la red virtual de Azure NetApp Files

A continuación, conecte la red virtual AKS (vnet) con Azure NetApp Files vnet siguiendo estos pasos:

1. En el cuadro de búsqueda de la parte superior del portal de Azure, escriba redes virtuales.
2. Haga clic en vnet aks- vnet-name y, a continuación, escriba peerings en el campo de búsqueda.
3. Haga clic en +Agregar e introduzca la información proporcionada en la siguiente tabla:

Campo	Valor o descripción
Nombre de enlace de relación entre iguales	aks-vnet-name_to_anf
SubscriptionId	Suscripción de la red virtual de Azure NetApp Files a la que se está creando una relación de paridad
Partner de vnet peering	Red virtual de Azure NetApp Files



Deje todas las secciones que no sean asteriscos por defecto

4. Haga clic en AGREGAR o en Aceptar para agregar la conexión a la red virtual.

Si desea más información, visite ["Crear, cambiar o eliminar una conexión de red virtual entre iguales"](#).

### Trident

Trident es un proyecto de código abierto que NetApp mantiene para el almacenamiento persistente en contenedores para aplicaciones. Trident se ha implementado como una controladora de aprovisionamiento externa que se ejecuta como un "pod", supervisando volúmenes y automatizando totalmente el proceso de aprovisionamiento.

Trident de NetApp permite una integración sin problemas con K8S mediante la creación y el montaje de volúmenes persistentes para almacenar conjuntos de datos de entrenamiento y modelos entrenados. Esta funcionalidad facilita a los científicos e ingenieros de datos el uso de K8 sin los problemas de almacenar y gestionar manualmente conjuntos de datos. Trident también elimina la necesidad que tienen los científicos de datos de aprender a gestionar nuevas plataformas de datos a medida que integra las tareas relacionadas con la gestión de datos a través de la integración lógica de las API.

## Instale Trident

Para instalar el software Trident, realice los pasos siguientes:

1. "Primero instale el timón".
2. Descargue y extraiga el instalador de Trident 21.01.1.

```
wget
https://github.com/NetApp/trident/releases/download/v21.01.1/trident-
installer-21.01.1.tar.gz
tar -xf trident-installer-21.01.1.tar.gz
```

3. Cambie el directorio a. trident-installer.

```
cd trident-installer
```

4. Copiar tridentctl a un directorio del sistema \$PATH.

```
cp ./tridentctl /usr/local/bin
```

5. Instale Trident en el clúster K8s con Helm:

- a. Cambiar el directorio al directorio del timón.

```
cd helm
```

- b. Instale Trident.

```
helm install trident trident-operator-21.01.1.tgz --namespace trident
--create-namespace
```

- c. Compruebe el estado de Trident pods de la manera habitual K8s:

```
kubect1 -n trident get pods
```

- d. Si todos los pods están ya en funcionamiento, se instala Trident y se podrá seguir avanzando.

## Configure el back-end de Azure NetApp Files y la clase de almacenamiento

Para configurar el back-end de Azure NetApp Files y la clase de almacenamiento, complete los siguientes pasos:

1. Vuelva al directorio inicial.



```
cd ~
```

2. Clone el "repositorio de proyectos" lane-detection-SCNN-horovod.
3. Vaya a la trident-config directorio.

```
cd ./lane-detection-SCNN-horovod/trident-config
```

4. Crear un principio de servicio de Azure (el principio del servicio es cómo Trident se comunica con Azure para acceder a sus recursos de Azure NetApp Files).

```
az ad sp create-for-rbac --name
```

El resultado debería ser como el ejemplo siguiente:

```
{
  "appId": "xxxxx-xxxx-xxxx-xxxx-xxxxxxxxxxxx",
  "displayName": "netapptrident",
  "name": "http://netapptrident",
  "password": "xxxxxxxxxxxxxxxxxxxxxxxxxxxxxxxx",
  "tenant": "xxxxxxxx-xxxx-xxxx-xxxx-xxxxxxxxxxxx"
}
```

5. Cree Trident backend json archivo.
6. Con el editor de texto preferido, complete los siguientes campos de la tabla siguiente dentro de la anf-backend.json archivo.

Campo	Valor
SubscriptionId	Su ID de suscripción de Azure
ID de tenantID	Su ID de inquilino de Azure (de la salida de az ad sp en el paso anterior)
ID del Cliente	Su AppID (de la salida de az ad sp en el paso anterior)
ClientSecret	Su contraseña (de la salida de az ad sp en el paso anterior)

El archivo debería tener el siguiente ejemplo:

```
{
  "version": 1,
  "storageDriverName": "azure-netapp-files",
  "subscriptionID": "fakec765-4774-fake-ae98-a721add4fake",
  "tenantID": "fakef836-edc1-fake-bff9-b2d865eefake",
  "clientID": "fake0f63-bf8e-fake-8076-8de91e57fake",
  "clientSecret": "SECRET",
  "location": "westeurope",
  "serviceLevel": "Standard",
  "virtualNetwork": "anf-vnet",
  "subnet": "default",
  "nfsMountOptions": "vers=3,proto=tcp",
  "limitVolumeSize": "500Gi",
  "defaults": {
    "exportRule": "0.0.0.0/0",
    "size": "200Gi"
  }
}
```

7. Indique a Trident que cree el back-end de Azure NetApp Files en la `trident` espacio de nombres, utilizar `anf-backend.json` como el archivo de configuración de la siguiente manera:

```
tridentctl create backend -f anf-backend.json -n trident
```

8. Cree la clase de almacenamiento:

- a. Los usuarios de K8 aprovisionan volúmenes mediante el uso de EVs que especifican una clase de almacenamiento por nombre. Indique a K8S que cree una clase de almacenamiento `azurenetaappfiles` Que hará referencia al back-end de Azure NetApp Files creado en el paso anterior utilizando lo siguiente:

```
kubectl create -f anf-storage-class.yaml
```

- b. Compruebe que la clase de almacenamiento se haya creado mediante el siguiente comando:

```
kubectl get sc azurenetaappfiles
```

El resultado debería ser como el ejemplo siguiente:

NAME	PROVISIONER	RECLAIMPOLICY	VOLUMEBINDINGMODE	ALLOWVOLUMEEXPANSION	AGE
azurenetaappfiles	csi.trident.netapp.io	Delete	Immediate	false	98s

## Ponga en marcha y configure componentes snapshot para volúmenes en AKS

Si el clúster no viene preinstalado con los componentes Snapshot de volumen correctos, puede realizar la

instalación manual de estos componentes ejecutando los siguientes pasos:



AKS 1.18.14 no tiene una controladora Snapshot instalada previamente.

1. Instale los CRD de la versión beta de instantánea utilizando los siguientes comandos:

```
kubectl create -f https://raw.githubusercontent.com/kubernetes-csi/external-snapshotter/release-3.0/client/config/crd/snapshot.storage.k8s.io_volumesnapshotclasses.yaml
kubectl create -f https://raw.githubusercontent.com/kubernetes-csi/external-snapshotter/release-3.0/client/config/crd/snapshot.storage.k8s.io_volumesnapshotcontents.yaml
kubectl create -f https://raw.githubusercontent.com/kubernetes-csi/external-snapshotter/release-3.0/client/config/crd/snapshot.storage.k8s.io_volumesnapshots.yaml
```

2. Instale el controlador Snapshot con los siguientes documentos de GitHub:

```
kubectl apply -f https://raw.githubusercontent.com/kubernetes-csi/external-snapshotter/release-3.0/deploy/kubernetes/snapshot-controller/rbac-snapshot-controller.yaml
kubectl apply -f https://raw.githubusercontent.com/kubernetes-csi/external-snapshotter/release-3.0/deploy/kubernetes/snapshot-controller/setup-snapshot-controller.yaml
```

3. Configurar K8s volumesnapshotclass: Antes de crear una instantánea de volumen, a. "[clase de snapshot de volumen](#)" debe estar configurado. Cree una clase de copia Snapshot de volumen para Azure NetApp Files y utilícela para crear versiones ML con la tecnología Snapshot de NetApp. Cree volumesnapshotclass netapp-csi-snapclass y configúrelo como predeterminado "volumesnapshotclass" como tal:

```
kubectl create -f netapp-volume-snapshot-class.yaml
```

El resultado debería ser como el ejemplo siguiente:

```
volumesnapshotclass.snapshot.storage.k8s.io/netapp-csi-snapclass created
```

4. Compruebe que la clase de copia Snapshot de volumen se haya creado con el siguiente comando:

```
kubectl get volumesnapshotclass
```

El resultado debería ser como el ejemplo siguiente:

NAME	DRIVER	DELETIONPOLICY	AGE
netapp-csi-snapclass	csi.trident.netapp.io	Delete	63s

## EJECUCIÓN:instalación de IA

Para instalar RUN:AI, realice los siguientes pasos:

1. "Instalar EJECTAR:clúster AI en AKS".
2. Vaya a [app.runai.ai](https://app.runai.ai), haga clic en Crear nuevo proyecto y asigne un nombre a la detección de carriles. Creará un espacio de nombres en un clúster K8s a partir de `runai-` seguido del nombre del proyecto. En este caso, el espacio de nombres creado sería la detección de pistas en ejecución.

New Project

Basics

Node Affinity

Time Limit

Basics

Project Name ⓘ

lane-detection

Assigned GPUs

3

Over-quota for project

☒ Allow over-quota

Save Cancel

3. "INSTALAR RUN:AI CLI".
4. En el terminal, establezca la detección de carriles como UNA EJECUCIÓN predeterminada: Proyecto de IA mediante el siguiente comando:

```
`runai config project lane-detection`
```

El resultado debería ser como el ejemplo siguiente:

```
Project lane-detection has been set as default project
```

5. Crear ClusterRole y ClusterRoleBinding para el espacio de nombres del proyecto (por ejemplo, lane-detection) por lo tanto, la cuenta de servicio predeterminada que pertenece a runai-lane-detection el espacio de nombres tiene permiso para ejecutar volumesnapshot operaciones durante la ejecución de trabajos:
  - a. Enumere los espacios de nombres para comprobarlo runai-lane-detection existe usando este comando:

```
kubectl get namespaces
```

El resultado debería aparecer como el ejemplo siguiente:

NAME	STATUS	AGE
default	Active	130m
kube-node-lease	Active	130m
kube-public	Active	130m
kube-system	Active	130m
runai	Active	4m44s
runai-lane-detection	Active	13s
trident	Active	102m

6. Crear función de clúster netappsnapshot Y ClusterRoleBinding netappsnapshot con los siguientes comandos:

```
`kubectl create -f runai-project-snap-role.yaml`  
`kubectl create -f runai-project-snap-role-binding.yaml`
```

### Descargue y procese el conjunto de datos de TuSimple como RUN:AI job

El proceso para descargar y procesar el conjunto de datos TuSimple COMO UNA EJECUCIÓN: El trabajo de IA es opcional. Se trata de los siguientes pasos:

1. Cree y empuje la imagen del docker, o omita este paso si desea utilizar una imagen del docker existente (por ejemplo, muneer7589/download-tusimple:1.0)
  - a. Cambie al directorio principal:

```
cd ~
```

- b. Vaya al directorio de datos del proyecto lane-detection-SCNN-horovod:

```
cd ../lane-detection-SCNN-horovod/data
```

- c. Modificar `build_image.sh` script de shell y cambiar el repositorio de docker a la suya. Por ejemplo, sustituir `muneer7589` con el nombre de repositorio de docker. También puede cambiar el nombre y LA ETIQUETA de la imagen del docker (por ejemplo `download-tusimple` y `1.0`):

```
#!/bin/bash
#
# A simple script to build the Docker image.
#
# $ build_image.sh
set -ex

IMAGE=muneer7589/download-tusimple
TAG=1.0

# Build image
echo "Building image: "$IMAGE
docker build . -f Dockerfile \
  --tag "${IMAGE}:${TAG}"
echo "Finished building image: "$IMAGE

# Push image
echo "Pushing image: "$IMAGE
docker push "${IMAGE}:${TAG}"
echo "Finished pushing image: "$IMAGE
```

- d. Ejecute el script para crear la imagen de docker y empújela al repositorio de docker mediante los siguientes comandos:

```
chmod +x build_image.sh
./build_image.sh
```

2. Envíe el TRABAJO RUN: AI para descargar, extraer, preprocesar y almacenar el conjunto de datos de detección de carriles TuSimple en un `pvc` Creado dinámicamente por Trident de NetApp:

- a. Use los siguientes comandos para enviar LA EJECUCIÓN: AI job:

```
runai submit
--name download-tusimple-data
--pvc azurenetafiles:100Gi:/mnt
--image muneer7589/download-tusimple:1.0
```

b. Introduzca la información de la siguiente tabla para enviar EL trabajo RUN:AI:

Campo	Valor o descripción
-name	Nombre del trabajo
-pvc	PVC del formato [StorageClassName]:Size:ContainerMountPath en el envío de trabajos anterior, está creando un PVC basado en demanda usando Trident con archivos azurenetappfiles de la clase de almacenamiento. La capacidad de volumen persistente aquí es 100Gi y está montada en path /mnt.
-imagen	Imagen de Docker que se utilizará al crear el contenedor para este trabajo

El resultado debería ser como el ejemplo siguiente:

```
The job 'download-tusimple-data' has been submitted successfully
You can run `runai describe job download-tusimple-data -p lane-detection` to check the job status
```

c. Enumere los trabajos ENVIADOS RUN:AI.

```
runai list jobs
```

```
Showing jobs for project lane-detection
NAME          STATUS      AGE  NODE          IMAGE                                     TYPE  PROJECT      USER          GPUs Allocated (Requested)
PODs Running (Pending) SERVICE URL(S)
download-tusimple-data ContainerCreating 1m   aks-agentpool-34613062-vmss00000a muneer7589/download-tusimple:1.0 Train lane-detection veronmartina 0 (0)
1 (0)
```

d. Compruebe los registros de trabajos enviados.

```
runai logs download-tusimple-data -t 10
```

```
751150K ..... 6% 16.2M 20m37s
751200K ..... 6% 11.1M 20m37s
751250K ..... 6% 12.5M 20m36s
751300K ..... 6% 11.3M 20m36s
751350K ..... 6% 15.2M 20m36s
751400K ..... 6% 10.5M 20m36s
751450K ..... 6% 15.2M 20m36s
751500K ..... 6% 14.1M 20m36s
751550K ..... 6% 24.3M 20m36s
751600K ..... 6% 26.3M 20m36s
```

e. Enumere la pvc creado. Utilice esto pvc comando para la formación en el siguiente paso.

```
kubectl get pvc | grep download-tusimple-data
```

El resultado debería ser como el ejemplo siguiente:

- a. Compruebe el trabajo EN RUN: AI UI (o. `app.run.ai`).

Job Name	Status	User	Project	Total Run Time	Creation Time	Type	GPU Utilization	Used CPU
download-tusimple-data	Running	verronma...	lane-detection	00:07:11	03/03/21, 2:51PM	Train	-	0.00
build1	Deleted	root	lane-detection	00:01:56	03/01/21, 10:18...	Interactive	-	-
download-tusimple-data	Deleted	root	lane-detection	-	03/01/21, 9:58AM	Train	-	-
download-tusimple-data	Deleted	root	lane-detection	-	03/01/21, 10:03...	Train	-	-
download-tusimple-data	Deleted	root	lane-detection	00:02:55	03/01/21, 10:24...	Train	-	-
download-tusimple-data	Deleted	root	lane-detection	-	03/01/21, 10:30...	Train	-	-
download-tusimple-data	Deleted	root	lane-detection	00:13:17	03/01/21, 11:41...	Train	-	-
download-tusimple-data-1	Deleted	verronma...	lane-detection	-	02/26/21, 5:30PM	Train	-	-

### Realice la formación de detección de carriles distribuidos con Horovod

El entrenamiento de detección de carriles distribuidos con Horovod es un proceso opcional. Sin embargo, estos son los pasos implicados:

1. Cree y empuje la imagen del docker o omita este paso si desea utilizar la imagen del docker existente (por ejemplo, `muneer7589/dist-lane-detection:3.1`) :

- a. Cambie al directorio inicial.

```
cd ~
```

- b. Vaya al directorio del proyecto `lane-detection-SCNN-horovod`.

```
cd ./lane-detection-SCNN-horovod
```

- c. Modifique el `build_image.sh` script de shell y cambie el repositorio de docker a la suya (por ejemplo, reemplace `muneer7589` con el nombre del repositorio del docker). También puede cambiar el nombre y LA ETIQUETA de la imagen del docker (`dist-lane-detection y. 3.1, for example`).



```
#!/bin/bash
#
# A simple script to build the distributed Docker image.
#
# $ build_image.sh
set -ex

IMAGE=muneer7589/dist-lane-detection
TAG=3.0

# Build image
echo "Building image: "$IMAGE
docker build . -f Dockerfile \
  --tag "${IMAGE}:${TAG}"
echo "Finished building image: "$IMAGE

# Push image
echo "Pushing image: "$IMAGE
docker push "${IMAGE}:${TAG}"
echo "Finished pushing image: "$IMAGE
```

- d. Ejecute la secuencia de comandos para crear la imagen de docker y empújela al repositorio de docker.

```
chmod +x build_image.sh
./build_image.sh
```

## 2. Enviar LA CARRERA: Al job para llevar a cabo la formación distribuida (MPI):

- a. Uso de envío DE LA EJECUCIÓN: La inteligencia artificial para la creación automática de la RVP en el paso anterior (para la descarga de datos) solo le permite tener acceso a RWO, que no permite que varios POD o nodos accedan al mismo RVP para el entrenamiento distribuido. Actualice el modo de acceso a ReadWriteMany y utilice el parche Kubernetes para hacerlo.
- b. En primer lugar, ejecute el siguiente comando para obtener el nombre del volumen de la RVP:

```
kubectl get pvc | grep download-tusimple-data
```

```
root@ai-w-gpu-2:/mnt/ai_data/anf_runai/lane-detection-SCNN-horovod# kubectl get pvc | grep download-tusimple-data
pvc-download-tusimple-data-0   Bound   pvc-bb03b74d-2c17-40c4-a445-79f3de8d16d5   100Gi   RWX   azurenetappfiles   2d4h
```

- c. Aplicar parches al modo de acceso de volumen y actualización a ReadWriteMany (reemplace el nombre del volumen por el suyo en el siguiente comando):

```
kubectl patch pv pvc-bb03b74d-2c17-40c4-a445-79f3de8d16d5 -p
'{"spec":{"accessModes":["ReadWriteMany"]}}'
```

- d. Envíe el TRABAJO RUN: Al MPI para la ejecución del trabajo de formación distribuida utilizando la información de la tabla siguiente:

```

runai submit-mpi
--name dist-lane-detection-training
--large-shm
--processes=3
--gpu 1
--pvc pvc-download-tusimple-data-0:/mnt
--image muneer7589/dist-lane-detection:3.1
-e USE_WORKERS="true"
-e NUM_WORKERS=4
-e BATCH_SIZE=33
-e USE_VAL="false"
-e VAL_BATCH_SIZE=99
-e ENABLE_SNAPSHOT="true"
-e PVC_NAME="pvc-download-tusimple-data-0"

```

Campo	Valor o descripción
nombre	Nombre del trabajo de formación distribuido
gran km	Monte un dispositivo /dev/sm de gran tamaño. Es un sistema de archivos compartidos montado en la RAM y proporciona memoria compartida lo suficientemente grande para que varios trabajadores de CPU procesen y carguen lotes en la RAM de la CPU.
procesos	Número de procesos de formación distribuidos
gpu	Número de GPU/procesos que se van a asignar para la tarea en esta tarea, hay tres procesos de trabajo de GPU (--process=3), cada uno asignado con una única GPU (--gpu 1)
rvp	Utilice el volumen persistente existente (pvc-download-tusimple-data-0) creado por el trabajo anterior (download-tusimple-data) y se monta en la ruta /mnt
imagen	Imagen de Docker que se utilizará al crear el contenedor para este trabajo
Defina las variables de entorno que se van a establecer en el contenedor	
USE_WORKERS	Al establecer el argumento en true, se activa la carga de datos multiproceso
NÚM_TRABAJADORES	Número de procesos de trabajo del cargador de datos
TAMAÑO_LOTE	Tamaño de lote de entrenamiento
VALOR_USO	Establecer el argumento en true permite la validación
VAL_BATCH_SIZE	Tamaño del lote de validación

Campo	Valor o descripción
ENABLE_SNAPSHOT	Establecer el argumento en true permite tomar instantáneas de datos y modelos entrenados para el control de versiones EN ML
NOMBRE_PVC	Nombre de la rvp de la que se va a realizar una instantánea. En la presentación de trabajos anterior, está tomando una instantánea de pvc-download-tusimple-data-0, que consta de conjuntos de datos y modelos entrenados

El resultado debería ser como el ejemplo siguiente:

```
The job 'dist-lane-detection-training' has been submitted successfully
You can run `runai describe job dist-lane-detection-training -p lane-detection` to check the job status
```

e. Enumera el trabajo enviado.

```
runai list jobs
```

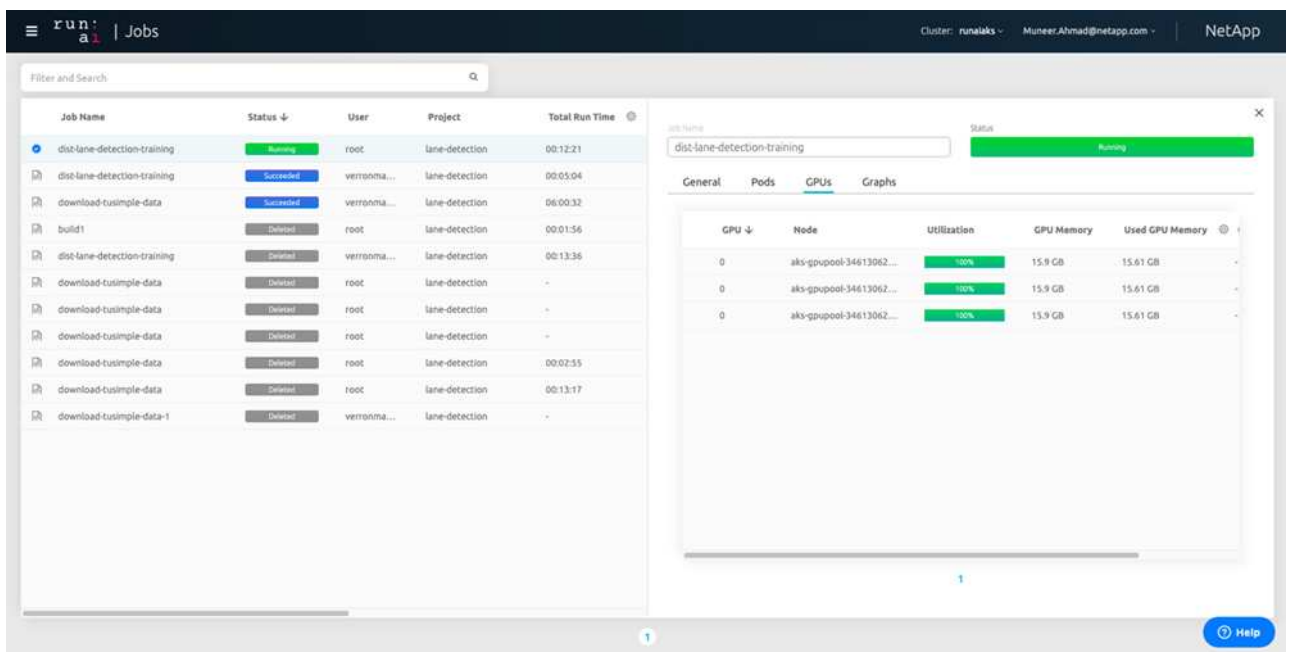
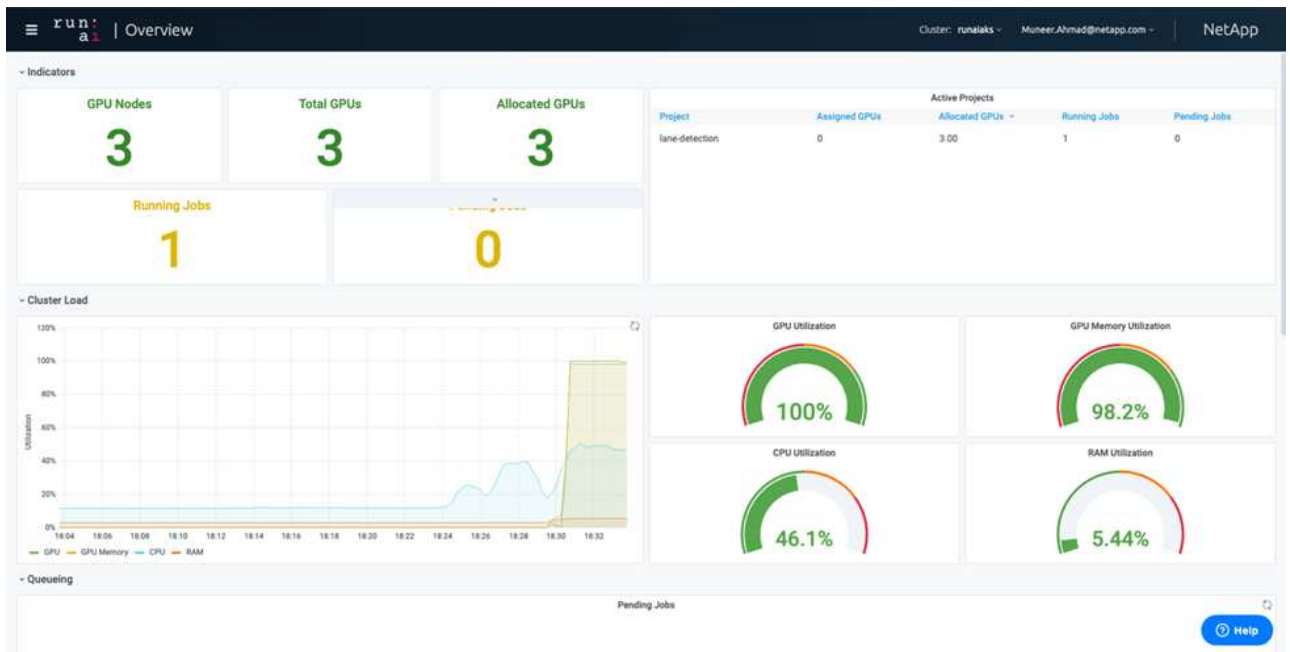
NAME	SERVICE URL(S)	STATUS	AGE	NODE	IMAGE	TYPE	PROJECT	USER	GPUs Allocated (Requested)	PODs
download-tusimple-data		Succeeded	1d		muneer7589/download-tusimple:1.0	Train	lane-detection	verronmartina	0 (0)	0 (0)
dist-lane-detection-training		Init:0/1	2m	<multiple>	muneer7589/dist-lane-detection:3.1	Train	lane-detection	root	3 (3)	4 (0)

f. Registros de trabajos enviados:

```
runai logs dist-lane-detection-training
```

```
root@ai-w-gpu-2:~/runai# runai logs dist-lane-detection-training
Running with 3 workers
2021-03-04 17:29:23.158449: I tensorflow/stream_executor/platform/default/dso_loader.cc:48] Successfully opened dynamic library libcudart.so.10.1
+ POD_NAME=dist-lane-detection-training-worker-0
+ [ d = - ]
+ shift
+ /opt/kube/kubect1 cp /opt/kube/hosts dist-lane-detection-training-worker-0:/etc/hosts_of_nodes
+ POD_NAME=dist-lane-detection-training-worker-2
+ [ d = - ]
+ shift
+ /opt/kube/kubect1 cp /opt/kube/hosts dist-lane-detection-training-worker-2:/etc/hosts_of_nodes
+ POD_NAME=dist-lane-detection-training-worker-1
```

g. Compruebe la tarea de entrenamiento EN EJECUCIÓN: Interfaz gráfica de usuario de IA (o app.runai.ai): RUN: Consola de IA, como se puede ver en las siguientes figuras. La primera figura detalla tres GPU asignadas para el trabajo de entrenamiento distribuido, repartidos en tres nodos en AKS, y la segunda EJECUCIÓN: Trabajos de IA:



h. Una vez finalizada la formación, compruebe la copia de Snapshot de NetApp que se creó y vinculado con RUN: Trabajo de IA.

```
runai logs dist-lane-detection-training --tail 1
```

```
[1,0]<stdout>:Snapshot snap-pvc-download-tusimple-data-0-dist-lane-detection-training-launcher-2021-03-05-16-23-42 created in namespace runai-lane-detection
```

```
kubectl get volumesnapshots | grep download-tusimple-data-0
```

## Restaura datos de la copia Snapshot de NetApp

Para restaurar datos de la copia Snapshot de NetApp, complete los siguientes pasos:

1. Cambie al directorio inicial.

```
cd ~
```

2. Vaya al directorio del proyecto lane-detection-SCNN-horovod.

```
cd ../lane-detection-SCNN-horovod
```

3. Modificar `restore-snapshot-pvc.yaml` y actualícelo `dataSource` name Campo de la copia Snapshot desde la que desea restaurar datos. También puede cambiar el nombre del PVC al que se restaurarán los datos, en este ejemplo su `restored-tusimple`.

```
apiVersion: v1
kind: PersistentVolumeClaim
metadata:
  name: restored-tusimple
spec:
  storageClassName: azurenetappfiles
  dataSource:
    name: snap-pvc-download-tusimple-data-0-dist-lane-detection-training-launcher-2021-03-05-16-23-42
    kind: VolumeSnapshot
    apiGroup: snapshot.storage.k8s.io
  accessModes:
    - ReadWriteMany
  resources:
    requests:
      storage: 100Gi
```

4. Cree una nueva RVP mediante `restore-snapshot-pvc.yaml`.

```
kubectl create -f restore-snapshot-pvc.yaml
```

El resultado debería ser como el ejemplo siguiente:

```
persistentvolumeclaim/restored-tusimple created
```

5. Si desea utilizar los datos recién restaurados para el entrenamiento, el envío de trabajos sigue siendo el mismo que antes; sólo sustituya el `PVC_NAME` con el restaurado `PVC_NAME` al enviar el trabajo de formación, como se muestra en los siguientes comandos:

```
runai submit-mpi
--name dist-lane-detection-training
--large-shm
--processes=3
--gpu 1
--pvc restored-tusimple:/mnt
--image muneer7589/dist-lane-detection:3.1
-e USE_WORKERS="true"
-e NUM_WORKERS=4
-e BATCH_SIZE=33
-e USE_VAL="false"
-e VAL_BATCH_SIZE=99
-e ENABLE_SNAPSHOT="true"
-e PVC_NAME="restored-tusimple"
```

### Evaluación del rendimiento

Para mostrar la escalabilidad lineal de la solución, se han realizado pruebas de rendimiento para dos supuestos: Una GPU y tres GPU. La asignación de GPU, la utilización de la GPU y la memoria, se han capturado diferentes métricas de uno y tres nodos durante el entrenamiento en el conjunto de datos de detección de carriles TuSimple. Los datos se incrementan cinco veces sólo por analizar la utilización de los recursos durante los procesos de entrenamiento.

La solución permite a los clientes comenzar con un conjunto de datos pequeño y unas pocas GPU. Cuando aumentan la cantidad de datos y la demanda de GPU, los clientes pueden escalar horizontalmente de forma dinámica los terabytes del nivel estándar y escalar rápidamente hasta el nivel Premium para obtener el cuádruple de rendimiento por terabyte sin necesidad de mover datos. Este proceso se explica más detalladamente en la sección, "[Niveles de servicio de Azure NetApp Files](#)".

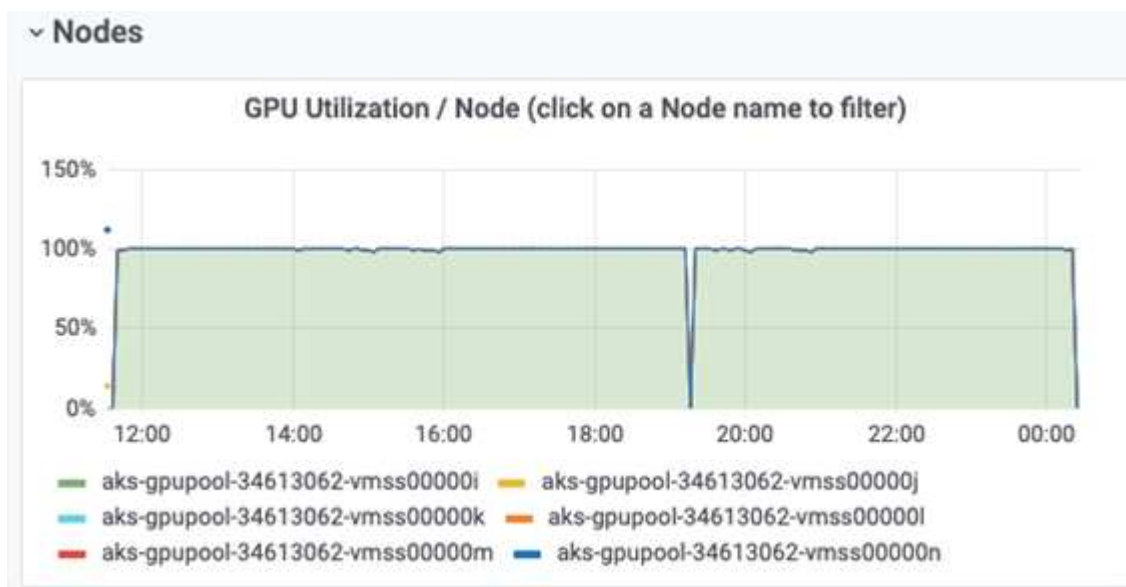
El tiempo de procesamiento en una GPU era de 12 horas y 45 minutos. El tiempo de procesamiento en tres GPU en tres nodos era de aproximadamente 4 horas y 30 minutos.

Las cifras que se muestran a lo largo del resto de este documento muestran ejemplos de rendimiento y escalabilidad basados en las necesidades empresariales individuales.

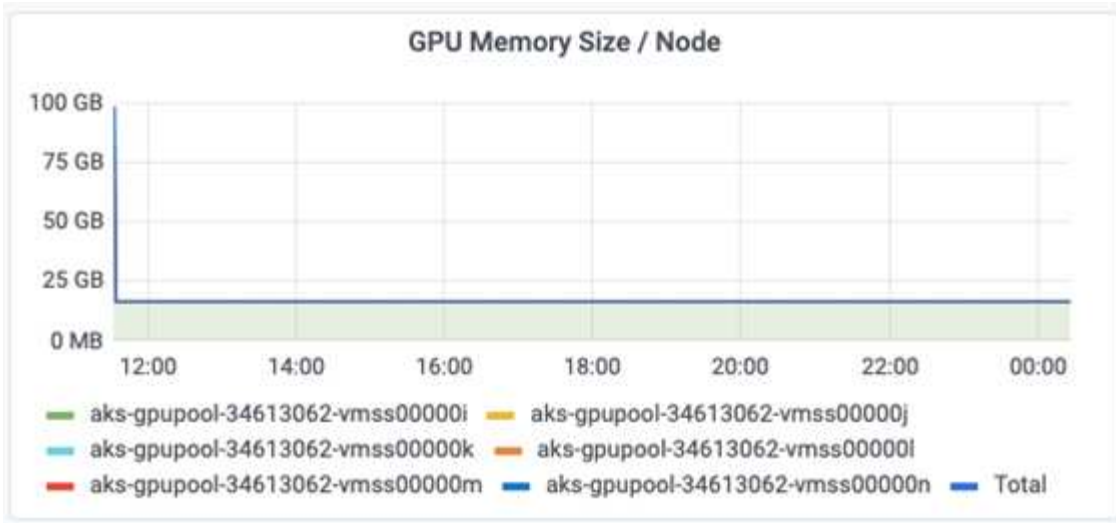
La siguiente figura muestra la asignación de 1 GPU y la utilización de memoria.



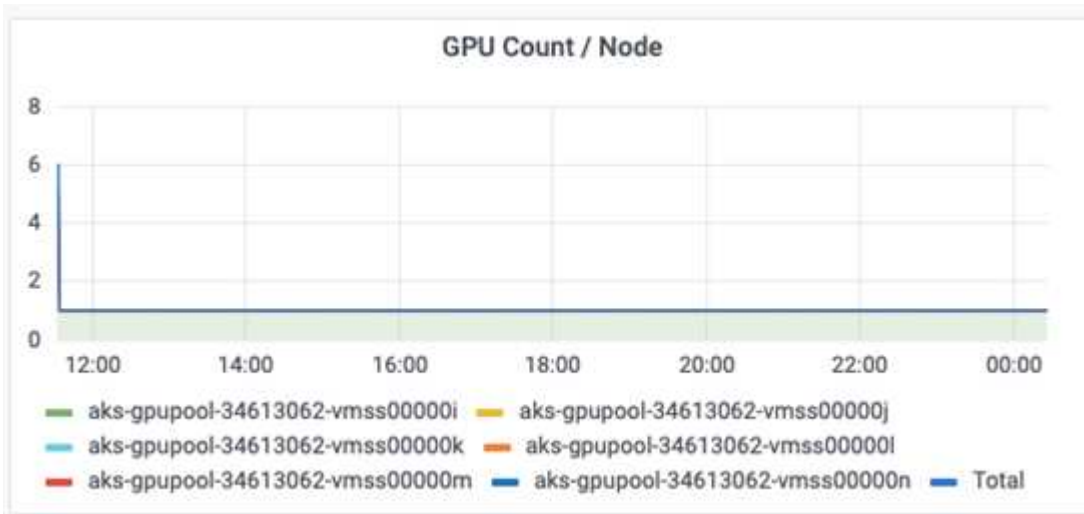
La siguiente figura ilustra el uso de GPU de un solo nodo.



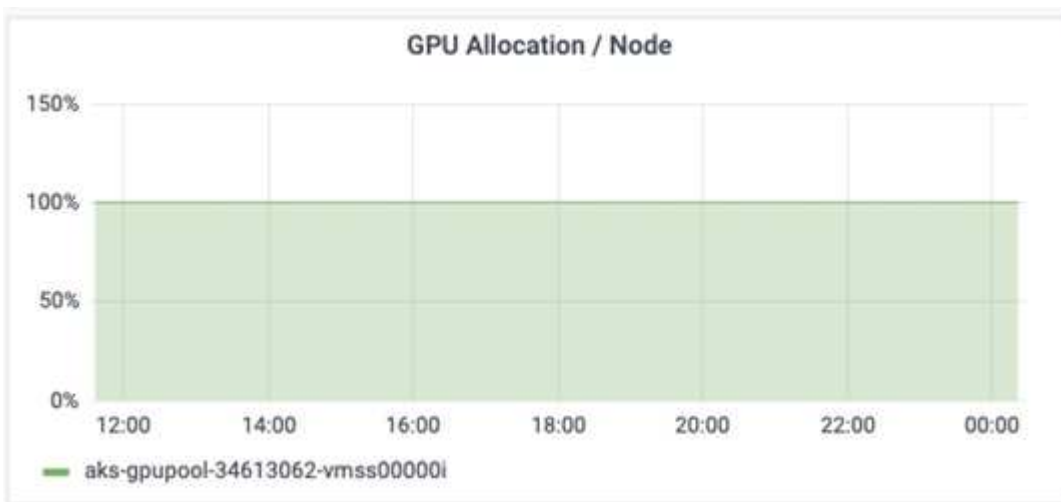
La siguiente figura ilustra el tamaño de la memoria de un solo nodo (16 GB).



La siguiente figura muestra el número de GPU de nodo único (1).

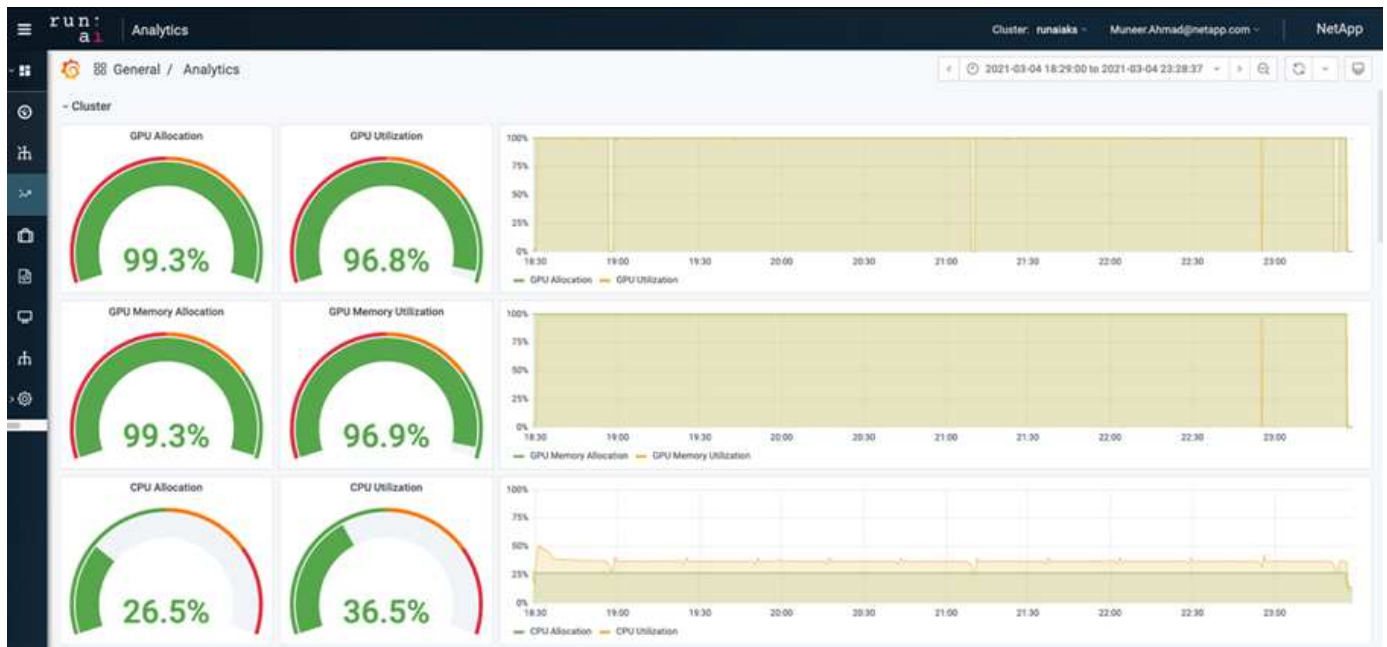


La siguiente figura muestra la asignación de GPU de un solo nodo (%).

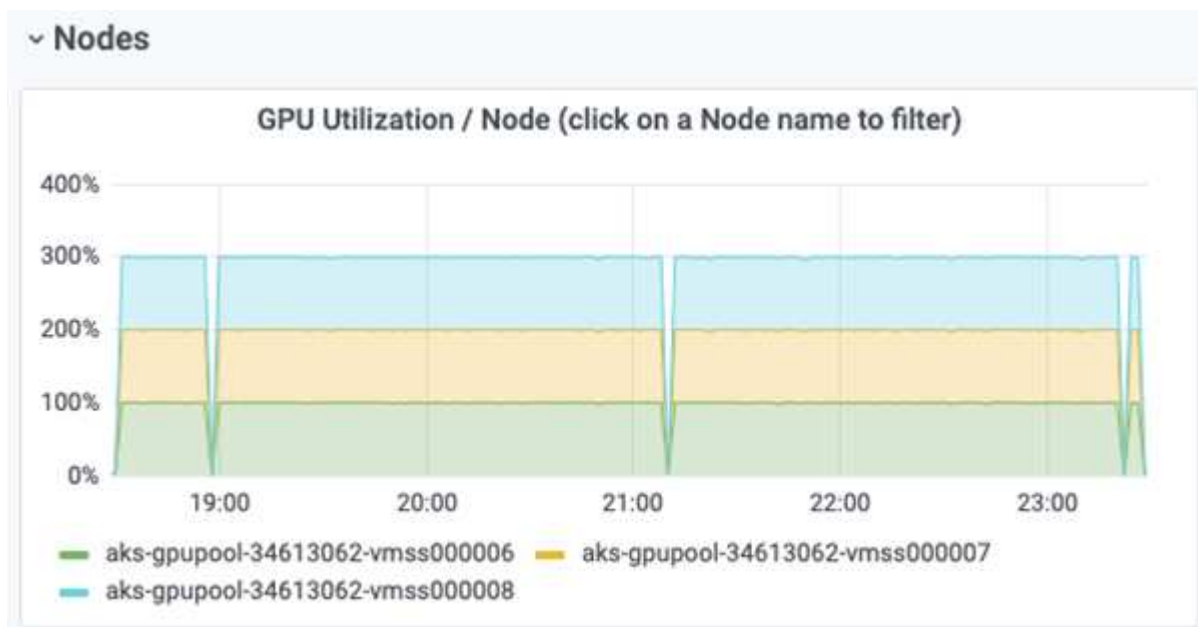


La siguiente figura muestra tres GPU en tres nodos: Asignación de GPU y memoria.

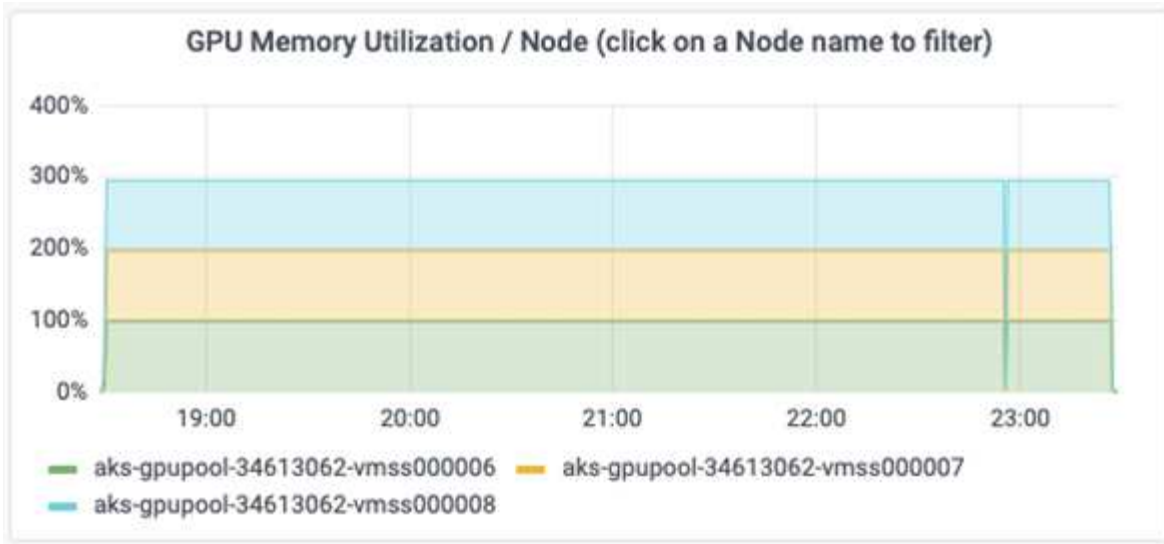




La siguiente figura muestra tres GPU en la utilización de tres nodos (%).



La siguiente figura muestra tres GPU en una utilización de la memoria de tres nodos (%).



### Niveles de servicio de Azure NetApp Files

Es posible cambiar el nivel de servicio de un volumen existente si se mueve el volumen a otro pool de capacidad que utiliza "nivel de servicio" se desea para el volumen. Este cambio de nivel de servicio existente del volumen no requiere la migración de los datos. Además, no afecta el acceso al volumen.

### Cambie dinámicamente el nivel de servicio de un volumen

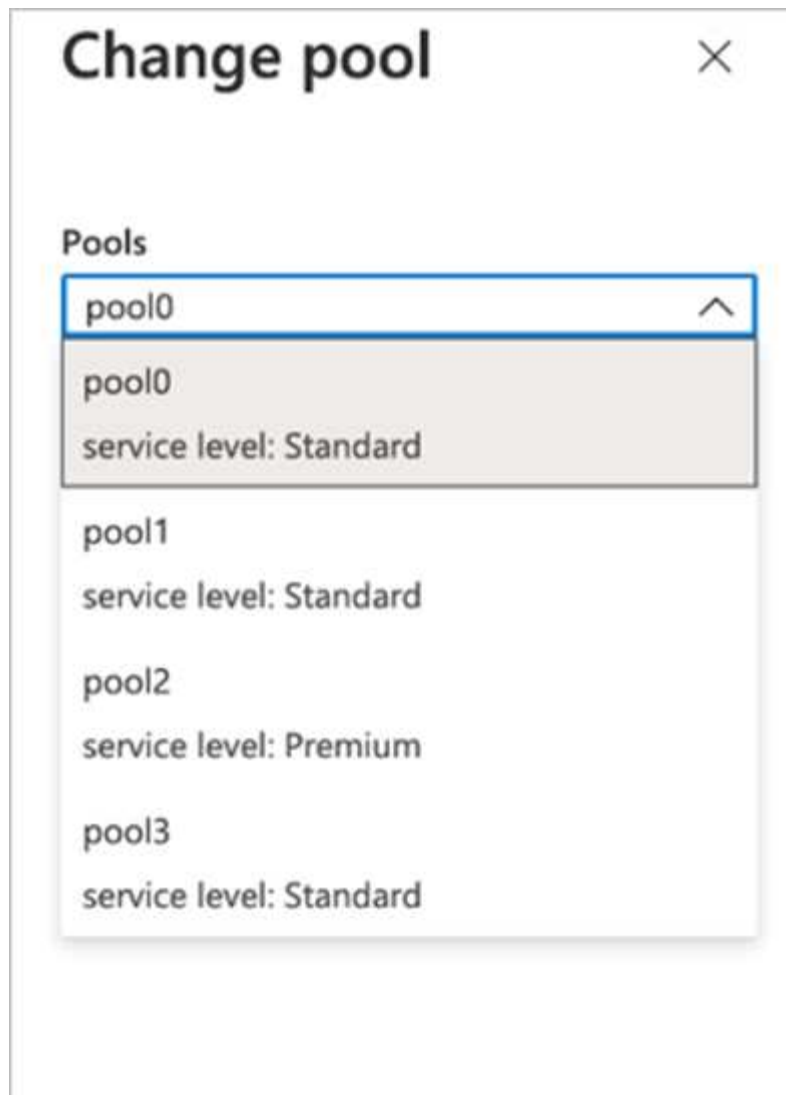
Para cambiar el nivel de servicio de un volumen, siga estos pasos:

1. En la página Volumes, haga clic con el botón derecho en el volumen cuyo nivel de servicio desea cambiar. Seleccione Cambiar pool.

NFSv3	10.28.254.4:/norootfor	Standard	pool0	...
NFSv4.1	NAS-735a.docs.lab:/for	Premium		
NFSv4.1	NAS-735a.docs.lab:/krt	Premium		
NFSv3	10.28.254.4:/moveme0	Premium		
NFSv3	10.28.254.4:/placeholder	Premium		

Resize
 Edit
 Change pool
 Delete

2. En la ventana Cambiar pool, seleccione el pool de capacidad al que desea mover el volumen. A continuación, haga clic en Aceptar.



### Automatizar el cambio de nivel de servicio

El cambio de nivel de servicio dinámico se encuentra actualmente en la vista previa pública, pero no está habilitado de forma predeterminada. Para activar esta función en la suscripción a Azure, siga estos pasos proporcionados en el documento “[Cambie dinámicamente el nivel de servicio de un volumen](#)”.

- También puede utilizar los siguientes comandos para Azure: CLI. Para obtener más información sobre cómo cambiar el tamaño del pool de Azure NetApp Files, visite ["volumen de archivos de az netapparchivo: Gestione los recursos de volúmenes del Azure NetApp Files \(ANF\)"](#).

```
az netappfiles volume pool-change -g mygroup
--account-name myacname
-pool-name mypoolname
--name myvolname
--new-pool-resource-id mynewresourceid
```

- La set- aznetappfilesvolumepool El cmdlet que se muestra aquí puede cambiar el pool de un volumen Azure NetApp Files. Para obtener más información sobre el cambio del tamaño del pool de volúmenes y Azure PowerShell, visite ["Cambiar el pool de un volumen Azure NetApp Files"](#).

```
Set-AzNetAppFilesVolumePool
-ResourceGroupName "MyRG"
-AccountName "MyAnfAccount"
-PoolName "MyAnfPool"
-Name "MyAnfVolume"
-NewPoolResourceId 7d6e4069-6c78-6c61-7bf6-c60968e45fbf
```

## Conclusión

NetApp and RUN: IA se ha asociado en la creación de este informe técnico para mostrar las funcionalidades únicas de la Azure NetApp Files junto con LA EJECUCIÓN: Plataforma de IA para simplificar la orquestación de las cargas de trabajo de IA. Este informe técnico proporciona una arquitectura de referencia para optimizar el proceso de canalizaciones de datos y orquestación de cargas de trabajo para la formación en la detección de carriles distribuidos.

En conclusión, con respecto a la formación distribuida a escala (especialmente en un entorno de cloud público), la orquestación de recursos y el componente de almacenamiento son una parte fundamental de la solución. Asegurarse de que la gestión de datos nunca afecta al procesamiento de varias GPU, por lo que se traduce en una utilización óptima de los ciclos de la GPU. Por lo tanto, hacer que el sistema sea lo más rentable posible para fines de capacitación distribuidos a gran escala.

El Data Fabric que ofrece NetApp supera el reto al permitir a los científicos e ingenieros de datos conectarse entre sí tanto en las instalaciones como en el cloud disponer de datos síncronos sin realizar ninguna intervención manual. En otras palabras, el tejido de datos suaviza el proceso de gestión del flujo de trabajo de IA expandiéndose por varias ubicaciones. También facilita la disponibilidad de datos basada en demanda al acercar los datos a la computación y realizar análisis, formación y validación donde y cuando sea necesario. Esta funcionalidad no solo permite la integración de datos, sino también la protección y seguridad de toda la canalización de datos.

## Información adicional

Si quiere más información sobre el contenido de este documento, consulte los siguientes documentos o sitios web:

- Conjunto de datos: TuSimple

["https://github.com/TuSimple/tusimple-benchmark/tree/master/doc/lane\\_detection"](https://github.com/TuSimple/tusimple-benchmark/tree/master/doc/lane_detection)

- Arquitectura de redes de aprendizaje profundo: Red neuronal convolucional espacial

["https://arxiv.org/abs/1712.06080"](https://arxiv.org/abs/1712.06080)

- Marco de entrenamiento de aprendizaje profundo distribuido: Horovod

["https://horovod.ai/"](https://horovod.ai/)

- EJECUCIÓN: Solución de orquestación de contenedores de IA: EJECUCIÓN: Introducción del producto de IA

["https://docs.run.ai/home/components/"](https://docs.run.ai/home/components/)

- EJECUCIÓN: Documentación de instalación de IA

["https://docs.run.ai/Administrator/Cluster-Setup/cluster-install/#step-3-install-runai"](https://docs.run.ai/Administrator/Cluster-Setup/cluster-install/#step-3-install-runai)

["https://docs.run.ai/Administrator/Researcher-Setup/cli-install/#runai-cli-installation"](https://docs.run.ai/Administrator/Researcher-Setup/cli-install/#runai-cli-installation)

- Envío de trabajos EN EJECUCIÓN: Interfaz de línea de comandos de IA

["https://docs.run.ai/Researcher/cli-reference/runai-submit/"](https://docs.run.ai/Researcher/cli-reference/runai-submit/)

["https://docs.run.ai/Researcher/cli-reference/runai-submit-mpi/"](https://docs.run.ai/Researcher/cli-reference/runai-submit-mpi/)

- Recursos de cloud para Azure: Azure NetApp Files

["https://docs.microsoft.com/azure/azure-netapp-files/"](https://docs.microsoft.com/azure/azure-netapp-files/)

- Servicio Azure Kubernetes

["https://azure.microsoft.com/services/kubernetes-service/-features"](https://azure.microsoft.com/services/kubernetes-service/-features)

- SKU de Azure VM

["https://azure.microsoft.com/services/virtual-machines/"](https://azure.microsoft.com/services/virtual-machines/)

- Azure VM con SKU de GPU

["https://docs.microsoft.com/azure/virtual-machines/sizes-gpu"](https://docs.microsoft.com/azure/virtual-machines/sizes-gpu)

- Trident de NetApp

["https://github.com/NetApp/trident/releases"](https://github.com/NetApp/trident/releases)

- Data Fabric con tecnología de NetApp

["https://www.netapp.com/data-fabric/what-is-data-fabric/"](https://www.netapp.com/data-fabric/what-is-data-fabric/)

- Documentación de productos de NetApp

["https://www.netapp.com/support-and-training/documentation/"](https://www.netapp.com/support-and-training/documentation/)

## **TR-4841: Sistema operativo de IA para el cloud híbrido con almacenamiento en caché de datos**

Rick Huang, David Arnette, NetApp Yochay Ettun, cnvrg.io

El crecimiento exponencial de los datos y el crecimiento exponencial DE LOS ML y la IA se han convergido para crear una economía de zettabytes con retos exclusivos de desarrollo e implementación.

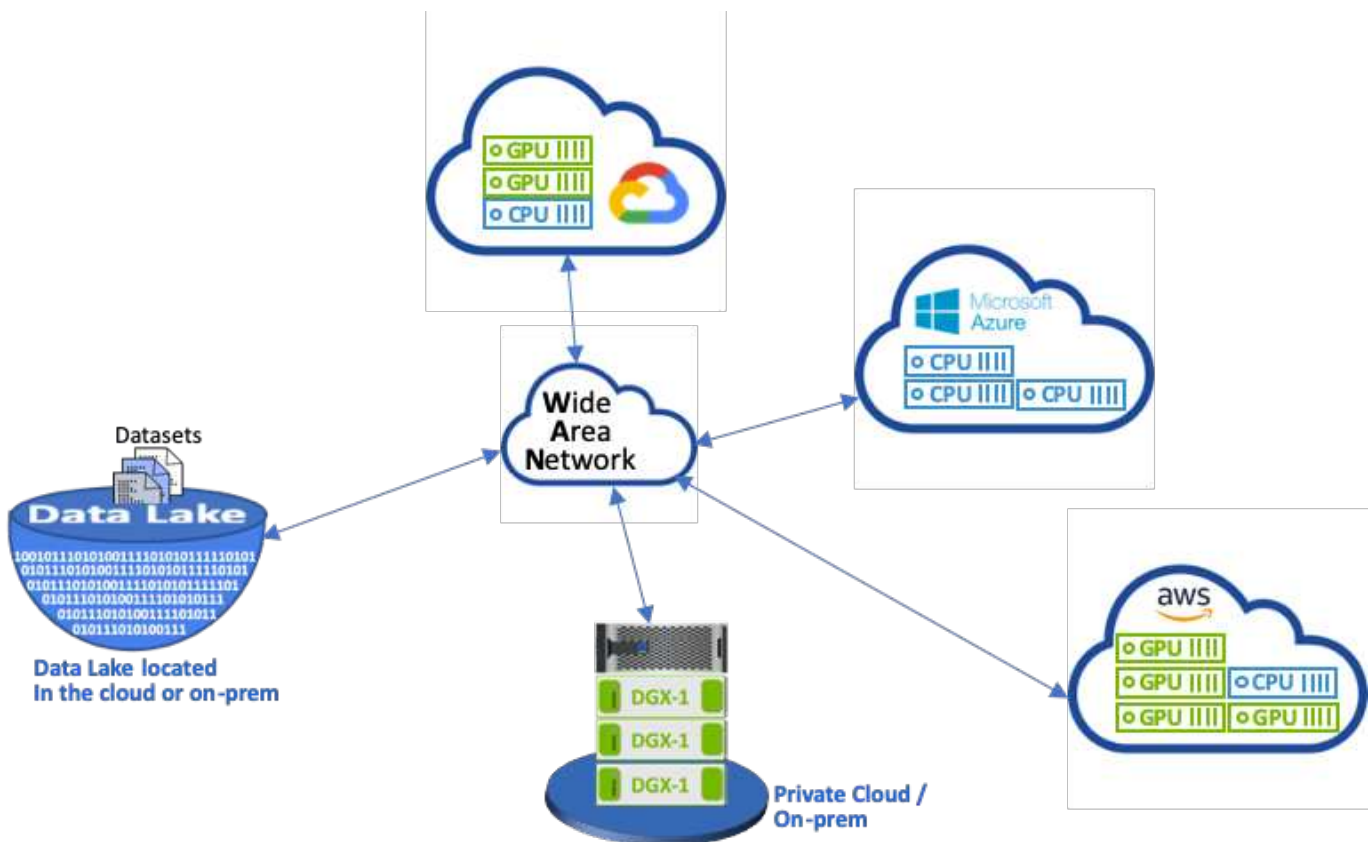
A pesar de que los modelos DE ML consumen gran cantidad de datos y requieren un almacenamiento de datos de alto rendimiento cerca de los recursos informáticos, en la práctica no es tan sencillo implementar este modelo, especialmente con instancias de cloud híbrido y computación elástica. Normalmente, se almacenan cantidades masivas de datos en lagos de datos de bajo coste, en los que los recursos informáticos de IA de alto rendimiento, como las GPU, no pueden acceder a ellos de forma eficiente. Este problema se

agrava en una infraestructura de cloud híbrido, en la que algunas cargas de trabajo funcionan en el cloud y otras están ubicadas en las instalaciones o en un entorno diferente de informática de alto rendimiento.

En este documento, presentamos una nueva solución que permite a los profesionales DE TECNOLOGÍA y a los ingenieros de datos crear una verdadera plataforma de IA de cloud híbrido con un concentrador de datos con topología que permite a los científicos de datos crear de forma instantánea y automática una caché de sus conjuntos de datos cerca de sus recursos informáticos, dondequiera que se encuentren. Como resultado, no solo se puede realizar el entrenamiento con modelos de alto rendimiento, sino que se obtienen ventajas adicionales, como la colaboración de varios profesionales de la IA, que tienen acceso inmediato a las cachés de los conjuntos de datos, las versiones y los linajes dentro de un concentrador de versiones del conjunto de datos.

### Descripción general de casos de uso y declaración de problemas

Los conjuntos de datos y las versiones de conjuntos de datos normalmente se encuentran en un lago de datos, como el almacenamiento basado en objetos StorageGRID de NetApp, que ofrece un coste reducido y otras ventajas operativas. Los científicos de datos extraen estos conjuntos de datos y los ingenieros en varios pasos para prepararlos para el entrenamiento con un modelo específico, a menudo creando varias versiones en el proceso. Al siguiente paso, el científico de datos debe seleccionar recursos informáticos optimizados (GPU, instancias de CPU de gama alta, un clúster local, etc.) para ejecutar el modelo. En la siguiente figura, se muestra la falta de proximidad del conjunto de datos en un entorno de computación DE ML.



Sin embargo, deben ejecutarse varios experimentos de entrenamiento en paralelo en distintos entornos informáticos, cada uno de los cuales requiere una descarga del conjunto de datos del lago de datos, lo cual es un proceso costoso y largo. No se garantiza la proximidad del conjunto de datos a entorno de computación

(especialmente para un cloud híbrido). Además, otros miembros del equipo que llevan a cabo sus propios experimentos con el mismo conjunto de datos deben pasar por el mismo arduo proceso. Más allá del obvio y lento acceso a los datos, los retos incluyen las dificultades para seguir las versiones de los conjuntos de datos, compartir conjuntos de datos, colaboración y reproducibilidad.

### Requisitos del cliente

Los requisitos del cliente pueden variar para lograr una ejecución DE ML de alto rendimiento mientras se usan los recursos de forma eficiente. Por ejemplo, los clientes pueden requerir lo siguiente:

- Acceso rápido a los conjuntos de datos desde cada instancia de computación que ejecuta el modelo de entrenamiento sin incurrir en costosas descargas y complejidades de acceso a los datos
- El uso de cualquier instancia de computación (GPU o CPU) en el cloud o en las instalaciones, sin preocuparse de la ubicación de los conjuntos de datos
- Aumenta la eficiencia y la productividad ejecutando múltiples experimentos de entrenamiento en paralelo con diferentes recursos informáticos en el mismo conjunto de datos sin retrasos innecesarios y latencia de los datos
- Costes mínimos de instancias de computación
- Reproducibilidad mejorada con herramientas para mantener registros de los conjuntos de datos, su linaje, versiones y otros detalles de metadatos
- Colaboración y uso compartido mejorados para que cualquier miembro autorizado del equipo pueda acceder a los conjuntos de datos y realizar experimentos

Para implementar el almacenamiento de conjuntos de datos en caché con el software de gestión de datos ONTAP de NetApp, los clientes deben realizar las siguientes tareas:

- Configure y establezca el almacenamiento NFS más cercano a los recursos de computación.
- Determinar qué conjunto de datos y versión almacenar en caché.
- Supervise la memoria total comprometida con los conjuntos de datos almacenados en caché y cuánto almacenamiento NFS está disponible para las confirmaciones adicionales de la caché (por ejemplo, gestión de caché).
- La caducidad de los conjuntos de datos de la caché si no se han utilizado en un momento determinado. El valor predeterminado es un día; hay otras opciones de configuración disponibles.

### Descripción general de la solución

En esta sección se revisa una canalización de ciencia de datos convencional y sus inconvenientes. También presenta la arquitectura de la solución de almacenamiento en caché del conjunto de datos propuesta.

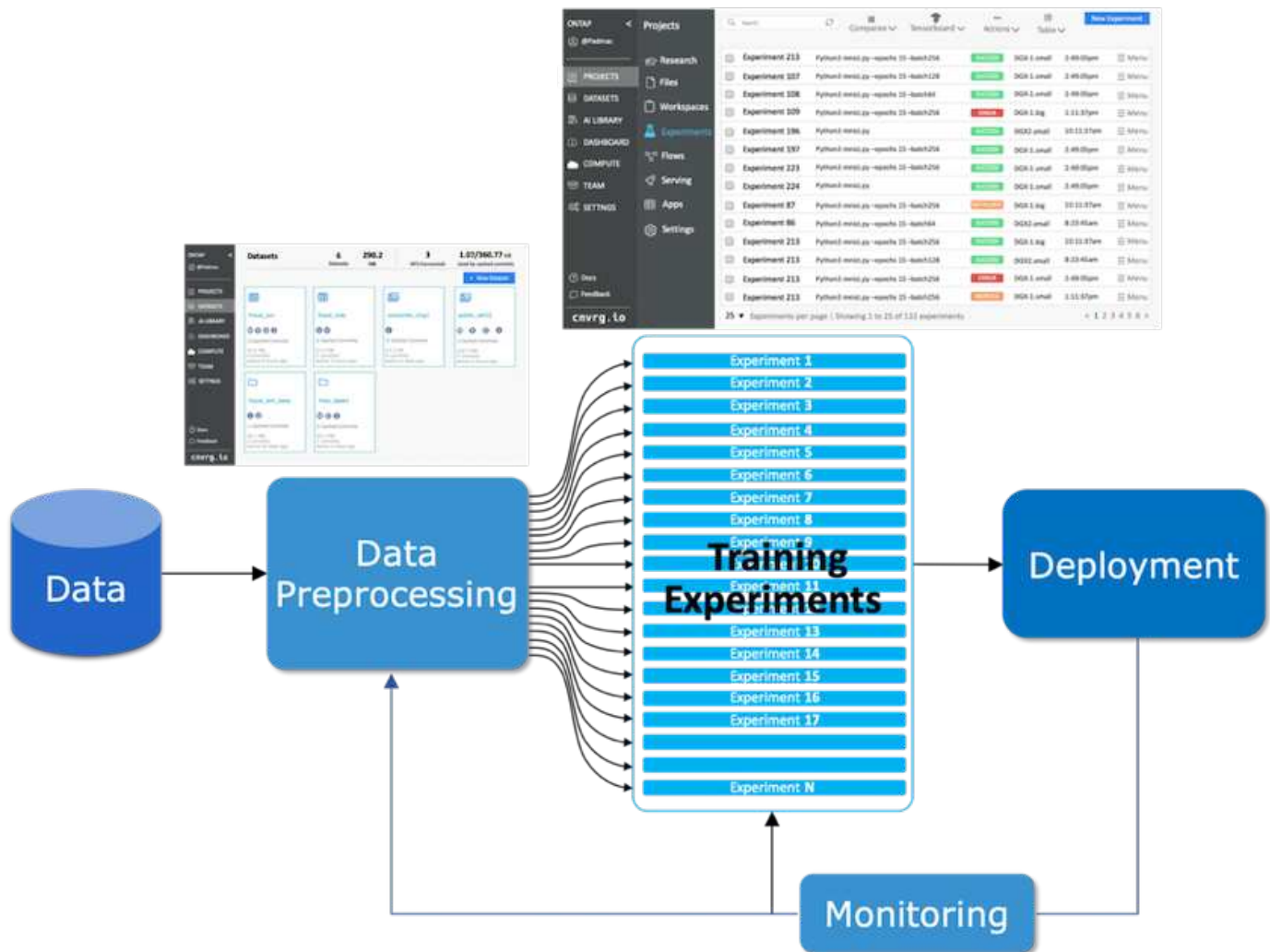
#### Canalización convencional de ciencia de datos e inconvenientes

Una secuencia típica de desarrollo e implementación de modelos ML implica pasos iterativos que incluyen lo siguiente:

- Ingestión de datos
- Procesamiento previo de los datos (crear varias versiones de los conjuntos de datos)
- Ejecución de múltiples experimentos que implican la optimización de hiperparámetros, diferentes modelos, etc.

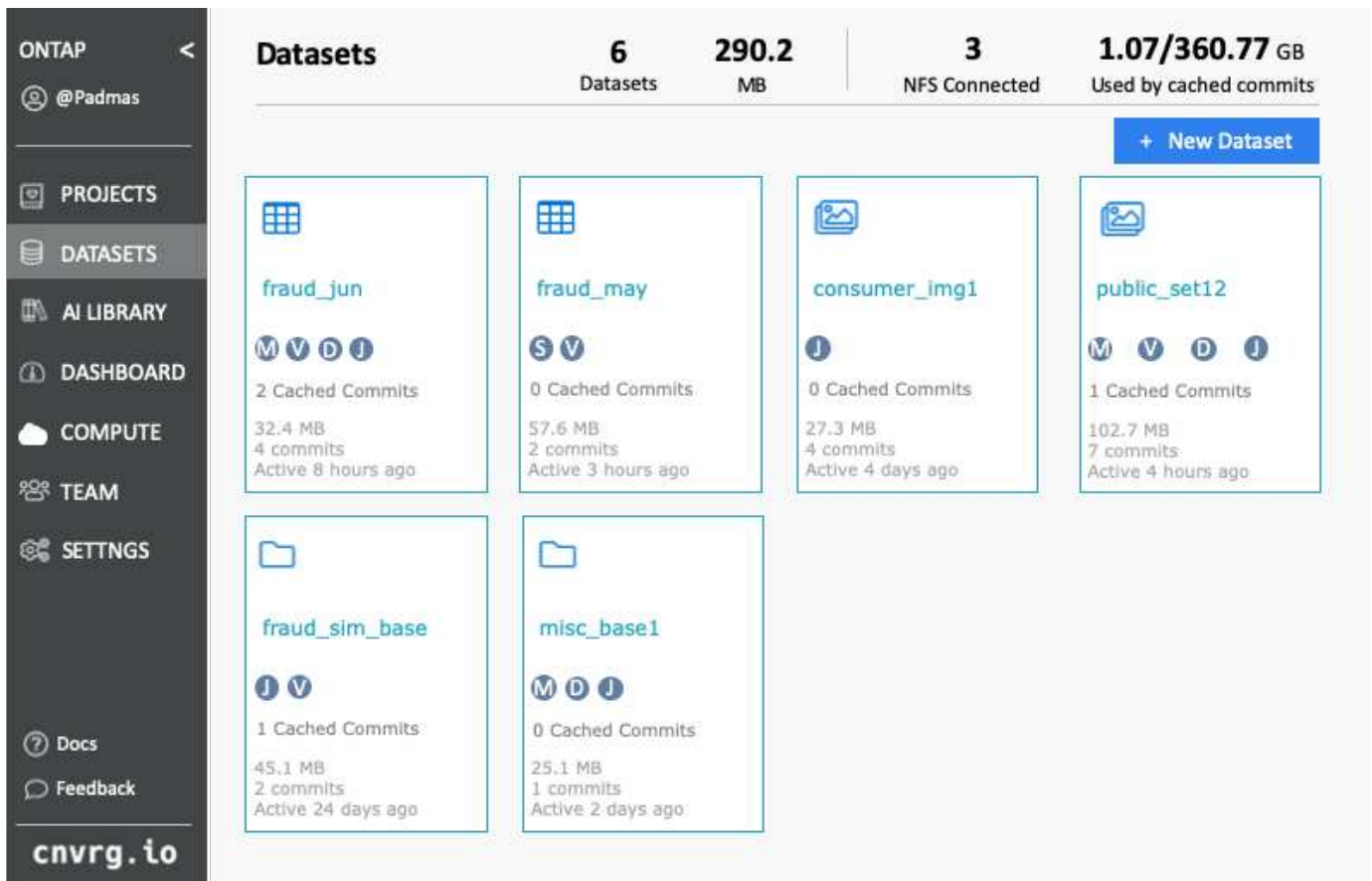


- Puesta en marcha
- Monitoringcvrg.io ha desarrollado una plataforma completa para automatizar todas las tareas, desde la investigación hasta la implementación. En la siguiente figura se muestra una pequeña muestra de capturas de pantalla del panel relativas a la canalización.



Es muy común tener múltiples conjuntos de datos en juego desde los repositorios públicos y los datos privados. Además, es probable que cada conjunto de datos tenga varias versiones como resultado de la limpieza del conjunto de datos o la ingeniería de funciones. Se necesita un panel que proporcione un concentrador de conjuntos de datos y un concentrador de versiones para garantizar que las herramientas de colaboración y consistencia estén disponibles para el equipo, como se puede ver en la siguiente figura.





El siguiente paso de la canalización es el entrenamiento, que requiere varias instancias paralelas de modelos de entrenamiento, cada uno asociado con un conjunto de datos y una determinada instancia de computación. La vinculación de un conjunto de datos con un cierto experimento con una determinada instancia informática es un reto, ya que es posible que algunos experimentos se realicen mediante instancias de GPU de Amazon Web Services (AWS), mientras que otras instancias de DGX-1 o DGX-2 se llevan a cabo en las instalaciones. Es posible que se ejecuten otros experimentos en servidores de CPU en GCP, mientras que la ubicación del conjunto de datos no está muy cerca de los recursos informáticos que realizan el entrenamiento. Una proximidad razonable tendría conectividad de 10 GbE completa o más baja latencia desde el almacenamiento del conjunto de datos a la instancia de computación.

Es una práctica común que los científicos de datos descarguen el conjunto de datos a la instancia de computación que realice el entrenamiento y ejecute el experimento. Sin embargo, este enfoque presenta varios problemas posibles:

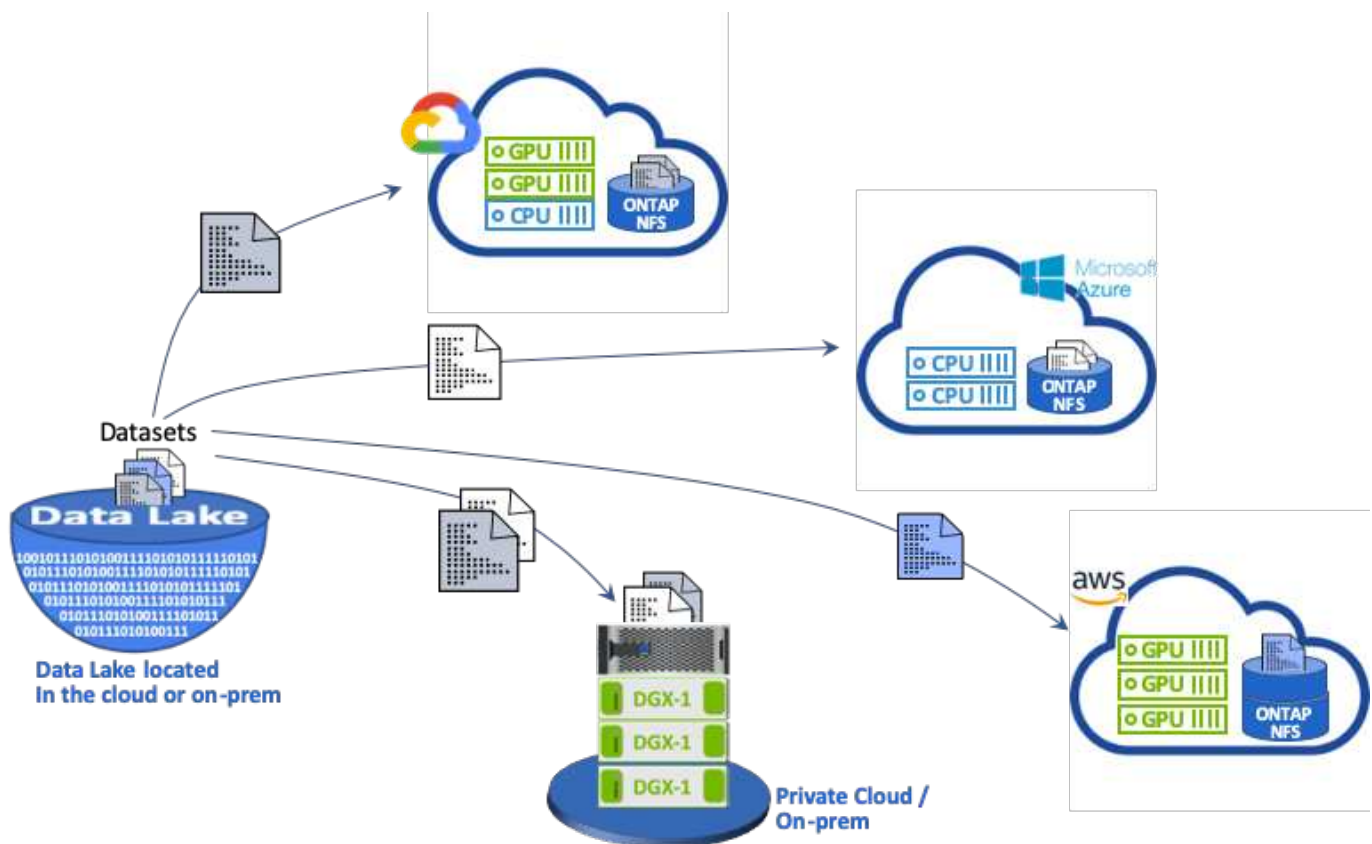
- Cuando el científico de datos descarga el conjunto de datos en una instancia informática, no hay garantías de que el almacenamiento informático integrado tenga un alto rendimiento (un ejemplo de un sistema de alto rendimiento sería la solución NVMe A800 de ONTAP AFF).
- Cuando el conjunto de datos descargado reside en un nodo de computación, el almacenamiento puede convertirse en un cuello de botella cuando se ejecutan los modelos distribuidos en varios nodos (a diferencia del almacenamiento distribuido de alto rendimiento de ONTAP de NetApp).
- La siguiente iteración del experimento de entrenamiento podría realizarse en una instancia de computación diferente debido a conflictos de cola o prioridades, creando de nuevo una distancia significativa de la red desde el conjunto de datos hasta la ubicación de computación.
- Otros miembros del equipo que ejecutan experimentos de entrenamiento en el mismo clúster informático no pueden compartir este conjunto de datos; cada uno realiza la descarga (costosa) del conjunto de datos desde una ubicación arbitraria.

- Si se necesitan otros conjuntos de datos o versiones del mismo conjunto de datos para las siguientes tareas de entrenamiento, los científicos de datos deben volver a realizar la descarga (costosa) del conjunto de datos en la instancia informática que realiza la versión training. NetApp y cnvrg.io han creado una nueva solución de almacenamiento en caché de conjuntos de datos que elimina estos obstáculos. Esta solución crea una ejecución acelerada de la canalización DE ML mediante el almacenamiento en caché de conjuntos de datos activos en el sistema de almacenamiento de alto rendimiento de ONTAP. Con NFS de ONTAP, los conjuntos de datos se almacenan en caché una vez (y solo una) en una estructura de datos con tecnología de NetApp (como AFF A800), que se encuentra junto con el sistema informático. Dado que el almacenamiento de alta velocidad NFS de ONTAP de NetApp puede dar servicio a nodos de computación DE varios ML, el rendimiento de los modelos de entrenamiento se optimiza, lo que permite ahorrar costes, productividad y eficiencia operativa a la organización.

### Arquitectura de la solución

Esta solución de NetApp y cnvrg.io proporciona almacenamiento en caché de conjuntos de datos, como se muestra en la siguiente figura. El almacenamiento en caché de conjuntos de datos permite a los científicos de datos elegir la versión deseada de conjuntos de datos o conjuntos de datos y moverlos a la caché NFS de ONTAP, que se encuentra cerca del clúster de computación DE ML. Ahora, el científico de datos puede realizar varios experimentos sin incurrir en retrasos ni descargas. Además, todos los ingenieros de colaboración pueden utilizar el mismo conjunto de datos con el clúster de computación conectado (con la libertad de elegir cualquier nodo) sin descargar adicionales del lago de datos. A los científicos de datos se les ofrece un panel que realiza un seguimiento y supervisa todos los conjuntos de datos y versiones y ofrece una vista de los conjuntos de datos que se almacenan en caché.

La plataforma cnvrg.io detecta automáticamente conjuntos de datos antiguos que no se han utilizado durante un cierto tiempo y los desaloja de la caché, por lo que mantiene un espacio libre de la caché NFS para conjuntos de datos utilizados con mayor frecuencia. Es importante tener en cuenta que el almacenamiento en caché del conjunto de datos con ONTAP funciona en el cloud y en las instalaciones, por lo que proporciona la máxima flexibilidad.



## Conceptos y componentes

En esta sección se tratan conceptos y componentes asociados al almacenamiento en caché de datos en un flujo DE trabajo DE ML.

### Aprendizaje automático

LA ML se está convirtiendo rápidamente en un factor esencial para muchas empresas y organizaciones de todo el mundo. Por lo tanto, los equipos DE TECNOLOGÍA y DevOps se enfrentan ahora al reto de estandarizar cargas de trabajo DE ML y aprovisionar cloud, recursos informáticos en las instalaciones y recursos informáticos híbridos que dan soporte a los flujos de trabajo dinámicos e intensivos que requieren las tareas de ML y las canalizaciones.

### Aprendizaje automático basado en contenedores y Kubernetes

Los contenedores son instancias aisladas del espacio de usuario que se ejecutan sobre un kernel de sistema operativo host compartido. La adopción de contenedores está aumentando rápidamente. Los contenedores ofrecen muchos de los mismos beneficios de uso de pruebas de espacio que las máquinas virtuales (VM). Sin embargo, debido a que se eliminan las capas de hipervisor y de sistema operativo «guest» de las que dependen las máquinas virtuales, los contenedores son mucho más ligeros.

Los contenedores también permiten el empaquetado eficiente de dependencias de aplicaciones, tiempos de ejecución, etc. directamente en una aplicación. El formato de embalaje de contenedor más utilizado es el contenedor Docker. Una aplicación que se haya contenedor en el formato de contenedor Docker se puede ejecutar en cualquier máquina que pueda ejecutar contenedores Docker. Esto es cierto incluso si las dependencias de la aplicación no están presentes en la máquina, porque todas las dependencias están empaquetadas en el propio contenedor. Para obtener más información, visite la ["Sitio web de Docker"](#).

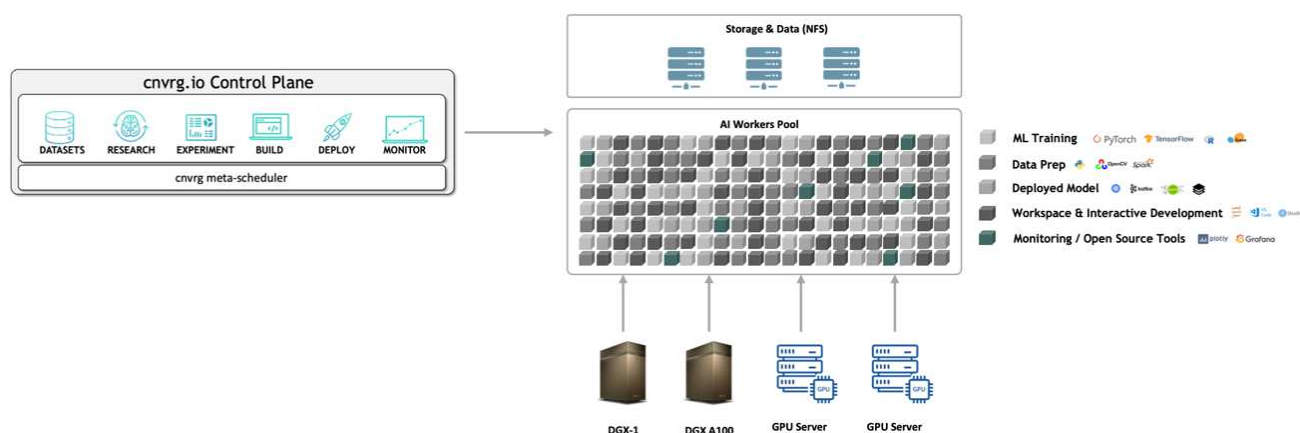
Kubernetes, el popular orquestador de contenedores, permite a los científicos de datos lanzar trabajos y canalizaciones flexibles basados en contenedores. También permite a los equipos de infraestructura gestionar y supervisar cargas de trabajo DE ML en un único entorno gestionado y nativo del cloud. Para obtener más información, visite la ["Sitio web de Kubernetes"](#).

### cnvrg.io

Cnvrg.io es un sistema operativo de IA que transforma la forma en la que las empresas gestionan, escalan y aceleran la IA y el desarrollo científico de datos de la investigación a la producción. La plataforma de código primero está desarrollada por científicos de datos para científicos de datos y ofrece flexibilidad para ejecutarse en las instalaciones o en el cloud. Gracias a la gestión de modelos, MLOPS y soluciones continuas DE ML, cnvrg.io aporta una tecnología de primera línea a los equipos de ciencia de datos para que puedan dedicar menos tiempo a DevOps y centrarse en la auténtica magia, gracias a los algoritmos. Desde que utiliza cnvrg.io, los equipos de distintos sectores han obtenido más modelos de producción, lo que da como resultado un aumento del valor empresarial.

### Cnvrg.io Meta-Scheduler

cnvrg. i/o tiene una arquitectura única que permite A LOS DEPARTAMENTOS DE TECNOLOGÍA e ingenieros conectar distintos recursos informáticos al mismo plano de control y que cnvrg.io gestiona tareas DE ML en todos los recursos. Esto significa que puede conectar varios clústeres de Kubernetes en las instalaciones, servidores de VM y cuentas de cloud, y ejecutar cargas de trabajo DE ML en todos los recursos, como se muestra en la siguiente figura.



## Cnvrg.io almacenamiento en caché de datos

cnvrg.io permite a los científicos de datos definir versiones de conjuntos de datos calientes y fríos con su tecnología de almacenamiento en caché de datos. De forma predeterminada, los conjuntos de datos se almacenan en una base de datos de almacenamiento de objetos centralizada. A continuación, los científicos de datos pueden almacenar en caché una versión de datos específica en el recurso de computación seleccionado para ahorrar tiempo en la descarga y, por tanto, aumentar EL desarrollo DE ML y la productividad. Los conjuntos de datos que se almacenan en la caché y no se utilizan durante unos días se borran automáticamente del NFS seleccionado. El almacenamiento en caché y el borrado de la caché se pueden realizar con un solo clic; no se requiere codificación, NI trabajo de DevOps.

## Cnvrg.io fluye y canalizaciones ML

Cnvrg.io fluye es una herramienta para construir tuberías DE PRODUCCIÓN ML. Cada componente de un flujo es un script/código que se ejecuta en una computación seleccionada con una imagen de Docker base. Este diseño permite a los científicos e ingenieros de datos crear una única canalización que puede ejecutar tanto en las instalaciones como en el cloud. cnvrg.io garantiza que los datos, los parámetros y los artefactos se mueven entre los diferentes componentes. Además, se supervisa y se sigue cada flujo para obtener ciencia de datos reproducibles al 100%.

## NÚCLEO cnvrg.io

El NÚCLEO cnvrg.io es una plataforma gratuita para que la comunidad de ciencia de datos pueda ayudar a los científicos de datos a centrarse más en la ciencia de datos y menos en DevOps. La infraestructura flexible DE CORE aporta a los científicos de datos el control de usar cualquier idioma, marco de IA o entorno informático, ya sea en las instalaciones o en el cloud, para poder hacer lo que mejor hacen o crear algoritmos. El NÚCLEO cnvrg.io se puede instalar fácilmente con un único comando en cualquier clúster de Kubernetes.

## ONTAP AI de NetApp

ONTAP AI es una arquitectura de referencia de centro de datos para cargas de trabajo DE APRENDIZAJE profundo (DL) y ML que utiliza sistemas de almacenamiento AFF de NetApp y sistemas DGX de NVIDIA con GPU Tesla V100. ONTAP AI se basa en el protocolo de archivos NFS estándar del sector en Ethernet de 100 GB y proporciona a los clientes una infraestructura DE APRENDIZAJE PROFUNDO DE alto rendimiento QUE utiliza tecnologías estándar para el centro de datos para reducir los gastos generales de implementación y administración. Con una red y protocolos estandarizados, ONTAP AI se integra en entornos de cloud híbrido a la vez que mantiene la coherencia y la simplicidad operativas. Como solución de infraestructura prevalidada, ONTAP AI reduce el tiempo y el riesgo de la puesta en marcha y la sobrecarga de la administración de forma significativa, lo que permite a los clientes lograr una rentabilidad de la inversión más rápida.

## DeepOps de NVIDIA

DeepOps es un proyecto de código abierto de NVIDIA que, con Ansible, automatiza la puesta en marcha de clústeres de servidores de GPU de acuerdo con las prácticas recomendadas. DeepOps es modular y se puede utilizar para realizar varias tareas de puesta en marcha. En este documento y en el ejercicio de validación descrito, DeepOps se utiliza para poner en marcha un clúster de Kubernetes que consta de nodos de trabajo de servidor GPU. Para obtener más información, visite la ["Sitio web DeepOps"](#).

## Trident de NetApp

Trident es un orquestador de almacenamiento de código abierto desarrollado y mantenido por NetApp que simplifica en gran medida la creación, la gestión y el consumo de almacenamiento persistente para cargas de trabajo de Kubernetes. Trident, en sí misma, una aplicación nativa de Kubernetes, se ejecuta directamente en un clúster de Kubernetes. Con Trident, los usuarios de Kubernetes (desarrolladores, científicos de datos, administradores de Kubernetes, etc.) pueden crear, gestionar e interactuar con volúmenes de almacenamiento persistente en el formato Kubernetes estándar, con el que ya están familiarizados. Al mismo tiempo, pueden aprovechar las funciones avanzadas de gestión de datos de NetApp y un Data Fabric con tecnología de NetApp. Trident elimina las complejidades del almacenamiento persistente y facilita el consumo. Para obtener más información, visite la ["Sitio web de Trident"](#).

## StorageGRID de NetApp

StorageGRID de NetApp es una plataforma de almacenamiento de objetos definida por software diseñada para satisfacer estas necesidades proporcionando un almacenamiento sencillo y similar al cloud a el que los usuarios pueden acceder mediante el protocolo S3. StorageGRID es un sistema de escalado horizontal diseñado para admitir varios nodos en sitios conectados a Internet, independientemente de la distancia. Con el motor de políticas inteligente de StorageGRID, los usuarios pueden elegir objetos de codificación de borrado en todos los sitios para lograr resiliencia geográfica o replicación de objetos entre sitios remotos para minimizar la latencia de acceso WAN. StorageGRID proporciona un excelente lago de datos de almacenamiento de objetos primarios de cloud privado en esta solución.

## Cloud Volumes ONTAP de NetApp

El software de gestión de datos Cloud Volumes ONTAP de NetApp proporciona control, protección y eficiencia para los datos de usuarios con la flexibilidad de proveedores de cloud público como AWS, Google Cloud Platform y Microsoft Azure. Cloud Volumes ONTAP es un software para la gestión de datos nativo del cloud, integrado en el software de almacenamiento ONTAP de NetApp, que proporciona a los usuarios una plataforma de almacenamiento universal superior que cubre sus necesidades de datos en el cloud. Disponer de un mismo software de almacenamiento en el cloud y en las instalaciones proporciona a los usuarios el valor de una estructura de datos sin necesidad de formar al personal INFORMÁTICO en todos los métodos nuevos para gestionar los datos.

Para los clientes interesados en modelos de puesta en marcha de cloud híbrido, Cloud Volumes ONTAP puede proporcionar las mismas funcionalidades y un rendimiento líder en la mayoría de clouds públicos para proporcionar una experiencia de usuario fluida y coherente en cualquier entorno.

## Requisitos de hardware y software

En esta sección se tratan los requisitos tecnológicos de la solución de IA de ONTAP.

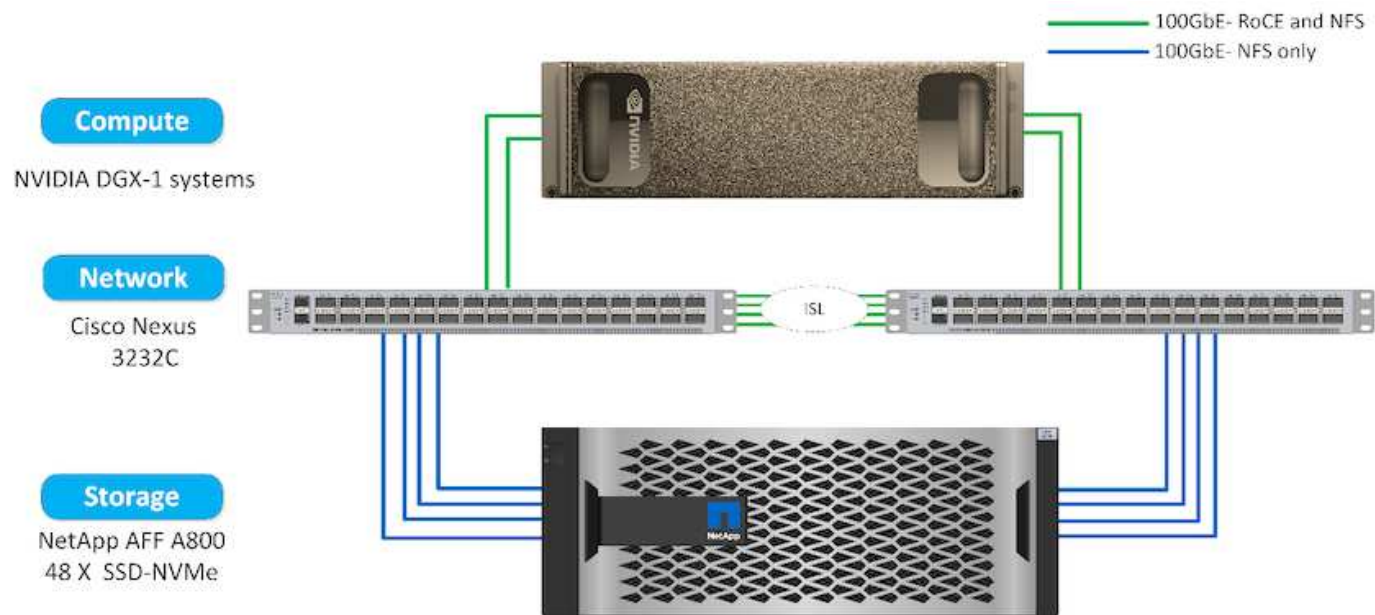
### Requisitos de hardware

Aunque los requisitos de hardware dependen de cargas de trabajo específicas de los clientes, es posible poner en marcha IA de ONTAP a cualquier escala para la ingeniería de datos, el entrenamiento de modelos y la inferencia de producción desde una única GPU hasta configuraciones a escala de rack para operaciones



DE APRENDIZAJE AUTOMÁTICO o aprendizaje profundo a gran escala. Si quiere más información sobre ONTAP AI, consulte ["Sitio web de IA de ONTAP"](#).

Esta solución se validó utilizando un sistema DGX-1 para computación, un sistema de almacenamiento AFF A800 de NetApp y Cisco Nexus 3232C para conectividad de red. El AFF A800 utilizado en esta validación puede admitir hasta 10 sistemas DGX-1 para la mayoría de cargas de trabajo DE ML/DL. En la siguiente figura se muestra la topología ONTAP AI utilizada para entrenar el modelo en esta validación.



Para ampliar esta solución a un cloud público, Cloud Volumes ONTAP puede ponerse en marcha junto con los recursos informáticos de la GPU del cloud e integrarse en un tejido de datos del cloud híbrido que permita a los clientes utilizar los recursos que sean apropiados para una carga de trabajo en concreto.

Requisitos de software

En la siguiente tabla se muestran las versiones de software específicas que se utilizan en la validación de esta solución.

Componente	Versión
Ubuntu	18.04.4 LTS
SO DGX DE NVIDIA	4.4.0
DeepOps de NVIDIA	20.02.1
Kubernetes	1.15
Timón	3.1.0
cnvrg.io	3.0.0
ONTAP de NetApp	9.6P4

Para esta validación de soluciones, se puso en marcha Kubernetes como un clúster de un solo nodo en el sistema DGX-1. Para las puestas en marcha a gran escala, se deben poner en marcha nodos maestros de Kubernetes independientes para proporcionar una alta disponibilidad de los servicios de gestión y reservar valiosos recursos DGX para las cargas de trabajo DE APRENDIZAJE AUTOMÁTICO y aprendizaje profundo.

## Detalles de validación y puesta en marcha de la solución

En las siguientes secciones se tratan los detalles de la validación y puesta en marcha de la solución.

### Puesta en marcha de IA de ONTAP

La puesta en marcha de ONTAP AI requiere la instalación y configuración de hardware de almacenamiento, computación y redes. Este documento no cubre las instrucciones específicas para la puesta en marcha de la infraestructura de IA de ONTAP. Para obtener información detallada sobre la implementación, consulte ["NVA-1121-PUESTA en MARCHA: ONTAP AI de NetApp, con tecnología NVIDIA"](#).

Para esta validación de soluciones, se creó y montó un solo volumen en el sistema DGX-1. A continuación, ese punto de montaje se montó en los contenedores para que los datos sean accesibles para el entrenamiento. Para puestas en marcha a gran escala, Trident de NetApp automatiza la creación y el montaje de volúmenes para eliminar la sobrecarga administrativa y permitir la gestión de recursos por parte del usuario final.

### Puesta en marcha de Kubernetes

Para poner en marcha y configurar su clúster de Kubernetes con NVIDIA DeepOps, realice las siguientes tareas desde un host de salto de implementación:

1. Descargue NVIDIA DeepOps siguiendo las instrucciones de ["Página de Inicio"](#) En el sitio de NVIDIA DeepOps GitHub.
2. Implemente Kubernetes en su clúster siguiendo las instrucciones que se indican en ["Guía de puesta en marcha de Kubernetes"](#) En el sitio de NVIDIA DeepOps GitHub.



Para que la puesta en marcha de DeepOps Kubernetes funcione, debe haber el mismo usuario en todos los nodos maestro y de trabajador de Kubernetes.

Si la implementación falla, cambie el valor de `kubect1_localhost` a falso in `deepops/config/group_vars/k8s-cluster.yml` y repita el paso 2. La Copy `kubect1 binary to ansible host` tarea, que sólo se ejecuta cuando el valor de `kubect1_localhost` Es cierto que se basa en el módulo `Fetch Ansible`, que tiene problemas conocidos con el uso de memoria. Estos problemas de uso de la memoria a veces pueden provocar un error en la tarea. Si la tarea falla debido a un problema de memoria, el resto de la operación de implementación no se completa correctamente.

Si la implementación se completa correctamente después de cambiar el valor de `kubect1_localhost` para `false`, a continuación, debe copiar manualmente la `kubect1 binary` Desde un nodo maestro de Kubernetes a un host de salto de implementación. Puede encontrar la ubicación del `kubect1 binary` en un nodo maestro específico ejecutando el `which kubect1` comando directamente en ese nodo.

### Puesta en marcha de cnvrg.io

#### Ponga en marcha EL NÚCLEO cnvrg mediante Helm

Timón es la forma más sencilla de poner en marcha cnvrg rápidamente mediante cualquier clúster, en las instalaciones, MiniClube o en cualquier clúster de cloud (como AKS, EKS y GKE). En esta sección se describe cómo se instaló cnvrg en una instancia de las instalaciones (DGX-1) con Kubernetes instalado.

## Requisitos previos

Antes de completar la instalación, debe instalar y preparar las siguientes dependencias en su equipo local:

- Kubectl
- Timón 3.x
- Clúster de Kubernetes 1.15 o posterior

## Despliegue mediante Helm

1. Para descargar los gráficos cnvrg del timón más actualizados, ejecute el siguiente comando:

```
helm repo add cnvrg https://helm.cnvrg.io
helm repo update
```

2. Antes de poner en marcha cnvrg, necesita la dirección IP externa del clúster y el nombre del nodo en el que se va a implementar cnvrg. Para poner en marcha cnvrg en un clúster de Kubernetes en las instalaciones, ejecute el siguiente comando:

```
helm install cnvrg cnvrg/cnvrg --timeout 1500s --wait \ --set
global.external_ip=<ip_of_cluster> \ --set global.node=<name_of_node>
```

3. Ejecute el `helm install` comando. Todos los servicios y sistemas se instalan automáticamente en el clúster. El proceso puede tardar hasta 15 minutos.
4. La `helm install` command puede tardar hasta 10 minutos. Cuando finalice la implementación, vaya a la dirección URL de su cnvrg recién implementado o agregue el nuevo clúster como un recurso dentro de su organización. La `helm` Command le informa de la URL correcta.

```
Thank you for installing cnvrg.io!
Your installation of cnvrg.io is now available, and can be reached via:
Talk to our team via email at
```

5. Cuando el estado de todos los contenedores se ejecuta o se completa, cnvrg se ha implementado correctamente. Debería ser similar a la siguiente salida de ejemplo:



NAME	READY	STATUS	RESTARTS	AGE	
cnvrg-app-69fbb9df98-6xrgf		1/1	Running	0	2m
cnvrg-sidekiq-b9d54d889-5x4fc		1/1	Running	0	2m
controller-65895b47d4-s96v6		1/1	Running	0	2m
init-app-vs-config-wv9c4		0/1	Completed	0	9m
init-gateway-vs-config-2zbpp		0/1	Completed	0	9m
init-minio-vs-config-cd2rg		0/1	Completed	0	9m
minio-0		1/1	Running	0	2m
postgres-0		1/1	Running	0	2m
redis-695c49c986-kcvt9		1/1	Running	0	2m
seeder-wh655		0/1	Completed	0	2m
speaker-5sghr		1/1	Running	0	2m

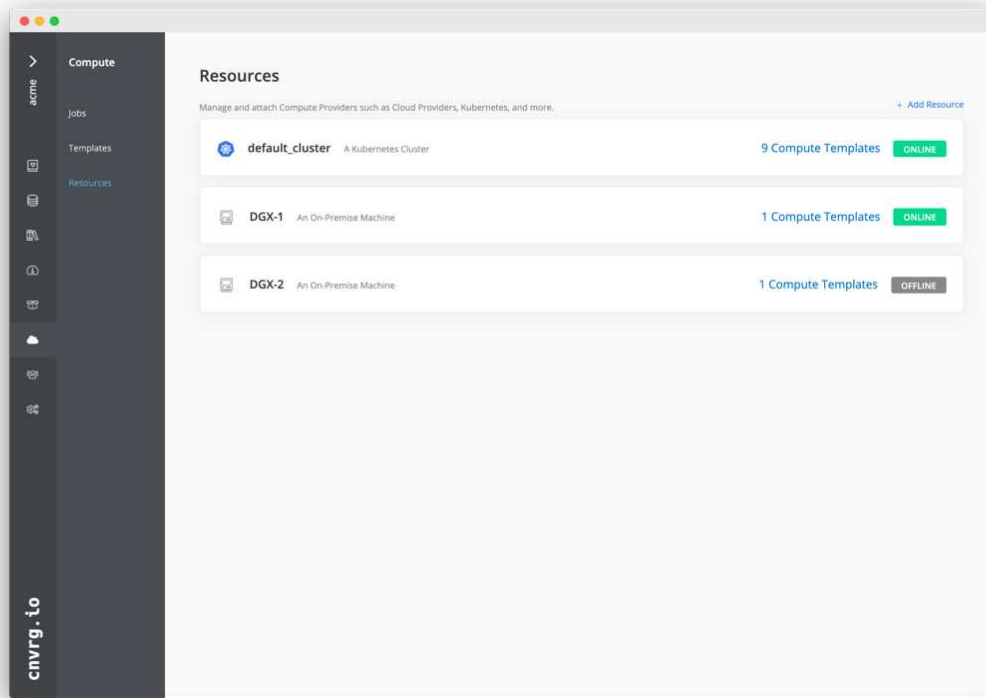
## Formación de modelos de visión computarizada con ResNet50 y el conjunto de datos de rayos X torácicos

Cnvrg.io AI OS se puso en marcha en una configuración de Kubernetes en una arquitectura de IA ONTAP de NetApp basada en el sistema NVIDIA DGX. Para la validación, se utilizó el conjunto de datos de radiografía de tórax de los NIH que consistía en imágenes desidentificadas de radiografías de tórax. Las imágenes estaban en formato PNG. Los datos fueron proporcionados por el Centro clínico de los NIH y están disponibles a través del ["Sitio de descarga de NIH"](#). Utilizamos una muestra de 250 GB de los datos con 627, 615 imágenes en 15 clases.

El conjunto de datos se cargó en la plataforma cnvrg y se almacenó en caché en una exportación NFS desde el sistema de almacenamiento AFF A800 de NetApp.

## Configure los recursos de computación

La arquitectura cnvrg y la funcionalidad de programación de metadatos permiten que los ingenieros y profesionales DE TECNOLOGÍA adjunte diferentes recursos informáticos a una única plataforma. En nuestra configuración, utilizamos el mismo cnvrg de clúster que se puso en marcha para ejecutar cargas de trabajo de aprendizaje profundo. Si necesita conectar clústeres adicionales, utilice la GUI, tal como se muestra en la siguiente captura de pantalla.

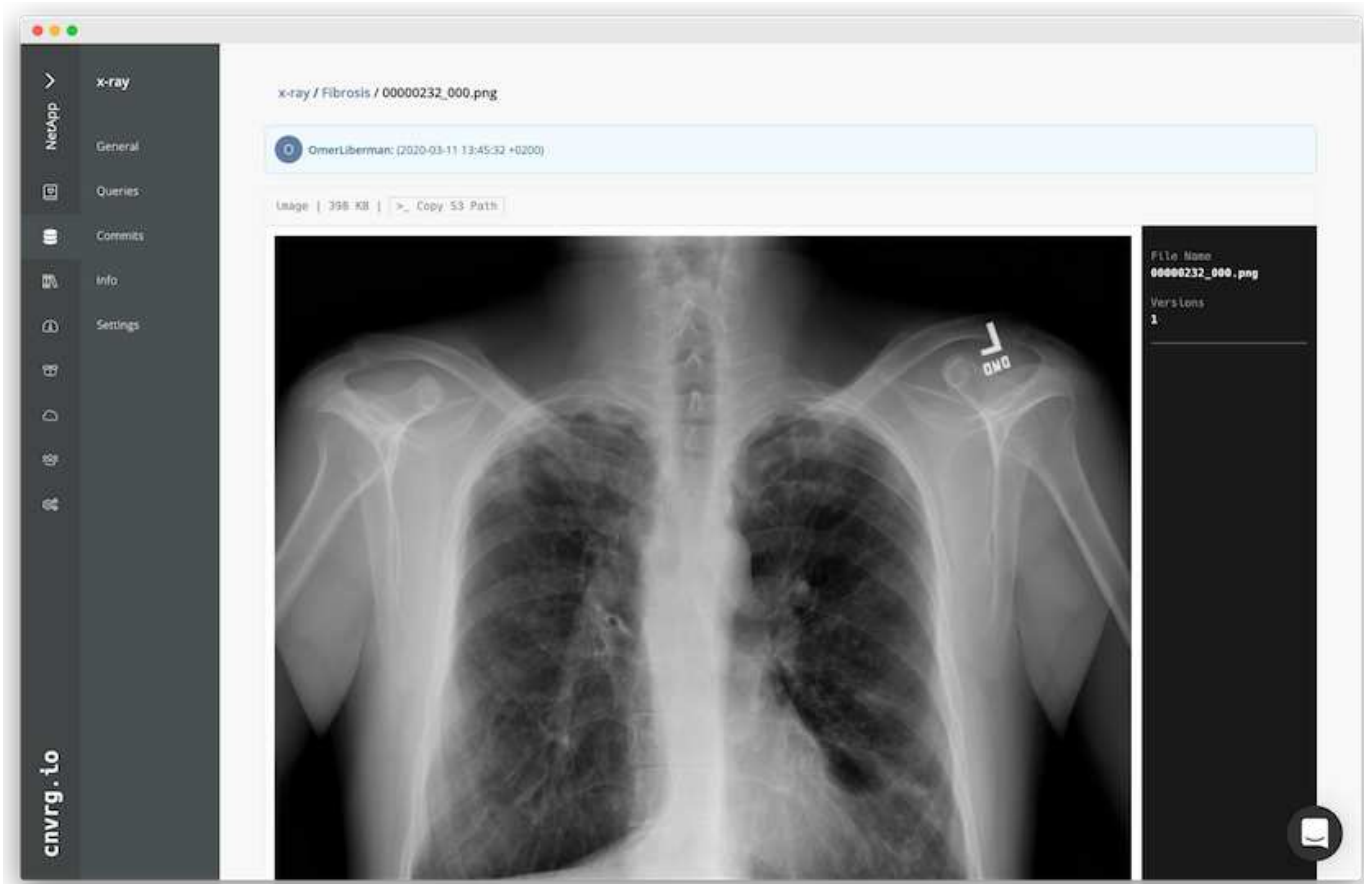


## Cargar datos

Para cargar datos en la plataforma cnvrg, puede utilizar la interfaz gráfica de usuario o la interfaz de línea de comandos cnvrg. En el caso de conjuntos de datos de gran tamaño, NetApp recomienda el uso de la CLI porque es una herramienta sólida, escalable y fiable capaz de gestionar un gran número de archivos.

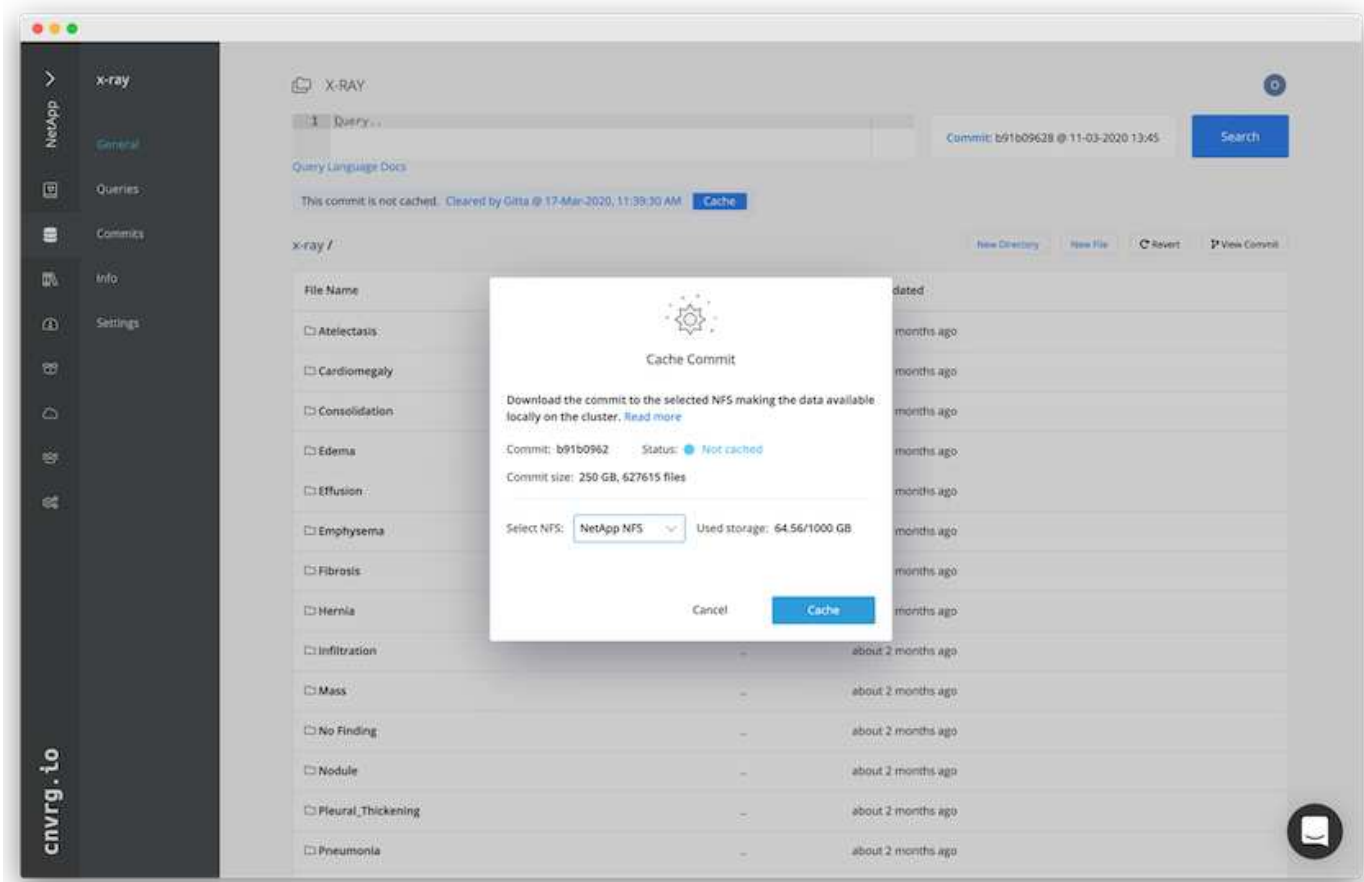
Para cargar datos, realice los siguientes pasos:

1. Descargue el "[CLI de cnvrg](#)".
2. desplácese hasta el directorio de rayos x.
3. Inicialice el conjunto de datos en la plataforma con el `cnvrg data init` comando.
4. Cargue todo el contenido del directorio en el lago de datos central con el `cnvrg data sync` Command. una vez que los datos se cargan en el almacén de objetos central (StorageGRID, S3 u otros), puede navegar por la GUI. La siguiente figura muestra un archivo PNG cargado de imagen de fibrosis torácica. Además, cnvrg introduce los datos de modo que cualquier modelo que cree pueda reproducirse en la versión de los datos.



## Datos de Cach

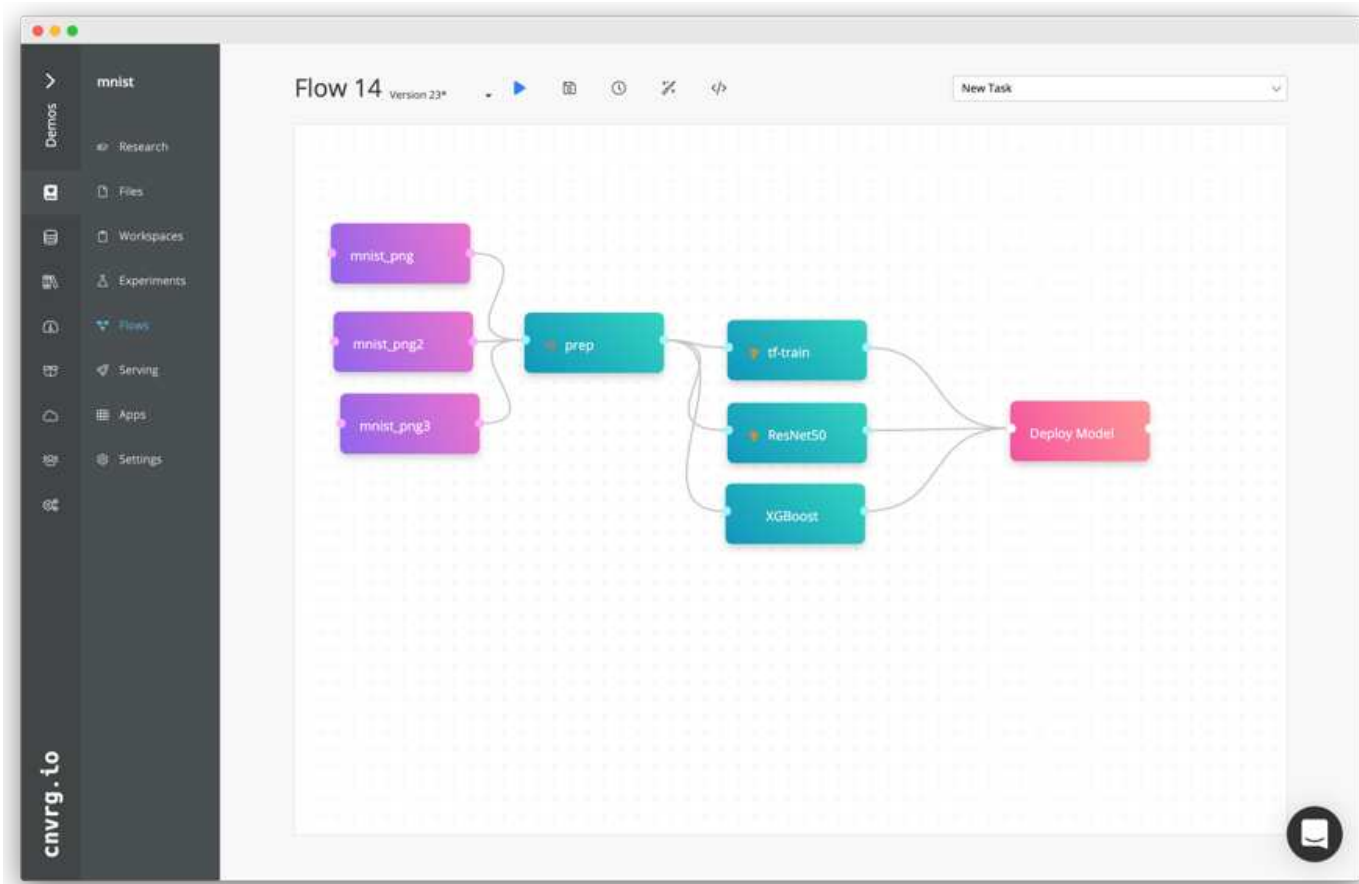
Para realizar el entrenamiento con más rapidez y evitar la descarga de más de 600 000 archivos para cada entrenamiento y experimento, utilizamos la función de almacenamiento en caché de datos después de cargar los datos inicialmente en el almacén de objetos del lago de datos central.



Después de que los usuarios hagan clic en caché, cnvrg descarga los datos de su registro específico del almacén de objetos remoto y los almacena en caché en el volumen NFS de ONTAP. Una vez que finalice, los datos estarán disponibles para el entrenamiento instantáneo. Además, si los datos no se utilizan durante unos días (para el entrenamiento de modelos o la exploración, por ejemplo), cnvrg borra automáticamente la memoria caché.

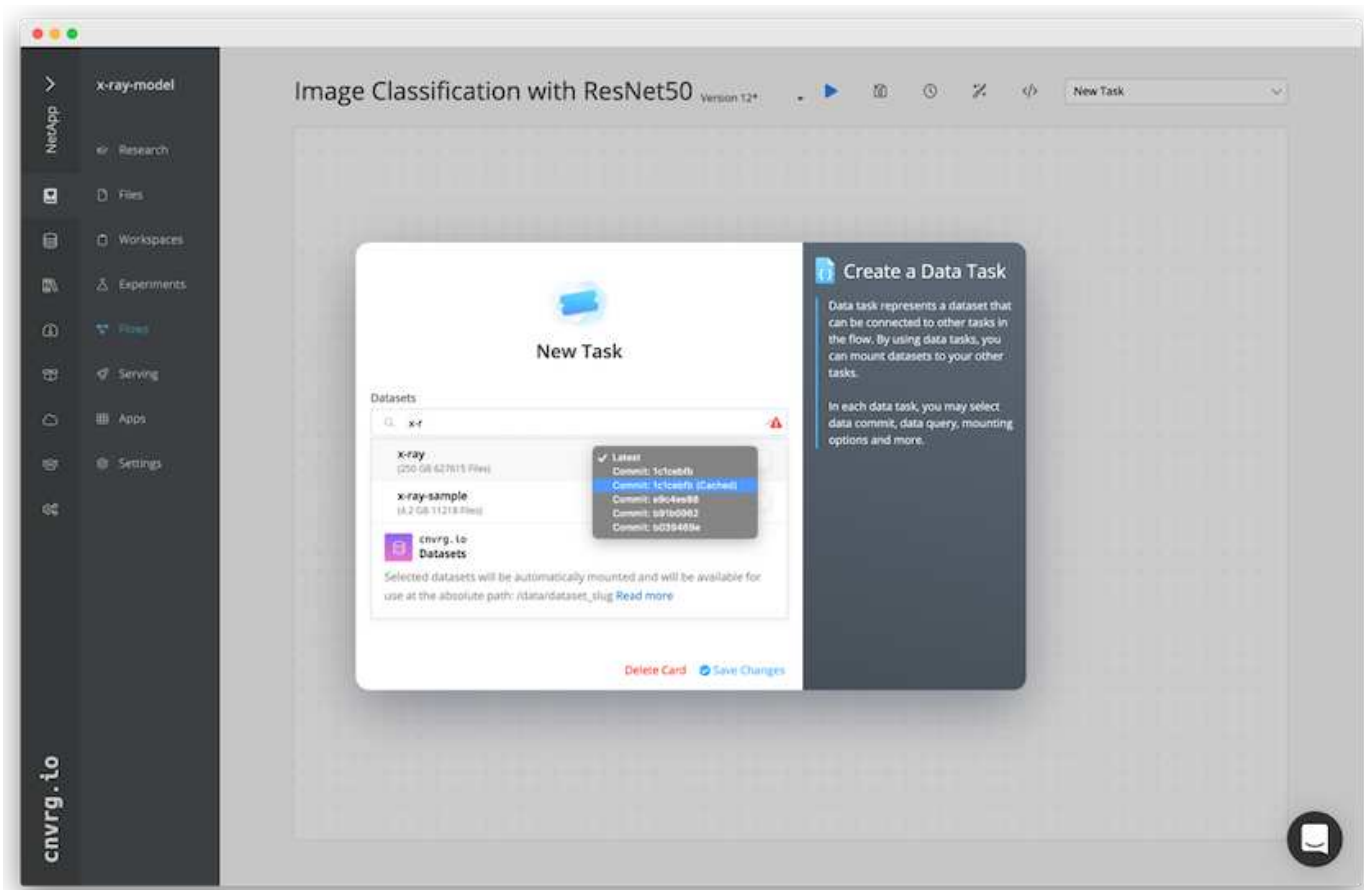
### Cree un ML de canalización con datos en caché

Cnvrg fluye le permite construir fácilmente tuberías ML de producción. Los flujos son flexibles, pueden funcionar para cualquier tipo de caso de uso DE ML y pueden crearse a través de la GUI o el código. Cada componente de un flujo puede ejecutarse en un recurso de computación diferente con una imagen de Docker diferente, por lo que es posible crear canalizaciones de APRENDIZAJE AUTOMÁTICO optimizadas y de cloud híbrido.



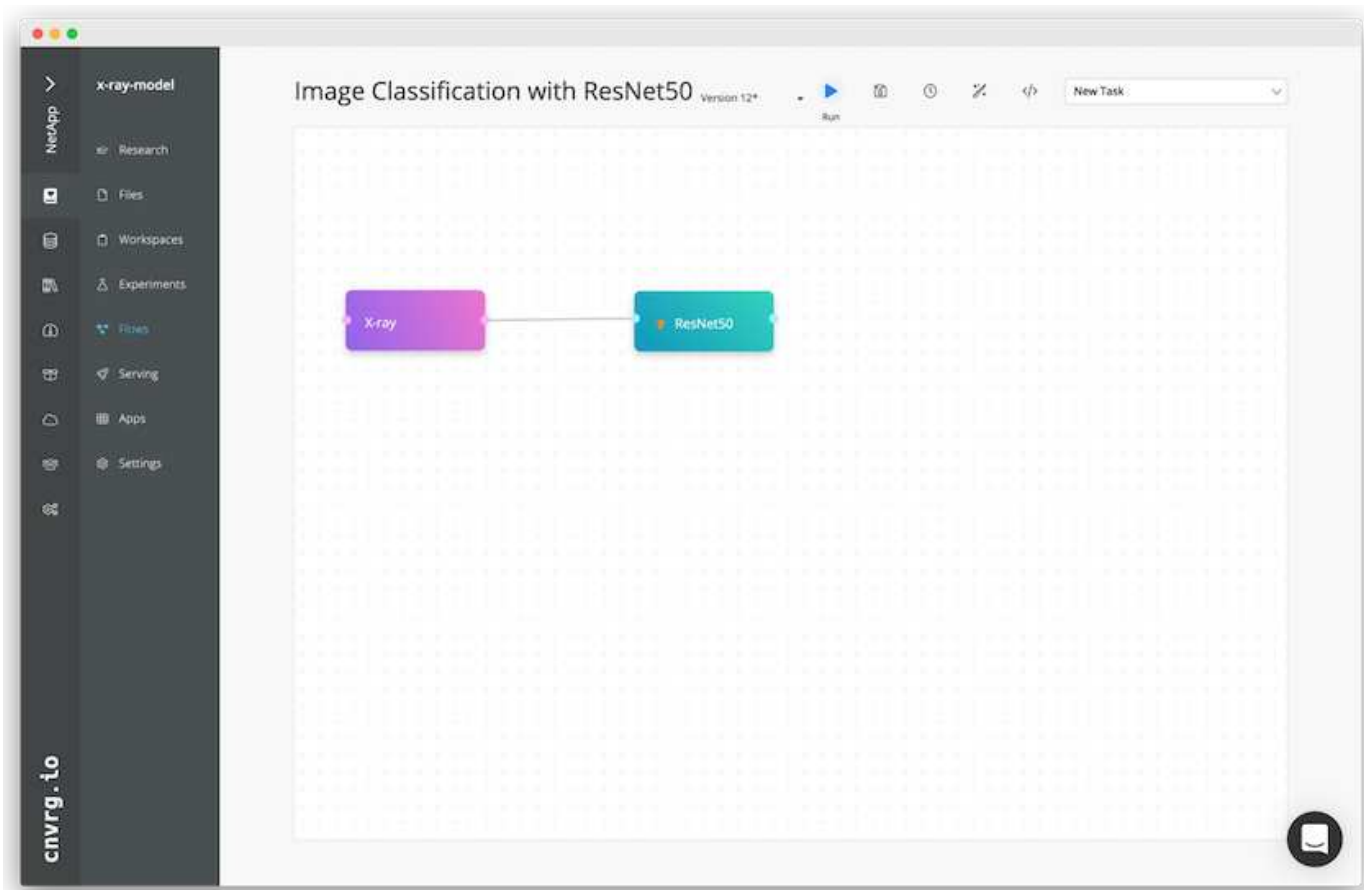
### Cómo construir el flujo de la radiografía de tórax: Ajuste de los datos

Hemos añadido nuestro conjunto de datos a un flujo recién creado. Al agregar el conjunto de datos, puede seleccionar la versión específica (Commit) e indicar si desea la versión en caché. En este ejemplo, hemos seleccionado la confirmación almacenada en caché.



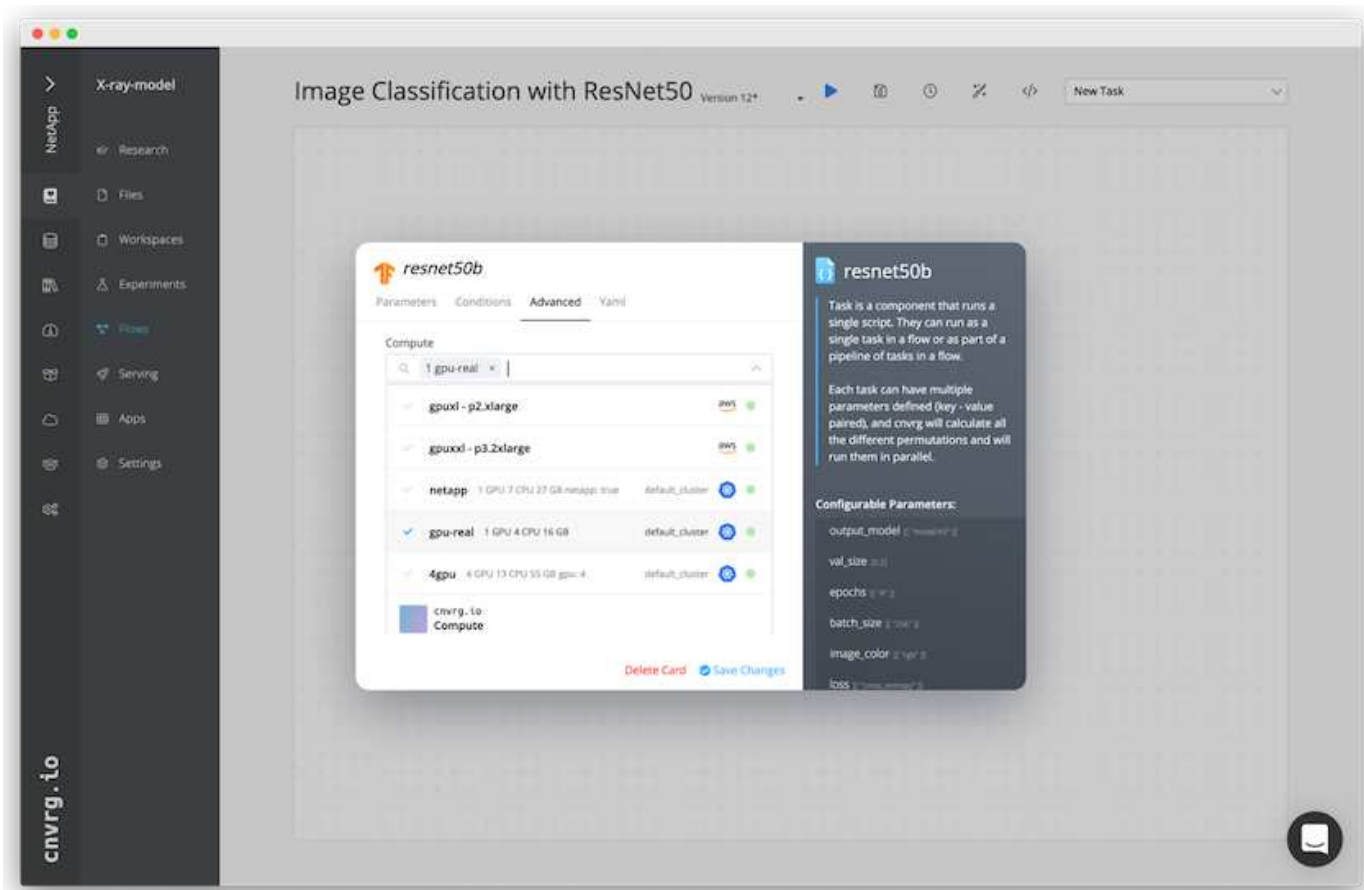
## Cómo construir el flujo de rayos X del tórax: Ajuste del modelo de entrenamiento: ResNet50

En la canalización, puede agregar cualquier tipo de código personalizado que desee. En cnvrg, también existe la biblioteca de IA, una colección de componentes DE ML reutilizables. En la biblioteca de IA existen algoritmos, scripts, orígenes de datos y otras soluciones que se pueden usar en cualquier flujo DE APRENDIZAJE profundo o DE ML. En este ejemplo, hemos seleccionado el módulo ResNet50 premontado. Se utilizaron parámetros predeterminados como batch\_size:128, épocas:10 y más. Estos parámetros pueden verse en los documentos de la Biblioteca de IA. La siguiente captura de pantalla muestra el nuevo flujo con el conjunto de datos de rayos X conectado a ResNet50.



## Defina el recurso de computación para ResNet50

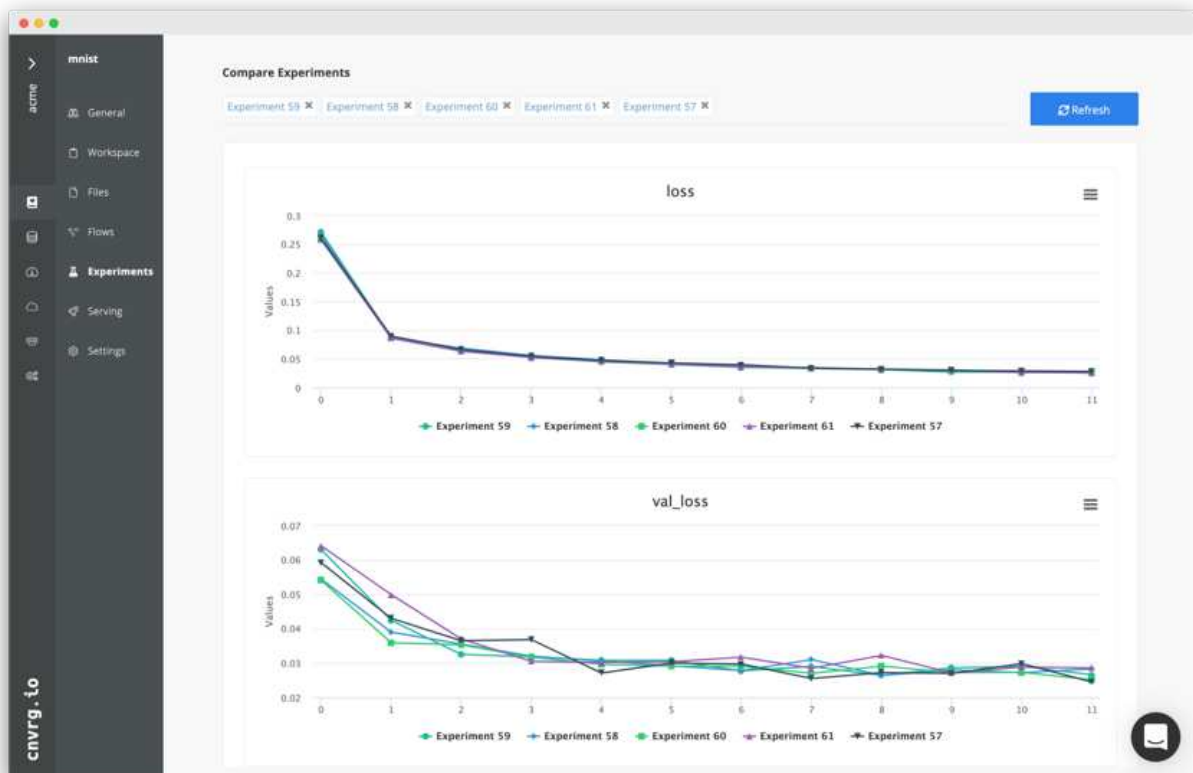
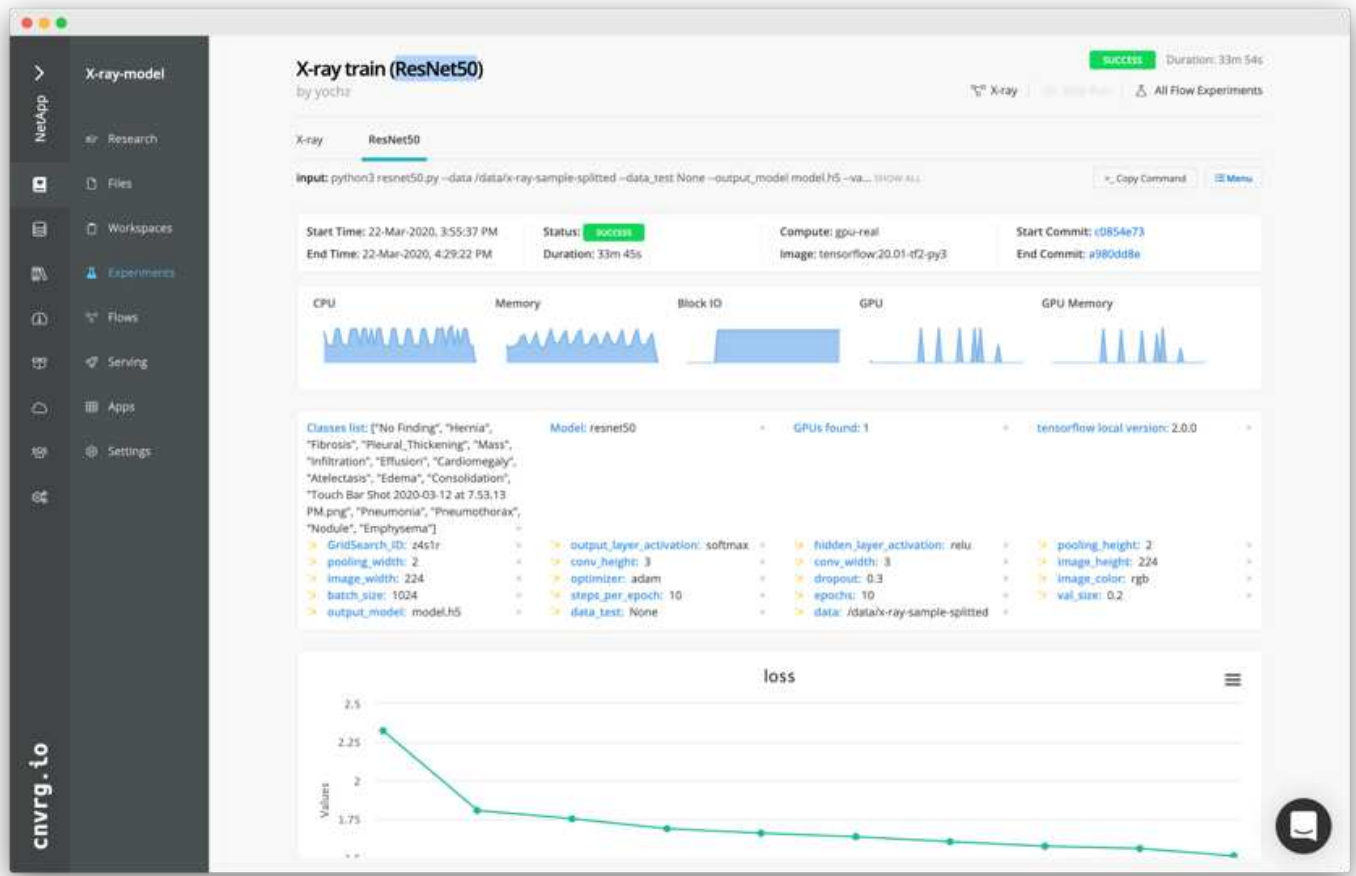
Cada algoritmo o componente en flujos cnvrg puede ejecutarse en una instancia de computación diferente, con una imagen de Docker diferente. En nuestra configuración, queríamos ejecutar el algoritmo de entrenamiento en los sistemas DGX de NVIDIA con la arquitectura de IA ONTAP de NetApp. En la siguiente figura, hemos seleccionado `gpu-real`, que es una plantilla de cálculo y una especificación para nuestro clúster local. También creamos una cola de plantillas y seleccionamos varias plantillas. De esta manera, si el `gpu-real` no se puede asignar el recurso (si, por ejemplo, otros científicos de datos lo están utilizando), puede habilitar la explosión automática en el cloud añadiendo una plantilla de proveedor de cloud. La siguiente captura de pantalla muestra el uso de `gpu-real` como nodo de computación para ResNet50.



## Seguimiento y seguimiento de resultados

Después de ejecutar un flujo, cnvrg activa el motor de seguimiento y supervisión. Cada ejecución de un flujo se documenta y actualiza automáticamente en tiempo real. Hiperparámetros, métricas, uso de recursos (utilización de GPU, etc.), versión de código, artefactos, registros Y así sucesivamente están automáticamente disponibles en la sección experimentos, como se muestra en las dos capturas de pantalla siguientes.





## Conclusión

NetApp y cnvrg.io se han asociado para ofrecer a los clientes una solución completa de gestión de datos para el desarrollo de software ML y DL. ONTAP AI proporciona almacenamiento y computación de alto rendimiento para cualquier escala de operación y el software cnvrg.io optimiza los flujos de trabajo de ciencia de datos y mejora la utilización de recursos.

## Reconocimientos

- Mike Oglesby, Ingeniero Técnico de Marketing de NetApp
- Santosh Rao, director técnico sénior, NetApp

## Dónde encontrar información adicional

Si quiere obtener más información sobre el contenido de este documento, consulte los siguientes recursos:

- Nvrg.io ( "<https://cnvrg.io>"):
  - NÚCLEO Cnvrg (plataforma DE ML libre)  
<https://cnvrg.io/platform/core>
  - Documentos de Cnvrg  
["https://app.cnvrg.io/docs"](https://app.cnvrg.io/docs)
- Servidores DGX-1 de NVIDIA:
  - Servidores DGX-1 de NVIDIA  
<https://www.nvidia.com/en-us/data-center/dgx-1/>
  - GPU de núcleo tensor NVIDIA Tesla V100  
<https://www.nvidia.com/en-us/data-center/tesla-v100/>
  - GPU CLOUD DE NVIDIA (NGC)  
<https://www.nvidia.com/en-us/gpu-cloud/>
- Sistemas AFF de NetApp:
  - Especificaciones técnicas de AFF  
<https://www.netapp.com/us/media/d-3582.pdf>
  - FlashAdvantage de NetApp para AFF  
<https://www.netapp.com/us/media/ds-3733.pdf>
  - Documentación de ONTAP 9.x.  
<http://mysupport.netapp.com/documentation/productlibrary/index.html?productID=62286>

- Informe técnico de NetApp FlexGroup

<https://www.netapp.com/us/media/tr-4557.pdf>

- Almacenamiento persistente de NetApp para contenedores:

- Trident de NetApp

<https://netapp.io/persistent-storage-provisioner-for-kubernetes/>

- Matriz de interoperabilidad de NetApp:

- Herramienta de matriz de interoperabilidad de NetApp

<https://mysupport.netapp.com/matrix/#welcome>

- Redes de IA de ONTAP:

- Switches Cisco Nexus 3232C

<https://www.cisco.com/c/en/us/products/switches/nexus-3232c-switch/index.html>

- Switches Mellanox Spectrum serie 2000

[http://www.mellanox.com/page/products\\_dyn?product\\_family=251&mtag=sn2000](http://www.mellanox.com/page/products_dyn?product_family=251&mtag=sn2000)

- Marco Y herramientas DE ML:

- DALÍ

<https://github.com/NVIDIA/DALI>

- TensorFlow: Un marco de aprendizaje automático de código abierto para todos

<https://www.tensorflow.org/>

- Horovod: El marco de Uber para el aprendizaje automático distribuido de código abierto para TensorFlow

<https://eng.uber.com/horovod/>

- Habilitación de GPU en el ecosistema Container Runtime

<https://devblogs.nvidia.com/gpu-containers-runtime/>

- Docker

<https://docs.docker.com>

- Kubernetes

<https://kubernetes.io/docs/home/>

- DeepOps de NVIDIA

<https://github.com/NVIDIA/deepops>

- Kubeflow

<http://www.kubeflow.org/>

- Servidor de portátiles Jupyter

<http://www.jupyter.org/>

- Conjuntos de datos y pruebas:
  - Conjunto de datos de radiografía de tórax de los NIH

<https://nihcc.app.box.com/v/ChestXray-NIHCC>

- Xiaosong Wang, Yifan Peng, le Lu, Zhiyong Lu, Mohammadadi Bagheri, Ronald Summers, ChestX-ray8: Base de datos de rayos X en el pecho a escala de hospital y puntos de referencia sobre clasificación y localización de enfermedades comunes del tórax, IEEE CVPR, págs 3462-3471, 2017TR-4841-0620

## **Informe técnico TR-4732: Big Data Analytics datos en inteligencia artificial**

Karthikeyan Nagalingam, NetApp

Este documento describe cómo trasladar los datos de análisis de Big Data y los datos de computación de alto rendimiento a IA. La IA procesa datos de NFS a través de exportaciones NFS, mientras que los clientes suelen tener sus datos de IA en una plataforma de análisis de Big Data, como HDFS, Blob o S3, así como plataformas HPC como GPFS. Este documento proporciona directrices para trasladar datos de análisis de Big Data y datos de HPC a IA mediante NetApp XCP y NIPAM. También hablamos de las ventajas empresariales que supone trasladar datos de Big Data y de HPC a IA.

### **Conceptos y componentes**

#### **Almacenamiento de análisis de Big Data**

Los análisis de Big Data son el principal proveedor de almacenamiento para HDFS. Un cliente suele utilizar un sistema de archivos compatible con Hadoop (HCFS), como almacenamiento blob de Windows Azure, MapR File System (MapR-FS) y almacenamiento de objetos S3.

#### **Sistema de archivos paralelos general**

GPFS de IBM es un sistema de archivos empresariales que ofrece una alternativa a HDFS. GPFS proporciona flexibilidad para que las aplicaciones decidan el tamaño del bloque y el diseño de replicación, lo que proporciona un buen rendimiento y eficiencia.

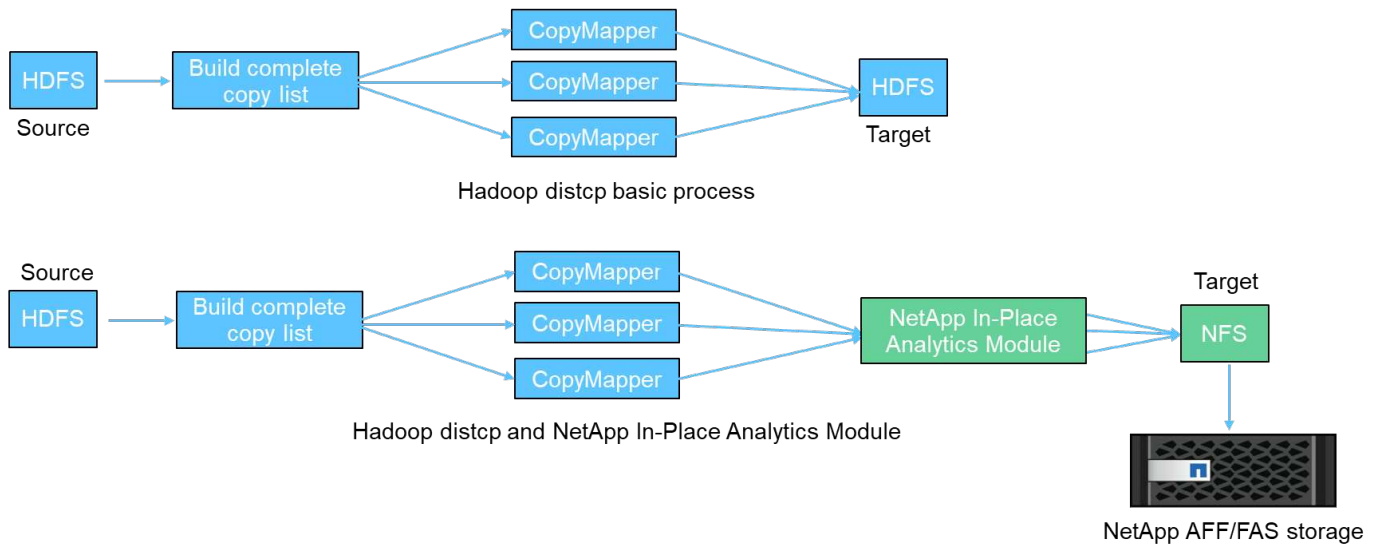
#### **Módulo de análisis in situ de NetApp**

El módulo de análisis in situ (NIPAM, in situ) de NetApp sirve como controlador para que los clústeres de Hadoop accedan a datos NFS. Consta de cuatro componentes: Un pool de conexión, un InputStream NFS, una caché de gestión de archivos y un OutputStream NFS. Para obtener más información, consulte "[TR-4382: Módulo de análisis in situ de NetApp.](#)"

#### **Copia distribuida de Hadoop**

La copia distribuida de Hadoop (DistCp) es una herramienta de copia distribuida que se usa para las tareas de adaptación entre clústeres y dentro de clústeres de gran tamaño. Esta herramienta utiliza MapReduce para la

distribución de datos, el manejo de errores y los informes. Amplía la lista de archivos y directorios e introduce las tareas de asignación para copiar los datos de la lista de origen. La siguiente imagen muestra la operación DistCp en HDFS y no HDFS.



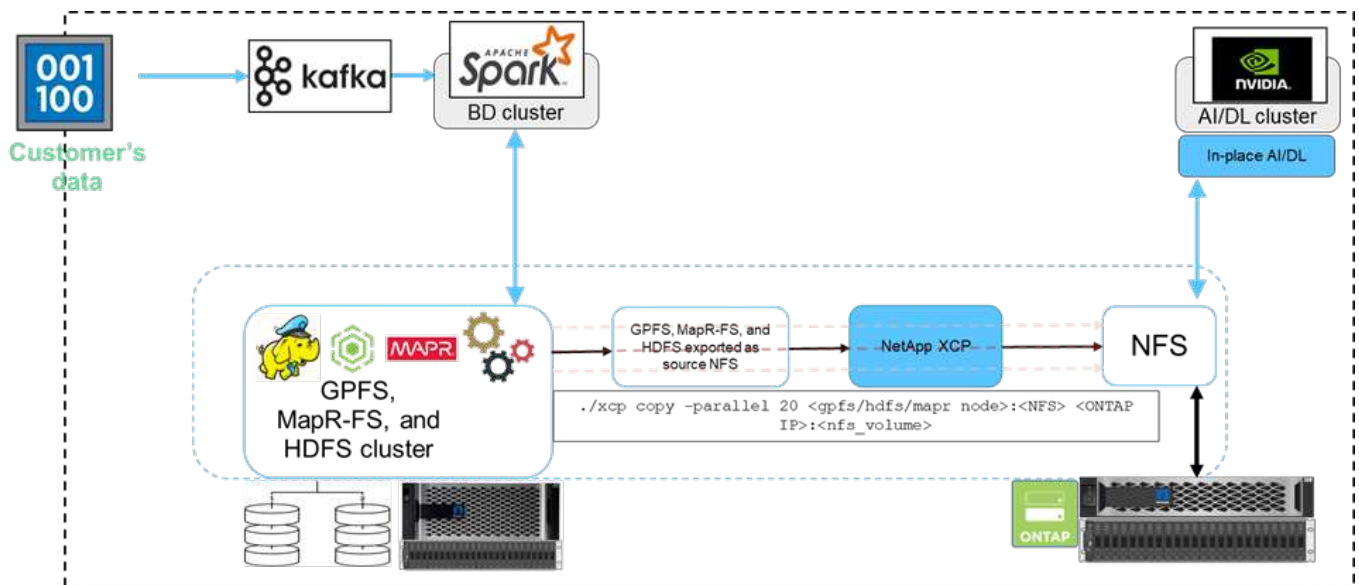
Hadoop DistCp mueve datos entre los dos sistemas HDFS sin necesidad de utilizar un controlador adicional. NetApp proporciona el controlador para sistemas que no son HDFS. En un destino NFS, NIPAM proporciona el controlador para copiar datos que Hadoop DistCp utiliza para comunicarse con destinos NFS al copiar datos.

### Cloud Volumes Service de NetApp

Cloud Volumes Service de NetApp es un servicio de archivos nativo del cloud con un rendimiento extremo. Este servicio ayuda a sus clientes a acelerar el plazo de comercialización mediante el rápido aumento y reducción de los recursos, así como el uso de las funciones de NetApp para mejorar la productividad y reducir el tiempo de inactividad del personal. Cloud Volumes Service es la alternativa adecuada para la recuperación ante desastres y sus backups en el cloud, ya que reduce el espacio global del centro de datos y consume menos almacenamiento en cloud público nativo.

### XCP de NetApp

NetApp XCP es un software cliente que permite una migración de datos de cualquiera a NetApp y de NetApp a NetApp rápida y fiable. Esta herramienta está diseñada para copiar una gran cantidad de datos NAS no estructurados de cualquier sistema NAS a una controladora de almacenamiento de NetApp. La herramienta de migración XCP utiliza un motor de transmisión de E/S multicanal y de varios canales que puede procesar muchas solicitudes en paralelo, como la migración de datos, listas de archivos o directorios y la creación de informes de espacio. Esta es la herramienta de migración de datos de NetApp predeterminada. Puede utilizar XCP para copiar datos de un clúster de Hadoop y de una HPC al almacenamiento NFS de NetApp. El siguiente diagrama muestra la transferencia de datos de un clúster Hadoop y HPC a un volumen NFS de NetApp mediante XCP.



## Copia y sincronización de NetApp BlueXP

La copia y sincronización de BlueXP de NetApp es un software como servicio de replicación de datos híbrida que transfiere y sincroniza datos de NFS, S3 y CIFS sin problemas y de forma segura entre el almacenamiento on-premises y el almacenamiento en la nube. Este software se utiliza para migración de datos, archivado, colaboración, análisis, etc. Una vez transferidos los datos, BlueXP Copy y Sync sincronizan continuamente los datos entre el origen y el destino. De cara al futuro, luego transfiere el delta. También protege los datos dentro de su propia red, en el cloud o en las instalaciones. Este software se basa en el modelo de pago por uso, que ofrece una solución rentable y ofrece capacidades de supervisión y generación de informes para su transferencia de datos.

## Inferencia de IA en el perímetro - NetApp con Lenovo ThinkSystem - Diseño de la solución

### TR-4886: Inferencia de IA en el Edge - NetApp con Lenovo ThinkSystem - Diseño de la solución

Sathish Thyagarajan, NetApp Miroslav Hodak, Lenovo

#### Resumen

Varios supuestos de aplicaciones emergentes, como los sistemas avanzados de asistencia al conductor (ADAS), el sector 4.0, las ciudades inteligentes y el Internet de las cosas (IoT), requieren el procesamiento de flujos de datos continuos con una latencia cercana a cero. Este documento describe una arquitectura de computación y almacenamiento para poner en marcha la inferencia de inteligencia artificial (IA) basada en GPU en controladoras de almacenamiento de NetApp y servidores Lenovo ThinkSystem en un entorno perimetral que satisface estos requisitos. Este documento también proporciona datos de rendimiento para las pruebas de rendimiento de inferencia MLPerf estándares del sector, por lo que evalúa diversas tareas de inferencia en servidores periféricos equipados con GPU T4 de NVIDIA. Investigamos el rendimiento de escenarios de inferencia multisequencia, sin conexión y con múltiples flujos, y mostramos que la arquitectura con un sistema de almacenamiento en red compartido rentable tiene un alto rendimiento y proporciona un punto central para la gestión de modelos y datos en servidores periféricos múltiples.

#### Introducción

Las empresas están generando cada vez más volúmenes masivos de datos en el extremo de la red. Con el fin de obtener el máximo valor a partir de sensores inteligentes y datos del Internet de las cosas, las

organizaciones buscan una solución de transmisión de eventos en tiempo real que permita el uso de la informática perimetral. Así, pues, cada vez se realizan más tareas informáticas exigentes fuera del perímetro y fuera de los centros de datos. La inferencia de IA es uno de los motores de esta tendencia. Los servidores periféricos proporcionan suficiente potencia computacional para estas cargas de trabajo, sobre todo cuando se utilizan aceleradores, pero el almacenamiento limitado suele ser un problema, especialmente en entornos con varios servidores. En este documento mostramos cómo puede poner en marcha un sistema de almacenamiento compartido en un entorno perimetral y cómo beneficia a las cargas de trabajo de inferencia de IA sin perjudicar el rendimiento.

Este documento describe una arquitectura de referencia para la inferencia de IA en el perímetro. Combina múltiples servidores Lenovo ThinkSystem Edge con un sistema de almacenamiento de NetApp para crear una solución fácil de poner en marcha y gestionar. Está previsto ser una guía de referencia para implementaciones prácticas en diversas situaciones, como la planta de fábrica con múltiples cámaras y sensores industriales, sistemas de punto de venta (POS) en transacciones de venta al por menor o sistemas de autoconducción completa (FSD) que identifican anomalías visuales en vehículos autónomos.

Este documento recoge las pruebas y validación de una configuración de computación y almacenamiento que consta de Lenovo ThinkSystem SE350 Edge Server y un sistema de almacenamiento AFF y EF-Series de NetApp de gama básica. Las arquitecturas de referencia proporcionan una solución eficiente y rentable para puestas en marcha de IA, a la vez que proporcionan servicios de datos completos, protección de datos integrada, escalabilidad fluida y almacenamiento de datos conectado al cloud con ONTAP de NetApp y el software para la gestión de datos SANtricity de NetApp.

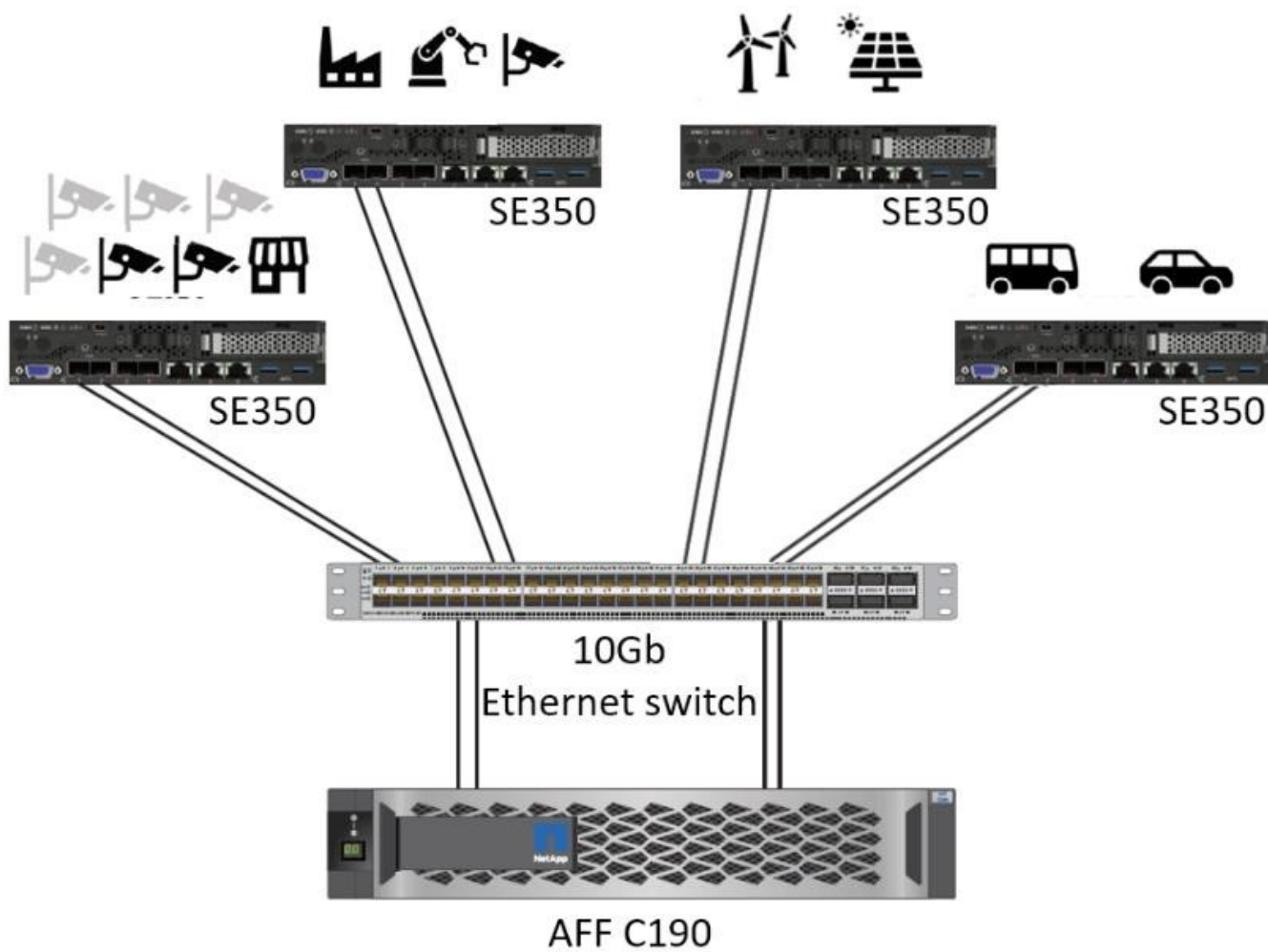
## **Público objetivo**

Este documento está dirigido a los siguientes destinatarios:

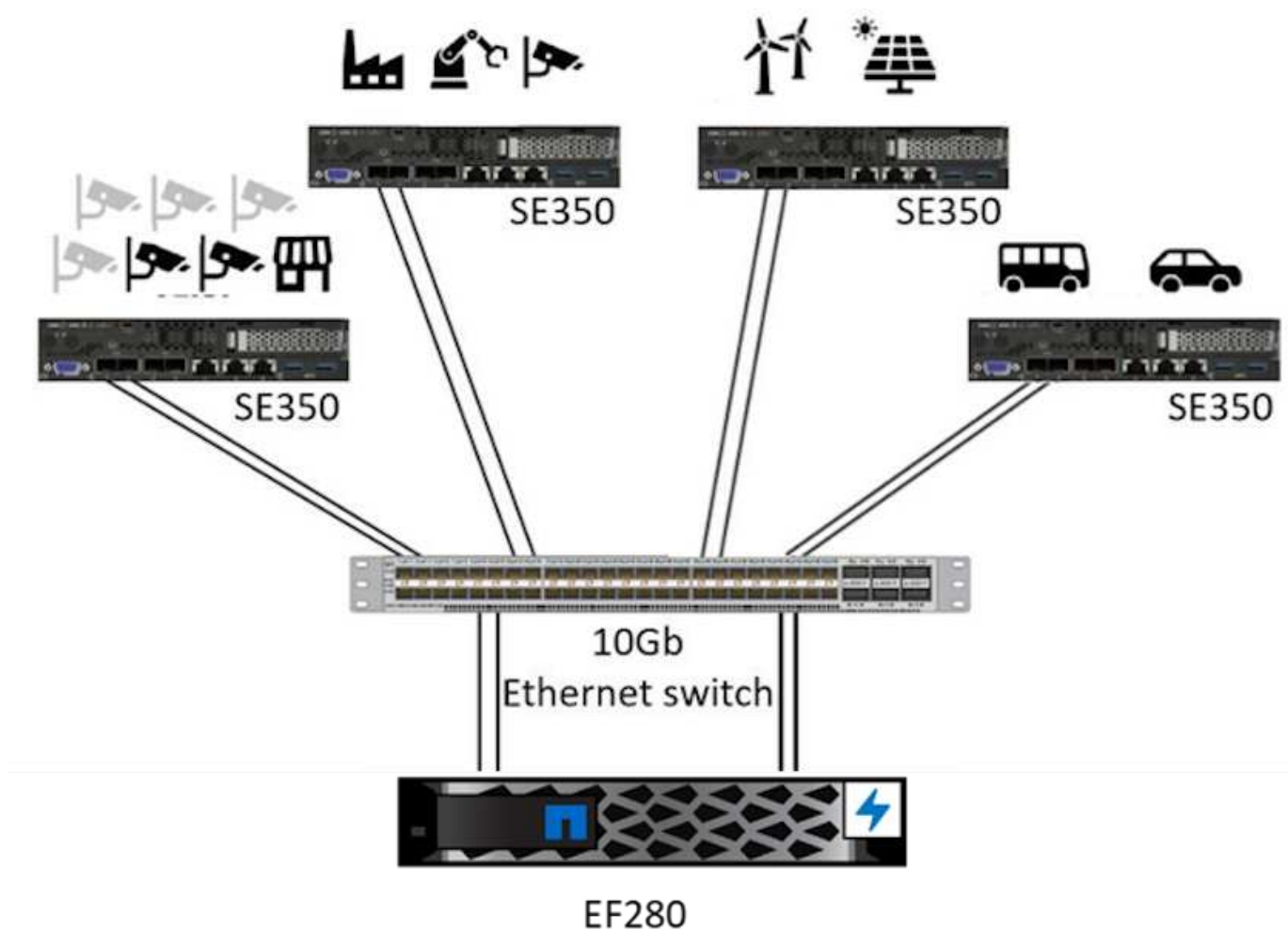
- Líderes de negocio y arquitectos empresariales que quieren aprovechar la IA en el perímetro.
- Científicos de datos, ingenieros de datos, investigadores DE IA/aprendizaje automático y desarrolladores de sistemas de IA.
- Arquitectos empresariales que diseñan soluciones para el desarrollo de modelos y aplicaciones de IA/ML.
- Científicos e ingenieros de IA que buscan formas eficientes de poner en marcha modelos DE aprendizaje profundo (DL) Y ML.
- Los administradores de dispositivos periféricos y los administradores de servidor perimetral son responsables de la puesta en marcha y la gestión de modelos de inferencia perimetrales.

## **Arquitectura de la solución**

Este servidor Lenovo ThinkSystem y la solución de almacenamiento ONTAP o SANtricity de NetApp están diseñados para gestionar la inferencia de IA en grandes conjuntos de datos mediante la potencia de procesamiento de GPU junto con CPU tradicionales. Esta validación demuestra un alto rendimiento y una gestión de datos óptima con una arquitectura que utiliza uno o varios servidores de borde Lenovo SR350 interconectados con un único sistema de almacenamiento AFF de NetApp, como se muestra en las siguientes dos figuras.

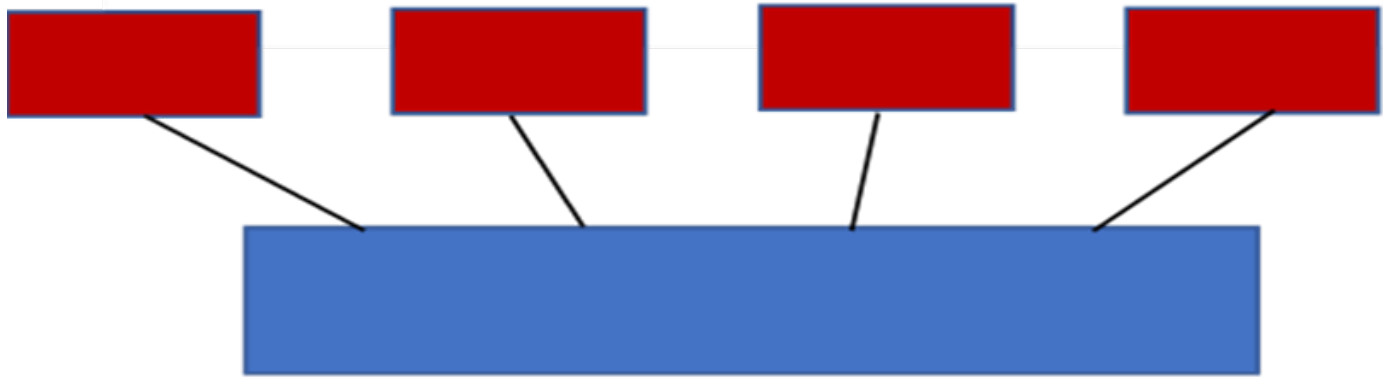






La descripción general de la arquitectura lógica de la siguiente figura muestra las funciones de los elementos de computación y almacenamiento de esta arquitectura. Específicamente, muestra lo siguiente:

- Dispositivos de computación periféricos que realizan inferencia en los datos que recibe de cámaras, sensores, etc.
- Un elemento de almacenamiento compartido que ofrece diversos objetivos:
  - Ofrece una ubicación central para los modelos de inferencia y otros datos que necesitan para realizar la inferencia. Los servidores informáticos acceden al almacenamiento directamente y utilizan los modelos de inferencia en toda la red sin necesidad de copiarlos de forma local.
  - Los modelos actualizados se empujan aquí.
  - Archiva los datos de entrada que los servidores perimetrales reciben para analizarlos posteriormente. Por ejemplo, si los dispositivos periféricos están conectados a las cámaras, el elemento de almacenamiento mantiene los vídeos capturados por las cámaras.



rojo	azul
Sistema de computación Lenovo	Sistema de almacenamiento AFF de NetApp
Dispositivos periféricos que realizan inferencia en las entradas de cámaras, sensores, etc.	Almacenamiento compartido que mantiene los modelos de inferencia y los datos de los dispositivos periféricos para su posterior análisis.

Esta solución de NetApp y Lenovo ofrece las siguientes ventajas clave:

- GPU acelera la computación en el extremo.
- Implementación de varios servidores periféricos, respaldados y gestionado desde un almacenamiento compartido.
- Protección de datos sólida para cumplir con los objetivos de punto de recuperación (RPO) y los objetivos de tiempo de recuperación (RTO) bajos sin pérdida de datos.
- Gestión de datos optimizada con copias Snapshot y clones de NetApp para optimizar los flujos de trabajo de desarrollo.

### Cómo utilizar esta arquitectura

Este documento valida el diseño y el rendimiento de la arquitectura propuesta. Sin embargo, no hemos probado ciertas piezas a nivel de software, como la gestión y sincronización de datos de contenedores, cargas de trabajo o modelos con cloud o centros de datos en las instalaciones, ya que son específicas de un escenario de puesta en marcha. Aquí existen varias opciones.

En el nivel de gestión de contenedores, la gestión de contenedores de Kubernetes es una buena opción y es compatible con una versión totalmente ascendente (Canonical) o con una versión modificada adecuada para implementaciones empresariales (Red Hat). La ["Plano de control de IA de NetApp"](#) Que utiliza Trident de NetApp y los recientemente añadidos ["Kit de herramientas de operaciones de datos de NetApp"](#) Proporciona trazabilidad, funciones de gestión de datos, interfaces y herramientas integradas para que los científicos e ingenieros de datos se integren con el sistema de almacenamiento de NetApp. Kubeflow, el kit DE herramientas ML para Kubernetes, proporciona funcionalidades de IA adicionales junto con compatibilidad para el control de versiones de modelos y KFServing en varias plataformas como TensorFlow Serving o NVIDIA Triton inferpensado Server. Otra opción es la plataforma NVIDIA EGX, que proporciona gestión de las cargas de trabajo junto con acceso a un catálogo de contenedores de inferencia de IA habilitados para GPU. Sin embargo, es posible que estas opciones requieran gran esfuerzo y experiencia para ponerlas en producción y es posible que requieran la ayuda de un proveedor de software independiente (ISV) o un asesor de terceros.

## Áreas de soluciones

La ventaja clave de la inferencia de IA y la computación perimetral es la capacidad de que los dispositivos calculen, procesen y analicen datos con un alto nivel de calidad sin latencia. Hay demasiados ejemplos de casos de uso de computación perimetral que describir en este documento, pero aquí hay algunos ejemplos destacados:

### Automóviles: Vehículos autónomos

La ilustración clásica de la informática avanzada se encuentra en los sistemas avanzados de asistencia al conductor (ADAS) en vehículos autónomos (AV). La IA en vehículos sin conductor debe procesar rápidamente una gran cantidad de datos procedentes de cámaras y sensores para garantizar su seguridad. Tomar demasiado tiempo para interpretar entre un objeto y un humano puede significar la vida o la muerte, por lo tanto poder procesar los datos lo más cerca posible del vehículo es crucial. En este caso, uno o varios servidores de computación periféricos se encarga de las entradas de cámaras, RADAR, LiDAR y otros sensores, mientras que el almacenamiento compartido contiene modelos de inferencia y almacena datos de entrada de los sensores.

### Atención sanitaria: Monitorización de pacientes

Uno de los mayores impactos de la IA y la informática perimetral es su capacidad para mejorar la supervisión continua de pacientes para enfermedades crónicas, tanto en las unidades de cuidados intensivos como en las unidades de cuidados intensivos (UCI). Los datos de los dispositivos periféricos que supervisan los niveles de insulina, la respiración, la actividad neurológica, el ritmo cardíaco y las funciones gastrointestinales requieren un análisis instantáneo de los datos que deben ser objeto de acciones inmediatas porque hay poco tiempo para actuar y salvar la vida de alguien.

### Venta al por menor: Pago sin cajero

La computación perimetral puede ayudar a los minoristas a reducir el tiempo de salida y aumentar el tráfico de pies. Los sistemas sin cajero admiten varios componentes, como los siguientes:

- Autenticación y acceso. Conectar el comprador físico a una cuenta validada y permitir el acceso al espacio de venta al por menor.
- Supervisión de inventario. Utilizar sensores, etiquetas RFID y sistemas de visión computarizada para confirmar la selección o desección de artículos por parte de los compradores.

Aquí, cada uno de los servidores perimetrales gestiona cada contador de retirada y el sistema de almacenamiento compartido sirve como punto de sincronización central.

### Servicios financieros: Seguridad humana en quioscos y prevención del fraude

Las organizaciones bancarias utilizan la IA y la informática perimetral para innovar y crear experiencias bancarias personalizadas. Los quioscos interactivos, mediante el análisis de datos en tiempo real y la inferencia de IA, permiten ahora a los cajeros automáticos no sólo ayudar a los clientes a retirar el dinero, sino también supervisar de forma proactiva los quioscos a través de las imágenes capturadas con las cámaras para identificar el riesgo para la seguridad humana o el comportamiento fraudulento. En este escenario, los servidores periféricos informáticos y los sistemas de almacenamiento compartido se conectan a quioscos y cámaras interactivos para ayudar a los bancos a recopilar y procesar datos con modelos de inferencia de IA.

### Fabricación: Industria 4.0

La cuarta revolución industrial (Industry 4.0) ha comenzado, junto con tendencias emergentes como Smart Factory e impresión 3D. Con el fin de prepararse para un futuro impulsado por los datos, la comunicación y el

Internet de las cosas (M2M) a gran escala están integrados para una mayor automatización sin necesidad de intervención humana. La fabricación ya está altamente automatizada y añadir características de IA es una continuación natural de la tendencia a largo plazo. La IA permite automatizar operaciones que se pueden automatizar con la ayuda de la visión computarizada y otras funcionalidades de IA. Puede automatizar el control de calidad o las tareas que se basan en la visión humana o en la toma de decisiones para realizar análisis más rápidos de materiales en líneas de ensamblaje en plantas de fabricación para ayudar a las plantas de fabricación a cumplir con los estándares ISO requeridos de gestión de la seguridad y la calidad. Aquí, cada servidor perimetral informático está conectado a una matriz de sensores que supervisan el proceso de fabricación, y cuando es necesario, los modelos de inferencia actualizados se ven empujados al almacenamiento compartido.

### **Telecomunicaciones: Detección de óxido, inspección de torre y optimización de la red**

La industria de las telecomunicaciones utiliza técnicas de IA y visión informática para procesar imágenes que detectan automáticamente el óxido e identifican las torres celulares que contienen corrosión y, por lo tanto, requieren una inspección adicional. El uso de imágenes con drones y modelos de IA para identificar regiones distintas de una torre para analizar la oxidación, las grietas superficiales y la corrosión ha aumentado en los últimos años. La demanda continúa creciendo para tecnologías de IA que permiten inspeccionar eficientemente la infraestructura de telecomunicaciones y las torres de células, evaluar periódicamente la degradación y reparar rápidamente cuando sea necesario.

Además, otro caso de uso emergente en las telecomunicaciones es el uso de algoritmos de IA y ML para predecir patrones de tráfico de datos, detectar dispositivos compatibles con 5G y automatizar y aumentar la gestión energética de múltiples entradas y salidas múltiples (MIMO). El hardware de MIMO se utiliza en las torres de radio para aumentar la capacidad de la red; sin embargo, esto viene con costos de energía adicionales. Los modelos ML para “modo de suspensión MIMO” implementados en las zonas de células pueden predecir el uso eficiente de las radios y ayudar a reducir los costes de consumo de energía para los operadores de redes móviles (MNO). Las soluciones de computación avanzada y de inferencia de IA ayudan a las MNO a reducir la cantidad de datos transmitidos hacia los centros de datos, reducir su TCO, optimizar las operaciones de red y mejorar el rendimiento general de los usuarios finales.

### **Información general de la tecnología**

Esta sección describe la base tecnológica de esta solución de IA.

#### **Sistemas AFF de NetApp**

Los sistemas de almacenamiento AFF de NetApp más innovadores permiten que las puestas en marcha de inferencia de IA en el perímetro satisfagan los requisitos de almacenamiento de la empresa con un rendimiento líder del sector, una flexibilidad superior, integración con el cloud y la mejor gestión de datos de su clase. Los sistemas AFF de NetApp han sido diseñados específicamente para flash y ayudan a acelerar, gestionar y proteger los datos esenciales para la empresa.

- Los sistemas de almacenamiento AFF de NetApp de gama básica se basan en hardware FAS2750 y medios flash SSD
- Dos controladoras en configuración de alta disponibilidad



Los sistemas de almacenamiento C190 de gama básica de AFF admiten las siguientes funciones:

- Un número máximo de unidades de estado sólido de 960 GB
- Dos configuraciones posibles:
  - Ethernet (10 GbE): 4 puertos 10GBASE-T (RJ-45)
  - Unificado (FC de 16 GB o 10 GbE): 4 puertos de adaptador de destino unificado 2 (UTA2)
- Una capacidad efectiva máxima de 50,5 TB



En el caso de cargas de trabajo NAS, un único sistema AFF C190 de gama básica admite un rendimiento de 4,4 GB/s para lecturas secuenciales y 230 000 IOPS para lecturas aleatorias pequeñas con latencias de 1 ms o menos.

#### AFF A220 de NetApp

NetApp también ofrece otros sistemas de almacenamiento de gama básica que proporcionan un mayor rendimiento y escalabilidad para las puestas en marcha a gran escala. En el caso de cargas de trabajo NAS, un único sistema AFF A220 de gama básica admite:

- Rendimiento de 6,2 Gbps para lecturas secuenciales
- 3375 000 IOPS para lecturas aleatorias pequeñas con latencias de 1 ms o menos
- Recuento máximo de unidades de SSD 144 TB, 3,8 TB o 7,6 TB
- AFF A220 es escalable a más de 1 PB de capacidad efectiva

#### AFF A250 de NetApp

- La capacidad efectiva máxima es de 35 PB con el escalado horizontal máximo de 2-24 nodos (12 parejas de alta disponibilidad).
- Ofrece un aumento del rendimiento de  $\geq 45\%$  con respecto a AFF A220
- 440 000 IOPS lecturas aleatorias a 1 ms
- Basado en el último lanzamiento de ONTAP de NetApp: ONTAP 9.8

- Aprovecha dos Ethernet de 25 GB para alta disponibilidad e interconexión de clúster

### Sistemas E-Series EF de NetApp

EF-Series es una familia de cabinas DE almacenamiento SAN all-flash de gama básica y media que pueden acelerar el acceso a sus datos y ayudarle a obtener valor de ella más rápido con el software SANtricity de NetApp. Estos sistemas ofrecen almacenamiento flash SAS y NVMe y le proporcionan IOPS asequibles hasta extremas, tiempos de respuesta inferiores a 100 microsegundos y un ancho de banda de hasta 44 Gbps, lo cual los convierte en ideales para cargas de trabajo mixtas y aplicaciones exigentes como la inferencia de IA y la computación de alto rendimiento (HPC).

La siguiente figura muestra el sistema de almacenamiento EF280 de NetApp.



### EF280 de NetApp

- Compatibilidad con FC de 32 GB/16 GB, iSCSI de 25 GB/10 GB y SAS de 12 GB
- La capacidad efectiva máxima es 96 unidades que totalizan 1,5 PB
- Rendimiento de 10 Gbps (lecturas secuenciales)
- 300K IOPS (lectura aleatoria)
- El sistema EF280 de NetApp es la cabina all-flash (AFA) más económica de la cartera de NetApp

### EF300 de NetApp

- 24 unidades SSD NVMe para una capacidad total de 367 TB
- Opciones de expansión con un total de 240 HDD NL-SAS, 92 SSD SAS o una combinación de ambos
- 100 GB NVMe/IB, NVMe/roce, Iser/IB y SRP/IB
- NVME DE 32 GB/FC, FCP
- iSCSI de 25 GB
- 20 Gbps (lecturas secuenciales)
- 670K IOPS (lecturas aleatorias)



Para obtener más información, consulte ["Especificaciones técnicas de las cabinas all-flash EF-Series de NetApp EF600, F300, EF570 y EF280"](#).

## ONTAP 9 de NetApp

ONTAP 9.8.1, la última generación del software de gestión del almacenamiento de NetApp, permite a las empresas modernizar su infraestructura y realizar la transición a un centro de datos preparado para el cloud. ONTAP ofrece las mejores capacidades de gestión de datos y permite la gestión y protección de los datos con un solo conjunto de herramientas, sin importar dónde residan. También puede mover los datos libremente a donde sea necesario: El perímetro, el núcleo o el cloud. ONTAP 9.8.1 incluye numerosas funciones que simplifican la gestión de datos, aceleran y protegen los datos esenciales y permiten disfrutar de funcionalidades de infraestructura de nueva generación en arquitecturas de cloud híbrido.

### Simplificar la gestión de los datos

La gestión de los datos es crucial para las operaciones TECNOLÓGICAS empresariales, de modo que se utilicen recursos apropiados para las aplicaciones y conjuntos de datos. ONTAP incluye las siguientes funciones para facilitar y simplificar las operaciones, y reducir el coste total de las operaciones:

- **Compactación de datos inline y deduplicación expandida.** la compactación de datos reduce el espacio perdido dentro de los bloques de almacenamiento, y la deduplicación aumenta significativamente la capacidad efectiva. Esto es aplicable a los datos almacenados localmente y a los datos organizados en niveles en el cloud.
- **Calidad de servicio (AQoS) mínima, máxima y adaptativa.** los controles de calidad de servicio (QoS) granulares ayudan a mantener los niveles de rendimiento para aplicaciones críticas en entornos altamente compartidos.
- **FabricPool de NetApp.** esta función ofrece una organización automática por niveles de datos inactivos a opciones de almacenamiento en cloud privado o público, como la solución de almacenamiento Amazon Web Services (AWS), Azure y StorageGRID de NetApp. Para obtener más información sobre FabricPool, consulte "[TR-4598](#)".

### Acelere y proteja sus datos

ONTAP 9 ofrece niveles superiores de rendimiento y protección de datos, y amplía estas capacidades de las siguientes maneras:

- **Rendimiento y menor latencia.** ONTAP ofrece el rendimiento más alto posible con la menor latencia posible.
- **Protección de datos.** ONTAP ofrece capacidades integradas de protección de datos con administración común en todas las plataformas.
- **Cifrado de volumen de NetApp (NVE).** ONTAP ofrece cifrado nativo a nivel de volumen con compatibilidad para gestión de claves incorporada y externa.
- **Multitenancy y autenticación multifactor.** ONTAP permite compartir recursos de infraestructura con los niveles más altos de seguridad.

### Infraestructura preparada para futuros retos

ONTAP 9 ayuda a satisfacer las exigentes y siempre cambiantes necesidades de la empresa con las siguientes funciones:

- **Escalado sencillo y operaciones no disruptivas.** ONTAP admite la adición sin interrupciones de capacidad a las controladoras existentes y a los clústeres de escalado horizontal. Los clientes pueden empezar a utilizar tecnologías punteras como NVMe y FC 32 GB, sin necesidad de realizar costosas migraciones de datos y sin cortes.
- **Conexión en cloud.** ONTAP es el software de gestión del almacenamiento con mejor conexión en cloud,



con opciones de almacenamiento definido por software (ONTAP Select) e instancias nativas del cloud (Cloud Volumes Service de NetApp) en todos los clouds públicos.

- **Integración con aplicaciones emergentes.** ONTAP ofrece servicios de datos de clase empresarial para plataformas y aplicaciones de última generación, como vehículos autónomos, ciudades inteligentes e Industria 4.0, utilizando la misma infraestructura que soporta las aplicaciones empresariales existentes.

### **SANtricity de NetApp**

SANtricity de NetApp está diseñado para ofrecer un rendimiento líder del sector, fiabilidad y simplicidad en las cabinas all-flash EF-Series y flash híbrido. Logre el máximo rendimiento y aprovechamiento de sus cabinas all-flash EF-Series y flash híbridas de E-Series para aplicaciones de cargas de trabajo pesadas, como análisis de datos, videovigilancia y backup y recuperación de datos. Con SANtricity, los ajustes de configuración, el mantenimiento, la expansión de la capacidad y otras tareas se pueden realizar mientras el almacenamiento sigue online. SANtricity también ofrece una protección de datos superior, supervisión proactiva y seguridad certificada: Todo accesible a través de la sencilla interfaz integrada de System Manager. Para obtener más información, consulte ["Especificaciones técnicas del software SANtricity para E-Series de NetApp"](#).

### **Rendimiento optimizado**

El software SANtricity optimizado para el rendimiento ofrece datos —con una alta tasa de IOPS, un alto rendimiento y una baja latencia— a todas sus aplicaciones de análisis de datos, videovigilancia y backup. Acelere el rendimiento en aplicaciones con baja latencia y una alta tasa de IOPS, y en aplicaciones con requisitos de amplio ancho de banda y un rendimiento alto.

### **Maximice el tiempo de actividad**

Complete todas las tareas de gestión mientras el almacenamiento sigue en línea. Modifique las configuraciones, realice tareas de mantenimiento o amplíe la capacidad sin interrumpir el flujo de I/O. Proporcione los mayores niveles de fiabilidad posibles con funciones automatizadas, opciones de configuración en línea, la tecnología de pools de discos dinámicos (DPP) de vanguardia, entre otras.

### **Esté tranquilo**

El software SANtricity ofrece una protección de datos superior, supervisión proactiva y seguridad certificada: Todo a través de la sencilla interfaz incluida de System Manager. Simplifique las tareas de gestión del almacenamiento. Obtenga la flexibilidad que necesita para realizar un mejor ajuste de todos los sistemas de almacenamiento E-Series. Gestione su sistema E-Series de NetApp en todo momento y en cualquier lugar. Nuestra interfaz integrada basada en web optimiza el flujo de trabajo de gestión.

### **Trident de NetApp**

["Trident"](#) De NetApp es un orquestador de almacenamiento dinámico de código abierto para Docker y Kubernetes que simplifica la creación, la gestión y el consumo de almacenamiento persistente. Trident, una aplicación nativa de Kubernetes, se ejecuta directamente dentro de un clúster de Kubernetes. Trident permite que los clientes implementen sin problemas imágenes de contenedores de DL en el almacenamiento de NetApp y proporciona una experiencia de clase empresarial para implementaciones de contenedores de IA. Los usuarios de Kubernetes (como desarrolladores DE ML y científicos de datos) pueden crear, gestionar y automatizar la orquestación y el clonado para aprovechar las funcionalidades de gestión de datos avanzadas de NetApp, impulsadas por la tecnología de NetApp.

### **Copia y sincronización de NetApp BlueXP**

["Copia y sincronización de BlueXP"](#) Es un servicio de NetApp que ofrece una sincronización de datos rápida y segura. Ya tenga que transferir archivos entre recursos compartidos de archivos NFS o SMB en las



instalaciones, NetApp StorageGRID, NetApp ONTAP S3, NetApp Cloud Volumes Service, Azure NetApp Files, Amazon Simple Storage Service (Amazon S3), Amazon Elastic File System (Amazon EFS), Azure Blob, Google Cloud Storage, o IBM Cloud Object Storage, BlueXP Copy and Sync mueve los archivos a donde los necesites de forma rápida y segura. Una vez transferidos los datos, estarán completamente disponibles para su uso tanto en origen como en destino. BlueXP Copy and Sync sincroniza los datos de forma continua en función de tu programación predefinida, moviendo solo los deltas, por lo que se reducen al mínimo el tiempo y el dinero que se invierten en la replicación de datos. Copia y sincronización de BlueXP es una herramienta de software como servicio (SaaS) extremadamente sencilla de configurar y usar. Las transferencias de datos activadas por BlueXP Copy and Sync se llevan a cabo por agentes de datos. Puedes poner en marcha agentes de datos de BlueXP Copy y Sync en AWS, Azure, Google Cloud Platform o en las instalaciones.

## **Servidores Lenovo ThinkSystem**

Los servidores Lenovo ThinkSystem incluyen hardware, software y servicios innovadores que resuelven los desafíos actuales de los clientes y ofrecen un enfoque de diseño modular, evolutivo y adecuado para su propósito para afrontar los desafíos del futuro. Estos servidores se capitalizan en las mejores tecnologías estándar del sector, junto con innovaciones diferenciadas de Lenovo, para proporcionar la mayor flexibilidad posible en servidores x86.

Las ventajas clave de la implementación de servidores Lenovo ThinkSystem incluyen:

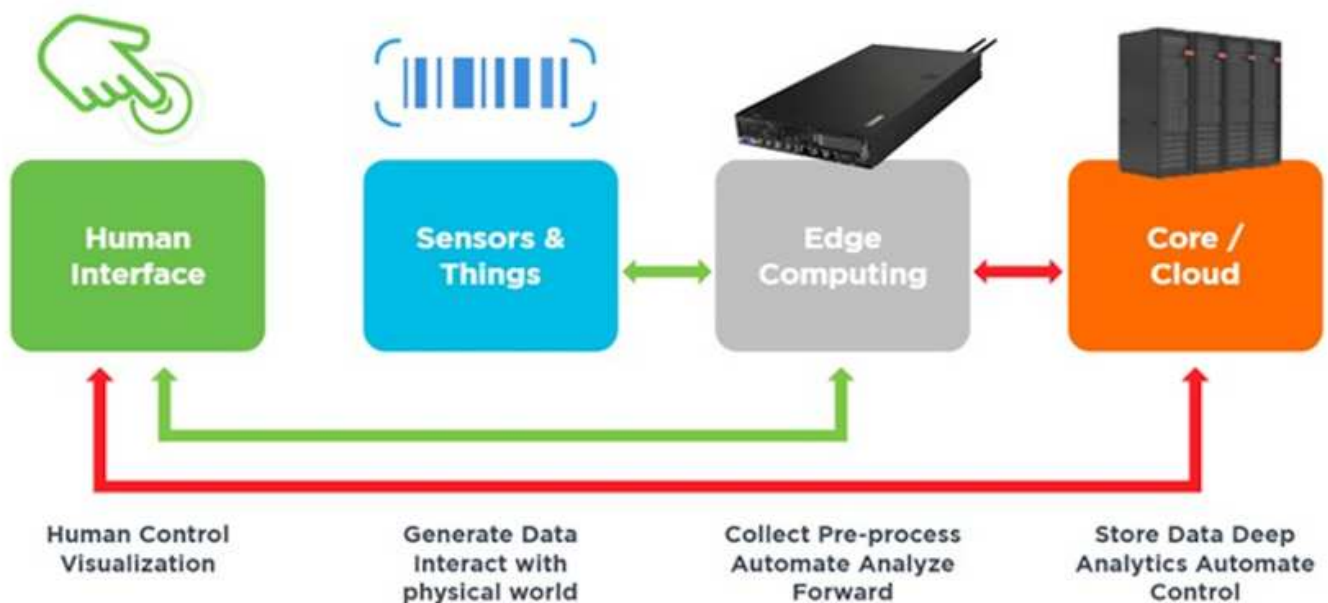
- Diseños modulares y altamente escalables que crecen a medida que lo hace su negocio
- Resiliencia líder en el sector para ahorrar horas de costosos tiempos de inactividad no programados
- Tecnologías flash rápidas para reducir las latencias, acelerar los tiempos de respuesta y gestionar los datos de forma más inteligente en tiempo real

En el ámbito de la IA, Lenovo está adoptando un enfoque práctico para ayudar a las empresas a comprender y adoptar las ventajas DEL APRENDIZAJE AUTOMÁTICO y la IA para sus cargas de trabajo. Los clientes de Lenovo pueden explorar y evaluar las ofertas de IA de Lenovo en los centros de innovación de IA de Lenovo para comprender por completo el valor de su caso de uso en particular. Con el fin de mejorar la rentabilidad de la inversión, este enfoque centrado en el cliente proporciona a los clientes una prueba de concepto para las plataformas de desarrollo de soluciones que están listas para usar y optimizadas para la IA.

## **Servidor Lenovo ThinkSystem SE350 Edge**

La computación perimetral permite analizar los datos de dispositivos de IoT en el extremo de la red antes de enviarlos al centro de datos o al cloud. El sistema ThinkSystem SE350 de Lenovo, como se muestra en la siguiente figura, está diseñado para los requisitos únicos de implementación en el perímetro, con un enfoque en flexibilidad, conectividad, seguridad y capacidad de gestión remota en un factor de forma compacto y reforzado con el medio ambiente.

Con el procesador Intel Xeon D con la flexibilidad necesaria para admitir la aceleración de las cargas de trabajo de IA perimetral, el SE350 se ha diseñado específicamente para afrontar los retos de las implementaciones de servidores en una gran variedad de entornos fuera del centro de datos.



## Rendim. MLPerf

MLPerf es el conjunto de pruebas de rendimiento líder del sector para evaluar el rendimiento de la IA. Cubre muchas áreas de IA aplicada, incluida la clasificación de imágenes, la detección de objetos, las imágenes médicas y el procesamiento del lenguaje natural (NLP). En esta validación, hemos utilizado cargas de trabajo de inferencia v0.7, que es la última iteración de la inferencia MLPerf al finalizar esta validación. La ["Inferencia del rendimiento ML0,7"](#) la suite incluye cuatro nuevas pruebas de rendimiento para centros de datos y sistemas periféricos:

- **BERT.** representación de encoder bidireccional de Transformers (BERT) ajustada para responder preguntas utilizando el conjunto de datos de escuadrón.
- **DLRM.** el modelo de recomendación de aprendizaje profundo (DLRM) es un modelo de personalización y recomendación que se entrena para optimizar las tarifas de clic (CTR).
- **3D U-Net.** la arquitectura 3D U-Net está entrenada en el conjunto de datos de segmentación del tumor cerebral (Brats).

- **RNN-T.** el transductor de red neural recurrente (RNN-T) es un modelo de reconocimiento automático de voz (ASR) que se entrena en un subconjunto de LibriSpeech. Los resultados y el código de la inferencia de MLPerf están disponibles y se liberan públicamente bajo la licencia de Apache. La inferencia MLPerf tiene una división Edge, que admite los siguientes escenarios:
- **Flujo único.** este escenario imita a los sistemas en los que la capacidad de respuesta es un factor crítico, como las consultas de IA sin conexión realizadas en smartphones. Las consultas individuales se envían al sistema y se registran los tiempos de respuesta. como resultado, se indica la latencia del percentil 90 de todas las respuestas.
- **MultiStream.** este punto de referencia es para sistemas que procesan la entrada de varios sensores. Durante la prueba, las consultas se envían a un intervalo de tiempo fijo. Se impone una limitación de calidad de servicio (latencia máxima permitida). La prueba informa del número de flujos que el sistema puede procesar mientras cumple la restricción QoS.
- **Fuera de línea.** este es el escenario más sencillo que cubre las aplicaciones de procesamiento por lotes y la métrica es el procesamiento en muestras por segundo. Todos los datos están disponibles para el sistema y el punto de referencia mide el tiempo que tarda en procesar todas las muestras.

Lenovo ha publicado puntuaciones de inferencia MLPerf para SE350 con T4, el servidor utilizado en este documento. Consulte los resultados en "<https://mlperf.org/inference-results-0-7/>" En la sección "Edge, Closed Division", en la entrada #0.7-145.

## Plan de pruebas

Este documento sigue a la inferencia MLPerf v0.7 "[codificación](#)", Inferencia MLPerf v1.1 "[codificación](#)", y. "[bases de datos](#)". Ejecutamos pruebas de rendimiento MLPerf diseñadas para la inferencia en el perímetro tal y como se definen en la tabla siguiente.

Zona	Tarea	Modelo	Conjunto de datos	Tamaño de QSL	Calidad	Limitación de latencia de múltiples flujos
Visión	Clasificación de imágenes	Resnet50v1.5	ImageNET (224 x 224)	1024	99% del FP32	50 ms
Visión	Detección de objetos (grande)	SSD-ResNet34	COCO (1200x1200)	64	99% del FP32	66 ms
Visión	Detección de objetos (pequeño)	SSD-MobileNetsv1	COCO (300 x 300)	256	99% del FP32	50 ms
Visión	Segmentación de imagen médica	UNET 3D	Brats 2019 (224x224x160 )	16	99% y 99.9% del FP32	n.a.
Voz	Voz a texto	RNNT	Llibrispeech dev-Clean	2513	99% del FP32	n.a.
Idioma	Procesamiento de idiomas	BERT	Escuadra v1.1	10833	99% del FP32	n.a.

En la siguiente tabla se presentan los escenarios de referencia de Edge.

Zona	Tarea	Escenarios
Visión	Clasificación de imágenes	Flujo único, sin conexión, multisequencia
Visión	Detección de objetos (grande)	Flujo único, sin conexión, multisequencia
Visión	Detección de objetos (pequeño)	Flujo único, sin conexión, multisequencia
Visión	Segmentación de imagen médica	Flujo único, sin conexión
Voz	Voz a texto	Flujo único, sin conexión
Idioma	Procesamiento de idiomas	Flujo único, sin conexión

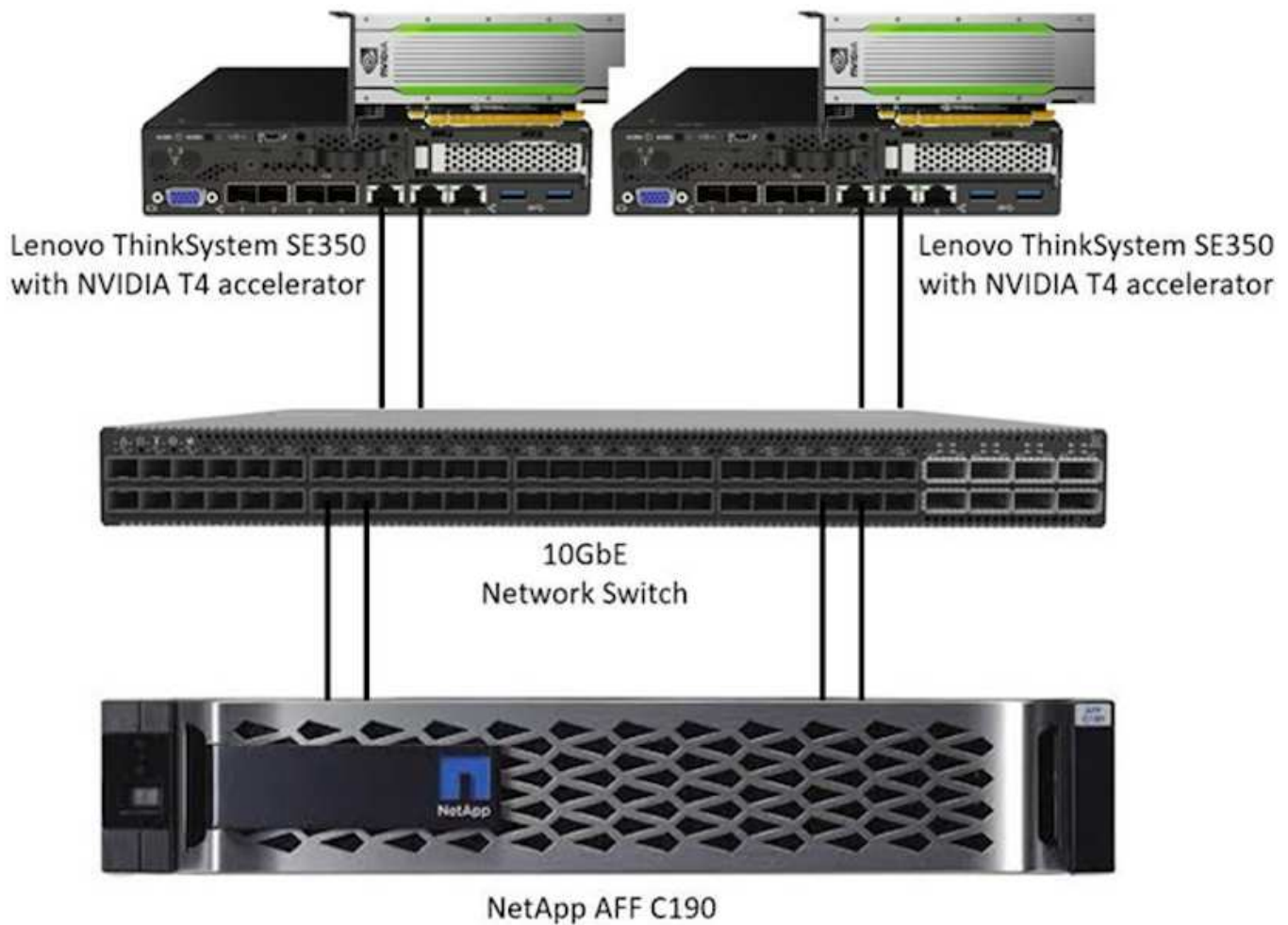
Realizamos estas pruebas de rendimiento utilizando la arquitectura de almacenamiento en red desarrollada en esta validación y comparamos los resultados de las ejecuciones locales en servidores perimetrales que antes se enviaban al MLPerf. La comparación consiste en determinar cuánto impacto tiene el almacenamiento compartido en el rendimiento de la inferencia.

### Configuración de prueba

La figura siguiente muestra la configuración de prueba. Utilizamos el sistema de almacenamiento AFF C190 de NetApp y dos servidores Lenovo ThinkSystem SE350 (cada uno con un acelerador NVIDIA T4). Estos componentes se conectan a través de un switch de red de 10 GbE. El almacenamiento en red contiene conjuntos de datos de validación y pruebas y modelos preentrenados. Los servidores proporcionan una funcionalidad computacional y se accede al almacenamiento mediante un protocolo NFS.

En este apartado se describen las configuraciones probadas, la infraestructura de red, el servidor SE350 y los detalles de aprovisionamiento de almacenamiento. En la siguiente tabla se enumeran los componentes básicos para la arquitectura de la solución.

Componentes de la solución	Detalles
Servidores Lenovo ThinkSystem	<ul style="list-style-type: none"> <li>• 2 servidores SE350 con una tarjeta GPU NVIDIA T4</li> </ul>
	<ul style="list-style-type: none"> <li>• Cada servidor contiene una CPU Intel Xeon D-2123IT con cuatro núcleos físicos que funcionan a 2,20 GHz y 128 GB de RAM</li> </ul>
Sistema de almacenamiento AFF de gama básica de NetApp (par de alta disponibilidad)	<ul style="list-style-type: none"> <li>• Software ONTAP 9 de NetApp</li> <li>• 24 unidades SSD de 960 GB</li> <li>• Protocolo NFS</li> <li>• Un grupo de interfaces por controladora, con cuatro direcciones IP lógicas para puntos de montaje</li> </ul>



La siguiente tabla enumera la configuración del almacenamiento: AFF C190 con 2 ranuras de 24 unidades 2RU.

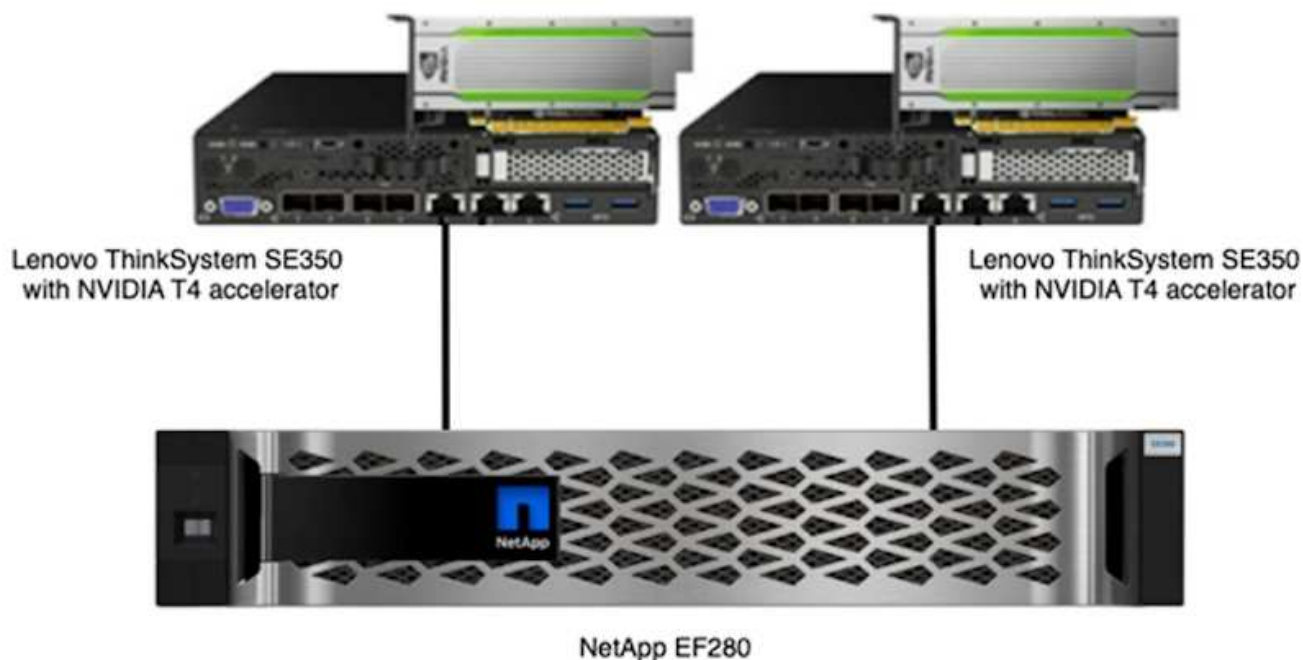
Controladora	Agregado	Volumen FlexGroup	Agregatesize	Volumen	Punto de montaje del sistema operativo
Control1	Agr1	/Netaptenovo_AI_fg	8,42 TIB	15 TB	/netapp_lenovo_fg
Control2	Agr2		8,42 TIB		

La carpeta /netappLenovo\_AI\_fg contiene los conjuntos de datos utilizados para la validación del modelo.

La siguiente figura muestra la configuración de prueba. Utilizamos el sistema de almacenamiento EF280 de NetApp y dos servidores Lenovo ThinkSystem SE350 (cada uno con un acelerador NVIDIA T4). Estos componentes se conectan a través de un switch de red de 10 GbE. El almacenamiento en red contiene conjuntos de datos de validación y pruebas y modelos preentrenados. Los servidores proporcionan una funcionalidad computacional y se accede al almacenamiento mediante un protocolo NFS.

La siguiente tabla enumera la configuración de almacenamiento para EF280.

Controladora	Grupo de volúmenes	Volumen	Volumen	Tamaño DDPsize	Método de conexión
Control1	DDP1	Volumen 1	8,42 TIB	16 TB	SE350-1 a iSCSI LUN 0
Control2		Volumen 2	8,42 TIB		SE350-2 a iSCSI LUN 1



## Procedimiento de prueba

Este apartado describe los procedimientos de prueba utilizados para validar esta solución.

### Configuración de inferencia de IA y sistema operativo

Para AFF C190, utilizamos Ubuntu 18.04 con controladores NVIDIA y docker con soporte para GPU de NVIDIA y usamos MLPerf "codificación" Disponible como parte de la presentación de Lenovo a la inferencia MLPerf v0.7.

Para EF280, utilizamos Ubuntu 20.04 con controladores NVIDIA y docker con soporte para las GPU de NVIDIA y MLPerf "codificación" Disponible como parte de la presentación de Lenovo a la inferencia MLPerf v1.1.

Para configurar la inferencia de IA, siga estos pasos:

1. Descargue los conjuntos de datos que requieren registro, el conjunto de datos ImageNET 2012 Validation, el conjunto de datos Criteo Terabyte y el conjunto de entrenamiento Brats 2019 y, a continuación, descomprima los archivos.
2. Cree un directorio de trabajo con al menos 1 TB y defina una variable ambiental `MLPERF_SCRATCH_PATH` referencia al directorio.

Debe compartir este directorio en el almacenamiento compartido para el caso de uso del almacenamiento de red o en el disco local cuando realice pruebas con datos locales.

3. Ejecute la Marca `prebuild` comando, que crea e inicia el contenedor docker para las tareas de inferencia necesarias.



Los siguientes comandos se ejecutan desde el contenedor docker en ejecución:

- Descargue modelos de IA preformados para tareas de inferencia de MLPerf: `make download_model`
- Descargue conjuntos de datos adicionales que se pueden descargar gratuitamente: `make download_data`
- Preprocesar los datos: `Make preprocess_data`
- Ejecución: `make build`.
- Cree motores de inferencia optimizados para la GPU en servidores informáticos: `make generate_engines`
- Para ejecutar cargas de trabajo de inferencia, ejecute el siguiente (un comando):

```
make run_harness RUN_ARGS="--benchmarks=<BENCHMARKS>  
--scenarios=<SCENARIOS>"
```

### Se ejecuta la inferencia de IA

Se ejecutaron tres tipos de ejecuciones:

- Inferencia de IA de un único servidor con almacenamiento local
- Inferencia de IA de un único servidor con el almacenamiento en red
- Inferencia de la IA en varios servidores con el almacenamiento en red

### Resultados de la prueba

Se llevaron a cabo una multitud de pruebas para evaluar el rendimiento de la arquitectura propuesta.

Existen seis cargas de trabajo diferentes (clasificación de imágenes, detección de objetos [pequeño], detección de objetos [grande], imágenes médicas, voz a texto, Y procesamiento de lenguaje natural [NLP]), que se puede ejecutar en tres escenarios diferentes: Sin conexión, flujo único y multisequencia.



El último escenario se implementa sólo para la clasificación de imágenes y la detección de objetos.

Esto permite 15 cargas de trabajo posibles, todas ellas probadas en tres configuraciones diferentes:

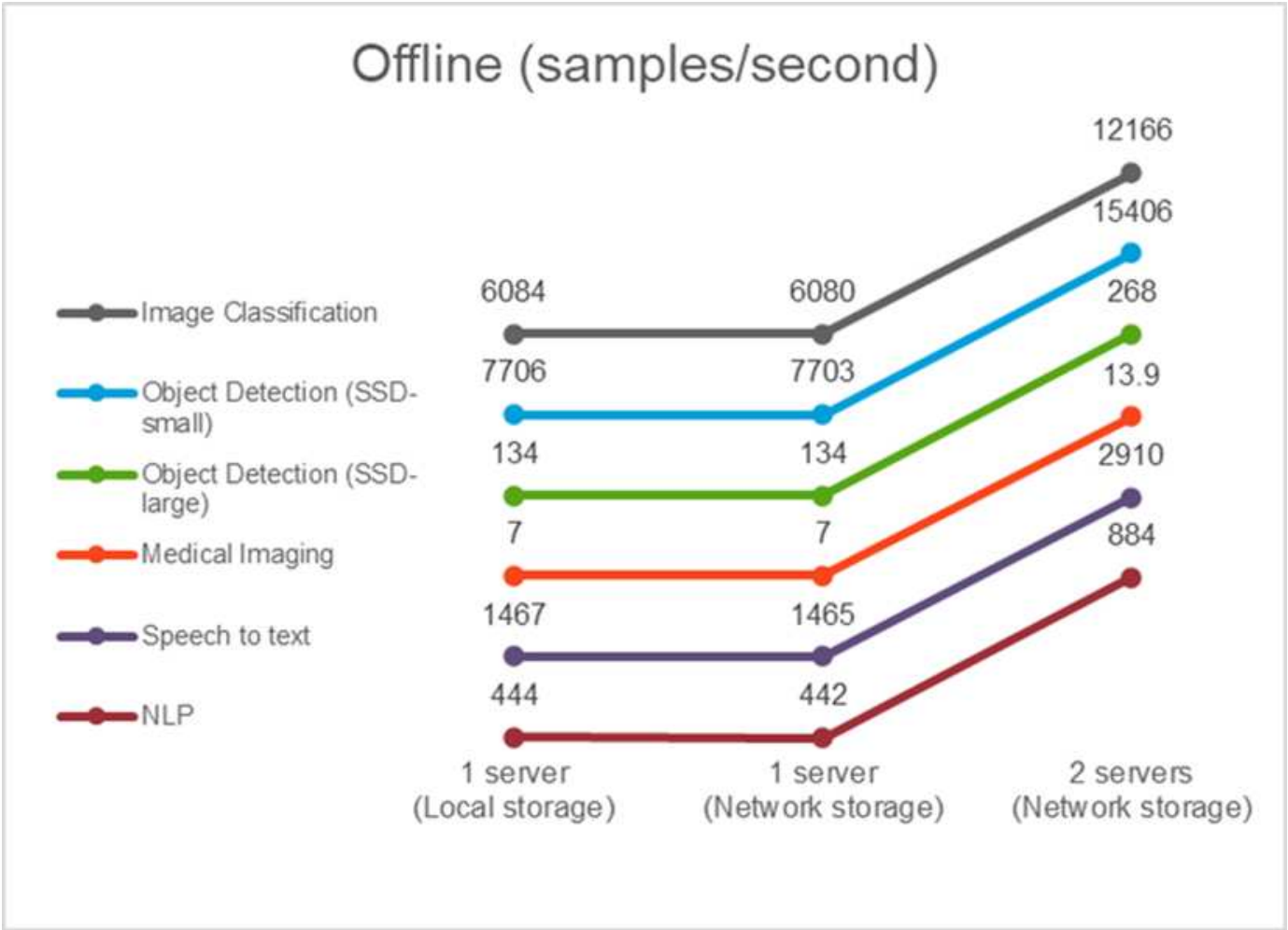
- Un único servidor/almacenamiento local
- Almacenamiento único de servidor/red
- Almacenamiento en red/multiservidor



Los resultados se describen en las siguientes secciones.

**Inferencia de la IA en un escenario sin conexión para AFF**

En este escenario, todos los datos estaban disponibles para el servidor y se midió el tiempo que se tardaba en procesar todas las muestras. Reportamos los anchos de banda en muestras por segundo como resultados de las pruebas. Cuando se utilizaron más de un servidor informático, hemos indicado el total de ancho de banda total en todos los servidores. En la siguiente figura se muestran los resultados de los tres casos de uso. Para el caso de dos servidores, reportamos un ancho de banda combinado de ambos servidores.

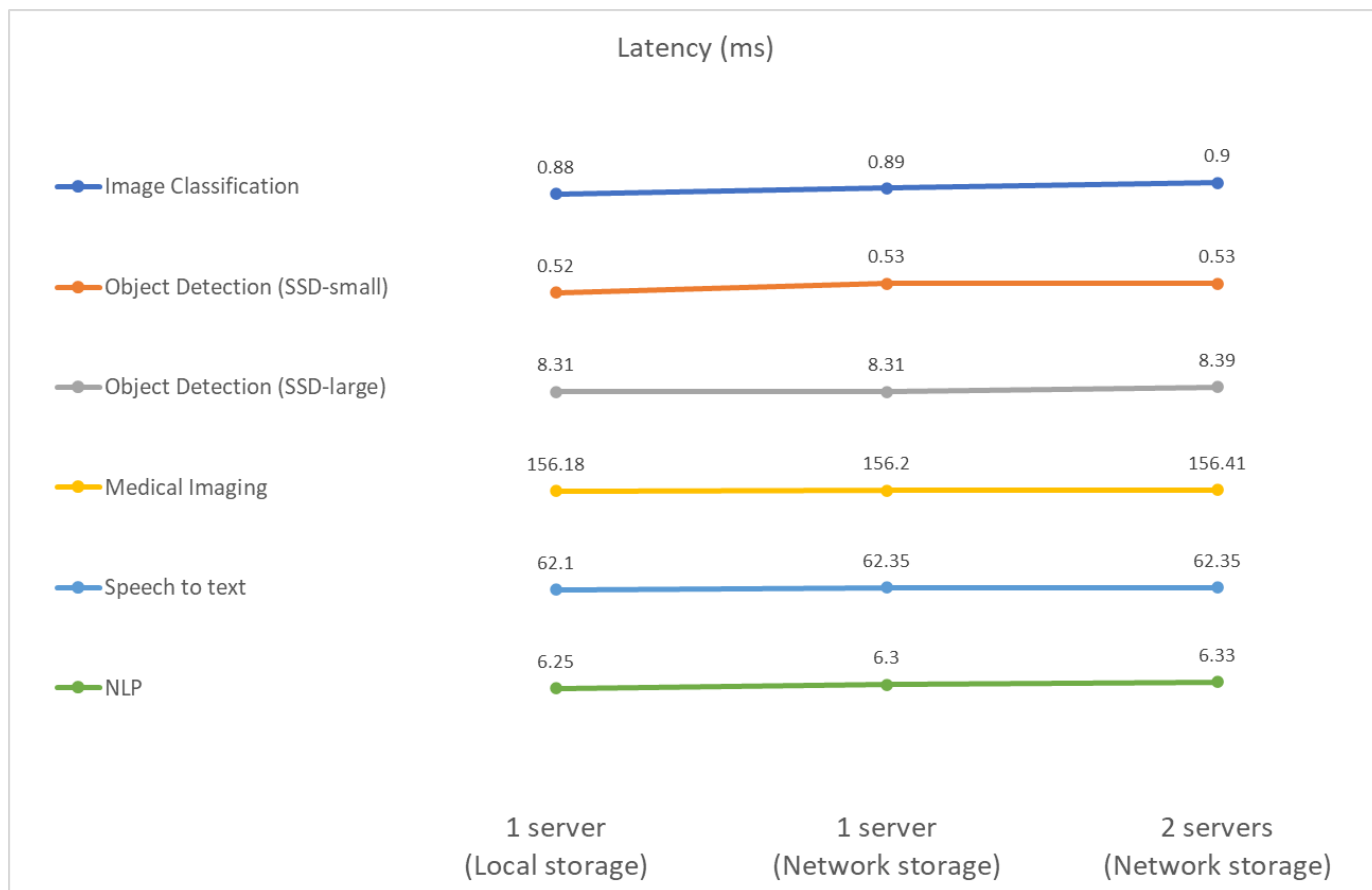


Los resultados muestran que el almacenamiento en red no afecta de manera negativa al rendimiento; el cambio es mínimo y, en algunas tareas, no se encuentra ninguno. Cuando se añade el segundo servidor, el ancho de banda total se duplica exactamente o, en el peor de los casos, el cambio es inferior al 1%.

**Inferencia de IA en un escenario de flujo único para AFF**

Esta prueba mide la latencia. Para el caso de varios servidores informáticos, reportamos la latencia media. Los resultados del conjunto de tareas se indican en la siguiente figura. Para el caso de dos servidores, hemos registrado la latencia media de ambos servidores.

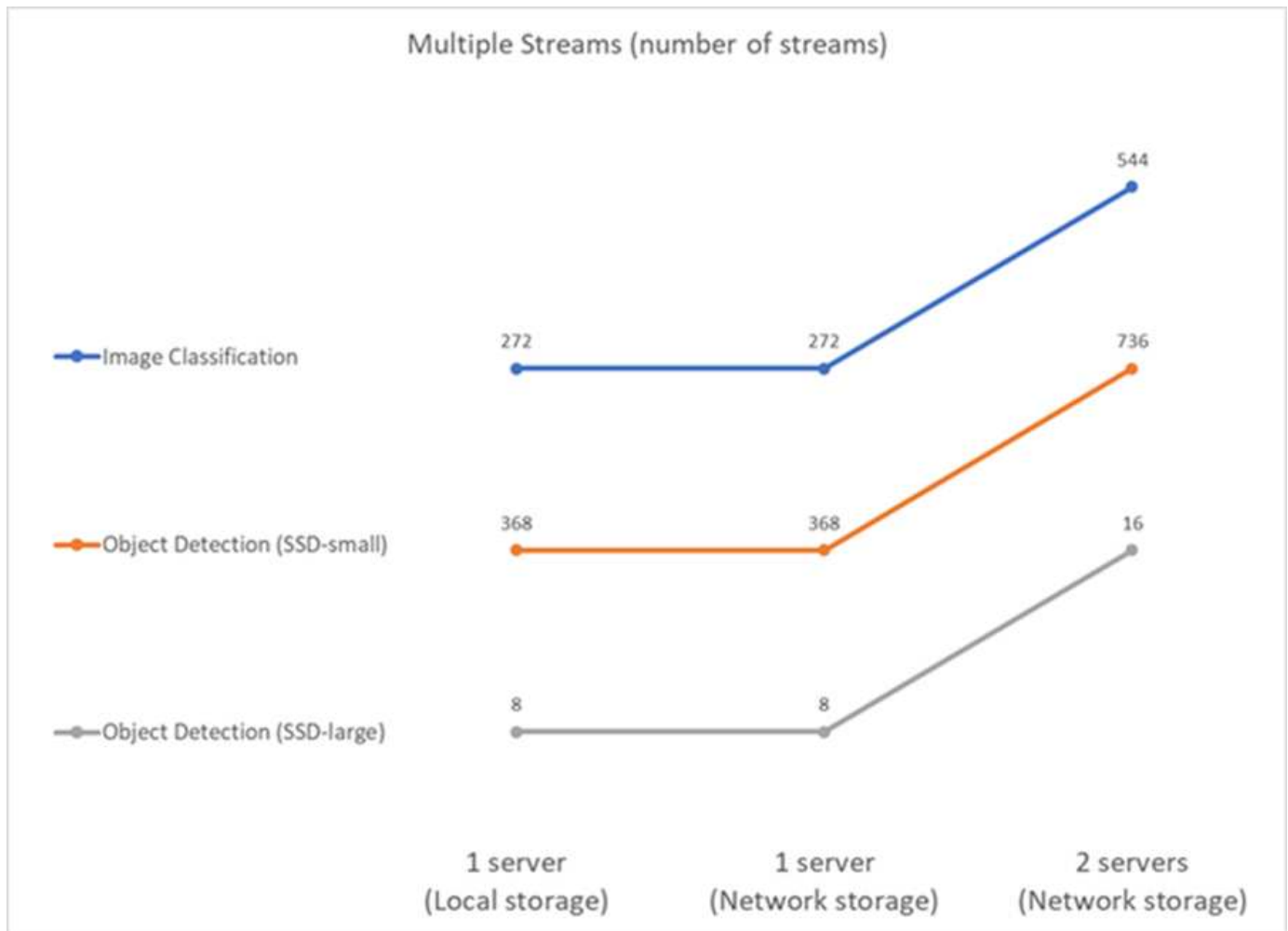




Los resultados, de nuevo, muestran que el almacenamiento de red es suficiente para manejar las tareas. La diferencia entre el almacenamiento local y el de red en un caso de servidor es mínima o ninguna. De igual modo, cuando dos servidores utilizan el mismo almacenamiento, la latencia de ambos servidores se mantiene igual o se modifica en una cantidad muy pequeña.

### Inferencia de IA en un escenario de transmisión múltiple para AFF

En este caso, el resultado es el número de flujos que el sistema puede manejar mientras se satisface la restricción QoS. Por lo tanto, el resultado siempre es un entero. En más de un servidor, reportamos el número total de flujos sumados en todos los servidores. No todas las cargas de trabajo admiten este escenario, pero hemos ejecutado las que lo hacen. Los resultados de nuestras pruebas se resumen en la siguiente figura. Para el caso de dos servidores, reportamos el número combinado de flujos de ambos servidores.



Los resultados muestran un rendimiento perfecto de la configuración; el almacenamiento local y en red ofrecen los mismos resultados y al añadir el segundo servidor se duplica el número de flujos que puede gestionar la configuración propuesta.

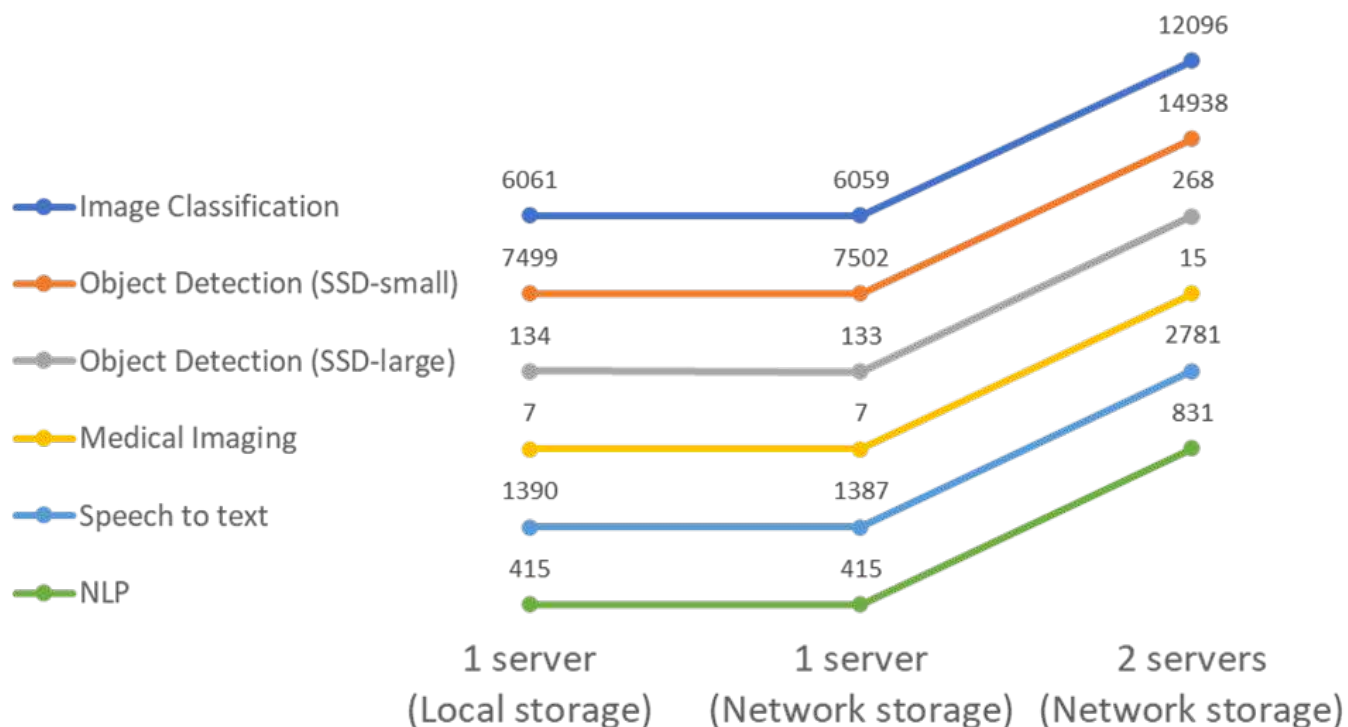
#### Resultados de la prueba para EF

Se llevaron a cabo una multitud de pruebas para evaluar el rendimiento de la arquitectura propuesta. Existen seis cargas de trabajo diferentes (clasificación de imágenes, detección de objetos [pequeño], detección de objetos [grande], imágenes médicas, voz a texto, Y procesamiento de lenguaje natural [NLP]), que se ejecutaron en dos escenarios diferentes: Offline y de flujo único. Los resultados se describen en las siguientes secciones.

#### Inferencia de la IA en un escenario sin conexión para EF

En este escenario, todos los datos estaban disponibles para el servidor y se midió el tiempo que se tardaba en procesar todas las muestras. Reportamos los anchos de banda en muestras por segundo como resultados de las pruebas. En las ejecuciones de un solo nodo reportamos la media de ambos servidores, mientras que en dos ejecuciones de servidor reportamos un ancho de banda total resumido en todos los servidores. Los resultados de los casos de uso se muestran en la siguiente figura.

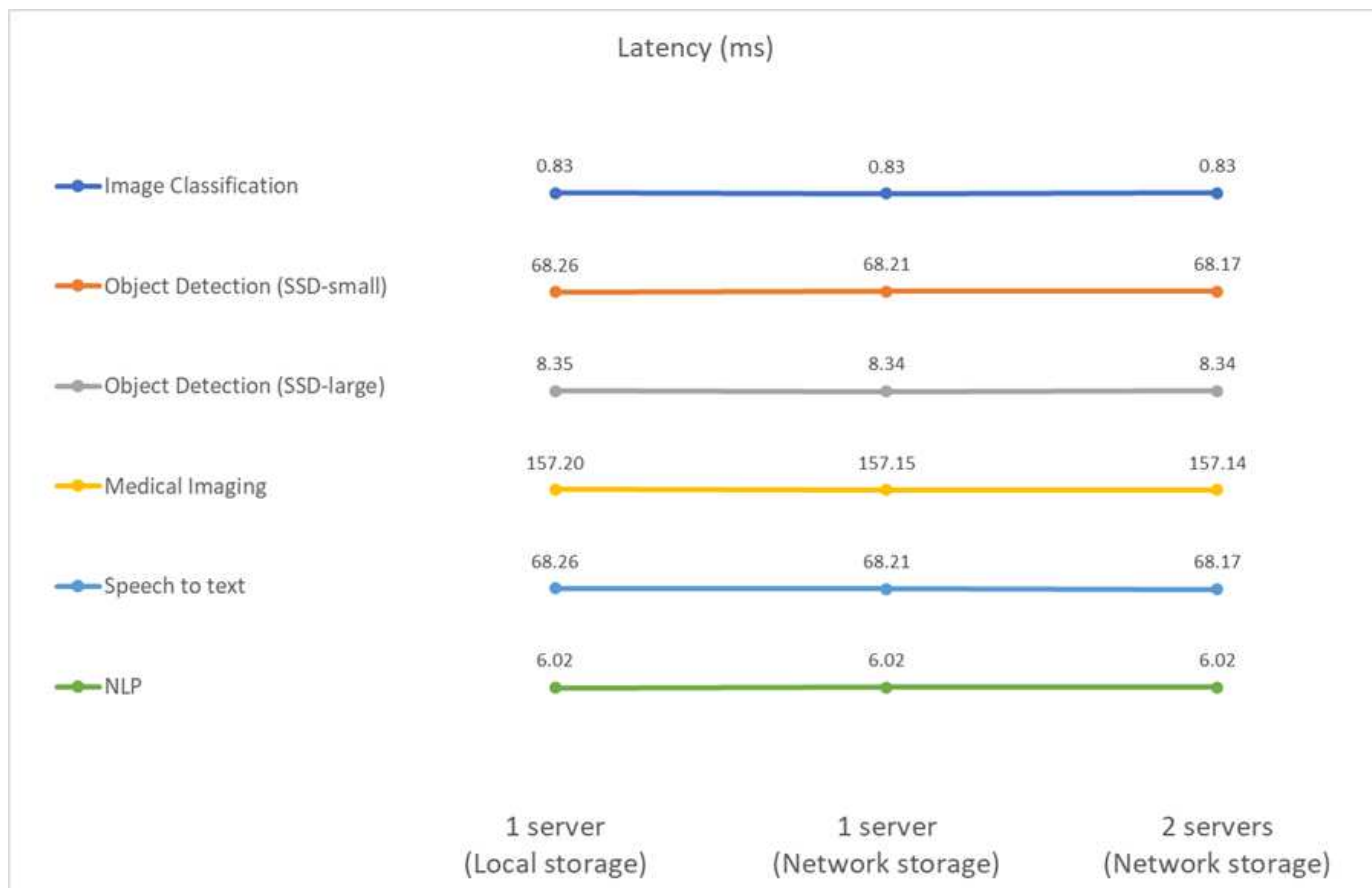
## Offline (samples/second)



Los resultados muestran que el almacenamiento en red no afecta de manera negativa al rendimiento; el cambio es mínimo y, en algunas tareas, no se encuentra ninguno. Cuando se añade el segundo servidor, el ancho de banda total se duplica exactamente o, en el peor de los casos, el cambio es inferior al 1%.

### Inferencia de IA en un escenario de flujo único para EF

Esta prueba mide la latencia. En todos los casos, reportamos una latencia media en todos los servidores involucrados en las ejecuciones. Los resultados para el conjunto de tareas se dan.



Los resultados muestran de nuevo que el almacenamiento en red es suficiente para manejar las tareas. La diferencia entre el almacenamiento local y el de red en un caso de servidor es mínima o ninguna. De igual modo, cuando dos servidores utilizan el mismo almacenamiento, la latencia de ambos servidores se mantiene igual o se modifica en una cantidad muy pequeña.

### Opciones de ajuste de tamaño de arquitectura

Puede ajustar la configuración utilizada para la validación y adaptarla a otros casos prácticos.

#### Servidor de computación

Utilizamos una CPU Intel Xeon D-2123IT, que es el nivel más bajo de CPU compatible con SE350, con cuatro núcleos físicos y 60 W TDP. Aunque el servidor no admite la sustitución de CPU, se puede pedir con una CPU más potente. La CPU más alta admitida es Intel Xeon D-2183IT con 16 núcleos, 100 W con 2,20 GHz. Esto aumenta considerablemente la capacidad computacional de la CPU. Mientras que la CPU no era un cuello de botella para ejecutar las cargas de trabajo de inferencia en sí, ayuda con el procesamiento de datos y otras tareas relacionadas con la inferencia. Actualmente, NVIDIA T4 es la única GPU disponible para casos de uso periféricos; por lo tanto, actualmente, no es posible actualizar o degradar el GPU.

#### Almacenamiento compartido

Para las pruebas y la validación, se utilizó el sistema AFF C190 de NetApp, que tiene una capacidad de almacenamiento máxima de 50,5 TB, una rendimiento de 4,4 GB/s para lecturas secuenciales y 230 000 IOPS para pequeñas lecturas aleatorias, con los fines de este documento y se ha demostrado ser apto para cargas de trabajo de inferencia periférica.

No obstante, si se requiere más capacidad de almacenamiento o velocidades de red más rápidas, deberá utilizar el o el sistema AFF A220 de NetApp "AFF A250 de NetApp" sistemas de almacenamiento. Además, el sistema EF280 de NetApp, que tiene una capacidad máxima de 1,5 PB, 10 Gbps de ancho de banda, también se utilizó para la validación de esta solución. Si prefiere más capacidad de almacenamiento con un mayor ancho de banda, "EF300 de NetApp" puede utilizarse.

## Conclusión

La automatización impulsada por IA y la informática perimetral es un enfoque líder que ayuda a las organizaciones empresariales a lograr la transformación digital y maximizar la eficiencia y la seguridad operativas. En el caso de los entornos periféricos, los datos se procesan con mayor rapidez, porque no tienen que desplazarse hacia y desde un centro de datos. Por lo tanto, el coste asociado al envío de datos a los centros de datos o al cloud se ve disminuido. La reducción de la latencia y el aumento de la velocidad pueden ser beneficiosas cuando las empresas deben tomar decisiones casi en tiempo real mediante modelos de inferencia de IA implementados en el perímetro.

Los sistemas de almacenamiento de NetApp ofrecen el mismo rendimiento o mejor que el almacenamiento SSD local y ofrecen las siguientes ventajas a los científicos de datos, ingenieros de datos, desarrolladores de IA/ML y responsables de la toma de decisiones EMPRESARIALES o TECNOLÓGICAS:

- Uso compartido de datos sin esfuerzo entre sistemas de IA, análisis y otros sistemas de negocio cruciales. Este uso compartido de datos reduce la sobrecarga de la infraestructura, mejora el rendimiento y optimiza la gestión de datos en toda la empresa.
- Almacenamiento y cálculo escalables de forma independiente para minimizar los costes y mejorar el uso de recursos.
- Flujos de trabajo de desarrollo y puesta en marcha optimizados que utilizan copias Snapshot integradas y clones para espacios de trabajo de usuario instantáneos con gestión eficiente del espacio, control de versiones integrado y una puesta en marcha automatizada.
- Protección de datos de clase empresarial para la recuperación ante desastres y la continuidad del negocio. La solución de NetApp y Lenovo presentada en este documento es una arquitectura flexible y de escalado horizontal ideal para puestas en marcha de inferencia de IA en el perímetro.

## Reconocimientos

- J.J. Falkanger, Sr. Director de Soluciones HPC e IA de Lenovo
- Dave Arnette, ingeniero técnico de marketing, NetApp
- Joey Parnell, responsable técnico de soluciones de IA de E-Series de NetApp
- Cody Harryman, ingeniero de control de calidad, NetApp

## Dónde encontrar información adicional

Para obtener más información sobre la información descrita en este documento, consulte los siguientes documentos y/o sitios web:

- Página de producto de las cabinas AFF a-Series de NetApp

["https://www.netapp.com/data-storage/aff-a-series/"](https://www.netapp.com/data-storage/aff-a-series/)

- Software de gestión de datos ONTAP de NetApp: Biblioteca de información de ONTAP 9

<http://mysupport.netapp.com/documentation/productlibrary/index.html?productID=62286>

- TR-4727: Introducción a EF-Series de NetApp

<https://www.netapp.com/pdf.html?item=/media/17179-tr4727pdf.pdf>

- Especificaciones técnicas del software SANtricity para E-Series de NetApp

<https://www.netapp.com/pdf.html?item=/media/19775-ds-3171-66862.pdf>

- Almacenamiento persistente de NetApp para contenedores: Trident de NetApp

["https://netapp.io/persistent-storage-provisioner-for-kubernetes/"](https://netapp.io/persistent-storage-provisioner-for-kubernetes/)

- Rendim. MLPerf

- ["https://mlcommons.org/en/"](https://mlcommons.org/en/)

- ["http://www.image-net.org/"](http://www.image-net.org/)

- ["https://mlcommons.org/en/news/mlperf-inference-v11/"](https://mlcommons.org/en/news/mlperf-inference-v11/)

- Copia y sincronización de NetApp BlueXP

["https://docs.netapp.com/us-en/occm/concept\\_cloud\\_sync.html#how-cloud-sync-works"](https://docs.netapp.com/us-en/occm/concept_cloud_sync.html#how-cloud-sync-works)

- Prueba de rendimiento de TensorFlow

["https://github.com/tensorflow/benchmarks"](https://github.com/tensorflow/benchmarks)

- Servidor Lenovo ThinkSystem SE350 Edge

["https://lenovopress.com/lp1168"](https://lenovopress.com/lp1168)

- Sistema de almacenamiento flash unificado Lenovo ThinkSystem DM5100F

["https://lenovopress.com/lp1365-thinksystem-dm5100f-unified-flash-storage-array"](https://lenovopress.com/lp1365-thinksystem-dm5100f-unified-flash-storage-array)

## **Artículo técnico WP-7328: La IA conversacional de NetApp con NVIDIA Jarvis**

Rick Huang, Sung-han Lin, NetApp Davide Onofrío, NVIDIA

La familia de sistemas NVIDIA DGX está compuesta por los primeros sistemas basados en inteligencia artificial (IA) integrados del mundo, creados expresamente para la IA empresarial. Los sistemas de almacenamiento AFF de NetApp proporcionan un rendimiento extremo y funcionalidades de gestión de datos de cloud híbrido líderes en el sector. NetApp y NVIDIA se han asociado para crear la arquitectura de referencia de IA ONTAP de NetApp, una solución lista para usar para cargas de trabajo de IA y aprendizaje automático (ML) que ofrece rendimiento, fiabilidad y soporte de clase empresarial.

Este whitepaper proporciona instrucciones direccionales a los clientes para crear sistemas de IA conversacionales que respaldan diferentes casos de uso en diversos sectores verticales. Incluye información sobre la implementación del sistema mediante NVIDIA Jarvis. Las pruebas se llevaron a cabo con una estación DGX de NVIDIA y un sistema de almacenamiento AFF A220 de NetApp.

El público objetivo de la solución incluye los siguientes grupos:

- Arquitectos empresariales que diseñan soluciones para el desarrollo de modelos de IA y software para casos de uso de IA conversacionales como un asistente de ventas al por menor virtual
- Científicos de datos que buscan formas eficientes de lograr los objetivos de desarrollo del modelado del lenguaje
- Ingenieros de datos encargados de mantener y procesar datos de texto, como preguntas de clientes y transcripciones de diálogo
- Responsables de la toma de decisiones tecnológicas y ejecutivas y líderes de negocio interesados en transformar la experiencia de IA conversacional y lograr el plazo de comercialización más rápido de iniciativas de IA

## **Descripción general de la solución**

### **NetApp ONTAP AI y copia y sincronización de BlueXP**

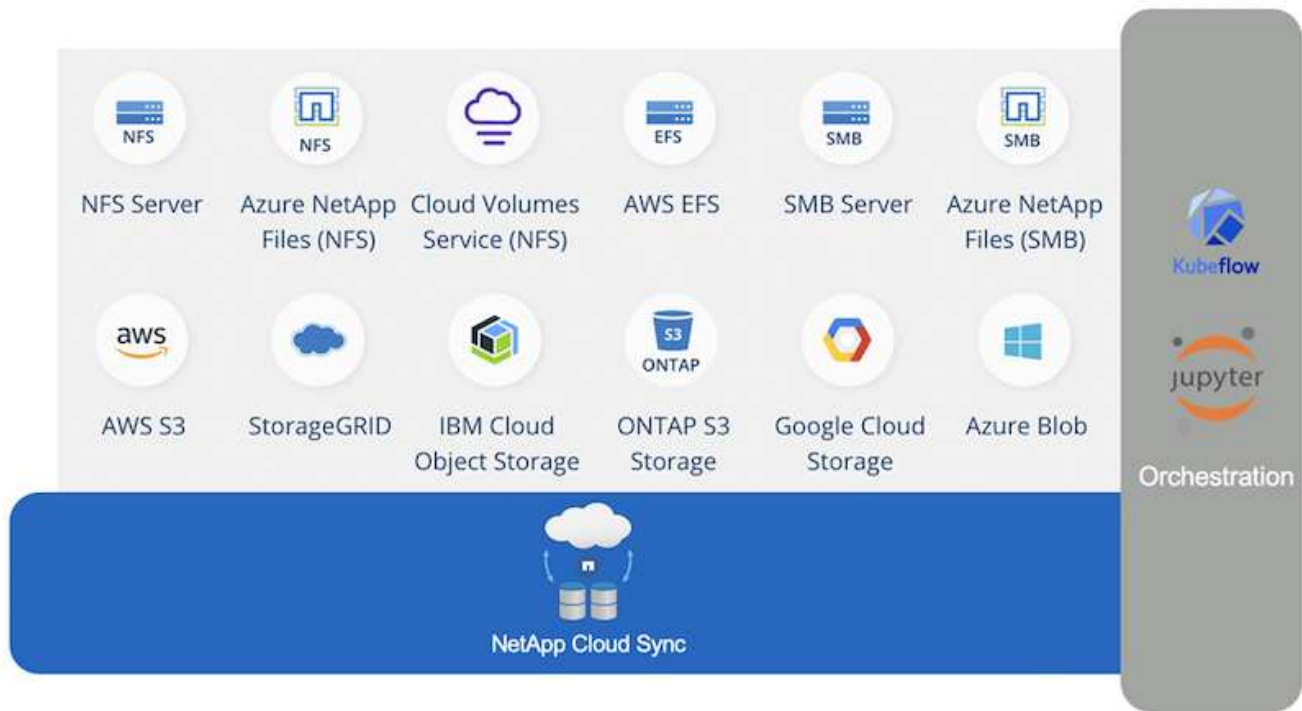
La arquitectura de ONTAP AI de NetApp, impulsada por los sistemas NVIDIA DGX y los sistemas de almacenamiento conectados al cloud de NetApp, fue desarrollada y verificada por NetApp y NVIDIA. Esta arquitectura de referencia proporciona a las organizaciones DE TI las siguientes ventajas:

- Elimina las complejidades de diseño
- Permite un escalado independiente de las capacidades de computación y almacenamiento
- Permite a los clientes empezar con poco y escalar sin problemas
- Ofrece una gama de opciones de almacenamiento para diferentes funcionalidades de rendimiento y rentabilidad. ONTAP AI integra a la perfección los sistemas de almacenamiento DGX y AFF A220 de NetApp con redes de vanguardia. Los sistemas ONTAP AI y DGX de NetApp simplifican las puestas en marcha de IA al eliminar complejidades y conjeturas en la fase de diseño. Los clientes pueden empezar poco a poco y aumentar sus sistemas de forma ininterrumpida, a la vez que gestionan de forma inteligente datos entre el perímetro, el núcleo y el cloud.

La función Copy and Sync de NetApp BlueXP te permite mover datos fácilmente a través de diversos protocolos, ya sea entre dos unidades NFS, dos unidades CIFS o una unidad de archivos y almacenamiento de Amazon S3, Amazon Elastic File System (EFS) o Azure Blob. El funcionamiento activo-activo permite seguir trabajando tanto con el origen como con el objetivo al mismo tiempo, sincronizando de forma incremental los cambios de datos cuando sea necesario. Al permitirte mover y sincronizar gradualmente los datos entre cualquier sistema de origen y de destino, ya sea on-premises o en la nube, BlueXP Copy and Sync ofrece una amplia variedad de nuevas formas de utilizar los datos. Es muy fácil llegar a la migración de datos entre sistemas locales, integración en el cloud y migración al cloud, o colaboración y análisis de datos. La siguiente figura muestra los orígenes y destinos disponibles.

En sistemas de IA conversacional, los desarrolladores pueden aprovechar la función Copy and Sync de BlueXP para archivar el historial de conversaciones desde el cloud a los centros de datos y permitir el entrenamiento sin conexión de modelos de procesamiento del lenguaje natural (NLP). Gracias a los modelos de formación que permiten reconocer más intenciones, el sistema de IA conversacional estará mejor equipado para gestionar preguntas más complejas por parte de los usuarios finales.

### **NVIDIA Jarvis Multimodal Framework**



"**NVIDIA Jarvis**" Es un marco integral para crear servicios de IA que puedan conversación. Incluye los siguientes servicios optimizados para GPU:

- Reconocimiento automático de voz (ASR)
- Comprensión del lenguaje natural (NLU)
- Integración con servicios logísticos específicos de dominio
- Texto a voz (TTS)
- Los servicios basados en Jarvis utilizan modelos de aprendizaje profundo de vanguardia para hacer frente a la compleja y desafiante tarea de IA conversacional en tiempo real. Para permitir la interacción natural y en tiempo real con un usuario final, los modelos deben completar el cálculo en menos de 300 milisegundos. Las interacciones naturales son difíciles, y requieren la integración sensorial multimodal. Los gasoductos son también complejos y requieren coordinación entre los servicios mencionados.

Jarvis es un marco de aplicaciones totalmente acelerado para la creación de servicios de IA conversacionales multimodales que usan una canalización de aprendizaje profundo integral. El marco Jarvis incluye modelos de IA conversacionales, herramientas y servicios integrales optimizados para tareas de habla, visión y NLU. Además de los servicios de IA, Jarvis le permite fusionar simultáneamente la visión, el audio y otras entradas de sensores para proporcionar capacidades como conversaciones multiusuario y multicontexto en aplicaciones como asistentes virtuales, diarización multiusuario y asistentes de centros de llamadas.

#### **NVIDIA Nemo**

"**NVIDIA Nemo**" Es un kit de herramientas Python de código abierto para crear, formar y ajustar modelos de IA de última generación acelerados por GPU y conversacionales que usan interfaces de programación de aplicaciones (API) fáciles de usar. Nemo ejecuta una computación de precisión mixta mediante núcleos tensores en GPU de NVIDIA y puede escalarse verticalmente hasta varias GPU fácilmente para ofrecer el máximo rendimiento de entrenamiento posible. Nemo se utiliza para crear modelos para aplicaciones ASR, NLP y TTS en tiempo real, como transcripciones de videollamadas, asistentes inteligentes de vídeo y soporte de centros de llamadas automatizados en diferentes sectores verticales, incluyendo servicios sanitarios, financieros, minoristas y telecomunicaciones.



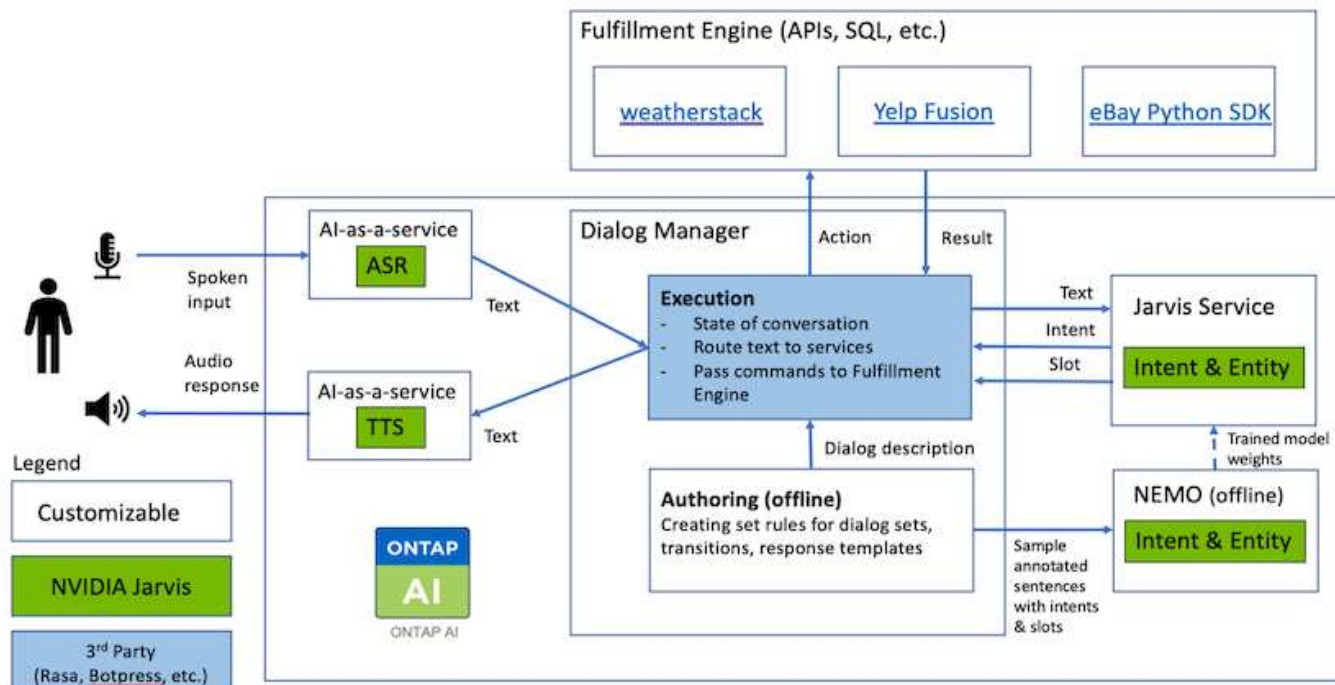
Utilizamos Nemo para entrenar modelos que reconocen intentos complejos de preguntas de usuario en el historial de conversaciones archivadas. Esta capacitación amplía las capacidades del asistente virtual al por menor más allá de lo que Jarvis apoya como se ofrece.

### Resumen del caso de uso de venta al por menor

Con NVIDIA Jarvis, creamos un asistente de venta al por menor virtual que acepta comentarios o comentarios y responde a preguntas sobre el tiempo, puntos de interés y precios de inventario. El sistema de IA conversacional puede recordar el flujo de conversaciones, por ejemplo, formular una pregunta de seguimiento si el usuario no especifica la ubicación para el clima o los puntos de interés. El sistema también reconoce entidades complejas como “comida tailandesa” o “memoria portátil”. Entiende preguntas de lenguaje natural como “¿lloverá la próxima semana en los Ángeles?” Una demostración del asistente virtual al por menor se puede encontrar en ["Personalizar Estados y flujos para el caso de uso de comercio minorista"](#).

### Tecnología de soluciones

La siguiente figura ilustra la arquitectura del sistema de IA conversacional propuesta. Puede interactuar con el sistema con la señal de voz o con la entrada de texto. Si se detecta una entrada hablada, Jarvis AI-as-a-service (AlaaS) realiza ASR para producir texto para Dialog Manager. Dialog Manager recuerda los estados de la conversación, enruta el texto a los servicios correspondientes y pasa los comandos al motor de ejecución. El servicio Jarvis NLP toma texto, reconoce intentos y entidades, y envía esos intentos y ranuras de entidad de vuelta a Dialog Manager, que luego envía Acción al motor de cumplimiento. El motor de ejecución consta de API o bases de datos SQL de terceros que responden a las consultas de los usuarios. Después de recibir el resultado de Logística Engine, Dialog Manager enruta el texto a Jarvis TTS AlaaS para producir una respuesta de audio para el usuario final. Podemos archivar el historial de conversaciones, anotar frases con intentos y ranuras para la formación de Nemo, de forma que el servicio de NLP mejore a medida que más usuarios interactúan con el sistema.



## Requisitos de hardware

Esta solución se validó con un sistema de almacenamiento DGX Station y uno AFF A220. Jarvis requiere una GPU T4 o V100 para realizar cálculos profundos de la red neuronal.

En la siguiente tabla se enumeran los componentes de hardware necesarios para implementar la solución según se ha probado.

Hardware subyacente	Cantidad
GPU T4 O V100	1
Estación DGX de NVIDIA	1

## Requisitos de software

En la siguiente tabla se enumeran los componentes de software necesarios para implementar la solución según se ha probado.

De NetApp	Versión u otra información
Software de gestión de datos ONTAP de NetApp	9.6
Firmware de switch Cisco NX-OS	7.0(3)I6(1)
SO DGX DE NVIDIA	4.0.4 - Ubuntu 18.04 LTS
NVIDIA Jarvis Framework	EA v0.2
NVIDIA Nemo	<a href="https://nvcv.io/nvidia/nemo:v0.10">nvcv.io/nvidia/nemo:v0.10</a>
Plataforma contenedora Docker	18.06.1-ce [e68fc7a]

## Descripción general

En esta sección se ofrece información detallada sobre la implantación del asistente de venta al por menor virtual.

### Despliegue de Jarvis

Puede registrarse para "[Programa Jarvis Early Access](#)" Para obtener acceso a contenedores Jarvis en NVIDIA GPU Cloud (NGC). Después de recibir credenciales de NVIDIA, puede implementar Jarvis siguiendo los pasos siguientes:

1. Firma a NGC.
2. Establezca la organización en NGC: `ea-2-jarvis`.
3. Localice Jarvis EA v0.2 activos: Contenedores Jarvis están en `Private Registry > Organization Containers`.
4. Seleccione Jarvis: Desplácese a `Model Scripts` y haga clic en `Jarvis Quick Start`
5. Compruebe que todos los activos funcionan correctamente.
6. Busque la documentación para crear sus propias aplicaciones: Los archivos PDF se pueden encontrar en `Model Scripts > Jarvis Documentation > File Browser`.

### Personalizar Estados y flujos para el caso de uso de comercio minorista

Puede personalizar los Estados y flujos de Dialog Manager para sus casos de uso específicos. En nuestro ejemplo de venta al por menor, tenemos los siguientes cuatro archivos yaml para dirigir la conversación según diferentes intentos.

S la siguiente lista de nombres de archivo y descripción de cada archivo:

- `main_flow.yml`: Define los principales flujos y estados de conversación y dirige el flujo a los otros tres archivos yaml cuando sea necesario.
- `retail_flow.yml`: Contiene estados relacionados con preguntas sobre puntos de interés o minoristas. El sistema proporciona la información de la tienda más cercana o el precio de un artículo determinado.
- `weather_flow.yml`: Contiene estados relacionados con las preguntas sobre el clima. Si no se puede determinar la ubicación, el sistema hace una pregunta de seguimiento para aclarar.
- `error_flow.yml`: Trata los casos en los que las intenciones del usuario no entran en los tres archivos yaml anteriores. Después de mostrar un mensaje de error, el sistema vuelve a enruta para aceptar preguntas de usuario. las siguientes secciones contienen las definiciones detalladas de estos archivos yaml.

#### main\_flow.yml

```
name: JarvisRetail
intent_transitions:
  jarvis_error: error
  price_check: retail_price_check
  inventory_check: retail_inventory_check
  store_location: retail_store_location
```

```

weather.weather: weather
weather.temperature: temperature
weather.sunny: sunny
weather.cloudy: cloudy
weather.snow: snow
weather.rainfall: rain
weather.snow_yes_no: snowfall
weather.rainfall_yes_no: rainfall
weather.temperature_yes_no: tempyesno
weather.humidity: humidity
weather.humidity_yes_no: humidity
navigation.startnavigationpoi: retail # Transitions should be context
and slot based. Redirecting for now.
navigation.geteta: retail
navigation.showdirection: retail
navigation.showmappoi: idk_what_you_talkin_about
nomatch.none: idk_what_you_talkin_about
states:
  init:
    type: message_text
    properties:
      text: "Hi, welcome to NARA retail and weather service. How can I
help you?"
    input_intent:
      type: input_context
      properties:
        nlp_type: jarvis
        entities:
          intent: dontcare
# This state is executed if the intent was not understood
dont_get_the_intent:
  type: message_text_random
  properties:
    responses:
      - "Sorry I didn't get that! Please come again."
      - "I beg your pardon! Say that again?"
      - "Are we talking about weather? What would you like to know?"
      - "Sorry I know only about the weather"
      - "You can ask me about the weather, the rainfall, the
temperature, I don't know much more"
    delay: 0
    transitions:
      next_state: input_intent
  idk_what_you_talkin_about:
    type: message_text_random
    properties:

```

```

    responses:
      - "Sorry I didn't get that! Please come again."
      - "I beg your pardon! Say that again?"
      - "Are we talking about retail or weather? What would you like to
know?"
      - "Sorry I know only about retail and the weather"
      - "You can ask me about retail information or the weather, the
rainfall, the temperature. I don't know much more."
    delay: 0
    transitions:
      next_state: input_intent
  error:
    type: change_context
    properties:
      update_keys:
        intent: 'error'
    transitions:
      flow: error_flow
  retail_inventory_check:
    type: change_context
    properties:
      update_keys:
        intent: 'retail_inventory_check'
    transitions:
      flow: retail_flow
  retail_price_check:
    type: change_context
    properties:
      update_keys:
        intent: 'check_item_price'
    transitions:
      flow: retail_flow
  retail_store_location:
    type: change_context
    properties:
      update_keys:
        intent: 'find_the_store'
    transitions:
      flow: retail_flow
  weather:
    type: change_context
    properties:
      update_keys:
        intent: 'weather'
    transitions:
      flow: weather_flow

```

```
temperature:
  type: change_context
  properties:
    update_keys:
      intent: 'temperature'
  transitions:
    flow: weather_flow
rainfall:
  type: change_context
  properties:
    update_keys:
      intent: 'rainfall'
  transitions:
    flow: weather_flow
sunny:
  type: change_context
  properties:
    update_keys:
      intent: 'sunny'
  transitions:
    flow: weather_flow
cloudy:
  type: change_context
  properties:
    update_keys:
      intent: 'cloudy'
  transitions:
    flow: weather_flow
snow:
  type: change_context
  properties:
    update_keys:
      intent: 'snow'
  transitions:
    flow: weather_flow
rain:
  type: change_context
  properties:
    update_keys:
      intent: 'rain'
  transitions:
    flow: weather_flow
snowfall:
  type: change_context
  properties:
    update_keys:
```

```

        intent: 'snowfall'
    transitions:
        flow: weather_flow
tempyesno:
    type: change_context
    properties:
        update_keys:
            intent: 'tempyesno'
    transitions:
        flow: weather_flow
humidity:
    type: change_context
    properties:
        update_keys:
            intent: 'humidity'
    transitions:
        flow: weather_flow
end_state:
    type: reset
    transitions:
        next_state: init

```

## retail\_flow.yml

```

name: retail_flow
states:
    store_location:
        type: conditional_exists
        properties:
            key: '{{location}}'
        transitions:
            exists: retail_state
            notexists: ask_retail_location
    retail_state:
        type: Retail
        properties:
        transitions:
            next_state: output_retail
    output_retail:
        type: message_text
        properties:
            text: '{{retail_status}}'
        transitions:
            next_state: input_intent
    ask_retail_location:

```

```

    type: message_text
    properties:
      text: "For which location? I can find the closest store near you."
    transitions:
      next_state: input_retail_location
input_retail_location:
  type: input_user
  properties:
    nlp_type: jarvis
    entities:
      slot: location
      require_match: true
    transitions:
      match: retail_state
      notmatch: check_retail_jarvis_error
output_retail_acknowledge:
  type: message_text_random
  properties:
    responses:
      - 'ok in {{location}}'
      - 'the store in {{location}}'
      - 'I always wanted to shop in {{location}}'
    delay: 0
  transitions:
    next_state: retail_state
output_retail_notlocation:
  type: message_text
  properties:
    text: "I did not understand the location. Can you please repeat?"
  transitions:
    next_state: input_intent
check_retail_jarvis_error:
  type: conditional_exists
  properties:
    key: '{{jarvis_error}}'
  transitions:
    exists: show_retail_jarvis_api_error
    notexists: output_retail_notlocation
show_retail_jarvis_api_error:
  type: message_text
  properties:
    text: "I am having troubled understanding right now. Come again on that?"
  transitions:
    next_state: input_intent

```



```

name: weather_flow
states:
  check_weather_location:
    type: conditional_exists
    properties:
      key: '{{location}}'
    transitions:
      exists: weather_state
      notexists: ask_weather_location
  weather_state:
    type: Weather
    properties:
    transitions:
      next_state: output_weather
  output_weather:
    type: message_text
    properties:
      text: '{{weather_status}}'
    transitions:
      next_state: input_intent
  ask_weather_location:
    type: message_text
    properties:
      text: "For which location?"
    transitions:
      next_state: input_weather_location
  input_weather_location:
    type: input_user
    properties:
      nlp_type: jarvis
      entities:
        slot: location
        require_match: true
    transitions:
      match: weather_state
      notmatch: check_jarvis_error
  output_weather_acknowledge:
    type: message_text_random
    properties:
      responses:
        - 'ok in {{location}}'
        - 'the weather in {{location}}'
        - 'I always wanted to go in {{location}}'
    delay: 0

```

```

    transitions:
      next_state: weather_state
output_weather_notlocation:
  type: message_text
  properties:
    text: "I did not understand the location, can you please repeat?"
  transitions:
    next_state: input_intent
check_jarvis_error:
  type: conditional_exists
  properties:
    key: '{{jarvis_error}}'
  transitions:
    exists: show_jarvis_api_error
    notexists: output_weather_notlocation
show_jarvis_api_error:
  type: message_text
  properties:
    text: "I am having troubled understanding right now. Come again on
that, else check jarvis services?"
  transitions:
    next_state: input_intent

```

## error\_flow.yml

```

name: error_flow
states:
  error_state:
    type: message_text_random
    properties:
      responses:
        - "Sorry I didn't get that!"
        - "Are we talking about retail or weather? What would you like to
know?"
        - "Sorry I know only about retail information or the weather"
        - "You can ask me about retail information or the weather, the
rainfall, the temperature. I don't know much more"
        - "Let's talk about retail or the weather!"
      delay: 0
    transitions:
      next_state: input_intent

```

**Conéctese a API de terceros como motor de cumplimiento**

Conectamos las siguientes API de terceros como motor de cumplimiento de normativas

para responder a las preguntas:

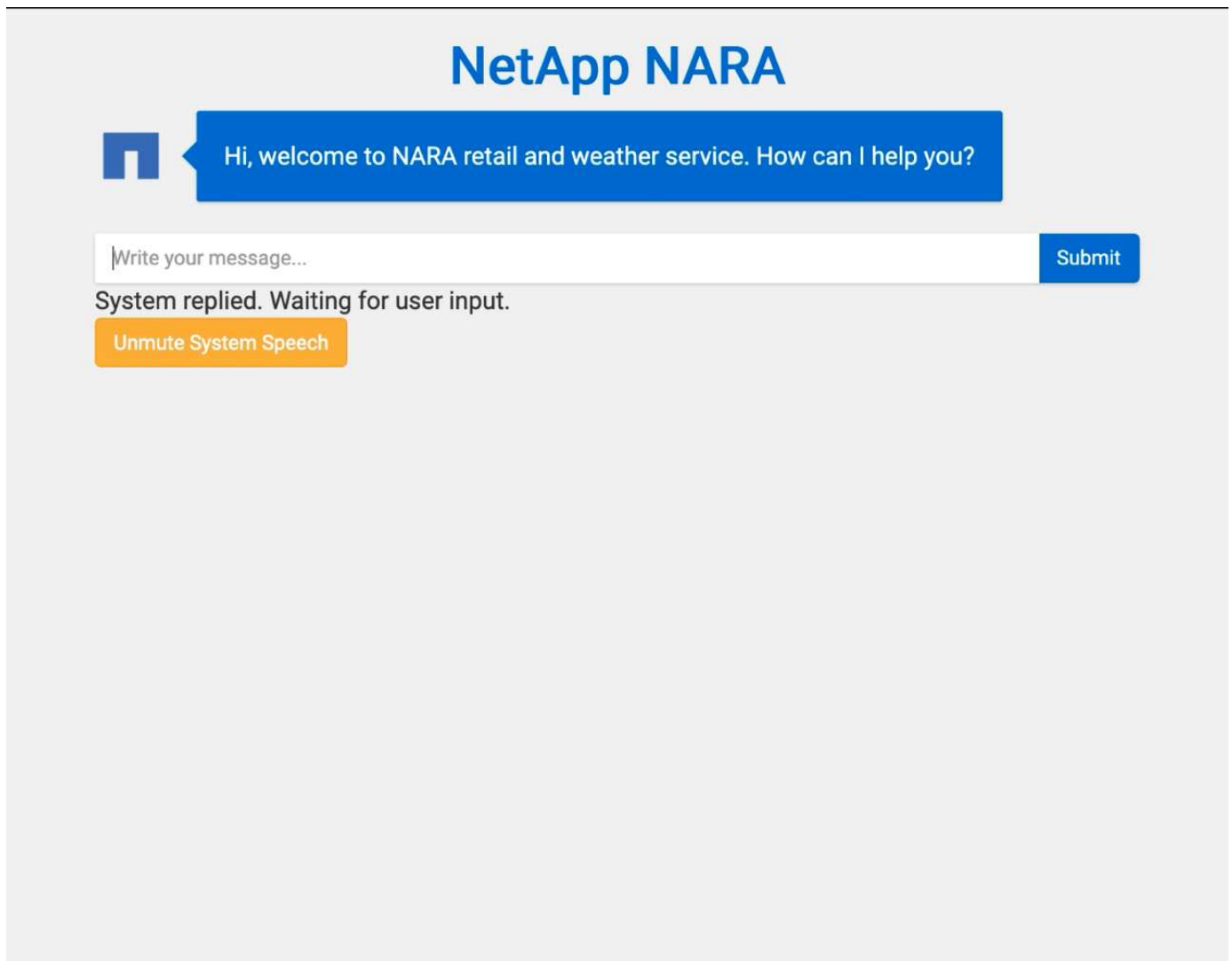
- "API de WeatherStack": regresa el clima, la temperatura, las precipitaciones y la nieve en un lugar determinado.
- "API de Fusion yelp": devuelve la información de la tienda más cercana en una ubicación determinada.
- "SDK de Python de eBay": devuelve el precio de un artículo determinado.

#### Demostración del asistente de venta al por menor de NetApp

Grabamos un vídeo de demostración del Asistente para minoristas de NetApp (NARA).

#### Demostración en vídeo de NARA

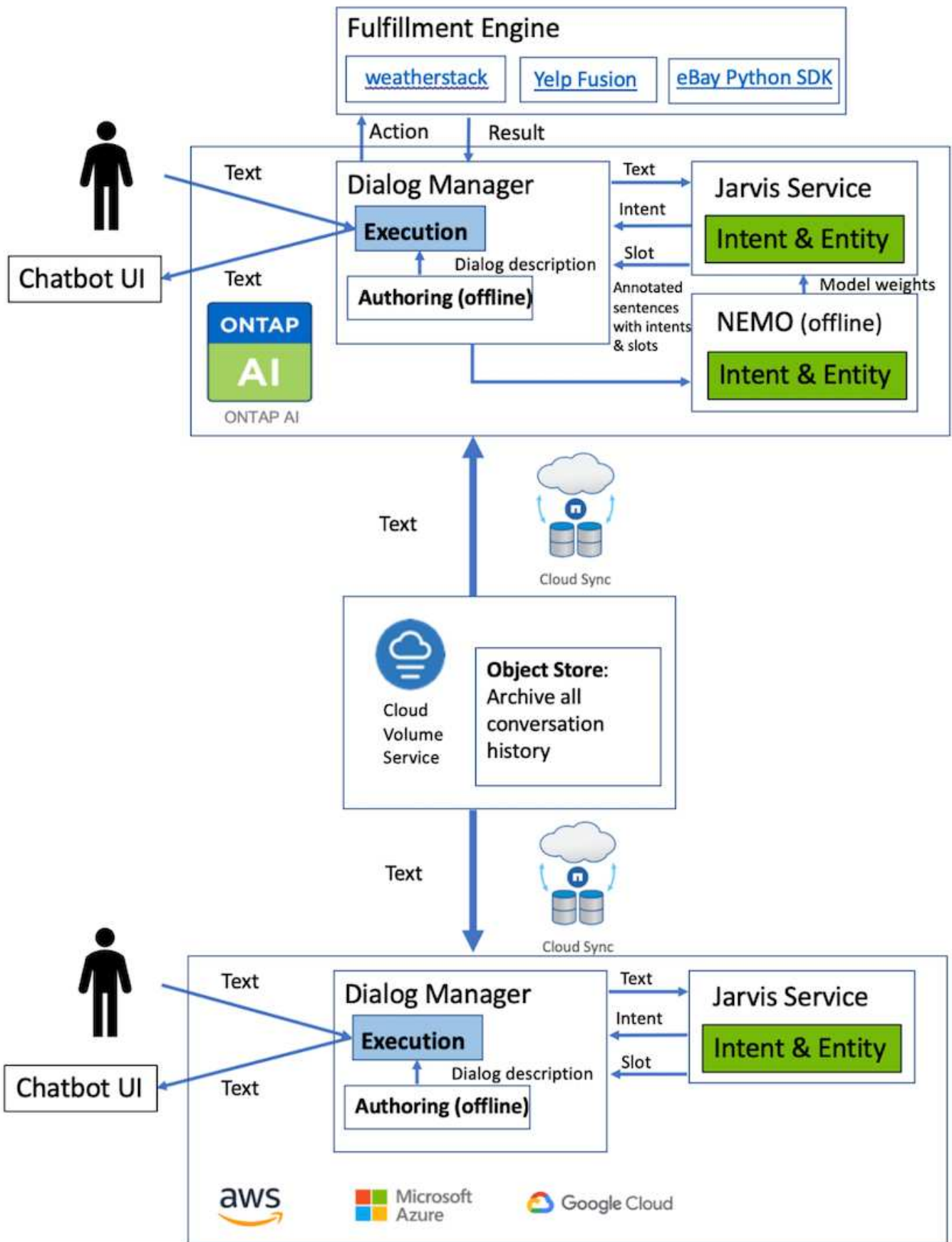
[Demostración en vídeo de NARA](#)



Utiliza la copia y sincronización de NetApp BlueXP para archivar el historial de conversaciones

Al volcar el historial de conversaciones en un archivo CSV una vez al día, podemos aprovechar la copia y sincronización de BlueXP para descargar los archivos de registro en un almacenamiento local. La siguiente figura muestra la arquitectura que indica que

Jarvis ha puesto en marcha on-premises y en nubes públicas, mientras utiliza la función Copy y Sync de BlueXP para enviar el historial de conversaciones para la formación de Nemo. En la sección encontrará más información sobre la formación de Nemo "[Expandir modelos de intención utilizando Nemo Training](#)".



NVIDIA Nemo es un kit de herramientas creado por NVIDIA para crear aplicaciones de IA conversacionales. Este kit de herramientas incluye colecciones de módulos preentrenados para ASR, NLP y TTS, lo que permite a investigadores y científicos de datos componer fácilmente arquitecturas complejas de redes neuronales y centrarse más en el diseño de sus propias aplicaciones.

Como se muestra en el ejemplo anterior, NARA sólo puede manejar un tipo limitado de preguntas. Esto se debe a que el modelo NLP pre-entrenado sólo entrena en este tipo de preguntas. Si queremos permitir QUE NARA pueda gestionar una gama más amplia de preguntas, debemos volver a formar este con nuestros propios conjuntos de datos. Por lo tanto, aquí mostramos cómo podemos utilizar Nemo para ampliar el modelo NLP para satisfacer los requisitos. Comenzamos convirtiendo el registro recolectado de NARA en el formato de Nemo, y luego entrenamos con el conjunto de datos para mejorar el modelo NLP.

### Modelo

Nuestro objetivo es permitir A NARA ordenar los elementos según las preferencias del usuario. Por ejemplo, podemos pedir A NARA que sugiera el restaurante de sushi con mejor calificación o que quiera QUE NARA busque los vaqueros con el precio más bajo. Para ello, utilizamos el modelo de detección de intención y relleno de ranuras proporcionado en Nemo como modelo de entrenamiento. Este modelo permite A NARA comprender la intención de buscar preferencias.

### Preparación de datos

Para entrenar el modelo, recopilamos el conjunto de datos para este tipo de preguntas y lo convertimos al formato Nemo. Aquí enumeramos los archivos que utilizamos para entrenar el modelo.

#### dict.intents.csv

Este archivo enumera todos los intentos que queremos que el Nemo entienda. Aquí tenemos dos intentos principales y una intención que sólo se utiliza para categorizar las preguntas que no encajan en ninguno de los intentos primarios.

```
price_check
find_the_store
unknown
```

#### dict.slots.csv

En este archivo se enumeran todas las ranuras que podemos etiquetar en nuestras preguntas de formación.

```
B-store.type
B-store.name
B-store.status
B-store.hour.start
B-store.hour.end
B-store.hour.day
B-item.type
```

B-item.name  
B-item.color  
B-item.size  
B-item.quantity  
B-location  
B-cost.high  
B-cost.average  
B-cost.low  
B-time.period\_of\_time  
B-rating.high  
B-rating.average  
B-rating.low  
B-interrogative.location  
B-interrogative.manner  
B-interrogative.time  
B-interrogative.personal  
B-interrogative  
B-verb  
B-article  
I-store.type  
I-store.name  
I-store.status  
I-store.hour.start  
I-store.hour.end  
I-store.hour.day  
I-item.type  
I-item.name  
I-item.color  
I-item.size  
I-item.quantity  
I-location  
I-cost.high  
I-cost.average  
I-cost.low  
I-time.period\_of\_time  
I-rating.high  
I-rating.average  
I-rating.low  
I-interrogative.location  
I-interrogative.manner  
I-interrogative.time  
I-interrogative.personal  
I-interrogative  
I-verb  
I-article  
O

## train.tsv

Este es el conjunto de datos de entrenamiento principal. Cada línea comienza con la pregunta que sigue a la lista de la categoría de intención en el archivo dict.intent.csv. La etiqueta se enumera a partir de cero.

## train\_slots.tsv

```
20 46 24 25 6 32 6
52 52 24 6
23 52 14 40 52 25 6 32 6
...
```

## Entrenar el modelo

```
docker pull nvcr.io/nvidia/nemo:v0.10
```

A continuación, utilizamos el siguiente comando para iniciar el contenedor. En este comando, limitamos el contenedor para usar una única GPU (ID de GPU = 1), ya que se trata de un ejercicio de entrenamiento de poco peso. También mapeamos nuestro espacio de trabajo local /Workspace/nemo/ a la carpeta dentro del contenedor /nemo.

```
NV_GPU='1' docker run --runtime=nvidia -it --shm-size=16g \
    --network=host --ulimit memlock=-1 --ulimit
stack=67108864 \
    -v /workspace/nemo:/nemo\
    --rm nvcr.io/nvidia/nemo:v0.10
```

Dentro del contenedor, si queremos empezar desde el modelo ORIGINAL BERT pre-entrenado, podemos usar el siguiente comando para iniciar el procedimiento de entrenamiento. data\_dir es el argumento para establecer la ruta de los datos de entrenamiento. dir\_trabajo le permite configurar dónde desea almacenar los archivos de punto de control.

```
cd examples/nlp/intent_detection_slot_tagging/
python joint_intent_slot_with_bert.py \
    --data_dir /nemo/training_data\
    --work_dir /nemo/log
```

Si contamos con nuevos conjuntos de datos de entrenamiento y queremos mejorar el modelo anterior, podemos utilizar el siguiente comando para continuar desde el punto que hemos detenido. checkpoint\_dir lleva la ruta a la carpeta de puntos de control anteriores.



```
cd examples/nlp/intent_detection_slot_tagging/
python joint_intent_slot_infer.py \
    --data_dir /nemo/training_data \
    --checkpoint_dir /nemo/log/2020-05-04_18-34-20/checkpoints/ \
    --eval_file_prefix test
```

## Inferencia del modelo

Se debe validar el rendimiento del modelo entrenado después de una serie determinada de épocas. El siguiente comando nos permite probar la consulta una por una. Por ejemplo, en este comando, queremos comprobar si nuestro modelo puede identificar adecuadamente la intención de la consulta `where can I get the best pasta`.

```
cd examples/nlp/intent_detection_slot_tagging/
python joint_intent_slot_infer_b1.py \
    --checkpoint_dir /nemo/log/2020-05-29_23-50-58/checkpoints/ \
    --query "where can i get the best pasta" \
    --data_dir /nemo/training_data/ \
    --num_epochs=50
```

A continuación, se muestra la salida de la inferencia. En el resultado, podemos ver que nuestro modelo entrenado puede predecir correctamente la intención `find_the_store`, y devolver las palabras clave en las que estamos interesados. Con estas palabras clave, permitimos A LA NARA buscar lo que los usuarios desean y realizar una búsqueda más precisa.

```
[NeMo I 2020-05-30 00:06:54 actions:728] Evaluating batch 0 out of 1
[NeMo I 2020-05-30 00:06:55 inference_utils:34] Query: where can i get the
best pasta
[NeMo I 2020-05-30 00:06:55 inference_utils:36] Predicted intent:      1
find_the_store
[NeMo I 2020-05-30 00:06:55 inference_utils:50] where      B-
interrogative.location
[NeMo I 2020-05-30 00:06:55 inference_utils:50] can        O
[NeMo I 2020-05-30 00:06:55 inference_utils:50] i          O
[NeMo I 2020-05-30 00:06:55 inference_utils:50] get        B-verb
[NeMo I 2020-05-30 00:06:55 inference_utils:50] the        B-article
[NeMo I 2020-05-30 00:06:55 inference_utils:50] best       B-rating.high
[NeMo I 2020-05-30 00:06:55 inference_utils:50] pasta     B-item.type
```

## Conclusión

Un auténtico sistema de IA conversacional se involucra en un diálogo similar al humano, comprende el contexto y proporciona respuestas inteligentes. Estos modelos de IA a menudo son enormes y muy complejos. Con las GPU de NVIDIA y el almacenamiento de

NetApp, se pueden formar y optimizar modelos masivos de lenguaje de última generación para ejecutar inferencia rápidamente. Se trata de una importante medida para poner fin a la compensación entre un modelo de IA que es rápido frente a uno que es grande y complejo. Los modelos de comprensión del lenguaje optimizados para GPU se pueden integrar en las aplicaciones de IA en sectores como el sanitario, el minorista y los servicios financieros. Esto hace posible que asistentes de voz digitales avanzados en altavoces inteligentes y líneas de servicio al cliente. Estos sistemas de IA conversacionales de alta calidad permiten a las empresas de sectores verticales proporcionar servicios personalizados que antes eran impensables al tratar con los clientes.

Jarvis permite el despliegue de casos de uso como asistentes virtuales, avatares digitales, sensor multimodal Fusion (CV con fusible ASR/NLP/TTS) o cualquier caso de uso independiente ASR/NLP/TTS/CV, como la transcripción. Hemos creado un asistente de venta al por menor virtual que puede responder a preguntas relacionadas con el tiempo, los puntos de interés y los precios del inventario. También demostramos cómo mejorar las funcionalidades de comprensión del lenguaje natural del sistema de IA conversacional archivando el historial de conversaciones con BlueXP Copy y Sync y entrenando modelos de Nemo en datos nuevos.

## Reconocimientos

Los autores agradecen las contribuciones que nuestros estimados colegas de NVIDIA han realizado a este documento técnico: Davide Onofrío, Alex Qi, Sicong Ji, Marty Jain y Robert Sohigian. Los autores también desean reconocer las contribuciones de miembros clave del equipo de NetApp: Santosh Rao, David Arnette, Michael Oglesby, Brent Davis, Andy Sayare, Erik Mulder, y Mike McNamara.

Nuestro más sincero agradecimiento a todos aquellos que nos ofrecieron sus valiosos comentarios y su experiencia para la creación de este informe.

## Dónde encontrar información adicional

Si quiere obtener más información sobre el contenido de este documento, consulte los siguientes recursos:

- NVIDIA DGX Station, GPU V100, GPU Cloud
  - Estación DGX de NVIDIA <https://www.nvidia.com/en-us/data-center/dgx-station/>["https://www.nvidia.com/en-us/data-center/dgx-station/"]
  - GPU de núcleo tensor NVIDIA V100 <https://www.nvidia.com/en-us/data-center/tesla-v100/>["https://www.nvidia.com/en-us/data-center/tesla-v100/"]
  - NVIDIA NGC <https://www.nvidia.com/en-us/gpu-cloud/>["https://www.nvidia.com/en-us/gpu-cloud/"]
- NVIDIA Jarvis Multimodal Framework
  - NVIDIA Jarvis <https://developer.nvidia.com/nvidia-jarvis>["https://developer.nvidia.com/nvidia-jarvis"]
  - Acceso temprano a NVIDIA Jarvis <https://developer.nvidia.com/nvidia-jarvis-early-access>["https://developer.nvidia.com/nvidia-jarvis-early-access"]
- NVIDIA Nemo
  - NVIDIA Nemo <https://developer.nvidia.com/nvidia-nemo>["https://developer.nvidia.com/nvidia-nemo"]

- Guía para desarrolladores <https://nvidia.github.io/NeMo/>["https://nvidia.github.io/NeMo/"]
- Sistemas AFF de NetApp
  - Especificaciones técnicas de AFF a-Series de NetApp <https://www.netapp.com/us/media/ds-3582.pdf>["https://www.netapp.com/us/media/ds-3582.pdf"]
  - Ventaja de NetApp Flash para All Flash FAS <https://www.netapp.com/us/media/ds-3733.pdf>["https://www.netapp.com/us/media/ds-3733.pdf"]
  - Biblioteca de información de ONTAP  
9 <http://mysupport.netapp.com/documentation/productlibrary/index.html?productID=62286>["http://mysupport.netapp.com/documentation/productlibrary/index.html?productID=62286"]
  - Informe técnico de NetApp ONTAP FlexGroup Volumes <https://www.netapp.com/us/media/tr-4557.pdf>["https://www.netapp.com/us/media/tr-4557.pdf"]
- ONTAP AI de NetApp
  - Guía de diseño de ONTAP AI con DGX-1 y Cisco Networking <https://www.netapp.com/us/media/nva-1121-design.pdf>["https://www.netapp.com/us/media/nva-1121-design.pdf"]
  - Guía de puesta en marcha de ONTAP AI con DGX-1 y Cisco Networking <https://www.netapp.com/us/media/nva-1121-deploy.pdf>["https://www.netapp.com/us/media/nva-1121-deploy.pdf"]
  - Guía de diseño de ONTAP AI con DGX-1 y Mellanox Networking <http://www.netapp.com/us/media/nva-1138-design.pdf>["http://www.netapp.com/us/media/nva-1138-design.pdf"]
  - Guía de diseño de ONTAP AI con DGX-2 <https://www.netapp.com/us/media/nva-1135-design.pdf>["https://www.netapp.com/us/media/nva-1135-design.pdf"]

## TR-4858: Solución de orquestación de NetApp con Run:AI

Rick Huang, David Arnette, Sung-han Lin, NetApp Yaron Goldberg, Run:AI

Los sistemas de almacenamiento AFF de NetApp proporcionan un rendimiento extremo y funcionalidades de gestión de datos de cloud híbrido líderes en el sector. NetApp y Run:AI se han asociado para demostrar las funcionalidades únicas de la solución de IA ONTAP de NetApp para cargas de trabajo de inteligencia artificial (IA) y aprendizaje automático (ML) que ofrecen rendimiento, fiabilidad y soporte de clase empresarial. Ejecutar la orquestación de la IA de las cargas de trabajo de IA añade una programación basada en Kubernetes y una plataforma de uso de recursos para ayudar a los investigadores a gestionar y optimizar la utilización de la GPU. Junto con los sistemas DGX de NVIDIA, la solución combinada de NetApp, NVIDIA y Run:AI proporciona una pila de infraestructura diseñada específicamente para cargas de trabajo de IA empresariales. Este informe técnico proporciona instrucciones direccionales a los clientes que crean sistemas de IA conversacionales que respaldan diversos casos prácticos y mercados verticales del sector. Incluye información sobre la puesta en marcha de Run:AI y un sistema de almacenamiento AFF A800 de NetApp y sirve como arquitectura de referencia para obtener el método más sencillo de poner en marcha con éxito iniciativas de IA de forma rápida.

El público objetivo de la solución incluye los siguientes grupos:

- Los arquitectos empresariales que diseñan soluciones para el desarrollo de modelos de IA y software para

casos prácticos basados en Kubernetes, como microservicios en contenedores

- Los científicos de datos que buscan formas eficientes de alcanzar objetivos de desarrollo de modelos eficientes en un entorno de clúster con varios equipos y proyectos
- Ingenieros de datos a cargo del mantenimiento y la ejecución de modelos de producción
- Responsables de la toma de decisiones de TECNOLOGÍA y ejecutivos y líderes empresariales que desean crear la experiencia óptima de utilización de recursos de clúster de Kubernetes y lograr el plazo de comercialización más rápido posible gracias a las iniciativas de IA

## Descripción general de la solución

### Plano de control de la IA y la IA con ONTAP de NetApp

La arquitectura de IA de ONTAP de NetApp, desarrollada y verificada por NetApp y NVIDIA, cuenta con la tecnología de sistemas NVIDIA DGX y los sistemas de almacenamiento conectados al cloud de NetApp. Esta arquitectura de referencia proporciona a las organizaciones DE TI las siguientes ventajas:

- Elimina las complejidades de diseño
- Permite un escalado independiente de las capacidades de computación y almacenamiento
- Permite a los clientes empezar con poco y escalar sin problemas
- Ofrece opciones de almacenamiento para distintos niveles de rendimiento y coste

ONTAP AI de NetApp integra perfectamente los sistemas DGX y los sistemas de almacenamiento AFF A800 de NetApp con una red de vanguardia. Los sistemas ONTAP AI y DGX de NetApp simplifican las puestas en marcha de IA al eliminar complejidades y conjeturas en la fase de diseño. Los clientes pueden empezar poco a poco y aumentar sus sistemas de forma ininterrumpida, a la vez que gestionan de forma inteligente datos entre el perímetro, el núcleo y el cloud.

El plano de control de IA de NetApp es una solución de gestión de experimentos y datos de IA, ML y aprendizaje profundo (DL) para científicos de datos e ingenieros de datos. A medida que las organizaciones aumentan el uso de la IA, deben hacer frente a numerosos retos, como la escalabilidad de la carga de trabajo y la disponibilidad de los datos. El plano de control de IA de NetApp responde a estos retos mediante funcionalidades como clonar rápidamente un espacio de nombres de datos del mismo modo que una repo de Git y definir e implementar flujos de trabajo de entrenamiento de IA que incorporen la creación casi instantánea de bases de datos y modelos para la trazabilidad y versionado. Con el plano de control de IA de NetApp, puede replicar datos sin problemas entre sitios y regiones y aprovisionar rápidamente espacios de trabajo de los portátiles Jupyter con acceso a conjuntos de datos masivos.

### Ejecutar:Plataforma de IA para orquestación de cargas de trabajo de IA

Ejecución:la IA ha creado la primera plataforma de orquestación y virtualización del mundo para la infraestructura de IA. Al abstraer las cargas de trabajo del hardware subyacente, Run:AI crea un pool compartido de recursos de GPU que se pueden aprovisionar de forma dinámica, lo que permite una orquestación eficiente de las cargas de trabajo de IA y un uso optimizado de las GPU. Los científicos de datos pueden consumir sin problemas grandes cantidades de potencia de GPU para mejorar y acelerar sus investigaciones mientras los equipos DE TECNOLOGÍA conservan el control centralizado entre sitios y la visibilidad en tiempo real sobre el aprovisionamiento de recursos, la cola y el uso. La plataforma Run:AI se ha creado sobre Kubernetes, por lo que permite una integración sencilla con los flujos de trabajo existentes de tecnología y ciencia de datos.

La plataforma Run:AI ofrece los siguientes beneficios:

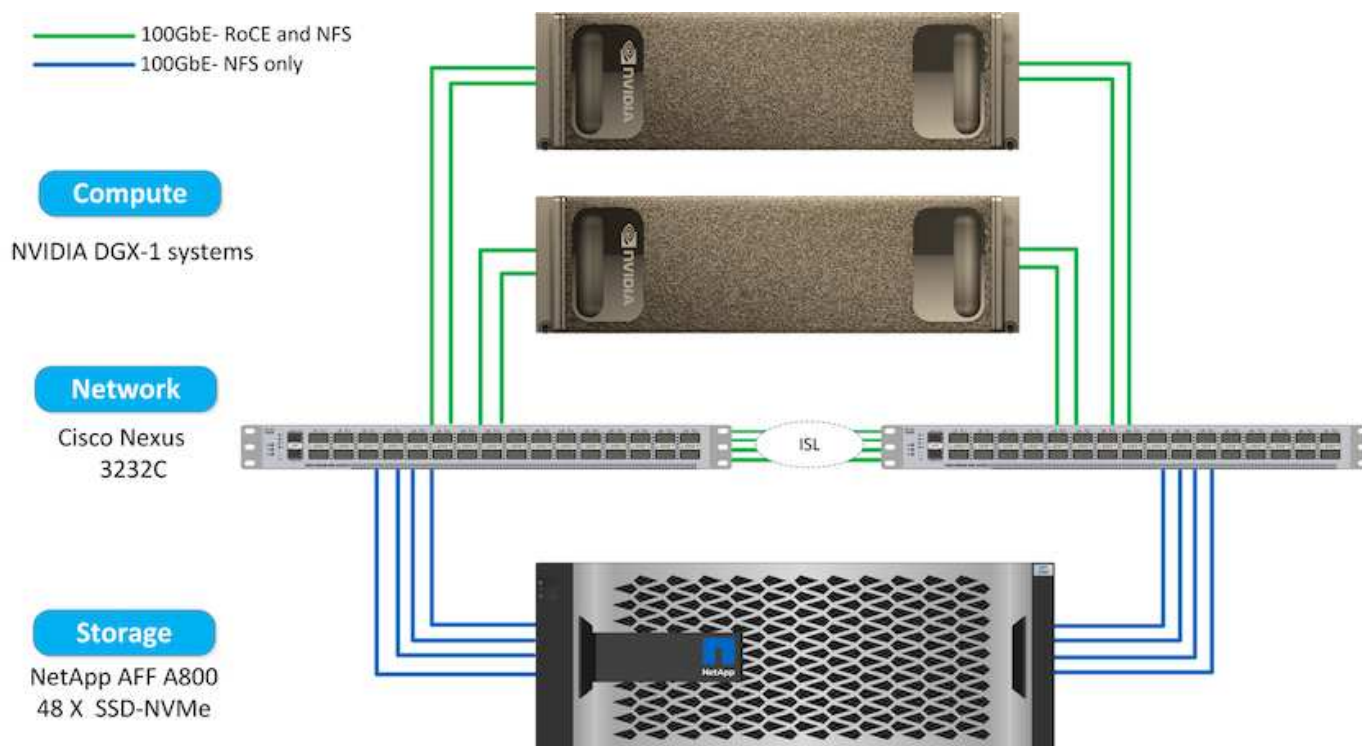
- **Plazo de innovación más rápido.** con los mecanismos de puesta en cola, priorización y agrupación de

recursos de IA junto con un sistema de almacenamiento de NetApp, los investigadores se eliminan de los problemas de gestión de la infraestructura y pueden centrarse exclusivamente en la ciencia de datos. Ejecución: los clientes de NetApp y IA aumentan la productividad ejecutando tantas cargas de trabajo como necesiten sin cuellos de botella en la canalización de datos o en la computación.

- **Aumento de la productividad del equipo.** Ejecutar: los algoritmos de justicia de IA garantizan que todos los usuarios y equipos obtengan su parte justa de los recursos. Es posible predefinir políticas sobre proyectos de prioridad y la plataforma permite la asignación dinámica de recursos de un usuario o equipo a otro, lo que ayuda a los usuarios a obtener un acceso puntual a los recursos de la GPU codiciados.
- **Uso mejorado de la GPU.** el Planificador Run:AI permite a los usuarios utilizar fácilmente GPU fraccionarias, GPU enteros y varios nodos de GPU para el entrenamiento distribuido en Kubernetes. De esta forma, las cargas de trabajo de IA se ejecutan según sus necesidades, no en función de la capacidad. Los equipos de ciencia de datos pueden realizar más experimentos de IA en la misma infraestructura.

## Tecnología de soluciones

Esta solución se implementó con un sistema AFF A800 de NetApp, dos servidores DGX-1 y dos switches Cisco Nexus 3232C 100 GbE. Cada servidor DGX-1 está conectado a los switches Nexus mediante cuatro conexiones 100 GbE que se utilizan para las comunicaciones entre GPU. Para ello se utilizan accesos remotos directos a la memoria (RDMA) sobre Ethernet convergente (roce). En estos enlaces también se producen las comunicaciones IP tradicionales para el acceso al almacenamiento NFS. Cada controladora de almacenamiento está conectada a los switches de red mediante cuatro enlaces 100 GbE. En la siguiente figura se muestra la arquitectura de la solución de IA de ONTAP utilizada en este informe técnico para todos los escenarios de pruebas.



### Hardware utilizado en esta solución

Esta solución se validó con la arquitectura de referencia ONTAP AI dos nodos DGX-1 y un sistema de almacenamiento AFF A800. Consulte ["NVA-1121"](#) para obtener más detalles sobre la infraestructura utilizada en esta validación.

En la siguiente tabla se enumeran los componentes de hardware necesarios para implementar la solución según se ha probado.

Hardware subyacente	Cantidad
Sistemas DGX-1	2
A800 de AFF	1
Switches Nexus 3232C	2

### Requisitos de software

Esta solución se validó con una puesta en marcha de Kubernetes básica con el operador Run:AI instalado. Kubernetes se puso en marcha usando el ["DeepOps de NVIDIA"](#) motor de puesta en marcha, que implementa todos los componentes necesarios para un entorno listo para la producción. DeepOps ya se implementa automáticamente ["Trident de NetApp"](#) Para la integración del almacenamiento persistente con el entorno k8S, se crearon las clases de almacenamiento predeterminadas, de modo que los contenedores aprovechan el almacenamiento del sistema de almacenamiento A800 de AFF. Para obtener más información sobre Trident con Kubernetes en ONTAP AI, consulte ["TR-4798"](#).

En la siguiente tabla se enumeran los componentes de software necesarios para implementar la solución según se ha probado.

De NetApp	Versión u otra información
Software de gestión de datos ONTAP de NetApp	9.6p4
Firmware de switch Cisco NX-OS	7.0(3)I6(1)
SO DGX DE NVIDIA	4.0.4 - Ubuntu 18.04 LTS
La versión de Kubernetes	1.17
Versión de Trident	20.04.0
Ejecución: CLI de IA	v2.1.13
Ejecutar:versión del operador de Kubernetes de orquestación de IA	1.0.39
Plataforma contenedora Docker	18.06.1-ce [e68fc7a]

Requisitos de software adicionales para la ejecución: Se puede encontrar IA en ["Ejecutar:requisitos previos del clúster de GPU de IA"](#).

### Óptimo uso de clúster y de la GPU con Run:AI

En las siguientes secciones se ofrecen detalles sobre la instalación de Run:AI, escenarios de prueba y resultados realizados en esta validación.

Validamos el funcionamiento y el rendimiento de este sistema mediante herramientas estándar del sector,

incluidas las pruebas de rendimiento TensorFlow. El conjunto de datos ImageNET se utilizó para entrenar ResNet-50, que es un famoso modelo de red neuronal convolucional (CNN) DL para la clasificación de imágenes. RESNET-50 ofrece un resultado de entrenamiento preciso con un tiempo de procesamiento más rápido, lo que nos permitió generar una demanda suficiente sobre el almacenamiento.

#### Ejecutar:instalación de AI

Para instalar Run:AI, lleve a cabo los siguientes pasos:

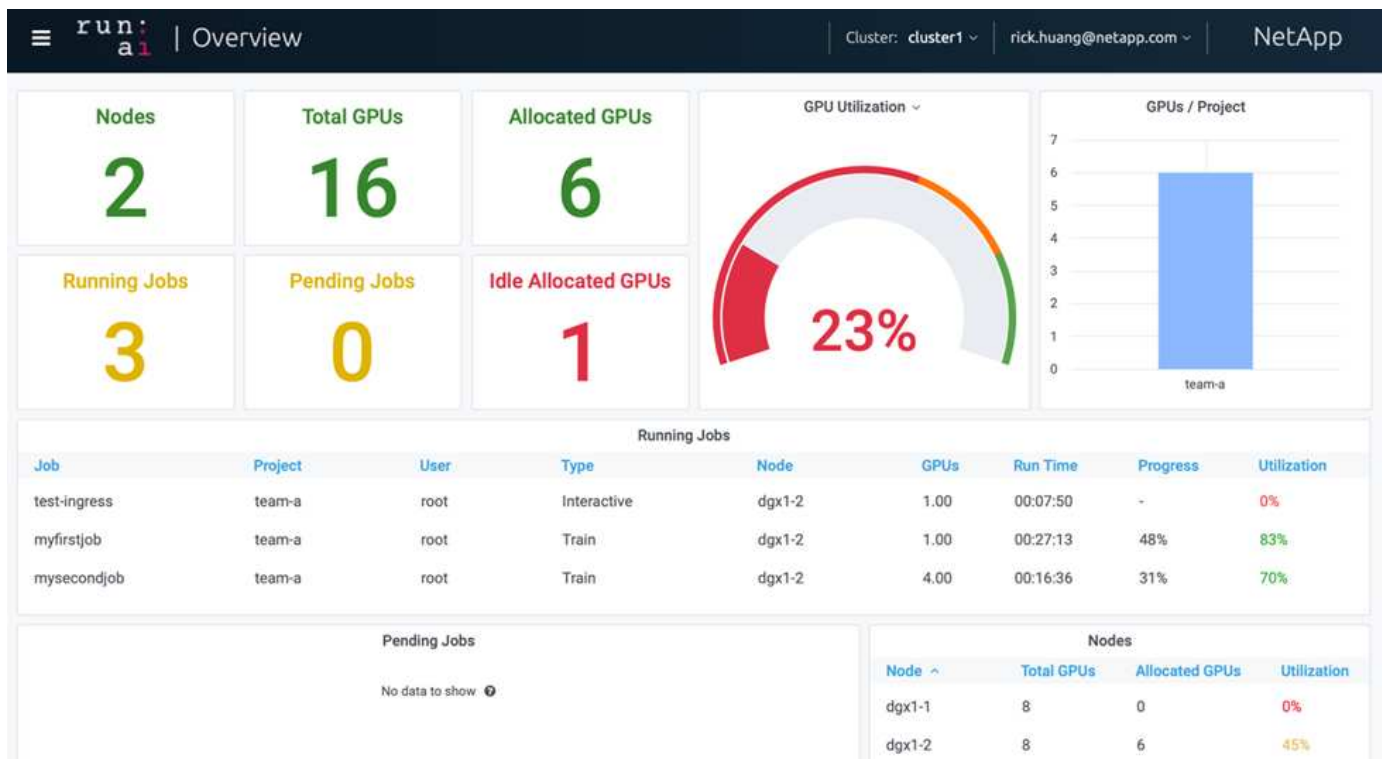
1. Instale el clúster de Kubernetes con DeepOps y configure la clase de almacenamiento predeterminada de NetApp.
2. Prepare los nodos de GPU:
  - a. Compruebe que los controladores de NVIDIA están instalados en los nodos de la GPU.
  - b. Compruebe que `nvidia-docker` se instala y configura como el tiempo de ejecución de docker predeterminado.
3. Ejecución de instalación:AI:
  - a. Inicie sesión en el ["Ejecución: IU de administración de IA"](#) para crear el clúster.
  - b. Descargue el creado `runai-operator-<clustername>.yaml` archivo.
  - c. Aplique la configuración del operador al clúster de Kubernetes.

```
kubectl apply -f runai-operator-<clustername>.yaml
```

4. Compruebe la instalación:
  - a. Vaya a ["https://app.run.ai/"](https://app.run.ai/).
  - b. Vaya a la consola Overview.
  - c. Compruebe que el número de GPU de la parte superior derecha refleja el número esperado de GPU y los nodos de GPU forman parte de la lista de servidores. Para obtener más información acerca de la puesta en marcha de Run:IA, consulte ["Instalar Run:IA en un clúster de Kubernetes en las instalaciones"](#) y.. ["Instalación de la CLI Run:AI"](#).

#### Ejecutar:Paneles de IA y vistas

Después de instalar Run:AI en su clúster de Kubernetes y configurar los contenedores correctamente, verá las siguientes consolas y vistas en ["https://app.run.ai/"](https://app.run.ai/) en su navegador, como se muestra en la siguiente figura.



Hay 16 GPU en total en el clúster proporcionados por dos nodos DGX-1. Puede ver el número de nodos, el total de GPU disponibles, las GPU asignadas con cargas de trabajo, el número total de trabajos en ejecución, los trabajos pendientes y las GPU asignadas inactivas. En el lado derecho, el diagrama de barras muestra las GPU por proyecto, que resume cómo usan los distintos equipos el recurso de clúster. En el medio se encuentra la lista de trabajos actualmente en ejecución con detalles de trabajo, incluidos el nombre del trabajo, el proyecto, el usuario, el tipo de trabajo, El nodo en el que se ejecuta cada trabajo, el número de GPU asignados para ese trabajo, el tiempo de ejecución actual del trabajo, el progreso del trabajo en porcentaje y el uso de la GPU para ese trabajo. Tenga en cuenta que el clúster está infrautilizado (uso de la GPU al 23%) porque solo hay tres trabajos en ejecución enviados por un único equipo (team-a).

En la siguiente sección, mostramos cómo crear varios equipos en la pestaña proyectos y asignar GPU para cada equipo con el fin de maximizar el uso del clúster y gestionar los recursos cuando hay muchos usuarios por clúster. Los escenarios de prueba imitan entornos empresariales en los que los recursos de memoria y GPU se comparten entre cargas de trabajo de entrenamiento, inferencia e interactivas.

### Creación de proyectos para equipos de ciencia de datos y asignación de GPU

Los investigadores pueden enviar cargas de trabajo a través de la CLI Run:AI, Kubeflow o procesos similares. Para agilizar la asignación de recursos y crear prioridades, Run:AI introduce el concepto de proyectos. Los proyectos son entidades de cuota que asocian un nombre de proyecto con la asignación y las preferencias de la GPU. Se trata de una forma sencilla y cómoda de gestionar varios equipos de ciencia de datos.

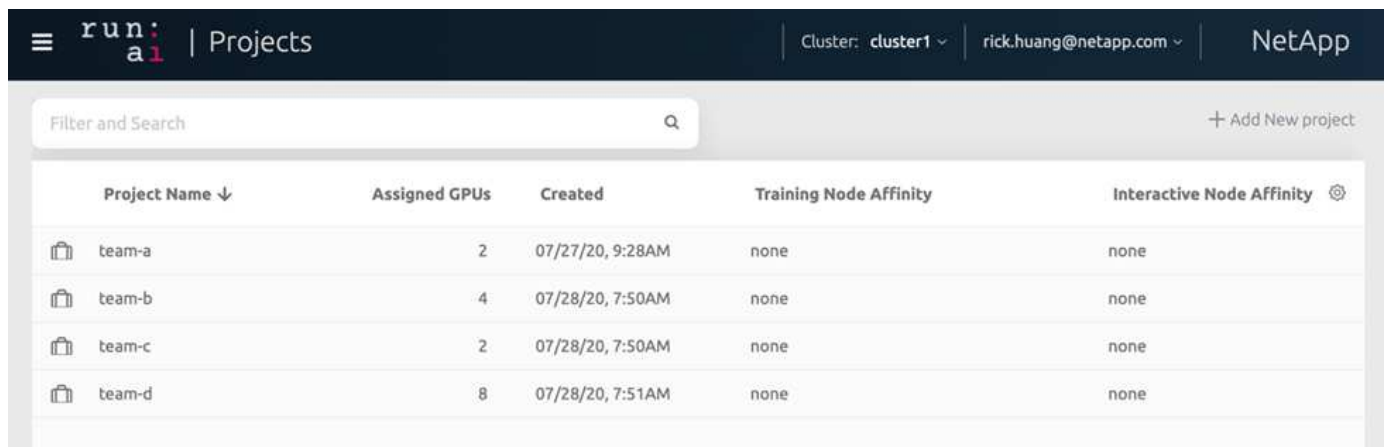
Un investigador que presenta una carga de trabajo debe asociar un proyecto con una solicitud de carga de trabajo. El programador Run:AI compara la solicitud con las asignaciones actuales y el proyecto, y determina si la carga de trabajo puede asignarse recursos o si debe permanecer en estado pendiente.

Como administrador del sistema, puede establecer los siguientes parámetros en la ficha Ejecutar:proyectos AI:



- **Proyectos modelo.** establecer un proyecto por usuario, establecer un proyecto por equipo de usuarios y establecer un proyecto por proyecto de organización real.
- \* Cuotas del proyecto.\* cada proyecto está asociado con una cuota de GPU que se pueden asignar para este proyecto al mismo tiempo. Se trata de una cuota garantizada en el sentido de que se garantiza que los investigadores que utilizan este proyecto obtengan este número de GPU, independientemente del estado del clúster. Por lo general, la suma de la asignación de proyectos debe ser igual al número de GPU del clúster. Más allá de eso, un usuario de este proyecto puede recibir un exceso de cuota. Mientras no se utilicen las GPU, un investigador que usa este proyecto puede obtener más GPU. Demostramos escenarios de prueba de exceso de cuota y consideraciones de equidad en ["Utilización óptima del clúster gracias a la asignación de GPU por encima de una cuota"](#), ["Equidad en la asignación de recursos básicos"](#), y ["Justicia por exceso de cuotas"](#).
- Cree un proyecto nuevo, actualice un proyecto existente y elimine un proyecto existente.
- **Limite los trabajos para que se ejecuten en grupos de nodos específicos.** Puede asignar proyectos específicos para que se ejecuten solo en nodos específicos. Esto resulta útil cuando el equipo de proyecto necesita hardware especializado, por ejemplo, con suficiente memoria. Como alternativa, un equipo de proyecto podría ser el propietario de hardware específico que se adquirió con un presupuesto especializado o cuando sea necesario dirigir las cargas de trabajo de creación o interactivas para trabajar en hardware más débil y dirigir la formación de mayor duración o las cargas de trabajo sin supervisión a nodos más rápidos. Para ver los comandos para agrupar nodos y establecer la afinidad de un proyecto específico, consulte ["Ejecute:Documentación de IA"](#).
- **Limitar la duración de los trabajos interactivos.** Los investigadores se olvidan frecuentemente de cerrar trabajos interactivos. Esto podría conducir a una pérdida de recursos. Algunas organizaciones prefieren limitar la duración de los trabajos interactivos y cerrarlos automáticamente.

La siguiente figura muestra la vista proyectos con cuatro equipos creados. A cada equipo se le asigna un número diferente de GPU para dar cuenta de diferentes cargas de trabajo, con un número total de GPU igual al de las GPU disponibles totales en un clúster que consta de dos DGX-1s.



Project Name ↓	Assigned GPUs	Created	Training Node Affinity	Interactive Node Affinity ⓘ
team-a	2	07/27/20, 9:28AM	none	none
team-b	4	07/28/20, 7:50AM	none	none
team-c	2	07/28/20, 7:50AM	none	none
team-d	8	07/28/20, 7:51AM	none	none

### Enviar trabajos en ejecución:AI CLI

En esta sección se ofrecen detalles sobre los comandos básicos Run:AI que puede utilizar para ejecutar cualquier trabajo con Kubernetes. Se divide en tres partes según el tipo de carga de trabajo. Las cargas de trabajo de IA/ML/AP se pueden dividir en dos tipos genéricos:

- **Sesiones de formación desatendida.** Con estos tipos de cargas de trabajo, el científico de datos prepara una carga de trabajo de ejecución automática y la envía para su ejecución. Durante la ejecución, el cliente puede examinar los resultados. Este tipo de carga de trabajo se utiliza a menudo en la producción o

cuando el desarrollo de modelos se encuentra en una etapa en la que no se requiere intervención humana.

- **Sesiones interactivas de construcción.** Con este tipo de cargas de trabajo, el científico de datos abre una sesión interactiva con Bash, Jupyter Notebook, PyCharm remoto o IDE similares y accede directamente a los recursos de la GPU. Incluimos un tercer escenario para ejecutar cargas de trabajo interactivas con puertos conectados a fin de revelar un puerto interno al usuario del contenedor.

## Cargas de trabajo de formación desatendida

Después de configurar los proyectos y asignar GPU, puede ejecutar cualquier carga de trabajo de Kubernetes usando el siguiente comando en la línea de comandos:

```
$ runai project set team-a runai submit hyper1 -i gcr.io/run-ai-demo/quickstart -g 1
```

Este comando inicia un trabajo de entrenamiento desatendido para el equipo a con una asignación de una única GPU. El trabajo se basa en una imagen de docker de muestra, `gcr.io/run-ai-demo/quickstart`. Nosotros nombramos el trabajo `hyper1`. A continuación, puede supervisar el progreso del trabajo ejecutando el siguiente comando:

```
$ runai list
```

En la siguiente figura se muestra el resultado del `runai list` comando. Los Estados típicos que puede ver incluyen los siguientes:

- `ContainerCreating`. El contenedor docker se está descargando del repositorio en la nube.
- `Pending`. El trabajo está a la espera de ser programado.
- `Running`. El trabajo se está ejecutando.

```
~> runai list
Showing jobs for project team-a
NAME      STATUS  AGE  NODE                                     IMAGE                                     TYPE  PROJECT  USER  GPUs
hyper1    Running  11s  gke-dev-yaron1-gpu-4-pool-154f511d-5nk5 gcr.io/run-ai-demo/quickstart          Train team-a  yaron  1
```

Para obtener un estado adicional en su trabajo, ejecute el siguiente comando:

```
$ runai get hyper1
```

Para ver los registros del trabajo, ejecute el `runai logs <job-name>` comando:

```
$ runai logs hyper1
```

En este ejemplo, debería ver el registro de una sesión DL en ejecución, incluyendo la época de entrenamiento actual, ETA, valor de la función de pérdida, precisión y tiempo transcurrido para cada paso.

Puede ver el estado del clúster en la interfaz de usuario Run:AI en "<https://app.run.ai/>". En Paneles >

Descripción general, puede supervisar el uso de la GPU.

Para detener esta carga de trabajo, ejecute el siguiente comando:

```
$ runai delte hyper1
```

Este comando detiene la carga de trabajo de entrenamiento. Puede verificar esta acción ejecutando `runai list` de nuevo. Para obtener información detallada, consulte ["iniciar cargas de trabajo de formación sin supervisión"](#).

### Cargas de trabajo de compilación interactiva

Después de configurar proyectos y asignar GPU, puede ejecutar una carga de trabajo de compilación interactiva utilizando el siguiente comando en la línea de comandos:

```
$ runai submit build1 -i python -g 1 --interactive --command sleep --args infinity
```

El trabajo se basa en un pitón de imagen de Docker de muestra. Nombramos el edificio de trabajo 1.



La `-- interactive` indicador significa que el trabajo no tiene inicio ni fin. Es responsabilidad del investigador cerrar el trabajo. El administrador puede definir un límite de tiempo para los trabajos interactivos después de que el sistema los termine.

La `--g 1` El indicador asigna una única GPU a este trabajo. El comando y el argumento proporcionados son `--command sleep--args infinity`. Debe proporcionar un comando o el contenedor se inicia y, a continuación, sale inmediatamente.

Los siguientes comandos funcionan de forma similar a los comandos descritos en [Cargas de trabajo de formación desatendida](#):

- `runai list`: Muestra el nombre, el estado, la edad, el nodo, la imagen, Proyecto, usuario y GPU para trabajos.
- `runai get build1`: Muestra el estado adicional en la generación de trabajos 1.
- `runai delete build1`: Detiene la carga de trabajo interactiva build1. para obtener un shell bash en el contenedor, el siguiente comando:

```
$ runai bash build1
```

Esto proporciona un shell directo en el equipo. A continuación, los científicos de datos pueden desarrollar o afinar sus modelos dentro del contenedor.

Puede ver el estado del clúster en la interfaz de usuario Run:AI en ["https://app.run.ai"](https://app.run.ai). Para obtener información detallada, consulte ["iniciar y utilizar cargas de trabajo de compilación interactivas"](#).

## Cargas de trabajo interactivas con puertos conectados

Como extensión de las cargas de trabajo de compilación interactiva, puede revelar puertos internos al usuario del contenedor al iniciar un contenedor con la CLI Run:AI. Esto resulta útil para entornos cloud, trabajar con Jupyter Notebooks o conectarse a otros microservicios. "Entrada" Permite el acceso a los servicios de Kubernetes desde fuera del clúster de Kubernetes. Puede configurar el acceso mediante la creación de una colección de reglas que definan qué conexiones entrantes llegan a qué servicios.

Para mejorar la gestión del acceso externo a los servicios de un clúster, sugerimos que instalen los administradores de clúster "Entrada" Y configurar LoadBalancer.

Para utilizar Ingress como tipo de servicio, ejecute el siguiente comando para establecer el tipo de método y los puertos al enviar la carga de trabajo:

```
$ runai submit test-ingress -i jupyter/base-notebook -g 1 \
--interactive --service-type=ingress --port 8888 \
--args="--NotebookApp.base_url=test-ingress" --command=start-notebook.sh
```

Una vez que el contenedor se haya iniciado correctamente, ejecute `runai list` para ver la SERVICE URL (S) Con el que acceder al Cuaderno de Jupyter. La dirección URL está compuesta por el punto final de entrada, el nombre del trabajo y el puerto. Por ejemplo, consulte <https://10.255.174.13/test-ingress-8888>.

Para obtener información detallada, consulte "iniciar una carga de trabajo de compilación interactiva con puertos conectados".

### Se obtiene un uso elevado del clúster

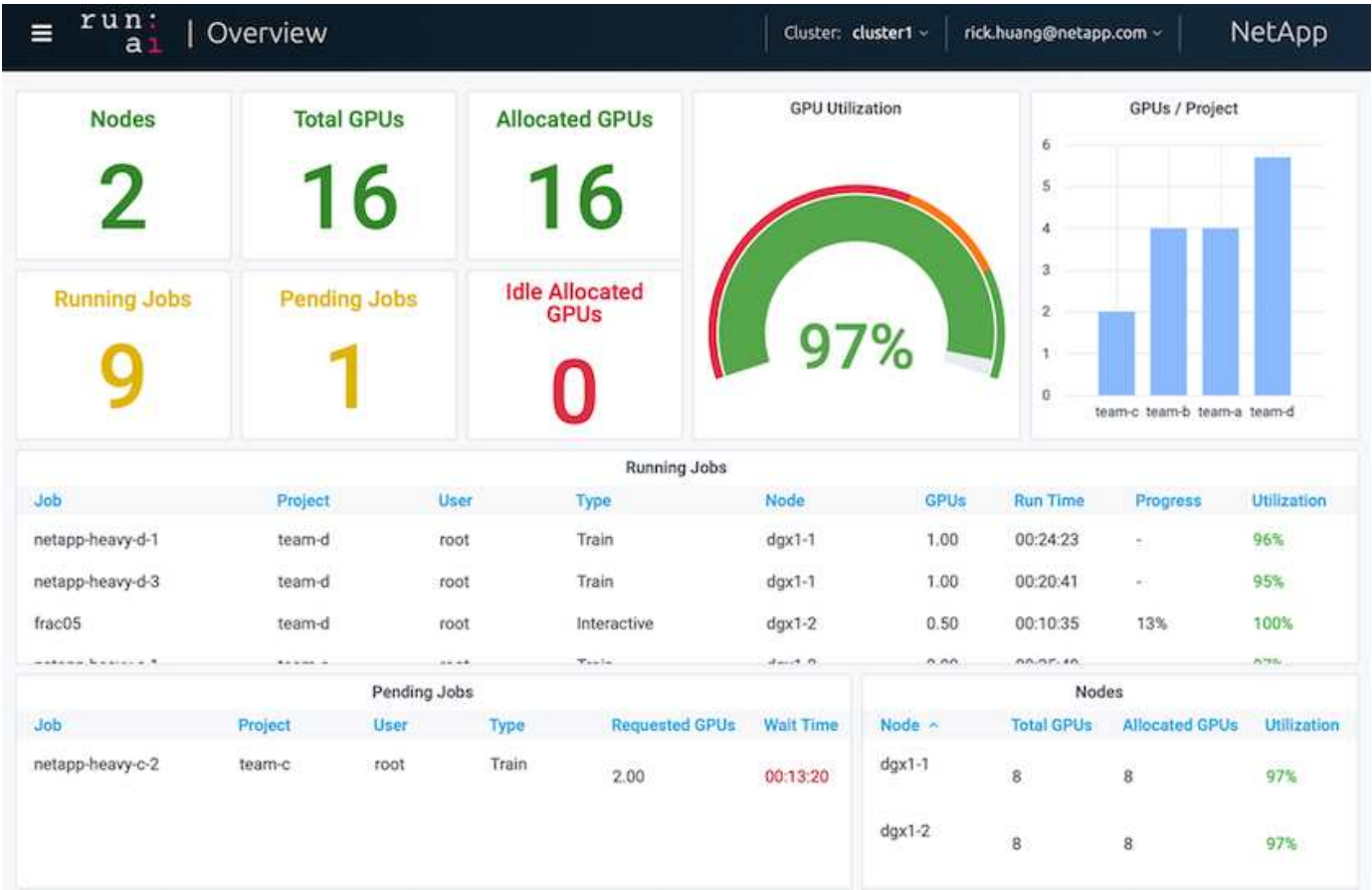
En esta sección, emulamos un escenario realista en el que cuatro equipos de ciencia de datos envían sus propias cargas de trabajo para demostrar la solución de orquestación Run:AI que logra un uso elevado del clúster mientras mantiene la priorización y el equilibrio de los recursos de la GPU. Empezamos utilizando la prueba de rendimiento ResNet-50 descrita en el apartado "RESNET-50 con resumen de la prueba de rendimiento de conjunto de datos ImageNET":

```
$ runai submit netapp1 -i netapp/tensorflow-tf1-py3:20.01.0 --local-image
--large-shm -v /mnt:/mnt -v /tmp:/tmp --command python --args
"/netapp/scripts/run.py" --args "--
dataset_dir=/mnt/mount_0/dataset/imagenet/imagenet_original/" --args "--
num_mounts=2" --args "--dgx_version=dgx1" --args "--num_devices=1" -g 1
```

Ejecutamos la misma prueba de rendimiento ResNet-50 como en la "NVA-1121". Usamos la bandera `--local-image` para contenedores que no residen en el repositorio de docker público. Montamos los directorios `/mnt` y `/tmp` En el nodo DGX-1 del host a `/mnt` y `/tmp` al contenedor, respectivamente. El conjunto de datos se encuentra en NetApp AFFA800 con el `dataset_dir` argumento que apunta al directorio. Ambas `--num_devices=1` y `-g 1` Significa que asignamos una GPU para este trabajo. El primero es un argumento para el `run.py` script, mientras que el último es un indicador para el `runai submit` comando.

La siguiente figura muestra un panel de información general del sistema con un uso del 97 % de la GPU y las

dieciséis GPU disponibles asignadas. Puede ver fácilmente cuántas GPU se asignan a cada equipo en el gráfico de barras GPU/proyecto. El panel trabajos en ejecución muestra los nombres de los trabajos en ejecución actuales, el proyecto, el usuario, el tipo, el nodo, Las GPU consumidas, tiempo de ejecución, progreso y detalles de uso. Se muestra una lista de las cargas de trabajo que están en cola con el tiempo de espera en trabajos pendientes. Finalmente, el recuadro Nodes ofrece cifras de GPU y utilización de nodos DGX-1 individuales en el clúster.



**Asignación de GPU fraccionaria para las cargas de trabajo menos exigentes o interactivas**

Cuando investigadores y desarrolladores trabajan en sus modelos, ya sea en las etapas de desarrollo, ajuste de hiperparámetros o depuración, estas cargas de trabajo suelen requerir menos recursos computacionales. Por lo tanto, es más eficiente aprovisionar GPU y memoria fraccionarias de modo que la misma GPU se pueda asignar simultáneamente a otras cargas de trabajo. Ejecutar:la solución de orquestación de IA proporciona un sistema de uso compartido de GPU fraccionario para las cargas de trabajo en contenedores en Kubernetes. El sistema admite cargas de trabajo que ejecutan programas CUDA y se adapta especialmente a tareas de IA ligeras como la inferencia y la creación de modelos. El sistema de GPU fraccionaria proporciona a los equipos de ciencia de datos e ingeniería de IA la capacidad de ejecutar varias cargas de trabajo simultáneamente en una única GPU. De este modo, las empresas pueden ejecutar más cargas de trabajo, como visión informática, reconocimiento de voz y procesamiento de lenguaje natural en el mismo hardware, con lo que se reducen los costes.

Ejecutar: el sistema de GPU fraccionaria de IA crea de manera efectiva GPU lógicas virtualizadas con su propia memoria y espacio de computación que los contenedores pueden utilizar y acceder como si fueran procesadores independientes. De este modo, es posible ejecutar varias cargas de trabajo en contenedores en paralelo y en la misma GPU sin interferir entre sí. La solución es transparente, sencilla y portátil y no requiere ningún cambio en los contenedores en sí.

Una usecase típica podría ver dos a ocho trabajos ejecutándose en la misma GPU, lo que significa que podría trabajar ocho veces más con el mismo hardware.

Para el trabajo `frac05` perteneciente al proyecto `team-d` En la siguiente figura podemos comprobar que el número de GPU asignadas era de 0.50. Esto es verificado por el `nvidia-smi` Comando, que muestra que la memoria de la GPU disponible para el contenedor era de 16,255 MB: La mitad de las 32 GB por GPU V100 en el nodo DGX-1.

```
root@run-deploy:~# runai bash frac05 -p team-d
root@frac05-0:/workload# nvidia-smi
Tue Jul 28 15:17:03 2020
```

NVIDIA-SMI 450.51.05 Driver Version: 450.51.05 CUDA Version: 11.0									
GPU	Name	Persistence-MI	Bus-Id	Disp.A	Volatile	Uncorr.	ECC		
Fan	Temp	Perf	Pwr:Usage/Cap	Memory-Usage	GPU-Util	Compute	M.		
							MIG	M.	
0	Tesla V100-SXM2...	On	00000000:07:00.0	Off			0		
N/A	57C	P0	240W / 300W	15525MiB / 16255MiB	100%	Default			
						N/A			

Processes:							
GPU	GI	CI	PID	Type	Process name	GPU Memory	
	ID	ID				Usage	
0	N/A	N/A	156	C	python3	15525MiB	

#### Utilización óptima del clúster gracias a la asignación de GPU por encima de una cuota

En esta sección y en las secciones ["Equidad en la asignación de recursos básicos"](#), y ["Justicia por exceso de cuotas"](#), Hemos diseñado escenarios de pruebas avanzados para demostrar las capacidades de orquestación Run:AI para la administración de cargas de trabajo complejas, la programación preventiva automática y el aprovisionamiento de GPU con exceso de cuota. Hemos hecho esto para lograr un uso elevado de los recursos de clúster y optimizar la productividad de los equipos científicos de datos a nivel empresarial en un entorno de IA de ONTAP.

Para estas tres secciones, defina los siguientes proyectos y cuotas:



Proyecto	Cuota
equipo a	4
equipo-b	2
equipo-c	2
equipo d	8

Además, utilizamos los siguientes contenedores para estas tres secciones:

- Portátil Jupyter: `jupyter/base-notebook`
- Ejecutar:AI Quickstart: `gcr.io/run-ai-demo/quickstart`

Definimos los siguientes objetivos para este escenario de prueba:

- Muestre la simplicidad del aprovisionamiento de recursos y cómo los recursos se abstraen de los usuarios
- Mostrar cómo los usuarios pueden aprovisionar fácilmente fracciones de una GPU y un número entero de GPU
- Muestre cómo el sistema elimina los cuellos de botella de computación al permitir que equipos o usuarios hagan uso de su cuota de recursos si hay GPU libres en el clúster
- Muestre cómo se eliminan los cuellos de botella en la canalización de datos utilizando la solución de NetApp cuando se ejecutan tareas con un gran procesamiento como el contenedor de NetApp
- Muestra cómo se ejecutan los diversos tipos de contenedores mediante el sistema
  - Portátil Jupyter
  - Ejecute:contenedor de IA
- Muestra una alta utilización cuando el clúster está lleno

Para obtener información detallada sobre la secuencia de comandos real ejecutada durante la prueba, consulte ["Detalles de las pruebas para la sección 4.8"](#).

Cuando se envían las 13 cargas de trabajo, puede ver una lista con los nombres de los contenedores y las GPU asignadas, como se muestra en la siguiente figura. Contamos con siete cursos de formación y seis trabajos interactivos, que simulan cuatro equipos de ciencia de datos, cada uno con sus propios modelos en ejecución o en desarrollo. Para trabajos interactivos, los desarrolladores individuales están utilizando Jupyter Notebooks para escribir o depurar su código. Por lo tanto, es adecuado aprovisionar fracciones de GPU sin usar demasiados recursos de clúster.

```
root@run-deploy:~# runai list -A
NAME          STATUS  AGE  NODE  IMAGE                                     TYPE      PROJECT  USER  GPUs  CREATED BY CLI  SERVICE URL(S)
b-4-gg        Running  2m   dgx1-2  gcr.io/run-ai-demo/quickstart           Train     team-b   root  2     true           http://10.61.218.134/a-1-1-jupyter,
c-5-g          Running  2m   dgx1-2  gcr.io/run-ai-demo/quickstart           Train     team-c   root  1     true           https://10.61.218.134/a-1-1-jupyter
c-4-gg        Running  2m   dgx1-1  gcr.io/run-ai-demo/quickstart           Train     team-c   root  2     true           http://10.61.218.134/a-1-1-jupyter,
b-3-g          Running  2m   dgx1-1  gcr.io/run-ai-demo/quickstart           Train     team-b   root  1     true           https://10.61.218.134/a-1-1-jupyter
c-3-g02        Running  2m   dgx1-1  gcr.io/run-ai-demo/quickstart           Interactive team-c   root  0.2   true           http://10.61.218.134/a-1-1-jupyter,
d-1-gggg       Running  2m   dgx1-2  gcr.io/run-ai-demo/quickstart           Train     team-d   root  4     true           https://10.61.218.134/a-1-1-jupyter
c-2-g03        Running  2m   dgx1-1  gcr.io/run-ai-demo/quickstart           Interactive team-c   root  0.3   true           http://10.61.218.134/a-1-1-jupyter,
c-1-g05        Running  2m   dgx1-1  gcr.io/run-ai-demo/quickstart           Interactive team-c   root  0.5   true           https://10.61.218.134/a-1-1-jupyter
a-2-gg         Running  3m   dgx1-1  gcr.io/run-ai-demo/quickstart           Train     team-a   root  2     true           http://10.61.218.134/a-1-1-jupyter,
b-2-g04        Running  3m   dgx1-2  gcr.io/run-ai-demo/quickstart           Interactive team-b   root  0.4   true           https://10.61.218.134/a-1-1-jupyter
a-1-g          Running  3m   dgx1-1  gcr.io/run-ai-demo/quickstart           Train     team-a   root  1     true           http://10.61.218.134/a-1-1-jupyter,
b-1-g06        Running  3m   dgx1-2  gcr.io/run-ai-demo/quickstart           Interactive team-b   root  0.6   true           https://10.61.218.134/a-1-1-jupyter
a-1-1-jupyter  Running  3m   dgx1-1  jupyter/base-notebook                   Interactive team-a   root  1     true           http://10.61.218.134/a-1-1-jupyter,
```

Los resultados de este escenario de prueba muestran lo siguiente:

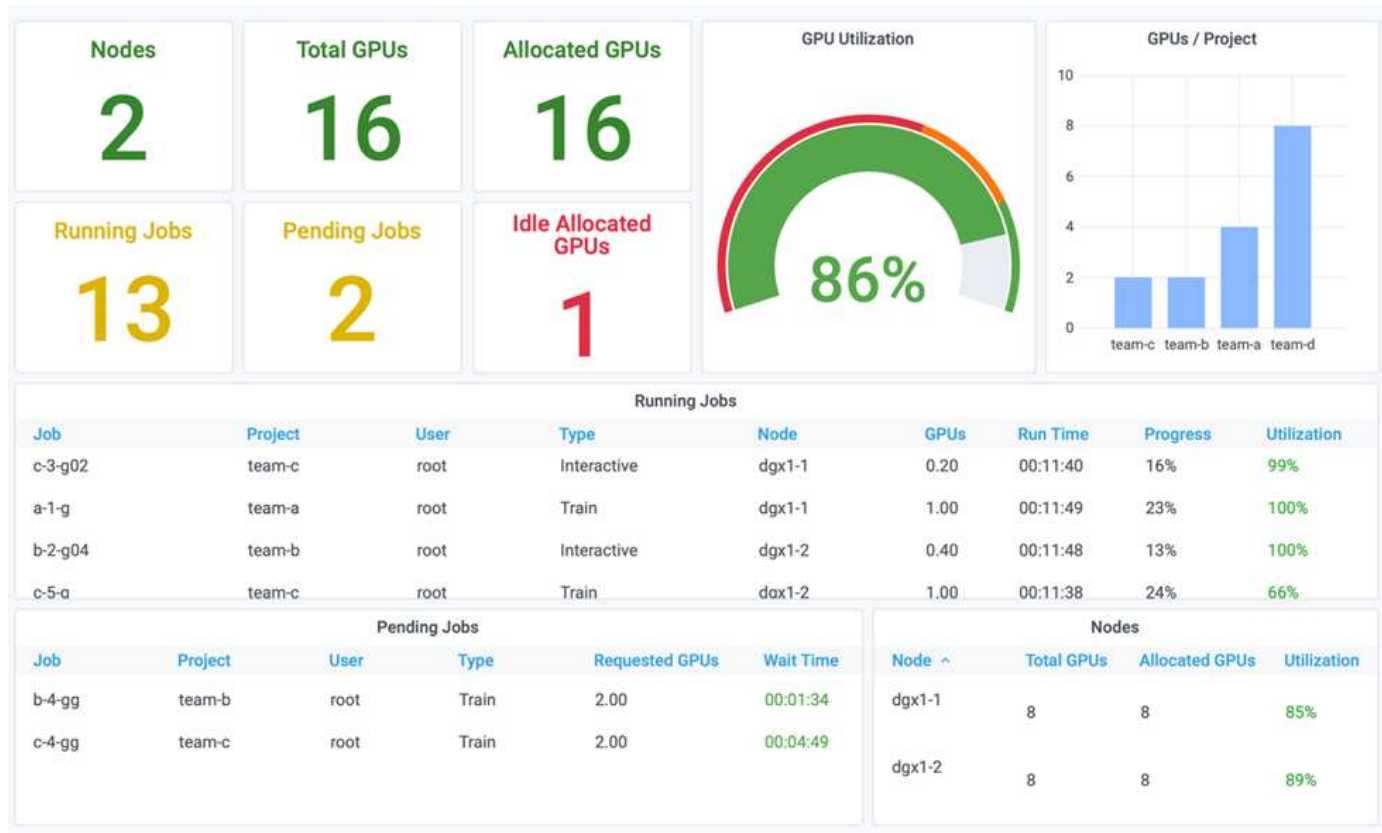
- El clúster debe estar lleno: Se utilizan 16/16 GPU.
- Uso elevado de clúster.
- Más experimentos que las GPU debido a la asignación fraccionaria.
- team-d no está utilizando toda su cuota; por lo tanto, team-b y.. team-c Puede utilizar más GPU para sus experimentos, lo que acelera el plazo de innovación.

#### Equidad en la asignación de recursos básicos

En esta sección, lo mostramos cuando team-d Solicita más GPU (están por debajo de su cuota), el sistema coloca en pausa las cargas de trabajo de team-b y.. team-c y los mueve a un estado pendiente de manera equitativa.

Para obtener información detallada, incluidos los envíos de trabajos, las imágenes contenedoras utilizadas y las secuencias de comandos ejecutadas, consulte la sección ["Detalles de las pruebas para la sección 4.9"](#).

La siguiente figura muestra el uso resultante del clúster, las GPU asignadas por equipo y los trabajos pendientes debido al equilibrio de carga automático y a la programación preventiva. Podemos observar que cuando el número total de GPU solicitado por todas las cargas de trabajo de equipo supera el total de GPU disponibles en el clúster, el algoritmo de integridad interno de Run:AI pone en pausa un trabajo por cada una team-b y.. team-c porque han cumplido con su cuota de proyecto. Esto proporciona una utilización elevada del clúster global, mientras que los equipos de ciencia de datos siguen trabajando con las limitaciones de recursos definidas por un administrador.



Los resultados de este escenario de prueba demuestran lo siguiente:

- **Equilibrio de carga automático.** el sistema equilibra automáticamente la cuota de las GPU, de modo que cada equipo utiliza ahora su cuota. Las cargas de trabajo en pausa pertenecen a equipos que se



encontraban por encima de su cuota.

- **Pausa de uso compartido justo.** el sistema elige detener la carga de trabajo de un equipo que estaba por encima de su cuota y luego detener la carga de trabajo del otro equipo. Ejecutar: la IA tiene algoritmos internos de justicia.

### Justicia por exceso de cuotas

En esta sección, ampliamos el escenario en el que varios equipos envían cargas de trabajo y superan su cuota. De esta manera, mostramos cómo el algoritmo de equidad de Run:AI asigna recursos de clúster según la proporción de cuotas predefinidas.

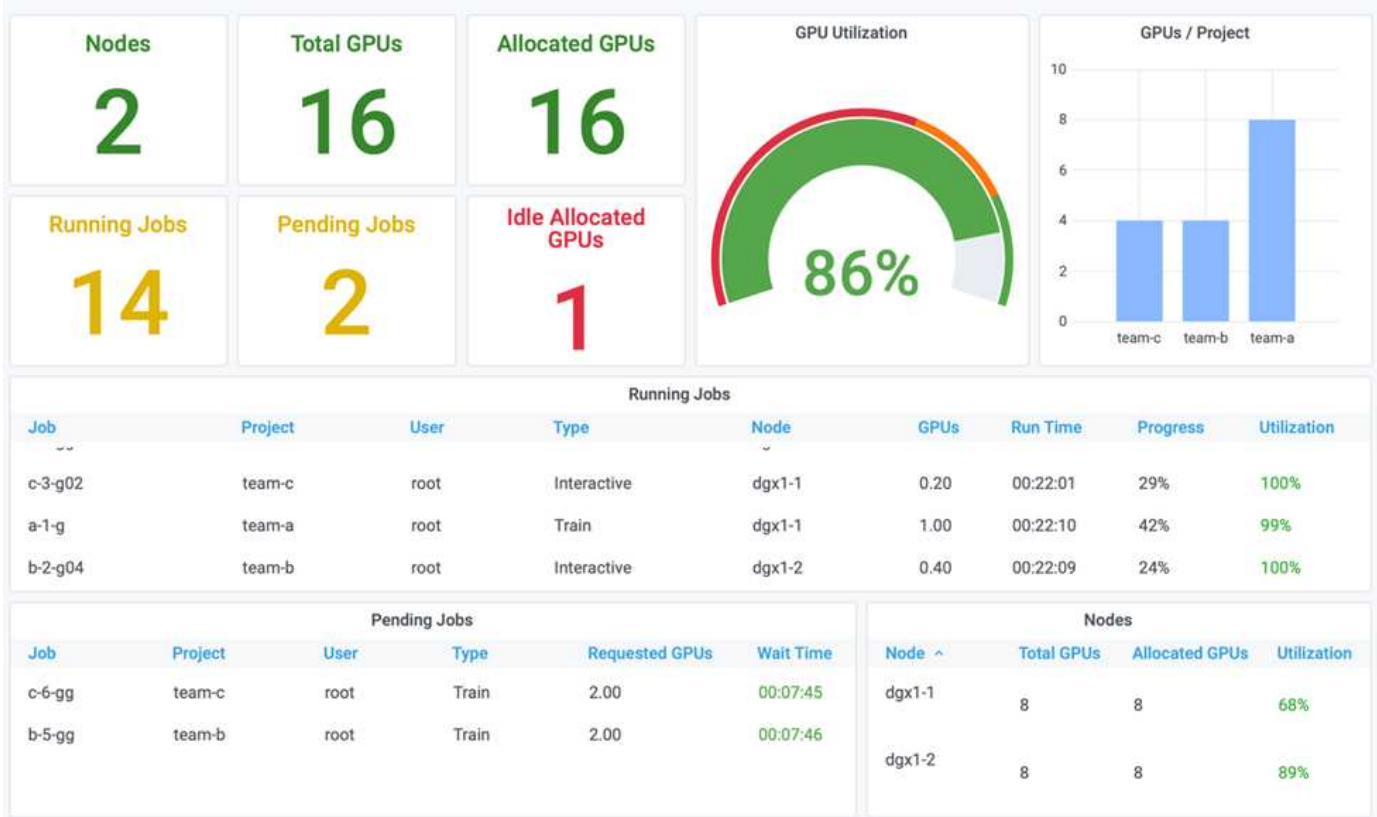
Objetivos de este escenario de prueba:

- Muestra el mecanismo de cola cuando varios equipos solicitan GPU a través de su cuota.
- Muestre cómo distribuye el sistema una cuota del clúster entre varios equipos que se encuentran por encima de su cuota en función de la proporción entre sus cuotas, de modo que el equipo con la cuota mayor obtenga una parte mayor de la capacidad de reserva.

Al final de ["Equidad en la asignación de recursos básicos"](#), hay dos cargas de trabajo en cola: una para `team-b` y uno para `team-c`. En esta sección, ponemos en cola las cargas de trabajo adicionales.

Para obtener información detallada, incluidos los envíos de trabajos, las imágenes contenedoras utilizadas y las secuencias de comandos ejecutadas, consulte ["Detalles de la prueba para la sección 4.10"](#).

Cuando todos los trabajos se someten según la sección ["Detalles de la prueba para la sección 4.10"](#), el panel del sistema muestra eso `team-a`, `team-b`, y `team-c`. Todos tienen más GPU que su cuota predefinida. `team-a` Ocupa cuatro GPU más que su cuota de software predefinida (cuatro), mientras que `team-b` y `team-c`. Cada uno ocupa dos GPU más que su cuota blanda (dos). La proporción de GPU asignadas por encima de una cuota es igual a la de su cuota predefinida. Esto se debe a que el sistema utilizó la cuota preconfigurada como referencia de prioridad y se aprovisionó de manera acorde cuando varios equipos soliciten más GPU, por lo que superó su cuota. Este equilibrio de carga automático proporciona justicia y priorización cuando los equipos de ciencia de datos empresariales están involucrados activamente en el desarrollo y la producción de modelos de IA.



Los resultados de este escenario de prueba muestran lo siguiente:

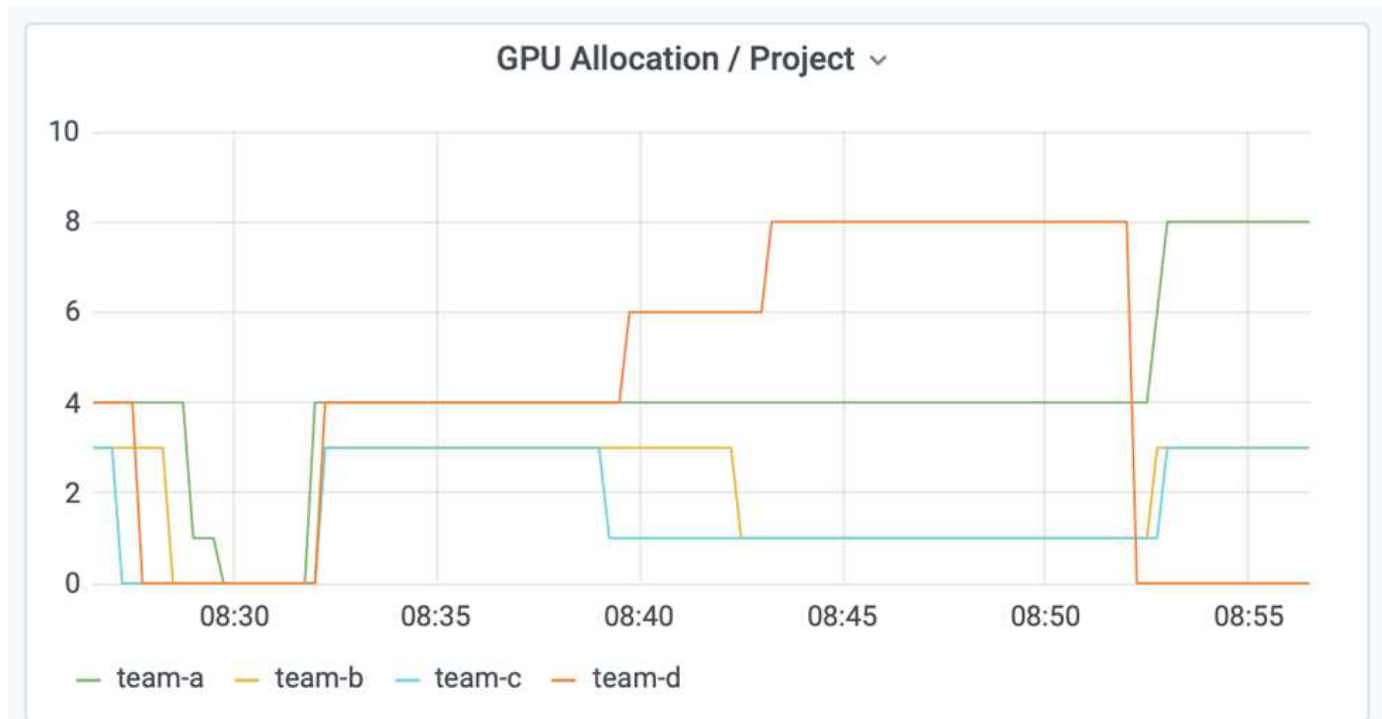
- El sistema empieza a deponer las cargas de trabajo de otros equipos en cola.
- El orden de la depuesta en cola se decide de acuerdo a los algoritmos de justicia, tal que team-b y team-c Obtenga la misma cantidad de GPU de cuota superior (ya que tienen una cuota similar) y team-a Obtiene una cantidad doble de GPU, ya que su cuota es dos veces superior a la de team-b y team-c.
- Toda la asignación se realiza automáticamente.

Por lo tanto, el sistema debe estabilizarse en los siguientes estados:

Proyecto	GPU asignadas	Comentar
equipo a	8/4	Cuatro GPU a lo largo de la cuota. Cola vacía.
equipo-b	4/2	Dos GPU por encima de la cuota. Una carga de trabajo en cola.
equipo-c	4/2	Dos GPU por encima de la cuota. Una carga de trabajo en cola.
equipo d	0/8	No utiliza GPU en absoluto, sin cargas de trabajo en cola.

La siguiente figura muestra la asignación de GPU por proyecto a lo largo del tiempo en el panel Run:AI Analytics de las secciones "Utilización óptima del clúster gracias a la asignación de GPU por encima de una cuota", "Equidad en la asignación de recursos básicos", y "Justicia por exceso de cuotas". Cada línea de la figura indica el número de GPU aprovisionadas para un equipo de ciencias de datos en cualquier momento. Podemos observar que el sistema asigna dinámicamente las GPU de acuerdo con las cargas de trabajo

enviadas. Esto permite a los equipos revisar las cuotas cuando hay GPU disponibles en el clúster y, después, adelantarse a los trabajos según sea necesario, antes de alcanzar, por fin, un estado estable para los cuatro equipos.



#### Guardar datos en un volumen persistente aprovisionado por Trident

Trident de NetApp es un proyecto de código abierto totalmente compatible diseñado para ayudarle a satisfacer las sofisticadas demandas de persistencia de sus aplicaciones en contenedores. Puede leer y escribir datos en un volumen persistente (VP) de Kubernetes aprovisionado por Trident con la ventaja añadida de la organización en niveles de datos, el cifrado, la tecnología Snapshot de NetApp, el cumplimiento de normativas y el alto rendimiento que ofrece el software de gestión de datos ONTAP de NetApp.

#### Reutilizar EVs en un espacio de nombres existente

En el caso de proyectos de IA de mayor tamaño, es posible que sea más eficiente que diferentes contenedores para leer y escribir datos en el mismo PV de Kubernetes. Para reutilizar una solicitud de volumen persistente de Kubernetes (PVC), el usuario ya debe haber creado una RVP. Consulte ["Documentación de Trident de NetApp"](#) Para obtener más detalles sobre la creación de una RVP. A continuación se muestra un ejemplo de reutilización de un PVC existente:

```
$ runai submit pvc-test -p team-a --pvc test:/tmp/pvc1mount -i gcr.io/run-ai-demo/quickstart -g 1
```

Ejecute el siguiente comando para ver el estado del trabajo `pvc-test` para el proyecto `team-a`:

```
$ runai get pvc-test -p team-a
```

Debe ver el montaje PV /tmp/pvc1 montado a. team-a trabajo pvc-test. De este modo, se pueden leer varios contenedores desde el mismo volumen, lo que resulta útil cuando hay varios modelos competidores en desarrollo o producción. Los científicos de datos pueden crear un conjunto de modelos y, a continuación, combinar los resultados de las predicciones por voto mayoritario u otras técnicas.

Utilice lo siguiente para acceder al shell del contenedor:

```
$ runai bash pvc-test -p team-a
```

A continuación, puede comprobar el volumen montado y acceder a los datos dentro del contenedor.

Esta capacidad de reutilización de RVP funciona con volúmenes FlexVol de NetApp y volúmenes ONTAP FlexGroup de NetApp, lo que permite a los ingenieros de datos disfrutar de opciones de gestión de datos más flexibles y sólidas para aprovechar su Data Fabric con tecnología de NetApp.

Conclusión

NetApp y Run:AI se han asociado en este informe técnico para mostrar las exclusivas funcionalidades de la solución ONTAP AI de NetApp junto con la plataforma Run:AI para simplificar la orquestación de las cargas de trabajo de IA. Los pasos anteriores proporcionan una arquitectura de referencia para optimizar el proceso de canalizaciones de datos y orquestación de cargas de trabajo para el aprendizaje profundo. Animamos a los clientes que deseen poner en marcha estas soluciones a NetApp y a Run:AI para obtener más información.

Detalles de las pruebas para la sección 4.8

Esta sección contiene los detalles de las pruebas de la sección "Utilización óptima del clúster gracias a la asignación de GPU por encima de una cuota".

Enviar trabajos en el siguiente orden:

Proyecto	Imagen	N.o de GPU	Total	Comentar
equipo a	Jupyter	1	1/4	–
equipo a	NetApp	1	2/4	–
equipo a	Ejecución: IA	2	4/4	Usando toda su cuota
equipo-b	Ejecución: IA	0.6	0.6/2	GPU fraccionaria
equipo-b	Ejecución: IA	0.4	1/2	GPU fraccionaria
equipo-b	NetApp	1	2/2	–
equipo-b	NetApp	2	4/2	Dos sobre cuota
equipo-c	Ejecución: IA	0.5	0.5/2	GPU fraccionaria
equipo-c	Ejecución: IA	0.3	0.8/2	GPU fraccionaria
equipo-c	Ejecución: IA	0.2	1/2	GPU fraccionaria

Proyecto	Imagen	N.o de GPU	Total	Comentar
equipo-c	NetApp	2	3/2	Uno sobre cuotas
equipo-c	NetApp	1	4/2	Dos sobre cuota
equipo d	NetApp	4	4/8	Utilizando la mitad de su cuota

Estructura de comandos:

```
$ runai submit <job-name> -p <project-name> -g <#GPUs> -i <image-name>
```

Secuencia de comandos real utilizada en la prueba:

```
$ runai submit a-1-1-jupyter -i jupyter/base-notebook -g 1 \
  --interactive --service-type=ingress --port 8888 \
  --args="--NotebookApp.base_url=team-a-test-ingress" --command=start
-notebook.sh -p team-a
$ runai submit a-1-g -i gcr.io/run-ai-demo/quickstart -g 1 -p team-a
$ runai submit a-2-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-a
$ runai submit b-1-g06 -i gcr.io/run-ai-demo/quickstart -g 0.6
--interactive -p team-b
$ runai submit b-2-g04 -i gcr.io/run-ai-demo/quickstart -g 0.4
--interactive -p team-b
$ runai submit b-3-g -i gcr.io/run-ai-demo/quickstart -g 1 -p team-b
$ runai submit b-4-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-b
$ runai submit c-1-g05 -i gcr.io/run-ai-demo/quickstart -g 0.5
--interactive -p team-c
$ runai submit c-2-g03 -i gcr.io/run-ai-demo/quickstart -g 0.3
--interactive -p team-c
$ runai submit c-3-g02 -i gcr.io/run-ai-demo/quickstart -g 0.2
--interactive -p team-c
$ runai submit c-4-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-c
$ runai submit c-5-g -i gcr.io/run-ai-demo/quickstart -g 1 -p team-c
$ runai submit d-1-gggg -i gcr.io/run-ai-demo/quickstart -g 4 -p team-d
```

En este punto, debe tener los siguientes estados:

Proyecto	GPU asignadas	Cargas de trabajo en cola
equipo a	4/4 (cuota suave/asignación real)	Ninguno
equipo-b	4/2	Ninguno
equipo-c	4/2	Ninguno
equipo d	4/8	Ninguno

Consulte la sección "[Mayor uso de clúster con la asignación de GPU over-quota](#)" para discusiones sobre el escenario de pruebas en curso.

#### Detalles de las pruebas para la sección 4.9

Esta sección contiene detalles de la prueba para la sección "[Equidad en la asignación de recursos básicos](#)".

Enviar trabajos en el siguiente orden:

Proyecto	N.o de GPU	Total	Comentar
equipo d	2	6/8	La carga de trabajo Team-b/c se pone en pausa y se mueve a. pending.
equipo d	2	8/8	Las cargas de trabajo de otros equipos (b/c) se pausan y pasan a. pending.

Consulte la siguiente secuencia de comandos ejecutada:

```
$ runai submit d-2-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-d$  
runai submit d-3-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-d
```

En este punto, debe tener los siguientes estados:

Proyecto	GPU asignadas	Cargas de trabajo en cola
equipo a	4/4	Ninguno
equipo-b	2/2	Ninguno
equipo-c	2/2	Ninguno
equipo d	8/8	Ninguno

Consulte la sección "[Equidad en la asignación de recursos básicos](#)" para un debate sobre el escenario de pruebas en curso.

#### Detalles de las pruebas para la sección 4.10

Esta sección contiene detalles de la prueba para la sección "[Justicia por exceso de cuotas](#)".

Enviar trabajos en el siguiente orden para team-a, team-b, y. team-c:

Proyecto	N.o de GPU	Total	Comentar
equipo a	2	4/4	1 carga de trabajo en cola

Proyecto	N.o de GPU	Total	Comentar
equipo a	2	4/4	2 cargas de trabajo en cola
equipo-b	2	2/2	2 cargas de trabajo en cola
equipo-c	2	2/2	2 cargas de trabajo en cola

Consulte la siguiente secuencia de comandos ejecutada:

```
$ runai submit a-3-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-a$
runai submit a-4-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-a$ runai
submit b-5-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-b$ runai
submit c-6-gg -i gcr.io/run-ai-demo/quickstart -g 2 -p team-c
```

En este punto, debe tener los siguientes estados:

Proyecto	GPU asignadas	Cargas de trabajo en cola
equipo a	4/4	Dos cargas de trabajo que solicitan GPU dos cada una
equipo-b	2/2	Dos cargas de trabajo que solicitan dos GPU cada una
equipo-c	2/2	Dos cargas de trabajo que solicitan dos GPU cada una
equipo d	8/8	Ninguno

A continuación, elimine todas las cargas de trabajo para team-d:

```
$ runai delete -p team-d d-1-gggg d-2-gg d-3-gg
```

Consulte la sección "[Justicia por exceso de cuotas](#)", para las discusiones sobre el escenario de pruebas en curso.

## Dónde encontrar información adicional

Si quiere obtener más información sobre el contenido de este documento, consulte los siguientes recursos:

- Sistemas DGX de NVIDIA
  - Sistema DGX-1 de NVIDIA <https://www.nvidia.com/en-us/data-center/dgx-1/>
  - GPU de núcleo tensor NVIDIA V100 <https://www.nvidia.com/en-us/data-center/tesla-v100/>
  - NVIDIA NGC <https://www.nvidia.com/en-us/gpu-cloud/>
- Ejecute: solución de orquestación de contenedores de IA

- Ejecución: Introducción de producto de IA <https://docs.run.ai/home/components/>
- Ejecución: documentación de instalación de IA <https://docs.run.ai/Administrator/Cluster-Setup/Installing-Run-AI-on-an-on-premise-Kubernetes-Cluster/>  
<https://docs.run.ai/Administrator/Researcher-Setup/Installing-the-Run-AI-Command-Line-Interface/>
- Enviar trabajos en Ejecutar:AI CLI <https://docs.run.ai/Researcher/Walkthroughs/Walkthrough-Launch-Unattended-Training-Workloads-/>  
<https://docs.run.ai/Researcher/Walkthroughs/Walkthrough-Start-and-Use-Interactive-Build-Workloads-/>
- Asignar fracciones de GPU en la CLI de Run:AI <https://docs.run.ai/Researcher/Walkthroughs/Walkthrough-Using-GPU-Fractions/>
- Plano de control de IA de NetApp
  - Informe técnico <https://www.netapp.com/us/media/tr-4798.pdf>
  - Demostración en formato breve [https://youtu.be/gfr\\_sO27Rvo](https://youtu.be/gfr_sO27Rvo)
  - Repositorio de GitHub [https://github.com/NetApp/kubeflow\\_jupyter\\_pipeline](https://github.com/NetApp/kubeflow_jupyter_pipeline)
- Sistemas AFF de NetApp
  - Especificaciones técnicas de AFF a-Series de NetApp <https://www.netapp.com/us/media/ds-3582.pdf>
  - Ventaja de NetApp Flash para All Flash FAS <https://www.netapp.com/us/media/ds-3733.pdf>
  - Biblioteca de información de ONTAP <http://mysupport.netapp.com/documentation/productlibrary/index.html?productID=62286>
  - Informe técnico de NetApp ONTAP FlexGroup Volumes <https://www.netapp.com/us/media/tr-4557.pdf>
- ONTAP AI de NetApp
  - Guía de diseño de ONTAP AI con DGX-1 y Cisco Networking <https://www.netapp.com/us/media/nva-1121-design.pdf>
  - Guía de puesta en marcha de ONTAP AI con DGX-1 y Cisco Networking <https://www.netapp.com/us/media/nva-1121-deploy.pdf>
  - Guía de diseño de ONTAP AI con DGX-1 y Mellanox Networking <http://www.netapp.com/us/media/nva-1138-design.pdf>
  - Guía de diseño de ONTAP AI con DGX-2 <https://www.netapp.com/us/media/nva-1135-design.pdf>

## **TR-4799-DESIGN: Arquitectura de referencia ONTAP AI de NetApp para cargas de trabajo de conducción autónoma**

David Arnette y Sung-han Lin, NetApp

La familia de sistemas NVIDIA DGX es la primera plataforma de inteligencia artificial (IA) integrada del mundo, específicamente diseñada para IA empresarial. Los sistemas de almacenamiento AFF de NetApp proporcionan un rendimiento extremo y funcionalidades de gestión de datos de cloud híbrido líderes en el sector. NetApp y NVIDIA se han asociado para crear la arquitectura de referencia de IA ONTAP de NetApp con el fin de ofrecer a los clientes una solución lista para usar que respalda cargas de trabajo DE IA y aprendizaje automático (ML) con rendimiento, fiabilidad y soporte de nivel empresarial.

["TR-4799-DESIGN: Arquitectura de referencia ONTAP AI de NetApp para cargas de trabajo de conducción autónoma"](#)



## **TR-4811: Arquitectura de referencia de IA ONTAP de NetApp para el sector sanitario: Imágenes de diagnóstico - diseño de la solución**

Rick Huang, Sung-han Lin, Sathish Thyagarajan, NetApp Jacci Cenci, NVIDIA

Esta arquitectura de referencia ofrece directrices para los clientes que crean una infraestructura de inteligencia artificial (IA) con sistemas NVIDIA DGX-2 y almacenamiento AFF de NetApp para casos de uso de sanidad. Incluye información sobre los flujos de trabajo de alto nivel utilizados en el desarrollo de modelos de aprendizaje profundo (DL) para imágenes de diagnóstico médico, casos de prueba validados y resultados. También incluye recomendaciones para el dimensionamiento de las puestas en marcha de los clientes.

["TR-4811: Arquitectura de referencia de IA ONTAP de NetApp para el sector sanitario: Imágenes de diagnóstico - diseño de la solución"](#)

## **TR-4807: Arquitectura de referencia de IA de ONTAP de NetApp para cargas de trabajo de servicios financieros - Diseño de la solución**

Karthikeyan Nagalingam, Sung-han Lin, NetApp Jacci Cenci, NVIDIA

Esta arquitectura de referencia ofrece directrices a los clientes que crean una infraestructura de inteligencia artificial con los sistemas NVIDIA DGX-1 y el almacenamiento AFF de NetApp para casos de uso con sectores financieros. Por ejemplo, incluye información sobre los flujos de trabajo de alto nivel utilizados en el desarrollo de modelos de aprendizaje profundo para los casos de pruebas y los resultados de servicios financieros. También incluye recomendaciones para el dimensionamiento de las puestas en marcha de los clientes.

["TR-4807: Arquitectura de referencia de IA de ONTAP de NetApp para cargas de trabajo de servicios financieros - Diseño de la solución"](#)

## **IA generativa y valor de NetApp**

Autor: Sathish Thyagarajan, NetApp

### **Sinopsis**

La demanda de inteligencia artificial (IA) generativa está impulsando la disrupción en todos los sectores y mejorando la creatividad empresarial y la innovación de productos. Muchas organizaciones usan la IA generativa para crear nuevas características de productos, mejorar la productividad de ingeniería y crear prototipos de aplicaciones impulsadas por IA que ofrezcan mejores resultados y experiencias de los consumidores. La IA generativa, como los Transformadores Generativos Pre-entrenados (GPT), utilizan redes neuronales para crear nuevo contenido, tan diverso como texto, audio y vídeo. Dada la escala extrema y los conjuntos de datos masivos que participan en los modelos de lenguaje grande (LLM), es crucial diseñar una infraestructura de IA robusta que aproveche las atractivas funciones de almacenamiento de datos de las opciones de puesta en marcha locales, híbridas y multicloud, y reducir los riesgos asociados a la movilidad de datos, la gobernanza y la protección de datos antes de que las empresas puedan diseñar soluciones de IA. En este documento se describen estas consideraciones y las capacidades de IA de NetApp® correspondientes que permiten una gestión de datos y un movimiento de datos sin interrupciones en la canalización de datos de

IA para el entrenamiento, el reciclaje, el ajuste y la inferencia de modelos de IA generativos.

## Resumen ejecutivo

Más recientemente, después del lanzamiento de ChatGPT, un spin-off de GPT-3 en noviembre de 2022, las nuevas herramientas de IA utilizadas para generar texto, código, imagen o incluso proteínas terapéuticas en respuesta a las indicaciones de los usuarios han ganado fama significativa. Esto indica que los usuarios pueden realizar una solicitud utilizando lenguaje natural y la IA interpretará y generará texto, como artículos de noticias o descripciones de productos que reflejen la solicitud del usuario o produzcan código, música, voz, efectos visuales y activos 3D mediante algoritmos entrenados en datos ya existentes. Como resultado, frases como la difusión estable, alucinaciones, ingeniería de avisos y alineación de valores están emergiendo rápidamente en el diseño de sistemas de IA. Estos modelos de aprendizaje automático (ML) autosupervisados o semisupervisados están cada vez más disponibles como modelos de base preformados (FM) a través de los proveedores de servicios cloud y otros proveedores de firma de IA, que están adoptando diversos establecimientos empresariales de diferentes sectores para una amplia gama de tareas descendentes de procesamiento del lenguaje natural (NLP). Como afirman firmas de analistas de investigación como McKinsey: “El impacto de la IA generativa en la productividad podría agregar billones de dólares en valor a la economía global”. Mientras que las empresas se están replanteando la IA como partners de pensamiento para los humanos y FMS se están ampliando a la vez a lo que las empresas e instituciones pueden hacer con la IA generativa, las oportunidades para gestionar volúmenes masivos de datos seguirán creciendo. Este documento presenta información introductoria sobre IA generativa y los conceptos de diseño con relación a las funcionalidades de NetApp que aportan valor a los clientes de NetApp, tanto en entornos híbridos como en las instalaciones o multicloud.

**Por lo tanto, ¿qué es lo que los clientes utilizan NetApp en sus entornos de IA?** NetApp ayuda a las organizaciones a satisfacer las complejidades creadas por el rápido crecimiento de datos y nube, la gestión de múltiples nubes y la adopción de tecnologías de nueva generación, como la IA. NetApp ha combinado varias funcionalidades en software de gestión inteligente de datos e infraestructuras de almacenamiento que han estado bien equilibradas con un alto rendimiento optimizado para cargas de trabajo de IA. Las soluciones de IA generativas como las LLM necesitan leer y procesar sus conjuntos de datos de origen desde el almacenamiento a la memoria numerosas veces para fomentar la inteligencia. NetApp ha sido una empresa líder en tecnologías de movilidad de datos, gobernanza de datos y seguridad de datos en el ecosistema del perímetro al núcleo y al cloud, sirviendo a los clientes empresariales a crear soluciones de IA a escala. NetApp, con una sólida red de partners ha ayudado a directores de datos, ingenieros de IA, arquitectos empresariales y científicos de datos en el diseño de una canalización de datos libre para la preparación, la protección de datos así como las responsabilidades estratégicas de gestión de datos de la formación e inferencia de modelos de IA, optimizando el rendimiento y la escalabilidad del ciclo de vida de la IA/ML. Tecnologías y capacidades de datos de NetApp, como NetApp® ONTAP AI® para la canalización de datos de aprendizaje profundo, NetApp® SnapMirror® para transportar datos de forma fluida y eficiente entre puntos finales de almacenamiento, y NetApp® FlexCache® para el renderizado en tiempo real. Cuando el flujo de datos cambia de batch a tiempo real y la ingeniería de datos ocurre en un momento oportuno, aporta valor a la implementación de modelos de IA generativa en tiempo real. A medida que las empresas de todo tipo adoptan nuevas herramientas de IA, deben hacer frente a desafíos en materia de datos desde el perímetro hasta el centro de datos y el cloud que demandan soluciones de IA escalables, responsables y explicables. Como referente en datos para cloud híbrido y multicloud, NetApp se compromete a crear una red de partners y soluciones conjuntas que pueden ayudar en todos los aspectos de la construcción de una canalización de datos y lagos de datos para el entrenamiento de modelos de IA generativos (entrenamiento previo), el ajuste, la inferencia basada en el contexto y la supervisión del deterioro de modelos de las LL.

## ¿Qué es la IA generativa?









La IA generativa está cambiando la forma en que creamos contenido, generamos nuevos conceptos de diseño y exploramos composiciones novedosas. Ilustra los marcos de redes neuronales como la Red Adversario Generativa (GAN), los Autocodificadores Variacionales (VAE) y los Transformadores Generativos Pre-entrenados (GPT), que pueden generar contenido nuevo, como texto, código, imágenes, audio, vídeo, etc. y

datos sintéticos. Los modelos basados en transformadores como Chat-GPT de OpenAI, Bard de Google, abrazar la FLORACIÓN DE Face y la llama de Meta han surgido como la tecnología fundamental que sustenta muchos avances en modelos lingüísticos grandes. Del mismo modo, Dall-E de OpenAI, Meta's CM3leon y Google's Imagen son ejemplos de modelos de difusión de texto a imagen que ofrecen a los clientes un grado de fotorealismo sin precedentes para crear imágenes nuevas y complejas desde cero o editar imágenes existentes para generar imágenes contextuales de alta calidad utilizando el aumento de conjuntos de datos y la síntesis de texto a imagen enlazando textuales y semántica visual. Los artistas digitales están empezando a aplicar una combinación de tecnologías de renderizado como Nerf (Neural Radiance Field) con IA generativa para convertir imágenes 2D estáticas en escenas 3D inmersivas. En general, las LLM se caracterizan ampliamente por cuatro parámetros: (1) Tamaño del modelo (normalmente en miles de millones de parámetros); (2) Tamaño del conjunto de datos de entrenamiento; (3) Costo de entrenamiento y (4) Rendimiento del modelo después del entrenamiento. Los LLMs también se dividen principalmente en tres arquitecturas de transformadores. (i) Modelos solo para codificadores. Por ejemplo, BERT (Google, 2018); (ii) Encoder-Decoder, por ejemplo, BART (Meta, 2020) y (iii) modelos solo Decoder. Por ejemplo, llama (Meta, 2023), Palm-E (Google, 2023). Según los requisitos del negocio, independientemente de la arquitectura que elija una compañía el número de parámetros del modelo (N) y el número de tokens (D) del conjunto de datos de entrenamiento generalmente determinan el coste base de entrenamiento (entrenamiento previo) o el ajuste de un LLM.

### **Casos prácticos empresariales y tareas de NLP descendentes**

Las empresas de todos los sectores están descubriendo cada vez más potencial para que la IA extraiga y produzca nuevas formas de valor a partir de los datos existentes para operaciones empresariales, ventas, marketing y servicios legales. Según la inteligencia de mercado de IDC (International Data Corporation) sobre casos de uso e inversiones de IA generativa global, la gestión del conocimiento en el desarrollo de software y el diseño de productos será la más afectada, seguida de la creación de historias para el marketing y la generación de código para los desarrolladores. En el sector sanitario, las organizaciones de investigación clínica están abriendo nuevos caminos en la medicina. Los modelos preentrenados como ProteinBERT incorporan anotaciones de Ontología Genética (GO) para diseñar rápidamente estructuras de proteínas para medicamentos médicos, lo que representa un hito significativo en el descubrimiento de fármacos, la bioinformática y la biología molecular. Las empresas de biotecnología han iniciado ensayos en humanos para la medicina generalizada descubierta por IA, que tiene como objetivo tratar enfermedades como la fibrosis pulmonar (FPI), una enfermedad pulmonar que causa cicatrices irreversibles del tejido pulmonar.

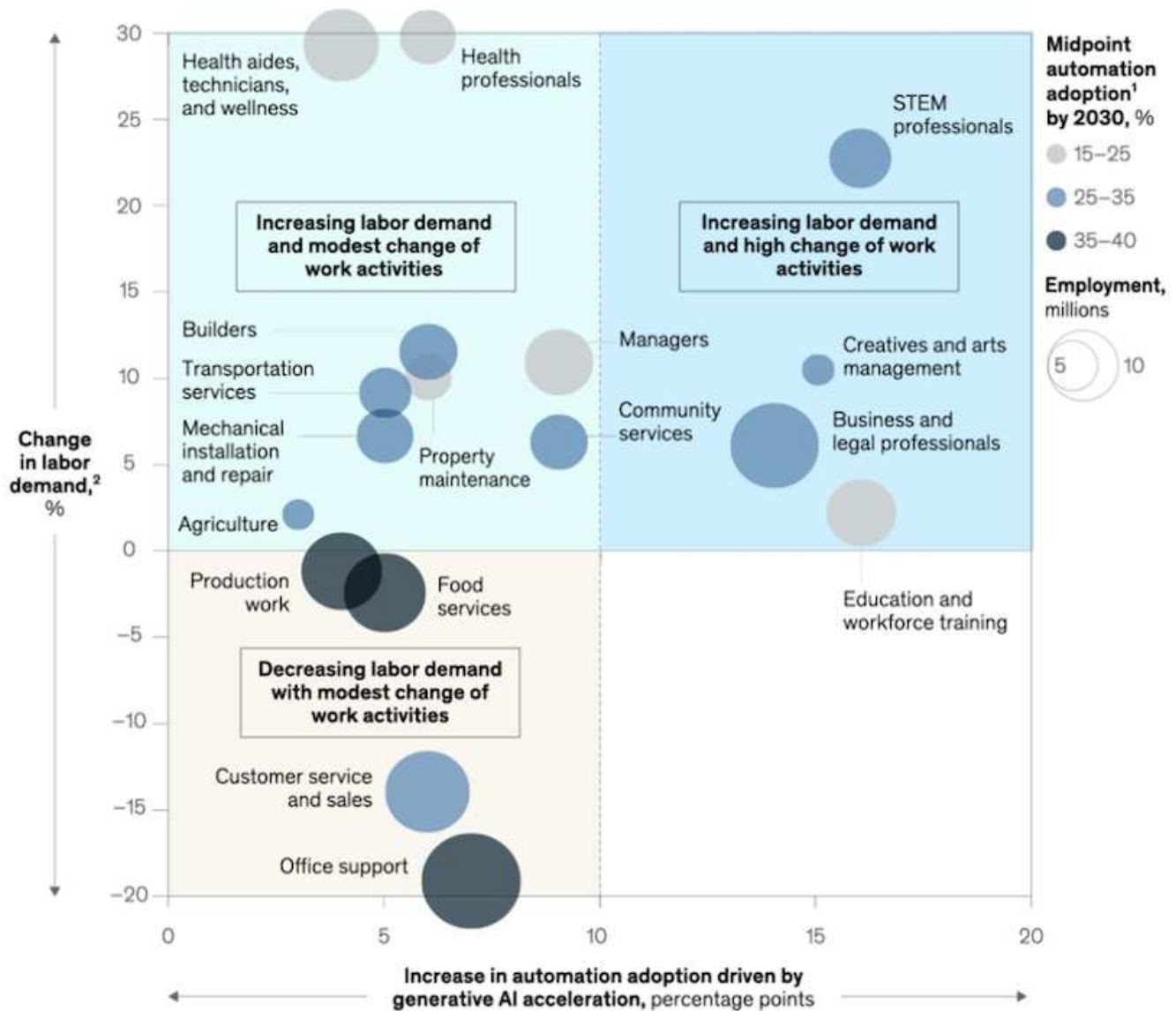
Figura 1: Casos prácticos que impulsan la IA generativa

 <p><b>Chatbots</b></p>	 <p><b>Drug discovery</b></p>
 <p><b>Text generation</b></p>	 <p><b>Genome model expression</b></p>
 <p><b>Image generation</b></p>	 <p><b>Classification</b></p>
 <p><b>Code generation</b></p>	 <p><b>Speech-to-Text</b></p>

El aumento en la adopción de la automatización impulsado por la IA generativa también está cambiando el suministro y la demanda de las actividades del trabajo para muchas ocupaciones. Según McKinsey, el mercado laboral estadounidense (diagrama a continuación) ha pasado por una rápida transición, que solo puede continuar teniendo en cuenta el impacto de la IA.

Fuente: McKinsey & Company

## Estimated labor demand change and generative AI automation acceleration by occupation, US, 2022–30



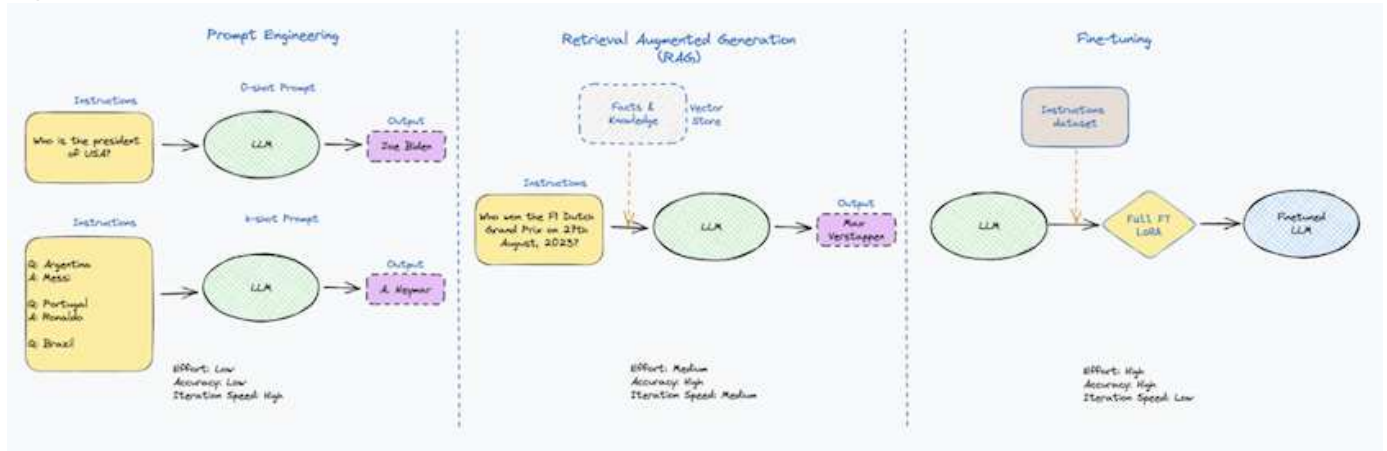
### Función del almacenamiento en IA generativa

Las LLMs dependen en gran medida del aprendizaje profundo, las GPU y la computación. Sin embargo, cuando el búfer de la GPU se llena, los datos deben escribirse rápidamente en el almacenamiento. Mientras que algunos modelos de IA son lo suficientemente pequeños para ejecutarse en la memoria, los LLM requieren un alto volumen de IOPS y un almacenamiento de alto rendimiento para proporcionar un acceso rápido a grandes conjuntos de datos, especialmente si se trata de miles de millones de tokens o millones de imágenes. Para un requisito de memoria de GPU típico de un LLM, la memoria necesaria para entrenar un modelo con 1 millones de parámetros podría llegar a 80GB @32-bit de precisión completa. En ese caso, la llama 2 de Meta, una familia de LLMs que varían en escala de 7 mil millones a 70 mil millones de parámetros, puede requerir 70x80, aprox 5600GB o 5,6TB GB de RAM de GPU. Además, la cantidad de memoria que necesita es directamente proporcional al número máximo de tokens que desea generar. Por ejemplo, si desea generar salidas de hasta 512 tokens (alrededor de 380 palabras), necesita "512MB". Puede parecer intrascendente, pero, si quieres correr lotes más grandes, comienza a sumar. Por lo tanto, lo hace muy caro para las organizaciones que entrenan o ajustan los LLM en la memoria, lo que convierte al almacenamiento en la piedra angular para la IA generativa.

### Tres enfoques principales para LLMs

Para la mayoría de las empresas, basado en las tendencias actuales, el enfoque para implementar LLMs se puede condensar en 3 escenarios básicos. Como se describe en un reciente ["Harvard Business Review"](#) artículo: (1) Capacitación (pre-entrenamiento) Un LLM desde cero – costoso y requiere habilidades expertas de IA/ML; (2) Ajuste de un modelo de base con datos empresariales – complejos, pero factibles; (3) Uso de la generación aumentada de recuperación (RAG) para consultar repositorios de documentos, API y bases de datos vectoriales que contienen datos de la empresa. Cada uno de ellos tiene desventajas entre el esfuerzo, la velocidad de iteración, la rentabilidad y la precisión de los modelos en sus implementaciones, que se utilizan para resolver diferentes tipos de problemas (diagrama a continuación).

Figura 3: Tipos de problemas



### Modelos de base

Un modelo de base (FM) también conocido como modelo base es un modelo de IA de gran tamaño (LLM) entrenado sobre grandes cantidades de datos sin etiquetar, utilizando la autosupervisión a escala y, por lo general, adaptado a una gran variedad de tareas posteriores del PLN. Dado que los datos de entrenamiento no son etiquetados por los humanos, el modelo emerge en lugar de ser explícitamente codificado. Esto significa que el modelo puede generar historias o una narrativa propia sin estar explícitamente programado para hacerlo. Por lo tanto, una característica importante de FM es la homogeneización, lo que significa que el mismo método se utiliza en muchos dominios. Sin embargo, con las técnicas de personalización y ajuste, los FMS integrados en los productos que aparecen en estos días no solo son buenos para generar texto, texto a imágenes y texto a código, sino también para explicar tareas específicas de dominio o depurar código. Por ejemplo, FMS como el Codex de OpenAI o el Code Llama de Meta pueden generar código en varios lenguajes de programación basados en descripciones de lenguaje natural de una tarea de programación. Estos modelos son competentes en más de una docena de lenguajes de programación, incluyendo Python, C#, JavaScript, Perl, Ruby, y SQL. Entienden la intención del usuario y generan código específico que logra la tarea deseada útil para el desarrollo de software, la optimización de código y la automatización de tareas de programación.

### Ajuste fino, especificidad de dominio y nueva formación

Una de las prácticas comunes con la implementación de LLM después de la preparación de datos y el procesamiento previo de los datos es seleccionar un modelo previamente entrenado que se haya entrenado en un conjunto de datos grande y diverso. En el contexto del ajuste fino, puede tratarse de un modelo de lenguaje grande de código abierto como ["Meta's Llama 2"](#) entrenado en 70 000 millones de parámetros y 2 billones de tokens. Una vez seleccionado el modelo preentrenado, el siguiente paso es ajustarlo en los datos específicos del dominio. Esto implica ajustar los parámetros del modelo y entrenarlo en los nuevos datos para adaptarse a un dominio y tarea específicos. Por ejemplo, BloombergGPT, un LLM propietario entrenado en una amplia gama de datos financieros que sirven a la industria financiera. Los modelos específicos de dominio diseñados y entrenados para una tarea específica generalmente tienen mayor precisión y rendimiento dentro de su alcance, pero baja transferibilidad entre otras tareas o dominios. Cuando el entorno empresarial y los

datos cambian durante un período, la precisión de predicción del FM podría comenzar a disminuir en comparación con su rendimiento durante las pruebas. Esto ocurre cuando resulta crucial volver a entrenar o ajustar el modelo. El reentrenamiento de modelos en IA/ML tradicional hace referencia a actualizar un modelo DE ML puesto en marcha con nuevos datos, por lo general se realizan para eliminar dos tipos de derivas que ocurren. (1) Derivación del concepto: Cuando el vínculo entre las variables de entrada y las variables de destino cambia con el tiempo, ya que la descripción de lo que queremos predecir los cambios, el modelo puede producir predicciones inexactas. (2) Derivación de datos: Se produce cuando las características de los datos de entrada cambian, como los cambios en los hábitos o el comportamiento del cliente a lo largo del tiempo y, por lo tanto, la incapacidad del modelo para responder a dichos cambios. De manera similar, la capacitación se aplica a FMS/LLMs, sin embargo, puede ser mucho más costosa (en millones de dólares), por lo tanto, no es algo que la mayoría de las organizaciones puedan considerar. Está bajo investigación activa, todavía emergiendo en el reino de LLMOps. Por lo tanto, en lugar de volver a entrenarse, cuando se produce la decadencia de modelos en FMS ajustados, las empresas pueden optar por el ajuste de nuevo (mucho más barato) con un conjunto de datos más nuevo. Para una perspectiva de costes, a continuación se muestra un ejemplo de una tabla de precios de modelo de Azure-OpenAI Services. Para cada categoría de tareas, los clientes pueden ajustar y evaluar modelos en conjuntos de datos específicos.

Fuente: Microsoft Azure

Model	Per 1000 token
Text-Ada	\$0.0001
GPT-3.5 Turbo	\$0.003
GPT-4	\$0.06
Text-Davinci	\$0.02
Model	Per 100 images
Dall-E	\$2

Ingeniería de pedidos e inferencia

La ingeniería rápida se refiere a los métodos efectivos de cómo comunicarse con LLMs para realizar las tareas deseadas sin actualizar los pesos del modelo. Tan importante como el entrenamiento y el ajuste de los modelos de IA es para las aplicaciones de PLN, la inferencia es igualmente importante, cuando los modelos entrenados responden a las indicaciones de los usuarios. Los requisitos del sistema de inferencia suelen estar mucho más en el rendimiento de lectura del sistema de almacenamiento de IA que alimenta datos desde los LLM a las GPU, ya que necesitan poder aplicar miles de millones de parámetros del modelo almacenados para producir la mejor respuesta.

LLMOps, Model Monitoring y Vectorstores

Al igual que las operaciones tradicionales de aprendizaje automático (MLOps), las operaciones de modelos de lenguaje grande (LLMOps) también requieren la colaboración de científicos de datos e ingenieros de DevOps con herramientas y prácticas recomendadas para la gestión de LLM en entornos de producción. Sin embargo, el flujo de trabajo y la pila de tecnología para LLMs podrían variar de alguna manera. Por ejemplo, LLM Pipelines creados usando marcos como LangChain unen varias llamadas de API LLM a puntos finales de incrustación externos, como vectorstores o bases de datos vectoriales. El uso de un punto final de incrustación y un almacén de vectores para conectores descendentes (como una base de datos vectorial) representa un desarrollo significativo en la forma en que se almacenan y se accede a los datos. A diferencia

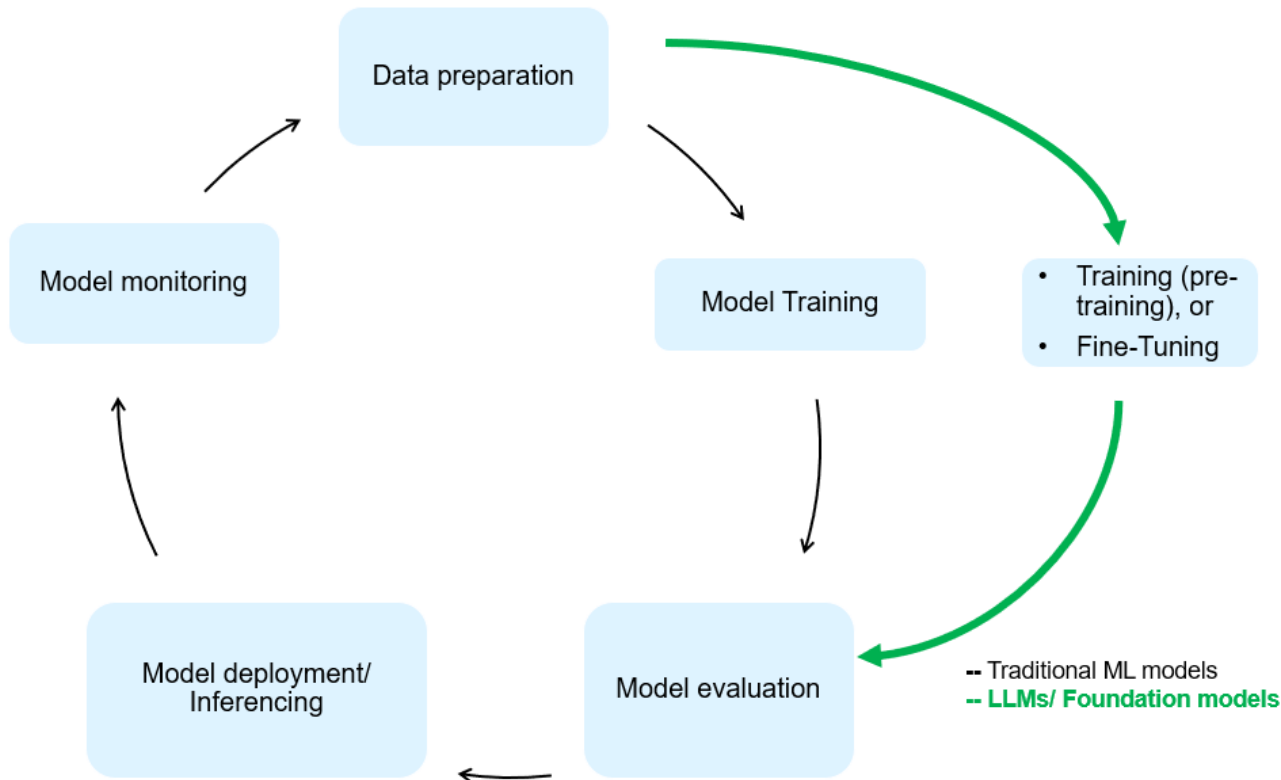
de los modelos tradicionales DE ML que se desarrollan desde cero, las LLM a menudo dependen del aprendizaje de transferencia, ya que estos modelos comienzan con FMS que están ajustados con nuevos datos para mejorar el rendimiento en un dominio más específico. Por lo tanto, es crucial que LLMOps ofrezca las capacidades de gestión de riesgos y monitoreo de deterioro de modelos.

Riesgos y Ética en la Era de la IA Generativa

“ChatGPT – Es elegante, pero todavía arroja tonterías.”– MIT Tech Review. La basura en-basura, siempre ha sido el caso difícil con la informática. La única diferencia con la IA generativa es que destaca en hacer que la basura sea altamente creíble, lo que lleva a resultados imprecisos. Las LLMs son propensas a inventar hechos para adaptarse a la narrativa que está construyendo. Por lo tanto, las empresas que ven la IA generativa como una gran oportunidad para reducir sus costes con equivalentes de IA necesitan detectar de manera eficiente falsificaciones profundas, reducir los sesgos y reducir los riesgos para mantener los sistemas honestos y éticos. Una canalización de datos con flujo libre con una sólida infraestructura de IA que admita la movilidad de datos, la calidad de los datos, la gobernanza y la protección de datos a través de cifrado integral y barreras de IA es eminente en el diseño de modelos de IA generativos responsables y explicables.

Situación de cliente y NetApp

Figura 3: Aprendizaje automático/flujo de trabajo de modelos de lenguaje grande



\*¿Estamos entrenando o ajustando? \* La pregunta de si (a) entrenar un modelo LLM desde cero, ajustar un FM pre-entrenado, o usar RAG para recuperar datos de repositorios de documentos fuera de un modelo de base y aumentar las indicaciones, y (b) ya sea mediante el aprovechamiento de LLM de código abierto (por ejemplo, Llama 2) o FMS propietario (por ejemplo, ChatGPT, Bard, AWS Bedrock) es una decisión estratégica para las organizaciones. Cada enfoque tiene una compensación entre la rentabilidad, la gravedad de los datos, las operaciones, la precisión de los modelos y la gestión de LLMs.

NetApp como empresa adopta la IA internamente en su cultura de trabajo y en su enfoque del diseño de productos y los esfuerzos de ingeniería. Por ejemplo, la protección autónoma frente al ransomware de NetApp se crea con la IA y el aprendizaje automático. Proporciona la detección temprana de las anomalías en el

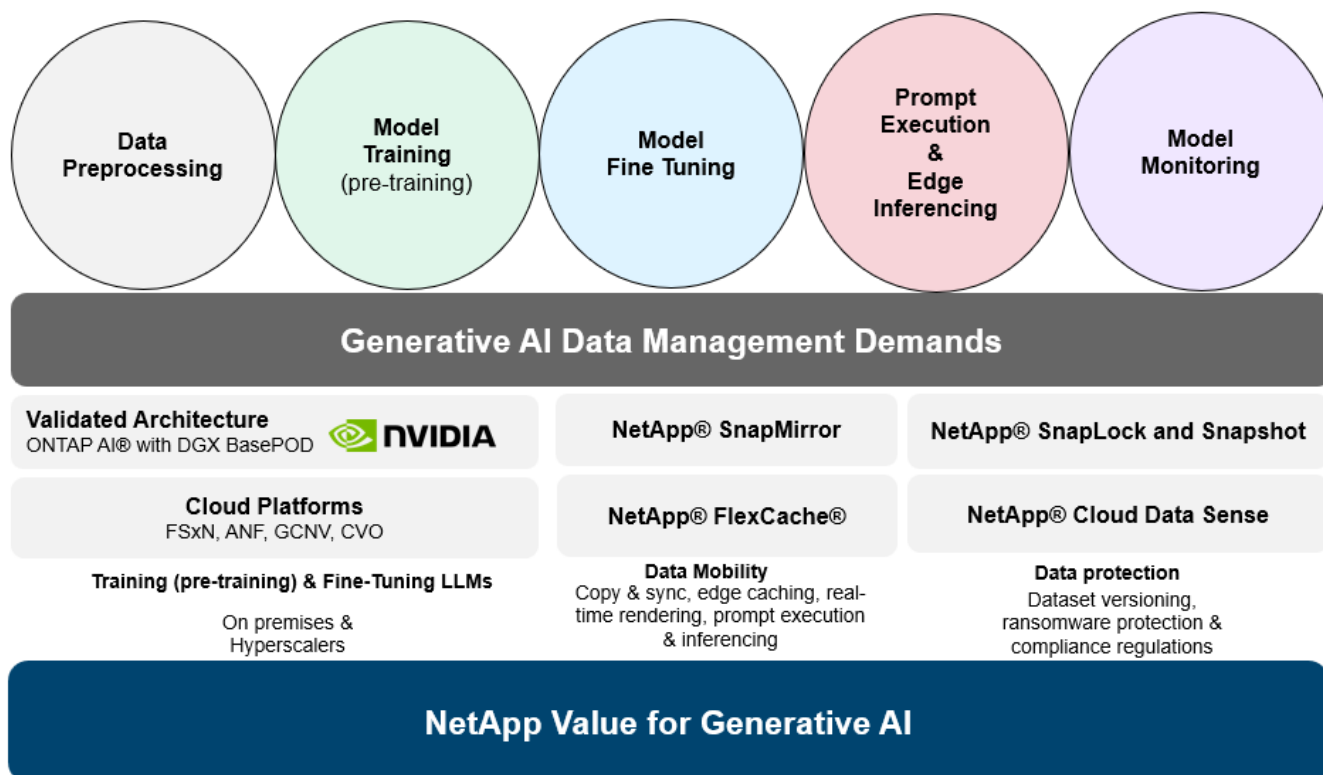


sistema de archivos para ayudar a identificar las amenazas antes de que afecten a las operaciones. En segundo lugar, NetApp utiliza la IA predictiva para sus operaciones empresariales, como la previsión de ventas e inventario y los bots conversacionales, para ayudar a los clientes en los servicios de soporte de productos de centro de llamadas, especificaciones técnicas, garantía y manuales de servicio, etc. En tercer lugar, NetApp aporta valor al cliente a la canalización de datos de IA y al flujo de trabajo DE ML/LLM a través de productos y soluciones que sirven a los clientes que crean soluciones de IA predictivas como previsión de demanda, imágenes médicas, análisis de sentimientos, y soluciones de IA generativas como Gans para la detección de anomalías en imágenes industriales en el sector de fabricación y la detección de fraudes y lavado de dinero en servicios bancarios y financieros con productos y capacidades de NetApp como NetApp® ONTAP AI®, NetApp® SnapMirror® y NetApp® FlexCache®.

## Funcionalidades de NetApp

El movimiento y la gestión de los datos en aplicaciones de IA generativa, como bot conversacional, generación de código, generación de imágenes o expresión de modelos genómicos, puede abarcar el perímetro, los centros de datos privados y el ecosistema multicloud híbrido. Por ejemplo, un robot de IA en tiempo real que ayuda a un pasajero a actualizar su boleto de avión a clase empresarial desde una aplicación de usuario final expuesta a través de API de modelos pre-entrenados como ChatGPT no puede lograr esa tarea por sí mismo, ya que la información del pasajero no está disponible públicamente en Internet. La API requiere acceso a la información personal del pasajero y a la información del boleto de la aerolínea que puede existir en un ecosistema híbrido o multicloud. Un escenario similar podría aplicarse a los científicos que comparten una molécula de fármaco y datos de pacientes a través de una aplicación de usuario final que utiliza LLMs para llevar a cabo ensayos clínicos en el descubrimiento de fármacos que involucran a instituciones de investigación biomédica de uno a varios. Los datos confidenciales que se pasan a FMS o LLMs pueden incluir información PII, información financiera, información de salud, datos biométricos, datos de ubicación, etc. datos de comunicaciones, comportamiento en línea e información legal. En tal caso de representación en tiempo real, ejecución rápida e inferencia periférica, hay movimiento de datos desde la aplicación de usuario final a extremos de almacenamiento a través de modelos LLM propios o de código abierto, hasta un centro de datos en las instalaciones o en plataformas de cloud público. En todos estos escenarios, la movilidad y la protección de datos son cruciales para las operaciones de IA que involucran LLM que dependen de grandes conjuntos de datos de entrenamiento y del movimiento de estos datos.

Figura 4: IA generativa: Canalización de datos de LLM



La cartera de infraestructura de almacenamiento, datos y servicios cloud de NetApp cuenta con la tecnología del software inteligente para la gestión de datos.

**Preparación de datos:** El primer pilar de la pila de tecnología LLM está en gran parte intacto de la pila de ML tradicional más antigua. El preprocesamiento de datos en la canalización de IA es necesario para normalizar y limpiar los datos antes del entrenamiento o ajuste. En este paso se incluyen conectores para ingerir datos dondequiera que estén en forma de un nivel Amazon S3 o en sistemas de almacenamiento en las instalaciones, como un almacén de archivos o de objetos como NetApp StorageGRID.

**NetApp® ONTAP** es la tecnología fundamental que sustenta las soluciones de almacenamiento críticas de NetApp en el centro de datos y la nube. ONTAP incluye varias funciones y funcionalidades de gestión y protección de datos, como la protección automática frente a ransomware contra ciberataques, funciones integradas de transporte de datos y funcionalidades de eficiencia del almacenamiento para una gama de arquitecturas en las instalaciones, híbridas y multiclouds en NAS, SAN, objetos, y las situaciones de almacenamiento definido por software (SDS) de las implementaciones de LLM.

**NetApp® ONTAP AI®** para el entrenamiento de modelos de aprendizaje profundo. NetApp® ONTAP® es compatible con NVIDIA GPU Direct Storage™ con el uso de NFS sobre RDMA para clientes de NetApp con el clúster de almacenamiento de ONTAP y nodos de computación NVIDIA DGX. Ofrece un rendimiento rentable para leer y procesar conjuntos de datos de origen desde el almacenamiento a la memoria numerosas veces para fomentar la inteligencia, lo que permite a las organizaciones formar, ajustar y escalar el acceso a las LLM.

**NetApp® FlexCache®** es una capacidad de almacenamiento en caché remoto que simplifica la distribución de archivos y almacena en caché solo los datos leídos activamente. Puede ser útil para la formación de LLM, reformaciones y ajustes precisos, lo que ofrece valor a los clientes con requisitos de negocio como la renderización en tiempo real e la inferencia de LLM.

**NetApp® SnapMirror** es una función de ONTAP que replica instantáneas de volumen entre dos sistemas ONTAP cualquiera. Esta función transfiere de forma óptima los datos que se encuentren en el perímetro a su centro de datos ubicado en las instalaciones o al cloud. SnapMirror puede utilizarse para mover datos de

forma segura y eficiente entre instalaciones y clouds de proveedores a hiperescala, cuando los clientes quieran desarrollar una IA generativa en clouds con RAG que contenga datos empresariales. Transfiere eficientemente solo los cambios, ahorrando ancho de banda y acelerando la replicación, ofreciendo así funciones esenciales de movilidad de datos durante las operaciones de entrenamiento, re-entrenamiento y ajuste de FMS o LLM.

**NetApp® SnapLock** aporta capacidad de disco inmutable en los sistemas de almacenamiento basados en ONTAP para el control de versiones de conjuntos de datos. La arquitectura microcore está diseñada para proteger los datos de los clientes con el motor FPolicy™ Zero Trust. NetApp garantiza que los datos de los clientes estén disponibles al resistir los ataques de denegación de servicio (DoS) cuando un atacante interactúa con un LLM de una forma que consume particularmente recursos.

**NetApp® Cloud Data Sense** ayuda a identificar, mapear y clasificar la información personal presente en los conjuntos de datos empresariales, promulgar políticas, cumplir con los requisitos de privacidad en las instalaciones o en la nube, ayudar a mejorar la postura de seguridad y cumplir con las regulaciones.

**Clasificación NetApp® BlueXP™**, con tecnología Cloud Data Sense. Los clientes pueden escanear, analizar, categorizar y actuar automáticamente sobre los datos en todo el conjunto de datos, detectar riesgos de seguridad, optimizar el almacenamiento y acelerar las puestas en marcha de cloud. Combina servicios de datos y almacenamiento a través de su plano de control unificado, los clientes pueden utilizar instancias de GPU para computación y entornos de multicloud híbridos para la organización en niveles de almacenamiento frío, así como archivos y backups.

**NetApp dualidad de archivo-objeto.** NetApp ONTAP permite el acceso de protocolo doble para NFS y S3. Con esta solución, los clientes pueden acceder a los datos NFS desde los portátiles de Amazon AWS SageMaker a través de bloques de S3 TB desde NetApp Cloud Volumes ONTAP. Esto ofrece flexibilidad a los clientes que necesitan un acceso fácil a fuentes de datos heterogéneas con la capacidad de compartir datos desde NFS y S3. Por ejemplo, afinar FMS como los modelos de generación de texto Llama 2 de Meta en SageMaker con acceso a cubos de objetos de archivo.

**El servicio NetApp® Cloud Sync** ofrece una forma sencilla y segura de migrar datos a cualquier destino, en la nube o en las instalaciones. Cloud Sync transfiere y sincroniza sin problemas los datos entre el almacenamiento en las instalaciones o en el cloud, NAS y almacenes de objetos.

**NetApp XCP** es un software cliente que permite migraciones de datos de cualquiera a NetApp y de NetApp a NetApp rápidas y fiables. XCP también proporciona la capacidad de mover datos masivos de forma eficiente desde sistema de archivos Hadoop HDFS a NFS de ONTAP, S3 o StorageGRID, y los análisis de archivos XCP proporcionan visibilidad sobre el sistema de archivos.

**NetApp® DataOps Toolkit** es una biblioteca de Python que facilita a los científicos de datos, DevOps e ingenieros de datos la realización de varias tareas de gestión de datos, como el aprovisionamiento, la clonación o la captura casi instantánea de un volumen de datos o un espacio de trabajo JupyterLab respaldado por un almacenamiento NetApp escalable de alto rendimiento.

**Seguridad de los productos de NetApp.** Los LLMs pueden revelar inadvertidamente datos confidenciales en sus respuestas, por lo tanto, una preocupación para los CISO que estudian las vulnerabilidades asociadas con las aplicaciones de IA que aprovechan los LLMs. Como se describe en OWASP (Open Worldwide Application Security Project), los problemas de seguridad como el envenenamiento de datos, la fuga de datos, la denegación de servicio y las inyecciones rápidas dentro de LLMs pueden afectar a las empresas de la exposición de los datos al acceso no autorizado que presta servicios a los atacantes. Los requisitos de almacenamiento de datos deben incluir comprobaciones de la integridad y snapshots inmutables de datos estructurados, semiestructurados y no estructurados. Las copias Snapshot de NetApp y SnapLock se utilizan para el control de versiones de conjuntos de datos. Proporciona un estricto control de acceso basado en roles (RBAC), así como protocolos seguros y cifrado estándar del sector para proteger los datos en reposo y en tránsito. Cloud Insights y Cloud Data Sense juntos ofrecen funcionalidades para ayudarte a identificar de

forma forense el origen de la amenaza y priorizar los datos que debes restaurar.

### **ONTAP AI con DGX BasePOD**

La arquitectura de referencia de IA NetApp® ONTAP® con NVIDIA DGX BasePOD es una arquitectura escalable para cargas de trabajo de aprendizaje automático (ML) e inteligencia artificial (IA). Durante la crucial fase de formación de los LLM, suelen copiarse datos a intervalos periódicos del almacenamiento de datos al clúster de formación. Los servidores empleados en esta fase utilizan las GPU para paralelizar los cálculos, lo que crea un tremendo apetito de datos. Satisfacer las necesidades brutas de ancho de banda E/S es esencial para poder mantener un uso de GPU elevado.

### **ONTAP AI con NVIDIA AI Enterprise**

NVIDIA AI Enterprise es una suite integral y nativa en el cloud de software de IA y análisis de datos optimizado, certificado y compatible con NVIDIA para ejecutarse en VMware vSphere con sistemas certificados por NVIDIA. Este software facilita la puesta en marcha, la gestión y el escalado simples y rápidos de las cargas de trabajo de IA en el entorno de cloud híbrido moderno. NVIDIA AI Enterprise, con la tecnología de NetApp y VMware, ofrece una excelente carga de trabajo de IA y gestión de datos en un paquete simplificado y conocido.

### **1P plataformas en la nube**

Las ofertas de almacenamiento en la nube totalmente gestionadas están disponibles de forma nativa en Microsoft Azure como Azure NetApp Files (ANF), en AWS como Amazon FSx para NetApp ONTAP (FSxN) y en Google como Google Cloud NetApp Volumes (GNCV). 1P es un sistema de archivos gestionado y de alto rendimiento que permite a los clientes ejecutar cargas de trabajo de IA de alta disponibilidad con seguridad de datos mejorada en nubes públicas, para ajustar LLMs/FMS con plataformas de ML nativas en la nube como AWS SageMaker, Azure-OpenAI Services y Vertex AI de Google.

### **Suite de soluciones para partners de NetApp**

Además de sus funciones, tecnologías y productos de datos básicos, NetApp también colabora estrechamente con una sólida red de partners de IA con el fin de ofrecer valor añadido a los clientes.

**Las Guardrails de NVIDIA** en los sistemas de IA sirven como salvaguardas para garantizar el uso ético y responsable de las tecnologías de IA. Los desarrolladores de IA pueden elegir definir el comportamiento de las aplicaciones impulsadas por LLM en temas específicos e impedir que participen en discusiones sobre temas no deseados. Guardrails, un kit de herramientas de código abierto, proporciona la capacidad de conectar un LLM a otros servicios, de forma fluida y segura para crear sistemas conversacionales LLM confiables, seguros y seguros.

**Domino Data Lab** proporciona herramientas versátiles de nivel empresarial para crear y producir IA generativa: Rápida, segura y económica, dondequiera que se encuentre en su viaje a IA. Con la plataforma Enterprise MLOps de Domino, los científicos de datos pueden utilizar las herramientas preferidas y todos sus datos, entrenar e implementar modelos fácilmente en cualquier lugar y administrar riesgos y de manera rentable, todo desde un único centro de control.

**Modzy para Edge AI.** NetApp® y Modzy se han asociado para ofrecer IA a escala a cualquier tipo de datos, incluidas imágenes, audio, texto y tablas. Modzy es una plataforma MLOps para implementar, integrar y ejecutar modelos de IA, ofrece a los científicos de datos las capacidades de supervisión de modelos, detección de deriva y explicabilidad, con una solución integrada para una inferencia LLM perfecta.

**Run:AI** y NetApp se han asociado para demostrar las capacidades únicas de la solución de IA de NetApp ONTAP con la plataforma de gestión de clústeres Run:AI para simplificar la orquestación de cargas de trabajo de IA. Divide y une automáticamente los recursos de GPU, diseñados para escalar sus canalizaciones de

procesamiento de datos a cientos de máquinas con marcos de integración integrados para Spark, Ray, DASK y Rapids.

## Conclusión

La IA generativa solo puede producir resultados efectivos cuando se entrena el modelo en remas de datos de calidad. Si bien las LLMs han logrado hitos notables, es fundamental reconocer sus limitaciones, desafíos de diseño y riesgos asociados con la movilidad de datos y la calidad de los datos. Las LLM dependen de conjuntos de datos de entrenamiento de gran tamaño y dispares procedentes de orígenes de datos heterogéneos. Los resultados inexactos o los resultados sesgados generados por los modelos pueden poner en peligro tanto a las empresas como a los consumidores. Estos riesgos pueden corresponder a las restricciones que surgen para las LLM potencialmente de los retos de gestión de datos asociados con la calidad, la seguridad de los datos y la movilidad de los mismos. NetApp ayuda a las organizaciones a cumplir las complejidades creadas por el rápido crecimiento de los datos, la movilidad de datos, la gestión multicloud y la adopción de IA. La infraestructura de IA a escala y la gestión de datos eficiente son cruciales para definir el éxito de las aplicaciones de IA como la IA generativa. Es crucial que los clientes cubran todas las situaciones de puesta en marcha sin comprometer su capacidad de expansión a medida que las empresas necesitan mantener el control de la rentabilidad, el gobierno de los datos y las prácticas éticas de IA. NetApp trabaja constantemente para ayudar a los clientes a simplificar y acelerar sus puestas en marcha de IA.

## TR-4785: Puesta en marcha de IA con E-Series de NetApp y BeeGFS

Nagalakshmi Raju, Daniel Landas, Nathan Swartz, Amine Bennani, NetApp

Las aplicaciones de inteligencia artificial (IA), aprendizaje automático (ML, por sus siglas en inglés) y aprendizaje profundo (DL, por sus siglas en inglés) implican conjuntos de datos de gran tamaño y cálculos elevados. Para ejecutar estas cargas de trabajo con éxito, necesita una infraestructura ágil que le permita escalar horizontalmente tanto los nodos de almacenamiento como los de computación sin problemas. Este informe incluye los pasos para poner en marcha un modelo de entrenamiento de IA en un modo distribuido, que permite un escalado horizontal fluido de los nodos de computación y almacenamiento. El informe también incluye varios indicadores de rendimiento para mostrar cómo una solución, que combina el almacenamiento E-Series de NetApp con el sistema de archivos en paralelo BeeGFS, proporciona una solución flexible, rentable y sencilla para cargas de trabajo de IA.

["TR-4785: Puesta en marcha de IA con E-Series de NetApp y BeeGFS"](#)

## Diseño de NVA-1150: Quantum StorNext con la guía de diseño de sistemas E-Series de NetApp

Ryan Rodine, NetApp

Este documento proporciona información detallada sobre cómo diseñar una solución de sistema de archivos paralelo StorNext con los sistemas de almacenamiento E-Series de NetApp. Esta solución abarca la cabina all-flash EF280 de NetApp, la cabina NVMe all-flash EF300 de NetApp, la cabina NVMe all-flash EF600 y el sistema híbrido E5760 de NetApp. Ofrece una caracterización del rendimiento basada en las pruebas comparativas de Frametest, una herramienta que se utiliza ampliamente para realizar pruebas en el sector del entretenimiento y los medios de comunicación.

## **NVA-1150-PUESTA en MARCHA: Quantum StorNext con la guía de puesta en marcha de sistemas E-Series de NetApp**

Ryan Rodine, NetApp

Este documento proporciona detalles sobre cómo poner en marcha una solución de sistema de archivos paralelo StorNext con sistemas de almacenamiento E-Series de NetApp. Esta solución abarca la cabina all-flash EF280 de NetApp, la cabina NVMe all-flash EF300 de NetApp, la cabina NVMe all-flash EF600 de NetApp y el sistema híbrido E5760 de NetApp. Ofrece una caracterización del rendimiento basada en las pruebas comparativas de Frametest, una herramienta que se utiliza ampliamente para realizar pruebas en el sector del entretenimiento y los medios de comunicación.

["NVA-1150-PUESTA en MARCHA: Quantum StorNext con la guía de puesta en marcha de sistemas E-Series de NetApp"](#)

## Información de copyright

Copyright © 2024 NetApp, Inc. Todos los derechos reservados. Imprimido en EE. UU. No se puede reproducir este documento protegido por copyright ni parte del mismo de ninguna forma ni por ningún medio (gráfico, electrónico o mecánico, incluidas fotocopias, grabaciones o almacenamiento en un sistema de recuperación electrónico) sin la autorización previa y por escrito del propietario del copyright.

El software derivado del material de NetApp con copyright está sujeto a la siguiente licencia y exención de responsabilidad:

ESTE SOFTWARE LO PROPORCIONA NETAPP «TAL CUAL» Y SIN NINGUNA GARANTÍA EXPRESA O IMPLÍCITA, INCLUYENDO, SIN LIMITAR, LAS GARANTÍAS IMPLÍCITAS DE COMERCIALIZACIÓN O IDONEIDAD PARA UN FIN CONCRETO, CUYA RESPONSABILIDAD QUEDA EXIMIDA POR EL PRESENTE DOCUMENTO. EN NINGÚN CASO NETAPP SERÁ RESPONSABLE DE NINGÚN DAÑO DIRECTO, INDIRECTO, ESPECIAL, EJEMPLAR O RESULTANTE (INCLUYENDO, ENTRE OTROS, LA OBTENCIÓN DE BIENES O SERVICIOS SUSTITUTIVOS, PÉRDIDA DE USO, DE DATOS O DE BENEFICIOS, O INTERRUPCIÓN DE LA ACTIVIDAD EMPRESARIAL) CUALQUIERA SEA EL MODO EN EL QUE SE PRODUJERON Y LA TEORÍA DE RESPONSABILIDAD QUE SE APLIQUE, YA SEA EN CONTRATO, RESPONSABILIDAD OBJETIVA O AGRAVIO (INCLUIDA LA NEGLIGENCIA U OTRO TIPO), QUE SURJAN DE ALGÚN MODO DEL USO DE ESTE SOFTWARE, INCLUSO SI HUBIEREN SIDO ADVERTIDOS DE LA POSIBILIDAD DE TALES DAÑOS.

NetApp se reserva el derecho de modificar cualquiera de los productos aquí descritos en cualquier momento y sin aviso previo. NetApp no asume ningún tipo de responsabilidad que surja del uso de los productos aquí descritos, excepto aquello expresamente acordado por escrito por parte de NetApp. El uso o adquisición de este producto no lleva implícita ninguna licencia con derechos de patente, de marcas comerciales o cualquier otro derecho de propiedad intelectual de NetApp.

Es posible que el producto que se describe en este manual esté protegido por una o más patentes de EE. UU., patentes extranjeras o solicitudes pendientes.

LEYENDA DE DERECHOS LIMITADOS: el uso, la copia o la divulgación por parte del gobierno están sujetos a las restricciones establecidas en el subpárrafo (b)(3) de los derechos de datos técnicos y productos no comerciales de DFARS 252.227-7013 (FEB de 2014) y FAR 52.227-19 (DIC de 2007).

Los datos aquí contenidos pertenecen a un producto comercial o servicio comercial (como se define en FAR 2.101) y son propiedad de NetApp, Inc. Todos los datos técnicos y el software informático de NetApp que se proporcionan en este Acuerdo tienen una naturaleza comercial y se han desarrollado exclusivamente con fondos privados. El Gobierno de EE. UU. tiene una licencia limitada, irrevocable, no exclusiva, no transferible, no sublicenciable y de alcance mundial para utilizar los Datos en relación con el contrato del Gobierno de los Estados Unidos bajo el cual se proporcionaron los Datos. Excepto que aquí se disponga lo contrario, los Datos no se pueden utilizar, desvelar, reproducir, modificar, interpretar o mostrar sin la previa aprobación por escrito de NetApp, Inc. Los derechos de licencia del Gobierno de los Estados Unidos de América y su Departamento de Defensa se limitan a los derechos identificados en la cláusula 252.227-7015(b) de la sección DFARS (FEB de 2014).

## Información de la marca comercial

NETAPP, el logotipo de NETAPP y las marcas que constan en <http://www.netapp.com/TM> son marcas comerciales de NetApp, Inc. El resto de nombres de empresa y de producto pueden ser marcas comerciales de sus respectivos propietarios.