



Dettagli sulla distribuzione e sulla convalida della soluzione

NetApp Solutions

NetApp
April 26, 2024

This PDF was generated from https://docs.netapp.com/it-it/netapp-solutions/ai/hcaios_ontap_ai_deployment.html on April 26, 2024. Always check docs.netapp.com for the latest.

Sommario

- Dettagli sulla distribuzione e sulla convalida della soluzione 1
 - Implementazione di ONTAP ai 1
 - Implementazione di Kubernetes 1
 - Implementazione di cnvrg.io 1

Dettagli sulla distribuzione e sulla convalida della soluzione

Le sezioni seguenti illustrano i dettagli dell'implementazione e della convalida della soluzione.

Implementazione di ONTAP ai

L'implementazione di ONTAP ai richiede l'installazione e la configurazione dell'hardware di rete, calcolo e storage. Le istruzioni specifiche per l'implementazione dell'infrastruttura ai di ONTAP esulano dall'ambito di questo documento. Per informazioni dettagliate sull'implementazione, vedere ["NVA-1121-DEPLOY: NetApp ONTAP ai, basato su NVIDIA"](#).

Per la convalida di questa soluzione, è stato creato un singolo volume e montato sul sistema DGX-1. Tale punto di montaggio è stato quindi montato sui container per rendere i dati accessibili per la formazione. Per implementazioni su larga scala, NetApp Trident automatizza la creazione e il montaggio dei volumi per eliminare i costi amministrativi e consentire la gestione delle risorse da parte dell'utente finale.

Implementazione di Kubernetes

Per implementare e configurare il cluster Kubernetes con NVIDIA DeepOps, eseguire le seguenti operazioni da un host di distribuzione jump:

1. Scaricare NVIDIA DeepOps seguendo le istruzioni sul ["Pagina introduttiva"](#) Sul sito NVIDIA DeepOps GitHub.
2. Implementare Kubernetes nel cluster seguendo le istruzioni sul ["Guida all'implementazione di Kubernetes"](#) Sul sito NVIDIA DeepOps GitHub.



Affinché l'implementazione di DeepOps Kubernetes funzioni, lo stesso utente deve esistere su tutti i nodi master e worker di Kubernetes.

Se l'implementazione non riesce, modificare il valore di `kubectl_localhost` a `false` in `deepops/config/group_vars/k8s-cluster.yml` e ripetere il punto 2. Il Copy `kubectl` binary to `ansible host` attività, che viene eseguita solo quando il valore di `kubectl_localhost` È vero, si basa sul modulo `fetch` Ansible, che presenta problemi noti di utilizzo della memoria. A volte, questi problemi di utilizzo della memoria possono causare un errore nell'attività. Se l'operazione non riesce a causa di un problema di memoria, il resto dell'operazione di implementazione non viene completata correttamente.

Se l'implementazione viene completata correttamente dopo aver modificato il valore di `kubectl_localhost` a `false`, quindi è necessario copiare manualmente `kubectl` binary Da un nodo master Kubernetes all'host di salto per l'implementazione. È possibile trovare la posizione di `kubectl` binary su un nodo master specifico eseguendo `which kubectl` comando direttamente su quel nodo.

Implementazione di cnvrg.io

Implementare il CORE cnvrg utilizzando Helm

Helm è il modo più semplice per implementare rapidamente cnvrg utilizzando qualsiasi cluster, on-premise, Minikube o qualsiasi cluster cloud (come AKS, EKS e GKE). Questa sezione descrive come cnvrg è stato installato su un'istanza on-premise (DGX-1) con Kubernetes installato.

Prerequisiti

Prima di completare l'installazione, è necessario installare e preparare le seguenti dipendenze sul computer locale:

- Kubectl
- Timone 3.x
- Kubernetes cluster 1.15+

Implementazione con Helm

1. Per scaricare i grafici di comando più aggiornati, eseguire il seguente comando:

```
helm repo add cnvrg https://helm.cnvrg.io
helm repo update
```

2. Prima di implementare cnvrg, è necessario disporre dell'indirizzo IP esterno del cluster e del nome del nodo su cui verrà implementato cnvrg. Per implementare cnvrg in un cluster Kubernetes on-premise, eseguire il seguente comando:

```
helm install cnvrg cnvrg/cnvrg --timeout 1500s --wait \ --set
global.external_ip=<ip_of_cluster> \ --set global.node=<name_of_node>
```

3. Eseguire `helm install` comando. Tutti i servizi e i sistemi vengono installati automaticamente sul cluster. Il processo può richiedere fino a 15 minuti.
4. Il `helm install` il comando può richiedere fino a 10 minuti. Una volta completata l'implementazione, accedere all'URL del cnvrg appena distribuito o aggiungere il nuovo cluster come risorsa all'interno dell'organizzazione. Il `helm` il comando indica l'URL corretto.

```
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. Quando lo stato di tutti i container è in esecuzione o completo, cnvrg è stato implementato correttamente. Dovrebbe essere simile al seguente output di esempio:

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-kcbt9		1/1	Running	0	2m
seeder-wh655		0/1	Completed	0	2m
speaker-5sqhr		1/1	Running	0	2m

Formazione sul modello di visione artificiale con ResNet50 e il set di dati radiologici Chest

Il sistema operativo ai cnvrg.io è stato implementato su una configurazione Kubernetes su un'architettura NetApp ONTAP ai basata sul sistema NVIDIA DGX. Per la convalida, abbiamo utilizzato il set di dati radiologici NIH Chest costituito da immagini anonimizzate dei raggi X del torace. Le immagini erano in formato PNG. I dati sono stati forniti dal NIH Clinical Center e sono disponibili tramite ["Sito di download NIH"](#). Abbiamo utilizzato un campione di 250 GB dei dati con 627, 615 immagini in 15 classi.

Il set di dati è stato caricato sulla piattaforma cnvrg ed è stato memorizzato nella cache di un'esportazione NFS dal sistema di storage NetApp AFF A800.

Impostare le risorse di calcolo

L'architettura cnvrg e la funzionalità di meta-scheduling consentono a tecnici e professionisti IT di collegare diverse risorse di calcolo a una singola piattaforma. Nella nostra configurazione, abbiamo utilizzato lo stesso cluster cnvrg implementato per l'esecuzione dei carichi di lavoro di deep-learning. Se è necessario collegare altri cluster, utilizzare la GUI, come mostrato nella seguente schermata.



Caricare i dati

Per caricare i dati sulla piattaforma cnvrg, è possibile utilizzare la GUI o la CLI cnvrg. Per i set di dati di grandi dimensioni, NetApp consiglia di utilizzare CLI perché si tratta di uno strumento potente, scalabile e affidabile in grado di gestire un gran numero di file.

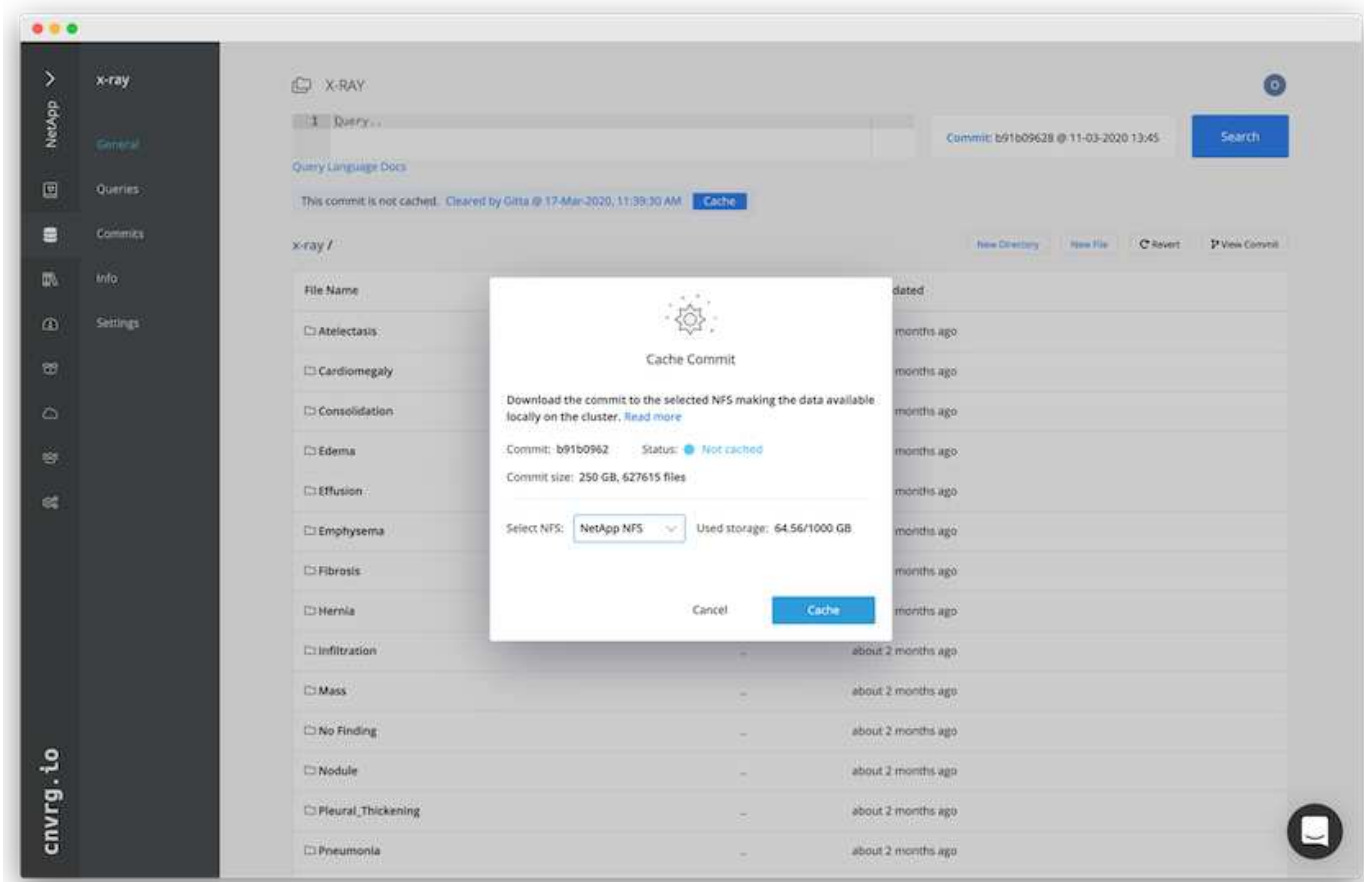
Per caricare i dati, attenersi alla seguente procedura:

1. Scaricare il "CLI cnvrg".
2. accedere alla directory dei raggi x.
3. Inizializzare il set di dati nella piattaforma con `cnvrg data init` comando.
4. Caricare tutti i contenuti della directory nel data Lake centrale con `cnvrg data sync` Command. una volta caricati i dati nell'archivio centrale di oggetti (StorageGRID, S3 o altri), è possibile navigare con la GUI. La figura seguente mostra un file PNG di immagine della fibrosi a raggi X del torace caricato. Inoltre, cnvrg consente di eseguire la versione dei dati in modo che qualsiasi modello creato possa essere riprodotto fino alla versione dei dati.



Dati di cach

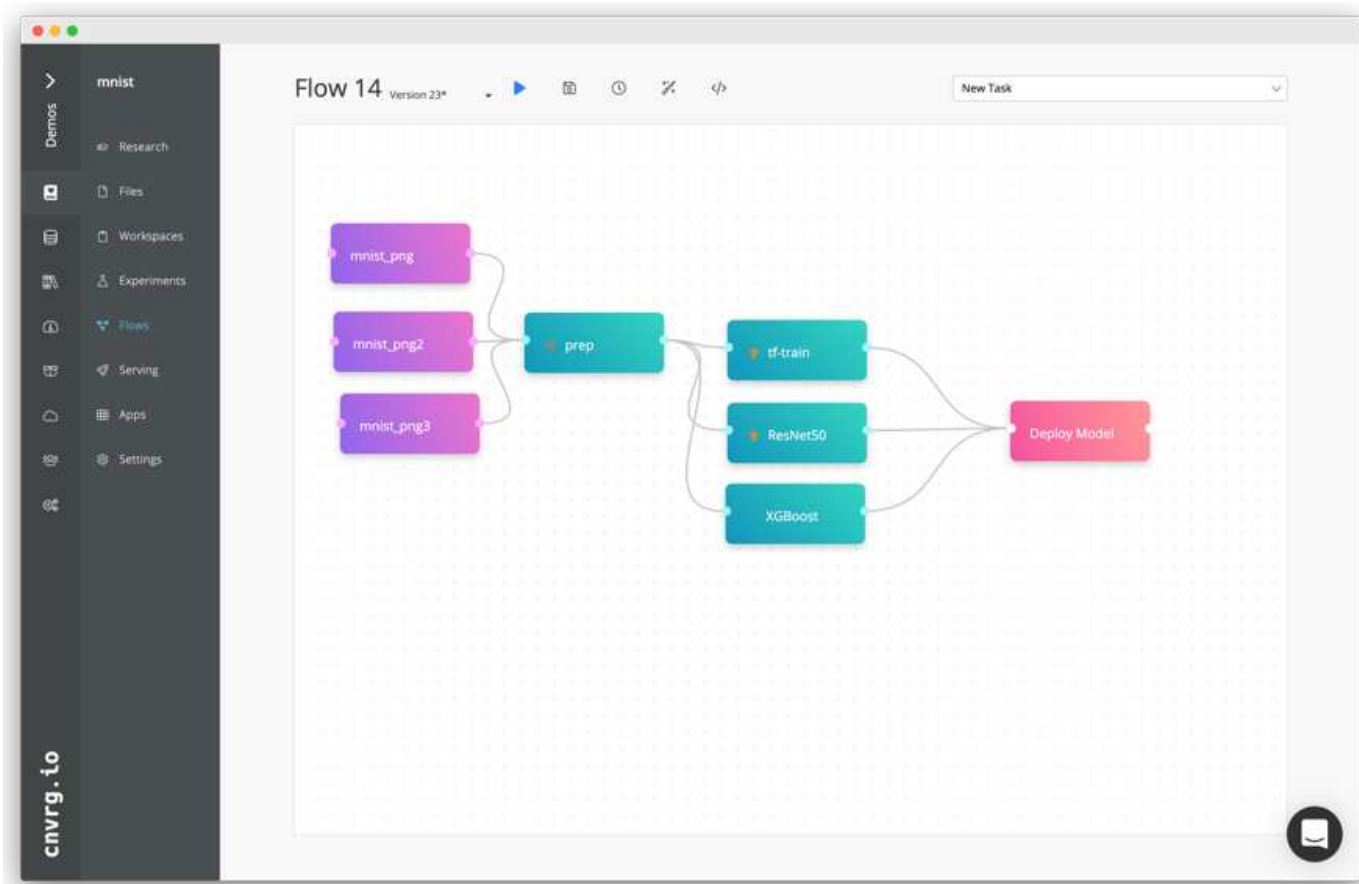
Per accelerare il training ed evitare il download di oltre 600.000 file per ciascun modello di training ed esperimento, abbiamo utilizzato la funzionalità di caching dei dati dopo che i dati sono stati inizialmente caricati nell'archivio centrale di oggetti data-Lake.



Dopo che gli utenti hanno fatto clic su cache, cnvrg scarica i dati nel relativo commit specifico dall'archivio remoto di oggetti e li memorizza nella cache del volume NFS di ONTAP. Al termine, i dati saranno disponibili per il training istantaneo. Inoltre, se i dati non vengono utilizzati per alcuni giorni (ad esempio, per la formazione o l'esplorazione del modello), cnvrg cancella automaticamente la cache.

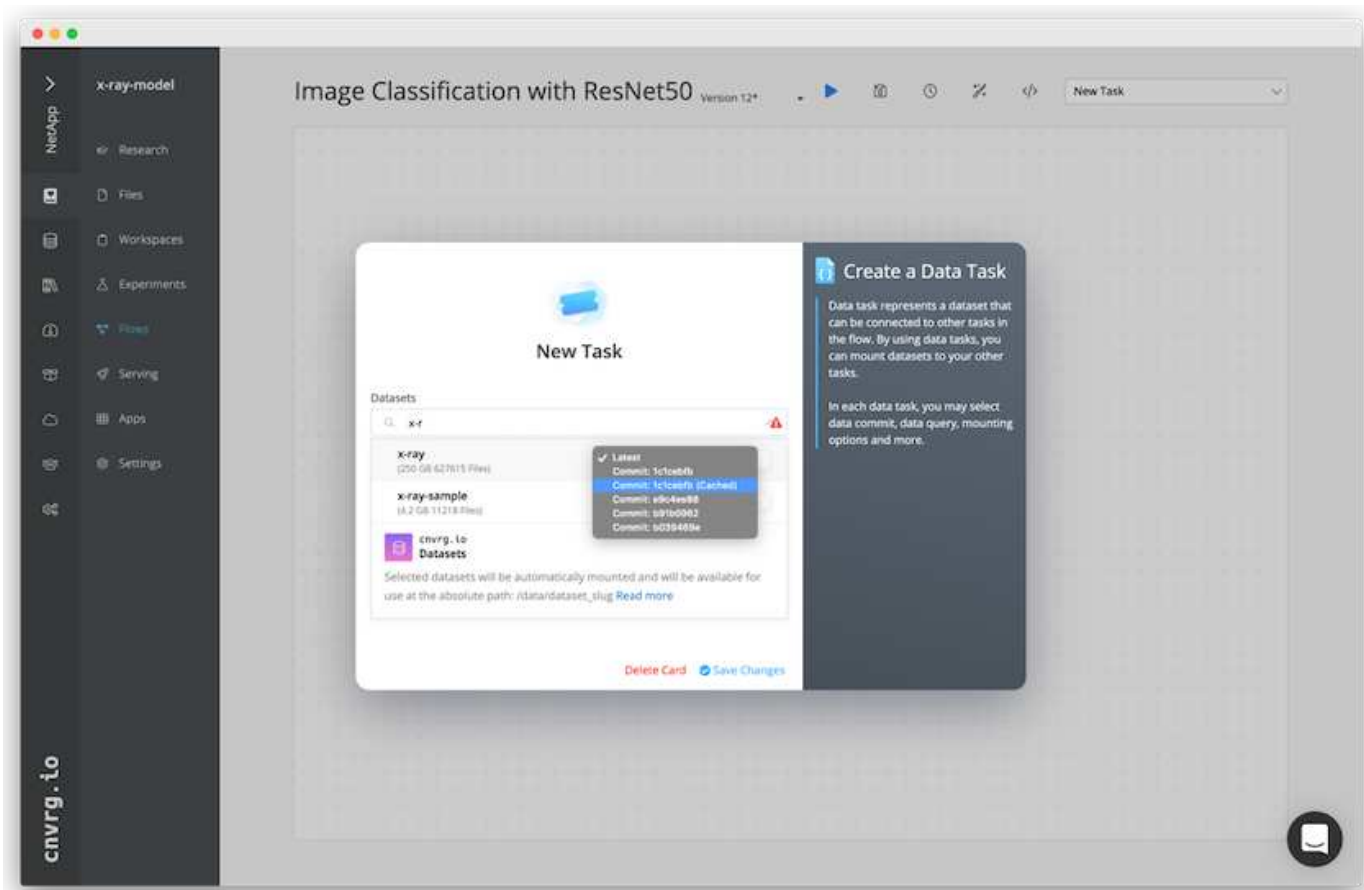
Crea una pipeline ML con i dati memorizzati nella cache

Cnvrg Flows consente di creare facilmente pipeline ML di produzione. I flussi sono flessibili, possono funzionare per qualsiasi tipo di caso d'utilizzo DI ML e possono essere creati attraverso la GUI o il codice. Ogni componente di un flusso può essere eseguito su una diversa risorsa di calcolo con un'immagine Docker diversa, il che rende possibile la creazione di cloud ibrido e pipeline ML ottimizzate.



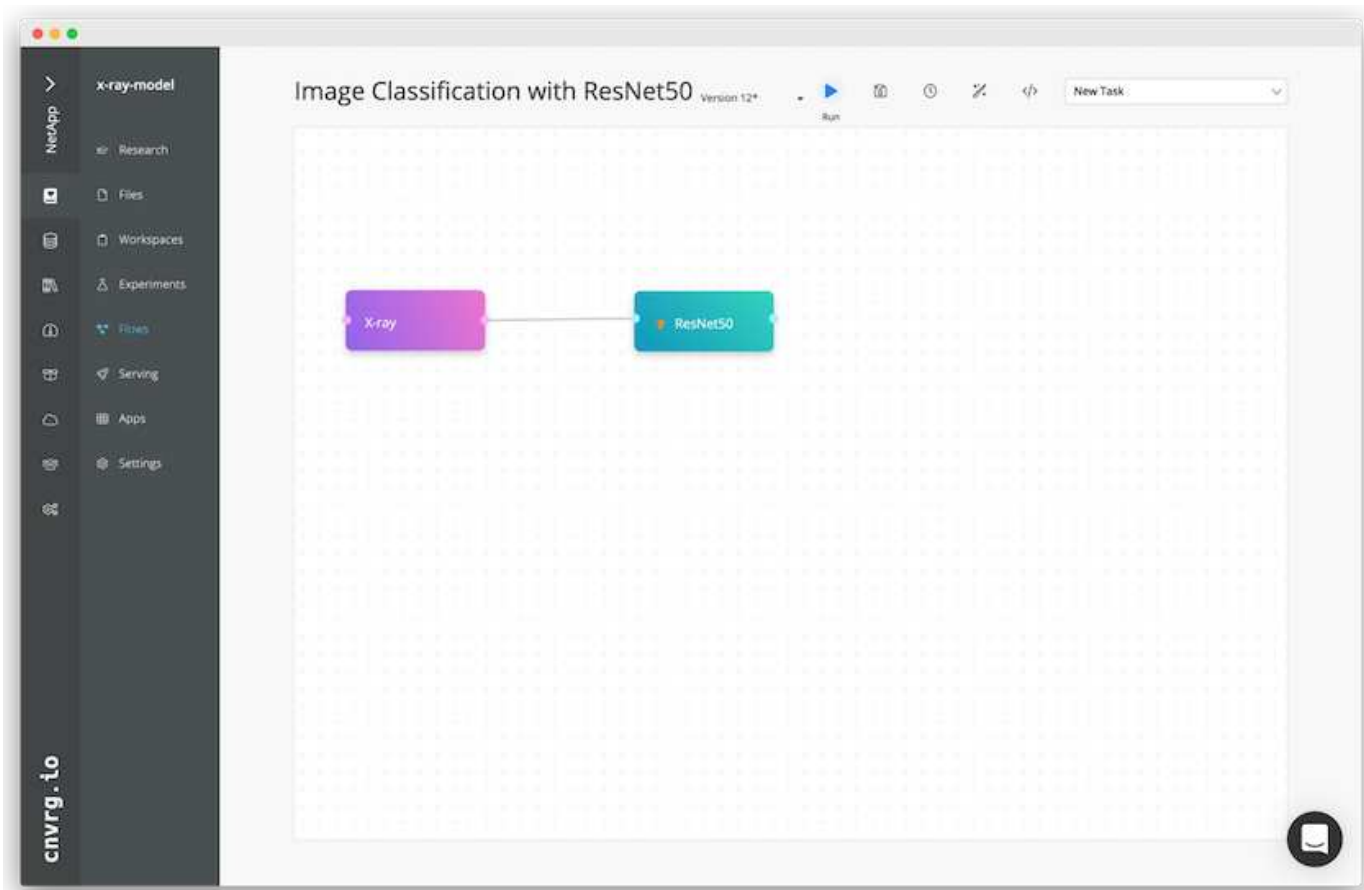
Creazione del flusso di raggi X del torace: Impostazione dei dati

Abbiamo aggiunto il nostro set di dati a un flusso appena creato. Quando si aggiunge il dataset, è possibile selezionare la versione specifica (commit) e indicare se si desidera la versione memorizzata nella cache. In questo esempio, è stato selezionato il commit memorizzato nella cache.



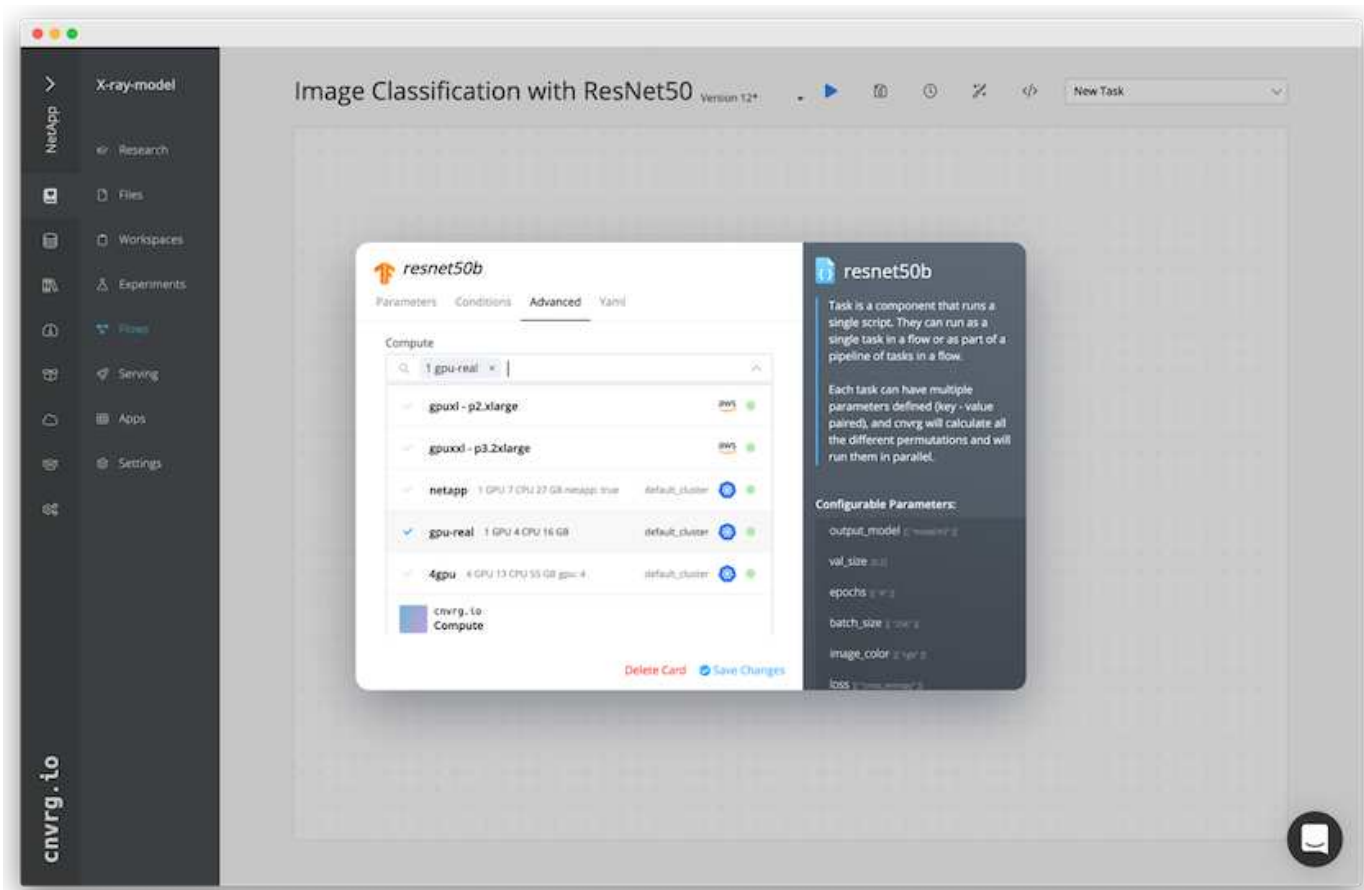
Creazione del flusso di raggi X del torace: Impostazione del modello di training: ResNet50

Nella pipeline, è possibile aggiungere qualsiasi tipo di codice personalizzato desiderato. In cnvrg è disponibile anche la libreria ai, una raccolta di componenti ML riutilizzabili. Nella libreria ai sono presenti algoritmi, script, origini dati e altre soluzioni che possono essere utilizzate in qualsiasi ML o flusso di deep learning. In questo esempio, è stato selezionato il modulo ResNet50 preinstallato. Abbiamo utilizzato parametri predefiniti come batch_size:128, epochs:10 e molto altro ancora. Questi parametri possono essere visualizzati nei documenti della ai Library. La seguente schermata mostra il nuovo flusso con il set di dati radiologici collegato a ResNet50.



Definire la risorsa di calcolo per ResNet50

Ogni algoritmo o componente nei flussi cnvrg può essere eseguito su un'istanza di calcolo diversa, con un'immagine Docker diversa. Nella nostra configurazione, volevamo eseguire l'algoritmo di training sui sistemi NVIDIA DGX con l'architettura NetApp ONTAP ai. Nella figura seguente, è stato selezionato `gpu-real`, che è un modello di calcolo e una specifica per il nostro cluster on-premise. Abbiamo anche creato una coda di modelli e selezionato più modelli. In questo modo, se il `gpu-real` non è possibile allocare le risorse (se, ad esempio, altri data scientist le stanno utilizzando), quindi è possibile attivare la diffusione automatica del cloud aggiungendo un modello di cloud provider. La seguente schermata mostra l'utilizzo di `gpu-real` come nodo di calcolo per ResNet50.



Monitoraggio e monitoraggio dei risultati

Una volta eseguito un flusso, cnvrg attiva il motore di monitoraggio e tracciamento. Ogni esecuzione di un flusso viene documentata e aggiornata automaticamente in tempo reale. Hyperparameters, metriche, utilizzo delle risorse (utilizzo della GPU e altro ancora), versione del codice, artefatti, log, E così via sono disponibili automaticamente nella sezione Experiments (esperimenti), come mostrato nelle due schermate seguenti.



Informazioni sul copyright

Copyright © 2024 NetApp, Inc. Tutti i diritti riservati. Stampato negli Stati Uniti d'America. Nessuna porzione di questo documento soggetta a copyright può essere riprodotta in qualsiasi formato o mezzo (grafico, elettronico o meccanico, inclusi fotocopie, registrazione, nastri o storage in un sistema elettronico) senza previo consenso scritto da parte del detentore del copyright.

Il software derivato dal materiale sottoposto a copyright di NetApp è soggetto alla seguente licenza e dichiarazione di non responsabilità:

IL PRESENTE SOFTWARE VIENE FORNITO DA NETAPP "COSÌ COM'È" E SENZA QUALSIVOGLIA TIPO DI GARANZIA IMPLICITA O ESPRESSA FRA CUI, A TITOLO ESEMPLIFICATIVO E NON ESAUSTIVO, GARANZIE IMPLICITE DI COMMERCIALIZZABILITÀ E IDONEITÀ PER UNO SCOPO SPECIFICO, CHE VENGONO DECLINATE DAL PRESENTE DOCUMENTO. NETAPP NON VERRÀ CONSIDERATA RESPONSABILE IN ALCUN CASO PER QUALSIVOGLIA DANNO DIRETTO, INDIRETTO, ACCIDENTALE, SPECIALE, ESEMPLARE E CONSEGUENZIALE (COMPRESI, A TITOLO ESEMPLIFICATIVO E NON ESAUSTIVO, PROCUREMENT O SOSTITUZIONE DI MERCI O SERVIZI, IMPOSSIBILITÀ DI UTILIZZO O PERDITA DI DATI O PROFITTI OPPURE INTERRUZIONE DELL'ATTIVITÀ AZIENDALE) CAUSATO IN QUALSIVOGLIA MODO O IN RELAZIONE A QUALUNQUE TEORIA DI RESPONSABILITÀ, SIA ESSA CONTRATTUALE, RIGOROSA O DOVUTA A INSOLVENZA (COMPRESA LA NEGLIGENZA O ALTRO) INSORTA IN QUALSIASI MODO ATTRAVERSO L'UTILIZZO DEL PRESENTE SOFTWARE ANCHE IN PRESENZA DI UN PREAVVISO CIRCA L'EVENTUALITÀ DI QUESTO TIPO DI DANNI.

NetApp si riserva il diritto di modificare in qualsiasi momento qualunque prodotto descritto nel presente documento senza fornire alcun preavviso. NetApp non si assume alcuna responsabilità circa l'utilizzo dei prodotti o materiali descritti nel presente documento, con l'eccezione di quanto concordato espressamente e per iscritto da NetApp. L'utilizzo o l'acquisto del presente prodotto non comporta il rilascio di una licenza nell'ambito di un qualche diritto di brevetto, marchio commerciale o altro diritto di proprietà intellettuale di NetApp.

Il prodotto descritto in questa guida può essere protetto da uno o più brevetti degli Stati Uniti, esteri o in attesa di approvazione.

LEGENDA PER I DIRITTI SOTTOPOSTI A LIMITAZIONE: l'utilizzo, la duplicazione o la divulgazione da parte degli enti governativi sono soggetti alle limitazioni indicate nel sottoparagrafo (b)(3) della clausola Rights in Technical Data and Computer Software del DFARS 252.227-7013 (FEB 2014) e FAR 52.227-19 (DIC 2007).

I dati contenuti nel presente documento riguardano un articolo commerciale (secondo la definizione data in FAR 2.101) e sono di proprietà di NetApp, Inc. Tutti i dati tecnici e il software NetApp forniti secondo i termini del presente Contratto sono articoli aventi natura commerciale, sviluppati con finanziamenti esclusivamente privati. Il governo statunitense ha una licenza irrevocabile limitata, non esclusiva, non trasferibile, non cedibile, mondiale, per l'utilizzo dei Dati esclusivamente in connessione con e a supporto di un contratto governativo statunitense in base al quale i Dati sono distribuiti. Con la sola esclusione di quanto indicato nel presente documento, i Dati non possono essere utilizzati, divulgati, riprodotti, modificati, visualizzati o mostrati senza la previa approvazione scritta di NetApp, Inc. I diritti di licenza del governo degli Stati Uniti per il Dipartimento della Difesa sono limitati ai diritti identificati nella clausola DFARS 252.227-7015(b) (FEB 2014).

Informazioni sul marchio commerciale

NETAPP, il logo NETAPP e i marchi elencati alla pagina <http://www.netapp.com/TM> sono marchi di NetApp, Inc. Gli altri nomi di aziende e prodotti potrebbero essere marchi dei rispettivi proprietari.