



## 解決策の導入と検証の詳細 NetApp Solutions

NetApp  
April 10, 2024

This PDF was generated from [https://docs.netapp.com/ja-jp/netapp-solutions/ai/hcaios\\_ontap\\_ai\\_deployment.html](https://docs.netapp.com/ja-jp/netapp-solutions/ai/hcaios_ontap_ai_deployment.html) on April 10, 2024. Always check docs.netapp.com for the latest.

# 目次

解決策の導入と検証の詳細 .....	1
ONTAP AI の導入 .....	1
Kubernetes の導入 .....	1
cnvrg.io の展開 .....	1

# 解決策の導入と検証の詳細

以降のセクションでは、解決策の導入と検証の詳細について説明します。

## ONTAP AI の導入

ONTAP AI を導入するには、ネットワーク、コンピューティング、ストレージハードウェアのインストールと設定が必要です。ONTAP AI インフラの導入手順については、本ドキュメントでは説明していません。導入の詳細については、["NVA-1121-deploy : NetApp ONTAP AI、Powered by NVIDIA"](#)を参照してください。

この解決策検証では、1つのボリュームを作成して DGX-1 システムにマウントしました。その後、そのマウントポイントをコンテナにマウントして、トレーニング用のデータにアクセスできるようにしました。大規模な環境では、NetApp Trident によってボリュームの作成とマウントが自動化されるため、管理上のオーバーヘッドが発生しないとともに、エンドユーザによるリソースの管理が可能になります。

## Kubernetes の導入

NVIDIA DeepOps で Kubernetes クラスタを導入および設定するには、導入ジャンプホストから次のタスクを実行します。

1. の手順に従って、NVIDIA DeepOps をダウンロードします "[「はじめに」 ページ](#)" NVIDIA DeepOps GitHub サイトで入手できます。
2. の手順に従って、クラスタに Kubernetes を導入します。 "[Kubernetes 導入ガイド](#)" NVIDIA DeepOps GitHub サイトで入手できます。



DeepOps Kubernetes 環境を使用するには、Kubernetes マスターノードとワーカーノードがすべて同じユーザである必要があります。

配備に失敗した場合は 'kubectl\_localhost' の値を 'deepops/config/group\_vars/k8s-cluster.yml' で false に変更し、手順 2 を繰り返します。Copy kubectl binary to Ansible host タスクは 'kubectl\_localhost' の値が true の場合にのみ実行され、メモリ使用に関する既知の問題がある FETCH Ansible モジュールに依存します。このようなメモリ使用の問題により、原因がタスクを失敗させることがあります。メモリ問題が原因でタスクが失敗した場合は、以降の導入処理は正常に完了しません。

「kubectl\_localhost」の値を「false」に変更した後で展開が正常に完了した場合、「kubectl binary」を Kubernetes マスターノードから配備ジャンプホストに手動でコピーする必要があります。特定のマスター・ノード上で 'kubectl binary' の場所を確認するには 'which kubectl' コマンドをそのノード上で直接実行します。

## cnvrg.io の展開

### Helm を使用して cnvrg コアを導入します

Helm は、任意のクラスタ、オンプレミス、Minikube、または任意のクラウドクラスタ（AKS、EKS、GKE など）を使用して、cnvrg を迅速に導入する最も簡単な方法です。このセクションでは、Kubernetes がインストールされたオンプレミス（DGX-1）インスタンスに cnvrg がインストールされた方法について説明します。

## 前提条件

インストールを完了する前に、ローカルマシンに次の依存関係をインストールして準備する必要があります。

- Kubectl のように入力する
- Helm 3.x
- Kubernetes クラスタ 1.15 以降

## Helm を使用して展開します

1. 最新の cnvrg Helm チャートをダウンロードするには、次のコマンドを実行します。

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

2. cnvrg を導入する前に、クラスタの外部 IP アドレス、および cnvrg を導入するノードの名前が必要です。オンプレミスの Kubernetes クラスタに cnvrg を導入するには、次のコマンドを実行します。

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

3. 「helm install」コマンドを実行します。すべてのサービスとシステムがクラスタに自動的にインストールされます。この処理には最大 15 分かかることがあります。
4. 「helm install」コマンドの所要時間は最大 10 分です。展開が完了したら、新しく展開した cnvrg の URL に移動するか、新しいクラスタを組織内のリソースとして追加します。「helm」コマンドは正しい URL を通知します。

```
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. すべてのコンテナのステータスが「Running」または「Complete」の場合、cnvrg は正常に展開されています。次のような出力が表示されます。

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

## ResNet50 および胸部 X 線を使用したコンピュータビジョンモデルトレーニング データセット

NVIDIA DGX システムを基盤とする NetApp ONTAP AI アーキテクチャ上の Kubernetes セットアップに、cnvrg.io AI OS が導入されました。検証には、胸部 X 線の匿名画像からなる NIH 胸部 X 線データセットを使用しました。画像は PNG 形式でした。このデータは NIH クリニカルセンタおよびによって提供された は、から使用できます ["NIH ダウンロードサイト"](#)。250 GB のサンプルデータを 15 クラスの 627、615 イメージで使用しました。

データセットは cnvrg プラットフォームにアップロードされ、NetApp AFF A800 ストレージシステムからの NFS エクスポートにキャッシュされました。

## コンピューティングリソースをセットアップする

cnvrg アーキテクチャおよびメタスケジューリング機能により、エンジニアおよび IT プロフェッショナルは、異なるコンピューティングリソースを 1 つのプラットフォームに接続できます。今回のセットアップでは、ディープラーニングワークロードの実行用に導入されたクラスター cnvrg を使用しました。追加のクラスターを接続する必要がある場合は、次のスクリーンショットに示すように、GUI を使用してください。

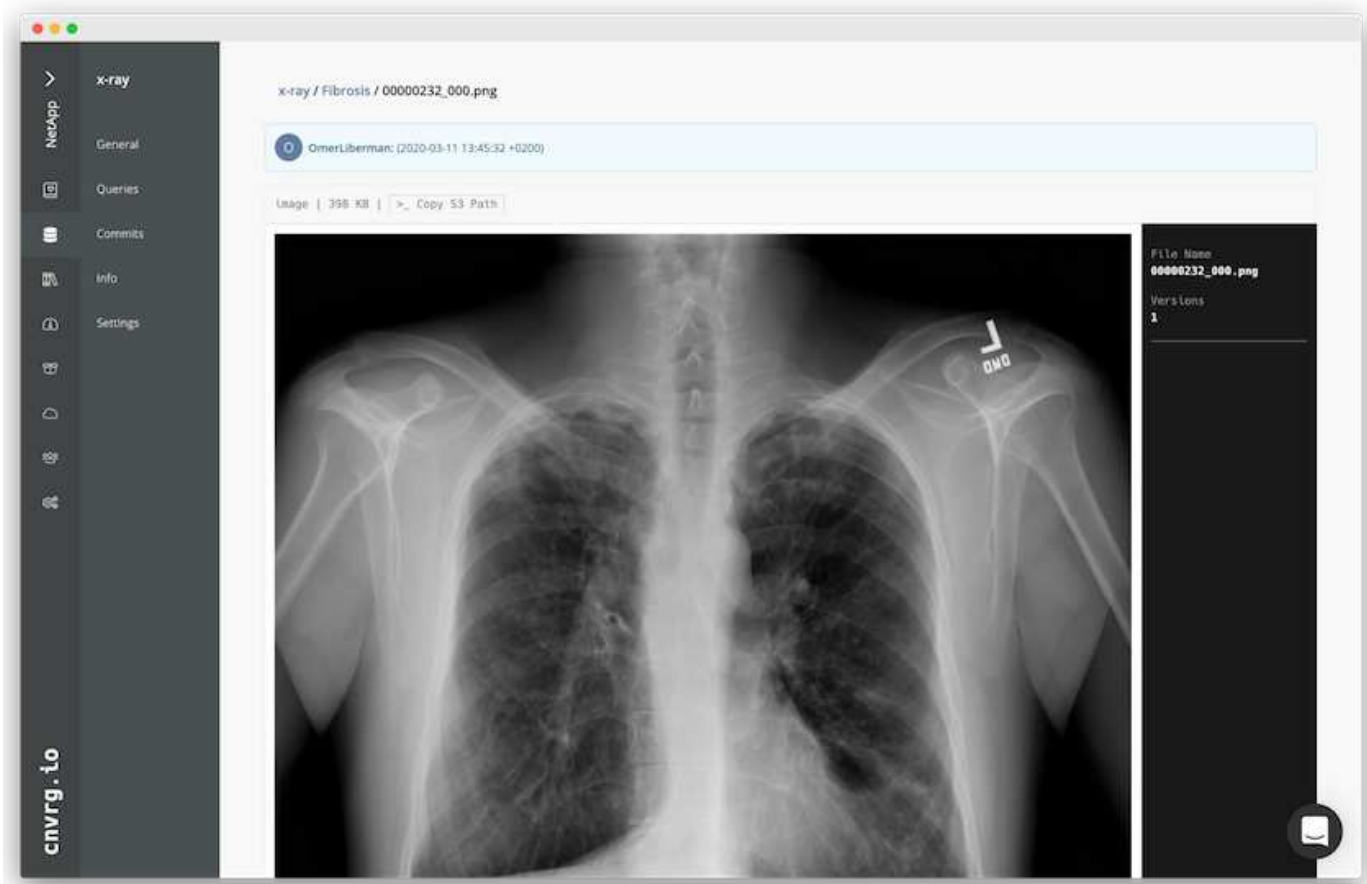


## データをロードします

cnvrg プラットフォームにデータをアップロードするには、GUI または cnvrg CLI を使用します。大規模なデータセットの場合は、CLI の使用を推奨します。CLI は、多数のファイルを処理できる、拡張性と信頼性に優れた強力なツールです。

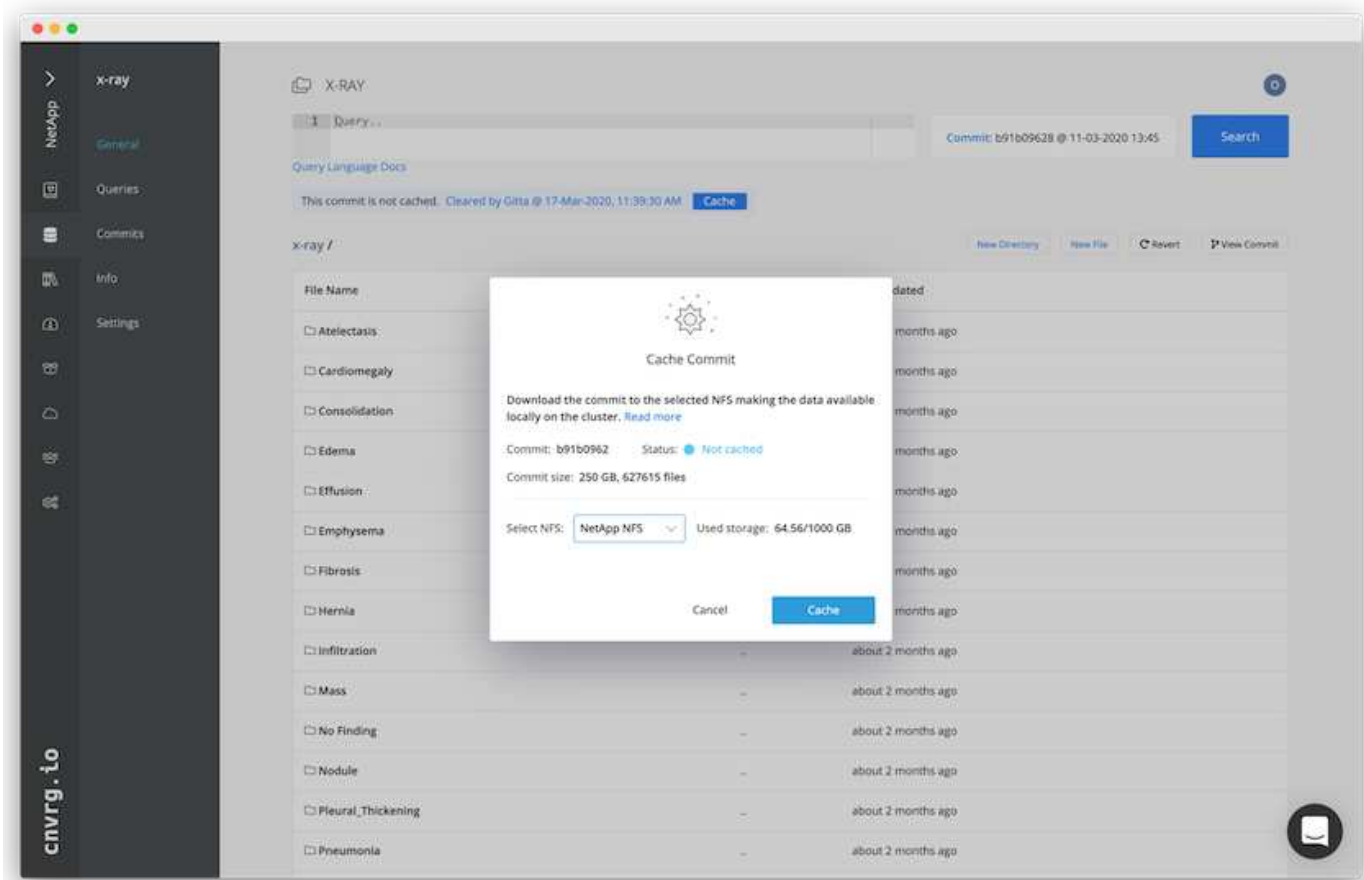
データをアップロードするには、次の手順を実行します。

1. をダウンロードします **"cnvrg CLI"**。
2. X 線ディレクトリに移動します。
3. 「cnvrg data init」コマンドを使用して、プラットフォーム内のデータセットを初期化します。
4. 「cnvrg data sync」コマンドを使用して、ディレクトリのすべての内容を中央のデータレイクにアップロードします。データが中央のオブジェクトストア（StorageGRID、S3、またはその他）にアップロードされたら、GUI で参照できます。次の図は、ロードされた胸部 X 線線維症画像 PNG ファイルを示しています。さらに、cnvrg は、ビルドしたすべてのモデルをデータバージョンに複製できるように、データをバージョン化します。



## マッシュデータ

トレーニングを高速化し、モデルのトレーニングや実験ごとに 600k 以上のファイルをダウンロードしないようにするために、データを最初に中央のデータレイクオブジェクトストアにアップロードしたあとにデータキャッシュ機能を使用しました。



ユーザーが Cache をクリックすると、cnvrg はリモートオブジェクトストアから特定のコミットでデータをダウンロードし、ONTAP NFS ボリュームにキャッシュします。完了すると、データをすぐにトレーニングに利用できるようになります。さらに、データが数日間使用されていない場合（たとえば、モデルのトレーニングや探索など）、cnvrg は自動的にキャッシュをクリアします。

## キャッシュデータで ML パイプラインを構築

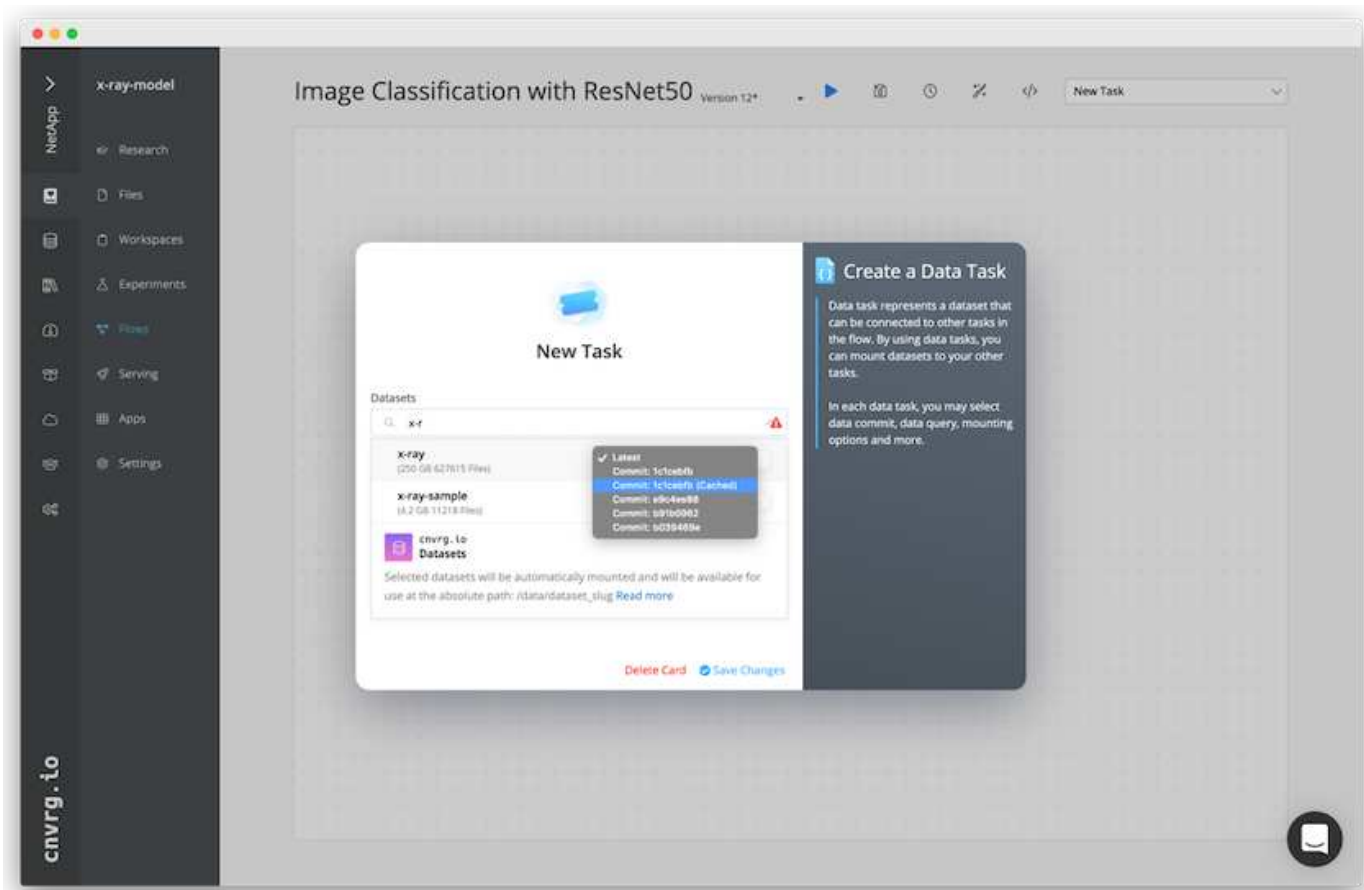
cnvrg フローを使用すると、本番 ML パイプラインを簡単に構築できます。フローは柔軟性が高く、あらゆる種類の ML ユースケースに対応し、GUI またはコードを使用して作成できます。フロー内の各コンポーネントは、異なる Docker イメージを使用して異なるコンピューティングリソース上で実行できるため、ハイブリッドクラウドを構築し、ML パイプラインを最適化できます。





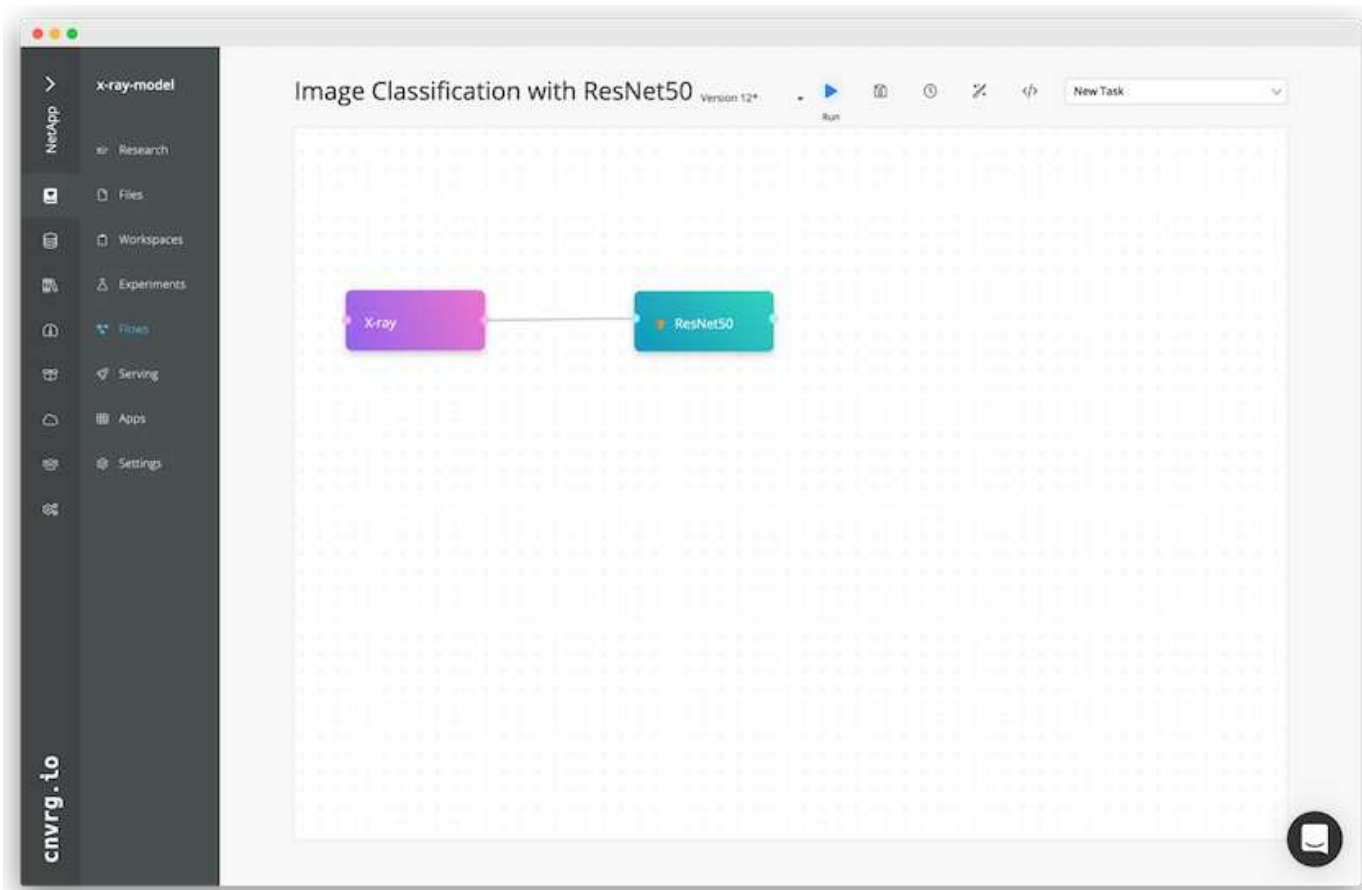
## 胸部 X 線フローの構築：データの設定

新しく作成したフローにデータセットを追加しました。データセットを追加する際には、特定のバージョン（commit）を選択し、キャッシュされたバージョンが必要かどうかを指定できます。この例では、キャッシュされたコミットを選択しました。



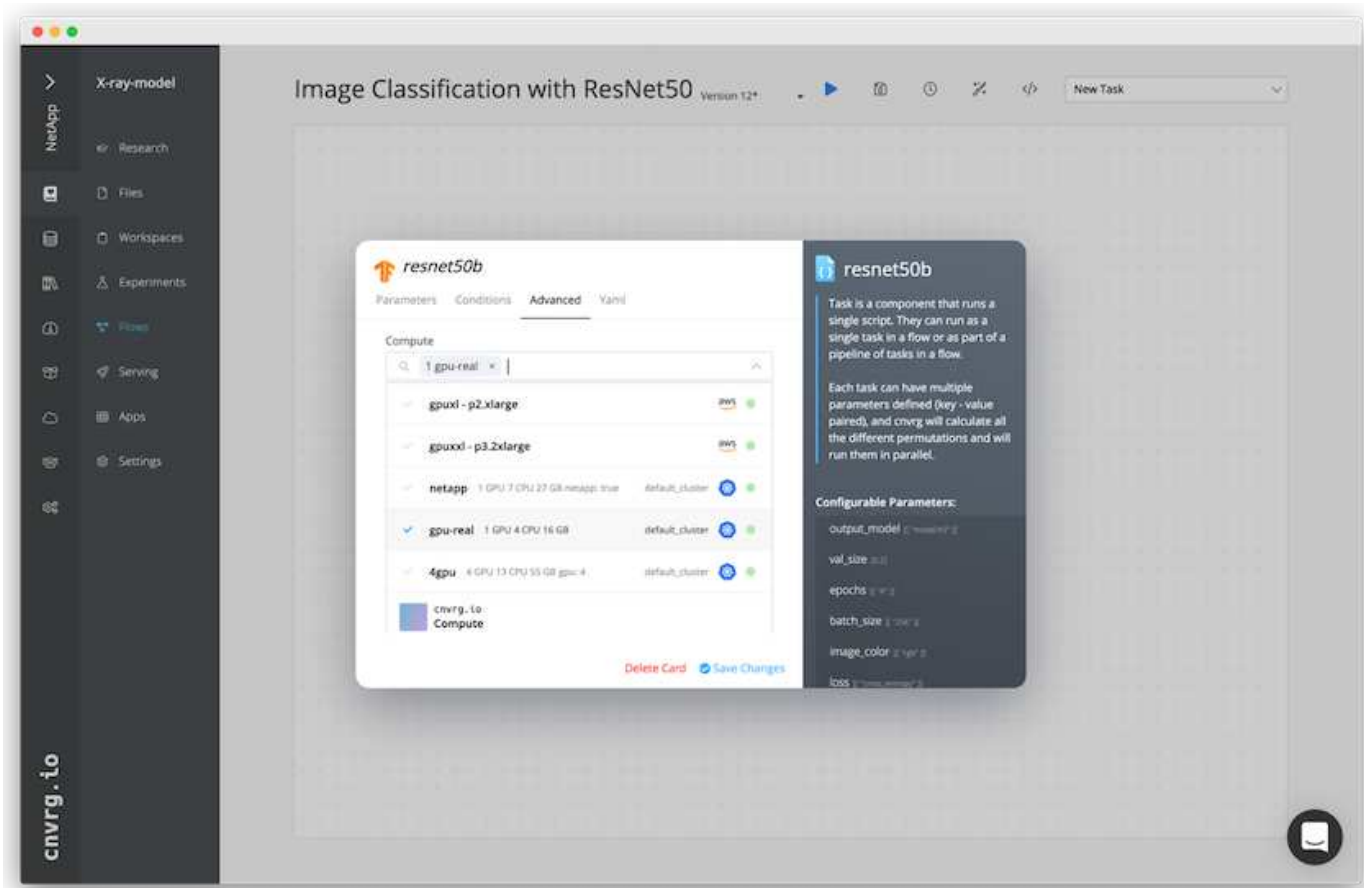
## 胸部 X 線フローの構築：トレーニングモデルの設定： **ResNet50**

パイプラインでは、任意の種類のカスタムコードを追加できます。cnvrg には、再利用可能な ML コンポーネントコレクションである AI ライブラリもあります。AI ライブラリには、アルゴリズム、スクリプト、データソースなど、あらゆる ML やディープラーニングフローで利用できるソリューションがあります。この例では、ResNet50 の事前ビルドモジュールを選択しました。batch\_size : 128、epochs : 10 などのデフォルトパラメータを使用しました。これらのパラメータは AI ライブラリのドキュメントで確認できます。次のスクリーンショットは、X 線データセットが ResNet50 に接続された新しいフローを示しています。



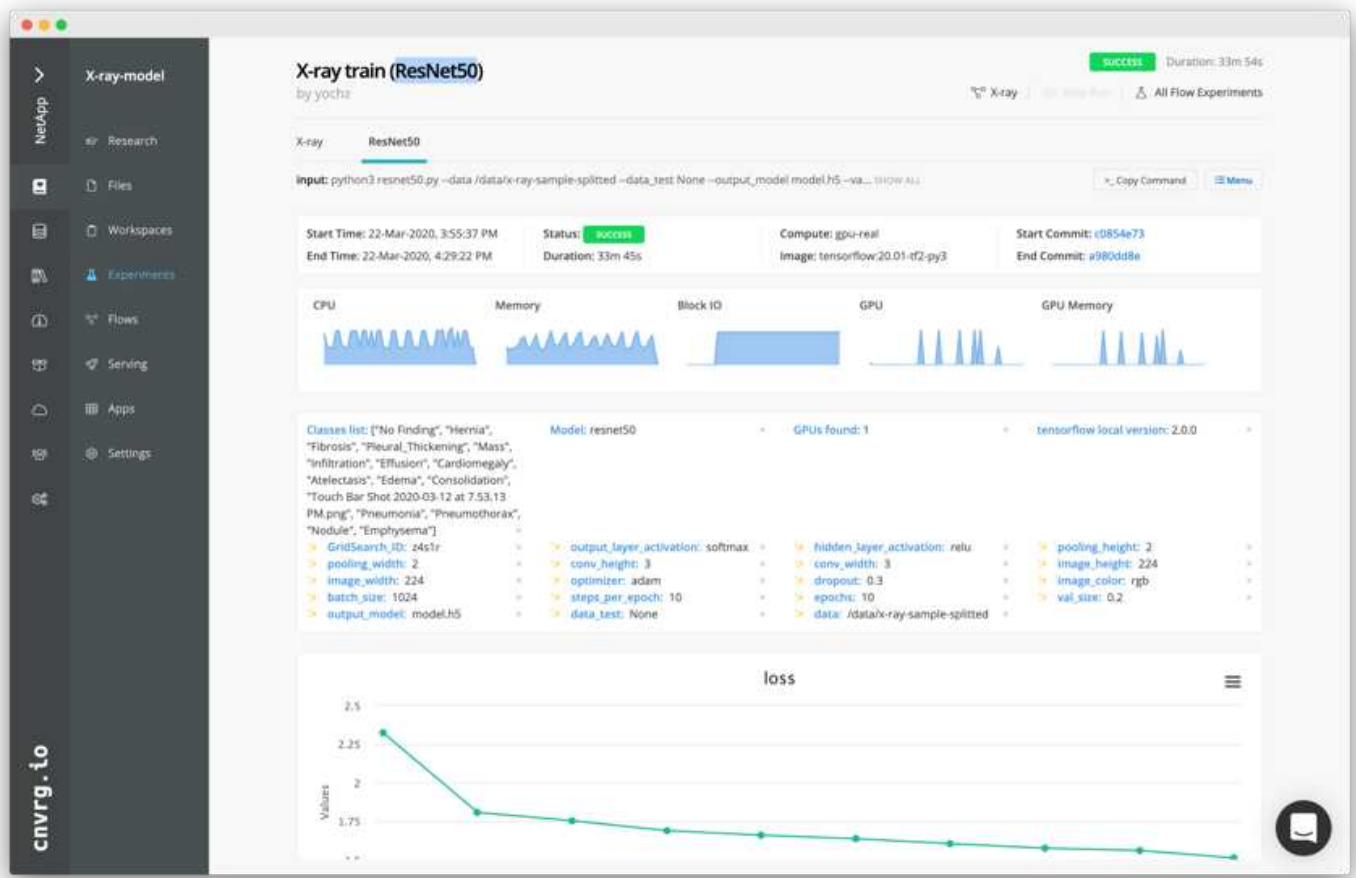
## ResNet50 の計算リソースを定義します

cnvrg フロー内の各アルゴリズムまたはコンポーネントは、異なる Docker イメージを使用して、異なるコンピューティングインスタンス上で実行できます。セットアップでは、NetApp ONTAP AI アーキテクチャを採用した NVIDIA DGX システムでトレーニングアルゴリズムを実行したいと考えていました。次の図では、「GPU - REAL」を選択しました。これは、オンプレミスクラスタのコンピューティングテンプレートであり、仕様です。また、テンプレートのキューを作成し、複数のテンプレートを選択しました。このようにして 'GPU 実数のリソースを割り当てることができない場合 (たとえば '他のデータ・サイエンティストがリソースを使用している場合) は 'クラウド・プロバイダ・テンプレートを追加して '自動クラウド・バーストを有効にできます次のスクリーンショットは、ResNet50 のコンピューティングノードとしての GPU 実数の使用を示しています。



## 結果の追跡と監視

フローが実行されると、cnvrg はトラッキングおよびモニタリングエンジンをトリガーします。フローの各実行は自動的に文書化され、リアルタイムで更新されます。ハイパーパラメータ、指標、リソース使用率（GPU 利用率など）、コードバージョン、アーティファクト、ログ また、次の 2 つのスクリーンショットに示すように、[テスト] セクションで自動的に使用できるようになります。



## 著作権に関する情報

Copyright © 2024 NetApp, Inc. All Rights Reserved. Printed in the U.S. このドキュメントは著作権によって保護されています。著作権所有者の書面による事前承諾がある場合を除き、画像媒体、電子媒体、および写真複写、記録媒体、テープ媒体、電子検索システムへの組み込みを含む機械媒体など、いかなる形式および方法による複製も禁止します。

ネットアップの著作物から派生したソフトウェアは、次に示す使用許諾条項および免責条項の対象となります。

このソフトウェアは、ネットアップによって「現状のまま」提供されています。ネットアップは明示的な保証、または商品性および特定目的に対する適合性の暗示的保証を含み、かつこれに限定されないいかなる暗示的な保証も行いません。ネットアップは、代替品または代替サービスの調達、使用不能、データ損失、利益損失、業務中断を含み、かつこれに限定されない、このソフトウェアの使用により生じたすべての直接的損害、間接的損害、偶発的損害、特別損害、懲罰的損害、必然的損害の発生に対して、損失の発生の可能性が通知されていたとしても、その発生理由、根拠とする責任論、契約の有無、厳格責任、不法行為（過失またはそうでない場合を含む）にかかわらず、一切の責任を負いません。

ネットアップは、ここに記載されているすべての製品に対する変更を随時、予告なく行う権利を保有します。ネットアップによる明示的な書面による合意がある場合を除き、ここに記載されている製品の使用により生じる責任および義務に対して、ネットアップは責任を負いません。この製品の使用または購入は、ネットアップの特許権、商標権、または他の知的所有権に基づくライセンスの供与とはみなされません。

このマニュアルに記載されている製品は、1つ以上の米国特許、その他の国の特許、および出願中の特許によって保護されている場合があります。

権利の制限について：政府による使用、複製、開示は、DFARS 252.227-7013（2014年2月）およびFAR 5252.227-19（2007年12月）のRights in Technical Data -Noncommercial Items（技術データ - 非商用品目に関する諸権利）条項の(b)(3)項、に規定された制限が適用されます。

本書に含まれるデータは商用製品および / または商用サービス（FAR 2.101の定義に基づく）に関係し、データの所有権はNetApp, Inc.にあります。本契約に基づき提供されるすべてのネットアップの技術データおよびコンピュータ ソフトウェアは、商用目的であり、私費のみで開発されたものです。米国政府は本データに対し、非独占的かつ移転およびサブライセンス不可で、全世界を対象とする取り消し不能の制限付き使用权を有し、本データの提供の根拠となった米国政府契約に関連し、当該契約の裏付けとする場合にのみ本データを使用できます。前述の場合を除き、NetApp, Inc.の書面による許可を事前に得ることなく、本データを使用、開示、転載、改変するほか、上演または展示することはできません。国防総省にかかる米国政府のデータ使用权については、DFARS 252.227-7015(b)項（2014年2月）で定められた権利のみが認められます。

## 商標に関する情報

NetApp、NetAppのロゴ、<http://www.netapp.com/TM>に記載されているマークは、NetApp, Inc.の商標です。その他の会社名と製品名は、それを所有する各社の商標である場合があります。