



デベロッパーガイド

AWS Deep Learning AMIs



AWS Deep Learning AMIs: デベロッパーガイド

Copyright © 2025 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon の商標およびトレードドレスは Amazon 以外の製品およびサービスに使用することはできません。また、お客様に誤解を与える可能性がある形式で、または Amazon の信用を損なう形式で使用することもできません。Amazon が所有していないその他のすべての商標は Amazon との提携、関連、支援関係の有無にかかわらず、それら該当する所有者の資産です。

Table of Contents

DLAMI とは	1
本ガイドについて	1
前提条件	1
ユースケースの例	1
機能	2
事前インストールされたフレームワーク	2
事前インストールされた GPU ソフトウェア	3
モデルの処理および可視化	3
DLAMI のリリースノート	4
Base DLAMI	4
シングルフレームワークの DLAMI	5
マルチフレームワークの DLAMI	6
入門	7
DLAMI の選択	7
CUDA のインストール環境およびフレームワークのバインド	8
基本	9
Conda	9
アーキテクチャ	11
OS	11
インスタンスの選択	12
料金	13
利用可能なリージョン	13
GPU	14
CPU	15
Inferentia	15
Trainium	16
設定	17
DLAMI ID の検索	17
インスタンスの起動	19
インスタンスへの接続	21
Jupyter の設定	21
サーバーの保護	22
サーバーの起動	23
クライアントの接続	23

ログイン	25
クリーンアップ	28
DLAMI の使用	29
Conda DLAMI	29
Deep Learning AMI with Conda の概要	29
DLAMI にログインする	30
TensorFlow 環境を開始する	30
PyTorch Python 3 環境に切り替える	31
環境を削除する	32
Base DLAMI	32
Deep Learning Base AMI の使用	32
CUDA のバージョンの設定	33
Jupyter ノートブック	33
インストールされたチュートリアルを操作する	34
Jupyter で環境を切り替える	34
チュートリアル	35
フレームワークのアクティブ化	35
Elastic Fabric Adapter	39
GPU のモニタリングおよび最適化	52
AWS 推論	62
ARM64 DLAMI	84
推論	87
モデル提供	88
DLAMI のアップグレード	92
DLAMI のアップグレード	92
ソフトウェアの更新	93
リリースの通知	94
セキュリティ	96
データ保護	97
Identity and Access Management	98
アイデンティティを使用した認証	98
ポリシーを使用したアクセスの管理	101
Amazon EMR での IAM	104
コンプライアンス検証	104
耐障害性	105
インフラストラクチャセキュリティ	105

モニタリング	106
使用状況の追跡	106
DLAMI サポートポリシー	107
DLAMI サポートFAQs	107
どのフレームワークバージョンにセキュリティパッチが適用されますか?	108
セキュリティパッチはどのオペレーティングシステムで取得されますか?	108
新しいフレームワークバージョンがリリースされると、どのイメージが AWS 公開されますか?	108
どのイメージに新しい SageMaker AI/AWS 機能が追加されますか?	108
サポート対象フレームワークの表では、現在のバージョンはどのように定義されていますか?	109
サポート対象テーブルにないバージョンを実行している場合はどうなりますか?	109
DLAMIs フレームワークバージョンの以前のパッチバージョンをサポートしていますか? ..	109
サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにありますか?	109
新しいイメージはどのくらいの頻度でリリースされますか?	110
ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか?	110
新しいパッチが適用されたフレームワークバージョン、または更新されたフレームワークバージョンが利用可能になった場合はどうなりますか?	110
フレームワークのバージョンを変更せずに依存関係は更新されますか?	110
使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか?	110
アクティブにメンテナンスされなくなったフレームワークバージョンのイメージにはパッチが適用されますか?	112
古いフレームワークバージョンを使用するにはどうすればよいですか?	112
フレームワークとそのバージョンでのサポート変更に関する最新情報を得るにはどうすればよいですか?	112
Anaconda リポジトリを使用するには商用ライセンスが必要ですか?	112
重要な変更点	113
DLAMI の NVIDIA ドライバーの変更に関するよくある質問	113
何が変わったのですか?	113
この変更が行われた理由は何ですか?	114
この変更により影響を受けるのはどの DLAMI ですか?	115
ユーザーにとってこの変更にはどのような意味がありますか?	115
新しい DLAMI で失われる機能はありますか?	115
この変更は Deep Learning Containers に影響しますか?	116
関連情報	117

非推奨の機能	118
ドキュメント履歴	120
.....	cxxiii

とは AWS Deep Learning AMIs

AWS Deep Learning AMIs (DLAMI) は、クラウドでの深層学習に使用できるカスタマイズされたマシンイメージを提供します。DLAMIs は、小さな CPU のみのインスタンスから最新のハイパワーマルチ GPU インスタンスまで、さまざまな Amazon Elastic Compute Cloud (Amazon EC2) インスタンスタイプ AWS リージョン でほとんどので使用できます。DLAMI は、[NVIDIA CUDA](#)、[NVIDIA cuDNN](#)、および最も人気が高い深層学習フレームワークの最新リリースを使用するように事前に設定されたとうえで提供されます。

本ガイドについて

本ガイドのコンテンツは、DLAMI の起動と使用に役立ちます。また、トレーニングと推論の両方で深層学習の一般的なユースケースについてもいくつか説明します。その他、目的に応じた適切な AMI の選択や、好まれるインスタンスの種類についても取り上げます。

さらに DLAMI には、サポートされているフレームワークが提供するチュートリアルがいくつか含まれています。このガイドでは、各フレームワークをアクティブ化し、開始するのに適切なチュートリアルを見つける方法について説明します。また、分散トレーニング、デバッグ、Inferentia と AWS Trainium AWS の使用、その他の主要な概念に関するチュートリアルもあります。ブラウザでチュートリアルを実行するように Jupyter Notebook サーバーを設定する手順については、「[DLAMI インスタンスでの Jupyter Notebook サーバーの設定](#)」を参照してください。

前提条件

DLAMI を正常に実行するには、コマンドラインツールと基本的な Python についての知識があることが推奨されます。

DLAMI ユースケースの例

(AWS Deep Learning AMIs DLAMI) の一般的なユースケースの例を次に示します。

深層学習に関する学習 - DLAMI は、機械学習および深層学習フレームワークの習得や教育にうってつけです。DLAMI により、各フレームワークのインストール環境のトラブルシューティングや、同じコンピュータ上でフレームワークを協調させる作業に悩まなくて済むようになります。DLAMI には Jupyter Notebook が含まれています。そのため、機械学習や深層学習が初めての方でも、フレームワークが提供するチュートリアルを簡単に実行できます。

アプリ開発 - アプリの開発を担当しており、深層学習によりアプリで AI の最新テクノロジーを利用することに関心をお持ちの方にとって、DLAMI は、申し分のない試験台になります。各フレームワークには、ディープラーニングを開始する方法についてのチュートリアルが付属しています。これらのチュートリアルの多くには Model Zoo があり、ニューラルネットワークを自分で作成したり、モデルトレーニングを実行したりせずに、簡単にディープラーニングを試すことができます。また、画像検出アプリケーションをわずか数分で作成する方法や、chatbot 用の音声認識アプリを作成する方法が例として示されています。

機械学習およびデータ分析 - データサイエンティストの方、または深層学習でデータを処理することに関心をお持ちの方は、フレームワークの多くで R と Spark がサポートされていることに気付くでしょう。スケーラブルなパーソナライズ用データ処理システムと予測システムを構築するまで、シンプルな回帰を実行する方法についてのチュートリアルが用意されています。

調査 - 新しいフレームワークを試したり、新しいモデルをテストしたり、新しいモデルをトレーニングしたりしたいと考えている調査担当者にとって、DLAMI とスケール用の AWS 機能は、複数のトレーニングノードの面倒なインストールと管理の複雑さを軽減できます。

Note

より多くの GPU (最大 8 個) を持つ大きなインスタンスにインスタンスタイプをアップグレードすることを最初に選択する場合もあるかもしれませんが、DLAMI インスタンスのクラスターを作成して、水平方向にスケーリングすることも可能です。クラスター構築の詳細については、[DLAMI の関連情報](#) を参照してください。

DLAMI の機能

AWS Deep Learning AMIs (DLAMI) の機能には、プリインストールされた深層学習フレームワーク、GPU ソフトウェア、モデルサーバー、モデル可視化ツールなどがあります。

事前インストールされたフレームワーク

現在、DLAMI には、主に 2 つのタイプがあります。また、それぞれについて、オペレーティングシステム (OS) とソフトウェアのバージョンに関連する他のバリエーションがあります。

- [Deep Learning AMI with Conda](#) - conda パッケージ、および独立した Python 環境を使用して個別にインストールされるフレームワーク。
- [Deep Learning Base AMI](#) - フレームワークはインストールされておらず、[NVIDIA CUDA](#) と他の依存関係のみがインストールされています。

Deep Learning AMI with Conda は、conda 環境を使用して各フレームワークを分離します。そのため、各フレームワークを自由に切り替えることができ、依存関係の競合を心配する必要はありません。Deep Learning AMI with Conda は、以下のフレームワークをサポートしています。

- PyTorch
- TensorFlow 2

Note

DLAMI では、Apache MXNet、Microsoft Cognitive Toolkit (CNTK)、Caffe、Caffe2、Theano、Chainer、Keras などの深層学習フレームワークはサポートされなくなりました。

事前インストールされた GPU ソフトウェア

CPU のみのインスタンスを使用する場合でも、DLAMI では、[NVIDIA CUDA](#) と [NVIDIA cuDNN](#) が使用されます。インストールされたソフトウェアは、インスタンスタイプに関係なく同じです。GPU 専用ツールは、GPU が 1 つ以上あるインスタンスでのみ機能することに注意してください。インスタンスタイプの詳細については、「[DLAMI インスタンスタイプの選択](#)」を参照してください。

CUDA の詳細については、「[CUDA のインストール環境およびフレームワークのバインド](#)」を参照してください。

モデルの処理および可視化

Deep Learning AMI with Conda には、TensorFlow 用とモデル可視化用の TensorBoard 用のモデルサーバーがあらかじめインストールされています。詳細については、「[TensorFlow Serving](#)」を参照してください。

DLAMI のリリースノート

ここでは、現在サポートされているすべての AWS Deep Learning AMIs (DLAMI) オプションの詳細なリリースノートを確認できます。

サポートが終了した DLAMI フレームワークのリリースノートについては、「[DLAMI フレームワーク サポートポリシー](#)」ページの「Unsupported Framework Release Notes Archive」セクションを参照してください。

Note

AWS Deep Learning AMIs には、セキュリティパッチの夜間リリース頻度があります。これらの増分セキュリティパッチは公式リリースノートには含まれていません。

Base DLAMI

GPU

- X86
 - [AWS Deep Learning Base AMI \(Amazon Linux 2023\)](#)
 - [AWS Deep Learning Base AMI \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning Base AMI \(Ubuntu 20.04\)](#)
 - [AWS Deep Learning Base AMI \(Amazon Linux 2\)](#)
- ARM64
 - [AWS Deep Learning Base ARM64 AMI \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning Base ARM64 AMI \(Amazon Linux 2\)](#)
 - [AWS Deep Learning Base ARM64 AMI \(Amazon Linux 2023\)](#)

Qualcomm

- X86
 - [AWS Deep Learning Base Qualcomm AMI \(Amazon Linux 2\)](#)

AWS Neuron

- [Neuron DLAMI ガイド](#)を参照してください。

シングルフレームワークの DLAMI

PyTorch 固有の AMI

GPU

- X86
 - [AWS Deep Learning AMI GPU PyTorch 2.6 \(Amazon Linux 2023\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.6 \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.5 \(Amazon Linux 2023\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.5 \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.4 \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.3 \(Ubuntu 20.04\)](#)
 - [AWS Deep Learning AMI GPU PyTorch 2.3 \(Amazon Linux 2\)](#)
- ARM64
 - [AWS Deep Learning ARM64 AMI GPU PyTorch 2.6 \(Amazon Linux 2023\)](#)
 - [AWS Deep Learning ARM64 AMI GPU PyTorch 2.6 \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning ARM64 AMI GPU PyTorch 2.5 \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning ARM64 AMI GPU PyTorch 2.4 \(Ubuntu 22.04\)](#)
 - [AWS Deep Learning ARM64 AMI GPU PyTorch 2.3 \(Ubuntu 22.04\)](#)

AWS Neuron

- [Neuron DLAMI ガイド](#)を参照してください。

TensorFlow 固有の AMI

GPU

- X86
 - [AWS Deep Learning AMI GPU TensorFlow 2.18 \(Amazon Linux 2023\)](#)
 - [AWS Deep Learning AMI GPU TensorFlow 2.18 \(Ubuntu 22.04\)](#)

- [AWS Deep Learning AMI GPU TensorFlow 2.17 \(Ubuntu 22.04\)](#)

AWS Neuron

- [Neuron DLAMI ガイド](#)を参照してください。

マルチフレームワークの DLAMI

Tip

機械学習フレームワークを1つだけ使用する場合は、[シングルフレームワークの DLAMI](#)をお勧めします。

GPU

- X86
 - [AWS Deep Learning AMI \(Amazon Linux 2\)](#)

AWS Neuron

- [Neuron DLAMI ガイド](#)を参照してください。

DLAMI の使用開始

このガイドには、最適な DLAMI を選択する方法やユースケースと予算に適したインスタンスタイプを選択する方法についてのヒントのほか、興味深いカスタムセットアップを説明している [DLAMI の関連情報](#) が記載されています。

Amazon EC2 を初めて使用する場合 AWS は、 から始めます [Deep Learning AMI with Conda](#)。Amazon EC2 や Amazon EMR、Amazon EFS、Amazon S3 などの他の AWS サービスに精通しており、分散トレーニングや推論が必要なプロジェクトにこれらのサービスを統合することに関心がある場合は、ユースケースに適している [DLAMI の関連情報](#) かどうかを確認してください。

[DLAMI の選択](#) を参照し、アプリケーションに最適なインスタンスタイプの概念を理解することをお勧めします。

次のステップ

[DLAMI の選択](#)

DLAMI の選択

GPU DLAMI [リリースノートに記載されているように、さまざまな DLAMI オプション](#) が用意されています。ユースケースに適切な DLAMI を選択できるように、開発されたハードウェアのタイプまたは機能別にイメージをグループ化しています。トップレベルのグループは次のとおりです。

- DLAMI タイプ： Base、Single-Framework、Multi-Framework (Conda DLAMI)
- コンピューティングアーキテクチャ： x86 ベース、Arm64-based [AWS Graviton](#)
- プロセッサタイプ： [GPU](#)、[CPU](#)、[Inferentia](#)、[Trainium](#)
- SDK: [CUDA](#)、[AWS Neuron](#)
- OS: Amazon Linux、Ubuntu

このガイドの残りのトピックでは、より詳細な情報を説明します。

トピック

- [CUDA のインストール環境およびフレームワークのバインド](#)
- [Deep Learning Base AMI](#)
- [Deep Learning AMI with Conda](#)

- [DLAMI アーキテクチャのオプション](#)
- [DLAMI のオペレーティングシステムのオプション](#)

次回の予定

[Deep Learning AMI with Conda](#)

CUDA のインストール環境およびフレームワークのバインド

深層学習は、非常に最先端のものですが、各フレームワークは、"安定した" バージョンを提供しています。これらの安定したバージョンは、CUDA または cuDNN の最新の実装および機能とともに動作しない場合があります。ユースケースと必要な機能に基づいて、フレームワークを選択できます。よくわからない場合は、最新の Deep Learning AMI with Conda を使用してください。これには、CUDA を利用したすべてのフレームワークの公式 pip バイナリが含まれており、各フレームワークでサポートされる最新バージョンが使用されます。最新のバージョンが必要な場合で、深層学習環境をカスタマイズするには、Deep Learning Base AMI を使用してください。

「[安定性とリリース候補](#)」で当社のガイドを参照してください。

DLAMI with CUDA の選択

[Deep Learning Base AMI](#) には、使用可能なすべての CUDA バージョンシリーズがあります

[Deep Learning AMI with Conda](#) には、使用可能なすべての CUDA バージョンシリーズがあります

Note

MXNet、CNTK、Caffe、Caffe2、Theano、Chainer、Keras Conda 環境は AWS Deep Learning AMIs に含まれなくなりました。

特定のフレームワークのバージョン番号については、[DLAMI のリリースノート](#)を参照してください。

この DLAMI タイプを選択するか、次回の予定オプションを使ってさまざまな DLAMI の詳しい情報を確認してください。

いずれかのバージョンの CUDA を選択してから、そのバージョンを含む DLAMI の詳細なリストを付録で確認するか、次回の予定オプションで別の DLAMI の詳細を確認してください。

次回の予定

[Deep Learning Base AMI](#)

関連トピック

- CUDA のバージョンを切り替える方法については、「[Deep Learning Base AMI の使用](#)」チュートリアルを参照してください。

Deep Learning Base AMI

Deep Learning Base AMI は、深層学習のための真っ白なキャンバスのようなものです。特定のフレームワークをインストールするまでに必要なすべてが含まれており、CUDA のバージョンも選択できます。

Base DLAMI を選択する理由

最先端のディープラーニングプロジェクトへの参加を目指すプロジェクト協力者は、この AMI グループを有効活用できます。また、環境を展開する際に、最新の NVIDIA ソフトウェアがインストールされ動作していることを確信できるため、インストールするフレームワークとバージョンの選択に集中できます。

この DLAMI タイプを選択するか、次回の予定オプションを使ってさまざまな DLAMI の詳しい情報を確認してください。

次回の予定

[DLAMI with Conda](#)

関連トピック

- [Deep Learning Base AMI の使用](#)

Deep Learning AMI with Conda

Conda DLAMI は conda 仮想環境を使用しており、マルチフレームワークまたは単一フレームワークのいずれかの DLAMI があります。これらの環境は、インストールした異なるフレームワークの独立性が維持され、フレームワーク間の切り替えが合理化されるように設定されます。これは、DLAMI が提供するすべてのフレームワークを学習したり、試したりするのに最適です。ほとんどのユーザーにとって、新しい Deep Learning AMI with Conda は申し分のないものとなります。

フレームワークからの最新バージョンで頻繁に更新され、最新の GPU ドライバとソフトウェアが導入されることとなります。これらは通常、AWS Deep Learning AMIs ほとんどのドキュメントではと呼ばれます。これらの DLAMIs、Ubuntu 20.04、Ubuntu 22.04、Amazon Linux 2、Amazon Linux 2023 オペレーティングシステムをサポートしています。オペレーティングシステムのサポートは、アップストリーム OS からのサポートによって異なります。

安定性とリリース候補

Conda AMI は、各フレームワークの最新の正式リリース版の最適化バイナリを使用します。リリース候補および実験的機能は対象外です。この最適化は、C5、および C4 CPU インスタンスタイプにおけるトレーニングおよび推論を高速化する高速化テクノロジー (MKL DNN など) をサポートするフレームワークに依存します。また、バイナリは高度な Intel 指示セット (AVX、AVX-2、SSE4.1、SSE4.2 などを含む) をサポートするようにコンパイルされています。これによって、Intel CPU アーキテクチャ上のベクターおよび浮動小数点オペレーションが高速化されます。さらに、GPU インスタンスタイプでは、CUDA および cuDNN は最新の公式リリースがサポートするバージョンで更新されます。

Deep Learning AMI with Conda は、Amazon EC2 インスタンスにフレームワークの最初のアクティベーションから最適化されたフレームワークのバージョンを自動的にインストールします。詳細については、「[Deep Learning AMI with Conda の使用](#)」を参照してください。

カスタムまたは最適化されたビルドオプションを使用してソースからインストールする場合は、[Deep Learning Base AMI](#) の方が適しています。

Python 2 の廃止

Python オープンソースコミュニティは、Python 2 のサポートを 2020 年 1 月 1 日に正式に終了しました。TensorFlow と PyTorch コミュニティは、TensorFlow 2.1 および PyTorch 1.4 リリースが Python 2 をサポートする最後のリリースになると発表しました。Python 2 Conda 環境を含む DLAMI (v26、v25 など) の以前のリリースは、引き続き利用できます。ただし、以前に公開された DLAMI のバージョンで Python 2 Conda 環境の更新を提供するのは、それらのバージョンのオープンソースコミュニティによって公開されたセキュリティ修正がある場合のみです。TensorFlow および PyTorch フレームワークの今後のバージョンを利用する DLAMI リリースには、Python 2 Conda 環境は含められません。

CUDA サポート

特定の CUDA バージョン番号は [GPU DLAMI リリースノート](#) で確認できます。

次回の予定

[DLAMI アーキテクチャのオプション](#)

関連トピック

- Deep Learning AMI with Conda の使用に関するチュートリアルについては、[Deep Learning AMI with Conda の使用](#)のチュートリアルを参照してください。

DLAMI アーキテクチャのオプション

AWS Deep Learning AMIsは、x86 ベースまたは Arm64 ベースの [AWS Graviton2](#) のアーキテクチャで提供されています。

ARM64 GPU DLAMI の使用開始方法については、「[ARM64 DLAMI](#)」を参照してください。使用可能なインスタンスタイプの詳細については、[DLAMI インスタンスタイプの選択](#)を参照してください。

次回の予定

[DLAMI のオペレーティングシステムのオプション](#)

DLAMI のオペレーティングシステムのオプション

DLAMI は、以下のオペレーティングシステムで提供されています。

- Amazon Linux 2
- Amazon Linux 2023
- Ubuntu 20.04
- Ubuntu 22.04

古いバージョンのオペレーティングシステムは、廃止された DLAMI でも使用できます。DLAMI の廃止について詳しくは、「[DLAMI の廃止](#)」を参照してください。

DLAMI を選択する前に、必要なインスタンスタイプを評価し、AWS リージョンを特定してください。

次回の予定

[DLAMI インスタンスタイプの選択](#)

DLAMI インスタンスタイプの選択

一般的に、DLAMI のインスタンスタイプを選択するときは、次の点を考慮してください。

- 深層学習を初めて導入する場合は、GPU が 1 個搭載のインスタンスがニーズに合っている可能性があります。
- 予算重視の場合は、CPU のみのインスタンスを使用できます。
- 深層学習モデル推論のために高いパフォーマンスとコスト効率を最適化する場合は、AWS Inferentia チップでインスタンスを使用できます。
- Arm64 ベースの CPU アーキテクチャを備えた高いパフォーマンスの GPU インスタンスをお探しなら、G5g インスタンスタイプを使用できます。
- 推論と予測用に事前トレーニング済みモデルを実行することに興味がある場合は、[Amazon Elastic Inference](#) を Amazon EC2 インスタンスにアタッチできます。Amazon Elastic Inference では、GPU の一部を利用したアクセラレーターにアクセスできます。
- 大容量の推論サービスの場合は、大容量のメモリを搭載した 1 つの CPU インスタンス、またはそのようなインスタンスのクラスターが、より適したソリューションになることがあります。
- 大量のデータまたは大きなバッチサイズを持つ大規模なモデルを使用している場合は、より多くのメモリを持つより大きなインスタンスが必要になります。モデルを GPU のクラスターに配布することもできます。バッチサイズを小さくした場合、メモリの少ないインスタンスを使用すると改善されることがあります。これは、精度とトレーニング速度に影響を与える可能性があります。
- 大規模で高レベルのノード間通信を必要とする NVIDIA Collective Communications Library (NCCL) を使用した機械学習アプリケーションの実行に関心がある場合は、[Elastic Fabric Adapter \(EFA\)](#) を使用できます。

インスタンスの詳細については、[EC2 インスタンスタイプ](#)を参照してください。

次のトピックでは、インスタンスタイプに関する考慮事項について説明します。

Important

ディープラーニング AMI には、ドライバ、ソフトウェアやツールキット (NVIDIA Corporation によって開発、所有あるいは提供) が含まれています。これらの NVIDIA ドライバ、ソフトウェア、ツールキットは、NVIDIA ハードウェアが含まれている Amazon EC2 インスタンスでのみ使用することにユーザーが同意するものとします。

トピック

- [DLAMI の価格設定](#)
- [DLAMI で利用可能なリージョン](#)
- [推奨 GPU インスタンス](#)
- [推奨 CPU インスタンス](#)
- [推奨 Inferentia インスタンス](#)
- [推奨 Trainium インスタンス](#)

DLAMI の価格設定

DLAMI に含まれている深層学習フレームワークは無料で、それぞれに独自のオープンソースライセンスが付与されています。DLAMI に含まれているソフトウェアは無料ですが、基盤となる Amazon EC2 インスタンスハードウェアの料金を支払う必要があります。

Amazon EC2 インスタンスタイプには無料とされているものがあります。そうした無料のインスタンスのいずれかで DLAMI を実行できます。これは、そのインスタンスの容量だけを使う場合は、DLAMI の使用が完全に無料という意味です。CPU コア数を増やす、ディスク容量を増やす、RAM を増やす、GPU を 1 つ以上使用するなど、より強力なインスタンスが必要な場合は、無料利用枠のインスタンスクラス以外のインスタンスが必要になります。

インスタンスの選択と料金の詳細については、「[Amazon EC2 料金表](#)」を参照してください。

DLAMI で利用可能なリージョン

各リージョンでサポートされるインスタンスタイプの範囲はそれぞれ異なり、多くの場合、インスタンスタイプのコストもリージョンごとにわずかに違います。DLAMI を利用できないリージョンはありますが、選択したリージョンに DLAMI をコピーすることは可能です。詳細については、[AMI のコピー](#)を参照してください。リージョンの選択リストを確認し、必ず自社または自社の顧客に近いリージョンを選択してください。複数の DLAMI を使用する予定があり、クラスターを作成する可能性がある場合は、クラスター内のすべてのノードに同じリージョンを使う必要があります。

リージョンの詳細については、「[Amazon EC2 service endpoints](#)」を参照してください。

次回の予定

[推奨 GPU インスタンス](#)

推奨 GPU インスタンス

GPU インスタンスは、深層学習の大半の目的に推奨されます。新しいモデルのトレーニングは CPU インスタンス上よりも GPU インスタンス上の方が高速に実行できます。複数の GPU インスタンスがある場合や、トレーニングを複数の GPU インスタンスに分散した場合は、ほぼ直線的な拡張性を得ることができます。

以下のインスタンスタイプで DLAMI がサポートされています。GPU インスタンスタイプのオプションとその使用方法の詳細については、[EC2 インスタンスタイプ](#)を参照し、高速コンピューティングを選択してください。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 P5e インスタンス](#)には、NVIDIA Tesla H200 GPU が最大 8 個搭載されます。
- [Amazon EC2 P5 インスタンス](#)には、NVIDIA Tesla H100 GPU が最大 8 個搭載されます。
- [Amazon EC2 P4 インスタンス](#)には、NVIDIA Tesla A100 GPU が最大 8 個搭載されます。
- [Amazon EC2 P3 インスタンス](#)には、NVIDIA Tesla V100 GPU が最大 8 個搭載されます。
- [Amazon EC2 G3 インスタンス](#)には、NVIDIA Tesla M60 GPU が最大 4 個搭載されます。
- [Amazon EC2 G4 インスタンス](#)には、NVIDIA T4 GPU が最大 4 個搭載されます。
- [Amazon EC2 G5 インスタンス](#)には、NVIDIA A10G GPU が最大 8 個搭載されます。
- [Amazon EC2 G6 インスタンス](#)には、NVIDIA L4 GPU が最大 8 個搭載されます。
- [Amazon EC2 G6e インスタンス](#)には、NVIDIA L40S Tensor Core GPU が最大 8 個搭載されます。
- [Amazon EC2 G5g インスタンス](#)には、Arm64 ベースの [AWS Graviton2 プロセッサ](#)が搭載されません。

DLAMI インスタンスでは、GPU プロセスをモニタリングおよび最適化するためのツールが提供されています。GPU プロセスのモニタリングの詳細については、[GPU のモニタリングおよび最適化](#)を参照してください。

G5g インスタンスの操作に関する具体的なチュートリアルについては、[ARM64 DLAMI](#)を参照してください。

次回の予定

[推奨 CPU インスタンス](#)

推奨 CPU インスタンス

予算に制約がある場合、ディープラーニングを学習する目的の場合、または予測サービスの実行が唯一の目的である場合は、CPU のカテゴリに数多くの低価格な選択肢があります。一部のフレームワークは、C5 (一部のリージョンでのみ使用可) CPU インスタンスタイプにおけるトレーニングと推定を高速化する Intel の MKL DNN を活用しています。CPU インスタンスタイプの詳細については、[EC2 インスタンスタイプ](#) を参照し、コンピューティング最適化を選択してください。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 C5 インスタンス](#)には最大 72 個のインテル vCPU があります。C5 インスタンスは、科学モデリング、バッチ処理、分散分析、ハイパフォーマンスコンピューティング (HPC)、機械学習/深層学習による推論などに優れています。

次回の予定

[推奨 Inferentia インスタンス](#)

推奨 Inferentia インスタンス

AWS Inferentia インスタンスは、深層学習モデル推論ワークロードに高いパフォーマンスとコスト効率を提供するように設計されています。具体的には、Inf2 インスタンスタイプは AWS Inferentia チップと [AWS Neuron SDK](#) を使用します。これは、TensorFlow や PyTorch などの一般的な機械学習フレームワークと統合されています。

お客様は Inf2 インスタンスを使用して、検索、レコメンデーションエンジン、コンピュータビジョン、音声認識、自然言語処理、パーソナライゼーション、不正検出などの大規模な機械学習推論アプリケーションをクラウド内で低コストで実行できます。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 Inf2 インスタンス](#)には、最大 16 AWS Inferentia チップと 100 Gbps のネットワークスループットがあります。

AWS Inferentia DLAMIs 「」を参照してください [DLAMI を使用した AWS Inferentia チップ](#)。

次回の予定

[推奨 Trainium インスタンス](#)

推奨 Trainium インスタンス

AWS Trainium インスタンスは、深層学習モデル推論ワークロードに高いパフォーマンスとコスト効率を提供するように設計されています。具体的には、Trn1 インスタンスタイプは AWS Trainium チップと [AWS Neuron SDK](#) を使用します。これは TensorFlow や PyTorch などの一般的な機械学習フレームワークと統合されています。

お客様は Trn1 インスタンスを使用して、検索、レコメンデーションエンジン、コンピュータビジョン、音声認識、自然言語処理、パーソナライゼーション、不正検出などの大規模な機械学習推論アプリケーションをクラウド内で低コストで実行できます。

Note

モデルのサイズは、インスタンスを選択する際の要因となります。モデルがインスタンスの使用可能な RAM を超えている場合は、アプリケーション用に十分なメモリを持つ別のインスタンスタイプを選択します。

- [Amazon EC2 Trn1 インスタンス](#)には、最大 16 個の AWS Trainium チップと 100 Gbps のネットワークスループットがあります。

DLAMI インスタンスの設定

[DLAMI を選択し](#)、使用する [Amazon Elastic Compute Cloud \(Amazon EC2\) インスタンスタイプ](#) を選択したら、新しい DLAMI インスタンスを設定する準備が整います。

まだ DLAMI および EC2 インスタンスタイプを選択していない場合は、「[DLAMI の使用開始](#)」を参照してください。

トピック

- [DLAMI の ID の検索](#)
- [DLAMI インスタンスの起動](#)
- [DLAMI インスタンスへの接続](#)
- [DLAMI インスタンスでの Jupyter Notebook サーバーの設定](#)
- [DLAMI インスタンスのクリーンアップ](#)

DLAMI の ID の検索

DLAMI にはそれぞれ一意の識別子 (ID) があります。Amazon EC2 コンソールを使用して DLAMI インスタンスを起動する場合、DLAMI ID を使用して、使用する DLAMI を検索することもできます。AWS Command Line Interface (AWS CLI) を使用して DLAMI インスタンスを起動する場合、この ID が必要です。

選択した DLAMI の ID は、Amazon EC2 または Parameter Store の一機能用の AWS CLI コマンドを使用して確認できます AWS Systems Manager。のインストールと設定の手順については AWS CLI、[「ユーザーガイド」の AWS CLI](#)「の開始方法」を参照してください。AWS Command Line Interface

Using Parameter Store

`ssm get-parameter` を使用して DLAMI ID をを見つけるには

次の `ssm get-parameter` コマンドでは、`--name` オプションの場合、パラメータの名前形式は `/aws/service/deeplearning/ami/$architecture/$ami_type/latest/ami-id` です。この名前形式では、`architecture` は `x86_64` または `arm64` のいずれかになります。DLAMI 名を取得し、「`deep`」、「`learning`」、「`ami`」のキーワードを削除して、`ami_type` を指定します。AMI 名は [DLAMI のリリースノート](#) にあります。

⚠ Important

このコマンドを使用するには、使用する AWS Identity and Access Management (IAM) プリンシパルに アクセス `ssm:GetParameter` 許可が必要です。IAM プリンシパルの詳細については、「IAM ユーザーガイド」の「IAM ロール」にある「[その他のリソース](#)」セクションを参照してください。

```
aws ssm get-parameter --name /aws/service/deeplearning/ami/x86_64/base-oss-  
nvidia-driver-ubuntu-22.04/latest/ami-id \  
--region us-east-1 --query "Parameter.Value" --output text
```

出力は次の例に類似したものになります:

```
ami-09ee1a996ac214ce7
```

ℹ Tip

現在サポートされている DLAMI フレームワークの一部については、[DLAMI のリリースノート](#) でより具体的な `ssm get-parameter` コマンドの例を確認できます。選択した DLAMI のリリースノートへのリンクを選択し、その ID クエリをリリースノートで検索します。

Using Amazon EC2 CLI

`ec2 describe-images` を使用して DLAMI ID をを見つけるには

次の [ec2 describe-images](#) コマンドで、フィルター `Name=name` の値に DLAMI 名を入力します。特定のフレームワークのリリースバージョンを指定したり、バージョン番号を疑問符 (?) に置き換えて最新のリリースを取得したりできます。

```
aws ec2 describe-images --region us-east-1 --owners amazon \  
--filters 'Name=name,Values=Deep Learning Base OSS Nvidia Driver GPU AMI (Ubuntu  
22.04) ????????' 'Name=state,Values=available' \  
--query 'reverse(sort_by(Images, &CreationDate))[:1].ImageId' --output text
```

出力は次の例に類似したものになります:


```
ami-09ee1a996ac214ce7
```

i Tip

選択した DLAMI に固有の ec2 describe-images コマンドの例については、「[DLAMI のリリースノート](#)」を参照してください。選択した DLAMI のリリースノートへのリンクを選択し、その ID クエリをリリースノートで検索します。

次のステップ

[DLAMI インスタンスの起動](#)

DLAMI インスタンスの起動

DLAMI インスタンスの起動に使用する DLAMI の [ID が見つかったら](#)、インスタンスを起動する準備が整います。起動するには、Amazon EC2 コンソールまたは AWS Command Line Interface () を使用できますAWS CLI。

i Note

このチュートリアルでは、Deep Learning Base OSS Nvidia Driver GPU AMI (Ubuntu 22.04) に固有の内容に言及する可能性があります。異なる DLAMI を選択している場合でも、このガイドに従うことができます。

EC2 console

i Note

ハイパフォーマンスコンピューティング (HPC) と機械学習アプリケーションを加速するために、Elastic Fabric Adapter (EFA) を使用して DLAMI インスタンスを起動できます。固有の説明については、「[EFA を使用した AWS Deep Learning AMIs インスタンスの起動](#)」を参照してください。

1. [EC2 コンソール](#)を開きます。

2. 一番上のナビゲーション AWS リージョン で現在の を書き留めます。目的のリージョンでない場合、このオプションを変更して続行します。詳細については、Amazon Web Services 全般のリファレンスの「[Amazon EC2 service endpoints](#)」を参照してください。
3. [Launch Instance] (インスタンスの起動) を選択します。
4. インスタンスの名前を入力し、適切な DLAMI を選択します。
 - a. [自分の全 AMI] で既存の DLAMI を検索するか、[クイックスタート] を選択します。
 - b. DLAMI ID で検索します。オプションを参照し、目的の項目を選択します。
5. インスタンスタイプを選択します。DLAMI の推奨インスタンスファミリーは、[DLAMI のリリースノート](#) で確認できます。DLAMI インスタンスタイプに関する一般的な推奨事項については、「[DLAMI インスタンスタイプの選択](#)」を参照してください。
6. [Launch Instance] (インスタンスの起動) を選択します。

AWS CLI

- を使用するには AWS CLI、使用する DLAMI の ID、AWS リージョン および EC2 インスタンスタイプ、およびセキュリティトークン情報が必要です。その後、[ec2 run-instances](#) AWS CLI コマンドを使用してインスタンスを起動できます。

のインストールと設定の手順については AWS CLI、[「ユーザーガイド」の AWS CLI](#)「の開始方法」を参照してください。AWS Command Line Interface コマンド例を含む詳細については、「[AWS CLIの Amazon EC2 インスタンスの起動、リスト、終了](#)」を参照してください。

Amazon EC2 コンソールまたは を使用してインスタンスを起動したら AWS CLI、インスタンスの準備ができるまで待ちます。この待機時間は通常わずか数分です。インスタンスの状態は [Amazon EC2 コンソール](#) で確認できます。詳細については、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 インスタンスのステータスチェック](#)」を参照してください。

次のステップ

[DLAMI インスタンスへの接続](#)

DLAMI インスタンスへの接続

[DLAMI インスタンスを起動](#)してインスタンスを実行したら、SSH を使用してクライアント (Windows、macOS、または Linux) からそのインスタンスに接続できます。手順については、「Amazon EC2 ユーザーガイド」の「[SSH を使用した Linux インスタンスへの接続](#)」を参照してください。

ログイン後に Jupyter Notebook サーバーを設定する場合は、SSH ログインコマンドのコピーを手元に保管してください。Jupyter ウェブページに接続する際に、このコマンドのバリエーションを使用します。

次のステップ

[DLAMI インスタンスでの Jupyter Notebook サーバーの設定](#)

DLAMI インスタンスでの Jupyter Notebook サーバーの設定

Jupyter Notebook サーバーを使用すると、DLAMI インスタンスから Jupyter Notebook を作成して実行できます。Jupyter ノートブックを使用すると、AWS インフラストラクチャを使用し、DLAMI に組み込まれたパッケージにアクセスしながら、トレーニングと推論のための機械学習 (ML) 実験を実行できます。Jupyter Notebook の詳細については、[Jupyter のユーザードキュメント](#)のウェブサイトで「Jupyter Notebook」を参照してください。

Jupyter Notebook サーバーを設定するには、次の作業を実行する必要があります。

- DLAMI インスタンスで Jupyter Notebook サーバーを設定する。
- Jupyter Notebook サーバーに接続するようにクライアントを設定する。Windows クライアント、macOS クライアント、および Linux クライアントの設定手順が用意されています。
- Jupyter Notebook サーバーにログインして、設定をテストする。

これらのステップを完了するには、以下のトピックの手順に従います。Jupyter Notebook サーバーを設定したら、DLAMI に付属している Notebook チュートリアルサンプルを実行できます。詳細については、「[Jupyter ノートブックチュートリアルを実行する](#)」を参照してください。

トピック

- [DLAMI インスタンスでの Jupyter Notebook サーバーの保護](#)
- [DLAMI インスタンスでの Jupyter Notebook サーバーの起動](#)

- [DLAMI インスタンスで Jupyter Notebook サーバーにクライアントを接続する](#)
- [DLAMI インスタンスでの Jupyter Notebook サーバーへのログイン](#)

DLAMI インスタンスでの Jupyter Notebook サーバーの保護

Jupyter Notebook サーバーを安全に保つには、パスワードを設定し、サーバーの SSL 証明書を作成することをお勧めします。パスワードと SSL を設定するには、まず [DLAMI インスタンスに接続](#)し、次の手順に従います。

Jupyter Notebook サーバーを保護するには

1. Jupyter にはパスワードユーティリティが用意されています。次のコマンドを実行して、プロンプトに任意のパスワードを入力します。

```
$ jupyter notebook password
```

出力は次のようになります。

```
Enter password:
Verify password:
[NotebookPasswordApp] Wrote hashed password to /home/ubuntu/.jupyter/
jupyter_notebook_config.json
```

2. 自己署名 SSL 証明書を作成します。プロンプトに従って、該当するローカリティを入力します。プロンプトを空白のままにする場合は、`.`を入力する必要があります。この回答は証明書の機能に影響を及ぼしません。

```
$ cd ~
$ mkdir ssl
$ cd ssl
$ openssl req -x509 -nodes -days 365 -newkey rsa:2048 -keyout mykey.key -out
mycert.pem
```

Note

サードパーティが署名した定期的 SSL 証明書を作成して、ブラウザでセキュリティ警告が発生しないようにすることもできます。このプロセスは非常に複雑になります。詳細について

は、Jupyter Notebook ユーザードキュメントの「[Securing a notebook server](#)」を参照してください。

次のステップ

[DLAMI インスタンスでの Jupyter Notebook サーバーの起動](#)

DLAMI インスタンスでの Jupyter Notebook サーバーの起動

[Jupyter Notebook サーバーをパスワードと SSL で保護](#)したら、サーバーを起動できます。DLAMI インスタンスにログインし、以前に作成した SSL 証明書を使用する次のコマンドを実行します。

```
$ jupyter notebook --certfile=~/.ssl/mycert.pem --keyfile ~/.ssl/mykey.key
```

サーバーが起動していれば、クライアントコンピュータから SSH トンネルを通じてサーバーに接続できます。サーバーを実行すると、サーバーが実行中であることを確認するメッセージが Jupyter から出力されます。トンネルを作成するまでは機能しないため、このとき表示されるローカルホスト URL 経由でサーバーにアクセスできるという吹き出しは無視してください。

Note

Jupyter ウェブインターフェイスを使用してフレームワークを切り替えるとき、Jupyter は自動的に環境に切り替えを処理します。詳細については、「[Jupyter で環境を切り替える](#)」を参照してください。

次のステップ

[DLAMI インスタンスで Jupyter Notebook サーバーにクライアントを接続する](#)

DLAMI インスタンスで Jupyter Notebook サーバーにクライアントを接続する

[DLAMI インスタンスで Jupyter Notebook サーバーを起動](#)したら、サーバーに接続するように Windows、macOS、または Linux の各クライアントを設定します。接続したら、サーバー上のワークスペース内にノートブックを作成して、そのノートブックにアクセスしたり、サーバー上で深層学習コードを実行したりできます。

前提条件

次が手元に揃っていることを確認します。これらは、SSH トンネルを設定する際に必要です。

- Amazon EC2 インスタンスのパブリック DNS 名。詳細については、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 インスタンスホスト名のタイプ](#)」を参照してください。
- プライベートキーファイルのキーペア。キーペアへのアクセスに関する詳細は、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 のキーペアと Amazon EC2 インスタンス](#)」を参照してください。

Windows、macOS、または Linux の各クライアントから接続する

Windows、macOS、または Linux の各クライアントから DLAMI インスタンスに接続するには、クライアントのオペレーティングシステムの手順に従います。

Windows

SSH を使用して Windows クライアントから DLAMI インスタンスに接続するには

1. PuTTY などの Windows 用 SSH クライアントを使用します。手順については、「Amazon EC2 ユーザーガイド」の「[PuTTY を使用して Linux インスタンスに接続する](#)」を参照してください。その他の SSH 接続オプションについては、「[Connect to your Linux instance using SSH](#)」を参照してください。
2. (オプション) 実行中の Jupyter サーバーへの SSH トンネルを作成します。Windows クライアントに Git Bash をインストールし、macOS および Linux の各クライアントの接続手順に従います。

macOS or Linux

SSH を使用して Linux または macOS の各クライアントから DLAMI インスタンスに接続するには

1. ターミナルを開きます。
2. ローカルポート 8888 に対するすべてのリクエストをリモート Amazon EC2 インスタンスのポート 8888 に転送するために、次のコマンドを実行します。Amazon EC2 インスタンスにアクセスするキーの場所と Amazon EC2 インスタンスのパブリック DNS 名を置き換えて、コマンドを更新します。Amazon Linux AMI の場合、ユーザー名は ubuntu の代わりに ec2-user であることに注意してください。

```
$ ssh -i ~/mykeypair.pem -N -f -L 8888:localhost:8888 ubuntu@ec2-###-##-##-###.compute-1.amazonaws.com
```

このコマンドを実行すると、Jupyter Notebook サーバーを実行しているリモート Amazon EC2 インスタンスとクライアントの間にトンネルが開通します。

次のステップ

[DLAMI インスタンスでの Jupyter Notebook サーバーへのログイン](#)

DLAMI インスタンスでの Jupyter Notebook サーバーへのログイン

[クライアントを DLAMI インスタンスの Jupyter Notebook サーバーに接続](#)したら、サーバーにログインできます。

ブラウザでサーバーにログインするには

1. ブラウザのアドレスバーに次の URL を入力するか、次のリンクをクリックします。<https://localhost:8888>
2. 自己署名 SSL 証明書を使用すると、ブラウザは警告を表示し、このウェブサイトの閲覧を続行しないように求められます。



Your connection is not private

Attackers might be trying to steal your information from **localhost** (for example, passwords, messages, or credit cards). [Learn more](#)

NET::ERR_CERT_AUTHORITY_INVALID

Help improve Safe Browsing by sending some [system information and page content](#) to Google.
[Privacy policy](#)



Back to safety

この設定は自身が行ったものであるため、安全に続行できます。ブラウザによっては、「高度な設定」、「詳細の表示」などのボタンが表示されることがあります。



Your connection is not private

Attackers might be trying to steal your information from **localhost** (for example, passwords, messages, or credit cards). [Learn more](#)

NET::ERR_CERT_AUTHORITY_INVALID

Help improve Safe Browsing by sending some [system information and page content](#) to Google.
[Privacy policy](#)

Hide advanced

Back to safety

This server could not prove that it is **localhost**; its security certificate is not trusted by your computer's operating system. This may be caused by a misconfiguration or an attacker intercepting your connection.

[Proceed to localhost \(unsafe\)](#)

これをクリックして「localhostに進む」リンクをクリックします。接続に成功すると、Jupyter Notebook サーバーのウェブページが表示されます。この時点で、前に設定したパスワードを入力するように求められます。

これで、DLAMI インスタンスで実行中の Jupyter Notebook サーバーにアクセスできるようになりました。新しいノートブックを作成することも、用意されている [チュートリアル](#) を実行することもできます。

DLAMI インスタンスのクリーンアップ

DLAMI インスタンスが不要になった場合は、Amazon EC2 でそのインスタンスを停止または終了させることによって予想外の料金の発生を避けられます。

インスタンスを停止する場合は、インスタンスをそのまま保持して、後で再度使用するときに開始できます。設定、ファイル、およびその他の非揮発性情報は Amazon Simple Storage Service (Amazon S3) のボリュームに格納されます。インスタンスが停止している間は、ボリュームの保持に対して S3 料金が発生しますが、コンピューティングリソースには料金が発生しません。再度インスタンスを起動すると、データを含むそのストレージボリュームがマウントされます。

インスタンスを終了した場合は、インスタンスは消去され、再度起動することはできません。もちろん、終了したインスタンスのコンピューティングリソースに対してそれ以上の料金は発生しません。ただし、データは引き続き Amazon S3 に残り、引き続き S3 料金が発生する可能性があります。終了したインスタンスに関連して、今後いかなる料金も発生させないためには、Amazon S3 のストレージボリュームも削除する必要があります。手順については、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 インスタンスを終了する](#)」を参照してください。

stopped や terminated などの Amazon EC2 インスタンスの状態の詳細については、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 インスタンスの状態変更](#)」を参照してください。

DLAMI の使用

トピック

- [Deep Learning AMI with Conda の使用](#)
- [Deep Learning Base AMI の使用](#)
- [Jupyter ノートブックチュートリアルを実行する](#)
- [チュートリアル](#)

以降のセクションでは、Deep Learning AMI with Conda を使用して環境を切り替える方法、各フレームワークからサンプルコードを実行する方法、および Jupyter を実行して、さまざまなノートブックチュートリアルを試す方法について説明します。

Deep Learning AMI with Conda の使用

トピック

- [Deep Learning AMI with Conda の概要](#)
- [DLAMI にログインする](#)
- [TensorFlow 環境を開始する](#)
- [PyTorch Python 3 環境に切り替える](#)
- [環境を削除する](#)

Deep Learning AMI with Conda の概要

Conda は、Windows、macOS、および Linux で稼働するオープンソースのパッケージ管理システムおよび環境管理システムです。Conda では、パッケージとその依存関係を迅速にインストール、実行、および更新できます。また、ローカルコンピュータ上で簡単に環境を作成、保存、ロードし、環境を切り替えることができます。

Deep Learning AMI with Conda は、深層学習環境を簡単に切り替えられるように設定されています。次に示す手順では、conda での基本的なコマンドを説明します。これらのコマンドは、フレームワークの基本インポートが機能していることや、フレームワークでいくつかの簡単な操作を実行できることの確認にも役立ちます。その後、DLAMI で提供されているより詳細なチュートリアルや、各フレームワークのプロジェクトサイトにあるフレームワークの例に進むことができます。

DLAMI にログインする

サーバーにログインすると、サーバーの「その日のメッセージ (MOTD)」が表示され、さまざまなディープラーニングフレームワークを切り替えるための各種の Conda コマンドが示されます。以下に MOTD の例を示します。新しいバージョンの DLAMI がリリースされる事により、場合によっては特定の MOTD が異なる場合があります。

```
=====
AMI Name: Deep Learning OSS Nvidia Driver AMI (Amazon Linux 2) Version 77
Supported EC2 instances: G4dn, G5, G6, Gr6, P4d, P4de, P5
  * To activate pre-built tensorflow environment, run: 'source activate
tensorflow2_p310'
  * To activate pre-built pytorch environment, run: 'source activate
pytorch_p310'
  * To activate pre-built python3 environment, run: 'source activate python3'

NVIDIA driver version: 535.161.08

CUDA versions available: cuda-11.7 cuda-11.8 cuda-12.0 cuda-12.1 cuda-12.2

Default CUDA version is 12.1

Release notes: https://docs.aws.amazon.com/dlami/latest/devguide/appendix-ami-release-notes.html
AWS Deep Learning AMI Homepage: https://aws.amazon.com/machine-learning/amis/
Developer Guide and Release Notes: https://docs.aws.amazon.com/dlami/latest/devguide/what-is-dlami.html
Support: https://forums.aws.amazon.com/forum.jspa?forumID=263
For a fully managed experience, check out Amazon SageMaker at https://aws.amazon.com/sagemaker
=====
```

TensorFlow 環境を開始する

Note

初めて Conda 環境を起動する際には、ロードするまで辛抱して待機してください。Deep Learning AMI with Conda は、EC2 インスタンスにフレームワークの最初のアクティベーションから最適化されたフレームワークのバージョンを自動的にインストールします。継続する遅延はありません。

1. Python 3 用の TensorFlow 仮想環境を有効化します。

```
$ source activate tensorflow2_p310
```

2. iPython ターミナルを起動します。

```
(tensorflow2_p310)$ ipython
```

3. クイック TensorFlow プログラムを実行します。

```
import tensorflow as tf
hello = tf.constant('Hello, TensorFlow!')
sess = tf.Session()
print(sess.run(hello))
```

"Hello, Tensorflow!" が表示されます。

次回の予定

[Jupyter ノートブックチュートリアルを実行する](#)

PyTorch Python 3 環境に切り替える

iPython コンソールをまだ使用している場合は、`quit()` を使用して、環境を切り替える準備を整えます。

- Python 3 用の PyTorch 仮想環境を有効化します。

```
$ source activate pytorch_p310
```

PyTorch コードをテストする

インストール環境をテストするには、配列を作成および出力する PyTorch コードを Python で作成します。

1. iPython ターミナルを起動します。

```
(pytorch_p310)$ ipython
```

2. PyTorch をインポートします。

```
import torch
```

サードパーティー製パッケージに関する警告メッセージが表示される場合があります。このメッセージは無視できます。

3. 各要素をランダムに初期化して、5x3 行列を作成します。配列を出力します。

```
x = torch.rand(5, 3)
print(x)
```

結果を確認します。

```
tensor([[0.3105, 0.5983, 0.5410],
        [0.0234, 0.0934, 0.0371],
        [0.9740, 0.1439, 0.3107],
        [0.6461, 0.9035, 0.5715],
        [0.4401, 0.7990, 0.8913]])
```

環境を削除する

DLAMI にスペースが足りない場合、使用していない Conda パッケージをアンインストールすることを選択できます。

```
conda env list
conda env remove --name <env_name>
```

Deep Learning Base AMI の使用

Deep Learning Base AMI の使用

Base AMI には、独自でカスタマイズしたディープラーニング環境をデプロイするための GPU ドライバの基盤プラットフォームとアクセラレーションライブラリが含まれています。AMI はデフォルトで、任意の 1 つの NVIDIA CUDA バージョン環境で設定されています。CUDA のさまざまなバージョン間で切り替えを行うこともできます。これを行う方法については、次の手順を参照してください。

CUDA のバージョンの設定

NVIDIA の `nvcc` プログラムを実行して、CUDA バージョンを検証できます。

```
nvcc --version
```

次の Bash コマンドを使用して、特定の CUDA バージョンを選択して検証できます。

```
sudo rm /usr/local/cuda
sudo ln -s /usr/local/cuda-12.0 /usr/local/cuda
```

詳細については、「[Base DLAMI リリースノート](#)」を参照してください。

Jupyter ノートブックチュートリアルを実行する

各深層学習プロジェクトにはチュートリアルと例が付属しています。ほとんどの場合、それらはすべての DLAMI で実行することができます。[Deep Learning AMI with Conda](#) を選択した場合は、厳選されたチュートリアルが既にセットアップされ、すぐに試すことができるという利点があります。

Important

DLAMI にインストールされている Jupyter ノートブックチュートリアルを実行するには、[DLAMI インスタンスでの Jupyter Notebook サーバーの設定](#)の内容を実行する必要があります。

Jupyter サーバー実行されたら、ウェブブラウザからチュートリアルを実行できます。Deep Learning AMI with Conda を実行している場合や、Python 環境をセットアップしている場合は、Jupyter ノートブックのインターフェイスから Python カーネルを切り替えることができます。フレームワーク固有のチュートリアルを実行する前に、適切なカーネルを選択してください。これに関するその他の例が Deep Learning AMI with Conda のユーザー向けに提供されています。

Note

多くのチュートリアルでは追加の Python モジュールが必要になりますが、それらが DLAMI にセットアップされていない可能性があります。"xyz module not found" などのエラーが発生した場合は、DLAMI にログインし、上記を参考に環境をアクティブ化して、必要なモジュールをインストールしてください。

i Tip

ディープラーニングのチュートリアルや例の多くは、1つ以上の GPU を必要とします。GPU を使用しないインスタンスタイプの場合、例のコードを実行するために、コードを一部変更することが必要になる可能性があります。

インストールされたチュートリアルを操作する

Jupyter サーバーにログインしてチュートリアルディレクトリが表示されると (Deep Learning AMI with Conda 上のみ)、各フレームワークの名前別に整理されたチュートリアルのフォルダを確認できます。フレームワークのリストが表示されない場合は、現在の DLAMI で、そのフレームワークのチュートリアルが提供されていないことを示しています。フレームワークの名前をクリックし、リストされているチュートリアルを確認して、チュートリアルをクリックして起動します。

Deep Learning AMI with Conda でノートブックを初めて実行すると、使用する環境を指定するように求められます。選択のためのリストがプロンプトで表示されます。それぞれの環境は、以下のパターンに従って命名されています。

Environment (conda_framework_python-version)

たとえば、"Environment (conda_mxnet_p36)" と表示された場合は、環境に MXNet と Python 3 が存在しています。別のバリエーションとして "Environment (conda_mxnet_p27)" があります。この場合は環境に MXNet と Python 2 が存在しています。

i Tip

アクティブ化される CUDA のバージョンについて不明な点がある場合は、最初に DLAMI にログインする際の MOTD で確認できます。

Jupyter で環境を切り替える

別のフレームワークのチュートリアルを試す場合は、現在実行中のカーネルを検証する必要があります。この情報は、Jupyter のインターフェイスの右上、ログアウトボタンの真下に表示されます。開いているノートブックでカーネルを変更するには、Jupyter のメニュー項目を [Kernel]、[Change Kernel] の順にクリックして、実行中のノートブックに適合する環境をクリックします。

カーネルに変更が加わると、これまでに実行したすべての項目の状態が消去されるため、この時点ですべてのセルを再度実行することが必要になります。

Tip

フレームワークの切り替えは、興味深くさまざまな情報も入手できる操作ですが、メモリ不足が発生する可能性があります。エラーが表示されるようになったら、実行中の Jupyter サーバーが表示されているターミナルウィンドウを確認してください。参考になるメッセージやエラーログをここで確認でき、メモリ不足のエラーもここに表示されます。この問題を修正するには、Jupyter サーバーのホームページにアクセスして、バックグラウンドで実行中でメモリを占有していると思われる各チュートリアルで、[Running] タブ、[Shutdown] の順にクリックします。

チュートリアル

Deep Learning AMI with Conda のソフトウェアの使用法のチュートリアルを以下に示します。

トピック

- [フレームワークのアクティブ化](#)
- [Elastic Fabric Adapter を使用した分散トレーニング](#)
- [GPU のモニタリングおよび最適化](#)
- [DLAMI を使用した AWS Inferentia チップ](#)
- [ARM64 DLAMI](#)
- [推論](#)
- [モデル提供](#)

フレームワークのアクティブ化

以下に示しているのは、Deep Learning AMI with Conda にインストールされている深層学習フレームワークです。フレームワークをクリックすると、そのフレームワークをアクティブ化する方法を参照できます。

トピック

- [PyTorch](#)
- [TensorFlow 2](#)

PyTorch

PyTorch の有効化

フレームワークの安定版 Conda パッケージがリリースされると、DLAMI でテストされ事前にインストールされます。テストされていない最新のナイトリービルドを実行する場合は、手動で「[PyTorch のナイトリービルド \(実験段階\) をインストールする](#)」ことができます。

現在インストールされているフレームワークをアクティブ化するには、使用する Deep Learning AMI with Conda に関する以下の手順に従います。

Python 3 ベースの PyTorch を CUDA と MKL-DNN で使用する場合は、以下のコマンドを実行します。

```
$ source activate pytorch_p310
```

iPython ターミナルを起動します。

```
(pytorch_p310)$ ipython
```

クイック PyTorch プログラムを実行します。

```
import torch
x = torch.rand(5, 3)
print(x)
print(x.size())
y = torch.rand(5, 3)
print(torch.add(x, y))
```

最初のランダムな配列が出力された後、そのサイズが出力され、次に別のランダムな配列が追加で出力されます。

PyTorch のナイトリービルド (実験段階) をインストールする

ナイトリービルドから PyTorch をインストールする方法

最新の PyTorch ビルドを Deep Learning AMI with Conda の PyTorch Conda 環境のいずれかまたは両方にインストールできます。

- (Python 3 用オプション) - Python 3 PyTorch 環境を有効化します。

```
$ source activate pytorch_p310
```

2. 残りの手順は、pytorch_p310 環境を使用していることを前提としています。現在インストールされている PyTorch を削除します。

```
(pytorch_p310)$ pip uninstall torch
```

3.
 - (GPU インスタンス用オプション) - PyTorch と CUDA.0 の最新ナイトリービルドをインストールします:

```
(pytorch_p310)$ pip install torch_nightly -f https://download.pytorch.org/whl/nightly/cu100/torch_nightly.html
```

- (CPU インスタンス用オプション) - GPU を持たないインスタンスの場合は PyTorch の最新ナイトリービルドをインストールします:

```
(pytorch_p310)$ pip install torch_nightly -f https://download.pytorch.org/whl/nightly/cpu/torch_nightly.html
```

4. 最新のナイトリービルドが正常にインストールされたことを確認するには、IPython ターミナルを起動して PyTorch のバージョンを確認します。

```
(pytorch_p310)$ ipython
```

```
import torch
print (torch.__version__)
```

出力は 1.0.0.dev20180922 のように表示されます。

5. PyTorch ナイトリービルドが MNIST のサンプルで正常に機能することを確認するには、PyTorch のサンプルリポジトリからテストスクリプトを実行します。

```
(pytorch_p310)$ cd ~
(pytorch_p310)$ git clone https://github.com/pytorch/examples.git pytorch_examples
(pytorch_p310)$ cd pytorch_examples/mnist
(pytorch_p310)$ python main.py || exit 1
```

他のチュートリアル

他のチュートリアルおよび例については、フレームワークの公式ドキュメントである [PyTorch ドキュメント](#)と [PyTorch](#) のウェブサイトを参照してください。

TensorFlow 2

このチュートリアルでは、Deep Learning AMI with Conda (DLAMI on Conda) を実行しているインスタンスで TensorFlow 2 を有効化し、TensorFlow 2 プログラムを実行する方法を示します。

フレームワークの安定版 Conda パッケージがリリースされると、DLAMI でテストされ事前にインストールされます。

TensorFlow 2 のアクティブ化

Conda の DLAMI 上で TensorFlow を実行するには

1. TensorFlow 2 を有効化するには、DLAMI with Conda の Amazon Elastic Compute Cloud (Amazon EC2) インスタンスを開きます。
2. Python 3 ベースの TensorFlow 2 および Keras 2 を CUDA 10.1 および MKL-DNN で使用する場合は、以下のコマンドを実行します。

```
$ source activate tensorflow2_p310
```

3. iPython ターミナルを起動します。

```
(tensorflow2_p310)$ ipython
```

4. TensorFlow 2 プログラムを実行して、正常に作動していることを確認します。

```
import tensorflow as tf
hello = tf.constant('Hello, TensorFlow!')
tf.print(hello)
```

Hello, TensorFlow! が画面に表示される必要があります。

他のチュートリアル

その他のチュートリアルや例については、[TensorFlow Python API](#) の TensorFlow ドキュメントを参照するか、[TensorFlow](#) ウェブサイトを参照してください。

Elastic Fabric Adapter を使用した分散トレーニング

[Elastic Fabric Adapter](#) (EFA) は、ハイパフォーマンスコンピューティング (HPC) アプリケーションを高速化するために DLAMI インスタンスにアタッチできるネットワークデバイスです。EFA を使用すると、AWS クラウドが提供するスケーラビリティ、柔軟性、伸縮性で、オンプレミス HPC クラスターのアプリケーションパフォーマンスを実現できます。

以下のトピックでは、EFA と DLAMI の使用を開始する方法について説明します。

Note

この [Base GPU DLAMI リスト](#) からお使いの DLAMI を選択してください。

トピック

- [EFA を使用した AWS Deep Learning AMIs インスタンスの起動](#)
- [DLAMI での EFA の使用](#)

EFA を使用した AWS Deep Learning AMIs インスタンスの起動

最新版 Base DLAMI は EFA と共に使用する準備ができており、GPU インスタンス用の必要なドライバー、カーネルモジュール、Libfabric、OpenMPI および [NCCL OFI プラグイン](#) が付属しています。

サポートされている Base DLAMI の CUDA バージョンは、[リリースノート](#) で確認できます。

メモ:

- EFA 上の mpirun を使用して NCCL アプリケーションを実行する際、EFA がサポートされているインストレーションへの完全なパスを以下のように指定する必要があります。

```
/opt/amazon/openmpi/bin/mpirun <command>
```

- アプリケーションで EFA を使用するには、[DLAMI での EFA の使用](#) に示すように、mpirun コマンドに FI_PROVIDER="efa" を追加します。

トピック

- [EFA 対応のセキュリティグループを準備する](#)
- [インスタンスの起動](#)
- [EFA 添付ファイルの確認](#)

EFA 対応のセキュリティグループを準備する

EFA には、セキュリティグループ自体とのインバウンドおよびアウトバウンドのトラフィックをすべて許可するセキュリティグループが必要です。詳細については、「[EFA のドキュメント](#)」を参照してください。

1. Amazon EC2 コンソールの <https://console.aws.amazon.com/ec2/> を開いてください。
2. ナビゲーションペインで [セキュリティグループ] を選択して、[セキュリティグループの作成] を選択します。
3. [セキュリティグループの作成] ウィンドウで、以下を行います。
 - [セキュリティグループ名] に、EFA-enabled security group のような、分かりやすいセキュリティグループ名を入力してください。
 - (オプション) [説明] に、セキュリティグループの簡単な説明を入力してください。
 - [VPC] で、EFA 対応のインスタンスを起動する VPC を選択します。
 - [作成] を選択します。
4. 作成したセキュリティグループを選択し、[説明] タブで [グループ ID] をコピーします。
5. [Inbound] (インバウンド) タブおよび [Outbound] (アウトバウンド) タブで、次の手順を実行します。
 - [編集] を選択します。
 - [Type] で、[All traffic] を選択します。
 - [ソース] で [カスタム] を選択します。
 - コピーしたセキュリティグループ ID をフィールドに貼り付けます。
 - [保存] を選択します。

6. 「[Linux インスタンスのインバウンドトラフィックの承認](#)」を参照するインバウンドトラフィックを有効にします。このステップを抜かすと、DLAMI インスタンスと通信できなくなります。

インスタンスの起動

の EFA AWS Deep Learning AMIs は現在、次のインスタンスタイプとオペレーティングシステムでサポートされています。

- P3dn: Amazon Linux 2、Ubuntu 20.04
- P4d, P4de: Amazon Linux 2、Amazon Linux 2023、Ubuntu 20.04、Ubuntu 22.04
- P5, P5e, P5en: Amazon Linux 2、Amazon Linux 2023、Ubuntu 20.04、Ubuntu 22.04

以下のセクションでは、EFA 対応の DLAMI インスタンスを起動する方法について説明します。EFA 対応のインスタンスの起動の詳細については、[クラスタープレイスメントグループで EFA 対応のインスタンスを起動する](#)を参照してください。

1. Amazon EC2 コンソール (<https://console.aws.amazon.com/ec2/>) を開きます。
2. [インスタンスの作成] を選択してください。
3. AMI を選択ページで、[DLAMI リリースノートページ](#)にある、サポートされている DLAMI を選択します。
4. [Choose an Instance Type] (インスタンスタイプの選択) ページで、次のいずれかのサポート対象のインスタンスタイプを選択し、[Next: Configure Instance Details] (次の手順: インスタンスの詳細の設定) を選択します。サポートされているインスタンスのリストについては、[EFA と MPI の開始方法](#)のリンクを参照してください。
5. [Configure Instance Details] ページで以下の操作を実行します。
 - [インスタンス数] に、起動する EFA 対応のインスタンスの数を入力します。
 - [ネットワーク] および [サブネット] で、インスタンスを起動する VPC およびサブネットを選択します。
 - [オプション][プレイスメントグループ] で、[プレイスメントグループにインスタンスを追加] を選択します。最適なパフォーマンスを得るには、プレイスメントグループ内でインスタンスを起動します。
 - [オプション][プレイスメントグループ名] で、[新しいプレイスメントグループに追加] を選択し、プレイスメントグループの分かりやすい名前を入力して、[プレイスメントグループ戦略] で [クラスター] を選択します。

- 必ず、このページで [Elastic Fabric Adapter] を有効にしてください。このオプションが無効になっている場合は、選択したインスタンスタイプに対応するサブネットに変更します。
 - [ネットワークインターフェイス] セクションの [eth0] で、[新しいネットワークインターフェイス] を選択してください。必要に応じて、プライマリ IPv4 アドレスと 1 つ以上のセカンダリ IPv4 アドレスを指定できます。関連付けられている IPv6 CIDR ブロックを持つサブネットにインスタンスを起動する場合は、必要に応じて、プライマリ IPv6 アドレスと 1 つ以上のセカンダリ IPv6 アドレスを指定することができます。
 - [次の手順: ストレージの追加] を選択します。
6. [ストレージの追加] ページで、AMI で指定されたボリュームに加えてインスタンスにアタッチするボリューム (例: ルートデバイスのボリューム) を指定し、[Next: Add Tags (次へ: タグの追加)] を選択します。
 7. [Add Tags] ページで、ユーザーフレンドリーな名前などを使ってインスタンスのタグを指定し、[Next: Configure Security Group] を選択します。
 8. [セキュリティグループの設定] ページの [セキュリティグループの割り当て] で、[既存のセキュリティグループを選択する] を選択し、前に作成したセキュリティグループを選択します。
 9. [Review and Launch] を選択してください。
 10. [インスタンス作成の確認] ページで設定を確認し、[起動] を選択してキーペアを選択し、インスタンスを起動します。

EFA 添付ファイルの確認

コンソールから

インスタンスを起動したら、AWS コンソールでインスタンスの詳細を確認します。これを行うには、EC2 コンソールでインスタンスを選択し、ページ下のペインにある [Description (説明)] タブを確認します。[Network Interfaces: eth0 (ネットワークインターフェイス: eth0)] というパラメータを探し、eth0 をクリックするとポップアップが開きます。[Elastic Fabric Adapter] が有効になっていることを確認します。

EFA が有効になっていない場合は、次のいずれかの方法でこれを修正できます。

- EC2 インスタンスを終了し、同じ手順で新しいインスタンスを起動します。EFA が添付されていることを確認します。
- 既存のインスタンスに EFA を添付します。
 1. EC2 コンソールで、[Network Interfaces (ネットワークインターフェイス)] に移動します。

2. [Create a Network Interface (ネットワークインターフェイスを作成)] をクリックします。
3. インスタンスが入っている同じサブネットを選択します。
4. 必ず、[Elastic Fabric Adapter] を有効にし、[Create] (作成) をクリックします。
5. [EC2 Instances (EC2 インスタンス)] タブに戻り、インスタンスを選択します。
6. [Actions: Instance State] (アクション: インスタンスの状態) に移動し、EFA にアタッチする前にインスタンスを停止します。
7. [Actions (アクション)] から、[Networking: Attach Network Interface (ネットワーキング: ネットワークインターフェイスの接続)] を選択します。
8. 先ほど作成したインターフェイスを選択し、[Attach (接続)] をクリックします。
9. インスタンスを再起動します。

インスタンスから

DLAMI にすでに、次のテストスクリプトが存在します。これを実行して、カーネルモジュールが正しくロードされていることを確認します。

```
$ fi_info -p efa
```

出力は以下のようになります。

```
provider: efa
  fabric: EFA-fe80::e5:56ff:fe34:56a8
  domain: efa_0-rdm
  version: 2.0
  type: FI_EP_RDM
  protocol: FI_PROTO_EFA
provider: efa
  fabric: EFA-fe80::e5:56ff:fe34:56a8
  domain: efa_0-dgrm
  version: 2.0
  type: FI_EP_DGRAM
  protocol: FI_PROTO_EFA
provider: efa;ofi_rxd
  fabric: EFA-fe80::e5:56ff:fe34:56a8
  domain: efa_0-dgrm
  version: 1.0
  type: FI_EP_RDM
  protocol: FI_PROTO_RXD
```

セキュリティグループ構成の確認

DLAMI にすでに、次のテストスクリプトが存在します。これを実行して、作成したセキュリティグループが正しく設定されていることを確認します。

```
$ cd /opt/amazon/efa/test/  
$ ./efa_test.sh
```

出力は以下のようになります。

```
Starting server...  
Starting client...  
bytes  #sent  #ack  total  time  MB/sec  usec/xfer  Mxfers/sec  
64      10     =10   1.2k   0.02s  0.06    1123.55    0.00  
256     10     =10   5k     0.00s  17.66   14.50     0.07  
1k      10     =10   20k    0.00s  67.81   15.10     0.07  
4k      10     =10   80k    0.00s  237.45  17.25     0.06  
64k     10     =10   1.2m   0.00s  921.10  71.15     0.01  
1m      10     =10   20m    0.01s  2122.41 494.05    0.00
```

応答しなかったり完了しない場合は、セキュリティグループに正しいインバウンド/アウトバウンドルールが設定されていることを確認します。

DLAMI での EFA の使用

次の項では、EFA を使用して、AWS Deep Learning AMIs でマルチノードアプリケーションを実行する方法について説明します。

EFA でマルチノードアプリケーションを実行

ノードのクラスター全体でアプリケーションを実行するには、次の構成が必要です。

トピック

- [パスワードレス SSH の有効化](#)
- [Hosts ファイルの作成](#)
- [NCCL テスト](#)

パスワードレス SSH の有効化

クラスター内の 1 つのノードをリーダーノードとして選択します。残りのノードはメンバーノードと呼ばれます。

1. リーダーノードで、RSA キーペアを生成します。

```
ssh-keygen -t rsa -N "" -f ~/.ssh/id_rsa
```

2. リーダーノードのプライベートキーの許可を変更します。

```
chmod 600 ~/.ssh/id_rsa
```

3. パブリックキー ~/.ssh/id_rsa.pub をコピーし、クラスター内のメンバーノードの ~/.ssh/authorized_keys に追加します。
4. これで、プライベート IP を使用して、リーダーノードからメンバーノードに直接ログインできるようになります。

```
ssh <member private ip>
```

5. strictHostKeyChecking を無効にし、リーダーノードの ~/.ssh/config ファイルに以下を追加することで、リーダーノードでエージェント転送を有効化します。

```
Host *  
    ForwardAgent yes  
Host *  
    StrictHostKeyChecking no
```

6. Amazon Linux 2 インスタンスでは、リーダーノードで次のコマンドを実行して、設定ファイルへの正しいアクセス許可を付与します。

```
chmod 600 ~/.ssh/config
```

Hosts ファイルの作成

リーダーノードで、クラスター内のノードを識別する Hosts ファイルを作成します。Hosts ファイルには、クラスター内の各ノードのエントリが必要です。~/hosts ファイルを作成し、次のようにプライベート IP を使用して各ノードを追加します。

```
localhost slots=8  
<private ip of node 1> slots=8  
<private ip of node 2> slots=8
```

NCCL テスト

Note

これらのテストは、EFA バージョン 1.38.0 と OFI NCCL プラグイン 1.13.2 を使用して実行されています。

以下は、複数のコンピューティングノードで機能とパフォーマンスの両方をテストするために Nvidia が提供する NCCL テストのサブセットです

サポートされているインスタンス: P3dn, P4, P5, P5e, P5en

パフォーマンステスト

P4d.24xlarge でのマルチノード NCCL パフォーマンステスト

EFA で NCCL パフォーマンステストを確認するには、公式の [NCCL-Tests Repo](#) で実施可能な標準 NCCL パフォーマンステストを実施します。DLAMI には、CUDA XX.X 用に構築されたこのテストが付属しています。EFA を使用して独自のスクリプトを実行することもできます。

独自のスクリプトを作成する場合は、次のガイダンスを参照してください。

- EFA で NCCL アプリケーションを実行中に、例に示すように、mpirun への完全なパスを使用します。
- クラスター内のインスタンスの数と GPU の数に基づいて、パラメータ np と N を変更します。
- NCCL_DEBUG=INFO フラグを追加し、ログの EFA 使用状況が、[Selected Provider is EFA] になっていることを確認します。
- 検証用に解析するトレーニングログの場所を設定します。

```
TRAINING_LOG="testEFA_$(date +%N).log"
```

いずれかのメンバーノードでコマンド `watch nvidia-smi` を使用し、GPU の使用状況を監視します。次の `watch nvidia-smi` コマンドは汎用 CUDA xx.x バージョン用のコマンドで、インスタンスのオペレーティングシステムによって異なります。スクリプト内の CUDA バージョンを置き換えることで、Amazon EC2 インスタンスで使用可能な任意の CUDA バージョンに対するコマンドを実行できます。

- Amazon Linux 2、Amazon Linux 2023:

```
$ /opt/amazon/openmpi/bin/mpirun -n 16 -N 8 \
-x NCCL_DEBUG=INFO --mca pml ^cm \
-x LD_LIBRARY_PATH=/usr/local/cuda-xx.x/efa/lib:/usr/local/cuda-xx.x/lib:/usr/
local/cuda-xx.x/lib64:/usr/local/cuda-xx.x:/opt/amazon/efa/lib64:/opt/amazon/openmpi/
lib64:$LD_LIBRARY_PATH \
--hostfile hosts --mca btl tcp,self --mca btl_tcp_if_exclude lo,docker0 --bind-to
none \
/usr/local/cuda-xx.x/efa/test-cuda-xx.x/all_reduce_perf -b 8 -e 1G -f 2 -g 1 -c 1 -n
100 | tee ${TRAINING_LOG}
```

- Ubuntu 20.04、Ubuntu 20.04:

```
$ /opt/amazon/openmpi/bin/mpirun -n 16 -N 8 \
-x NCCL_DEBUG=INFO --mca pml ^cm \
-x LD_LIBRARY_PATH=/usr/local/cuda-xx.x/efa/lib:/usr/local/cuda-xx.x/lib:/usr/
local/cuda-xx.x/lib64:/usr/local/cuda-xx.x:/opt/amazon/efa/lib:/opt/amazon/openmpi/
lib:$LD_LIBRARY_PATH \
--hostfile hosts --mca btl tcp,self --mca btl_tcp_if_exclude lo,docker0 --bind-to
none \
/usr/local/cuda-xx.x/efa/test-cuda-xx.x/all_reduce_perf -b 8 -e 1G -f 2 -g 1 -c 1 -n
100 | tee ${TRAINING_LOG}
```

出力は次のようになります。

```
# nThread 1 nGpus 1 minBytes 8 maxBytes 1073741824 step: 2(factor) warmup iters: 5
iters: 100 agg iters: 1 validation: 1 graph: 0
#
# Using devices
# Rank 0 Group 0 Pid 33378 on ip-172-31-42-25 device 0 [0x10] NVIDIA A100-
SXM4-40GB
# Rank 1 Group 0 Pid 33379 on ip-172-31-42-25 device 1 [0x10] NVIDIA A100-
SXM4-40GB
# Rank 2 Group 0 Pid 33380 on ip-172-31-42-25 device 2 [0x20] NVIDIA A100-
SXM4-40GB
# Rank 3 Group 0 Pid 33381 on ip-172-31-42-25 device 3 [0x20] NVIDIA A100-
SXM4-40GB
# Rank 4 Group 0 Pid 33382 on ip-172-31-42-25 device 4 [0x90] NVIDIA A100-
SXM4-40GB
# Rank 5 Group 0 Pid 33383 on ip-172-31-42-25 device 5 [0x90] NVIDIA A100-
SXM4-40GB
```


64	16	float	sum	-1	178.8	0.00	0.00	0
178.7	0.00	0.00	0					
128	32	float	sum	-1	178.2	0.00	0.00	0
177.8	0.00	0.00	0					
256	64	float	sum	-1	178.6	0.00	0.00	0
178.8	0.00	0.00	0					
512	128	float	sum	-1	177.2	0.00	0.01	0
177.1	0.00	0.01	0					
1024	256	float	sum	-1	179.2	0.01	0.01	0
179.3	0.01	0.01	0					
2048	512	float	sum	-1	181.3	0.01	0.02	0
181.2	0.01	0.02	0					
4096	1024	float	sum	-1	184.2	0.02	0.04	0
183.9	0.02	0.04	0					
8192	2048	float	sum	-1	191.2	0.04	0.08	0
190.6	0.04	0.08	0					
16384	4096	float	sum	-1	202.5	0.08	0.15	0
202.3	0.08	0.15	0					
32768	8192	float	sum	-1	233.0	0.14	0.26	0
232.1	0.14	0.26	0					
65536	16384	float	sum	-1	238.6	0.27	0.51	0
235.1	0.28	0.52	0					
131072	32768	float	sum	-1	237.2	0.55	1.04	0
236.8	0.55	1.04	0					
262144	65536	float	sum	-1	248.3	1.06	1.98	0
247.0	1.06	1.99	0					
524288	131072	float	sum	-1	309.2	1.70	3.18	0
307.7	1.70	3.20	0					
1048576	262144	float	sum	-1	408.7	2.57	4.81	0
404.3	2.59	4.86	0					
2097152	524288	float	sum	-1	613.5	3.42	6.41	0
607.9	3.45	6.47	0					
4194304	1048576	float	sum	-1	924.5	4.54	8.51	0
914.8	4.58	8.60	0					
8388608	2097152	float	sum	-1	1059.5	7.92	14.85	0
1054.3	7.96	14.92	0					
16777216	4194304	float	sum	-1	1269.9	13.21	24.77	0
1272.0	13.19	24.73	0					
33554432	8388608	float	sum	-1	1642.7	20.43	38.30	0
1636.7	20.50	38.44	0					
67108864	16777216	float	sum	-1	2446.7	27.43	51.43	0
2445.8	27.44	51.45	0					
134217728	33554432	float	sum	-1	4143.6	32.39	60.73	0
4142.4	32.40	60.75	0					

```

268435456    67108864    float    sum    -1    7351.9    36.51    68.46    0
7346.7    36.54    68.51    0
536870912    134217728    float    sum    -1    13717    39.14    73.39    0
13703    39.18    73.46    0
1073741824    268435456    float    sum    -1    26416    40.65    76.21    0
26420    40.64    76.20    0
...
# Out of bounds values : 0 OK
# Avg bus bandwidth    : 15.5514

```

検証テスト

EFA テストが有効な結果を返したことを検証するには、次のテストを使用して確認します。

- EC2 インスタンスメタデータを使用してインスタンスタイプを取得します。

```

TOKEN=$(curl -X PUT "http://169.254.169.254/latest/api/token" -H "X-aws-ec2-metadata-token-ttl-seconds: 21600")
INSTANCE_TYPE=$(curl -H "X-aws-ec2-metadata-token: $TOKEN" -v http://169.254.169.254/latest/meta-data/instance-type)

```

- [パフォーマンステスト](#) を実行します。
- 以下のパラメータを設定します。

```

CUDA_VERSION
CUDA_RUNTIME_VERSION
NCCL_VERSION

```

- 結果を次のように検証します。

```

RETURN_VAL=`echo $?`
if [ ${RETURN_VAL} -eq 0 ]; then

    # [0] NCCL INFO NET/OFI Initializing aws-ofi-nccl 1.13.2-aws
    # [0] NCCL INFO NET/OFI Using CUDA driver version 12060 with runtime 12010

    # cudaDriverVersion 12060 --> This is max supported cuda version by nvidia driver
    # NCCL version 2.23.4+cuda12.5 --> This is NCCL version compiled with cuda version

    # Validation of logs

```



```

grep "NET/OFI Configuring AWS-specific options" ${TRAINING_LOG} || { echo "AWS-
specific options text not found"; exit 1; }
grep "busbw" ${TRAINING_LOG} || { echo "busbw text not found"; exit 1; }
grep "Avg bus bandwidth " ${TRAINING_LOG} || { echo "Avg bus bandwidth text not
found"; exit 1; }
grep "NCCL version $NCCL_VERSION" ${TRAINING_LOG} || { echo "Text not found: NCCL
version $NCCL_VERSION"; exit 1; }
if [[ ${INSTANCE_TYPE} == "p4d.24xlarge" ]]; then
    grep "NET/Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Text not found: NET/
Libfabric/0/GDRDMA"; exit 1; }
    grep "NET/OFI Selected Provider is efa (found 4 nics)" ${TRAINING_LOG} ||
{ echo "Selected Provider is efa text not found"; exit 1; }
    elif [[ ${INSTANCE_TYPE} == "p4de.24xlarge" ]]; then
        grep "NET/Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Avg bus bandwidth
text not found"; exit 1; }
        grep "NET/OFI Selected Provider is efa (found 4 nics)" ${TRAINING_LOG} ||
{ echo "Avg bus bandwidth text not found"; exit 1; }
        elif [[ ${INSTANCE_TYPE} == "p5.48xlarge" ]]; then
            grep "NET/Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Avg bus bandwidth
text not found"; exit 1; }
            grep "NET/OFI Selected Provider is efa (found 32 nics)" ${TRAINING_LOG} ||
{ echo "Avg bus bandwidth text not found"; exit 1; }
            elif [[ ${INSTANCE_TYPE} == "p5e.48xlarge" ]]; then
                grep "NET/Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Avg bus bandwidth
text not found"; exit 1; }
                grep "NET/OFI Selected Provider is efa (found 32 nics)" ${TRAINING_LOG} ||
{ echo "Avg bus bandwidth text not found"; exit 1; }
                elif [[ ${INSTANCE_TYPE} == "p5en.48xlarge" ]]; then
                    grep "NET/Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Avg bus bandwidth
text not found"; exit 1; }
                    grep "NET/OFI Selected Provider is efa (found 32 nics)" ${TRAINING_LOG} ||
{ echo "Avg bus bandwidth text not found"; exit 1; }
                    elif [[ ${INSTANCE_TYPE} == "p3dn.24xlarge" ]]; then
                        grep "NET/OFI Selected Provider is efa (found 4 nics)" ${TRAINING_LOG} ||
{ echo "Selected Provider is efa text not found"; exit 1; }
                    fi
                echo "***** check_efa_nccl_all_reduce passed for cuda
version ${CUDA_VERSION} *****"
            else
                echo "***** check_efa_nccl_all_reduce failed for cuda
version ${CUDA_VERSION} *****"
            fi

```

- ベンチマークデータにアクセスするには、マルチノード all_reduce テストからテーブル出力の最終行を解析します。

```
benchmark=$(sudo cat ${TRAINING_LOG} | grep '1073741824' | tail -n1 | awk -F " "
'{{print $12}}' | sed 's/ //' | sed 's/ 5e-07//')
if [[ -z "${benchmark}" ]]; then
  echo "benchmark variable is empty"
  exit 1
fi

echo "Benchmark throughput: ${benchmark}"
```

GPU のモニタリングおよび最適化

次のセクションでは、GPU の最適化とモニタリングオプションについて説明します。このセクションは、モニタリング、監督、事前処理、トレーニングが伴う一般的なワークフローと同様に編成されています。

- [モニタリング](#)
 - [CloudWatch を使用して GPU をモニタリングする](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

モニタリング

DLAMI には、いくつかの GPU モニタリングツールがプリインストールされています。このガイドでは、ダウンロードしてインストールするために利用できるツールについても言及されています。

- [CloudWatch を使用して GPU をモニタリングする](#) - Amazon CloudWatch で GPU 使用状況統計を報告するプリインストールされたユーティリティ。
- [nvidia-smi CLI](#) - 全体的な GPU コンピューティングおよびメモリ使用率をモニタリングするユーティリティ。これは (AWS Deep Learning AMIs DLAMI) にプリインストールされています。
- [NVML C ライブラリ](#) - GPU モニタリングおよび管理機能に直接アクセスできる C ベースの API。これは、内部の nvidia-smi CLI によって使用され、DLAMI にプリインストールされています。また、それらの言語での開発を容易にするため、Python および Perl がバインドされていま

す。DLAMI にプリインストールされた gpumon.py ユーティリティは、[nvidia-ml-py](#) の pynvml パッケージを使用しています。

- [NVIDIA DCGM](#) - クラスタ管理ツール。開発者ページにアクセスし、このツールをインストールして設定する方法を確認してください。

i Tip

NVIDIA の開発者ブログで、DLAMI にインストールされている CUDA ツールの使用方法に関する最新情報を確認してください。

- (「[Monitoring TensorCore utilization using Nsight IDE and nvprof](#)」)。

CloudWatch を使用して GPU をモニタリングする

GPU で DLAMI を使用すると、トレーニングや推論中にその使用状況を追跡する方法が必要になることがあります。これは、データパイプラインの最適化や深層学習ネットワークのチューニングに役立ちます。

CloudWatch を使用して GPU メトリクスを設定するには、次の 2 つの方法があります。

- [AWS CloudWatch エージェントを使用してメトリクスを設定する \(推奨\)](#)
- [プリインストールされた gpumon.py スクリプトを使用してメトリクスを設定する](#)

AWS CloudWatch エージェントを使用してメトリクスを設定する (推奨)

DLAMI を[統合 CloudWatch エージェント](#)と統合して GPU メトリクスを設定し、Amazon EC2 高速インスタンスの GPU コプロセスの使用状況をモニタリングします。

DLAMI で [GPU メトリクス](#)を設定するには、次の 4 つの方法があります。

- [最小限の GPU メトリクスを設定する](#)
- [部分的な GPU メトリクスを設定する](#)
- [使用可能なすべての GPU メトリクスを設定する](#)
- [カスタム GPU メトリクスを設定する](#)

更新とセキュリティパッチの詳細については、「[AWS CloudWatch エージェントのセキュリティパッチ適用](#)」を参照してください。

前提条件

開始するには、Amazon EC2 インスタンスの IAM 許可を設定して、インスタンスが CloudWatch にメトリクスをプッシュできるようにする必要があります。詳しい手順については、「[CloudWatch エージェントで使用する IAM ロールとユーザーを作成する](#)」を参照してください。

最小限の GPU メトリクスを設定する

dlami-cloudwatch-agent@minimal systemd サービスを使用して最小限の GPU メトリクスを設定します。このサービスは以下のメトリクスを設定します。

- utilization_gpu
- utilization_memory

事前設定済みの最小限の GPU メトリクス向けの systemd サービスは以下の場所にあります。

```
/opt/aws/amazon-cloudwatch-agent/etc/dlami-amazon-cloudwatch-agent-minimal.json
```

以下のコマンドで systemd サービスを有効にして起動します。

```
sudo systemctl enable dlami-cloudwatch-agent@minimal
sudo systemctl start dlami-cloudwatch-agent@minimal
```

部分的な GPU メトリクスを設定する

dlami-cloudwatch-agent@partial systemd サービスを使用して部分的な GPU メトリクスを設定します。このサービスは以下のメトリクスを設定します。

- utilization_gpu
- utilization_memory
- memory_total
- memory_used
- memory_free

事前設定済みの部分的な GPU メトリクス向けの systemd サービスは以下の場所にあります。

```
/opt/aws/amazon-cloudwatch-agent/etc/dlami-amazon-cloudwatch-agent-partial.json
```

以下のコマンドで systemd サービスを有効にして起動します。

```
sudo systemctl enable dlami-cloudwatch-agent@partial
sudo systemctl start dlami-cloudwatch-agent@partial
```

使用可能なすべての GPU メトリクスを設定する

dlami-cloudwatch-agent@all systemd サービスを使用して使用可能なすべての GPU メトリクスを設定します。このサービスは以下のメトリクスを設定します。

- utilization_gpu
- utilization_memory
- memory_total
- memory_used
- memory_free
- temperature_gpu
- power_draw
- fan_speed
- pcie_link_gen_current
- pcie_link_width_current
- encoder_stats_session_count
- encoder_stats_average_fps
- encoder_stats_average_latency
- clocks_current_graphics
- clocks_current_sm
- clocks_current_memory
- clocks_current_video

事前設定済みの使用可能なすべての GPU メトリクス向けの systemd サービスは以下の場所にあります。

```
/opt/aws/amazon-cloudwatch-agent/etc/dlami-amazon-cloudwatch-agent-all.json
```

以下のコマンドで systemd サービスを有効にして起動します。

```
sudo systemctl enable dlami-cloudwatch-agent@all
sudo systemctl start dlami-cloudwatch-agent@all
```

カスタム GPU メトリクスを設定する

事前設定済みのメトリクスが要件を満たさない場合は、カスタム CloudWatch エージェント設定ファイルを作成できます。

カスタム設定を作成する

カスタム設定ファイルを作成するには、「[CloudWatch エージェント設定ファイルを手動で作成または編集する](#)」の詳細手順を参照してください。

この例では、スキーマ定義が `/opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.json` にあると仮定します。

カスタムファイルを使用してメトリクスを設定する

以下のコマンドを実行して、カスタムファイルに従って CloudWatch エージェントを設定します。

```
sudo /opt/aws/amazon-cloudwatch-agent/bin/amazon-cloudwatch-agent-ctl \
-a fetch-config -m ec2 -s -c \
file:/opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.json
```

AWS CloudWatch エージェントのセキュリティパッチ適用

新しくリリースされた DLAMIs は、利用可能な最新の AWS CloudWatch エージェントセキュリティパッチで設定されます。以下のセクションを参照して、お使いのオペレーティングシステムに応じて、現在の DLAMI を最新のセキュリティパッチで更新してください。

Amazon Linux 2

を使用して yum、Amazon Linux 2 DLAMI の最新の AWS CloudWatch エージェントセキュリティパッチを取得します。

```
sudo yum update
```

Ubuntu

Ubuntu で DLAMI の最新の AWS CloudWatch セキュリティパッチを取得するには、Amazon S3 ダウンロードリンクを使用して AWS CloudWatch エージェントを再インストールする必要があります。

```
wget https://s3.region.amazonaws.com/amazoncloudwatch-agent-region/ubuntu/arm64/latest/amazon-cloudwatch-agent.deb
```

Amazon S3 ダウンロードリンクを使用して AWS CloudWatch エージェントをインストールする方法の詳細については、[「サーバーでの CloudWatch エージェントのインストールと実行」](#)を参照してください。

プリインストールされた **gpumon.py** スクリプトを使用してメトリクスを設定する

gpumon.py というユーティリティは、DLAMI にプリインストールされています。これは、CloudWatch と統合されており、GPU ごとの使用状況 (GPU メモリ、GPU 温度、GPU 電源) のモニタリングをサポートしています。このスクリプトは、モニタリングしたデータを CloudWatch に定期的に送信します。スクリプトでいくつかの設定を変更することで、CloudWatch に送信されるデータの詳細度レベルを設定できます。ただし、スクリプトを開始する前に、CloudWatch がメトリクスを受け取るように設定する必要があります。

CloudWatch を使用した GPU モニタリングを設定して実行する方法

1. IAM ユーザーを作成するか、既存のユーザーを変更し、メトリクスを CloudWatch に発行するためのポリシーを追加します。新しいユーザーを作成する場合、認証情報をメモしてください。次のステップで必要になります。

検索する IAM ポリシーは「cloudwatch:PutMetricData」です。追加されるポリシーは次のようになります。

```
{
  "Version": "2012-10-17",
  "Statement": [
    {
      "Action": [
        "cloudwatch:PutMetricData"
      ],
      "Effect": "Allow",
      "Resource": "*"
    }
  ]
}
```

```
    }
  ]
}
```

 Tip

IAM ユーザーの作成と CloudWatch のポリシーの追加の詳細については、[CloudWatch のドキュメント](#)を参照してください。

- DLAMI で、[AWS 構成](#)を実行し、IAM ユーザー認証情報を指定します。

```
$ aws configure
```

- 実行する前に、gpumon ユーティリティにいくつかの変更が必要になる場合があります。gpumon ユーティリティと README は次のコードブロックに定義された場所にあります。gpumon.py スクリプトの詳細については、[スクリプトの Amazon S3 の場所](#)を参照してください。

```
Folder: ~/tools/GPUCloudWatchMonitor
Files:  ~/tools/GPUCloudWatchMonitor/gpumon.py
        ~/tools/GPUCloudWatchMonitor/README
```

オプション:

- インスタンスが us-east-1 でない場合、gpumon.py でリージョンを変更します。
 - CloudWatch namespace や、store_reso によるレポートの期間などの他のパラメータを変更します。
- 現在、スクリプトでは Python 3 のみがサポートされています。希望するフレームワークの Python 3 環境を有効化するか、DLAMI の一般的な Python 3 環境を有効化します。

```
$ source activate python3
```

- gpumon ユーティリティをバックグラウンドで実行します。

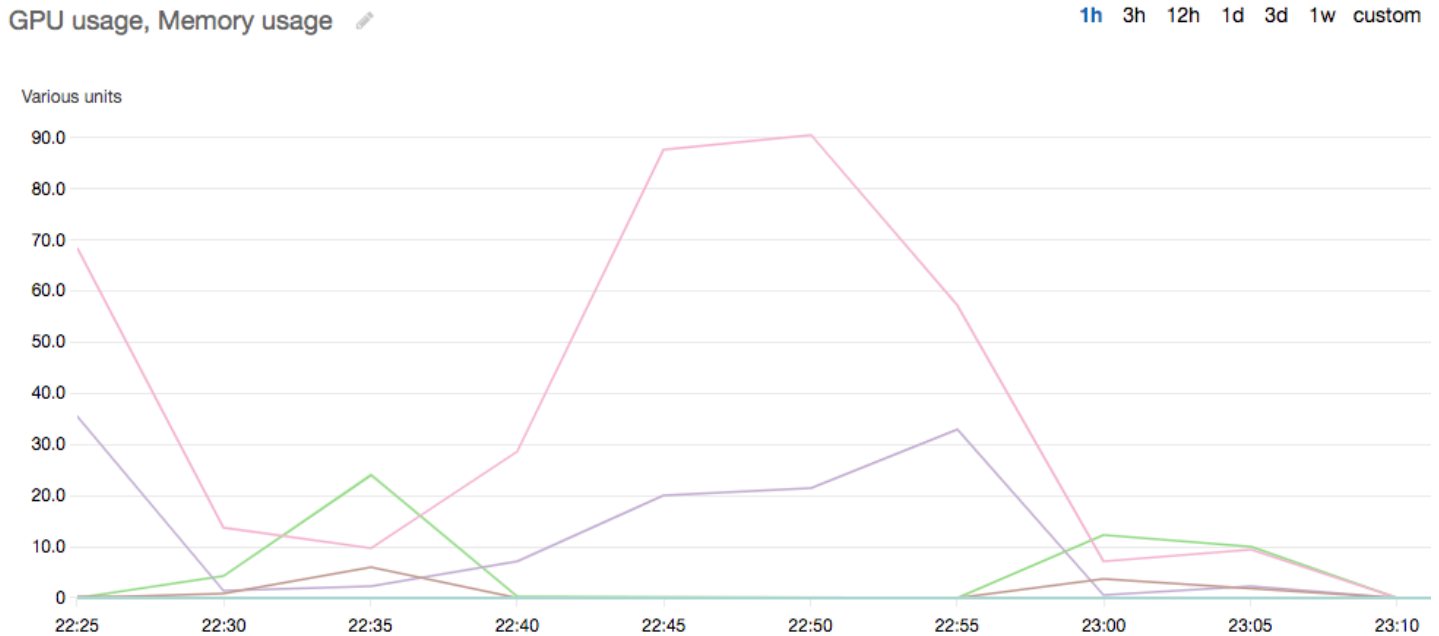
```
(python3)$ python gpumon.py &
```

- ブラウザを開いて <https://console.aws.amazon.com/cloudwatch/> にアクセスし、メトリクスを選択します。名前空間は「DeepLearningTrain」になります。

i Tip

gpumon.py を変更することで名前空間を変更できます。store_reso を調整することで、レポートの間隔を調整することもできます。

次に示すのは、p2.8xlarge インスタンス上のトレーニングジョブをモニタリングする gpumon.py の実行について報告する CloudWatch グラフの例です。



必要に応じて、GPU モニタリングおよび最適化に関する他のトピックも参照できます。

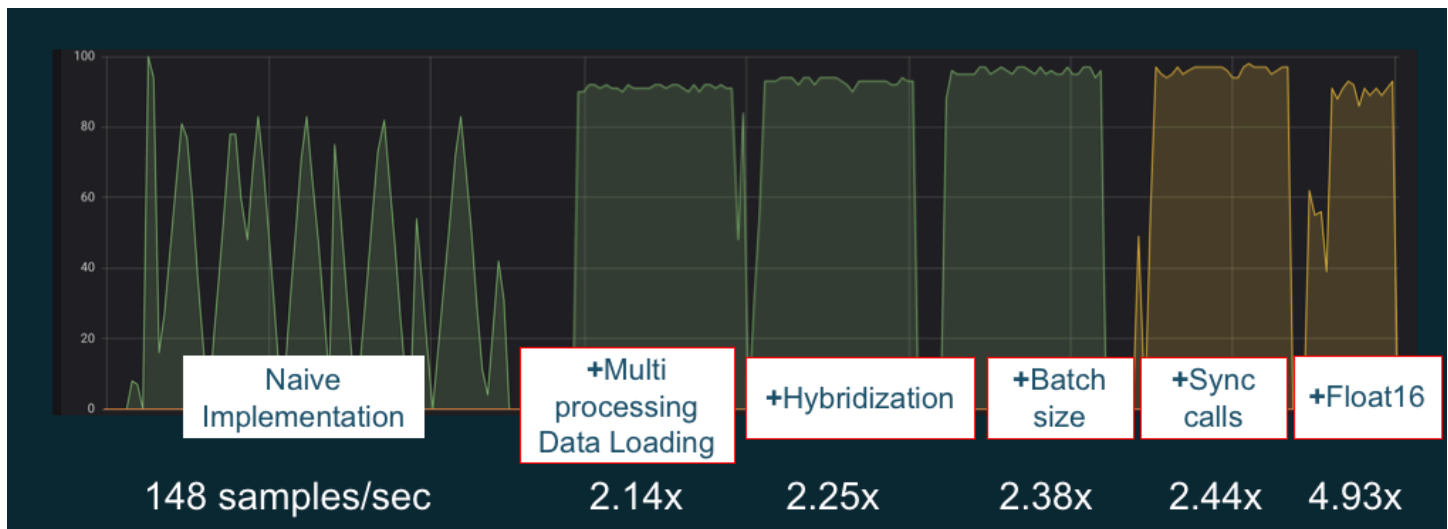
- [モニタリング](#)
 - [CloudWatch を使用して GPU をモニタリングする](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

最適化

GPU を最大限に活用するため、データパイプラインを最適化し、深層学習ネットワークをチューニングできます。次の図が示すように、ニューラルネットワークのネイティブまたは基本実装では、GPU が一貫性なく使用され、潜在能力が十分に引き出されない可能性があります。事前処理と

データのロードを最適化すると、CPU から GPU へのボトルネックを減らすことができます。ハイブリダイゼーションを使用し (フレームワークでサポートされている場合)、バッチサイズを調整してコールを同期することで、ニューラルネットワーク自体を調整することができます。ほとんどのフレームワークでは、多精度 (float16 または int8) トレーニングを使用することもできます。スループットの向上に劇的な効果が及ぶ可能性があります。

以下の図は、さまざまな最適化を適用した場合の累積的なパフォーマンス向上を示しています。結果は、処理するデータと最適化するネットワークによって異なります。



GPU パフォーマンスの最適化の例。グラフの出典: [MXNet と Gluon のパフォーマンスのヒント](#)

次のガイドでは、DLAMI を使用し、GPU パフォーマンスの向上に役立つオプションについて紹介します。

トピック

- [前処理](#)
- [トレーニング](#)

前処理

トランスフォーメーションやオーグメンテーションによるデータの事前処理は、多くの場合、CPU バウンド処理であり、パイプライン全体でボトルネックになる可能性があります。フレームワークには、画像処理用の組み込み演算子がありますが、DALI (データオーグメンテーションライブラリ) はフレームワークの組み込みオプションよりパフォーマンスが上回ることを実証しています。

- NVIDIA データオーグメンテーションライブラリ (DALI): DALI は、データオーグメンテーションを GPU にオフロードします。DLAMI にはプリインストールされていませんが、DLAMI や他の

Amazon Elastic Compute Cloud インスタンスに DALI をインストールするか、サポートされているフレームワークコンテナをロードすることにより、DALI にアクセスできます。詳細については、NVIDIA ウェブサイトの [DALI プロジェクトページ](#) を参照してください。ユースケースの例とコードサンプルのダウンロードについては、[SageMaker 前処理トレーニングパフォーマンス](#) のサンプルを参照してください。

- [nvJPEG](#): C プログラマー向け GPU 加速 JPEG デコーダーライブラリ。1つのイメージまたはバッチのデコードに加えて、深層学習に共通する後続のトランスフォーメーション操作もサポートされています。nvJPEG には DALI が組み込まれています。または、[NVIDIA ウェブサイトの nvjpeg ページ](#) からダウンロードし、別個に使用することもできます。

必要に応じて、GPU モニタリングおよび最適化に関する他のトピックも参照できます。

- [モニタリング](#)
 - [CloudWatch を使用して GPU をモニタリングする](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

トレーニング

混合精度トレーニングでは、メモリの量が同じでより大規模なネットワークをデプロイしたり、単精度または倍精度ネットワークと比較してメモリの使用量を減らしたりすることができます。これにより、コンピューティングパフォーマンスが向上します。また、複数ノードに分散したトレーニングでは重要な要素である、少量かつ高速なデータ転送というメリットも得られます。混合精度トレーニングを利用するには、データキャストと損失スケールを調整する必要があります。混合精度をサポートするフレームワークでこれを行う方法について説明しているガイドを以下に示します。

- [NVIDIA Deep Learning SDK](#) - MXNet、PyTorch、および TensorFlow の混合精度実装について説明している NVIDIA ウェブサイト上のドキュメント。

Tip

必ず、ウェブサイトで選択したフレームワークを確認し、「混合精度」または「fp16」を検索して最新の最適化手法を参照してください。以下の混合精度ガイドが役に立つ可能性があります。

- [Mixed-precision training with TensorFlow \(ビデオ\)](#) - NVIDIA ブログサイト上。
- [Mixed-precision training using float16 with MXNet](#) - MXNet ウェブサイト上のよくある質問記事。
- [NVIDIA Apex: a tool for easy mixed-precision training with PyTorch](#) - NVIDIA ウェブサイトのブログ記事。

必要に応じて、GPU モニタリングおよび最適化に関する他のトピックも参照できます。

- [モニタリング](#)
 - [CloudWatch を使用して GPU をモニタリングする](#)
- [最適化](#)
 - [前処理](#)
 - [トレーニング](#)

DLAMI を使用した AWS Inferentia チップ

AWS Inferentia は、によって設計されたカスタム機械学習チップ AWS で、高性能な推論予測に使用できます。このチップを使用するには、Amazon Elastic Compute Cloud インスタンスをセットアップし、AWS Neuron ソフトウェア開発キット (SDK) を使用して Inferentia チップを呼び出します。お客様に最高の Inferentia エクスペリエンスを提供するために、Neuron が AWS Deep Learning AMIs (DLAMI) に組み込まれています。

以下のトピックでは、Inferentia と DLAMI の使用を開始する方法について説明します。

内容

- [AWS Neuron を使用した DLAMI インスタンスの起動](#)
- [AWS Neuron での DLAMI の使用](#)

AWS Neuron を使用した DLAMI インスタンスの起動

最新の DLAMI は Inferentia AWS で使用できるようになり、AWS Neuron API パッケージが付属しています。DLAMI インスタンスを起動するには、[DLAMI の起動と設定](#)を参照してください。DLAMI を作成したら、以下のステップを使用して、AWS Inferentia チップと AWS Neuron リソースがアクティブであることを確認します。

内容

- [インスタンスの確認](#)
- [AWS Inferentia デバイスの識別](#)
- [リソースの使用状況の表示](#)
- [Neuron Monitor \(neuron-monitor\) の使用](#)
- [Neuron ソフトウェアのアップグレード](#)

インスタンスの確認

インスタンスを使用する前に、インスタンスが正しくセットアップされ、Neuron で設定されていることを確認します。

AWS Inferentia デバイスの識別

インスタンス上の Inferentia デバイスの数を確認するには、次のコマンドを使用します。

```
neuron-ls
```

インスタンスに Inferentia デバイスがアタッチされている場合、出力は次のようになります。

```
+-----+-----+-----+-----+-----+
| NEURON | NEURON | NEURON | CONNECTED | PCI      |
| DEVICE | CORES  | MEMORY | DEVICES   | BDF      |
+-----+-----+-----+-----+-----+
| 0      | 4      | 8 GB   | 1         | 0000:00:1c.0 |
| 1      | 4      | 8 GB   | 2, 0      | 0000:00:1d.0 |
| 2      | 4      | 8 GB   | 3, 1      | 0000:00:1e.0 |
| 3      | 4      | 8 GB   | 2         | 0000:00:1f.0 |
+-----+-----+-----+-----+-----+
```

示されている出力は Inf1.6xLarge インスタンスから取得されたものであり、次の列が含まれていません。

- NEURON DEVICE: NeuronDevice に割り当てられた論理 ID。この ID は、異なる NeuronDevice を使用するように複数のランタイムを設定する場合に使用されます。
- NEURON CORES: NeuronDevice に存在する NeuronCore の数。
- NEURON MEMORY: NeuronDevice 内の DRAM メモリの量。

- CONNECTED DEVICES: NeuronDevice に接続されている他の NeuronDevice。
- PCI BDF: NeuronDevice の PCI バスデバイス機能 (BDF) ID。

リソースの使用状況の表示

neuron-top コマンドを使用して、NeuronCore と vCPU の使用率、メモリ使用率、ロードされたモデル、Neuron アプリケーションに関する有用な情報を表示します。引数を指定せずに neuron-top を起動すると、NeuronCore を使用するすべての機械学習アプリケーションのデータが表示されます。

```
neuron-top
```

アプリケーションが 4 つの NeuronCore を使用している場合、出力は以下の画像のようになります。



Neuron ベースの推論アプリケーションを監視および最適化するためのリソースの詳細については、[Neuron Tools](#) を参照してください。

Neuron Monitor (neuron-monitor) の使用

Neuron Monitor は、システムで実行されている Neuron ランタイムからメトリクスを収集し、収集したデータを JSON 形式で標準出力にストリーミングします。これらのメトリクスは、設定ファイルを指定して設定するメトリクスグループに編成されます。Neuron Monitor の詳細については、[Neuron Monitor User Guide](#) を参照してください。

Neuron ソフトウェアのアップグレード

DLAMI 内で Neuron SDK ソフトウェアを更新する方法については、AWS 「Neuron [セットアップガイド](#)」を参照してください。

次のステップ

[AWS Neuron での DLAMI の使用](#)

AWS Neuron での DLAMI の使用

AWS Neuron SDK の一般的なワークフローは、以前にトレーニングした機械学習モデルをコンパイルサーバーでコンパイルすることです。その後、実行のためにアーティファクトを Inf1 インスタンスに配布します。AWS Deep Learning AMIs (DLAMI) には、Inferentia を使用する Inf1 インスタンスで推論をコンパイルして実行するために必要なものがすべてプリインストールされています。

以下のセクションでは、Inferentia とともに DLAMI を使用する方法について説明します。

内容

- [TensorFlow-Neuron と AWS Neuron Compiler の使用](#)
- [AWS Neuron TensorFlow Serving の使用](#)
- [MXNet-Neuron と AWS Neuron Compiler の使用](#)
- [MXNet-Neuron モデルサービングの使用](#)
- [PyTorch-Neuron と AWS Neuron Compiler の使用](#)

TensorFlow-Neuron と AWS Neuron Compiler の使用

このチュートリアルでは、AWS Neuron コンパイラを使用して Keras ResNet-50 モデルをコンパイルし、保存済みモデルとして SavedModel 形式でエクスポートする方法を示します。このフォーマットは、一般的な TensorFlow モデル交換可能な形式です。また、入力例を使用して Inf1 インスタンスで推論を実行する方法も学習します。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [Resnet50 コンパイル](#)
- [ResNet50 推論](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して、TensorFlow-Neuron Conda 環境をアクティブにします。

```
source activate aws_neuron_tensorflow_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

Resnet50 コンパイル

次の内容を含む **tensorflow_compile_resnet50.py** という Python スクリプトを作成します。この Python スクリプトは、Keras ResNet50 モデルをコンパイルし、保存されたモデルとしてエクスポートします。

```
import os
import time
import shutil
import tensorflow as tf
import tensorflow.neuron as tfn
import tensorflow.compat.v1.keras as keras
from tensorflow.keras.applications.resnet50 import ResNet50
from tensorflow.keras.applications.resnet50 import preprocess_input
```



```
# Create a workspace
WORKSPACE = './ws_resnet50'
os.makedirs(WORKSPACE, exist_ok=True)

# Prepare export directory (old one removed)
model_dir = os.path.join(WORKSPACE, 'resnet50')
compiled_model_dir = os.path.join(WORKSPACE, 'resnet50_neuron')
shutil.rmtree(model_dir, ignore_errors=True)
shutil.rmtree(compiled_model_dir, ignore_errors=True)

# Instantiate Keras ResNet50 model
keras.backend.set_learning_phase(0)
model = ResNet50(weights='imagenet')

# Export SavedModel
tf.saved_model.simple_save(
    session          = keras.backend.get_session(),
    export_dir       = model_dir,
    inputs           = {'input': model.inputs[0]},
    outputs          = {'output': model.outputs[0]})

# Compile using Neuron
tfn.saved_model.compile(model_dir, compiled_model_dir)

# Prepare SavedModel for uploading to Inf1 instance
shutil.make_archive(compiled_model_dir, 'zip', WORKSPACE, 'resnet50_neuron')
```

次のコマンドを使用してモデルをコンパイルします。

```
python tensorflow_compile_resnet50.py
```

コンパイル処理には数分かかります。完了すると、出力は次のようになります。

```
...
INFO:tensorflow:fusing subgraph neuron_op_d6f098c01c780733 with neuron-cc
INFO:tensorflow:Number of operations in TensorFlow session: 4638
INFO:tensorflow:Number of operations after tf.neuron optimizations: 556
INFO:tensorflow:Number of operations placed on Neuron runtime: 554
INFO:tensorflow:Successfully converted ./ws_resnet50/resnet50 to ./ws_resnet50/
resnet50_neuron
```

...

コンパイル後、保存されたモデルは **ws_resnet50/resnet50_neuron.zip** に圧縮されます。以下のコマンドを使用して、モデルを解凍し、推論用のサンプルイメージをダウンロードします。

```
unzip ws_resnet50/resnet50_neuron.zip -d .
curl -O https://raw.githubusercontent.com/awslabs/mxnet-model-server/master/docs/images/kitten_small.jpg
```

ResNet50 推論

次の内容を含む **tensorflow_infer_resnet50.py** という Python スクリプトを作成します。このスクリプトは、以前にコンパイルされた推論モデルを使用して、ダウンロードしたモデルに対して推論を実行します。

```
import os
import numpy as np
import tensorflow as tf
from tensorflow.keras.preprocessing import image
from tensorflow.keras.applications import resnet50

# Create input from image
img_sgl = image.load_img('kitten_small.jpg', target_size=(224, 224))
img_arr = image.img_to_array(img_sgl)
img_arr2 = np.expand_dims(img_arr, axis=0)
img_arr3 = resnet50.preprocess_input(img_arr2)
# Load model
COMPILED_MODEL_DIR = './ws_resnet50/resnet50_neuron/'
predictor_inferentia = tf.contrib.predictor.from_saved_model(COMPILED_MODEL_DIR)
# Run inference
model_feed_dict={'input': img_arr3}
infa_rslts = predictor_inferentia(model_feed_dict);
# Display results
print(resnet50.decode_predictions(infa_rslts["output"], top=5)[0])
```

次のコマンドを使用して、モデルに対して推論を実行します。

```
python tensorflow_infer_resnet50.py
```

出力は次のようになります。

```
...
[('n02123045', 'tabby', 0.6918919), ('n02127052', 'lynx', 0.12770271), ('n02123159',
'tiger_cat', 0.08277027), ('n02124075', 'Egyptian_cat', 0.06418919), ('n02128757',
'snow_leopard', 0.009290541)]
```

次のステップ

[AWS Neuron TensorFlow Serving の使用](#)

AWS Neuron TensorFlow Serving の使用

このチュートリアルでは、TensorFlow Serving で使用する保存済みモデルをエクスポートする前に、グラフを作成し、AWS Neuron コンパイルステップを追加する方法を示します。TensorFlow Serving は、ネットワーク全体で推論をスケールアップすることを可能にする提供システムです。Neuron TensorFlow Serving は通常の TensorFlow Serving と同じ API を使用します。唯一の違いは、保存されたモデルを Inferentia AWS 用にコンパイルする必要があり、エントリポイントは `tensorflow_model_server_neuron` という名前の別のバイナリであることです。バイナリは、`/usr/local/bin/tensorflow_model_server_neuron` にあり、DLAMI にあらかじめインストールされています。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [保存したモデルのコンパイルとエクスポート](#)
- [保存したモデルの提供](#)
- [モデルサーバーへの推論リクエストを生成する](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して、TensorFlow-Neuron Conda 環境をアクティブにします。

```
source activate aws_neuron_tensorflow_p36
```

現在の Conda 環境を終了する必要がある場合は、次のコマンドを実行します。

```
source deactivate
```

保存したモデルのコンパイルとエクスポート

以下の内容が含まれた Python スクリプト `tensorflow-model-server-compile.py` を作成します。このスクリプトは、Neuron を使用してグラフを作成し、コンパイルします。次に、コンパイルされたグラフを保存されたモデルとしてエクスポートします。

```
import tensorflow as tf
import tensorflow.neuron
import os

tf.keras.backend.set_learning_phase(0)
model = tf.keras.applications.ResNet50(weights='imagenet')
sess = tf.keras.backend.get_session()
inputs = {'input': model.inputs[0]}
outputs = {'output': model.outputs[0]}

# save the model using tf.saved_model.simple_save
model_dir = "./resnet50/1"
tf.saved_model.simple_save(sess, model_dir, inputs, outputs)

# compile the model for Inferentia
neuron_model_dir = os.path.join(os.path.expanduser('~'), 'resnet50_inf1', '1')
tf.neuron.saved_model.compile(model_dir, neuron_model_dir, batch_size=1)
```

次のコマンドを使用してモデルをコンパイルします。

```
python tensorflow-model-server-compile.py
```

出力は次のようになります。

```
...
INFO:tensorflow:fusing subgraph neuron_op_d6f098c01c780733 with neuron-cc
INFO:tensorflow:Number of operations in TensorFlow session: 4638
INFO:tensorflow:Number of operations after tf.neuron optimizations: 556
INFO:tensorflow:Number of operations placed on Neuron runtime: 554
INFO:tensorflow:Successfully converted ./resnet50/1 to /home/ubuntu/resnet50_inf1/1
```

保存したモデルの提供

モデルをコンパイルしたら、次のコマンドを使用して、保存したモデルに `tensorflow_model_server_neuron` バイナリを提供できます。

```
tensorflow_model_server_neuron --model_name=resnet50_inf1 \
  --model_base_path=$HOME/resnet50_inf1/ --port=8500 &
```

出力は次のようになります。コンパイルされたモデルは、推論の準備のためにサーバーによって Inferentia デバイスの DRAM にステージングされます。

```
...
2019-11-22 01:20:32.075856: I external/org_tensorflow/tensorflow/cc/saved_model/
loader.cc:311] SavedModel load for tags { serve }; Status: success. Took 40764
microseconds.
2019-11-22 01:20:32.075888: I tensorflow_serving/servables/tensorflow/
saved_model_warmup.cc:105] No warmup data file found at /home/ubuntu/resnet50_inf1/1/
assets.extra/tf_serving_warmup_requests
2019-11-22 01:20:32.075950: I tensorflow_serving/core/loader_harness.cc:87]
Successfully loaded servable version {name: resnet50_inf1 version: 1}
2019-11-22 01:20:32.077859: I tensorflow_serving/model_servers/
server.cc:353] Running gRPC ModelServer at 0.0.0.0:8500 ...
```

モデルサーバーへの推論リクエストを生成する

次の内容で `tensorflow-model-server-infer.py` という Python スクリプトを作成します。このスクリプトはサービスフレームワークである gRPC を介して推論を実行してください。

```
import numpy as np
import grpc
import tensorflow as tf
from tensorflow.keras.preprocessing import image
from tensorflow.keras.applications.resnet50 import preprocess_input
from tensorflow_serving.apis import predict_pb2
from tensorflow_serving.apis import prediction_service_pb2_grpc
from tensorflow.keras.applications.resnet50 import decode_predictions

if __name__ == '__main__':
    channel = grpc.insecure_channel('localhost:8500')
    stub = prediction_service_pb2_grpc.PredictionServiceStub(channel)
    img_file = tf.keras.utils.get_file(
        "./kitten_small.jpg",
        "https://raw.githubusercontent.com/aws-labs/mxnet-model-server/master/docs/
images/kitten_small.jpg")
    img = image.load_img(img_file, target_size=(224, 224))
    img_array = preprocess_input(image.img_to_array(img)[None, ...])
    request = predict_pb2.PredictRequest()
    request.model_spec.name = 'resnet50_inf1'
    request.inputs['input'].CopyFrom(
        tf.contrib.util.make_tensor_proto(img_array, shape=img_array.shape))
    result = stub.Predict(request)
    prediction = tf.make_ndarray(result.outputs['output'])
    print(decode_predictions(prediction))
```

次のコマンドで gRPC を使用して、モデルの推論を実行します。

```
python tensorflow-model-server-infer.py
```

出力は次のようになります。

```
[(['n02123045', 'tabby', 0.6918919), ('n02127052', 'lynx', 0.12770271), ('n02123159',
'tiger_cat', 0.08277027), ('n02124075', 'Egyptian_cat', 0.06418919), ('n02128757',
'snow_leopard', 0.009290541)]]
```

MXNet-Neuron と AWS Neuron Compiler の使用

MXNet-Neuron コンパイル API は、AWS Inferentia デバイスで実行できるモデルグラフをコンパイルする方法を提供します。

この例では、API を使用して ResNet-50 モデルをコンパイルし、それを使用して推論を実行します。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [Resnet50 コンパイル](#)
- [ResNet50 推論](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して、MXNet-Neuron Conda 環境をアクティブにします。

```
source activate aws_neuron_mxnet_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

Resnet50 コンパイル

次の内容で `mxnet_compile_resnet50.py` という Python スクリプトを作成します。このスクリプトは、MXNet-Neuron コンパイル Python API を使用して、ResNet-50 モデルをコンパイルします。

```
import mxnet as mx
import numpy as np

print("downloading...")
```

```
path='http://data.mxnet.io/models/imagenet/'
mx.test_utils.download(path+'resnet/50-layers/resnet-50-0000.params')
mx.test_utils.download(path+'resnet/50-layers/resnet-50-symbol.json')
print("download finished.")

sym, args, aux = mx.model.load_checkpoint('resnet-50', 0)

print("compile for inferentia using neuron... this will take a few minutes...")
inputs = { "data" : mx.nd.ones([1,3,224,224], name='data', dtype='float32') }

sym, args, aux = mx.contrib.neuron.compile(sym, args, aux, inputs)

print("save compiled model...")
mx.model.save_checkpoint("compiled_resnet50", 0, sym, args, aux)
```

次のコマンドを使用してモデルをコンパイルします。

```
python mxnet_compile_resnet50.py
```

コンパイルには数分かかります。コンパイルが終了すると、次のファイルが現在のディレクトリに表示されます。

```
resnet-50-0000.params
resnet-50-symbol.json
compiled_resnet50-0000.params
compiled_resnet50-symbol.json
```

ResNet50 推論

次の内容で **mxnet_infer_resnet50.py** という Python スクリプトを作成します。このスクリプトは、サンプルイメージをダウンロードし、それを使用して、コンパイルされたモデルを持つ推論を実行します。

```
import mxnet as mx
import numpy as np

path='http://data.mxnet.io/models/imagenet/'
mx.test_utils.download(path+'synset.txt')
```



```
fname = mx.test_utils.download('https://raw.githubusercontent.com/awsmlabs/mxnet-model-server/master/docs/images/kitten_small.jpg')
img = mx.image.imread(fname)

# convert into format (batch, RGB, width, height)
img = mx.image.imresize(img, 224, 224)
# resize
img = img.transpose((2, 0, 1))
# Channel first
img = img.expand_dims(axis=0)
# batchify
img = img.astype(dtype='float32')

sym, args, aux = mx.model.load_checkpoint('compiled_resnet50', 0)
softmax = mx.nd.random_normal(shape=(1,))
args['softmax_label'] = softmax
args['data'] = img
# Inferentia context
ctx = mx.neuron()

exe = sym.bind(ctx=ctx, args=args, aux_states=aux, grad_req='null')
with open('synset.txt', 'r') as f:
    labels = [l.rstrip() for l in f]

exe.forward(data=img)
prob = exe.outputs[0].asnumpy()
# print the top-5
prob = np.squeeze(prob)
a = np.argsort(prob)[::-1]
for i in a[0:5]:
    print('probability=%f, class=%s' %(prob[i], labels[i]))
```

次のコマンドを使用して、コンパイルされたモデルで推論を実行します。

```
python mxnet_infer_resnet50.py
```

出力は次のようになります。

```
probability=0.642454, class=n02123045 tabby, tabby cat
probability=0.189407, class=n02123159 tiger cat
probability=0.100798, class=n02124075 Egyptian cat
```

```
probability=0.030649, class=n02127052 lynx, catamount  
probability=0.016278, class=n02129604 tiger, Panthera tigris
```

次のステップ

[MXNet-Neuron モデルサービングの使用](#)

MXNet-Neuron モデルサービングの使用

このチュートリアルでは、事前にトレーニングされた MXNet モデルを使用して、マルチモデル サーバー (MMS) でリアルタイムのイメージ分類を実行する方法を学習します。MMS は、Machine Learning や深層学習フレームワークを使用してトレーニングされた深層学習モデルを提供するための、柔軟で使いやすいツールです。このチュートリアルには、AWS Neuron を使用したコンパイルステップと、MXNet を使用した MMS の実装が含まれています。

Neuron SDK の詳細については、[AWS Neuron SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [コード例のダウンロード](#)
- [モデルのコンパイル](#)
- [推論の実行](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して、MXNet-Neuron Conda 環境をアクティブにします。

```
source activate aws_neuron_mxnet_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

コード例のダウンロード

この例を実行するには、次のコマンドを使用してコード例をダウンロードします。

```
git clone https://github.com/aws-labs/multi-model-server
cd multi-model-server/examples/mxnet_vision
```

モデルのコンパイル

次の内容で `multi-model-server-compile.py` という Python スクリプトを作成します。このスクリプトは、ResNet50 モデルを Inferentia デバイスタarget にコンパイルします。

```
import mxnet as mx
from mxnet.contrib import neuron
import numpy as np

path='http://data.mxnet.io/models/imagenet/'
mx.test_utils.download(path+'resnet/50-layers/resnet-50-0000.params')
mx.test_utils.download(path+'resnet/50-layers/resnet-50-symbol.json')
mx.test_utils.download(path+'synset.txt')

nn_name = "resnet-50"

#Load a model
sym, args, auxs = mx.model.load_checkpoint(nn_name, 0)

#Define compilation parameters# - input shape and dtype
inputs = {'data' : mx.nd.zeros([1,3,224,224], dtype='float32')}

# compile graph to inferentia target
csym, cargs, cauxs = neuron.compile(sym, args, auxs, inputs)

# save compiled model
mx.model.save_checkpoint(nn_name + "_compiled", 0, csym, cargs, cauxs)
```

モデルをコンパイルするには、次のコマンドを使用します。

```
python multi-model-server-compile.py
```

出力は次のようになります。

```
...
```

```
[21:18:40] src/nnvm/legacy_json_util.cc:209: Loading symbol saved by previous version
v0.8.0. Attempting to upgrade...
[21:18:40] src/nnvm/legacy_json_util.cc:217: Symbol successfully upgraded!
[21:19:00] src/operator/subgraph/build_subgraph.cc:698: start to execute partition
graph.
[21:19:00] src/nnvm/legacy_json_util.cc:209: Loading symbol saved by previous version
v0.8.0. Attempting to upgrade...
[21:19:00] src/nnvm/legacy_json_util.cc:217: Symbol successfully upgraded!
```

次の内容を含むと `signature.json` という名前のファイルを作成し、入力名と形状を設定します。

```
{
  "inputs": [
    {
      "data_name": "data",
      "data_shape": [
        1,
        3,
        224,
        224
      ]
    }
  ]
}
```

次のコマンドを使用して `synset.txt` ファイルを呼び出すことができます。このファイルは、ImageNet 予測クラスの名前の一覧です。

```
curl -O https://s3.amazonaws.com/model-server/model_archive_1.0/examples/
squeeze_net_v1.1/synset.txt
```

`model_server_template` フォルダ内のテンプレートに続くカスタムサービスクラスを作成します。次のコマンドを使用して、テンプレートを現在の作業ディレクトリにコピーします。

```
cp -r ../model_service_template/* .
```

次のように `mxnet_model_service.py` モジュールを編集して、`mx.cpu()` コンテキストを `mx.neuron()` コンテキストに置き換えます。MXNet-Neuron は `NDArray` および `Gluon API` をサポートしていないため、`model_input` の不要なデータコピーをコメントアウトする必要があります。

```
...
self.mxnet_ctx = mx.neuron() if gpu_id is None else mx.gpu(gpu_id)
...
#model_input = [item.as_in_context(self.mxnet_ctx) for item in model_input]
```

次のコマンドを使用して、モデルアーカイバでモデルをパッケージ化します。

```
cd ~/multi-model-server/examples
model-archiver --force --model-name resnet-50_compiled --model-path mxnet_vision --
handler mxnet_vision_service:handle
```

推論の実行

マルチモデルサーバーを起動し、次のコマンドを使用して RESTful API を使用するモデルをロードします。neuron-rtd がデフォルト設定で実行されていることを確認します。

```
cd ~/multi-model-server/
multi-model-server --start --model-store examples > /dev/null # Pipe to log file if you
  want to keep a log of MMS
curl -v -X POST "http://localhost:8081/models?
  initial_workers=1&max_workers=4&synchronous=true&url=resnet-50_compiled.mar"
sleep 10 # allow sufficient time to load model
```

以下のコマンドでサンプルイメージを使用して推論を実行します。

```
curl -O https://raw.githubusercontent.com/awslabs/multi-model-server/master/docs/
  images/kitten_small.jpg
curl -X POST http://127.0.0.1:8080/predictions/resnet-50_compiled -T kitten_small.jpg
```

出力は次のようになります。

```
[
  {
    "probability": 0.6388034820556641,
    "class": "n02123045 tabby, tabby cat"
  },
  {
    "probability": 0.16900072991847992,
    "class": "n02123159 tiger cat"
  },
  {
```

```
    "probability": 0.12221276015043259,
    "class": "n02124075 Egyptian cat"
  },
  {
    "probability": 0.028706775978207588,
    "class": "n02127052 lynx, catamount"
  },
  {
    "probability": 0.01915954425930977,
    "class": "n02129604 tiger, Panthera tigris"
  }
]
```

テスト後にクリーンアップするには、RESTful API を使用して削除コマンドを発行し、次のコマンドを使用してモデルサーバーを停止します。

```
curl -X DELETE http://127.0.0.1:8081/models/resnet-50_compiled

multi-model-server --stop
```

以下の出力が表示されます。

```
{
  "status": "Model \"resnet-50_compiled\" unregistered"
}
Model server stopped.
Found 1 models and 1 NCGs.
Unloading 10001 (MODEL_STATUS_STARTED) :: success
Destroying NCG 1 :: success
```

PyTorch-Neuron と AWS Neuron Compiler の使用

PyTorch-Neuron コンパイル API は、AWS Inferentia デバイスで実行できるモデルグラフをコンパイルする方法を提供します。

トレーニング済みモデルは、Inf1 インスタンスにデプロイする前に、Inferentia ターゲットにコンパイルする必要があります。以下のチュートリアルでは、torchvision ResNet50 モデルをコンパイルし、保存された TorchScript モジュールとしてエクスポートします。このモデルは、推論を実行するために使用されます。

便宜上、このチュートリアルではコンパイルと推論の両方に Inf1 インスタンスを使用します。実際には、c5 インスタンスファミリーなどの別のインスタンスタイプを使用してモデルをコンパイルで

きます。その後、コンパイルされたモデルを Inf1 推論サーバーにデプロイする必要があります。詳細については、[AWS Neuron PyTorch SDK のドキュメント](#)を参照してください。

内容

- [前提条件](#)
- [Conda 環境のアクティブ化](#)
- [Resnet50 コンパイル](#)
- [ResNet50 推論](#)

前提条件

このチュートリアルを使用する前に、[AWS Neuron を使用した DLAMI インスタンスの起動](#) の設定ステップを完了しておく必要があります。また、深層学習および DLAMI の使用にも精通している必要があります。

Conda 環境のアクティブ化

次のコマンドを使用して、PyTorch-Neuron Conda 環境をアクティブにします。

```
source activate aws_neuron_pytorch_p36
```

現在の Conda 環境を終了するには、次のコマンドを実行します。

```
source deactivate
```

Resnet50 コンパイル

次の内容で **pytorch_trace_resnet50.py** という Python スクリプトを作成します。このスクリプトは、PyTorch-Neuron コンパイル Python API を使用して、ResNet-50 モデルをコンパイルします。

Note

torchvision と torch パッケージのバージョン間には、torchvision モデルのコンパイル時に注意すべき依存関係があります。これらの依存関係ルールは、pip を介して管理できます。torchvision==0.6.1 は torch==1.5.1 リリースと対応し、torchvision==0.8.2 は torch==1.7.1 リリースに対応しています。

```
import torch
import numpy as np
import os
import torch_neuron
from torchvision import models

image = torch.zeros([1, 3, 224, 224], dtype=torch.float32)

## Load a pretrained ResNet50 model
model = models.resnet50(pretrained=True)

## Tell the model we are using it for evaluation (not training)
model.eval()
model_neuron = torch.neuron.trace(model, example_inputs=[image])

## Export to saved model
model_neuron.save("resnet50_neuron.pt")
```

コンパイルスクリプトを実行します。

```
python pytorch_trace_resnet50.py
```

コンパイルには数分かかります。コンパイルが完了すると、コンパイルされたモデルは `resnet50_neuron.pt` としてローカルディレクトリに保存されます。

ResNet50 推論

次の内容で **pytorch_infer_resnet50.py** という Python スクリプトを作成します。このスクリプトは、サンプルイメージをダウンロードし、それを使用して、コンパイルされたモデルを持つ推論を実行します。

```
import os
import time
import torch
import torch_neuron
import json
import numpy as np

from urllib import request

from torchvision import models, transforms, datasets
```



```
## Create an image directory containing a small kitten
os.makedirs("./torch_neuron_test/images", exist_ok=True)
request.urlretrieve("https://raw.githubusercontent.com/aws-labs/mxnet-model-server/
master/docs/images/kitten_small.jpg",
                    "./torch_neuron_test/images/kitten_small.jpg")

## Fetch labels to output the top classifications
request.urlretrieve("https://s3.amazonaws.com/deep-learning-models/image-models/
imagenet_class_index.json", "imagenet_class_index.json")
idx2label = []

with open("imagenet_class_index.json", "r") as read_file:
    class_idx = json.load(read_file)
    idx2label = [class_idx[str(k)][1] for k in range(len(class_idx))]

## Import a sample image and normalize it into a tensor
normalize = transforms.Normalize(
    mean=[0.485, 0.456, 0.406],
    std=[0.229, 0.224, 0.225])

eval_dataset = datasets.ImageFolder(
    os.path.dirname("./torch_neuron_test/"),
    transforms.Compose([
        transforms.Resize([224, 224]),
        transforms.ToTensor(),
        normalize,
    ])
)

image, _ = eval_dataset[0]
image = torch.tensor(image.numpy()[np.newaxis, ...])

## Load model
model_neuron = torch.jit.load( 'resnet50_neuron.pt' )

## Predict
results = model_neuron( image )

# Get the top 5 results
top5_idx = results[0].sort()[1][-5:]

# Lookup and print the top 5 labels
```

```
top5_labels = [idx2label[idx] for idx in top5_idx]

print("Top 5 labels:\n {}".format(top5_labels) )
```

次のコマンドを使用して、コンパイルされたモデルで推論を実行します。

```
python pytorch_infer_resnet50.py
```

出力は次のようになります。

```
Top 5 labels:
['tiger', 'lynx', 'tiger_cat', 'Egyptian_cat', 'tabby']
```

ARM64 DLAMI

AWS ARM64 GPU DLAMIsは、深層学習ワークロードに高いパフォーマンスとコスト効率を提供するように設計されています。具体的には、G5g インスタンスタイプは Arm64-based [AWS Graviton2 プロセッサ](#)を搭載しています。これは、[AWS Graviton2 プロセッサ](#)によってゼロから構築され、お客様がクラウドでワークロードを実行する方法に最適化されています。AWS ARM64 GPU DLAMIsは、Docker、NVIDIA Docker、NVIDIA Driver、CUDA、CuDNN、NCCL、および TensorFlow や PyTorch などの一般的な機械学習フレームワークで事前設定されています。

G5g インスタンスタイプでは、Graviton2 の価格およびパフォーマンスの利点を活用することで、GPU アクセラレーション機能を備えた x86 ベースのインスタンスと比較して、GPU アクセラレーションを利用した深層学習モデルを大幅に低いコストでデプロイできます。

ARM64 DLAMI を選択する

選択した ARM64 DLAMI を備えた [G5g インスタンス](#)を起動します。

DLAMI を起動するステップバイステップの手順については、[DLAMI を起動および設定する](#)を参照してください。

最新の ARM64 DLAMI のリストについては、[DLAMI のリリースノート](#)を参照してください。

使用を開始する

以下のトピックでは、ARM64 DLAMI の使用を開始する方法について説明します。

内容

- [ARM64 GPU PyTorch DLAMI の使用](#)

ARM64 GPU PyTorch DLAMI の使用

AWS Deep Learning AMIs は Arm64 プロセッサベースの GPUs で使用できるようになり、PyTorch 用に最適化されています。ARM64 GPU PyTorch DLAMI には、深層学習のトレーニングと推論のユースケース用に [PyTorch](#)、[TorchVision](#)、および [TorchServe](#) で事前設定された Python 環境が含まれています。

内容

- [PyTorch Python 環境を確認する](#)
- [PyTorch でトレーニングサンプルを実行する](#)
- [PyTorch で推論サンプルを実行する](#)

PyTorch Python 環境を確認する

G5g インスタンスに接続し、次のコマンドを使用して Base Conda 環境を有効化します。

```
source activate base
```

コマンドプロンプトは、PyTorch、TorchVision、およびその他のライブラリが含まれた Base Conda 環境で作業していることが示されます。

```
(base) $
```

PyTorch 環境のデフォルトのツールパスを確認します。

```
(base) $ which python
(base) $ which pip
(base) $ which conda
(base) $ which mamba
>>> import torch, torchvision
>>> torch.__version__
>>> torchvision.__version__
>>> v = torch.autograd.Variable(torch.randn(10, 3, 224, 224))
>>> v = torch.autograd.Variable(torch.randn(10, 3, 224, 224)).cuda()
```

```
>>> assert isinstance(v, torch.Tensor)
```

PyTorch でトレーニングサンプルを実行する

サンプル MNIST トレーニングジョブを実行します。

```
git clone https://github.com/pytorch/examples.git
cd examples/mnist
python main.py
```

出力は以下のようになります。

```
...
Train Epoch: 14 [56320/60000 (94%)]    Loss: 0.021424
Train Epoch: 14 [56960/60000 (95%)]    Loss: 0.023695
Train Epoch: 14 [57600/60000 (96%)]    Loss: 0.001973
Train Epoch: 14 [58240/60000 (97%)]    Loss: 0.007121
Train Epoch: 14 [58880/60000 (98%)]    Loss: 0.003717
Train Epoch: 14 [59520/60000 (99%)]    Loss: 0.001729
Test set: Average loss: 0.0275, Accuracy: 9916/10000 (99%)
```

PyTorch で推論サンプルを実行する

次のコマンドを使用して、事前トレーニング済みの densenet161 モデルをダウンロードし、TorchServe を使用して推論を実行します。

```
# Set up TorchServe
cd $HOME
git clone https://github.com/pytorch/serve.git
mkdir -p serve/model_store
cd serve

# Download a pre-trained densenet161 model
wget https://download.pytorch.org/models/densenet161-8d451a50.pth >/dev/null

# Save the model using torch-model-archiver
torch-model-archiver --model-name densenet161 \
  --version 1.0 \
  --model-file examples/image_classifier/densenet_161/model.py \
  --serialized-file densenet161-8d451a50.pth \
  --handler image_classifier \
```

```
--extra-files examples/image_classifier/index_to_name.json \  
--export-path model_store  
  
# Start the model server  
torchserve --start --no-config-snapshots \  
  --model-store model_store \  
  --models densenet161=densenet161.mar &> torchserve.log  
  
# Wait for the model server to start  
sleep 30  
  
# Run a prediction request  
curl http://127.0.0.1:8080/predictions/densenet161 -T examples/image_classifier/  
kitten.jpg
```

出力は以下のようになります。

```
{  
  "tiger_cat": 0.4693363308906555,  
  "tabby": 0.4633873701095581,  
  "Egyptian_cat": 0.06456123292446136,  
  "lynx": 0.0012828150065615773,  
  "plastic_bag": 0.00023322898778133094  
}
```

次のコマンドを使用して densenet161 モデルの登録を解除し、サーバーを停止します。

```
curl -X DELETE http://localhost:8081/models/densenet161/1.0  
torchserve --stop
```

出力は以下のようになります。

```
{  
  "status": "Model \"densenet161\" unregistered"  
}  
TorchServe has stopped.
```

推論

このセクションでは、DLAMI のフレームワークとツールを使用して推論を実行する方法についてチュートリアルを提供します。

推論ツール

- [TensorFlow Serving](#)

モデル提供

Deep Learning AMI with Conda にインストールされているモデル提供のオプションを以下に示します。いずれかのオプションをクリックして、使用方法をご覧ください。

トピック

- [TensorFlow Serving](#)
- [TorchServe](#)

TensorFlow Serving

[TensorFlow Serving](#) は、機械学習モデル向けの柔軟で高パフォーマンスの処理システムです。

tensorflow-serving-api には、単一のフレームワーク DLAMI がプリインストールされています。テンソルフローサービングを使用するには、まず TensorFlow 環境をアクティブ化します。

```
$ source /opt/tensorflow/bin/activate
```

次に、任意のテキストエディタを使用して、以下の内容のスクリプトを作成します。このスクリプトに test_train_mnist.py という名前を付けます。このスクリプトは、イメージを分類するニューラルネットワーク機械学習モデルをトレーニングおよび評価する [TensorFlow チュートリアル](#) から参照されます。

```
import tensorflow as tf
mnist = tf.keras.datasets.mnist

(x_train, y_train),(x_test, y_test) = mnist.load_data()
x_train, x_test = x_train / 255.0, x_test / 255.0

model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dropout(0.2),
```

```
tf.keras.layers.Dense(10, activation='softmax')
])

model.compile(optimizer='adam',
              loss='sparse_categorical_crossentropy',
              metrics=['accuracy'])

model.fit(x_train, y_train, epochs=5)
model.evaluate(x_test, y_test)
```

サーバーの場所、ポート、ハスキーの写真のファイル名をパラメータとして渡してスクリプトを実行します。

```
$ /opt/tensorflow/bin/python3 test_train_mnist.py
```

このスクリプトは出力に時間がかかることがあるため、少し待ちます。トレーニングが完了すると、以下が表示されます。

```
I0000 00:00:1739482012.389276    4284 device_compiler.h:188] Compiled cluster using
XLA! This line is logged at most once for the lifetime of the process.
1875/1875 [=====] - 24s 2ms/step - loss: 0.2973 - accuracy:
0.9134
Epoch 2/5
1875/1875 [=====] - 3s 2ms/step - loss: 0.1422 - accuracy:
0.9582
Epoch 3/5
1875/1875 [=====] - 3s 1ms/step - loss: 0.1076 - accuracy:
0.9687
Epoch 4/5
1875/1875 [=====] - 3s 2ms/step - loss: 0.0872 - accuracy:
0.9731
Epoch 5/5
1875/1875 [=====] - 3s 1ms/step - loss: 0.0731 - accuracy:
0.9771
313/313 [=====] - 0s 1ms/step - loss: 0.0749 - accuracy:
0.9780
```

その他の機能と例

TensorFlow Serving について詳しくお知りになりたい場合は、[TensorFlow ウェブサイト](#)を参照してください。

TorchServe

TorchServe は、PyTorch からエクスポートされた深層学習モデルを供給するための柔軟なツールです。TorchServe には、Deep Learning AMI with Conda がプリインストールされています。

TorchServe の使用方法の詳細については、[Model Server for PyTorch のドキュメント](#)を参照してください。

トピック

TorchServe でイメージ分類モデルを供給する

このチュートリアルでは、TorchServe でイメージ分類モデルを供給する方法を説明します。PyTorch で提供されている Densenet-161 モデルを使用します。サーバーが実行されると、予測リクエストをリッスンします。イメージをアップロードすると (この例では猫の画像)、サーバーはモデルがトレーニングされたクラスから最適な 5 つの一致するクラスの予測を返します。

TorchServe でイメージ分類モデル例を供給するには

1. Deep Learning AMI with Conda v34 以降の Amazon Elastic Compute Cloud (Amazon EC2) インスタンスに接続します。
2. `pytorch_p310` 環境をアクティブ化します。

```
source activate pytorch_p310
```

3. TorchServe リポジトリをクローンし、モデルを保存するディレクトリを作成します。

```
git clone https://github.com/pytorch/serve.git
mkdir model_store
```

4. モデルアーカイバを使用してモデルをアーカイブします。extra-files パラメータは、TorchServe リポジトリからのファイルを使用するので、必要に応じてパスを更新します。モデルアーカイバの詳細については、[Torch Model archiver for TorchServe](#) を参照してください。

```
wget https://download.pytorch.org/models/densenet161-8d451a50.pth
torch-model-archiver --model-name densenet161 --version 1.0 --model-file ./
serve/examples/image_classifier/densenet_161/model.py --serialized-file
densenet161-8d451a50.pth --export-path model_store --extra-files ./serve/examples/
image_classifier/index_to_name.json --handler image_classifier
```


5. TorchServe を実行してエンドポイントを開始します。> /dev/null を追加すると、ログ出力が抑止されます。

```
torchserve --start --ncs --model-store model_store --models densenet161.mar > /dev/null
```

6. 猫の画像をダウンロードして TorchServe 予測エンドポイントに送信します。

```
curl -O https://s3.amazonaws.com/model-server/inputs/kitten.jpg
curl http://127.0.0.1:8080/predictions/densenet161 -T kitten.jpg
```

予測エンドポイントは次のような上位 5 つの予測に類似する予測を JSON で返します。ここでは、エジプシャンマウが含まれている可能性が 47%、続いてトラネコが含まれている可能性が 46% となっています。

```
{
  "tiger_cat": 0.46933576464653015,
  "tabby": 0.463387668132782,
  "Egyptian_cat": 0.0645613968372345,
  "lynx": 0.0012828196631744504,
  "plastic_bag": 0.00023323058849200606
}
```

7. テストが完了したら、以下のようにサーバーを停止します。

```
torchserve --stop
```

その他の例

TorchServe には、DLAMI インスタンスで実行できるさまざまな例があります。それらの例は [TorchServe プロジェクトリポジトリの例のページ](#) で見つかります。

詳細情報

Docker で TorchServe を設定する方法や、最新の TorchServe 機能など、TorchServe のその他のドキュメントについては、GitHub で [TorchServe プロジェクトページ](#) を参照してください。

DLAMI のアップグレード

ここでは、DLAMI のアップグレードについての情報および DLAMI 上のソフトウェアを更新するためのヒントを示します。

オペレーティングシステムやその他のインストールされているソフトウェアは、パッチや更新が利用可能になり次第適用して、常に最新の状態を保ってください。

Amazon Linux または Ubuntu を使用している場合は、DLAMI にログインしたときに、更新が利用可能であれば通知され、更新の手順が表示されます。Amazon Linux のメンテナンスの詳細については、[インスタンスソフトウェアの更新](#)を参照してください。Ubuntu インスタンスの場合は、[Ubuntu の公式ドキュメント](#)を参照してください。

Windows では、Windows Update でソフトウェアとセキュリティの更新を定期的にチェックします。必要に応じて、更新が自動的に適用されるようにします。

Important

Meltdown と Spectre の脆弱性、およびオペレーティングシステムにパッチを適用して対処する方法については、[「セキュリティ情報 AWS 2018-013」](#)を参照してください。

トピック

- [DLAMI の新しいバージョンへのアップグレード](#)
- [ソフトウェア更新のヒント](#)
- [新しいアップデートの通知を受け取る](#)

DLAMI の新しいバージョンへのアップグレード

DLAMI のシステムイメージは、新しい深層学習フレームワークのリリース、CUDA やその他のソフトウェア更新、およびパフォーマンスチューニングを活用するため、定期的に更新されています。DLAMI をしばらく使用していて、更新を活用したい場合は、新しいインスタンスを起動する必要があります。また、データセット、チェックポイント、またはその他の重要なデータを手動で移行する必要があります。代わりに、Amazon EBS を使用してデータを保持し、新しい DLAMI にアタッチできます。このようにして頻繁にアップグレードし、データの移行にかかる時間を最小限に抑えることができます。

Note

DLAMI 間で Amazon EBS ボリュームをアタッチおよび移行する場合は、DLAMI と新しいボリュームの両方が同じアベイラビリティゾーンにある必要があります。

1. Amazon EC2 コンソールを使用して、新しい Amazon EBS ボリュームを作成します。詳細な手順については、[Amazon EBS ボリュームの作成](#)を参照してください。
2. 新しく作成した Amazon EBS ボリュームを既存の DLAMI にアタッチします。詳細な手順については、[Amazon EBS ボリュームのアタッチ](#)を参照してください。
3. データセット、チェックポイント、設定ファイルなどのデータを転送します。
4. DLAMI を起動します。詳細な手順については、[DLAMI インスタンスの設定](#)を参照してください。
5. 古い DLAMI から Amazon EBS ボリュームをデタッチします。詳細な手順については、[Amazon EBS ボリュームのデタッチ](#)を参照してください。
6. Amazon EBS ボリュームを新しい DLAMI にアタッチします。ステップ 2 から手順に従ってボリュームをアタッチします。
7. データが新しい DLAMI で利用できることを確認したら、古い DLAMI を停止して削除します。詳細なクリーンアップ手順については、「[DLAMI インスタンスのクリーンアップ](#)」を参照してください。

ソフトウェア更新のヒント

時として、DLAMI のソフトウェアを手動で更新する必要がある場合があります。一般的には、pip を使用して Python パッケージを更新することをお勧めします。また、Deep Learning AMI with Conda の Conda 環境内のパッケージの更新にも pip を使用してください。特定のフレームワークまたはソフトウェアのアップグレードとインストールの手順については、それぞれのウェブサイトを参照してください。

Note

パッケージが正常に更新されることは保証できません。互換性のない依存関係を持つ環境にパッケージを更新しようとする、更新が失敗する可能性があります。そのような場合は、ライブラリのメンテナンス担当者に連絡して、パッケージの依存関係を更新できるかどうかを確認してください。また、更新を許可するように環境を変更することもできます。ただ

し、このような変更を行う場合は、既存のパッケージを削除または更新する必要があり、この環境の安定性が保証されなくなります。

AWS Deep Learning AMIs には、多くの Conda 環境と多くのパッケージがプリインストールされています。プリインストールされているパッケージの数が多いため、互換性が保証されているパッケージのセットを見つけることは困難です。「環境が矛盾しています。パッケージプランをよく確認してください」という警告が表示されることがあります。DLAMI では、DLAMI が提供する環境が正しいことを保証しますが、ユーザーがインストールしたパッケージが正しく機能することは保証できません。

新しいアップデートの通知を受け取る

Note

AWS Deep Learning AMIs には、セキュリティパッチのリリース頻度が毎週あります。これらの増分セキュリティパッチについてはリリースの通知が送信されますが、公式リリースノートには含まれない場合があります。

新しい DLAMI がリリースされると毎回通知を受け取ることができます。通知は、以下のトピックを使用して [Amazon SNS](#) で公開されます。

```
arn:aws:sns:us-west-2:767397762724:dlami-updates
```

新しい DLAMI が公開されると、メッセージがここに投稿されます。このメッセージには、AMI のバージョン、メタデータ、リージョン AMI ID が含まれます。

これらのメッセージは複数の方法で受信できます。次の方法を使用することをお勧めします。

1. [\[Amazon SNS コンソール\]](#) を開きます。
2. ナビゲーションバーで、必要に応じて AWS リージョンを米国西部 (オレゴン) に変更します。サブスクライブする SNS 通知が作成されたリージョンを選択する必要があります。
3. ナビゲーションペインで、[サブスクリプション]、[サブスクリプションの作成] の順に選択します。
4. [サブスクリプションの作成] ダイアログボックスで、次の操作を行います。

- a. [トピックの ARN] には、次の Amazon リソースネーム (ARN) をコピーして貼り付けます:
arn:aws:sns:us-west-2:767397762724:dlami-updates
 - b. プロトコルで、[Amazon SQS、AWS Lambda、E メール、E メール JSON] から 1 つ選択します。
 - c. [エンドポイント] には、通知の受信に使用するリソースの E メールアドレスまたは Amazon リソースネーム (ARN) を入力します。
 - d. [Create subscription] を選択します。
5. 「AWS 通知 - サブスクリプションの確定」という件名の確認メールを受け取ります。メールを開いて [サブスクリプションを確定] を選択して受信登録を完了します。

のセキュリティ AWS Deep Learning AMIs

でのクラウドセキュリティが最優先事項 AWS です。お客様は AWS、最もセキュリティの影響を受けやすい組織の要件を満たすように構築されたデータセンターとネットワークアーキテクチャを活用できます。

セキュリティは、AWS とユーザーの間で共有される責任です。[責任共有モデル](#)では、これをクラウドのセキュリティおよびクラウド内のセキュリティと説明しています。

- クラウドのセキュリティ – AWS は、AWS のサービス で実行されるインフラストラクチャを保護する責任を担います AWS クラウド。AWS また、は、お客様が安全に使用できるサービスも提供します。サードパーティーの監査者は、[AWS コンプライアンスプログラム](#)コンプライアンスプログラムの一環として、当社のセキュリティの有効性を定期的にテストおよび検証。が適用されるコンプライアンスプログラムの詳細については AWS Deep Learning AMIs、「[コンプライアンスプログラムAWS による対象範囲内のサービスコンプライアンスプログラム](#)」を参照してください。
- クラウド内のセキュリティ – お客様の責任は、使用する によって決まり AWS のサービス ます。また、お客様は、お客様のデータの機密性、企業の要件、および適用可能な法律および規制などの他の要因についても責任を担います。

このドキュメントは、DLAMI を使用する際に責任共有モデルを適用する方法を理解するのに役立ちます 以下のトピックでは、セキュリティおよびコンプライアンスの目的を達成するために DLAMI を設定する方法を示します。また、DLAMI リソースのモニタリングと保護 AWS のサービス に役立つ他の の使用方法についても説明します。

詳細については、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 のセキュリティ](#)」を参照してください。

トピック

- [でのデータ保護 AWS Deep Learning AMIs](#)
- [の ID とアクセスの管理 AWS Deep Learning AMIs](#)
- [のコンプライアンス検証 AWS Deep Learning AMIs](#)
- [の耐障害性 AWS Deep Learning AMIs](#)
- [のインフラストラクチャセキュリティ AWS Deep Learning AMIs](#)
- [AWS Deep Learning AMIs インスタンスのモニタリング](#)

でのデータ保護 AWS Deep Learning AMIs

責任 AWS [共有モデル](#)、でのデータ保護に適用されます AWS Deep Learning AMIs。このモデルで説明されているように、AWS はすべての を実行するグローバルインフラストラクチャを保護する責任があります AWS クラウド。ユーザーは、このインフラストラクチャでホストされるコンテンツに対する管理を維持する責任があります。また、使用する「AWS のサービス」のセキュリティ設定と管理タスクもユーザーの責任となります。データプライバシーの詳細については、[データプライバシーに関するよくある質問](#)を参照してください。欧州でのデータ保護の詳細については、AWS セキュリティブログに投稿された [AWS 責任共有モデルおよび GDPR](#) のブログ記事を参照してください。

データ保護の目的で、認証情報を保護し AWS アカウント、AWS IAM Identity Center または AWS Identity and Access Management (IAM) を使用して個々のユーザーを設定することをお勧めします。この方法により、それぞれのジョブを遂行するために必要な権限のみが各ユーザーに付与されます。また、次の方法でデータを保護することもお勧めします:

- 各アカウントで多要素認証 (MFA) を使用します。
- SSL/TLS を使用して AWS リソースと通信します。TLS 1.2 が必須で、TLS 1.3 をお勧めします。
- で API とユーザーアクティビティのログ記録を設定します AWS CloudTrail。CloudTrail 証跡を使用して AWS アクティビティをキャプチャする方法については、「AWS CloudTrail ユーザーガイド」の [CloudTrail 証跡の使用](#) を参照してください。
- AWS 暗号化ソリューションと、内のすべてのデフォルトのセキュリティコントロールを使用します AWS のサービス。
- Amazon Macie などの高度な管理されたセキュリティサービスを使用します。これらは、Amazon S3 に保存されている機密データの検出と保護を支援します。
- コマンドラインインターフェイスまたは API AWS を介して にアクセスするときに FIPS 140-3 検証済み暗号化モジュールが必要な場合は、FIPS エンドポイントを使用します。利用可能な FIPS エンドポイントの詳細については、「[連邦情報処理規格 \(FIPS\) 140-3](#)」を参照してください。

お客様の E メールアドレスなどの極秘または機密情報を、タグ、または [名前] フィールドなどの自由形式のテキストフィールドに含めないことを強くお勧めします。これは、コンソール、API、または SDK を使用して DLAMI AWS CLI または他の AWS のサービスを使用する場合も同様です。AWS SDKs タグ、または名前に使用される自由記述のテキストフィールドに入力したデータは、請求または診断ログに使用される場合があります。外部サーバーに URL を提供する場合、そのサーバーへのリクエストを検証できるように、認証情報を URL に含めないことを強くお勧めします。

の ID とアクセスの管理 AWS Deep Learning AMIs

AWS Identity and Access Management (IAM) は、管理者が AWS リソースへのアクセスを安全に制御 AWS のサービス するのに役立つです。IAM 管理者は、誰が認証 (サインイン) され、DLAMI リソースを使用する認可を受ける (許可がある) ことができるかを制御します。IAM は、追加料金なしで使用できる AWS のサービス です。

Identity and Access Management の詳細については、「[Amazon EC2 の Identity and Access Management](#)」を参照してください。

トピック

- [アイデンティティを使用した認証](#)
- [ポリシーを使用したアクセスの管理](#)
- [Amazon EMR での IAM](#)

アイデンティティを使用した認証

認証は、ID 認証情報 AWS を使用して にサインインする方法です。として、IAM ユーザーとして AWS アカウントのルートユーザー、または IAM ロールを引き受けることによって、認証 (にサインイン AWS) される必要があります。

ID ソースを介して提供された認証情報を使用して、フェデレーテッド ID AWS として にサインインできます。AWS IAM Identity Center (IAM Identity Center) ユーザー、会社のシングルサインオン認証、Google または Facebook 認証情報は、フェデレーテッド ID の例です。フェデレーテッド ID としてサインインする場合、IAM ロールを使用して、前もって管理者により ID フェデレーションが設定されています。フェデレーション AWS を使用して にアクセスすると、間接的にロールを引き受けることとなります。

ユーザーの種類に応じて、AWS Management Console または AWS アクセスポータルにサインインできます。へのサインインの詳細については AWS、AWS サインイン ユーザーガイドの「[へのサインイン方法 AWS アカウント](#)」を参照してください。

AWS プログラムで にアクセスする場合、は Software Development Kit (SDK) とコマンドラインインターフェイス (CLI) AWS を提供し、認証情報を使用してリクエストを暗号化して署名します。AWS ツールを使用しない場合は、自分でリクエストに署名する必要があります。リクエストに自分で署名する推奨方法の使用については、「IAM ユーザーガイド」の「[API リクエストに対する AWS Signature Version 4](#)」を参照してください。

使用する認証方法を問わず、追加セキュリティ情報の提供をリクエストされる場合もあります。例えば、AWS では、多要素認証 (MFA) を使用してアカウントのセキュリティを向上させることをお勧めします。詳細については、「AWS IAM Identity Center ユーザーガイド」の「[多要素認証](#)」および「IAM ユーザーガイド」の「[IAM のAWS 多要素認証](#)」を参照してください。

AWS アカウント ルートユーザー

を作成するときは AWS アカウント、アカウント内のすべての およびリソースへの AWS のサービス 完全なアクセス権を持つ 1 つのサインインアイデンティティから始めます。この ID は AWS アカウント ルートユーザーと呼ばれ、アカウントの作成に使用した E メールアドレスとパスワードでサインインすることでアクセスできます。日常的なタスクには、ルートユーザーを使用しないことを強くお勧めします。ルートユーザーの認証情報は保護し、ルートユーザーでしか実行できないタスクを実行するときに使用します。ルートユーザーとしてサインインする必要があるタスクの完全なリストについては、IAM ユーザーガイドの「[ルートユーザー認証情報が必要なタスク](#)」を参照してください。

IAM ユーザーとグループ

[IAM ユーザー](#)は、単一のユーザーまたはアプリケーションに対して特定のアクセス許可 AWS アカウント を持つ 内のアイデンティティです。可能であれば、パスワードやアクセスキーなどの長期的な認証情報を保有する IAM ユーザーを作成する代わりに、一時的な認証情報を使用することをお勧めします。ただし、IAM ユーザーでの長期的な認証情報が必要な特定のユースケースがある場合は、アクセスキーをローテーションすることをお勧めします。詳細については、「IAM ユーザーガイド」の「[長期的な認証情報を必要とするユースケースのためにアクセスキーを定期的にローテーションする](#)」を参照してください。

[IAM グループ](#)は、IAM ユーザーの集団を指定するアイデンティティです。グループとしてサインインすることはできません。グループを使用して、複数のユーザーに対して一度に権限を指定できます。多数のユーザーグループがある場合、グループを使用することで権限の管理が容易になります。例えば、IAMAdmins という名前のグループを設定して、そのグループに IAM リソースを管理する許可を与えることができます。

ユーザーは、ロールとは異なります。ユーザーは 1 人の人または 1 つのアプリケーションに一意に関連付けられますが、ロールはそれを必要とする任意の人が引き受けるようになっています。ユーザーには永続的な長期の認証情報がありますが、ロールでは一時認証情報が提供されます。詳細については、「IAM ユーザーガイド」の「[IAM ユーザーに関するユースケース](#)」を参照してください。

IAM ロール

[IAM ロール](#)は、特定のアクセス許可 AWS アカウント を持つ 内の ID です。これは IAM ユーザーに似ていますが、特定のユーザーには関連付けられていません。で IAM ロールを一時的に引き受けるには AWS Management Console、[ユーザーから IAM ロール \(コンソール\) に切り替える](#) ことができます。ロールを引き受けるには、または AWS API オペレーションを AWS CLI 呼び出すか、カスタム URL を使用します。ロールを使用する方法の詳細については、「IAM ユーザーガイド」の「[ロールを引き受けるための各種方法](#)」を参照してください。

IAM ロールと一時的な認証情報は、次の状況で役立ちます:

- フェデレーションユーザーアクセス - フェデレーティッド ID に許可を割り当てるには、ロールを作成してそのロールの許可を定義します。フェデレーティッド ID が認証されると、その ID はロールに関連付けられ、ロールで定義されている許可が付与されます。フェデレーションのロールについては、「IAM ユーザーガイド」の「[サードパーティー ID プロバイダー \(フェデレーション\) 用のロールを作成する](#)」を参照してください。IAM Identity Center を使用する場合は、許可セットを設定します。アイデンティティが認証後にアクセスできるものを制御するため、IAM Identity Center は、権限セットを IAM のロールに関連付けます。アクセス許可セットの詳細については、「AWS IAM Identity Center User Guide」の「[Permission sets](#)」を参照してください。
- 一時的な IAM ユーザー権限 - IAM ユーザーまたはロールは、特定のタスクに対して複数の異なる権限を一時的に IAM ロールで引き受けることができます。
- クロスアカウントアクセス - IAM ロールを使用して、自分のアカウントのリソースにアクセスすることを、別のアカウントの人物 (信頼済みプリンシパル) に許可できます。クロスアカウントアクセス権を付与する主な方法は、ロールを使用することです。ただし、一部の では AWS のサービス、(プロキシとしてロールを使用する代わりに) リソースに直接ポリシーをアタッチできます。クロスアカウントアクセスにおけるロールとリソースベースのポリシーの違いについては、「IAM ユーザーガイド」の「[IAM でのクロスアカウントのリソースへのアクセス](#)」を参照してください。
- クロスサービスアクセス — 一部の は、他の の機能 AWS のサービス を使用します AWS のサービス。例えば、あるサービスで呼び出しを行うと、通常そのサービスによって Amazon EC2 でアプリケーションが実行されたり、Amazon S3 にオブジェクトが保存されたりします。サービスでは、呼び出し元プリンシパルの許可、サービスロール、またはサービスリンクロールを使用してこれを行う場合があります。
- 転送アクセスセッション (FAS) - IAM ユーザーまたはロールを使用して でアクションを実行すると AWS、プリンシパルと見なされます。一部のサービスを使用する際に、アクションを実行することで、別のサービスの別のアクションがトリガーされることがあります。FAS は、 を呼び出すプリンシパルのアクセス許可と AWS のサービス、ダウンストリームサービス AWS の

サービスへのリクエストのリクエストをリクエストするを使用します。FAS リクエストは、サービスが他の AWS のサービス または リソースとのやり取りを完了する必要があるリクエストを受け取った場合にのみ行われます。この場合、両方のアクションを実行するためのアクセス許可が必要です。FAS リクエストを行う際のポリシーの詳細については、「[転送アクセスセッション](#)」を参照してください。

- サービスロール - サービスがユーザーに代わってアクションを実行するために引き受ける [IAM ロール](#)です。IAM 管理者は、IAM 内からサービスロールを作成、変更、削除することができます。詳細については、「IAM ユーザーガイド」の「[AWS のサービスに許可を委任するロールを作成する](#)」を参照してください。
- サービスにリンクされたロール - サービスにリンクされたロールは、にリンクされたサービスロールの一種です AWS のサービス。サービスは、ユーザーに代わってアクションを実行するロールを引き受けることができます。サービスにリンクされたロールは に表示され AWS アカウント、サービスによって所有されます。IAM 管理者は、サービスリンクロールのアクセス許可を表示できますが、編集することはできません。
- Amazon EC2 で実行されているアプリケーション - IAM ロールを使用して、EC2 インスタンスで実行され、AWS CLI または AWS API リクエストを行うアプリケーションの一時的な認証情報を管理できます。これは、EC2 インスタンス内でのアクセスキーの保存に推奨されます。AWS ロールを EC2 インスタンスに割り当て、そのすべてのアプリケーションで使用できるようにするには、インスタンスにアタッチされたインスタンスプロファイルを作成します。インスタンスプロファイルにはロールが含まれ、EC2 インスタンスで実行されるプログラムは一時的な認証情報を取得できます。詳細については、「IAM ユーザーガイド」の「[Amazon EC2 インスタンスで実行されるアプリケーションに IAM ロールを使用して許可を付与する](#)」を参照してください。

ポリシーを使用したアクセスの管理

でアクセスを制御する AWS には、ポリシーを作成し、ID AWS または リソースにアタッチします。ポリシーは AWS、アイデンティティまたはリソースに関連付けられているときにアクセス許可を定義する のオブジェクトです。 は、プリンシパル (ユーザー、ルートユーザー、またはロールセッション) がリクエストを行うときに、これらのポリシー AWS を評価します。ポリシーでの権限により、リクエストが許可されるか拒否されるかが決まります。ほとんどのポリシーは JSON ドキュメント AWS として に保存されます。JSON ポリシードキュメントの構造と内容の詳細については、IAM ユーザーガイドの [JSON ポリシー概要](#)を参照してください。

管理者は JSON AWS ポリシーを使用して、誰が何にアクセスできるかを指定できます。つまり、どのプリンシパルがどのリソースに対してどのような条件下でアクションを実行できるかということです。

デフォルトでは、ユーザーやロールに権限はありません。IAM 管理者は、リソースで必要なアクションを実行するための権限をユーザーに付与する IAM ポリシーを作成できます。その後、管理者はロールに IAM ポリシーを追加し、ユーザーはロールを引き受けることができます。

IAM ポリシーは、オペレーションの実行方法を問わず、アクションの許可を定義します。例えば、iam:GetRole アクションを許可するポリシーがあるとします。そのポリシーを持つユーザーは、AWS Management Console、AWS CLI または AWS API からロール情報を取得できます。

アイデンティティベースのポリシー

アイデンティティベースポリシーは、IAM ユーザーグループ、ユーザーのグループ、ロールなど、アイデンティティにアタッチできる JSON 許可ポリシードキュメントです。これらのポリシーは、ユーザーとロールが実行できるアクション、リソース、および条件をコントロールします。アイデンティティベースポリシーの作成方法については、「IAM ユーザーガイド」の「[カスタマー管理ポリシーでカスタム IAM アクセス許可を定義する](#)」を参照してください。

アイデンティティベースのポリシーは、さらにインラインポリシーまたはマネージドポリシーに分類できます。インラインポリシーは、単一のユーザー、グループ、またはロールに直接埋め込まれています。管理ポリシーは、内の複数のユーザー、グループ、ロールにアタッチできるスタンドアロンポリシーです AWS アカウント。管理ポリシーには、AWS 管理ポリシーとカスタマー管理ポリシーが含まれます。マネージドポリシーまたはインラインポリシーのいずれかを選択する方法については、「IAM ユーザーガイド」の「[管理ポリシーとインラインポリシーのいずれかを選択する](#)」を参照してください。

リソースベースのポリシー

リソースベースのポリシーは、リソースに添付する JSON ポリシードキュメントです。リソースベースのポリシーには例として、IAM ロールの信頼ポリシーや Amazon S3 バケットポリシーがあげられます。リソースベースのポリシーをサポートするサービスでは、サービス管理者はポリシーを使用して特定のリソースへのアクセスを制御できます。ポリシーがアタッチされているリソースの場合、指定されたプリンシパルがそのリソースに対して実行できるアクションと条件は、ポリシーによって定義されます。リソースベースのポリシーでは、[プリンシパルを指定する](#)必要があります。プリンシパルには、アカウント、ユーザー、ロール、フェデレーティッドユーザー、またはを含めることができます AWS のサービス。

リソースベースのポリシーは、そのサービス内にあるインラインポリシーです。リソースベースのポリシーでは、IAM の AWS マネージドポリシーを使用できません。

アクセスコントロールリスト (ACL)

アクセスコントロールリスト (ACL) は、どのプリンシパル (アカウントメンバー、ユーザー、またはロール) がリソースにアクセスするための許可を持つかを制御します。ACL はリソースベースのポリシーに似ていますが、JSON ポリシードキュメント形式は使用しません。

Amazon S3、AWS WAF、および Amazon VPC は、ACLs。ACL の詳細については、「Amazon Simple Storage Service デベロッパーガイド」の「[アクセスコントロールリスト \(ACL\) の概要](#)」を参照してください。

その他のポリシータイプ

AWS は、一般的でない追加のポリシータイプをサポートします。これらのポリシータイプでは、より一般的なポリシータイプで付与された最大の権限を設定できます。

- **アクセス許可の境界** - アクセス許可の境界は、アイデンティティベースポリシーによって IAM エンティティ (IAM ユーザーまたはロール) に付与できる権限の上限を設定する高度な機能です。エンティティにアクセス許可の境界を設定できます。結果として得られる権限は、エンティティのアイデンティティベースポリシーとそのアクセス許可の境界の共通部分になります。Principal フィールドでユーザーまたはロールを指定するリソースベースのポリシーでは、アクセス許可の境界は制限されません。これらのポリシーのいずれかを明示的に拒否した場合、権限は無効になります。アクセス許可の境界の詳細については、「IAM ユーザーガイド」の「[IAM エンティティのアクセス許可の境界](#)」を参照してください。
- **サービスコントロールポリシー (SCPs)** - SCPs は、の組織または組織単位 (OU) の最大アクセス許可を指定する JSON ポリシーです AWS Organizations。AWS Organizations は、ビジネスが所有する複数のをグループ化して一元管理するためのサービス AWS アカウントです。組織内のすべての機能を有効にすると、サービスコントロールポリシー (SCP) を一部またはすべてのアカウントに適用できます。SCP は、各を含むメンバーアカウントのエンティティのアクセス許可を制限します AWS アカウントのルートユーザー。Organizations と SCP の詳細については、「AWS Organizations ユーザーガイド」の「[サービスコントロールポリシー \(SCP\)](#)」を参照してください。
- **リソースコントロールポリシー (RCP)** - RCP は、所有する各リソースにアタッチされた IAM ポリシーを更新することなく、アカウント内のリソースに利用可能な最大数のアクセス許可を設定するために使用できる JSON ポリシーです。RCP は、メンバーアカウントのリソースに対するアクセス許可を制限し、組織に属しているかどうかにかかわらず AWS アカウントのルートユーザー、を含む ID に対する有効なアクセス許可に影響を与える可能性があります。RCP AWS のサービスをサポートする のリストを含む Organizations と RCPs [「リソースコントロールポリシー \(RCPs\)」](#) を参照してください。AWS Organizations

- セッションポリシー - セッションポリシーは、ロールまたはフェデレーションユーザーの一時的なセッションをプログラムで作成する際にパラメータとして渡す高度なポリシーです。結果としてセッションの権限は、ユーザーまたはロールのアイデンティティベースポリシーとセッションポリシーの共通部分になります。また、リソースベースのポリシーから権限が派生する場合があります。これらのポリシーのいずれかを明示的に拒否した場合、権限は無効になります。詳細については、「IAM ユーザーガイド」の「[セッションポリシー](#)」を参照してください。

複数のポリシータイプ

1つのリクエストに複数のタイプのポリシーが適用されると、結果として作成される権限を理解するのがさらに難しくなります。複数のポリシータイプが関係する場合にリクエストを許可するかどうかが AWS を決定する方法については、IAM ユーザーガイドの「[ポリシー評価ロジック](#)」を参照してください。

Amazon EMR での IAM

Amazon EMR で IAM を使用して、ユーザー、AWS リソース、グループ、ロール、ポリシーを定義できます。AWS のサービス これらのユーザーとロールがアクセスできるものを制御することもできます。

Amazon EMR での IAM の詳細については、「[Amazon EMR のAWS Identity and Access Management](#)」を参照してください。

のコンプライアンス検証 AWS Deep Learning AMIs

サードパーティーの監査者は、複数のコンプライアンスプログラムの一環として AWS Deep Learning AMIs のセキュリティと AWS コンプライアンスを評価します。サポートされているコンプライアンスプログラムの詳細については、「[Amazon EC2 のコンプライアンス検証](#)」を参照してください。

特定のコンプライアンスプログラム AWS のサービスの範囲内の のリストについては、「[コンプライアンスAWS プログラムによる対象範囲内のサービスコンプライアンス](#)」を参照してください。一般的な情報については、[AWS 「コンプライアンスプログラム」](#)を参照してください。

を使用して、サードパーティーの監査レポートをダウンロードできます AWS Artifact。詳細については、「[AWS Artifact でレポートをダウンロードする](#)」、「」を参照してください。

DLAMI を使用する場合のお客様のコンプライアンス責任は、お客様のデータの機密性、貴社のコンプライアンス目的、適用される法律および規制によって決まります。は、コンプライアンスに役立つ以下のリソース AWS を提供します。

- [「セキュリティ & コンプライアンス クイックリファレンスガイド」](#) – これらのデプロイガイドには、アーキテクチャ上の考慮事項の説明と、AWS でセキュリティとコンプライアンスに重点を置いたベースライン環境をデプロイするための手順が記載されています。
- [AWS コンプライアンスリソース](#) – このワークブックとガイドのコレクションは、お客様の業界と地域に適用される場合があります。
- [「デベロッパーガイド」の AWS Config 「ルールによるリソースの評価」](#) – この AWS Config サービスは、リソース設定が内部プラクティス、業界ガイドライン、および規制にどの程度準拠しているかを評価します。AWS Config
- [AWS Security Hub](#) – これにより AWS のサービス、内のセキュリティ状態を包括的に把握できます AWS。Security Hub は、セキュリティコントロールを使用して AWS リソースを評価し、セキュリティ業界標準とベストプラクティスに対するコンプライアンスをチェックします。

の耐障害性 AWS Deep Learning AMIs

AWS グローバルインフラストラクチャは、AWS リージョン およびアベイラビリティゾーンを中心に構築されています。は、低レイテンシー、高スループット、および高度に冗長なネットワークに接続された、物理的に分離および分離された複数のアベイラビリティゾーン AWS リージョンを提供します。アベイラビリティゾーンでは、ゾーン間で中断することなく自動的にフェイルオーバーするアプリケーションとデータベースを設計および運用することができます。アベイラビリティゾーンは、従来の単一または複数のデータセンターインフラストラクチャよりも可用性が高く、フォールトトレラントで、スケーラブルです。

AWS リージョン およびアベイラビリティゾーンの詳細については、[AWS 「グローバルインフラストラクチャ」](#) を参照してください。

データの耐障害性とバックアップのニーズをサポートするための Amazon EC2 機能については、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 の耐障害性](#)」を参照してください。

のインフラストラクチャセキュリティ AWS Deep Learning AMIs

のインフラストラクチャセキュリティ AWS Deep Learning AMIs は Amazon EC2 によってサポートされています。詳細については、「[Amazon EC2 ユーザーガイド](#)」の「Amazon EC2 でのインフラストラクチャセキュリティ」を参照してください。

AWS Deep Learning AMIs インスタンスのモニタリング

モニタリングは、AWS Deep Learning AMIs インスタンスやその他の AWS ソリューションの信頼性、可用性、パフォーマンスを維持する上で重要な部分です。DLAMI インスタンスには、GPU 使用統計を Amazon CloudWatch に報告するユーティリティを含む、いくつかの GPU モニタリングツールが付属しています。詳細については、「[GPU のモニタリングおよび最適化](#)」および「Amazon EC2 ユーザーガイド」の「[Amazon EC2 リソースのモニタリング](#)」を参照してください。

DLAMI インスタンスの使用状況の追跡をオプトアウトする

次の AWS Deep Learning AMIs オペレーティングシステムディストリビューションには、AWS がインスタンスタイプ、インスタンス ID、DLAMI タイプ、および OS 情報を収集できるようにするコードが含まれています。

Note

AWS は、DLAMI 内で使用するコマンドなど、DLAMI に関するその他の情報を収集または保持しません。

- Amazon Linux 2
- Amazon Linux 2023
- Ubuntu 20.04
- Ubuntu 22.04

使用状況の追跡をオプトアウトするには

選択すると、新しい DLAMI インスタンスの使用状況の追跡をオプトアウトできます。オプトアウトするには、起動時に Amazon EC2 インスタンスにタグを追加する必要があります。タグでは、キー `OPT_OUT_TRACKING` を使用し、関連した値を `true` に設定する必要があります。詳細については、「Amazon EC2 ユーザーガイド」の「[Amazon EC2 リソースのタグ付け](#)」を参照してください。

DLAMI サポートポリシー

ここでは、AWS Deep Learning AMIs (DLAMI) のサポートポリシーの詳細を確認できます。

AWS 現在サポートされている DLAMI フレームワークとオペレーティングシステムのリストについては、[DLAMI サポートポリシーページ](#)を参照してください。以下の用語は、サポートポリシーページとこのページに記載されているすべての DLAMIs に適用されます。

- 現行バージョンでは、フレームワークのバージョンを x.y.z 形式で指定します。この形式では、x はメジャーバージョン、y はマイナーバージョン、z はパッチバージョンを指します。例えば、TensorFlow 2.10.1 では、メジャーバージョンは 2、マイナーバージョンは 10、パッチバージョンは 1 です。
- パッチの終了は、特定のフレームワークまたはオペレーティングシステムのバージョン AWS をサポートする期間を指定します。

特定の DLAMI の詳細は、「[DLAMI のリリースノート](#)」を参照してください。

DLAMI サポートFAQs

- [どのフレームワークバージョンにセキュリティパッチが適用されますか？](#)
- [セキュリティパッチはどのオペレーティングシステムで取得されますか？](#)
- [新しいフレームワークバージョンがリリースされると、どのイメージが AWS 公開されますか？](#)
- [どのイメージに新しい SageMaker AI/AWS 機能が追加されますか？](#)
- [サポート対象フレームワークの表では、現在のバージョンはどのように定義されていますか？](#)
- [サポート対象テーブルにないバージョンを実行している場合はどうなりますか？](#)
- [DLAMIs フレームワークバージョンの以前のパッチバージョンをサポートしていますか？](#)
- [サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにありますか？](#)
- [新しいイメージはどのくらいの頻度でリリースされますか？](#)
- [ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか？](#)
- [新しいパッチが適用されたフレームワークバージョン、または更新されたフレームワークバージョンが利用可能になった場合はどうなりますか？](#)

- [フレームワークのバージョンを変更せずに依存関係は更新されますか？](#)
- [使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか？](#)
- [アクティブにメンテナンスされなくなったフレームワークバージョンのイメージにはパッチが適用されますか？](#)
- [古いフレームワークバージョンを使用するにはどうすればよいですか？](#)
- [フレームワークとそのバージョンでのサポート変更に関する最新情報を得るにはどうすればよいですか？](#)
- [Anaconda リポジトリを使用するには商用ライセンスが必要ですか？](#)

どのフレームワークバージョンにセキュリティパッチが適用されますか？

フレームワークバージョンがサポート[AWS Deep Learning AMIs ポリシーテーブル](#)のサポートされているフレームワークバージョンにある場合、セキュリティパッチが取得されます。

セキュリティパッチはどのオペレーティングシステムで取得されますか？

オペレーティングシステムがサポート[AWS Deep Learning AMIs ポリシー表](#)のサポートされているオペレーティングシステムのバージョンにリストされている場合、セキュリティパッチが取得されます。

新しいフレームワークバージョンがリリースされると、どのイメージが AWS 公開されますか？

TensorFlow と PyTorch の新しいバージョンがリリースされると、すぐに新しい DLAMI が公開されます。これには、フレームワークのメジャーバージョン、メジャー/マイナーバージョン、メジャー/マイナー/パッチバージョンが含まれます。また、新しいバージョンのドライバーやライブラリが利用可能になった場合も、イメージが更新されます。イメージメンテナンスの詳細については、「[使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか？](#)」を参照してください。

どのイメージに新しい SageMaker AI/AWS 機能が追加されますか？

新機能は通常、PyTorch と TensorFlow 向けの最新バージョンの DLAMI でリリースされます。新しい SageMaker AI または AWS 機能の詳細については、特定のイメージのリリースノートを参照してください。使用可能な DLAMI のリストについては、「[DLAMI のリリースノート](#)」を参照してください。

い。イメージメンテナンスの詳細については、「[使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか?](#)」を参照してください。

サポート対象フレームワークの表では、現在のバージョンはどのように定義されていますか？

[AWS Deep Learning AMIs サポートポリシー表](#)の最新バージョンは、が GitHub で利用可能 AWS にする最新のフレームワークバージョンを参照しています。各最新リリースには、DLAMI のドライバー、ライブラリ、関連パッケージの更新が含まれています。イメージメンテナンスの詳細については、「[使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか?](#)」を参照してください。

サポート対象テーブルにないバージョンを実行している場合はどうなりますか？

[AWS Deep Learning AMIs サポートポリシーテーブル](#)にないバージョンを実行している場合、最新のドライバー、ライブラリ、関連パッケージがない可能性があります。up-to-dateバージョンについては、任意の最新の DLAMI を使用して、サポートされているフレームワークまたはオペレーティングシステムのいずれかにアップグレードすることをお勧めします。使用可能な DLAMI のリストについては、「[DLAMI のリリースノート](#)」を参照してください。

DLAMIs フレームワークバージョンの以前のパッチバージョンをサポートしていますか？

いいえ。サポート [AWS Deep Learning AMIs ポリシー表](#)に記載されているように、最初の GitHub リリースから 365 日後にリリースされた各フレームワークの最新バージョンのメジャーバージョンの最新のパッチバージョンをサポートしています。詳細については、「[サポート対象テーブルにないバージョンを実行している場合はどうなりますか?](#)」を参照してください。

サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにありますか？

最新のフレームワークバージョンで DLAMI を使用するには、CLI AWS または SSM パラメータを使用して [DLAMI ID](#) を取得し、[E2C2コンソール](#)を使用して DLAMI を起動します。AWS Deep Learning AMIs ID を取得するためのサンプル CLI AWS または SSM パラメータコマンドについては、DLAMI リリースノートページの[単一フレームワーク DLAMI リリースノート](#)を参照してください。選択したフレームワークバージョンは、サポート [AWS Deep Learning AMIs ポリシー表](#)のサポートされているフレームワークバージョンにリストされている必要があります。

新しいイメージはどのくらいの頻度でリリースされますか？

最新のパッチバージョンを提供することは、AWS の最優先事項です。パッチを適用したイメージを、できるだけ早く定期的に作成しています。新たにパッチが適用されたフレームワークのバージョン (例: TensorFlow 2.9 から TensorFlow 2.9.1) や新しいマイナーリリースバージョン (例: TensorFlow 2.9 から TensorFlow 2.10) がないかチェックして、できるだけ早く公開するようにしています。既存のバージョンの TensorFlow が新しいバージョンの CUDA と共にリリースされると、AWS は新しい CUDA バージョンをサポートする、そのバージョンの TensorFlow 用の新しい DLAMI をリリースします。

ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか？

いいえ。DLAMI のパッチ更新は「インプレース」更新ではありません。

新しい EC2 インスタンスをオンにし、ワークロードとスクリプトを移行してから、以前のインスタンスをオフにする必要があります。

新しいパッチが適用されたフレームワークバージョン、または更新されたフレームワークバージョンが利用可能になった場合はどうなりますか？

DLAMI の変更を通知するには、関連する DLAMI の通知をサブスクライブしてください。[「新しい更新に関する通知の受信」](#)を参照してください。

フレームワークのバージョンを変更せずに依存関係は更新されますか？

AWS はフレームワークのバージョンを変更せずに依存関係を更新します。ただし、依存関係の更新によって互換性が失われる場合、AWS は別のバージョンでイメージ作成します。更新された依存関係の情報については、必ず「[DLAMI のリリースノート](#)」を確認してください。

使用しているフレームワークバージョンに対する有効なサポートはいつ終了しますか？

DLAMI イメージはイミュータブルです。一度作成されると変更されません。フレームワークバージョンの有効なサポートが終了する主な理由は 4 つあります。

- [フレームワークバージョン \(パッチ\) のアップグレード](#)
- [AWS セキュリティパッチ](#)

- [パッチ終了日 \(エージングアウト\)](#)
- [依存関係のサポート終了](#)

Note

バージョンパッチのアップグレードやセキュリティパッチは頻繁に行われるため、DLAMI のリリースノートページを頻繁に確認し、変更があった場合はアップグレードすることをお勧めします。

フレームワークバージョン (パッチ) のアップグレード

TensorFlow 2.7.0 に基づく DLAMI ワークロードがあり、GitHub で TensorFlow リリースバージョン 2.7.1 を使用している場合、は TensorFlow 2.7.1 の新しい DLAMI を AWS リリースします。TensorFlow 2.7.1 の新しいイメージがリリースされると、2.7.0 の以前のイメージはアクティブにメンテナンスされなくなります。TensorFlow 2.7.0 の DLAMI には、以降のパッチは適用されません。TensorFlow 2.7 の DLAMI リリースノートページが最新の情報に更新されます。マイナーパッチごとの個別のリリースノートページはありません。

パッチアップグレードにより作成された新しい DLAMI には、新しい [AMI ID](#) が割り当てられます。

AWS セキュリティパッチ

TensorFlow 2.7.0 のイメージに基づくワークロードがあり、セキュリティパッチ AWS を作成する場
合、新しいバージョンの DLAMI が TensorFlow 2.7.0 用にリリースされます。TensorFlow 2.7.0 の
以前のバージョンのイメージは、アクティブにメンテナンスされなくなります。詳細については、
「[ワークロードの実行中にインスタンスにインプレースでパッチが適用されますか?](#)」を参照してく
ださい。最新の DLAMI を検索する手順については、「[サポートされるフレームワークバージョン用の最新のパッチ適用済みイメージはどこにありますか?](#)」を参照してください。

パッチアップグレードにより作成された新しい DLAMI には、新しい [AMI ID](#) が割り当てられます。

パッチ終了日 (エージングアウト)

DLAMI のパッチの終了日は、GitHub のリリース日から 365 日後です。

[マルチフレームワーク DLAMI](#) の場合、いずれかのフレームワークバージョンが更新されると、更新されたバージョンを含む新しい DLAMI が必要になります。古いフレームワークバージョンの DLAMI は、アクティブにメンテナンスされなくなります。

⚠ Important

フレームワークのメジャー更新がある場合は例外とします。例えば、TensorFlow 1.15 が TensorFlow 2.0 に更新された場合、GitHub のリリース日から 2 年間、またはオリジンフレームワークのメンテナンスチームがサポートを終了してから 6 か月後のいずれか早い方の日付まで、TensorFlow 1.15 の最新バージョンを引き続きサポートします。

依存関係のサポート終了

Python 3.6 で TensorFlow 2.7.0 DLAMI イメージでワークロードを実行していて、そのバージョンの Python がサポート終了とマークされている場合、Python 3.6 をベースとするすべての DLAMI イメージはアクティブにメンテナンスされなくなります。同様に、Ubuntu 16.04 のような OS バージョンがサポート終了とマークされている場合、Ubuntu 16.04 に依存するすべての DLAMI イメージはアクティブにメンテナンスされなくなります。

アクティブにメンテナンスされなくなったフレームワークバージョンのイメージにはパッチが適用されますか？

いいえ。アクティブにメンテナンスされていないイメージには新しいリリースは適用されません。

古いフレームワークバージョンを使用するにはどうすればよいですか？

古いフレームワークバージョンの DLAMI を使用するには、[DLAMI ID](#) を取得し、その ID を使用して [EC2 コンソール](#) を使用して DLAMI を起動します。AMI ID を取得するための AWS CLI コマンドについては、[単一フレームワーク DLAMI リリースノートのリリースノート](#) ページを参照してください。

フレームワークとそのバージョンでのサポート変更に関する最新情報を得るにはどうすればよいですか？

「[DLAMI リリースノート](#)」の [AWS Deep Learning AMIs フレームワークサポートポリシーの表](#) で、DLAMI フレームワークとバージョンの最新情報を参照してください。

Anaconda リポジトリを使用するには商用ライセンスが必要ですか？

Anaconda は特定のユーザー向けの商用ライセンスモデルに移行しました。アクティブにメンテナンスされている DLAMI は、Anaconda チャンネルから公開されているオープンソースバージョンの Conda ([conda-forge](#)) に移行されました。

DLAMI に対する NVIDIA ドライバーの重要な変更点

2023 年 11 月 15 日、は DLAMI が使用する NVIDIA ドライバーに関連する AWS Deep Learning AMIs (DLAMI) に重要な変更 AWS を加えました。DLAMIs 変更内容と、DLAMI の使用に影響するかどうかについては「[DLAMI の NVIDIA ドライバーの変更に関するよくある質問](#)」を参照してください。

DLAMI の NVIDIA ドライバーの変更に関するよくある質問

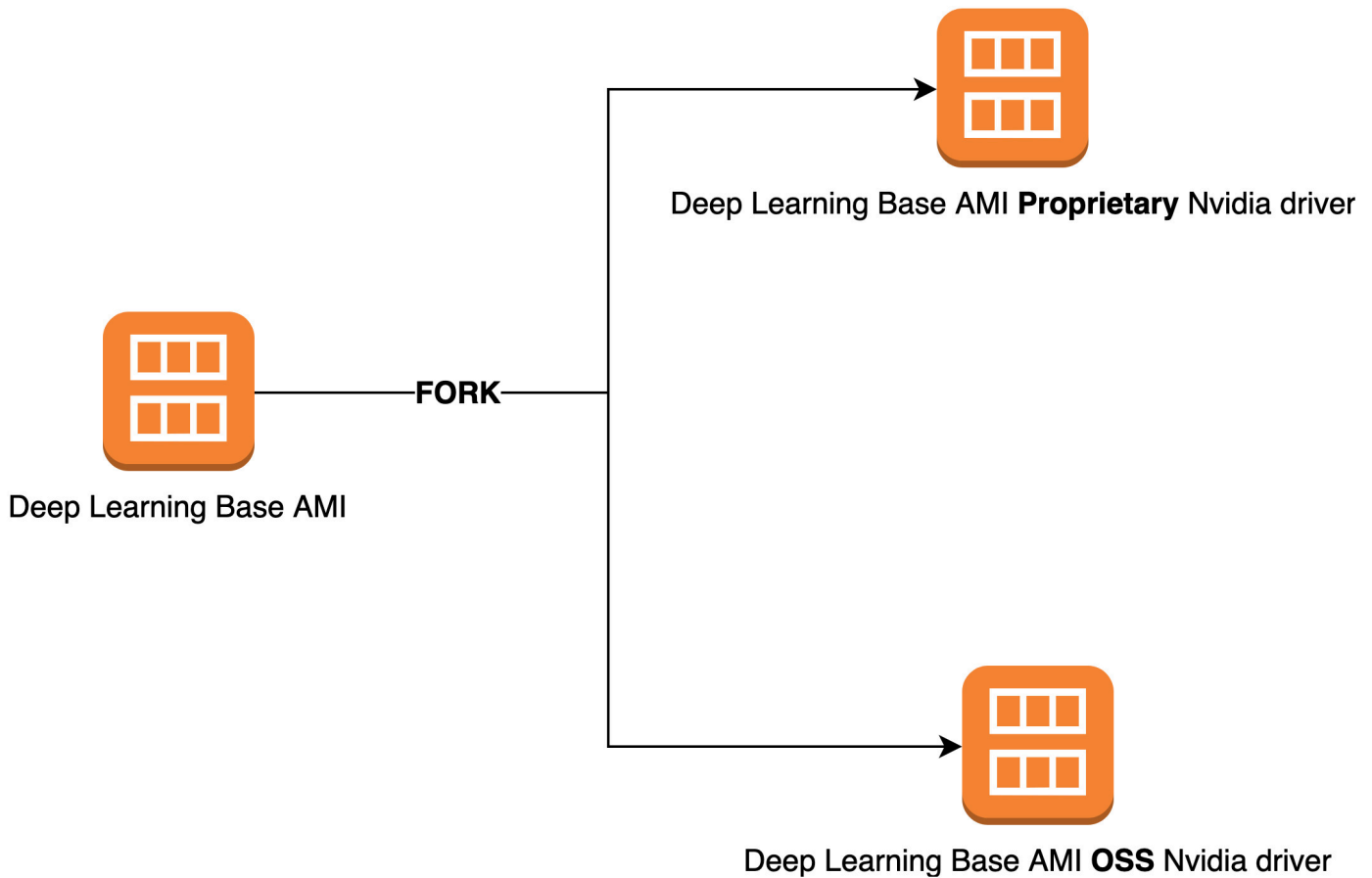
- [何が変わったのですか？](#)
- [この変更が行われた理由は何ですか？](#)
- [この変更により影響を受けるのはどの DLAMI ですか？](#)
- [ユーザーにとってこの変更にはどのような意味がありますか？](#)
- [新しい DLAMI で失われる機能はありますか？](#)
- [この変更は Deep Learning Containers に影響しますか？](#)

何が変わったのですか？

AWS は DLAMI を次の 2 つのグループに分割しました。

- NVIDIA 独自のドライバーを使用する DLAMI (P3、P3dn、G3 をサポート)
- NVIDIA OSS ドライバーを使用する DLAMI (G4dn、G5、P4、P5 をサポート)

その結果、2 つのカテゴリそれぞれに、新しい名前と新しい AMI ID を持つ新しい DLAMI を作成しました。これらの DLAMI には互換性がありません。つまり、あるグループの DLAMI は、他のグループがサポートするインスタンスをサポートしません。例えば、P5 をサポートする DLAMI は G3 をサポートせず、G3 をサポートする DLAMI は P5 をサポートしません。



この変更が行われた理由は何ですか？

以前、NVIDIA GPU 用の DLAMI には NVIDIA 独自のカーネルドライバーが含まれていました。しかし、アップストリームの Linux カーネルコミュニティにより、NVIDIA GPU ドライバーなどの独自のカーネルドライバーが他のカーネルドライバーと通信できないようにする変更が受け入れられました。この変更により、P4 および P5 シリーズのインスタンスの GPUDirect RDMA が無効になります。これは、GPU が分散トレーニングに EFA を効率的に使用できるようにするメカニズムです。その結果、DLAMI は OpenRM ドライバー (NVIDIA のオープンソースドライバー) を使用するようになり、G4dn、G5、P4、P5 をサポートするためにオープンソース EFA ドライバーにリンクされました。しかし、この OpenRM ドライバーは古いインスタンス (P3 や G3 など) をサポートしていません。そのため AWS では、今後も両方のインスタンスタイプをサポートする最新で高パフォーマンスの安全な DLAMI を提供するために、DLAMI を 2 つのグループに分割しました。1 つには OpenRM ドライバー (G4dn、G5、P4、P5 をサポート) が、もう 1 つには従来の独自のドライバー (P3、P3dn、G3 をサポート) があります。

この変更により影響を受けるのはどの DLAMI ですか？

この変更はすべての DLAMI に影響します。

ユーザーにとってこの変更にはどのような意味がありますか？

サポートされている Amazon Elastic Compute Cloud (Amazon EC2) インスタンスタイプで実行されている限り、引き続きすべての DLAMI は機能、パフォーマンス、セキュリティを提供します。DLAMI がサポートする EC2 インスタンスタイプを確認するには、その DLAMI のリリースノートを確認してから、サポートされている EC2 インスタンスを探します。現在サポートされている DLAMI オプションのリストとリリースノートへのリンクについては、「[DLAMI のリリースノート](#)」を参照してください。

さらに、現在の DLAMIs を呼び出すには、正しい AWS Command Line Interface (AWS CLI) コマンドを使用する必要があります。

P3、P3dn、G3 をサポートするベース DLAMI の場合は、次のコマンドを使用します。

```
aws ec2 describe-images --region us-east-1 --owners amazon \  
--filters 'Name=name,Values=Deep Learning Base Proprietary Nvidia Driver AMI (Amazon Linux 2) Version ??.' 'Name=state,Values=available' \  
--query 'reverse(sort_by(Images, &CreationDate))[:1].ImageId' --output text
```

G4dn、G5、P4、P5 をサポートするベース DLAMI の場合は、次のコマンドを使用します。

```
aws ec2 describe-images --region us-east-1 --owners amazon \  
--filters 'Name=name,Values=Deep Learning Base OSS Nvidia Driver AMI (Amazon Linux 2) Version ??.' 'Name=state,Values=available' \  
--query 'reverse(sort_by(Images, &CreationDate))[:1].ImageId' --output text
```

新しい DLAMI で失われる機能はありますか？

いいえ、機能が失われることはありません。現在の DLAMI は、サポートされている EC2 インスタンスタイプで実行すれば、これまでの DLAMI のすべての機能、パフォーマンス、セキュリティを提供します。

この変更は Deep Learning Containers に影響しますか？

いいえ。この変更は NVIDIA ドライバーが含まれていないため、AWS 深層学習コンテナには影響しませんでした。ただし、Deep Learning Containers は、必ず基盤となるインスタンスと互換性のある AMI で実行してください。

DLAMI の関連情報

DLAMI の関連情報を含むその他のリソースは、AWS Deep Learning AMIs デベロッパーガイドの外部にあります。で AWS re:Post、DLAMI に関する他のお客様からの質問を確認するか、独自の質問をしてください。AWS Machine Learning ブログやその他の AWS ブログで、DLAMI に関する公式投稿をお読みください。

AWS re:Post

[タグ : AWS Deep Learning AMIs](#)

AWS ブログ

- [AWS Machine Learning ブログ | カテゴリ : AWS Deep Learning AMIs](#)
- [AWS Machine Learning ブログ | Amazon EC2 C5 および P3 インスタンスで最適化された TensorFlow 1.6 を使用したトレーニングの高速化](#)
- [AWS Machine Learning ブログ | Machine Learning プラクティショナー AWS Deep Learning AMIs 向けの新機能](#)
- [AWS Partner Network \(APN\) ブログ | 新しいトレーニングコースが利用可能: での Machine Learning と深層学習の概要 AWS](#)
- [AWS ニュースブログ | を使用した深層学習へのジャーニー AWS](#)

DLAMI の非推奨の機能

次の表に、AWS Deep Learning AMIs (DLAMI) の非推奨機能、非推奨になった日付、および非推奨になった理由の詳細を示します。

機能	日付	詳細
Ubuntu 16.04	10/07/2021	Ubuntu Linux 16.04 LTS は、2021 年 4 月 30 日に 5 年間の LTS ウィンドウが終了し、ベンダーによってサポートされなくなりました。2021 年 10 月から、新規リリースでの Deep Learning Base AMI (Ubuntu 16.04) に対する更新はなくなりました。以前のリリースは、引き続き利用可能です。
Amazon Linux	10/07/2021	Amazon Linux は 2020 年 12 月から サポート終了 です。2021 年 10 月から、新規リリースでの Deep Learning AMI (Amazon Linux) に対する更新はなくなりました。Deep Learning AMI (Amazon Linux) の以前のリリースは引き続き利用できます。
Chainer	2020 年 7 月 1 日	Chainer は、2019 年 12 月時点での メジャーリリースのサポート終了 を発表しました。そのため、2020 年 7 月以降は、DLAMI に Chainer Conda 環境は含ま

機能	日付	詳細
		<p>れなくなります。これらの環境を含む DLAMI の以前のリリースは引き続き利用できます。これらのフレームワークのオープンソースコミュニティによってセキュリティ修正が公開されている場合にのみ、これらの環境の更新を提供します。</p>
Python 3.6	2020 年 6 月 15 日	<p>お客様の要望により、新しい TF/MX/PT リリースに対しては Python 3.7 に移行中です。</p>
Python 2	2020 年 1 月 1 日	<p>Python オープンソースコミュニティは、Python 2 のサポートを正式に終了しました。</p> <p>TensorFlow、PyTorch、MXNet コミュニティもまた、TensorFlow 1.15、TensorFlow 2.1、PyTorch 1.4、MxNet 1.6.0の各リリースが、Python 2 をサポートする最後のリリースとなると発表しています。</p>

DLAMI のドキュメント履歴

次の表は、最新の DLAMI リリースと AWS Deep Learning AMIs デベロッパーガイドの関連する変更の履歴を示しています。

最新の変更

変更	説明	日付
TensorFlow Serving を使用した MNIST モデルのトレーニング	Tensorflow サービスを使用して MNIST モデルをトレーニングする例。	2025 年 2 月 14 日
ARM64 DLAMI	は Arm64 プロセッサベースの GPUs でイメージをサポートする AWS Deep Learning AMIs になりました。	2021 年 11 月 29 日
TensorFlow 2	Deep Learning AMI with Conda には、CUDA 10 の TensorFlow 2 が付属するようになりました。	2019 年 12 月 3 日
AWS 推論	Deep Learning AMI は、Inferentia AWS ハードウェアと AWS Neuron SDK をサポートするようになりました。	2019 年 12 月 3 日
ナイトリービルドからの PyTorch のインストール	Deep Learning AMI with Conda で PyTorch をアンインストールしてから PyTorch のナイトリービルドをインストールする方法を説明するチュートリアルが追加されました。	2018 年 9 月 25 日

[Conda チュートリアル](#)

MOTD の例では、最新リリースを反映するように更新されました。

2018 年 7 月 23 日

以前の変更

次の表は、2018 年 7 月より前の DLAMI リリースと関連する変更の履歴を示しています。

変更	説明	日付
TensorFlow (Horovod を使用)	TensorFlow および Horovod を使用した ImageNet のトレーニングのチュートリアルを追加します。	2018 年 6 月 6 日
アップグレードガイド	アップグレードガイドを追加しました。	2018 年 5 月 15 日
新しいリージョンと新しい 10 分間チュートリアル	追加された新しいリージョン: 米国西部 (北カリフォルニア)、南米、カナダ (中部)、欧州 (ロンドン)、および欧州 (パリ)。また、「Deep Learning AMI の使用開始」という 10 分間チュートリアルの最初のリリースが追加されました。	2018 年 4 月 26 日
Chainer チュートリアル	マルチ GPU、単一の GPU および CPU モデルで Chainer を使用するためのチュートリアルが追加されました。CUDA A 統合は、複数のフレームワークで CUDA 8 から CUDA 9 にアップグレードされました。	2018 年 2 月 28 日

変更	説明	日付
Linux AMI v3.0 に加えて、MXNet Model Server、TensorFlow Serving、および TensorBoard の導入	MXNet Model Server v0.1.5、TensorFlow Serving v1.4.0、および TensorBoard v0.4.0 を使用した新しいモデルおよび可視化処理機能を Conda AMI で使用するチュートリアルを追加しました。AMI およびフレームワーク CUDA 機能については、Conda および CUDA の概要で説明されています。最新のリリースノートは https://aws.amazon.com/releasenotes/ に移動しました。	2018 年 1 月 25 日
Linux AMI v2.0	Base、Source、および Conda AMI が NCCL 2.1 で更新されました。Source および Conda AMI が MXNet v1.0、PyTorch 0.3.0、および Keras 2.0.9 で更新されました。	2017 年 12 月 11 日
2 つの Windows AMI オプションを追加	Windows 2012 R2 および 2016 AMI をリリース: AMI セレクションガイドおよびリリースノートに追加。	2017 年 11 月 30 日
初回のドキュメントリリース	変更されたトピック/セクションへのリンクを含む変更の詳細な説明。	2017 年 11 月 15 日

翻訳は機械翻訳により提供されています。提供された翻訳内容と英語版の間で齟齬、不一致または矛盾がある場合、英語版が優先します。