



开发人员指南

AWS Deep Learning AMIs



AWS Deep Learning AMIs: 开发人员指南

Copyright © 2024 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon 的商标和商业外观不得用于任何非 Amazon 的商品或服务，也不得以任何可能引起客户混淆、贬低或诋毁 Amazon 的方式使用。所有非 Amazon 拥有的其他商标均为各自所有者的财产，这些所有者可能附属于 Amazon、与 Amazon 有关联或由 Amazon 赞助，也可能不是如此。

Table of Contents

DLAMI 是什么？	1
关于本指南	1
先决条件	1
使用案例示例	1
特征	2
预装的框架	2
预装GPU的软件	3
模型处理和可视化	3
DLAMIs 的发布说明	4
基础 DLAMIs	4
单框架 DLAMIs	5
多框架 DLAMIs	6
入门	7
选择一个 DLAMI	7
CUDA 安装和框架绑定	8
基础	9
Conda	9
架构	10
OS	10
选择实例	11
定价	12
区域可用性	12
GPU	12
CPU	13
Inferentia	14
Trainium	14
设置	16
查找DLAMI身份证	16
启动 实例	18
连接到 实例	19
设置 Jupyter	20
保护服务器	20
启动服务器	21
连接客户端	22

登录	23
清理	25
使用 DLAMI	27
Conda DLAMI	27
AMI使用 Conda 进行深度学习简介	27
登录到你的 DLAMI	28
启动 TensorFlow 环境	28
切换到 PyTorch Python 3 环境	29
删除环境	30
基础 DLAMI	30
使用深度学习基础 AMI	30
配置CUDA版本	30
Jupyter 笔记本	31
导航已安装的教程	32
通过 Jupyter 切换环境	32
教程	32
激活框架	33
Elastic Fabric Adapter	36
GPU监控和优化	49
AWS 推论	58
ARM64 DLAMI	80
推理	83
模型处理	83
升级 DLAMI	89
DLAMI升级	89
软件更新	90
发布通知	90
安全性	92
数据保护	92
身份和访问管理	93
使用身份进行身份验证	94
使用策略管理访问	96
IAM与亚马逊合作 EMR	98
合规性验证	98
恢复能力	98
基础结构安全性	99

监控	99
使用情况跟踪	99
框架支持策略	101
DLAMI框架支持 FAQs	101
哪些框架版本会获得安全补丁？	102
AWS 发布新框架版本时会发布哪些镜像？	102
哪些图像获得了新的 SageMaker AI/AWS 功能？	102
“支持的框架”表中是如何定义当前版本的？	102
如果我运行的版本不在“支持的框架”表中，该怎么办？	102
是否DLAMIs支持以前版本的 TensorFlow？	102
如何找到支持的框架版本的最新补丁映像？	102
多长时间发布一次新映像？	103
运行工作负载时，能在我的实例上以替代方式安装补丁吗？	103
如果有新的补丁或更新的框架版本可用，会发生什么呢？	103
是否可在不更改框架版本的情况下更新依赖项？	103
对我的框架版本的主动支持何时结束？	103
对于框架版本不再主动维护的映像，会为其安装补丁吗？	105
如何使用旧框架版本？	105
如何保持框架及其版本 up-to-date的支持变更？	105
是否需要商业许可证才能使用 Anaconda 存储库？	105
重要更改	106
DLAMINVIDIA司机更换 FAQs	106
更改了哪些内容？	106
为什么需要进行此更改？	107
这一变化影响DLAMIs了哪些？	107
这对您意味着什么？	108
较新的版本会失去功能DLAMIs吗？	108
这一变化是否影响了 Deep Learning Containers？	108
相关信息	109
已弃用功能	110
文档历史记录	112
.....	CXV

什么是 AWS Deep Learning AMIs ?

AWS Deep Learning AMIs (DLAMI) 提供可用于云端深度学习的自定义机器映像。DLAMIs它们大多适用于各种亚马逊弹性计算云 (AmazonEC2) 实例类型，从CPU仅限小型的实例到最新的高性能多实例。GPU它们预DLAMIs先配置了 [NVIDIAACUDADNN](#)和 [NVIDIAcu](#) 以及最受欢迎的深度学习框架的最新版本。

关于本指南

中的内容可以帮助您启动和使用DLAMIs。该指南涵盖了几个可同时用于训练和推理的常见深度学习使用案例。它还涵盖了如何根据自己的AMI目的选择正确的实例以及您可能更喜欢的实例类型。

此外，还DLAMIs包括其支持的框架提供的几个教程。本指南可以向您展示如何激活每个框架，并找到相应的入门教程。它还有关于分布式训练、调试、使用 AWS Inferentia 和 AWS Trainium 以及其他关键概念的教程。有关如何设置 Jupyter Notebook 服务器以在浏览器中运行教程的说明，请参阅[在实例上设置 Jupyter 笔记本服务器 DLAMI](#)。

先决条件

要成功运行DLAMIs，我们建议您熟悉命令行工具和基本的 Python。

DLAMI 使用案例示例

以下是 AWS Deep Learning AMIs (DLAMI) 的一些常见用例的示例。

学习深度学习 — DLAMI 是学习或教授机器学习和深度学习框架的绝佳选择。这DLAMIs消除了对每个框架的安装进行故障排除并让它们在同一台计算机上运行所带来的麻烦。DLAMIs其中包括 Jupyter 笔记本，便于运行框架为机器学习和深度学习新手提供的教程。

应用程序开发 — 如果你是一名应用程序开发者，有兴趣使用深度学习让你的应用程序利用人工智能的最新进展，那么DLAMI这是你的理想试验台。每个框架都附带了关于如何开始使用深度学习的教程，其中许多教程都有模型动物园，可以让您轻松试用深度学习，您不需要自己创建神经网络或进行模型训练。一些示例向您展示如何在几分钟内构建映像检测应用程序，或如何为您自己的聊天自动程序构建语音识别应用程序。

机器学习和数据分析：如果您是数据科学家，或者您有兴趣通过深度学习处理数据，则您会发现许多框架都支持 R 和 Spark。您将找到关于进行简单回归的教程，以及为个性化和预测系统构建可扩展的数据处理系统的教程。

研究 — 如果你是一名想要尝试新框架、测试新模型或训练新模型的研究人员，那么DLAMI扩展 AWS 能力可以减轻繁琐的安装和管理多个训练节点所带来的痛苦。

Note

虽然您最初的选择可能是将您的实例类型升级到具有更多实例的大型实例GPU（最多 8 个），但您也可以通过创建DLAMI实例集群进行横向扩展。有关集群构建的更多信息，请参阅[相关的信息 DLAMI](#)。

DLAMI的特征

AWS Deep Learning AMIs (DLAMI) 的功能包括预装的深度学习框架、GPU软件、模型服务器和模型可视化工具。

预装的框架

当前，DLAMI与操作系统 (OS) 和软件版本相关的其他变体有两种主要形式：

- [AMI使用 Conda 进行深度学习](#)：在单独的 Python 环境中使用 conda 程序包单独安装的框架。
- [深度学习基础 AMI](#)— 未安装任何框架；仅[NVIDIA CUDA](#)安装其他依赖项。

Conda AMI 的深度学习使用conda环境来隔离每个框架，因此您可以随意在它们之间切换，而不必担心它们的依赖关系会发生冲突。AMI使用 Conda 进行深度学习支持以下框架：

- PyTorch
- TensorFlow 2

Note

DLAMI不再支持以下深度学习框架：Apache MXNet、微软认知工具包 (CNTK)、Caffe、Caffe、Caffe2、Theano、Chainer 和 Keras。

预装GPU的软件

即使你使用的是CPU仅限实例，DLAMIs也会有 [NVIDIA CUDA](#) and [NVIDIA cuDNN](#)。无论实例类型是什么，安装的软件都是相同的。请记住，GPU特定工具仅适用于至少有一个GPU工具的实例。有关实例类型的更多信息，请参阅[选择DLAMI实例类型](#)。

有关 CUDA 的更多信息，请参阅 [CUDA 安装和框架绑定](#)。

模型处理和可视化

AMI使用 Conda 进行深度学习预装了模型服务器 TensorFlow，用于模型可视化和模型 TensorBoard 可视化。有关更多信息，请参阅 [TensorFlow 服务](#)。

DLAMIs 发行说明

在这里，您可以找到当前支持的所有选项 AWS Deep Learning AMIs (DLAMI) 的详细发行说明。

有关我们不再支持的DLAMI框架的发行说明，请参阅 [Framework Support Policy 页面的不支持的DLAMI框架](#)发行说明存档部分。

Note

AWS Deep Learning AMIs 他们有每晚发布安全补丁的节奏。我们不会在官方发布说明中包含这些增量安全补丁。

基础 DLAMIs

GPU

- X86
 - [AWS 深度学习基础AMI \(亚马逊 Linux 2023 \)](#)
 - [AWS 深度学习基础 AMI \(Ubuntu 22.04\)](#)
 - [AWS 深度学习基础 AMI \(Ubuntu 20.04\)](#)
 - [AWS 深度学习基础AMI \(亚马逊 Linux 2 \)](#)
- ARM64
 - [AWS 深度学习基础 ARM64 AMI \(Ubuntu 22.04\)](#)
 - [AWS 深度学习基础 ARM64AMI \(亚马逊 Linux 2 \)](#)

AWS Neuron

- X86
 - [AWS 深度学习基础AMI神经元 \(亚马逊 Linux 2 \)](#)
 - [AWS 深度学习基础AMI神经元 \(Ubuntu 20.04\)](#)

Qualcomm

- X86

- [AWS 深度学习基础高通AMI \(亚马逊 Linux 2 \)](#)

单框架 DLAMIs

PyTorch-特定的 AMIs

GPU

- X86
 - [AWS 深度学习 AMI GPU PyTorch 2.5 \(Ubuntu 22.04\)](#)
 - [AWS 深度学习 AMI GPU PyTorch 2.5 \(亚马逊 Linux 2023 \)](#)
 - [AWS 深度学习 AMI GPU PyTorch 2.4 \(Ubuntu 22.04\)](#)
 - [AWS 深度学习 AMI GPU PyTorch 2.3 \(Ubuntu 20.04\)](#)
 - [AWS 深度学习 AMI GPU PyTorch 2.3 \(亚马逊 Linux 2 \)](#)
 - [AWS 深度学习 AMI GPU PyTorch 2.2 \(Ubuntu 20.04\)](#)
 - [AWS 深度学习 AMI GPU PyTorch 2.2 \(亚马逊 Linux 2\)](#)
- ARM64
 - [AWS 深度学习 ARM64 AMI GPU PyTorch 2.5 \(Ubuntu 22.04\)](#)
 - [AWS 深度学习 ARM64 AMI GPU PyTorch 2.4 \(Ubuntu 22.04\)](#)
 - [AWS 深度学习 ARM64 AMI GPU PyTorch 2.3 \(Ubuntu 22.04\)](#)
 - [AWS 深度学习 ARM64 AMI GPU PyTorch 2.2 \(Ubuntu 20.04\)](#)

AWS Neuron

- X86
 - [AWS 深度学习 AMI Neuron PyTorch 1.13 \(亚马逊 Linux 2\)](#)
 - [AWS 深度学习AMI神经元 PyTorch 1.13 \(Ubuntu 20.04\)](#)

TensorFlow-特定的 AMIs

GPU

- X86
 - [AWS 深度学习 AMI GPU TensorFlow 2.18 \(亚马逊 Linux 2023 \)](#)

- [AWS 深度学习 AMI GPU TensorFlow 2.18 \(Ubuntu 22.04\)](#)
- [AWS 深度学习 AMI GPU TensorFlow 2.17 \(Ubuntu 22.04\)](#)
- [AWS 深度学习 AMI GPU TensorFlow 2.16 \(亚马逊 Linux 2\)](#)
- [AWS 深度学习 AMI GPU TensorFlow 2.16 \(Ubuntu 20.04\)](#)

AWS Neuron

- X86
 - [AWS 深度学习 AMI Neuron TensorFlow 2.10 \(亚马逊 Linux 2\)](#)
 - [AWS 深度学习 AMI 神经元 TensorFlow 2.10 \(Ubuntu 20.04\)](#)

多框架 DLAMIs

Tip

如果您只使用一个机器学习框架，那么我们建议使用[单一框架DLAMI](#)。

GPU

- X86
 - [AWS 深度学习 AMI \(亚马逊 Linux 2\)](#)

AWS Neuron

- X86
 - [AWS 深度学习 AMI 神经元 \(亚马逊 Linux 2023\)](#)
 - [AWS 深度学习 AMI 神经元 \(Ubuntu 22.04\)](#)

开始使用 DLAMI

本指南包含有关选择适合您的实例、选择适合您的用例和预算的实例类型以及描述[相关的信息 DLAMI](#)可能感兴趣的自定义设置的提示。DLAMI

如果您不熟悉使用 AWS 或使用 Amazon EC2，请从[AMI使用 Conda 进行深度学习](#)。如果您熟悉 Amazon EC2 和其他 AWS 服务，例如亚马逊 EMREFS、亚马逊或 Amazon S3，并且有兴趣将这些服务集成[相关的信息 DLAMI](#)到需要分布式训练或推理的项目中，那么请查看是否适合您的用例。

我们建议您查看[选择一个 DLAMI](#)，以了解最适合您的应用程序的实例类型。

后续步骤

[选择一个 DLAMI](#)

选择一个 DLAMI

如[GPU DLAMI 发行说明](#)中所述，我们提供了一系列 DLAMI 选项。为了帮助您 DLAMI 为自己的用例选择正确的镜像，我们按开发镜像所针对的硬件类型或功能对镜像进行分组。我们的顶层分组是：

- DLAMI 类型：基本、单框架、多框架 (Conda) DLAMI
- 计算架构：[基于 x86、基于 arm64 的 Graviton AWS](#)
- 处理器类型：GPU <https://docs.aws.amazon.com/dlami/latest/devguide/gpu/>、[Inferentia CPU、Trainium](#)
- SDK：[CUDA](#)，[AWS 神经元](#)
- 操作系统：亚马逊 Linux、Ubuntu

本指南中的其他主题有助于进一步给您提供信息和进一步让您了解详情。

主题

- [CUDA 安装和框架绑定](#)
- [深度学习基础 AMI](#)
- [AMI 使用 Conda 进行深度学习](#)
- [DLAMI 架构选项](#)
- [DLAMI 操作系统选项](#)

后续步骤

[AMI使用 Conda 进行深度学习](#)

CUDA 安装和框架绑定

虽然深度学习都非常先进，但每个框架都提供“稳定”版本。这些稳定版本可能不适用于最新的 CUDA 或 cuDNN 实施和功能。您的用例和所需功能可以帮助您选择框架。如果您不确定，就请使用最新的带 Conda 的深度学习 AMI。它为所有使用 CUDA 的框架提供了官方 pip 二进制文件，同时使用每个框架均支持的最新版本。如果您需要最新版本，并要自定义深度学习环境，就请使用深度学习基础 AMI。

如需进一步指导，请在 [稳定版本与候选版本](#) 上查看我们的指南。

选择带 CUDA 的 DLAMI

[深度学习基础 AMI](#) 具有所有可用的 CUDA 版本系列

[AMI使用 Conda 进行深度学习](#) 具有所有可用的 CUDA 版本系列

Note

我们不再在 AWS Deep Learning AMIs 中包含 MXNet、CNTK、Caffe、Caffe2、Theano、Chainer 或 Keras Conda 环境。

有关具体框架版本号，请参阅 [DLAMIs 发行说明](#)。

选择这个 DLAMI 类型，或通过后续步骤选项了解有关其他 DLAMI 的更多信息。

选择一个 CUDA 版本并在附录中查看具有该版本的完整 DLAMI 列表，或者使用后续步骤选项来了解不同 DLAMI 的更多信息。

后续步骤

[深度学习基础 AMI](#)

相关主题

- 有关在 CUDA 版本之间切换的说明，请参阅 [使用深度学习基础 AMI 教程](#)。

深度学习基础 AMI

深度学习基础 AMI 就像一张用于深度学习的空白画布。它包含您需要的一切，直到安装特定的框架，并具有您选择的 CUDA 版本。

为什么要选择基础 DLAMI

该 AMI 团队对于那些想要分享深度学习项目和建立最新版本的项目贡献者非常有用。适用于那些希望推出自己环境且有信心安装和使用最新 NVIDIA 软件的人员，从而他们可以专注于选择要安装的框架和版本。

选择这个 DLAMI 类型，或通过后续步骤选项了解有关其他 DLAMI 的更多信息。

后续步骤

[使用 Conda 的 DLAMI](#)

相关主题

- [使用深度学习基础 AMI](#)

AMI使用 Conda 进行深度学习

Conda DLAMI 使用conda虚拟环境，它们存在于多框架或单一框架中。DLAMIs这些环境配置为单独安装不同的框架，并简化框架之间的切换。这对了解和体验 DLAMI 必须提供的所有框架很有好处。大多数用户发现全新 Conda 深度学习非常适合他们。

它们经常使用框架中的最新版本进行更新，并具有最新的GPU驱动程序和软件。在大多数文档 [AWS Deep Learning AMIs](#) 中，它们通常被称为。它们DLAMIs支持 Ubuntu 20.04、Ubuntu 22.04、亚马逊 Linux 2、亚马逊 Linux 2023 操作系统。操作系统支持取决于上游操作系统的支持。

稳定版本与候选版本

Conda AMIs 使用每个框架中最新正式版本的优化二进制文件。不会有候选版本和实验性功能。优化取决于框架对英特尔等加速技术的支持 MKLDNN，这些技术可以加快对 C5 和 C4 实例类型的训练和推理。CPU编译二进制文件还支持高级英特尔指令集，包括但不限于 AVX-2 AVX、SSE4 .1 和 SSE4 .2。它们可以加速英特尔CPU架构上的矢量和浮点运算。此外，对于GPU实例类型，CUDA和 cu DNN 将使用最新官方版本支持的任何版本进行更新。

AMI使用 Conda 进行深度学习会在框架首次激活时自动为您的 Amazon EC2 实例安装最优化的框架版本。有关更多信息，请参阅[在 Conda 中使用深度学习](#)。

如果要使用自定义或优化的版本选项从源安装，[深度学习基础 AMI](#) 可能是您更好的选择。

Python 2 弃用

Python 开源社区已于 2020 年 1 月 1 日正式结束对 Python 2 的支持。TensorFlow 和 PyTorch 社区已经宣布，TensorFlow 2.1 和 PyTorch 1.4 版本是最后一个支持 Python 2 的版本。包含 Python 2 Conda 环境的 DLAMI (v26、v25 等) 的先前版本继续可用。但是，只有开源社区发布了针对先前发布的版本的 Python 2 Conda 环境的安全修复后，我们才会为这些 DLAMI 版本提供 Python 2 Conda 环境的更新。DLAMI TensorFlow 和 PyTorch 框架的最新版本不包含 Python 2 Conda 环境。

CUDA Support

具体 CUDA 版本号可以在 [GPU DLAMI 发行说明](#) 中找到。

后续步骤

[DLAMI 架构选项](#)

相关主题

- 有关在 Conda 中使用深度学习 AMI 的教程，请参阅 [在 Conda 中使用深度学习教程](#)。

DLAMI 架构选项

AWS Deep Learning AMIs 附带基于 x86 或基于 Arm64 的 [AWS Graviton2](#) 架构。

有关 ARM64 GPU DLAMI 入门的信息，请参阅 [ARM64 DLAMI](#)。有关可用实例类型的更多详细信息，请参阅 [选择 DLAMI 实例类型](#)。

后续步骤

[DLAMI 操作系统选项](#)

DLAMI 操作系统选项

DLAMIs 在以下操作系统中提供。

- Amazon Linux 2
- Amazon Linux 2023
- Ubuntu 20.04

- Ubuntu 22.04

旧版本的操作系统在已弃用DLAMIs版本上可用。有关DLAMI弃用的更多信息，请参阅[“弃用” DLAMI](#)。在选择之前DLAMI，请评估您需要的实例类型并确定您的 AWS 区域。

后续步骤

[选择DLAMI实例类型](#)

选择DLAMI实例类型

更笼统地说，在为选择实例类型时，请考虑以下几点DLAMI。

- 如果您不熟悉深度学习，那么具有单一的实例GPU可能适合您的需求。
- 如果您注重预算，则可以使用CPU仅限于预算的实例。
- 如果您希望优化深度学习模型推断的高性能和成本效益，则可以使用带有 AWS Inferentia 芯片的实例。
- 如果您正在寻找具有基于 ARM64 CPU 架构的高性能GPU实例，则可以使用 g5G 实例类型。
- 如果您有兴趣运行预训练模型进行推理和预测，则可以将 Amazon [Elastic Inference 附加到您的亚马逊实例](#)。EC2 Amazon Elastic Inference 让你只需几分之一即可访问加速器。GPU
- 对于大容量推理服务，具有大量内存的单个CPU实例或此类实例的集群可能是更好的解决方案。
- 如果您使用的是具有大量数据或较大批处理大小的大型模型，那么您需要具有更多内存的大型实例。您也可以将模型分发到一个集群GPUs。您可能会发现，如果减小批处理大小，则使用内存较少的实例将是更好的选择。这可能会影响您的准确性和训练速度。
- 如果您有兴趣使用NVIDIA集体通信库 (NCCL) 运行机器学习应用程序，需要大规模进行高水平的节点间通信，则可能需要使用 [Elastic Fabric Adapter \(EFA\)](#)。

有关实例的更多详细信息，请参阅[实例类型](#)。

以下主题提供有关实例类型注意事项的信息。

Important

深度学习AMIs包括NVIDIA公司开发、拥有或提供的驱动程序、软件或工具包。您同意仅在包含NVIDIA硬件的 Amazon EC2 实例上使用这些NVIDIA驱动程序、软件或工具包。

主题

- [DLAMI 的定价](#)
- [DLAMI 区域可用性](#)
- [推荐GPU实例](#)
- [推荐CPU实例](#)
- [推荐的 Inferentia 实例](#)
- [推荐的 Trainium 实例](#)

DLAMI 的定价

中包含的深度学习框架DLAMI是免费的，每个框架都有自己的开源许可证。尽管中包含的软件DLAMI是免费的，但您仍然需要为底层 Amazon EC2 实例硬件付费。

有些 Amazon EC2 实例类型被标记为免费。可以在其中一个免费实例DLAMI上运行。这意味着，当您仅使用该DLAMI实例的容量时，完全可以免费使用。如果您需要一个具有更多CPU内核、更多磁盘空间RAM、更多或一个或多个的更强大的实例GPUs，那么您需要一个不属于免费套餐实例类的实例。

有关实例选择和定价的更多信息，请参阅 [Amazon EC2 定价](#)。

DLAMI 区域可用性

每个区域支持不同的实例类型范围，并且通常在不同的区域中，一种实例类型的成本略有不同。DLAMIs并非在每个地区都可用，但可以复制DLAMIs到您选择的区域。有关更多信息，[AMI请参阅复制](#)。请注意区域选择列表，并确保您选择一个靠近您或您客户的区域。如果您计划使用多个集群DLAMI并可能创建集群，请务必对集群中的所有节点使用相同的区域。

有关区域的更多信息，请访问 [Amazon EC2 服务终端节点](#)。

后续步骤

[推荐GPU实例](#)

推荐GPU实例

对于大多数深度学习目的，我们建议使用GPU实例。在实例上训练新模型比在GPU实例上训练新模型更快。CPU当您有多GPU实例或在多个实例上使用分布式训练时，您可以进行亚线性扩展。GPUs

以下实例类型支持DLAMI。有关GPU实例类型选项及其用途的信息，请参阅[实例类型](#)并选择加速计算。

Note

应将模型大小作为选择实例的一个因素。如果您的模型超过了实例的可用容量RAM，请选择其他具有足够内存的实例类型，以供应用程序使用。

- [亚马逊 EC2 P5e 实例](#)最多有 8 个 NVIDIA Tesla H200。GPUs
- [亚马逊 EC2 P5 实例](#)最多有 8 个 NVIDIA Tesla H GPUs 100。
- [亚马逊 EC2 P4 实例](#)最多有 8 个 NVIDIA 特斯拉 A GPUs 100。
- [亚马逊 EC2 P3 实例](#)最多有 8 个 NVIDIA Tesla V GPUs 100。
- [亚马逊 EC2 G3 实例](#)最多有 4 个 NVIDIA Tesla M GPUs 60。
- [亚马逊 EC2 G4 实例](#)最多有 4 个 NVIDIA T GPUs 4。
- [亚马逊 EC2 G5 实例](#)最多有 8 个 NVIDIA A GPUs 10G。
- [亚马逊 EC2 G6 实例](#)最多有 8 个 NVIDIA L GPUs 4。
- [亚马逊 EC2 G6e 实例](#)最多有 8 个 NVIDIA L40S 张量核心。GPUs
- [亚马逊 EC2 G5G 实例具有基于 ARM64 的 Gravit AWS on2 处理器。](#)

DLAMI实例提供了用于监控和优化GPU流程的工具。有关监控GPU流程的更多信息，请参阅[GPU监控和优化](#)。

有关使用 G5g 实例的特定教程，请参阅 [ARM64 DLAMI](#)。

后续步骤

[推荐CPU实例](#)

推荐CPU实例

无论您是预算有限，正在学习深度学习，还是只想运行预测服务，该CPU类别中都有许多经济实惠的选择。一些框架利用英特尔的优势 MKLDNN，它可以加快 C5（并非在所有地区都可用）CPU实例类型的训练和推理。有关CPU实例类型的信息，请参阅并选择计算优化。

Note

应将模型大小作为选择实例的一个因素。如果您的模型超过了实例的可用容量RAM，请选择其他具有足够内存的实例类型，以供应用程序使用。

- [Amazon EC2 C5 实例](#)最多有 72 个英特尔 vCPUs。C5 实例擅长科学建模、批处理、分布式分析、高性能计算 (HPC) 以及机器和深度学习推理。

后续步骤

[推荐的 Inferentia 实例](#)

推荐的 Inferentia 实例

AWS Inferentia 实例旨在为深度学习模型推理工作负载提供高性能和成本效益。具体而言，Inf2 实例类型使用 AWS Inferentia 芯片和 Ne [AWS uron SDK](#)，后者与流行的机器学习框架（例如和）集成。TensorFlow PyTorch

客户使用 Inf2 实例之后，能够以最低的云端成本来运行大规模的机器学习推理应用程序，例如搜索、推荐引擎、计算机视觉、语音识别、自然语言处理、个性化和欺诈检测。

Note

应将模型大小作为选择实例的一个因素。如果您的模型超过了实例的可用容量 RAM，请选择其他具有足够内存的实例类型，以供应用程序使用。

- [Amazon EC2 Inf2 实例](#)最多有 16 个 AWS 推理芯片和 100 Gbps 的网络吞吐量。

有关开始使用 AWS Inferentia 的更多信息 DLAMIs，请参阅。[AWS 推理芯片带有 DLAMI](#)

后续步骤

[推荐的 Trainium 实例](#)

推荐的 Trainium 实例

AWS Trainium 实例旨在为深度学习模型推理工作负载提供高性能和成本效益。具体而言，Trn1 实例类型使用 T AWS rainium 芯片和 Ne [AWS uron SDK](#)，后者与流行的机器学习框架（例如和）集成。TensorFlow PyTorch

客户使用 Trn1 实例之后，能够以最低的云端成本来运行大规模的机器学习推理应用程序，例如搜索、推荐引擎、计算机视觉、语音识别、自然语言处理、个性化和欺诈检测。

Note

应将模型大小作为选择实例的一个因素。如果您的模型超过了实例的可用容量RAM，请选择其他具有足够内存的实例类型，以供应用程序使用。

- [Amazon EC2 Trn1 实例](#)最多有 16 个 T AWS trainium 芯片和 100 Gbps 的网络吞吐量。

设置实DLAMI例

选择DLAMI并[选择要使用的亚马逊弹性计算云 \(AmazonEC2\) 实例类型](#)后，就可以设置新DLAMI实例了。

如果您尚未选择DLAMI和EC2实例类型，请参阅[开始使用 DLAMI](#)。

主题

- [正在查找 a 的 ID DLAMI](#)
- [启动 DLAMI 实例](#)
- [连接到实DLAMI例](#)
- [在实例上设置 Jupyter 笔记本服务器 DLAMI](#)
- [清理实DLAMI例](#)

正在查找 a 的 ID DLAMI

每个都DLAMI有一个唯一的标识符 (ID)。当您使用 Amazon EC2 控制台启动DLAMI实例时，您可以选择使用该 DLAMI ID 来搜索要使用的实例。DLAMI使用 AWS Command Line Interface (AWS CLI) 启动DLAMI实例时，需要此 ID。

您可以使用亚马逊EC2或 Parameter Store DLAMI 的 AWS CLI 命令找到自己选择的商品的 ID，该功能为 AWS Systems Manager。有关安装和配置的说明 AWS CLI，请参阅《AWS Command Line Interface 用户指南》AWS CLI中的[“入门”](#)。

Using Parameter Store

要查找DLAMI身份证，请使用 `ssm get-parameter`

在以下[`ssm get-parameter`](#)命令中，对于该`--name`选项，参数名称的格式为`/aws/service/deeplearning/ami/$architecture/$ami_type/latest/ami-id`。在这种名称格式中，`architecture`可以是`x86_64`或`arm64`。`ami_type`通过取DLAMI名字并删除关键字“deep”、“learning”和“ami”来指定。AMI名字可以在中找到[DLAMIs 发行说明](#)。

⚠ Important

要使用此命令，您使用的 AWS Identity and Access Management (IAM) 主体必须具有 `ssm:GetParameter` 权限。有关 IAM 委托人的更多信息，请参阅《IAM 用户指南》中 IAM 角色的 [其他资源](#) 部分。

- ```
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
```

该输出值应该类似于以下内容：

```
ami-09ee1a996ac214ce7
```

**i Tip**

对于当前支持的某些 DLAMI 框架，您可以在中找到更具体的示例 `ssm get-parameter` 命令 [DLAMIs 发行说明](#)。选择您所选发行说明的链接 DLAMI，然后在发行说明中找到其 ID 查询。

**Using Amazon EC2 CLI**

要查找 DLAMI 身份证，请使用 `ec2 describe-images`

在以下 `ec2 describe-images` 命令中，为筛选 `Name=name` 器的值输入 DLAMI 名称。可以为给定的框架指定发布版本，也可以通过将版本号替换为问号 ( ? ) 来获取最新版本。

- ```
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
```

该输出值应该类似于以下内容：

```
ami-09ee1a996ac214ce7
```

i Tip

有关特定于您选择的ec2 describe-images命令DLAMI的示例，请参阅[DLAMIs 发行说明](#)。选择您所选发行说明的链接DLAMI，然后在发行说明中找到其 ID 查询。

后续步骤

[启动 DLAMI 实例](#)

启动 DLAMI 实例

找到要用于启动DLAMI实例的 [DLAMI ID](#) 后，就可以启动该实例了。要启动它，您可以使用 Amazon EC2 控制台或 AWS Command Line Interface (AWS CLI)。

i Note

在本演练中，我们可能会特别提及深度学习基础 OSS Nvidia驱动程序 GPUAMI (Ubuntu 22.04)。即使您选择其他 DLAMI，您也应能按照本指南进行操作。

EC2 console

i Note

要加速高性能计算 (HPC) 和机器学习应用程序，您可以使用 Elastic Fabric Adapter (EFA) 启动您的DLAMI实例。有关具体说明，请参阅[使用启动 AWS Deep Learning AMIs 实例 EFA](#)。

1. 打开控制[EC2台](#)。
2. 在最上面的导航栏 AWS 区域 中记下你当前的状态。如果这不是您所需的区域，请在继续操作之前更改此选项。有关更多信息，请参阅中的[亚马逊EC2服务终端节点Amazon Web Services 一般参考](#)。
3. 选择启动实例。
4. 为您的实例输入一个名称DLAMI，然后选择适合您的名称。

- a. 在“我的”DLAMI 中查找现有内容AMIs或选择“快速入门”。
 - b. 按 DLAMI ID 搜索。浏览这些选项，然后选中您的选择。
5. 选择一个实例类型。您可以在中找到适合您的DLAMI推荐实例系列[DLAMIs 发行说明](#)。有关 DLAMI实例类型的一般建议，请参阅[选择DLAMI实例类型](#)。
 6. 选择启动实例。

AWS CLI

- 要使用 AWS CLI，您必须拥有要使用的 DLAMI ID、AWS 区域 和EC2实例类型以及您的安全令牌信息。然后，您可以使用[ec2 run-instances](#) AWS CLI 命令启动实例。

有关安装和配置说明 AWS CLI，请参阅《AWS Command Line Interface 用户指南》AWS CLI中的“入门”。有关更多信息，包括命令示例，请参阅[启动、列出和关闭 Amazon EC2 实例 AWS CLI](#)。

使用 Amazon EC2 控制台或启动实例后 AWS CLI，等待实例准备就绪。这通常仅需要几分钟时间。您可以在 [Amazon EC2 控制台](#) 中验证实例的状态。有关更多信息，请参阅亚马逊EC2用户指南中的[亚马逊EC2实例状态检查](#)。

后续步骤

[连接到实例DLAMI例](#)

连接到实例DLAMI例

在[启动DLAMI实例并且该实例正在运行](#)之后，您可以使用从客户端（Windows、macOS 或 Linux）连接到该实例。SSH有关说明，请参阅亚马逊EC2用户指南SSH中的[使用连接到您的 Linux 实例](#)。

如果你想在SSH登录后设置 Jupyter Notebook 服务器，请随身携带登录命令的副本。要连接到 Jupyter 网页，可以使用该命令的变体。

后续步骤

[在实例上设置 Jupyter 笔记本服务器 DLAMI](#)

在实例上设置 Jupyter 笔记本服务器 DLAMI

使用 Jupyter 笔记本服务器，您可以从您的实例创建和运行 Jupyter 笔记本。DLAMI 使用 Jupyter 笔记本，您可以在使用 AWS 基础架构和访问内置软件包的同时，进行机器学习 (ML) 实验进行训练和推理。DLAMI 有关 Jupyter Notebook 的更多信息，请参阅 Jupyter 用户文档网站上的 [The Jupyter Notebook](#)。

要设置 Jupyter Notebook 服务器，您必须：

- 在您的 DLAMI 实例上配置 Jupyter 笔记本服务器。
- 配置客户端以连接到 Jupyter Notebook 服务器。我们提供适用于 Windows、macOS 和 Linux 客户端的配置说明。
- 通过登录 Jupyter Notebook 服务器，测试设置。

要完成这些步骤，请遵循以下主题中的说明。设置 Jupyter Notebook 服务器后，您可以运行中附带的示例笔记本教程。DLAMIs 有关更多信息，请参阅 [运行 Jupyter 笔记本电脑教程](#)。

主题

- [在实例上保护 Jupyter 笔记本服务器 DLAMI](#)
- [在实例上启动 Jupyter 笔记本服务器 DLAMI](#)
- [在实例上将客户端连接到 Jupyter Notebook 服务器 DLAMI](#)
- [在实例上登录 Jupyter 笔记本服务器 DLAMI](#)

在实例上保护 Jupyter 笔记本服务器 DLAMI

为了保障 Jupyter Notebook 服务器的安全，我们建议您为服务器设置密码并创建 SSL 证书。要配置密码和 SSL，请先 [连接到您的 DLAMI 实例](#)，然后按照以下说明进行操作。

保护 Jupyter Notebook 服务器

1. Jupyter 提供了一个密码实用工具。运行以下命令，在命令提示符处输入您的首选密码。

```
$ jupyter notebook password
```

输出类似如下：

```
Enter password:
```

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

2. 创建自签名SSL证书。按照提示填写您认为适当的区域。如果要提示留空，则必须输入。您的答案将不会影响证书的功能性。

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

Note

您可能有兴趣创建一个由第三方签名且不会导致浏览器向您发出安全警告的常规SSL证书。此过程涉及内容较多。有关更多信息，请参阅 Jupyter Notebook 用户文档中的 [Securing a notebook server](#)。

后续步骤

[在实例上启动 Jupyter 笔记本服务器 DLAMI](#)

在实例上启动 Jupyter 笔记本服务器 DLAMI

使用[密码保护 Jupyter Notebook 服务器后SSL](#)，就可以启动服务器了。登录您的DLAMI实例并运行以下命令，该命令使用您之前创建的SSL证书。

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

服务器启动后，您现在可以从客户端计算机通过SSH隧道连接到服务器。当服务器运行时，您将看到来自 Jupyter 的一些输出，这些输出证实服务器正在运行。此时，请忽略关于您可以通过本地主机访问服务器的标注URL，因为在您创建隧道之前，这不会起作用。

Note

当您使用 Jupyter Web 界面切换框架时，Jupyter 将处理切换环境。有关更多信息，请参阅 [通过 Jupyter 切换环境](#)。

后续步骤

[在实例上将客户端连接到 Jupyter Notebook 服务器 DLAMI](#)

在实例上将客户端连接到 Jupyter Notebook 服务器 DLAMI

在[您的DLAMI实例上启动 Jupyter Notebook 服务器](#)后，将您的 Windows、macOS 或 Linux 客户端配置为连接到该服务器。当您连接时，可以在工作区中的服务器上创建和访问 Jupyter Notebook，并在服务器上运行深度学习代码。

先决条件

确保你有以下内容，你需要这些东西来设置SSH隧道：

- 您的 Amazon EC2 实例的公共DNS名称。有关更多信息，请参阅[亚马逊EC2用户指南中的亚马逊EC2实例主机名类型](#)。
- 私有密钥文件的密钥对。有关访问密钥对的更多信息，请参阅[亚马逊EC2用户指南中的亚马逊EC2密钥对和亚马逊EC2实例](#)。

从 Windows、macOS 或 Linux 客户端进行连接

要从 Windows、macOS 或 Linux 客户端连接到您的DLAMI实例，请按照客户端操作系统的说明进行操作。

Windows

要从 Windows 客户端连接到您的DLAMI实例，请使用 SSH

1. 使用适用于 Windows 的SSH客户端，例如 Pu TTY。有关说明，请参阅亚马逊EC2用户指南 TTY中的[使用 Pu 连接您的 Linux 实例](#)。有关其他SSH连接选项，请参阅[使用连接您的 Linux 实例SSH](#)。
2. （可选）创建通往正在运行的 Jupyter 服务器的SSH隧道。在 Windows 客户端上安装 Git Bash，然后按照 macOS 和 Linux 客户端的连接说明进行操作。

macOS or Linux

要从 macOS 或 Linux 客户端连接到您的DLAMI实例，请使用 SSH

1. 打开终端。

2. 运行以下命令将本地端口 8888 上的所有请求转发到远程 Amazon EC2 实例上的端口 8888。通过替换访问亚马逊实例的密钥位置和亚马逊 EC2 实例的公共 DNS 名称来更新命令。EC2 注意，对于 Amazon Linux AMI，用户名 `ec2-user` 取而代之 `ubuntu`。

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

此命令在您的客户端和运行 Jupyter Notebook 服务器的远程 Amazon EC2 实例之间打开一条隧道。

后续步骤

[在实例上登录 Jupyter 笔记本服务器 DLAMI](#)

在实例上登录 Jupyter 笔记本服务器 DLAMI

将[客户端连接到 DLAMI 实例上的 Jupyter Notebook 服务器](#)后，您可以登录该服务器。

在浏览器中登录服务器

1. 在浏览器的地址栏中输入以下内容 URL，或点击此链接：<https://localhost:8888>
2. 使用自签名 SSL 证书后，您的浏览器将警告您并提示您避免继续访问该网站。



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

由于这是您自己设置的，所以可以安全地继续。根据您的浏览器情况，有可能出现“高级”、“显示详细信息”等类似按钮。



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 \(unsafe\)](#)

单击此消息，然后单击“前进到本地主机”链接。如果连接成功，则会看到 Jupyter Notebook 服务器网页。此时，系统将要求您提供先前设置的密码。

现在，您可以访问实例上运行的 Jupyter 笔记本服务器。DLAMI您可以创建新的笔记本或运行提供的[教程](#)。

清理实例DLAMI例

当您不再需要您的DLAMI实例时，可以在 Amazon 上将其停止或终止EC2，以免产生意外费用。

如果您停止某个实例，可以保留它，稍后当您想再次使用它时再启动它。配置、文件和其它非易失性信息存储在 Amazon Simple Storage Service (Amazon S3) 上的卷中。当实例停止时，您会因保留卷而产生 S3 费用，但不会因计算资源而产生费用。当您再次启动实例时，它将装载包含数据的存储卷。

如果您终止实例，实例将消失，您将无法再次启动它。当然，对于终止的实例，您不会再为计算资源支付任何费用。但是，您的数据仍存储在 Amazon S3 上，并且您可能会继续支付 S3 费用。为了防止与终止的实例相关的所有进一步费用，还必须删除 Amazon S3 上的存储卷。有关说明，请参阅[亚马逊EC2用户指南中的终止亚马逊EC2实例](#)。

有关亚马逊EC2实例状态（例如stopped和terminated）的更多信息，请参阅[亚马逊EC2用户指南中的亚马逊EC2实例状态更改](#)。

使用 DLAMI

主题

- [在 Conda 中使用深度学习](#)
- [使用深度学习基础 AMI](#)
- [运行 Jupyter 笔记本电脑教程](#)
- [教程](#)

以下各节介绍如何使用AMI带有 Conda 的深度学习来切换环境、运行每个框架的示例代码以及运行 Jupyter，以便您可以尝试不同的笔记本教程。

在 Conda 中使用深度学习

主题

- [AMI使用 Conda 进行深度学习简介](#)
- [登录到你的 DLAMI](#)
- [启动 TensorFlow 环境](#)
- [切换到 PyTorch Python 3 环境](#)
- [删除环境](#)

AMI使用 Conda 进行深度学习简介

Conda 是一个开源程序包管理系统和环境管理系统，在 Windows、macOS 和 Linux 上运行。Conda 快速安装、运行和更新程序包及其依赖项。Conda 可轻松创建、保存、加载和切换本地计算机上的环境。

AMI使用 Conda 进行深度学习的配置可以让您在深度学习环境之间轻松切换。以下说明为您介绍与 conda 相关的一些基本命令。它们还可以帮助您验证框架的基本导入正常运行，并且您可以使用框架运行一些简单操作。然后，您可以继续查看随 DLAMI 提供的更全面的教程，或者每个框架的项目站点上提供的框架示例。

登录到你的 DLAMI

登录服务器后，您将看到一则服务器“每日消息”(MOTD)，其中描述了可用于在不同的深度学习框架之间切换的各种 Conda 命令。以下是一个示例 MOTD。随着新版本的发布，您的具体情况 MOTD DLAMI 可能会有所不同。

```
=====
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
=====
```

启动 TensorFlow 环境

Note

在启动您的第一个 Conda 环境时，请在其加载期间耐心等待。AMI 使用 Conda 进行深度学习会在框架首次激活时自动为您的 EC2 实例安装最优化的框架版本。您不应期望后续的延迟。

1. 激活 Python 3 的 TensorFlow 虚拟环境。

```
$ source activate tensorflow2_p310
```

2. 启动终端 iPython 端。

```
(tensorflow2_p310)$ ipython
```

3. 运行一个快速 TensorFlow 程序。

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

您应该会看到“Hello, Tensorflow!”

后续步骤

[运行 Jupyter 笔记本电脑教程](#)

切换到 PyTorch Python 3 环境

如果您仍在 iPython 控制台中，请使用 `quit()`，然后准备好切换环境。

- 激活 Python 3 的 PyTorch 虚拟环境。

```
$ source activate pytorch_p310
```

测试一些 PyTorch 代码

要测试您的安装，请使用 Python 编写用于创建和打印数组的 PyTorch 代码。

1. 启动终端 iPython 端。

```
(pytorch_p310)$ ipython
```

2. 导入 PyTorch。

```
import torch
```

您可能会看到一条关于第三方软件包的警告消息。您可以忽略它。

3. 创建一个 5x3 矩阵，将元素随机初始化。打印数组。

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

验证结果。

```
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]])
```

删除环境

注意：如果您用尽了 DLAMI 上的空间，则可以选择卸载不用的 Conda 软件包：

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

使用深度学习基础 AMI

使用深度学习基础 AMI

Base AMI 带有 GPU 驱动程序和加速库的基础平台，用于部署您自己的自定义深度学习环境。默认情况下 AMI，使用任意一个 NVIDIA CUDA 版本的环境进行配置。您也可以在不同的版本之间切换 CUDA。有关如何执行此操作，请参阅以下说明。

配置 CUDA 版本

您可以通过运行的 `nvcc` 程序 NVIDIA 来验证 CUDA 版本。

```
nvcc --version
```

您可以使用以下 `bash` 命令选择和验证特定 CUDA 版本：

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

有关更多信息，请参阅 [Base DLAMI 版本说明](#)。

运行 Jupyter 笔记本电脑教程

每个深度学习项目的源代码都附带了教程和示例，大多数情况下，它们可以在任何 DLAMI 上运行。如果您选择了 [AMI使用 Conda 进行深度学习](#)，那么您将获得一些已经建立并准备好尝试的精选教程的额外好处。

Important

要运行安装在 DLAMI 上的 Jupyter 笔记本电脑教程，您需要[在实例上设置 Jupyter 笔记本服务器 DLAMI](#)。

Jupyter 服务器运行后，您可以通过 Web 浏览器运行这些教程。如果你正在AMI使用 Conda 运行深度学习或者已经设置了 Python 环境，则可以从 Jupyter 笔记本界面切换 Python 内核。在尝试运行特定于框架的教程之前，请选择合适的内核。为使用 Conda 进行深度学习AMI的用户提供了更多示例。

Note

许多教程需要一些额外 Python 模块，它们在您的 DLAMI 上可能尚未设置。如果您收到类似 "xyz module not found" 的错误，请登录到 DLAMI，激活环境（如上所述），然后安装必要的模块。

Tip

深度学习教程和示例通常依赖于一个或多个GPUs。如果您的实例类型没有GPU，则可能需要更改示例的某些代码才能使其运行。

导航已安装的教程

登录 Jupyter 服务器并可以看到教程目录（仅限AMI使用 Conda 的深度学习）后，将显示按每个框架名称排列的教程文件夹。如果您没有看到某个框架，则表明在您当前 DLAMI 上该框架的教程不可用。单击框架名称以查看列出的教程，然后单击一个教程，将其启动。

当你第一次使用 Conda 在深AMI度学习上运行笔记本时，它会想知道你想使用哪个环境。它会提示您从列表中进行选择。每个环境都根据以下模式命名：

Environment (conda_framework_python-version)

例如，你可能会看到Environment (conda_mxnet_p36)，这表示环境中 MXNet Python 3。另一个变体是Environment (conda_mxnet_p27)，这表示环境中 MXNet Python 2。

Tip

如果您担心哪个版本CUDA处于活动状态，则可以通过首次登录MOTD时查看该版本DLAMI。

通过 Jupyter 切换环境

如果您决定尝试一个不同框架的教程，一定要验证当前正在运行的内核。此信息可以在 Jupyter 界面的右上方看到，就在注销按钮的下方。您可以在任何打开的笔记本电脑上更改内核，方法是依次单击 Kernel、Change Kernel 菜单项，然后单击正运行的笔记本电脑适合的环境。

此时您需要重新运行任何单元，因为内核中的更改将会擦除之前运行的任何内容的状态。

Tip

在框架之间进行切换可能很有趣而且很有教育意义，但是可能导致耗尽内存。如果您开始遇到错误，请查看运行 Jupyter 服务器的终端窗口。这里有一些有用的消息和错误记录，你可能会看到一个 out-of-memory 错误。要解决这一问题，您可以转到 Jupyter 服务器的主页中，单击 Running 选项卡，然后为每个教程单击 Shutdown，因为这些教程可能仍然在后台运行，导致耗尽了所有内存。

教程

以下是有关如何在 Conda 软件中使用深AMI度学习的教程。

主题

- [激活框架](#)
- [使用 Elastic Fabric Adapter 进行分布式训练](#)
- [GPU监控和优化](#)
- [AWS 推理芯片带有 DLAMI](#)
- [ARM64 DLAMI](#)
- [推理](#)
- [模型处理](#)

激活框架

以下是使用 Conda 在深度学习上安装的深AMI度学习框架。单击某个框架即可了解如何将其激活。

主题

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

PyTorch

正在激活 PyTorch

当框架的稳定的 Conda 程序包发布时，它会在 DLAMI 上进行测试并预安装。如果您希望运行最新的、未经测试的每日构建版本，您可以手动[Install PyTorch 的夜间构建 \(实验版\)](#)。

要激活当前安装的框架，请按照AMI使用 Conda 进行深度学习中的以下说明进行操作。

对于带有CUDA和 MKL-的 Pyth PyTorch on 3DNN，请运行以下命令：

```
$ source activate pytorch_p310
```

启动终端 iPython 端。

```
(pytorch_p310)$ ipython
```

运行一个快速 PyTorch 程序。

```
import torch
```

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

您应该会看到系统输出初始随机数组，然后输出大小，然后添加另一个随机数组。

Install PyTorch 的夜间构建 (实验版)

如何 PyTorch 从夜间版本中安装

您可以在使用 Conda 进行深度学习AMI的任一或两 PyTorch 个 Conda 环境中安装最新 PyTorch 版本。

1. • (Python 3 的选项) -激活 Python 3 PyTorch 环境 :

```
$ source activate pytorch_p310
```

2. 其余步骤假定您使用的是 pytorch_p310 环境。移除当前安装的 PyTorch :

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

3. • (GPU实例可选) -使用 CUDA .0 安装最新的夜 PyTorch 间版本 :

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

- (CPU实例可选) -为不GPU带以下选项的 PyTorch 实例安装最新的夜间版本 :

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

4. 要验证您是否已成功安装最新的夜间版本，请启动IPython终端并检查的版本。 PyTorch

```
(pytorch_p310)$ ipython
```

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

输出应类似于以下内容 : 1.0.0.dev20180922

5. 要验证 PyTorch 夜间构建是否适用于该MNIST示例，您可以从 PyTorch的示例存储库中运行测试脚本：

```
(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
```

更多教程

有关更多教程和示例，请参阅该框架的官方[PyTorch 文档](#)、[文档](#)和[PyTorch](#)网站。

TensorFlow 2

本教程介绍如何在AMI使用 Conda 运行深度学习的实例上激活 TensorFlow 2 (DLAMI在 Conda 上) 并运行 TensorFlow 2 程序。

当框架的稳定的 Conda 程序包发布时，它会在 DLAMI 上进行测试并预安装。

正在激活 TensorFlow 2

和 Conda DLAMI 一起跑步 TensorFlow

1. 要激活 TensorFlow 2，请DLAMI使用 Conda 打开的亚马逊弹性计算云 (AmazonEC2) 实例。
2. 对于 TensorFlow 带有 CUDA 10.1 和 MKL-的 Python 3 上的 2 和 Keras 2DNN，请运行以下命令：

```
$ source activate tensorflow2_p310
```

3. 启动终端 iPython 端：

```
(tensorflow2_p310)$ ipython
```

4. 运行 TensorFlow 2 程序以验证它是否正常运行：

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

Hello, TensorFlow! 应显示在您的屏幕上。

更多教程

有关更多教程和示例，请参阅 [TensorFlow Python TensorFlow](#) 文档API或[TensorFlow](#)访问网站。

使用 Elastic Fabric Adapter 进行分布式训练

[弹性结构适配器](#) (EFA) 是一种网络设备，您可以将其连接到您的DLAMI实例以加速高性能计算 (HPC) 应用程序。EFA借助 AWS 云提供的可扩展性、灵活性和弹性，使您能够实现本地HPC集群的应用程序性能。

以下主题将向您展示如何开始将 EFA 与 DLAMI 结合使用。

Note

DLAMI从这个[基础GPU DLAMI列表](#)中选择你的

主题

- [使用启动 AWS Deep Learning AMIs 实例 EFA](#)
- [在 DLAMI 上使用 EFA](#)

使用启动 AWS Deep Learning AMIs 实例 EFA

最新DLAMI的 Base 已准备就绪，EFA并附带了所需的驱动程序、内核模块、libfabric、openmpi 和[NCCLOFI插件](#)。GPU

您可以在[发行说明DLAMI](#)中找到支持的 Base CUDA 版本。

注意：

- 在mpirun上运行NCCL应用程序时EFA，必须将EFA支持的安装的完整路径指定为：

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

- 要使您的应用程序能够使用 EFA，请将 FI_PROVIDER="efa" 添加到 mpirun 命令，如 [在 DLAMI 上使用 EFA](#) 中所示。

主题

- [准备启用 EFA 的安全组](#)
- [启动实例](#)
- [验证 EFA 附件](#)

准备启用 EFA 的安全组

EFA 需要使用一个安全组，以允许进出安全组本身的所有入站和出站流量。有关更多信息，请参阅 [EFA 文档](#)。

1. 打开 Amazon EC2 控制台，网址为 <https://console.aws.amazon.com/ec2/>。
2. 在导航窗格中，选择安全组，然后选择创建安全组。
3. 在创建安全组窗口中，执行以下操作：
 - 对于安全组名称，请输入一个描述性的安全组名称，例如 EFA-enabled security group。
 - (可选) 对于描述，请输入安全组的简要描述。
 - 对于 VPC，VPC 选择您打算在其中启动 EFA 已启用实例的。
 - 选择创建。
4. 选择您创建的安全组，然后在描述 选项卡上复制组 ID。
5. 在入站和出站选项卡上，执行以下操作：
 - 选择编辑。
 - 对于类型，请选择所有流量。
 - 对于 Source，选择 Custom。
 - 将您复制的安全组 ID 粘贴到该字段中。
 - 选择保存。
6. 启用入站流量，请参考 [授权您 Linux 实例的入站流量](#)。如果您跳过此步骤，您将无法与您的 DLAMI 实例进行通信。

启动实例

EFAon 上 AWS Deep Learning AMIs 目前支持以下实例类型和操作系统：

- p3dn：亚马逊 Linux 2、Ubuntu 20.04
- p4d、p4de：亚马逊 Linux 2、亚马逊 Linux 2023、Ubuntu 20.04、Ubuntu 22.04

- P5、p5e、p5eN : 亚马逊 Linux 2、亚马逊 Linux 2023、Ubuntu 20.04、Ubuntu 22.04

以下部分说明如何启动EFA已启用的DLAMI实例。有关启动EFA已启用实例的更多信息，请参阅[已EFA启用的实例启动到集群置放群组](#)。

1. 打开 Amazon EC2 控制台，网址为 <https://console.aws.amazon.com/ec2/>。
2. 选择启动实例。
3. 在“选择一个AMI页面”上，选择“[DLAMI发行说明](#)”页面上DLAMI找到的支持版本
4. 在选择实例类型页面上，选择以下支持的实例类型之一，然后选择下一步：配置实例详细信息。有关支持的实例列表，请参阅此链接：[开始使用EFA和 MPI](#)
5. 在配置实例详细信息页面中，执行以下操作：
 - 对于 Number of instances (实例的数量)，请输入要启动的启用了 EFA 的实例数量。
 - 对于网络和子网，选择要在其中启动实例的VPC和子网。
 - [可选] 对于置放群组，请选择将实例添加到置放群组。为获得最佳性能，请在置放群组中启动实例。
 - [可选] 对于置放群组名称，请选择添加到新的置放群组，输入置放群组的描述性名称，然后对于置放群组策略选择集群。
 - 请务必在此页面上启用“Elastic Fabric Adapter”。如果禁用此选项，请将子网更改为支持所选实例类型的子网。
 - 在网络接口部分中，为设备 eth0 选择新网络接口。您可以选择指定一个主IPv4地址和一个或多个辅助IPv4地址。如果您要在具有关联IPv6CIDR区块的子网中启动实例，则可以选择指定一个主IPv6地址和一个或多个辅助IPv6地址。
 - 选择下一步：添加存储。
6. 在 Add Storage (添加存储) 页面上，除了 AMI 指定的卷 (如根设备卷) 以外，还要指定要附加到实例的卷，然后选择 Next: Add Tags (下一步：添加标签)。
7. 在添加标签页面上，为实例指定标签 (例如，便于用户识别的名称)，然后选择下一步：配置安全组。
8. 在配置安全组页面上，对于分配安全组选择一个现有的安全组，然后选择先前创建的安全组。
9. 选择审核并启动。
10. 在核查实例启动页面上，检查这些设置，然后选择启动以选择一个密钥对并启动您的实例。

验证 EFA 附件

通过控制台

启动实例后，请在 AWS 控制台中查看实例详细信息。为此，请在 EC2 控制台中选择实例，然后查看页面下方窗格中的描述选项卡。找到参数“Network Interfaces: eth0”，然后单击 eth0，这将弹出一个弹出窗口。确保已启用“Elastic Fabric Adapter”。

如果未启用 EFA，您可以通过以下任一方式解决此问题：

- 使用相同的步骤终止 EC2 实例并启动一个新实例。确保已附加 EFA。
- 将 EFA 附加到现有实例。
 1. 在 EC2 控制台中，转到网络接口。
 2. 单击“创建虚拟网络接口”。
 3. 选择您的实例所在的相同子网。
 4. 确保启用“Elastic Fabric Adapter”并单击“创建”。
 5. 返回 EC2 实例选项卡并选择您的实例。
 6. 转到“操作：实例状态”并在附加 EFA 之前停止实例。
 7. 从“操作”中，选择“网络连接：连接网络接口”。
 8. 选择您刚刚创建的界面，然后单击“附加”。
 9. 重新启动您的实例。

通过实例

以下测试脚本已存在于 DLAMI 中。运行它以确保内核模块正确加载。

```
$ fi_info -p efa
```

您的输出应类似于以下内容。

```
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

```

验证安全组配置

以下测试脚本已存在于 DLAMI 中。运行它以确保您创建的安全组配置正确。

```

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

```

您的输出应类似于以下内容。

```

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

```

如果它停止响应或未完成，请确保您的安全组具有正确的入站/出站规则。

在 DLAMI 上使用 EFA

以下部分描述如何使用 EFA 在 AWS Deep Learning AMIs 上运行多节点应用程序。

使用 EFA 运行多节点应用程序

要跨节点集群运行应用程序，需要以下配置

主题

- [启用无密码 SSH](#)

- [创建主机文件](#)
- [NCCL测试](#)

启用无密码 SSH

选择集群中的一个节点作为领导节点。其余节点称为成员节点。

1. 在领导节点上，生成RSA密钥对。

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

2. 更改领导节点上私有密钥的权限。

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

3. 复制公有密钥 `~/.ssh/id_rsa.pub` 并将其附加到集群中成员节点的 `~/.ssh/authorized_keys` 之后。
4. 现在，您应该能够使用私有 ip 从领导节点直接登录到成员节点。

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

5. 通过将以下内容添加到领导节点上的 `~/.ssh/config` 文件中，禁用 `strictHostKey`检查并在领导节点上启用代理转发：

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

6. 在 Amazon Linux 2 实例上，在领导节点上运行以下命令，为配置文件提供正确的权限：

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

创建主机文件

在领导节点上，创建主机文件以标识集群中的节点。主机文件必须针对集群中的每个节点都有一个条目。创建文件 `~/hosts` 并使用私有 IP 添加每个节点，如下所示：

```
localhost slots=8
```

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

NCCL测试

Note

这些测试已使用EFA版本 1.30.0 和OFINCCCL插件 1.7.4 运行。

下面列出了Nvidia提供的NCCL测试子集，用于在多个计算节点上测试功能和性能

支持的实例：P3dn、P4、P5

功能测试

NCCL消息传输多节点测试

`nccl_message_transfer` 是一个简单的测试，旨在确保插件按预期运行。NCCL OFI该测试验证了NCCL的连接建立和数据传输APIs的功能。使用运行NCCL应用程序时，请确保使用示例中所示的 `mpirun` 的完整路径。EFA根据集群GPUs中的实例数量和实例数量更改参数`np`。有关更多信息，请参阅[AWS OFINCCCL文档](#)。

以下 `nccl_message_transfer` 测试适用于通用 `xx.x` 版本。CUDA您可以通过替换脚本中的CUDA版本来运行适用于您的 Amazon EC2 实例中任何可用CUDA版本的命令。

```
$/opt/amazon/openmpi/bin/mpirun -n 2 -N 1 --hostfile hosts \  
-x LD_LIBRARY_PATH=/usr/local/cuda-xx.x/efa/lib:/usr/local/cuda-xx.x/lib:/usr/  
local/cuda-xx.x/lib64:/usr/local/cuda-xx.x:$LD_LIBRARY_PATH \  
--mca btl tcp,self --mca btl_tcp_if_exclude lo,docker0 --bind-to none \  
opt/aws-ofi-nccl/tests/nccl_message_transfer
```

您的输出应与以下内容类似。您可以查看输出以查看EFA它是否被用作OFI提供程序。

```
INFO: Function: nccl_net_ofi_init Line: 1069: NET/OFI Selected Provider is efa (found 4  
 nics)  
INFO: Function: nccl_net_ofi_init Line: 1160: NET/OFI Using transport protocol SENDRECV  
INFO: Function: configure_ep_inorder Line: 261: NET/OFI Setting  
 FI_OPT_EFA_SENDRECV_IN_ORDER_ALIGNED_128_BYTES not supported.  
INFO: Function: configure_nccl_proto Line: 227: NET/OFI Setting NCCL_PROTO to "simple"  
INFO: Function: main Line: 86: NET/OFI Process rank 1 started. NCCLNet device used on  
 ip-172-31-13-179 is AWS Libfabric.
```

```

INFO: Function: main Line: 91: NET/OFI Received 4 network devices
INFO: Function: main Line: 111: NET/OFI Network supports communication using CUDA
  buffers. Dev: 3
INFO: Function: main Line: 118: NET/OFI Server: Listening on dev 3
INFO: Function: main Line: 131: NET/OFI Send connection request to rank 1
INFO: Function: main Line: 173: NET/OFI Send connection request to rank 0
INFO: Function: main Line: 137: NET/OFI Server: Start accepting requests
INFO: Function: main Line: 141: NET/OFI Successfully accepted connection from rank 1
INFO: Function: main Line: 145: NET/OFI Send 8 requests to rank 1
INFO: Function: main Line: 179: NET/OFI Server: Start accepting requests
INFO: Function: main Line: 183: NET/OFI Successfully accepted connection from rank 0
INFO: Function: main Line: 187: NET/OFI Rank 1 posting 8 receive buffers
INFO: Function: main Line: 161: NET/OFI Successfully sent 8 requests to rank 1
INFO: Function: main Line: 251: NET/OFI Got completions for 8 requests for rank 0
INFO: Function: main Line: 251: NET/OFI Got completions for 8 requests for rank 1

```

性能测试

在 p4d.24xLarge 上进行多节点NCCL性能测试

要检查NCCL性能EFA，请运行官方测试存储库中提供的标准NCCL[NCCL性能测试](#)。DLAMI附带这个已经为 CUDA XX.X 构建的测试。同样，你可以用它运行自己的脚本。EFA

构建您自己的脚本时，请参阅以下指南：

- 使用运行NCCL应用程序时，使用示例中所示的 mpirun 的完整路径。EFA
- 根据集群GPU中的实例数量更改参数 np 和 N。
- 添加 NCCL_DEBUG = INFO 标志，并确保日志将EFA使用情况显示为“所选提供商是EFA”。
- 设置要解析的训练日志位置以进行验证

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

watch nvidia-smi在任何成员节点上使用该命令来监控GPU使用情况。以下watch nvidia-smi命令适用于通用 CUDA xx.x 版本，具体取决于您的实例的操作系统。您可以通过替换脚本中的CUDA版本来运行适用于您的 Amazon EC2 实例中任何可用CUDA版本的命令。

- Amazon Linux 2 :

```
$ /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 :

```
$ /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}
```

您的输出应与以下内容类似 :

```
# 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 9591 on ip-172-31-4-37 device 0 [0x10] NVIDIA A100-SXM4-40GB
# Rank 1 Group 0 Pid 9592 on ip-172-31-4-37 device 1 [0x10] NVIDIA A100-SXM4-40GB
# Rank 2 Group 0 Pid 9593 on ip-172-31-4-37 device 2 [0x20] NVIDIA A100-SXM4-40GB
# Rank 3 Group 0 Pid 9594 on ip-172-31-4-37 device 3 [0x20] NVIDIA A100-SXM4-40GB
# Rank 4 Group 0 Pid 9595 on ip-172-31-4-37 device 4 [0x90] NVIDIA A100-SXM4-40GB
# Rank 5 Group 0 Pid 9596 on ip-172-31-4-37 device 5 [0x90] NVIDIA A100-SXM4-40GB
# Rank 6 Group 0 Pid 9597 on ip-172-31-4-37 device 6 [0xa0] NVIDIA A100-SXM4-40GB
# Rank 7 Group 0 Pid 9598 on ip-172-31-4-37 device 7 [0xa0] NVIDIA A100-SXM4-40GB
# Rank 8 Group 0 Pid 10216 on ip-172-31-13-179 device 0 [0x10] NVIDIA A100-
SXM4-40GB
# Rank 9 Group 0 Pid 10217 on ip-172-31-13-179 device 1 [0x10] NVIDIA A100-
SXM4-40GB
# Rank 10 Group 0 Pid 10218 on ip-172-31-13-179 device 2 [0x20] NVIDIA A100-
SXM4-40GB
# Rank 11 Group 0 Pid 10219 on ip-172-31-13-179 device 3 [0x20] NVIDIA A100-
SXM4-40GB
```

```

# Rank 12 Group 0 Pid 10220 on ip-172-31-13-179 device 4 [0x90] NVIDIA A100-
SXM4-40GB
# Rank 13 Group 0 Pid 10221 on ip-172-31-13-179 device 5 [0x90] NVIDIA A100-
SXM4-40GB
# Rank 14 Group 0 Pid 10222 on ip-172-31-13-179 device 6 [0xa0] NVIDIA A100-
SXM4-40GB
# Rank 15 Group 0 Pid 10223 on ip-172-31-13-179 device 7 [0xa0] NVIDIA A100-
SXM4-40GB
ip-172-31-4-37:9591:9591 [0] NCCL INFO Bootstrap : Using ens32:172.31.4.37
ip-172-31-4-37:9591:9591 [0] NCCL INFO NET/Plugin: Failed to find ncclCollNetPlugin_v6
symbol.
ip-172-31-4-37:9591:9591 [0] NCCL INFO NET/Plugin: Failed to find ncclCollNetPlugin
symbol (v4 or v5).
ip-172-31-4-37:9591:9591 [0] NCCL INFO cudaDriverVersion 12020
NCCL version 2.18.5+cuda12.2
...
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Initializing aws-ofi-nccl 1.7.4-aws
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Using CUDA runtime version 11070
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Configuring AWS-specific options
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Using CUDA runtime version 11070
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Configuring AWS-specific options
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Setting provider_filter to efa
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Setting FI_EFA_FORK_SAFE environment
variable to 1
ip-172-31-4-37:9024:9062 [6] NCCL INFO NET/OFI Disabling NVLS support due to NCCL
version 21602
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Setting provider_filter to efa
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Setting FI_EFA_FORK_SAFE environment
variable to 1
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Disabling NVLS support due to NCCL
version 21602
ip-172-31-4-37:9020:9063 [2] NCCL INFO NET/OFI Running on p4d.24xlarge platform,
Setting NCCL_TOPO_FILE environment variable to /opt/aws-ofi-nccl/share/aws-ofi-nccl/
xml/p4d-24x1-topo.xml
...
-----some output truncated-----
#                                     out-of-place
#           in-place
#   size      count      type  redop  root    time  algbw  busbw #wrong
#   time      algbw      busbw #wrong
#   (us)      (GB/s)    (GB/s)
#           (B)      (elements)
#           (us)      (GB/s)    (GB/s)
#           0         0         float    sum     -1     11.02   0.00   0.00   0
11.04   0.00   0.00   0

```

	0	0	0	float	sum	-1	11.01	0.00	0.00	0
11.00	0.00	0.00	0							
	0	0	0	float	sum	-1	11.02	0.00	0.00	0
11.02	0.00	0.00	0							
	0	0	0	float	sum	-1	11.01	0.00	0.00	0
11.00	0.00	0.00	0							
	0	0	0	float	sum	-1	11.02	0.00	0.00	0
11.02	0.00	0.00	0							
	256	4	0	float	sum	-1	632.7	0.00	0.00	0
628.2	0.00	0.00	0							
	512	8	0	float	sum	-1	627.4	0.00	0.00	0
629.6	0.00	0.00	0							
	1024	16	0	float	sum	-1	632.2	0.00	0.00	0
631.7	0.00	0.00	0							
	2048	32	0	float	sum	-1	631.0	0.00	0.00	0
634.2	0.00	0.00	0							
	4096	64	0	float	sum	-1	623.3	0.01	0.01	0
633.6	0.01	0.01	0							
	8192	128	0	float	sum	-1	635.1	0.01	0.01	0
633.5	0.01	0.01	0							
	16384	256	0	float	sum	-1	634.8	0.03	0.02	0
637.0	0.03	0.02	0							
	32768	512	0	float	sum	-1	647.9	0.05	0.05	0
636.8	0.05	0.05	0							
	65536	1024	0	float	sum	-1	658.9	0.10	0.09	0
667.0	0.10	0.09	0							
	131072	2048	0	float	sum	-1	671.9	0.20	0.18	0
662.9	0.20	0.19	0							
	262144	4096	0	float	sum	-1	692.1	0.38	0.36	0
685.1	0.38	0.36	0							
	524288	8192	0	float	sum	-1	715.3	0.73	0.69	0
696.6	0.75	0.71	0							
	1048576	16384	0	float	sum	-1	734.6	1.43	1.34	0
729.2	1.44	1.35	0							
	2097152	32768	0	float	sum	-1	785.9	2.67	2.50	0
794.5	2.64	2.47	0							
	4194304	65536	0	float	sum	-1	837.2	5.01	4.70	0
837.6	5.01	4.69	0							
	8388608	131072	0	float	sum	-1	929.2	9.03	8.46	0
931.4	9.01	8.44	0							
	16777216	262144	0	float	sum	-1	1773.6	9.46	8.87	0
1772.8	9.46	8.87	0							
	33554432	524288	0	float	sum	-1	2110.2	15.90	14.91	0
2116.1	15.86	14.87	0							

```

    67108864      1048576      float      sum      -1      2650.9      25.32      23.73      0
2658.1  25.25  23.67      0
    134217728      2097152      float      sum      -1      3943.1      34.04      31.91      0
3945.9  34.01  31.89      0
    268435456      4194304      float      sum      -1      7216.5      37.20      34.87      0
7178.6  37.39  35.06      0
    536870912      8388608      float      sum      -1      13680      39.24      36.79      0
13676   39.26  36.80      0
[ 1073741824      16777216      float      sum      -1      25645      41.87      39.25      0
25497   42.11  39.48      0 ] <- Used For Benchmark
...
# Out of bounds values : 0 OK
# Avg bus bandwidth    : 7.46044

```

验证测试

要验证EFA测试是否返回了有效结果，请使用以下测试进行确认：

- 使用实例元数据获取EC2实例类型：

```

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)

```

- 运行[性能测试](#)
- 设置以下参数

```

CUDA_VERSION
CUDA_RUNTIME_VERSION
NCCL_VERSION

```

- 验证结果，如下所示：

```

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

    # Information on how the version come from logs
    #
    # ip-172-31-27-205:6427:6427 [0] NCCL INFO cudaDriverVersion 12020
    # NCCL version 2.16.2+cuda11.8
    # ip-172-31-27-205:6427:6820 [0] NCCL INFO NET/OFI Initializing aws-ofi-nccl
    1.7.1-aws

```

```

# ip-172-31-27-205:6427:6820 [0] NCCL INFO NET/OFI Using CUDA runtime version
11060

# cudaDriverVersion 12020 --> This is max supported cuda version by nvidia
driver
# NCCL version 2.16.2+cuda11.8 --> This is NCCL version compiled with cuda
version
# Using CUDA runtime version 11060 --> This is selected cuda version

# Validation of logs
grep "NET/OFI Using CUDA runtime version ${CUDA_RUNTIME_VERSION}" ${TRAINING_LOG}
|| { echo "Runtime cuda text not found"; exit 1; }
grep "NET/OFI Initializing aws-ofi-nccl" ${TRAINING_LOG} || { echo "aws-ofi-nccl
is not working, please check if it is installed correctly"; exit 1; }
grep "NET/OFI Configuring AWS-specific options" ${TRAINING_LOG} || { echo "AWS-
specific options text not found"; exit 1; }
grep "Using network AWS Libfabric" ${TRAINING_LOG} || { echo "AWS Libfabric 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/AWS Libfabric/0/GDRDMA" ${TRAINING_LOG} || { echo "Text not found:
NET/AWS 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; }
    grep "aws-ofi-nccl/xml/p4d-24x1-topo.xml" ${TRAINING_LOG} || { echo "Topology
file not found"; exit 1; }
elif [[ ${INSTANCE_TYPE} == "p4de.24xlarge" ]]; then
    grep "NET/AWS 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; }
    grep "aws-ofi-nccl/xml/p4de-24x1-topo.xml" ${TRAINING_LOG} || { echo
"Topology file not found"; exit 1; }
elif [[ ${INSTANCE_TYPE} == "p5.48xlarge" ]]; then
    grep "NET/AWS 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; }

```

```

    grep "aws-ofi-nccl/xml/p5.48x1-topo.xml" ${TRAINING_LOG} || { echo "Topology
file 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

```

- 要访问基准数据，可以解析多节点 all_reduce 测试的表输出的最后一行：

```

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监控和优化

以下部分将指导您完成GPU优化和监控选项。本部分的组织方式如同一个典型工作流程一样，其中包含监控监督预处理和训练。

- [监控](#)
 - [GPUs使用监视器 CloudWatch](#)
- [优化](#)
 - [预处理](#)
 - [训练](#)

监控

您预DLAMI装了多个GPU监控工具。本指南还将介绍可用于下载和安装的工具。

- [GPUs使用监视器 CloudWatch](#)-预装的实用程序，可向 Amazon CloudWatch 报告GPU使用情况统计信息。
- [nvidia-smi CLI](#)-用于监控整体GPU计算和内存利用率的实用程序。它已预先安装在你的 AWS Deep Learning AMIs (DLAMI) 上。
- [NVMLC 库](#)-基于 C 的库API，用于直接访问GPU监控和管理功能。nvidia-smi 在幕后使用它，并已 CLI预先安装在你的。DLAMI它还具有 Python 和 Perl 绑定以方便采用这些请求进行开发。你上预安装的 gpumon.py 实用程序DLAMI使用的是来自的 pynvml 软件包。[nvidia-ml-py](#)
- [NVIDIADCGM](#)-集群管理工具。请访问开发人员页面，了解如何安装和配置此工具。

Tip

请查看开发NVIDIA者博客，了解有关使用您安装的CUDA工具的最新信息DLAMI：

- [使用 Nsight IDE 和 nvprof 监控 TensorCore 利用率。](#)

GPUs使用监视器 CloudWatch

当你将你的DLAMI与一起使用时，GPU你可能会发现你正在寻找在训练或推理期间跟踪其使用情况的方法。这对于优化您的数据管道以及调整深度学习网络非常有用。

有两种方法可以配置GPU指标 CloudWatch：

- [使用 AWS CloudWatch 代理配置指标 \(推荐\)](#)
- [使用预安装 gpumon.py 脚本来配置指标](#)

使用 AWS CloudWatch 代理配置指标 (推荐)

将您DLAMI与[统一 CloudWatch 代理](#)集成，以配置GPU指标并监控 Amazon EC2 加速实例中GPU协处理的利用率。

有四种方法可以用你的配置[GPU指标](#)DLAMI：

- [配置最低GPU指标](#)
- [配置部分GPU指标](#)
- [配置所有可用GPU指标](#)
- [配置自定义GPU指标](#)

有关更新和安全补丁的信息，请参阅 [为代理安装安全补 AWS CloudWatch 丁](#)

先决条件

首先，您必须配置允许您的EC2实例向推送指标的 Amazon 实例IAM权限 CloudWatch。有关详细步骤，请参阅[创建用于 CloudWatch 代理的IAM角色和用户](#)。

配置最低GPU指标

使用该dlami-cloudwatch-agent@minimalsystemd服务配置最低GPU指标。此服务配置以下指标：

- utilization_gpu
- utilization_memory

您可以在以下位置找到用于最低预配置GPU指标的systemd服务：

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

使用以下命令启用并启动 systemd 服务：

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

配置部分GPU指标

使用该dlami-cloudwatch-agent@partialsystemd服务配置部分GPU指标。此服务配置以下指标：

- utilization_gpu
- utilization_memory
- memory_total
- memory_used
- memory_free

您可以在以下位置找到部分预配置GPU指标的systemd服务：

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

使用以下命令启用并启动 `systemd` 服务：

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

配置所有可用GPU指标

使用该 `dlami-cloudwatch-agent@all` `systemd` 服务配置所有可用GPU指标。此服务配置以下指标：

- `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`

您可以在以下位置找到所有可用的预配置GPU指标的 `systemd` 服务：

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

使用以下命令启用并启动 `systemd` 服务：

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

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

配置自定义GPU指标

如果预配置的指标不符合您的要求，则可以创建自定义 CloudWatch 代理配置文件。

创建自定义配置文件

要创建自定义配置文件，请参阅[手动创建或编辑 CloudWatch 代理配置文件](#)中的详细步骤。

在此示例中，假设架构定义位于 `/opt/aws/amazon-cloudwatch-agent/etc/amazon-cloudwatch-agent.json`。

使用您的自定义文件来配置指标

运行以下命令根据您的自定义文件配置 CloudWatch 代理：

```
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
```

为代理安装安全补 AWS CloudWatch 丁

新发布DLAMIs的服务器配置了最新的可用 AWS CloudWatch 代理安全补丁。根据您的选择的操作系统，请参阅以下章节，DLAMI使用最新的安全补丁更新您的最新安全补丁。

Amazon Linux 2

用于获取yum适用于亚马逊 Linux 2 的最新 AWS CloudWatch 代理安全补丁DLAMI。

```
sudo yum update
```

Ubuntu

要使用 Ubuntu 获取最新 AWS CloudWatch 安全补丁，必须使用 Amazon S3 下载链接重新安装 AWS CloudWatch 代理。DLAMI

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

有关使用 Amazon S3 下载链接安装 AWS CloudWatch 代理的更多信息，请参阅在[服务器上安装和运行 CloudWatch 代理](#)。

使用预安装 `gpumon.py` 脚本来配置指标

一个名为 `gpumon.py` 的实用工具已预安装在您的 DLAMI 上。它集成 CloudWatch 并支持对每次 GPU 使用情况的监控：GPU 内存、GPU 温度和 GPU 功率。该脚本会定期将监控的数据发送到 CloudWatch。您可以通过更改脚本中的一些设置来配置要发送到 CloudWatch 的数据的粒度级别。但是，在启动脚本之前，您需要先进行设置 CloudWatch 才能接收指标。

如何使用设置和运行 GPU 监控 CloudWatch

1. 创建 IAM 用户，或修改现有用户，使其具有向其发布指标的策略 CloudWatch。如果创建新用户，请记住凭证，因为您将在下一步中需要这些凭证。

要搜索的 IAM 策略是“`cloudwatch:PutMetricData`”。要添加的策略如下所示：

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

Tip

有关创建 IAM 用户和为添加策略的更多信息 CloudWatch，请参阅 [CloudWatch 文档](#)。

2. 在您的 DLAMI，运行 [AWS 配置](#) 并指定 IAM 用户凭证。

```
$ aws configure
```

3. 您可能需要先对 `gpumon` 实用工具进行一些修改，然后再运行该工具。您可以在以下代码块 README 中定义的位置找到 `gpumon` 实用程序。有关 `gpumon.py` 脚本的更多信息，请参阅 [脚本的 Amazon S3 位置](#)。

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

```
~/tools/GPUCloudWatchMonitor/README
```

选项：

- 如果您的实例位于 us-east-1 中，请在 gpumon.py NOT 中更改区域。
 - 更改其他参数，例如 CloudWatchnamespace 或报告周期 store_reso。
4. 目前，该脚本仅支持 Python 3。激活您的首选框架的 Python 3 环境或激活 DLAMI 的一般 Python 3 环境。

```
$ source activate python3
```

5. 在后台中运行 gpumon 实用工具。

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

6. 打开您的浏览器前往 <https://console.aws.amazon.com/cloudwatch/>，然后选择指标。它将有一个命名空间 “DeepLearningTrain”。

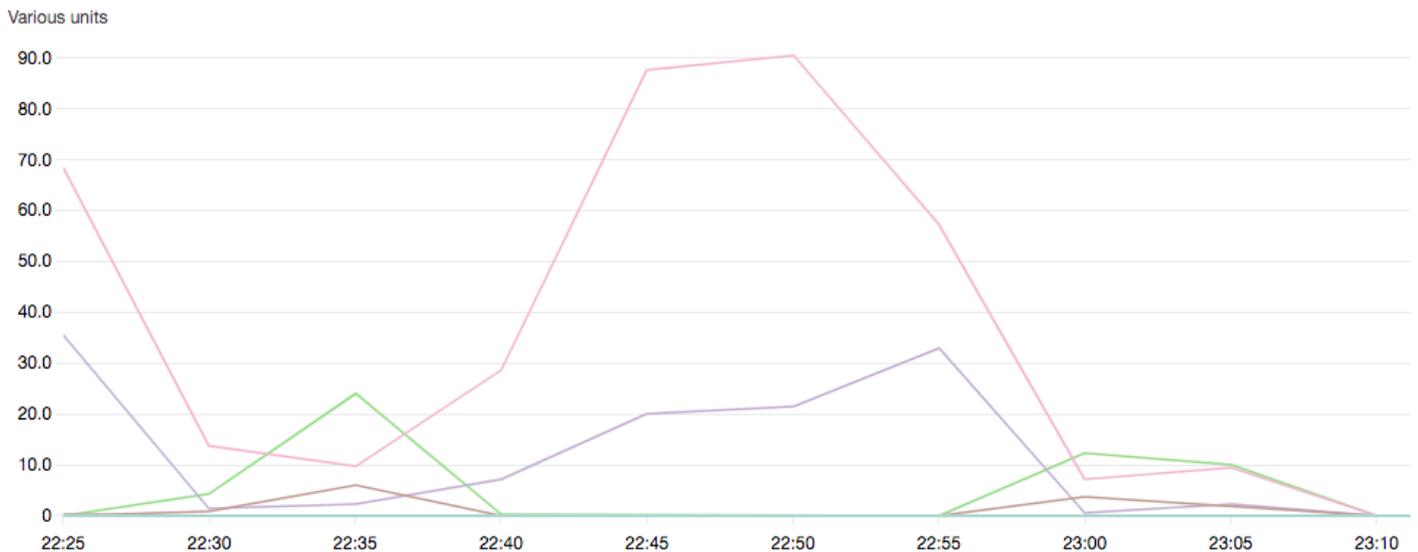
 Tip

您可以修改 gpumon.py 来更改该命名空间。您也可以通过调整 store_reso 来修改报告间隔。

以下是一个示例 CloudWatch 图表，报告了 gpumon.py 在监视 p2.8xlarge 实例上的训练作业的情况。

GPU usage, Memory usage 

1h 3h 12h 1d 3d 1w custom



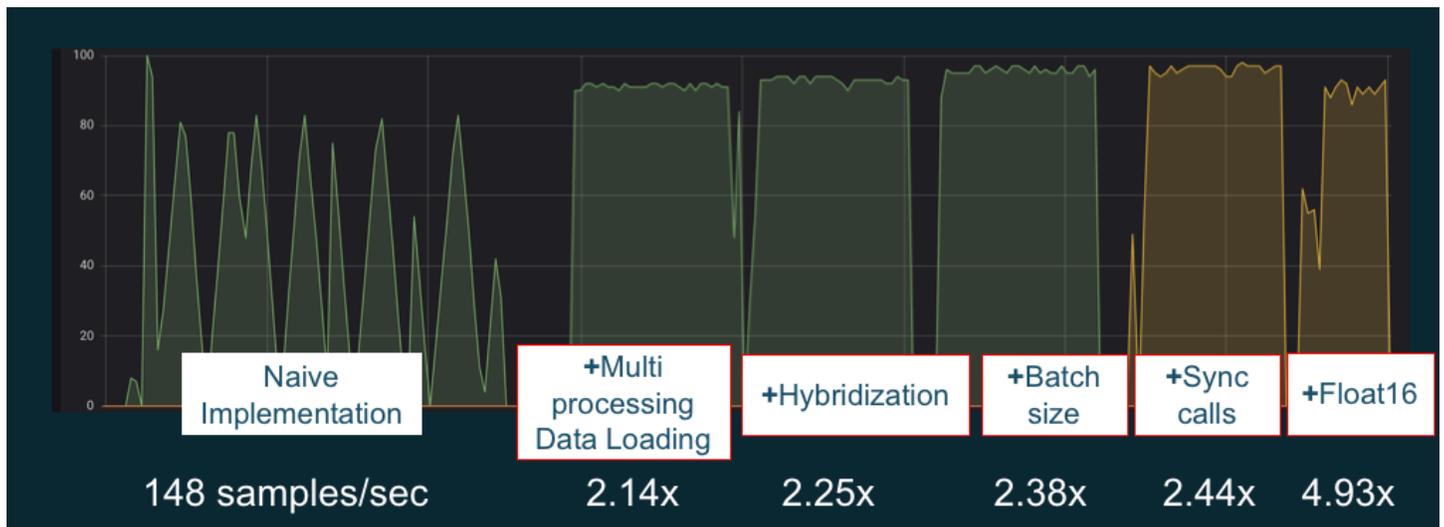
您可能对有关GPU监控和优化的其他主题感兴趣：

- [监控](#)
 - [GPUs使用监视器 CloudWatch](#)
- [优化](#)
 - [预处理](#)
 - [训练](#)

优化

要充分利用您的数据GPUs，您可以优化数据管道并调整深度学习网络。如下图所述，神经网络的天真或基本实现可能会使用不GPU一致的东西，也无法充分发挥其潜力。当你优化预处理和数据加载时，你可以减少瓶颈从你的CPU到你的GPU您可以通过使用混合化（该框架支持时）、调整批大小并同步调用来调整神经网络本身。您也可以在大多数框架中使用多精度（float16 或 int8）训练，从而可以显著提高吞吐量。

以下图表显示应用不同优化时的累积性能提升。您的结果将取决于要处理的数据和要优化的网络。



GPU性能优化示例。图表来源：[使用 MXNet Gluon 的性能技巧](#)

以下指南介绍了适用于您的选项DLAMI，可帮助您提高GPU性能。

主题

- [预处理](#)
- [训练](#)

预处理

通过转换或增强进行数据预处理通常是一个CPU受限的过程，这可能是整个管道中的瓶颈。框架具有用于图像处理的内置运算符，但是DALI（数据增强库）的性能比框架的内置选项有所提高。

- NVIDIA数据增强库 (DALI): DALI 将数据增强卸载到 GPU。它未预装在上DLAMI，但您可以通过在您的或其他亚马逊弹性计算云实例上安装DLAMI或加载支持的框架容器来访问它。详情请参阅NVIDIA网站上的[DALI项目页面](#)。有关示例用例和下载代码示例，请参阅 [SageMaker 预处理训练](#)性能示例。
- nvJPEG：面向 GPU C 程序员的加速JPEG解码器库。它支持解码单个图像或批次，以及深度学习中常见的后续转换操作。nv JPEG 内置了DALI，或者你可以从[NVIDIA网站的 nvjpeg](#) 页面下载并单独使用。

您可能对有关GPU监控和优化的其他主题感兴趣：

- [监控](#)
 - [GPUs使用监视器 CloudWatch](#)
- [优化](#)

- [预处理](#)
- [训练](#)

训练

利用混合精度训练，您可以使用相同的内存量部署更大的网络，或者减少内存使用量（与您的单精度或双精度网络相比），并且您将看到计算性能增加。您还将受益于更小且更快的数据传输，这在多节点分布式训练中是一个重要因素。要利用混合精度训练，您需要调整数据转换和损失比例。以下是介绍如何针对支持混合精度的框架执行此操作的指南。

- [NVIDIA深度学习 SDK](#)-NVIDIA 网站上描述了MXNet、PyTorch和的混合精度实现的文档。
TensorFlow

Tip

请务必针对您选择的框架检查网站，并且搜索“混合精度”或“fp16”，了解最新的优化方法。下面是可能对您有帮助的一些混合精度指南：

- [混合精度训练 TensorFlow（视频）](#)-在NVIDIA博客网站上。
- [使用 float16 进行混合精度训练，以及网站上MXNet的一FAQ篇文章](#)。MXNet
- [NVIDIAApex：一款用于轻松进行混合精度训练的工具 PyTorch](#)——网站上的一篇博客文章。NVIDIA

您可能对有关GPU监控和优化的其他主题感兴趣：

- [监控](#)
 - [GPUs使用监视器 CloudWatch](#)
- [优化](#)
 - [预处理](#)
 - [训练](#)

AWS 推理芯片带有 DLAMI

AWS Inferentia 是一款由其设计的自定义机器学习芯片 AWS ，可用于高性能的推理预测。要使用该芯片，请设置一个 Amazon Elastic Compute Cloud 实例，然后使用 Ne AWS uron 软件开发套件

(SDK) 调用 Inferentia 芯片。为了向客户提供最佳的 Inferentia 体验，Neuron 已内置在 () 中。AWS Deep Learning AMIs DLAMI

以下主题介绍如何开始使用 Inferentia 和 DLAMI

内容

- [启动带有 Neuron DLAMI on 的实例](#)
- [与 Neuron DLAMI on 一起使用](#)

启动带有 Neuron DLAMI on 的实例

最新版本已准备 DLAMI 好与 AWS Inferentia 一起使用，并附带 Neuron API on 软件包。要启动实例 DLAMI，请参阅[启动和配置 DLAMI](#)。完成后 DLAMI，请按照此处的步骤确保您的 AWS 推理芯片和 AWS 神经元资源处于活动状态。

内容

- [验证您的实例](#)
- [识别 AWS 推理设备](#)
- [查看资源使用量](#)
- [使用 Neuron Monitor \(Neuron 监视器 \)](#)
- [升级 Neuron 软件](#)

验证您的实例

在使用您的实例之前，验证该实例是否已针对 Neuron 进行正确的设置和配置。

识别 AWS 推理设备

要确定实例上的 Inferentia 设备数量，请使用以下命令：

```
neuron-ls
```

如果您的实例已附加了 Inferentia 设备，则输出将如下所示：

```
+-----+-----+-----+-----+-----+
| 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

提供的输出取自 Inf1.6xlarge 实例，包括以下各列：

- NEURONDEVICE：分配给逻辑 ID NeuronDevice。此 ID 用于将多个运行时配置为使用不同的 NeuronDevices 运行时。
- NEURONCORES：中 NeuronCores 存在的数量 NeuronDevice。
- NEURONMEMORY：中的 DRAM 内存量 NeuronDevice。
- CONNECTEDDEVICES：其他 NeuronDevices 连接到 NeuronDevice。
- PCIBDF：的 PCI 总线设备功能 (BDF) ID NeuronDevice。

查看资源使用量

使用 `neuron-top` 命令查看有关 NeuronCore 和 v CPU 利用率、内存使用情况、加载的模型和 Neuron 应用程序的有用信息。不 `neuron-top` 带参数启动将显示所有使用的机器学习应用程序的数据 NeuronCores。

```
neuron-top
```

当应用程序使用四时 NeuronCores，输出应类似于下图：

```

neuron-top
Neuroncore Utilization
NC0 NC1 NC2 NC3
ND0 [ 100%] [ 100%] [ 100%] [ 100%]
ND1 [ 0.00%] [ 0.00%] [ 0.00%] [ 0.00%]
ND2 [ 0.00%] [ 0.00%] [ 0.00%] [ 0.00%]
ND3 [ 0.00%] [ 0.00%] [ 0.00%] [ 0.00%]

vCPU and Memory Info
System vCPU Usage [ 8.69%, 9.47%] Runtime vCPU Usage [ 3.22%, 5.30%]
Runtime Memory Host [ 2.5MB/ 46.0GB] Runtime Memory Device 198.3MB

Loaded Models
Model ID Host Memory Device Memory
[-] ND 0
[-] NC0 10001 638.5KB 49.6MB
[-integ-tests/out-test7_resnet50_v2_fp16_b1_tpb1_tf
[+] NC1 638.5KB 49.6MB
[+] NC2 638.5KB 49.6MB
[+] NC3 638.5KB 49.6MB

Neuron Apps [1]:inference app 1 [2]:inference app 2 [3]:inference app 3 [4]:inference app 4
q: quit arrows: move tree selection enter: expand/collapse tree item x: expand/collapse entire tree a/d: previous/next tab 1-9: select tab

```

有关用于监控和优化基于 Neuron 的推理应用程序的资源的更多信息，请参阅 [Neuron 工具](#)。

使用 Neuron Monitor (Neuron 监视器)

Neuron Monitor 从系统上运行的 Neuron 运行时收集指标，并将收集的数据以格式流式传输到 stdout。JSON 这些指标按指标组进行组织，您可以通过提供配置文件进行配置。有关 Neuron Monitor 的更多信息，请参阅 [Neuron Monitor 用户指南](#)。

升级 Neuron 软件

有关如何在其中更新 Neuron SDK 软件的信息 DLAMI，请参阅《AWS 神经元 [设置指南](#)》。

下一个步骤

[与 Ne AWS ur DLAMI on 一起使用](#)

与 Ne AWS ur DLAMI on 一起使用

Ne AWS uron 的典型工作流程SDK是在编译服务器上编译之前训练过的机器学习模型。之后，将构件分发给 Inf1 实例以供执行。AWS Deep Learning AMIs (DLAMI) 预装了在使用 Inferentia 的 Inf1 实例中编译和运行推理所需的一切。

以下各节介绍如何使用 Inferent DLAMI ia。

内容

- [使用 TensorFlow-Neuron 和 Neuron 编译器 AWS](#)
- [使用 AWS 神经元服务 TensorFlow](#)
- [使用 MXNet-Neuron 和 Neuron 编译器 AWS](#)
- [使用MXNet神经元模型服务](#)
- [使用 PyTorch-Neuron 和 Neuron 编译器 AWS](#)

使用 TensorFlow-Neuron 和 Neuron 编译器 AWS

本教程演示如何使用 Ne AWS uron 编译器编译 Keras ResNet -50 模型并将其以格式导出为已保存的模型。SavedModel 这种格式是典型的 TensorFlow 模型可互换格式。您还可以学习如何使用示例输入，在 Inf1 实例上运行推理过程。

有关神经元的更多信息SDK，请参阅 Ne [AWS uron 文档](#)。SDK

内容

- [前提条件](#)
- [激活 Conda 环境](#)
- [Resnet50 编译](#)
- [ResNet50 推论](#)

前提条件

使用本教程之前，您应已完成 [启动带有 Ne AWS ur DLAMI on 的实例](#) 中的设置步骤。您还应该熟悉深度学习知识以及如何使用 DLAMI。

激活 Conda 环境

使用以下命令激活 TensorFlow-Neuron conda 环境：

```
source activate aws_neuron_tensorflow_p36
```

要退出当前 Conda 环境，请运行以下命令：

```
source deactivate
```

Resnet50 编译

创建一个名为 **tensorflow_compile_resnet50.py** 的 Python 脚本，其中包含以下内容。此 Python 脚本编译 Keras ResNet 50 模型并将其导出为已保存的模型。

```
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')
```

使用以下命令编译该模型：

```
python tensorflow_compile_resnet50.py
```

编译过程将需要几分钟时间。完成后，您的输出应与以下内容类似：

```
...
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
...
```

编译后，保存的模型将进行压缩，放置在 `ws_resnet50/resnet50_neuron.zip` 中。使用以下命令对模型解压缩，并下载用于推理的示例图像：

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

ResNet50 推论

创建一个名为 `tensorflow_infer_resnet50.py` 的 Python 脚本，其中包含以下内容。此脚本使用先前编译的推理模型，对下载的模式运行推理过程。

```
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])
```

使用以下命令对模型运行推理过程：

```
python tensorflow_infer_resnet50.py
```

您的输出应与以下内容类似：

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

下一个步骤

[使用 AWS 神经元服务 TensorFlow](#)

使用 AWS 神经元服务 TensorFlow

本教程展示了在导出保存的模型以用 AWS 于 Serving 之前，如何构造图形并添加 Neuron 编译步骤 TensorFlow。TensorFlow Serving 是一种服务系统，允许您在网络上扩大推理规模。Neuro TensorFlow n Serving 的使用方式与普通 TensorFlow 服务API相同。唯一的区别是，必须为 AWS Inferentia 编译保存的模型，并且入口点是名为的不同二进

制文件。tensorflow_model_server_neuron二进制文件位于 /usr/local/bin/tensorflow_model_server_neuron 中，并已预安装在 DLAMI 中。

有关神经元的更多信息 SDK，请参阅 [Neuron SDK](#)。

内容

- [前提条件](#)
- [激活 Conda 环境](#)
- [编译和导出保存的模型](#)
- [处理保存的模型](#)
- [生成发送给模型服务器的推理请求](#)

前提条件

使用本教程之前，您应已完成 [启动带有 Neuron 的 DLAMI 的实例](#) 中的设置步骤。您还应该熟悉深度学习知识以及如何使用 DLAMI。

激活 Conda 环境

使用以下命令激活 TensorFlow-Neuron conda 环境：

```
source activate aws_neuron_tensorflow_p36
```

如果需要退出当前 Conda 环境，请运行：

```
source deactivate
```

编译和导出保存的模型

创建一个名为 tensorflow-model-server-compile.py 的 Python 脚本，其中包含以下内容。该脚本构造一个图形并使用 Neuron 对其进行编译。然后，它将编译的图形导出为保存的模型。

```
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
modeldir = "./resnet50/1"
tf.saved_model.simple_save(sess, modeldir, inputs, outputs)

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

使用以下命令编译该模型：

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

您的输出应与以下内容类似：

```
...
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
```

处理保存的模型

当模型编译完成后，您可以使用以下命令，通过 `tensorflow_model_server_neuron` 二进制文件处理保存的模型：

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

您的输出应与以下内容类似。服务器将编译后的模型暂存在 Inferentia 设备中，为推理做准备。DRAM

```
...
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 ...
```

生成发送给模型服务器的推理请求

创建一个名为 `tensorflow-model-server-infer.py` 的 Python 脚本，其中包含以下内容。该脚本通过 `g` (RPC服务框架) 运行推理。

```
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))
```

使用RPC带有以下命令的 g 在模型上运行推理：

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

您的输出应与以下内容类似：

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

使用 MXNet-Neuron 和 Neuron 编译器 AWS

MXNet-Neuron 编译API提供了一种编译模型图的方法，您可以在 AWS Inferentia 设备上运行该模型。

在此示例中，您使用编译 ResNet -50 模型并使用它来运行推理。API

有关神经元的更多信息SDK，请参阅 [Neuron 文档](#)。SDK

内容

- [前提条件](#)
- [激活 Conda 环境](#)
- [Resnet50 编译](#)
- [ResNet50 推论](#)

前提条件

使用本教程之前，您应已完成 [启动带有 Neuron DLAMI on 的实例](#) 中的设置步骤。您还应该熟悉深度学习知识以及如何使用 DLAMI。

激活 Conda 环境

使用以下命令激活 MXNet-Neuron conda 环境：

```
source activate aws_neuron_mxnet_p36
```

要退出当前 Conda 环境，请运行：

```
source deactivate
```

Resnet50 编译

创建一个名为 **mxnet_compile_resnet50.py** 的 Python 脚本，其中包含以下内容。该脚本使用 MXNet-Neuron 编译 Python 编译 Python API 来编译 ResNet -50 模型。

```
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)
```

使用以下命令编译该模型：

```
python mxnet_compile_resnet50.py
```

编译需要几分钟。当编译完成后，以下文件将出现在您的当前目录中：

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

ResNet50 推论

创建一个名为 `mxnet_infer_resnet50.py` 的 Python 脚本，其中包含以下内容。此脚本会下载一个示例映像，然后使用该映像对已编译的模型运行推理过程。

```
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]))
```

使用以下命令对已编译模型运行推理过程：

```
python mxnet_infer_resnet50.py
```

您的输出应与以下内容类似：

```
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
```

下一个步骤

[使用MXNet神经元模型服务](#)

使用MXNet神经元模型服务

在本教程中，您将学习使用预训练的MXNet模型通过多模型服务器 (MMS) 执行实时图像分类。MMS 是一款灵活的 easy-to-use 工具，用于提供使用任何机器学习或深度学习框架训练的深度学习模型。本教程包括使用 AWS Neuron 的编译步骤和MMS使用的MXNet实现。

有关神经元的更多信息SDK，请参阅 [Neuron SDK](#)，请参阅 [AWS Neuron 文档](#)。SDK

内容

- [前提条件](#)
- [激活 Conda 环境](#)
- [下载示例代码](#)
- [编译模型](#)
- [运行推理](#)

前提条件

使用本教程之前，您应已完成 [启动带有 Neuron DLAMI 的实例](#) 中的设置步骤。您还应该熟悉深度学习知识以及如何使用 DLAMI。

激活 Conda 环境

使用以下命令激活 MXNet-Neuron conda 环境：

```
source activate aws_neuron_mxnet_p36
```

要退出当前 Conda 环境，请运行：

```
source deactivate
```

下载示例代码

要运行本示例，请使用以下命令下载示例代码：

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

编译模型

创建一个名为 multi-model-server-compile.py 的 Python 脚本，其中包含以下内容。此脚本将 ResNet 50 模型编译为 Inferentia 设备目标。

```
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)
```

要编译模型，请使用以下命令：

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

您的输出应与以下内容类似：

```
...
[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!
```

创建一个名为 `signature.json` 的文件，其中包含以下内容，以便配置输入名称和形状：

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

使用以下命令下载 `synset.txt` 文件。此文件是 ImageNet 预测类的名称列表。

```
curl -O https://s3.amazonaws.com/model-server/model_archive_1.0/examples/
squeezenet_v1.1/synset.txt
```

基于 `model_server_template` 文件夹中的模板，创建自定义服务类。使用以下命令，将模板复制到您的当前工作目录中：

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

编辑 `mxnet_model_service.py` 模块，将 `mx.cpu()` 上下文替换为 `mx.neuron()` 上下文，如下所示。你还需要注释掉不必要的数据副本，`model_input` 因为 MXNet-Neuron 不支持 `NDArray` 和 `Gluon`。APIs

```
...
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]
```

使用以下命令，通过模型归档程序对模型进行打包：

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

运行推理

启动多模型服务器并使用以下 RESTful API 命令加载使用它的模型。确保 `neuron-rtd` 正在使用默认设置运行。

```
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
```

通过以下命令，使用示例图像运行推理：

```
curl -O https://raw.githubusercontent.com/aws-labs/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
```

您的输出应与以下内容类似：

```
[
```

```
[
  {
    "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"
  }
]
```

要在测试结束后进行清理，请通过发出删除命令RESTfulAPI并使用以下命令停止模型服务器：

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

multi-model-server --stop
```

您应看到以下输出：

```
{
  "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
```

使用 PyTorch-Neuron 和 Neuron 编译器 AWS

PyTorch-Neuron 编译API提供了一种编译模型图的方法，您可以在 AWS Inferentia 设备上运行该模型。

经过训练的模型必须先编译为 Inferentia 目标，才能部署在 Inf1 实例上。以下教程编译 torchvision ResNet 50 模型并将其导出为已保存的模块。TorchScript 此模型随后将用于运行推理。

为方便起见，本教程使用 Inf1 实例进行编译和推理。在实际操作中，您也可以使用其他实例类型来编译模型，例如 c5 实例系列。然后，您必须将已编译的模型部署到 Inf1 推理服务器中。有关更多信息，请参阅 [AWS Neuron PyTorch SDK 文档](#)。

内容

- [前提条件](#)
- [激活 Conda 环境](#)
- [Resnet50 编译](#)
- [ResNet50 推论](#)

前提条件

使用本教程之前，您应已完成 [启动带有 Neuron DLAMI 的实例](#) 中的设置步骤。您还应该熟悉深度学习知识以及如何使用 DLAMI。

激活 Conda 环境

使用以下命令激活 PyTorch-Neuron conda 环境：

```
source activate aws_neuron_pytorch_p36
```

要退出当前 Conda 环境，请运行：

```
source deactivate
```

Resnet50 编译

创建一个名为 **pytorch_trace_resnet50.py** 的 Python 脚本，其中包含以下内容。该脚本使用 PyTorch-Neuron 编译 Python API 来编译 ResNet -50 模型。

Note

在编译 torchvision 模型时，您需要注意：torchvision 与 torch 软件包的版本之间存在依赖关系。这些依赖关系规则可以通过 pip 进行管理。Torchvision==0.6.1 与 torch==1.5.1 版本匹配，而 torchvision==0.8.2 与 torch==1.7.1 版本匹配。

```
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")
```

运行编译脚本。

```
python pytorch_trace_resnet50.py
```

编译需要几分钟。编译完成后，已编译的模型将以 `resnet50_neuron.pt` 的形式保存在本地目录中。

ResNet50 推论

创建一个名为 **pytorch_infer_resnet50.py** 的 Python 脚本，其中包含以下内容。此脚本会下载一个示例映像，然后使用该映像对已编译的模型运行推理过程。

```
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) )
```

使用以下命令对已编译模型运行推理过程：

```
python pytorch_infer_resnet50.py
```

您的输出应与以下内容类似：

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

ARM64 DLAMI

AWS ARM64GPUDLAMIs旨在为深度学习工作负载提供高性能和成本效益。具体而言，g5G 实例类型采用基于 ARM64 的 [G AWS raviton2 处理器](#)，该处理器是从头开始构建的，AWS 并针对客户在云中运行工作负载的方式进行了优化。AWS ARM64GPUDLAMIs已预先配置了 Docker、NVIDIA Docker、D NVIDIA river DNN、NCCL、Cu，以及流行的机器学习框架，例如和。CUDA TensorFlow PyTorch

借助 g5G 实例类型，您可以利用 Graviton2 的价格和性能优势，与带加速功能的 x86 实例相比，部署 GPU加速深度学习模型的成本要低得多。GPU

选择一个 ARM64 DLAMI

使用您选择的[实例启动 g5G 实例](#)。ARM64 DLAMI

有关启动的 step-by-step说明DLAMI，请参阅[启动和配置DLAMI](#)。

有关最新版本的列表 ARM64DLAMIs，请参阅的[发行说明DLAMI](#)。

开始使用

以下主题向您展示了如何开始使用ARM64DLAMI。

内容

- [使用 ARM64 GPU PyTorch DLAMI](#)

使用 ARM64 GPU PyTorch DLAMI

已准备好在基 AWS Deep Learning AMIs 于 Arm64 处理器的情况下使用 GPUs，并针对以下方面进行了优化。PyTorchARM64GPU PyTorch DLAMI 包括预先配置了、和的 Python 环境 [PyTorch](#) [TorchVision](#)，[TorchServe](#) 用于深度学习训练和推理用例。

内容

- [验证 PyTorch Python 环境](#)
- [使用运行训练示例 PyTorch](#)
- [使用运行推理示例 PyTorch](#)

验证 PyTorch Python 环境

使用以下命令来连接您的 G5g 实例并激活基础 Conda 环境：

```
source activate base
```

您的命令提示符应表明您正在基本 Conda 环境中工作，该环境包含 PyTorch TorchVision、和其他库。

```
(base) $
```

验证 PyTorch 环境的默认刀具路径：

```
(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)
```

使用运行训练示例 PyTorch

运行示例 MNIST 训练作业：

```
git clone https://github.com/pytorch/examples.git
```

```
cd examples/mnist
python main.py
```

您的输出应类似于以下内容：

```
...
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%)
```

使用运行推理示例 PyTorch

使用以下命令下载预训练的 densenet161 模型并使用以下命令运行推理： TorchServe

```
# 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
```

您的输出应类似于以下内容：

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

使用以下命令来注销 densenet161 模型并停止服务器：

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

您的输出应类似于以下内容：

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

推理

本节提供有关如何使用框架和工具运行推理DLAMI的教程。

推理工具

- [TensorFlow 服务](#)

模型处理

以下是 Conda 深度学习中安装AMI的模型服务选项。单击其中一个选项可了解如何使用该选项。

主题

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

TensorFlow 服务

TensorFlow Serving 是一款适用于机器学习模型的灵活、高性能的服务系统。

预装了 tensorflow-serving-api Conda 的 Deep Learning AMI 深度学习！您将找到用于训练、导出和提供 MNIST 模型的示例脚本 `~/examples/tensorflow-serving/`。

要运行这些示例中的任何一个，请先使用 Conda 连接到您的 Deep Learning AMI 深度学习并激活 TensorFlow 环境。

```
$ source activate tensorflow2_p310
```

现在，将目录更改至服务示例脚本文件夹。

```
$ cd ~/examples/tensorflow-serving/
```

处理预训练的 Inception 模型

以下是您可尝试为不同的模型（如 Inception）提供服务的示例。作为一般规则，您需要将可维护模型和客户端脚本下载到您的 DLAMI。

使用 Inception 模型处理和测试推理

1. 下载该模型。

```
$ curl -O https://s3-us-west-2.amazonaws.com/tf-test-models/INCEPTION.zip
```

2. 解压缩模型。

```
$ unzip INCEPTION.zip
```

3. 下载一张哈士奇的照片。

```
$ curl -O https://upload.wikimedia.org/wikipedia/commons/b/b5/Siberian_Husky_bieyed_Flickr.jpg
```

4. 启动服务器。请注意，对于 Amazon Linux，您必须将用于 `model_base_path` 的目录从 `/home/ubuntu` 更改为 `/home/ec2-user`。

```
$ tensorflow_model_server --model_name=INCEPTION --model_base_path=/home/ubuntu/
examples/tensorflow-serving/INCEPTION/INCEPTION --port=9000
```

5. 对于在前台运行的服务器，您需要启动另一个终端会话才能继续。打开一个新的终端并 TensorFlow 使用激活 `source activate tensorflow2_p310`。然后，使用您的首选文本编辑器创建具有以下内容的脚本。将它命名为 `inception_client.py`。此脚本将映像文件名用作参数，并从预训练模型中获得预测结果。

```
from __future__ import print_function

import grpc
import tensorflow as tf
import argparse

from tensorflow_serving.apis import predict_pb2
from tensorflow_serving.apis import prediction_service_pb2_grpc

parser = argparse.ArgumentParser(
    description='TF Serving Test',
    formatter_class=argparse.ArgumentDefaultsHelpFormatter
)
parser.add_argument('--server_address', default='localhost:9000',
                    help='Tensorflow Model Server Address')
parser.add_argument('--image', default='Siberian_Husky_bi-eyed_Flickr.jpg',
                    help='Path to the image')
args = parser.parse_args()

def main():
    channel = grpc.insecure_channel(args.server_address)
    stub = prediction_service_pb2_grpc.PredictionServiceStub(channel)
    # Send request
    with open(args.image, 'rb') as f:
        # See prediction_service.proto for gRPC request/response details.
        request = predict_pb2.PredictRequest()
        request.model_spec.name = 'INCEPTION'
        request.model_spec.signature_name = 'predict_images'

        input_name = 'images'
        input_shape = [1]
        input_data = f.read()
        request.inputs[input_name].CopyFrom(
```

```
tf.make_tensor_proto(input_data, shape=input_shape))

result = stub.Predict(request, 10.0) # 10 secs timeout
print(result)

print("Inception Client Passed")

if __name__ == '__main__':
    main()
```

6. 现在，运行将服务器位置和端口以及哈士奇照片的文件名作为参数传递的脚本。

```
$ python3 inception_client.py --server=localhost:9000 --image Siberian_Husky_bi-
eyed_Flickr.jpg
```

训练和服务MNIST模型

对于本教程，我们将导出模型并随后通过 `tensorflow_model_server` 应用程序处理它。最后，您可以使用示例客户端脚本测试模型服务器。

运行将训练和导出MNIST模型的脚本。作为脚本的唯一参数，您需要为其提供一个文件夹位置以保存该模型。现在，我们可以把它放入 `mnist_model` 中。该脚本将会为您创建此文件夹。

```
$ python mnist_saved_model.py /tmp/mnist_model
```

请耐心等待，因为此脚本可能需要一段时间才能提供输出。当训练完成并最终导出模型后，您应该看到以下内容：

```
Done training!
Exporting trained model to mnist_model/1
Done exporting!
```

下一步是运行 `tensorflow_model_server` 以处理导出的模型。

```
$ tensorflow_model_server --port=9000 --model_name=mnist --model_base_path=/tmp/
mnist_model
```

为您提供了一个客户端脚本来测试服务器。

要对其进行测试，您将需要打开一个新的终端窗口。

```
$ python mnist_client.py --num_tests=1000 --server=localhost:9000
```

更多功能和示例

如果您有兴趣了解有关 TensorFlow 服务的更多信息，请[TensorFlow 访问该网站](#)。

您也可以使用 [Amazon Elastic Inference TensorFlow](#) 提供服务。如需更多信息，请查看有关如何将 [Elastic Inference 与 TensorFlow 服务结合使用的指南](#)。

TorchServe

TorchServe 是一款灵活的工具，用于提供已从中导出的深度学习模型 PyTorch。TorchServe 预装了带有 Conda 的深度学习。

有关使用的更多信息 TorchServe，[请参阅 PyTorch 文档模型服务器](#)。

主题

在上提供图像分类模型 TorchServe

本教程介绍如何使用提供图像分类模型 TorchServe。它使用提供的 DenseNet -161 模型。PyTorch 服务器运行后，它会监听预测请求。在这种情况下，如果您上传图像（一张小猫的图像），服务器会返回在其上训练该模型的类中匹配的前 5 个类的预测。

在上提供图像分类模型示例 TorchServe

1. 使用 Conda v34 或更高版本使用深度学习AMI连接到亚马逊弹性计算云 (AmazonEC2) 实例。
2. 激活 pytorch_p310 环境。

```
source activate pytorch_p310
```

3. 克隆 TorchServe 存储库，然后创建一个目录来存储您的模型。

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

4. 使用模型存档程序来存档模型。该 extra-files 参数使用 TorchServe 存储库中的文件，因此如有必要，请更新路径。有关模型存档器的更多信息，请参阅 [Torch 模型存档器](#)。TorchServe

```
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. 运行 TorchServe 以启动终端节点。添加 `> /dev/null` 会使日志输出静音。

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

6. 下载小猫的图像并将其发送到 TorchServe 预测端点：

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

预测端点返回的预测值与以下前五名预测JSON类似，其中图像包含埃及猫的概率为47%，其次是有虎斑猫的概率为46%。

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

7. 当您完成测试时，停止服务器：

```
torchserve --stop
```

其他示例

TorchServe 提供了可以在您的DLAMI实例上运行的各种示例。您可以在 [TorchServe项目存储库示例页面上](#)查看它们。

更多信息

有关更多 TorchServe 文档，包括如何 TorchServe使用 Docker 进行设置和最新 TorchServe 功能，请参阅[上的 TorchServe GitHub项目页面](#)。

升级 DLAMI

在此处，您将找到有关升级 DLAMI 的信息以及有关在 DLAMI 上更新软件的提示。

一旦有修补程序和更新推出便立即应用，从而始终保持操作系统和其他已安装软件为最新。

如果您使用的是 Amazon Linux 或 Ubuntu，则当您登录到您 DLAMI 时，便会收到更新通知（如果有）并且会看到更新说明。有关 Amazon Linux 维护的更多信息，请参阅[更新实例软件](#)。对于 Ubuntu 实例，请参阅官方 [Ubuntu 文档](#)。

在 Windows 上，定期检查 Windows Update 有无软件和安全更新。如果您愿意，可以自动应用更新。

Important

有关 Meltdown 和 Spectre 漏洞以及如何修补操作系统以解决这些漏洞的信息，请参阅[安全公告-2018-013](#)。AWS

主题

- [升级到新版 DLAMI](#)
- [有关软件更新的提示](#)
- [有新的更新时收到通知](#)

升级到新版 DLAMI

DLAMI 的系统映像会定期更新，以利用新的深度学习框架版本、CUDA 其他软件更新和性能调整。如果您已使用 DLAMI 一段时间并想要利用更新，则需要启动新实例。您还必须手动传输任何数据集、检查点或其他宝贵的数据。相反，您可以使用 Amazon EBS 来保留您的数据并将其附加到新的数据中 DLAMI。通过这种方法，您可以经常升级，同时最大限度地减少转换数据所需的时间。

Note

在连接和移动 Amazon EBS 卷时 DLAMIs，必须将 DLAMIs 和新卷都放在同一个可用区中。

1. 使用 Amazon EC2 console 创建新的亚马逊 EBS 卷。有关详细说明，请参阅[创建 Amazon EBS 卷](#)。

2. 将新创建的 Amazon EBS 卷附加到现有卷上 DLAMI。有关详细说明，请参阅[附加 Amazon EBS 卷](#)。
3. 传输您的数据，如数据集、检查点和配置文件。
4. 启动 DLAMI。有关详细指导，请参阅[设置实DLAMI例](#)。
5. 将 Amazon EBS 卷与您的旧 DLAMI 卷分离。有关详细说明，请参阅[分离 Amazon EBS 卷](#)。
6. 将 Amazon EBS 卷附加到您的新卷上 DLAMI。请按照步骤 2 中的相关说明附加卷。
7. 在确认数据可用于新的 DLAMI 上时，停止并终止旧的 DLAMI。有关更详细的清理说明，请参阅[清理实DLAMI例](#)。

有关软件更新的提示

有时，您可能想要在 DLAMI 上手动更新软件。通常建议您使用 pip 来更新 Python 软件包。你还应该使用 pip 在 Conda 的深度学习 AMI 上在 Conda 环境中更新软件包。有关升级和安装说明，请参阅特定框架或软件的网站。

Note

我们不能保证软件包更新一定成功。尝试在依赖项不兼容的环境中升级软件包可能会导致安装失败。在这种情况下，您应该联系库维护人员，看看是否可以更新软件包依赖项。或者，您可以尝试以允许更新的方式来修改环境。但是，这种修改可能意味着删除或更新现有软件包，这意味着我们无法再保证此环境的稳定性。

它预 AWS Deep Learning AMIs 装了许多 Conda 环境和许多软件包。由于预装软件包数量众多，要找到一组保证兼容的软件包非常困难。您可能会看到一条警告“环境不一致，请仔细检查套餐计划”。DLAMI 确保所有 DLAMI 提供的环境都正确无误，但不能保证任何用户安装的软件包都能正常运行。

有新的更新时收到通知

Note

AWS 深度学习 AMIs 每周都会发布安全补丁。将针对这些增量安全补丁发送发布通知，尽管它们可能未包含在官方发布说明中。

每当有新 DLAMI 内容发布时，您都可以收到通知。通知 SNS 使用以下主题在 [Amazon](#) 上发布。

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

当有新消息发布时，会在此DLAMI处发布消息。的版本、元数据和区域 AMI ID AMI 将包含在消息中。

可以使用几种不同的方法接收这些消息。我们建议您使用以下方法。

1. 打开 [Amazon SNS 控制台](#)。
2. 如有必要，在导航栏中将 AWS 区域更改为美国西部（俄勒冈）。您必须选择创建您订阅的SNS通知的区域。
3. 在导航窗格中，依次选择订阅、创建订阅。
4. 对于 Create subscription 对话框，执行以下操作：
 - a. 对于主题 ARN，复制并粘贴以下 Amazon 资源名称 (ARN)：**arn:aws:sns:us-west-2:767397762724:dlami-updates**
 - b. 对于协议，请从 [Amazon SQS、AWS Lambda、Email、Email-JSON] 中选择一个
 - c. 对于 Endpoint，输入您将用于接收通知的资源的电子邮件地址或 Amazon 资源名称 (ARN)。
 - d. 选择创建订阅。
5. 您将收到一封主题为 AWS 通知 - 订阅确认 的确认电子邮件。打开该电子邮件，选择确认订阅来完成订阅。

安全性 AWS Deep Learning AMIs

云安全 AWS 是重中之重。作为 AWS 客户，您可以受益于专为满足大多数安全敏感型组织的要求而构建的数据中心和网络架构。

安全是双方共同承担 AWS 的责任。[责任共担模式](#)将其描述为云的安全性和云中的安全性：

- 云安全 — AWS 负责保护在云 AWS 服务 中运行的基础架构 AWS Cloud。AWS 还为您提供可以安全使用的服务。作为[AWS 合规计划](#)的一部分，第三方审计师定期测试和验证我们安全的有效性。要了解适用的合规计划 AWS Deep Learning AMIs，请参阅按合规计划划分的[范围内的AWS服务按合规计划](#)。
- 云端安全 — 您的责任由您 AWS 服务 使用的内容决定。您还需要对其它因素负责，包括您的数据的敏感性、您的公司的要求以及适用的法律法规。

此文档将帮助您了解如何在使用 DLAMI 时应用责任共担模式。以下主题说明如何配置 DLAMI 以实现您的安全性和合规性目标。您还将学习如何使用其他方法 AWS 服务 来帮助您监控和保护您的DLAMI 资源。

有关更多信息，请参阅《[亚马逊EC2用户指南](#)》EC2中的“[亚马逊安全](#)”。

主题

- [中的数据保护 AWS Deep Learning AMIs](#)
- [的身份和访问管理 AWS Deep Learning AMIs](#)
- [合规性验证 AWS Deep Learning AMIs](#)
- [韧性在 AWS Deep Learning AMIs](#)
- [中的基础设施安全 AWS Deep Learning AMIs](#)
- [监控 AWS Deep Learning AMIs 实例](#)

中的数据保护 AWS Deep Learning AMIs

分 AWS [担责任模型](#)适用于中的数据保护 AWS Deep Learning AMIs。如本模型所述 AWS，负责保护运行所有内容的全球基础架构 AWS Cloud。您负责维护对托管在此基础结构上的内容的控制。您还负责您所使用的 AWS 服务 的安全配置和管理任务。有关数据隐私的更多信息，请参阅[数据隐私FAQ](#)。有关欧洲数据保护的信息，请参阅[责任AWS 共担模型和AWS安全GDPR](#)博客上的博客文章。

出于数据保护目的，我们建议您保护 AWS 账户凭据并使用 AWS IAM Identity Center 或 AWS Identity and Access Management (IAM) 设置个人用户。这样，每个用户只获得履行其工作职责所需的权限。还建议您通过以下方式保护数据：

- 对每个账户使用多重身份验证 (MFA)。
- 使用SSL/TLS与 AWS 资源通信。我们需要 TLS 1.2，建议使用 TLS 1.3。
- 使用API进行设置和用户活动记录 AWS CloudTrail。有关使用 CloudTrail 跟踪捕获 AWS 活动的信息，请参阅《AWS CloudTrail 用户指南》中的[使用跟 CloudTrail 踪](#)。
- 使用 AWS 加密解决方案以及其中的所有默认安全控件 AWS 服务。
- 使用高级托管安全服务（例如 Amazon Macie），它有助于发现和保护存储在 Amazon S3 中的敏感数据。
- 如果您在 AWS 通过命令行界面或访问时需要 FIPS 140-3 经过验证的加密模块API，请使用端点。FIPS有关可用FIPS端点的更多信息，请参阅[联邦信息处理标准 \(FIPS\) 140-3](#)。

强烈建议您切勿将机密信息或敏感信息（如您客户的电子邮件地址）放入标签或自由格式文本字段（如名称字段）。这包括您使用DLAMI或以其他 AWS 服务方式使用控制台时API、AWS CLI、或 AWS SDKs。在用于名称的标签或自由格式文本字段中输入的任何数据都可能会用于计费或诊断日志。如果您URL向外部服务器提供，我们强烈建议您不要在中包含凭据信息，URL以验证您对该服务器的请求。

的身份和访问管理 AWS Deep Learning AMIs

AWS Identity and Access Management (IAM) AWS 服务可以帮助管理员安全地控制对 AWS 资源的访问权限。IAM管理员控制谁可以通过身份验证（登录）和授权（拥有权限）使用DLAMI资源。IAM无需支付额外费用即可使用。AWS 服务

有关身份和访问管理的更多信息，请参阅 [Amazon 的身份和访问管理EC2](#)。

主题

- [使用身份进行身份验证](#)
- [使用策略管理访问](#)
- [IAM与亚马逊合作 EMR](#)

使用身份进行身份验证

身份验证是您 AWS 使用身份凭证登录的方式。您必须以 AWS 账户根用户、IAM 用户身份或通过担任 IAM 角色进行身份验证（登录 AWS）。

您可以使用通过身份源提供的凭据以 AWS 联合身份登录。AWS IAM Identity Center（IAM 身份中心）用户、贵公司的单点登录身份验证以及您的 Google 或 Facebook 凭据就是联合身份的示例。在您以联合身份登录时，您的管理员以前使用 IAM 角色设置了身份联合验证。当您使用联合访问 AWS 时，您就是在间接扮演一个角色。

根据您的用户类型，您可以登录 AWS Management Console 或 AWS 访问门户。有关登录的更多信息 AWS，请参阅《AWS 登录 用户指南》中的[如何登录到您 AWS 账户](#)的。

如果您 AWS 以编程方式访问，则会 AWS 提供软件开发套件 (SDK) 和命令行接口 (CLI)，以便使用您的凭据对请求进行加密签名。如果您不使用 AWS 工具，则必须自己签署请求。有关使用推荐的方法自行签署请求的更多信息，请参阅《IAM 用户指南》中的[API 请求 AWS 签名版本 4](#)。

无论使用何种身份验证方法，您可能需要提供其他安全信息。例如，AWS 建议您使用多重身份验证 (MFA) 来提高账户的安全性。要了解更多信息，请参阅用户指南中的[多因素身份验证](#)和 AWS IAM Identity Center 用户指南 IAM 中的[AWS 多因素身份验证](#)。IAM

AWS 账户 root 用户

创建时 AWS 账户，首先要有一个登录身份，该身份可以完全访问账户中的所有资源 AWS 服务和资源。此身份被称为 AWS 账户 root 用户，使用您创建账户时使用的电子邮件地址和密码登录即可访问该身份。强烈建议您不要使用根用户执行日常任务。保护好根用户凭证，并使用这些凭证来执行仅根用户可以执行的任务。有关需要您以根用户身份登录的任务的完整列表，请参阅《IAM 用户指南》中的[需要根用户凭证的任务](#)。

IAM 用户和组

[IAM 用户](#)是您内部 AWS 账户对个人或应用程序具有特定权限的身份。在可能的情况下，我们建议使用临时证书，而不是创建拥有密码和访问密钥等长期凭证的 IAM 用户。但是，如果您有需要 IAM 用户长期凭证的特定用例，我们建议您轮换访问密钥。有关更多信息，请参阅《IAM 用户指南》中的[对于需要长期凭证的使用场景定期轮换访问密钥](#)。

[IAM 组](#)是指定一个 IAM 用户集合的身份。您不能使用组的身份登录。您可以使用组来一次性为多个用户指定权限。如果有大量用户，使用组可以更轻松地管理用户权限。例如，您可以拥有一个名为的群组，IAMAdmins 并授予该群组管理 IAM 资源的权限。

用户与角色不同。用户唯一地与某个人员或应用程序关联，而角色旨在让需要它的任何人代入。用户具有永久的长期凭证，而角色提供临时凭证。要了解更多信息，请参阅 [《IAM用户指南》中的IAM用户用例](#)。

IAM 角色

[IAM角色](#)是您内部具有特定权限 AWS 账户 的身份。它类似于 IAM 用户，但未与特定人员关联。要在中临时扮演角色 AWS Management Console，可以[从用户切换到IAM角色（控制台）](#)。您可以通过调用 AWS CLI 或 AWS API操作或使用自定义操作来代入角色URL。有关使用角色的方法的更多信息，请参阅 [《IAM用户指南》中的代入角色的方法](#)。

具有临时凭证的 IAM 角色在以下情况下很有用：

- 联合用户访问 – 要向联合身份分配权限，请创建一个角色，并为该角色定义权限。当联合身份进行身份验证时，该身份将与角色相关联并被授予由此角色定义的权限。有关用于联合身份验证的角色的信息，请参阅 [《IAM用户指南》中的为第三方身份提供商（联合）创建角色](#)。如果您使用 IAM Identity Center，则需要配置权限集。为了控制您的身份在进行身份验证后可以访问的内容，IAM Identity Center 会将权限集关联到中的IAM角色。有关权限集的信息，请参阅 [《AWS IAM Identity Center 用户指南》中的权限集](#)。
- 临时IAM用户权限-IAM 用户或角色可以代入一个IAM角色，为特定任务临时获得不同的权限。
- 跨账户存取 - 您可以使用 IAM 角色允许其他账户中的某个人（可信任主体）访问您账户中的资源。角色是授予跨账户存取权限的主要方式。但是，对于某些资源 AWS 服务，您可以将策略直接附加到资源（而不是使用角色作为代理）。要了解角色和基于资源的跨账户访问策略之间的区别，请参阅IAM用户指南[IAM中的跨账户资源访问权限](#)。
- 跨服务访问 — 有些 AWS 服务 使用其他 AWS 服务服务中的功能。例如，当您在服务中拨打电话时，该服务通常会在 Amazon 中运行应用程序EC2或在 Amazon S3 中存储对象。服务可能会使用发出调用的主体的权限、使用服务角色或使用服务相关角色来执行此操作。
- 转发访问会话 (FAS)-当您使用IAM用户或角色在中执行操作时 AWS，您被视为委托人。当你使用某些服务时，你可能会执行一个操作，然后在不同的服务中启动另一个操作。FAS使用调用委托人的权限 AWS 服务以及 AWS 服务 向下游服务发出请求的请求。FAS只有当服务收到需要与其他 AWS 服务 或资源交互才能完成的请求时，才会发出请求。在这种情况下，您必须具有执行这两项操作的权限。有关提出FAS请求时的政策详情，请参阅[转发访问会话](#)。
- 服务角色-服务[IAM角色](#)是服务代替您执行操作的角色。IAM 管理员可以在 IAM 中创建、修改和删除服务角色。有关更多信息，请参阅 [《IAM用户指南》AWS 服务中的创建角色以向委派权限](#)。
- 服务相关角色-服务相关角色是一种与服务相关联的服务角色。AWS 服务服务可以代入代表您执行操作的角色。服务相关角色出现在您的中 AWS 账户，并且归服务所有。IAM 管理员可以查看但不能编辑服务相关角色的权限。

- 在 Amazon 上运行的应用程序 EC2 — 您可以使用 IAM 角色管理在 EC2 实例上运行并发出 AWS CLI 或 AWS API 请求的应用程序的临时证书。这优先于在 EC2 实例中存储访问密钥。要为 EC2 实例分配 AWS 角色并使其可供其所有应用程序使用，您需要创建一个附加到该实例的实例配置文件。实例配置文件包含角色，并使 EC2 实例上运行的程序能够获得临时凭证。有关更多信息，请参阅 IAM 用户指南中的[使用 IAM 角色向在 Amazon EC2 实例上运行的应用程序授予权限](#)。

使用策略管理访问

您可以 AWS 通过创建策略并将其附加到 AWS 身份或资源来控制中的访问权限。策略是其中的一个对象 AWS，当与身份或资源关联时，它会定义其权限。AWS 在委托人（用户、root 用户或角色会话）发出请求时评估这些策略。策略中的权限确定是允许还是拒绝请求。大多数策略都以 JSON 文档的 AWS 形式存储在中。有关 JSON 策略文档结构和内容的更多信息，请参阅 [《IAM 用户指南》中的 JSON 策略概述](#)。

管理员可以使用 AWS JSON 策略来指定谁有权访问什么。也就是说，哪个主体可以对什么资源执行操作，以及在什么条件下执行。

默认情况下，用户和角色没有权限。要授予用户对所需资源执行操作的权限，IAM 管理员可以创建 IAM 策略。管理员随后可以向角色添加 IAM 策略，用户可以代入这些角色。

IAM 策略定义操作的权限，无论您使用哪种方法执行操作。例如，假设您有一个允许 `iam:GetRole` 操作的策略。拥有该策略的用户可以从 AWS Management Console AWS CLI、或获取角色信息 AWS API。

基于身份的策略

基于身份的策略是可以附加到身份（例如 IAM 用户、用户组或角色）的 JSON 权限策略文档。这些策略控制用户和角色可在何种条件下对哪些资源执行哪些操作。要了解如何创建基于身份的策略，请参阅 IAM 用户指南中的[使用客户托管策略定义自定义 IAM 权限](#)。

基于身份的策略可以进一步归类为内联策略或托管式策略。内联策略直接嵌入单个用户、组或角色中。托管策略是独立的策略，您可以将其附加到中的多个用户、群组和角色 AWS 账户。托管策略包括 AWS 托管策略和客户托管策略。要了解如何在托管策略或内联策略之间进行[选择，请参阅《IAM 用户指南》中的在托管策略和内联策略之间进行选择](#)。

基于资源的策略

基于资源的 JSON 策略是您附加到资源的策略文档。基于资源的策略的示例包括 IAM 角色信任策略和 Amazon S3 桶策略。在支持基于资源的策略的服务中，服务管理员可以使用它们来控制对特定资源的

访问。对于在其中附加策略的资源，策略定义指定主体可以对该资源执行哪些操作以及在什么条件下执行。您必须在基于资源的策略中[指定主体](#)。委托人可以包括账户、用户、角色、联合用户或 AWS 服务。

基于资源的策略是位于该服务中的内联策略。您不能在基于资源的策略IAM中使用 AWS 托管策略。

访问控制列表 (ACLs)

访问控制列表 (ACLs) 控制哪些委托人 (账户成员、用户或角色) 有权访问资源。ACLs与基于资源的策略类似，尽管它们不使用JSON策略文档格式。

Amazon S3 AWS WAF、和亚马逊VPC就是支持的服务示例ACLs。要了解更多信息ACLs，请参阅《亚马逊简单存储服务开发者指南》中的[访问控制列表 \(ACL\) 概述](#)。

其他策略类型

AWS 支持其他不太常见的策略类型。这些策略类型可以设置更常用的策略类型向您授予的最大权限。

- 权限边界-权限边界是一项高级功能，您可以在其中设置基于身份的策略可以向IAM实体 (IAM用户或角色) 授予的最大权限。您可为实体设置权限边界。这些结果权限是实体基于身份的策略及其权限边界的交集。在 Principal 中指定用户或角色的基于资源的策略不受权限边界限制。任一项策略中的显式拒绝将覆盖允许。有关权限边界的更多信息，请参阅《IAM 用户指南》中的[IAM 实体的权限边界](#)。
- 服务控制策略 (SCPs)-SCPs 是为中的组织或组织单位 (OU) 指定最大权限的JSON策略 AWS Organizations。AWS Organizations 是一项用于对您的企业拥有的多 AWS 账户 项进行分组和集中管理的服务。如果您启用组织中的所有功能，则可以将服务控制策略 (SCPs) 应用于您的任何或所有帐户。对成员账户中的实体 (包括每个实体) 的权限进行了SCP限制 AWS 账户根用户。有关 Organization SCPs s 和的更多信息，请参阅《AWS Organizations 用户指南》中的[服务控制策略](#)。
- 资源控制策略 (RCPs) — RCPs 这些JSON策略可用于设置账户中资源的最大可用权限，而无需更新附加到您拥有的每项资源的IAM策略。这会RCP限制成员账户中资源的权限，并可能影响身份 (包括身份) 的有效权限 AWS 账户根用户，无论这些身份是否属于您的组织。有关 Organizations 的更多信息RCPs，包括 AWS 服务 该支持的列表RCPs，请参阅《AWS Organizations 用户指南》中的[资源控制策略 \(RCPs\)](#)。
- 会话策略：会话策略是当您以编程方式为角色或联合用户创建临时会话时作为参数传递的高级策略。结果会话的权限是用户或角色的基于身份的策略和会话策略的交集。权限也可以来自基于资源的策略。任一项策略中的显式拒绝将覆盖允许。有关更多信息，请参阅《IAM 用户指南》中的[会话策略](#)。

多个策略类型

当多个类型的策略应用于一个请求时，生成的权限更加复杂和难以理解。要了解在涉及多种策略类型时如何 AWS 确定是否允许请求，请参阅IAM用户指南中的[策略评估逻辑](#)。

IAM与亚马逊合作 EMR

您可以IAM与 Amazon 一起使用EMR来定义用户、AWS 资源、群组、角色和策略。您还可以控制AWS 服务 这些用户和角色可以访问哪些用户。

有关在亚马逊上使用的更多信息EMR，请参阅 IAM Amazon f [AWS Identity and Access Management or Amazon EMR](#)。

合规性验证 AWS Deep Learning AMIs

AWS Deep Learning AMIs 作为多个合规计划的一部分，第三方审计师对安全性和 AWS 合规性进行评估。有关支持的合规计划的信息，请参阅 [Amazon 合规性验证EC2](#)。

有关特定合规计划范围 AWS 服务 内的列表，请参阅按合规计划划分的[范围内的AWSAWS 服务按合规计划](#)。有关一般信息，请参阅[AWS 合规计划AWS](#)。

您可以使用下载第三方审计报告 AWS Artifact。有关更多信息，请参阅 [Downloading Reports in AWS Artifact](#)。

您在使用DLAMI时的合规责任取决于您的数据的敏感性、贵公司的合规目标以及适用的法律和法规。AWS 提供了以下资源来帮助实现合规性：

- [安全性与合规性快速入门指南](#) - 这些部署指南讨论了架构注意事项，并提供了在 AWS上部署基于安全性和合规性的基准环境的步骤。
- [AWS 合规资源AWS](#) — 此工作簿和指南集可能适用于您所在的行业和所在地区。
- [使用AWS Config 开发人员指南中的 AWS Config 规则评估资源](#) — 该 AWS Config 服务评估您的资源配置在多大程度上符合内部实践、行业准则和法规。
- [AWS Security Hub](#)— 这 AWS 服务 提供了您内部安全状态的全面视图 AWS。Security Hub 使用安全控制来评估您的 AWS 资源，并检查您是否符合安全行业标准和最佳实践。

韧性在 AWS Deep Learning AMIs

AWS 全球基础设施是围绕 AWS 区域 可用区构建的。AWS 区域 提供多个物理分隔和隔离的可用区，这些可用区通过低延迟、高吞吐量和高度冗余的网络连接。利用可用区，您可以设计和操作在可用区之

间无中断地自动实现失效转移的应用程序和数据库。与传统的单个或多个数据中心基础结构相比，可用区具有更高的可用性、容错性和可扩展性。

有关 AWS 区域 和可用区的更多信息，请参阅[AWS 全球基础设施](#)。

有关有助于满足您的数据弹性和备份需求的 Amazon EC2 功能的信息，请参阅《[亚马逊EC2用户指南](#)》[EC2中的 Amazon 弹性](#)。

中的基础设施安全 AWS Deep Learning AMIs

的基础设施安全由 Amazon 提供支持EC2。AWS Deep Learning AMIs 有关更多信息，请参阅《[亚马逊EC2用户指南](#)》[EC2中的 Amazon 基础设施安全](#)。

监控 AWS Deep Learning AMIs 实例

监控是维护您的 AWS Deep Learning AMIs 实例和其他 AWS 解决方案的可靠性、可用性和性能的重要组成部分。您的DLAMI实例附带多种GPU监控工具，包括向 Amazon 报告GPU使用情况统计数据的实用工具 CloudWatch。有关更多信息[GPU监控和优化](#)，请参阅《[亚马逊EC2用户指南](#)》中的“[监控亚马逊EC2资源](#)”。

选择退出实例的使用情况跟踪 DLAMI

以下 AWS Deep Learning AMIs 操作系统发行版包括 AWS 允许收集实例类型、实例 ID、DLAMI类型和操作系统信息的代码。

Note

AWS 不会收集或保留有关的任何其他信息DLAMI，例如您在中使用的命令DLAMI。

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

选择退出使用情况跟踪

如果您愿意，您可以选择退出新DLAMI实例的使用情况跟踪。要选择退出，您必须在启动期间向 Amazon EC2 实例添加标签。标签应使用密钥 `OPT_OUT_TRACKING`，并将关联值设置为 `true`。有关更多信息，请参阅 [《亚马逊EC2用户指南》中的为您的亚马逊EC2资源添加标签](#)。

DLAMI框架支持政策

在这里，您可以找到 AWS Deep Learning AMIs (DLAMI) 框架的支持政策的详细信息。

有关 AWS 当前支持的 DLAMI 框架的列表，请参阅 Framework Supp [DLAMIort Policy](#) 页面。在该页面上的表中，请记住以下内容：

- 当前版本以 x.y.z 格式指定框架版本。在这种格式中，x 表示主要版本，y 表示次要版本，z 表示补丁版本。例如，对于 TensorFlow 2.10.1，主版本为 2，次要版本为 10，补丁版本为 1。
- 补丁结束指定该框架版本 AWS 支持多长时间。

有关具体内容的详细信息 DLAMIs，请参阅 [DLAMIs 发行说明](#)。

DLAMI框架支持 FAQs

- [哪些框架版本会获得安全补丁？](#)
- [AWS 发布新框架版本时会发布哪些镜像？](#)
- [哪些图像获得了新的 SageMaker AI/AWS 功能？](#)
- [“支持的框架”表中是如何定义当前版本的？](#)
- [如果我运行的版本不在“支持的框架”表中，该怎么办？](#)
- [是否 DLAMIs 支持以前版本的 TensorFlow？](#)
- [如何找到支持的框架版本的最新补丁映像？](#)
- [多长时间发布一次新映像？](#)
- [运行工作负载时，能在我的实例上以替代方式安装补丁吗？](#)
- [如果有新的补丁或更新的框架版本可用，会发生什么呢？](#)
- [是否可在不更改框架版本的情况下更新依赖项？](#)
- [对我的框架版本的主动支持何时结束？](#)
- [对于框架版本不再主动维护的映像，会为其安装补丁吗？](#)
- [如何使用旧框架版本？](#)
- [如何保持框架及其版本 up-to-date 的支持变更？](#)
- [是否需要商业许可证才能使用 Anaconda 存储库？](#)

哪些框架版本会获得安全补丁？

如果框架版本在 [AWS Deep Learning AMIs 框架支持策略表](#) 中标记为已支持，它就会获得安全补丁。

AWS 发布新框架版本时会发布哪些镜像？

我们会在 TensorFlow 和 PyTorch 的新版本发布后 DLAMIs 不久发布新版本。这包括框架的主要版本、主要的次要版本和 major-minor-patch 版本。当新版本的驱动程序和库可用时，我们也会更新映像。有关映像维护的更多信息，请参阅 [对我的框架版本的主动支持何时结束？](#)。

哪些图像获得了新的 SageMaker AI/AWS 功能？

新功能通常在最新版本的 for DLAMIs PyTorch 和中发布 TensorFlow。有关新 SageMaker AI 或 AWS 功能的详细信息，请参阅特定图像的发行说明。有关可用版本的列表 DLAMIs，请参阅的 [发行说明 DLAMI](#)。有关映像维护的更多信息，请参阅 [对我的框架版本的主动支持何时结束？](#)。

“支持的框架”表中是如何定义当前版本的？

Frame AWS work S [AWS Deep Learning AMIs support Policy 表](#) 中的当前版本是指在上提供的最新框架版本 GitHub。每个最新版本都包括中驱动程序、库和相关软件包的更新 DLAMI。有关映像维护的信息，请参阅 [对我的框架版本的主动支持何时结束？](#)

如果我运行的版本不在“支持的框架”表中，该怎么办？

如果您运行的版本不在 [AWS Deep Learning AMIs 框架支持策略表](#) 中，则可能没有最新的驱动程序、库和相关包。要获得更多 up-to-date 版本，我们建议您使用自己选择的最新 DLAMI 版本升级到支持的框架之一。有关可用版本的列表 DLAMIs，请参阅的 [发行说明 DLAMI](#)。

是否 DLAMIs 支持以前版本的 TensorFlow？

不是。如 Framework Support P [olicy 表所述，我们支持](#) 每个框架最新主要版本的最新补丁版本，自其首次 GitHub 发布起 365 天后 AWS Deep Learning AMIs 发布。有关更多信息，请参阅 [如果我运行的版本不在“支持的框架”表中，该怎么办？](#)

如何找到支持的框架版本的最新补丁映像？

要在最新框架版本中 DLAMI 使用，请检索 [DLAMI ID](#)，然后使用它 DLAMI 通过 [EC2 控制台](#) 启动。有关检索 AWS Deep Learning AMIs ID 的示例 AWS CLI 命令，请参阅 DLAMI 发行说明页面 [单框架 DLAMI 发行说明](#)。您选择的框架版本必须在 [AWS Deep Learning AMIs 框架支持策略表](#) 中标记为已支持。

多长时间发布一次新映像？

提供更新的补丁版本是我们的首要任务。我们通常会尽早创建安装了补丁的映像。我们会监控新修补的框架版本（例如 TensorFlow 2.9 到 TensorFlow 2.9.1）和新的次要发行版本（例如 TensorFlow 2.9 到 TensorFlow 2.10），并尽早提供它们。当发布带有新版本的 TensorFlow 现有版本时CUDA，我们会DLAMI针对该版本发布支持新CUDA版本的新版本。TensorFlow

运行工作负载时，能在我的实例上以替代方式安装补丁吗？

不是。的补丁更新不DLAMI是“就地”更新。

您必须开启新EC2实例，迁移工作负载和脚本，然后关闭之前的实例。

如果有新的补丁或更新的框架版本可用，会发生什么呢？

请定期查看发布说明页面以获取您的映像。我们鼓励您在新的补丁或更新的框架可用时将框架升级。有关可用版本的列表DLAMIs，请参阅的[发行说明DLAMI](#)。

是否可在不更改框架版本的情况下更新依赖项？

我们在不更改框架版本的情况下更新依赖项。但是，如果依赖项更新导致不兼容，我们会创建不同版本的映像。请务必查看[发行说明以DLAMI获取](#)更新的依赖项信息。

对我的框架版本的主动支持何时结束？

DLAMI图像是不可变的。一旦创建，就不会改变。结束对框架版本的主动支持涉及四个主要原因：

- [框架版本（补丁）升级](#)
- [AWS 安全补丁](#)
- [补丁结束日期（已过期）](#)
- [依赖关系 end-of-support](#)

Note

由于版本补丁升级和安全补丁的频率很高，我们建议您DLAMI经常查看发行说明页面，并在进行更改时进行升级。

框架版本（补丁）升级

如果您的DLAMI工作负载基于 TensorFlow 2.7.0，并且在 2.7.1 版本上 TensorFlow 发布 GitHub，则会发布 2.7.1 AWS 版本的新DLAMI工作负载。TensorFlow 2.7.1 版本的新镜像发布后，将不再主动维护之前的 TensorFlow 2.7.0 镜像。版本DLAMI为 TensorFlow 2.7.0 的版本不会收到更多补丁。然后，TensorFlow 2.7 的DLAMI发行说明页面将使用最新信息进行更新。没有为每个次要补丁提供单独的发布说明页面。

由于补丁升级而DLAMIs创建的新用户将使用新 [AMIID](#) 指定。

AWS 安全补丁

如果您的工作负载基于 TensorFlow 2.7.0 版本的映像并 AWS 制作了安全补丁，则会发布适用于 TensorFlow 2.7.0 的新版本。DLAMI TensorFlow 2.7.0 版图像的先前版本已不再活跃维护。有关更多信息，请参阅[运行工作负载时，能在我的实例上以替代方式安装补丁吗？](#)有关查找最新版本步骤 DLAMI，请参阅 [如何找到支持的框架版本的最新补丁映像？](#)

由于补丁升级而DLAMIs创建的新用户将使用新 [AMIID](#) 指定。

补丁结束日期（已过期）

DLAMIs在 GitHub 发布日期 365 天后，他们的补丁结束日期。

对于[多框架 DLAMIs](#)，当其中一个框架版本更新时，需要更新版本的新DLAMI版本。DLAMI使用旧框架版本的，不再积极维护。

Important

当有重大框架更新时，我们会例外处理。例如。如果 TensorFlow 1.15 更新到 TensorFlow 2.0，那么我们将在自 GitHub 发布之日起两年内继续支持最新版本的 TensorFlow 1.15，或者在 Origin 框架维护团队取消支持后的六个月内（以较早的日期为准）。

依赖关系 end-of-support

如果您正在使用 Python 3.6 在 TensorFlow 2.7.0 DLAMI 映像上运行工作负载，并且该版本的 Python 已标记为 end-of-support，则所有基于 Python 3.6 的DLAMI图像都将不再被主动维护。同样，如果标记了像 Ubuntu 16.04 这样的操作系统版本 end-of-support，则所有依赖于 Ubuntu 16.04 的DLAMI映像都将不再被主动维护。

对于框架版本不再主动维护的映像，会为其安装补丁吗？

不会。不再主动维护的图像就不会有新版本。

如何使用旧框架版本？

要在DLAMI较旧的框架版本中使用，请检索 [DLAMIID](#)，然后使用它DLAMI通过[EC2控制台](#)启动。有关检索 AMI ID 的 AWS CLI命令，请参阅[单框架发行说明中的DLAMI发行说明](#)页面。

如何保持框架及其版本 up-to-date的支持变更？

使用 [Fr up-to-date amework Support Policy 表 \(DLAMI发行说明 \)](#) 继续关注DLAMI AWS Deep Learning AMIs [框架](#)和版本。

是否需要商业许可证才能使用 Anaconda 存储库？

Anaconda 转向了针对某些用户的商业许可模式。积极维护DLAMIs已从Anaconda频道迁移到公开可用的开源版本的Conda ([conda-forge](#))。

重要的NVIDIA驱动程序变更为 DLAMIs

2023 年 11 月 15 日，对与DLAMIs使用的NVIDIA驱动程序相关的 AWS Deep Learning AMIs (DLAMI) AWS 进行了重要更改。有关更改内容以及更改是否会影响您的使用情况的信息DLAMIs，请参阅[DLAMINVIDIA司机更换 FAQs](#)。

DLAMINVIDIA司机更换 FAQs

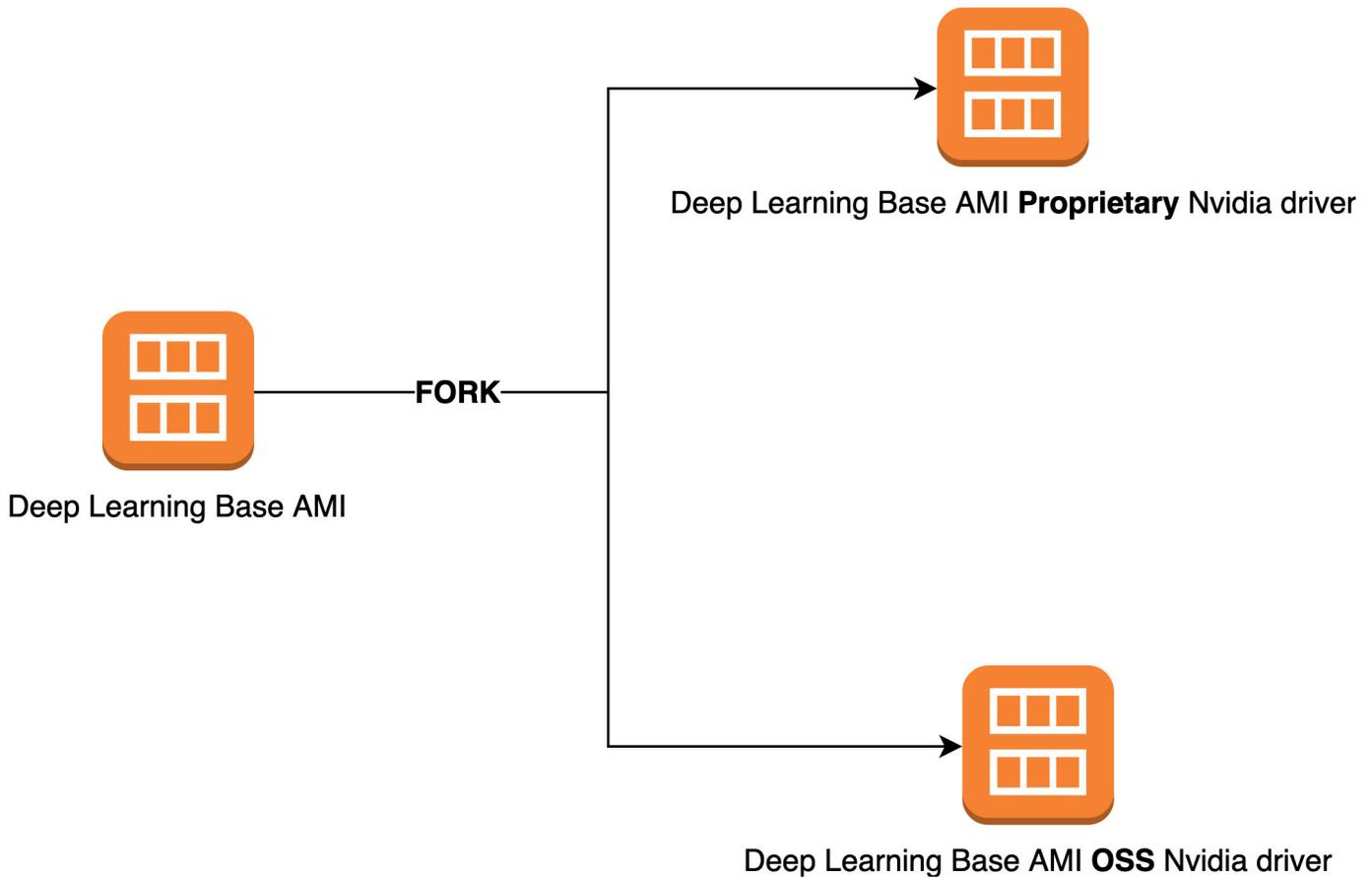
- [更改了哪些内容？](#)
- [为什么需要进行此更改？](#)
- [这一变化影响DLAMIs了哪些？](#)
- [这对您意味着什么？](#)
- [较新的版本会失去功能DLAMIs吗？](#)
- [这一变化是否影响了 Deep Learning Containers？](#)

更改了哪些内容？

我们分DLAMIs成两个独立的小组：

- DLAMIs使用NVIDIA专有驱动程序（支持 P3、p3dn、G3）
- DLAMIs使用NVIDIAOSS驱动程序（支持 g4dn、G5、P4、P5）

因此，我们用新名称和新名称分别DLAMIs为这两个类别创建了新名称AMIIDs。DLAMIs它们不可互换。也就是说，DLAMIs来自一个组的实例不支持另一个组支持的实例。例如DLAMI，支持 P5 的不支持 G3，支持 G3 DLAMI 的则不支持 P5。



为什么需要进行此更改？

以前，DLAMIs 包含来自 NVIDIA GPU 的专有内核驱动程序。但是，上游 Linux 内核社区接受了一项变更，该变更将专有内核驱动程序（例如 NVIDIA GPU 驱动程序）与其他内核驱动程序的通信隔离开来。此更改在 P4 和 P5 系列实例上禁用 GPU Direct RDMA，该机制允许 GPU 有效地用于分布式训练。因此，DLAMIs 现在使用 OpenRM 驱动程序（NVIDIA 开源驱动程序），该驱动程序与开源 EFA 驱动程序链接以支持 G4dn、G5、P4 和 P5。但是，此 OpenRM 驱动程序不支持较旧的实例（例如 P3 和 G3）。因此，为了确保我们继续提供支持这两种实例类型的最新、高性能和安全性 DLAMIs，我们将 DLAMIs 分为两组：一组使用 OpenRM 驱动程序（支持 G4dn、G5、P4 和 P5），另一组使用较旧的专有驱动程序（支持 P3、p3dn 和 G3）。

这一变化影响 DLAMIs 了哪些？

这一变化影响了所有人 DLAMIs。

这对您意味着什么？

只要您在支持的亚马逊弹性计算云 (AmazonEC2) 实例类型上运行，它们都DLAMIs将继续提供功能、性能和安全性。要确定DLAMI支持的EC2实例类型，请查看版本说明DLAMI，然后查找支持的EC2实例。有关当前支持的DLAMI选项列表及其发行说明的链接，请参阅[DLAMIs 发行说明](#)。

此外，您必须使用正确的 AWS Command Line Interface (AWS CLI) 命令来调用当前DLAMIs。

对于支持 P3、p3dn 和 G3 的基础DLAMIs，请使用以下命令：

```
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
```

对于支持 g4dN、G5、P4 和 P5 的基础DLAMIs，请使用以下命令：

```
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
```

较新的版本会失去功能DLAMIs吗？

否，不损失任何功能。当前版本DLAMIs提供前一个版本的所有功能、性能和安全性DLAMIs，前提是你支持的EC2实例类型上运行它们。

这一变化是否影响了 Deep Learning Containers？

不，此更改并未影响 Dee AWS p Learning Containers，因为它们不包括NVIDIA驱动程序。但是，请务必在与底层实例兼容的深度容器上AMIs运行 Deep Learning Containers。

相关的信息 DLAMI

您可以在《AWS Deep Learning AMIs 开发者指南》DLAMI之外找到包含相关信息的其他资源。在 AWS re:Post，查看其他客户提出的问题或自己提问。DLAMI在 Mach AWS ine Learning AWS 博客和其他博客上，阅读相关的官方帖子DLAMI。

AWS re:Post

[标签：AWS Deep Learning AMIs](#)

AWS 博客

- [AWS Machine Learning 博客 | 分类：AWS Deep Learning AMIs](#)
- [AWS Machine Learning 博客 | 在 Amazon EC2 C5 和 P3 实例上优化了 TensorFlow 1.6，训练速度更快](#)
- [AWS Machine Learning 博客 | AWS Deep Learning AMIs 面向机器学习从业者的新内容](#)
- [AWS Partner Network \(APN\) 博客 | 推出新的培训课程：Machine Learning 和深度学习简介 AWS](#)
- [AWS 新闻博客 | 深度学习之旅 AWS](#)

的已弃用功能 DLAMI

下表列出了 AWS Deep Learning AMIs (DLAMI) 的已弃用功能、我们弃用它们的日期以及有关我们为何弃用这些功能的详细信息。

特征	Date	详细信息
Ubuntu 16.04	10/07/2021	Ubuntu Linux 16.04的五年期LTS限LTS已于2021年4月30日结束，不再得到其供应商的支持。截至2021年10月，新版本中不再有深度学习基础AMI (Ubuntu 16.04) 的更新。先前的版本将继续可用。
Amazon Linux	10/07/2021	截至 2020 年 12 月 end-of-life ，亚马逊 Linux 已上市。截至2021年10月，新版本中不再有深度学习AMI (Amazon Linux) 的更新。先前版本的深度学习AMI (Amazon Linux) 将继续推出。
Chainer	2020 年 7 月 1 日	Chainer 宣布自 2019 年 12 月起 停用主要版本 。因此，从 2020 年 7 月起，我们将不再在 DLAMI 中包含 Chainer Conda 环境。包含这些环境的先前版本的 DLAMI 将继续可用。仅在开源社区针对这些框架发布安全修补程序时，我们才会为这些环境提供更新。

特征	Date	详细信息
Python 3.6	2020 年 6 月 15 日	由于客户的要求，我们正在迁移到 Python 3.7 来TF/MX/PT发布新版本。
Python 2	2020 年 1 月 1 日	<p>Python 开源社区已正式结束对 Python 2 的支持。</p> <p>TensorFlow PyTorch、和MXNet社区还宣布，TensorFlow 1.15、TensorFlow 2.1、1. PyTorch 4 和 MXNet 1.6.0 版本将是最后一个支持 Python 2 的版本。</p>

DLAMI 的文档历史记录

下表提供了《AWS Deep Learning AMIs 开发者指南》的最新DLAMI版本和相关变更的历史记录。

最近更改

变更	说明	日期
ARM64 DLAMI	AWS Deep Learning AMIs 现在支持基于 Arm64 处理器的 GPUs 图像。	2021 年 11 月 29 日
TensorFlow 2	Conda AMI 的深度学习现在支持 TensorFlow 2 和 CUDA 10。	2019 年 12 月 3 日
AWS 推论	深度学习 AMI 现在支持 AWS Inferentia 硬件和神经元。 AWS SDK	2019 年 12 月 3 日
将 TensorFlow 服务与初始模型配合使用	为 TensorFlow 服务添加了在盗梦空间模型中使用推理的示例，无论有没有使用 Elastic Inference，都可以使用 Elastic Inference。	2018 年 11 月 28 日
PyTorch 从夜间版本中安装	添加了一个教程，其中介绍了如何卸载 PyTorch，然后使用 Conda 在 Deep Learning PyTorch AMI 上安装每晚版本的。	2018 年 9 月 25 日
Conda 教程	该示例 MOTD 已更新，以反映最新版本。	2018 年 7 月 23 日

之前的更改

下表提供了 2018 年 7 月之前的早期 DLAMI 版本和相关变更的历史记录。

更改	描述	日期
TensorFlow 和 Horovod	添加了 ImageNet 使用 TensorFlow 和 Horovod 进行训练的教程。	2018 年 6 月 6 日
升级指南	添加了升级指南。	2018 年 5 月 15 日
新区域和新的 10 分钟教程	添加的新区域：美国西部 (加利福尼亚北部)、南美洲、加拿大 (中部)、欧洲 (伦敦) 和欧洲 (巴黎)。此外，还发布了名为“深度学习入门AMI”的 10 分钟教程。	2018 年 4 月 26 日
Chainer 教程	添加了在多模式GPU、单GPU/CPU模式和模式下使用 Chainer 的教程。CUDA多个框架的集成已从 CUDA 8 升级到 CUDA 9。	2018 年 2 月 28 日
Linux AMIs v3.0，以及MXNet 模型服务器、Servin TensorFlow 和 TensorBoard	添加了 Conda 的教程，AMIs 其中包含使用模型服务器 v0.1.5、Serving v1.4.0 和 v0.4.0 的新MXNet模型和可视化 TensorFlow 服务功能。TensorBoard AMI以及 Conda 和CUDA概述中描述的框架 CUDA功能。最新发行说明迁移到 https://aws.amazon.com/releases/notes/	2018 年 1 月 25 日
Linux AMIs v2.0	Base、Source 和 Conda AMIs 更新为 NCCL 2.1。Source 和 Conda AMIs 更新了 MXNet v1.0、PyTorch 0.3.0 和 Keras 2.0.9。	2017 年 12 月 11 日

更改	描述	日期
增加了两个 Windows AMI 选项	Windows 2012 R2 和 2016 已 AMIs发布：已添加到AMI选择指南中并添加到发行说明中。	2017 年 11 月 30 日
初始文档版本	更改的详细说明，带有指向所更改主题/章节的链接。	2017 年 11 月 15 日

本文属于机器翻译版本。若本译文内容与英语原文存在差异，则一律以英文原文为准。