

NVIDIA® JETSON AGX ORIN™ DEVELOPER KIT

Next-Level AI Performance for Next-Gen Robotics



Reviewer's Guide

Introducing NVIDIA® Jetson AGX Orin™	3
Best in Class Performance	6
Up to 8X Higher AI Performance	6
Significant Performance Speedups for Real-world AI Workloads	6
Performance on Vision AI and Conversational AI Models	8
Faster Development with the NVIDIA AI Software Platform	8
Appendix	12
Getting Started with the Jetson AGX Orin Developer Kit	12
Booting up your Jetson AGX Orin Developer Kit	12
Exploring the Developer Kit	12
Running Inference Benchmarks	14
NVIDIA TAO	16
NVIDIA RIVA	17
Additional Resources	19
Pytorch Wheels	19
Pytorch Container	19
NVIDIA Jetson AGX Orin Specs	20
NVIDIA Contact Information	21
NVIDIA North/Latin America Public Relations	21
NVIDIA Europe Public Relations	22
NVIDIA Asia/Pacific Public Relations	23

Introducing NVIDIA® Jetson AGX Orin™

NVIDIA Accelerated Computing has revolutionized the growth of the AI industry with high performance GPUs, SoCs, and optimized software stacks that are widely used across industries spanning cloud, data centers, and edge computing. NVIDIA Jetson platform was specifically created to deliver the benefits of AI computing to embedded edge devices.

The Jetson AGX Xavier launched in 2018 raised the bar on AI compute performance by delivering up to 32 TOPS of AI performance and enabling new AI-powered edge applications. NVIDIA Jetson AGX Xavier, combined with NVIDIA's comprehensive, optimized, and scalable AI software stack, are powering many commercial autonomous robots used in food delivery, crop harvesting, warehouse goods movement, robot assisted surgery, autonomous industrial inspections, and many more. The strong performance of Jetson AGX Xavier coupled with the powerful NVIDIA AI software stack has resulted in the rapid growth of the Jetson developer community to over one million registered developers and has enabled over 6000 customers to develop and deploy NVIDIA Jetson-based AI solutions.

The pace of development and deployment of AI-powered autonomous machines and robots continues to grow rapidly, and the next generation of applications require tremendous AI compute performance to handle multi-modal AI applications that need to run concurrently in real-time. As human-robot interactions increase in retail spaces, food delivery, hospitals, warehouses, factory floors and other commercial applications, autonomous robots will need to concurrently perform 3D perception, natural language understanding, path planning, obstacle avoidance, pose estimation, and many more autonomous actions that not just require significant AI performance, but also highly accurate and trained neural models for each application.

The NVIDIA Jetson AGX Orin is the highest-performing and newest member of the NVIDIA Jetson family. It delivers tremendous performance, class-leading energy efficiency, and is backed by the comprehensive NVIDIA AI software stack to power the next generation of demanding edge AI applications. The NVIDIA® Orin System-on-Chip (SoC) based on the NVIDIA Ampere GPU architecture with **2048 CUDA cores, 64 Tensor Cores, and 2 Deep Learning Accelerator (DLA) engines** delivers up to **275 TOPS** of raw AI performance. NVIDIA Jetson AGX Orin delivers the following key benefits:

- Up to 8X the raw AI compute performance and up to 2X the energy efficiency of Jetson AGX Xavier
- Same pin-out and footprint as the Jetson AGX Xavier, enabling customers to easily upgrade existing designs to Jetson AGX Orin
- Industry-leading AI performance enabling significant speedups for real world computer vision and conversational AI workloads
- Powerful NVIDIA AI software stack with support for SDKs such as NVIDIA JetPack, NVIDIA RIVA, NVIDIA DeepStream, NVIDIA Isaac, NVIDIA TAO, and others.

The Jetson AGX Orin Developer Kit contains everything needed for developers to get up and running quickly. The Jetson AGX Orin Developer Kit is priced at **\$1999** and is available for purchase through NVIDIA authorized distributors worldwide.

Developer kit contents

- Jetson AGX Orin module with heat sink and reference carrier board
- 802.11ac/abgn Wireless Network Interface Controller
- Power adapter and USB Type-C cord
- USB Type-C to USB Type-A cord
- Quick Start and Support Guide



Figure 1 Jetson AGX Orin Developer Kit

Jetson AGX Orin Developer Kit features

MODULE:	
GPU	NVIDIA Ampere architecture with 2048 NVIDIA® CUDA® cores and 64 Tensor cores
CPU	12-core Arm Cortex-A78AE v8.2 64-bit CPU 3MB L2 + 6MB L3
DL Accelerator	2x NVDLA v2.0
Vision Accelerator	PVA v2.0
Memory	32GB 256-bit LPDDR5 204.8 GB/s
Storage	64GB eMMC 5.1
Power	15W to 60W

Table 1 Jetson AGX Orin Developer Kit Module Specs

REFERENCE CARRIER BOARD:	
Camera	16 lane MIPI CSI-2 connector
PCIe	x16 PCIe slot supporting x8 PCIe Gen4
M.2 Key M	x4 PCIe Gen 4
M.2 Key E	x1 PCIe Gen 4, USB 2.0, UART, I2S
USB	Type C: 2x USB 3.2 Gen2 with USB-PD support Type A: 2x USB 3.2 Gen2, 2x USB 3.2 Gen1 Micro-B: USB 2.0
Networking	RJ45 (up to 10 GbE)
Display	DisplayPort 1.4a (+MST)
microSD slot	UHS-1 cards up to SDR104 mode
Others	40-pin header (I2C, GPIO, SPI, CAN, I2S, UART, DMIC) 12-pin automation header 10-pin audio panel header 10-pin JTAG header 4-pin fan header 2-pin RTC battery backup connector DC power jack Power, Force Recovery, and Reset buttons
Dimensions	110mm x 110mm x 71.65mm (Height includes feet, carrier board, module, and thermal solution)

Table 2 Jetson AGX Orin Developer Kit Carrier Board Specs

Best in Class Performance

Up to 8X Higher AI Performance

The power-efficient Jetson AGX Orin System-on-Module (SoM) delivers up to 275 TOPS¹ of AI performance within a 60-Watt power budget, an 8X improvement over the 32 TOPS delivered by Jetson AGX Xavier. For designs requiring lower power profiles, customers can tune their designs for power profiles ranging from 15W to 60W. Jetson AGX Orin also delivers up to 9X the DLA performance of Xavier providing even higher energy efficiency for inferencing applications that purely run on the DLA cores. Jetson AGX Orin also delivers up to 1.5X higher CPU performance and up to 1.5X higher DRAM bandwidth that helps to reduce bottlenecks and latencies when running multiple concurrent inferencing applications.

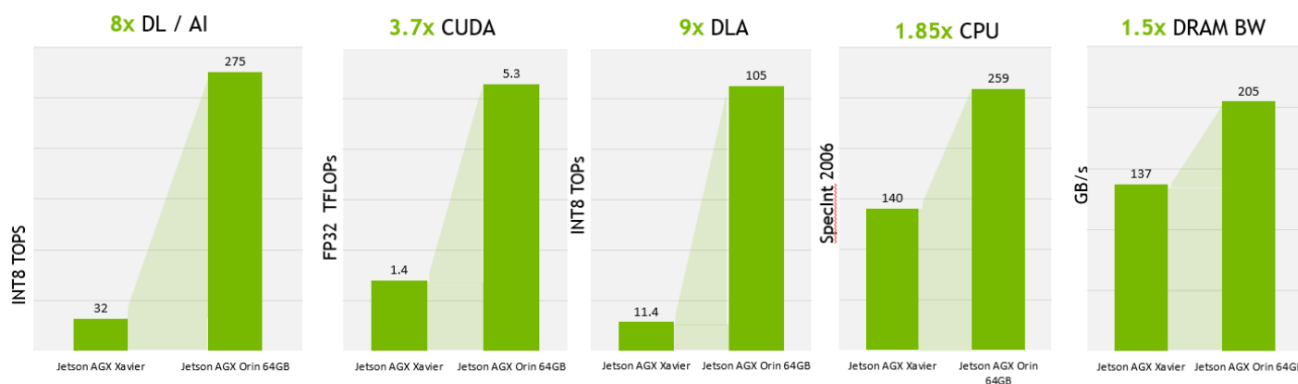


Figure 2 Jetson AGX Orin delivers 8x the AI performance of Jetson AGX Xavier

Significant Performance Speedups for Real-world AI Workloads

NVIDIA Jetson AGX Orin delivers significant performance speedups for various real-world AI workloads as measured by industry accepted benchmarks and performance on widely used neural networks.

The industry standard MLPerf benchmark defined by a consortium of industry leaders including NVIDIA, Google, Meta and others measures performance on common AI tasks for image classification, object detection and natural language processing using widely used neural networks such as ResNet-50, Mobilenet-v2, and GNMT. NVIDIA Jetson AGX Xavier is the currently at the top of the leaderboard for this benchmark, beating the competition by significant margins. Read the blog [here](#)

NVIDIA Jetson AGX Orin will further increase this lead and even beat the current leader, Jetson AGX Xavier by a significant margin. The results of this benchmark will be revealed on April 6th after it is reviewed and confirmed by the MLPerf benchmark organization.

¹ INT8 TOPS

Jetson AGX Orin delivers almost 3.5X the performance of Jetson AGX Xavier on popular pre-trained neural networks that are used for object detection, action recognition, pose estimation, and others. The new Sparsity feature introduced in the Ampere GPU architecture and available on Jetson AGX Orin, will help further boost the performance advantage. The continual optimizations in the underlying JetPack software, neural networks, and the adoption of Sparse networks are expected to further improve the performance lead of Jetson AGX Orin in real world workloads to up to 5X. Learn more about accelerating inferencing with sparsity [here](#).

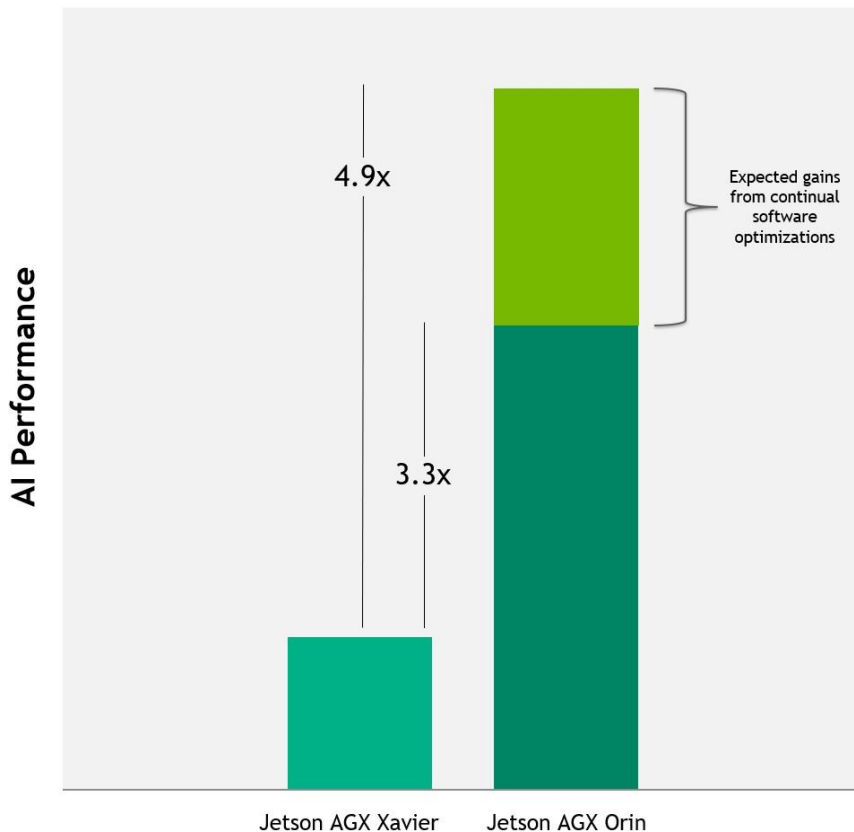


Figure 3 Jetson AGX Orin delivers more than 3x the AI performance of Jetson AGX Xavier²

² Relative performance gain represents the geometric mean of performance gains measured across a wide variety of production-ready pre-trained neural networks and inference models used in MLPerf.

Performance on Vision AI and Conversational AI Models










	Jetson AGX Xavier (Inferences/sec)	Jetson AGX Orin (Inferences/sec)
 PeopleNet	196	536
 Action Recognition 2D	471	1577
 Action Recognition 3D	32	105
 LPR	1190	4118
 Dashcam Net	671	1908
 Bodypose Net	172	559
 ASR: Citrinet 1024	34	113
 NLP: BERT-base	94	186
 TTS: Fastpitch-HifiGAN	9	32

Table 3 Jetson AGX Orin Performance on widely used Vision and Conversational AI models

Instructions for running the above benchmarks on Jetson AGX Orin and Jetson AGX Xavier are available in [Running Inference Benchmarks](#) section of Appendix

Faster Development with the NVIDIA AI Software Platform

The class-leading performance and energy efficiency of Jetson AGX Orin is backed by the same powerful NVIDIA AI platform that is deployed in GPU-accelerated data centers, hyperscale servers, and powerful AI workstations. The NVIDIA AI platform brings optimized AI tools, libraries, and NVIDIA SDKs such as CUDA, CuDNN, TensorRT, DeepStream, RIVA, TAO and Isaac to the Jetson platform, enabling developers to seamlessly train AI applications on powerful cloud GPUs and deploy trained networks on Jetson-powered AI edge devices.

Tools like **NVIDIA Omniverse Replicator** for synthetic data generation help in creating high quality datasets to boost model training, **NVIDIA Train-Adapt-Optimize (TAO)** and **Pre-Trained Models (PTM)** cuts down development time by up to 10X and provides an easier and faster way to accelerate

training and quickly create highly accurate production ready models. SDKs that are deployed in powerful datacenters, such as RIVA for accelerated conversational AI and Deepstream for accelerated vision AI, also run on the Jetson AGX Orin. NVIDIA Isaac SDK delivers the GPU-accelerated solutions for Jetson-based robotic applications. Developers can quickly and easily build advanced robots such as the one showcased in the [Orion demo video](#) using the software tools available on the NVIDIA AI platform.

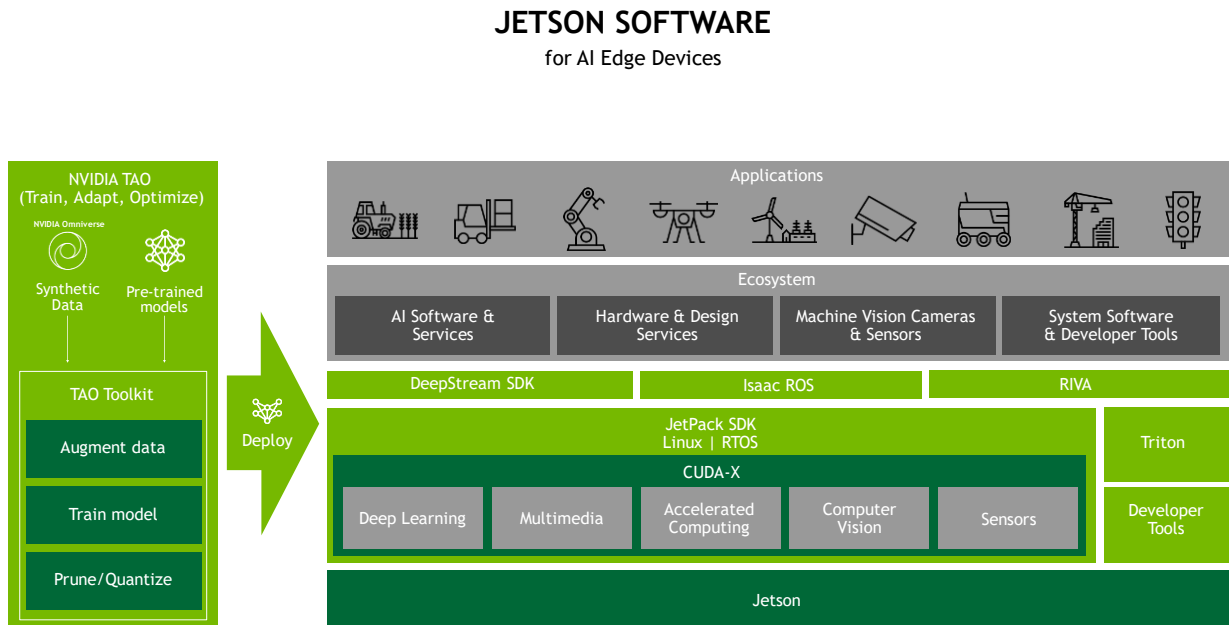


Figure 4 NVIDIA AI Software Stack for Jetson

NVIDIA JetPack is the foundational SDK for the Jetson edge AI platform. JetPack SDK provides a full development environment for hardware-accelerated AI-at-the-edge development. JetPack SDK provides a board support package, with bootloader, Linux kernel, Ubuntu desktop environment, and a complete set of libraries for acceleration of GPU computing, multimedia, graphics, and computer vision.

Jetson AGX Orin is powered by our latest release of JetPack 5.0 whose highlights include:

- Latest compute stack with latest versions of CUDA 11 and TensorRT 8
- Linux Kernel 5.10
- Ubuntu 20.04 based root file system
- UEFI for CPU bootloader
- OP-TEE for Trusted Execution Environment
- Hardware root of trust, Secureboot, Disk Encryption, Secure Storage and other security features
- Over-the-Air Updates to safely update any Jetson module deployed in the field

Please refer to [Booting up your Jetson AGX Orin Developer Kit](#) section in Appendix to set up your development kit with JetPack 5.0 SDK.

[NVIDIA TAO \(Train-Adapt-Optimize\)](#) is a framework that lets developers create custom, production-ready models, in hours rather than months, without AI expertise or large training datasets. The NVIDIA TAO Toolkit abstracts away the AI/deep learning framework complexity, letting you fine-tune on high-quality NVIDIA pre-trained AI models with only a fraction of the data compared to training from scratch. Customers can use the TAO Toolkit to fine-tune and optimize on a wide variety of use cases from computer vision and automatic speech recognition to speech synthesis and natural language understanding.

Data collection and annotation is an expensive and laborious process. Simulation can help bridge the need for data. NVIDIA **Omniverse Replicator** uses simulation to generate synthetic data that is an order of magnitude faster and cheaper to create than real data. With Omniverse Replicator you can quickly create diverse, massive and accurate datasets for training AI models.

We have provided you with a demo that will enable you to use TAO to train a model in the cloud and deploy the trained model on Jetson using DeepStream. Please refer to [NVIDIA TAO](#) section in Appendix for instructions on running the demo.

[NVIDIA Isaac ROS GEMs](#) are hardware-accelerated packages that make it easier for ROS developers to build high-performance solutions on NVIDIA hardware. [NVIDIA Isaac Sim](#), powered by Omniverse, is a scalable robotics simulation application. It includes Replicator - a tool to generate diverse synthetic datasets for training perception models. Isaac Sim is also a tool that powers photorealistic, physically accurate virtual environments to develop, test, and manage AI-based robots.

[NVIDIA RIVA](#) is an SDK for building GPU-accelerated conversational AI applications. RIVA includes **state of the art pre-trained models** for Automatic Speech Recognition (ASR) and Text-To-Speech (TTS). These pre-trained models are highly accurate and can be **easily customized** using the TAO Toolkit to improve accuracy on desired domains, accents, languages and use cases. NVIDIA RIVA speech models are optimized for TensorRT to **deliver high inferencing performance and low latencies** on Jetson AGX Orin.

We have provided you with a RIVA ASR demo which is a dictation application that showcases the performance of Jetson AGX Orin and the accuracy of the pre-trained speech recognition neural networks. Please refer to [NVIDIA RIVA](#) section in Appendix for more details and instructions to run these demos.

[DeepStream](#) is an SDK for rapidly developing and deploying Vision AI applications and services. DeepStream offers hardware acceleration beyond inference as it offers **hardware accelerated plugins** for end-to-end AI pipeline acceleration. It offers state-of-the-art throughput. Developers can also bring their own TensorFlow, PyTorch, or ONNX models and deploy them using DeepStream.

[NVIDIA Fleet Command](#)™ is a cloud service that securely deploys, manages, and scales AI applications across distributed edge infrastructure. Purpose-built for AI, Fleet Command is a **turnkey solution for AI lifecycle management**. It removes the complexity of building and maintaining an edge software platform by offering streamlined deployments, OTA updates, and detailed monitoring capabilities. Layered security protocols protect intellectual property and application insights from cloud to edge. With Fleet Command, organizations can go from **zero to AI in minutes**. Fleet Command on Jetson is coming soon.

All the above software technologies combined with performance of Jetson AGX Orin will enable developers to build multi-modal applications that run multiple concurrent neural networks for use cases such as personal assistant robots as demonstrated in this [video](#). The video showcases the Orion robot built by NVIDIA on the Jetson AGX Orin platform using the neural networks and SDKs that are in the sample demos provided to you.

Appendix

Getting Started with the Jetson AGX Orin Developer Kit

You can access this Reviewers Guide from <https://developer.nvidia.com/jetson-agx-orin-review>

Booting up your Jetson AGX Orin Developer Kit

The Jetson AGX Orin Developer Kit comes pre-flashed with Jetson Linux BSP. Just connect the power supply, keyboard, mouse, display and power on. Also make sure to connect Jetson to the internet over the built-in Wi-Fi or ethernet before you proceed further.

The first boot on this “private” preview pre-flashed image will take couple of minutes to show the initial configuration dialogue. The developer kits made which will be available for public after announcement in GTC 2022 will have a release quality pre-flashed BSP image with a much faster first boot.

Go through the simple initial setup process and skip installing chromium web browser when asked. There is a known issue in this private preview image.

Once the initial configuration is complete and developer kit has booted to desktop, you can install the latest BSP and JetPack components using the below command:

```
sudo add-apt-repository "deb https://repo.download.nvidia.com/jetson/jetson-50/t234 r34.0 main"
sudo add-apt-repository "deb https://repo.download.nvidia.com/jetson/jetson-50/common r34.0
main"
sudo apt dist-upgrade
```

During upgrade, the script will ask your choice few times, please enter “Y” for all choices.

Please do a hard reboot of the system by long pressing the power button OR by disconnecting the power supply and connecting it back again. After the system is rebooted and you have logged into the desktop, please install Jetpack components using:

```
sudo apt install nvidia-jetpack
```

After installing JetPack components, reboot the system and boot into the latest JetPack.

Exploring the Developer Kit

The developer kit is running JetPack 5.0 (pre-release) which has following components:

- CUDA 11.4
- cuDNN 8.3.2
- TensorRT 8.4.0
- OpenCV 4.5.4
- Vulkan 1.3
- VPI 2.0
- Nsight Systems 2021.5
- Nsight Graphics 2021.5

On top right of the desktop, there is a power profile selector. When clicked, it provides a drop-down menu of all software defined power modes for the developer kit. For the purposes of reviewing the developer kit, running benchmarks and the provided sample demos, please confirm that Jetson is in the MaxN power profile

JetPack comes with [multiple samples built in](#). These samples provide a preview and capabilities of different JetPack components. JetPack includes the following samples which can be compiled on the developer kit. We encourage you to run these samples to learn about the JetPack components.

JetPack component	Sample locations on reference filesystem
TensorRT	<code>/usr/src/tensorrt/samples/</code>
cuDNN	<code>/usr/src/cudnn_samples_v8/</code>
CUDA	<code>/usr/local/cuda-11.4/samples/</code>
MM API	<code>/usr/src/jetson_multimedia_api</code>
VPI	<code>/opt/nvidia/vpi2/samples/</code>

Running Inference Benchmarks

Download the tar ball named "benchmarks.tar.gz" required for this benchmarking from [here](#) and move it inside the Review Directory. And untar the tar ball:

```
mkdir $HOME/Review
```

```
cd $HOME/Review
```

Move benchmarks.tar.gz to inside the \$HOME/Review and then untar it:

```
tar -xvf benchmarks.tar.gz
```

Setup the requirements for running the benchmark by doing:

```
cd benchmarks
```

```
bash installBenchmarks.sh
```

The above script will download additional content from the internet. Please make sure that the devkit is online when you run the script

For a clean measurement, please reboot the system and then start benchmarking.

Vision Model Benchmarking

Run Vision Model Benchmarks using:

```
bash launchVisionBenchmark.sh
```

You will get the output similar to below:

Model	Name	FPS
0	peoplenet	531.319971
1	action_recog_2d	1576.913747
2	action_recog_3d	108.784537
3	dashcamnet	1905.534769
4	bodyposenet	560.921805
5	lpr_us	4118.400131

Conversational AI Model Benchmarking

Run Vision Model Benchmarks using:

```
bash launchConvBenchmark.sh
```

You will get the output similar to below. Highlighted in yellow are the latencies of ASR, NLP and TTS. You can convert that latency to FPS using 1000/latency.

-----Starting ASR Benchmark

Loading eval dataset...

filename: /benchdata/1272-135031-0000.wav

Done loading 1 files

Latencies (ms):

Median	90th	95th	99th	Avg
9.2084	9.3835	9.5279	14.065	9.4069

Intermediate latencies (ms):

Median	90th	95th	99th	Avg
9.2061	9.3615	9.4755	9.6975	9.3387

Final latencies (ms):

Median	90th	95th	99th	Avg
14.065	14.131	14.131	14.131	14.041

Run time: 54.441 sec.

Total audio processed: 54.425 sec.

Throughput: 0.99971 RTFX

-----Starting NLP Benchmark

Done sending 1000 requests

Done processing 1000 responses

Run time: 6.80158s

Total sequences processed: 1000

Throughput: 147.025 seq/sec

Latencies:	Median	90	95	99	Avg
	3.39	3.49	3.51	3.57	3.45

-----Starting TTS Benchmark

Latencies:

First audio - average: 0.024089 (This is in seconds. Convert to milliseconds using 1000/latency in seconds)

First audio - P90: 0.0248866

First audio - P95: 0.0249044

First audio - P99: 0.0249245

Chunk - average: 0.00473509

Chunk - P90: 0.00488245

Chunk - P95: 0.00489093

Chunk - P99: 0.0058794

Throughput (RTF): 50.8159

At the end stop benchmarking using:

```
bash stopBenchmark.sh
```

NVIDIA TAO

NVIDIA provides pre-trained, production-ready neural network models on the NVIDIA GPU Cloud (NGC). For each model, NVIDIA provides a deployable model and a trainable model. The deployable model is production ready and optimized to run on inferencing pipeline while the trainable model is intended for fine tuning using custom data with NVIDIA TAO.

In this demo, you will experience:

- Deploying production ready high accuracy, optimized models on Jetson
- Using Train-Adapt-Optimize to train a trainable model on cloud and then deploy the trained model on Jetson.

To setup the demo, Create a directory and name it as "Review" in your home directory if already not created.

```
mkdir $HOME/Review
```

Download the tar ball named "TAO-PTM.tar.gz" required for this demo from [here](#) and move it inside the Review Directory. And untar the tar ball:

```
cd $HOME/Review
```

Move "TAO-PTM.tar.gz" to inside the \$HOME/Review and then untar it:

```
tar -xvf TAO-PTM.tar.gz
```

Install the demo using:


```
cd $HOME/Review/TAO-PTM
```

```
bash installTAOPTM.sh
```

Installation will take couple of minutes.

After installation is complete, start the jupyter-lab by following the commands below:

```
cd $HOME/Review/TAO-PTM/
```

```
export PATH="$HOME/.local/bin/:$PATH"
```

```
jupyter-lab
```

Click on the link in the end of the output of the above command and it will open Jupyter notebook on the browser.

Deploy Pre-Trained Models on Jetson

In this section, you will deploy the Peoplenet deployable model using DeepStream and see it in action. In the Jupyter notebook on the browser, navigate (navigation pane on the left) to *PreTrainedModel* directory and open the notebook named *peoplenet_workflow.ipynb*. You can read about the PeopleNet model: the model architecture, training data and accuracy in the notebook.

Train using NVIDIA TAO and Deploy on Jetson

In this section, you will start from a trainable PeopleNet model and train in the cloud for an additional detection class. After training, you will download the trained model to Jetson and deploy using DeepStream. In the Jupyter notebook on the browser, navigate (navigation pane on the left) to *TrainAdaptOptimize* directory and open the notebook named *tao_workflow.ipynb*

For your review, we have provided minimal set of data for training and training for 500 epochs. The quality of detections benefits from the amount and quality of dataset. Since we have kept the data set minimal in order to shorten the training time, the model generated will not be of a production quality.

NVIDIA RIVA

In this demo, you will experience NVIDIA RIVA ASR in action. You will need a headset with mic which we have shipped to you in the reviewers package. Please connect the USB headset with mic to the developer kit before starting on the demos below.

Create a directory and name it as "Review" in your home directory if already not done.

```
mkdir $HOME/Review
```

Download the tar ball named "Riva.tar.gz" required for this demo from [here](#) and move it inside the Review Directory.

```
cd $HOME/Review
```

Move "Riva.tar.gz" to inside the \$HOME/Review and then untar it:

```
tar -xvf Riva.tar.gz
```

Install the demo using:

```
cd $HOME/Review/Riva
```

```
bash installRivaASR.sh
```

Once the installation is done, you can experience the following RIVA transcription demo:

Please remove the USB Camera and only have the USB Headset with mic connected. And then reboot. We require this step since we have hardcoded the demo to use the USB Headset with mic we have provided.

RIVA ASR Demo

First demo will showcase RIVA Automatic Speech Recognition (ASR) with a transcription use case. The demo uses Citrinet model architecture which is trained with ASR Set 3.0 - 16700 Hours dataset. The model is running totally locally on Jetson in this demo.

Launch the demo using below command:

```
cd $HOME/Review/Riva/
```

```
bash launchASR.sh
```

It will take approximately 3 minutes to load. Once loaded you will see a clear screen and at this point you can start talking. You will see the transcription happening in real time.

We encourage you to put our ASR to test by giving dictation to both RIVA ASR and a competing solution in the market and compare the accuracy and performance.

After you are done, press ctrl-c to stop the transcription, and then execute the following command.

```
bash stopRiva.sh
```

Additional Resources

Pytorch Wheels

We have hosted Pytorch wheels here:

<https://nvidia.box.com/shared/static/19je2l0ppy1fpq4mw1a5gsbb5y9fopy7.whl>

For installation, please do the following steps:

```
wget
https://nvidia.box.com/shared/static/19je2l0ppy1fpq4mw1a5gsbb5y9fopy7.whl -O
torch-1.10.0-cp38-cp38-linux_aarch64.whl

sudo apt-get install python3-pip libopenblas-base libopenmpi-dev

pip3 install Cython

pip3 install numpy torch-1.10.0-cp38-cp38-linux_aarch64.whl
```

Pytorch Container

Please pull the pytorch container by following steps below:

Step 1: Log into Nvidia GPU Cloud

```
sudo docker login nvcr.io
```

User username as \$oauthtoken and password as

“NzRkbGFham5hYmFjZWdxcDdwY2c1b3VyaXQ6MmWQ3NmVkNjQtNGIyNS00MGFkLTlhYzEtYzNjMjkzZTc2OTZi” like shown below.

```
docker login nvcr.io

Username: $oauthtoken
Password:
NzRkbGFham5hYmFjZWdxcDdwY2c1b3VyaXQ6MmWQ3NmVkNjQtNGIyNS00MGFkLTlhYzEtYzNjMjkzZTc2OTZi
```

Step 2: Pull and start the container

```
sudo docker run -it --rm --runtime nvidia --network host nvcr.io/ea-
linux4tegra/14t-pytorch:r34.0.1-pt1.10-py3
```

NVIDIA Jetson AGX Orin Specs

	Jetson AGX Xavier	Jetson AGX Orin
AI Performance	32 TOPS (INT8)	275 TOPS (INT8 with Sparsity) 138 TOPS (INT8)
GPU	512-core NVIDIA Volta GPU with 64 Tensor Cores	2048-core NVIDIA Ampere GPU with 64 Tensor Cores
DL Accelerator	2x NVDLA	2x NVDLA v2
Vision Accelerator	2x PVA v1	PVA v2
CPU	8-core NVIDIA Carmel Arm®v8.2 64-bit CPU 8MB L2 + 4MB L3	12-core NVIDIA Arm® Cortex A78AE v8.2 64-bit CPU 3MB L2 + 6MB L3
Memory	Up to 64GB 256-bit LPDDR4x @ 2133MHz 137 GB/s	Up to 64 GB 256-bit LPDDR5 @ 3200MHz 204.8 GB/s
Storage	32GB eMMC 5.1	64 GB eMMC 5.1
Video Encode	4x 4K60 8x 4K30 16x 1080p60 32x 1080p30 (H.265) H.264, VP9	2x 4K60 4x 4K30 8x 1080p60 16x 1080p30 (H.265) H.264, AV1
Video Decode	2x 8K30 6x 4K60 12x 4K30 26x 1080p60 52x 1080p30 (H.265) H.264, VP9	1x 8K30 3x 4K60 7x 4K30 11x 1080p60 22x 1080p30 (H.265) H.264, VP9, AV1
Camera	16 lanes MIPI CSI-2 (36 Virtual Channels) 8 lanes SLVS-EC D-PHY 40Gbps / C-PHY 62 Gbps	16 lanes MIPI CSI-2 (16 Virtual Channels*) D-PHY 2.1 40Gbps / C-PHY 2.0 164Gbps
PCI Express	16 lanes PCIe Gen 4 1 x8 + 1 x4 + 1 x2 + 2 x1	22 lanes PCIe Gen 4 Up to 2 x8, 1 x4, 2 x1
Ethernet	1 Gbe RGMII	1 Gbe RGMII 4x 10Gbe XFI
Mechanical	100mm x 87mm 699 pin connector	100mm x 87mm 699 pin connector
Power	10W to 30W	15W to 60W

Table 4 NVIDIA Jetson AGX Orin compared to NVIDIA Jetson AGX Xavier

NVIDIA Contact Information

Any questions when reviewing the Developer Kit? Email JAOReviewersTeam@nvidia.com

NVIDIA North/Latin America Public Relations

<p>David Pinto PR Manager, Autonomous Machines Office: 408 566 6950 dpinto@nvidia.com</p>	<p>Sridhar Ramaswamy Director, Enterprise and SHIELD Technical Marketing Cell: 510 545 3774 sramaswamy@nvidia.com</p>
<p>Michael Lim Director, Analyst Relations Office: 408 486 2376 mlim@nvidia.com</p>	

NVIDIA Europe Public Relations

<p>Jens Neuschäfer Senior PR Manager, Enterprise Europe Office: +49 89 6283 50015 Mobile: +49 173 6282912 jneuschafer@nvidia.com NVIDIA GmbH Haus 1 West, 3rd Floor Flössergasse 2 81369 Munich, Germany</p>	<p>Rick Napier Senior Technical Product Manager Northern Europe Office: +44 (118) 9184378 Mobile: +44 (7917) 630172 rnapier@nvidia.com NVIDIA UK 100 Brook Drive Green Park Reading RG2 6UJ</p>
---	--

NVIDIA Asia/Pacific Public Relations

<p>Jeff Yen Director, Technical Marketing, APAC Office : +886 987 263 193 jyen@nvidia.com NVIDIA 8, Kee Hu Road, Neihu Taipei 114 TAIWAN</p>	<p>Melody Tu Senior PR / Marketing Manager, APAC Office: +886 2 6605 5856 metu@nvidia.com NVIDIA TASA (TW/AU/SEA) 8, Kee Hu Road, Neihu Taipei 114 TAIWAN</p>
<p>Searching Shi Sr. Technical Marketing Manager, China Office: +86 75586919016 Email: seshi@nvidia.com 5F BLOCK 8 VISEEN BUSINESSPARK 9 HIGH-TECH 9TH SOUTH ROAD SHENZHEN HI-TECH IND. PARK SHENZHEN, GUANGDONG Shenzhen 518057 China</p>	<p>Alex Liu PR/Marketing Manager, China Office: +86 1058661510 Email: alliu@nvidia.com Fortune Financial Center Level 40, Units 05-2,06 Building #5, Middle Road, East 3rd Ring Chaoyang District, Beijing, China 100000</p>
<p>Kyle Kim Sr. Technical Marketing Manager, Korea Office: +82 2 6001 7186 kylek@nvidia.com NVIDIA Korea #2101, COEX Trade Tower, 159-1 Samsung-dong, Kangnam-gu, Seoul 135-729 KOREA</p>	<p>Sunny Lee Marketing Director, Korea Office: +82 2 6001 7123 slee@nvidia.com NVIDIA Korea #2101, COEX Trade Tower, 159-1 Samsung-dong Kangnam-gu, Seoul 135-729 KOREA</p>
<p>Kaori Nakamura Head of Public Relations, Japan Office : +81 3 6743 8712 knakamura@nvidia.com ATT New Tower 13F 2-11-7 Akasaka, Minato-ku, Tokyo 107-0052 , JAPAN</p>	<p>Masaki Sawai Technical Marketing Manager, Japan Office: +81 3 6743 8717 Email: msawai@nvidia.com ATT New Tower 13F 2-11-7 Akasaka, Minato-ku, Tokyo 107-0052, JAPAN</p>

John Gillooly

Technical Marketing Manager, Asia Pacific South

Office: +65 8286 8727

Email: igillooly@nvidia.com

SINGAPORE

Titus Su

Technical Marketing Manager, TASA

Office: +886 2 6605 5430

Email: tisu@nvidia.com

8, Kee Hu Road, Neihu

Taipei 114, TAIWAN

Notice

ALL INFORMATION PROVIDED IN THIS REVIEWER'S GUIDE, INCLUDING COMMENTARY, OPINION, NVIDIA DESIGN SPECIFICATIONS, REFERENCE BOARDS, FILES, DRAWINGS, DIAGNOSTICS, LISTS, AND OTHER DOCUMENTS (TOGETHER AND SEPARATELY, "MATERIALS") ARE BEING PROVIDED "AS IS." NVIDIA MAKES NO WARRANTIES, EXPRESSED, IMPLIED, STATUTORY, OR OTHERWISE WITH RESPECT TO MATERIALS, AND EXPRESSLY DISCLAIMS ALL IMPLIED WARRANTIES OF NONINFRINGEMENT, MERCHANTABILITY, AND FITNESS FOR A PARTICULAR PURPOSE.

Information furnished is believed to be accurate and reliable. However, NVIDIA Corporation assumes no responsibility for the consequences of use of such information or for any infringement of patents or other rights of third parties that may result from its use. No license is granted by implication or otherwise under any patent or patent rights of NVIDIA Corporation. Specifications mentioned in this publication are subject to change without notice. This publication supersedes and replaces all information previously supplied. NVIDIA Corporation products are not authorized for use as critical components in life support devices or systems without express written approval of NVIDIA Corporation.

Trademarks

NVIDIA, the NVIDIA logo, GeForce, Tegra, Tesla and Jetson are trademarks and/or registered trademarks of NVIDIA Corporation in the U.S. and other countries. All rights reserved. Other company and product names may be trademarks of the respective companies with which they are associated.

Copyright

© 2022 NVIDIA Corporation. All rights reserved.