



The Vision Behind MLPerf

A Community-driven ML Benchmark Suite for Software Frameworks
and Hardware Accelerators in Cloud and Edge Computing

Prof. Vijay Janapa Reddi

MLPerf Inference Chair

ARM Research Summit

September 17th, 2019



HARVARD

**School of Engineering
and Applied Sciences**

MLPerf is the Work of Many

Founding leads: Peter Bailis (Stanford), Greg Diamos (Baidu), Peter Mattson (Google), David Patterson (UC Berkeley / Google), Gu-Yeon Wei (Harvard), Matei Zaharia (Stanford)

Training chairs: Victor Bittorf (Google), Paulius Micikevicius (NVIDIA), Andy Hock (Cerebras)

Inference chairs: Christine Cheng (Intel), David Kanter (RWI), Vijay Janapa Reddi (Harvard), Carole-Jean Wu (Facebook), Guenther Schmuelling (Microsoft), Hanlin Tang (Intel), Bing Yu (MediaTek)

Many others see mlperf.org/about

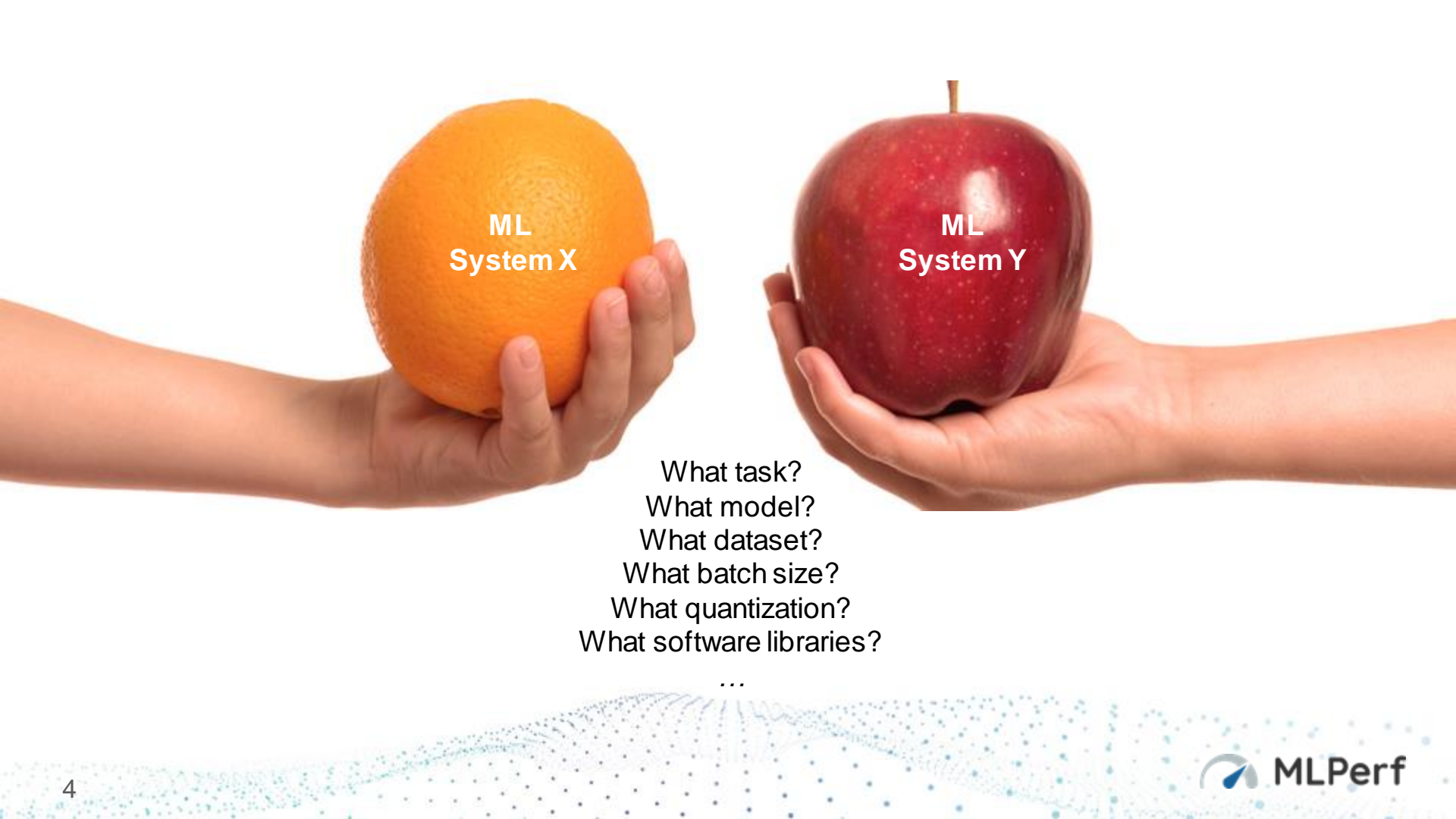
The New York Times

*Big Bets on A.I. Open a New
Frontier for Chip Start-Ups*

ML hardware is projected to be a
~\$60B industry in 2025.

(Tractica.com \$66.3B, Marketsandmarkets.com: \$59.2B)

than \$100 million from investors. Venture capitalists invested more than \$1.5 billion in chip start-ups last year, nearly doubling the investments made two years ago, according to the research firm CB Insights.



ML
System X

ML
System Y

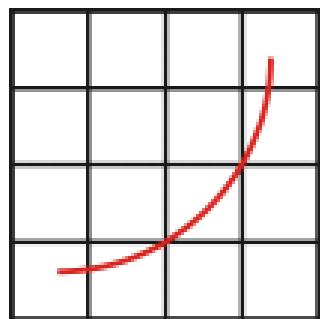
What task?
What model?
What dataset?
What batch size?
What quantization?
What software libraries?

...

Agenda

- ✓ *Why ML needs a benchmark suite?*
 - **Are there lessons we can borrow?**
 - What is MLPerf?
 - How does MLPerf curate a benchmark?
 - What is the “science” behind the curation?
 - Where are we heading now?
 - What comes next for MLPerf?

Yes! Look to successful history in benchmarks.



spec[®]

TPC[™]

SPEC Impact

- Settled **arguments in the marketplace** (grow the pie)
- Resolved internal **engineering debates** (better investments)
- Became a **standard** from research to practice in the industry
- Needed to **revise regularly** to maintain usefulness:
SPEC89, SPEC92, SPEC95, SPEC2000, SPEC2006, SPEC2017

SPEC fueled the Golden Age of microprocessor design.

Can we start a new **Golden Age** for ML Systems?

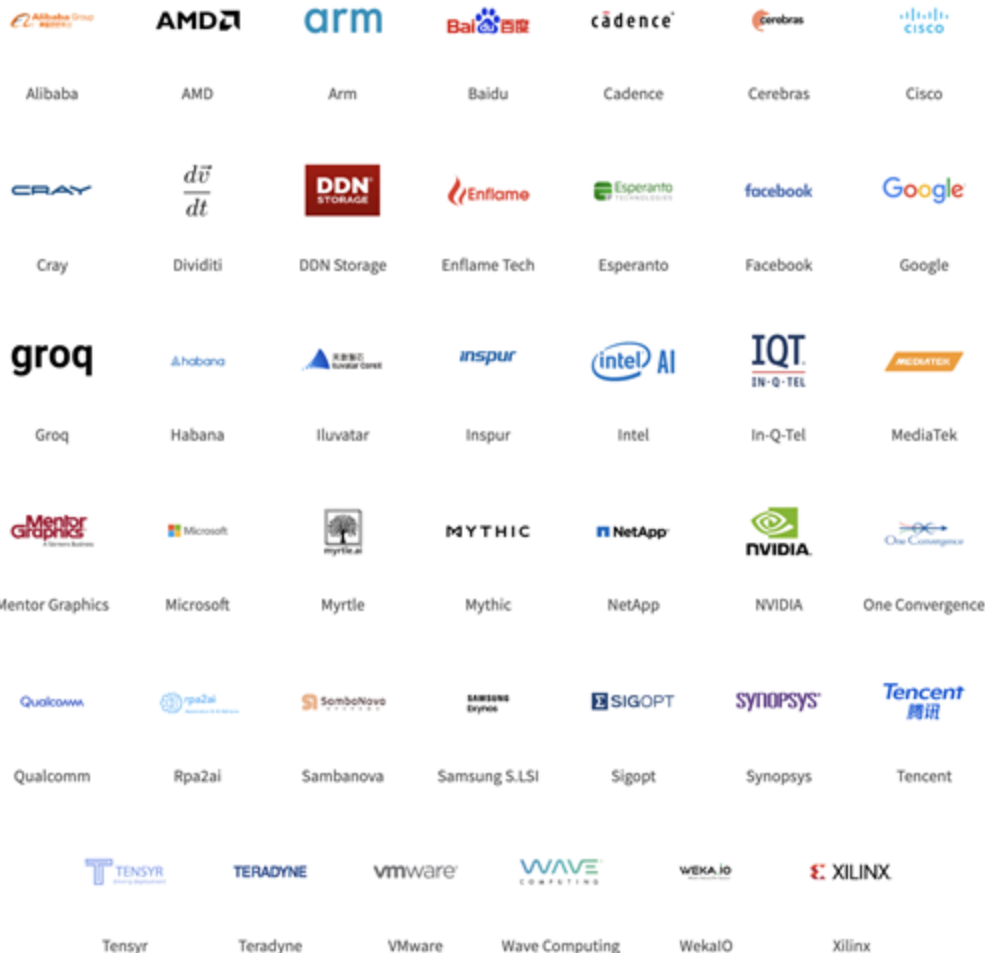
Agenda

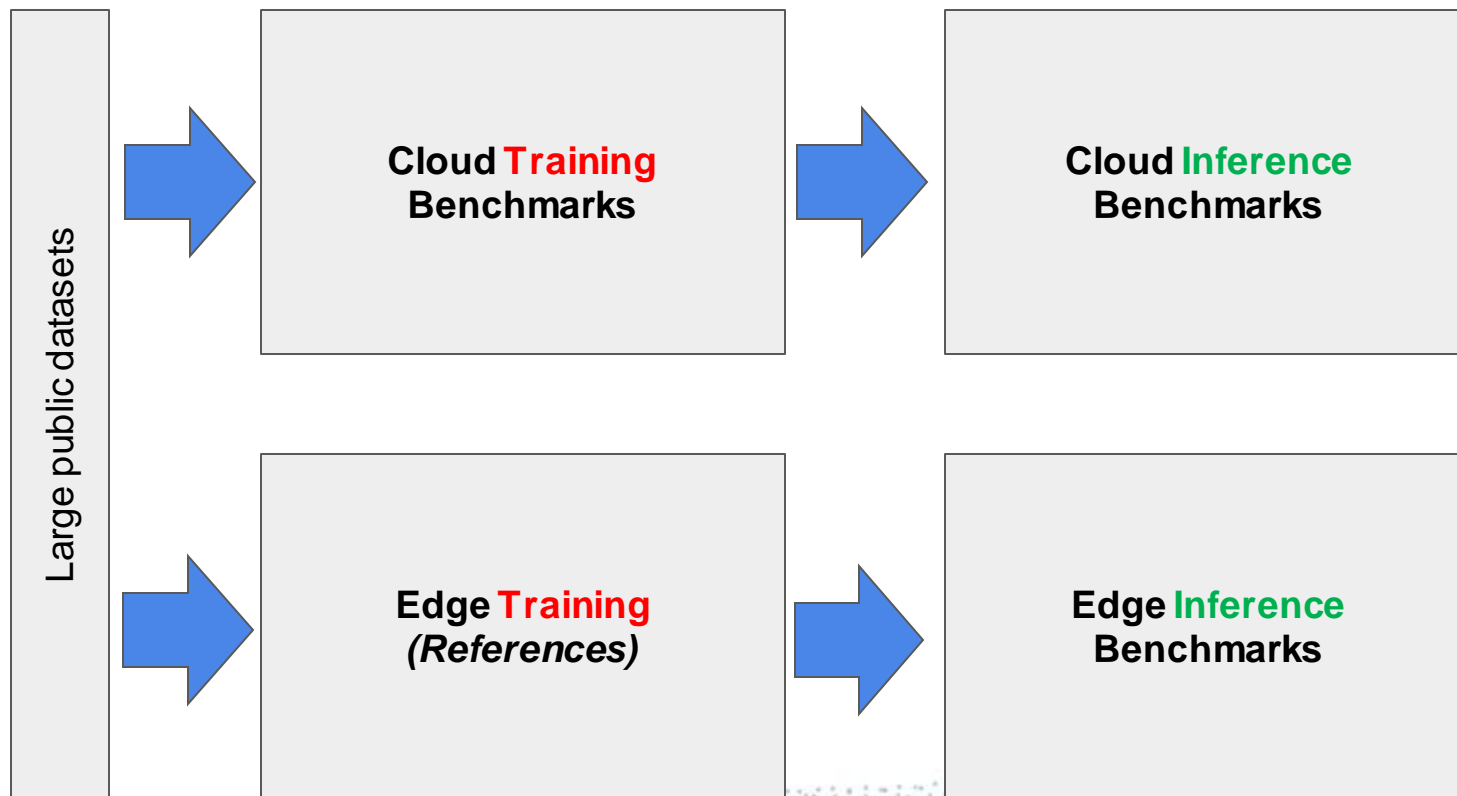
- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ **What is MLPerf?**
 - How does MLPerf curate a benchmark?
 - What is the “science” behind the curation?
 - Where are we heading now?
- What comes next for MLPerf?

What is MLPerf?

A machine learning performance benchmark suite with broad industry and academic support.

Contributions by researchers from





Agenda

- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ *What is MLPerf?*
 - **How does MLPerf curate a benchmark?**
 - What is the “science” behind the curation?
 - Where are we heading now?
- What comes next for MLPerf?

The ML Landscape is Large

Area	Vision	Language	Audio	Commerce	Action / RL	Other
Problem	Image Classification Object Detection / Segmentation Face ID HealthCare (Radiology) Video Detection Self-Driving	Translation Language Model Word Embedding	Speech Recognition Text-to-Speech Question Answering Keyword Spotting Language Modeling Chatbots Speaker ID Graph embeddings Content ID	Rating Recommendations Sentiment Analysis Next-action Healthcare (EHR) Fraud detection Anomaly detection Time series prediction Large scale regression	Games Go Robotics Health Care Bioinformatics	GANs 3D point clouds Word embeddings
Models	ResNet-50 TF Object Detection Detectron	BERT Transformer OpenNMT	Deep Speech 2 SQuAD Explorer	Neural Collaborative Filtering CNNs	DQN PPO A2C	
Datasets	ImageNet COCO	WMT English- German	LibriSpeech SQuAD LM-Benchmark	MovieLens-20M Amazon IMDB	Atari Go Chess Grasping	
Metrics	COCO mAP Prediction accuracy	BLEU	WER Perplexity	Prediction accuracy	Prediction accuracy Win/Loss	

ML Benchmark Design Choices

Big Questions	Training	Inference
1. Benchmark definition	What is the definition of a benchmark task?	
2. Benchmark selection	Which benchmark task to select?	
3. Metric definition	What is the measure of performance in ML systems?	
4. Implementation equivalence	How do submitters run on different hardware/software systems?	
5. Issues specific to training or inference	Which hyperparameters can submitters tune?	Quantization, calibration, and/or retraining?
	Reduce result variance?	
6. Results	Do we normalize and/or summarize results?	

Decision Making Guidelines

Model Range	Example	Principle
Maturity: Lowest common denominator, most widely used, or most advanced?	Image recognition: AlexNet, ResNet, or EfficientNet?	Cutting edge, not bleeding edge
Variety: What broad kind of deep neural network to choose?	Translation: GNMT with RNN vs. Transformer with Attention	Try and ensure coverage at a whole suite level
Complexity: Less or more weights?	Object detection: SSD vs. Mask R-CNN? Resolution?	Survey end-users and anticipate market demand
Practicality	Feasibility: Is there a public dataset?	Good for now > perfect.

MLPerf **Training** Benchmarks 0.5v, 0.6v

Area	Benchmark	Model	Dataset
Vision	Image classification	ResNet-50 v1.5	ImageNet
	Object detection	Mask R-CNN	COCO
		Mask-R-CNN	COCO
Language/Audio	Translation	Transformer	WMT Eng-Germ
		NMT	WMT Eng-Germ
Commerce	Recommendation	NCF	MovieLens-20M
Action	Reinforcement Learning	Mini-go	Go

MLPerf **Inference** Benchmarks 0.5v

Area	Benchmark	Model	Dataset
Vision	Image Classification	MobileNet v1 ResNet50	ImageNet (224x224) ImageNet (224x224)
	Object Detection	SSD-MobileNet v1 SSD-ResNet34	MS-COCO (300x300) MS-COCO (1200x1200)
Language	Translation	Google NMT	WMT Eng-Germ

Agenda

- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ *What is MLPerf?*
 - ✓ *How does MLPerf curate a benchmark?*
 - **What is the “science” behind the curation?**
 - Where are we heading now?
 - What comes next for MLPerf?

Agenda

- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ *What is MLPerf?*
 - ✓ *How does MLPerf curate a benchmark?*
 - **What is the “science” behind the curation?**
 - **Training**
 - Inference
 - Results
 - Where are we heading now?
 - What comes next for MLPerf?

Challenges in Benchmarking ML Systems

Training

- **Hyperparameters**
- **Metric**
- Scale
- Power
- Cost
- Variance
- On-premise vs. cloud
- ...

Inference

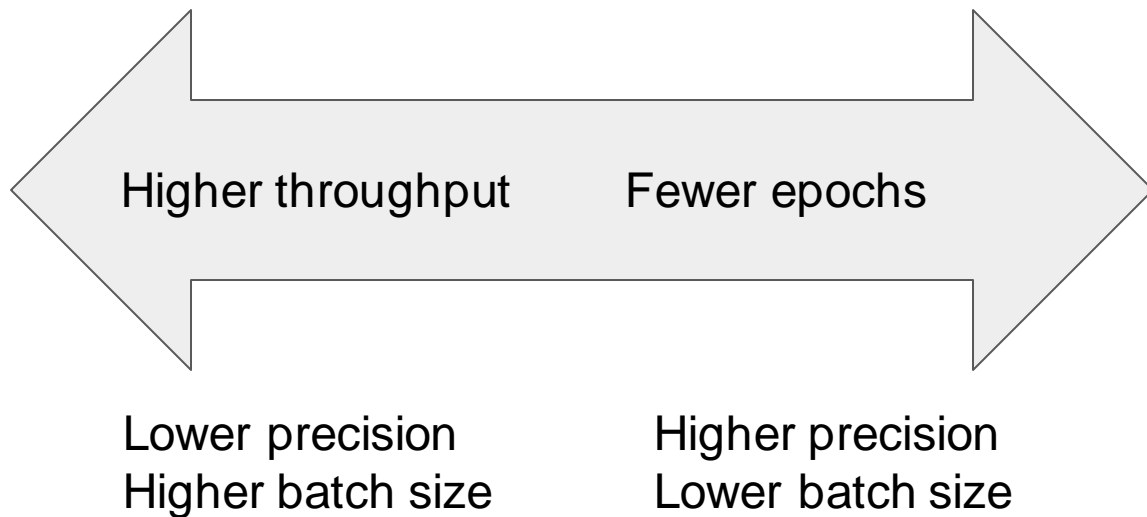
- **Scenarios**
- **Quantizations**
- Pruning
- Scores
- **Power**
- ...

Training **Metric**: Throughput vs. Time-to-Train

Throughput (samples / sec)
Easy / cheap to measure

Can increase throughput at
cost of total time to train!

Time-to-train (end-to-end)
Time to solution!



Training **Hyperparameter Tuning**

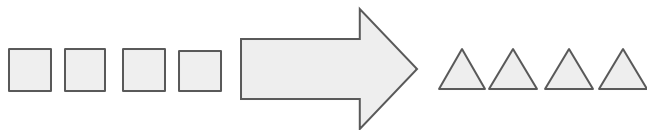
- Different **system sizes** \Rightarrow different **batch sizes** \Rightarrow different **hyperparameters**
- But, some working hyperparameters are better than others
- Finding good hyperparameters is expensive and not the point of the benchmark, **we do not want a hyperparameter tuning competition!**

- Solution v0.5, v0.6:
 - Hyperparameter “borrowing” during review process

Agenda

- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ *What is MLPerf?*
 - ✓ *How does MLPerf curate a benchmark?*
 - **What is the “science” behind the curation?**
 - Training
 - **Inference**
 - Results
 - Where are we heading now?
 - What comes next for MLPerf?

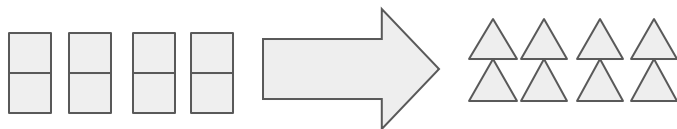
Inference Scenarios & Metrics to Measure



Single stream

(e.g. cell phone augmented vision)

Latency

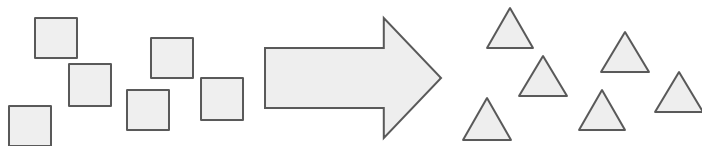


Multiple stream

(e.g. multiple camera driving assistance)

Number streams

subject to latency bound

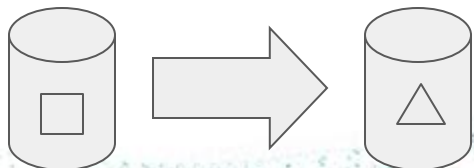


Server

(e.g. translation app)

QPS

subject to latency bound



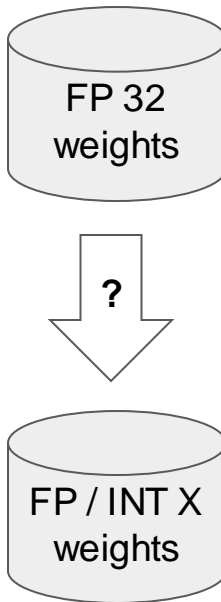
Offline

(e.g. photo sorting app)

Throughput

Inference Quantization and Retraining

- Quantization is key to efficient inference, but do not want a quantization contest – some can do it better than others
- Can the submitters quantize?
 - **Yes**, but must be principled: describe reproducible method
- Can the submitters calibrate?
 - **Yes**, but must use a fixed set of calibration data
- Can the submitters retrain?
 - **No**, not a retraining contest. But, provide retrained 8 bit.



Agenda

- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ *What is MLPerf?*
 - ✓ *How does MLPerf curate a benchmark?*
 - **What is the “science” behind the curation?**
 - Training
 - Inference
 - **Results**
 - Where are we heading now?
 - What comes next for MLPerf?

Results **Normalization** and/or **Scale**

- Do you present only the results?

System
Foo
Bar

Results lack scale information.

If so, an inefficient larger system can look better than an efficient smaller system.

Need supplemental normalization and/or scaling information

MLPerf provides some scale information

Current: Number of chips

Planned: Power

Results Summarization (or Not)?

- Should we have a single MLPerf score that summarizes all results?

System	ResNet	GNMT
Foo	3m	4m
Bar	1m	6m

- Summarized results **Pros**
 - Easy to communicate
Do it consistently
- Summarized results **Cons**
 - Oversimplifies
 - Some vendors only submit subsets
 - Users care about different subsets

MLPerf doesn't summarize.

We recommend weighted geometric mean.

Agenda

- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ *What is MLPerf?*
 - ✓ *How does MLPerf curate a benchmark?*
 - ✓ *What is the “science” behind the curation?*
 - **Where are we heading now?**
- What comes next for MLPerf?

NVIDIA Sets Six Records in AI Performance

Tensor Core GPUs achieve best performance across every MLPerf benchmark submitted;

AI & MACHINE LEARNING

MLPerf benchmark establishes that Google Cloud offers the most accessible

MLPERF RESULTS VALIDATE CPUS FOR DEEP LEARNING TRAINING



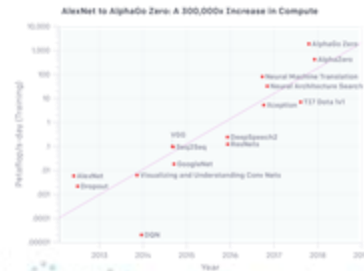
Inference Submissions Due Soon (Oct 11th)

Closed division submissions

- Requires using the specified model
- Enables apples-to-apples comparison
- Simplifies work for HW groups

Open division submissions

- Open division allows using any model
- Encourages innovation
- Ensures Closed division does not stagnate



Agenda

- ✓ *Why ML needs a benchmark suite?*
- ✓ *Are there lessons we can borrow?*
- ✓ *What is MLPerf?*
 - ✓ *How does MLPerf curate a benchmark?*
 - ✓ *What is the “science” behind the curation?*
 - ✓ *Where are we heading now?*
- **What comes next for MLPerf?**

ML Commons

We are creating a non-profit called MLCommons to “accelerate ML innovation and increase its positive impact on society.”

Datasets x Benchmarks x Best Practices

The People's Speech
100k hours of human speech covering the half of the 2020 world population
© Microsoft, Mozilla Common Voice, and others.
[Link]
[Link]

Objective: MLCommons aims to boost speech recognition accuracy for the entire world population by creating a new speech dataset with 100,000 hours of transcribed speech training to improve and diversify voice for machine translation, writing, and reading. Such datasets and general models are the foundation of advanced speech recognition technology available to the entire world's population by the mid-2020s. Existing systems are biased for the most common languages of the world: common languages such as English and Mandarin. However, the real distribution of human speech is very diverse. There are approximately one thousand spoken languages with more than ten million speakers each, many regional accents, and significant variation among dialects and environments. We will focus on boosting accuracy for the half of the distribution by targeting more speech diversity. We will also include specific measures that emerging populations of countries may require. We will also focus on the most diverse linguistic groups to ensure the dataset is representative of the global population and to support the needs of underserved and non-mainstream languages.

Key Deliverables: We are expected to participate in the creation of the dataset, open challenge dataset which has led to significant improvements in speech recognition accuracy by focusing on the community and other metrics. We also plan to release new data sets that are as diverse and as high quality as possible, including systems that have never been publicly available, new applications, and data from multiple, non-mainstream languages, such as the Linguistic Data Consortium, Mozilla Common Voice, LibriSpeech, Common Voice, and LibriVox. We plan to release general purpose, state-of-the-art, open-source, and high-quality technology to collect and label a new set of multilingual speech. Emerging generalist speech and language modeling technologies such as GPT-4, GPT-3.5, and GPT-4o have enabled capabilities to generate realistic speech and language. We will use these technologies to augment our dataset. We will also use state-of-the-art speech recognition to enable the creation of high performance compact systems that are near state-of-the-art speech recognition accuracy.

High Dataset Development: A small core team will rapidly iterate on the creation of real data from real world conditions by studying and analyzing patterns as well as incorporating new research. Accuracy of our real-world datasets based on new data sources will be evaluated on these real sets. We will quickly release an initial set of publicly available 100,000 hours to enable leading edge training systems. Subsequent releases will be targeted for the next iteration on new data sets and significant improvements in real world accuracy.

High Performance: We are excited with various approaches to track the world's top researchers. The use of general purpose and high-quality datasets will be a key factor in the success of our research. We will also release general purpose speech and language modeling systems, and we will release technology to assemble our dataset.

Speech Recognition System: We will evaluate the impact of new data sets on the state-of-the-art speech recognition system and release a new system with a significant accuracy improvement over the current state-of-the-art. We will also release a new system with a significant accuracy improvement over the current state-of-the-art. We will also release a new system with a significant accuracy improvement over the current state-of-the-art.

During this project we will develop more accurate models sometimes and in cooperation with corporations such as the MIT Open Speech Initiative (OSI) and Google.

High Performance Training System: Training a speech recognition system on 100,000 hours of speech will require a new high performance system. We will release the system in the next 100 days of development and we will release a new system with a significant accuracy improvement over the current state-of-the-art. We will also release a new system with a significant accuracy improvement over the current state-of-the-art. We will also release a new system with a significant accuracy improvement over the current state-of-the-art.

Call For Participation: We are seeking contributors to the project who can donate their ML time. See [Link] for details on how to participate.



More at **MLPerf.org**, or contact info@mlperf.org

