

2020  
embedded  
**VISION**  
summit®

# MLPerf: An Industry Standard Benchmark Suite for Machine Learning

Carole-Jean Wu

[carolejeanwu@fb.com](mailto:carolejeanwu@fb.com)

Facebook, ASU

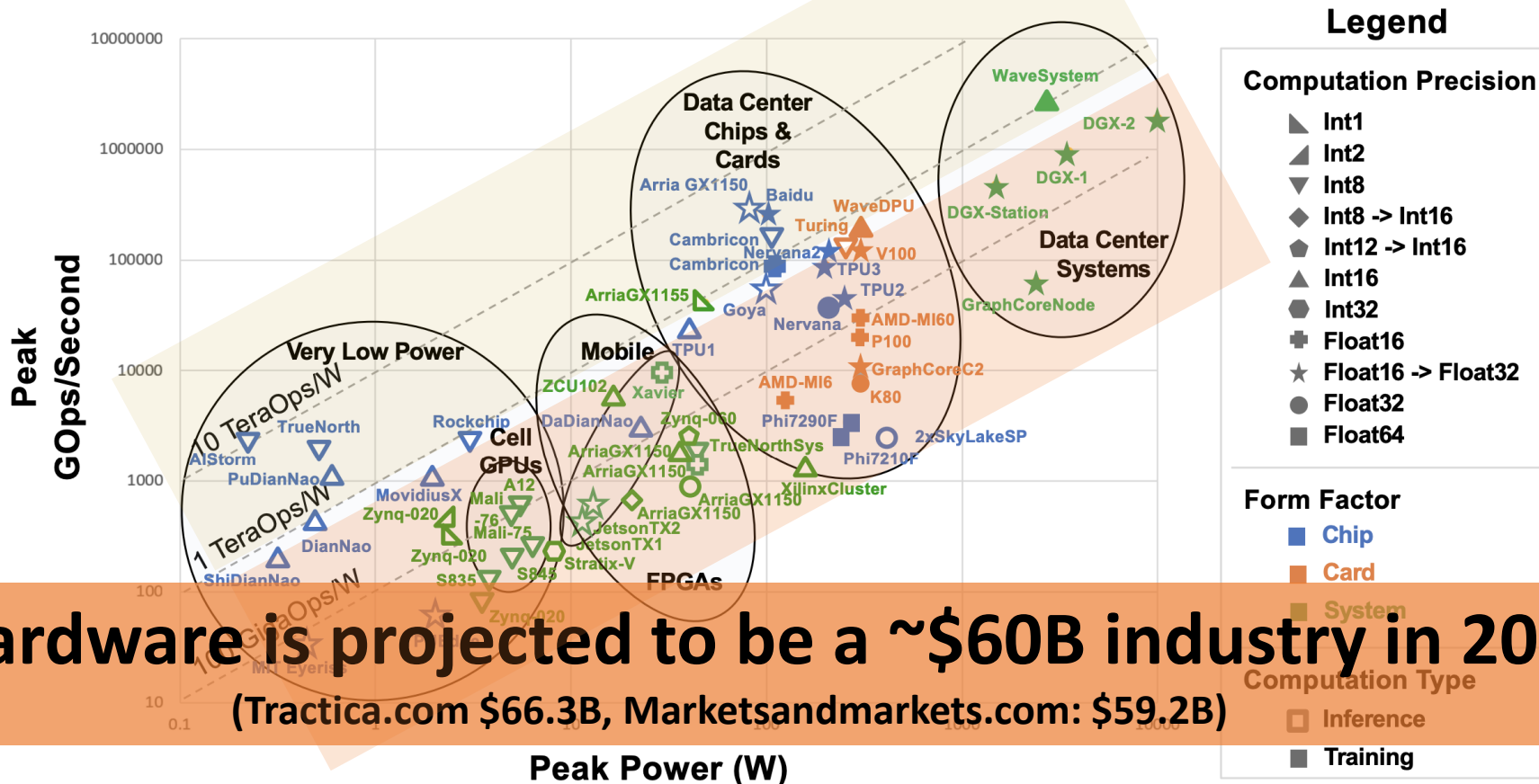
MLPerf Inference Chair and

Recommendation Benchmark Advisory Board Chair

(Work by Many People in the MLPerf Community)



# Deep Learning is Fueling the HW Renaissance



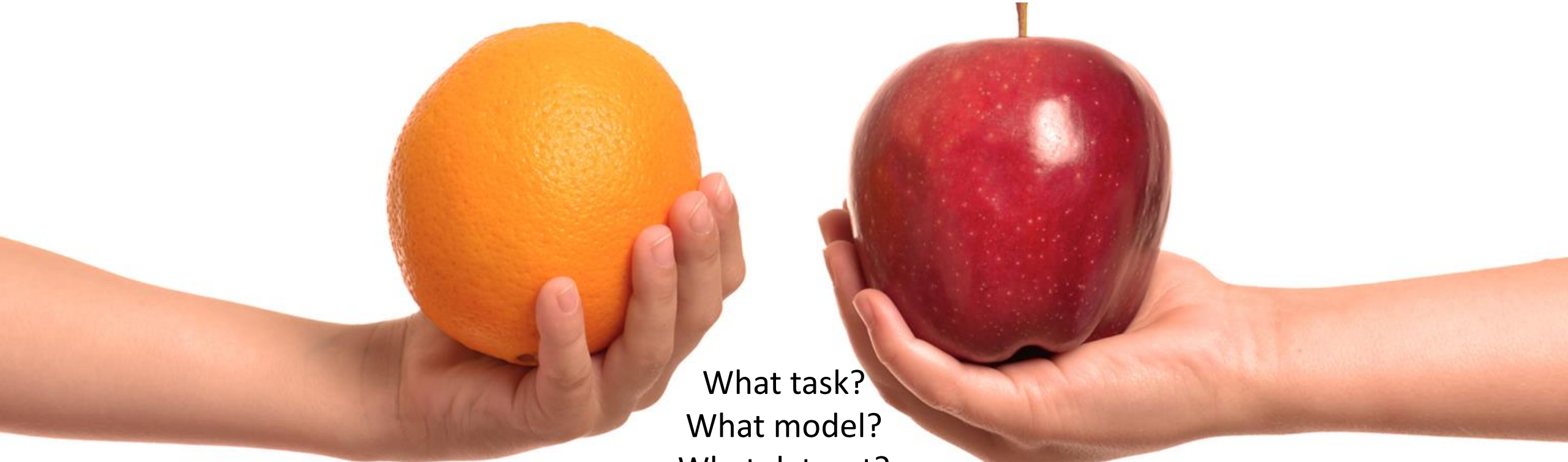
**ML hardware is projected to be a ~\$60B industry in 2025.**

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

Survey and Benchmarking of AI Accelerators. Reuther et al. MIT Lincoln Lab Supercomputing Center. Arxiv-2019

# How Do We Compare AI Hardware?

*Systems*



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

...

***“What get measured, gets improved.”***

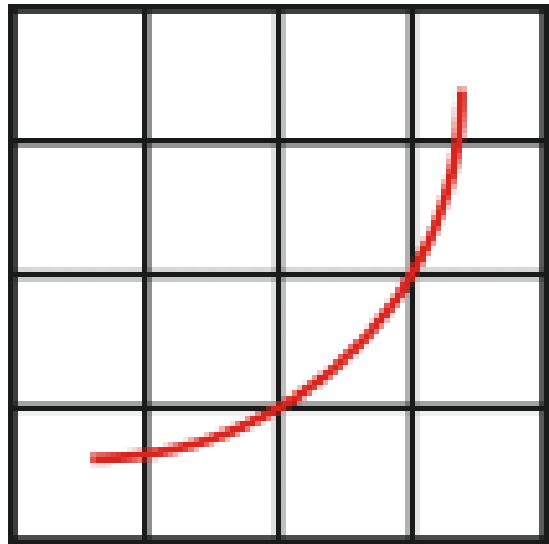
***— Peter Drucker***

Benchmarking aligns research with development,  
engineering with marketing,  
and competitors across the industry in  
pursuit of a clear objective

- *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?
- What comes next for MLPerf?
- How can we contribute to MLPerf?

# Are There **Lessons** We Can Borrow?

# Successful History in Benchmarks



**spec<sup>®</sup>**

**TPC<sup>™</sup>**

- Settled **arguments in the marketplace** (grow the pie)
- Resolved internal **engineering debates** (better investments)
- **Cooperative** nonprofit corporation with 22 members
- **Universities** join at modest cost and help drive innovation
- Became **standard** in marketplace, papers, and textbooks
- Needed to **revise regularly** to maintain usefulness:  
SPEC89, SPEC92, SPEC95, SPEC2000, SPEC2006, SPEC2017

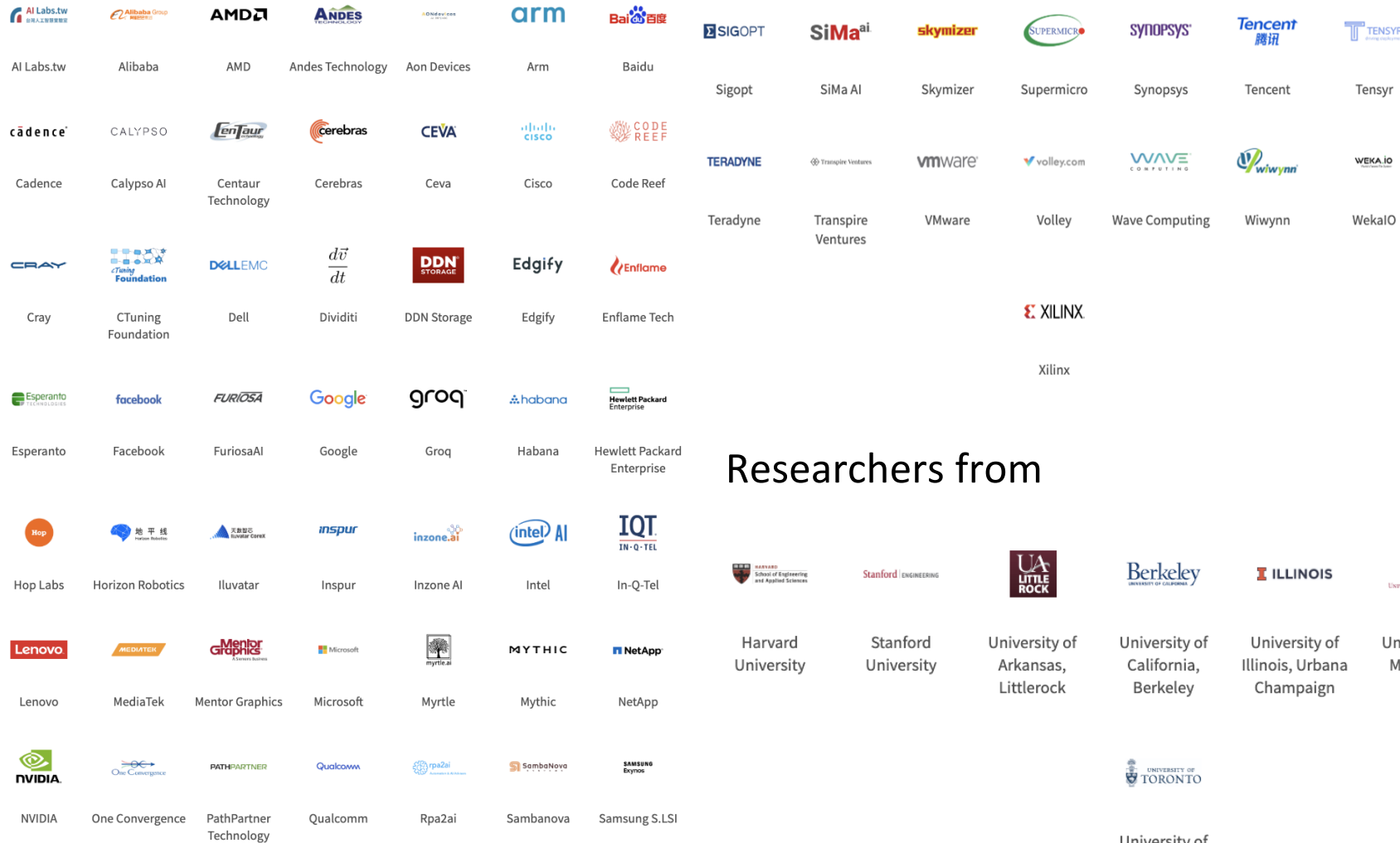
⇒ **Fueled the golden age of microprocessor design**



# How Do We Start a New Golden Age for ML Systems?

- *Why ML needs a benchmark suite?*
  - *Are there lessons we can borrow?*
- **What is MLPerf?**
  - How does MLPerf curate a benchmark?
  - Training
  - Inference
- What comes next for MLPerf?
- How can we contribute to MLPerf?

# MLPerf is an ML Performance Benchmarking Effort with Wide Industry and Academic Support



## Researchers from



# ML Benchmark Design Overview

Big Questions	Training	Inference
1. Benchmark definition	What is a benchmark task?	
2. Benchmark selection	Which benchmark tasks?	
3. Metric definition	What is performance?	
4. Implementation equivalence	How do submitters run on very 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. Presentation	Do we normalize and/or summarize results?	

# Training Benchmark Definition

Dataset



Train a  
model

Target Quality<sup>1</sup>

**75.9%**

Do we specify the model?

1. Target quality set by experts in area, raised as SOTA improves

# MLPerf Training Benchmark Suite (v0.7)

Benchmark	Data set	Model	Quality Threshold
Image classification	ImageNet (Deng et al., 2009)	ResNet-50 v1.5 (MLPerf, 2019b)	75.9% Top-1 accuracy
Object detection (lightweight)	COCO 2017 (Lin et al., 2014)	SSD-ResNet-34 (Liu et al., 2016)	23.0% mAP
Instance segmentation and object detection (heavyweight)	COCO 2017 (Lin et al., 2014)	Mask R-CNN (He et al., 2017a)	37.7 Box min AP, 33.9 Mask min AP
NLP	Wikipedia 01/01/2020	BERT	0.712 Mask-LM accuracy
Translation (nonrecurrent)	WMT17 EN-DE (WMT, 2017)	Transformer (Vaswani et al., 2017)	25.0 BLEU
Recommendation	Criteo Terabyte CTR	DLRM (Naumov et al., 2019)	0.8025 AUC
Reinforcement learning	Go (19 x 19 Board)	MiniGo (MLPerf, 2019a)	50% win rate

# Training Metric: Time to Reach Quality Target

- Quality target is **specific for each benchmark and close to state-of-the-art**
  - Updated w/ each release to keep up with the SOTA
- Time includes preprocessing and validation over of N runs
- MLPerf provides the reference implementations that achieve quality target

- *Why ML needs a benchmark suite?*
  - *Are there lessons we can borrow?*
- **What is MLPerf?**
  - *How does MLPerf curate a benchmark?*
  - *Training*
  - **Inference**
- What comes next for MLPerf?
- How can we contribute to MLPerf?



# Inference Benchmark Definition

Input



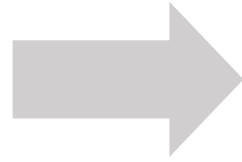
Trained  
ResNet

Result  
(with required quality, e.g. 75.1%)

“Leopard”

Do you specify the model? Closed Division does, Open Division does not.

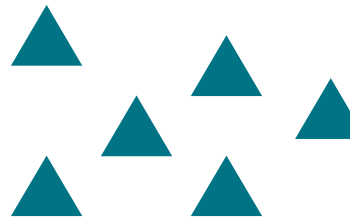
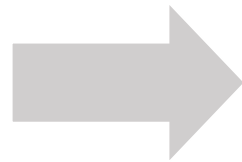
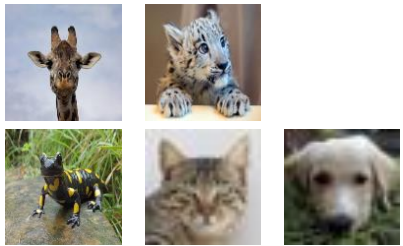
# But How is Inference Really Used?



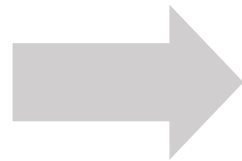
Single stream  
(e.g. cell phone augmented vision)



Multiple stream  
(e.g. multiple camera driving assistance)



Server  
(e.g. translation app)

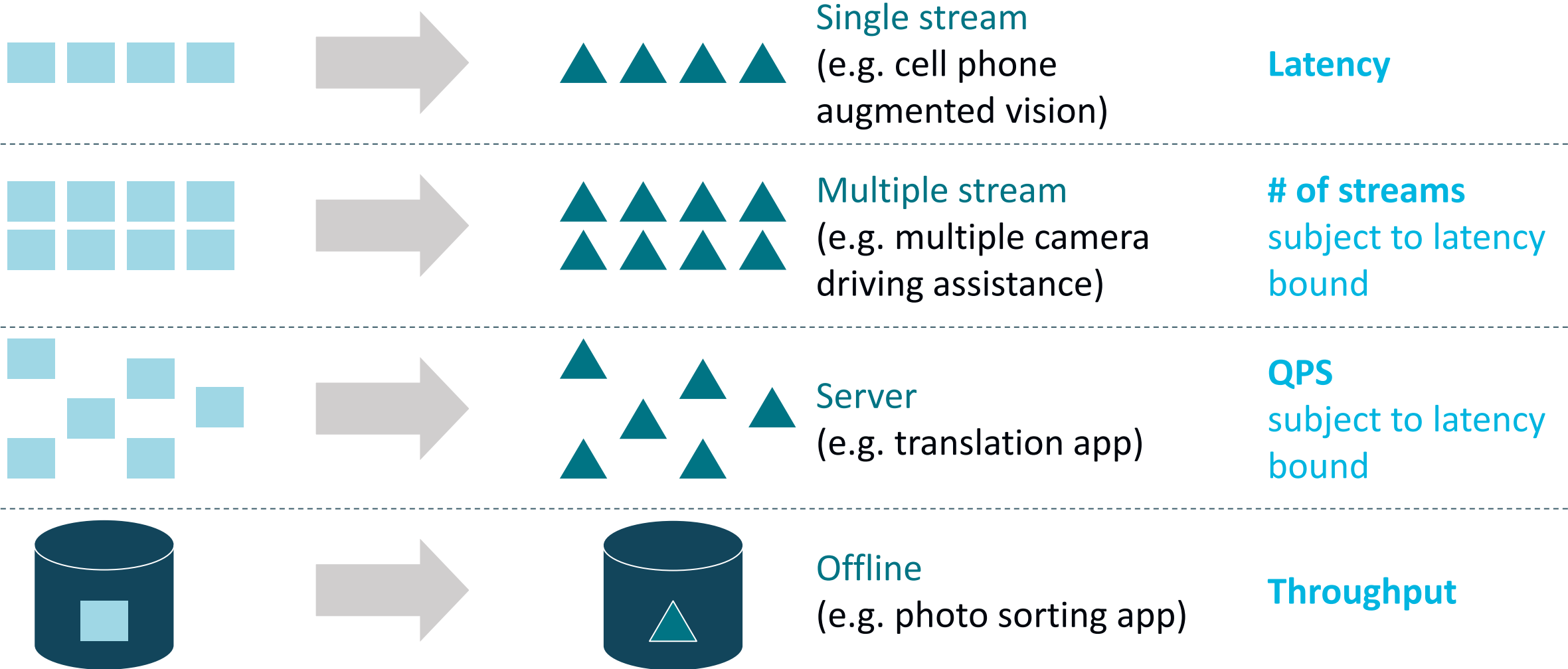


Offline  
(e.g. photo sorting app)

# MLPerf Inference Benchmark Suite (v0.5)

Area	Task	Model	Dataset	Quality	Server latency constraint	Multi-Stream latency constraint
Vision	Image classification	Resnet50-v1.5	ImageNet (224x224)	99% of FP32 (76.46%)	15 ms	50 ms
Vision	Image classification	MobileNets-v1 224	ImageNet (224x224)	98% of FP32 (71.68%)	10 ms	50 ms
Vision	Object detection	SSD-ResNet34	COCO (1200x1200)	99% of FP32 (0.20 mAP)	100 ms	66 ms
Vision	Object detection	SSD-MobileNets-v1	COCO (300x300)	99% of FP32 (0.22 mAP)	10 ms	50 ms
Language	Machine translation	GNMT	WMT16	99% of FP32 (23.9 BLEU)	250 ms	100 ms

# Inference Metric: One Metric for Each Scenario



# Inference Benchmark Suite v0.5

Area	Task	Model	Dataset	Quality	Server latency constraint	Multi-Stream latency constraint
Vision	Image classification	Resnet50-v1.5	ImageNet (224x224)	99% of FP32 (76.46%)	15 ms	50 ms
Vision	Image classification	MobileNets-v1 224	ImageNet (224x224)	98% of FP32 (71.68%)	10 ms	50 ms
Vision	Object detection	SSD-ResNet34	COCO (1200x1200)	99% of FP32 (0.20 mAP)	100 ms	66 ms
Vision	Object detection	SSD-MobileNets-v1	COCO (300x300)	99% of FP32 (0.22 mAP)	10 ms	50 ms
Language	Machine translation	GNMT	WMT16	99% of FP32 (23.9 BLEU)	250 ms	100 ms

# MLPerf Inference Benchmark Suite v0.7

	Neural Network
Vision	ResNet-50 v1.5
	SSD ResNet-34
	SSD MobileNet v1 (edge)
	3D UNET
Speech	RNN-T
Language	BERT Large
Commerce	DLRM (datacenter)

	MOBILE	Neural Network
Vision		MobileNet EdgeTPU
		SSD-MobileNet v2
		DeepLabv3
Language		Mobile-BERT

# MLPerf Training and Inference Call for Submissions

## Closed division submissions

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

## Open division submissions

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

# Timeline for Benchmark Result Submission

**Week -12**

Benchmark  
Freeze

- Discuss details on rules and models and implementations
- Develop the SW stack

**Week -6**

Code &  
Rule Freeze

- Clarify rules
- Fine-tune SW for HW
- Pass compliance check
- Sign CLA
- Release SW

**Week 0**

Submission  
Deadline

- Review submission results
- Prepare result release

**Week 4**

Result  
Publication

**ATTEND WEEKLY SUBMITTER WORKING GROUP MEETINGS**



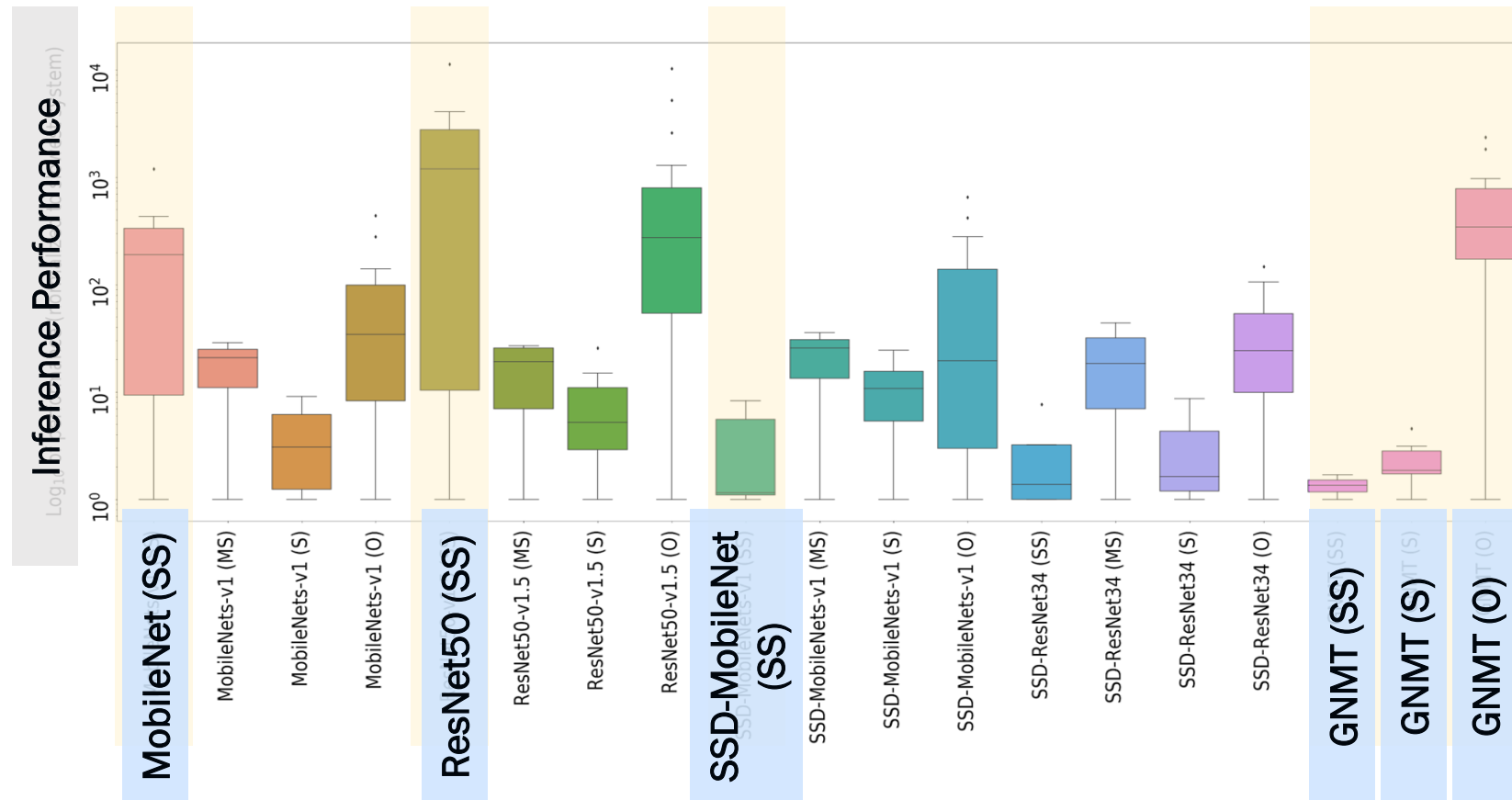
# MLPerf Inference Results

<https://mlperf.org/inference-results/>

ID	Submissions		Benchmarks				ResNet-50 v1.5			
	Submitter	System	Scenarios				Stream	MultiS	Server	
			MobileNet-v1	Stream	MultiS	Server				Offline
CATEGORY: Available										
Inf-0.5-1	Alibaba Cloud	Alibaba Cloud T4						17,473.60		
Inf-0.5-2	Dell EMC	Dell EMC R740					67,124.18	71,214.50		20,742.8
Inf-0.5-3	Dell EMC	Dell EMC R740xd with 2nd generation Intel® Xeon® Scalable Processor								
Inf-0.5-4	Dell EMC	Dell EMC R740xd with 2nd generation Intel® Xeon® Scalable Processor								
Inf-0.5-5	dividiti	Raspberry Pi 4 (rpi4)				394.34			1,916.65	
Inf-0.5-6	dividiti	Raspberry Pi 4 (rpi4)				103.60			448.31	
Inf-0.5-7	dividiti	Linaro HiKey960 (hikey960)				121.11			518.07	
Inf-0.5-8	dividiti	Linaro HiKey960 (hikey960)				50.77			203.99	
Inf-0.5-9	dividiti	Linaro HiKey960 (hikey960)				143.07			494.90	
Inf-0.5-10	dividiti	Huawei Mate 10 Pro (mate10pro)				74.20			354.13	
Inf-0.5-11	dividiti	Huawei Mate 10 Pro (mate10pro)				111.60			494.92	
Inf-0.5-12	dividiti	Firefly-RK3399 (firefly)				120.56			695.11	
Inf-0.5-13	dividiti	Firefly-RK3399 (firefly)				106.49			447.90	
Inf-0.5-14	dividiti	Firefly-RK3399 (firefly)				80.12			391.02	
Inf-0.5-15	Google	Cloud TPU v3								16,014.2
Inf-0.5-16	Google	2x Cloud TPU v3								
Inf-0.5-17	Google	4x Cloud TPU v3								
Inf-0.5-18	Google	8x Cloud TPU v3								
Inf-0.5-19	Google	16x Cloud TPU v3								
Inf-0.5-20	Google	32x Cloud TPU v3								
Inf-0.5-21	Habana Labs	HL-102-Goya PCI-board							0.24	700.00
Inf-0.5-22	Intel	Intel® Xeon® Platinum 9200 processors								
Inf-0.5-23	Intel	Intel® Xeon® Platinum 9200 processors				0.49	27,244.81	29,203.30	1.37	4,850

ML

# MLPerf Inference v0.5 Presented Almost 600 Results Across a Wide Range of Platform Scales



Normalized performance on log<sub>10</sub> scale for models and scenarios. Results are normalized to the slowest system and show up to a 10,000X range in performance.

- *Why ML needs a benchmark suite?*
- *Are there lessons we can borrow?*
- *What is MLPerf?*
  - *How does MLPerf curate a benchmark?*
  - *Training*
  - *Inference*
- **What comes next for MLPerf?**
- How can we contribute to MLPerf?

# MLPerf Focuses for 2020

Evolve the benchmark suites fairly

Improve efficiency information

Reduce result sparsity

Reduce benchmarking cost

Inference (Datacenter/Edge)

Inference (Mobile)



# Conclusion

- **Benchmarking ML Systems is hard** due to the fragmented ecosystem
- MLPerf is a **community-driven ML benchmark** for the HW/SW industry
- The benchmark suite helps level the playing field, **enabling ML system comparison**
  - Defines Tasks, Scenarios, Datasets, Methods
  - Establish clear set of metrics and divisions
  - Allows for hardware/software flexibility

MLPERF TRAINING BENCHMARK

Peter Mattson<sup>1</sup>, Christine Cheng<sup>2</sup>, Cody Coleman<sup>3</sup>, Greg Diamos<sup>4</sup>, Paulius Micikevicius<sup>5</sup>, David Patterson<sup>1,6</sup>, Hanlin Tang<sup>7</sup>, Gu-Yeon Wei<sup>8</sup>, Peter Bailis<sup>9</sup>, Victor Bittorf<sup>1</sup>, David Brooks<sup>4</sup>, Dhaos Chen<sup>1</sup>, Debajyoti Dutta<sup>4</sup>, Udit Gupta<sup>1</sup>, Kim Hazelwood<sup>1</sup>, Andrew Hock<sup>10</sup>, Xinyuan Huang<sup>8</sup>, Bill Jia<sup>1</sup>, Daniel Kang<sup>3</sup>, David Kanter<sup>11</sup>, Naveen Kumar<sup>1</sup>, Jeffery Liao<sup>12</sup>, Guokai Ma<sup>2</sup>, Deepak Narayanan<sup>3</sup>, Tayo Oguntibi<sup>1</sup>, Gennady Pekhimenko<sup>13</sup>, Lillian Pentecost<sup>7</sup>, Vijay Janapa Reddi<sup>1</sup>, Taylor Robie<sup>1</sup>, Tom St. John<sup>14</sup>, Carole-Jean Wu<sup>9</sup>, Lingjie Xu<sup>15</sup>, Cliff Young<sup>1</sup>, Matei Zaharia<sup>1</sup>

MLPerf Inference Benchmark

Vijay Janapa Reddi,<sup>1</sup> Christine Cheng,<sup>1</sup> David Kanter,<sup>3</sup> Peter Mattson,<sup>5</sup> Guenther Schmuelling,<sup>4</sup> Carole-Jean Wu,<sup>1</sup> Brian Anderson,<sup>1</sup> Maximilien Breughe,<sup>16</sup> Mark Charlebois,<sup>11</sup> William Chou,<sup>11</sup> Ramesh Chukka,<sup>1</sup> Cody Coleman,<sup>11</sup> Sam Davis,<sup>17</sup> Pan Deng,<sup>18</sup> Greg Diamos,<sup>19</sup> Jared Duke,<sup>1</sup> Dave Fick,<sup>20</sup> J. Scott Gardner,<sup>21</sup> Itay Hubara,<sup>22</sup> Sachin Iyengar,<sup>23</sup> Thomas B. Jablin,<sup>3</sup> Jeff Jiao,<sup>24</sup> Tom St. John,<sup>25</sup> Pankaj Kanwar,<sup>1</sup> David Lee,<sup>26</sup> Jeffery Liao,<sup>27</sup> Anton Lokhtin,<sup>28</sup> Francisco Massa,<sup>29</sup> Peng Meng,<sup>30</sup> Paulius Micikevicius,<sup>31</sup> Colin Osborne,<sup>32</sup> Gennady Pekhimenko,<sup>33</sup> Arun Tejasvee Raghunath Rajan,<sup>34</sup> Dilip Sequeira,<sup>35</sup> Ashish Srivastava,<sup>36</sup> Fei Sun,<sup>37</sup> Hanlin Tang,<sup>38</sup> Michael Thomson,<sup>39</sup> Frank Wei,<sup>40</sup> Ephrem Wu,<sup>41</sup> Lingjie Xu,<sup>42</sup> Koichi Yamada,<sup>43</sup> Bing Yu,<sup>44</sup> George Yuan,<sup>45</sup> Aaron Zhong,<sup>46</sup> Peizhao Zhang,<sup>47</sup> Yuchen Zhou<sup>48</sup>

<sup>1</sup>Harvard University <sup>2</sup>Intel <sup>3</sup>Real World Insights <sup>4</sup>Google <sup>5</sup>Microsoft <sup>6</sup>Facebook <sup>7</sup>NVIDIA <sup>8</sup>Qualcomm <sup>9</sup>Stanford University <sup>10</sup>Myrtle <sup>11</sup>Landing AI <sup>12</sup>Mythic <sup>13</sup>Advantage Engineering <sup>14</sup>Habana Labs <sup>15</sup>Alibaba T-Head <sup>16</sup>Facebook (formerly at MediaTek) <sup>17</sup>OPPO (formerly at Synopsys) <sup>18</sup>dividiti <sup>19</sup>Arm <sup>20</sup>Google <sup>21</sup>University of California, Berkeley <sup>22</sup>Intel <sup>23</sup>Google <sup>24</sup>Google <sup>25</sup>Google <sup>26</sup>Google <sup>27</sup>Google <sup>28</sup>Google <sup>29</sup>Google <sup>30</sup>Google <sup>31</sup>Google <sup>32</sup>Google <sup>33</sup>Google <sup>34</sup>Google <sup>35</sup>Google <sup>36</sup>Google <sup>37</sup>Google <sup>38</sup>Google <sup>39</sup>Google <sup>40</sup>Google <sup>41</sup>Google <sup>42</sup>Google <sup>43</sup>Google <sup>44</sup>Google <sup>45</sup>Google <sup>46</sup>Google <sup>47</sup>Google <sup>48</sup>Google

MLPerf: An Industry Standard Benchmark Suite for Machine Learning Performance

Peter Mattson  
Google Brain

Vijay Janapa Reddi  
Harvard University

Christine Cheng  
Intel

Cody Coleman  
Stanford University

Greg Diamos  
Landing AI

David Kanter  
Real World Technologies

Paulius Micikevicius  
NVIDIA

David Patterson  
Google Brain and University of California, Berkeley

Guenther Schmuelling  
Microsoft Azure AI Infrastructure

Hanlin Tang  
Intel

Gu-Yeon Wei  
Harvard University

Carole-Jean Wu  
Facebook and Arizona State University

**Abstract**—In this article, we describe the design choices behind MLPerf, a machine learning performance benchmark that has become an industry standard. The first two rounds of the MLPerf Training benchmark helped drive improvements to software-stack performance and scalability, showing a 1.3× speedup in the top 16-chip results despite higher quality targets and a 5.5× increase in system scale. The first round of MLPerf Inference received over 500 benchmark results from 14 different organizations, showing growing adoption.

# MLPerf is the Work of Many

Aaron Zhong

David Patterson

Jared Duke

Peter Bailis

Abid Muslim

Debajyoti Pal

Jeff Jiao

Peter Baldwin

Andrew Hock

Debo Dutta

Jeffery Liao

Peter Mattson

Ankur Ankur

Deepak Narayanan

Jonah Alben

Ramesh Chukka

Anton Lokhmotov

Dehao Chen

Jonathan Cohen

Sachin Idgunji

Arun Rajan

Dilip Sequeira

Kim Hazelwood

Sam Davis

Ashish Sirasao

Ephrem Wu

Koichi Yamada

Sarah Bird

Atsushi Ike

Fei Sun

Lillian Pentecost

Sergey Serebryakov

Bill Jia

Francisco Massa

Lingjie Xu

Steve Farrell

Bing Yu

Frank Wei

Mark Charlebois

Taylor Robie

Brian Anderson

Gennady Pekhimenko

Masafumi Yamazaki

Tayo Oguntebi

Carole-Jean Wu

George Yuan

Matei Zaharia

Thomas B. Jablin

Christine Cheng

Greg Diamos

Maximilien Breughe

Tom St. John

Cliff Young

Gu-Yeon Wei

Michael Thomson

Tsuguchika Tabaru

Cody Coleman

Guenther Schmuelling

Naveen Kumar

Udit Gupta

Colin Osborne

Guokai Ma

Pan Deng

Victor Bittorf

Daniel Kang

Hanlin Tang

Pankaj Kanwar

Vijay Janapa Reddi

Dave Fick

Itay Hubara

Paulius Micikevicius

William Chou

David Brooks

J. Scott Gardner

Peizhao Zhang

Xinyuan Huang

David Kanter

Jacob Balma

Peng Meng

Yuchen Zhou

# MLPerf Needs Your Help!

- Join the discussion community at [mlperf.org](https://mlperf.org)
- Help us by joining a working group
  - Cloud training and inference
  - Mobile Inference, HPC
  - On-premises scale, submitters, special topics
  - Help us design submission criteria, to include the data you want
- Propose new benchmarks and data sets
- Address challenging “special topic” issues
- Submit your benchmark results!

## info@mlperf.org

### Latest Inference Results October 19, 2020





Thank You

# MLPerf: An Industry Standard Benchmark Suite for Machine Learning



Carole-Jean Wu  
carolejeanwu@fb.com

MLPerf Inference Chair  
Recommendation Benchmark Advisory Board Chair