

tinyML[®] Summit

Enabling Ultra-low Power Machine Learning at the Edge

February 12-13, 2020

Burlingame, California



www.tinyML.org



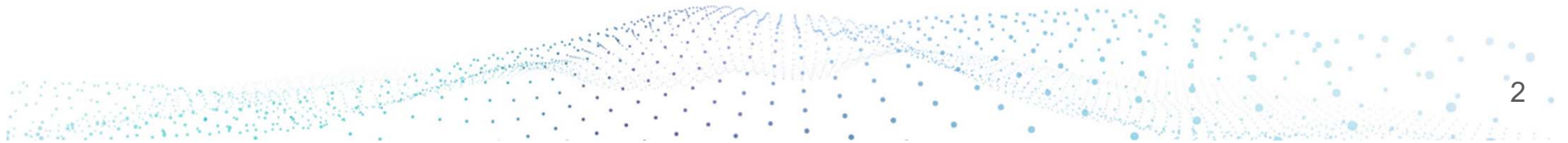
Benchmarking Ultra-low Power Machine Learning Systems

Prof. Vijay Janapa Reddi

Harvard University

MLPerf Inference Chair

*(representing the work of **many** people!)*





What is MLPerf ?

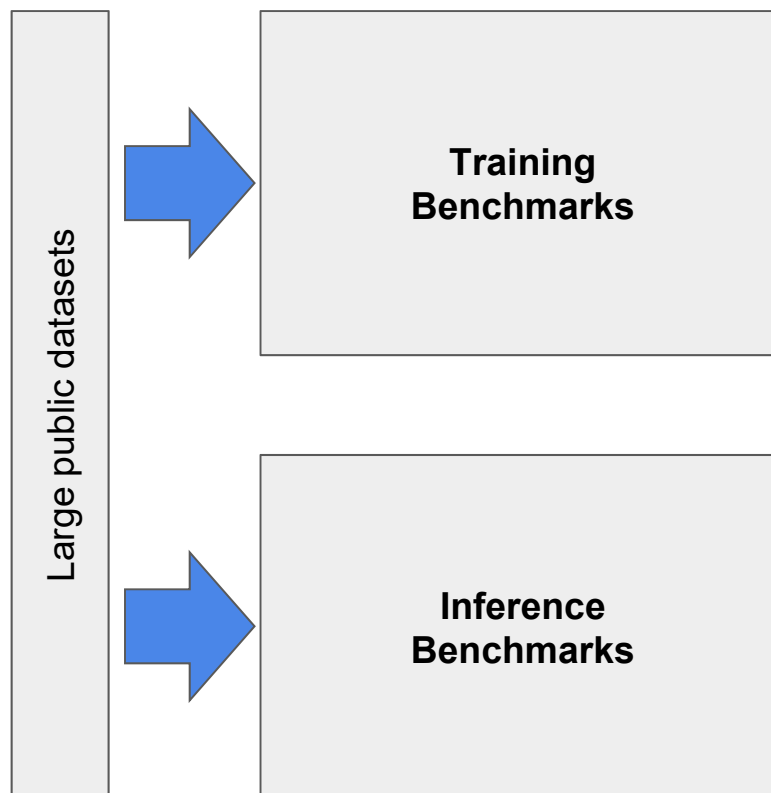


A Community-driven ML Benchmark Suite



1,000+ members, 50+ organizations, 8+ universities





1. Identify a set of **ML tasks** and models
2. Identify real-world **scenarios** to emulate
3. Outline the **rules** for benchmarking
4. Define a clear set of evaluation **metrics**
5. Collect **results** to publish

MLPerf Goals

- **Enforce replicability** to ensure reliable results
- Use **representative workloads**, reflecting production use-cases
- **Encourage innovation** to improve the state-of-the-art of ML
- Accelerate progress in ML via **fair and useful measurement**
- Serve both the **commercial and research communities**
- Keep **benchmarking affordable** (so that all can play)

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

Inference v0.5 Results

MLPerf Inference v0.5 Results

November 6th, 2019

Any use of the MLPerf results and site must comply with the [MLPerf Terms of Use](#).

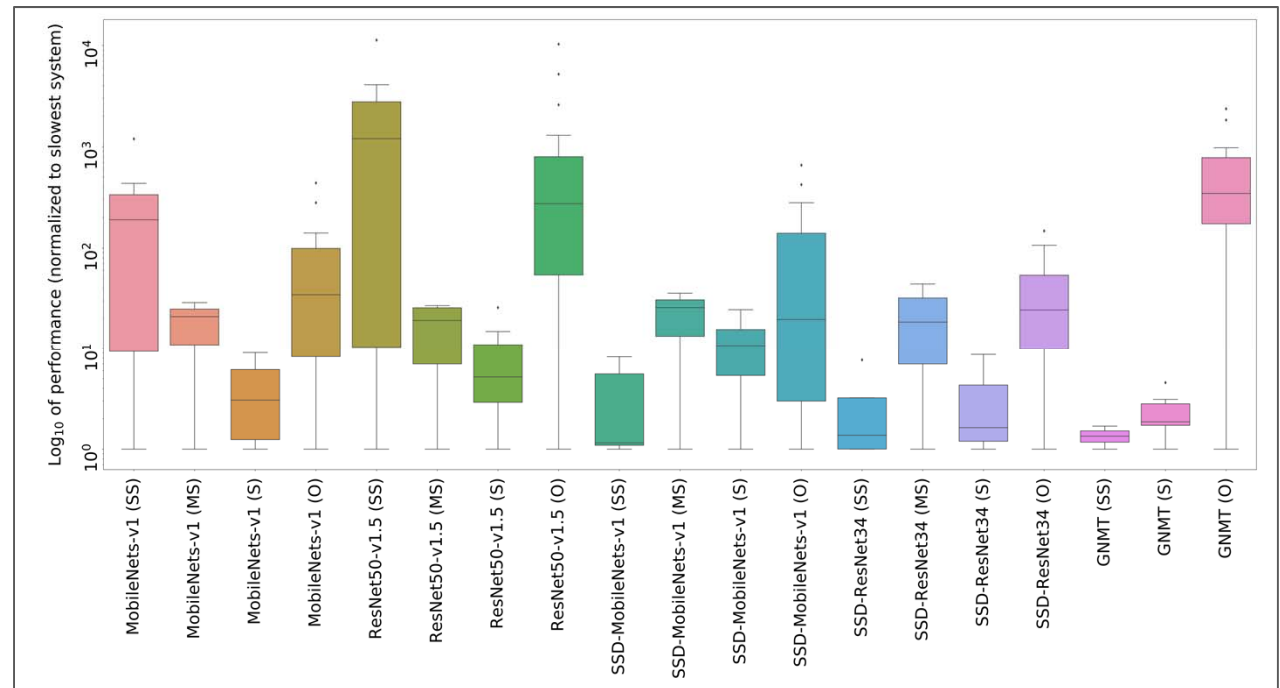
You may wish to read the [Inference Overview](#) to better understand the results.

- Closed Division Performance
- Open Division Performance

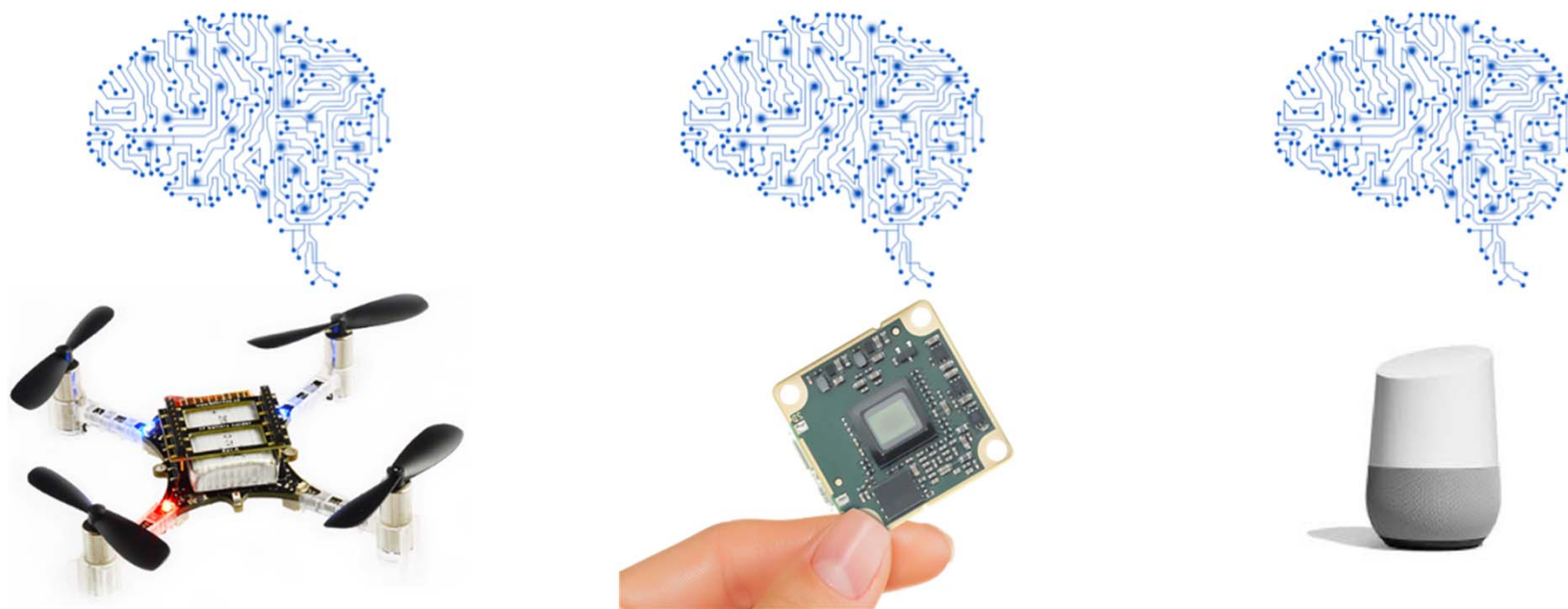
Closed Division Times																	
ID	Submitter	System	Benchmark results (Single Stream in milliseconds, MultiStream in no. streams, Server in QPS, Offline in inputs/second)														
			Image classification						Object detection								
			ImageNet				ImageNet		COCO				COCO				
			MobileNet-v1			ResNet-50 v1.5			SSD w/ MobileNet-v1				SSD w/				
Stream			MultiS			Server			Offline			Stream			MultiS		
CATEGORY: Available																	
Inf-0.5-1	Alibaba Cloud	Alibaba Cloud T4					17,473.60					5,540.10		7,431.20			
Inf-0.5-2	Dell EMC	Dell EMC R740			67,124.18		71,214.50			20,742.83		22,438.00		28,293.31	30,407.90		
Inf-0.5-3	Dell EMC	Dell EMC R740xd with 2nd generation Intel® Xeon® Scalable Processor										1.54		3,744.24			
Inf-0.5-4	Dell EMC	Dell EMC R740xd with 2nd generation Intel® Xeon® Scalable Processor										1.69		4,266.46			
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							

Inference Results

- 600+ inference results
- Over 30 systems submitted
- 10,000x difference in performance

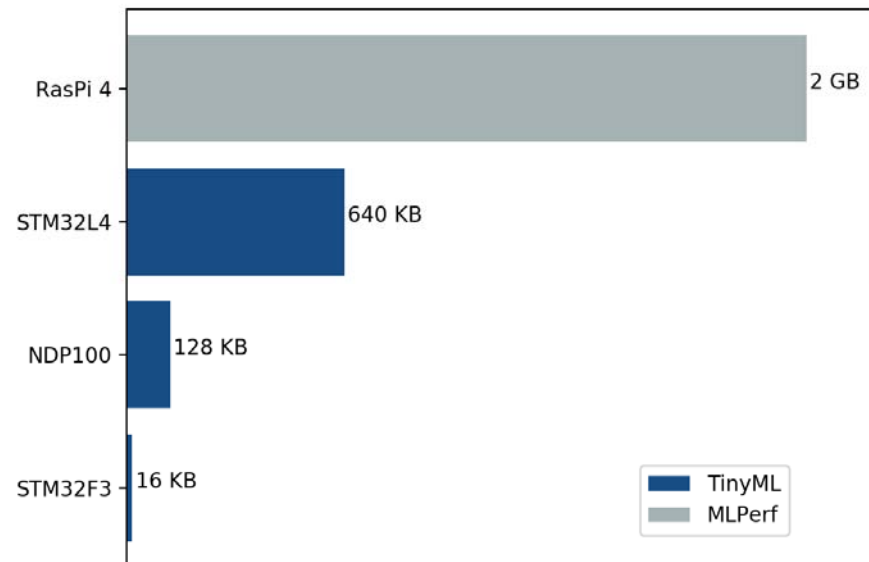


What about TinyML systems?



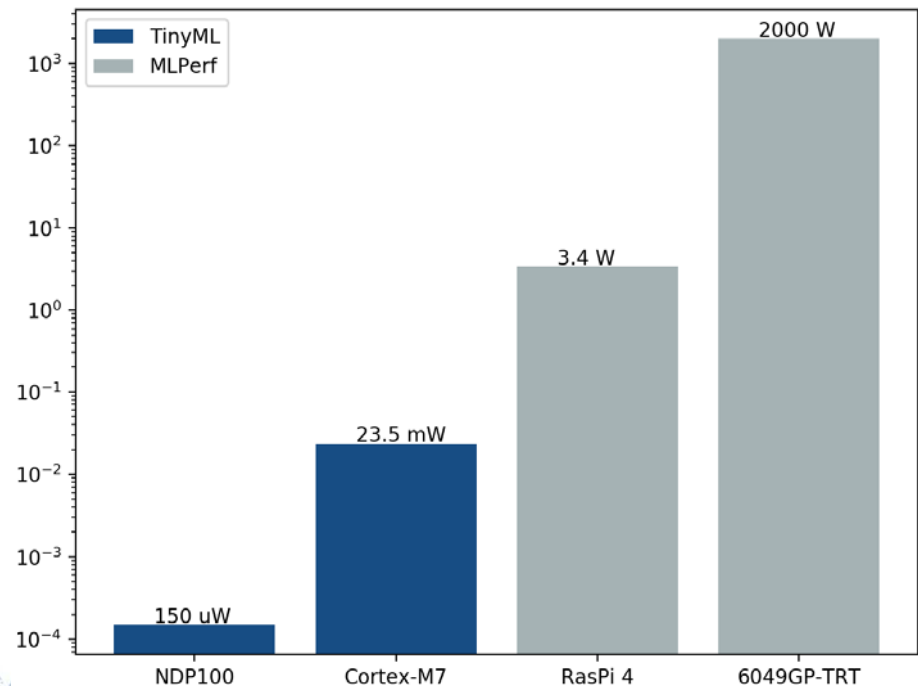
TinyML Challenges for ML Benchmarking

- Resources are extremely condensed in the tinyML devices
- Need to define a methodology to load the SUT for evaluation
- **How can we come up with a methodology that works across many different systems?**



TinyML Challenges for ML Benchmarking

- Power is optional in MLPerf
- MLPerf power working group is trying to develop a specification
- But power is a first-order design constraint in TinyML devices
- **How to define a power spec?**



TinyML Challenges for ML Benchmarking

- Plethora of techniques
 - Quantization
 - Sparsity
 - Pruning
 - Retraining
 - ...
- What are the rules that apply for tinyML?
- Which of the optimizations should be allowed while still enabling fair comparisons?

tinyMLPerf Working Group Members



HARVARD
John A. Paulson
School of Engineering
and Applied Sciences



OctoML



W
UNIVERSITY of
WASHINGTON



ASU
Arizona State



Microsoft



Reality AI



AONdevices

Explainable AI at the Edge

AI | DSP | ASIC



Red Hat

cadence



TinyML Tasks for ML Benchmarking

Task Category	Use Case	Model Type	Datasets
Audio	Audio Wake Words Context Recognition Control Words Keyword Detection	DNN CNN RNN LSTM	Speech Commands Audioset ExtraSensory Freesound DCASE
Image	Visual Wake Words Object Detection Gesture Recognition Object Counting Text Recognition	DNN CNN SVM Decision Tree KNN Linear	Visual Wake Words CIFAR10 MNIST ImageNet DVS128 Gesture
Physiological / Behavioral Metrics	Segmentation Anomaly Detection Forecasting Activity Detection	DNN Decision Tree SVM Linear	Physionet HAR DSA Opportunity
Industry Telemetry	Sensing Predictive Maintenance Motor Control	DNN Decision Tree SVM Linear Naive Bayes	UCI Air Quality UCI Gas UCI EMG NASA's PCoE

tinyMLPerf Benchmark Design Choices



Big Questions	Inference
1. Benchmark definition	What is the definition of a benchmark task?

tinyMLPerf Benchmark Design Choices



Big Questions	Inference
1. Benchmark definition	What is the definition of a benchmark task?
2. Benchmark selection	Which benchmark task to select?

tinyMLPerf Benchmark Design Choices



Big Questions	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?

tinyMLPerf Benchmark Design Choices



Big Questions	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?

tinyMLPerf Benchmark Design Choices



Big Questions	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	Quantization, calibration, and/or retraining?
	Reduce result variance?



tinyMLPerf Benchmark Design Choices

Big Questions	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	Quantization, calibration, and/or retraining?
	Reduce result variance?
6. Results	Do we normalize and/or summarize results?



ML Benchmark Design Choices: Examples

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



ML Benchmark Design Choices: Examples

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



ML Benchmark Design Choices: Examples

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 and anticipate market demand



ML Benchmark Design Choices: Examples

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 and anticipate market demand
Practicality: Availability of datasets?	Feasibility: Is there a public dataset?	Good now > perfect.

tinyMLPerf Benchmark Strawperson

Task Category	Use Case	Dataset
Audio	Audio Wake Words	Speech Commands
Visual	Visual Wake Words	Google's VWW dataset
Behavioral	Anomaly Detection	Physionet, HAR, DSA, Opportunity



tinyMLPerf

Benchmarking Resource-constrained Machine Learning Systems Colby Banbury, Max Lam, Vijay Janapa Reddi, David Kanter, Amin Fazl, Xinyuan Huang, Danilo Pietro Pau and the tinyMLPerf working group Harvard University, MLPerf, Samsung Semiconductor, Inc., Cisco Systems, STMicroelectronics



Abstract

Advancements in ultra-low-power machine learning (tinyML) hardware promises to unlock an entirely new class of intelligent applications. However, the complexity and dynamics of the field obscure the measurement of progress and make application design decisions intractable. In order to enable the continued innovation, a fair, replicable and robust method of comparison is needed. Since progress is often the result of increased hardware capability, a reliable tinyML hardware benchmark is required.

To fulfill the need, we have created a community effort to extend the scope of the existing MLPerf benchmarking suite to include tinyML devices. With the help of over 75 member organizations, this group, dubbed tinyMLPerf, has begun the process of developing a benchmarking suite.

Existing Benchmarks

Existing benchmarks do not represent ML workloads or they are too large to fit on tinyML constrained processors.

BENCHMARK	ML?	POWER?	TINY?
COREMARK	x	✓	✓
MLMARK	✓	x	x
MLPERF INFERENCE	✓	✓	x
TINYML REQUIREMENTS	✓	✓	✓

Survey of tinyML Use Cases, Models, and Datasets

The landscape of tinyML use cases is large and wide. We surveyed many state of the art use cases to determine the scope of a representative tinyMLPerf benchmark.

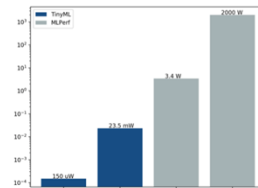
Input Type	Use Case	Model Type	Dataset
Audio	Audio Wake Words	DNN	Speech Commands
	Context Recognition	CNN	AudioSet
	Control Words Keyword Detection	RNN LSTM	ExtraSensory Freesound DCASE
Image	Visual Wake Words	DNN	Visual Wake Words
	Object Detection	CNN	CIFAR10
	Gesture Recognition	SVM	MNIST
	Object Counting Text Recognition	Decision Tree KNN Linear	ImageNet DVS128 Gesture
Physiological / Behavioral Metrics	Segmentation	DNN	Physionet
	Anomaly Detection Forecasting Activity Detection	Decision Tree SVM Linear	HAR DSA Opportunity
Industry Telemetry	Sensing	DNN	UCI Air Quality
	Predictive Maintenance Motor Control	Decision Tree SVM Linear Naive Bayes	UCI Gas UCI EMG NASA's PCoE

Challenges: Energy

An ideal tinyML benchmark would profile the energy efficiency of each system. Unfortunately, there are many challenges in fairly measuring energy usage:

- Maintaining the accuracy of energy measurement across the diverse range of processors, silicon technologies and memory architectures.
 - Determining the scope of the measurement.
 - Memories (RAM, FLASH)?
 - Peripherals? PLL?
 - Pre/Post processing? Interfaces ?
- Measuring energy consumption without significant work or alterations to the SUT.
- Preventing energy measurements from impacting the other metrics

Scope of tinyMLPerf vs. MLPerf Inference: Power Envelope



tinyML Systems consume drastically less power than traditional ML systems, yet still cover a large scope.

Challenges: Model Infancy

Despite the nascency of the field, tinyML systems are already diverse. This poses a number of challenges:

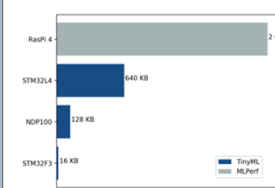
- Lack of standardization makes collecting metrics harder to formalize.
- Novel architectures have drastically different constraints and topologies.
- System requirements vary significantly across use cases.
- Performances are difficult to normalize.
- Different manufacturing technologies jeopardize comparisons with a fairly acceptable methodology.

Challenges: Memory

Memory constraints are one of the primary motivating factors for the creation of a tinyML specific benchmark. However, memory constraints add additional developmental challenges:

- Traditional benchmarks use NN models that are far too large in weights and activations.
- The overhead of the benchmark is more significant factor, pushing the need for non-intrusive inspection of key metric indicators.
- The System Under Test cannot hold the entire testing set, w/out involving host communication.
- Software (e.g. RTOS, drivers, built-in libraries) will require further discrimination.

Scope of tinyMLPerf vs. MLPerf Inference: Memory Envelope



Limited memory is a significant constraint for tinyML systems and the degree of which can vary widely.

Challenges: Processors Heterogeneity

tinyML is still a new field. It creates an opportunity to foster growth through community efforts (with industry support) but also poses a number of challenges for developing a robust benchmark that features industry acceptance and consensus:

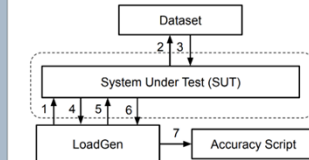
- Widely accepted tinyML neural network models.
- Large open source tinyML datasets.
- Frameworks are still evolving and few de-facto industry standards have become popular therefore model portability/interoperability is evolving:
 - e.g. tensorflow, keras, pytorch, mxnet, caffe etc, associated interoperable file formats (e.g. tflite, keras, onnx, nnet).

Use Cases Selected for v0.1

The criteria for the preliminary selection was to select three use cases that represented the scope of tinyML in terms of input type, size, neural network model type, and maturity. Model selection is still in progress.

Use Case	Dataset
Audio Wake Words	Speech Commands
Visual Wake Words	Google's VWW dataset
Anomaly Detection	Physionet, HAR, DSA, Opportunity

LoadGen and System Under Test



(workflow is subject to change)

Working Group Member Organizations



Join Us!

tinyMLPerf
Help us create
a tinyML Benchmark!



Join <https://groups.google.com/forum/#!members/mlperf-tiny>

Summary

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



Machine learning performance benchmarks are a number that improve training has high variance, with the same bin overcomes these performance and

30 Oct 2019

arXiv:1910.01500v2 [cs.LG]

MLPERF TRAINING BENCHMARK

Peter Mattson¹ Christine Cheng² Cody Coleman³ Greg Diamos⁴ Paulius Micikevicius⁵ David Patterson^{1,6} Hanlin Tang⁷ Gu-Yeon Wei⁸ Peter Bailis⁹ Victor Bittorf⁷ David Brooks⁷ Dehao Chen¹ Debojyoti Dutta⁸ Uditi Gupta⁷ Kim Hazelwood⁷ Andrew Hoek¹⁰ Xiyan Huang⁷ Bill Jia⁹ Daniel Kang⁷ David Kanter¹¹ Naven Kumar¹ Jeffery Liao¹² Guokai Ma⁷ Deepak Narayanan⁷ Tayo Oguntebi¹ Gennady Pekhimenko¹³ Lillian Pentecost⁷ Vijay Janapa Reddi⁷ Taylor Robie¹ Tom St. John¹⁴ Carole-Jean Wu⁹ Lingjie Xu¹⁵ Cliff Young⁷ Matei Zaharia¹

Machine learning performance benchmarks are a number that improve training has high variance, with the same bin overcomes these performance and

4 Nov 2019

arXiv:submit/2913405 [cs.LG]

MLPERF INFERENCE BENCHMARK

Vijay Janapa Reddi¹ Christine Cheng² David Kanter³ Peter Mattson⁴ Guenther Schumling⁵ Carole-Jean Wu⁶ Brian Anderson⁷ Maximilien Breughe⁷ Mark Charlebois⁸ William Chou⁹ Ramesh Chukka⁷ Cody Coleman¹⁰ Sam Davis¹⁰ Greg Diamos¹⁰ Jared Duke¹⁰ Dave Fick¹⁰ J. Scott Gardner¹³ Huy Habara¹⁴ Sachin Idnani¹⁷ Thomas B. Jablin¹⁸ Jeff Jiao¹⁵ Tom St. John¹⁶ Pankaj Kanwar¹ David Lee¹⁷ Jeffery Liao¹⁸ Anton Likhomstov¹⁹ Francisco Massa⁷ Paulius Micikevicius⁷ Colin Osborne²⁰ Gennady Pekhimenko²¹ Arun Tejusve Raghunath Rajan⁷ Dilip Sequeira²² Ashish Srivastava²³ Fei Sun²³ Hanlin Tang⁷ Michael Thomson²⁴ Frank Wei²⁵ Ephrem Wu²² Lingjie Xu²² Koichi Yamada² Bing Yu¹¹ George Yuan¹¹ Aaron Zhang²⁶ Yuchen Zhou²⁶

Machine-learning (ML) hardware and software system demand is burgeoning. Driven by ML applications, the number of different ML inference systems has exploded. Over 100 organizations are building ML inference chips, and the systems that incorporate existing models span at least three orders of magnitude in power consumption and four orders of magnitude in performance; they range from embedded devices to data-center solutions. Fueling the hardware are a dozen or more software frameworks and libraries. The myriad combinations of ML hardware and ML software make assessing ML-system performance in an architecture-neutral, representative, and reproducible manner challenging. There is a clear need for industry-wide standard ML benchmarking and evaluation criteria. MLPerf Inference answers that call. Driven by more than 30 organizations as well as more than 200 ML engineers and practitioners, MLPerf implements a set of rules and practices to ensure comparability across systems with wildly differing architectures. In this paper, we present the method and design principles of the initial MLPerf Inference release. The first call for submissions garnered more than 600 inference-performance measurements from 14 organizations, representing 44 systems that show a wide range of capabilities.

1 INTRODUCTION

Machine learning (ML) powers a variety of applications from computer vision (He et al., 2016; Goodfellow et al., 2014; Liu et al., 2016; Krizhevsky et al., 2012) and natural-language processing (Vaswani et al., 2017; Devlin et al., 2018) to self-driving cars (Xu et al., 2018; Badrinarayanan et al., 2017) and autonomous robotics (Levine et al., 2018). These applications are deployed at large scale and require substantial investment to optimize inference performance. Although training of ML models has been a development bottleneck and a considerable expense (Amodei & Hernandez, 2018), inference has become a critical workload, since models can serve as many as 200 trillion queries and perform over 6 billion translations a day (Lee et al., 2019b).

¹Harvard University ²Intel ³Real World Insights ⁴Google ⁵Microsoft ⁶NVIDIA ⁷Qualcomm ⁸Stanford University ⁹Myrtle ¹⁰Landing AI ¹¹Mythic ¹²Advantage Engineering ¹³Hakana Labs ¹⁴Alibaba T-Head ¹⁵Intel ¹⁶MediaTek ¹⁷Synopsys ¹⁸dividiti ¹⁹Arm ²⁰University of Toronto ²¹Xilinx ²²Alibaba (was Facebook) ²³Centaur Technology ²⁴Alibaba Cloud ²⁵General Motors. MLPerf Inference is the product of (1) individuals that led the benchmarking effort from the various organizations and (2) submitters that produced the first set of benchmark results. It takes both groups to have a successful industry benchmark. We credit the submitters and their organizations in the acknowledgments. Send correspondence email to vjr@eecs.harvard.edu and see mlperf.org.

To address these growing computational demands, hardware, software, and system developers have focused on inference performance for a variety of use cases by designing optimized ML hardware and software systems. Estimates indicate that over 100 companies are producing or are on the verge of producing optimized inference chips. By comparison, only about 20 companies target training.

Each system takes a unique approach to inference and presents a trade-off between latency, throughput, power, and model quality. For example, quantization and reduced precision are powerful techniques for improving inference latency, throughput, and power efficiency at the expense of accuracy (Han et al., 2015; 2016). After training with floating-point numbers, compressing model weights enables better performance by decreasing memory-bandwidth requirements and increasing computational throughput (e.g., by using wider vectors). Similarly, many weights can be removed to boost sparsity, which can reduce the memory footprint and the number of operations (Han et al., 2015; Molchanov et al., 2016; Li et al., 2016). Support for these techniques varies among systems, however, and these optimizations can drastically reduce final model quality. Hence, the field needs an ML inference benchmark that can quantify these trade-offs in an architecturally neutral, representative, and reproducible manner.



Join us!

Google group: <https://groups.google.com/forum/#!members/mlperf-tiny>

vj@eecs.harvard.edu
cbanbury@g.harvard.edu

Copyright Notice

The presentation(s) in this publication comprise the proceedings of tinyML[®] Summit 2020. The content reflects the opinion of the authors and their respective companies. This version of the presentation may differ from the version that was presented at the tinyML Summit. The inclusion of presentations in this publication does not constitute an endorsement by tinyML Foundation or the sponsors.

There is no copyright protection claimed by this publication. However, each presentation is the work of the authors and their respective companies and may contain copyrighted material. As such, it is strongly encouraged that any use reflect proper acknowledgement to the appropriate source. Any questions regarding the use of any materials presented should be directed to the author(s) or their companies.

tinyML is a registered trademark of the tinyML Foundation.

www.tinyML.org