

---

# MLPerf Tiny Benchmark

---

Colby Banbury\* Vijay Janapa Reddi\* Peter Torelli† Jeremy Holleman‡|| Nat Jeffries§

Csaba Kiraly¶ Pietro Montino\* David Kanter\*\* Sebastian Ahmed†† Danilo Pau‡‡

Urmish Thakker<sup>I</sup> Antonio Torrini<sup>II</sup> Peter Warden<sup>§</sup> Jay Cordaro <sup>‡</sup>Giuseppe Di Guglielmo<sup>III</sup>

Javier Duarte<sup>IV</sup> Stephen Gibellini<sup>‡</sup> Videet Parekh<sup>V</sup> Honson Tran<sup>V</sup> Nhan Tran<sup>VI</sup>

Niu Wenxu<sup>VII</sup> Xu Xuesong<sup>VII</sup>

## Abstract

1       Advancements in ultra-low-power *tiny* machine learning (TinyML) systems  
2       promise to unlock an entirely new class of smart applications. However, con-  
3       tinued progress is limited by the lack of a widely accepted and easily reproducible  
4       benchmark for these systems. To meet this need, we present MLPerf Tiny, the  
5       first industry-standard benchmark suite for ultra-low-power tiny machine learning  
6       systems. The benchmark suite is the collaborative effort of more than 50 orga-  
7       nizations from industry and academia and reflects the needs of the community.  
8       MLPerf Tiny measures the accuracy, latency, and energy of machine learning  
9       inference to properly evaluate the tradeoffs between systems. Additionally, MLPerf  
10       Tiny implements a modular design that enables benchmark submitters to show the  
11       benefits of their product, regardless of where it falls on the ML deployment stack,  
12       in a fair and reproducible manner. The suite features four benchmarks: keyword  
13       spotting, visual wake words, image classification, and anomaly detection.

14

## 15   1 Introduction

16   Machine learning (ML) inference on the edge is an increasingly attractive prospect due to its potential  
17   for increasing energy efficiency [4], privacy, responsiveness, and autonomy of edge devices. Thus  
18   far, the field edge ML has predominantly focused on mobile inference, but in recent years, there  
19   have been major strides towards expanding the scope of edge ML to ultra-low-power devices.  
20   The field, known as “TinyML” [1], achieves ML inference under a milliWatt, and thereby breaks  
21   the traditional power barrier preventing widely distributed machine intelligence. By performing  
22   inference on-device, and near-sensor, TinyML enables greater responsiveness and privacy while  
23   avoiding the energy cost associated with wireless communication, which at this scale is far higher  
24   than that of compute [5]. Furthermore, the efficiency of TinyML enables a class of smart, battery-  
25   powered, always-on applications that can revolutionize the real-time collection and processing of  
26   data. Deploying advanced ML applications at this scale requires the co-optimization of each layer of  
27   the ML deployment stack to achieve the maximum efficiency. Due to this complex optimization, the

---

\*Harvard University, †EEMBC, ‡Syntiant ||UNC Charlotte §Google ¶Digital Catapult \*VoiceMed  
\*\*MLCommons ††Qualcomm ‡‡STMicroelectronics <sup>I</sup>SambaNova Systems <sup>II</sup>Silicon Labs <sup>III</sup>Columbia <sup>IV</sup>UCSD  
<sup>V</sup>Latent AI <sup>VI</sup>Fermilab <sup>VII</sup>Peng Cheng Labs

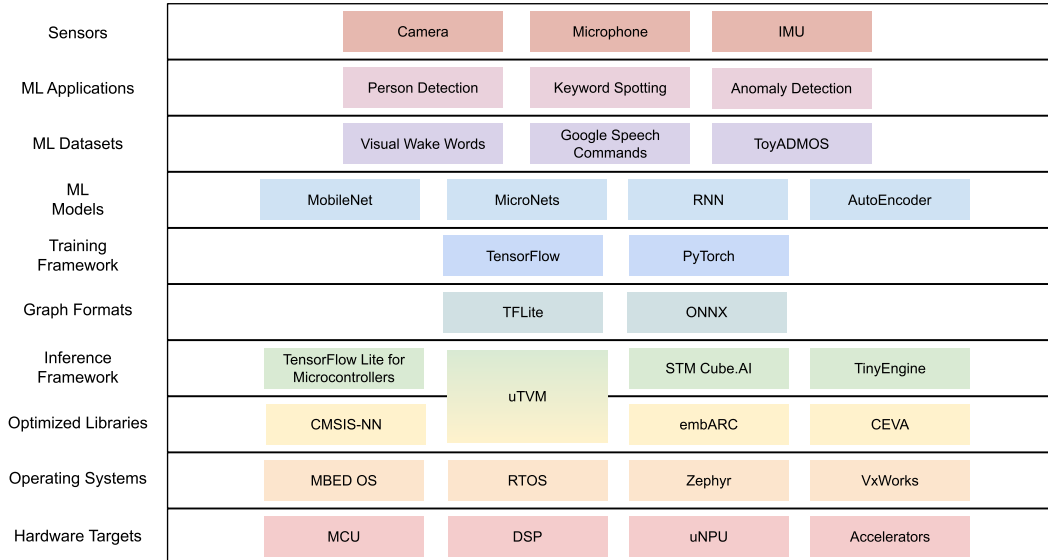


Figure 1: Summary of the Tiny Machine Learning Stack. There is diversity at every level which makes standardization for benchmarking challenging.

28 direct comparison of solutions is challenging and the impact of individual optimizations is difficult to  
 29 measure. In order to enable the continued innovation, a fair and reliable method of comparison is  
 30 needed.

31 In this paper, we present MLPerf Tiny, an open-source benchmark suite for TinyML systems. The  
 32 MLPerf Tiny inference benchmark suite provides a set of four standard benchmarks, selected by more  
 33 than 50 organizations in academia and industry. In order to capture the tradeoffs inherent to TinyML,  
 34 the benchmark suite measures latency, energy, and accuracy. MLPerf Tiny is designed with flexibility  
 35 and modularity in mind to support hardware and software users alike and provides complete reference  
 36 implementations to act as open-source community baselines.

## 37 2 Challenges

38 TinyML systems present a number of unique challenges to the design of a performance benchmark that  
 39 can be used to measure and quantify performance differences between various systems systematically.  
 40 We discuss the four primary obstacles and postulate how they might be overcome.

41 **Low Power** Power consumption is one of the defining features of TinyML systems. Therefore, a  
 42 useful benchmark should profile the energy efficiency of each device. However, there are many  
 43 challenges in fairly measuring energy consumption. TinyML devices can consume drastically  
 44 different amounts of power, which makes maintaining accuracy across the range of devices difficult.  
 45 Also, determining what falls under the scope of the power measurement is difficult to determine when  
 46 data paths and pre-processing steps can vary significantly between devices. Other factors like chip  
 47 peripherals and underlying firmware can impact the measurements.

48 **Limited Memory** While traditional ML systems like smartphones cope with resource constraints in  
 49 the order of a few GBs, tinyML systems are typically coping with resources that are two orders of  
 50 magnitude smaller. Traditional ML benchmarks use inference models that have drastically higher  
 51 peak memory requirements (in the order of gigabytes) than TinyML devices can provide. This also  
 52 complicates the deployment of a benchmarking suite as any overhead can make the benchmark too  
 53 big to fit. Individual benchmarks must also cover a wide range of devices; therefore, multiple levels  
 54 of quantization and precision should be represented in the benchmarking suite. Finally, a variety of  
 55 benchmarks should be chosen such that the diversity of the field is supported.

56 **Hardware Heterogeneity** Despite its nascency, TinyML systems are already diverse in their perfor-  
 57 mance, power, and capabilities. Devices range from general-purpose MCUs to novel architectures,  
 58 like in event-based neural processors [2] or memory compute citekim20191. This heterogeneity poses

59 a number of challenges as the system under test (SUT) will not necessarily include otherwise standard  
60 features, like a system clock or debug interface. Furthermore, creating a standard interface while  
61 minimizing porting effort is a key challenge. Today’s state-of-the-art benchmarks are not designed  
62 to handle these challenges readily. They need reengineering to be flexible enough to handle the  
63 hardware heterogeneity that is commonplace in the TinyML ecosystem.

64 **Software Heterogeneity** TinyML systems are often tightly coupled with their inference stack and  
65 deployment tools. To achieve the highest efficiency, users often develop their own tool chains to  
66 optimally deploy and execute a model on their hardware systems. This becomes increasingly critical  
67 on systems with multiple compute units, like accelerators and DSPs. This poses a challenge when  
68 designing a benchmark for these systems because any restriction, posed by the benchmark, on the  
69 inference stack, would negatively impact performance on these systems and result in unrepresentative  
70 results. Therefore we must balance optimality with portability, and comparability with representative-  
71 ness. A TinyML benchmark should support many options for model deployment and not impose any  
72 restrictions that may unfairly impact users with a different deployment stack.

73 **Cross-product** Figure 1 illustrates the diversity of options at every level in the TinyML stack. Each  
74 option and each layer has its own impact on performance and TinyML software users can provide an  
75 improvement to the overall system at any layer. A TinyML benchmark should enable these users to  
76 demonstrate the performance benefits of their solution in a controlled setting.

### 77 3 Related Work

78 There are a few ML related hardware benchmarks, however, none that accurately represent the  
79 performance of TinyML workloads on tiny hardware.

80 **EEMBC(r)’s CoreMark(r)** benchmark [9] has become the standard benchmark for MCU-class  
81 devices due to its ease of implementation and use of real algorithms. Yet, CoreMark does not profile  
82 full programs, nor does it accurately represent machine learning inference workloads.

83 **EEMBC’s MLMark(r)** benchmark [19] addresses these issues by using actual ML inference work-  
84 loads. However, the supported models are far too large for MCU-class devices and are not representa-  
85 tive of TinyML workloads. They require far too much memory (GBs) and have significant runtimes.  
86 Additionally, while CoreMark supports power measurements with the EEMBC ULPMark(tm)-CM  
87 benchmark, MLMark does not, which is critical for a TinyML benchmark.

88 **MLPerf**, a community-driven benchmarking effort, has recently introduced a benchmarking suite for  
89 ML inference [18] and has plans to add power measurements. However, much like MLMark, the  
90 current MLPerf inference benchmark precludes MCUs and other resource-constrained platforms due  
91 to a lack of small benchmarks and compatible implementations. As Table 1 summarizes, there is a  
92 clear and distinct need for a TinyML benchmark that caters to the unique needs of ML workloads,  
93 makes power a first-class citizen and prescribes a methodology that suits TinyML.

### 94 4 Benchmarks

95 All machine learning benchmarks fall somewhere on the continuum between low level and application  
96 level evaluation. Low level benchmarks target kernels that are core to many ML workloads, like  
97 matrix multiply, but they gloss critical elements like memory bandwidth or model level optimizations.  
98 On the other hand, application level benchmarks can obscure the target of the benchmark behind  
99 other stages of the application pipeline. MLPerf Tiny specifically targets model inference and does  
100 not include pre- or post-processing in the measurement window. Table 1 shows the v0.5 benchmarks.

101 In this section, we describe the benchmarks that form the MLPerf Tiny benchmark suite. Each  
102 benchmark targets a specific usecase and specifies a dataset, model, and quality target. Additionally,  
103 each benchmark has a reference implementation that includes training scripts, pre-trained models,  
104 and C code implementations. The reference implementations run the reference models in the TFLite  
105 format using TFLite for Microcontrollers (TFLM) [7] on the NUCLEO-L4R5ZI board. A known-  
106 good snapshot of the TFLM runtime is used to ensure stability, and the reference implementation is  
107 built using a bare-metal MBED project with the GCC-ARM toolchain. The reference implementations  
108 are open source and available on GitHub at <https://github.com/mlcommons/tiny>

Use Case	Dataset (Input Size)	Model (TFLite Model Size)	Quality Target (Metric)
Keyword Spotting	Speech Commands (49x10)	DS-CNN (52.5 KB)	90% (Top-1)
Visual Wake Words	VWW Dataset (96x96)	MobileNetV1 (325 KB)	80% (Top-1)
Image Classification	CIFAR10 (32x32)	ResNet (96 KB)	85% (Top-1)
Anomaly Detection	ToyADMOS (5*128)	FC-AutoEncoder (270 KB)	.85 (AUC)

Table 1: MLPerf Tiny v0.5 Inference Benchmarks.

#### 109 4.1 Visual Wake Words

110 **Rational** Tiny image processing models are becoming increasingly widespread, primarily for simple  
 111 image classification tasks. The Visual Wakewords challenge [6] tasked submitters with detecting  
 112 whether at least one person is in an image. This task is directly relevant to smart doorbell and  
 113 occupancy applications, and the reference network provided as part of the challenge fits on most  
 114 32-bit embedded microcontrollers.

115 **Dataset** The Visual Wakewords Challenge uses the MSCOCO 2014 dataset [15] as the training,  
 116 validation and test datasets for all person detection models. The dataset is preprocessed to train on  
 117 image which contain at least one person occupying more than 2.5% of the source image. Images are  
 118 also resized to 96x96 for model training.

119 **Model** The Visual Wakewords model is a MobilenetV1 [11] which takes 96x96 input images with an  
 120 alpha of 0.25 and two output classes (person and no person). The TFLM model is 325KB in size.

121 **Quality Target** Based on validation and training accuracy, the model reaches about 86% accuracy  
 122 across the preprocessed MSCOCO 2014 test dataset. In order to accommodate changes in accuracy  
 123 due to quantization and rounding differences between platforms, submissions to the closed category  
 124 should reach at least 80% accuracy across the same dataset.

#### 125 4.2 Image Classification

126 **Rational** Novel machine vision techniques offer a cross-industry potential for breakthroughs and  
 127 advances in autonomous and low power embedded solutions. Compact vision systems performing  
 128 image classification at low cost, high efficiency, low latency and high performances are pervading  
 129 manufacturing, IoT devices, and autonomous agents and vehicles. New hardware platforms, algo-  
 130 rithms and development tools form a wide variety of deep embedded vision systems requiring a point  
 131 of reference for scientific and industrial evaluation.

132 **Dataset** CIFAR-10 [14] is a labeled subset of the 80 Million Tiny Images dataset [20]. The low  
 133 resolution of the images make CIFAR-10 the most suitable source of data for training tiny image  
 134 classification models. It consists of 60000 32x32x3 RGB images, with 6000 images per class. The  
 135 10 different classes represent airplanes, cars, birds, cats, deers, dogs, frogs, horses, ships and trucks.  
 136 The dataset is divided into five training batches and one testing batch, each with 10000 images. A  
 137 significant amount of prior work in TinyML has used CIFAR-10 as a target dataset [8], by continuing  
 138 this trend, we create a point of reference in the benchmark suite that can be used to relate future  
 139 results to historical data points.

140 **Model** The Image Classification model is a customized ResNetv1 [10] which takes 32x32x3 input  
 141 images and outputs a probability vector of size 10. The custom model is made of fewer residual  
 142 stacks than the official ResNet: three compared to four. Moreover, the first convolutional layer is not  
 143 followed by the pooling layer due to the low resolution of the input data. The number of convolution  
 144 filters and the convolution strides dimension are lower as well compared to the official ResNet. The  
 145 TFLite model for IC is 96KB in size and fits on most 32-bit embedded microcontrollers.

146 **Quality Target** A set of 200 images from the CIFAR-10 test set are selected to evaluate the perfor-  
 147 mances of the IC reference implementation. The performance evaluation test is performed with the  
 148 benchmark framework software, i.e. the runner (see Appendix). The model reaches 86.5% accuracy

149 across the 200 testing raw images. To accommodate minor differences in quantization and various  
150 other optimizations, we set the quality target to 85% top-1 accuracy.

### 151 4.3 Keyword Spotting

152 **Rational** Recognition of specific words and brief phrases, known as keyword spotting, is one  
153 of the primary use cases of ultra-low-power machine learning. Voice is an important mode of  
154 human-machine interaction. Wakeword detection, a specific case of keyword spotting wherein a  
155 detector continuously monitors for a specific word or phrase (e.g. “Hey Siri”™, “Hey Google”™,  
156 “Alexa”™) in order to enable a larger processor, requires continuous operation and therefore low  
157 power consumption. For example, a 100 mA current drain would deplete a typical phone battery in  
158 one day without any other activity. Command phrases such as “volume up” or “turn left” are the  
159 provide a simple and natural interface to embedded devices. Both cases require low latency and low  
160 power consumption and must typically run on small, low-cost devices.

161 **Dataset** We used the Speech Commands v2 dataset[21], a collection of 105,829 utterances collected  
162 from 2,618 speakers with a variety of accents. It is freely available for download under the Creative  
163 Commons BY license. The dataset contains 30 words and a collection of background noises and is  
164 divided into training, validation, and test subsets such that any individual speaker only appears in one  
165 subset. Following the typical usage of this dataset, we used 10 words and combined the background  
166 noise with the remaining 20 words to approximate an open set labeled “unknown,” which, along  
167 with “silence”, results in 12 output classes. This approach exercises a combination command-phrase  
168 systems’ need for multiple words with wakeword systems’ need for an open set of background noise.

169 **Model** For the closed division model, we used the small depthwise-separable CNN described in [22].  
170 We chose this model because with 38.6K parameters, it fits within the available memory of most  
171 microcontrollers and similarly-scaled devices while achieving accuracy of 92.2% in our experiments.  
172 The model utilizes standard layers that can be expected of most neural network hardware. For  
173 this version of the benchmark, we opted to exclude feature extraction from the measurement and  
174 provide three choices of pre-computed features for the open division. Feature extraction is typically  
175 a small fraction of the overall compute cost, so the impact on the measurements is minor. While  
176 the pre-selected feature choices somewhat limit innovation in the overall KWS system, the features  
177 chosen represent the features most commonly used in KWS systems.

178 **Quality Target** To evaluate accuracy on the device under test, we randomly selected 1000 utterances  
179 from the speech commands. The quantized reference model demonstrated 91.6% accuracy on the  
180 full test set and 91.7% on the 1000-utterance subset. To allow for slight variations in quantization  
181 strategies, we set an accuracy requirement of 90%.

### 182 4.4 Anomaly Detection

183 **Rational** Anomaly detection is broadly the task of separating normal samples from anomalous  
184 samples and has applications in a large number of fields. It’s unsupervised variant, which is used in  
185 our benchmark, is of particular importance to industrial use cases such as early detection of machine  
186 anomalies, where failure types are many and failure events can be rare, and thus data is often only  
187 available on normal operation for training. Anomalies can be observed in various data sources (or  
188 combinations thereof), including power, temperature, vibration, with the most intuitive being the  
189 detection of anomalies based on sound captured by microphones. The most important distinctive  
190 features of this benchmark are the use of unsupervised learning and an autoencoder (AE).

191 **Dataset** For this benchmark we have selected to use the dataset from the DCASE2020 [13] compe-  
192 tition, which is itself a combination of two publicly available sources: the ToyADMOS [12], and  
193 the MIMII [17] datasets. The DCASE dataset contains data for six machine types (slide rail, fan,  
194 pump, valve, toy-car, toy-conveyor), and competition rules allowed for the training of separate models  
195 for each. Benchmarking would not benefit from this type of complexity, thus we choose to use  
196 the Toy-car machine type only. For training, normal sound samples of seven different toy-cars are  
197 provided, each having 1000 samples mixing the machine sound with environmental noise.

198 **Model** While several models are openly available from the competition, the competition itself did not  
199 focus on TinyML, and only a few of these fit our target model size. Among these, we have selected  
200 the AE based model also used in the reference implementation of DCASE2020 for multiple reasons.  
201 First, it is the reference for a large part of the literature around the audio anomaly detection problem

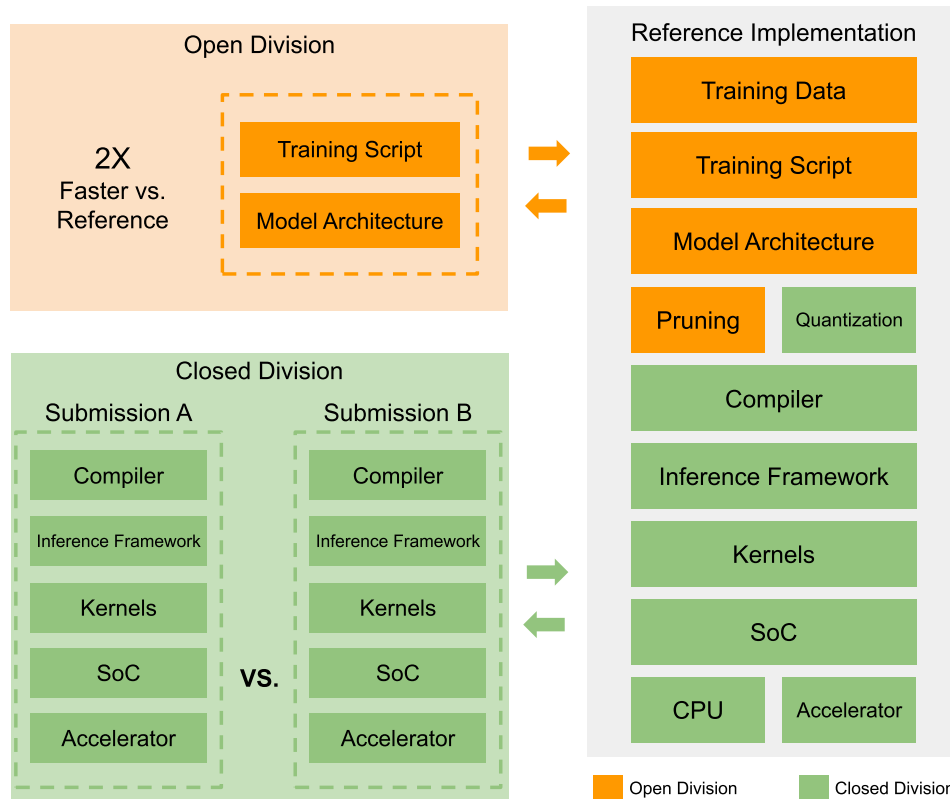


Figure 2: The modular design of MLPerf Tiny enables both the direct comparison of solutions and the demonstration of an improvement over the reference. The reference implementations are fully implemented solutions that allow individual components to be swapped out. The components in green can be modified in either division and the orange components can only be modified in the open division. The reference implementations also act as the baseline the results.

202 and thus a well known baseline. Second, it enhances the benchmark suite with a new model type  
 203 based entirely on FC layers. The model has input and output sizes of 640. Both the encoder and  
 204 decoder are made of four 128 unit FC layers with BatchNorm applied during the training and with  
 205 ReLU activation, while the bottleneck layer is of size 8. The model itself is not applied directly on  
 206 the 10 second audio. The audio is pre-processed into a log-mel-spectrogram with 128 bands and a  
 207 frame size of 32 ms. Then, the model is used repeatedly over a sliding window of five frames (hence  
 208 the 640 input size), and the MSE of the resulting reconstruction error is averaged over the central 6.4  
 209 second part of the spectrogram providing an overall anomaly score.

210 **Quality Target** The output of the autoencoder is an anomaly score, which assumes a further threshold  
 211 to be set before the binary normal/anomalous decision can be made. For this reason the parameterless  
 212 AUC-ROC (Area Under The Curve, Receiver Operating Characteristics) fits better for the quality  
 213 evaluation for this problem than metrics requiring the selection of a threshold. To ensure models  
 214 have generalization characteristics, we have selected an evaluation dataset composed of normal and  
 215 anomalous sounds from four different machines totalling 248 samples. On this set, the fp32 version  
 216 of the reference model achieves an AUC of 0.88 while the AUC after quantization is 0.86. Based on  
 217 these two numbers, we've set the threshold for the benchmark to AUC 0.85.

## 218 5 Run Rules

219 In this section we describe the modular design of the benchmark harness, as well as the closed and  
 220 open benchmark divisions, and the overall run rules of the benchmark suite to ensure reproducibility.

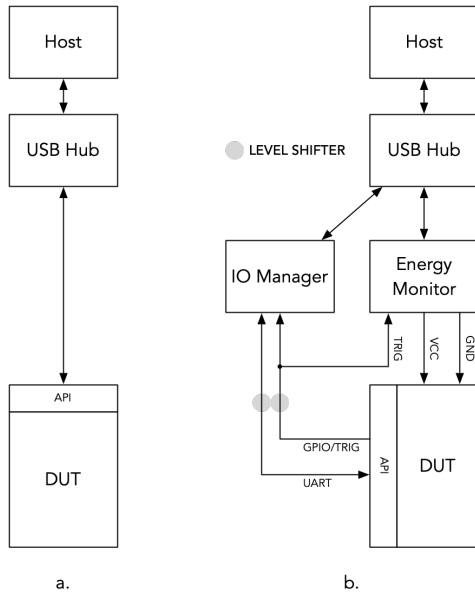


Figure 3: The two configuration modes of the benchmark framework for (a.) latency and accuracy measurement, or (b.) energy measurement.

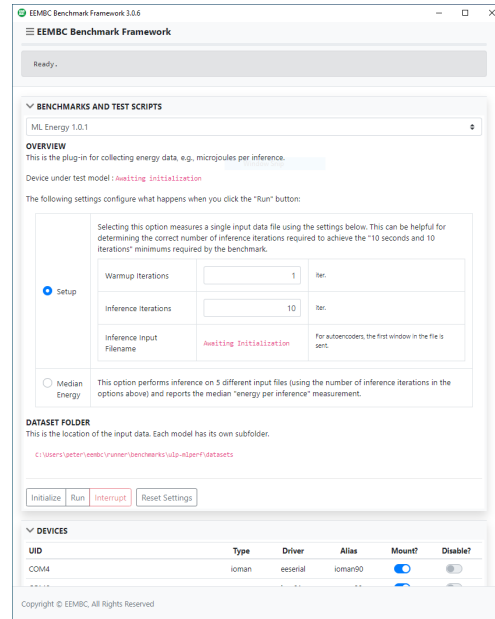


Figure 4: The graphical user interface (GUI) for the benchmark runner.

## 221 5.1 Modular Design

222 TinyML applications often require cross stack optimization to meet the tight constraints. Hardware and  
 223 software users can provide value by targeting specific components of the pipeline, like quantization,  
 224 or offer end-to-end solutions. Therefore, to enable users to demonstrate the competitive advantage of  
 225 their specific contribution, we employ a modular approach to benchmark design. Each benchmark in  
 226 the suite has a reference implementation that contains everything from training scripts to a reference  
 227 hardware platform. This reference implementation not only provides a baseline result, it can be  
 228 modified by a submitter wishing to show the performance of a single hardware or software component.  
 229 Figure 2 illustrates the modular components of a MLPerf Tiny reference implementation.

## 230 5.2 Closed and Open Divisions

231 MLPerf Tiny has two divisions for submitting results: a stricter closed division and a more flexible  
 232 open division. A submitting organization can submit to either or both divisions. This two division  
 233 design allows the benchmark to balance comparability and flexibility. Figure 2 illustrates which  
 234 components of the reference benchmark can be modified in which division. If only the components  
 235 shaded in green are modified then the submitter can submit to the open division. If any of the  
 236 components shaded in orange are modified then the submitter must submit to the open division.

237 The closed division enables a more direct comparison of systems. Submitters must use the same  
 238 models, datasets, and quality targets as the reference implementation. The closed division permits post  
 239 training quantization using the provided calibration datasets, but prohibits any retraining or weight  
 240 replacement. Figure 2 demonstrates how submissions to the closed division allow the comparison of  
 241 an inference stack to another in a controlled setting.

242 The open division is designed to broaden the scope of the benchmark and allow submitters to  
 243 demonstrate improvements to performance, energy, or accuracy at any stage in the machine learning  
 244 pipeline. The open division allow submitters to change the model, training scripts, and dataset. A  
 245 submission to the open division will still use the same test dataset to benchmark the accuracy but is  
 246 not required to meet the accuracy threshold, which allow submitters to demonstrate tradeoffs between  
 247 accuracy, latency, and energy consumption. Each submission to the open division must document  
 248 how it deviates from the reference implementation. Figure 2 demonstrates how submissions to the  
 249 open division allow users to demonstrate the specific value add of their product.

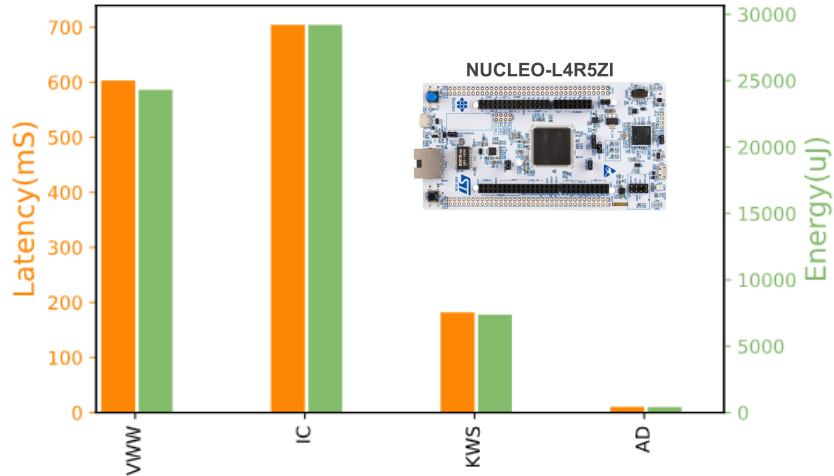


Figure 5: The energy and latency results of the reference implementations. Each reference implementation was run on the NUCLEO-L4R5ZI board which is shown in the top right.

### 250 5.3 Measurement Procedure

251 The platforms targeted by TinyML typically do not contain the resources required to run a complete  
 252 benchmark locally, e.g., file IO, standard input and output, interactivity, etc. As a result, a benchmark  
 253 framework (Figure 3) is required to facilitate controlling and measuring the device under test (DUT).  
 254 The benchmark framework used to implement the MLPerf Tiny benchmark is based on EEMBC’s  
 255 software development platform. The resulting program coordinates execution of the benchmark’s  
 256 behavioral model among the various components in the system as described in the Appendix. The  
 257 framework’s measurement procedure is as follows:

- 258 1. Latency – Perform this five times: download an input stimulus, load the input tensor  
 259 (converting the data as needed), run inference for a minimum of 10 seconds and 10 iterations,  
 260 measure the inferences per second (IPS); report the median IPS of the five runs as the score.
- 261 2. Accuracy – Perform a single inference on the entire set of validation inputs, and collect the  
 262 output tensor probabilities. The number of inputs vary, depending on the model. From the  
 263 individual results, compute the Top-1 percent and the AUC. Each model has a minimum  
 264 accuracy that must be met for the score to be valid.
- 265 3. Energy – Identical to latency, but in addition to measuring the number of inferences per  
 266 second, measure the total energy used in the timing window and compute micro-Joules per  
 267 inference. The same method of taking the median of five measurements is used.

268 Results are stored in a folder and can be reloaded again. An energy viewer allows the user to examine  
 269 the energy trace. Figure 5 illustrates the reference implementation results. The four benchmarks  
 270 cover a wide scope in terms of latency and energy and each reference meets the minimum accuracy.

## 271 6 Benchmark Assessment

272 In this section we assess the success of the benchmark at meeting the needs of the community.  
 273 The suite is designed to provide standardization and compatibility while enabling a diverse set of  
 274 organizations to submit results.

### 275 6.1 Submissions

276 The MLPerf Tiny benchmark suite accepts results from submitting organizations twice a year. The  
 277 submitting organizations must implement the benchmarks on their hardware/software stack and  
 278 submit their results and implementations at the deadline. All results are then transparently peer-  
 279 reviewed by the group of submitters and a review committee to ensure they conform to the rules.



Division	Dataset	Training	Model	Numerics	Framework	Hardware	Demonstrates
Closed	X	X	X	Int-8 PTQ	TFLite Micro	ARM MCU	Baseline Results
Closed	X	X	X	Int-8 PTQ	TFLite Micro	RISC-V MCU	RISC-V performance
Closed	X	X	X	FP-32 & Int-8 PTQ	LEIP Framework	RasPi 4	Hardware agnostic optimization
Closed	X	X	X	Int-8 PTQ	Syntiant TDK	NN Accelerator	Extreme hardware efficiency
Open	X	QKeras	✓	Int-6/8 QAT	HLS4ML	FPGA	Rapid end-to-end development

Table 2: Summary of the MLPerf Tiny v0.5 round of submissions. The submissions are diverse and illustrate the ability of the benchmark to demonstrate an advantage in a variety of ways. The X indicates no modification was made from the reference. The check mark indicates a modification. PTQ refers to post training quantization and QAT refers to quantization aware training.

280 The first round of submissions to MLPerf Tiny (which took place in June 2021) are summarized in  
281 Table 2. The results from the first round submissions are currently available at <https://mlcommons.org/en/inference-tiny-05/>. Each submission is publicly available on GitHub with instructions  
282 on how to reproduce each result at [https://github.com/mlcommons/tiny\\_results\\_v0.5](https://github.com/mlcommons/tiny_results_v0.5).  
283

## 284 6.2 Insights

285 The v0.5 submissions were diverse and demonstrated the desired flexibility of the benchmark suite.  
286 There were submissions to the open and closed divisions, as well as from hardware and software  
287 vendors. Each submitter had a specific element of their submission they wished to demonstrate and  
288 the benchmark suite was able to accommodate this diverse set of goals due to its modular design.

289 The submissions indicate general trends in TinyML. For example, the most common **numerical**  
290 **format** is 8-bit integer for inference as it offers a performance boost with little impact to the model  
291 accuracy. ML **frameworks** range from open source interpreters (TFLite Micro) to hardware specific  
292 inference compilers, indicating that there is still often a trade-off between optimization and portability.  
293 There were results collected on a wide variety of **hardware platforms**, including MCUs, accelerators,  
294 and FPGAs. The re-configurable hardware (FPGA) is able to utilize variable precision models for  
295 increased performance. The power consumption of these platforms ranged from  $\mu$ Watts to Watts.

296 More specifically, one organization used the benchmark to demonstrate that their SDK developer  
297 tools are hardware agnostic, easy to use and uniquely designed to optimize neural network inferences  
298 for compute, energy and memory, while preserving algorithmic accuracy. Another organization,  
299 a hardware vendor, used the benchmark to demonstrate outstanding performance on their NN  
300 accelerator, enabled by their novel microarchitecture design and optimized datapath, which avoids  
301 wait-states and maintains high computational efficiency. An academic research institute used the  
302 benchmark to demonstrate AI compute capability and potential applications for RISC-V-based AI  
303 micro-controllers. RISC-V is a free, open standard instruction set architecture (ISA) [3]. Finally, a  
304 multi-disciplinary team of scientists and engineers showcased their “hls4ml” (high-level synthesis  
305 for ML) open-source workflow. It was designed to enable researchers and engineers to codesign  
306 optimized neural networks for efficient dataflow architectures on a multitude of accelerator hardware  
307 platforms. Originally developed for the Large Hadron Collider to make ultrafast sub-microsecond ML  
308 inference, the workflow aims to serve the wider ML community to accelerate both the design process  
309 and neural network implementation across a broad range from low-power to high-performance  
310 devices.

311 These results indicate a snapshot of this emerging field but future submission rounds can be used  
312 to indicate the evolution of TinyML over time. For instance, despite a general trend in AI towards  
313 data-centric design, none of the submissions in the first round modified the training dataset. While  
314 we anticipate this will shift in future versions, TinyML design is still largely focused on models,  
315 frameworks, and hardware. To this end, as demonstrated, the benchmark meets a variety of needs.

## 316 7 Impact

317 The benchmarks have already acted as a standard set of tasks for TinyML research [4] and have been  
318 made into public projects on a TinyML development platform [16]. The benchmark will standardize  
319 the nascent field of TinyML and enable future progress through competition and comparability.  
320 TinyML as a field has the potential to democratize AI by removing the barrier of expensive hardware  
321 and can preserve privacy by keeping user’s data on the device that captures it. At the same time, the  
322 technology could be misused to more efficiently track and monitor unwilling individuals. Furthermore,  
323 because the devices themselves are inexpensive, TinyML can lead to increased electronic waste. By  
324 establishing a collaborative community, MLPerf Tiny can aid in the creation of standards for the  
325 responsible deployment of TinyML and mitigate the potential negative impacts of the technology.

## 326 8 Limitations

327 **New benchmarks & Long-term stability** MLPerf Tiny will continue to evolve to reflect the needs  
328 of the community. This will include new benchmarks that target new applications domains, such as  
329 wearables, medical devices, and environmental monitoring. While new benchmarks are being added  
330 and existing ones are evolved, it is also important to keep in mind that such a benchmark serves both  
331 as the comparison of state-of-the-art and for tracking historical progress. This latter goal requires  
332 benchmarks to be stable over time, eventually providing a multi-year perspective into the evolution of  
333 the field. Keeping the benchmark open source helps long-term reproducibility, while we also envision  
334 a subset of benchmarks to be kept long-term stable for this purpose. Also, thanks to MLCommons  
335 (mlcommons.org), a non-profit open engineering organization that hosts and develops the MLPerf  
336 benchmarks, the MLPerf Tiny benchmark will continue to be supported through future generations.

337 **Streaming inputs & Pre-processing** Time-domain tasks, such as the keyword spotting and anomaly  
338 detection benchmarks, typically involve continuously streaming inputs. Information from previous  
339 time steps can be exploited to improve the performance-efficiency tradeoff. However, limited  
340 bandwidth between the test runner (running on e.g. a PC) and the DUT makes it difficult to recreate a  
341 streaming scenario without adding delays incurred by data transfer. The choice of whether to include  
342 pre-processing in the measured benchmark can also have distorting effects on the results. Excluding  
343 feature extraction from measurement while allowing submitter-selected variations in feature extraction  
344 creates the possibility of a degenerate case where an entire model up to the penultimate layer is defined  
345 as “feature extraction”. Rigidly defined feature extraction precludes joint optimization over feature  
346 and model architecture that can be critical in the constrained systems targeted by this benchmark  
347 suite. Inclusion of feature extraction in the benchmark brings its own subtle complexities. Operating  
348 on a single discrete audio input in a non-streaming benchmark, one second of audio would involve  
349 one inference cycle and 40 feature extraction cycles, over-emphasizing the cost of feature extraction.  
350 In future versions, we aim to widen the scope of the benchmarks to include pre-processing as well.

351 **Extended coverage of layer types and model architectures** The closed division of the current  
352 benchmark suite includes models based mainly on FC and CNN layers in a variety of well known  
353 architectures, allowing only open division submissions to deviate from these. While our current suite  
354 provides a reasonable baseline that can be implemented on most HW platforms and openness to  
355 more experimental submissions, future benchmarks may include reference implementations using  
356 additional architectures e.g. RNNs.

## 357 9 Conclusion

358 The field of TinyML is poised to drive enormous growth within the IoT hardware and software  
359 industry. However, measuring the performance of these rapidly proliferating systems and comparing  
360 them in a meaningful way presents a considerable challenge; the complexity and dynamicity of the  
361 field obscure the measurement of progress and make embedded ML application and system design and  
362 deployment intractable. To enable more systematic development while fostering innovation, we need  
363 a fair, replicable, and robust method of evaluating TinyML systems. Developed as a collaboration  
364 between academia and industry, MLPerf Tiny is a suite of benchmarks for assessing the energy,  
365 latency, and accuracy of TinyML hardware, models, and runtimes. The benchmark suite and reference  
366 implementations are open-source and available at <https://github.com/mlcommons/tiny>.

## 367 **Acknowledgments**

368 MLPerf Tiny is the work of many individuals from multiple organizations. In this section, we  
369 acknowledge all those who helped produce the first set of results or supported the overall benchmark  
370 development. This work was also sponsored in part by the ADA (Applications Driving Architectures)  
371 Center, a JUMP Center co-sponsored by SRC and DARPA.

372 **Amazon** Amin Fazel

373 **Arizona State University** Jae-sun Seo

374 **California State Polytechnic University** Robert Hurtado

375 **CERN** Nicolo Ghielmetti

376 **Cisco Systems** Xinyuan Huang

377 **Columbia** Shvetank Prakash (Harvard)

378 **dividiti** Anton Lokhmotov

379 **Fermilab** Benjamin Hawks and Jules Muhizi (Harvard)

380 **Google** Ian Nappier, Dylan Zika, and David Patterson (UCB)

381 **Harvard** Max Lam and William Fu

382 **KU Leuven** Marian Verhelst (IMEC)

383 **ON Semiconductor** Jeffrey Dods

384 **Peng Cheng Lab** Yang Shazhou and Tang Hongwei

385 **Reality AI** Jeff Sieracki

386 **Renesas** Osama Neiroukh

387 **Silicon Labs** Dan Riedler

388 **STMicroelectronics** Mahdi Chtourou

389 **Synopsys** Dmitry Utyansky

390 **Syantiant** Mohammadreza Heydari and Rouzbeh Shirvani

391 **UCSD** Rushil Roy

392 **University of York** Poonam Yadav

## 393 **References**

394 [1] Tinyml foundation. URL <https://www.tinyml.org/>.

395 [2] Brainchip: Ai at the edge, Apr 2021. URL <https://brainchipinc.com/>.

396 [3] K. Asanović and D. A. Patterson. Instruction sets should be free: The case for risc-v. *EECS*  
397 *Department, University of California, Berkeley, Tech. Rep. UCB/EECS-2014-146*, 2014.

398 [4] C. Banbury, C. Zhou, I. Fedorov, R. Matas, U. Thakker, D. Gope, V. Janapa Reddi, M. Mattina,  
399 and P. Whatmough. Micronets: Neural network architectures for deploying tinyml applications  
400 on commodity microcontrollers. *Proceedings of Machine Learning and Systems*, 3, 2021.

401 [5] T. Bouguera, J.-F. Diouris, J.-J. Chaillout, R. Jaouadi, and G. Andrieux. Energy consumption  
402 model for sensor nodes based on lora and lorawan. *Sensors*, 18(7):2104, 2018.

403 [6] A. Chowdhery, P. Warden, J. Shlens, A. Howard, and R. Rhodes. Visual wake words dataset.  
404 *CoRR*, abs/1906.05721, 2019. URL <http://arxiv.org/abs/1906.05721>.

405 [7] R. David, J. Duke, A. Jain, V. J. Reddi, N. Jeffries, J. Li, N. Kreeger, I. Nappier, M. Natraj,  
406 S. Regev, et al. Tensorflow lite micro: Embedded machine learning on tinyml systems. *arXiv*  
407 *preprint arXiv:2010.08678*, 2020.

- 408 [8] I. Fedorov, R. P. Adams, M. Mattina, and P. Whatmough. Sparse: Sparse architecture search for  
409 cnns on resource-constrained microcontrollers. In *Advances in Neural Information Processing*  
410 *Systems* 32, pages 4978–4990. Curran Associates, Inc., 2019.
- 411 [9] S. Gal-On and M. Levy. Exploring coremark a benchmark maximizing simplicity and efficacy.  
412 *The Embedded Microprocessor Benchmark Consortium*, 2012.
- 413 [10] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In  
414 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–  
415 778, 2016.
- 416 [11] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and  
417 H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications.  
418 *arXiv preprint arXiv:1704.04861*, 2017.
- 419 [12] Y. Koizumi, S. Saito, H. Uematsu, N. Harada, and K. Imoto. Toyadmos: A dataset of miniature-  
420 machine operating sounds for anomalous sound detection. In *2019 IEEE Workshop on Ap-*  
421 *plications of Signal Processing to Audio and Acoustics (WASPAA)*, pages 313–317. IEEE,  
422 2019.
- 423 [13] Y. Koizumi, Y. Kawaguchi, K. Imoto, T. Nakamura, Y. Nikaido, R. Tanabe, H. Purohit, K. Sue-  
424 fusa, T. Endo, M. Yasuda, and N. Harada. Description and discussion on DCASE2020 challenge  
425 task2: Unsupervised anomalous sound detection for machine condition monitoring. In *arXiv*  
426 *e-prints: 2006.05822*, pages 1–4, June 2020. URL <https://arxiv.org/abs/2006.05822>.
- 427 [14] A. Krizhevsky, V. Nair, and G. Hinton. Cifar-10 (canadian institute for advanced research).  
428 2009. URL <http://www.cs.toronto.edu/~kriz/cifar.html>.
- 429 [15] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick.  
430 Microsoft coco: Common objects in context. In *European conference on computer vision*, pages  
431 740–755. Springer, 2014.
- 432 [16] L. Moreau. Tynymper - speech commands. URL [https://studio.edgeimpulse.com/](https://studio.edgeimpulse.com/public/29349/latest)  
433 [public/29349/latest](https://studio.edgeimpulse.com/public/29349/latest).
- 434 [17] H. Purohit, R. Tanabe, K. Ichige, T. Endo, Y. Nikaido, K. Suefusa, and Y. Kawaguchi. Miii  
435 dataset: Sound dataset for malfunctioning industrial machine investigation and inspection. *arXiv*  
436 *preprint arXiv:1909.09347*, 2019.
- 437 [18] V. J. Reddi, C. Cheng, D. Kanter, P. Mattson, G. Schmuelling, C.-J. Wu, B. Anderson,  
438 M. Breughe, M. Charlebois, W. Chou, R. Chukka, C. Coleman, S. Davis, P. Deng, G. Di-  
439 amos, J. Duke, D. Fick, J. S. Gardner, I. Hubara, S. Idgunji, T. B. Jablin, J. Jiao, T. S. John,  
440 P. Kanwar, D. Lee, J. Liao, A. Lokhmotov, F. Massa, P. Meng, P. Micikevicius, C. Osborne,  
441 G. Pekhimenko, A. T. R. Rajan, D. Sequeira, A. Sirasao, F. Sun, H. Tang, M. Thomson, F. Wei,  
442 E. Wu, L. Xu, K. Yamada, B. Yu, G. Yuan, A. Zhong, P. Zhang, and Y. Zhou. Mlperf inference  
443 benchmark, 2019.
- 444 [19] P. Torelli and M. Bangale. Measuring inference performance of machine-learning frameworks on  
445 edge-class devices with the mlmark benchmark. URL [https://www.eembc.org/techlit/](https://www.eembc.org/techlit/articles/MLMARK-WHITEPAPER-FINAL-1.pdf)  
446 [articles/MLMARK-WHITEPAPER-FINAL-1.pdf](https://www.eembc.org/techlit/articles/MLMARK-WHITEPAPER-FINAL-1.pdf).
- 447 [20] A. Torralba, R. Fergus, and W. T. Freeman. 80 million tiny images: A large data set for  
448 nonparametric object and scene recognition. *IEEE transactions on pattern analysis and machine*  
449 *intelligence*, 30(11):1958–1970, 2008.
- 450 [21] P. Warden. Speech commands: A dataset for limited-vocabulary speech recognition. *arXiv*  
451 *preprint arXiv:1804.03209*, 2018.
- 452 [22] Y. Zhang, N. Suda, L. Lai, and V. Chandra. Hello edge: Keyword spotting on microcontrollers.  
453 *arXiv preprint arXiv:1711.07128*, 2017.

## 454 **A Benchmark Framework**

455 The benchmark framework coordinates execution of the benchmark’s behavioral model among the  
456 various components in the system as described in the following sections.

### 457 **A.1 Framework Hardware**

458 In its most basic form, the framework consists of a host PC, a binary runner application GUI, the DUT,  
459 and a thin shim of firmware on the DUT that adheres to the framework communication protocol.

460 The framework for this particular benchmark provides two hardware configurations: one for measur-  
461 ing latency and accuracy, and the other for measuring energy (see Figure 3). The former mode is a  
462 simple connection between the host PC and the DUT through a serial port. The latter configuration  
463 expands the framework to include an electrical-isolation proxy (IO Manager), and an energy monitor  
464 (EMON), which supplies and measures energy consumption. Both configurations share the same  
465 process for initializing the DUT, loading the input test data, triggering the inference, and collecting  
466 results. The only difference is the energy configuration must initialize and manage the IO Manager  
467 and EMON hardware as well.

468 Two framework configurations are provided because energy configuration is a more complex setup,  
469 and energy scores may not be desired.

470 In the performance configuration, the Host PC talks directly to the DUT. In energy configuration,  
471 the IO Manager passes commands from the host software to the DUT through a serial port proxy.  
472 This proxy maintains a resilient serial connection state regardless of power cycling, and level shifters  
473 provide voltage domain conversion because the IO Manager GPIOs run at 5 Volts, whereas the DUT  
474 GPIO may vary. This configuration electrically isolates the DUT. For this framework version, the IO  
475 Manager is deployed as an Arduino UNO with its own custom firmware.

476 The energy monitor supplies and measures energy. The runner contains three EMON drivers for  
477 the following hardware: the STMicroelectronics LPM01A, the Jetperch Joulescope JS110, and the  
478 Keysight N6705. Regardless of the EMON used, it must supply one channel for the DUT and the  
479 other for the level shifters, the power used by level-shifters is not included in the total energy score  
480 because it is a cost associated with the framework and not the DUT. Only the power delivered to the  
481 core is measured. Furthermore, only one power supply is allowed to power the core; no additional  
482 energy supplies, such as batteries, supercapacitors, or energy harvesters may be used to defeat the  
483 measurement. Some energy monitors provide settings to increase the sample frequency or change the  
484 voltage. The runner GUI exposes these options to the user, so that the user can measure energy at a  
485 lower voltage. Note that the tradeoff between performance and energy can be dictated by the voltage,  
486 as higher voltage is often required to run at higher frequencies, and on-board SMPS conversion  
487 efficiency may vary considerably. This is one of the key insights of the benchmark: performance vs.  
488 energy tradeoffs.

489 The configurable sampling frequency is the rate at which data is returned from the EMON, which is  
490 much lower than the actual ADC sampling rate for all three devices. Since these internal sampling  
491 rates are typically well-within the time constant of the decoupling DUT’s power delivery decoupling  
492 capacitance, the rate has no impact on the score.

### 493 **A.2 Framework Firmware**

494 To provide consistency between the benchmark framework and the DUT implementation, a simple  
495 API is defined in C code that runs on the DUT. The DUT must provide two basic functions to interface  
496 to the host runner: a UART interface for receiving commands and returning status, and a timestamp  
497 function. In performance mode, the timestamp function is a local timer with a resolution of at least  
498 one millisecond (1 kHz); in energy mode, the timestamp function generates a falling edge GPIO  
499 signal which is used to trigger an external timer and synchronize data collection with the energy  
500 monitor. The API C code also provides functions to load data into a local buffer, translate and copy  
501 that data to the input tensor, trigger an inference, and read back the predictive results.

502 All of this functionality is partitioned into boilerplate code that does not change, or internally  
503 implemented, and code that has to be ported to the particular platform, called submitter implemented.

504 The submitter implemented code is the API that connects to the particular SDK, which may consist  
505 of optimized MCU libraries, or interface to hardware accelerators.

506 The API implemented by the submitter requires:

- 507 1. A timestamp function which prints a timestamp with a minimum resolution of one millisecon-  
508 ond, or generates an open-drain, GPIO toggle, depending on whether the DUT is configured  
509 latency or energy measurement.
- 510 2. UART Tx and Rx functionality, for communicating with the framework software.
- 511 3. A function for loading the input tensor with data sent down by the framework runner UI.
- 512 4. A function for performing a single inference.
- 513 5. A function for printing the prediction results.

514 The first two API functions provide external triggering and generic communication support; the latter  
515 three implement the minimum requirements to perform inference. Once these functions have been  
516 implemented, the framework software can successfully detect and communicate with the DUT.

517 The internally implemented firmware connects the API function calls in such a way that the framework  
518 software can send down an input dataset and instruct the firmware to perform a given number of  
519 iterations. The results are then extracted from the DUT (or the energy monitor) by the framework  
520 software.

### 521 **A.3 Framework Software**

522 The benchmark framework includes a Host PC GUI application called the runner. The runner allows  
523 the user to perform some basic setup and configuration, and then executes a pre-defined command  
524 script which activates the different hardware components required to execute the benchmark.

525 The runner is needed for three main reasons:

- 526 1. To provide a consistent benchmark interface to the user.
- 527 2. To standardize the execution of the benchmark, including measurements
- 528 3. To facilitate downloading the large number of input files required for accuracy measurements,  
529 since the typical target platform has less than a megabyte of flash memory.

530 The runner also provides visualization of the energy data, and feedback that can be useful for  
531 debugging framework connectivity.

532 The runner GUI, for the energy configuration, is shown in Figure 4.

533 The runner detects and initiates a handshake with detected hardware. After initialization, the test is  
534 started by clicking a button on the GUI, and a series of asynchronous commands are issued. These  
535 commands use the DUT API firmware to load input data and perform measurements. After the test  
536 completes, the runner collects the scores from the DUT and presents them to the user. The input  
537 stimuli files are located in a directory on the host PC, and are fed into the DUT depending on the  
538 measurement.

539 The measurement procedure is as follows:

- 540 1. Latency – Perform the following five times: download an input stimulus, load the input  
541 tensor (converting the data as needed), run inference for a minimum of 10 seconds and 10  
542 iterations, measure the inferences per second (IPS); report the median IPS of the five runs as  
543 the score.
- 544 2. Accuracy – Perform a single inference on the entire set of validation inputs, and collect the  
545 output tensor probabilities. The number of inputs vary, depending on the model. From the  
546 individual results, compute the Top-1 percent and the AUC. Each model has a minimum  
547 accuracy that must be met for the score to be valid.
- 548 3. Energy – Identical to latency, but in addition to measuring the number of inferences per  
549 second, measure the total energy used in the timing window and compute micro-Joules per  
550 inference. The same method of taking the median of five measurements is used.

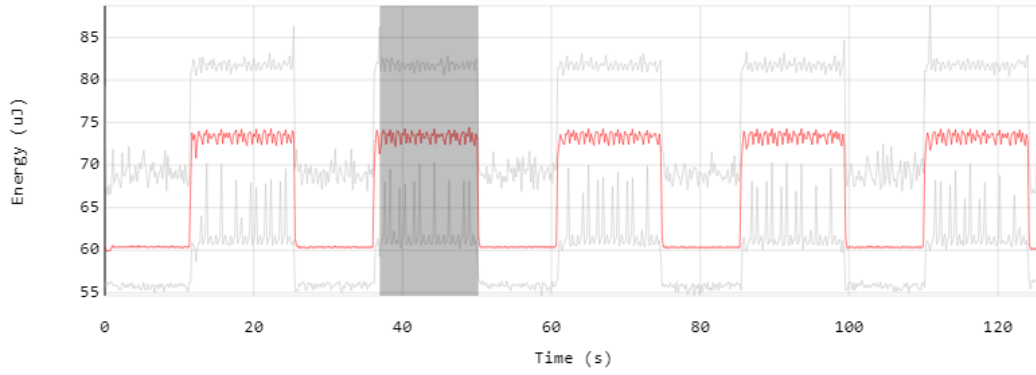


Figure 6: Energy Viewer

551 Results are stored in a folder and can be reloaded again. An energy viewer allows the user to examine  
552 the energy trace (Figure 6)

553 **Checklist**

- 554 1. For all authors...
- 555 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
- 556 contributions and scope? **[Yes]**
- 557 (b) Did you describe the limitations of your work? **[Yes]** See Section 8
- 558 (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See
- 559 Section 7
- 560 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
- 561 them? **[Yes]**
- 562 2. If you are including theoretical results...
- 563 (a) Did you state the full set of assumptions of all theoretical results? **[N/A]** No Theoretical
- 564 results
- 565 (b) Did you include complete proofs of all theoretical results? **[N/A]**
- 566 3. If you ran experiments (e.g. for benchmarks)...
- 567 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
- 568 mental results (either in the supplemental material or as a URL)? **[Yes]** See Section 4
- 569 and Conclusion
- 570 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
- 571 were chosen)? **[Yes]** All details are available on the linked github
- 572 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
- 573 ments multiple times)? **[No]** No, but the results of the benchmark are averaged over at
- 574 least 10 runs per input
- 575 (d) Did you include the total amount of compute and the type of resources used (e.g., type
- 576 of GPUs, internal cluster, or cloud provider)? **[No]** The models are small by their
- 577 nature so to reproduce the training portion of the reference implementations should
- 578 take only a few hours on an average laptop. Details provided in the GitHub repo link.
- 579 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 580 (a) If your work uses existing assets, did you cite the creators? **[Yes]** All datasets and
- 581 models are cited.
- 582 (b) Did you mention the license of the assets? **[No]** We provide citations to the datasets
- 583 and models which provide the licenses.
- 584 (c) Did you include any new assets either in the supplemental material or as a URL? **[Yes]**
- 585 Our benchmark is built on known-good reference models but some modifications have
- 586 been made. Section 4 outlines these changes and the models and training scripts are
- 587 available at the link provided in the conclusion.
- 588 (d) Did you discuss whether and how consent was obtained from people whose data you’re
- 589 using/curating? **[No]** The datasets are well recognized by the community.
- 590 (e) Did you discuss whether the data you are using/curating contains personally identifiable
- 591 information or offensive content? **[No]**
- 592 5. If you used crowdsourcing or conducted research with human subjects...
- 593 (a) Did you include the full text of instructions given to participants and screenshots, if
- 594 applicable? **[N/A]**
- 595 (b) Did you describe any potential participant risks, with links to Institutional Review
- 596 Board (IRB) approvals, if applicable? **[N/A]**
- 597 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 598 spent on participant compensation? **[N/A]**