

CS249r: Acceleration for ML: GPUs, TPUs and FPGAs



Oct 23



Course Logistics

Assignment – Why, What, How

- Assignment 1
 - **Why:** Vision (Camera) use cases + Dataset generation
 - **What:** End-to-end ML workflow
 - **How:** Frameworks (Edge Impulse) + Embedded Hardware (Nicla)
- Assignment 2
 - **Why:** Audio (Microphone) use cases
 - **What:** ML Pipeline (Data preprocessing + Model optimizations)
 - **How:** NAS (Edge Impulse) + Quantization/Pruning (TFLite)
- Assignment 3 - Model development
 - **Why:** Time series (IMU) use cases
 - **What:** ML frameworks programming (Tensorflow/TFLite/TFLM)
 - **How:** Bare metal implementation of end-to-end ML workflow
- Assignment 4 - Responsible AI
 - **Why:** Sustainability + Responsibility
 - **What:** Sustainability-aware system design and safe AI
 - **How:** TinyML Footprint + Adversarial Nibbler

Assignment 2

- Due tonight 11:59pm
- Any questions?
- Please check out the **#assignment2** slack channel for updates and answers to frequently asked questions

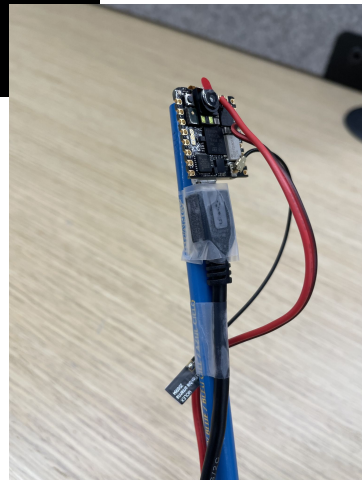
Assignment 3: Magic Nicla Wand

Due: November 6 at 11:59 pm

Objective:

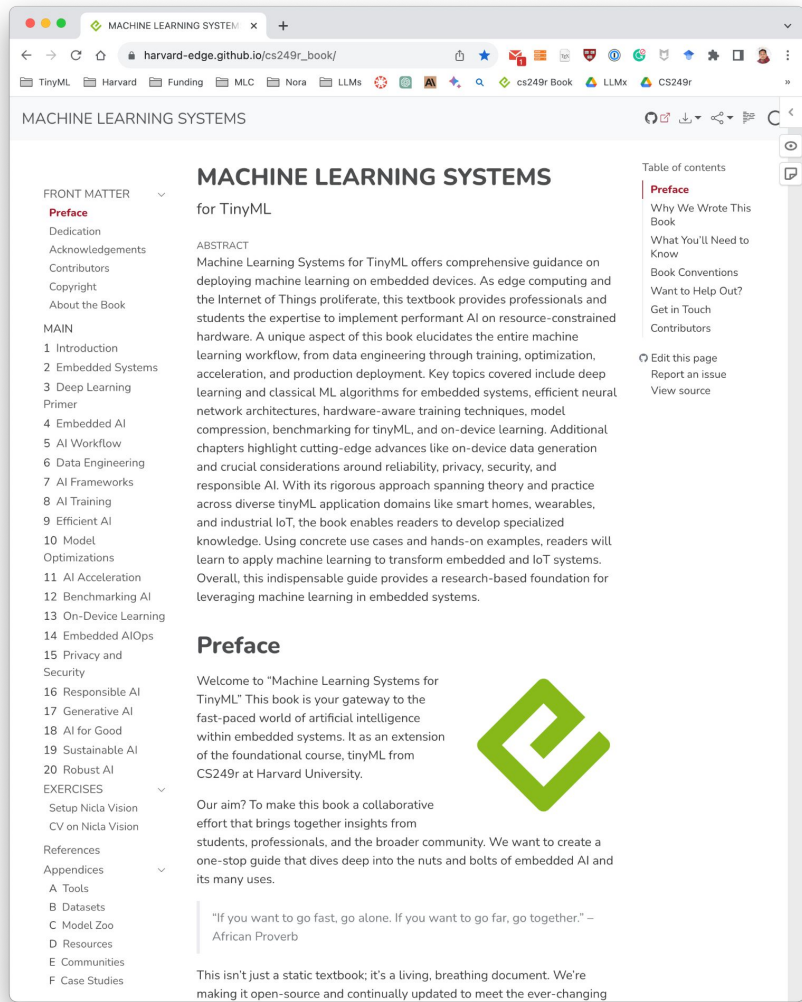
- Explore Tensorflow ecosystem (Tensorflow -> Tensorflow Lite -> Tensorflow Lite Micro)
- Model Optimization (quantization/pruning) using IMU data from Arduino Nicla Vision

Extra Credit: Deployment of model on Nicla



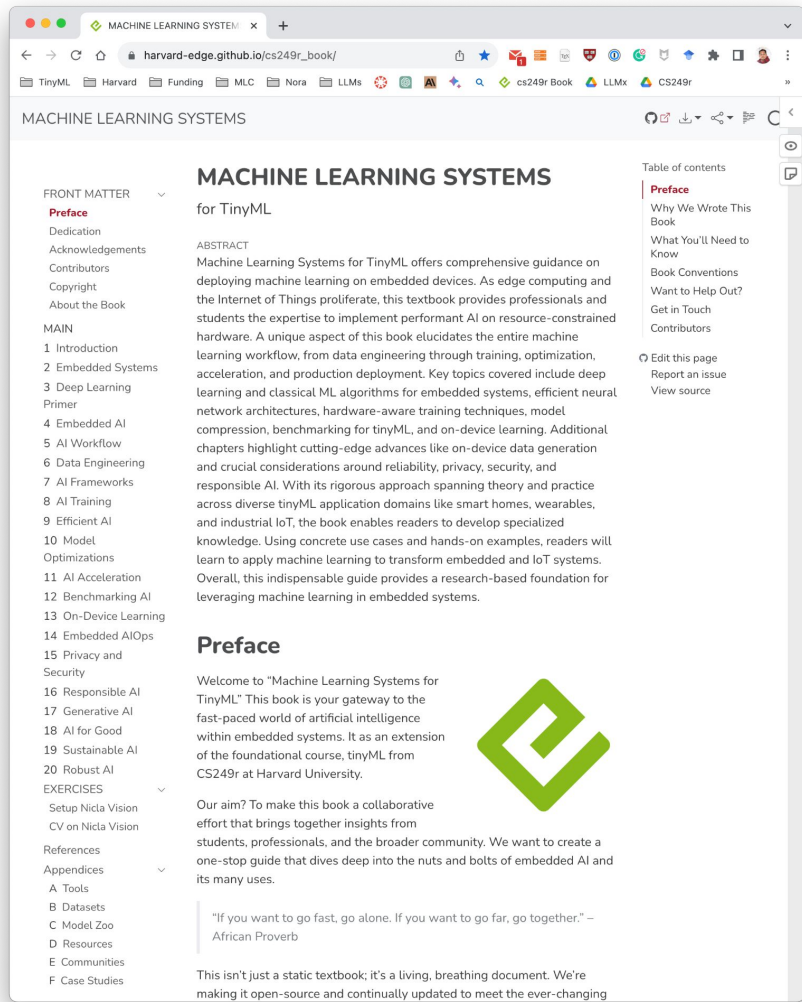
Scribing (again!)

- Week of topic
 - Meet with Matthew & VJ
 - Create a google doc
 - Share with staff
 - Iterate on rough outline with VJ/Matt
 - After 👍 from staff, put changes in a **single forked repo**
 - Create **one PR** with the entire chapter
 - Classmates will peer review using the PR



Scribing

- Grading (20% of your grade)
 - Part 1:
 - i. Paper review 10%
 - Content creation - 5%
 - Content curation - 5%
 - Part 2
 - i. Paper presentation 10%
 - References - 3%
 - Figures - 2%
 - Clarity - 5%



Scribing

- AI Frameworks
 - Available as a PR
 - Will be merged soon (EOD or tomorrow)

The screenshot shows a web browser displaying the GitHub page for the book "Machine Learning Systems for TinyML". The browser's address bar shows the URL "harvard-edge.github.io/cs249r_book/". The page has a sidebar on the left with a table of contents, including sections like "FRONT MATTER", "MAIN", and "EXERCISES". The main content area is titled "MACHINE LEARNING SYSTEMS for TinyML" and contains an "ABSTRACT" and a "Preface". The "Preface" section includes a green logo and a quote from an African Proverb. The right sidebar contains a "Table of contents" and links to "Edit this page", "Report an issue", and "View source".

MACHINE LEARNING SYSTEMS

for TinyML

ABSTRACT

Machine Learning Systems for TinyML offers comprehensive guidance on deploying machine learning on embedded devices. As edge computing and the Internet of Things proliferate, this textbook provides professionals and students the expertise to implement performant AI on resource-constrained hardware. A unique aspect of this book elucidates the entire machine learning workflow, from data engineering through training, optimization, acceleration, and production deployment. Key topics covered include deep learning and classical ML algorithms for embedded systems, efficient neural network architectures, hardware-aware training techniques, model compression, benchmarking for tinyML, and on-device learning. Additional chapters highlight cutting-edge advances like on-device data generation and crucial considerations around reliability, privacy, security, and responsible AI. With its rigorous approach spanning theory and practice across diverse tinyML application domains like smart homes, wearables, and industrial IoT, the book enables readers to develop specialized knowledge. Using concrete use cases and hands-on examples, readers will learn to apply machine learning to transform embedded and IoT systems. Overall, this indispensable guide provides a research-based foundation for leveraging machine learning in embedded systems.

Preface

Welcome to "Machine Learning Systems for TinyML." This book is your gateway to the fast-paced world of artificial intelligence within embedded systems. It as an extension of the foundational course, tinyML from CS249r at Harvard University.

Our aim? To make this book a collaborative effort that brings together insights from students, professionals, and the broader community. We want to create a one-stop guide that dives deep into the nuts and bolts of embedded AI and its many uses.

"If you want to go fast, go alone. If you want to go far, go together." – African Proverb

This isn't just a static textbook; it's a living, breathing document. We're making it open-source and continually updated to meet the ever-changing

Scribing (again!)

- This week (Oct 23)
 - Review Model optimizations chapter
- Next week
 - Review On-device learning chapter
 - Review of Benchmarking AI

The screenshot shows a web browser displaying the 'MACHINE LEARNING SYSTEMS for TinyML' page. The browser's address bar shows the URL 'harvard-edge.github.io/cs249r_book/'. The page has a sidebar on the left with a table of contents, including sections like 'FRONT MATTER', 'MAIN', 'EXERCISES', and 'References'. The main content area features the title 'MACHINE LEARNING SYSTEMS for TinyML' and an 'ABSTRACT' section. The abstract describes the book's focus on deploying machine learning on embedded devices, covering topics like training, optimization, acceleration, and production deployment. Below the abstract is a 'Preface' section with a green geometric logo. The preface welcomes readers to the book and states its aim to be a collaborative effort. A quote from an African proverb is also included. At the bottom, a note states that the book is a living document and will be updated.

MACHINE LEARNING SYSTEMS

MACHINE LEARNING SYSTEMS for TinyML

ABSTRACT

Machine Learning Systems for TinyML offers comprehensive guidance on deploying machine learning on embedded devices. As edge computing and the Internet of Things proliferate, this textbook provides professionals and students the expertise to implement performant AI on resource-constrained hardware. A unique aspect of this book elucidates the entire machine learning workflow, from data engineering through training, optimization, acceleration, and production deployment. Key topics covered include deep learning and classical ML algorithms for embedded systems, efficient neural network architectures, hardware-aware training techniques, model compression, benchmarking for tinyML, and on-device learning. Additional chapters highlight cutting-edge advances like on-device data generation and crucial considerations around reliability, privacy, security, and responsible AI. With its rigorous approach spanning theory and practice across diverse tinyML application domains like smart homes, wearables, and industrial IoT, the book enables readers to develop specialized knowledge. Using concrete use cases and hands-on examples, readers will learn to apply machine learning to transform embedded and IoT systems. Overall, this indispensable guide provides a research-based foundation for leveraging machine learning in embedded systems.

Preface

Welcome to "Machine Learning Systems for TinyML." This book is your gateway to the fast-paced world of artificial intelligence within embedded systems. It as an extension of the foundational course, tinyML from CS249r at Harvard University.

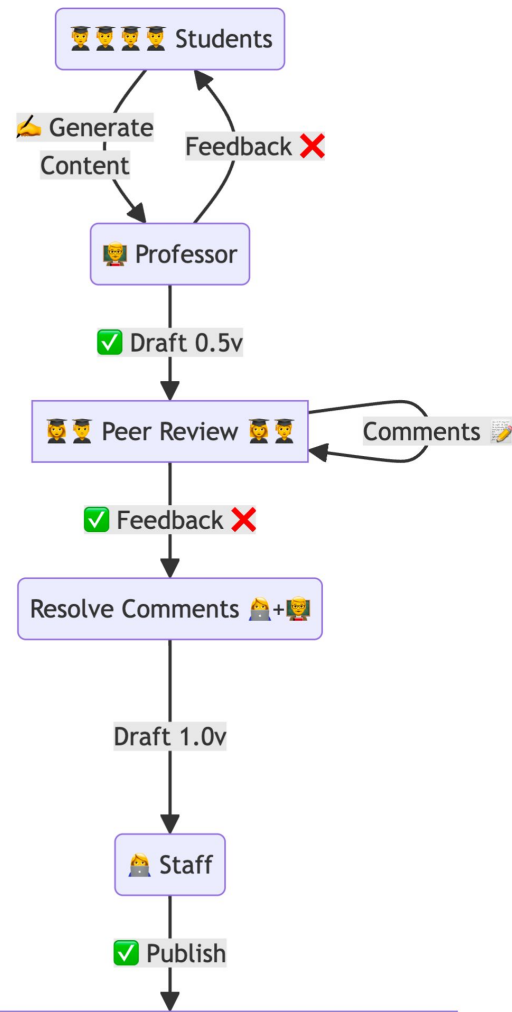
Our aim? To make this book a collaborative effort that brings together insights from students, professionals, and the broader community. We want to create a one-stop guide that dives deep into the nuts and bolts of embedded AI and its many uses.

"If you want to go fast, go alone. If you want to go far, go together." – African Proverb

This isn't just a static textbook; it's a living, breathing document. We're making it open-source and continually updated to meet the ever-changing

Assignment Schedule Updates


- Assignment 2
 - Due: October 23rd (Monday)
- Mid-Project Review
 - Due: October 30th (Monday)
- Assignment 3
 - Due: November 6th (Monday)
- Assignment 4 Part 1
 - Due: November 20th (Monday)
- Assignment 4 Part 2
 - Due: November 27th (Monday)
- Project Presentations
 - Due: December 4th (Monday)
- Final Report
 - Due: December 11th (Monday)



Projects

Mid-Project Presentation (**5 slides/3 mins**) – Oct 30th

Oct 30	MLOps	Daniel Situnayake, Head of Machine Learning at Edge Impulse	<p>Required</p> <ul style="list-style-type: none">• Edge Impulse (paper)• Hidden Technical Debt in Machine Learning (paper) <p>Optional</p> <ul style="list-style-type: none">• Data Cascades in High-Stakes AI (paper)• Pytorch RPC (paper)	None	Mid-Project Review Due
--------	-------	---	--	------	------------------------

A	B	C
	Paper Reading Table Leads	
	MLOps	Hidden Technical Debt
Table 1		Vijay Edupuganti
Table 2		Annie Landefeld
Table 3		Andrew Bass
Table 4		Aghyad Deeb
Table 5		Jared Ni

Project Update Slide Template

Put your slides in [here](#)

- Add your title slide
- Answer the questions with the template slides
- Stay on time



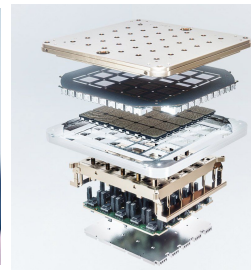
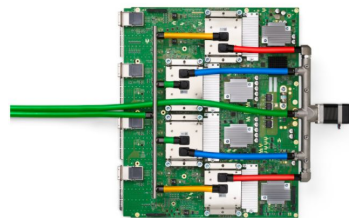
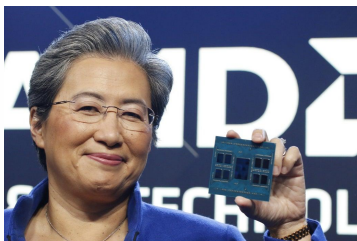
Lecture

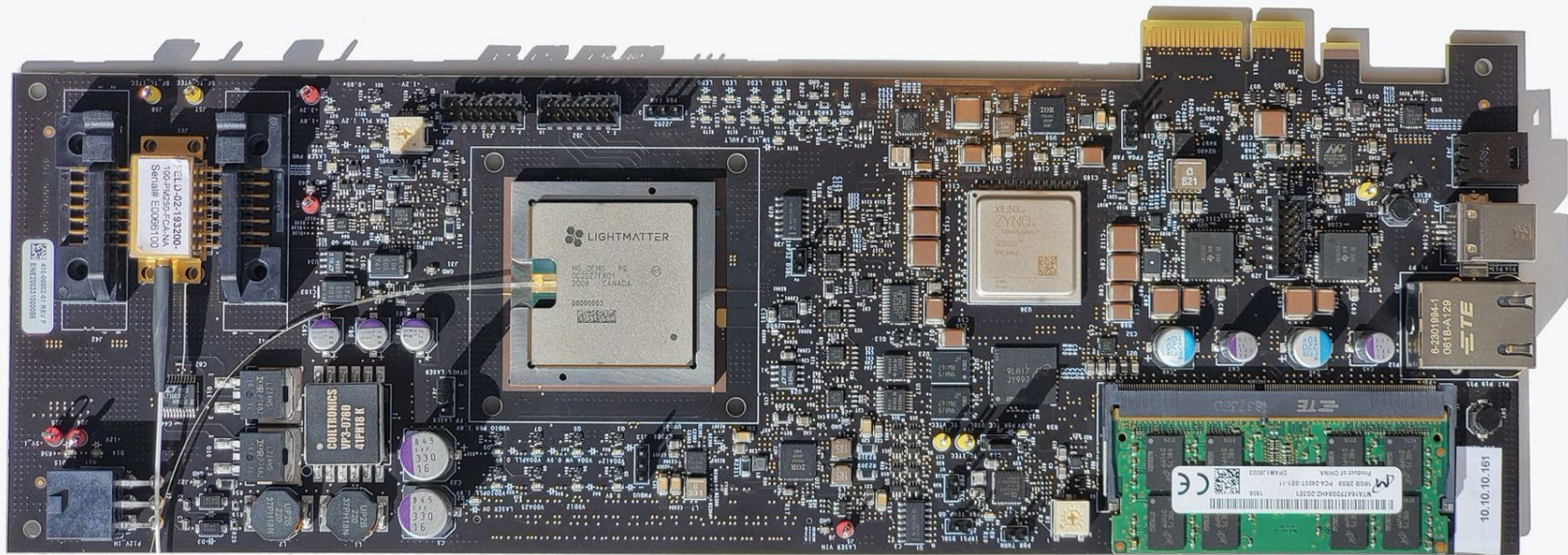
Course Topics

1. Overview and Introduction to Embedded Machine Learning
2. Data Engineering
3. Embedded Machine Learning Frameworks
4. Efficient Model Representation and Compression
- ~~5. Performance Metrics and Benchmarking of ML Systems~~
6. Learning on the Edge
7. Hardware Acceleration for Edge ML: GPUs, TPUs and FPGAs
8. Embedded MLOps
9. Secure and Privacy-Preserving On-Device ML
10. Responsible AI
11. Sustainability at the Edge
12. Generative AI at the Edge



Benchmarking ML Systems





AI Chip Landscape

S.T.

Tech Giants/System



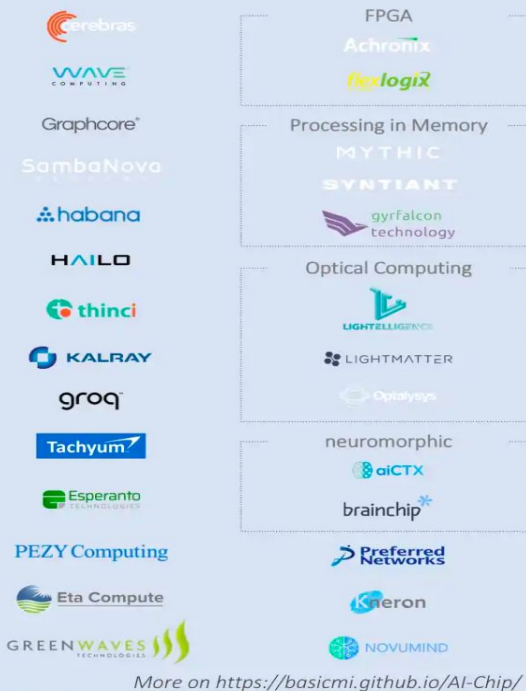
IC Vender/Fabless



Startup in China



Startup Worldwide



IP/Design Service

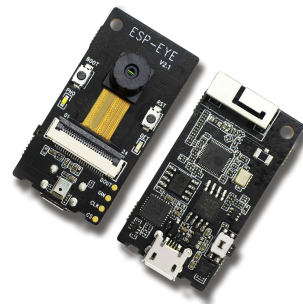
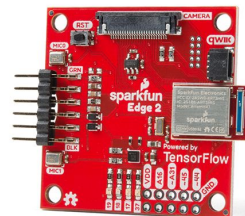
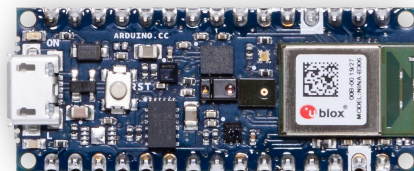


Compilers



Benchmarks







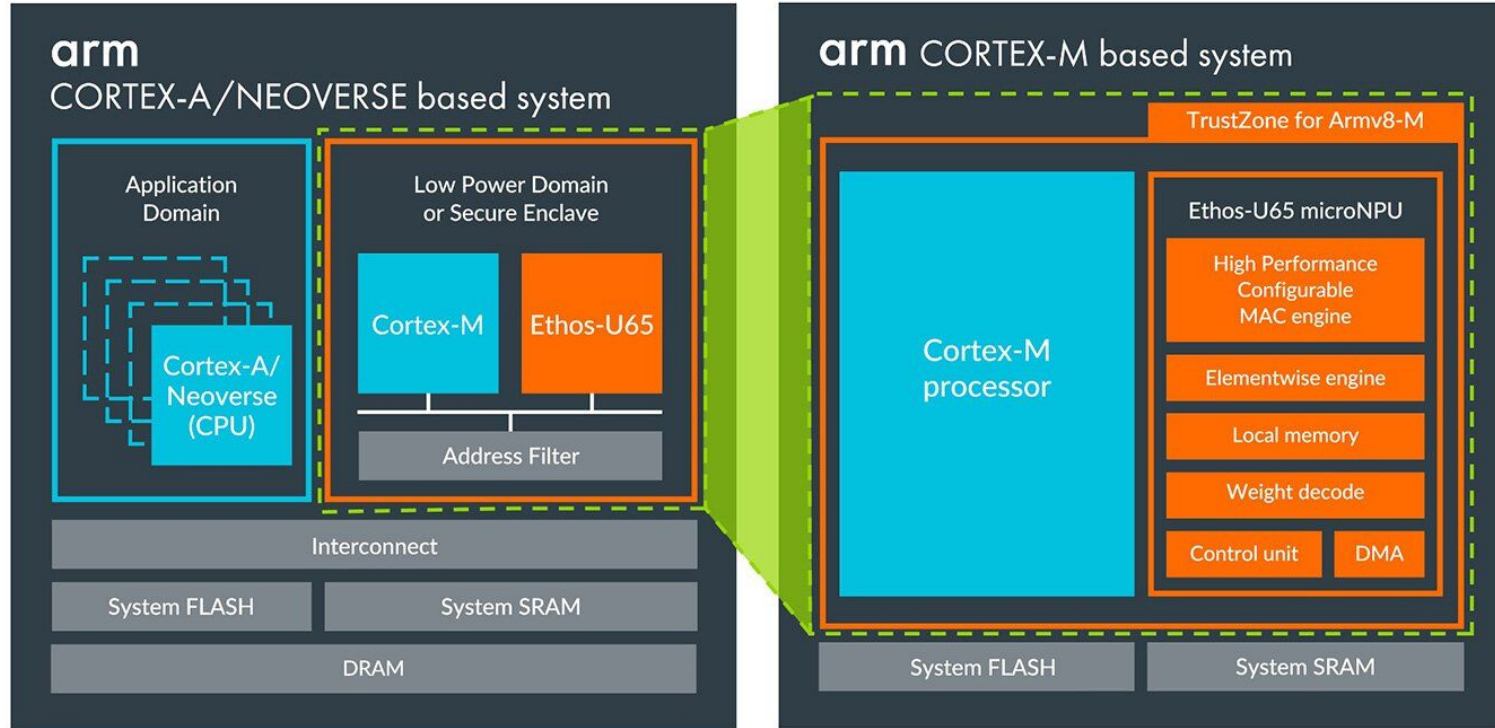
Board	MCU / ASIC	Clock	Memory	Sensors	Radio
Himax WE-I Plus EVB	HX6537-A 32-bit EM9D DSP	400 MHz	2MB flash 2MB RAM	Accelerometer, Mic, Camera	None
Arduino Nano 33 BLE Sense	32-bit nRF52840	64 MHz	1MB flash 256kB RAM	Mic, IMU, Temp, Humidity, Gesture, Pressure, Proximity, Brightness, Color	BLE
SparkFun Edge 2	32-bit ArtemisV1	48 MHz	1MB flash 384kB RAM	Accelerometer, Mic, Camera	BLE
Espressif EYE	32-bit ESP32-D0WD	240 MHz	4MB flash 520kB RAM	Mic, Camera	WiFi, BLE

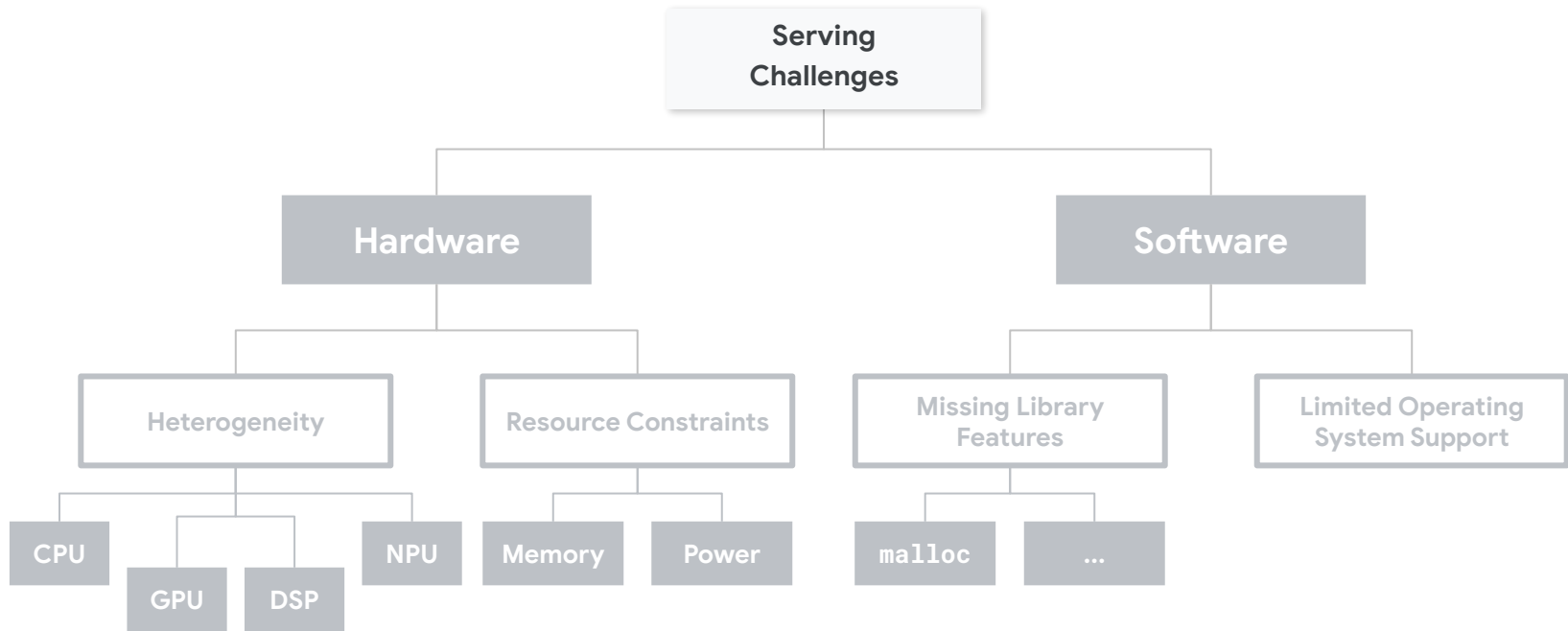
The diagram illustrates the progression of AI capabilities on ARM processors, categorized by processor family and performance level.

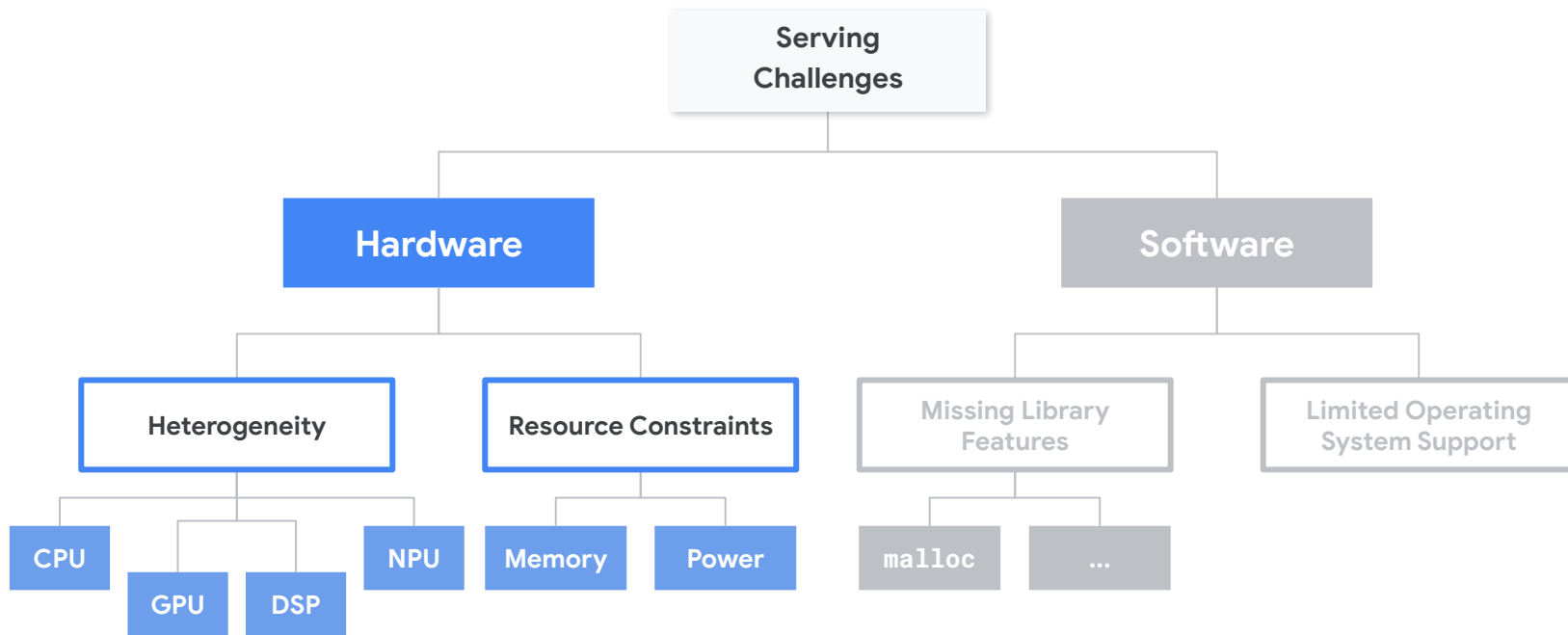
Processor Family	Capabilities
Cortex-M today	Vibration detection, Sensor fusion
Cortex-M55	Keyword detection, Anomaly detection, Object detection
Cortex-M and Ethos-U55	Gesture detection, Biometric awareness, Speech recognition
Cortex-A, Mali and Ethos-N	Object classification, Real-time recognition

TOP/s

Different Systems Applicable to Ethos-U





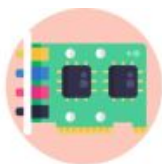




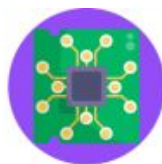
CPU



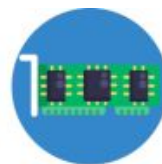
GPU



FPGA



**Flexible Dataflow
Accelerator**



**Fixed Dataflow
Accelerator***

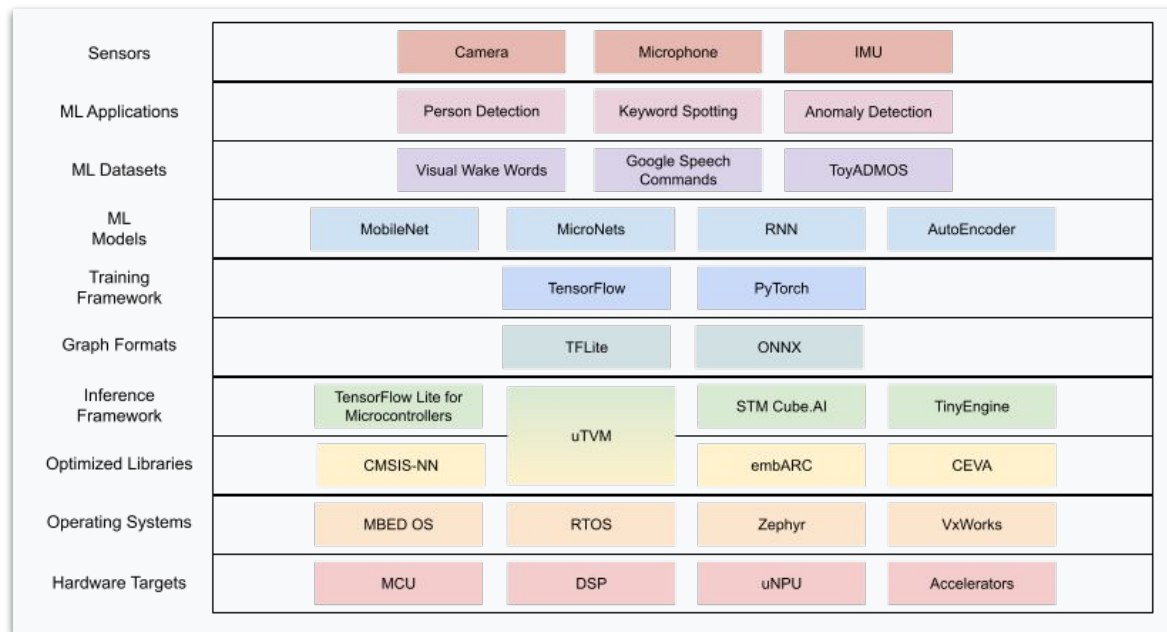


*** When the dataflow is optimal for a workload; Efficiency can drop under workload change**



TinyML System Stack is Complicated

- Machine learning system stack is **complicated**
- Many **different** models, datasets, models, frameworks, formats, compilers, libraries, operating systems, targets
- The **cross-product** makes it challenging to decipher system performance



Apples-to-apples comparison



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

...



bench·mark

/ˈben(t)SHmärk/

See definitions in:

All

Technology

Surveying

noun

1. a standard or point of reference against which things may be compared or assessed.
"a benchmark case"

Similar:

standard

point of reference

basis

gauge

criterion

specification



2. a surveyor's mark cut in a wall, pillar, or building and used as a reference point in measuring altitudes.

verb

evaluate or check (something) by comparison with a standard.
"we are **benchmarking** our performance **against** external criteria"

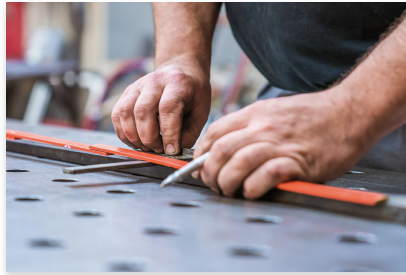
Definitions from Oxford Languages

Feedback

Benchmarking

Use to

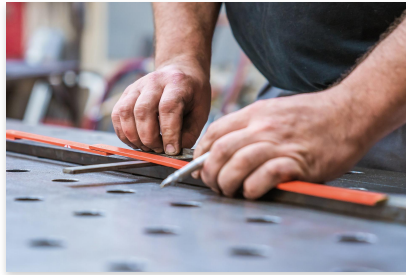
- **Compare** solutions



Benchmarking

Use to

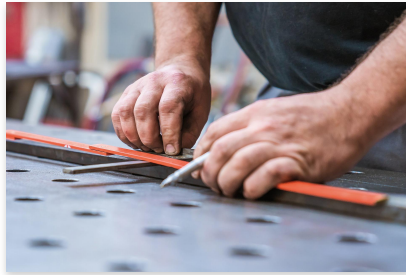
- **Compare** solutions
- **Inform** selection



Benchmarking

Use to

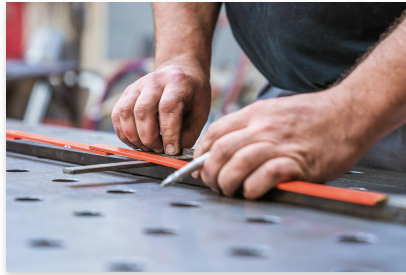
- **Compare** solutions
- **Inform** selection
- **Measure** and track progress



Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field



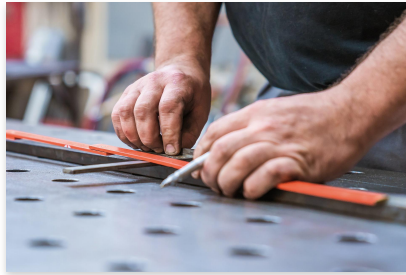
Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous



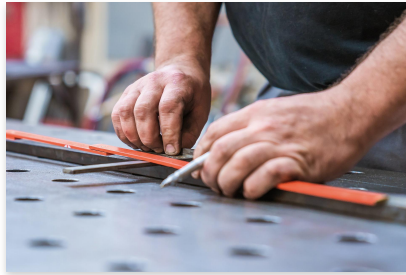
Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



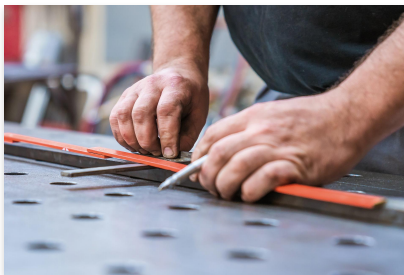
Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



Provides

- **Standardization** of use cases and workloads

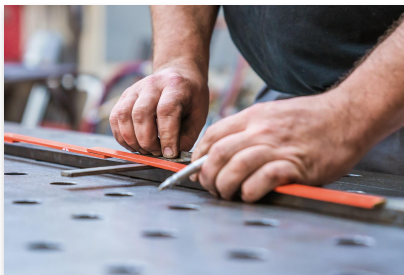
Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



Provides

- **Standardization** of use cases and workloads
- **Comparability** across heterogeneous HW/SW systems

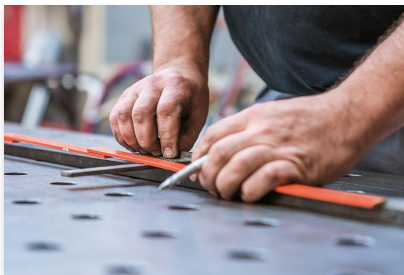
Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



Provides

- **Standardization** of use cases and workloads
- **Comparability** across heterogeneous HW/SW systems
- **Complex characterization** of system compromises

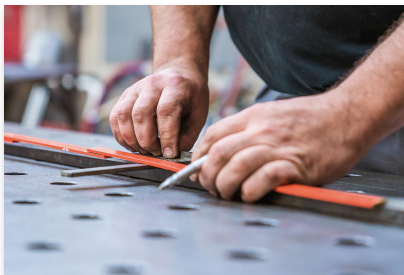
Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



Provides

- **Standardization** of use cases and workloads
- **Comparability** across heterogeneous HW/SW systems
- **Complex characterization** of system compromises
- **Verifiable and Reproducible** results

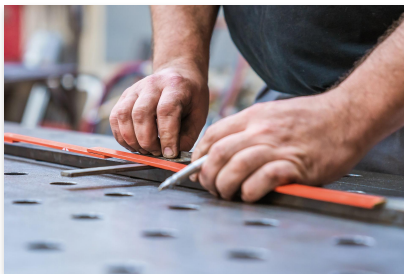
Benchmarking

Use to

- **Compare** solutions
- **Inform** selection
- **Measure** and track progress
- **Raise** the bar, **advance** the field

Requires

- **Methodology** that is both fair and rigorous
- **Community** support and consensus



Provides

- **Standardization** of use cases and workloads
- **Comparability** across heterogeneous HW/SW systems
- **Complex characterization** of system compromises
- **Verifiable and Reproducible** results



Goals



Enforce performance
result replicability to
ensure reliable results

Goals



Enforce **performance
result replicability** to
ensure reliable results



Use **representative
workloads**, reflecting
production use-cases

Goals



Enforce performance
result replicability to
ensure reliable results



Use representative
workloads, reflecting
production use-cases



Encourage innovation
to improve the
state-of-the-art of ML

Goals



Enforce **performance
result 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**

Goals



Enforce **performance**
result 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

Goals



Enforce **performance**
result 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 participate

Wide Array of ML Tasks

Task Category	Use Case
Audio	Audio Wake Words Context Recognition Control Words Keyword Detection
Image	Visual Wake Words Object Detection Gesture Recognition Object Counting Text Recognition
Physiological / Behavioral Metrics	Segmentation Anomaly Detection Forecasting Activity Detection
Industry Telemetry	Sensing Predictive Maintenance Motor Control

Wide Array of ML Tasks

Task Category	Use Case	Model Type
Audio	Audio Wake Words Context Recognition Control Words Keyword Detection	DNN CNN RNN LSTM
Image	Visual Wake Words Object Detection Gesture Recognition Object Counting Text Recognition	DNN CNN SVM Decision Tree KNN Linear
Physiological / Behavioral Metrics	Segmentation Anomaly Detection Forecasting Activity Detection	DNN Decision Tree SVM Linear
Industry Telemetry	Sensing Predictive Maintenance Motor Control	DNN Decision Tree SVM Linear Naive Bayes

Wide Array of ML Tasks

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

A Principled Approach to Subsetting

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

A Principled Approach to Subsetting

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

A Principled Approach to Subsetting

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?

A Principled Approach to Subsetting

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?

A Principled Approach to Subsetting

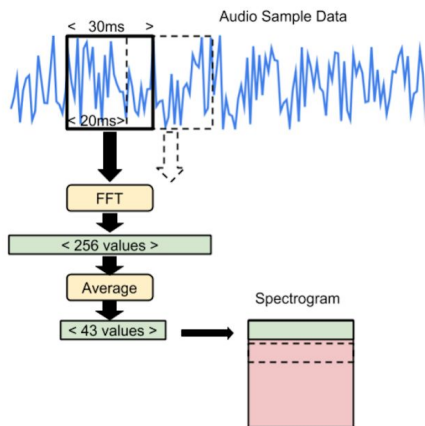
Big Questions	Inference
1. Benchmark definition	What is the definition of a benchmark task?
2. Benchmark selection	Which benchmark task to select?
3. Metric definition	What is the measure of “performance” in ML systems?
4. Implementation equivalence	How do submitters run on different hardware/software systems?
5. Issues with optimizations	Quantization, calibration, and/or retraining?

A Principled Approach to Subsetting

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 with optimizations	Quantization, calibration, and/or retraining?
6. Results	Do we normalize and/or summarize results?

MLPerf “Tiny” Tasks

Keyword Spotting



Warden, Pete. "Speech commands: A dataset for limited-vocabulary speech recognition." *arXiv preprint arXiv:1804.03209* (2018).

Visual Wake Words



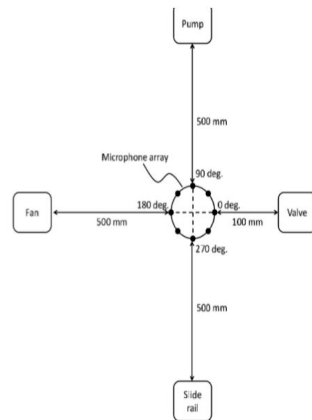
(a) 'Person'



(b) 'Not-person'

Chowdhery, Aakanksha, et al. "Visual wake words dataset." *arXiv preprint arXiv:1906.05721* (2019).

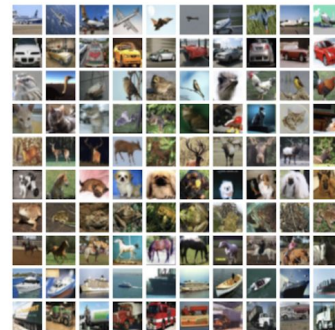
Anomaly Detection



Purohit, Harsh, et al. "MIMI dataset: Sound dataset for malfunctioning industrial machine investigation and inspection." *arXiv preprint arXiv:1909.09347* (2019).

Tiny Image Classification

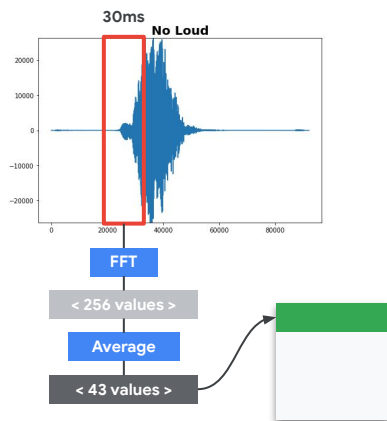
airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck



Krizhevsky, Alex, and Geoffrey Hinton. "Learning multiple layers of features from tiny images." (2009): 7.

MLPerf “Tiny” Tasks

Keyword Spotting



Warden, Pete. "Speech commands: A dataset for limited-vocabulary speech recognition." *arXiv preprint arXiv:1804.03209* (2018).

Visual Wake Words



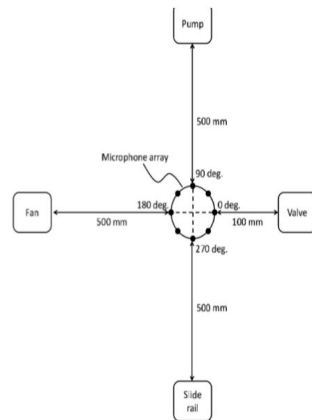
(a) 'Person'



(b) 'Not-person'

Chowdhery, Aakanksha, et al. "Visual wake words dataset." *arXiv preprint arXiv:1906.05721* (2019).

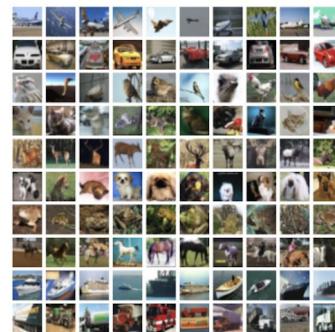
Anomaly Detection



Purohit, Harsh, et al. "MIMI dataset: Sound dataset for malfunctioning industrial machine investigation and inspection." *arXiv preprint arXiv:1909.09347* (2019).

Tiny Image Classification

airplane
automobile
bird
cat
deer
dog
frog
horse
ship
truck



Krizhevsky, Alex, and Geoffrey Hinton. "Learning multiple layers of features from tiny images." (2009): 7.

Problem
definition

Dataset
selection (public
domain)

Model selection

Model training
code

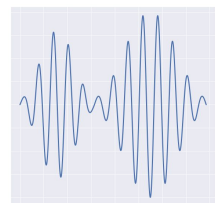
Derive "Tiny"
version:
Quantization

Embedded
implementation

Benchmarking
harness
integration

Deploy on
device

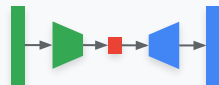
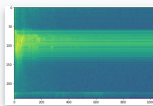
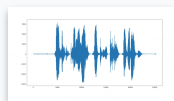
Example
benchmark run



Anomalous Sound
Detection System

Normal

Anomaly



FP32

INT8

Training
Code

ARM
mbed OS

Problem	AD
Model	FC-AE
Size	270 Kpar
Latency	10.4 ms/inf.
Accuracy	.86 AUC
Energy	516 μ J/inf.

Problem
definition

Dataset
selection (public
domain)

Model selection

Model training
code

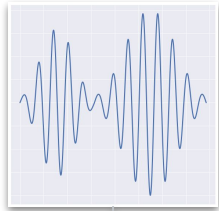
Derive "Tiny"
version:
Quantization

Embedded
implementation

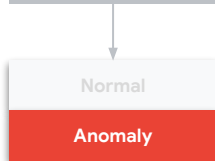
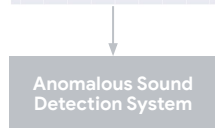
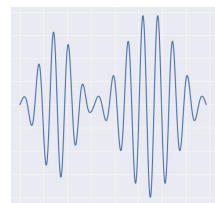
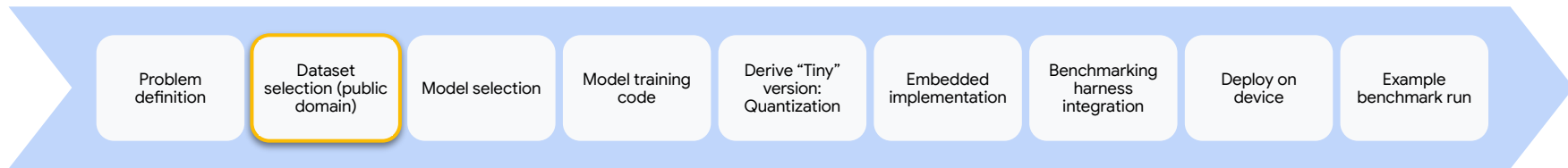
Benchmarking
harness
integration

Deploy on
device

Example
benchmark run



Anomalous Sound
Detection System



Problem
definition

Dataset
selection (public
domain)

Model selection

Model training
code

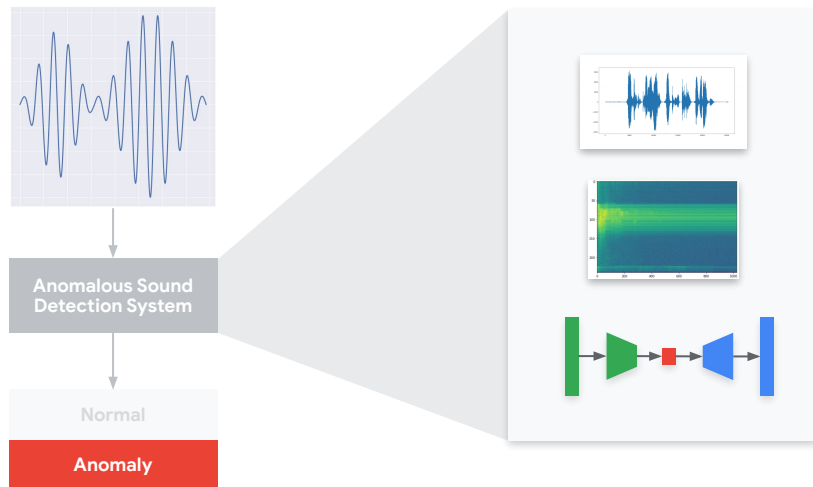
Derive "Tiny"
version:
Quantization

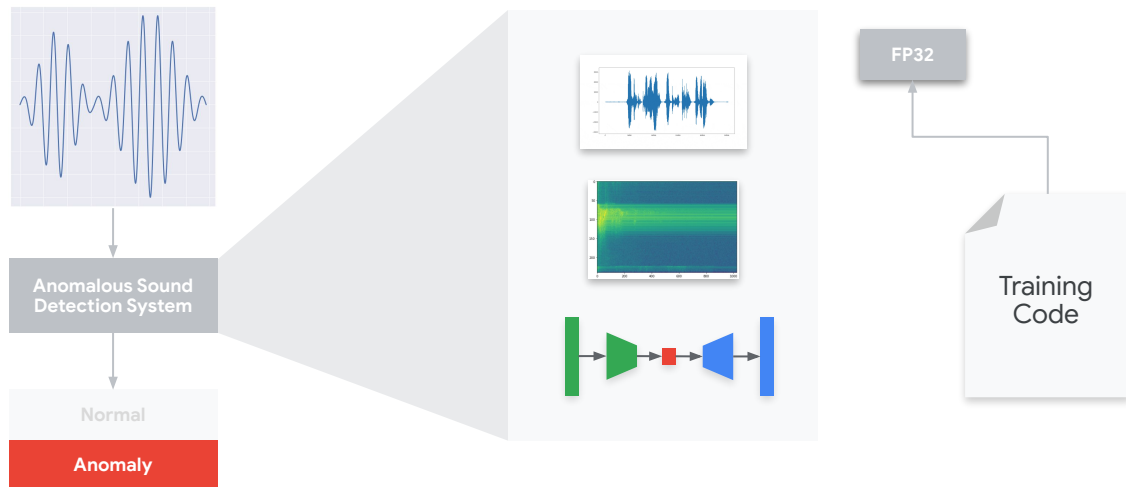
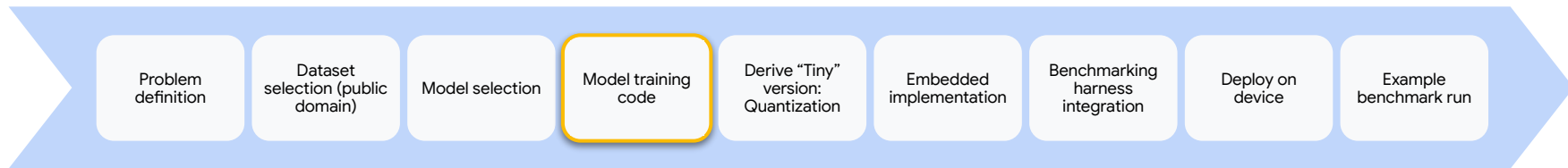
Embedded
implementation

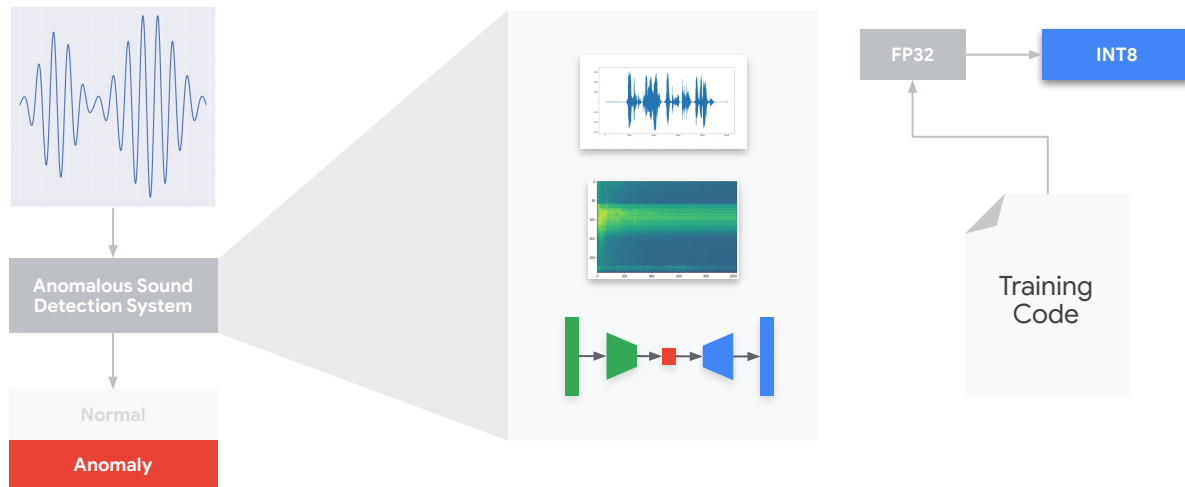
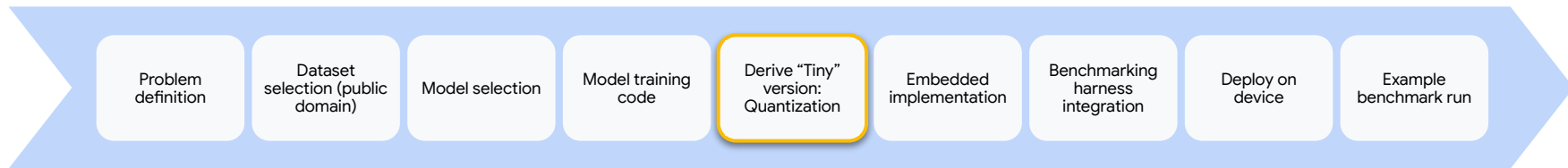
Benchmarking
harness
integration

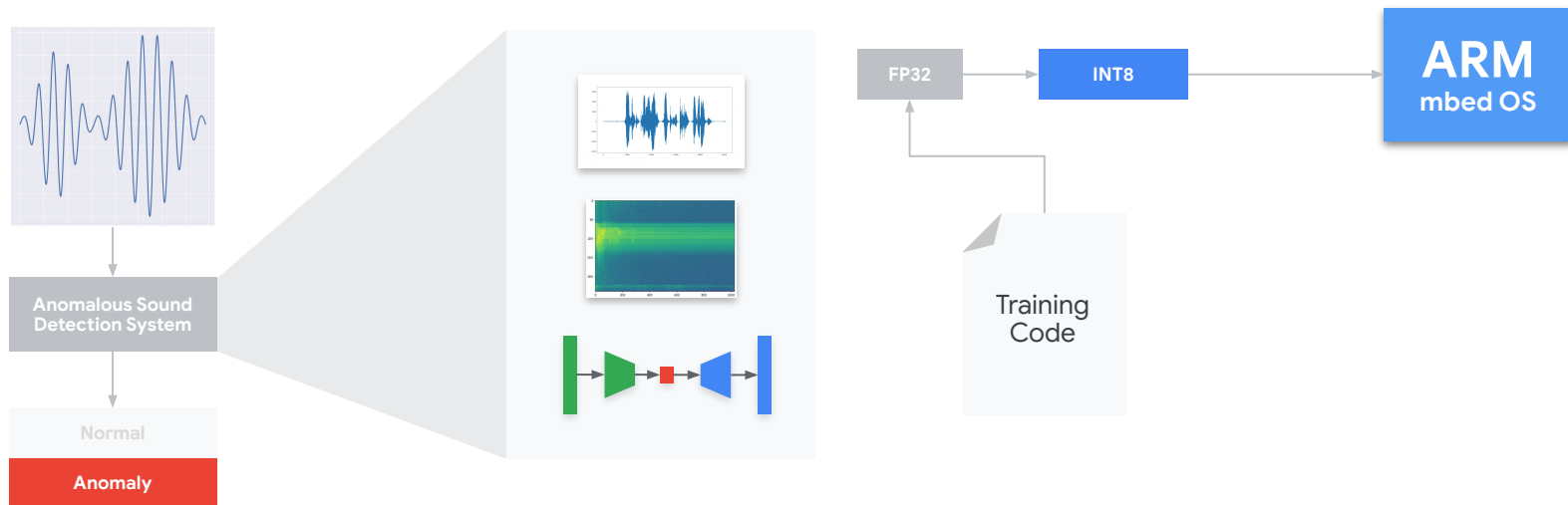
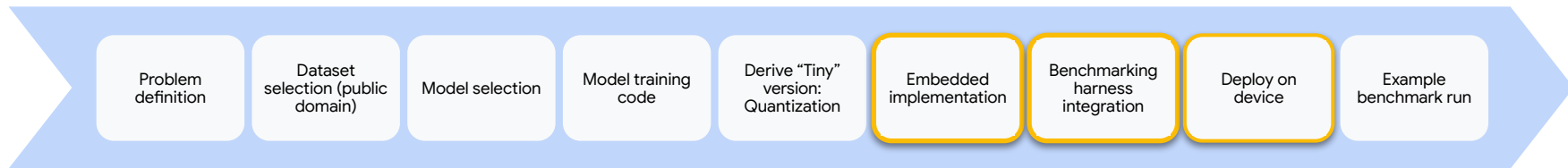
Deploy on
device

Example
benchmark run









Problem
definition

Dataset
selection (public
domain)

Model selection

Model training
code

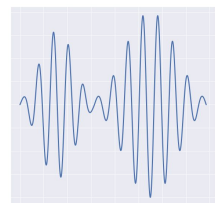
Derive "Tiny"
version:
Quantization

Embedded
implementation

Benchmarking
harness
integration

Deploy on
device

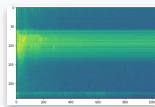
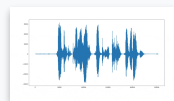
Example
benchmark run



Anomalous Sound
Detection System

Normal

Anomaly



FP32

INT8

Training
Code

ARM
mbed OS

Problem AD

Model	FC-AE
Size	270 Kpar
Latency	10.4 ms/inf.
Accuracy	.86 AUC
Energy	516 μ J/inf.

Metrics

Latency

Small fast dataset

Loop of inferences

No data-dependent
execution

```
Runtime requirements have been met.  
Performance results for window 10:  
# Inferences :      1000  
Runtime      :      10.524 sec.  
Throughput   :      95.020 inf./sec.  
Runtime requirements have been met.  
-----  
Median throughput is 95.019 inf./sec.  
-----
```



Metrics

Latency

Small fast dataset

Loop of inferences

No data-dependent execution

```
Runtime requirements have been met.  
Performance results for window 10:  
# Inferences :      1000  
Runtime      :      10.524 sec.  
Throughput   :      95.020 inf./sec.  
Runtime requirements have been met.  
-----  
Median throughput is 95.019 inf./sec.  
-----
```



a.

Accuracy

Evaluate on larger dataset

Top-1 accuracy & AUC

CLOSED: meet threshold
v.

OPEN: part of the metrics

Metrics

Latency

Small fast dataset

Loop of inferences

No data-dependent execution

```
Runtime requirements have been met.  
Performance results for window 10:  
# Inferences :      1000  
Runtime      :    10.524 sec.  
Throughput   :    95.020 inf./sec.  
Runtime requirements have been met.  
-----  
Median throughput is 95.019 inf./sec.
```



a.

Accuracy

Evaluate on larger dataset

Top-1 accuracy & AUC

CLOSED: meet threshold
v.

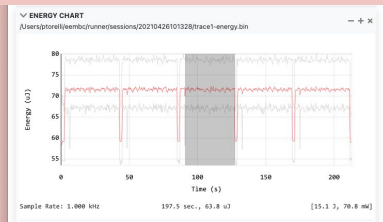
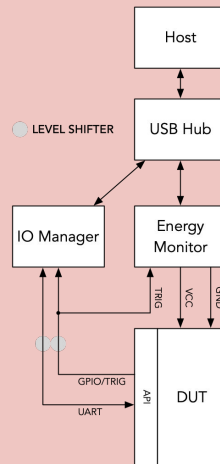
OPEN: part of the metrics

Energy

No
“cherry-picking”

Power Monitor
setup

Median result



V1.0 Results

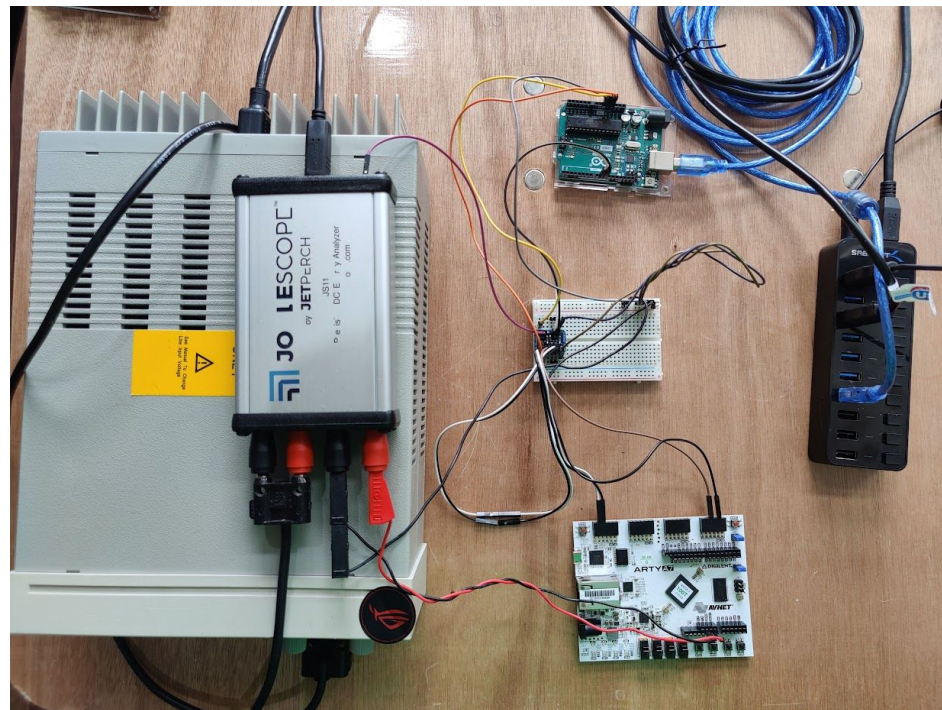
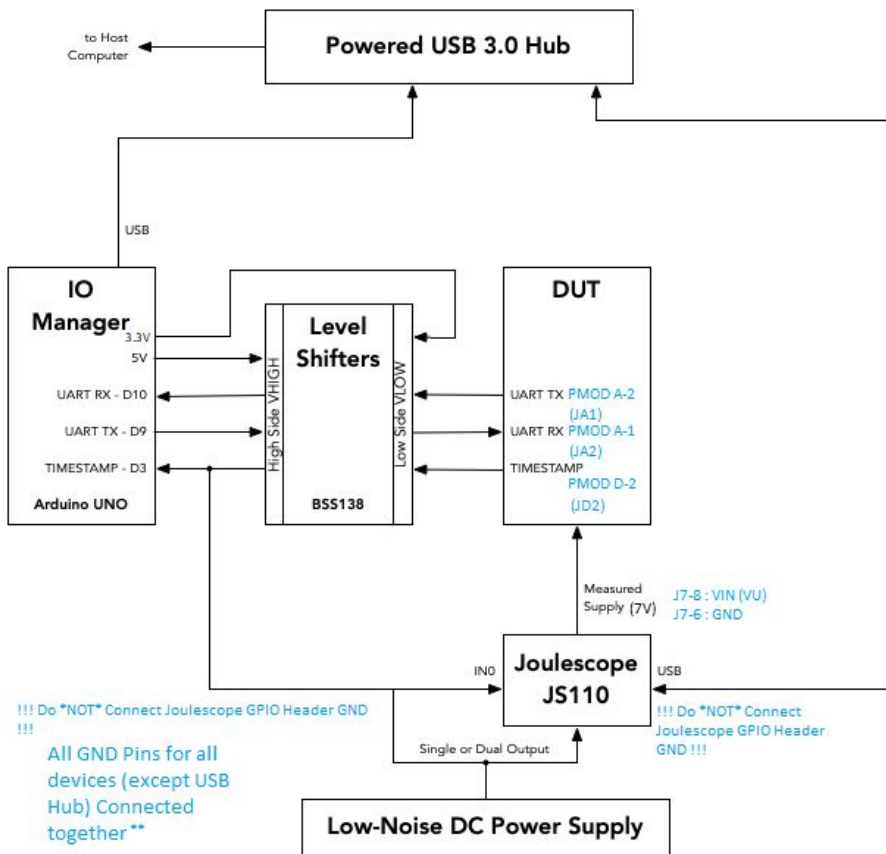
mlcommons.org/en/inference-tiny-10

Submitter	Board Name	SoC Name	Processor(s) & Number	Accelerator(s) & Number
Greenwaves Technologies	GAP9 EVK	GAP9	RISC-V Core (1+9)	NE16 (1)
Greenwaves Technologies	GAP9 EVK	GAP9	RISC-V Core (1+9)	NE16 (1)
OctoML	NRF5340DK	nRF5340	Arm® Cortex®-M33	
OctoML	NUCLEO-L4R5ZI	STM32L4R5ZIT6U	Arm® Cortex®-M4	
OctoML	NUCLEO-L4R5ZI	STM32L4R5ZIT6U	Arm® Cortex®-M4	
Plumerai	B_U585I_IOT02A	STM32U585	Arm® Cortex®-M33	
Plumerai	CY8CPROTO-062-4343w	PSoC 62 MCU	Arm® Cortex®-M4	
Plumerai	DISCO-F746NG	STM32F746	Arm® Cortex®-M7	
Plumerai	NUCLEO-L4R5ZI	STM32L4R5ZIT6U	Arm® Cortex®-M4	
Silicon Labs	xG24-DK2601B	EFR32MG24	Arm® Cortex®-M33	Silicon Labs MVP(1) (78 MHz, 1.8V)
STMicroelectronics	NUCLEO-H7A3ZI-Q	STM32H7A3ZIT6Q	Arm® Cortex®-M7	
STMicroelectronics	NUCLEO-L4R5ZI	STM32L4R5ZIT6U	Arm® Cortex®-M4	
STMicroelectronics	NUCLEO-U575ZI-Q	STM32U575ZIT6Q	Arm® Cortex®-M33	
Syntiant	NDP9120-EVL	NDP120	M0 + HiFi	Syntiant Core 2 (98MHz, 1.1V)
Syntiant	NDP9120-EVL	NDP120	M0 + HiFi	Syntiant Core 2 (30MHz, 0.9V)
Qualcomm Innovation Center	Next Generation Snapdragon Mobile Platform HDK	Next Generation Snapdragon Mobile Platform	Qualcomm Kryo CPU(1)	Qualcomm Sensing Hub(1)

Select Keyword Spotting Results

Submitter	SoC	Accelerator	Accuracy	Latency (mS)	Energy (uJ)
Syntiant	NDP120	Syntiant Core 2	90%	1.48	43.8
STMicroelectronics	STM32U575ZIT6Q	None	90%	44.2	1138.5

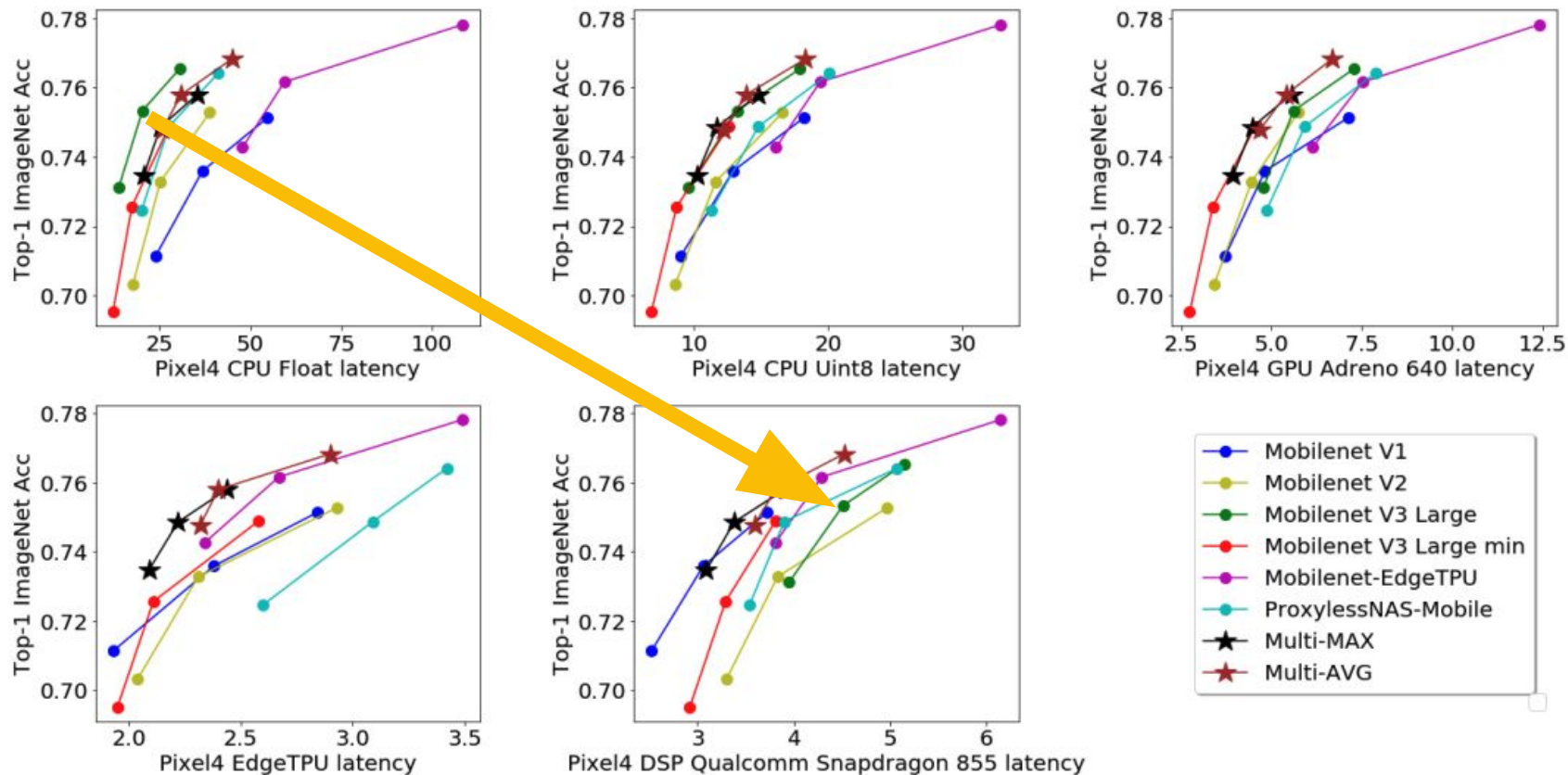
FPGA Energy Configuration



Hardware Lottery Problem

Source:

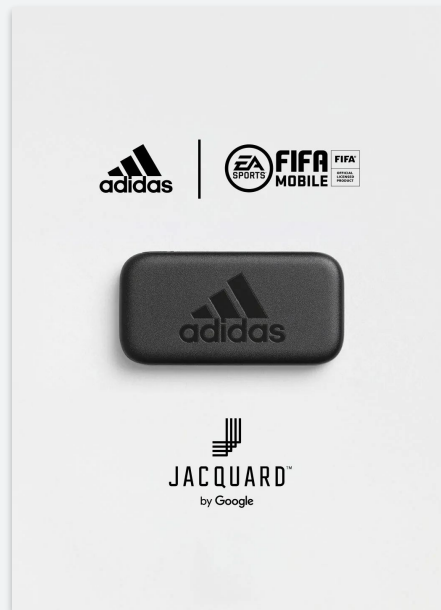
Chu, Grace, et al. "Discovering multi-hardware mobile models via architecture search." Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. 2021.



Emerging TinyML Use Cases

Example: Smart shoes

- Kicking
- Penalty kicking
- Passing
- Dribbling
- ...



Emerging TinyML Use Cases

Example: Augmented Reality

- Eye tracking
- Hand tracking
- Computer vision
- Superresolution
- ...



Toward Emerging Multi-DNN Models

Pipelined
DNNs



Keyword
Spotting

Speech
Processing

- Back-to-back execution
- Execution dependency

Toward Emerging Multi-DNN Models

Pipelined DNNs



Keyword Spotting

Speech Processing

- Back-to-back execution
- Execution dependency

Concurrent DNNs



Eye Tracking

Obstacle Detection

Video Processing

- Concurrent execution
- Execution deadline

Toward Emerging Multi-DNN Models

Pipelined DNNs



Keyword Spotting

Speech Processing

- Back-to-back execution
- Execution dependency

Concurrent DNNs



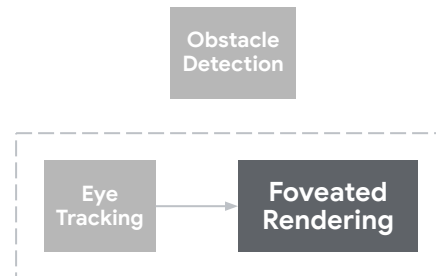
Eye Tracking

Obstacle Detection

Video Processing

- Concurrent execution
- Execution deadline

Concurrent & Pipelined DNNs



- Challenges from both pipelined and concurrent



Enforce **performance**
result 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

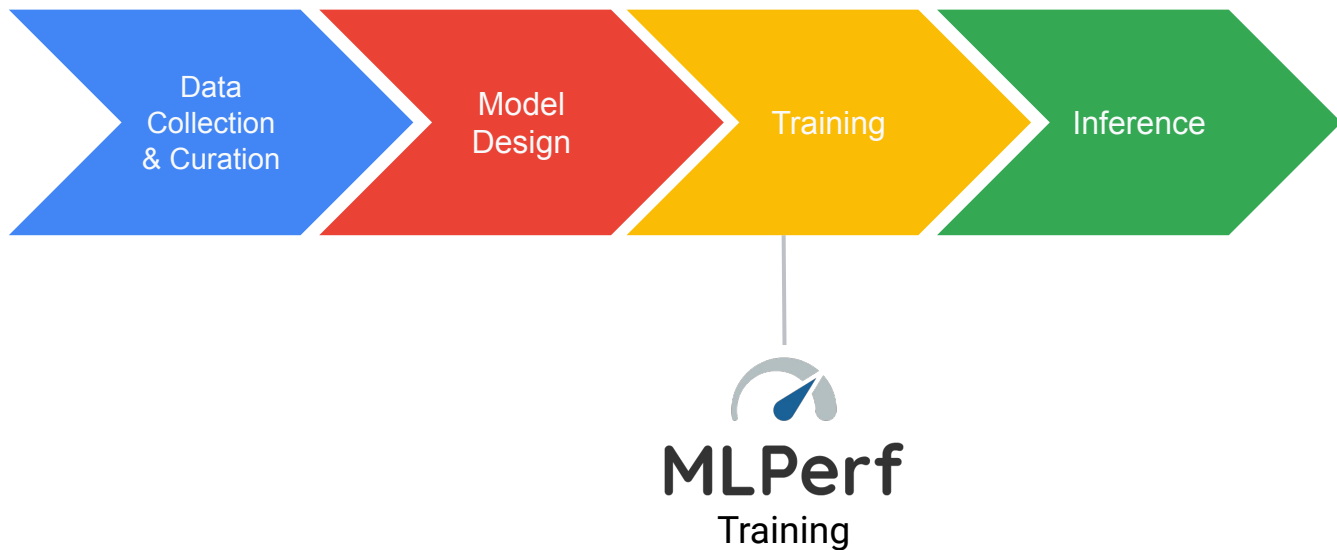


Keep **benchmarking**
affordable so that all
can participate

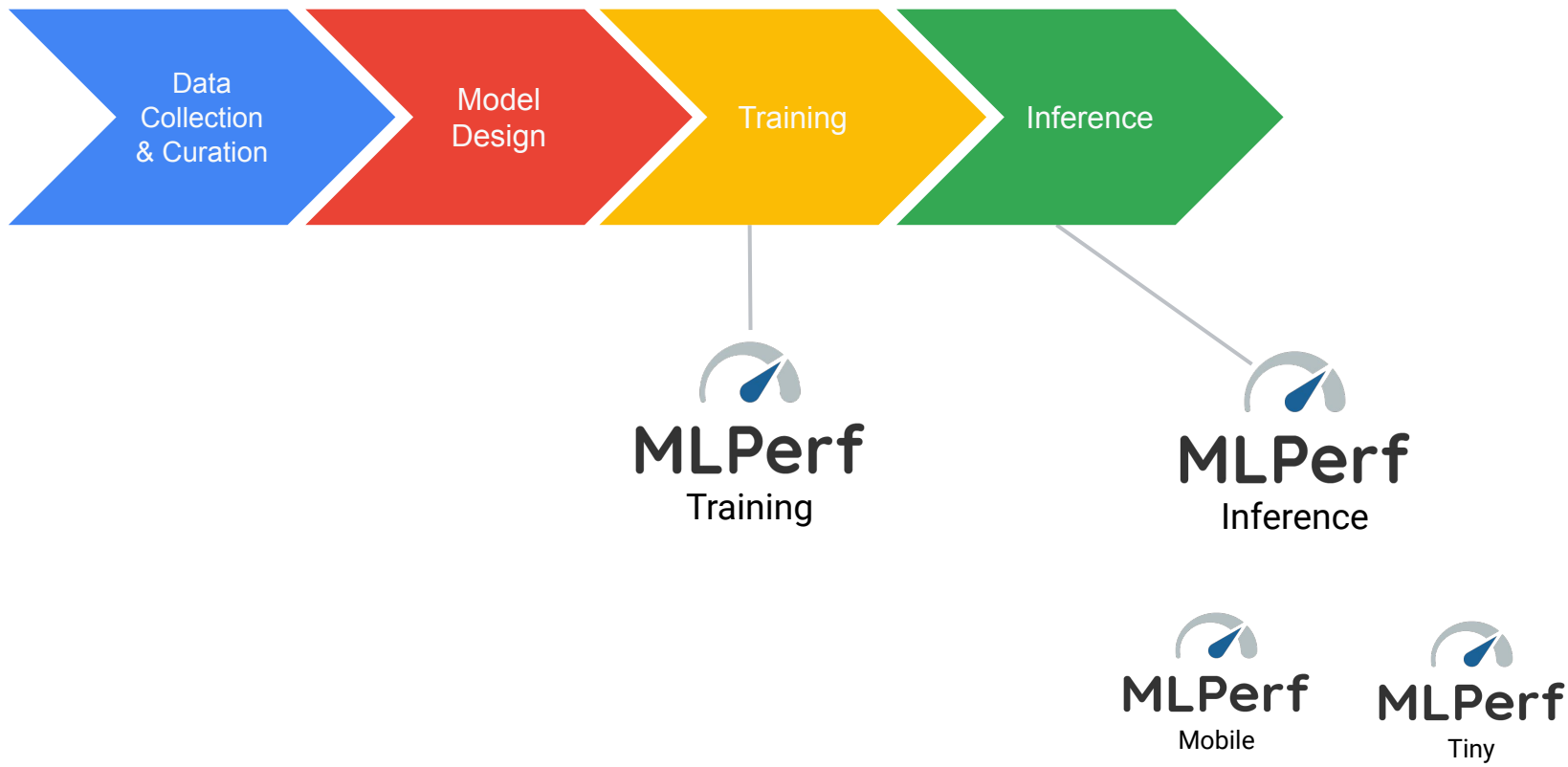
Landscape of AI benchmarking



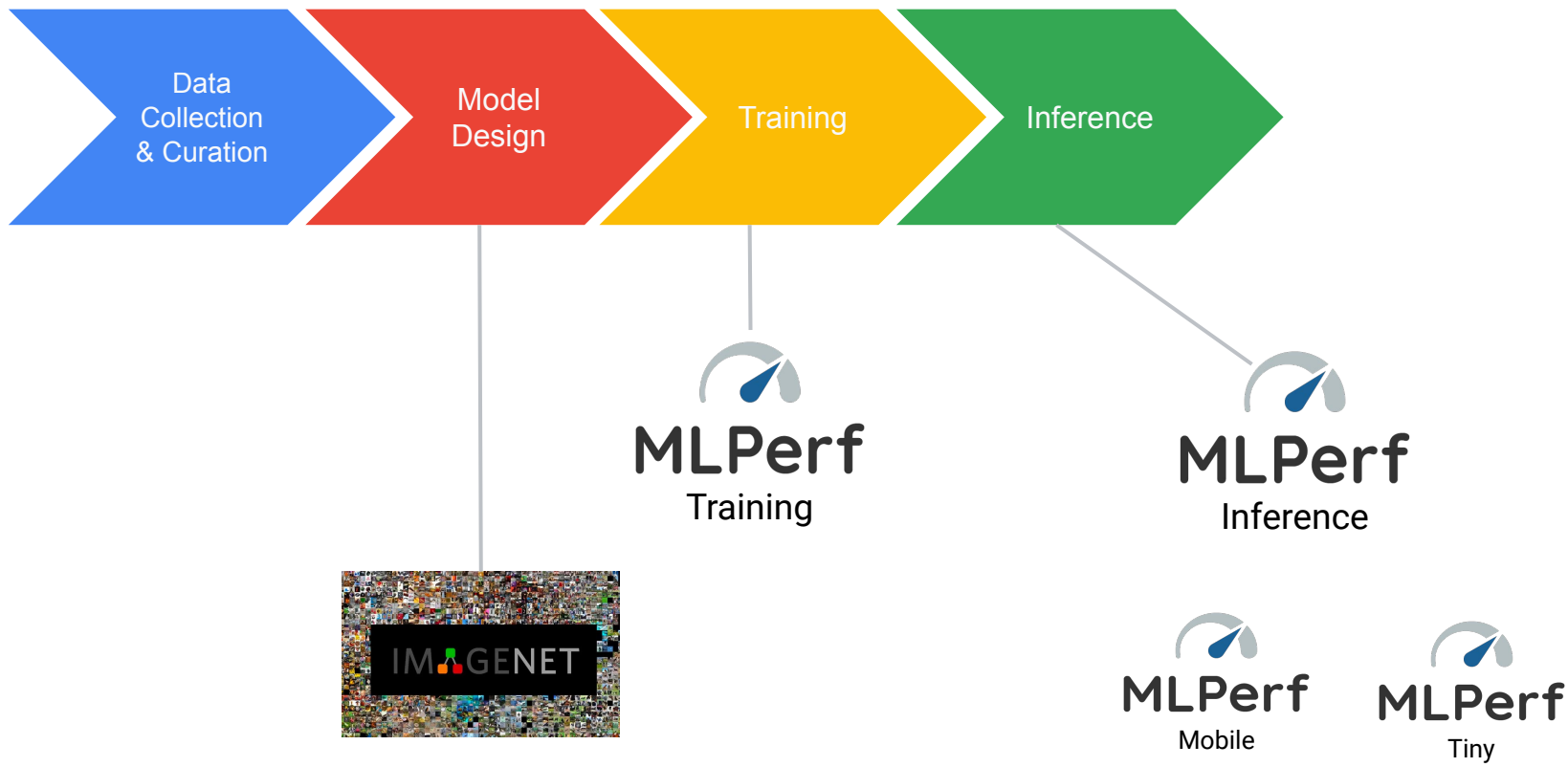
Landscape of AI benchmarking



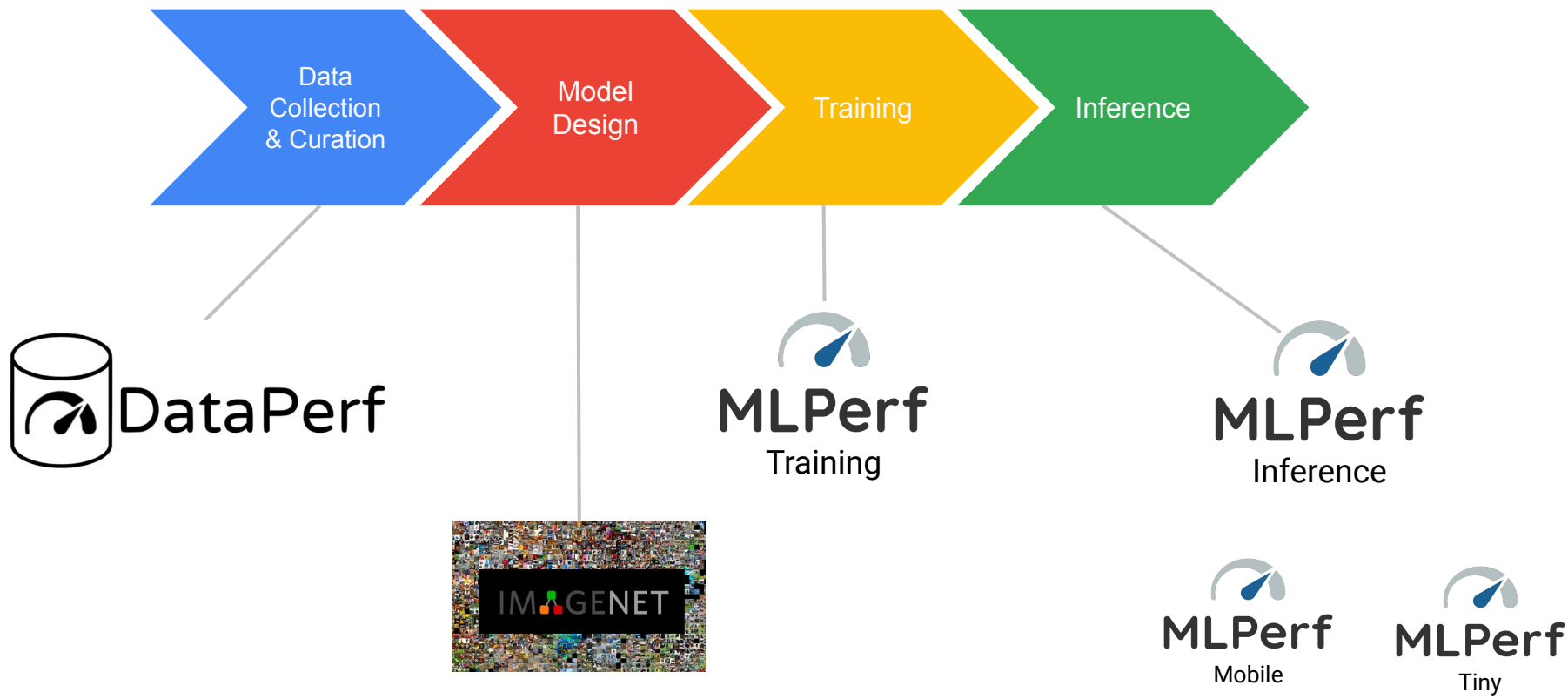
Landscape of AI benchmarking



Landscape of AI benchmarking



Landscape of AI benchmarking



MLPerf Tiny Benchmark

Colby Banbury* Vijay Janapa Reddi* Peter Torelli† Jeremy Holleman

Csaba Kiraly* Pietro Montino* David Kanter** Sebastian Ahmed††

Urmish Thakker‡ Antonio Torrini‡ Peter Warden§ Jay Cordaro †Gusep

Javier Duarte‡‡ Stephen Gibellini‡ Videet Parekh‡ Honson Tran‡‡

Niu Wenxu‡‡ Xu Xuesong‡‡

Abstract

Advancements in ultra-low-power *tiny* machine learning (TinyML) promise to unlock an entirely new class of smart applications. However, progress is limited by the lack of a widely accepted and easily reproducible benchmark for these systems. To meet this need, we present MLPerf Tiny, the first industry-standard benchmark suite for ultra-low-power tiny machine learning systems. The benchmark suite is the collaborative effort of more than 200 organizations from industry and academia and reflects the needs of the MLPerf Tiny community. It measures the accuracy, latency, and energy of machine learning inference to properly evaluate the tradeoffs between systems. Additionally, Tiny implements a modular design that enables benchmark submitters to benefit from their product, regardless of where it falls on the MLPerf Tiny deployment in a fair and reproducible manner. The suite features four benchmark categories: spotting, visual wake words, image classification, and anomaly detection.

1 Introduction

Machine learning (ML) inference on the edge is an increasingly attractive prospect for increasing energy efficiency [4], privacy, responsiveness, and autonomy of far, the field edge ML has predominantly focused on mobile inference, but it has been major strides towards expanding the scope of edge ML to ultra-low-power. The field, known as “TinyML” [1], achieves ML inference under a milliwatt (mW) traditional power barrier preventing widely distributed machine intelligence inference on-device, and near-sensor, TinyML enables greater responsiveness by avoiding the energy cost associated with wireless communication, which is often higher than that of compute [5]. Furthermore, the efficiency of TinyML enables a client device to power, always-on applications that can revolutionize the real-time collection of data. Deploying advanced ML applications at this scale requires the co-optimization of the ML deployment stack to achieve the maximum efficiency. Due to this complex

*Harvard University, †EEMBC, ‡Syntiant †UNC Charlotte ‡Google ‡Digital MLCommons †Qualcomm †STMicroelectronics †SambaNova Systems †Silicon Latent AI †Fermilab †Pang Cheng Labs

MLPerf Inference Benchmark

Vijay Janapa Reddi,* Christine Cheng,† David Kanter,‡ Guenther Schmuelling,† Carole-Jean Wu,† Brian Andersen,‡ Mark Charlebois,†† William Chou,†† Ramesh Chukka,†† Céline Deng,‡‡ Greg Diamos,‡‡ Jared Duke,‡‡ Dave Fick,‡‡ J. S. Sachin Idgunji,** Thomas B. Jablin,‡‡ Jeff Jiao,‡‡ Tom St. David Lee,‡‡ Jeffery Liao,‡‡ Anton Lokhtov,‡‡ Francis Paulius Mickiewicz,‡‡ Colin Osborne,‡‡ Gennady Pekhimenko,‡‡ Dilip Sequeira,‡‡ Ashish Srivastava,‡‡ Fei Sun,‡‡ Hanlin Tang,‡‡ Frank Wei,‡‡ Ephrem Wu,‡‡ Lingjie Xu,‡‡ Koichi George Yuan,** Aaron Zhong,‡‡ Peizhao Zhang,‡‡

*Harvard University †Intel ‡Real World Insights §Google ¶Microsoft ††Stanford University ‡‡Myrtle ‡‡Landing AI ‡‡Mythic ‡‡Advanta ‡‡Alibaba T-Head ‡‡Facebook (formerly at MediaTek) ‡‡OPPO (formerly at University of Toronto & Vector Institute) ‡‡Xilinx ‡‡Tesla ‡‡Centaur Technology ‡‡Alibaba Cloud ‡‡General Motors ‡‡Tencent ‡‡

Abstract—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 five 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. In this paper, we present our benchmarking method for evaluating ML inference systems. Driven by more than 30 organizations as well as more than 200 ML engineers and practitioners, MLPerf prescribes a set of rules and best practices to ensure comparability across systems with widely differing architectures. The first call for submissions garnered more than 600 reproducible inference-performance measurements from 14 organizations, representing over 30 systems that showcase a wide range of capabilities. The submissions attest to the benchmark’s flexibility and adaptability.

Index Terms—Machine Learning, Inference, Benchmarking

I. INTRODUCTION

Machine learning (ML) powers a variety of applications from computer vision ([20], [18], [34], [29]) and natural-language processing ([50], [16]) to self-driving cars ([55], [6]) and autonomous robotics [32]. Although ML-model training has been a development bottleneck and a considerable expense [4], inference has become a critical workload. Models can serve as many as 200 trillion queries and perform over 6 billion translations a day [31]. To address these growing computational demands, hardware, software, and system developers have focused on inference performance for a variety

of use cases by ware. Estimates specialized inference on 20 companies. Each ML system trading off latency result is many data sets, from inference engine performance not including but not, object detection automatic speed datations. Even many ML model of scenarios for continuously at multiple camera ML tasks have time processing model’s function specific, and this To quantify the that is architect

Both academic ML inference AIMatrix [3], E industry, as well and DAWNBe substantial computational developed with signers. As a machine learning

DataPerf:

Benchmarks for Data-Centric AI Development

Mark Mazumder¹, Colby Banbury¹, Xiaozhe Yao², Bojan Karlas², William Gaviria Rojas³, Sudnya Diamos³, Greg Diamos⁴, Lynn He⁵, Alicia Parrish⁹, Hannah Rose Kirk¹⁸, Jessica Quayle¹, Charvi Rastogi¹², Douwe Kiela^{10,22}, David Jurado²¹, David Kanter⁷, Rafael Mosquera^{7,21}, Juan Ciro^{7,21}, Lora Aroyo⁹, Bilge Acun⁹, Lingjiao Chen¹⁰, Mehul Smriti Raj⁹, Max Bartolo^{17,20}, Sabri Eyuboglu¹⁰, Amirata Ghorbani¹⁰, Emmett Goodman¹⁰, Oana Inel¹⁹, Tariq Kane^{3,9}, Christine R. Kirkpatrick¹¹, Tzu-Sheng Kuo¹², Jonas Mueller¹³, Tristan Thrush¹⁰, Joaquin Vanschoren¹⁴, Margaret Warren¹⁵, Adina Williams⁸, Serena Yeung¹⁰, Newsha Ardalani⁸, Praveen Paritosh⁷, Lilith Bat-Leah⁷, Ce Zhang⁷, James Zou¹⁰, Carole-Jean Wu⁸, Cody Coleman³, Andrew Ng^{4,5,10}, Peter Mattson⁹, and Vijay Janapa Reddi¹

¹Harvard University, ²ETH Zurich, ³Coactive.AI, ⁴Landing AI, ⁵DeepLearning.AI, ⁷MLCommons, ⁸Meta, ⁹Google, ¹⁰Stanford University, ¹¹San Diego Supercomputer Center, UC San Diego, ¹²Carnegie Mellon University, ¹³Cleanlab, ¹⁴Eindhoven University of Technology, ¹⁵Institute for Human and Machine Cognition, ¹⁶Kaggle, ¹⁷Cohere, ¹⁸University of Oxford, ¹⁹University of Zurich, ²⁰University College London, ²¹Factored, ²²Contextual AI

Abstract

Machine learning research has long focused on models rather than datasets, and prominent datasets are used for common ML tasks without regard to the breadth, difficulty, and faithfulness of the underlying problems. Neglecting the fundamental importance of data has given rise to inaccuracy, bias, and fragility in real-world applications, and research is hindered by saturation across existing dataset benchmarks. In response, we present DataPerf, a community-led benchmark suite for evaluating ML datasets and data-centric algorithms. We aim to foster innovation in data-centric AI through competition, comparability, and reproducibility. We enable the ML community to iterate on datasets, instead of just architectures, and we provide an open, online platform with multiple rounds of challenges to support this iterative development. The first iteration of DataPerf contains five benchmarks covering a wide spectrum of data-centric techniques, tasks, and modalities in vision, speech, acquisition, debugging, and diffusion prompting, and we support hosting new contributed benchmarks from the community. The benchmarks, online evaluation platform, and baseline implementations are open source, and the MLCommons Association will maintain DataPerf to ensure long-term benefits to academia and industry.

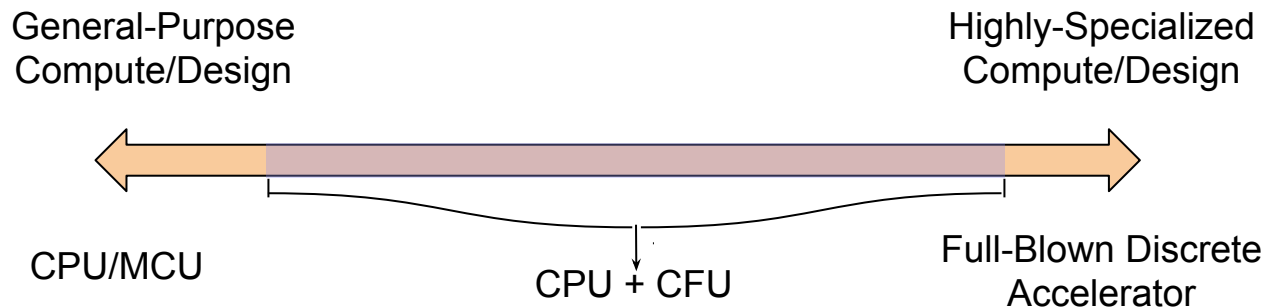
1 Introduction

Machine learning research has overwhelmingly focused on improving models rather than on improving datasets. Large public datasets such as ImageNet [14], Freebase [7], Switchboard [22], and SQuAD [44] serve as compasses for benchmarking model performance. Consequently, researchers eagerly adopt the largest existing dataset without fully considering its breadth, difficulty and fidelity to the underlying problem. Critically, better data quality [2] is increasingly necessary to improve generalization, avoid bias, and aid safety in data cascades [48]. Without high-quality training data models can exhibit performance discrepancies leading to reduced accuracy and persistent fairness

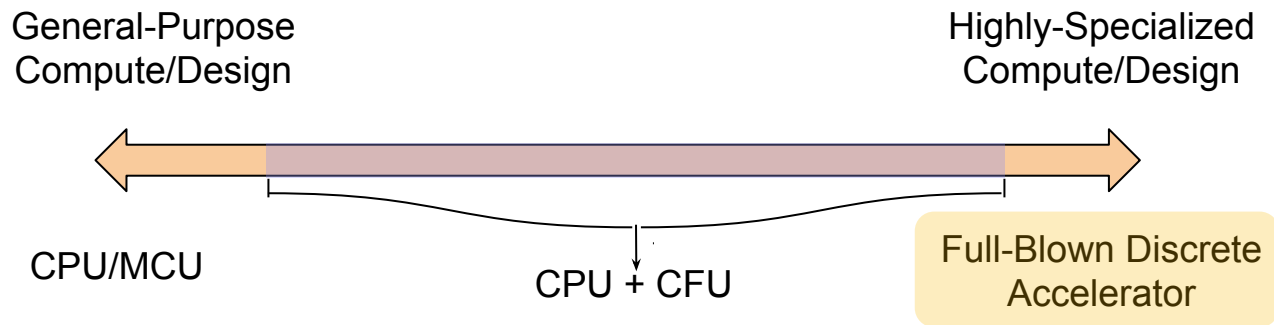


Paper Discussions

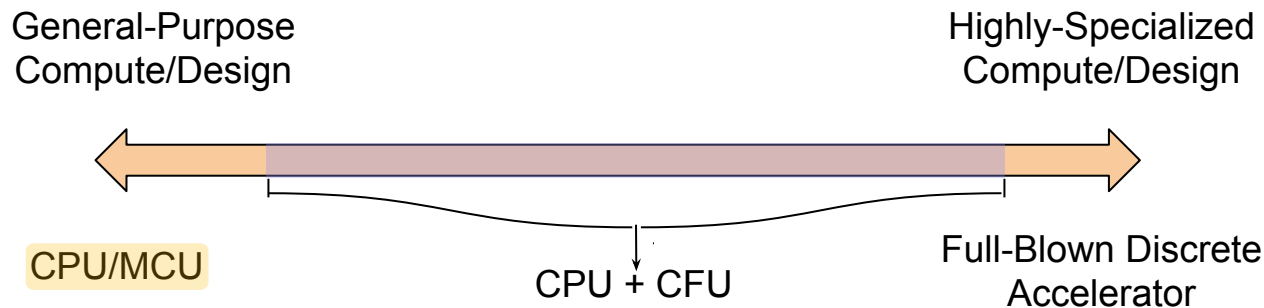
The Need for Agile and Full-Stack Frameworks



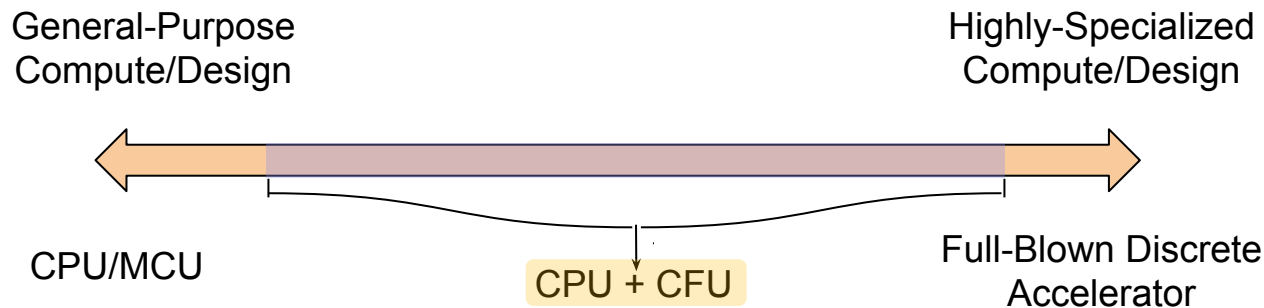
The Need for Agile and Full-Stack Frameworks



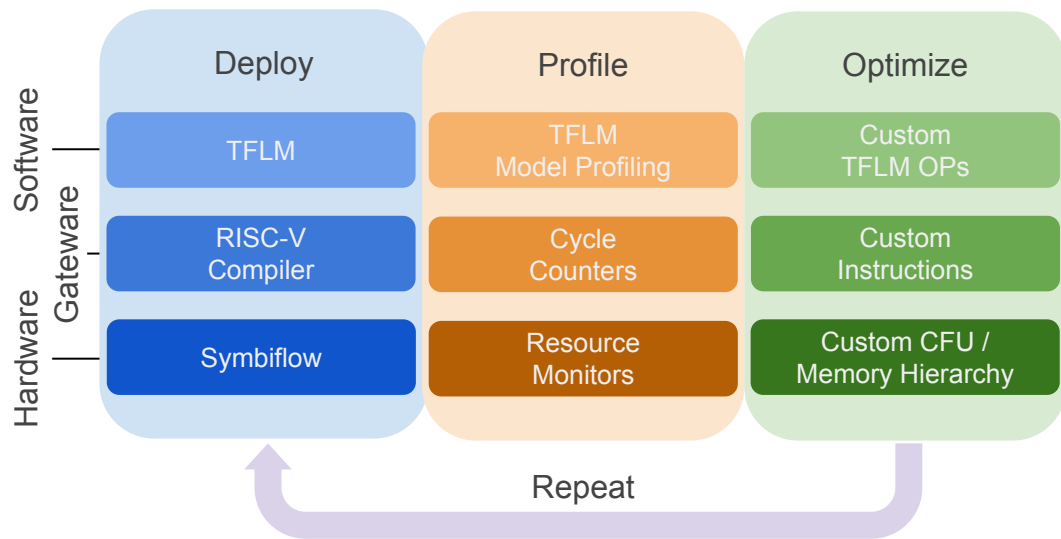
The Need for Agile and Full-Stack Frameworks



The Need for Agile and Full-Stack Frameworks

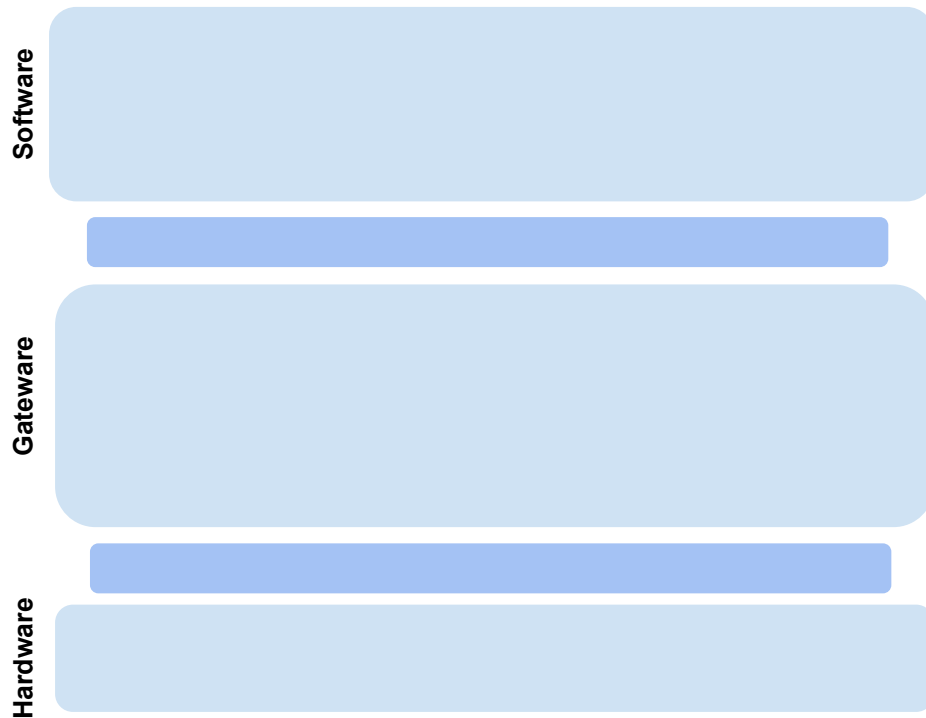


The Need for Agile and Full-Stack Frameworks



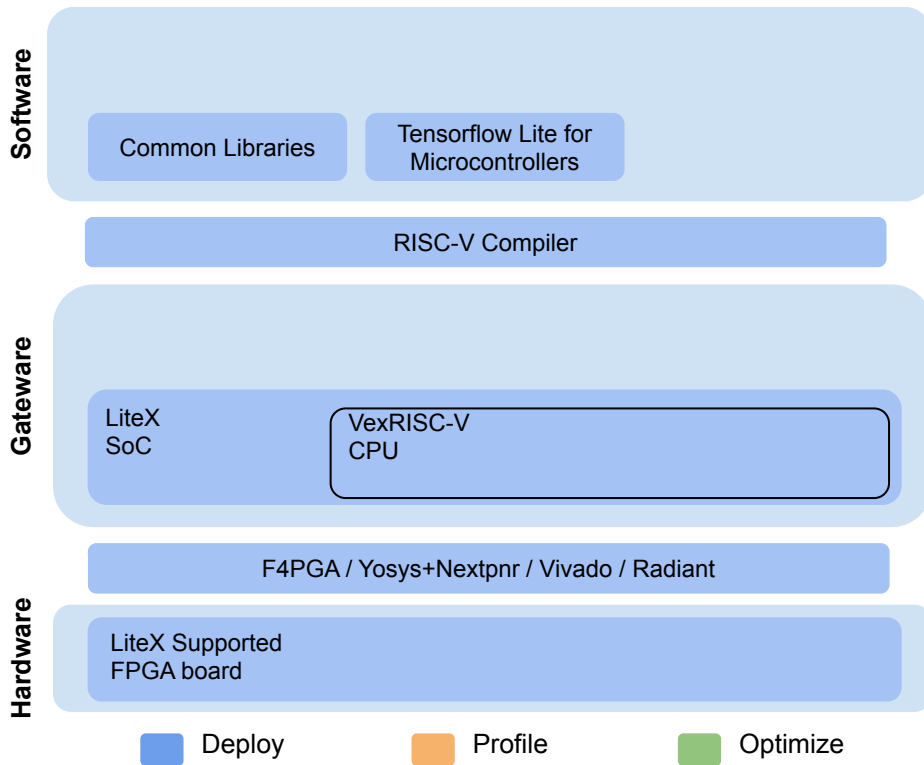
Agile Design Methodology

The Need for Agile and Full-Stack Frameworks



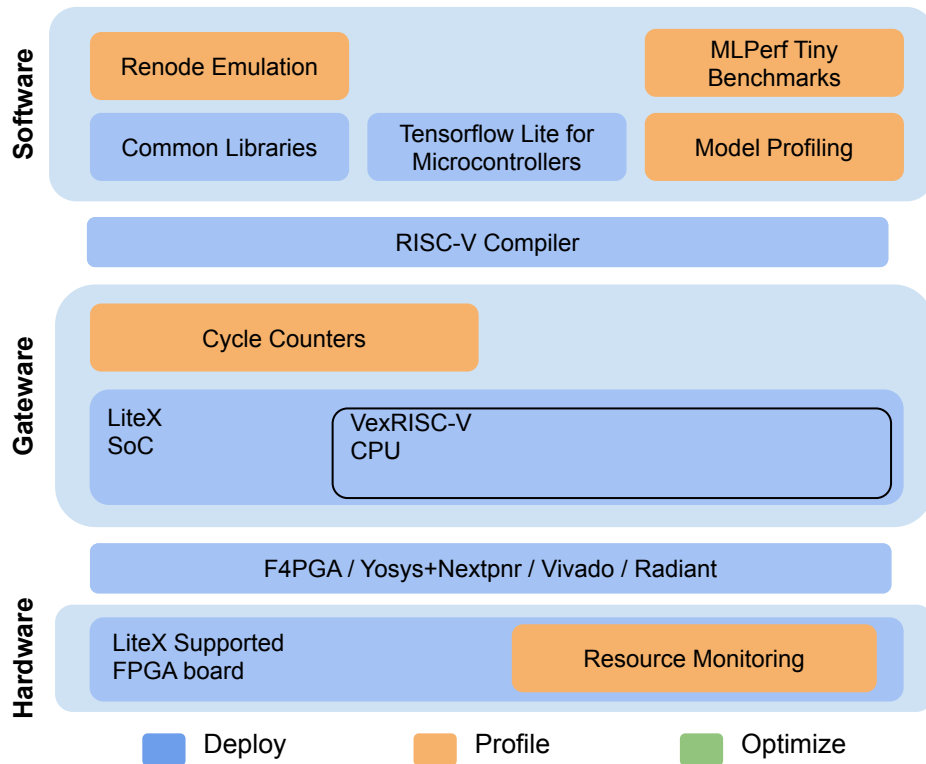
Full-Stack Open-Source Framework

The Need for Agile and Full-Stack Frameworks



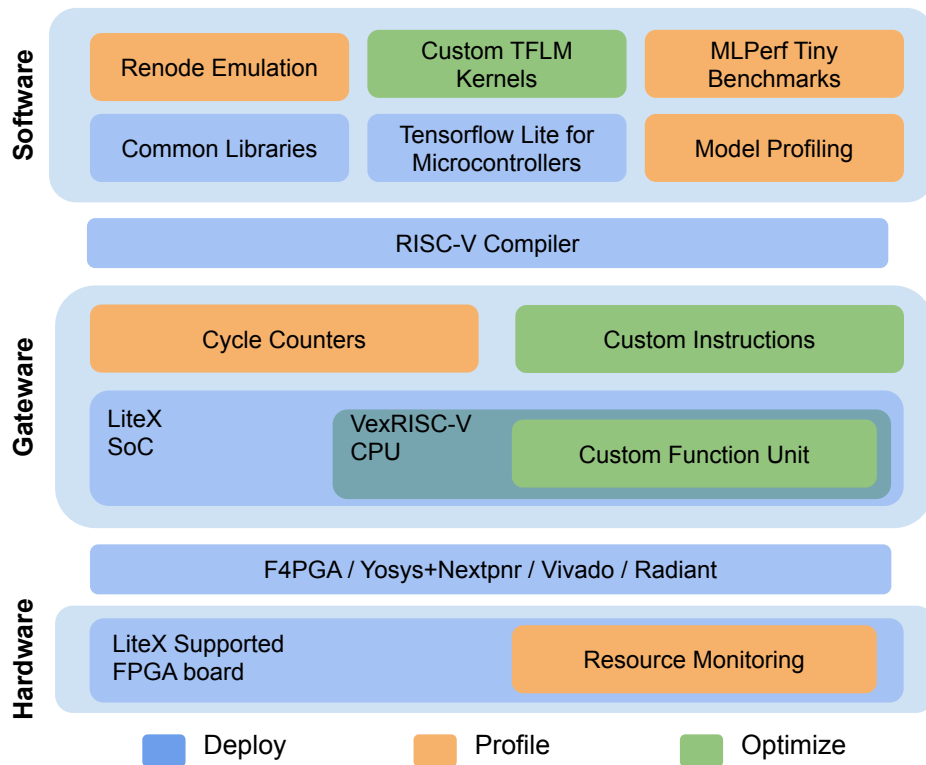
Full-Stack Open-Source Framework

The Need for Agile and Full-Stack Frameworks



Full-Stack Open-Source Framework

The Need for Agile and Full-Stack Frameworks



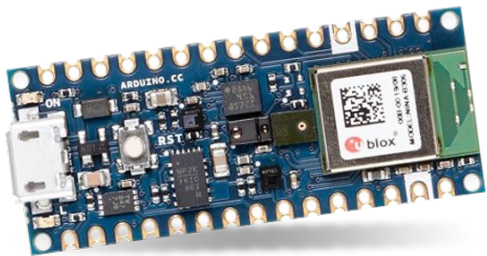
Full-Stack Open-Source Framework

Related Works

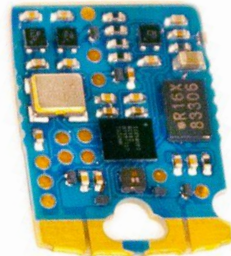
	Open Source	Full Stack	Full SoC	Tightly Coupled/ Specialized ISA	Fine-Grained Accelerated ML Ops	Hardware & Engineer In-The-Loop	Stock Compiler	Automated CPU↔Accelerator Design Space Exploration	TinyML Focus
CFU Playground	✓	✓	✓	✓	✓	✓	✓	✓	✓
Chipyard [20]	✓	✓	✓	✓	✓	✓	✓	✗	✗
Centrifuge [21]	✓	✓	✓	✓	✓	✓	✗	✗	✗
Embedded Scalable Platform [22]	✓	✓	✓	✗	✗	✓	✓	✗	✗
Gemmini [23]	✓	✓	✓	✗	✗	✗	✓	✗	✗
hls4ml [24]	✓	✗	✓	✗	✗	✓	✗	✗	✓
Deepburning [25]	✓	✓	✓	✗	✗	✗	✗	✗	✗
DNN-Weaver [26]	✓	✓	✗	✓	✗	✗	✗	✗	✗
DNN-Builder [27]	✓	✓	✗	✗	✗	✗	✗	✗	✗
FINN [28]	✓	✗	✗	✗	✗	✗	✗	✗	✗

TABLE III: Comparison of CFU Playground with open-source toolchains supporting custom hardware design for ML workloads. CFU Playground focuses on *open-source* development across the *full system* stack, while providing varying levels of flexibility for hardware and software (co-)design.

Accelerating TinyML on FPGAs

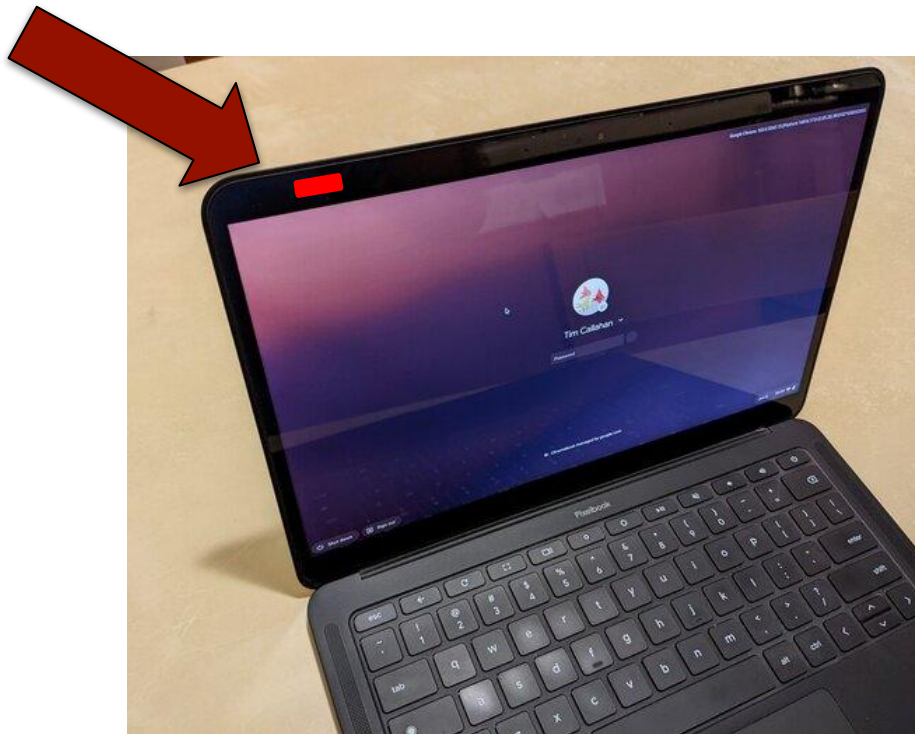


MCUs: KBs of RAM, Fixed/slow processor



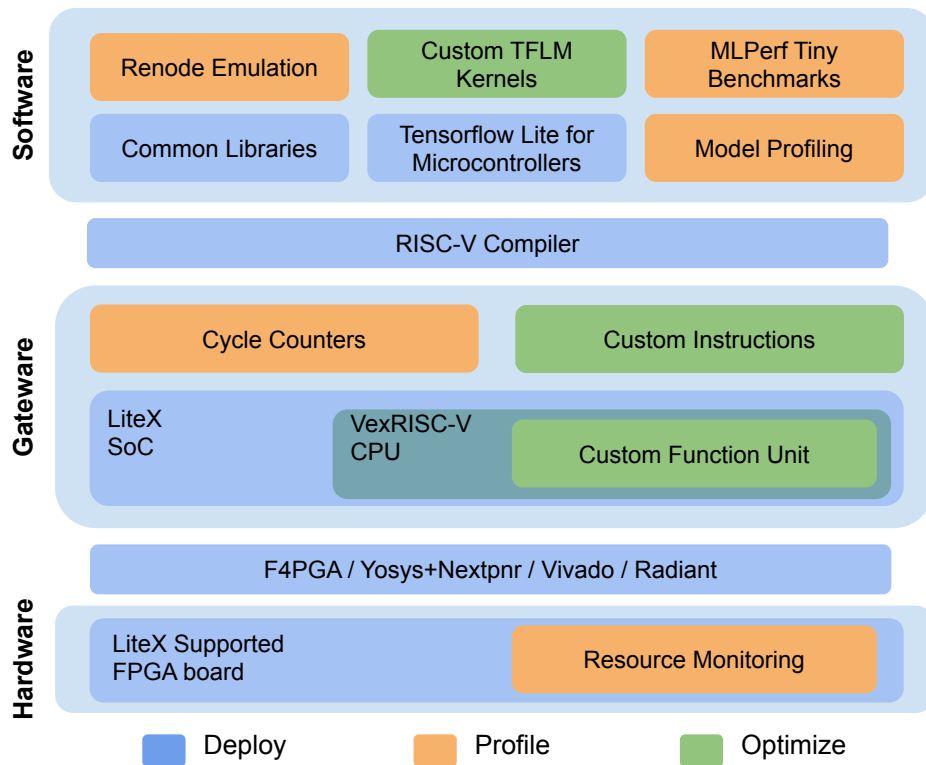
Specialized Hardware Customization (on FPGAs)

Real World Use Case



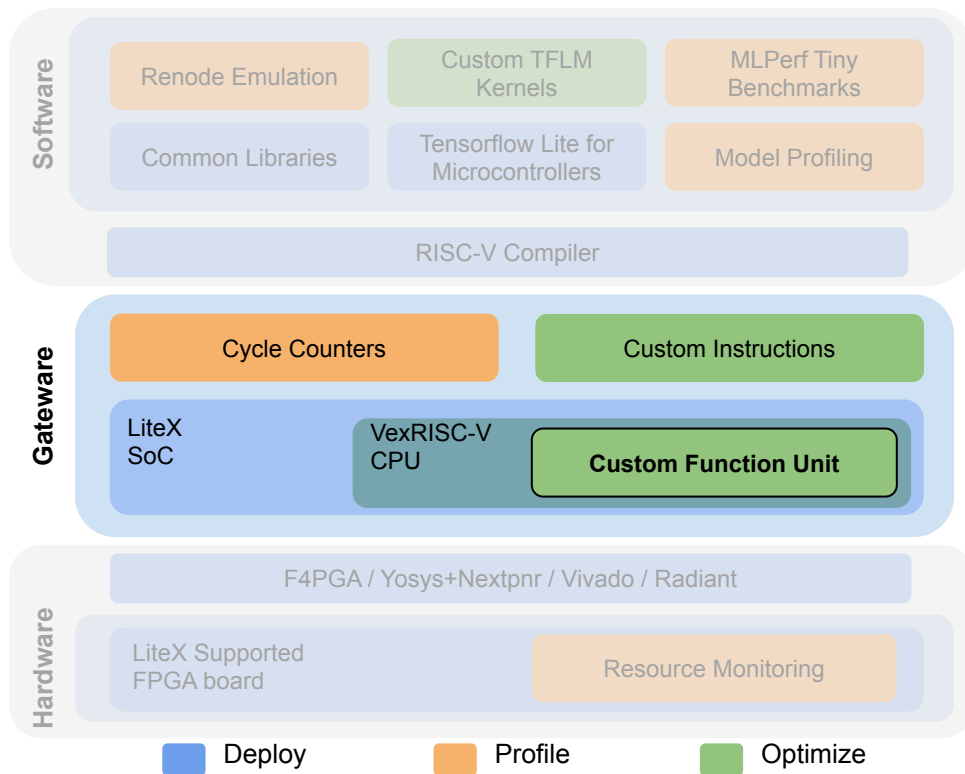
Chromebook Sensor Designed with CFU Playground

CFU Playground



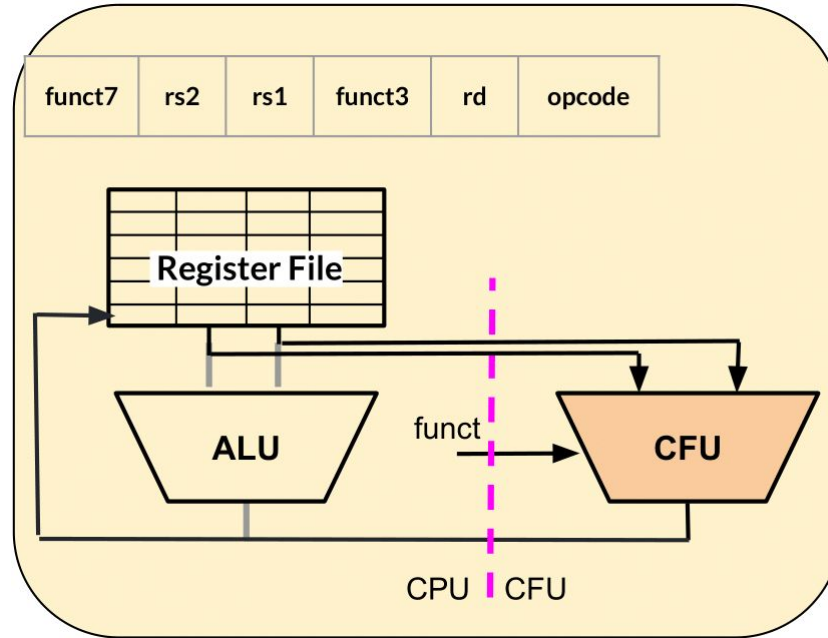
Full-Stack Open-Source Framework

CFU Playground



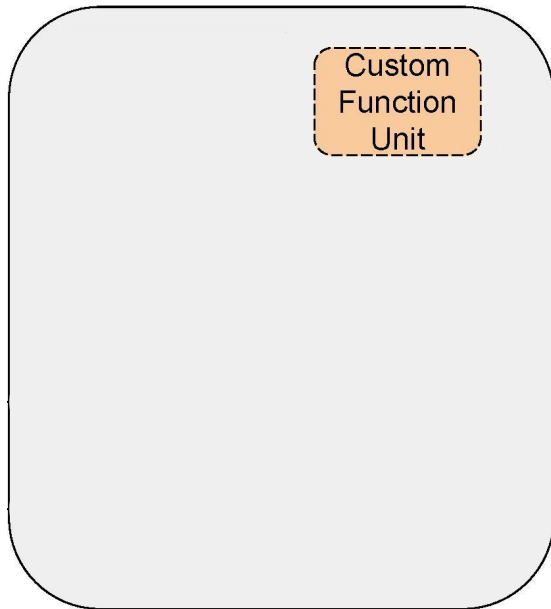
Full-Stack Open-Source Framework

Custom Function Units (CFU)

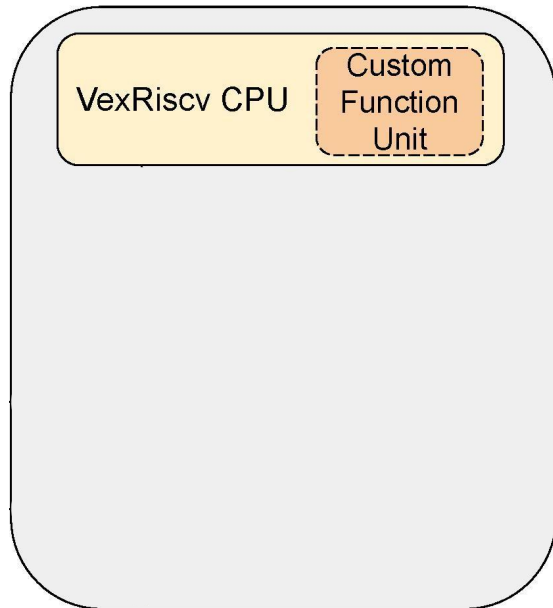


Acceleration via Custom Function Unit (CFU)

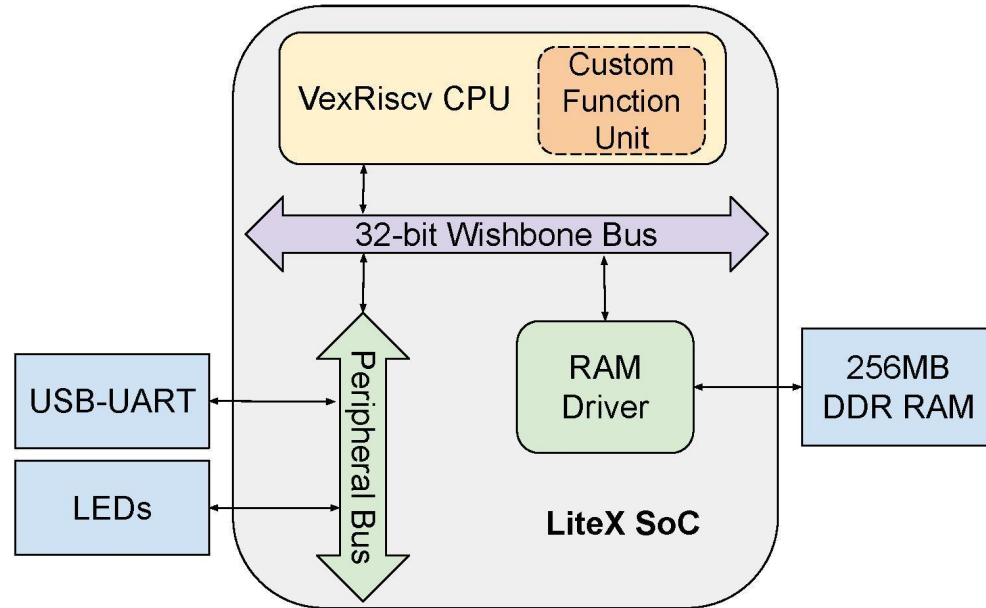
System On-Chip (SoC) Integration



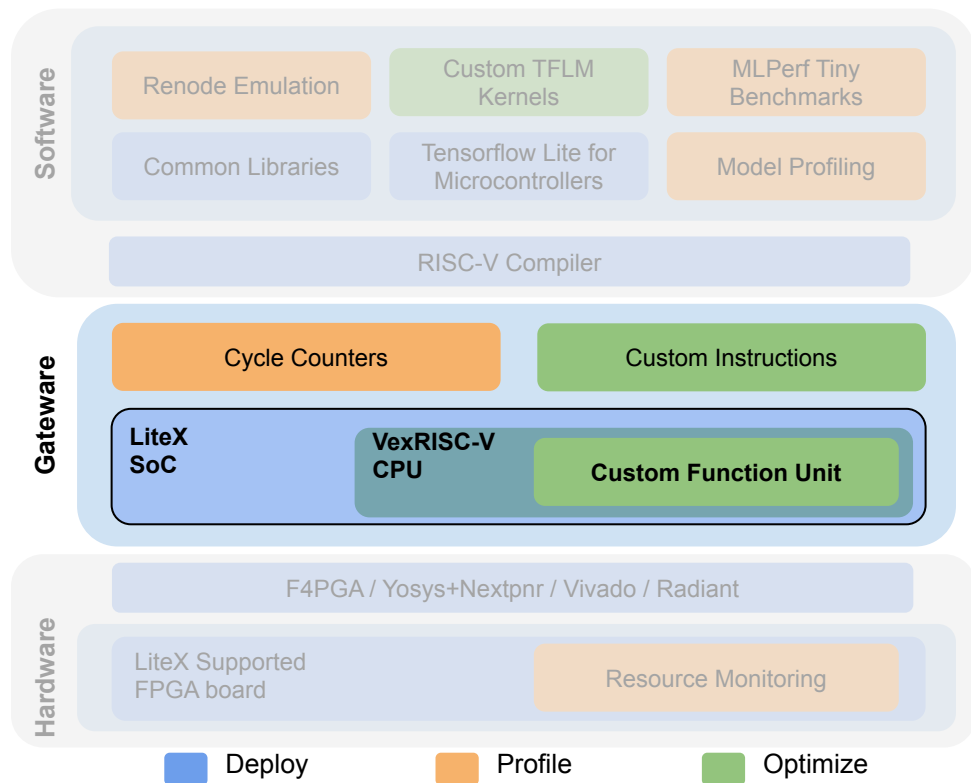
System On-Chip (SoC) Integration



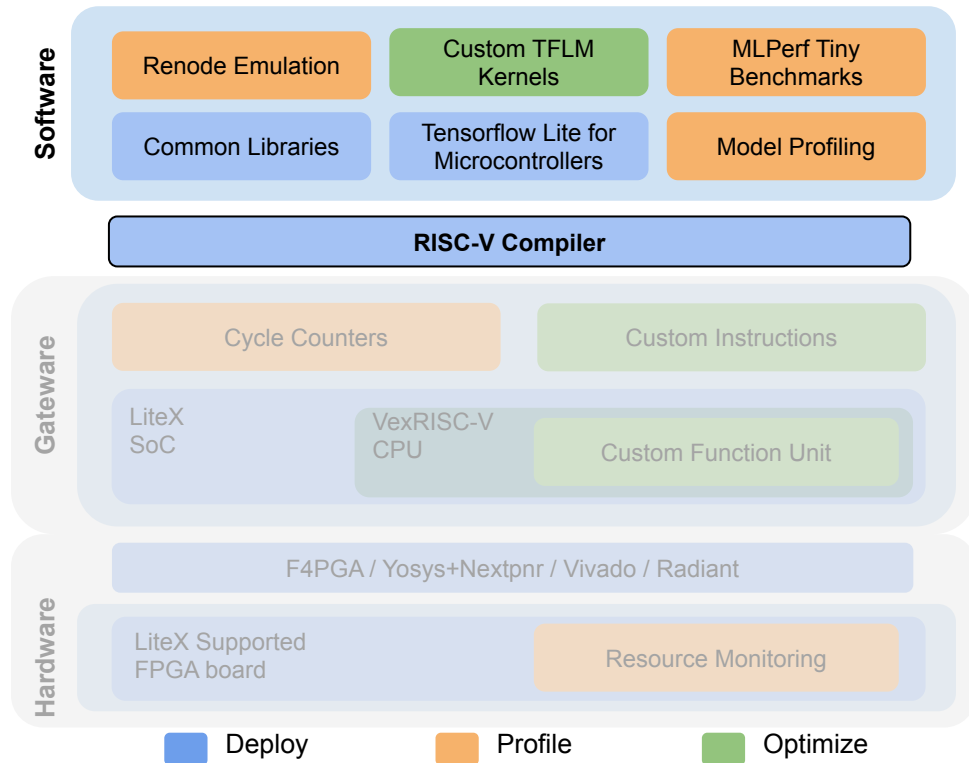
System On-Chip (SoC) Integration



System On-Chip (SoC) Integration



CFU Software Interface

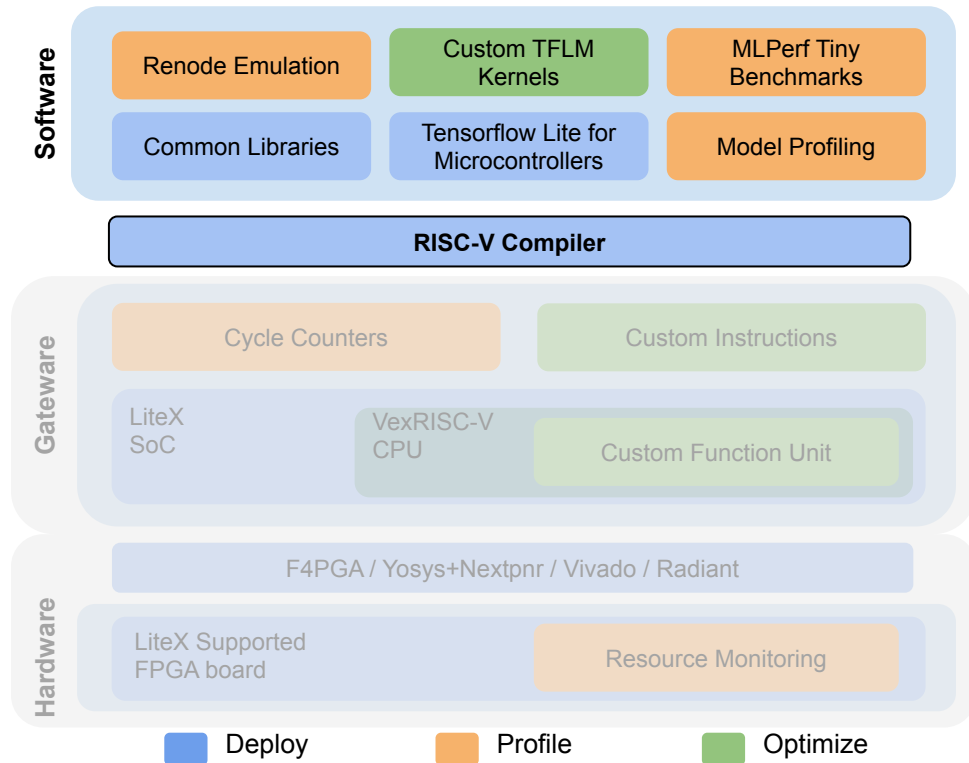


Access new instruction as C function call:
`rs1t = cfu_op(funcnt3, funct7, op1, op2);`

Compile-time constants

C / C++ variables / expressions

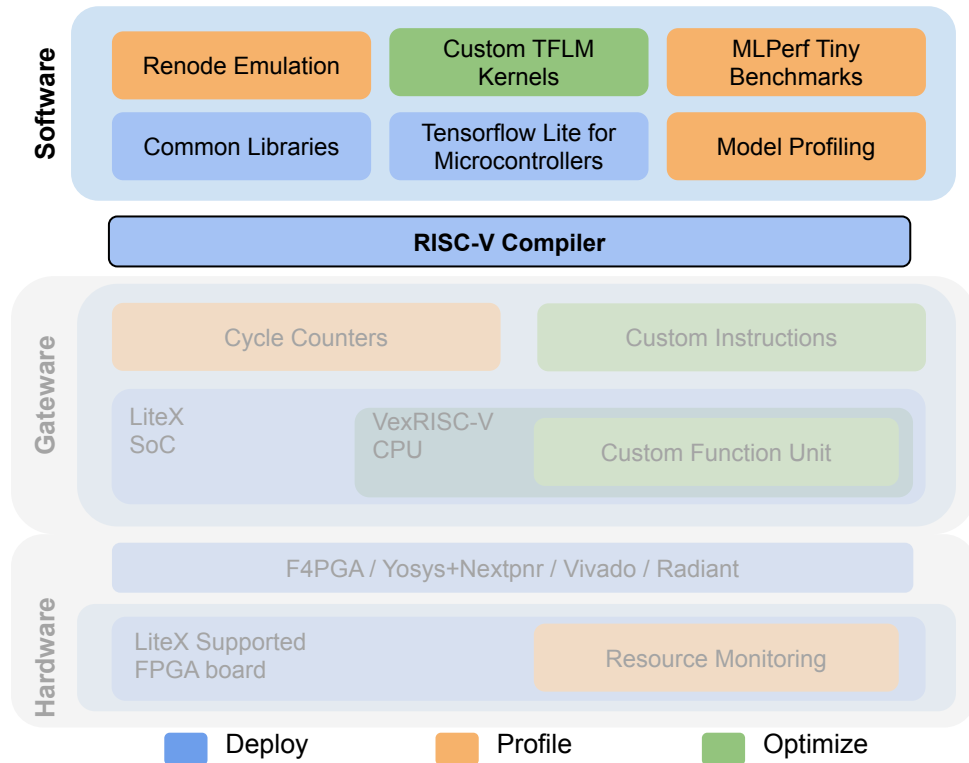
CFU Software Interface



Custom instruction macros intermix with plain C code

```
t1 = *x;  
t2 = cfu_op(0, 0, t1, b);  
t3 = cfu_op(1, 0, t2, b);  
*x = t3;
```

CFU Software Interface



Custom instruction macros intermix with plain C code

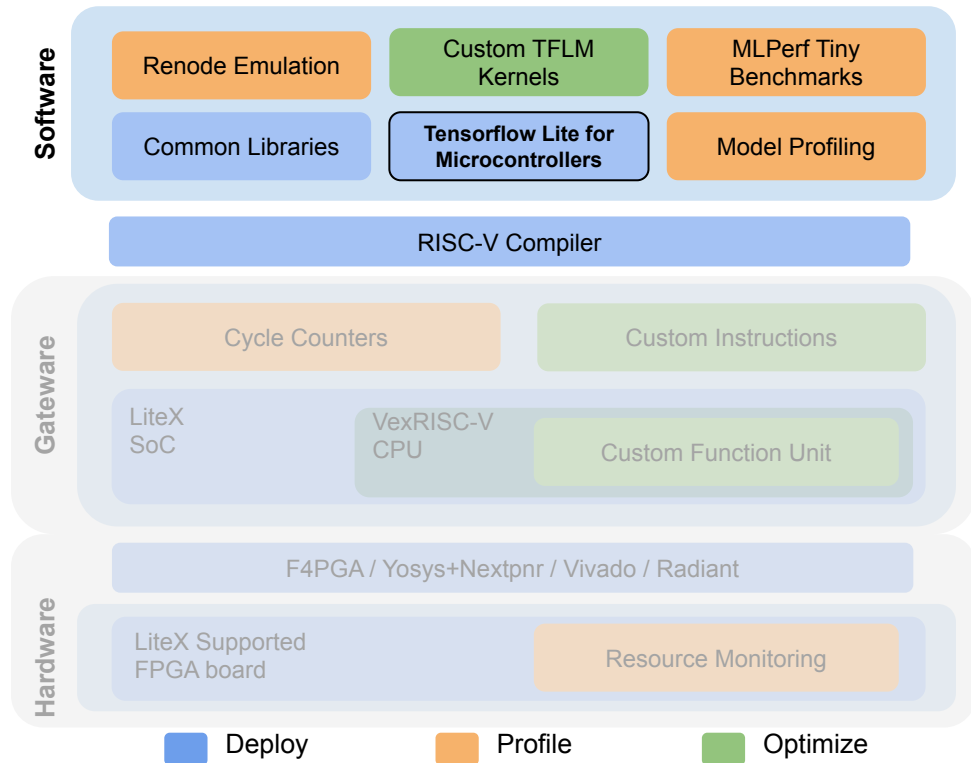
```
t1 = *x;  
t2 = cfu_op(0, 0, t1, b);  
t3 = cfu_op(1, 0, t2, b);  
*x = t3;
```

Compiled and disassembled:

400001a0:	00812783	lw	a5,8(sp)
400001a4:	00d7878b	cfu[0,0]	a5, a5, a3
400001a8:	00d7978b	cfu[0,1]	a5, a5, a3
400001ac:	00f12423	sw	a5,8(sp)

No overhead!

ML Deployment Framework



```
const int32_t input_offset = params.input_offset; // r = s(q - Z)

for (int batch = 0; batch < batches; ++batch) {
  for (int out_y = 0; out_y < output_height; ++out_y) {
    const int in_y_origin = (out_y * stride_height) - pad_height;
    for (int out_x = 0; out_x < output_width; ++out_x) {
      const int in_x_origin = (out_x * stride_width) - pad_width;
      for (int out_channel = 0; out_channel < output_depth; ++out_channel) {
        int32_t acc = 0;
        for (int filter_y = 0; filter_y < filter_height; ++filter_y) {
          const int in_y = in_y_origin + dilation_height_factor * filter_y;
          for (int filter_x = 0; filter_x < filter_width; ++filter_x) {
            const int in_x = in_x_origin + dilation_width_factor * filter_x;

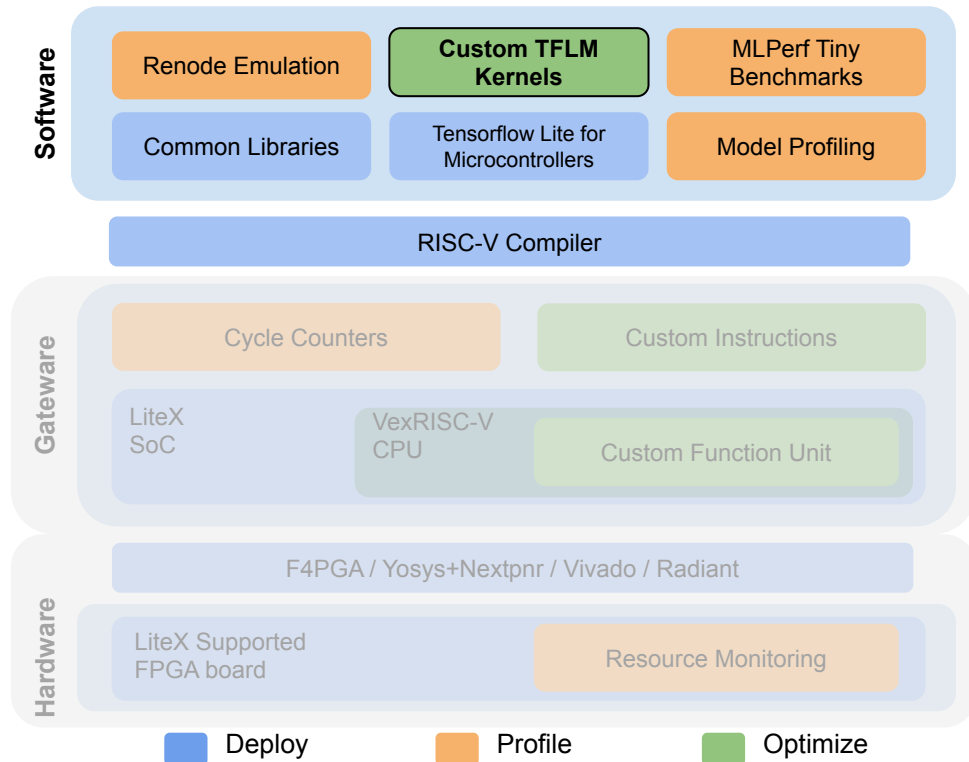
            // Zero padding by omitting the areas outside the image.
            const bool is_point_inside_image =
              (in_x >= 0) && (in_x < input_width) && (in_y >= 0) &&
              (in_y < input_height);

            if (!is_point_inside_image) {
              continue;
            }

            for (int in_channel = 0; in_channel < input_depth; ++in_channel) {
              int32_t input_val = input_data[Offset(input_shape, batch, in_y,
                                                       in_x, in_channel)];
              int32_t filter_val = filter_data[Offset(
                filter_shape, out_channel, filter_y, filter_x, in_channel)];

              acc += filter_val * (input_val + input_offset);
            }
          }
        }
      }
    }
  }
  (use acc)
}
```

Accelerated Kernels



```
const int32_t input_offset = params.input_offset; // r = s(q - Z)
```

```
// CFU: copy input_offset into the CFU
```

```
cfu_init_offset(input_offset);
```

```
for (int batch = 0; batch < batches; ++batch) {
  for (int out_y = 0; out_y < output_height; ++out_y) {
    const int in_y_origin = (out_y * stride_height) - pad_height;
    for (int out_x = 0; out_x < output_width; ++out_x) {
      const int in_x_origin = (out_x * stride_width) - pad_width;
      for (int out_channel = 0; out_channel < output_depth; ++out_channel) {
```

```
//int32_t acc = 0;
```

```
// CFU: set the CFU internal acc to ZERO
```

```
cfu_clear_acc();
```

```
for (int filter_y = 0; filter_y < filter_height; ++filter_y) {
  const int in_y = in_y_origin + dilation_height_factor * filter_y;
  for (int filter_x = 0; filter_x < filter_width; ++filter_x) {
    const int in_x = in_x_origin + dilation_width_factor * filter_x;
```

```
...
```

```
for (int in_channel = 0; in_channel < input_depth; ++in_channel) {
  int32_t input_val = input_data[Offset(input_shape, batch, in_y,
                                         in_x, in_channel)];
  int32_t filter_val = filter_data[Offset(
    filter_shape, out_channel, filter_y, filter_x, in_channel)];
```

```
// acc += filter_val * (input_val + input_offset);
```

```
// CFU: add-multiply-accumulate in the CFU
```

```
cfu_macc_with_offset(filter_val, input_val);
```

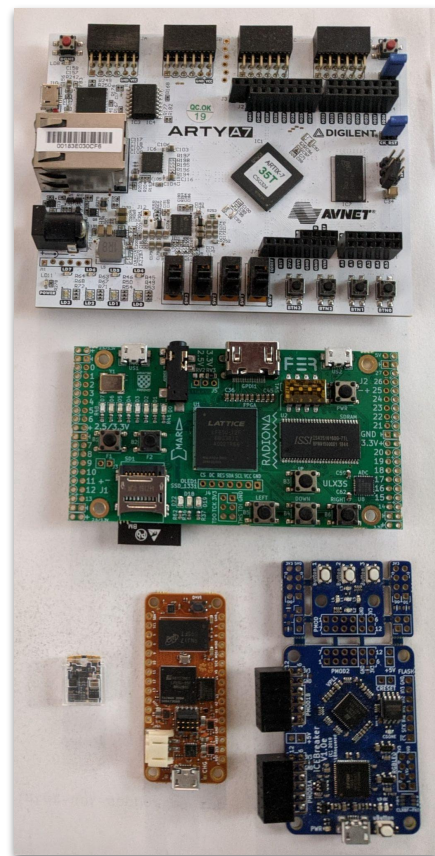
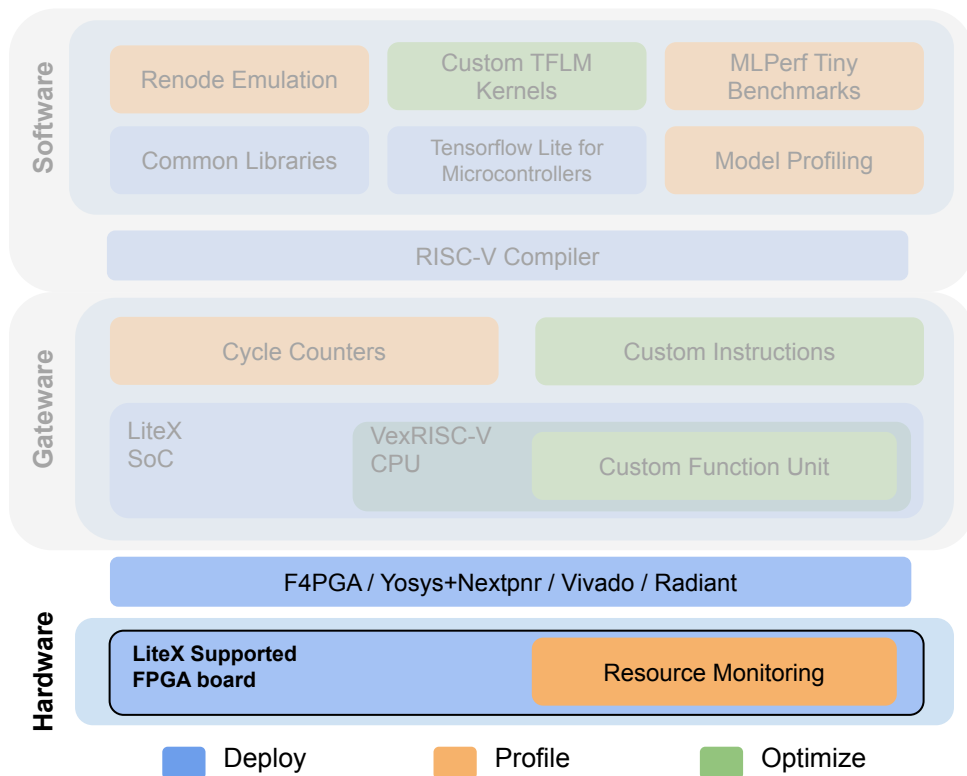
```
}
```

```
}
```

```
// CFU: retrieve final acc value from the CFU
```

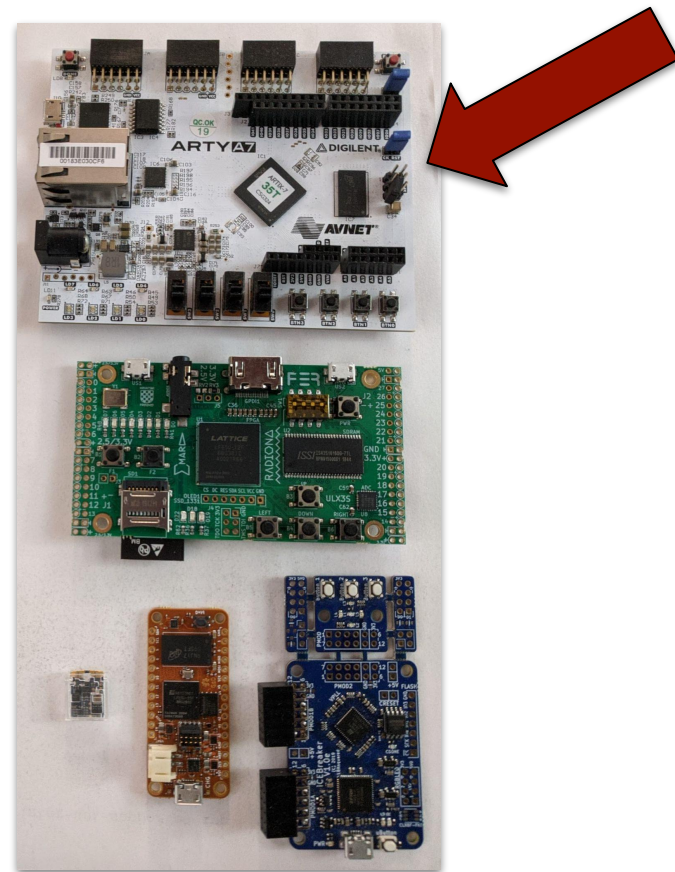
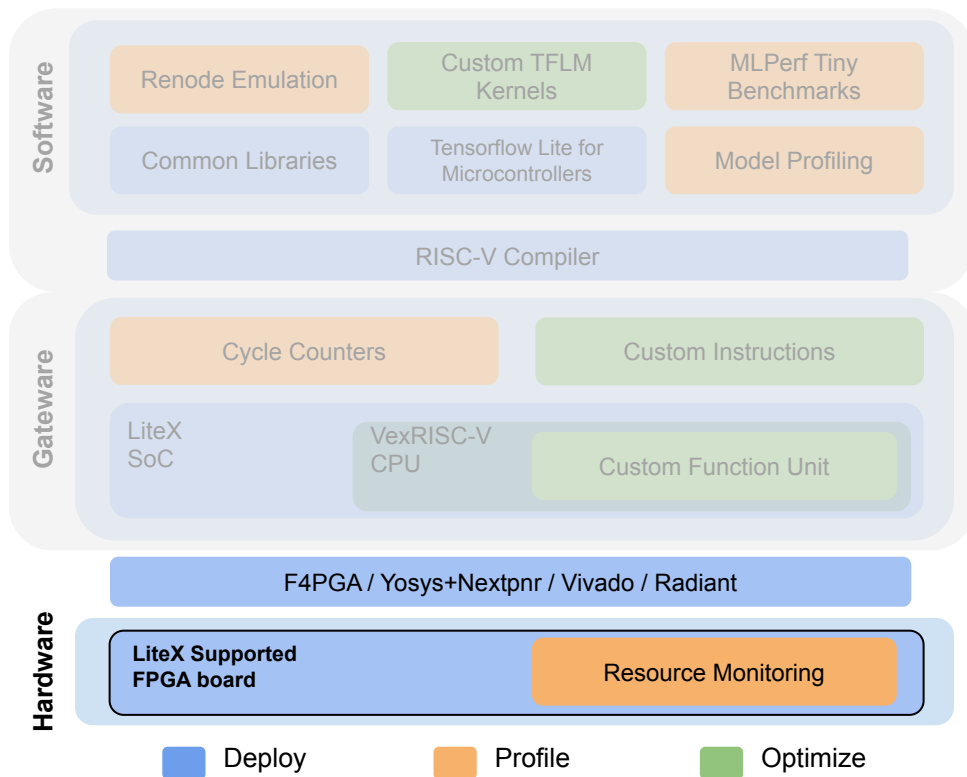
```
int32_t acc = cfu_get_acc();
```

Hardware In-The-Loop



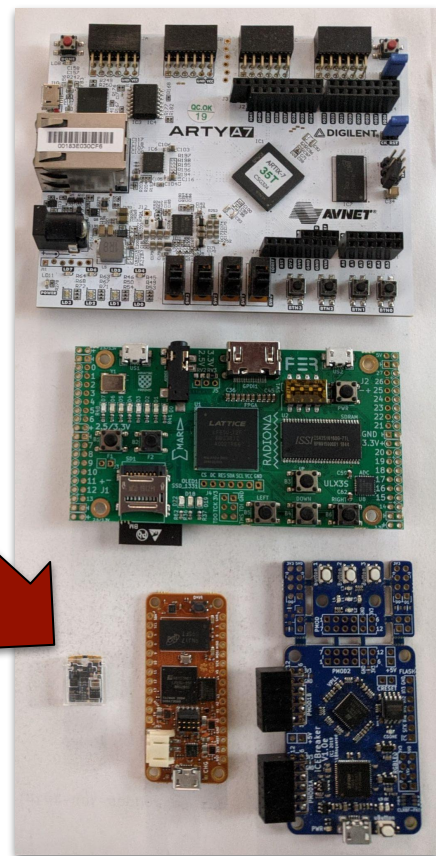
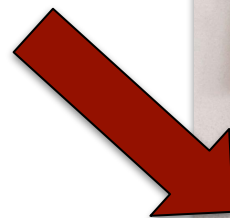
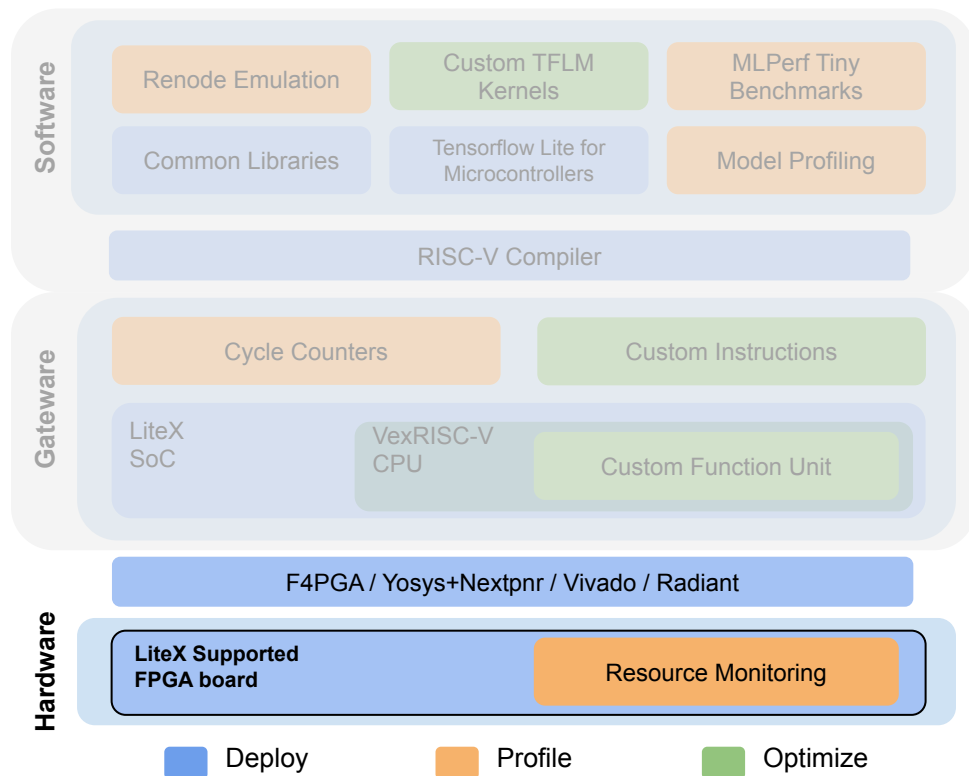
Diverse Family of FPGA Boards

Hardware In-The-Loop



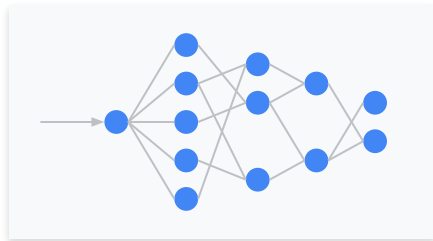
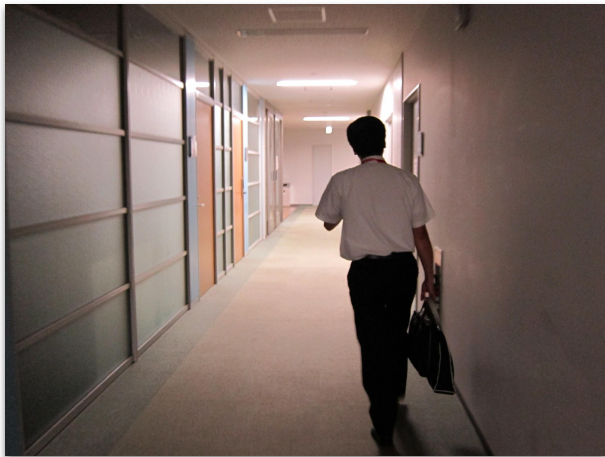
Diverse Family of FPGA Boards

Hardware In-The-Loop



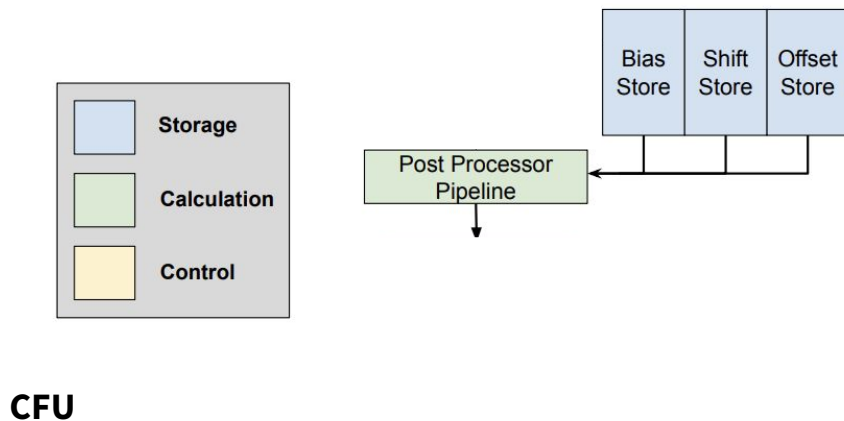
Diverse Family of FPGA Boards

FPGA Acceleration for Image Classification

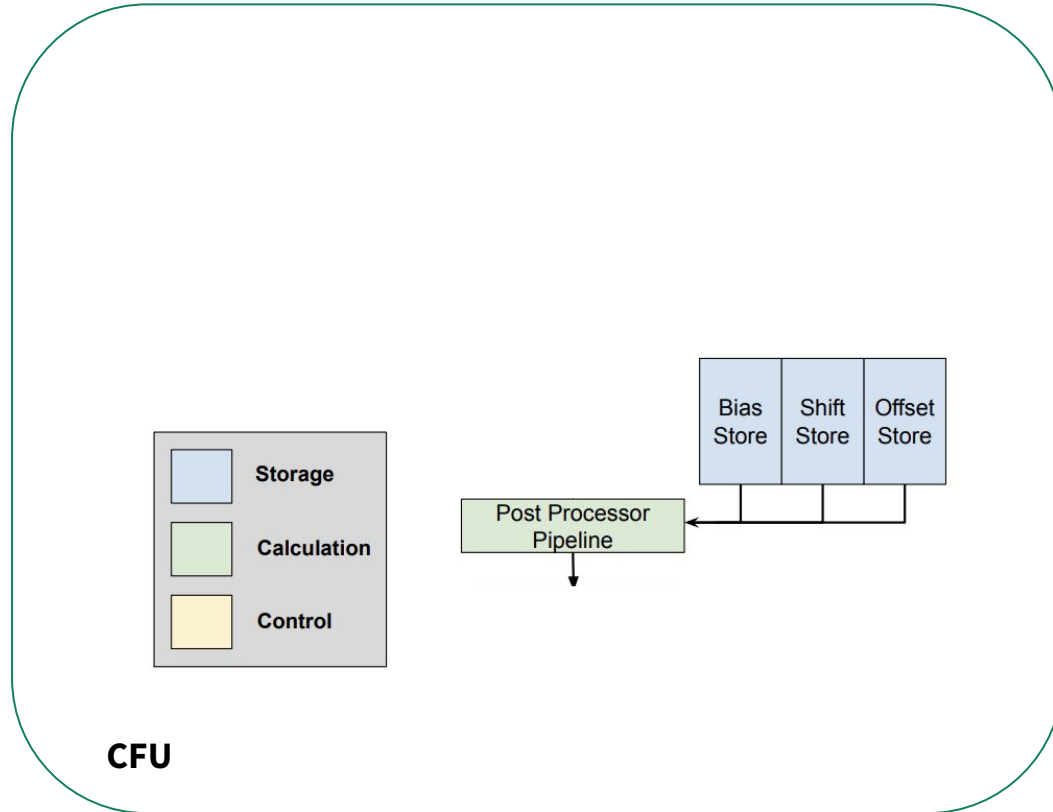


FPGA Acceleration for Image Classification

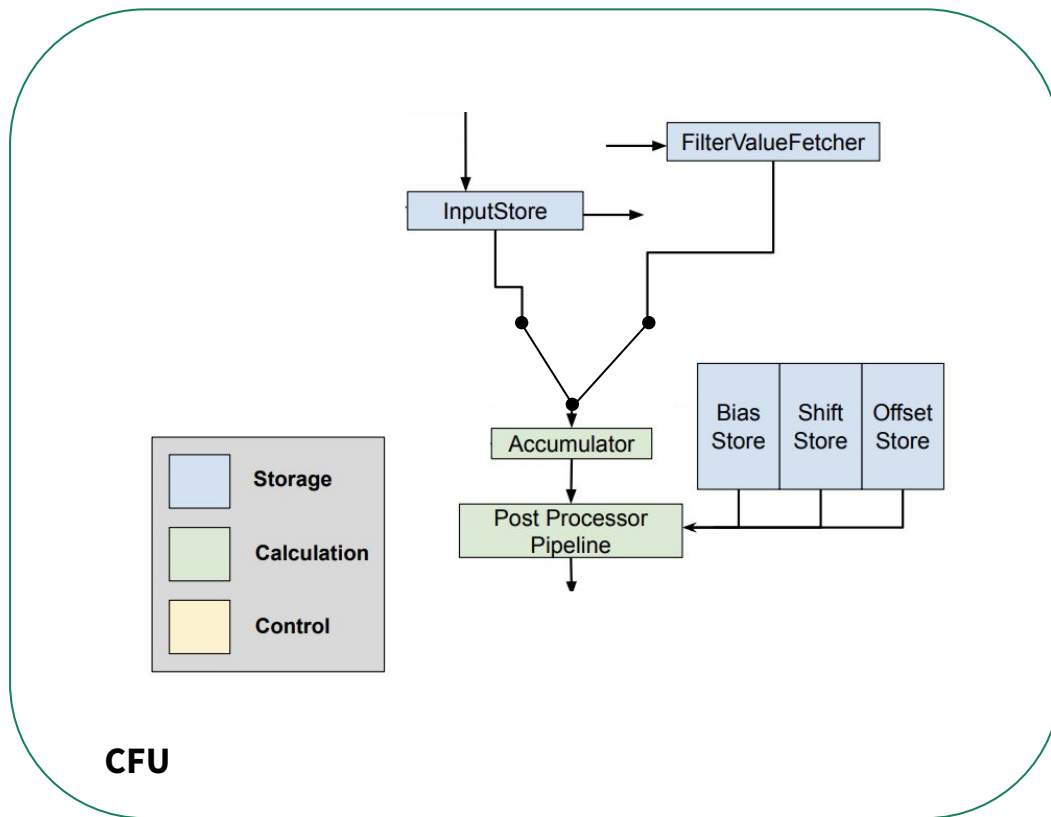
Start with Software Optimizations!



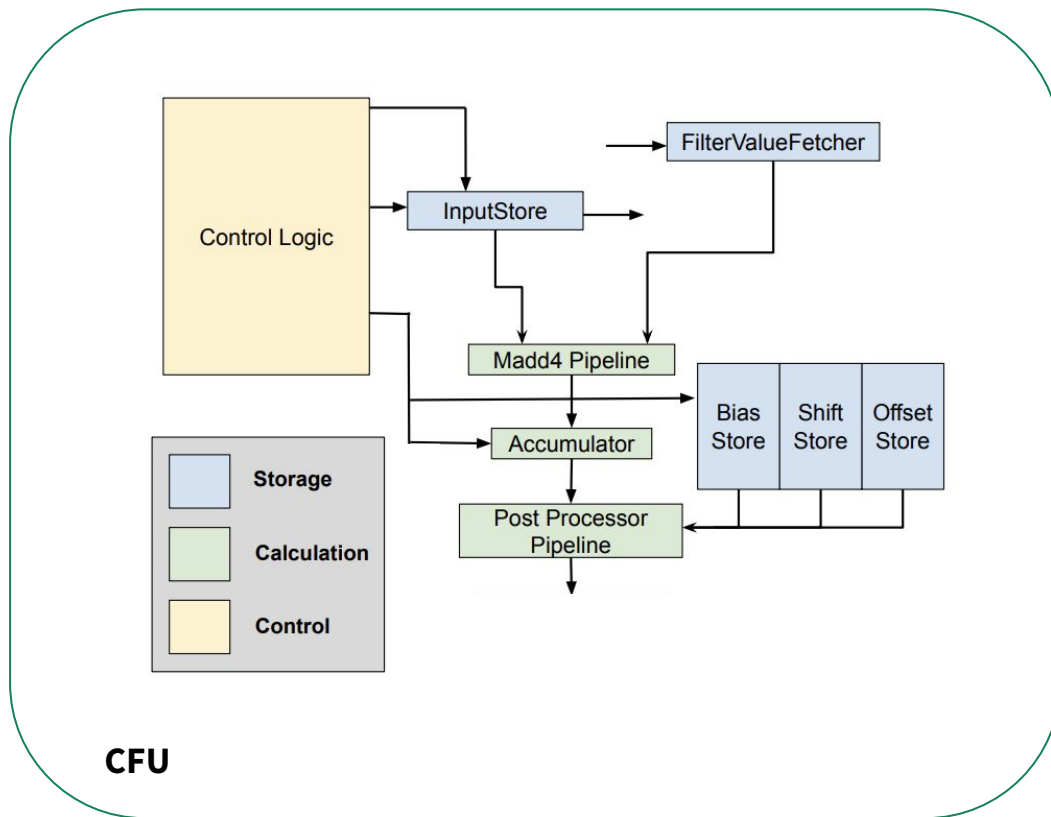
FPGA Acceleration for Image Classification



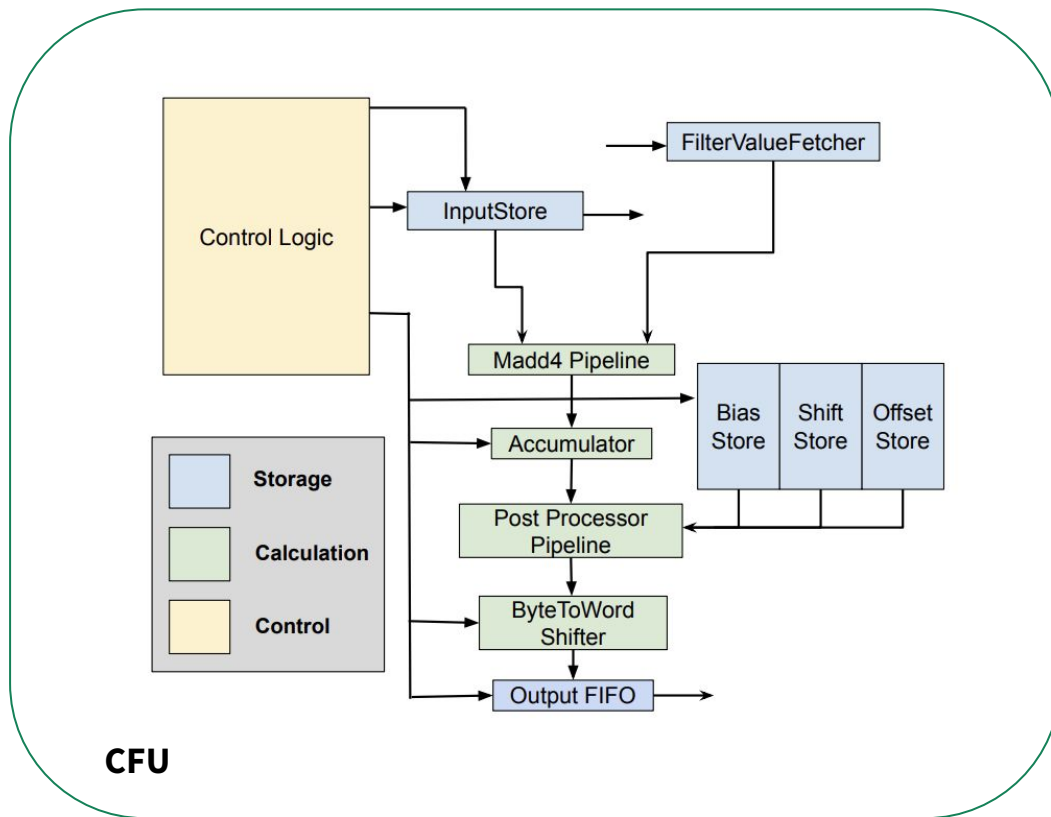
FPGA Acceleration for Image Classification



FPGA Acceleration for Image Classification

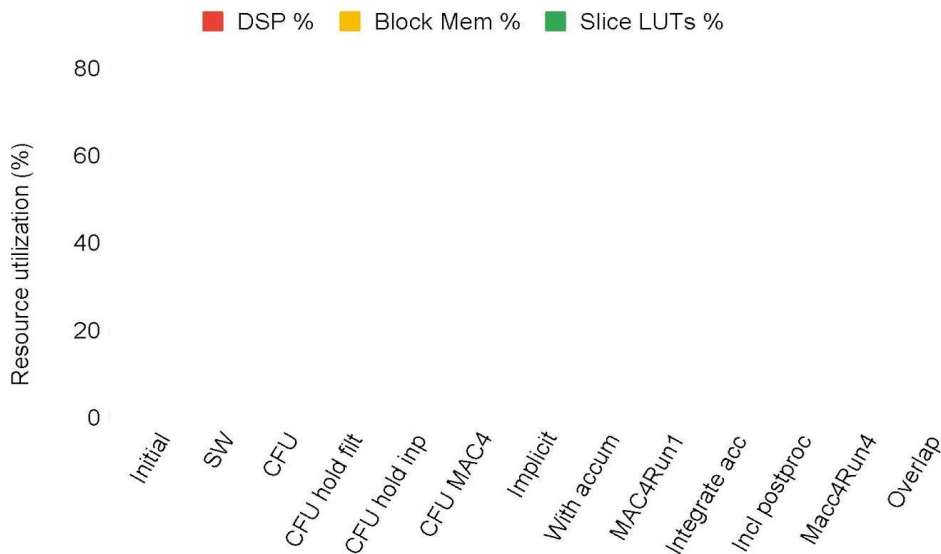


FPGA Acceleration for Image Classification



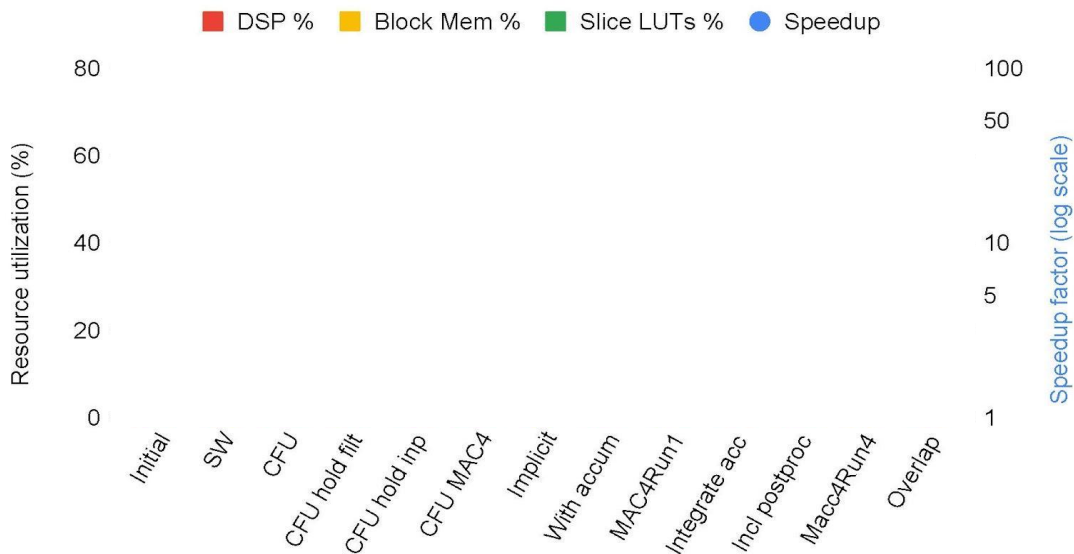
FPGA Acceleration for Image Classification

Image Classification on Arty



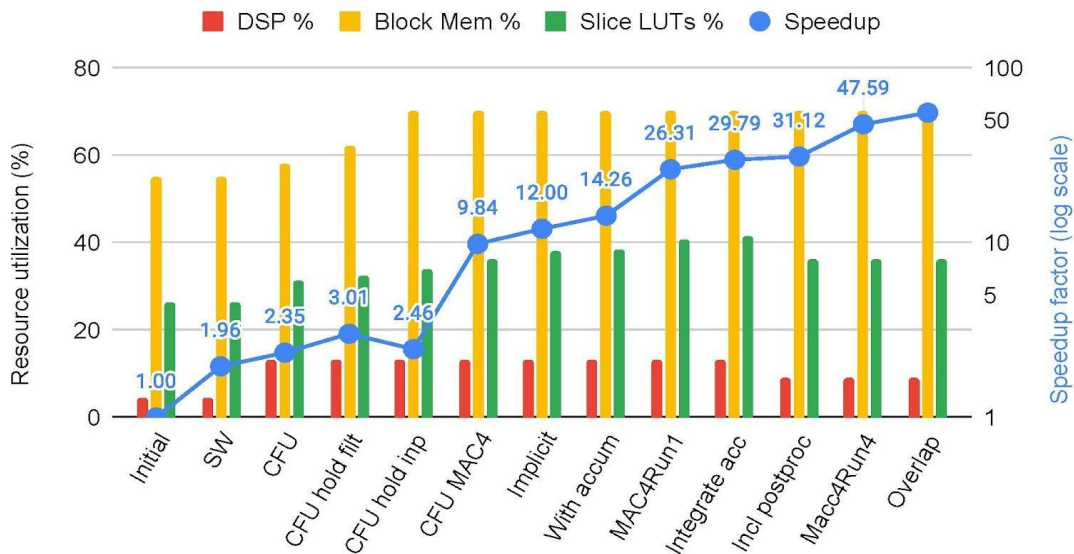
FPGA Acceleration for Image Classification

Image Classification on Arty



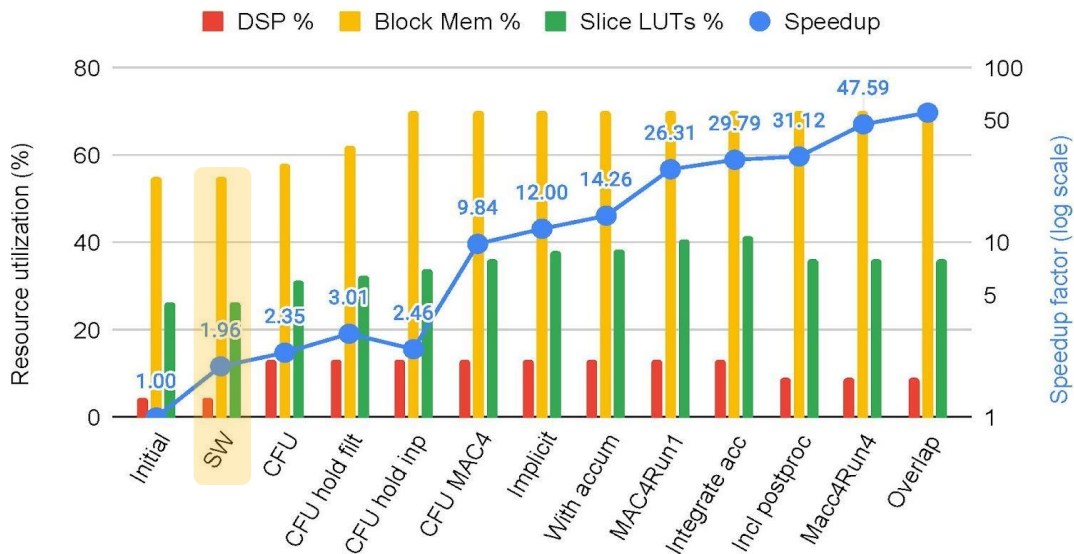
FPGA Acceleration for Image Classification

Image Classification on Arty



FPGA Acceleration for Image Classification

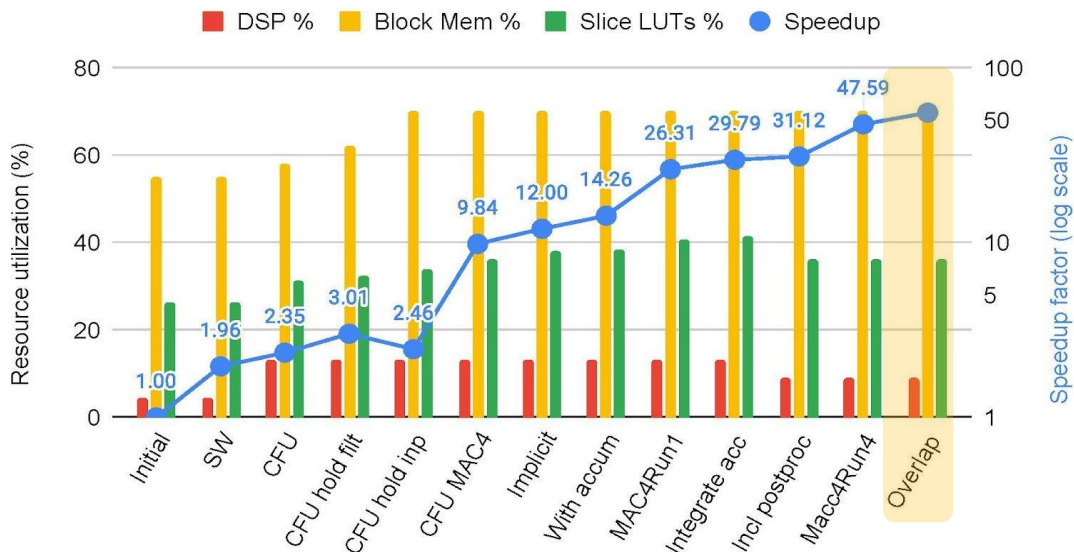
Image Classification on Arty



2x speedup from SW Optimizations

FPGA Acceleration for Image Classification

Image Classification on Arty



Total 55x speedup in 5 weeks

Human Presence Sensor

- **In Chromebook:**
 - An isolated camera+ML subsystem embedded in the display bezel
- **User features:**
 - Keep awake while present
 - Dim on leave
 - Wake on approach
 - Eavesdropper warning

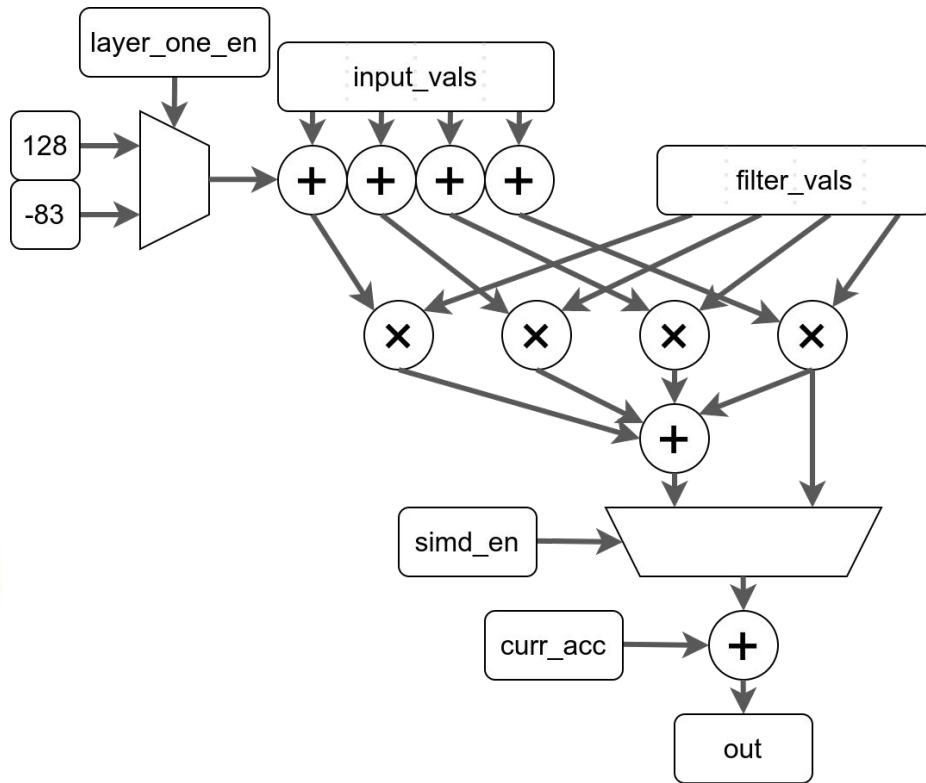


FPGA Acceleration for Keyword Spotting



FOMU FPGA

FPGA Acceleration for Keyword Spotting



75× speedup on model inference

How it started:

```
Running MLCommons Tiny V0.1 Keyword Spotting
Error_reporter OK!
Input: 490 bytes, 4 dims: 1 49 10 1

Tests for kws model
=====
0: Run with "down" input
1: Run with "go" input
2: Run with "left" input
g: Run golden tests (check for expected outputs)
x: eXit to previous menu
kws> █
```

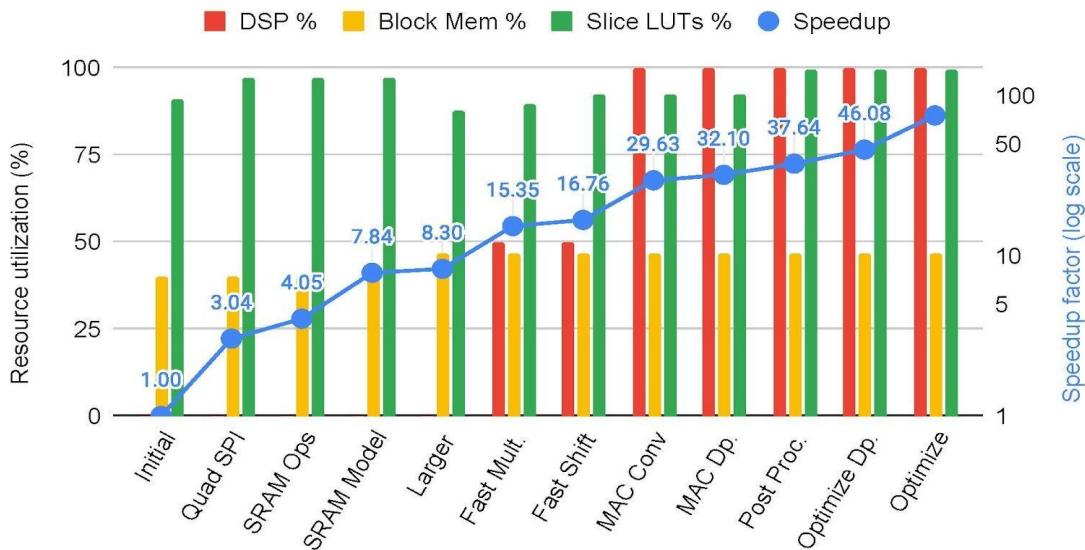
How it's going:

```
Running MLCommons Tiny V0.1 Keyword Spotting
Error_reporter OK!
Input: 490 bytes, 4 dims: 1 49 10 1

Tests for kws model
=====
0: Run with "down" input
1: Run with "go" input
2: Run with "left" input
g: Run golden tests (check for expected outputs)
x: eXit to previous menu
kws> █
```

FPGA Acceleration for Keyword Spotting

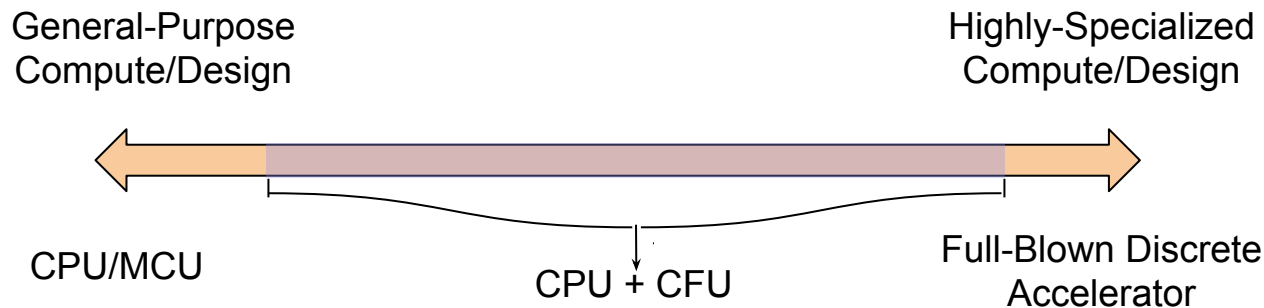
Keyword Spotting on FOMU



75x speedup in under 4 weeks

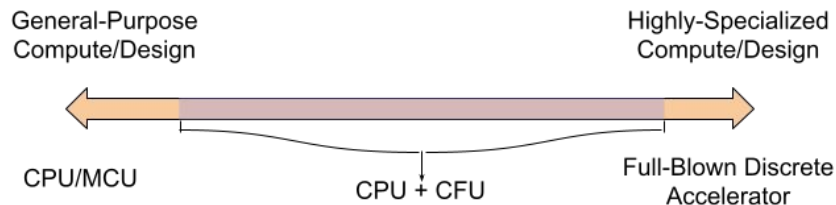
Design Space Exploration

(CFU) Accelerator vs (Soft) CPU

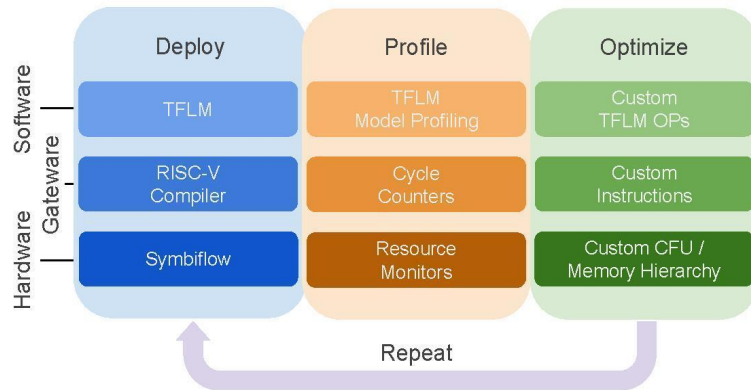


(Manual) Design Space Exploration

(CFU) Accelerator vs (Soft) CPU

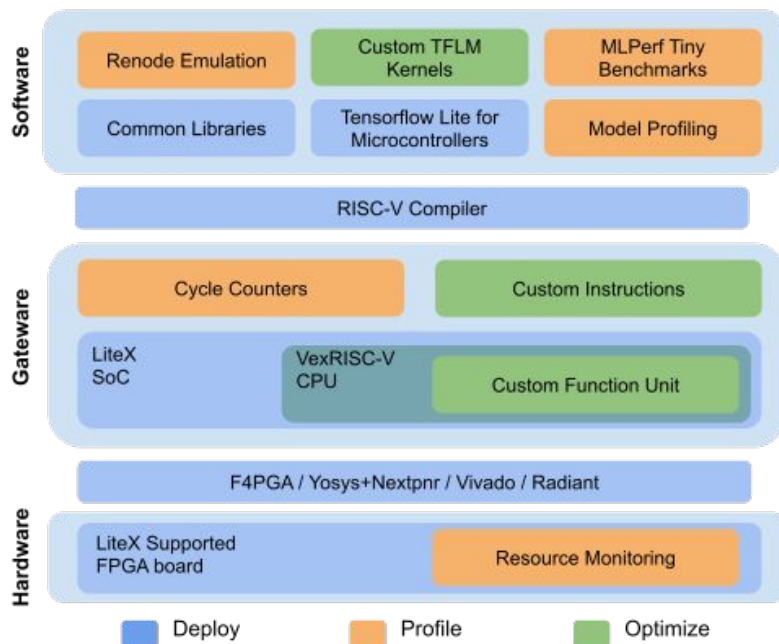


+



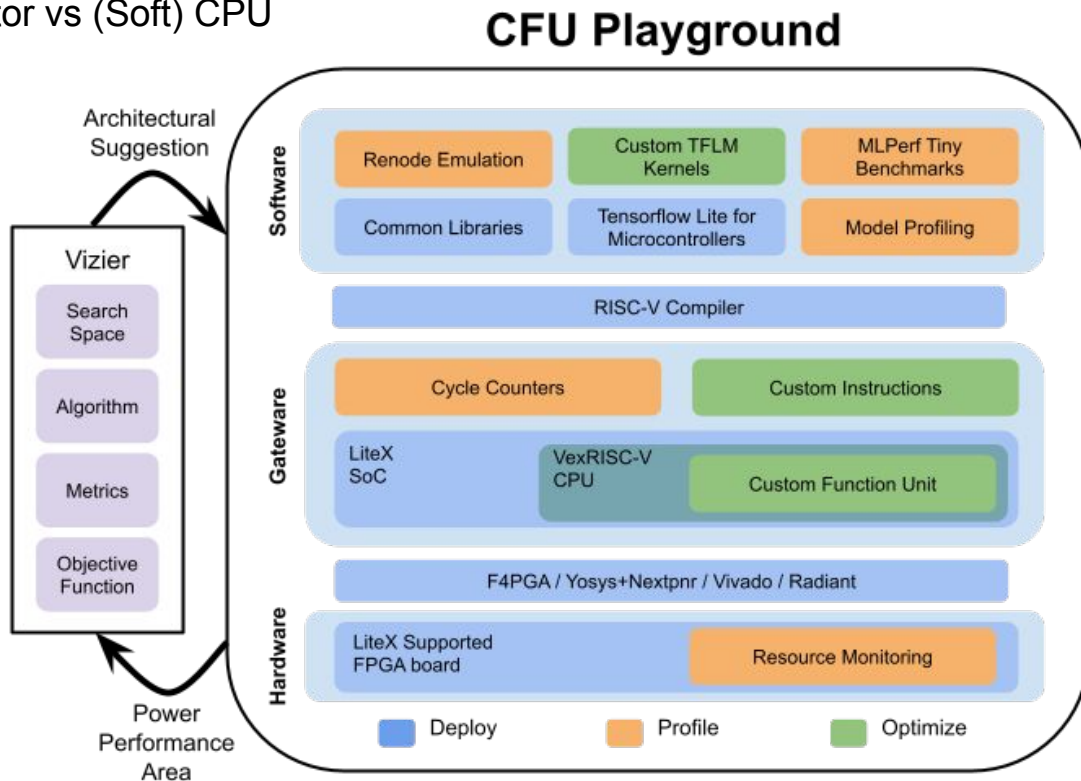
(Automated) Design Space Exploration

(CFU) Accelerator vs (Soft) CPU



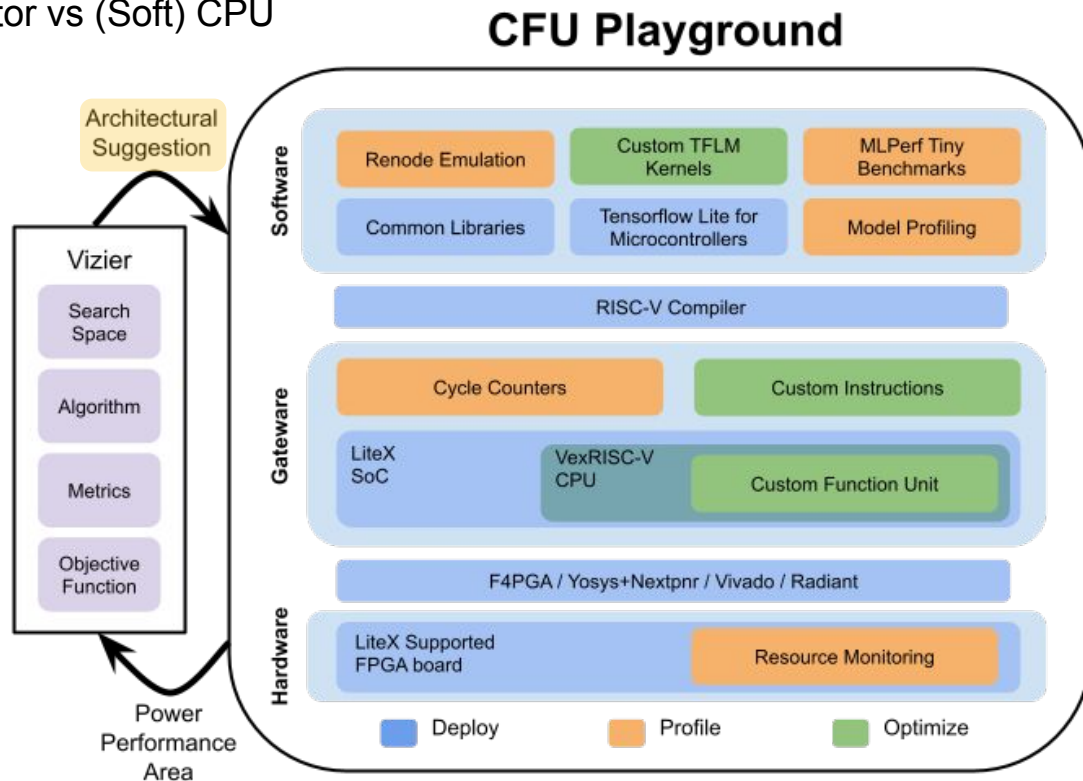
(Automated) Design Space Exploration

(CFU) Accelerator vs (Soft) CPU



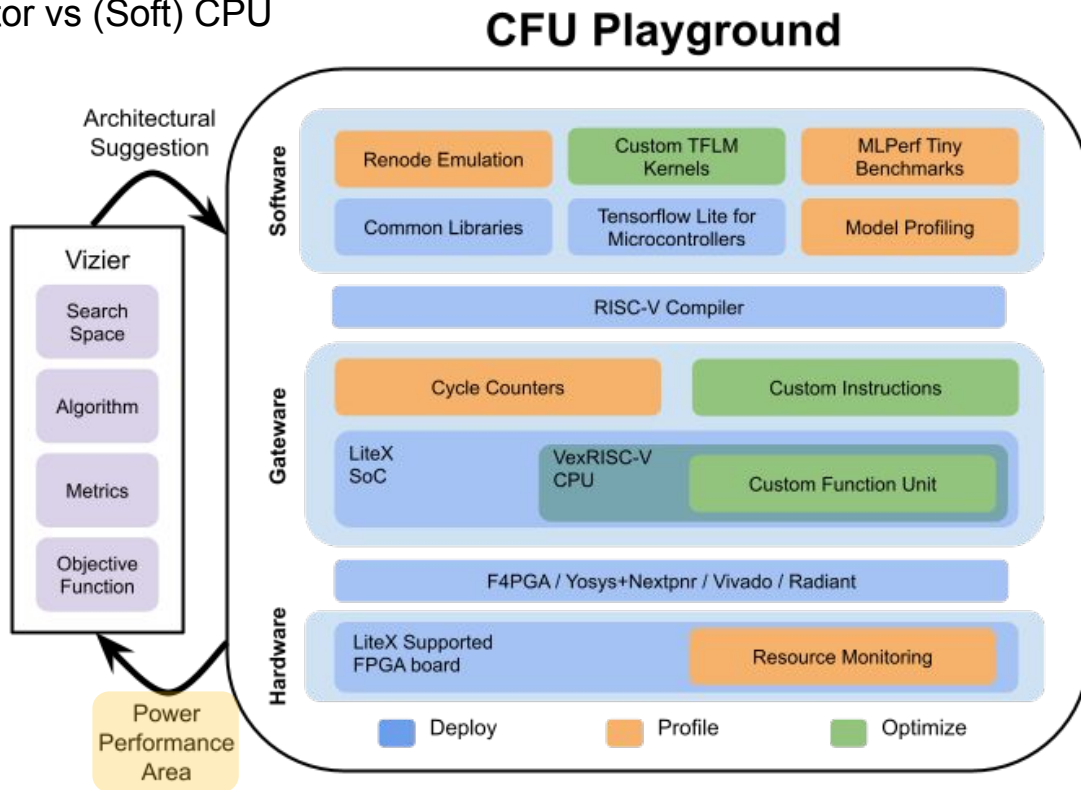
(Automated) Design Space Exploration

(CFU) Accelerator vs (Soft) CPU

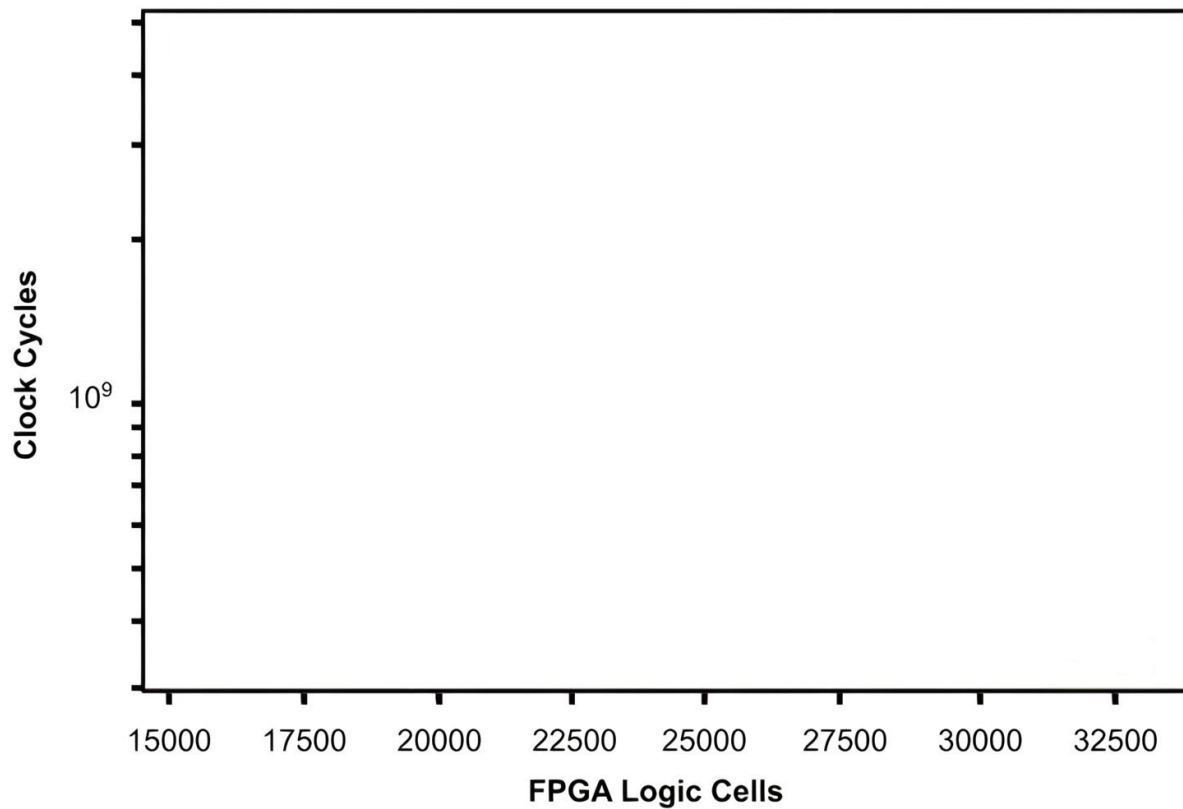


(Automated) Design Space Exploration

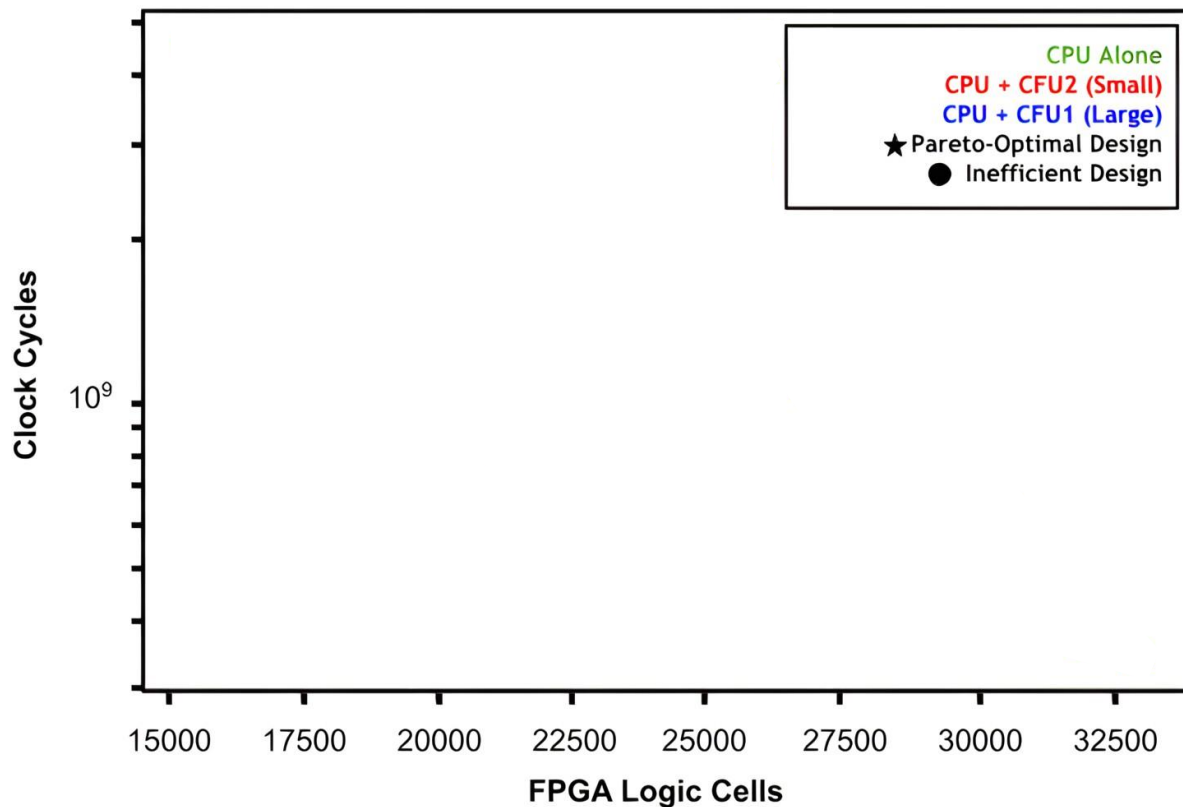
(CFU) Accelerator vs (Soft) CPU



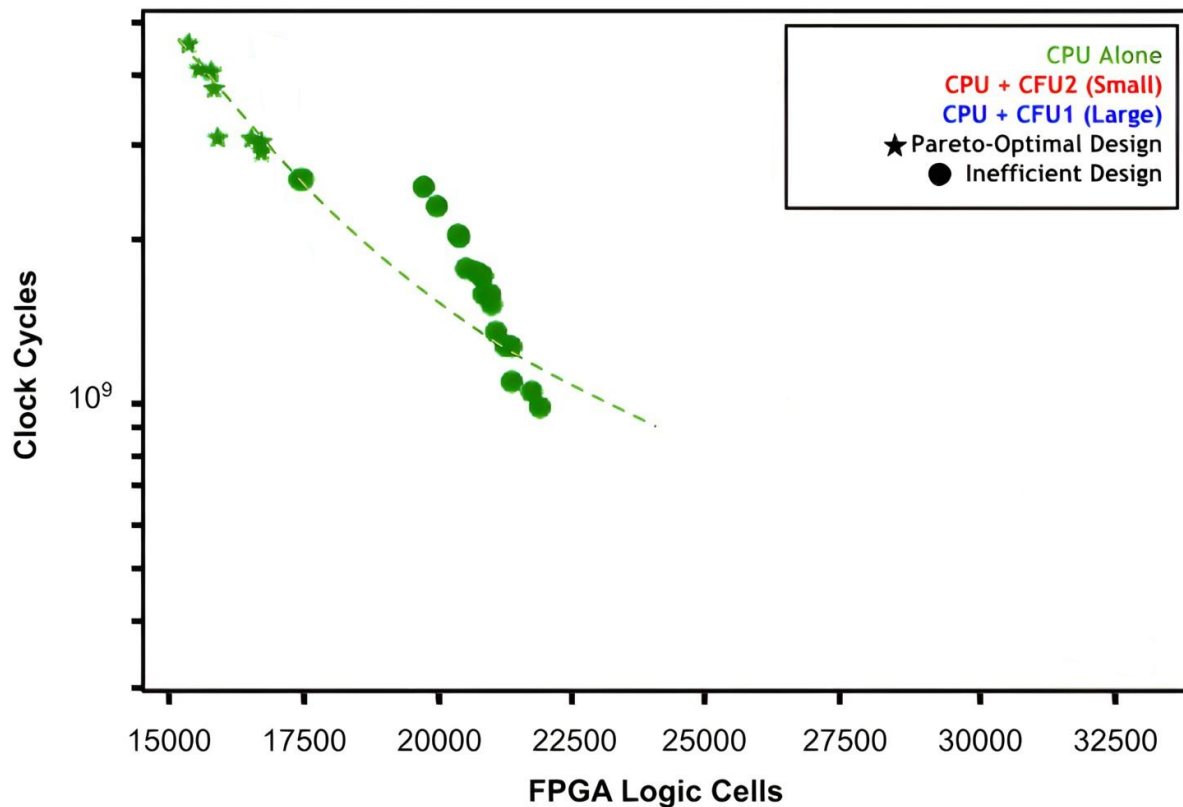
Design Space Exploration: CFU vs CPU



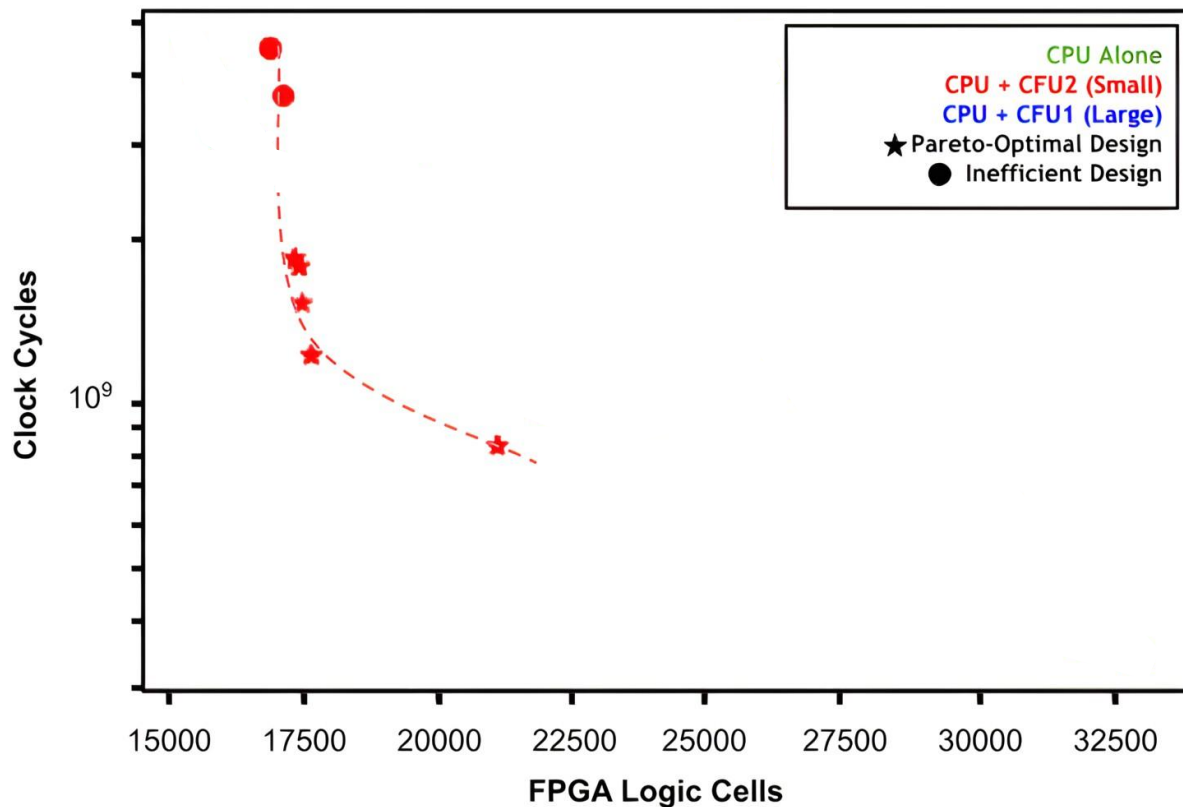
Design Space Exploration: CFU vs CPU



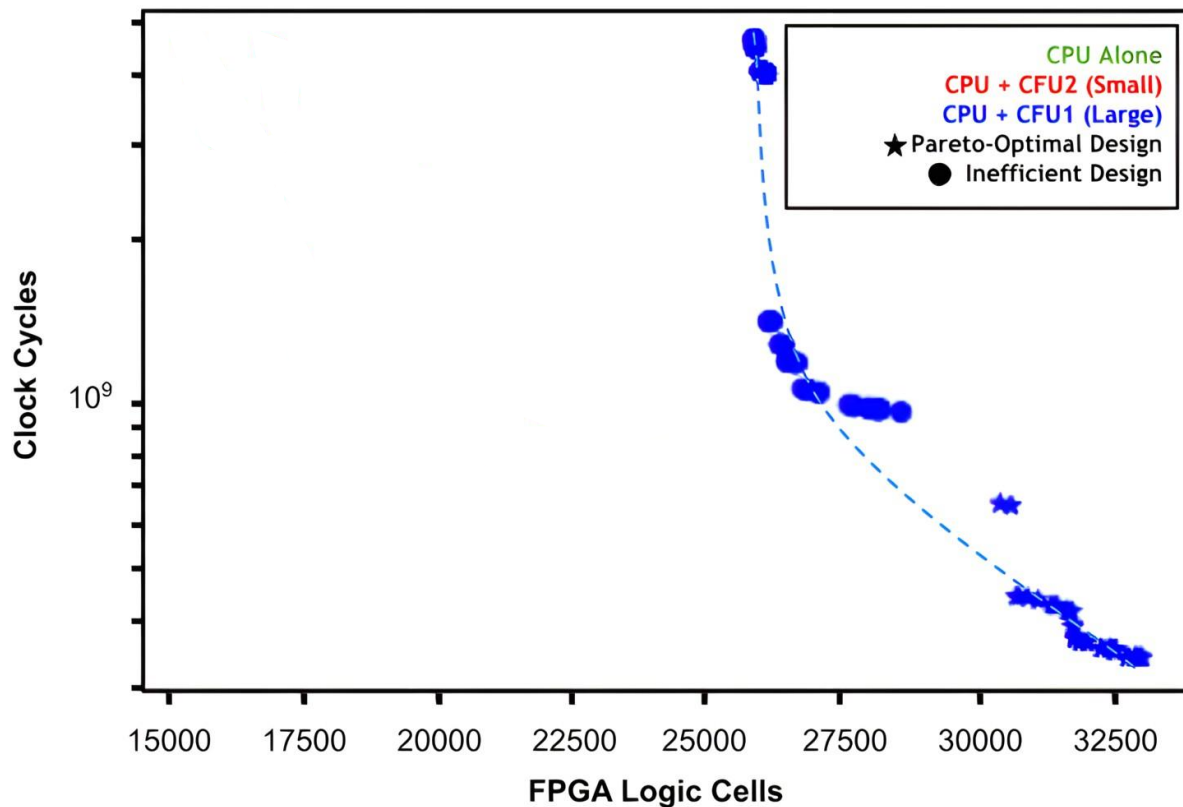
Design Space Exploration: CFU vs CPU



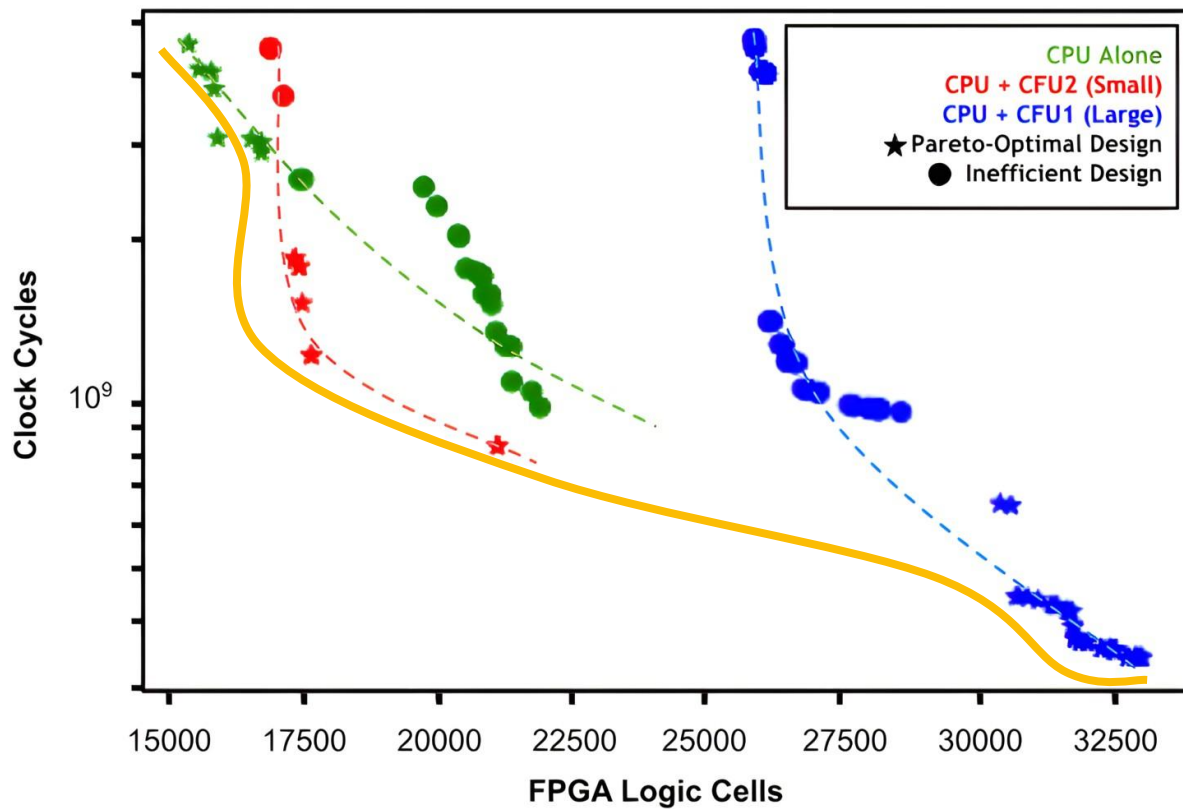
Design Space Exploration: CFU vs CPU



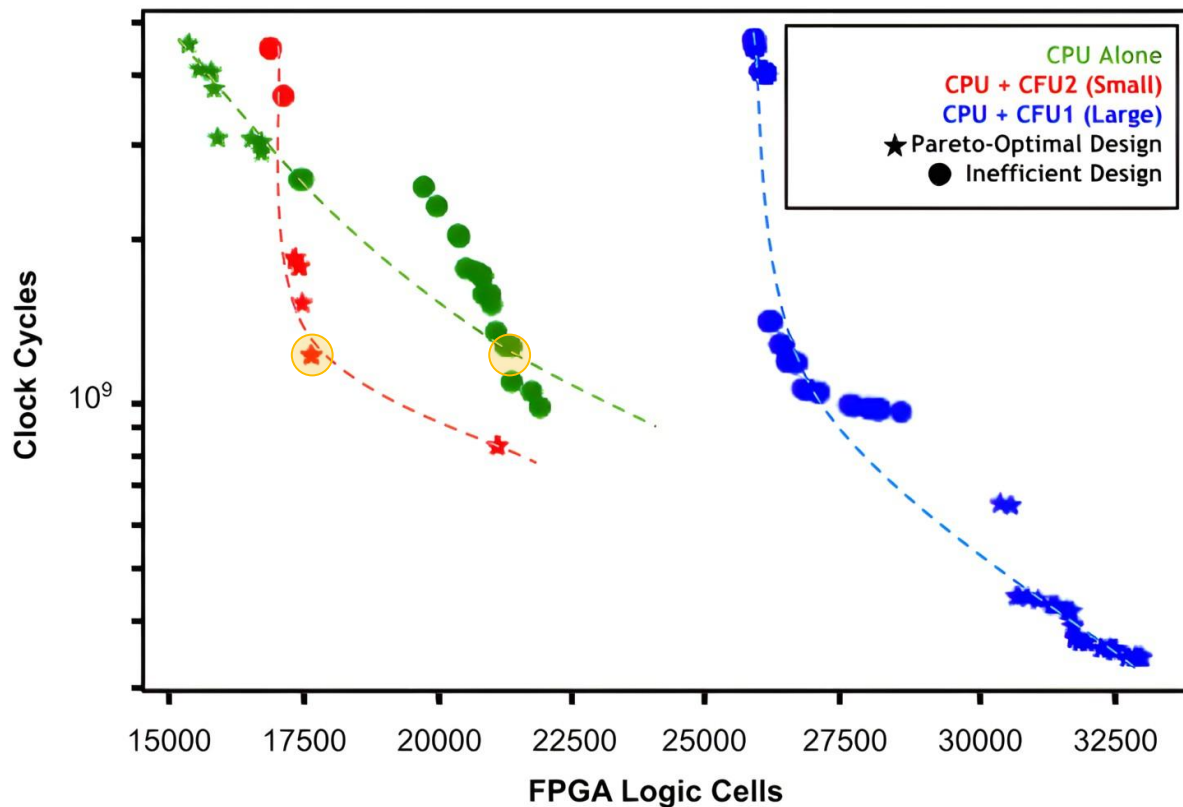
Design Space Exploration: CFU vs CPU



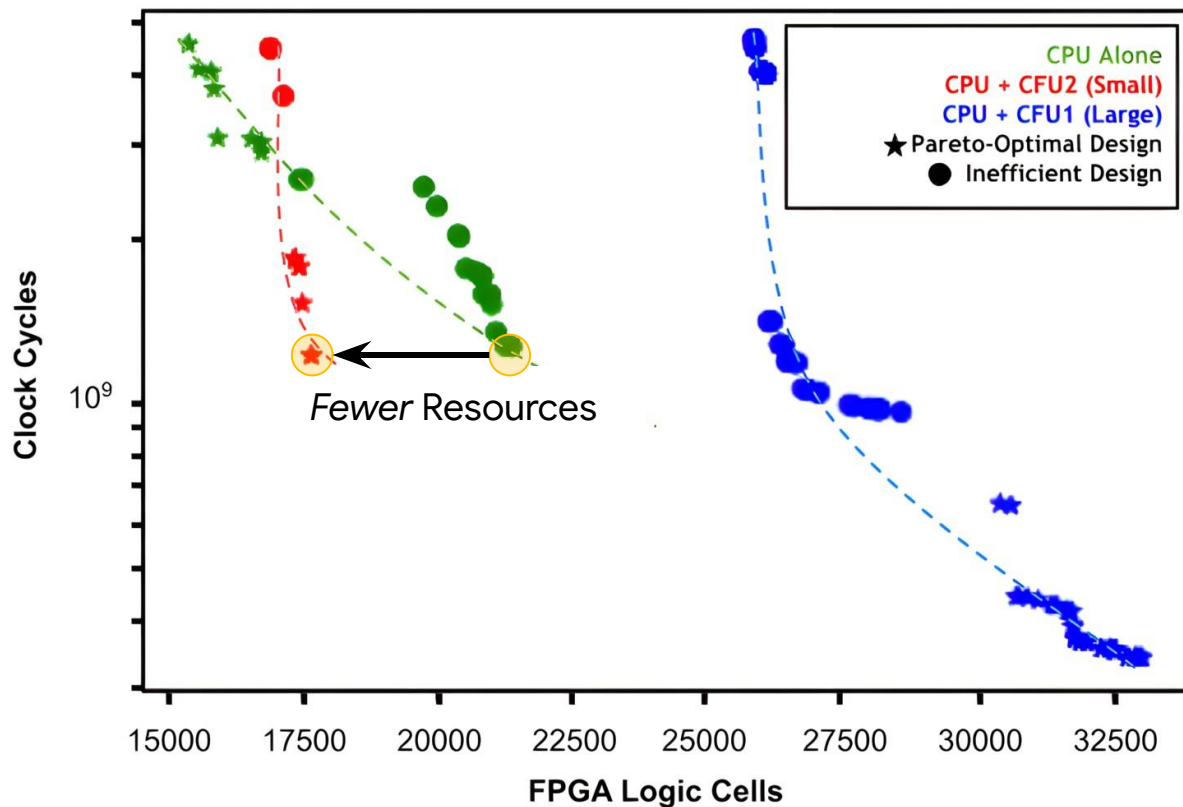
Design Space Exploration: CFU vs CPU



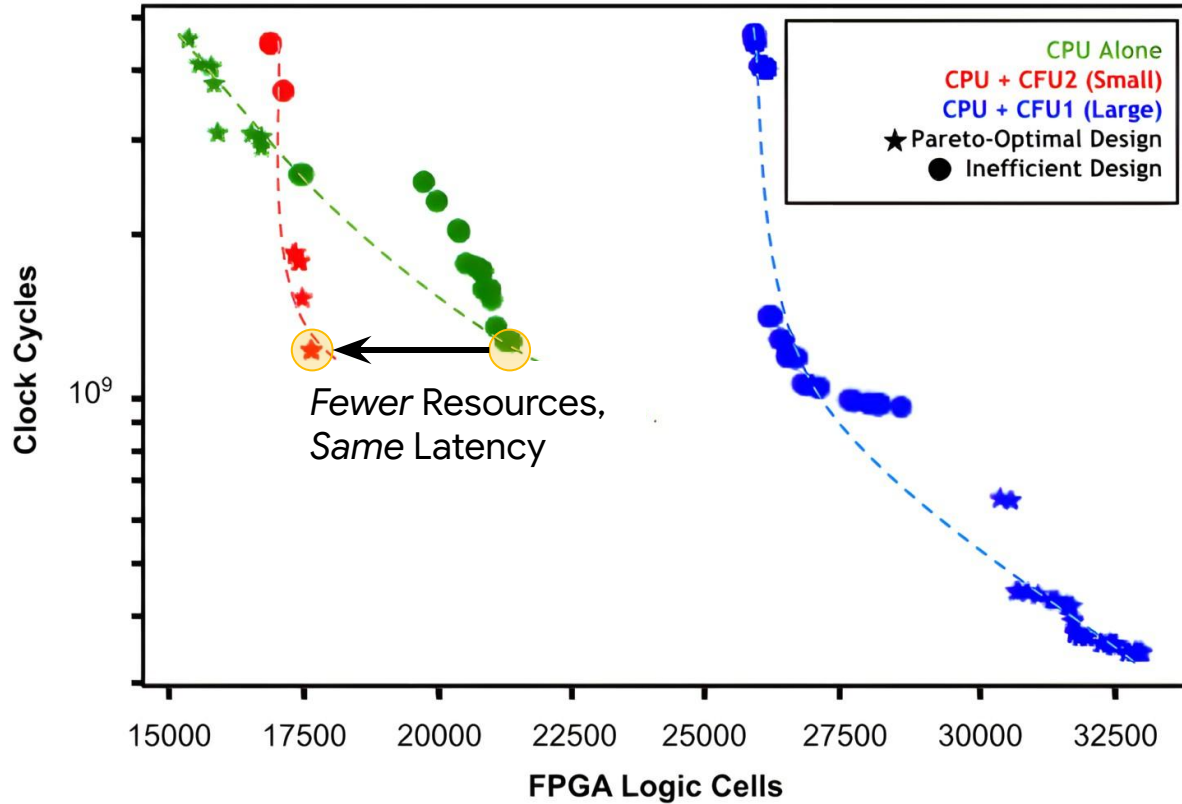
Design Space Exploration: CFU vs CPU



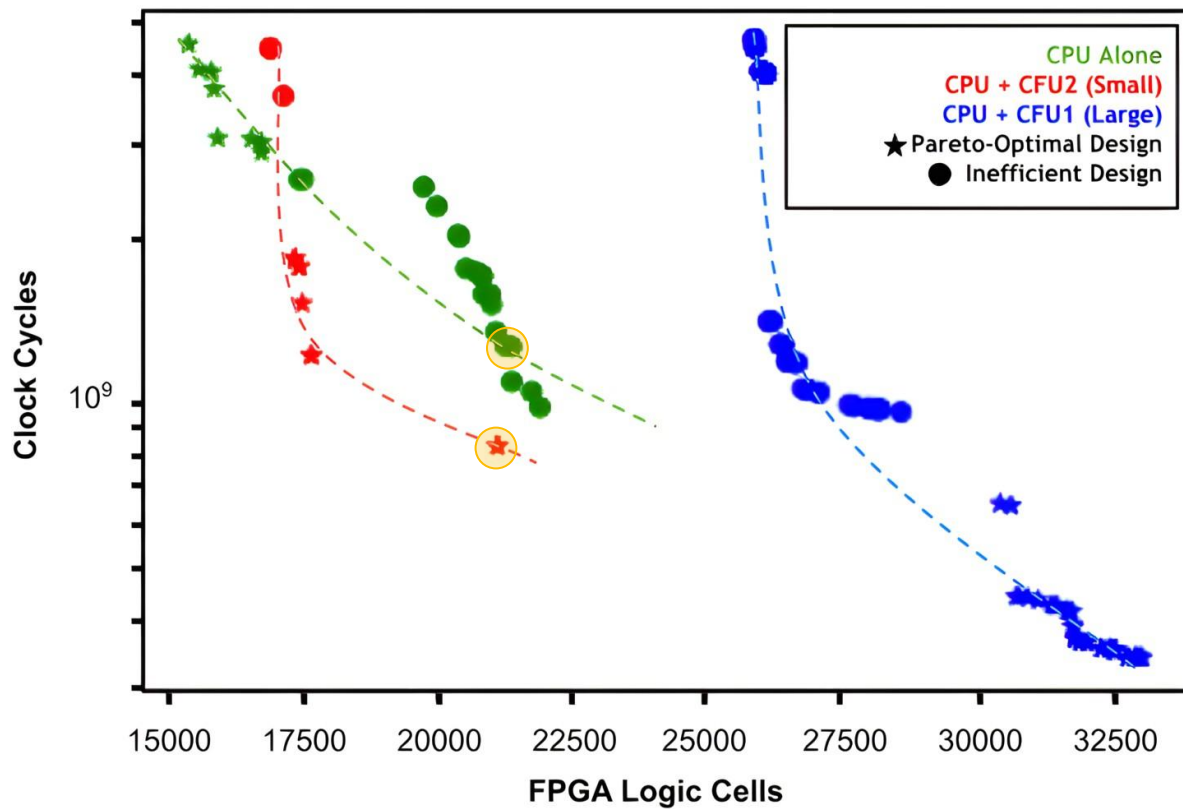
Design Space Exploration: CFU vs CPU



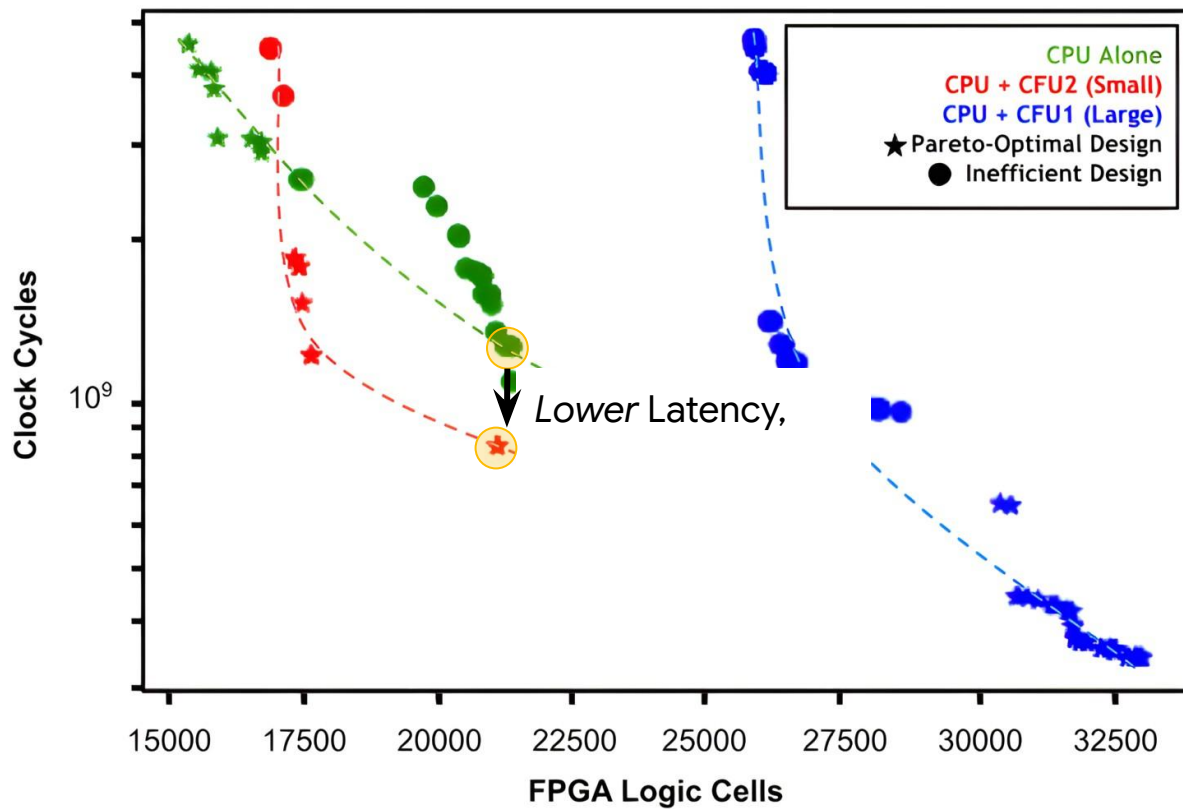
Design Space Exploration: CFU vs CPU



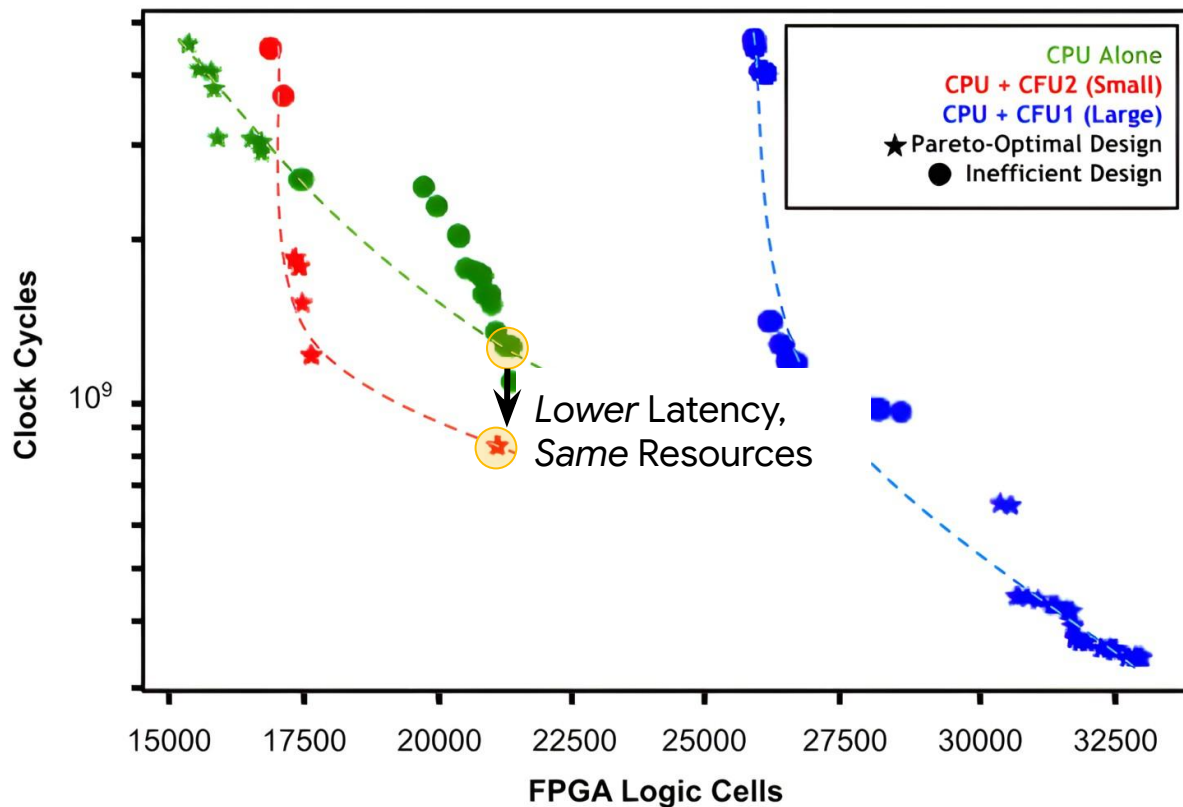
Design Space Exploration: CFU vs CPU



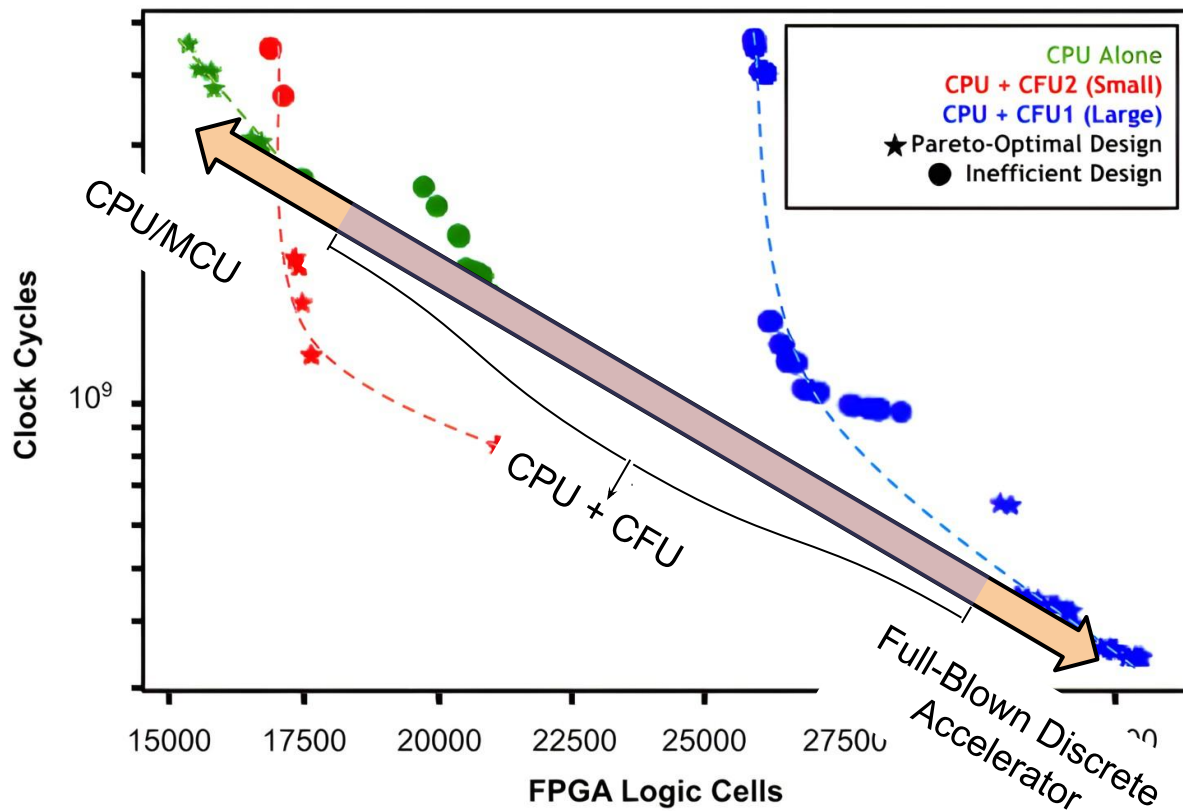
Design Space Exploration: CFU vs CPU



Design Space Exploration: CFU vs CPU



Design Space Exploration: CFU vs CPU



Key Takeaways

1. Full-stack framework that integrates open-source tools to facilitate community-driven research.
2. Agile methodology to iteratively design and evaluate tightly-coupled, bespoke TinyML accelerators.
3. Unique model-specific resource allocation trade-offs between CFU, CPU, and memory.
4. Automated design space exploration of the CPU paired with a CFU using Vizier.

CFU Playground: Full-Stack Open-Source Framework for Tiny Machine Learning (TinyML) Acceleration on FPGAs

Shvetank Prakash* Tim Callahan† Joseph Bushagour‡ Colby Banbury*
Alan V. Green† Pete Warden† Tim Ansell† Vijay Janapa Reddi*

*Harvard University †Google ‡Purdue University †Stanford University

Abstract—Need for the efficient processing of neural networks has given rise to the development of hardware accelerators. The increased adoption of specialized hardware has highlighted the need for more agile design flows for hardware-software co-design and domain-specific optimizations. In this paper, we present CFU Playground—a full-stack open-source framework that enables rapid and iterative design and evaluation of machine learning (ML) accelerators for embedded ML systems. Our tool provides a completely open-source end-to-end flow for hardware-software co-design on FPGAs and future systems research. This full-stack framework gives the users access to explore experimental and bespoke architectures that are customized and co-optimized for embedded ML. Our rapid, deploy-profile-optimization feedback loop lets ML hardware and software developers achieve significant returns out of a relatively small investment in customization. Using CFU Playground's design and evaluation loop, we show substantial speedups between 55× and 75×. The soft CPU coupled with the accelerator opens up a new, rich design space between the two components that we explore in an automated fashion using Vizier, an open-source black-box optimization service.

1. INTRODUCTION

Tiny machine learning (TinyML) is a fast-growing field at the intersection of ML algorithms and low-cost embedded systems. It enables on-device sensor data analytics (vision, audio, IMU, etc.) at ultra-low-power consumption. Processing data close to the sensor allows for an expansive new variety of always-on ML use-cases that preserve bandwidth, latency, and energy while improving responsiveness and maintaining privacy [1]. Given the need for energy efficiency when running ML on these embedded platforms, custom processor support and hardware accelerators for such systems could present the needed solutions. However, the field of ML is still in its infancy and fast-changing. Thus, it is desirable to avoid a massive non-recurring engineering (NRE) cost upfront, especially for low-cost embedded ML systems. Building ASICs is both costly and time-consuming. Moreover, since embedded systems are often task-specific, there is an opportunity to avoid general-purpose ML accelerators and instead explore task and model-specific ML acceleration methods. This setting presents the need for an agile design space exploration tool that allows us to adapt to the changing landscape of ML and hardware.

To enable holistic hardware-software co-design and evaluation of domain-specific performance optimizations easily,

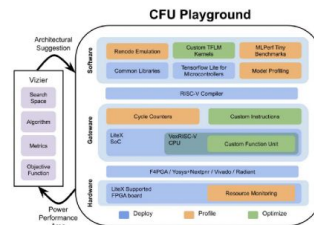


Fig. 1: CFU Playground allows users to design and evaluate model-specific ML enhancements to a “soft” CPU core. The Playground is wrapped around Vizier, an open-source black-box optimization service, to enable ML-driven design space exploration.

we present CFU Playground.¹ It is a full-stack open-source framework for iteratively (deploy→profile→optimize) exploring the design space of lightweight accelerators in an agile manner (Figure 1). The framework is unique in that it couples together various open-source software (TensorFlow Lite Micro/TFML, GCC), open-source RTL generation IP and toolkits (LiteX, VexRiscv, Migen, Amaranth), and open-source FPGA tools for synthesis, place, and route (yosys, nextpnr, F4PGA/SymbiFlow, etc.). By using open source for the entire stack, we enable the end-user to *customize* and *co-optimize hardware and software*, resulting in a specialized solution unencumbered by potential licensing restrictions and not tied to a particular FPGA, board, or vendor. CFU Playground yields large returns out of a relatively small investment in customized hardware and is useful for the long tail of low-volume applications.

Yet another novelty of CFU Playground is in its *ability to design custom function units (CFUs)* for distinct ML operations. CFUs represent a novel design space that balances ac-

¹CFU Playground is available at www.github.com/google/CFU-Playground.



Paper Discussions

Why are we reading these papers?

What are the important things we learn from these papers?

What can we compare and contrast with these papers?



Guest Speaker

Hardik Sharma

Hardik is the Director of Hardware Engineering at Bigstream. His research interests are domain specific hardware architectures for accelerating machine learning. He led the development of the first open-source FPGA-based hardware acceleration stack for DNNs at Georgia Tech.



[Website](#)

[Google Scholar](#)