

# **PULP-TrainLib: Enabling On-Device Training for RISC-V Multi-Core MCUs through Performance-Driven Autotuning**

---

**Davide Nadalini<sup>1,2</sup>, Manuele Rusci<sup>2</sup>, Giuseppe Tagliavini<sup>3</sup>,  
Leonardo Ravaglia<sup>2</sup>, Luca Benini<sup>2,4</sup>, and Francesco Conti<sup>2</sup>**

<sup>1</sup>DAUIN, Politecnico di Torino, Torino, Italy, [davide.nadalini@polito.it](mailto:davide.nadalini@polito.it)

<sup>2</sup>DEI, University of Bologna, Bologna, Italy,  
{[d.nadalini](mailto:d.nadalini),[manuele.rusci](mailto:manuele.rusci),[luca.benini](mailto:luca.benini),[f.conti](mailto:f.conti)}@unibo.it

<sup>3</sup>DISI, University of Bologna, Bologna, Italy, [giuseppe.tagliavini@unibo.it](mailto:giuseppe.tagliavini@unibo.it)

<sup>4</sup>IIS, ETH Zurich, Zurich, Switzerland, [lbenini@iis.ee.ethz.ch](mailto:lbenini@iis.ee.ethz.ch)

# OVERVIEW

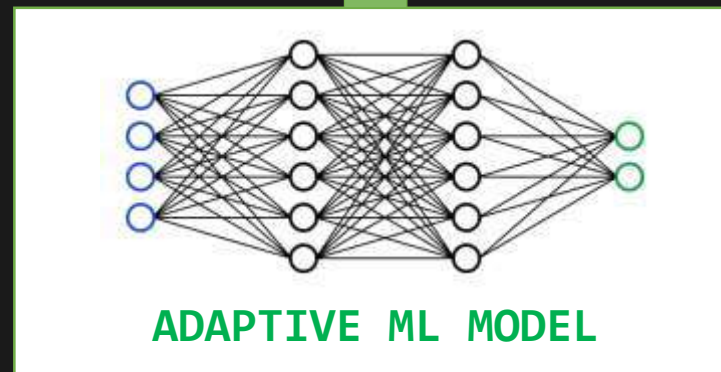
- Intro
  - Motivation (On-Device Training)
  - Contribution
- PULP-TrainLib
  - PULP-TrainLib
  - AutoTuner
- Results
  - In-Depth Analysis of AutoTuner's Results
  - Training TinyML Models
  - Comparison with the State-Of-Art
- Conclusion
- Q. & A.



# Introduction

# MOTIVATION

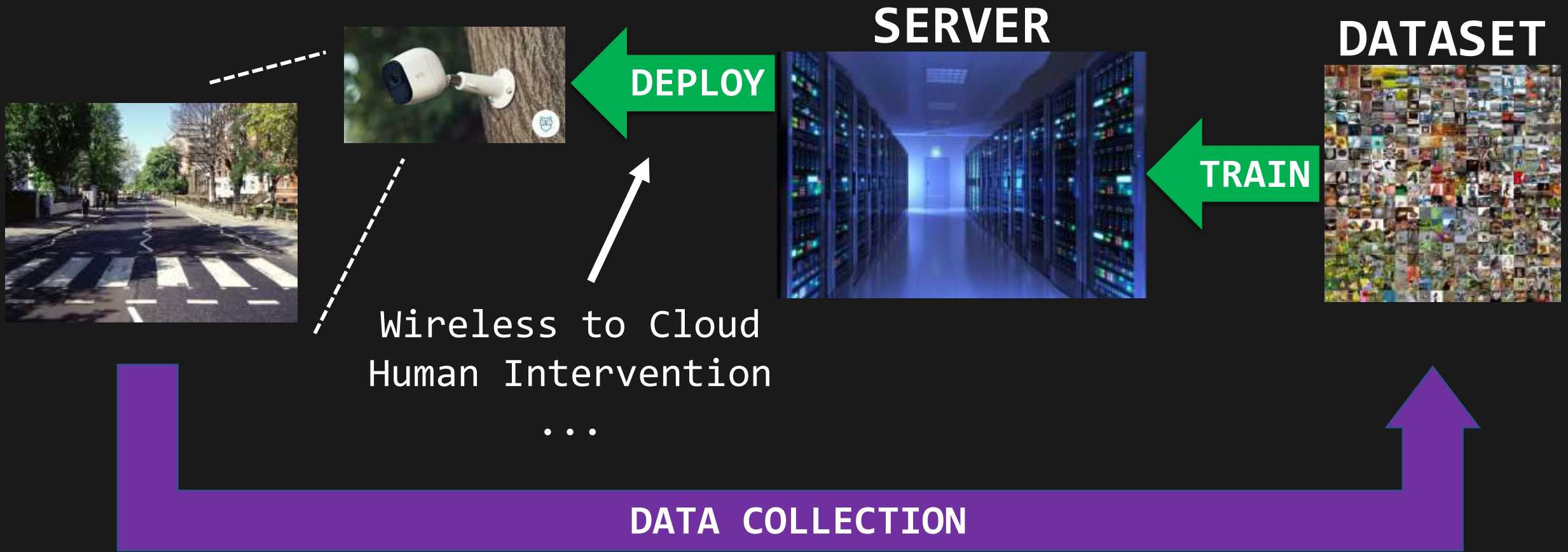
MINIVAN!



ADAPTIVE ML  
MODEL

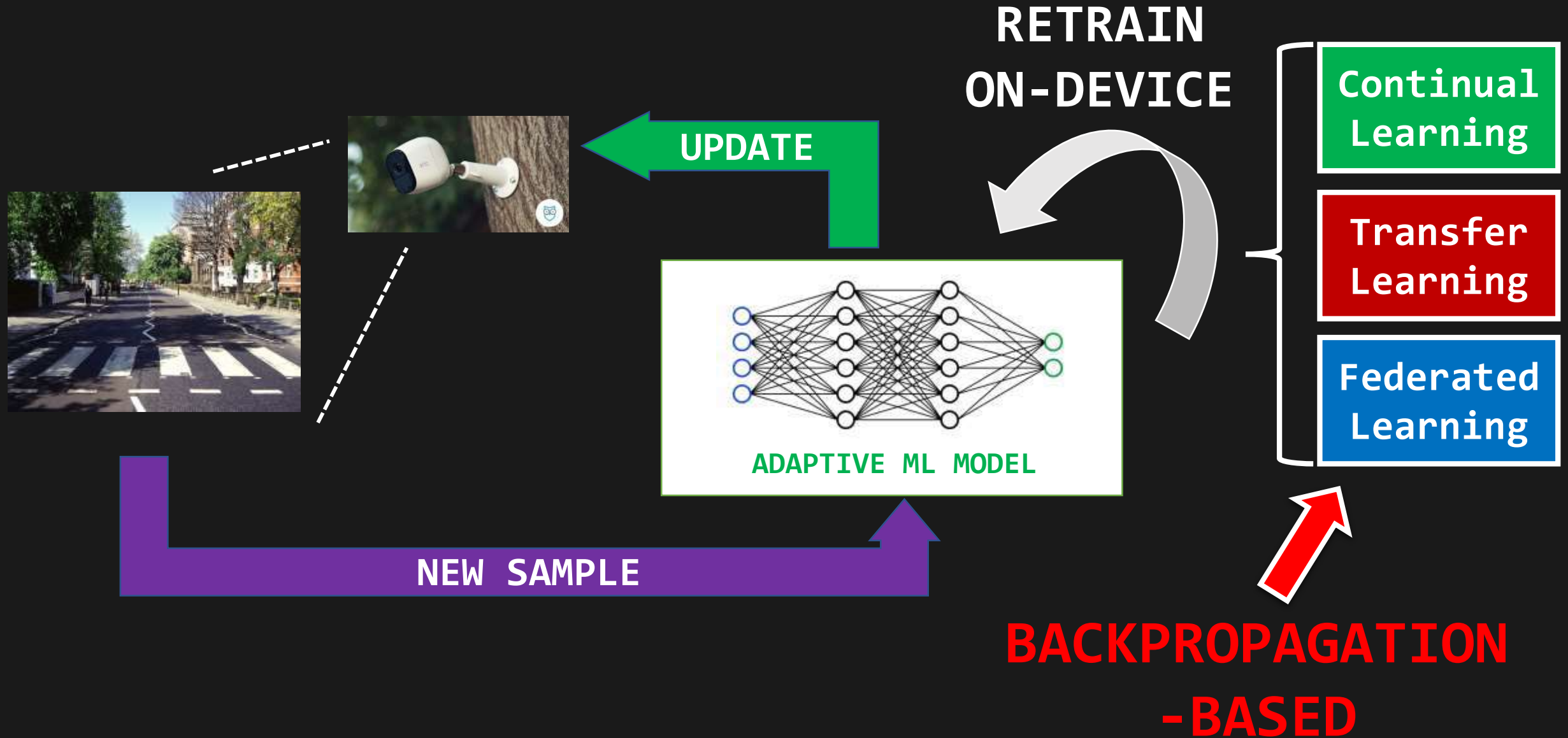
ML model's  
**parameters**  
updated based  
on  
environment-  
specific data

# STANDARD DNN TRAINING



**PRONE TO PRIVACY ISSUES, LIMITED SCALING, INCREASED POWER BUDGET**

# ON-DEVICE TRAINING



# TRAINING ON MCUs

SERVER



EDGE GPU



MCU



> 100 W

1-50 W

< 1 W

High Power  
High Performance

Low Power

Ultra-Low Power

State of the Art On-Device  
Training on MCUs:  
**AIfES for Arduino** (Fraunhofer IMS)

Suitable for IoT  
Edge Devices

<sup>1</sup>[https://github.com/Fraunhofer-IMS/AIfES\\_for\\_Arduino](https://github.com/Fraunhofer-IMS/AIfES_for_Arduino)

# CONTRIBUTION

**PULP-TrainLib**, the first Open-Source<sup>1</sup> training library for multicore RISC-V MCUs

**AutoTuner**, a HW-in-the-loop tool to optimize PULP-TrainLib

PULP-TrainLib + AutoTuner on an **8-Core** PULP Platform

<sup>1</sup><https://github.com/pulp-platform/pulp-trainlib>



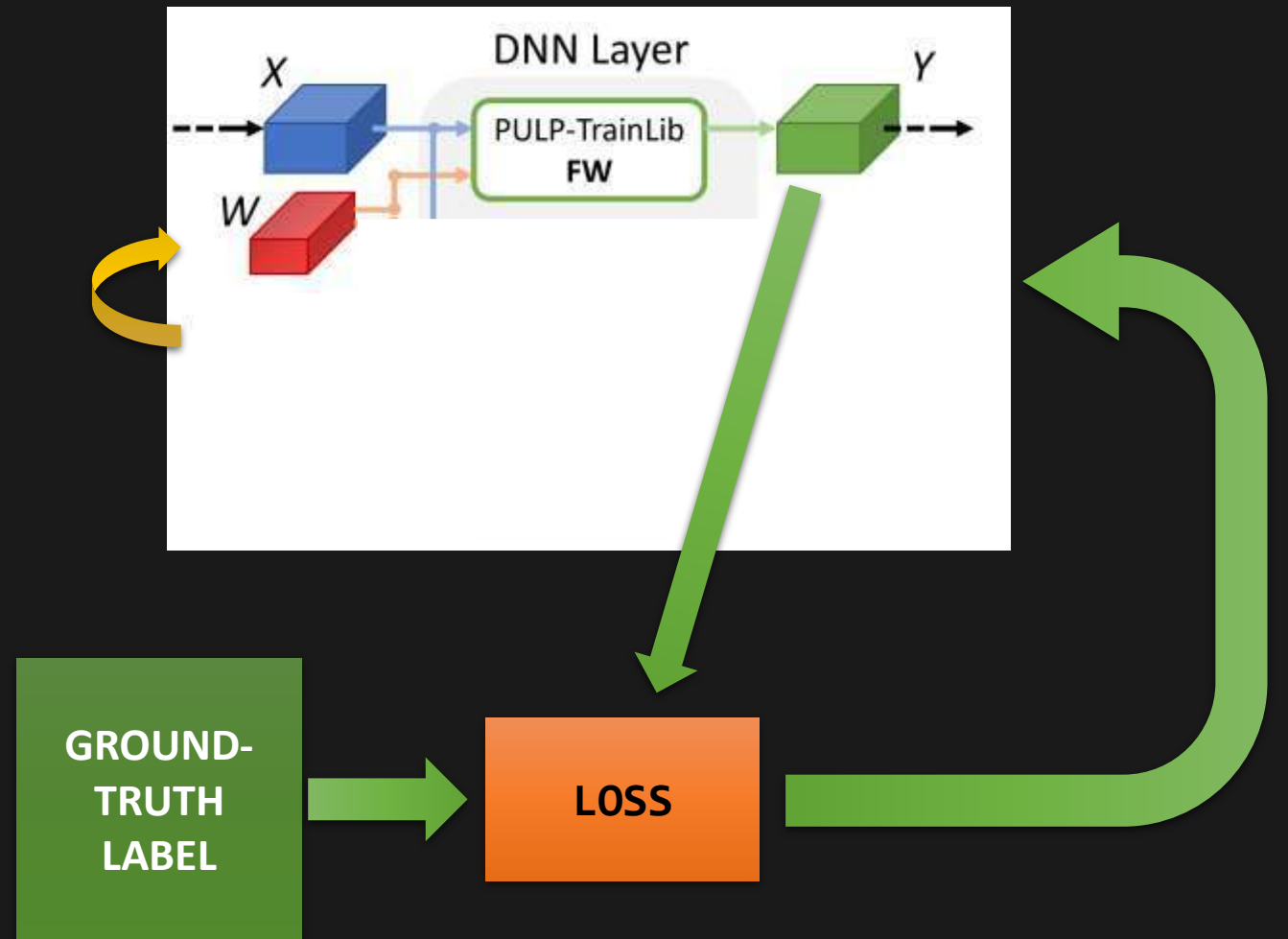


# PULP-TrainLib

# TRAINING WITH BACKPROPAGATION

GOAL: Update Deep Learning Model's  $W$

- Forward Pass
- Loss Function
- Weight Gradient
- Input Gradient
- Weight Update



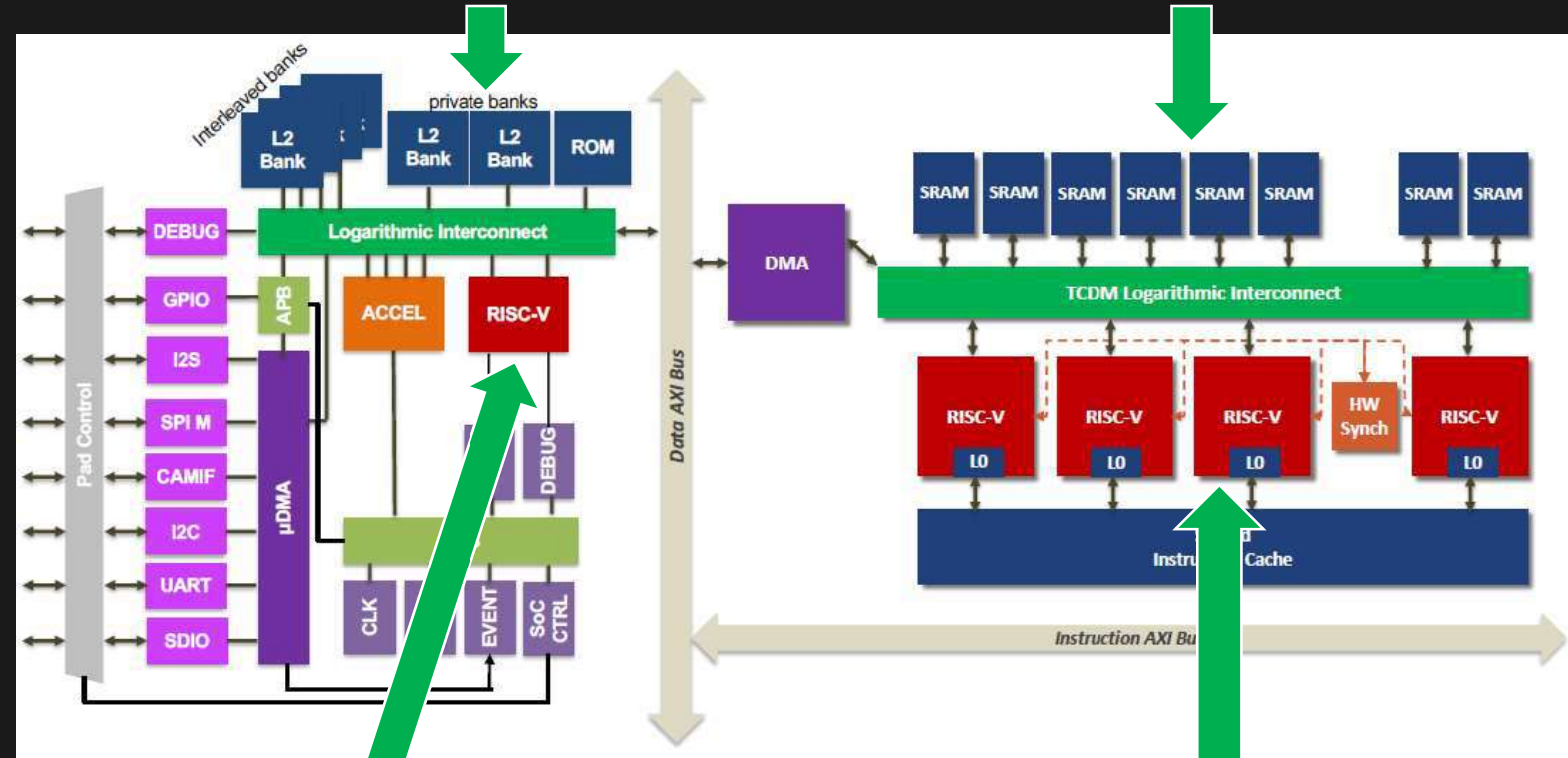
# THE PULP PLATFORM

**PULP** (Parallel Ultra-Low-Power): computational platform for energy-efficient and scalable edge computing, based on RISC-V cores.

**EMBODIMENTS:**  
**SoC (MCU), GVSoc**

512 kB - 2 MB  
SRAM (L2)

64 kB - 128 kB  
SRAM (L1)

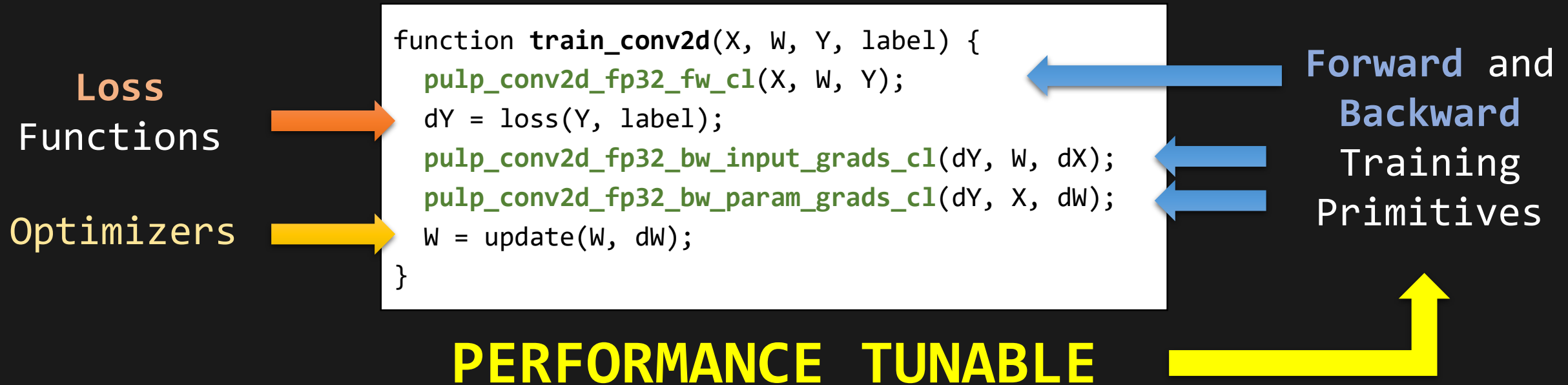


RISC-V MCU  
(Fabric Controller)

Cluster of 8 RISC-V Cores  
RV32IMFC-Xpulp

# PULP-TRAINLIB

The First SW DNN (Deep Neural Network) training library for multicore RISC-V MCUs

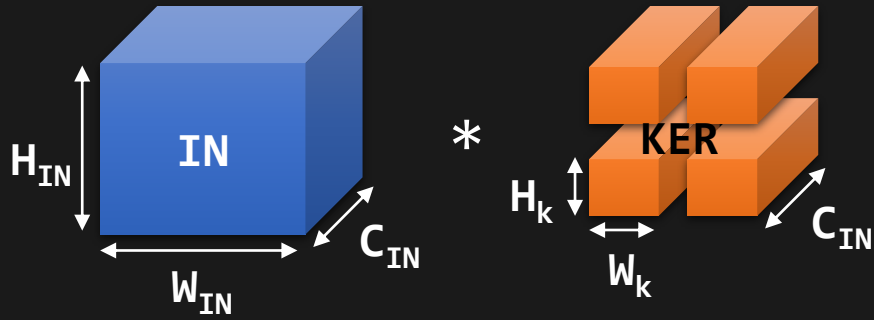


Floating Point FP32 SUPPORTED LAYERS: Fully-Connected, DepthWise Separable Convolution, 2D Convolution

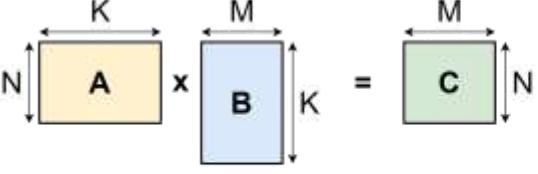
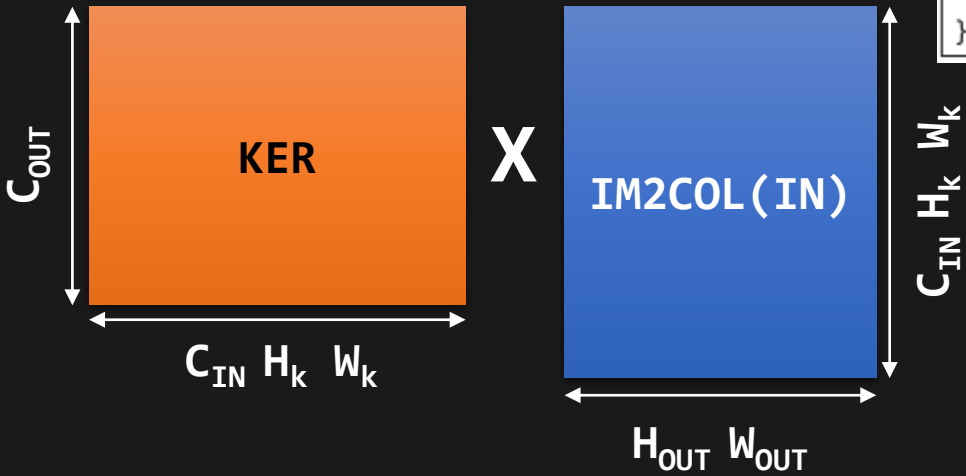
# PULP-TRAINLIB'S PRIMITIVES

UNROLLING

PARALLELIZATION



CONV2D FORWARD



```
function mm (A, B, C, sizes) {
  for i in N:
    for j in M:
      for k in K:
        // C[i, j] += A[i, k] * B[k, j]
        ld  rsA, Aik
        ld  rsB, Bkj
        fmac rdC, rsA, rsB
      st  rdC, addr
}
```

(a)

```
function mm_unroll_2x1 (A, B, C, sizes) {
  for i in N/2:
    for j in M:
      for k in K:
        // C[i, j] += A[i, k] * B[k, j]
        // C[i+1, j] += A[i+1, k] * B[k, j]
        ld  rsA1, Ai1k
        ld  rsA2, Ai2k
        ld  rsB, Bkj
        fmac rdC1, rsA1, rsB
        fmac rdC2, rsA2, rsB
      st  rdC1
      st  rdC2
}
```

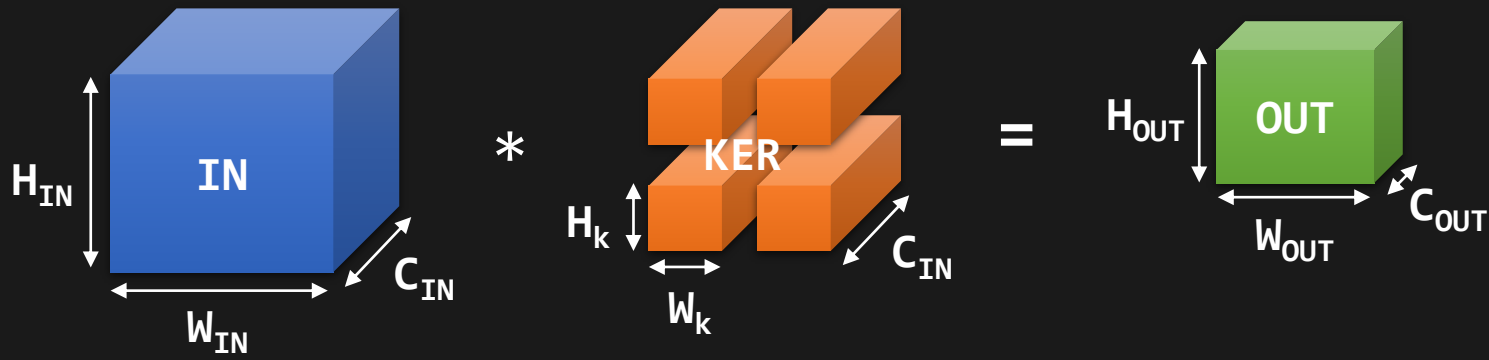
(b)

Split on N\_CORES

Table 1: PULP-TrainLib's Optimized MM Kernels

| MM Type       | Unrolling Factor | Parallelism | Description                       |
|---------------|------------------|-------------|-----------------------------------|
| mm            | -                | N or M      | Naive MM                          |
| mm_uJ         | J = 2            | N or M      | Unroll J inner products in K      |
| mm_unroll_1xV | U = 2, 4, 8      | N or M      | Unroll U columns of C             |
| mm_unroll_Ux1 | V = 2, 4, 8      | N or M      | Unroll V rows of C                |
| mm_unroll_UxV | U, V = 2, 4      | N or M      | Unroll U, V rows and columns of C |

# DEPLOYING A DNN TRAINING TASK



DEPLOYMENT  
ON MCU



TENSOR SIZE MAY  
EXCEED L1



TILING

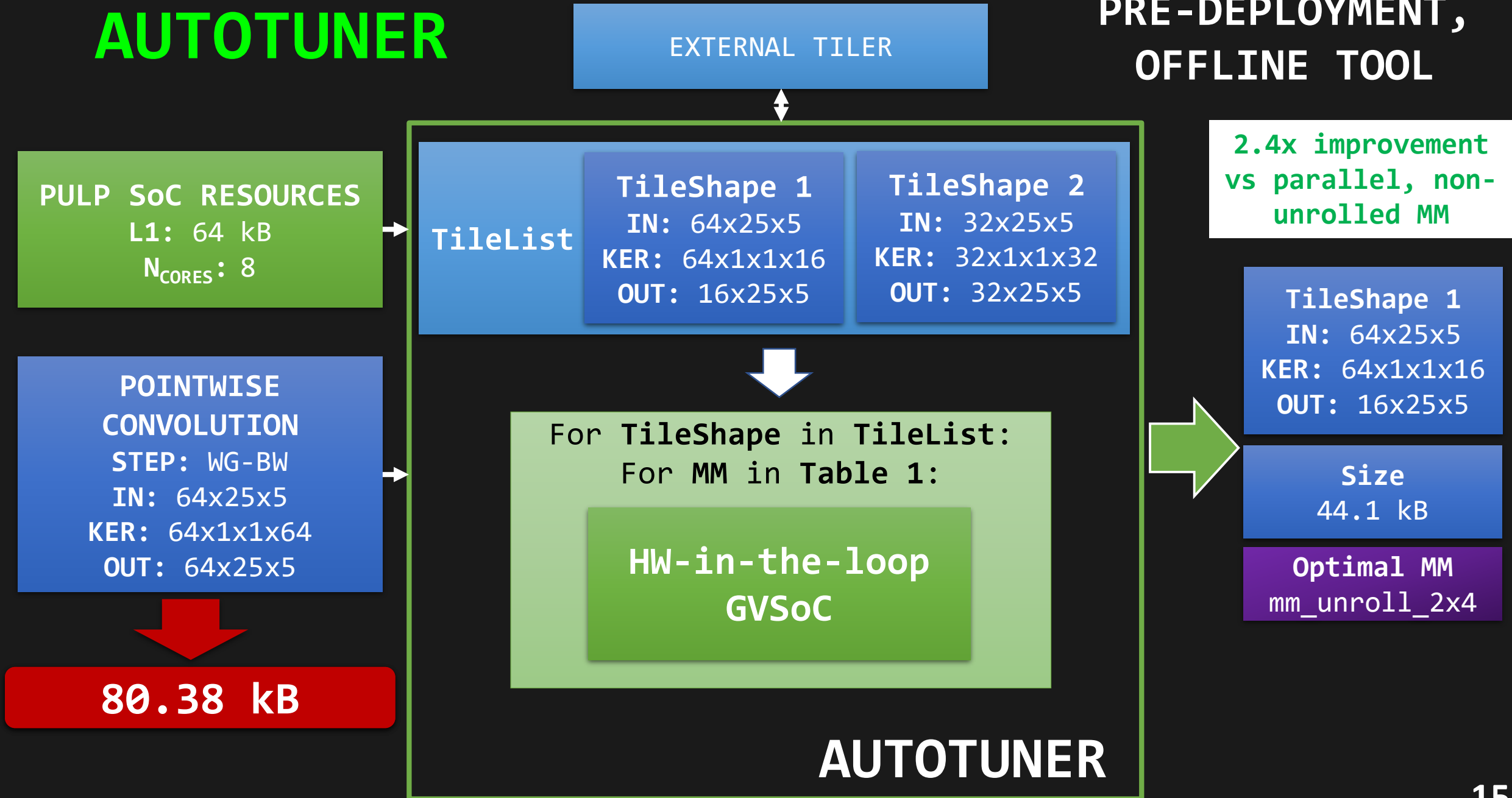
TILE SIZES NOT  
MULTIPLE OF  
UNROLLING FACTORS



OPTIMAL MM

# AUTOTUNER

PRE-DEPLOYMENT,  
OFFLINE TOOL

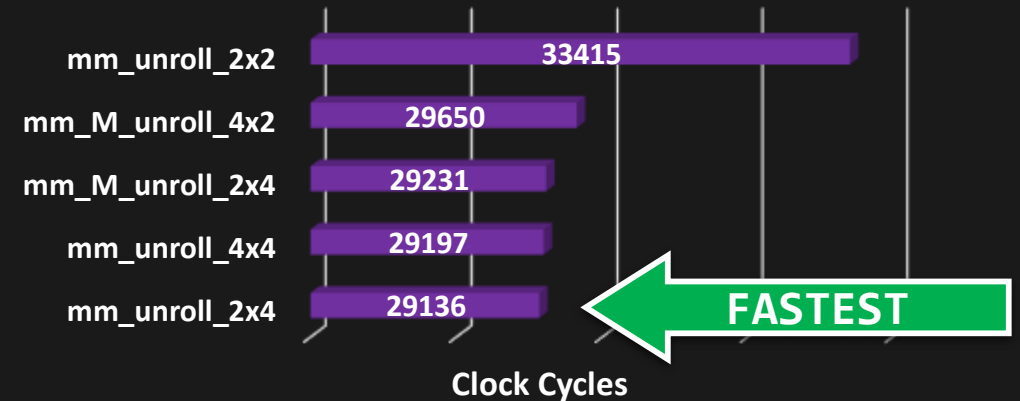
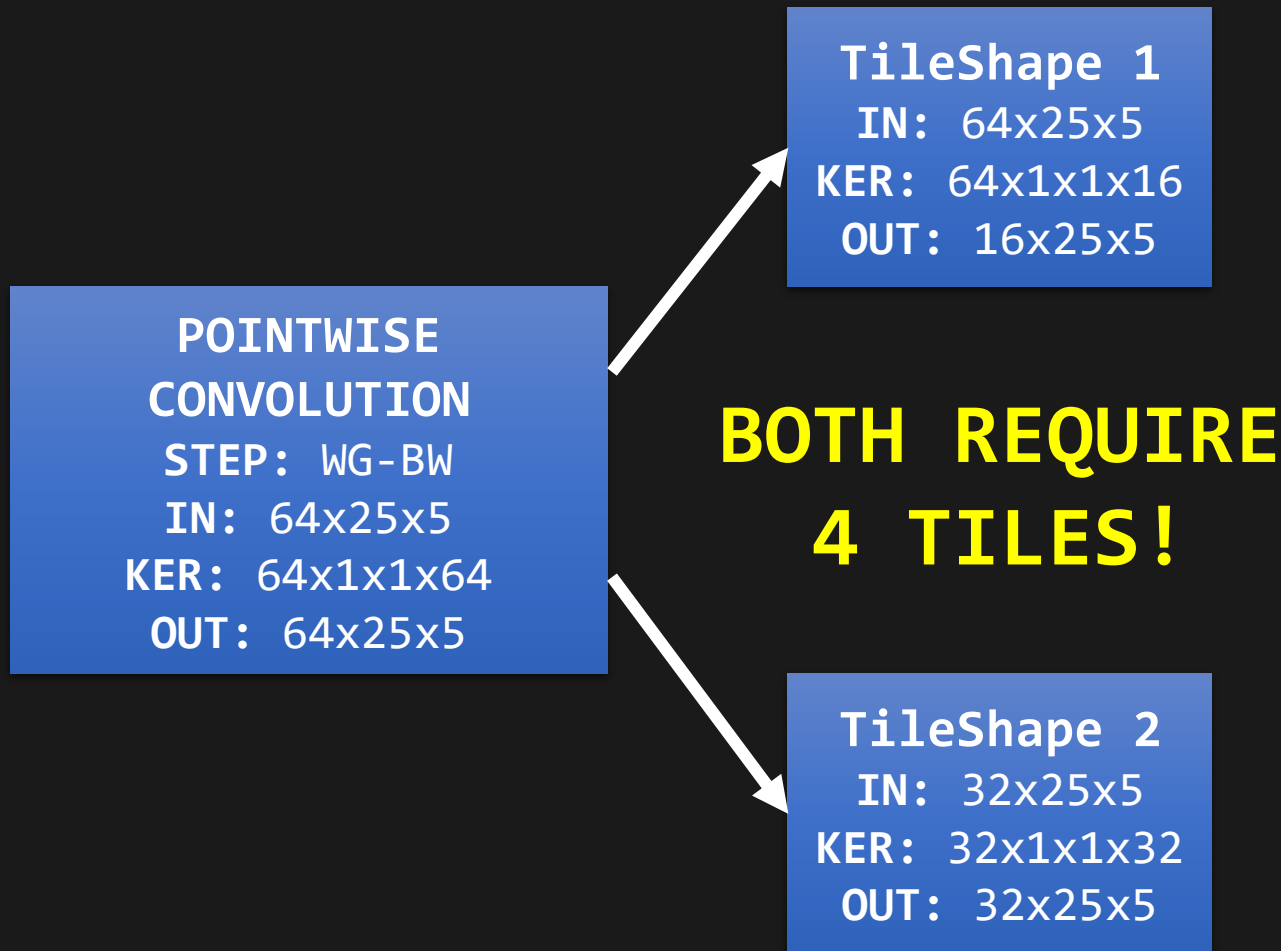




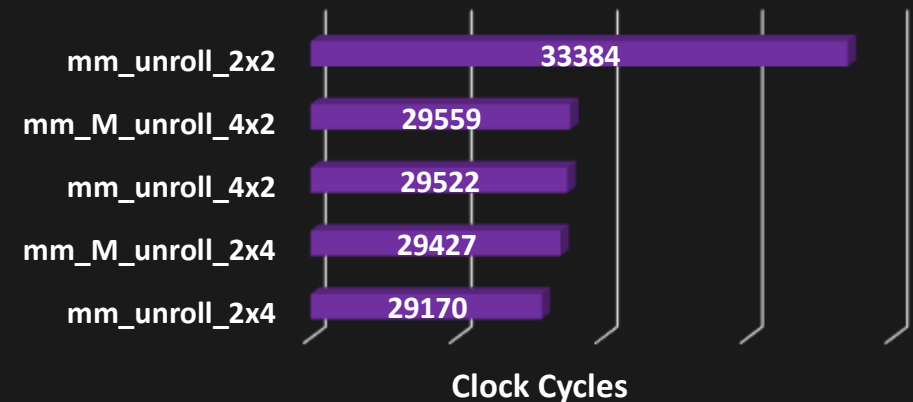
# Experimental Results



# AUTOTUNER - UNDER THE HOOD

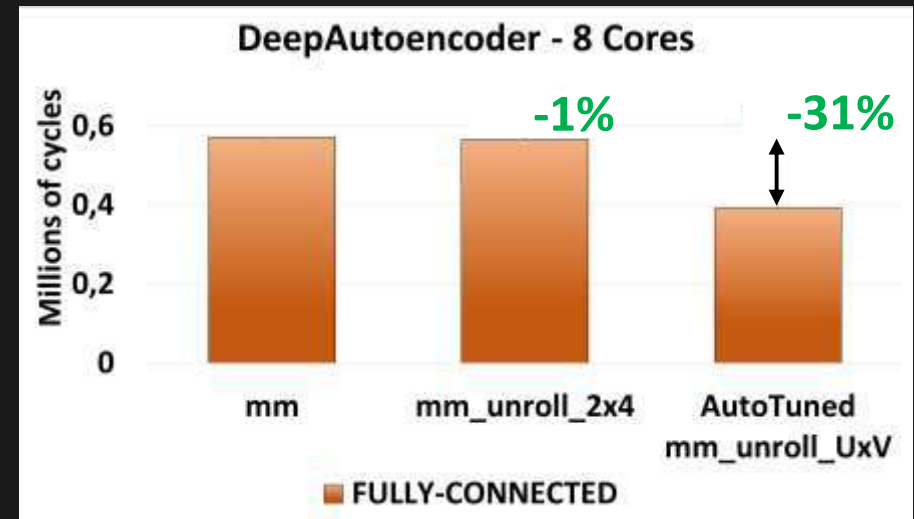
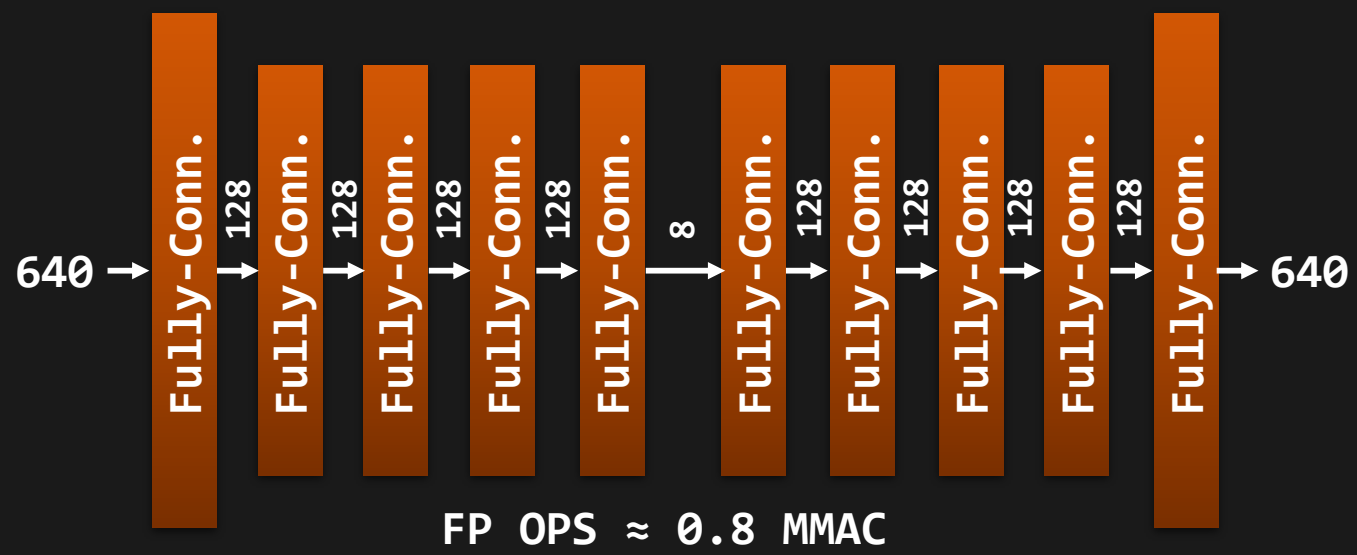
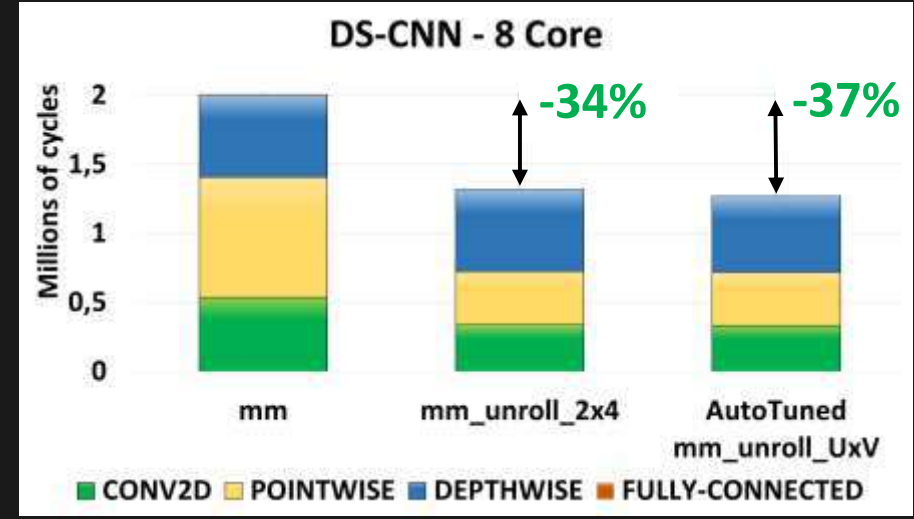
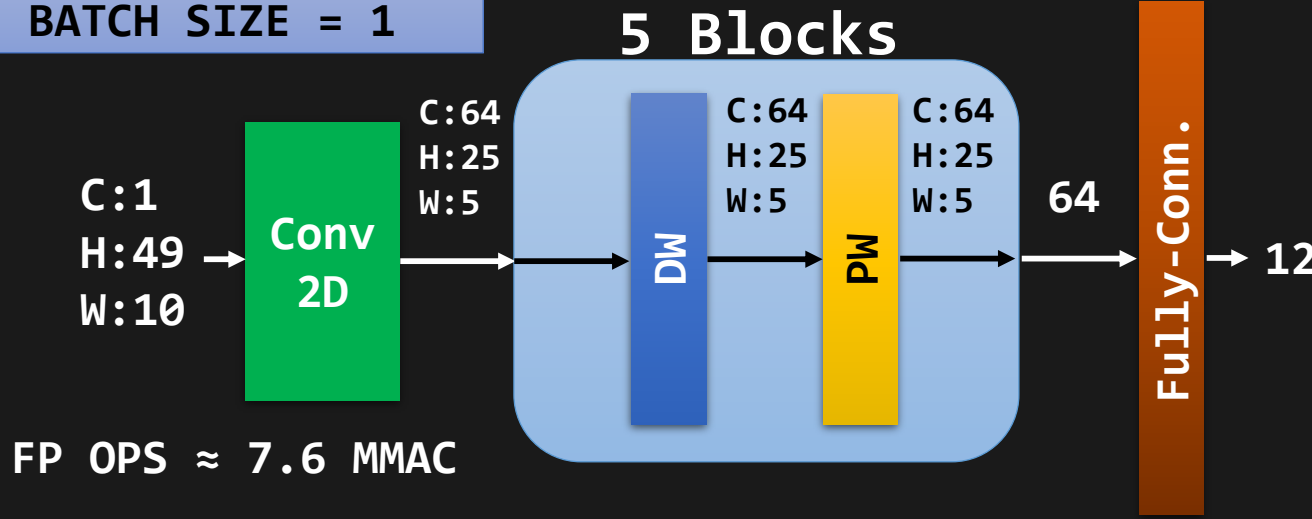


**PointWise WG-BW (with TileShape 1)**  
MAC/Clock: 4.39



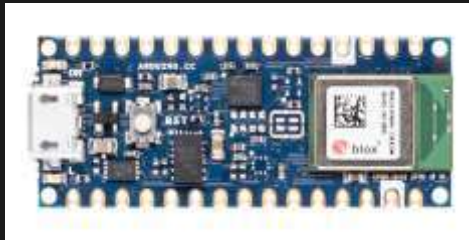
# COMPLETE TINYML MODELS' TRAINING

ONLINE LEARNING  
BATCH SIZE = 1



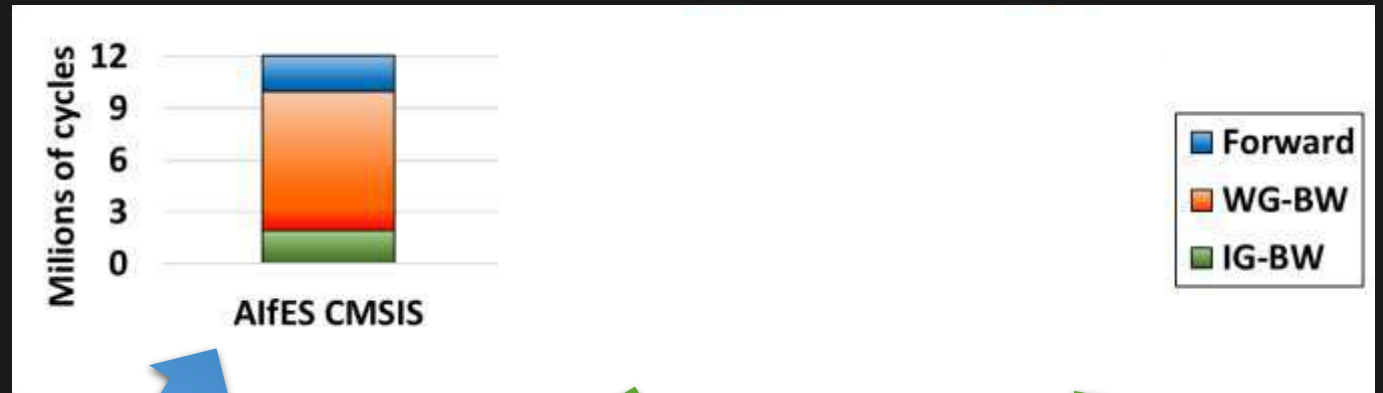
# COMPARISON WITH SOA - AIFES

AIFES: DNN on-device training/inference library for Arduino (with ARM MCUs)



Our target:  
**ARM CORTEX M4**  
VS  
GVSoc

## Training a Deep Autoencoder



AIFES core:  
ARM-CMSIS MM

PULP ISA  
(2x wrt ARM)  
+ AutoTuner

PARALLELISM  
+ AutoTuner  
(6.82x speedup  
wrt 1 RISC-V  
Core)



# Conclusion

# SUMMARY

**PULP-TrainLib**



**AutoTuner**

The first **Open-Source SW Library for On-Device Training**  
on RISC-V Multicore MCUs

Up to **4.39 MAC/clock** on **8 RISC-V cores** to execute FP32  
DNN training primitives (with AutoTuner)

**Almost real time training of TinyML models:**  
0.9 ms (Deep Autoencoder), 3 ms (DS-CNN)  
@450 MHz (PULP-based Vega SoC)

**30.7x** faster performance **wrt AlfES** on Cortex M4



THANK YOU FOR YOUR  
ATTENTION



Q. & A.