# PULP-TrainLib: On-Device Learning on Parallel Ultra-Low-Power MCUs

Davide Nadalini

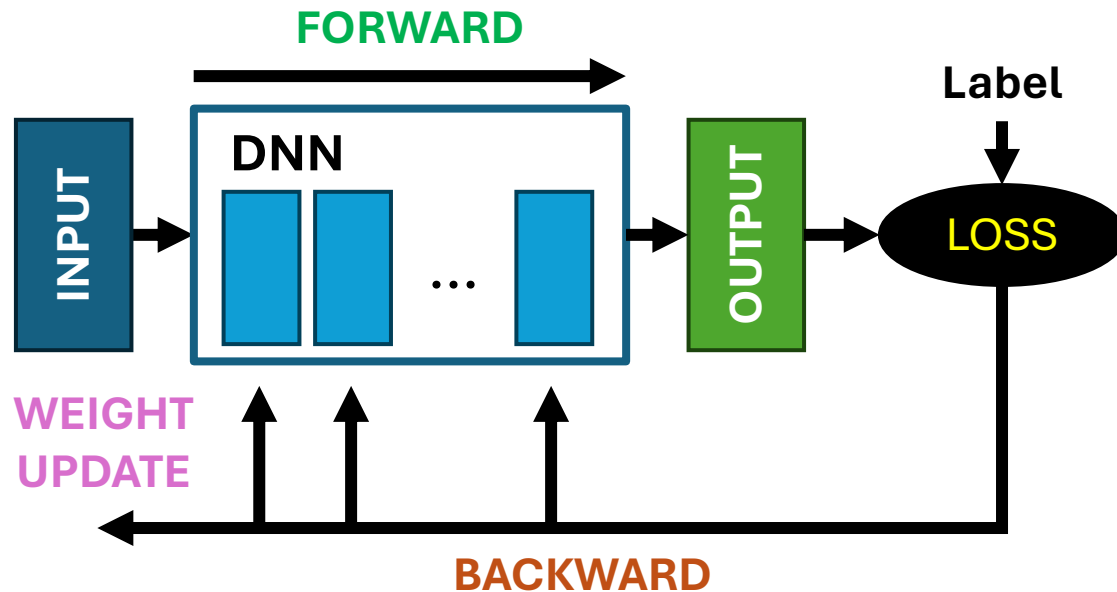d.nadalini@unibo.it

# Overview

- PULP-TrainLib

- GVSoC Simulator and Performance Metrics

- Analyzing and Generating On-Device Learning Code

- Parallelizing FP32 Code

- Optimizing FP16 Code with SIMD

- Comparison with the State-of-the-Art

# PULP-TrainLib

The first **latency-optimized open-source training library for RISC-V** multi-core **MCUs**

**FUNCTION SIGNATURE**
pulp_<layer/function>_<data_type>_<step>_<parallel?>();



**FORWARD**

**Label**

INPUT

**DNN**

...

OUTPUT

LOSS

**WEIGHT UPDATE**

**BACKWARD**

```
pulp_linear_fp32_fw_cl(&lin_args);

pulp_CrossEntropyLoss(&loss_args);

pulp_linear_fp32_bw_param_grads_cl(&lin_args);
pulp_linear_fp32_bw_input_grads_cl(&lin_args);

pulp_gradient_descent_fp32(&wgts);
```

# Repository Overview

PULP-TrainLib-Tutorial/

install_ub20.sh — Downloads, sets up requirements

pulp-trainlib/ — PULP-TrainLib + Code Generator
    tools/
        TrainLib_Deployer/
            TrainLib_Deployer.py

setup.sh — Sets the terminal up to compile your project

Ex01-TrainLib_Deployer/
    CNN_FP32/ — End-to-end ODL code
        <project_files>

**Download requirements & build PULP-SDK**

```
source install_ub20.sh
```

**IN EACH NEW TERMINAL RUN**

```
source setup.sh
```

**You can now compile you project!!**

4

# GVSoC simulator and performance metrics

```
user@PC:~/work$ source setup.sh
user@PC:~/work$ cd hello_world/
user@PC:~/work/hello_world/$ cat main.c

int main(void) {
    printf('Hello World!');
    return 0;
}


user@PC:~/work/hello_world/$ make clean all run


Hello World!
```
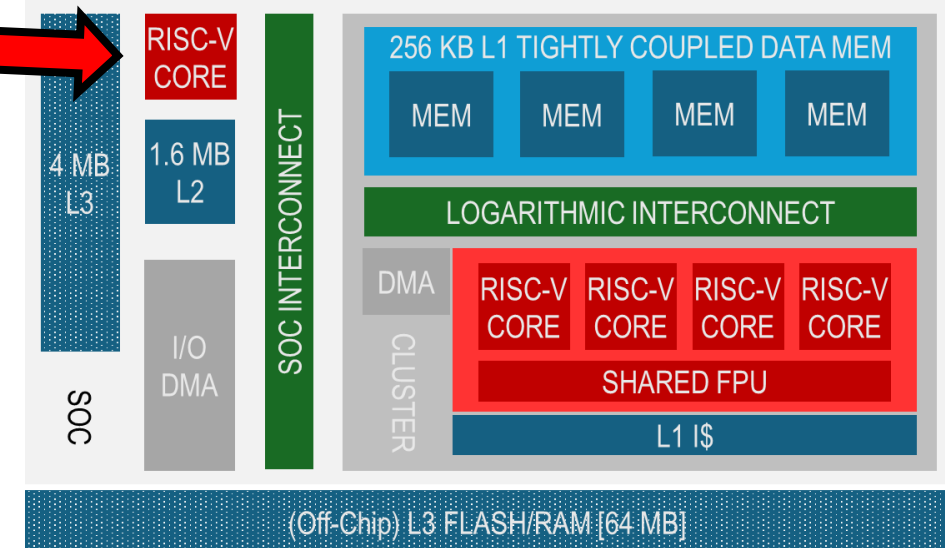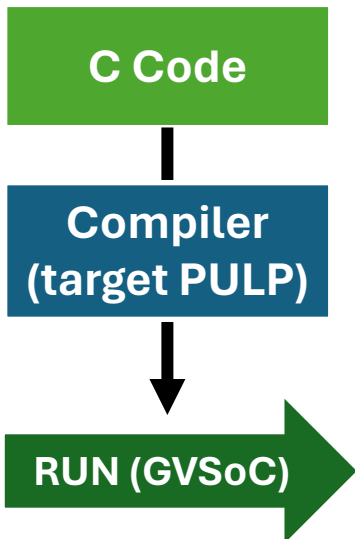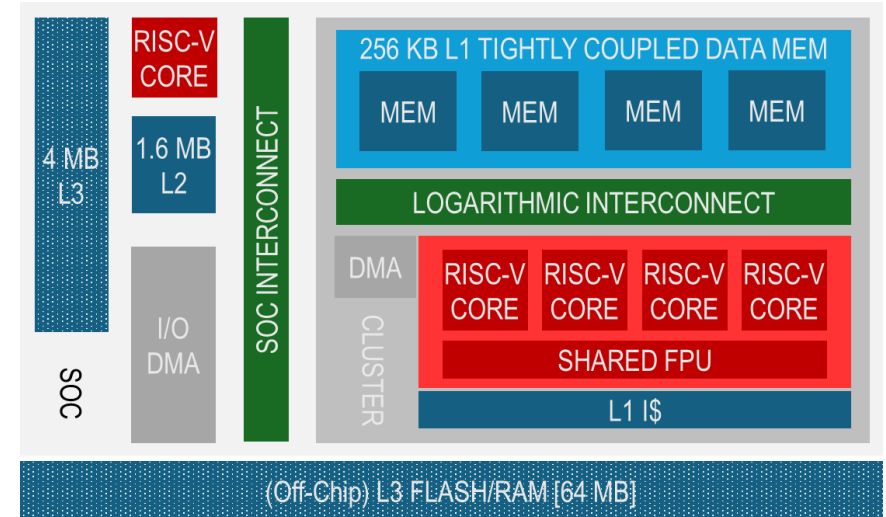
**Configurable Behavioural Simulator for the PULP Plaform**

**Included in PULP-SDK!**

# GVSoC simulator and performance metrics

**Sets up PULP-SDK & compiler toolchain**

```
user@PC:~/work$ source setup.sh
user@PC:~/work$ make clean all run
```



**C Code**

**Compiler (target PULP)**

**RUN (GVSoC)**

```
Layer 2 output:
-0.000023
....
-0.000006
Profiling performances...


[0] elapsed clock cycles = 750262      ①
[0] number of instructions = 578854    ②
[0] TCDM contentions = 0               ③
[0] load stalls = 164627              ④
[0] icache miss (clk cycles count) = 1947  ⑤
```
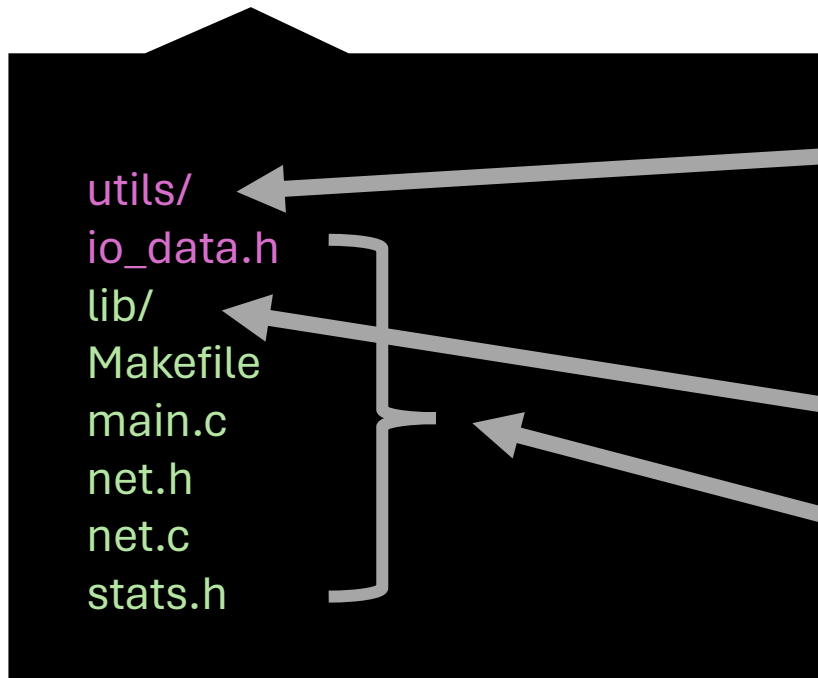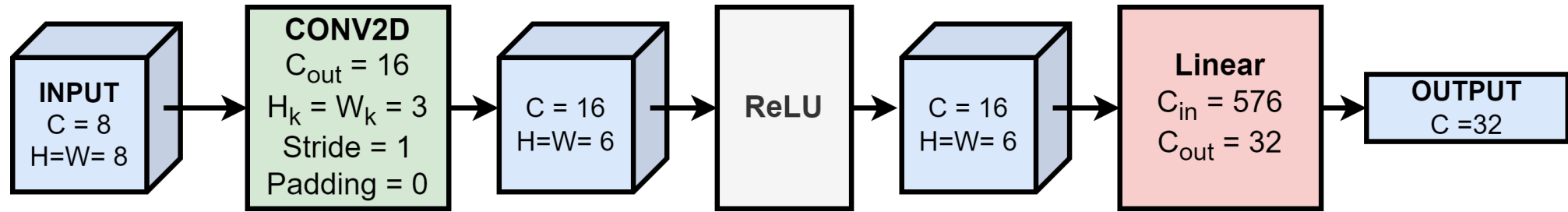
Program output

**PROFILING INFORMATION**
**Performance Counters (clock cycles)**

① **Latency of the program**

② **Number of instructions**

③ **Memory contentions in L1**

④ **Stalls while loading data from L1**

⑤ **Instruction cache misses**

# Training a model with PULP-TrainLib



INPUT
C = 8
H=W= 8

CONV2D
$C_{out}$ = 16
$H_k$ = $W_k$ = 3
Stride = 1
Padding = 0

C = 16
H=W= 6

ReLU

C = 16
H=W= 6

Linear
$C_{in}$ = 576
$C_{out}$ = 32

OUTPUT
C =32

utils/
io_data.h
lib/
Makefile
main.c
net.h
net.c
stats.h

**Test Golden Model data (PyTorch)**

**PULP-TrainLib**

**ODL C code**

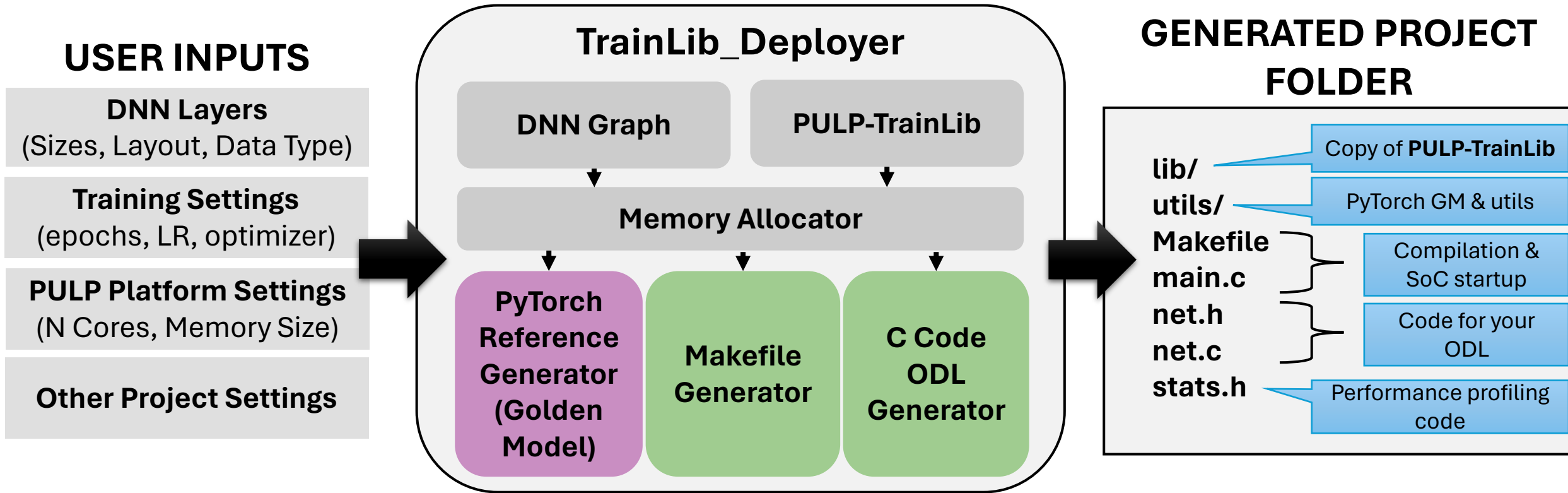**Set up environment variables**

```
source setup.sh
```

**Compile and run your ODL project**

```
cd Ex01-TrainLib_Deployer/CNN_FP32/
make clean get_golden all run
```

**DNN DATA TYPE = FP32, MEMORY = 207 kB**

# TrainLib-Deployer: ODL code generator



**USER INPUTS**

- **DNN Layers** (Sizes, Layout, Data Type)
- **Training Settings** (epochs, LR, optimizer)
- **PULP Platform Settings** (N Cores, Memory Size)
- **Other Project Settings**

**TrainLib_Deployer**

- DNN Graph
- PULP-TrainLib
- Memory Allocator
- PyTorch Reference Generator (Golden Model)
- Makefile Generator
- C Code ODL Generator

**GENERATED PROJECT FOLDER**

- lib/ → Copy of **PULP-TrainLib**
- utils/ → PyTorch GM & utils
- Makefile → Compilation & SoC startup
- main.c
- net.h → Code for your ODL
- net.c
- stats.h → Performance profiling code
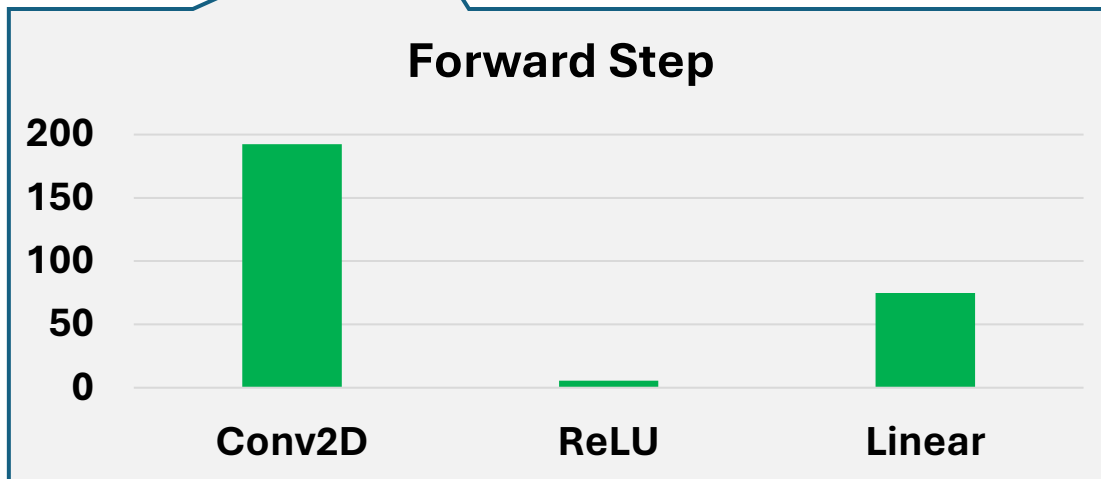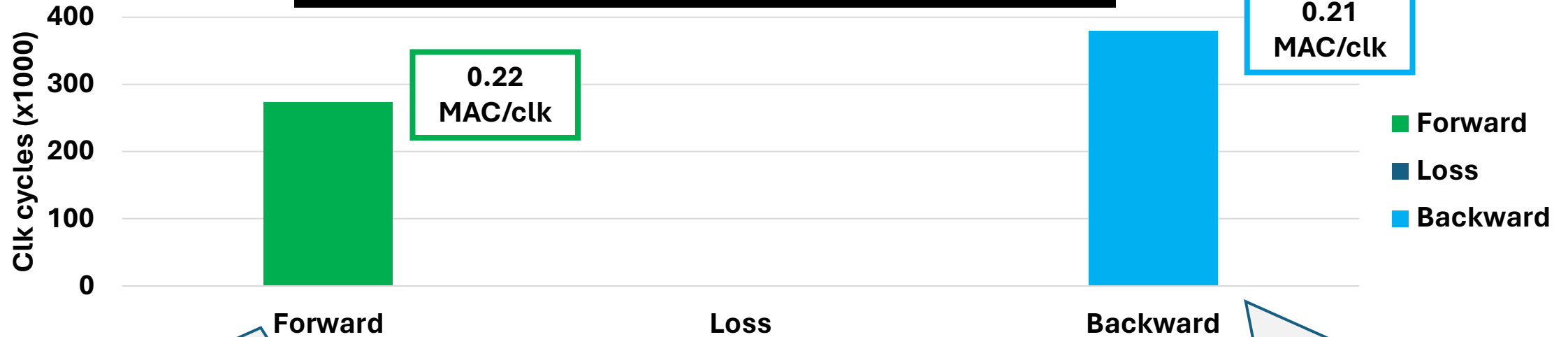
**Generate project with TrainLib_Deployer**

```
cd ../pulp-trainlib/tools/TrainLib_Deployer
python TrainLib_Deployer.py
```

# WHAT ABOUT THE CODE?
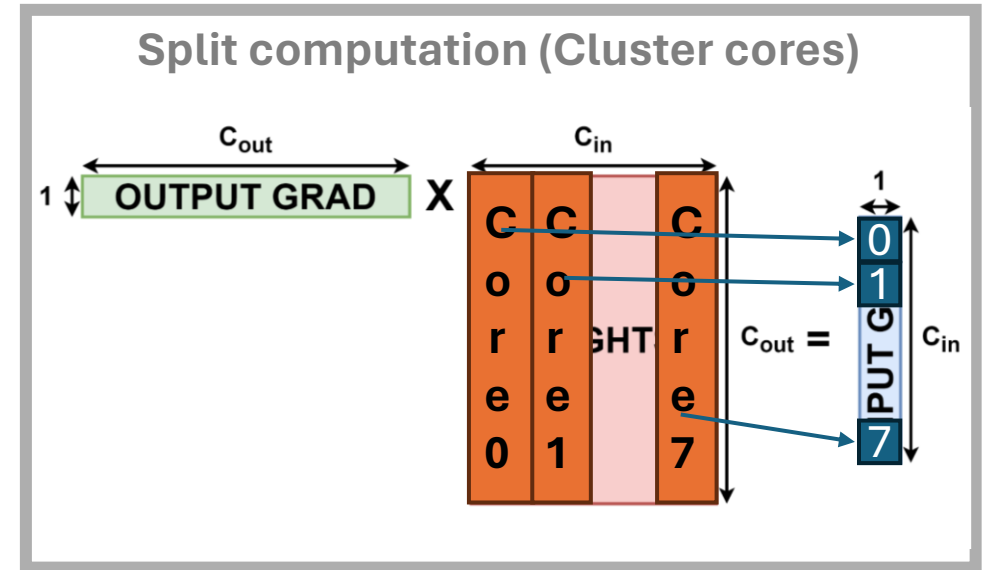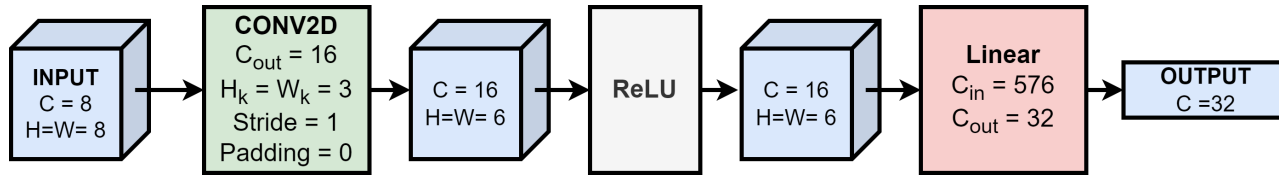
# ODL: latency of each step (FP32) – single core

```
cd Ex01-TrainLib_Deployer/CNN_FP32/
make clean get_golden all run NUM_CORES=1
```
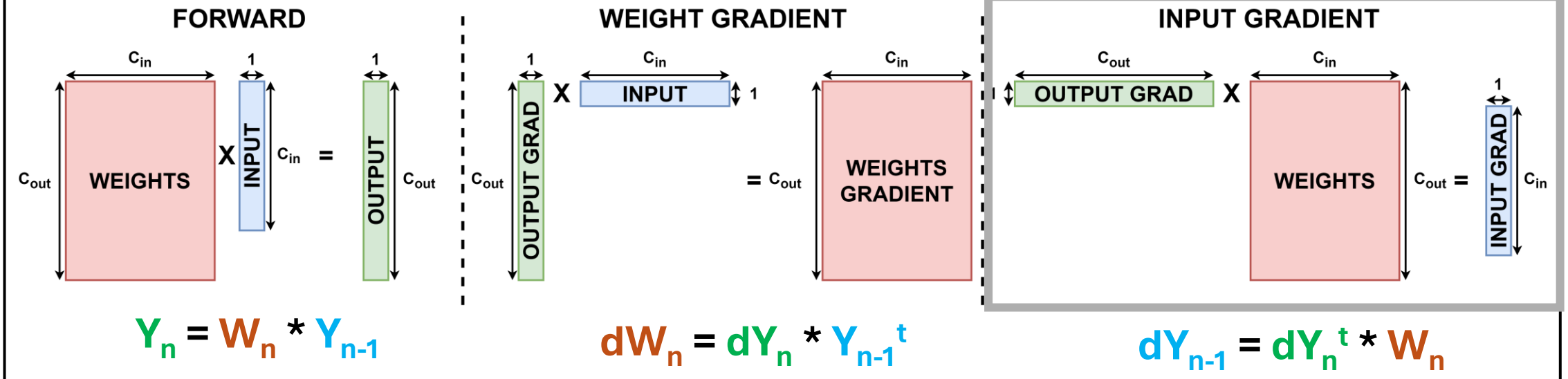
**Clk cycles (x1000)**

0.22 MAC/clk

0.21 MAC/clk

- Forward
- Loss
- Backward

Forward        Loss        Backward

## Forward Step

Conv2D        ReLU        Linear

## Backward Step

- Input Grad
- Weight Grad

Conv2D        ReLU        Linear

# Optimization 1: Parallelization

# Parallelizing ODL code



**Computations based on Linear Algebra**

$$Y_n = W_n * Y_{n-1}$$

$$dW_n = dY_n * Y_{n-1}^t$$

$$dY_{n-1} = dY_n^t * W_n$$

# WHAT ABOUT THE CODE?

# Parallelizing ODL code (cont'd)

```
cd Ex01-TrainLib_Deployer/CNN_FP32/
make clean get_golden all run NUM_CORES=8
```

**IDEAL: 8x speedup**

**REAL: TCDM contentions + parallelization overhead**

# Optimization 2: SIMD Vectorization (FP16)

# FP16 Optimization

## Pseudo-Assembly Instructions

**Floating-Point 32 (FP32)**

| 1 | 8 bits | 23 bits |
|---|--------|---------|

sign
exponent
mantissa

**Floating-Point 16 (FP16)***

| 1 | 5 bits | 10 bits |
|---|--------|---------|

sign
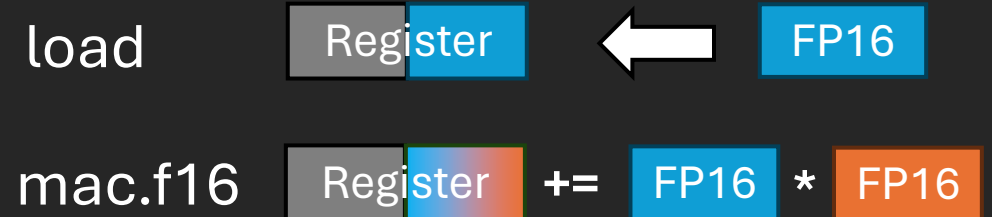exponent
mantissa

Memory Occupation (DNN):

~~207 kB (FP32)~~  →  **104 kB (FP16)**

## What about the latency?

*IEEE or BFLOAT (user-selectable, platform support for both)

### Non-Vectorized Instructions

load    Register  ⬅  FP16

mac.f16  Register  +=  FP16  *  FP16

### FP16 SIMD Vectorized Instructions

vload   Register  ⬅  FP16 FP16

vfmac   Register  +=  FP16  *  FP16
                  +=  FP16  *  FP16

vpack   Register  ⬅  Register
                     Register

# Linear Input Grad: Naive Vector-Matrix (FP16)

# FP16: maximizing SIMD performance



**Output Grad (O)** X

Row 0
Row 1
**Weights (W)** = **Input Grad (I)**

**L1 Memory**

**Row-Column Multiplication**

X X X

**Partial Products**

**SIMD Units:** we want **vectorized instructions**

**PSEUDO-ASSEMBLY (naive SIMD)**

```
vload     ra, O[k], O[k+1]          Loop K/2 times
load      rs2, W[k*M+j]
load      rs3, W[(k+1)*M+j]

vpack     rb, rs2, rs3
vfmac     rd1, ra, rb

add       rd, rd1, rd2
store     I[j], rd
```

**Vectorization overhead!!**

5 instructions / 2 MAC

18

# FP16: maximizing SIMD performance (cont'd)

# WHAT ABOUT THE CODE?

# FP16: maximizing SIMD performance (cont'd)

**Fully-Connected (Input Gradient)**

`cd Ex02-Fully-Connected-FP16/test_linear_fp16/`



Chart: Clock Cycles vs. vm_naive, vm_SIMD_naive, vm_T_SIMD

- **1.54x** — Not 2x!!! (single loads + pack)
- **1.91x** — Almost 2x!!
- **1.24x**

**Vector-Matrix (no SIMD)**

```
make clean get_golden all run
MATMUL_TYPE=0 TRANSPOSE_WEIGHTS=0
```

**Naive SIMD**

```
make clean get_golden all run
MATMUL_TYPE=1 TRANSPOSE_WEIGHTS=0
```

**SIMD with transposed Weights**

```
make clean get_golden all run
MATMUL_TYPE=2 TRANSPOSE_WEIGHTS=1
```

# Stacking Optimizations (Linear Input Grad)

Naive VM → Naive SIMD VM → Optimized SIMD VM → Parallel & Optimized SIMD VM

**Latency (Clock cycles, FP16)**

```
make clean get_golden all run MATMUL_TYPE=2
TRANSPOSE_WEIGHTS=1 NUM_CORES=8
```



**1.55x**

**1.91x**

**14.59x**

| Naive VM | Naive SIMD | Optimized SIMD | Parallel SIMD |
|---|---|---|---|
| No optimization | Only SIMD MAC, no vector loads | SIMD loads, MAC | SIMD loads, MAC, parallel execution |

# PULP-TrainLib: Further Optimizations

**Conv2D**

**Im2col + padding**

TOP = 4.55 MAC/clk
(22.1x wrt FP32)



**PointWise Convolution**

TOP = 4.36 MAC/clk
(21x wrt FP32)

**OPTIMIZATIONS:**
- FP32 Parallelization
- FP32 Unrolling + Parallelization
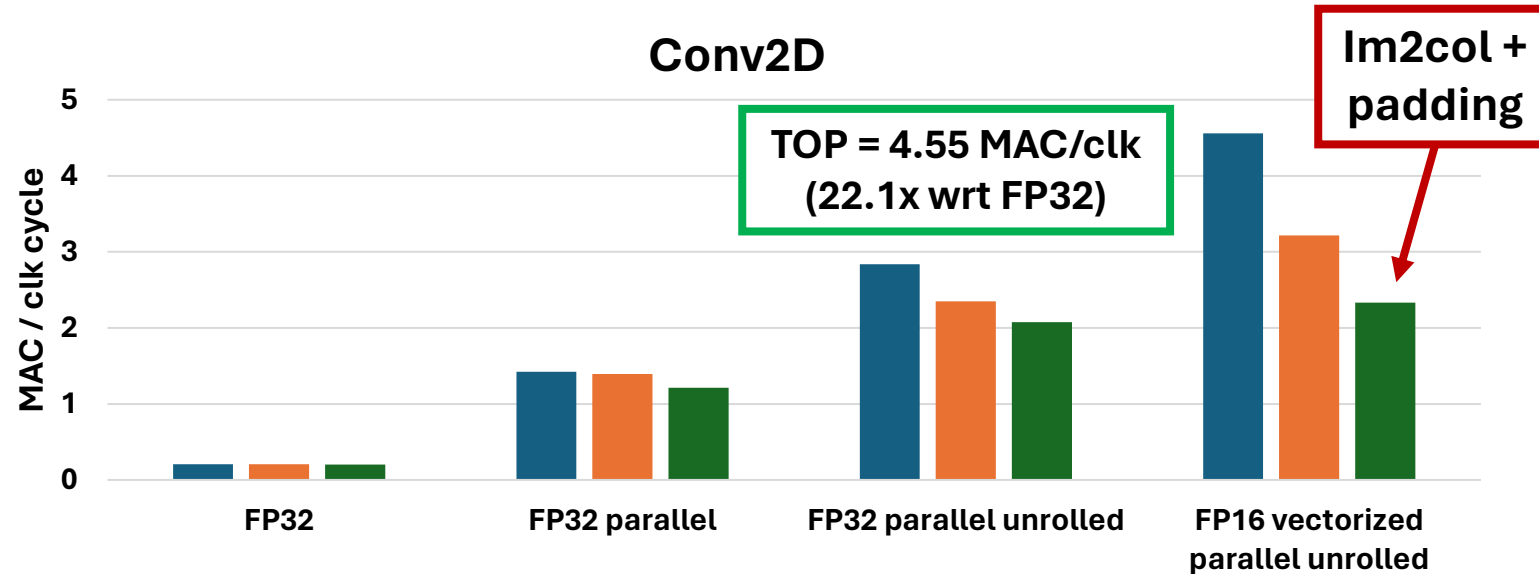- FP16 Vectorization + Unrolling + Parallelization

**SUPPORTED FEATURES:**
- **Layers**: Conv2D, PointWise, DepthWise, Linear
- **Activations**: ReLU, SoftMax
- **Losses**: MSE, CrossEntropy
- **Optimizers**: SGD
- **Pooling**: MaxPool, AvgPool

■ **Forward** ■ **Weight Grad** ■ **Input Grad**

# Comparison with the State-of-the-Art

| ODL Library / Method | Target Task | Target Device | Key Optimization | Retrainable Layers | Sparse Update | Data Type | Code Generation | Peak Training Performance (MAC/clk) |
|---|---|---|---|---|---|---|---|---|
| PULP-TrainLib [1, 2] | General Purpose | Multicore RISC-V MCUs | Parallelism, SIMD, unrolling | All (Convs, Linear) | No | FP32, FP16 | Manual DNN definition, auto memory management | 6.62[1] |
| TTE [3] | Image Classification / General Purpose | STM32 | Quantized Sparse Update | All (Convs, Linear) | Layers, Weight Channels, Biases | INT8, FP32 | Graph analysis with compile-time autodiff | ~ 0.1[2] |
| AIfES [5] | General Purpose | General Purpose | Unrolling (MatMul) | All (Convs, Linear) | No | FP32, INT32, INT8 | Graph Analysis (Pytorch, Keras) | 0.15[3] |

**WE ARE OPEN FOR CONTRIBUTION!**

[1]Profiled on hardware (Greenwaves GAP9)
[2]Extrapolation from data published in [3]
[3]Profiled on hardware (STM32L476RG)

# References

[1] Nadalini, D., Rusci, M., Tagliavini, G., Ravaglia, L., Benini, L. and Conti, F., 2022, July. **PULP-TrainLib: Enabling on-device training for RISC-V multi-core MCUs through performance-driven Autotuning**. In *International Conference on Embedded Computer Systems* (pp. 200-216). Cham: Springer International Publishing.

[2] Nadalini, D., Rusci, M., Benini, L. and Conti, F., 2023. **Reduced Precision Floating-Point Optimization for Deep Neural Network On-Device Learning on MicroControllers**. *arXiv preprint arXiv:2305.19167.*

[3] Lin, J., Zhu, L., Chen, W.M., Wang, W.C., Gan, C. and Han, S., 2022. **On-device training under 256kb memory**. *Advances in Neural Information Processing Systems*, *35*, pp.22941-22954.

[4] Wulfert, L., Kühnel, J., Krupp, L., Viga, J., Wiede, C., Gembaczka, P. and Grabmaier, A., 2024. **AIfES: A Next-Generation Edge AI Framework**. *IEEE Transactions on Pattern Analysis and Machine Intelligence.*

# Thank you for your attention

Davide Nadalini

d.nadalini@unibo.it