



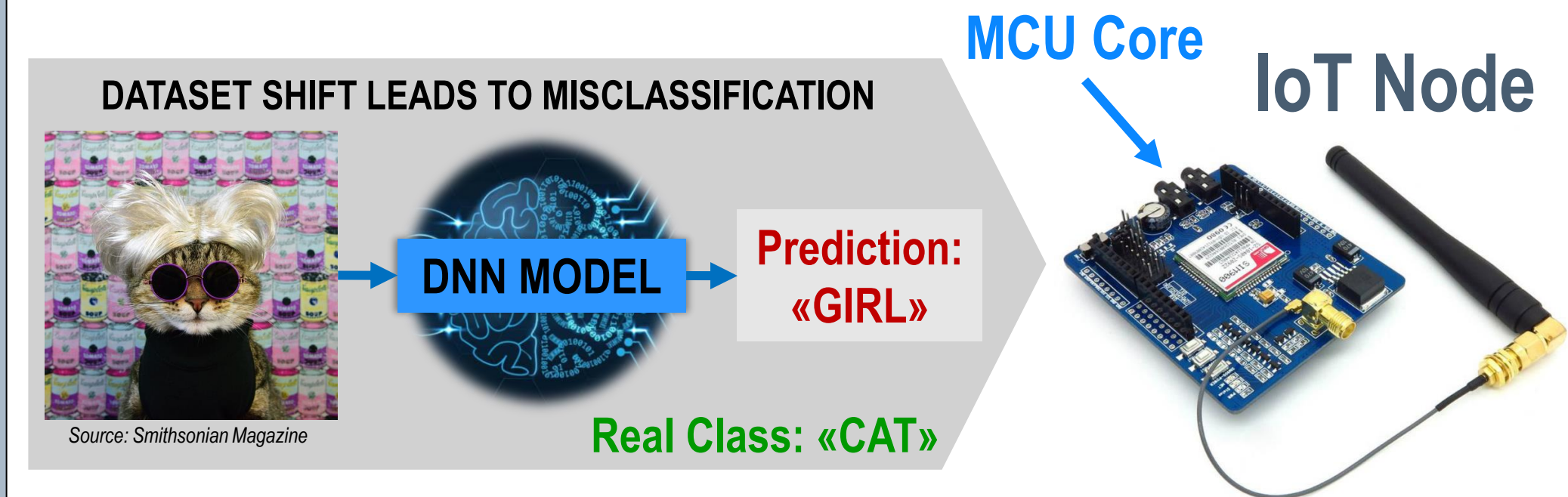
On-Device Learning (ODL) on RISC-V Multicore MCUs

Daide Nadalini^{1,2}, Manuele Rusci³, Luca Benini^{2,4} and Francesco Conti²

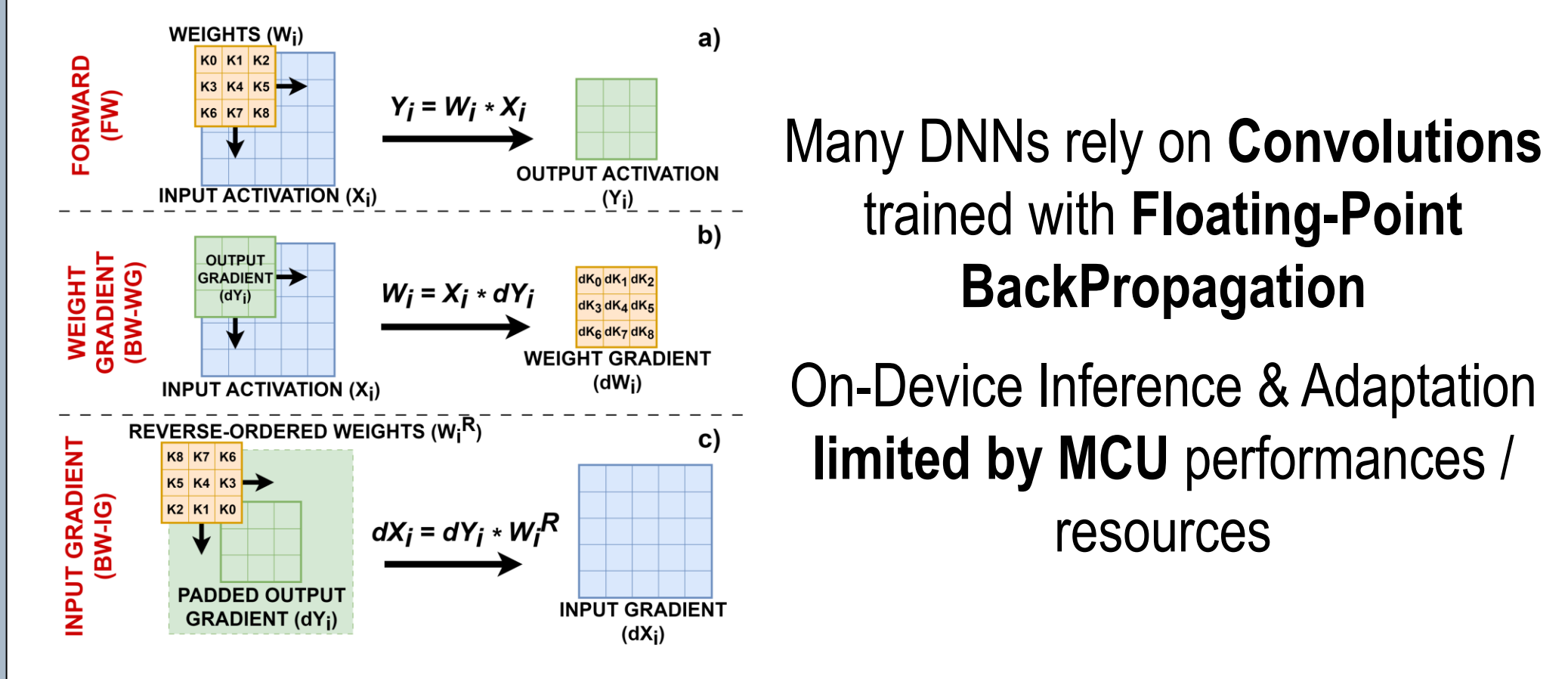
¹Politecnico di Torino, Torino, Italy, ²Università di Bologna, Bologna, Italy, ³Katholieke Universiteit Leuven, Leuven, Belgium, ⁴ETH Zurich, Zurich, Switzerland



Motivation



On-Device Learning (ODL): the ability to "update a model without data leaving your users' devices" (Tensorflow Lite)

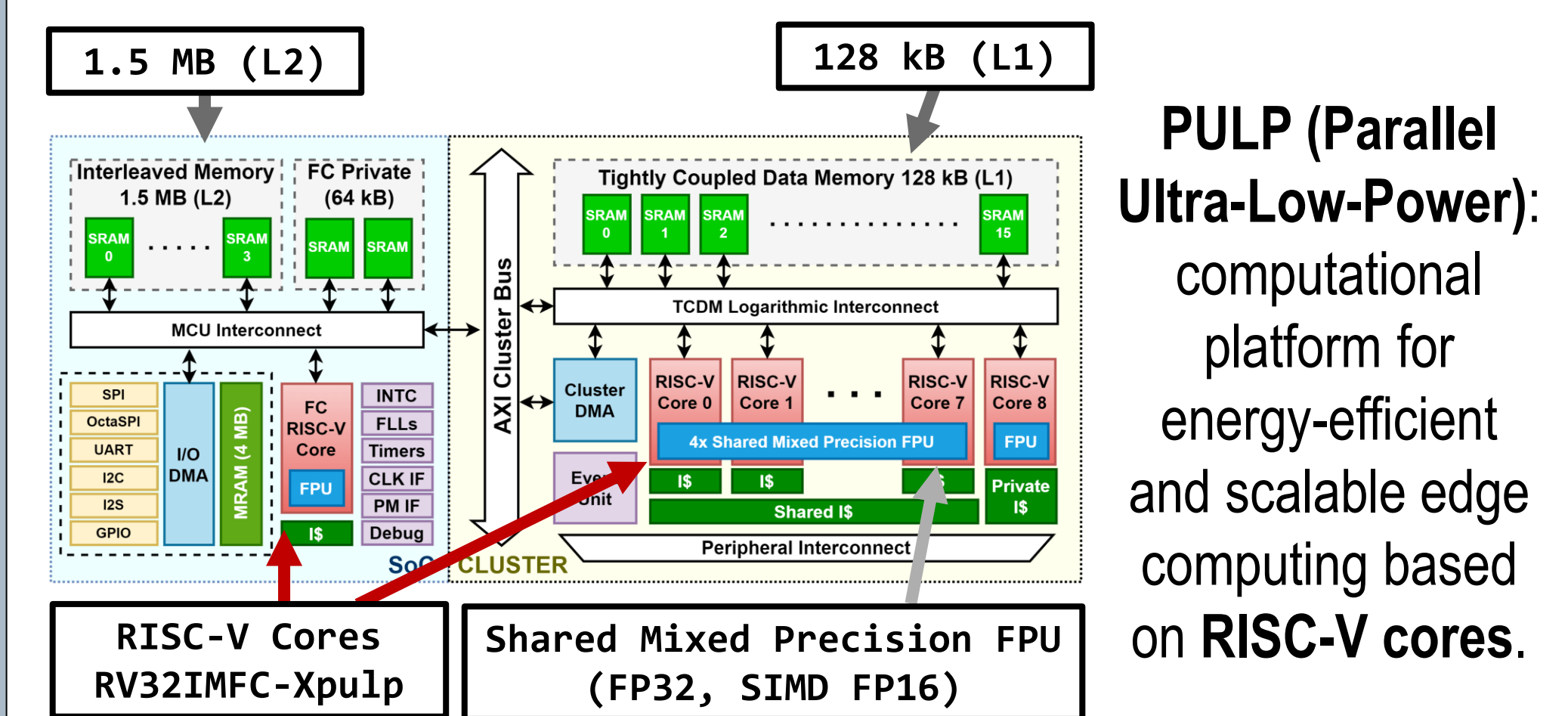


Many DNNs rely on Convolutions trained with Floating-Point BackPropagation

On-Device Inference & Adaptation limited by MCU performances / resources

Our Goal: enabling ODL on ultra-low-power MCUs

Our target: the PULP Platform¹



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

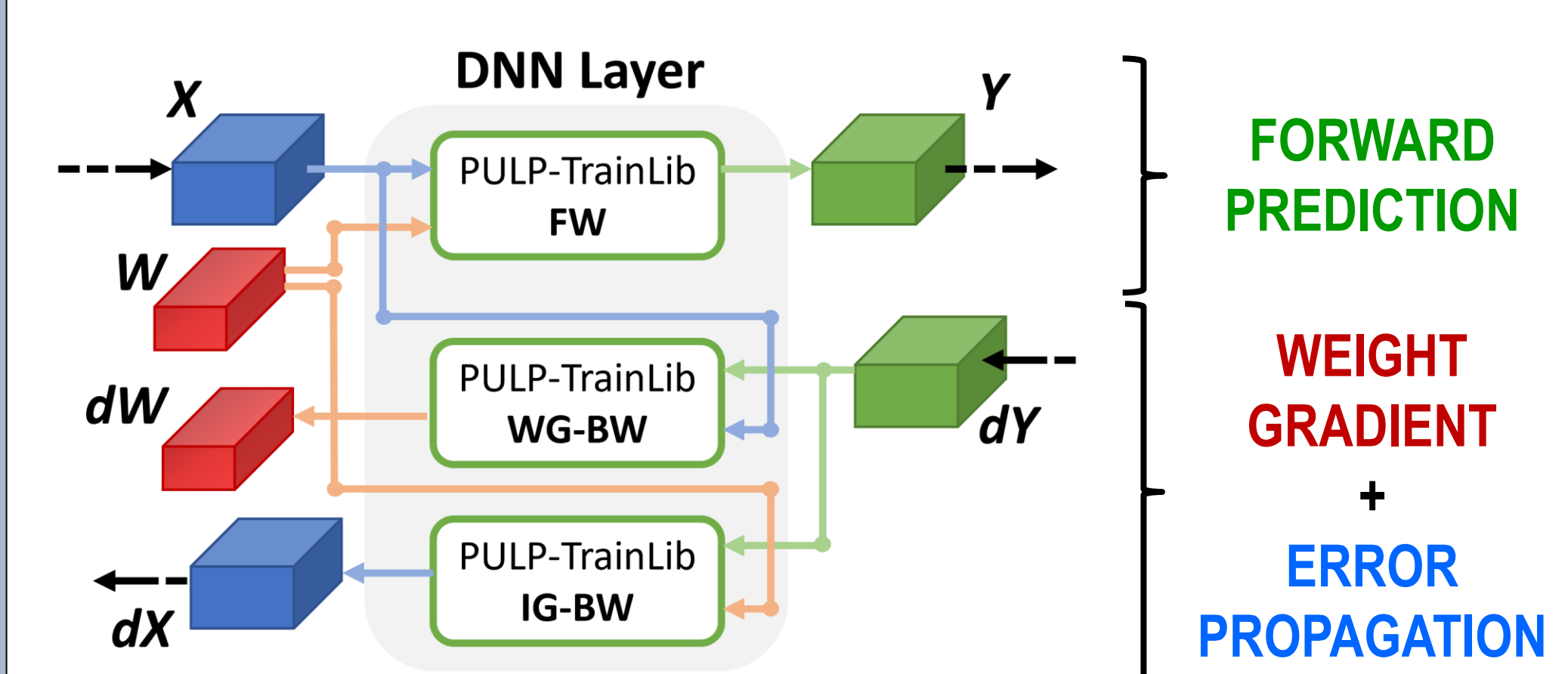
1) PULP-TrainLib: Enabling On-Device Training for RISC-V Multi-Core MCUs through Performance-Driven Autotuning²

Daide Nadalini^{1,2}, Manuele Rusci², Giuseppe Tagliavini², Leonardo Ravaglia², Luca Benini^{2,3}, and Francesco Conti²

¹Politecnico di Torino, ²Università di Bologna, ³ETH Zurich

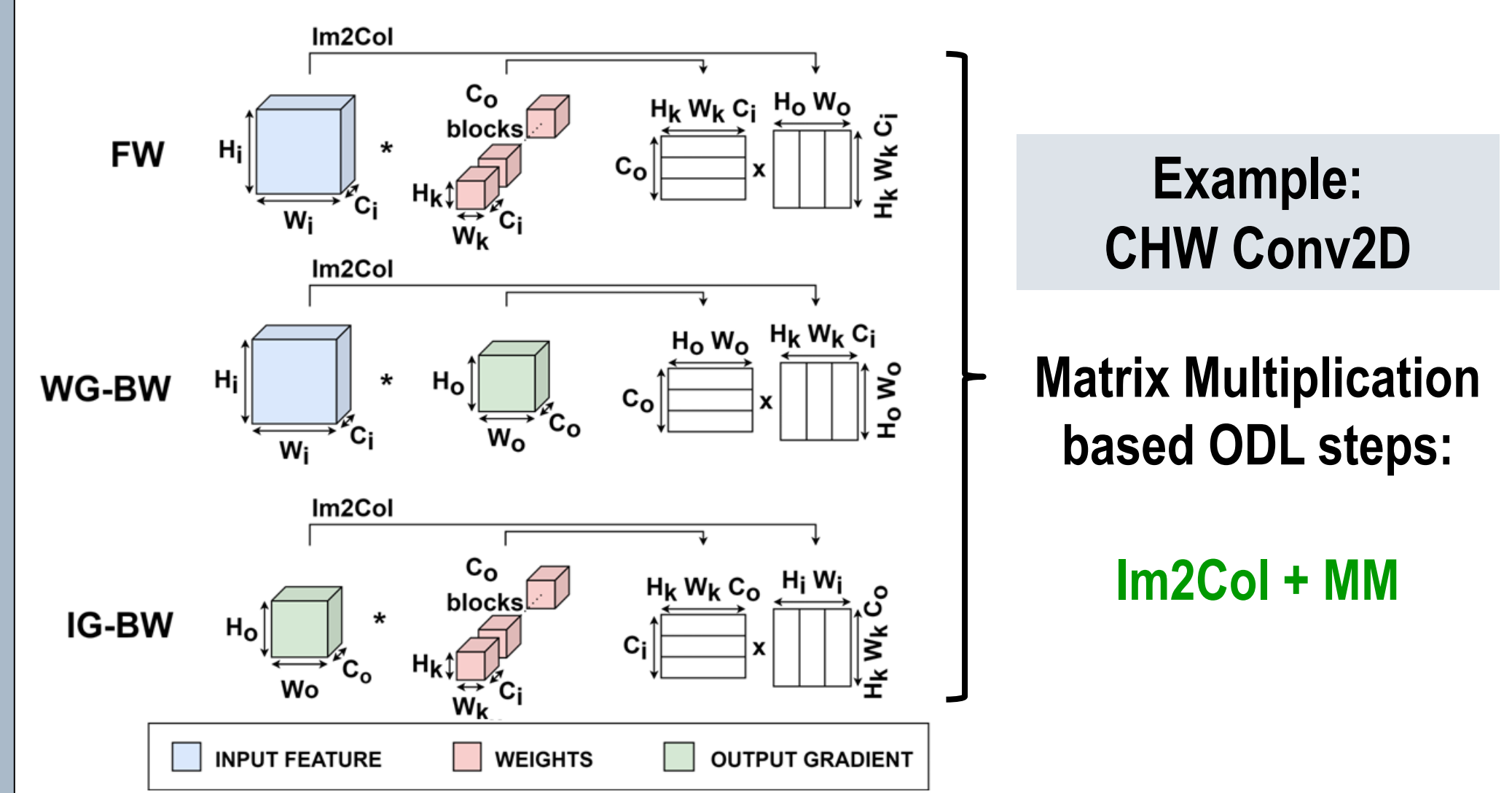
1) PULP-TrainLib, the first open-source³ training library for RISC-V multicore MCUs with Matrix-Multiplication (MM)-based performance-tunable Floating-Point (FP) primitives.

BackPropagation-based ODL primitives of CNN models



¹PULP Platform: <https://pulp-platform.org/>

²D. Nadalini, M. Rusci, G. Tagliavini, L. Ravaglia, L. Benini, and F. Conti, "PULP-TrainLib: Enabling On-Device Training for RISC-V Multi-Core MCUs through Performance-Driven Autotuning"

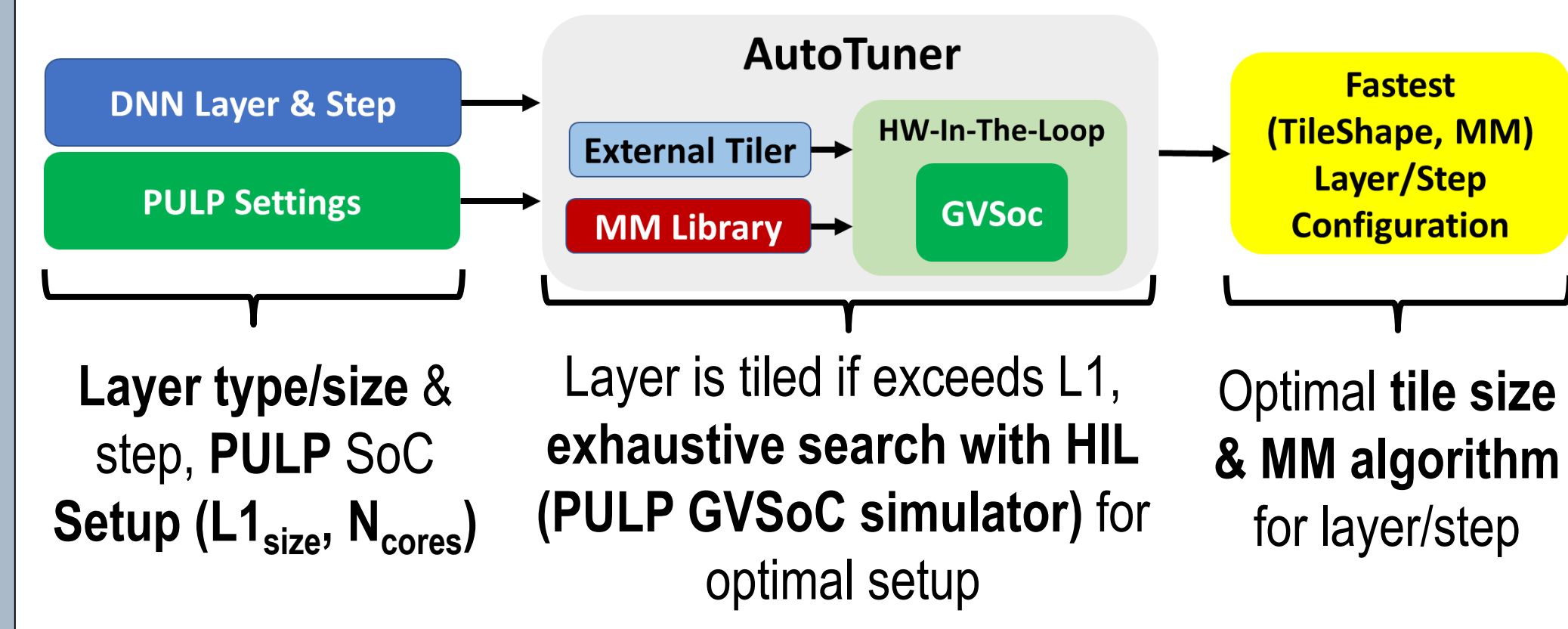


Example: CHW Conv2D
Matrix Multiplication based ODL steps:
Im2Col + MM

TABLE 1: PULP-TRAINLIB'S OPTIMIZED MM ALGORITHMS

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

2) AutoTuner, an HW-in-the-loop (HIL)-based tool to optimize PULP-TrainLib's primitives according to the DNN layer and step and the target PULP Platform settings

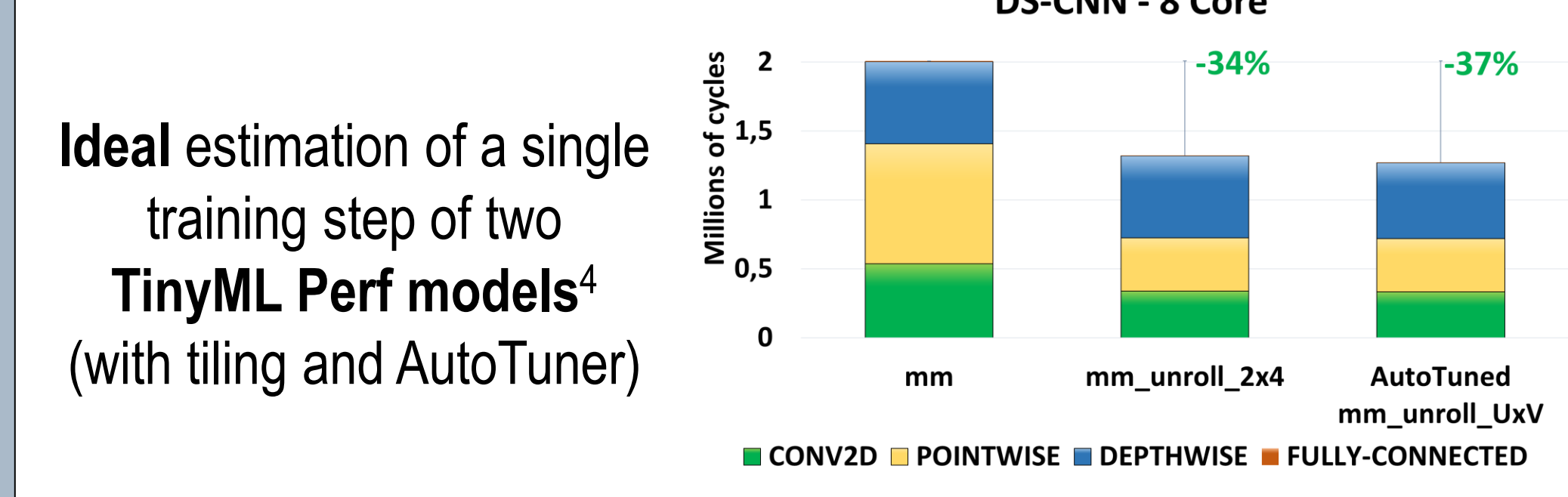
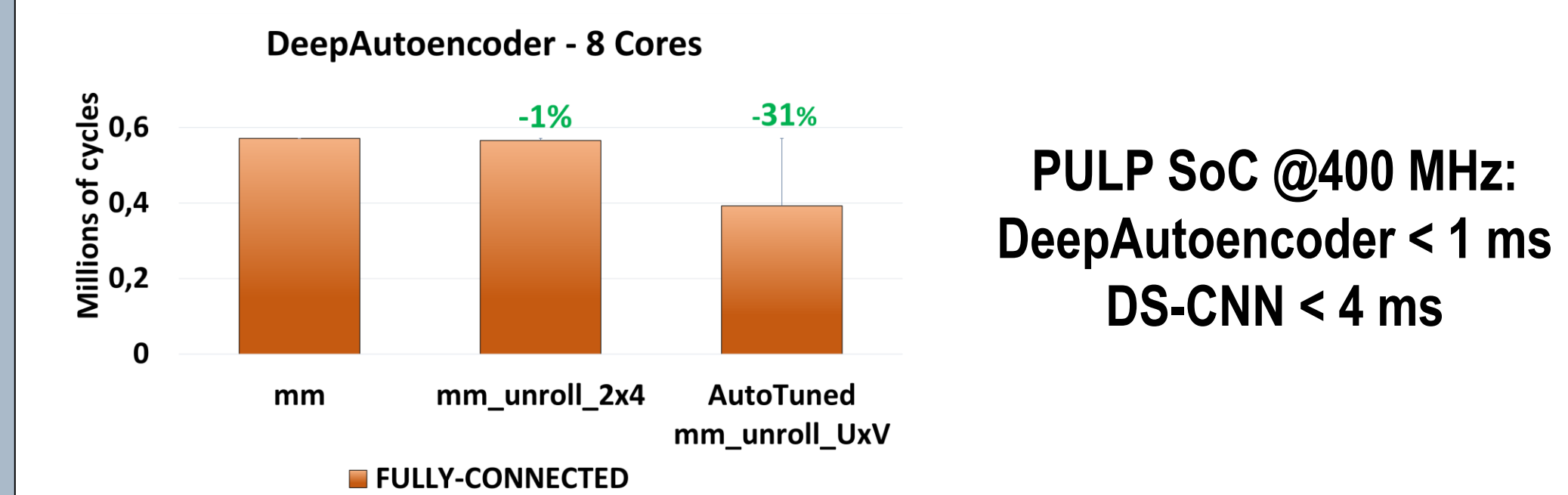


3) A detailed analysis of PULP-TrainLib and AutoTuner on an 8-Core PULP Platform

PULP GVSoc Setup: $L1_{SIZE} = 64\text{ kB}$, $N_{CORES} = \{1, 8\}$

	FW	WG-BW	IG-BW
Pointwise Convolution	0.10, 0.24, 0.53	1.87, 3.96	1.60, 3.50
STIM32	0.12, 0.24, 0.56	1.85, 4.39	1.60, 3.50
NAIVE	NAIVE	NAIVE	NAIVE
OPT	NAIVE	NAIVE	NAIVE

❖ Up to 2.4x speedup (autotuned vs one-size-fits-all, 4.39 MAC/clock on 8 RISC-V cores)
❖ 36.6x less latency vs unoptimized STM32

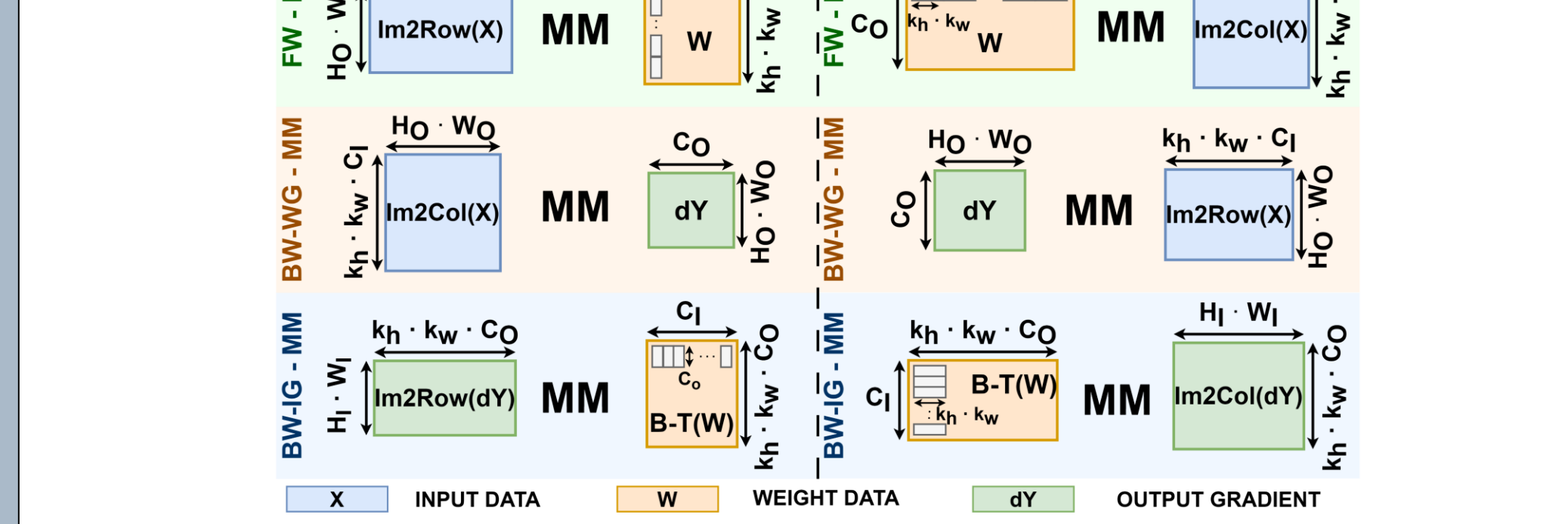


³PULP-TrainLib repository: <https://github.com/pulp-platform/pulp-trainlib>
⁴TinyMLPerf Benchmarks <https://github.com/mlcommons/tiny/tree/master/benchmark>

2) Reduced Precision Floating-Point Optimization for Deep Neural Network On-Device Learning on MicroControllers⁵

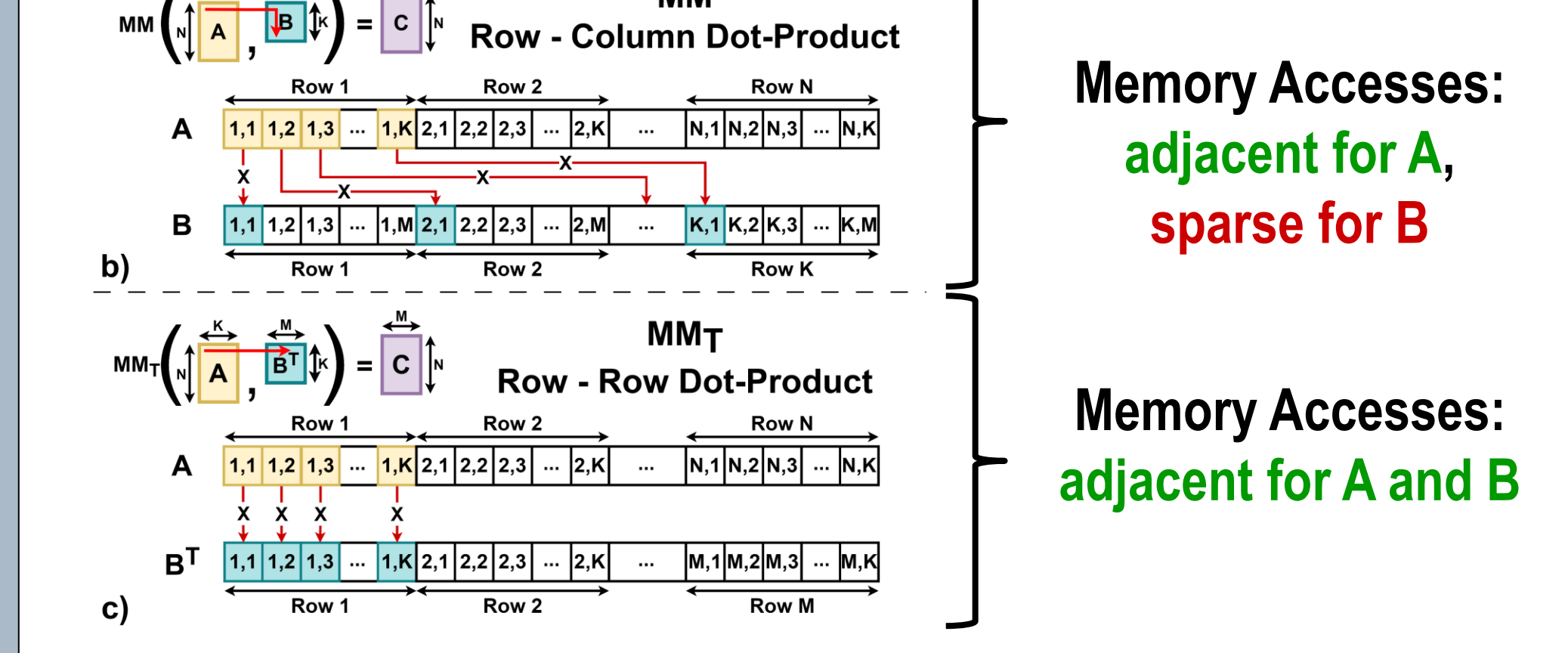
Daide Nadalini^{1,2}, Manuele Rusci³, Luca Benini^{1,4} and Francesco Conti¹

¹Università di Bologna, ²Politecnico di Torino, ³Katholieke Universiteit Leuven, ⁴ETH Zurich

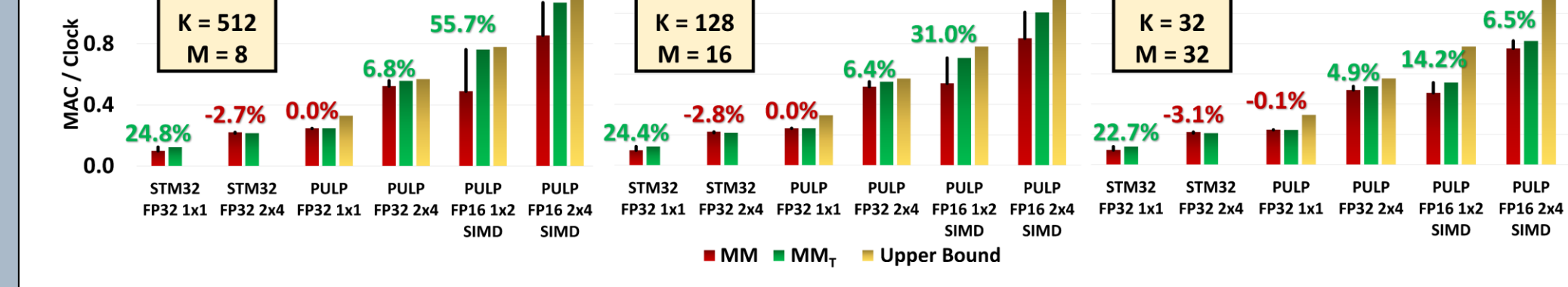


Main Idea: Row-Column MMs are inefficient with SIMD (Sparse load-stores)

1) Row-Column MM is replaced by MM_T (SIMD row-row MM)

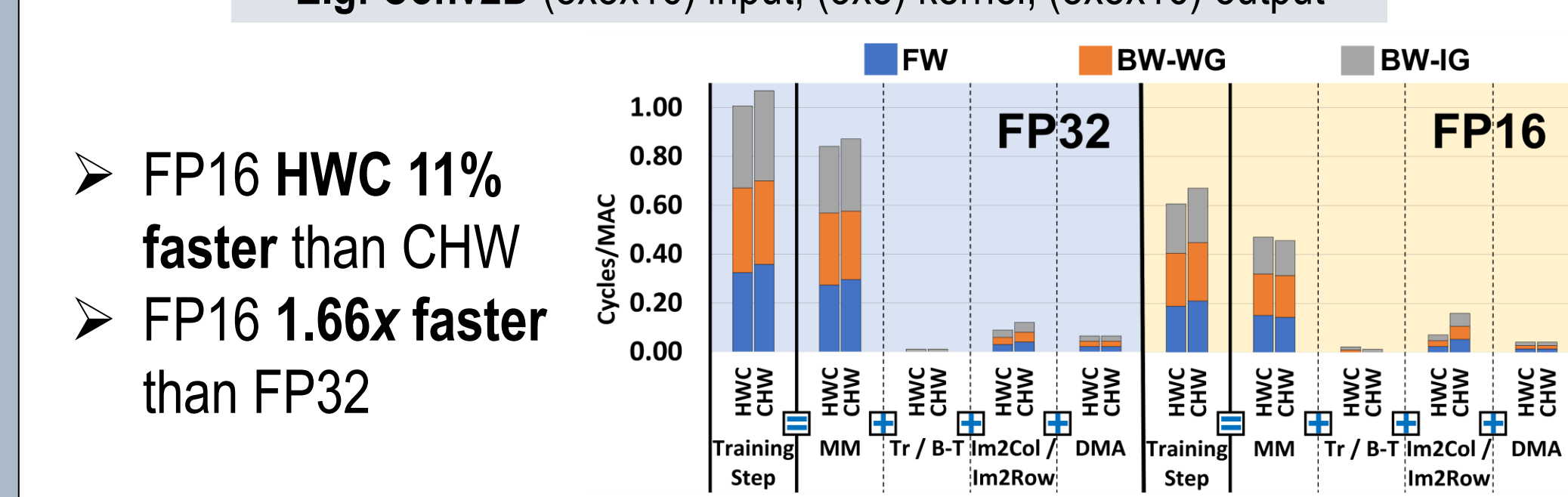
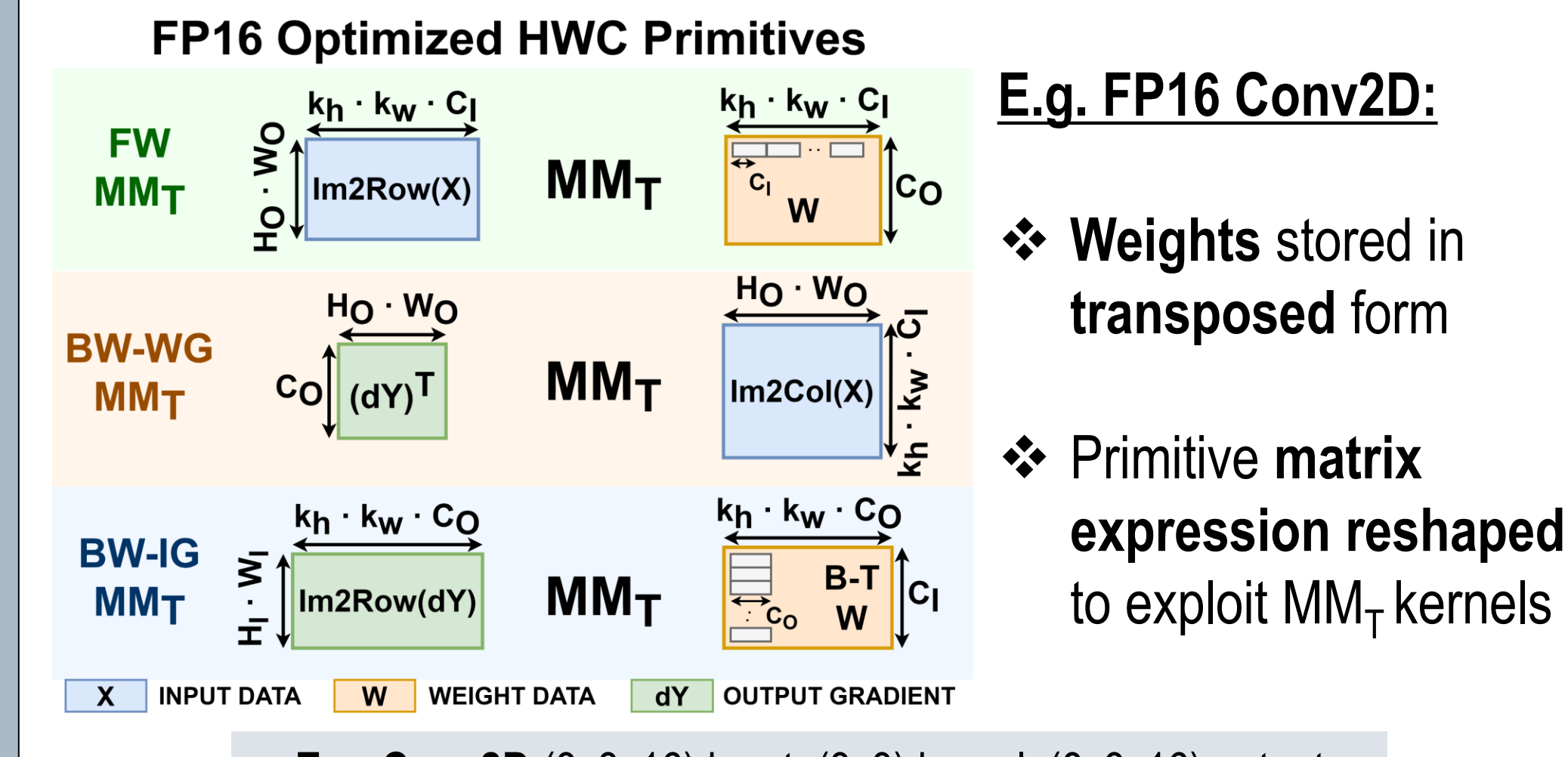


FP16 SIMD MM_T vs MM: loop unrolling on 1 Core



On 8 cores, up to 7.89 MAC/clock for MM_T (1.91x FP32 MM)

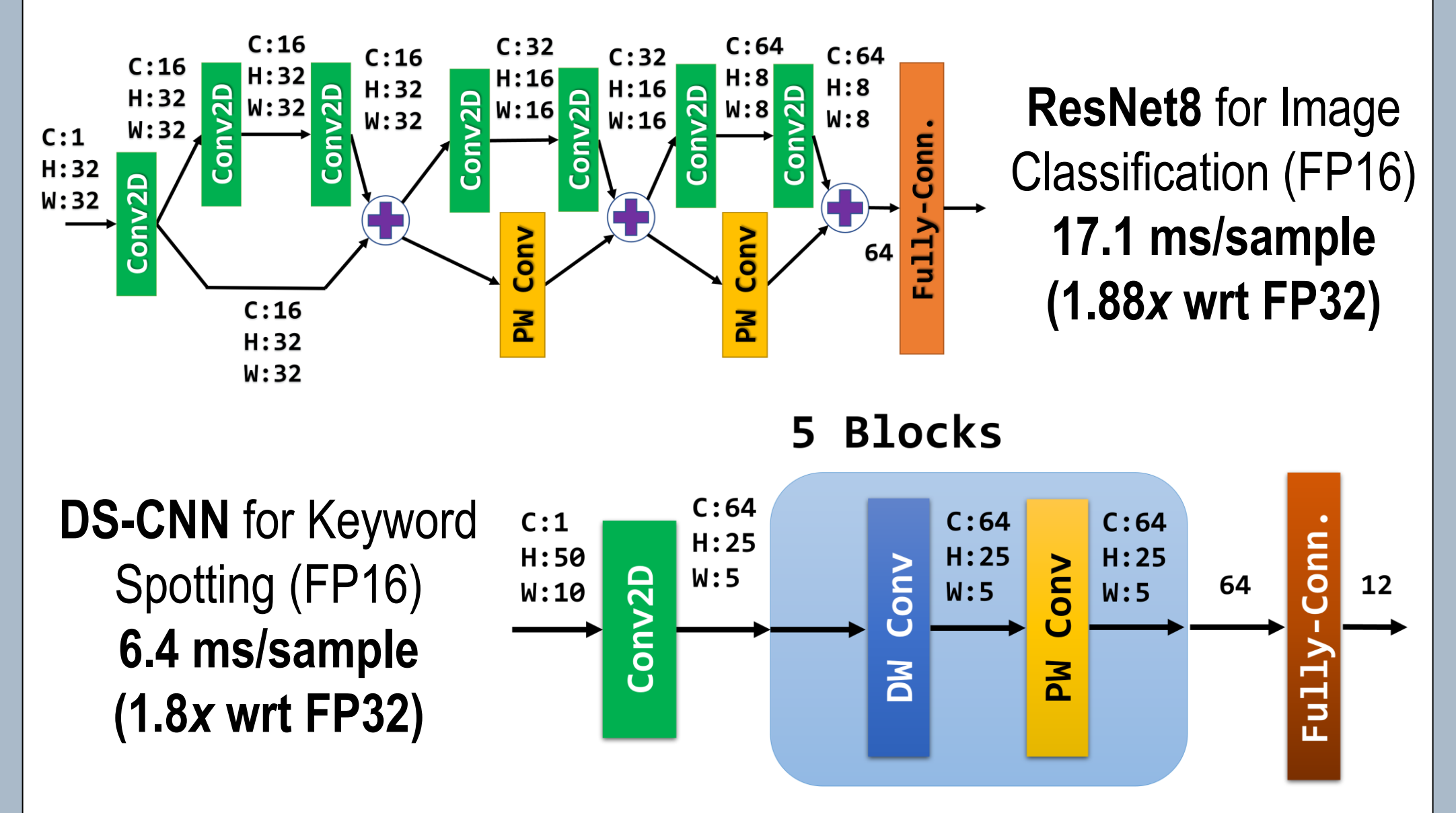
2) Primitive-level optimization of ODL training steps (BackPropagation-based)



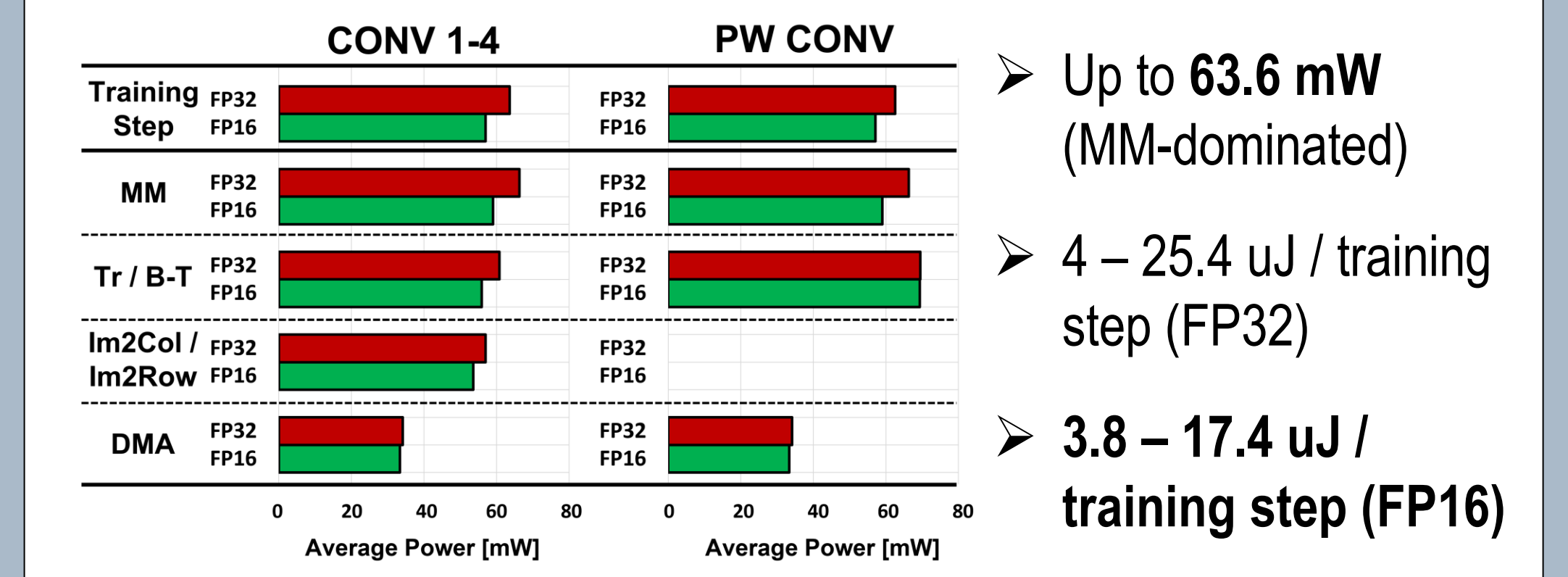
⁵D. Nadalini, M. Rusci, L. Benini, F. Conti, "Reduced Precision Floating-Point Optimization for Deep Neural Network On-Device Learning on MicroControllers"

3) FP16 optimization of full TinyML² models on a GreenWaves Technologies' GAP9 SoC @370MHz, 8 Cores

More accurate estimate: DMA transfers, Im2Col/Im2Row, real SoC (shared FPUs)



4) Power Analysis of ODL Primitives on the GAP9 SoC (FP32 vs FP16)



5) Real-time Continual Learning on MCU Targets (comparison with SoA)

Table 7: Latency and Energy Evaluation for Continual Learning on MCUs

	Platform	Latency (s)	Energy (J)
Aifetes[30]	STM32L4 @ 80 MHz	DW21	236655
		DW23	142157
		LIN27	102.4
PULP-TrainLib[31]	Greenwaves GAP9 @ 370 MHz	DW21	504.1
		DW23	303.8
		LIN27	0.89
This Work	Greenwaves GAP9 @ 370 MHz	DW21	308.5
		DW23	185.9
		LIN27	0.89

Continual Learning⁶ on a MobileNet: new class learned in 3-5 minutes with FP16 on GWT GAP9 (retraining the last 5-10 layers)

Conclusion

- Enabling ODL on MCUs allows real-time BackProp-based learning on IoT end nodes (Continual, Online, ..., Learning)
- Reduced Precision enables fast (1.66x FP32) backend for ODL with high enough precision
- Power consumption < 100 mW on Multi-Core RISC-V MCUs

⁶Latent Replay for Real-Time Continual Learning: <https://arxiv.org/abs/1912.01100>

