

ALMA MATER STUDIORUM Università di Bologna Towards Adaptive and Robust Tiny Machine Learning on Multi-Core Embedded RISC-V MCUs

> Manuele Rusci^{1,2} - <u>manuele.rusci@unibo.it</u> Francesco Conti¹ - <u>f.conti@unibo.it</u> Luca Benini^{1,3} - <u>luca.benini@unibo.it</u>

 ¹ Dipartimento di Ingegneria dell'Energia Elettrica e dell'Informazione (DEI) – Università di Bologna
 ² GreenWaves Technologies – Bologna
 ³ Integrated Systems Laboratory – ETH Zurich

Tiny Machine Learning for the IoT ecosystem



simple design server computational power

Network scalability Privacy of personal data TX Power consumption (100s mWs)

Miniaturized **low-power** and **low-cost** sensor nodes for cyber-physical systems



Tiny Machine Learning for the IoT ecosystem



Miniaturized **low-power** and **low-cost** sensor nodes for cyber-physical systems



Bringing Machine (Deep) Learning to Tiny Devices \rightarrow **Tiny Machine Learning**

Every **Smart Sensor** plays *data-to-information* conversion and transmits only low-bandwidth data to the server



TX Bandwidth and Power Reduction Need of Energy-Efficient TinyML devices, tools and design methods!

M. Rusci - IWES 10/12/2021

An (Open) HW perspective: the PULP platform

A **RISC-V based multi-core MCU-like platform**, with extended RISC-V ISA to efficiently support ML workloads, e.g. vectorized 8-bit MAC single-cycle instructions



PULP-NN: Accelerating DNN Inference on the PULP platform



- Peak Performance of **15.5 MAC /clk** on **8** cores
- □ **19.6x** and **30x** faster than CM7 and CM4 based MCUs using CMSIS-NN for CIFAR10 inference

Optimized GEMM on PULP by 1) maximizing **data reuse in** register file 2) exploiting SIMD MAC isntructions and 3) parallelism



Leveraging the PULP technology for real world applications

Automatic License Plate Detection and Recognition on GAP8



- Accuracy: 39% mAP for LP det. & > 99.13% for LP rec.
- Max recognition distance: 4m for detection and 2m for recognition
- □ **117mW** power envelope, **108 mJ** per inference.

73x less energy w.r.t. previous ALPR system.

FM388

<Shanahai>AFM388

#A FM388.

<Shanahai>AFM3883

<Zhejiang>AFM3883

<Hebei>PE3

<Beijing>E5V

Lossless compression of DNNs under severe memory constraints

On-chip memories are highly constrained (< 2MB) and off-chip memories present high costs (energy, \$).



Issue: 8-bit quantization (4x compression) is not sufficient to bring 'complex' models on MCUs

DNN Inference with Mixed-Low-Precision Quantization



- -2% wrt most accurate INT8 mobilenet (224_1.0)
- +8% wrt most accurate INT8 mobilenet fitting the memory (192_0.5)
- ➤ +2% wrt most accurate INT4 mobilenet (224_1.0)



Rule-based (Rusci, Mlsys 2020): rule-based bit reduction policy (from 8-bit to 2-bit) up to <u>fit</u> the memory constraints

Reinforcement-Learning (Rusci, ITEM 2020): same accuracy levels – the RL agent maximizes memory utilization

	Model	#Params (M)	FP32 Top1 Accuracy	Mixed Prec. Top1 Accuracy	Mixed Prec. Weight Compression	Mixed Prec. Activ. (average) Compression
	MobileNetV1 (224_1.0)	4.2	70.6%	68.4%	9x	5.1x
	MobileNetV1 (224_0.75)	2.6	68.4%	68.0%	5.4x	4.6x
	MobileNetV2 (224_1.0)	3.4	72.04%	55.1%	7.6x	5.1x
	MobileNetV2 (192_1.0)	3.4	70.7%	58.0%	7.4x	5.2x
	MobileNetV2 $(160_1.0)$	3.4	68.8%	67.5%	7.4x	4.4x
	MNasNet-A1 (224_1.0)	3.9	75.2%	62.6%	9.4x	4.7x
	MNasNet-A1 (224_0_35)	1.7	64.1%	64.1%*	4x (8-bit)	4x (8-bit)
* optimistic score: we assume a lossless 8-bit compression w.r.t. FP32						

(Rusci, Mlsys, 2020) Rusci et al, "Memory-driven mixed low precision quantization for enabling deep network inference on microcontrollers." *Mlsys* (2020).

(Rusci, ITEM 2020) Rusci, M. et al. "Leveraging Automated Mixed-Low-Precision Quantization for tiny edge microcontrollers." ITEM 2020.

Mixed-Precision Inference for MCUs

Low-bitwidth multi-precision is not natively supported by RISC-V ISA.

SW-emulation of 2-4-bit Mixed precision MAC adds overhead mixed-precision inference results slower than 8-bit inference

The **Mixed-Precision Inference Core (MPIC)** (Ottavi, 2020) added support for 4-, 2-bits and mixed precision operations from 16 down to 2 bits ➢ Novel SIMD DOT-Product unit





• Up to 11.7x better performance compared to Cortex M4

Demand for end-to-end tools...

But real-world is challenging!



- MCU's constraints prevent the deployment of robust "Big" models w/ high generalization capacity
- Training data differs from *real-world* sensor data
- More use-case specific data is always needed (e.g. custom keyword, different voice accent)
 New dataset collection and model retraining (from scratch)
- Takes 10's months for application development Not acceptable!

Class: "bottle"





Low-Power Camera

Google Open Images

Presently, the major issue towards the diffusion of smart sensor nodes relates to the lack of adaptation over time.

Continual Learning as the next paradigm for smart sensors: Adaptive TinyML

Continual Learning to train incrementally DL models **on device** starting from initial knowledge

• Latent replays to deal with catastrophic forgetting [Pellegrini 2019]



□ Quantize the latent replays (and the Frozen stage) to gain >4x compression at 1-3% accuracy loss (8-bit LR)

LR layer	FP32 baseline	FP32 + UINT-8	UINT-8 + UINT-8	FP32 + UINT-7	UINT-8 + UINT-7
27	72.7 ± 0.34	70.1 ± 0.54	69.2 ± 0.48	68.0 ± 0.63	67.8 ± 1.14
25	73.3 ± 0.58	70.9 ± 0.65	70.2 ± 0.67	66.2 ± 0.75	66.1 ± 0.94
23	75.0 ± 0.83	73.2 \pm 0.46	73.4 ± 0.66	71.1 \pm 0.63	69.9 ± 1.25
21	76.5 ± 0.63	74.9 ± 0.51	73.9 ± 1.67	72.7 ± 0.74	72.6 ± 1.30
19	77.7 ± 0.73	76.5 ± 0.48	76.0 ± 0.80	74.0 ± 0.57	75.2 ± 1.10

Accuracy on Core50 with 1500 MRs

- Running FP32 Backpropagation on a PULPbased system, up to 0.62 FMAC/clk with 8 cores and 512kB L1 memory
- □ **65x** faster than CM4-based device thanks to the higher frequency (22nm tech), 7.2x parallel speed-up and optimized ISA micro-architecture.

Ravaglia et al. "A TinyML Platform for On-Device Continual Learning with Quantized Latent Replays." JETCAS 2021

Conclusion

Robust Deep Learning on smart tiny device is getting real, thanks to a combination of factors:

- Energy-efficient but flexible HW
- SW stack to efficiently exploit the underlying HW
- DL-HW codesign to enable mixed-low precision on low-power processors
 - More tools are needed

Moving first steps into Adaptive TinyML:

- Proof of concept of On-Device Continual Learning on MCU-like platform
- More to be done... (application side, algorithm, SW, tools..)



ALMA MATER STUDIORUM Università di Bologna

Dr. Manuele Rusci

DEI – Università di Bologna

manuele.rusci@unibo.it

www.unibo.it