

**PULP PLATFORM** Open Source Hardware, the way it should be!

### Many shades of acceleration

An Open TinyML Platform Perspective

#### Luca Benini <lbenini@iis.ee.ethz.ch,luca.Benini@unibo.it>



Horizon 2020 European Union funding for Research & Innovation



Fonds national suisse Schweizerischer Nationalfonds Fondo nazionale svizzero Swiss National Science Foundation







Prof. of Digital Circuit and Systems @ ETHZ and UNIBO. h-index=110, 53'000+ citations, 1'000+ publications, fellow IEEE, ACM, Chief Architect in STMicroelectronics (2009-2012) Group of 100+ people



European Research Council





### CloudML → TinyML

#### **TinyML Opportunity**

Edge AI chips by device, 2020 and 2024 (millions of units) Speaker Wearable Smartphone Tablet Enterprise edge 2020 500 75 75 50 2024 1,000 100 150 100 250 200 400 600 800 1,000 1,200 1,400 1,600 1,800 0

Sources: MarketsandMarkets, Edge AI hardware market by device (smartphones, cameras, robots, automobiles, smart speakers, wearables, and smart mirrors), processor (CPU, GPU, ASIC, and others), power consumption, process, end user industry, and region—global forecast to 2024, April 4, 2019; Deloitte analysis.

Deloitte Insights | deloitte.com/insights

#### TinyML challenge Al capabilities in the power envelope of an MCU: **10-mW peak (1mW avg)**

### Al Workloads - DNNs

H Pham 2021(Google) arXiv:2003.10580v3



- 3

## Energy efficiency @ GOPS is the Challenge

ARM Cortex-M MCUs: M0+, M4, M7 (40LP, typ, 1.1V)\*



ETH ZUrich

### RI5CY – An Open MCU-class RISC-V Core for EE-AI

[Gautschi et al. TVLSI 2017]



### **PULP-NN: Xpulp ISA exploitation**



P↑ T↓↓↓ so, E=P\*T↓↓ Nice! But what about the GOPS? Faster+Superscalar is not efficient!

M7: 5.01 CoreMark/MHz-58.5 μW/MHz M4: 3.42 CoreMark/MHz-12.26 μW/MHz





### ML & Parallel, Near-threshold: a Marriage Made in Heaven

- As VDD decreases. operating speed decreases
- However efficiency increases  $\rightarrow$  more work done per Joule
- Until leakage effects start to dominate ürich
  - Put more units in parallel to get performance up and keep them busy with a parallel workload

ML is massively parallel and scales well (P/S  $\uparrow$  with NN size)





### Multiple RI5CY Cores (1-16)

RISC-V core RISC-V core RISC-V core RISC-V

**CLUSTER** 



### **Low-Latency Shared TCDM**



ETH zürich

### DMA for data transfers from/to L2



ETHZÜrich

### Shared instruction cache with private "loop buffer"



ETHZürich

### Results: RV32IMCXpulp vs RV32IMC

- 8-bit convolution
  - Open source DNN library
- 10x through xPULP
  - Extensions bring real speedup
- Near-linear speedup
  - Scales well for regular workloads
- 75x overall gain
  - Sub-byte: x2-4x better
  - Mixed precision supported (more later)



[Garofalo et al. arxiv.org/abs/1908.11263]

### An additional I/O controller is used for IO





[Burrello et al. arxiv.org/abs/2008.07127]

QuantLab **Quantization Laboratory** 

**NEMO NE**ural **M**inimization for pyt**O**rch

DORY Deployment Oriented to memoRY

**PULP-NN PULP Neural Network backend** 



### What's next? Sub-pJ/OP Accelerators



### **Tightly-coupled HW Compute Engine**



### Hardware Processing Engines (HWPEs)



#### HWPE efficiency vs. optimized RISC-V core

rich

- 1. Specialized datapath (e.g. systolic MAC) & internal storage (e.g. linebuffer, accum-regs)
- 2. Dedicated control (no I-fetch) with shadow registers (overlapped config-exec)
- 3. Specialized high-BW interco into L1 (on data-plane)



# $\mathbf{y}(k_{out}) = \text{binarize}_{\pm 1} \left( \mathbf{b}_{k_{out}} + \sum_{k_{in}} \left( \mathbf{W}(k_{out}, k_{in}) \otimes \mathbf{x}(k_{in}) \right) \right)$

$$\text{binarize}_{\pm 1}(t) = \text{sign}\left(\gamma \frac{t-\mu}{\sigma} + \beta\right)$$

$$\text{binarize}_{0,1}(t) = \begin{cases} 1 \text{ if } t \ge -\kappa/\lambda \doteq \tau, \text{ else } 0 & (\text{when } \lambda > 0) \\ 1 \text{ if } t \le -\kappa/\lambda \doteq \tau, \text{ else } 0 & (\text{when } \lambda < 0) \end{cases}$$

Extreme Quantization 
→ Binarization with XNOR nets

Sin	ary	product $\rightarrow$				XOR	
Α	В	out		A	В	out	
-1	-1	+1		0	0	1	
-1	+1	-1		0	1	0	
+1	-1	-1		1	0	0	
+1	+1	+1		1	1	1	

$$\mathbf{y}(k_{out}) = \text{binarize}_{0,1} \left( \sum_{k_{in}} \left( \mathbf{W}(k_{out}, k_{in}) \otimes \mathbf{x}(k_{in}) \right) \right)$$
  
hresholding  
Multi-bit accumulation

ETH ZUrich

### **XNE: XNOR Neural Engine**

[Di Mauro et al. TCAS I, 2020]





BINCONV: Binary dot-product and thresholding logic array

### **XNE Energy Efficiency**



L1 SCM, L2 high-density, low leakgage SRAM (activations), MRAM (weights)

#### But... Accuracy Loss is high even with retraining (10%+) Need flexible precision tuning!





- Many  $M \times N$  bits products...
- ... but one  $M \times N$  product is the superposition of  $M \times N$  1-bit products!

$$\mathbf{y}(k_{out}) = quant\left(\sum_{i=0..N}\sum_{k_{in}}\sum_{2^{i}2^{j}}\left(\mathbf{W}_{bin}(k_{out},k_{in})\otimes\mathbf{x}_{bin}(k_{in})\right)\right)$$
  
Q-bit output fmaps  
1-bit weights

One quantized NN can be emulated by superposition of power-of-2 weighted  $M \times N$  binary NN







### Mixed-Precision Quantized Networks – CMIX-NN

[Capotondi et al. TCAS II, 2020]

Apply minimum tensor-wise quantization to <u>fit</u> **memory constraints** with very-low accuracy drop



Only -2% wrt most accurate INT8 mobilenetV1 (224\_1.0) which does not fit on-chip

+8% wrt most accurate INT8 mobilenetV1 fitting on-chip (192\_0.5)

+7.5% wrt most accurate INT4 mobilenetV1 (224\_1.0) fitting on chip



# HW acceleration in perspective

Using 22FDX tech, NT@0.6V, High utilization, minimal IO & overhead

Energy-Efficient RV Core → 20pJ (8bit)

ISA-based 10-20x **→1-2pJ (8bit)** 



Configurable DP 10-20x  $\rightarrow$  50-100fJ (4bit)

Fully specialized DP 10-20x →5-10fJ (ternary)

LO 1088

Hzürich

\*See M. Scherer presentation, tinyML21 – sub 1fJ in 7nm

XPULPV2 &V3

HWCE, RBE, NE

**XNE, CUTIE\*** 

### Towards In-Sensor: Achieving sub-mW average power?

1mW average power with 10mW active power (10GOPS @ 1pJ/OP) → sub mW sleep



Duty cycling not acceptable when input events are asynchronous → watchful Sleep

Log(P)

Detect&Compress→1-10mW

Watchful sleep  $\rightarrow$  <1mW

Stream→ 100mW

### Need µW-range always-on Intelligence



ETH zürich

### **HD-Based smart Wake-Up Module**





### **HD-Based smart Wake-Up Module**



Y KANAG

### **HD-Based smart Wake-Up Module**



Y KINK

### Not Only CNNs: Hyper-Dimensional Computing



### **In-memory Hyperdimensional Computing**



ETHZürich

TER STOSLORUA

### HD-Based smart Wake-Up Module - Hypnos

[Eggiman et al. arxiv.org/abs/2102.02758]

github.com/pulp-platform/hypnos							
		Design (post P&R)					
		Tech	nology GF22		UHT		n
		Area		670kG	E		lo <sup>.</sup> lik
5		Max. Frequency		3 MHz			
ZULIC	f <sub>clk</sub>		32kHz		200kHz		
	max. sampling rate		150 SPS/Channel		1kSPS/Channel		
	P <sub>SWU, dynamic</sub>		0.99uW		6.21uW		
P <sub>SWU, leakage</sub>			0.7uW		0.7uW		
	P <sub>SPI, dynamic</sub>		1.28uW		8.00uW		
ER STU	P <sub>SWU, total</sub> Measure	<sub>WU, total</sub> Measured			14.9uW		

nplemented with west leakage cell orary (UHVT)

### All together in VEGA: Extreme Edge IoT Processor

[Rossi et al. ISSCC21]

- RISC-V cluster (8cores +1) 614GOPS/W @ 7.6GOPS (8bit DNNs), 79GFLOPS/W @ 1GFLOP (32bit FP appl)
- Multi-precision HWCE(4b/8b/16b) 3×3×3 MACs with normalization / activation: 32.2GOPS and 1.3TOPS/W (8bit)
  - 1.7 µW cognitive unit for autonomous wake-up from retentive sleep mode



### All together in VEGA: Extreme Edge IoT Processor

- RISC-V cluster (8cores +1) 614GOPS/W @ 7.6GOPS (8bit DNNs), 79GFLOPS/W @ 1GFLOP (32bit FP appl)
- Multi-precision HWCE(4b/8b/16b) 3×3×3 MACs with normalization / activation: 32.2GOPS and 1.3TOPS/W (8bit)
- 1.7 µW cognitive unit for autonomous wake-up from retentive sleep mode
- Fully-on chip DNN inference with 4MB MRAM



	Technology	22nm FDSOI		
	Chip Area	12mm <sup>2</sup>		
	SRAM	1.7 MB		
	MRAM	4 MB		
	VDD range	0.5V - 0.8V		
	VBB range	0V - 1.1V		
	Fr. Range	32 kHz - 450 MHz		
	Pow. Range	1.7 µW - 49.4 mW		



### Full DNN Energy (MobileNetV2)



### When you count mWatts, everything matters!

What about IO power? (Mem, Sensor)

- SPIs
  - I/O VDD=1.8V
  - fspi-max=50MHz,
  - Assuming duty-cycled operation @ various bandwidths
- ULP serial link (duty-cycled)
  - 10.2x less energy and 15.7x higher maximum BW compared to single SPI
  - 2.56x higher efficiency than the DDR Octal SPI @787Mbps
  - 5 → 3pJ/bit
  - However it's still 2mW@ 500Mbps
  - 3D integration: 0.15pJ/bit and below



[Okuhara et al. ISCAS20]





### Closing thoughts – Open Platform for TinyML

#### **PULP is an Open Platform**

- For science ... fundamental "research infrastructure" Reduce "getting up to speed" overhead for partners Enables fair and well controlled benchmarking
- For Business ... it is truly disruptive Reduces the NRE, faster innovation path for startups New business models (for profit and non-for profit) Exemplary collaboration with GF (Quentin, Vega...)

#### **Heterogeneous & Flexible**

- 1-3 orders of magnitude improvement (wrt to efficient RV) by acceleration
   ISA → Configurable → Fully customized + heterogeneous architectural combinations
- Focus on IO energy (memory, sensor) to achieve sub pJ/OP @ full platform 3D-IC technology is a key enabler



# Parallel Ultra Low Power

Luca Benini, Davide Rossi, Andrea Borghesi, Michele Magno, Simone Benatti, Francesco Conti, Francesco Beneventi, Daniele Palossi, Giuseppe Tagliavini, Antonio Pullini, Germain Haugou, Lukas Cavigelli, Manuele Rusci, Florian Glaser, Renzo Andri, Fabio Montagna, Bjoern Forsberg, Pasquale Davide Schiavone, Alfio Di Mauro, Victor Javier Kartsch Morinigo, Tommaso Polonelli, Fabian Schuiki, Stefan Mach, Andreas Kurth, Florian Zaruba, Manuel Eggimann, Philipp Mayer, Marco Guermandi, Xiaying Wang, Michael Hersche, Robert Balas, Antonio Mastrandrea, Matheus Cavalcante, Angelo Garofalo, Alessio Burrello, Gianna Paulin, Georg Rutishauser, Andrea Cossettini, Luca Bertaccini, Maxim Mattheeuws, Samuel Riedel, Sergei Vostrikov, Vlad Niculescu, Frank K. Gurkaynak, and many more that we forgot to mention



http://pulp-platform.org



@pulp\_platform