

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

Open, Parallel Ultra-Low Power Platforms for Extreme Edge Al

Luca Benini < Ibenini@iis.ee.ethz.ch>



European Commission Horizon 2020 European Union funding for Research & Innovation



FONDS NATIONAL SUISSE Schweizerischer Nationalfonds Fondo nazionale svizzero Swiss National Science Foundation







TH zürich

Cloud \rightarrow Edge \rightarrow Extreme Edge Al a.k.a. TinyML



E. Gousev, Qcomm research

#1 Customer Question on Amazon.com (out of 1,000+):

 I don't want any of my (private, personal) videos on any servers not in my control. Is this possible?

#2 Customer Question on Amazon.com (out of 1,000+):

water sense groups constraints and have from Asia Partial Works (A.

server an extension of the formation of the Witchel Watchel

2. How long does the battery charge last?

Extreme edge Al challenge Al capabilities in the power envelope of an MCU:

100mW peak (1mW avg)

9 March 2020

Te

Al Workloads from Cloud to Edge (Extreme?)







ETH zürich



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

3-cycle ALU-OP, 4-cyle MEM-OP→IPC loss: LD-use, Branch





9 March 2020

XPULP extensions: 25 kGE \rightarrow 40 kGE (1.6x)

PULP-NN: Xpulp ISA exploitation



Faster+Superscalar is not efficient!

Nice – But what about the GOPS? → M7: 5.01 CoreMark/MHz-58.5 µW/MHz M4: 3.42 CoreMark/MHz-12.26 µW/MHz

En zürich

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

- As VDD decreases, operating speed decreases
- However efficiency increases → more work done per Joule
- Until leakage effects start to dominate
- 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 个 with NN size)



Efficiency vs VDD chip01

ETH zürich



Multiple RI5CY Cores (1-16)







CLUSTER

9 March 2020

8



Low-Latency Shared TCDM



E H zürich



ç

DMA for data transfers from/to L2



ETH zürich

Shared instruction cache with private "loop buffer"



ETH zürich

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



n zürich 11

An additional I/O controller is used for IO



* 1:



EnHzürich



Nice, but what exactly is "open" in Open Source HW?

- Only the first stage of the silicon production pipeline can be open HW
 → RTL source code (in an HDL such as SystemVerilog)
- Later stages contain closed IP of various actors + tool licensing issues



PULP includes Cores+Interco+IO → and Open Platform



0 SIC **a**//a







PULP is silicon proven





























P S

Successful product development: GWT's GAP8

Two independent clock and voltage domains, from 0-133MHz/1V up to 0-250MHz/1.2V



What	Freq MHz	Exec Time ms 99.1		Cycles	Power mW
40nm Dual Issue MCU	216			21 400 000	60
GAP8 @1.0V	15.4	99.1	11 X	1 500 000	3.7
GAP8 @1.2V	17.5 17	8.7		1 500 000	70
GAP8 @1.0V w HWCE	4.7	99.1		460 000	0.8





New Application Frontiers: DroNET on NanoDrone



Pluggable PCB: PULP-Shield

- ~5g, 30x28mm
- GAP8 SoC
- 8 MB HDRAM
- 16 MB HFlash
- QVGA ULP HiMax camera
- Crazyflie 2.0
 nano-drone (27g)





Only onboard computation for autonomous flight + obstacle avoidance no human operator, no ad-hoc external signals, and no remote base-station!





- **.

Outlook: PULP-based nano UAV family



- FrontNET [1]: hovering in front of a freely moving user
- End-to-end, RESnet-like CNN
- Runs with DroneNET
- AI+control: 5-10% of tot. Power
- Collaboration with IDSIA (USI-SUPSI) Lugano, Switzerland



9 March 2020

[1] Mantegazza, Dario, et al. "Vision-based Control of a Quadrotor in User Proximity: Mediated vs End-to-End Learning Approaches." 2019 International Conference on Robotics and Automation (ICRA). IEEE, 2019.

n zürich

What's next? Sub-pJ/OP Accelerators, but flexibility needed!



asimovinstitute.org/neural-network-zoo

Hardware Processing Engines (HWPEs)



ETH zürich

2020

Hzürich

Sub-pJ/W Accelerator; Tightly-coupled HW Compute Engine



More Efficiency (2): Extreme Quantization

Model	Bit-width	Top-1 error	SOA INQ retraining
ResNet-18 ref	32	31.73%	
INQ	5	31.02%	
INQ	4	31.11%	
INQ	3	31.92%	
INQ	2 (ternary)	33.98%	2.2% loss \rightarrow 0% with 20% larger net

Low(er) precision: $8 \rightarrow 4 \rightarrow 2$



1 MAC Op = 2 Op (1 Op for the "sign-reverse", 1 Op for the add).





Thresholding

From +/-1 Binarization to XNORs

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

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

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

$$\mathbf{y}(k_{out}) = \operatorname{binarize}_{0,1} \left(\sum_{k_{in}} \left(\mathbf{W}(k_{out}, k_{in}) \otimes \mathbf{x}(k_{in}) \right) \right)$$

$$\mathbf{y}(k_{out}) = \operatorname{binarize}_{0,1} \left(\sum_{k_{in}} \left(\mathbf{W}(k_{out}, k_{in}) \otimes \mathbf{x}(k_{in}) \right) \right)$$

XOR

out

0 0 1

9 March 2020

Multi-bit accumulation

ETH zürich



XNE: XNOR Neural Engine



- Hzürich

BINCONV: Binary dot-product and thresholding logic array

XNE Energy Efficiency



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!

$$\mathbf{y}(k_{out}) = quant\left(\sum_{i=0..M}\sum_{j=0..N}\sum_{k_{in}}2^{i}2^{j}\left(\mathbf{W}_{bin}(k_{out},k_{in})\otimes\mathbf{x}_{bin}(k_{in})\right)\right)$$

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

Reconfigurable Binary Engine

$$\mathbf{y}(k_{out}) = quant\left(\sum_{i=0..M}\sum_{j=0..N}\sum_{k_{in}} 2^{i}2^{j}\left(\mathbf{W}_{bin}(k_{out},k_{in}) \otimes \mathbf{x}_{bin}(k_{in})\right)\right) \quad \text{e.g. a 3x3 conv}$$
with **N**=4 bits

9 March 2020

30



Entzürich

Reconfigurable Binary Engine

$$\mathbf{y}(k_{out}) = quant\left(\sum_{i=0..M}\sum_{j=0..N}\sum_{k_{in}} 2^{i} 2^{j} \left(\mathbf{W}_{bin}(k_{out},k_{in}) \otimes \mathbf{x}_{bin}(k_{in})\right)\right)$$



En zürich

 preload fmap buffer with input activations (up to 5x5 pixels, 32 channels, N=4 bits/pixel)



31

to/from L1 TCDM

Enh zürich

Reconfigurable Binary Engine

$$\mathbf{y}(k_{out}) = quant\left(\sum_{i=0..M}\sum_{j=0..N}\sum_{k_{in}} 2^{i}2^{j} \left(\mathbf{W}_{bin}(k_{out},k_{in}) \otimes \mathbf{x}_{bin}(k_{in})\right)\right)$$



- preload fmap buffer with input activations (up to 5x5 pixels, 32 channels, 4 bits/pixel)
- 2. each of the 81 blocks has as input one of the 25 px (32x4 bits)



ETH zürich

9 March 2020

3



Reconfigurable Binary Engine

 $\mathbf{y}(k_{out}) = \mathbf{quant}\left(\sum_{i=0..M}\sum_{j=0..N}\sum_{k_{in}} 2^{i} 2^{j} (\mathbf{W}_{\mathbf{bin}}(k_{out},k_{in}) \otimes \mathbf{x}_{\mathbf{bin}}(k_{in}))\right)$



to/from L1 TCDM



5. results are added columnwise and accumulated in 32-bit accumulator banks (9x32x32 bits) for many iterations until full accumulation

Reconfigurable Binary Engine

$$\mathbf{y}(k_{out}) = quant\left(\sum_{i=0..M}\sum_{j=0..N}\sum_{k_{in}}2^{i}2^{j}\left(\mathbf{W}_{bin}(k_{out},k_{in})\otimes\mathbf{x}_{bin}(k_{in})\right)\right)$$



- 5. results are added columnwise and accumulated in 32-bit accumulator banks (9x32x32 bits) for many iterations until full accumulation
- 6. after full accumulation, the accumulator values are **quantized** and **streamed out**

What about µW «sleep»?

Small always-on network \rightarrow

Triggers alarm and video capture/streaming for cloud-based forensics







Need µW-range always-on Intelligence



Not Only CNNs: Hyper-Dimensional Computing



Highly parallel, fault-tolerand binary operators, assoc-min-distance search

Merge storage & computation i.e. **In-memory computing**

9 March 2020

ETHZürich







ETH zürich

















ETH zürich



Entzürich









Entzürich



Enhzürich









Hzürich

Integrating Accelerators into PULP – The big picture



Academic open source \rightarrow Industrial open source



Rick O'Connor (OpenHW CEO, former RISC-V foundation director)

- OpenHW Group is a not-for-profit, global organization (EU,NA,Asia) driven by its members and individual contributors where HW and SW designers collaborate in the development of open-source cores, related IP, tools and SW such as the Core-V family of cores.
- OpenHW Group provides an infrastructure for hosting high quality open-source HW developments in line with industry best practices.



A vertical, application-focused open approach



- OpenTitan is the first open source silicon project building a transparent, high-quality reference design for silicon root of trust (RoT) chips.
- Founding Partners

ETH zürich



Google



Hzürich



nuvotor

Western Digital.



Feel the momentum!

Ibex RISC-V core, flash interface, communications ports, cryptography accelerators, and more.

Vibrant repository







EHZürich

HPC Vertical: The European Processor Initiative



Europe Needs its own Processors

- Processors now control almost every aspect of our lives
- Security (back doors etc.)
- Possible future restrictions on exports to EU due to increasing protectionism
- A competitive EU supply chain for HPC technologies will create jobs and growth in Europe
- Sovereignty (data, economical, embargo)



- High Performance General Purpose Processor for HPC
- High-performance RISC-V based accelerator
- Computing platform for autonomous cars
- Will also target the AI, Big Data and other markets in order to be economically sustainable

Closing thoughts... the open HW revolution

For science ... fundamental "research infrastructure" Community building: sharing of ideas, artefacts Fantastic tool for dissemination (more citations ③)! Reduce "getting up to speed" overhead for partners Enables fair and well controlled benchmarking For Business ... it is truly disruptive Reduces the NRE cost for silicon design Faster innovation path for startups New business models (for profit and non-for profit) Helps exchange of information across NDA walls Great for Marketing & Training **For society** ...long term sustained benefits More innovation, growth, jobs Personalized silicon vision "Moore-for-all" More Secure, safe, auditable HW



Posh Open Source Hardware (POSH): An open source System on Chip (SoC) design and verification ecosystem that enables cost effective design of ultra-complex SoCs

USA's electronics resurgence initiative



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

