

Plenty of Room at the Bottom?

Micropower Deep Learning for Cognitive Cyberphysical Systems

Luca Benini

http://www.pulp-platform.org











First, it was machine vision...

Now it's everywhere!





ETH Eidgenössische Technische Hochschule Zürich

Deep neural networks (DNNs) Swiss Federal Institute of Technology Zurich





Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

Key operation is dense M x V









CNN Computation: main kernel (per layer)





[®]GPUs are Great for Vanilla CNNs

Why?

Because they are good at matrix multiply \rightarrow 90% utilization is achievable (on lots of "cores")



Eidgenössische Technische Hochschule Züric

ETH

Biviss Federal Institute of Technology Zuth HW for deep Networks: Frenzy





Datacenter \rightarrow High-performance embedded \rightarrow mobile



CNN Workloads





Better networks are not necessarily more complex

DNNs Are Evolving Rapidly



Orders of magnitude compute effort an memory reduction with no loss in accuracy



What's Next?

Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

CPS Hierarchical Processing





Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

CPS Hierarchical Processing







Does it Matter?





CNN on MCUs?







Outline



Near Threshold Multiprocessing

- Non-Von Neumann Accelerators
- Aggressive Approximation
- From Frame-based to Event-based Processing
- Outlook and Conclusion







Near-Threshold Computing (NTC):

- **1.** Don't waste energy pushing devices in strong inversion
- 2. Recover performance with parallel execution

Swiss Federal Institute of Technolog Williprocessing





ULP (NT) Bottleneck: Memory

- "Standard" 6T SRAMs:
 - High VDDMIN
 - Bottleneck for energy efficiency
- Near-Threshold SRAMs (8T...)
 - Lower VDDMIN
 - Area/timing overhead (25%-50%)
 - High active energy
 - Low technology portability
- Standard Cell Memories (SCMs):
 - Wide supply voltage range
 - Lower read/write energy (2x 4x)
 - Easy technology portability
 - Major area overhead (4x)





ULP Memory Hierarchy





Near threshold FDSOI technology



Body bias: Highly effective knob for power & variability management!



Selective, fine-grained BB

- The cluster is partitioned in separate clock gating and body bias regions
- Body bias multiplexers (BBMUXes) control the well voltages of each region
- A Power Management Unit (PMU) automatically manages transitions between the operating modes
- Power modes of each region:
 - Boost mode: active + FBB
 - Normal mode: active + NO BB
 - Idle mode: clock gated + NO BB (in LVT) RBB (in RVT)



Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

The Evolution of the 'Species'



	PULPv1	PULPv2	PULPv3
# of cores	4	4	4
L2 memory	16 kB	64 kB	64 kB
TCDM	16kB SRAM	32kB SRAM 8kB SCM	32kB SRAM 16kB SCM
DVFS	no	yes	yes
I\$	4kB SRAM private	4kB SCM private	4kB SCM shared
DSP Extensions	no	no	yes
HW Synchronizer	no	no	yes
	PULPv1	PULPv2	PULPv3
Status	silicon proven	silicon proven	silicon proven
Technology	FD-SOI 28nm	FD-SOI 28nm flip-	FD-SOI 28nm
	conventional-well	well	conventional-well
Voltage range	0.45V - 1.2V	0.32V - 1.2V	0.4V - 0.7V
BB range	-1.8V - 0.9V	0.0V - 1.8V	-1.8V - 0.5V
Max freq.	475 MHz	850 MHz	200 MHz
Max perf.	1.9 GOPS	3.3 GOPS	1.8 GOPS
Peak en. eff.	60 GOPS/W	193 GOPS/W	385 GOPS/W

Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

The Evolution of the 'Species'







Extending RISC-V for CNNs

<32-bit precision -> SIMD2/4 opportunity

- 1. HW loops and Post modified LD/ST
- 2. Bit manipulations
- 3. Packed-SIMD ALU operations with dot product
- 4. Rounding and Normalizazion
- 5. Shuffle operations for vectors
 - V1 Baseline RISC-V RV32IMC HW loops
 - V2 Post modified Load/Store Mac
 - SIMD 2/4 + DotProduct + Shuffling
 - V3 Bit manipulation unit Lightweight fixed point

Small Power and Area overhead



Dot Product SIMD





Shuffle





PULP-CNN ISA-Extensions



Convolution Performance on PULP with 4 cores



5x5 convolution in only 6.6 cycles/pixel



Outline



- Near Threshold Multiprocessing
- Non-Von Neumann Accelerators
- Aggressive Approximation
- From Frame-based to Event-based Processing
- Outlook and Conclusion

Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich



Recovering silicon efficiency



Closing The Accelerator Efficiency Gap with Agile Customization



Computational Effort

- Computational effort
 - 10-class scene labeling on Stanford-BG
 - 7.5 GOp/frame for 320x240 image (#Op=2 × #MAC)
 - 260 GOp/frame for FHD
 - 1050 GOp/frame for 4k UHD





~90% workload is Conv

Origami CNN ASIC







- FP not needed: 12-bit signals sufficient
- Input to classification double-vs-12-bit accuracy loss < 0.5% (80.6% to 80.1%)</p>





Origami: The Architecture

- Keep moving window of all 8 input channels buffered on-chip
- 2. Perform all the convolutions
- Transmit result immediately after summing up

Maximize locality and data reuse (weights and pixels)





CNNs: typical workload



Example: ResNet-34

- classifies 224x224 images into 1000 classes
- ~ trained human-level performance
- ~ 21M parameters
- ~ 3.6G MAC operations

Scaling Origami to 28nm FDSOI

Performance for 10 fps: ~73 GOPS/s

Energy efficiency: ~2300 GOPS/W efficiency

0.4pj/OP

Origami core in 28nm FDSOI → 10 fps ResNet-34 with ~32mW



Outline



- Near Threshold Multiprocessing
- Non-Von Neumann Accelerators
- Aggressive Approximation
- From Frame-based to Event-based Processing
- Outlook and Conclusion



Pushing Further: YodaNN¹

- Approximation at the algorithmic side \rightarrow Binary weights
- BinaryConnect [Courbariaux, NIPS15], XOR NET [Rastegari, arXiv16]
 - Reduce weights to a binary value -1/+1
 - Stochastic Gradient Descent with Binarization in the Forward Path

$$w_{b,stoch} = \begin{cases} -1 & p_{-1} = \sigma(w) \\ 1 & p_1 = 1 - p_{-1} \end{cases} \qquad w_{b,det} = \begin{cases} -1 & w < 0 \\ 1 & w > 0 \end{cases}$$

- Learning large networks is still challenging, but starts to become feasible: ResNet-18 on ImageNet with 83.0% (binary-weight) vs. 89.2% (singleprecision) top-5 accuracy; and 60.8% vs. 69.3% top-1 accuracy
- Ultra-optimized HW is possible!
 - Major arithmetic density improvements: MAC → 2s compl. & Accum.
 - Area can be used for more energy-efficient weight storage
 - Storage reduction \rightarrow SCM memories for lower voltage \rightarrow E goes with $1/V^2$

¹After the Yedi Master from Star Wars - "Small in size but wise and powerful" cit. www.starwars.com

SoP-Unit Optimization





ImageBank

 $\frac{1000}{1000} = 2 \text{ Op (1 Op for the "sign-reverse", 1 Op for the add).}$

YODANN Energy Efficiency







Origami, YodaNN vs. Human

The «energy-efficient AI» challenge (e.g. Human vs. IBM Watson)

	Туре	Analog (bio)	Q2.9 Precision	Q2.9 Precision	Binary- Weight
	Network	human	ResNet-34	ResNet-18	ResNet-18
	Top-1 error [%]		21.53	30.7	39.2
	Top-5 error [%]	5.1	5.6	10.8	17.0
•	Hardware	Brain	Origami	Origami	YodaNN
	Energy-eff. [uJ/img]	100.000(*)	1086	543	31

*P_{brain} = 10W, 10% of the brain used for vision, trained human working at 10img/sec

- Game over for humans also in energy-efficient vision?
- Not yet! (object recognition is a super-simple task)





Heterogeneous PULP Cluster







Heterogeneous PULP Cluster

Hardware Convolutional Engine (HWCE) in the Cluster



HWCE CNN Performance

Cluster performance and energy efficiency on a 64x64 CNN layer (5x5 conv)



Scaled to ST FD-SOI 28nm @ Vdd=0.6V, f=115MHz



Outline



- Near Threshold Multiprocessing
- Non-Von Neumann Accelerators
- Aggressive Approximation
- From Frame-based to Event-based Processing
- Outlook and Conclusion

Back to System-Level



Smart Visual Sensor→ idle most of the time (nothing interesting to see)



- Event-Driven Computation, which occurs only when relevant events are detected by the sensor
- Event-based sensor interface to minimize IO energy (vs. Frame-based)
- Mixed-signal event triggering with an ULP imager, cochlea with internal processing AMS capability

A Neuromorphic Approach for doing *nothing* VERY well



GrainCam Imager





Analog internal image processing

- Contrast Extraction
- Motion Extraction, differencing two successive frames
- Background Subtraction with the reference image stored in pixel memory



Graincam Readout

120

100

80

60

40

20

Active

Idle



Readout modes:

- IDLE: readout the counter of asserted pixels
- Power Consumption [uW] ACTIVE: sending out the addresses of asserted pixels (address-coded representation), according raster scan order

Event-based sensing: output frame data bandwidth depends on the external context-activity





Power Management





System Design



M. **Rusci** et al. "A sub-mW IoT-endnode for always-on visual monitoring and smart triggering," in *IEEE Internet of Things Journal, 2017 (in print)*



Even-driven CNNs? Yes!

Binary Neural Networks reduce precision of weights and post-activation neurons to <u>1-bit precision</u> while leading to a limited performance drop





'Moving' pixel window

PE

ΡN

PO

Performing spatial filtering and binarization on the sensor die through mixed-signal sensing! \rightarrow in-sensor first stage of the binary NN!!



Per-pixel circitut for filtering and binarization





Event-Driven Binary Deep Network







Training challenge

Training Event-based Binarized Neural Network:

[ISSUE] Absence of huge amount of data for training

Modelling the "graincam filter" as a digital filter

Contrast	Va	$\max(p_E - p_O , p_N - p_o)$	Binary	$V_{0} = san(V_{c} - V_{th})$
Value	VC	$\max(p_E, p_O, p_N)$	Output	

Evaluation on **CIFAR-10** (10 classes, 45k training, 5k valid, 10k testing)

Baseline with RGB input	92%	
BNN with RGB input	86%	
Baseline with binary input	72%	
BNN with binary input	68%	
Model VGG-like with 12 Convolutional late Fully Connected Layers	ers and 3	
18% performance drop be input representation bu	cause o ut still	f

converges



Original RGB image Synthetic image Graincam image

Results





[2] http://podoce.dinf.usherbrooke.ca/



BNN implementation on PULP



$$o_{z}(x,y) = \frac{\rho(x,y) + b - \mu}{\sigma} \gamma + \beta \ge 0 \qquad \qquad \rho(x,y) \in N$$

if $\gamma \ge 0$ then $o_{z}(x,y) = \rho(x,y) \le \left[\mu - b - \frac{\beta * \sigma}{\gamma}\right]$ else $o_{z}(x,y) = \rho(x,y) \ge \left[\mu - b - \frac{\beta * \sigma}{\gamma}\right]$
inst logic operation and integer comparison!

ingle operation and integer companson!

Major opportunity for HW acceleration!



Preliminary Results

Scenario	BNN with RGB input	Event-based BNN
Image Sensor Power Consumption	1.1mW @30fps	$100\mu W @ 50 fps$
Image Size	632446 bits	8192 bits
Image Sensor Energy for frame capture	66.7 μJ	$2 \mu J$
Transfer Time (4bit SPI @50MHz)	3.1 msec	0.04 msec
Transfer Energy (8.9mW @0.7V)	$28 \ \mu J$	$2 \mu J$
BNN Execution Time (168MHz)	81.3 msec	75.3 msec
BNN Energy consumption (8.9mW @0.7V)	$725 \ \mu J$	671 μJ
Total System Energy for Classification	$820 \ \mu J$	$674 \ \mu J$

Statistics per frame	Frame-Based	Event-based
Idle (no motion)		
Sensor Power	1.1mW	$20\mu W$
Avg Sensor Data	19764 Bytes	-
Transfer Time	790μ sec	-
Processing Time	3.02 msec	-
Avg Processor Power	1.45mW	0.3mW (sleep)
Detection		
Sensor Power	1.1mW	$60 \mu W$
Avg Sensor Data	19764 Bytes	\sim 536 Bytes
Transfer Time	790μ sec	21.4μ sec
Processing Time	3.47 msec	187.6μ sec
Avg Processor Power	1.57mW	0.511mW
Classification		
Sensor Power	2mW	$60 \mu W$
Avg Sensor Data	79056 Bytes	1024 Bytes
Transfer Time	3.16 msec	41μ sec
Processing Time	81.3 msec	75.3 msec
Processor Energy	760 μ J	$677 \ \mu J$

84.6% vs. 81.6% Accuracy





Outline



- Near Threshold Multiprocessing
- Non-Von Neumann Accelerators
- Aggressive Approximation
- From Frame-based to Event-based Processing
- Outlook and Conclusion

System Integration: SoC & SiP



VivoSoC2 – SoC AFE integration



Multiple Biomedical Interfaces:

8 ExG channels with lead-off detection, Photoplethysmography (PPG), Bioimpedance monitoring, temperature sensing, 6 channels nerve blocker and stimulator, standard digital interfaces Full system in package concept



Vivosoc (130nm)

Low Cost Heterogeneous Integration:

CNN accelerator using highest density digital process for max GOPS/mW/mm². AFE A2D and GP processing in a good mixed-signal process, **Sensor and Memory** dies too...



Conclusions



- Near-sensor processing for the IoT
 - CNN are top performers for «sensor understanding» in the cloud today
- Energy efficiency requirements: pJ/OP and below
 - Technology scaling alone is not doing the job for us
 - Ultra-low power architecture and circuits are needed
 - Most promising technologies: 3D integration, low-leakage, non-volatile silicon-compatible mem-computing devices, but LOW maturity and HIGH cost → large-scale CNN engines will be early adopters
- CNNs can be taken into the ULP (mW power envelope) space
 - Non-von-Neumann acceleration
 - Very robust to low precision computations (deterministic and statistical)
 - fJ/OP is in sight!
- Open Source HW & SW approach → innovation ecosystem



Morale: Plenty of room at the bottom

Thanks!!!



www.pulp-platform.org www-micrel.deis.unibo.it/pulp-project iis-projects.ee.ethz.ch/index.php/PULP