

Applied Machine Learning Days @ EPFL

State of the art in hardware-accelerated neural networks

Frédéric Pétrot, Lorena Anghel, Liliana Andrade
Univ. Grenoble Alpes, CNRS, Grenoble INP^{*}, TIMA,
F-38000 Grenoble, France

🏠 tima.imag.fr/sls/people/petrot

✉ frederic.petrot@univ-grenoble-alpes.fr



^{*}Institute of Engineering Univ. Grenoble Alpes

The Brain: the Ultimate Autonomous System

- ▶ 1,2 to 1,4 kg, 1260 cm³
- ▶ Consumes between 15 and 30 Watts
- ▶ 86×10^9 neurones, $\approx 10^{12}$ synapses

The Brain: the Ultimate Autonomous System

- ▶ 1,2 to 1,4 kg, 1260 cm³
- ▶ Consumes between 15 and 30 Watts
- ▶ 86×10^9 neurones, $\approx 10^{12}$ synapses



J. Vitti and D. Silverman, "Bart the Genius", The Simpsons, 1990

The Brain: the Ultimate Autonomous System

- ▶ 1,2 to 1,4 kg, 1260 cm³
- ▶ Consumes between 15 and 30 Watts
- ▶ 86×10^9 neurones, $\approx 10^{12}$ synapses



J. Vitti and D. Silverman, "Bart the Genius", The Simpsons, 1990

Human Brain Project (EU Flagship)

- ▶ 16,000 neurones per 1 Watt chip
- ▶ 5.375 MW/brain
- ▶ 1.2×10^9 € from Europe: ≈ 1.4 cent of €/neurone

The Brain: the Ultimate Autonomous System

- ▶ 1,2 to 1,4 kg, 1260 cm³
- ▶ Consumes between 15 and 30 Watts
- ▶ 86×10^9 neurones, $\approx 10^{12}$ synapses
- ▶ (Energy comes from eggs, honey and mushrooms)



J. Vitti and D. Silverman, "Bart the Genius", The Simpsons, 1990

Human Brain Project (EU Flagship)

- ▶ 16,000 neurones per 1 Watt chip
- ▶ 5.375 MW/brain
- ▶ 1.2×10^9 € from Europe: ≈ 1.4 cent of €/neurone



The Brain: the Ultimate Autonomous System

- ▶ 1,2 to 1,4 kg, 1260 cm³
- ▶ Consumes between 15 and 30 Watts
- ▶ 86×10^9 neurones, $\approx 10^{12}$ synapses
- ▶ (Energy comes from eggs, honey and mushrooms)



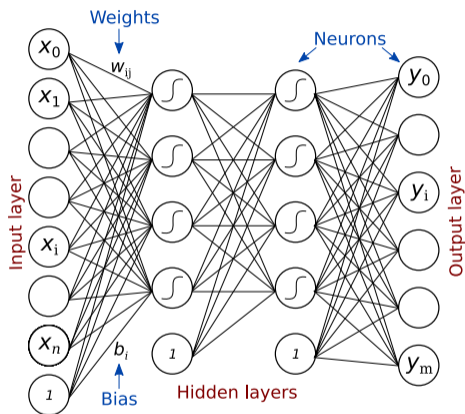
J. Vitti and D. Silverman, "Bart the Genius", The Simpsons, 1990

Human Brain Project (EU Flagship)

- ▶ 16,000 neurones per 1 Watt chip
- ▶ 5.375 MW/brain
- ▶ 1.2×10^9 € from Europe: ≈ 1.4 cent of €/neurone



Classical Deep Neural Network (DNN) Topology



$$u_i = \sum_{j=1}^n x_j w_{ij}$$

$$y_i = \text{act}(u_i + b_i)$$

Two phases

- ▶ Learning, off-line, floating point
- ▶ Inference, run-time, floating point

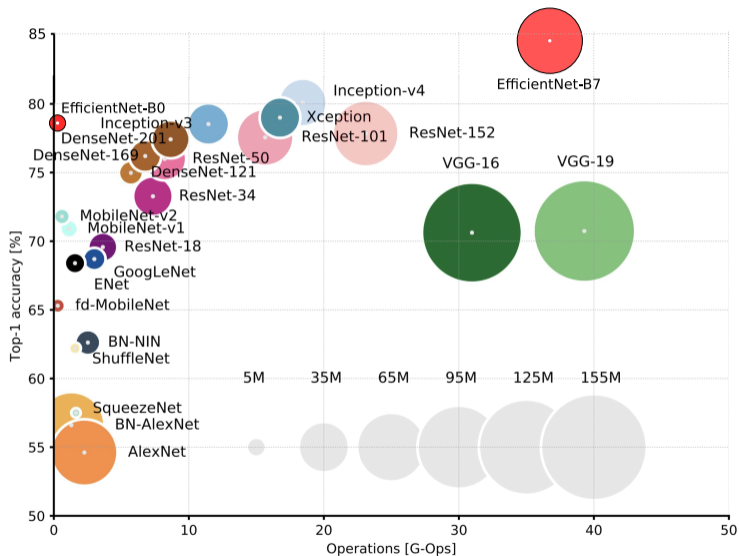
Weights

- ▶ Values “found” during learning
- ▶ Values used during inference

Issues

- ▶ Complex learning algorithms: GPU farms
- ▶ Limited computing power during inference

CNN Models: Accuracy, Operations and Weights



Our focus? Inference!

What's the interest, by the way?

Local computation!

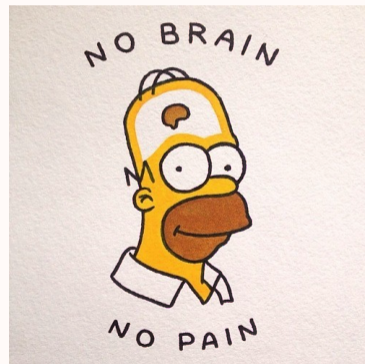
- ▶ Energy
 - No router, cloud server, ...
 - ⇒ Huge constraint in Edge Computing
 - ⇒ Worse in IoT
 - ⇒ Transmitting data costs energy
- ▶ Latency
 - Immediate response, no dead zone, no network reliability issue, ...
- ▶ Privacy
 - No storage in someone else's servers
 - Neither wire nor wireless sniffing possible

Constraints on HW Accelerated Neural Networks

- ▶ Accuracy needs depend on the application
- ▶ Silicon resources:
 - ⇒ Computations to perform
 - ⇒ Weights storage and access
- ▶ Energy efficiency
Typical constraints :
 - 10-100 μ W for wearables,
 - 10-100 mW for phones,
 - 1-10 W for plugged devices

Constraints on HW Accelerated Neural Networks

- ▶ Accuracy needs depend on the application
- ▶ Silicon resources:
 - ⇒ Computations to perform
 - ⇒ Weights storage and access
- ▶ Energy efficiency
Typical constraints :
 - 10-100 μ W for wearables,
 - 10-100 mW for phones,
 - 1-10 W for plugged devices



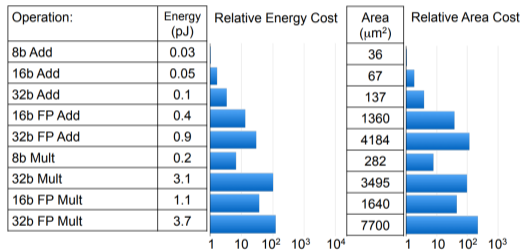
Inference involves a lot of computation...

- ▶ Elevated number of floating point (FP) operations

$$0.5G \leq \text{Nb of FLOPs} \leq 40G$$

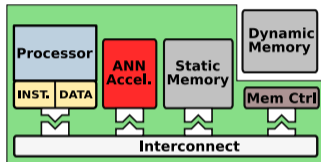
- ▶ Floating point operations are energy and area costly

(My 4 core-i7 PC ~120 GFLOPs \Rightarrow 30 GFLOPs/core.)



"Hardware Architectures for Deep Neural Networks", ISCA Tutorial, 2017

Inference involves a lot of memory access...



Operation:	Energy (pJ)	Relative Energy Cost
32b SRAM Read (8KB)	5	~10
32b DRAM Read	640	~10 ³

1 10 10² 10³ 10⁴

"Hardware Architectures for DNN", ISCA Tutorial, 2017

- ▶ Memory stores millions of (64-bit) weights
⇒ 4M (GoogLeNet), 60M (AlexNet), 130M (VGG)
- ▶ Memory access becomes the bottleneck
⇒ Each op needs 2 operands and produces a result
- ▶ An “elevated” power consumption is involved

Coping with GFLOPs and GBytes

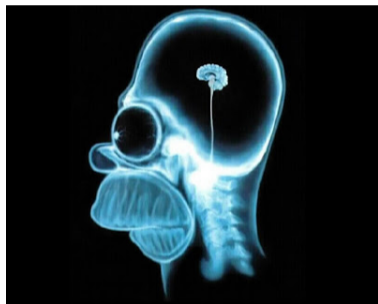
Alternatives: trade FLOPs for (some) accuracy loss

Simplify the operations

- ▶ Avoid sigmoid, batch normalization and stuff
- ▶ FP arithmetic is not HW friendly
⇒ Use data-types that are not 64-bit floats

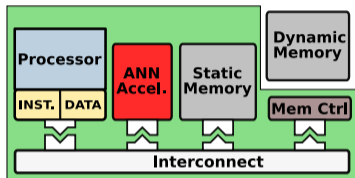
Alternatives: trade bytes for (some) accuracy loss

- ▶ Use “small” data types
- ▶ Integrate many memory cuts with processing elements and use them wisely
- ▶ Integrate computation into the memory itself

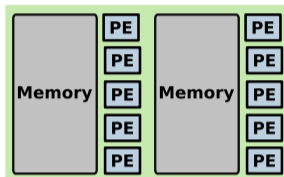


K. Usher, "The Dwarf in the Dirt", Bones, 2009
(Energy comes from donuts + beer)

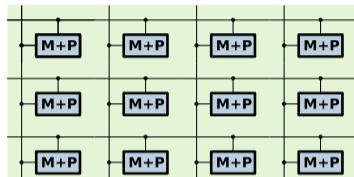
Typical Architectures for HW ANN



Exploit weight sparsity to optimize memory usage and weight placement



Use low precision/high efficiency computation along with on-chip memory storage of the weights



Integrate computation inside the memory itself, directly where the data is stored

Quantization: we need it anyway!

Quantization levels and accuracy...

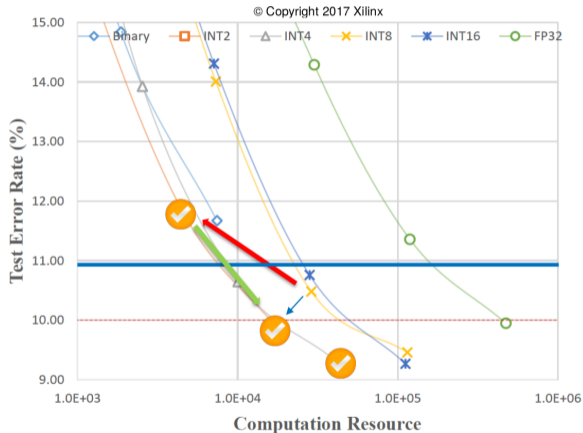
Just reducing precision,
reduce hardware cost &
increases error



Recuperate accuracy by
retraining & increasing
network size

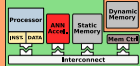


1b, 2b and 4b provide pareto
optimal solutions



Kees Vissers, "A Framework for Reduced Precision Neural Networks on FPGAs", MPSOC, 2017

Exploit Sparsity and Quantization



Custom hardware for sparse matrix-vector multiplication

Deep Compression Technique

Reduces storage requirements

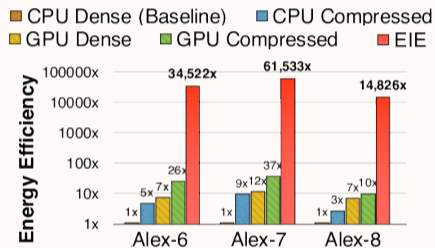
- ▶ Dedicated sparse matrix/vector representation
⇒ Eliminates redundant connections
- ▶ Quantizes weights

Quantization of AlexNet weights

- ▶ 256 shared weights (Conv layers) ⇒ 4 bits
- ▶ 35x of reduction (240MB ⇒ 6.9MB)

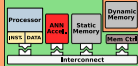
Weights stored into on-chip SRAM

⇒ 5 pJ/access (vs. 640 pJ/access off-chip DRAM)



Power efficiency

600 mW for Alexnet Fully-Connected layers



Acceleration using Low-Precision (ternary) weights

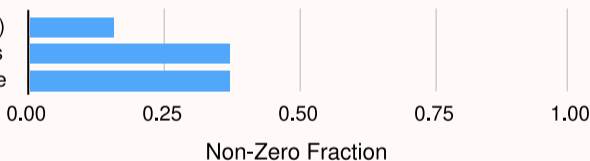
Only balanced **ternary weight** are used $\{-1, 0, +1\}$

- ▶ Floating point accumulations are kept
- ▶ Multipliers are not needed

Most of the FP operations operate on zero values

- ▶ Zero-skipping

FwdPass (Inference)
BWD Pass
Grad Update



Demonstrated highest accuracy

- ⇒ 93% on the ImageNet object classification challenge
- ⇒ Divide by 3 the number of FP operations

YodaNN: VLSI Implementation of binary-weights CNN Accelerator

Based on BinaryConnect approach

- ▶ Binary weights $\in \{-1, +1\}$
- ▶ 2's complement and multiplexers instead of multipliers
- ▶ Still full fledged adders: 12-bit activations

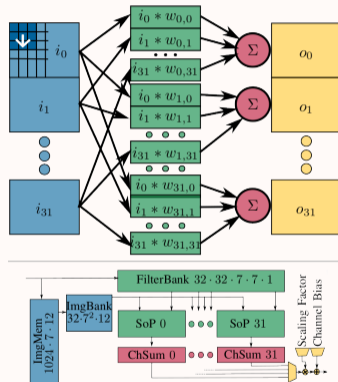
Large on-chip weights storage thanks to their size

- ▶ Latch-based standard cell memory

Flexible accelerator

- ▶ 7 kernel sizes supported

⇒ 61.2 TOP/s/W at 0.6V



FINN: Framework for building FPGA* accelerators

Mapping binarized neural networks to hardware

All values $\in \{-1, +1\}$

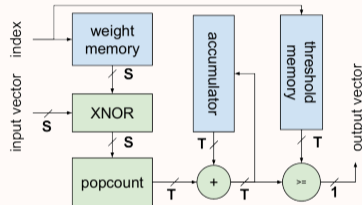
- ▶ Binary input activation
- ▶ Binary synapse weights
- ▶ Binary output activation

Weights kept in on-chip memory

- ⇒ Zynq-7000 FPGA technology
- ⇒ 80.1% accuracy for CIFAR-10
- ⇒ Total system power 25W

Y. Umuroglu et al., "FINN: A Framework for Fast, Scalable Binarized Neural Network Inference", FPGA, 2017

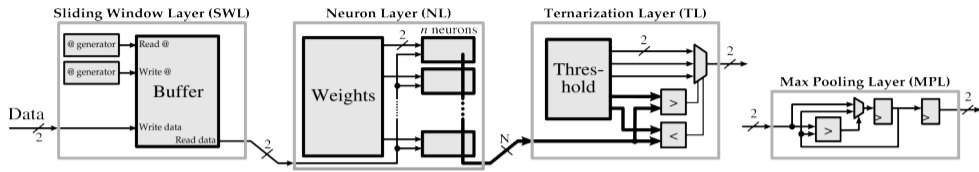
Convolution layer



- ▶ Dot-product between input vector and row of synaptic weight matrix
- ▶ Compares result to a threshold
- ▶ Produces single-bit output

* Field-Programmable Gate-Array: fine-grain reconfigurable hardware technology.

Ternary weights and ternary activations



FPGA Architecture for Ternary Neural Networks (TNN)

- ▶ Large-scale ternary CNN pipeline VGG-like (NN-64 or NN-128)
- ▶ Neuron layer (NL) → memory (ternary weights) + neurons
- ▶ Ternarization layer (TL) → ternary activations $\in \{-1, 0, +1\}$

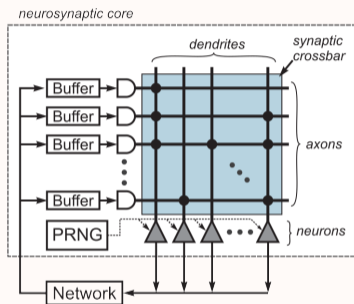
⇒ Error rate 13.29% for CIFAR-10 (vs. 19.9% in FINN)

⇒ Virtex-7 FPGA technology (VC709, Laaaaaaarge FPGA)

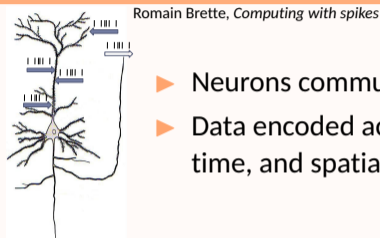
⇒ 1.62 TOP/s/W (vs 0.69 TOP/s/W in FINN)

⇒ Throughput > 60k fps

TrueNorth: Integrated Chip for Spiking Neural Networks



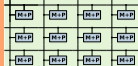
- ⇒ 4096 neuromorphic cores
 - ⇒ 1 million of digital neurons
 - ⇒ 256 millions of synapses
 - ⇒ 46 GSOP/s/W at 65 mW
- Asynchronous logic implementation



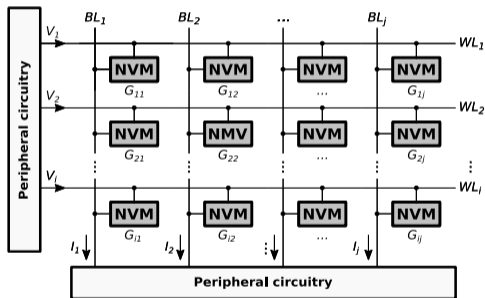
- ▶ Neurons communicate sending spikes
- ▶ Data encoded according to frequency, time, and spatial distribution of spikes

Non-Von Neumann architecture

- ▶ Neuromorphic core = 256 neurons (PE) + 64k synapses (memory)
- ▶ Memory and computation physically close to each other
- ▶ Reduction of power consumption



MAC operations using Non Volatile Memory (NVM)



$$I_j = V_1 \cdot G_{1j} + V_2 \cdot G_{2j} + \dots + V_i \cdot G_{ij}$$

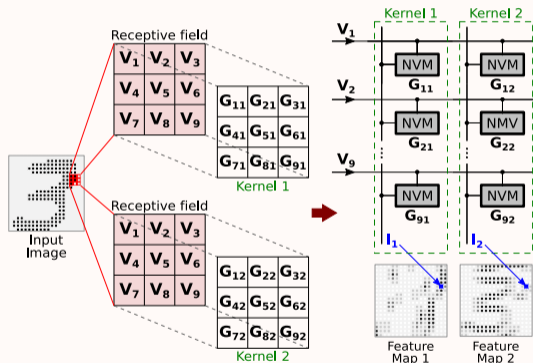
Computation accelerated by NVM arrays

- ▶ Synaptic weights are **not** stored in external memories
 - ⇒ Zero transfers between memory and processing elements
 - ⇒ Reduction of energy consumption

Arrays of resistive RAM devices

- ▶ Resistances vary according to voltages
- ▶ No CMOS access devices but complex peripheral circuitry
- ▶ Analog, intrinsically approximate, computations

Convolution Kernel in 12×12 Array

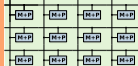


From receptive field to feature maps...

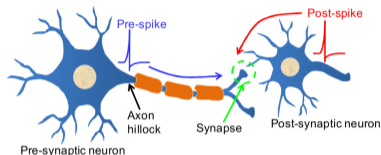
- ▶ Receptive field = row voltages
- ▶ Convolution kernel = column of resistive devices
- ▶ Convolution operation = column current

Interactive protocol programs kernels

Single demonstration on digits of MNIST



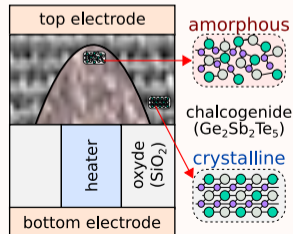
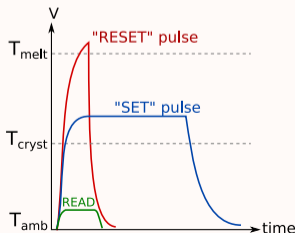
⇒ Artificial synapse using NVM



- ▶ Modeling synapses between neurons
- ▶ Input and output potentials fired between neurons (spikes)
- ▶ Synaptic connections are potentiated or depressed

Electronic Synapses Modeled by Phase Change Memory (PCM) Devices

- ▶ Programmed in different states (conductances)
- ▶ Compatible with CMOS components
- ▶ Scalable to nanometric dimensions



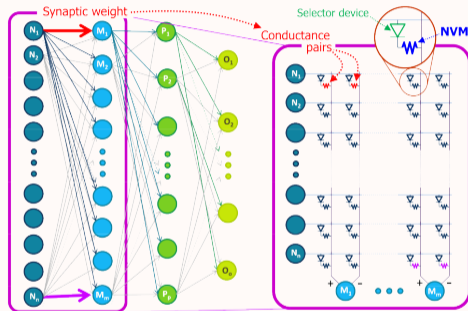
500x661 PCM Crossbar Array

Large scale implementation

- ▶ 3-layer perceptron
- ▶ 916 neurons
- ▶ 164865 synaptic connections

⇒ Accuracy: 82% (MNIST)

⇒ Low-power: at least 120x (vs. GPU)



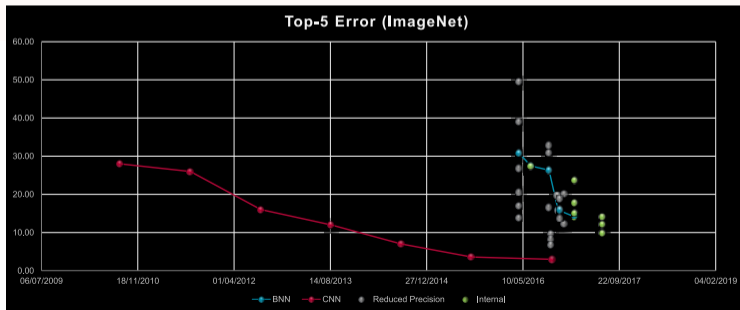
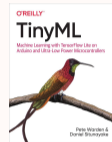
S. Burc et al., "Experimental Demonstration of Array-Level Learning with Phase Change Synaptic Devices", IEDM, 2013

G.W. Burr et al., "Experimental Demonstration and Tolerancing of a Large-Scale NN (165000 Synapses) using PCM as the Synaptic Weight Element", 2015

G.W. Burr et al., "Large-Scale Neural Networks Implemented with NVM as the Synaptic Weight Element: Comparative Performance Analysis", IEDM, 2015

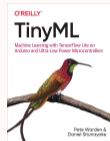
Classical digital (CMOS) architectures

- ▶ Quantization and compression is the way to go
- ▶ Few bits for weights and activations are enough in many cases
- ▶ Learning need to be changed according to bit-width
- ▶ Networks and HW architectures available today!



Classical digital (CMOS) architectures

- ▶ Quantization and compression is the way to go
- ▶ Few bits for weights and activations are enough in many cases
- ▶ Learning need to be changed according to bit-width
- ▶ Networks and HW architectures available today!



Processing in-memory approaches using emerging technologies

- ▶ Few experiments of NVM crossbar array implementations at large scale
- ▶ Comparison of energy gains is difficult
- ▶ Flexibility of the NVM to match a given ANN architecture questionable
- ▶ Promising ongoing research subject