LINE CONDUCTOR TECHNOLOGIES FOR ENERGY-EFFICIENT INTELLIGENCE

PETER.DEBACKER@IMEC.BE

CONFIDENTIAL



MACHINE LEARNING DEEP NEURAL NETWORKS



Training of CNN:

х1

x2

xЗ

x4

 Huge data set with known objects to determine value for millions of weights

Forward "guess"

Weight correction

Typically in data center on high-performance GPU



umec

MACHINE LEARNING DEEP NEURAL NETWORKS





x1 x2 xЗ x4 7 layers Activations Weights Multiply X1 W1 Accumulate $\int f(\sum_{i=1}^n W_i X_i)$ W2 X2 W3 X3

ເກາຍc

-

MACHINE LEARNING USING DEEP NEURAL NETWORKS >95% OF OPERATONS ARE MULTIPLY-ACCUMULATE

ReLU-4

ReLU-5

Pool-5

Conv-5



Conv-3

ReLU-3

Conv-4

- Convolutional (Conv) layers: feature extractors
- Fully-Connected (FC) layers: classifiers

Pool-2

Conv-2

ReLU-1

Conv-1

Pool-1

ReLU-2

- Remaining functions (pooling, threshold, normalization, ...) are critical for correct operation of DNN but do not dominate performance
- Number of layers and number of weights growing in pursuit of better classification accuracy.



Source: https://medium.com/towards-data-science/neural-network-architectures-156e5bad51ba and arXiv:1605.07678v4

ເຫາຍດ

DEEP NEURAL NETWORKS (CNN, LSTM)

MAINSTREAM DEEP LEARNING

- Convolutional Neural Networks: visual processing and object classification
 - Character recognition
 - Medical image analysis
 - People / face identification
 - ...
- LSTM
 - Time-series prediction, e.g. predict the weather or stocks
 - Action recognition, e.g. event detection in security cameras
 - Robotic movements, e.g. drive various joints based on sensory inputs
 - Speech recognition, e.g. ask Siri, Alexa, Cortana or Google assistant questions
 - Language models, e.g. Google translate
 - Anomaly detection, e.g. firewall applications in network traffic

SMART SERVICES REQUIRE ENERGY-EFFICIENT MACHINE LEARNING ENABLE SMART EDGE DEVICES ... GROWING DATACENTER INFERENCE WORKLOAD



HARDWARE PLATFORMS FOR MACHINE LEARNING

TODAY'S DOMINANT PLATFORMS TARGET HIGH-PERFORMANCE MARKET

- Focus on <u>training</u> of deep learning algorithms
- General purpose platforms (GPU, FPGA) support flexibility in still evolving field
- First application-specific IC (Google TPU) emerge



ເງຍ

NEURAL NET ACCELERATORS (ASIC) EMERGE FOR INFERENCE CONVENTIONAL DIGITAL CMOS

Power efficiency TOPSIN LOTOPSINN LOOTOPS/W 1000000 < 100 TOP/s/W tive Mobile ЬТ application Auton Google TPL LOOGOPSIN 100000 Tesla V100 specific HW **Binary NN** Tesla P40 accelerator large embedded .og(Speed) (Gop/s) 5nm (iN7) Tesla P100 Xavier LOGOPSIN Tesla P4 10000 technology memories ASIC aggressively Titan X 1GOPS/W GPU quantized NN 1000 Conventional D Power efficience 100 **FPGA** [Based on https://nicsefc.ee.tsinghua.edu.cn/projects/neural-network-accelerator/] 10 0.001 0.01 0.1 100 1 10 1000 Log (Power) (W)



OPPORTUNITIES FOR ANALOG IN-MEMORY COMPUTING CHALLENGES OF VON NEUMANN COMPUTING



- **Challenges:** Fundamental difficulties in scalability of memory bandwidth and capacity
- Emerging Opportunities
 - Performing calculation in memory
 - Non-volatile memory by limiting the data transfer between memory and CPU

Von-Neumann architectures

ເງຍ

ANALOG COMPUTE IN-MEMORY

ANALOG COMPUTE-IN-MEMORY ACCELERATORS FOR ML SUPPORTED BY NEW MEMORY TECHNOLOGY



- 1. Memristor (programmable resistor) stores weight value W_1 as an <u>analog quantity</u>: conductance Rw_1
- 2. Activation X₁ applied as <u>analog voltage</u>Vin₁
- 3. Ohm's law: memristor cell current $\sim X_1.W_1$
 - . Kirchoff's law: bit-line current $\sim \sum X_i W_i$





ANALOG COMPUTE-IN-MEMORY ACCELERATORS FOR ML SUPPORTED BY NEW MEMORY TECHNOLOGY



- Memristor (programmable resistor) stores weight value W₁ as an analog quantity: conductance Rw₁
- 2. Activation X₁ applied as analog voltage Vin₁
- 3. Ohm's law: memristor cell current $\sim X_1.W_1$
- 4. Kirchoff's law: bit-line current $\sim \sum X_i W_i$





ANALOG COMPUTE-IN-MEMORY ACCELERATORS FOR ML SUPPORTED BY NEW MEMORY TECHNOLOGY







CONFIDENTIAL



- Use memory array for massive parallel analog implementation of multiply-accumulate operations in DNN layer
- Memory array stores weights <u>and</u> implements a logic function (MAC) in analog fashion
 - \rightarrow compute-in-memory
 - \rightarrow computational memory
 - ightarrow neuromorphic computing

ເງຍ

POOLING, THRESHOLDING, ... IN DIGITAL DOMAIN DAC AND ADC REQUIRED





ເງງອ

Combine algorithm and architecture optimization with <u>semiconductor</u> <u>technology elements</u> to enable <u>energy-efficient</u> implementation of DNNs and beyond



SYSTEM DESIGN-TECHNOLOGY CO-OPTIMIZATION FOR ANALOG COMPUTE-IN-MEMORY

ANALOG COMPUTE-IN-MEMORY ACCELERATORS FOR ML CHALLENGES ...



- 1. Memristor (programmable resistor) stores weight value W_1 as an analog quantity: conductance Rw_1
- 2. Activation X₁ applied as analog voltage Vin₁
- 3. Ohm's law: memristor cell current $\sim X_1.W_1$
- 4. Kirchoff's law: bit-line current $\sim \sum X_i W_i$





CONFIDENTIAL

ANALOG COMPUTE-IN-MEMORY ACCELERATORS FOR ML CHALLENGES







HOW TO MAKE IT WORK? SYSTEM-TECHNOLOGY CO-OPTIMIZATION

- Deep neural networks
 - Perform inference with lower precision
 - Can tolerate some weight variation
- Neural network optimization
 - Train neural network with information of implementation limitations
 - Increase size of neural network to limit accuracy loss
- Architecture
 - Use multiple memory cells per weight for more precision
- Circuit and device
 - Tune devices for neural network operation





ເກາຍc

NEURAL NETWORK OPTIMIZATION FOR ACIM QUANTIZATION BRINGS PRECISION IN RANGE OF ACIM







- Accuracy loss due to quantization can be mitigated
 - Modified training procedure
 - Modified network structure:
 - Higher #operations
 - Simpler operations
 - No accuracy loss for key benchmark networks using
 - 3-levels for weight
 - 5-levels for activations
 - \rightarrow In range of analog storage elements and analog MAC



NEURAL NETWORK OPTIMIZATION FOR ACIM TRAIN WITH ACIM IMPERFECTIONS IN MIND

 Up to 10% noise on weights and activations, and 50% LSB noise on MAC can be mitigated by training.

Keyword spotting (KWS) Model	# param	Size	# Macs	#ACIM rows	#ACIM columns	Accuracy
LSTM (FP)	800K	3.2MB	79.2M	100-400	1600	93%
QLSTM	800K	200kB	79.2M	100-400	1600	93.2%
CNN (FP)	50K	200kB	3.5M	135	315	94.3%
QCNN	50K	12.5kB	3.5M	135	315	94.2%
ResNet-15 (Tang et al.2018)	238K	952kB	894M	405	585	95.8%

- All quantization with 3-level weight, 7-level activation, 15-level accumulation and additive noise
- CNN: ID dilated, 45 filters
- ResNet-15 is state-of-the-art for KWS; uses different preprocessing of data: filtering, more cepstral components, different window length and no input quantization layer.

Input: 39 features from MFCC					
FC (N=100) + BN					
ļ					
Quantize to 7 levels					
1					
ID-conv: 45*3 filter, 2-dilation - BN - ReLU					
ID-conv: 45*3 filter, 2-dilation - BN - ReLU					
ID-conv: 45*3 filter, 4-dilation - BN - ReLU					
ID-conv: 45*3 filter, 4-dilation - BN - ReLU					
ID-conv: 45*3 filter; 8-dilation - BN - ReLU					
ID-conv: 45*3 filter, 8-dilation - BN - ReLU					
ID-conv: 45*3 filter, 16-dilation - BN - ReLU					
Global-average pooling					
Softmax layer					

NEURAL NETWORK OPTIMIZATION FOR ACIM TRAIN WITH ACIM IMPERFECTIONS IN MIND

Up to 10% noise on weights and activations, and 50% LSB noise on MAC can be mitigated by training.

Keyword spotting (KWS) Model	# param	Size	# Macs	#ACIM rows	#ACIM columns	Accuracy
LSTM (FP)	800K	3.2MB	79.2M	100-400	1600	93%
QLSTM	800K	200kB	79.2M	100-400	1600	93.2%
CNN (FP)	50K	200kB	3.5M	135	315	94.3%
QCNN	50K	12.5kB	3.5M	135	315	94.2%
ResNet-15 (Tang et al.2018)	238K	952kB	894M	405	585	95.8%

- All quantization with 3-level weight, 7-level activation, 15-level accumulation and additive noise
- CNN: ID dilated, 45 filters

unec

- ResNet-15 is state-of-the-art for KWS; uses different preprocessing of data: filtering, more cepstral components, different window length and no input quantization layer.

	Input: 39 features from MFCC						
	ļ						
	FC (N=100) + BN						
	Quantize to 7 levels						
	Ļ						
10	D-conv: 45*3 filter, 2-dilation - BN -	ReLU					
10	D-conv: 45*3 filter, 2-dilation - BN -	ReLU					
10	0-conv: 45*3 filter, 4-dilation - BN -	ReLU					
10	0-conv: 45*3 filter, 4-dilation - BN -	ReLU					
IC	D-conv: 45*3 filter, 8-dilation - BN -	ReLU					
10	D-conv: 45*3 filter, 8-dilation - BN -	ReLU					
ID	-conv: 45*3 filter, 16-dilation - BN -	ReLU					
	Global-average pooling						
	Softmax laver						

Process variability can be accounted for during training



act[N-1]

act[0]-

MRAM based synaptic cell with 3-level weight encoding

M

colpair[0]

W

colpair[M-1]

-

W

AWA



ເງຍອ

CONCLUSIONS

CONCLUSIONS

- CNN and LSTM are the work horses for deep learning
 - Huge computational complexity (Mult-Acc) and storage requirements (Weights)
 - High-performance compute architectures dominant platform for development and training phase
 - GPU: I TOP/s/W, >100W
- Quantization of neural networks trade-off complexity vs. accuracy
 - Key for low-energy digital neural network inference accelerators
 - Binary neural network accelerators can approximate 100 TOP/s/W
 - Necessity for mixed-signal and analog compute-in-memory implementations
 - Inherent variability, stochastic behavior, noise limits precision of operations
 - Promise large gain in energy efficiency to >1,000 TOP/s/W
 - 3-level weights and 5b-activations provide required accuracy and are in range of Analog Compute-in-Memory
 - Several cell architectures (MRAM, PCM, DRAM, SRAM-based) can enable ACiM accelerators
 - Training procedures and network topologies are adapted to limit/eliminate accuracy loss



ເງຍ

CONFIDENTIAL

embracing a better life