

#### ALMA MATER STUDIORUM Università di Bologna

## A Fast-but-Gentle Introduction to Artificial Intelligence Acceleration

Dr. Francesco Conti, DEI & ARCES f.conti@unibo.it

### **Deep NN Timeline**

- 1940s: Neural networks were proposed
- 1960s: Deep neural networks were proposed
- 1989: Neural network for recognizing digits (LeNet)
- 1990s: Hardware for shallow neural nets
  - Example: Intel ETANN (1992) NVIDIA GPUs with CUDA available
- 2011: Breakthrough DNN-based speech recognition
  - Microsoft real-time speech translation
- 2012: DNNs for vision supplanting traditional ML
  - AlexNet for image classification
- 2014+: Rise of DNN accelerator research
  - Examples: Neuflow, DianNao, etc.



[Yann LeCun, ISSCC 2019]



### Deep ConvNets (depth inflation)





### Searching for "AI" on Google Image Search









[CC images, various sources]

### Searching for "Al" on Google Image Search





### Ok, not much information here!





[CC images, various sources]

### Searching for "Deep Neural Network" on Google Image Search









[CC images, various sources]

### Searching for "Deep Neural Network" on Google Image Search







### Much better! But still not crystal clear.



[CC images, various sources]

### Looking inside papers!



N×



[LSTM layer, image from Wikipedia, CC BY 4.0 Guillaume Chevalier]



[Inception-ResNet v2,

https://ai.googleblog.com/2016/08/improving-inception-and-image.html]



3x3 conv, 256

3x3 conv, 512, /2

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512

3x3 conv, 512 3x3 conv, 512

avg pool fc 1000

ALMA MATER STUDIORUM UNIVERSITÀ DI BOLOGNA











multiplications & additions





multiplications & additions



















1. how to represent data

Method	Data Range	Parameter Range	ImageNet Top1 Acc
Full-Precision	Real numbers (FP32)	Real numbers (FP32)	69.6
Linear Quantization	Integers 0 to 255	Integers -128 to +127	69.6
РАСТ	Integers 0 to 15	Integers -8 to +7	69.2
PACT	Integers 0 to 7	Integers -4 to +3	68.1
PACT	Integers 0 to 3	Integers -2 to +1	64.4
PACT	Integers 0 to 3	-1 or +1	62.9
XNOR-Net	-1 or +1	-1 or +1	51.2



1. how to **represent** data

Method	Data Range	Parameter Range	ImageNet Top1 Acc
Full-Precision	Real numbers (FP32)	Real numbers (FP32)	69.6
Linear Quantization	Integers 0 to 255	Integers -128 to +127	69.6
РАСТ	Integers 0 to 15	Integers -8 to +7	69.2
PACT	Integers 0 to 7	Integers -4 to +3	68.1
РАСТ	Integers 0 to 3	Integers -2 to +1	64.4
PACT	Integers 0 to 3	-1 or +1	62.9
XNOR-Net	-1 or +1	-1 or +1	51.2

#### 2. how to **compute** data transformations



multiplications & additions

i. dominated by simple arithmetic operations

ii. many-to-one → many possible compute orderings
 iii. independent operations → can be done in parallel / hierarchically
 iv. hardware complexity / speed /energy related to representation



1. how to **represent** data

Method	Data Range	Parameter Range	ImageNet Top1 Acc
Full-Precision	Real numbers (FP32)	Real numbers (FP32)	69.6
Linear Quantization	Integers 0 to 255	Integers -128 to +127	69.6
PACT	Integers 0 to 15	Integers -8 to +7	69.2
РАСТ	Integers 0 to 7	Integers -4 to +3	68.1
PACT	Integers 0 to 3	Integers -2 to +1	64.4
PACT	Integers 0 to 3	-1 or +1	62.9
XNOR-Net	-1 or +1	-1 or +1	51.2

#### 2. how to **compute** data transformations



multiplications & additions

3. where to store data and how to move them around



The main source of headaches for DL Accelerator Architects!



LMA MATER STUDIORUM Jniversità di Bologna

i. dominated by simple arithmetic operations
ii. many-to-one → many possible compute orderings
iii. independent operations → can be done in parallel / hierarchically
iv. hardware complexity / speed /energy related to representation

### **A Minimal Accelerator**

**COMPUTE:** 



#### **MEMORY:**

### Off-Chip>10<sup>-9</sup>J





Non-Von Neumann... see the other talk!





Off-Chip>10<sup>-9</sup>J





Worst Case: all memory R/W are DRAM accesses

AlexNet [NIPS 2012] has **724M** MACs → **2896M** DRAM accesses required

DRAM access 100-1000x less energyefficient than on-chip access!













Reduce memory cost by error resilience







Reduce memory cost by error resilience



Reduce transfer cost by data tiling





Reduce memory cost by error resilience

Reduce transfer cost by data tiling



### The Recipe for a Deep Learning Accelerator

- 1. Many **Multiply-Accumulate (MAC)** units to exploit parallelism
- 2. Flexible or Customized **on-chip memory organization** to keep as much data as possible on-chip, maximise its reuse...
- 3. Keep track of all **external memory** transfer overheads!



### The Recipe for a Deep Learning Accelerator

- 1. Many Multiply-Accumulate (MAC) units to exploit parallelism
- 2. Flexible or Customized **on-chip memory organization** to keep as much data as possible on-chip, maximise its reuse...
- 3. Keep track of all **external memory** transfer overheads!

*Custom architectures* 



Accelerator

Fabric

Intelligent Memory

16 SHAVE Vect

Myriad X VPU Architecture

#### Tesla FSF Chip



. D. 1080

#### Google Cloud TPU



#### Intel/Movidius Myriad X

Imaging Accelerato

Movidius

### The Recipe for a Deep Learning Accelerator









[WikiChip]

### The Queen of Deep Learning Accelerators: the GPU



### The Queen of Deep Learning Accelerators: the GPU



**GPU Tensor Cores** 

# **TENSOR CORE**

### Mixed Precision Matrix Math 4x4 matrices

**D** =

FP16 or FP32









FP16 or FP32

D = AB + C



### There's Plenty of Room at the Bottom

- The relationship between **data representation**, **network topology**, and **perf/energy/memory** is not yet fully explored (particularly for tiny devices)!
- The Von Neumann model could be suboptimal: is it possible to sidestep memory in doing Multiply-Adds?
- Will future sophisticated AI algorithms show the same "good" properties of DNNS: **regularity**, **parallelism**, **resilience**?



[Lin et al., MCUNet: Tiny Deep Learning on IoT Devices]





#### ALMA MATER STUDIORUM Università di Bologna

#### Prof. Francesco Conti

DEI – Università di Bologna

f.conti@unibo.it

www.unibo.it