



Reliable and Efficient hardware for **Trustworthy** and **Sustainable** Deep Neural Networks

Alberto Bosio

alberto.bosio@ec-lyon.fr



Institut des Nanotechnologies de Lyon UMR CNRS 5270

<http://inl.cnrs.fr>

Acknowledgments



RE-TRUSING Project, Grant number: ANR-21-CE24-0015
<https://inl.cnrs.fr/projects/re-trusting/>

AdequetDL Project, Grant number: ANR-18-CE23-0012



PEPR IA Adapting



Institut des Nanotechnologies de Lyon UMR CNRS 5270

<http://inl.cnrs.fr>

Outline

3

- Introduction
- Efficient HW accelerators
- Reliable HW accelerators
- Conclusions



Outline

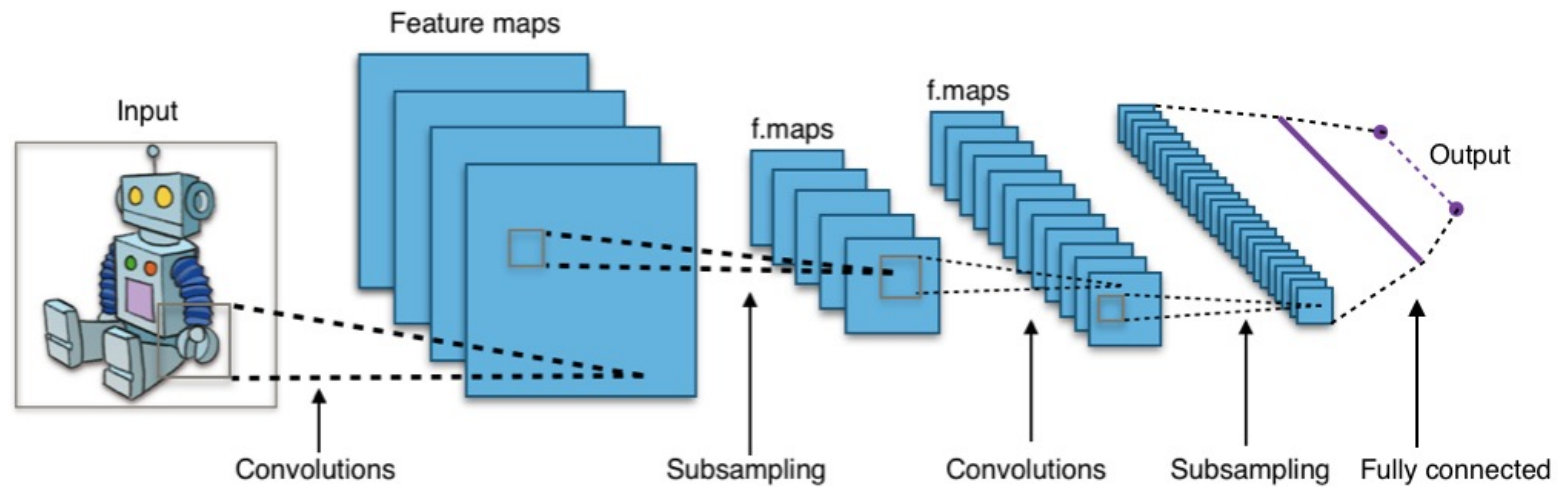
4

- Introduction
- Efficient HW accelerators
- Reliable HW accelerators
- Conclusions

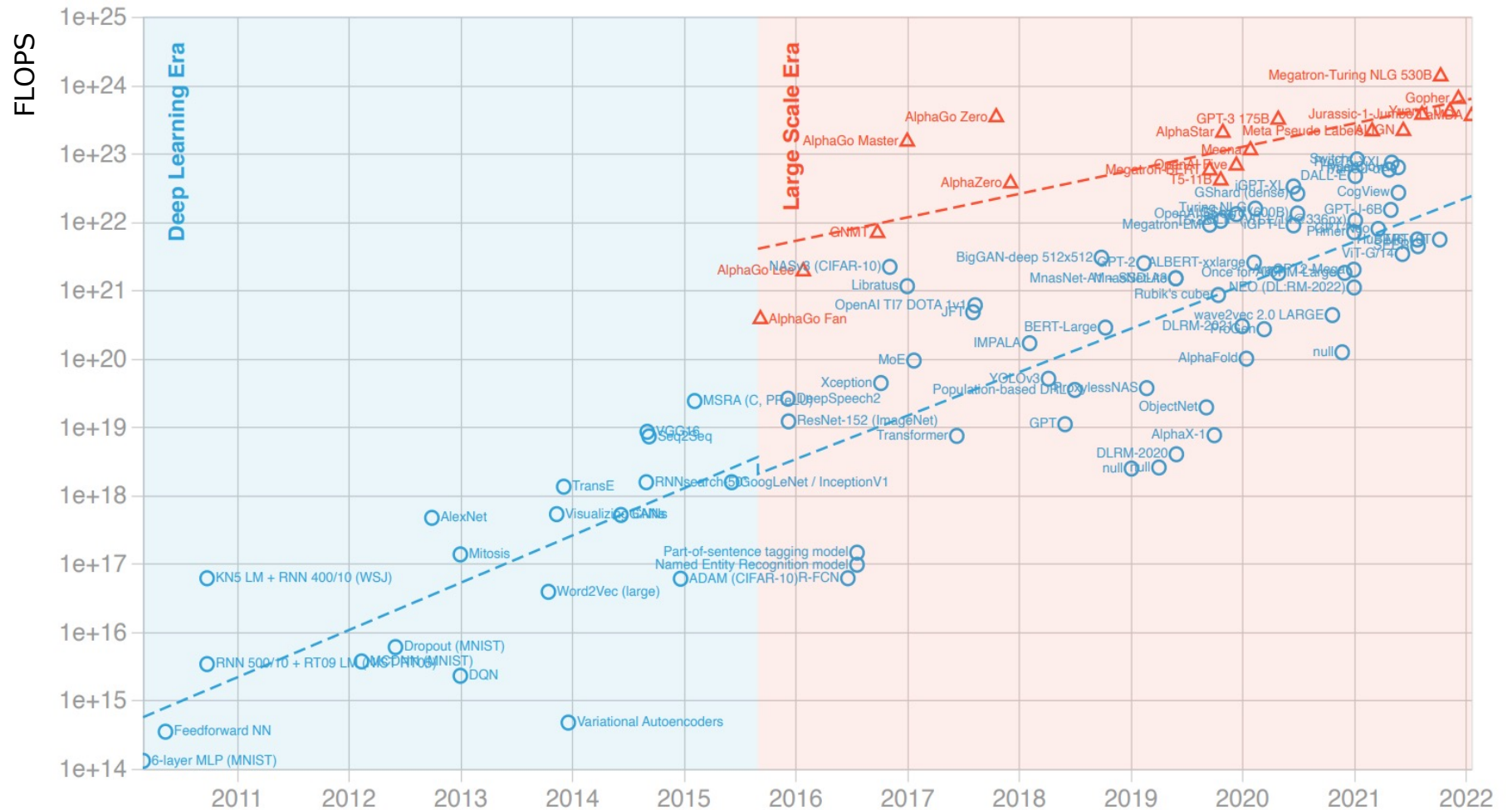


Context

- Deep Neural Networks

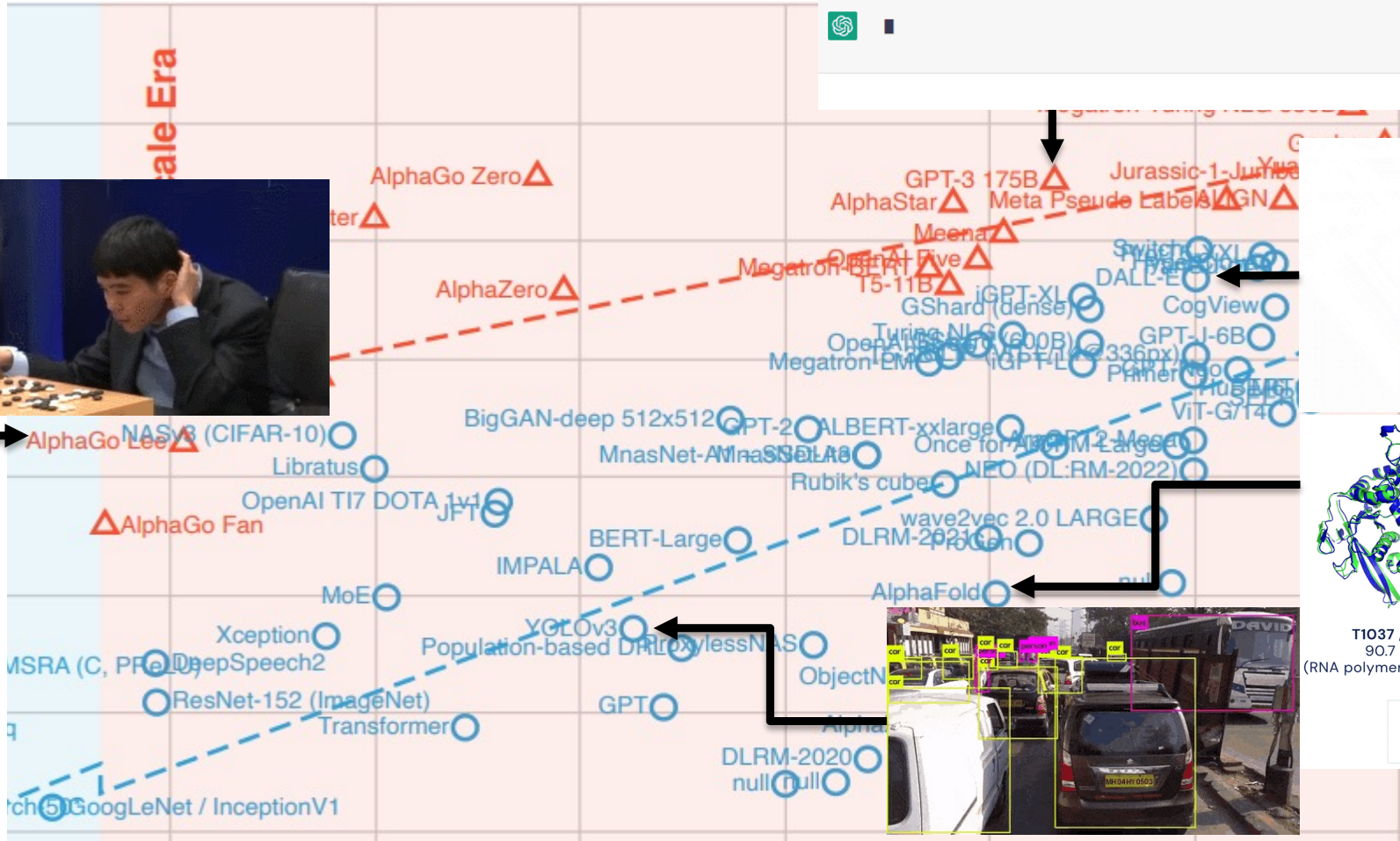
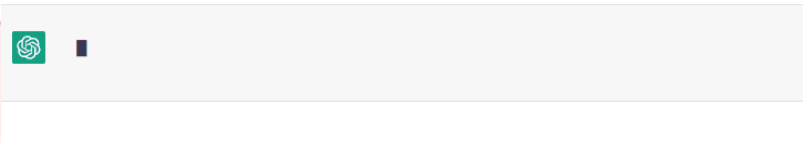


Context



Context

T Qu'est-ce que chatGPT ? Réponds-moi en deux phrases.

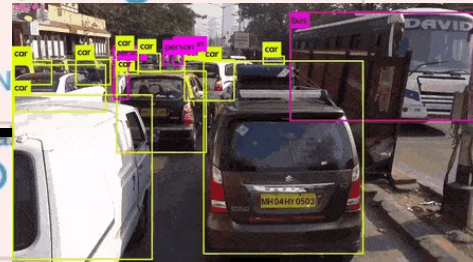


T1037 / 6vr4
90.7 GDT
(RNA polymerase domain)



T1049 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result
● Computational prediction



Why we talk about Trustworthy and Sustainable AI?

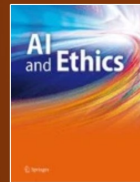
8

- **AI ethics**: “the study of ethical and societal issues facing developers, producers, consumers, citizens, policy makers, and civil society organizations.”

[Home](#) > [AI and Ethics](#) > [Article](#)

Sustainable AI: AI for sustainability and the sustainability of AI

Opinion Paper | [Open access](#) | [Published: 26 February 2021](#) | 1,213–218 (2021)



<https://link.springer.com/article/10.1007/s43681-021-00043-6>

Why we talk about Trustworthy and Sustainable AI?

9

- Waves:

1. Fanciful scenarios of robot uprisings

2. The problem of **explainability**



- The lack of equal representation in training data and the resulting biases in AI models (<https://www.theguardian.com/technology/2018/jan/12/google-racism-ban-gorilla-black-people>)
- Hardware malfunctions:
 - Intentional: Adv Attacks
 - Un-intentional: Hardware Faults

3. The **sustainable** development

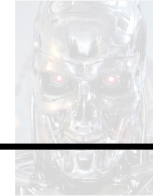
1. AlphaGo Zero generated 96 tonnes of CO2 over 40 days of research training which amounts to 1000 h of air travel or a carbon footprint of 23 American homes
2. Energy usage during ChatGPT's training has been estimated to be equivalent to that of an American household for over 700 years

Why we talk about Trustworthy and Sustainable AI?

10

- Waves:

1. Fanciful scenarios of robot uprisings



2. The problem of **explainability**

- The lack of equal representation in training data AI models (<https://www.theguardian.com/technology/2020/jan/12/google-racism-ban-gorilla-black-people>)
- Hardware malfunctions:
 - Intentional: Adv Attacks
 - Un-intentional: Hardware Faults



g biases in [1/12/google-](https://www.theguardian.com/technology/2020/jan/12/google-racism-ban-gorilla-black-people)

3. The **sustainable** development

1. AlphaGo Zero generated 96 tonnes of CO2 over 4 months, which amounts to 1000 h of air travel or a carbon footprint equivalent to 1000 American homes
2. Energy usage during ChatGPT's training has been estimated to be equivalent to that of an American household for over 700 years



h training
merican
equivalent

Outline

11

- Introduction
- **Efficient HW accelerators**
- Reliable HW accelerators
- Conclusions



State of the Art

Software

Application software

Virtual machines

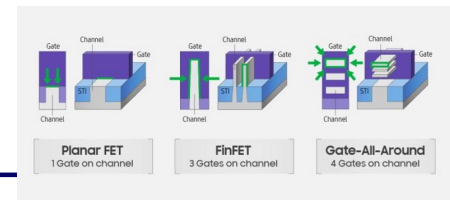
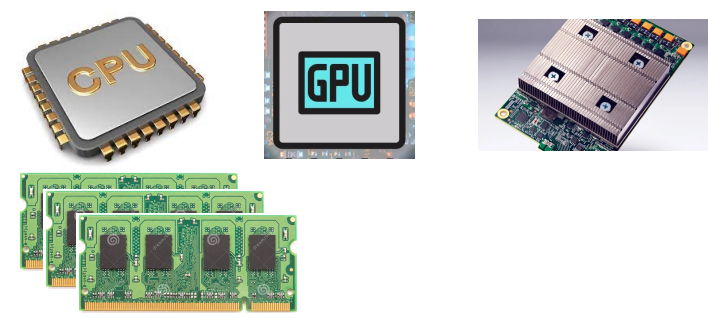
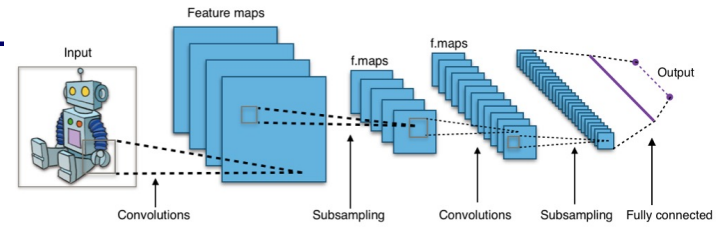
Operating system/Hypervisor

Firmware

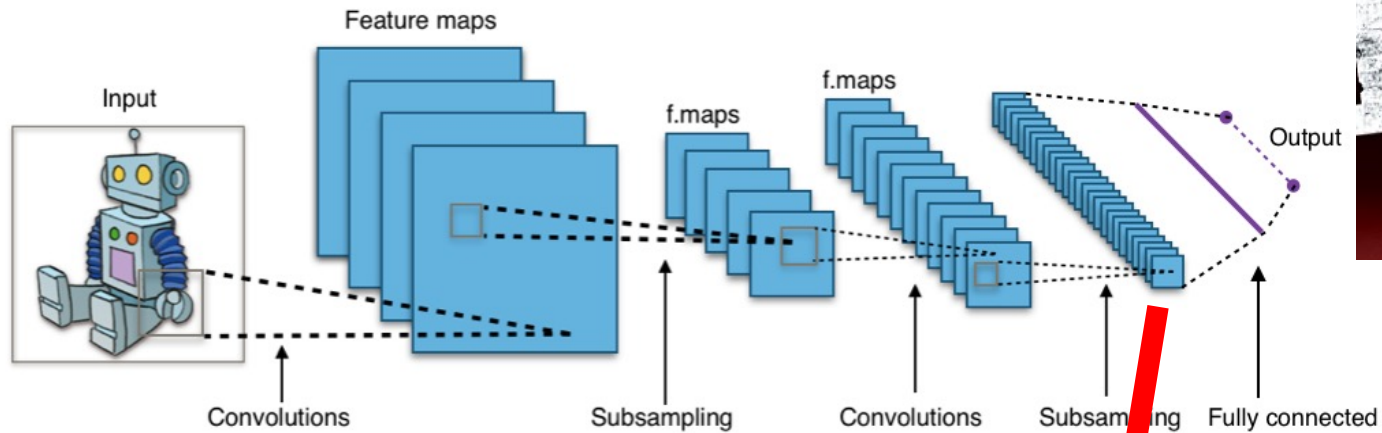
Hardware

Architecture

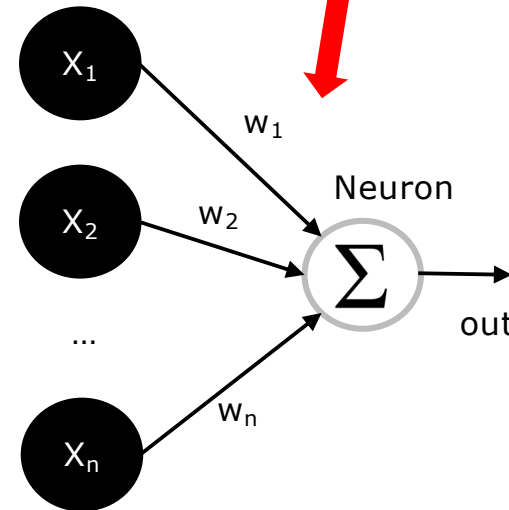
Technology

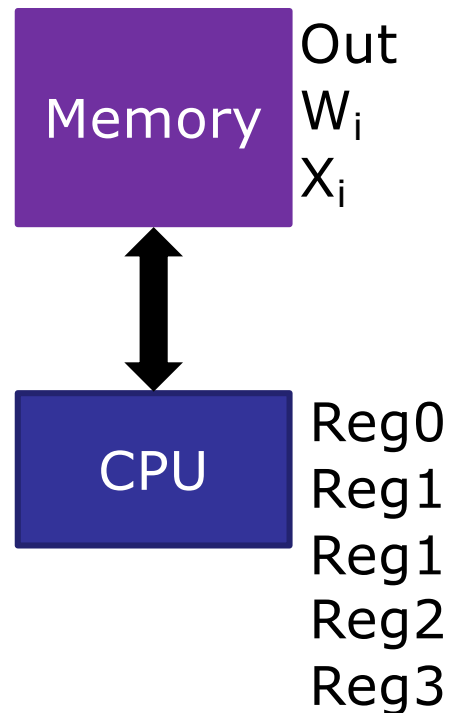


Why DNN are so complex?



```
for (I = 1 to N)  
  out += Wi * Xi
```





for (I = 1 to N)

```
out +=  $W_i$  *  $X_i$ 
```



assembly

```
Load  $W_i$ , reg0  
Load  $X_i$ , reg1  
Load out, reg2  
Mul reg3, reg1, reg0  
Add reg2, reg3, reg2  
Store reg2, out
```

6 memory accesses for instructions
4 memory accesses for data
2 operations

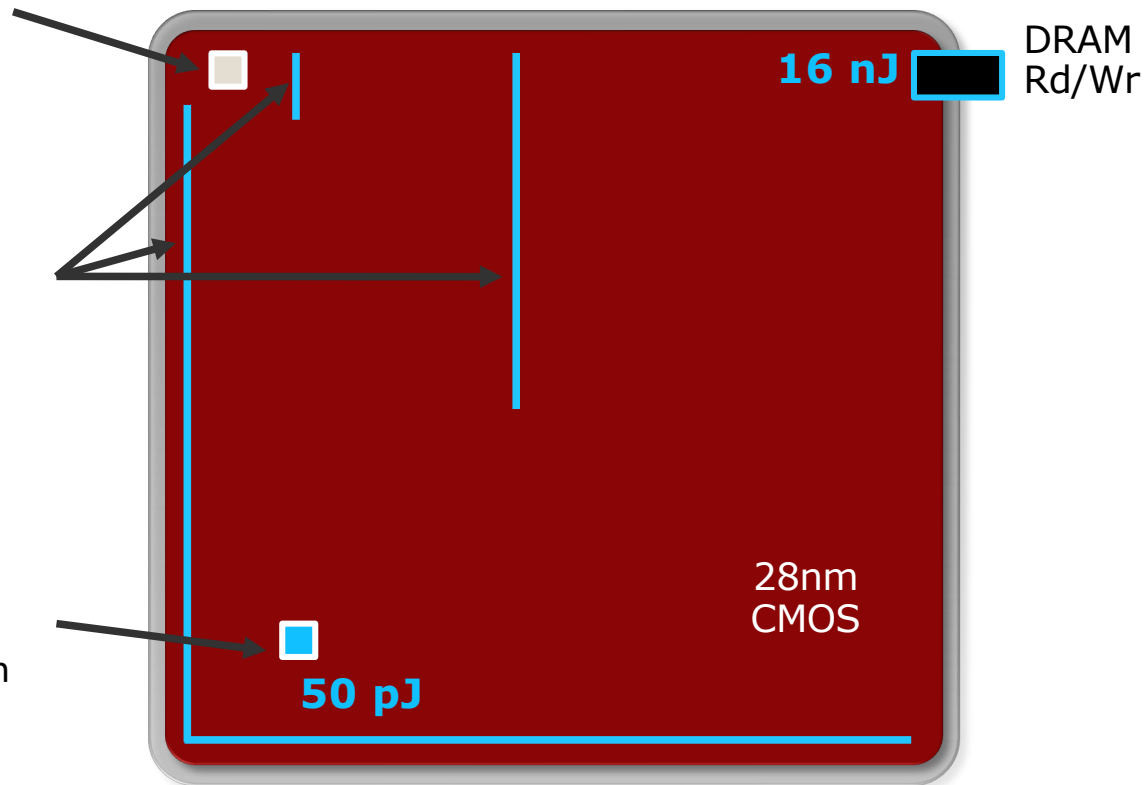
Energy Cost in a Processor

[Adapted from Dally, IPDPS'11]

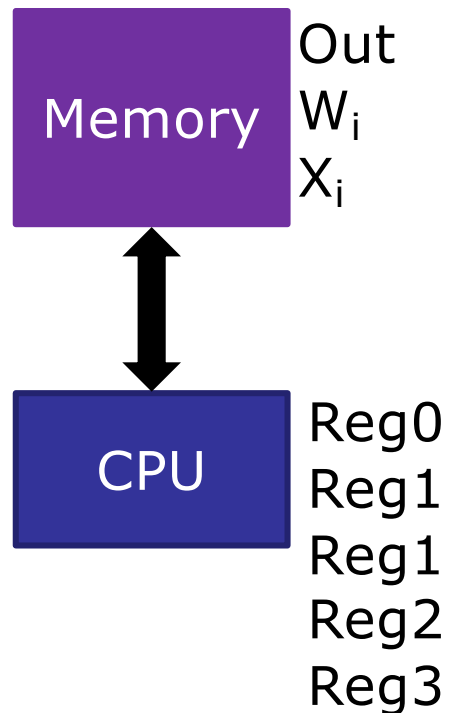
- 64-bit FPU: 20pJ/op
- 32-bit addition: 0.05pJ
- 16-bit multiply: 0.25pJ

- Wire energy
 - 32 bits: 40pJ/word/mm
 - 8 bits: 10pJ/word/mm

- Register-File
 - Depends on word-length



Courtesy of O. Sentieys



```
for (I = 1 to N)
```

```
  out +=  $W_i$  *  $X_i$ 
```



assembly

```
nJ Load  $W_i$ , reg0  
nJ Load  $X_i$ , reg1  
nJ Load out, reg2  
pJ Mul reg3, reg1, reg0  
pJ Add reg2, reg3, reg2  
nJ Store reg2, out
```

6 memory accesses for instructions
4 memory accesses for data
2 operations

We can target any abstraction level!

Software

Application software

Virtual machines

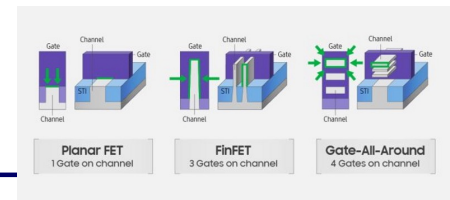
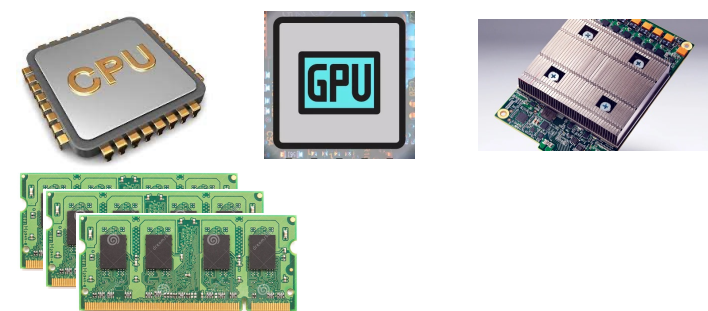
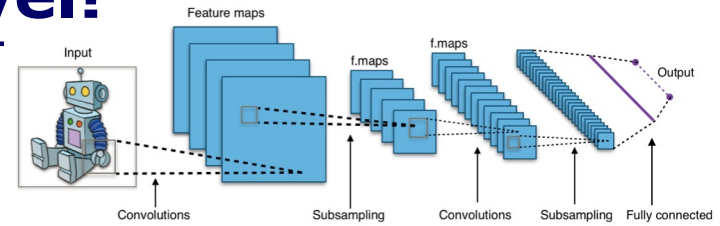
Operating system/Hypervisor

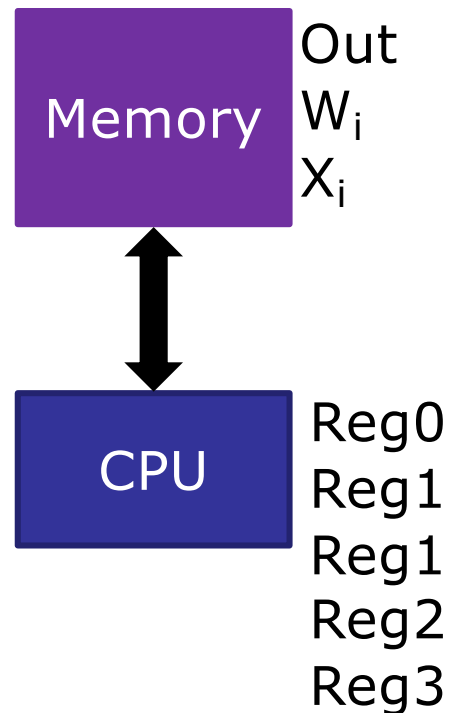
Firmware

Hardware

Architecture

Technology





```
for (I = 1 to N)
```

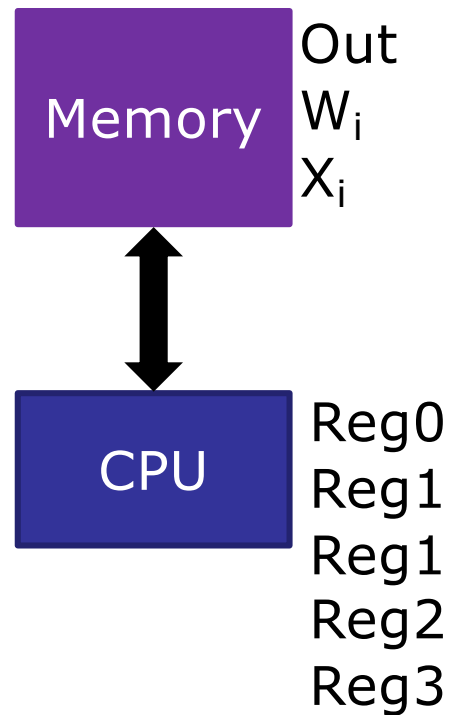
```
  out +=  $W_i$  *  $X_i$ 
```



assembly

```
Load  $W_i$ , reg0  
Load  $X_i$ , reg1  
Load out, reg2  
Mul reg3, reg1, reg0  
Add reg2, reg3, reg2  
Store reg2, out
```

6 memory accesses for instructions
4 memory accesses for data
2 operations



```
for (I = 1 to N)
```

```
  out +=  $W_i$  *  $X_i$ 
```



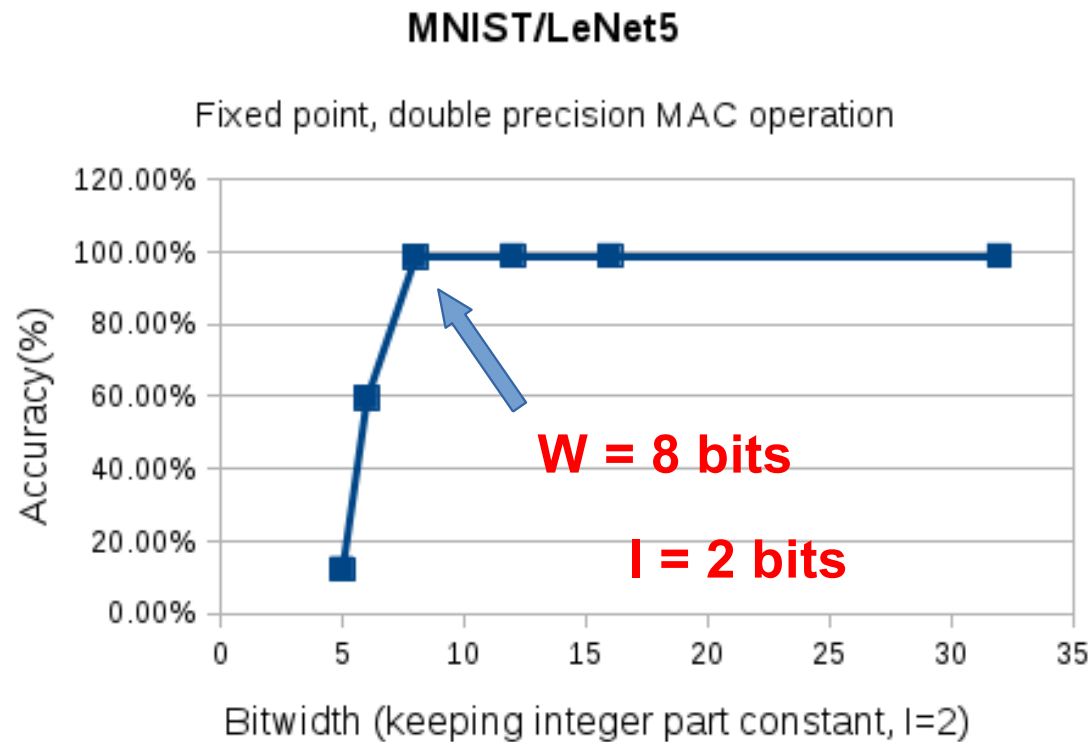
assembly

```
Load  $W_i$ , reg0  
Load  $X_i$ , reg1  
Load out, reg2  
Mul reg3, reg1, reg0  
Add reg2, reg3, reg2  
Store reg2, out
```

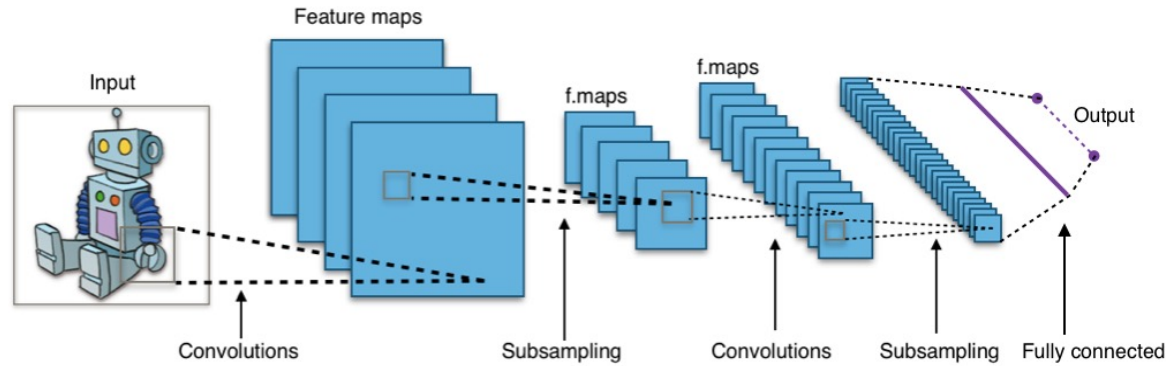
6 4 memory accesses for instructions
4 2 memory accesses for data
2 operations

Quantization

- 10k images, MNIST/LeNet-5
- Fixed-Point Arithmetic



Weight Sharing



5x5 Convolutional Kernel

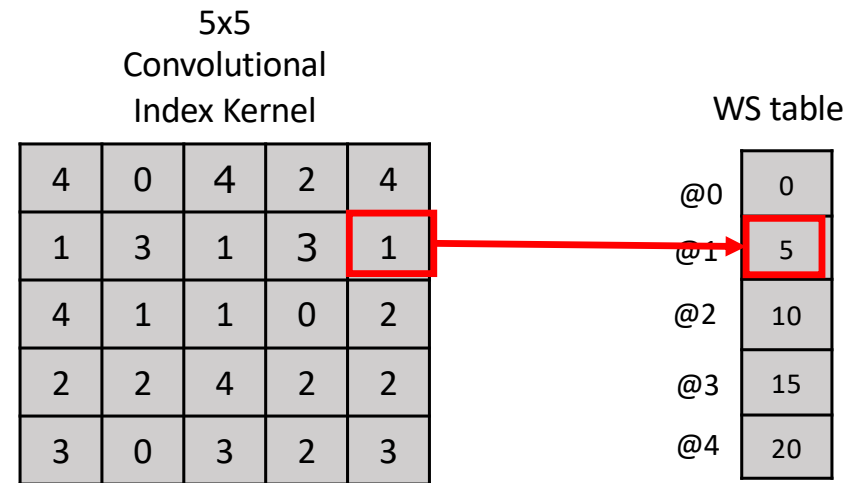
K-means Clustering

- 8 bits x W
- 200 bits

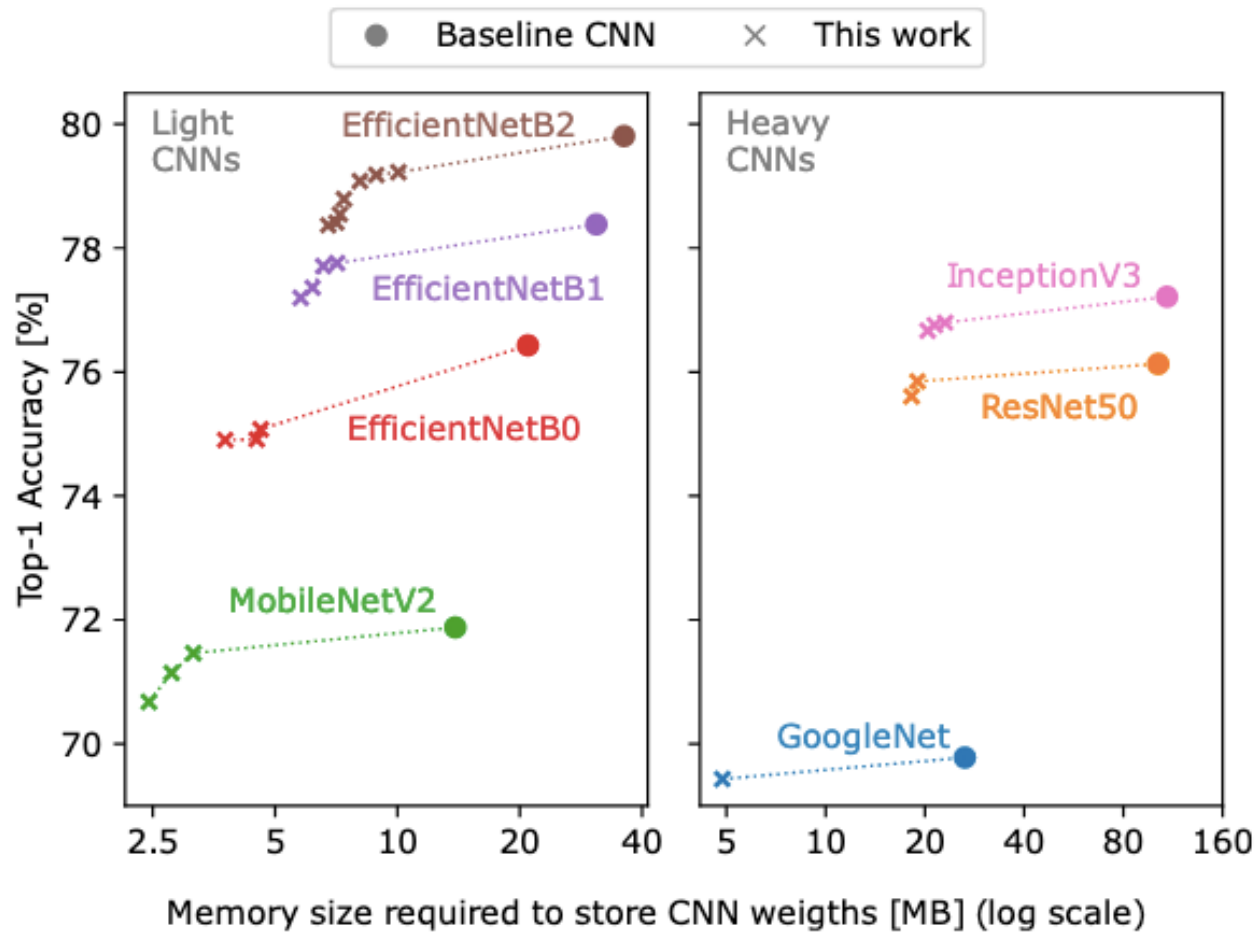
17	3	18	11	19
5	14	9	13	7
19	5	7	0	12
10	12	20	8	10
14	1	16	10	14

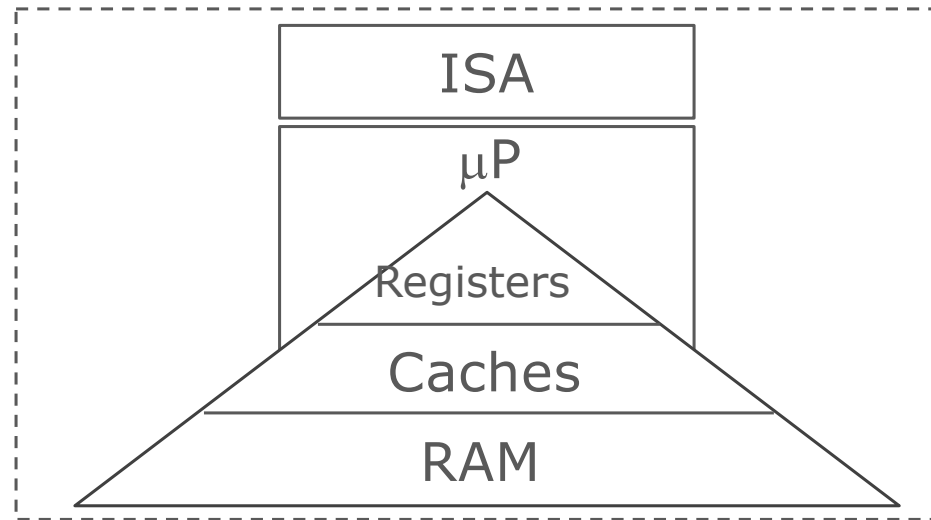


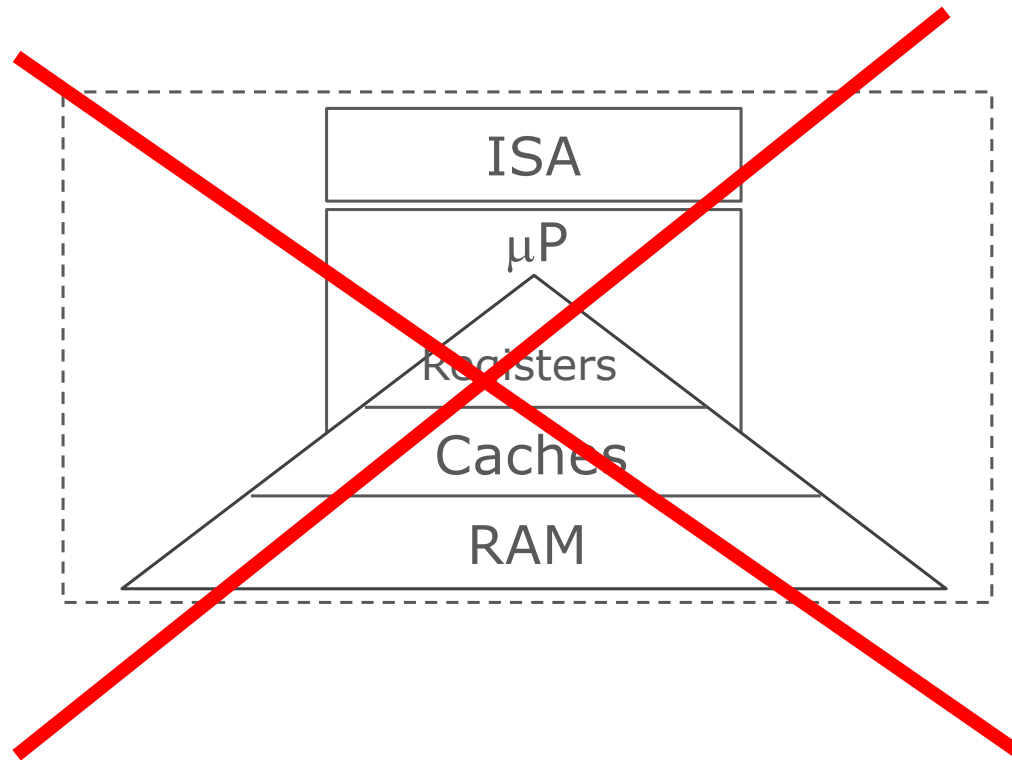
17	3	18	11	19
5	14	9	13	7
19	5	7	0	12
10	12	20	8	10
14	1	16	10	14



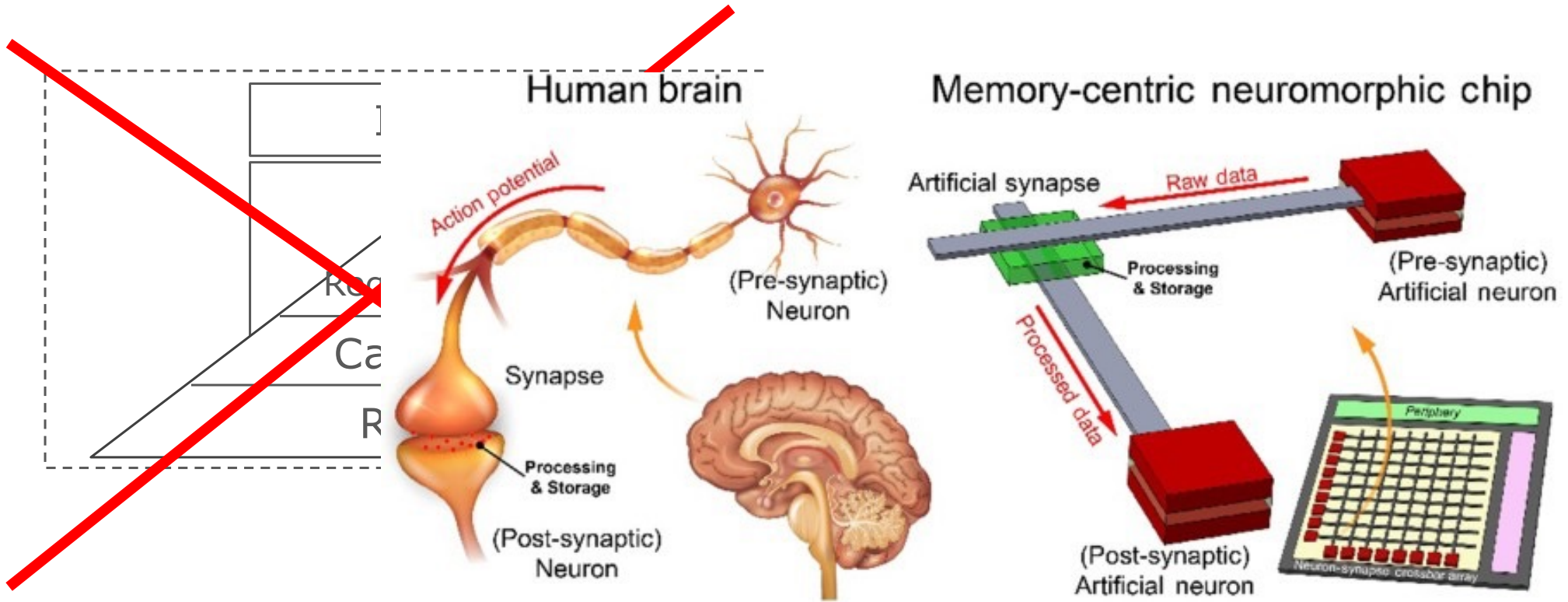
- 3 bits x W
 - 75 + 40 = 115 bits (instead of 200)
- ~42% bits reduction



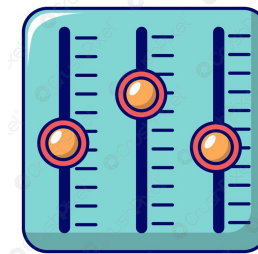
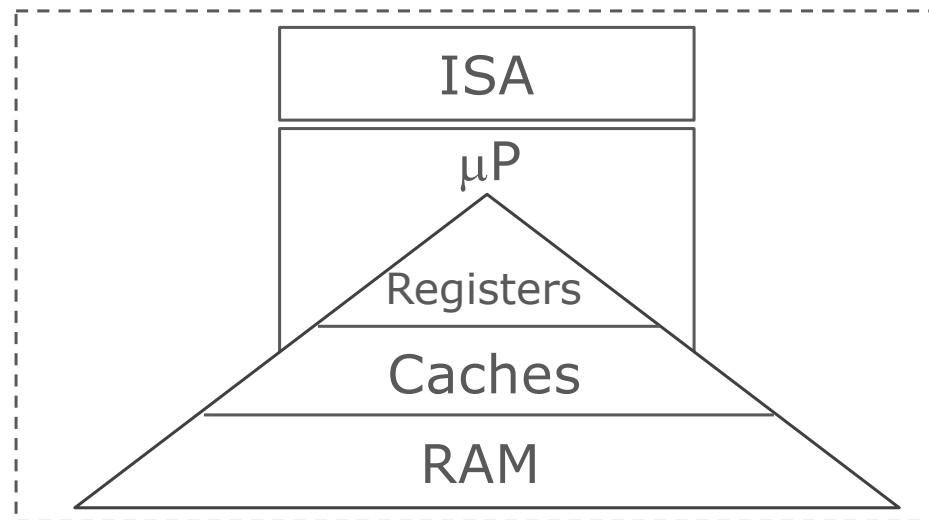




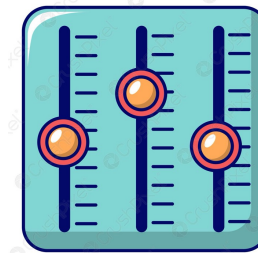
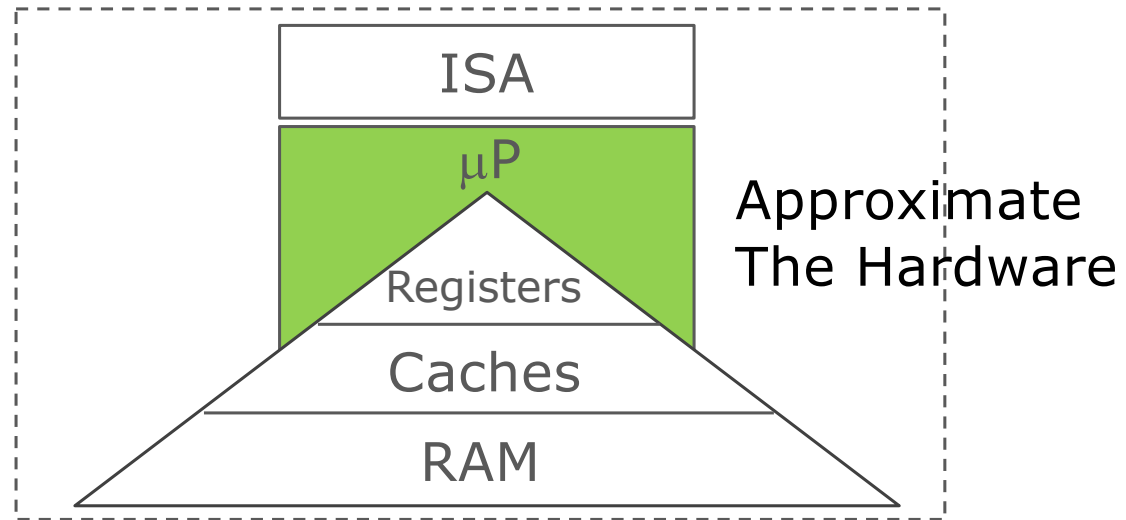
Hardware Level



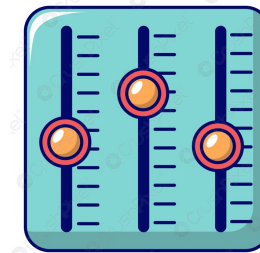
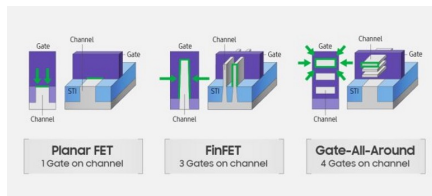
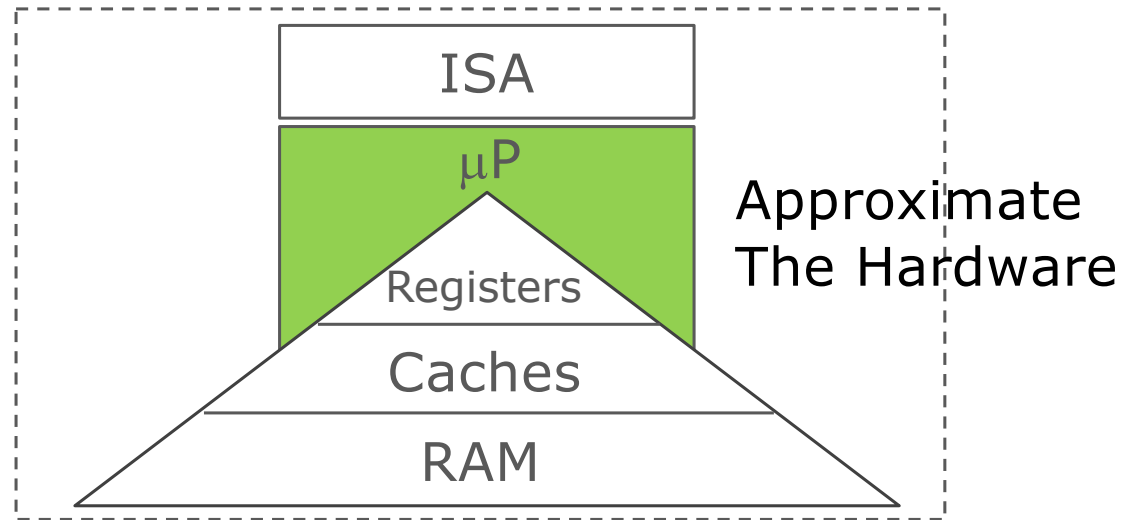
<https://link.springer.com/article/10.1007/s12274-021-3452-6>



Tune Vdd



Tune Vdd

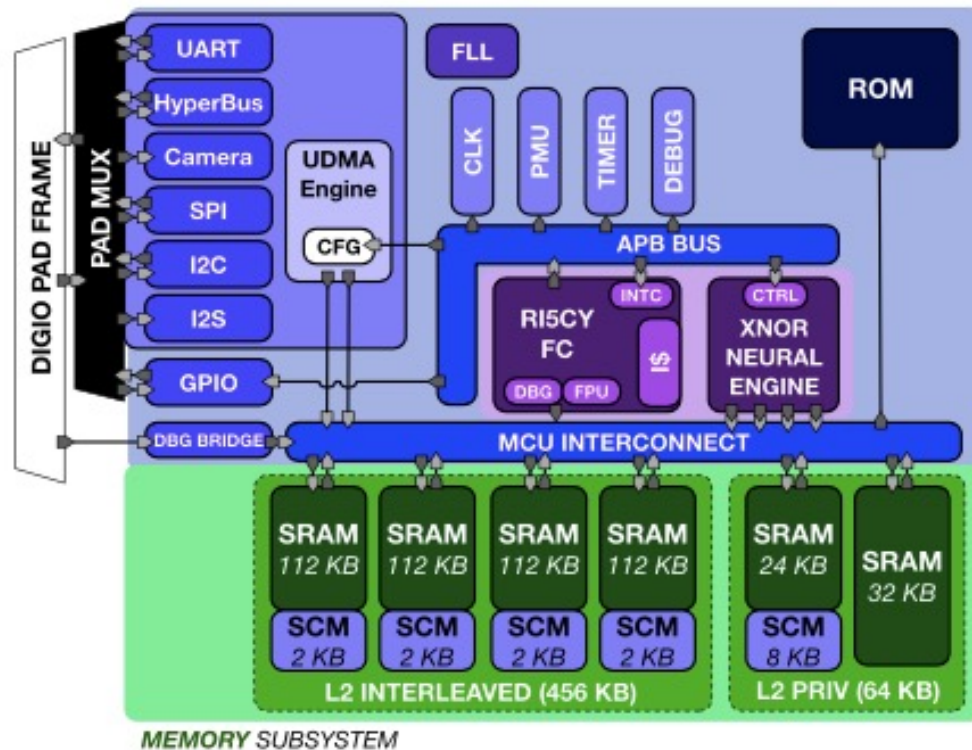


Tune Vdd

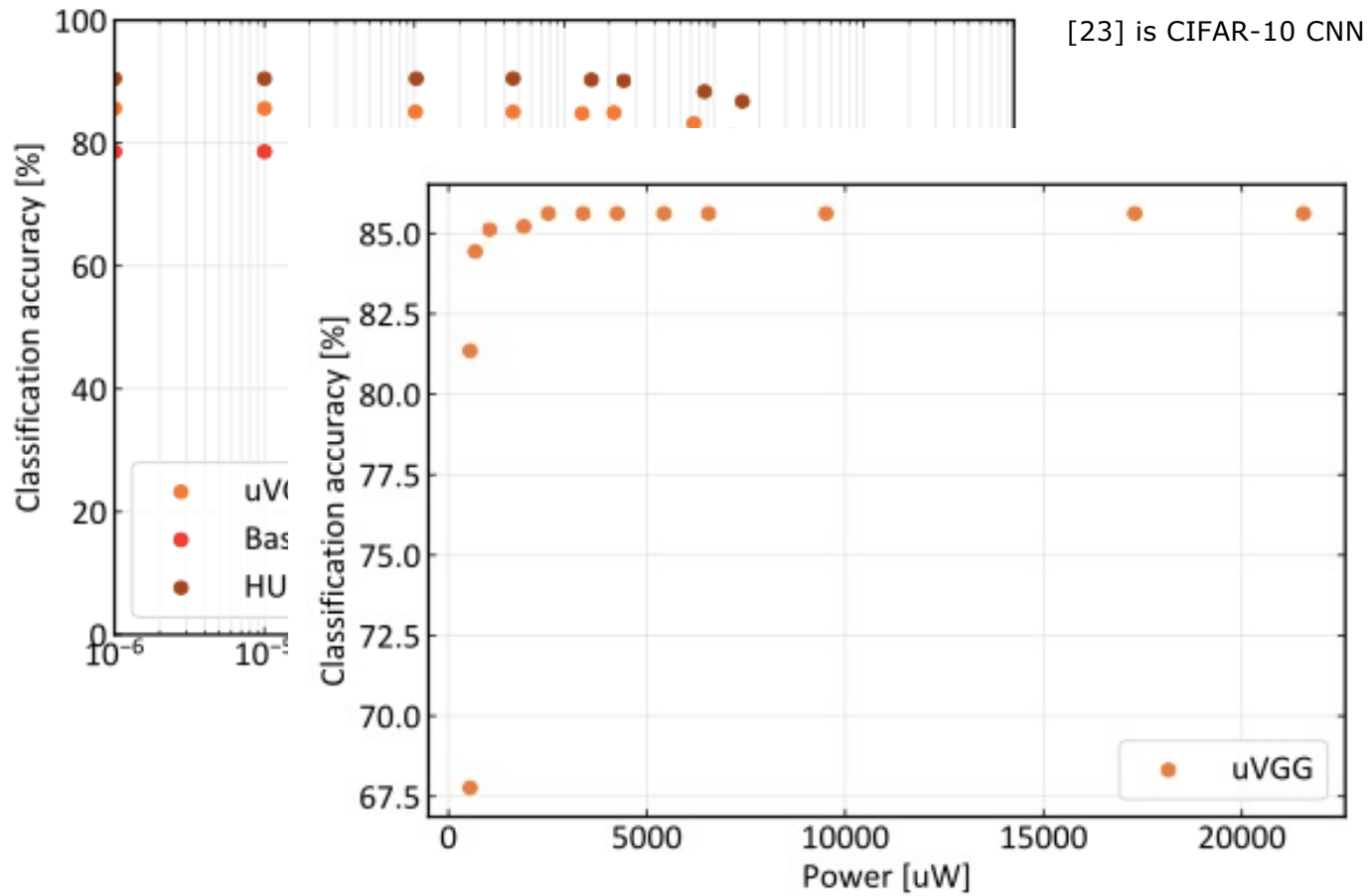
Improve the technology

Some examples: Over-Scaling

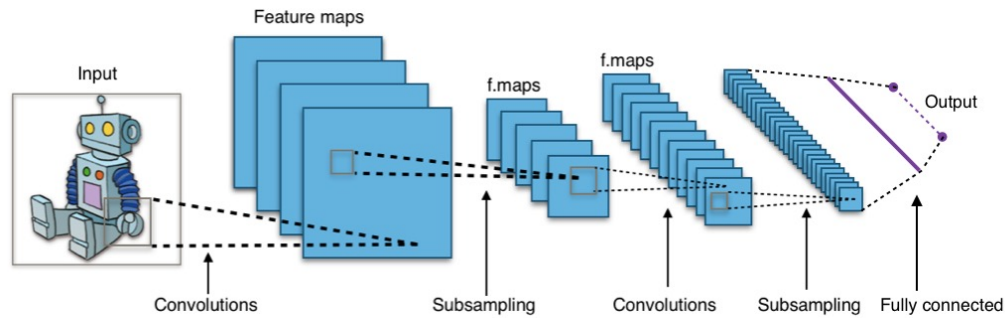
- Quentin SoC (based on PULPissimo system)



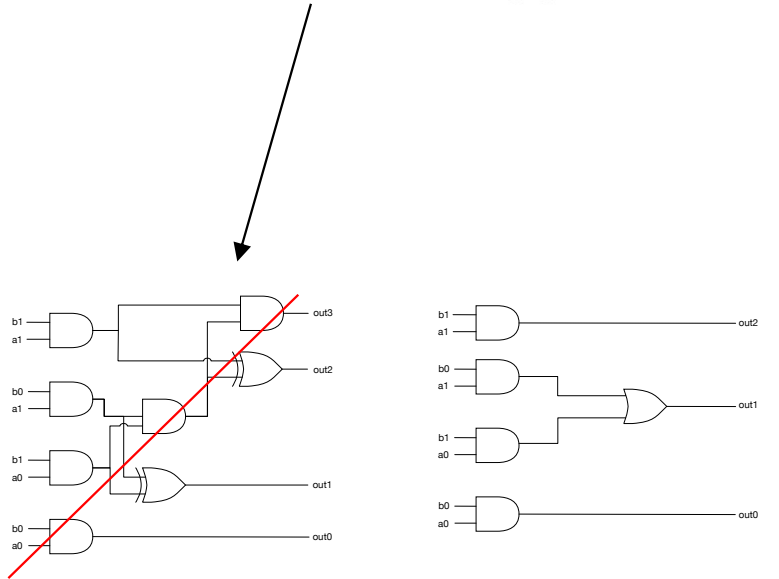
Over-Scaling:
reduce V_{dd} of SRAM



Functional Approximation

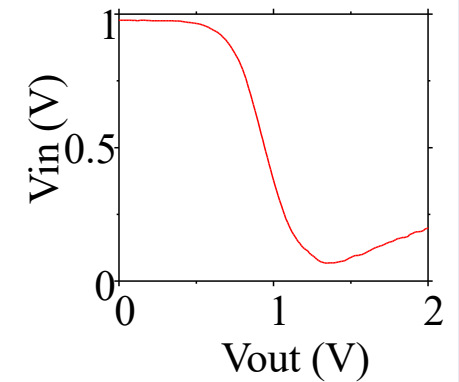
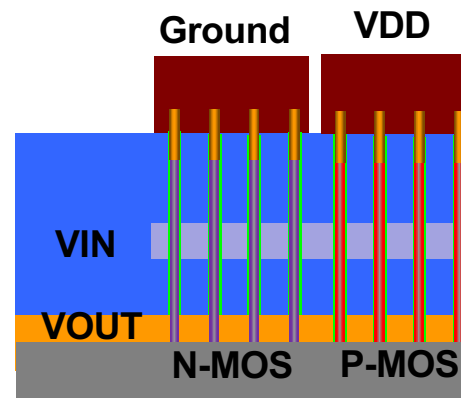
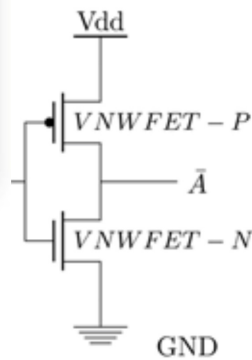
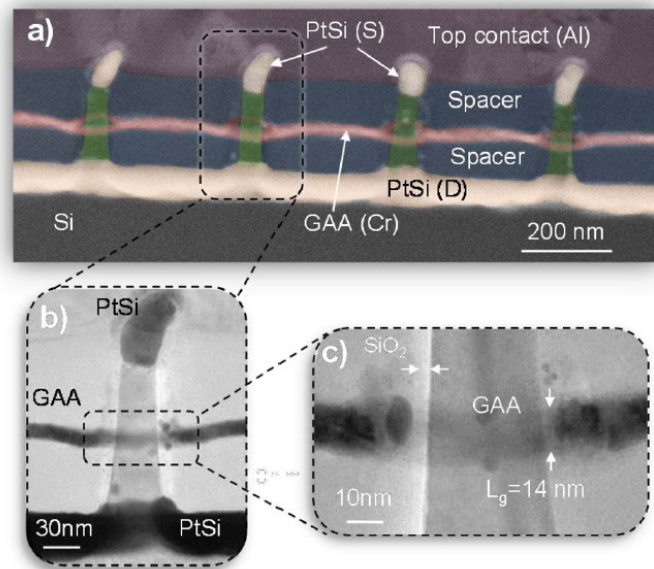


Multiplier	ER	Accuracy (%)	
		MNIST	SVHN
Exact	0	97.69	86.93
mul8-350	99.0	97.70	87.00
mul8-439	97.8	97.71	86.96
mul8-120	98.5	97.70	87.00
mul8-183	97.2	97.70	86.98
mul8-134	93.9	97.72	86.95



- Up to 71.45% more energy-efficient
- Up to 61.55% smaller

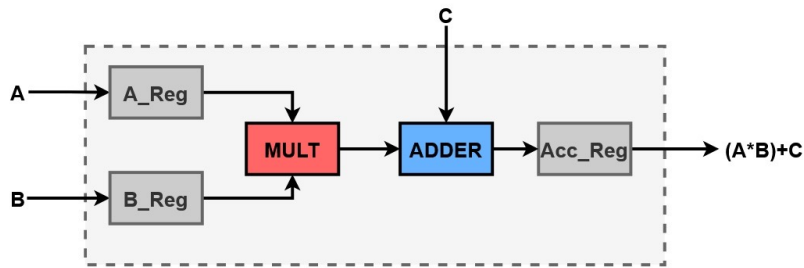
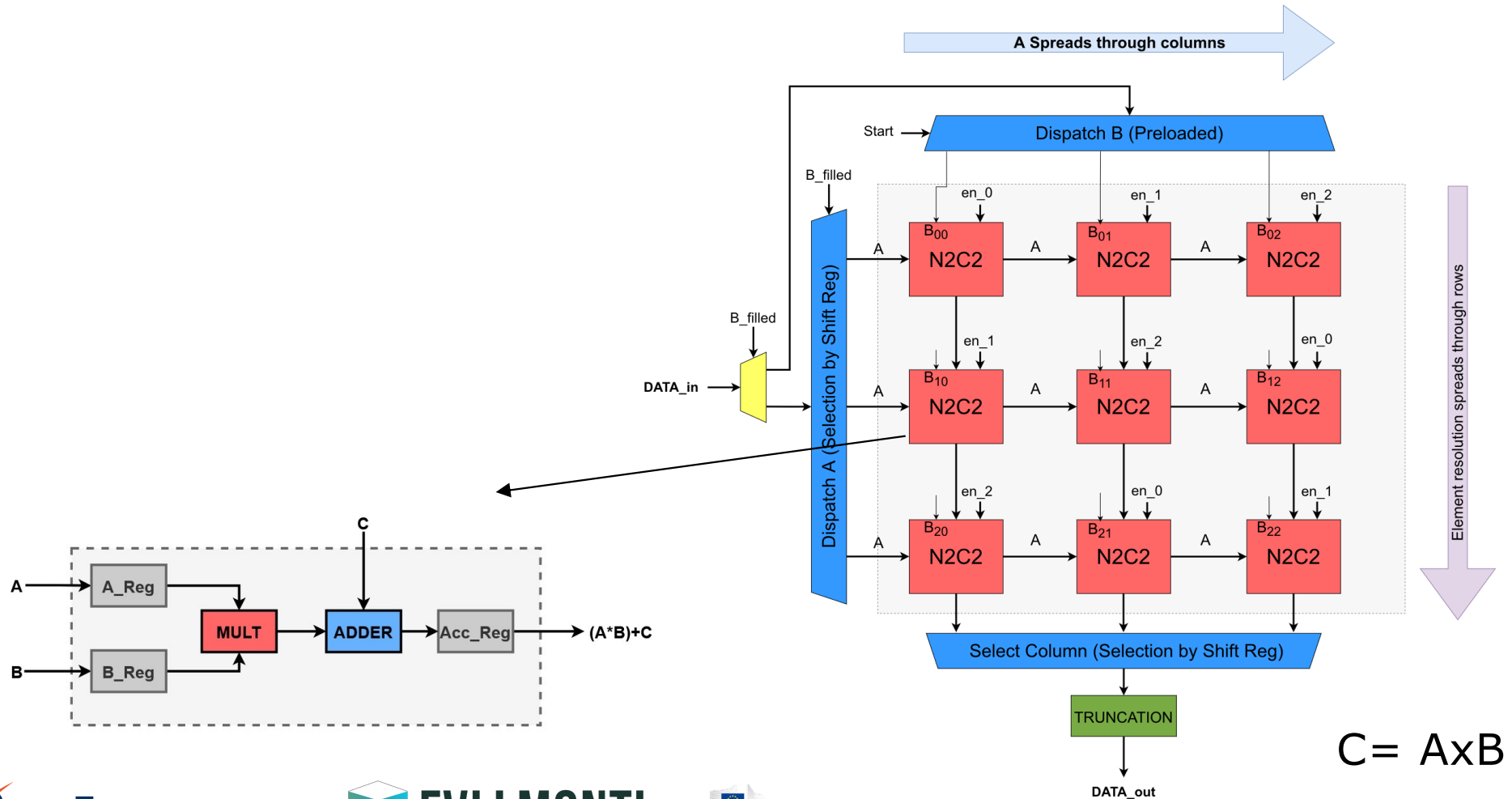
First hardware demo of 14 nm VGAA inverter



Y. Guerfi et al., Nanoscale research letters 11, 210 (2016)

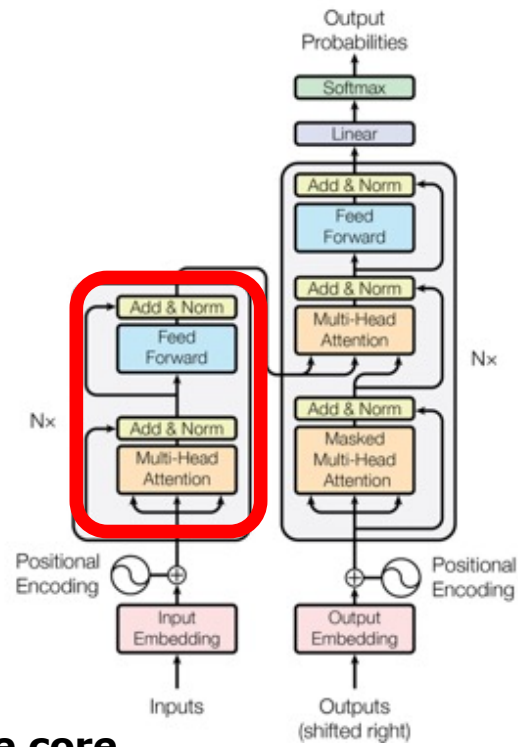
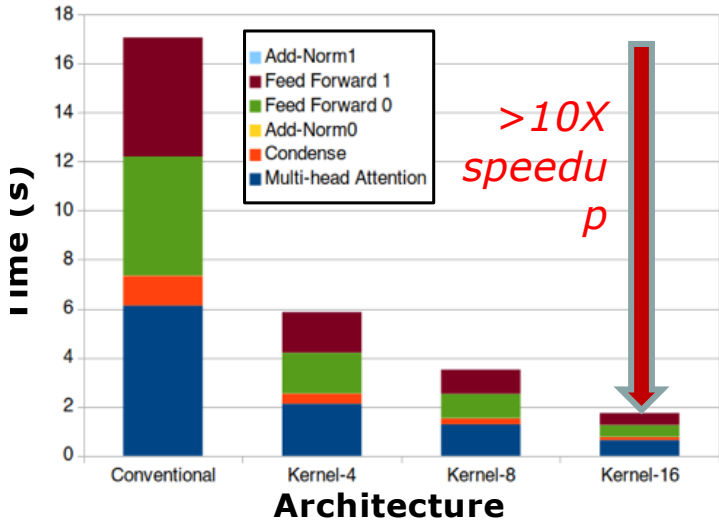
G. Larrieu et al., Solid-State Electronics 130, 9 (2017)

Systolic Array*



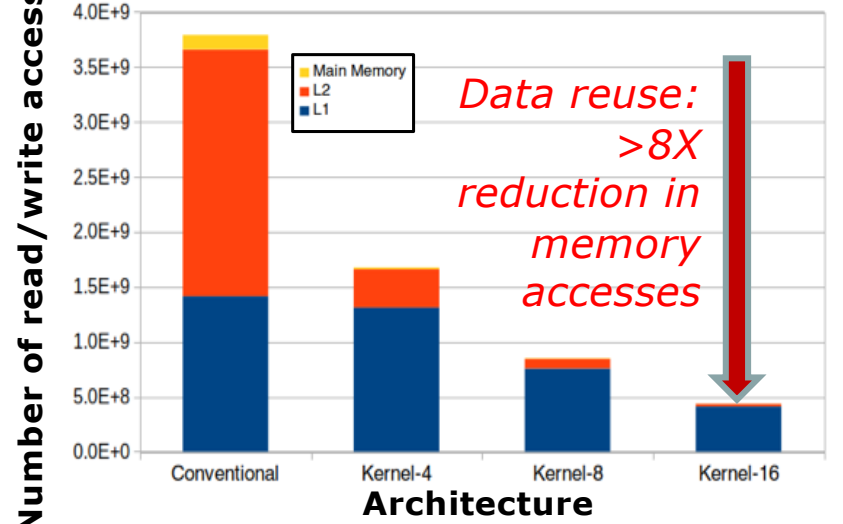
Performances

Inference time of a transformer block



Number of read/write accesses

Number of memory accesses



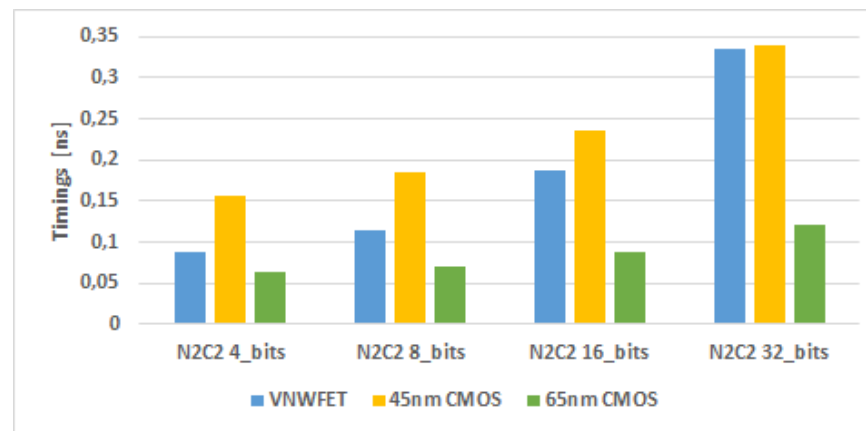
Kernel-N = NxN N²C² accelerator

Conventional =
 Non-accelerated ARM ISA - single core
 Clock frequency: 1GHz - Memory type:
 DDR4_2400_4x16
 Cache L1 size: 32KB - Cache L2 size: 1MB

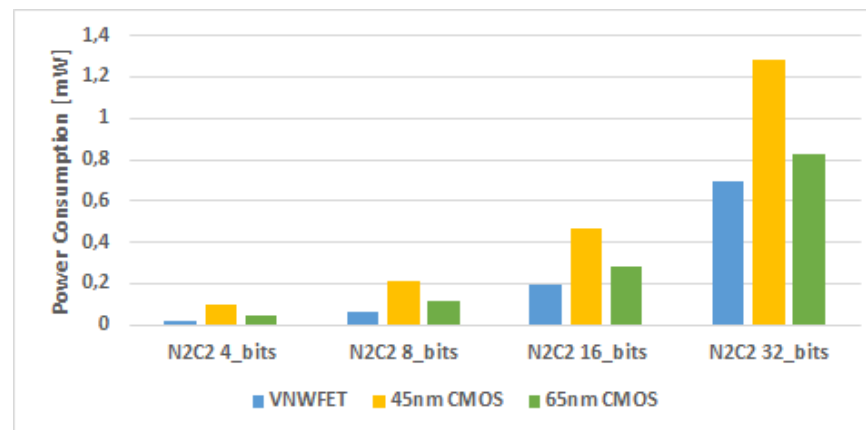


Tech library and comparison

Library	INV	NAND	NOR	XOR
VNWFET	4	3	3	3
45nm CMOS	4	1	1	1
65nm CMOS	20	15	15	10



Circuit	NVWFET	45nm CMOS	65nm CMOS
4-bits PE	230	530	221
8-bits PE	737	1674	705
16-bits PE	2553	5697	2321
32-bits PE	9395	20608	8492

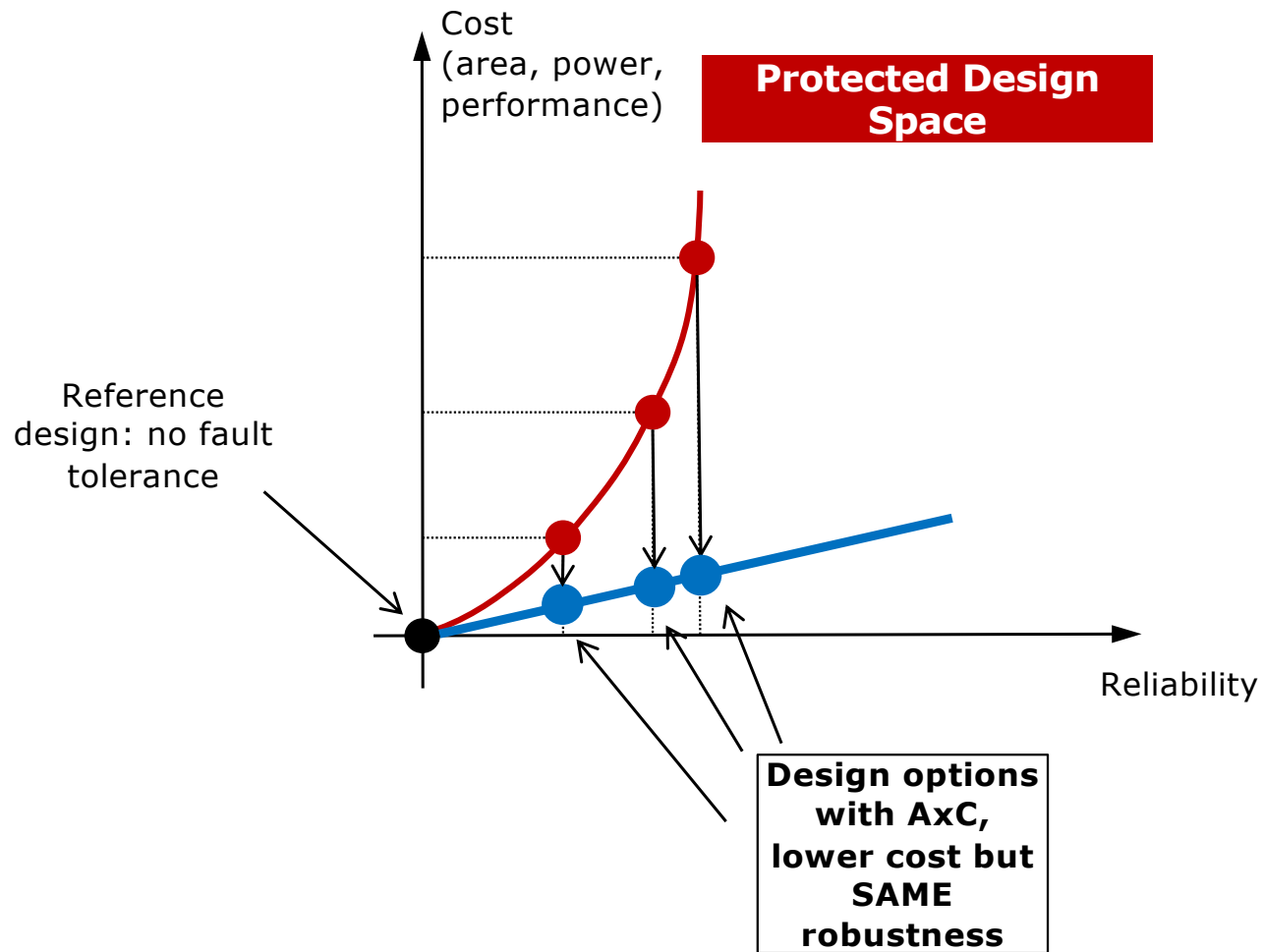


Outline

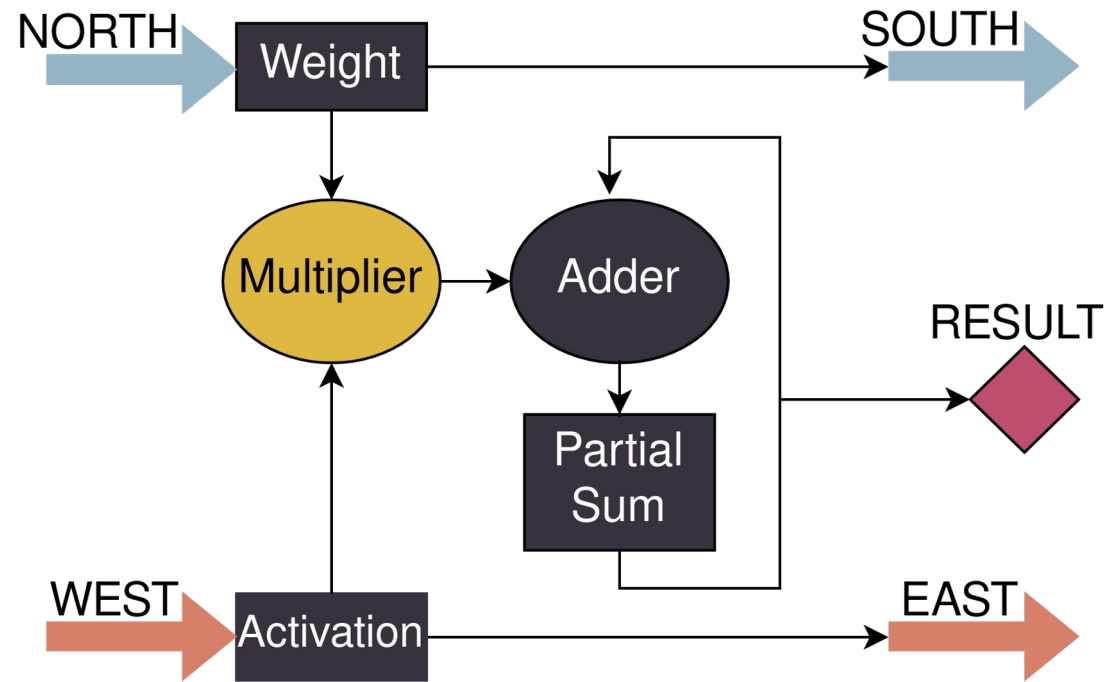
- Introduction
- Efficient HW accelerators
- **Reliable HW accelerators**
- Conclusions



Approximation is the key

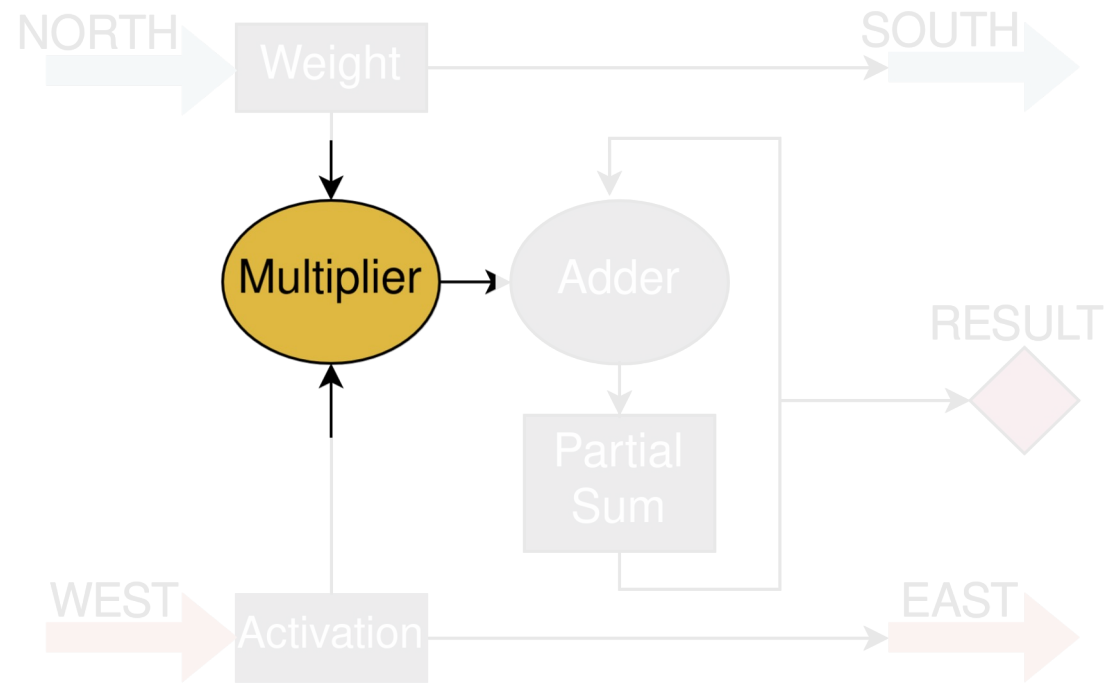


Approximate Systolic Array



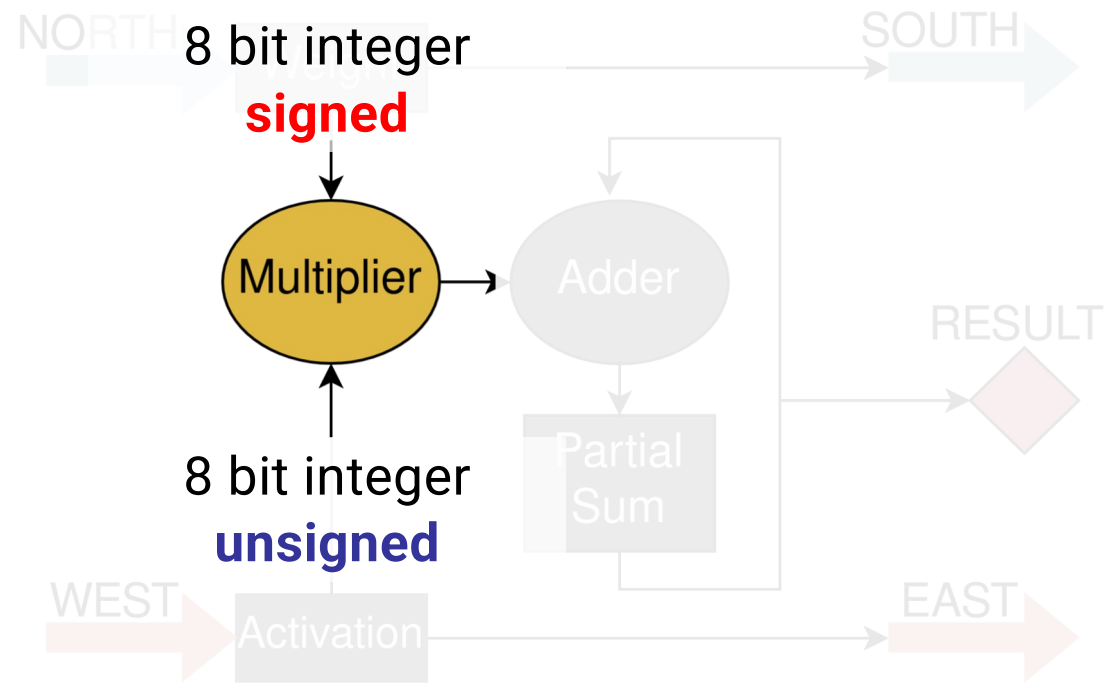
Processing Element Functional View

Approximate Systolic Array



Processing Element Functional View

Approximate Systolic Array



Processing Element Functional View

Used AxC Multipliers

Name	MAE	MAE-8
mul12s_2PT	0.000073	0.019
mul12s_2QH	0.0031	0.134
mul12s_2R5	0.0092	0.315
mul12s_34P	0.032	0.785
mul12s_2TE	0.19	6.080

MAE: Mean Absolute Error



https://ehw.fit.vutbr.cz/evoapproxlib/?folder=multiplers/12x12_signed

Institut des Nanotechnologies de Lyon UMR CNRS 5270

<http://inl.cnrs.fr>

Results

#	Bit-width	Multiplier	Energy Reduction [%]	Accuracy [%]	Injected Faults [%]	Masked [%]	Tolerable [%]	Critical [%]
Baseline	16	Precise	–	99.07%	10%	47.58%	29,18%	23.24%
1	8	Precise	50%	99.05%	19%	64.65%	23.01%	12.34%
2	8	mul12s_2PT	50.3%	99.08%	19%	63.9%	23.79%	12.3%
3	8	mul12s_2QH	51.21%	99.1%	19%	38.92%	44.85%	16.23%
4	8	mul12s_2R5	52%	99.06%	19%	26.96%	55.69%	17.34%
5	8	mul12s_34P	55%	98.24%	19%	74.16%	23.01%	2.83%
6	8	mul12s_2TE	55.6%	9.8%	19%	3.94%	27.5%	68.76%

Conclusions & Future works

44

- We need a holistic approach to achieve a HW-SW co-design methodology to design **sober** and **reliable** AI applications
- How to reach this goal?



Still a lot of work to reach a sober system

45



We have to build a novel flow

