

Christof Teuscher

ECE 410/510: Hardware for AI and ML

## Hardware for AI and ML Overview + Co-design

Portland State University  
Department of Electrical and Computer Engineering (ECE)

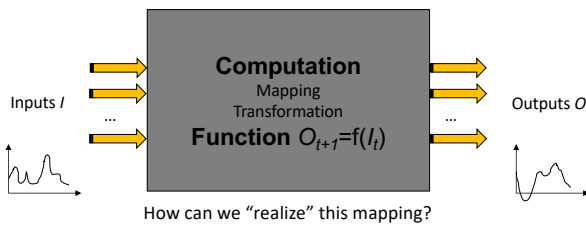
[www.teuscher-lab.com](http://www.teuscher-lab.com)  
[teuscher@pdx.edu](mailto:teuscher@pdx.edu)



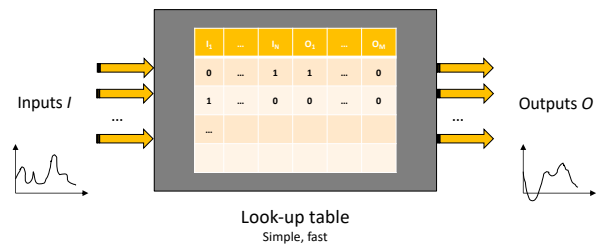
## What is a computer?

### What's the job of a computer architect / chip designer?

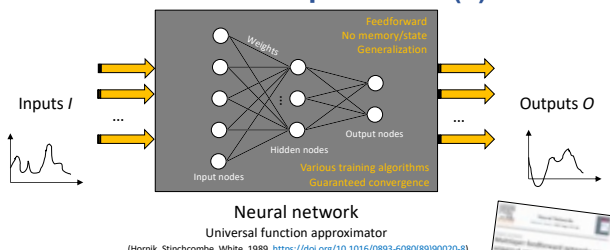
## What is computation? (1)



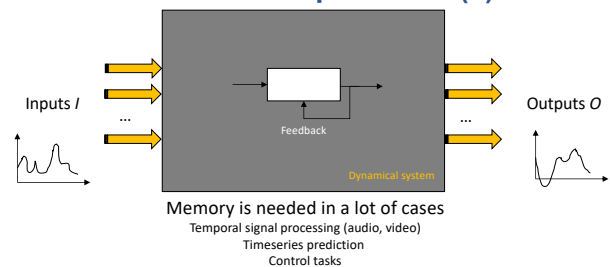
## What is computation? (2)



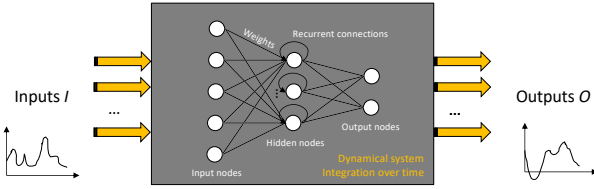
## What is computation? (3)



## What is computation? (4)



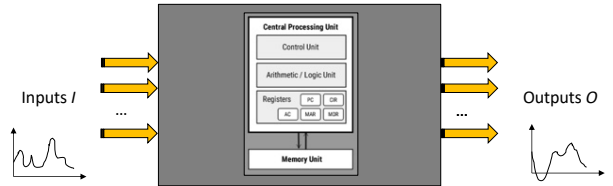
## What is computation? (5)



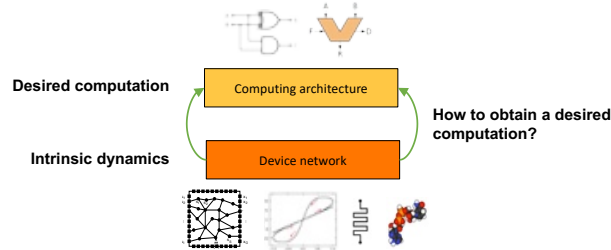
### Recurrent neural networks (RNN)

Training is challenging, various algorithms exist (backpropagation through time), convergence not guaranteed, many weights, slow learning. RNNs are also universal function approximators (Schaefer & Zimmerman, 2007, <https://doi.org/10.1142/S0129065707001111>)

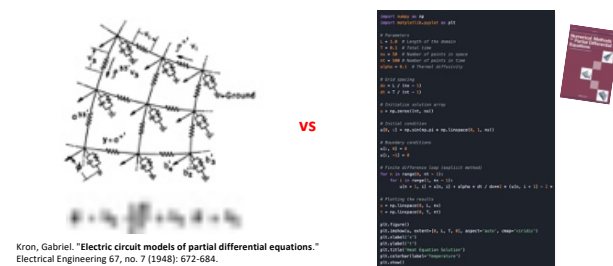
## What is computation? (6)



## Intrinsic vs designed computation (2)



## Intrinsic vs designed computation (3)



## Intrinsic vs designed computation (4)

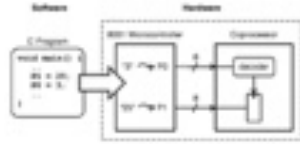


## What are the design trade-offs?

- **Performance**
- **Flexibility** (general purpose vs application-specific)
- **Energy efficiency**
- **Power density** (avoid burning chip, cost of cooling technology)
- **Design complexity**
- **Design cost**
- **Scalability**
- **Shrinking design schedules** (time-to-market)

## What is (HW/SW) co-design?

- The practice of taking the **best from software** design and the **best from hardware** design to solve design problems.
- HW/SW co-design deals with HW/SW **interfaces**.
- HW/SW co-design is the design of cooperating HW and SW components in a **single design effort**.
- Why design new HW if there is already one available to do the job? E.g., a RISC processor?
- From a flexibility and design effort perspective, as much as possible should be in software. [Easy, flexible, available libraries, etc.]



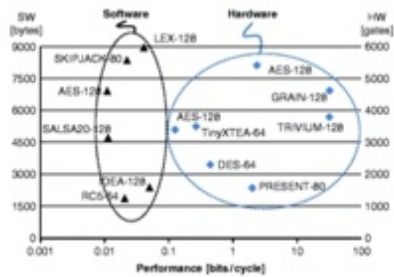
Schaumont, 2013

## Why is HW/SW co-design important for AI/ML

- Performance optimization:** Traditional hardware design assumes fixed software, while traditional software design assumes fixed hardware. Co-design breaks this barrier by simultaneously optimizing both, leading to significant performance improvements for AI workloads.
- Specialized acceleration:** AI algorithms have unique computation patterns (like matrix multiplications and convolutions) that benefit from custom hardware accelerators tailored to these specific operations.
- Memory bottlenecks:** AI models face severe memory bandwidth constraints. Co-designing hardware memory hierarchies with software that efficiently schedules operations can minimize data movement and maximize throughput.
- Energy efficiency:** By tailoring hardware precisely to the needs of ML workloads and optimizing software to take advantage of hardware capabilities, co-design dramatically reduces power consumption, which is critical for both data centers and edge devices.

## Cryptography on small embedded platforms

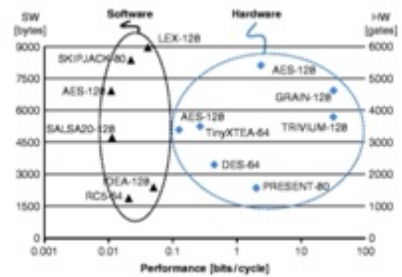
What do we learn here?



Schaumont, 2013

## Cryptography on small embedded platforms

Hardware crypto-architectures have, on the average, a higher relative performance compared to embedded processors.

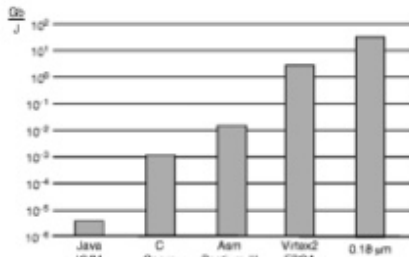


Schaumont, 2013

## Cryptography on small embedded platforms

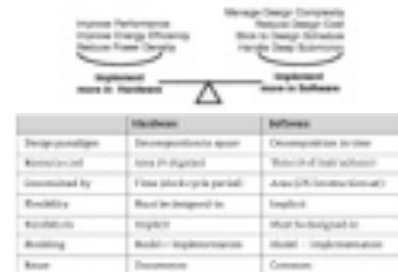
AES for different target platforms

What do we learn here?



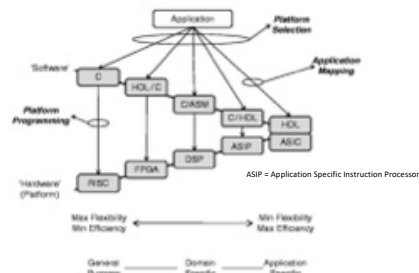
Schaumont, 2013

## Driving factors in HW/SW co-design + dualism



Schaumont, 2013

## The HW/SW co-design space

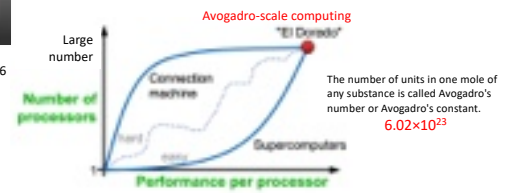


Schaumont, 2013

## The goals and drivers have changed



CM 1 & 2: Up to 65,536 1-bit processors

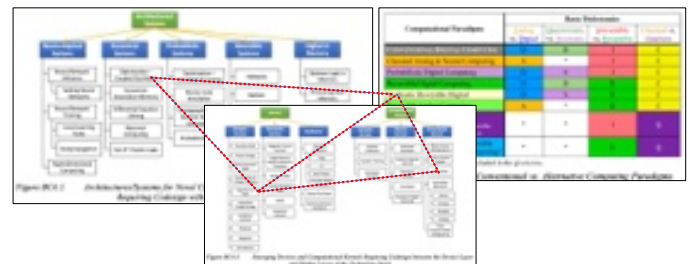


## The goals and drivers have changed



<https://royalsocietypublishing.org/doi/10.1098/rsta.2019.0061>

## Deep co-design across the entire stack



## What is Artificial Intelligence (AI)? What is Machine Learning (ML)?

## AI vs ML

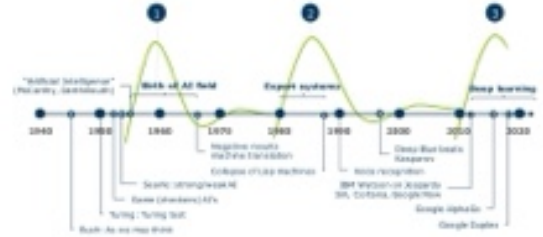
- **Machine Learning (ML)** is one technique used to achieve **Artificial Intelligence (AI)**.
- All ML is AI, but not all AI is ML.
- AI is the destination (e.g., intelligent machines), while ML is one path to get there (learning from data).
- **Artificial Intelligence (AI)**: Broader concept encompassing machines that can perform tasks requiring human intelligence. Ultimate goal: Artificial General Intelligence (AGI)
- **Machine Learning (ML)**: Subset of AI focused specifically on algorithms that improve through experience. Example: neural networks

## AI vs ML



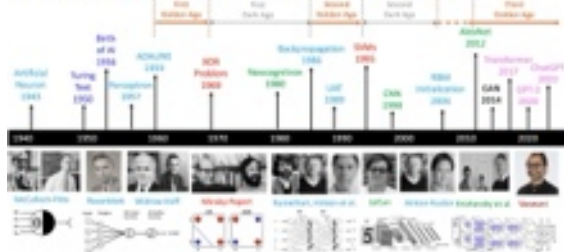
Source: Tignis

## A very brief history of AI



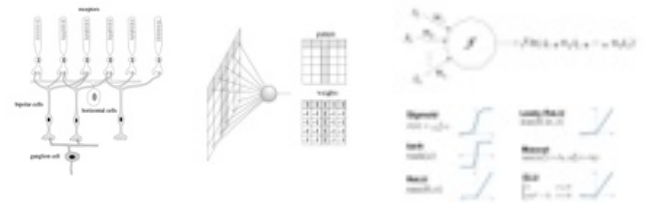
Strong AI = AGI = Artificial General Intelligence

## A very brief history of AI with deep learning



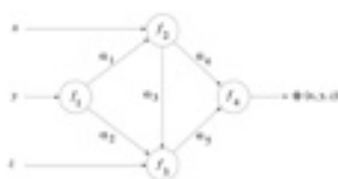
<https://medium.com/@lmpo/a-brief-history-of-ai-with-deep-learning-26f7948bc87b>

## An abstract neuron



Rojas, 1986

## A neural network



Rojas, 1986

## Special types of neural networks

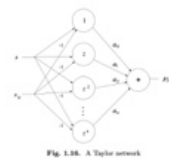


Fig. 1.16. A Taylor network.



Fig. 1.17. A Fourier network.

Figure 1.17 shows how a Fourier series can be implemented as a neural network. If the function  $F$  is to be developed as a Fourier series it has the form

$$F(x) = \sum_{n=1}^{\infty} [A_n \cos(n\pi x) + B_n \sin(n\pi x)] \quad (1.2)$$

Rojas, 1986



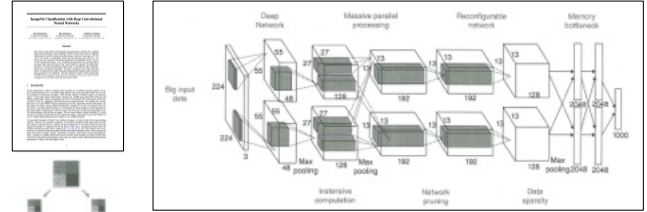
# Neuromorphic Computing with Large Scale Spiking Neural Networks

Neuromorphic Computing with Large Scale Spiking Neural Networks

Table 1: Comparison of leading neuromorphic hardware platforms in terms of neuron capacity, power efficiency, and scalability.

Hardware Platform	Technology	Neuron Capacity	Power Efficiency	Scalability
IBM TrueNorth	Digital	1M neurons	High	Medium
Intel Loihi	Digital	1M neurons	Very High	High
Spinnaker	Digital	1M neurons	Medium	High
BrainScaleC	Analog	1M neurons	Low	Low
TrueNorth	Hybrid	1M neurons	High	High

# AlexNet: A classical deep convolutional neural network



AlexNet paper: [https://papers.nips.cc/paper\\_files/paper/2012/file/c39862d3b9d6b76c8436e924a68c45b-Paper.pdf](https://papers.nips.cc/paper_files/paper/2012/file/c39862d3b9d6b76c8436e924a68c45b-Paper.pdf)

Liu & Law, Artificial Intelligence Hardware Design, 2021

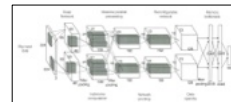
# Deep learning HW issues

Layer	Size	Filter	Depth	Stride	Padding	Number of parameters	Forward computation
Conv1 + ReLU	224x224x3	11x11	16	4		$(11 \times 11 \times 3 \times 16) \times 16 = 144,000$	$(11 \times 11 \times 3 \times 16) \times 16 \times 16 = 1,881,600$
Max pooling	16x16x16	3x3	2				
Norm	16x16x16						
Conv2 + ReLU	16x16x16	5x5	256	1	2	$(5 \times 5 \times 16 \times 256) \times 256 = 1,048,576$	$(5 \times 5 \times 16 \times 256) \times 256 \times 256 = 4,194,304$
Max pooling	8x8x256	3x3	2				
Norm	8x8x256						
Conv3 + ReLU	8x8x256	3x3	128	1	1	$(3 \times 3 \times 256 \times 128) \times 128 = 885,120$	$(3 \times 3 \times 256 \times 128) \times 128 \times 128 = 1,441,856$
Conv4 + ReLU	8x8x128	3x3	128	1	1	$(3 \times 3 \times 128 \times 128) \times 128 = 1,441,856$	$(3 \times 3 \times 128 \times 128) \times 128 \times 128 = 1,441,856$
Conv5 + ReLU	8x8x128	3x3	128	1	1	$(3 \times 3 \times 128 \times 128) \times 128 = 1,441,856$	$(3 \times 3 \times 128 \times 128) \times 128 \times 128 = 1,441,856$
Max pooling	4x4x128	3x3	2				
Dropout	4x4x128						
Dropout	4x4x128						
FC1 + ReLU	4096					$4096 \times 4096 = 16,777,216$	$4096 \times 4096 = 16,777,216$
Dropout	4096						
FC2 + ReLU	4096					$4096 \times 4096 = 16,777,216$	$4096 \times 4096 = 16,777,216$
Dropout	4096						
FC3 + ReLU	1000					$1000 \times 1000 = 1,000,000$	$1000 \times 1000 = 1,000,000$
Dropout	1000						
Softmax	1000						

FC = fully connected

Liu & Law, Artificial Intelligence Hardware Design, 2021

# Exploiting parallelism

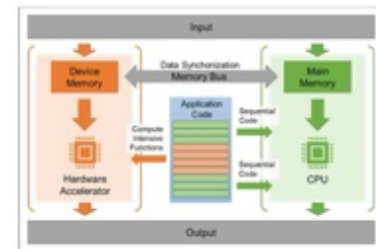


Does that come at a cost?



# Hardware for AI/ML

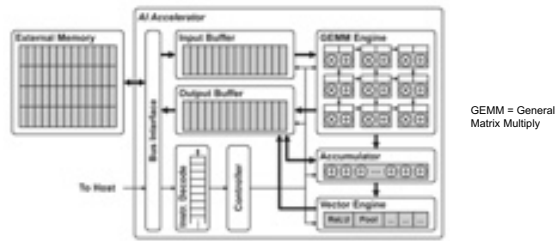
# General co-processing mechanism



Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shih; Artificial Intelligence and Hardware Accelerators, 2025



## General AI accelerator architecture



Mahra, Ashutosh; Cha, Jaekwang; Park, Hyunbin; Kim, Shihoo. Artificial Intelligence and Hardware Accelerators, 2025

