

Hardware for AI and ML: an Overview

Portland State University
 Department of Electrical and Computer Engineering (ECE)
www.teuscher-lab.com
teuscher@pdx.edu



teuscher•Lab
teuscher-lab.com

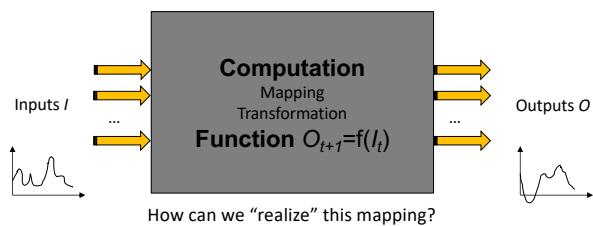
Portland State
 UNIVERSITY

What is a computer?

teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY

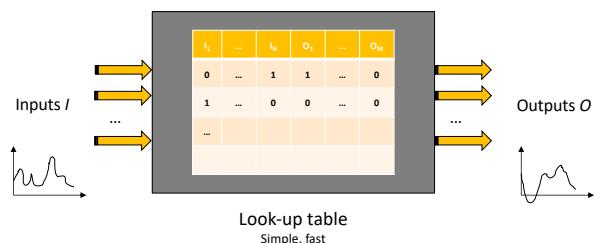
What is computation? (1)



teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY

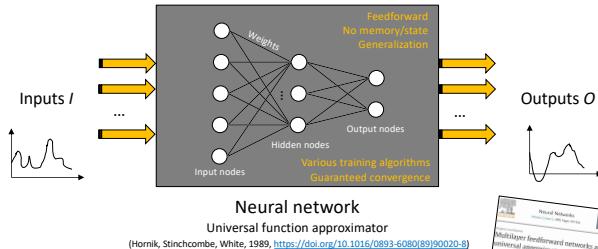
What is computation? (2)



teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY

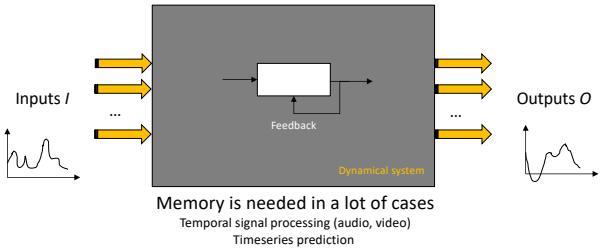
What is computation? (3)



teuscher•Lab
teuscher-lab.com

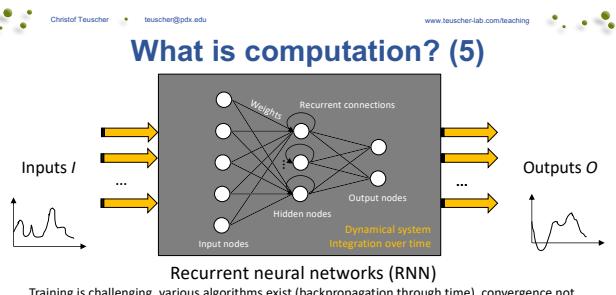
Portland State
 UNIVERSITY

What is computation? (4)



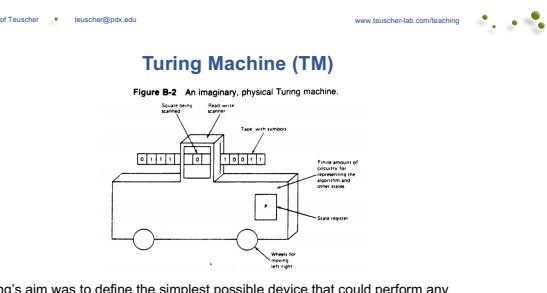
teuscher•Lab
teuscher-lab.com

Portland State
 UNIVERSITY



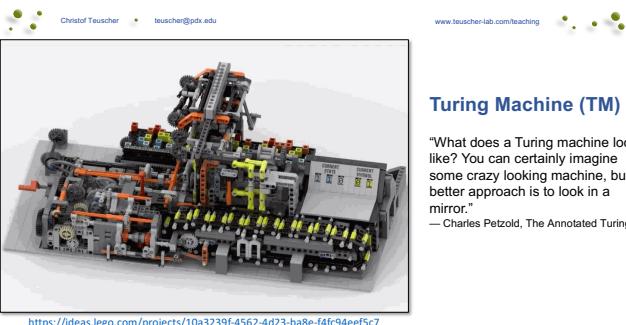
teuscher•Lab
teuscher-lab.com

Portland State UNIVERSITY



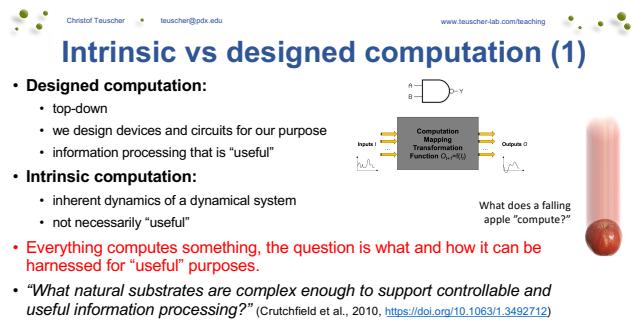
teuscher•Lab
teuscher-lab.com

Portland State UNIVERSITY



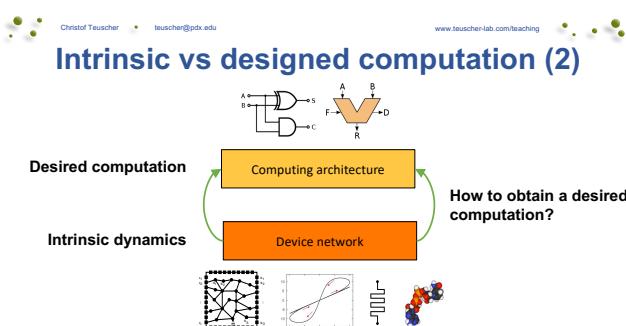
teuscher•Lab
teuscher-lab.com

Portland State UNIVERSITY



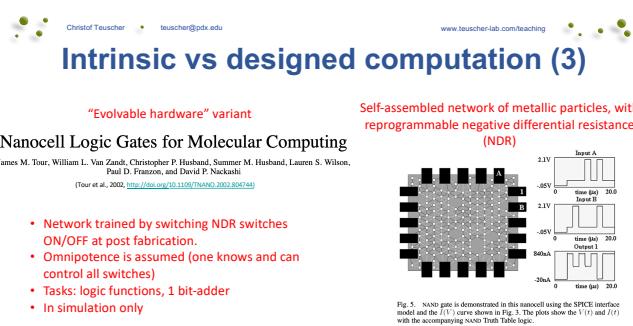
teuscher•Lab
teuscher-lab.com

Portland State UNIVERSITY



teuscher•Lab
teuscher-lab.com

Portland State UNIVERSITY



teuscher•Lab
teuscher-lab.com

Portland State UNIVERSITY

Intrinsic vs designed computation (4)

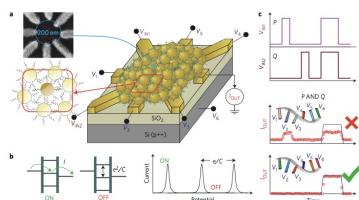
LETTERS
PUBLISHED ONLINE 20 SEPTEMBER 2015 | DOI: 10.1038/NNANO.2015.207

nature
nanotechnology

Evolution of a designless nanoparticle network into reconfigurable Boolean logic

S. K. Bose^a, C. P. Lawrence^b, Z. Liu^c, K. S. Makarewicz^c, R. M. J. van Damme^c, H. J. Broersma^c and W. G. van der Wiel^{a*}

(Bose et al., 2015, <https://doi.org/10.1038/nnano.2015.207>)



teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

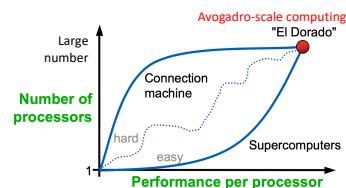
What kind of computing machinery would an extraterrestrial build?



teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

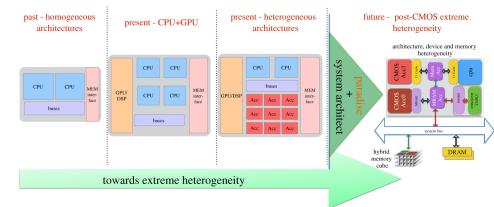
The goals and drivers have changed



teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

The goals and drivers have changed



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

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

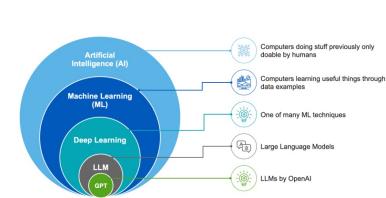
AI vs ML

- 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

teuscher•Lab
teuscher-lab.com

Portland State
UNIVERSITY

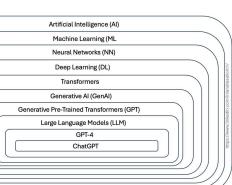
AI vs ML



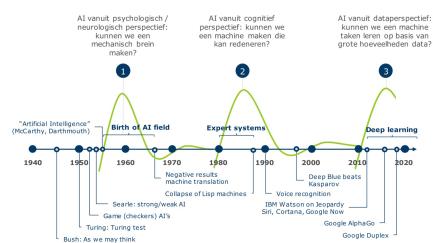
Source: Tignis

teuscher•Lab
teuscher-lab.com

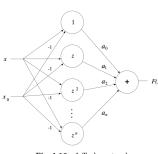
Portland State
UNIVERSITY



A very brief history of AI



Special types of neural networks



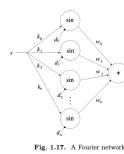
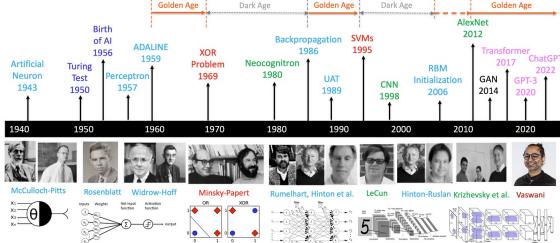
$$F(x) = a_0 + a_1(x - x_0) + a_2(x - x_0)^2 + \dots + a_n(x - x_0)^n + \dots,$$


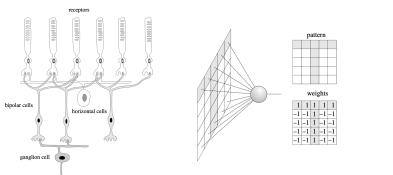
Figure 1.17 shows how a Fourier series can be implemented as a neural network. If the function $F(x)$ is to be developed as a Fourier series it has the form

$$F(x) = \sum_{i=0}^{\infty} (a_i \cos(ix) + b_i \sin(ix)). \quad (1.2)$$

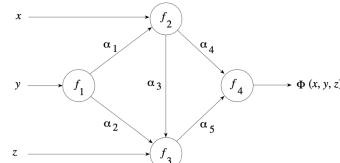
A very brief history of AI with deep learning



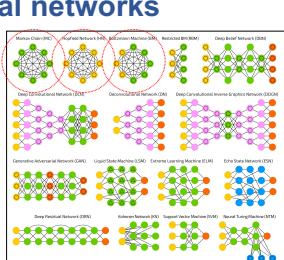
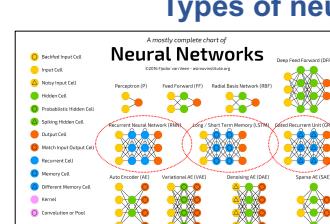
An abstract neuron

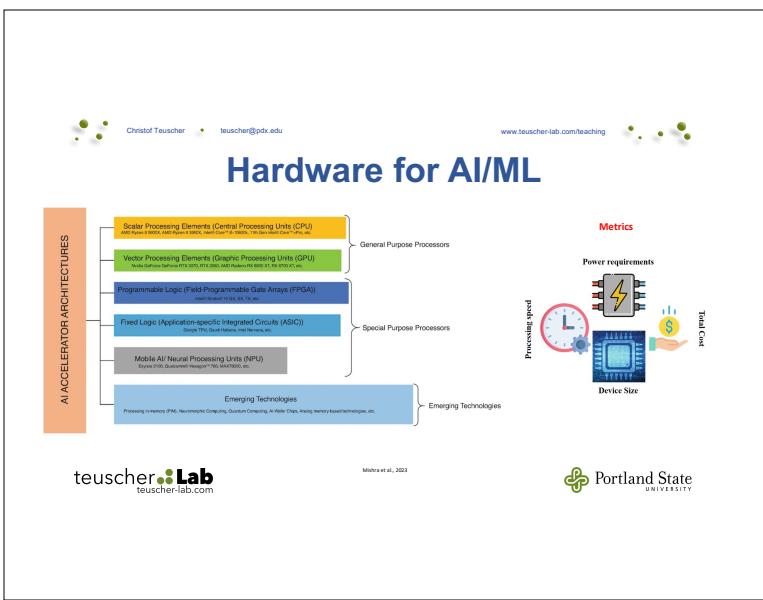
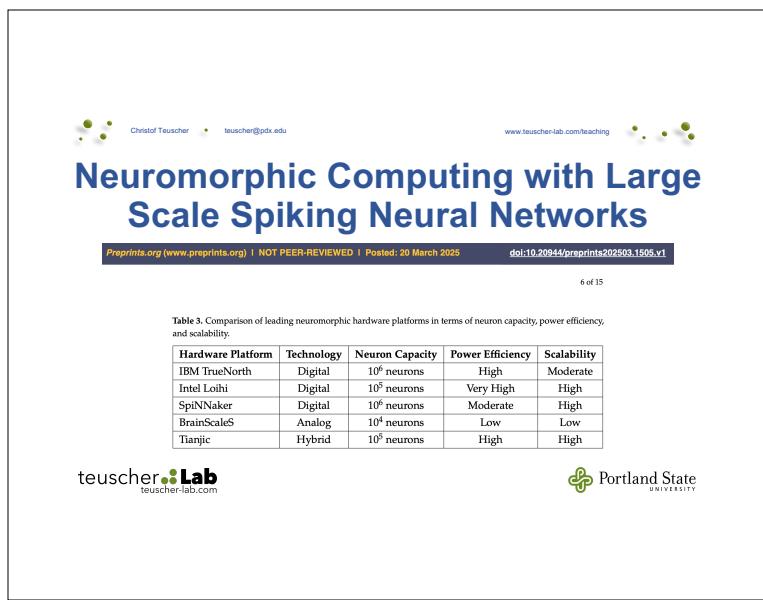
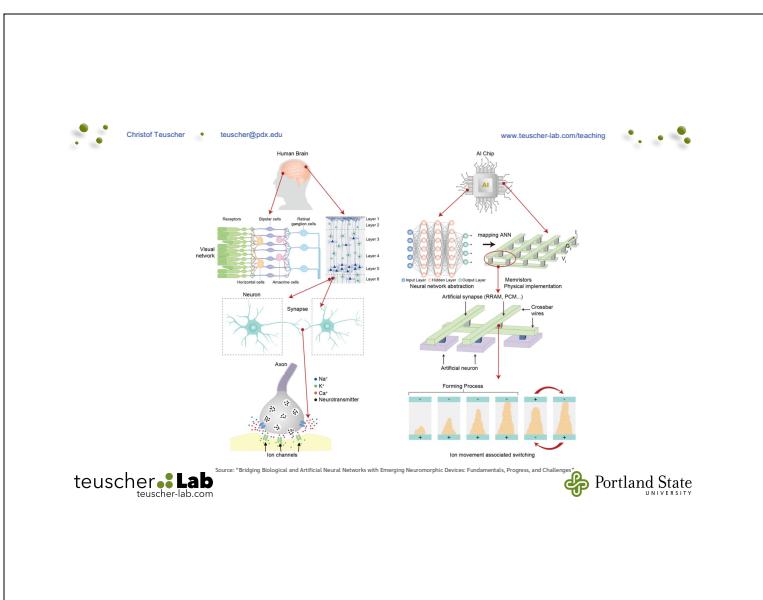
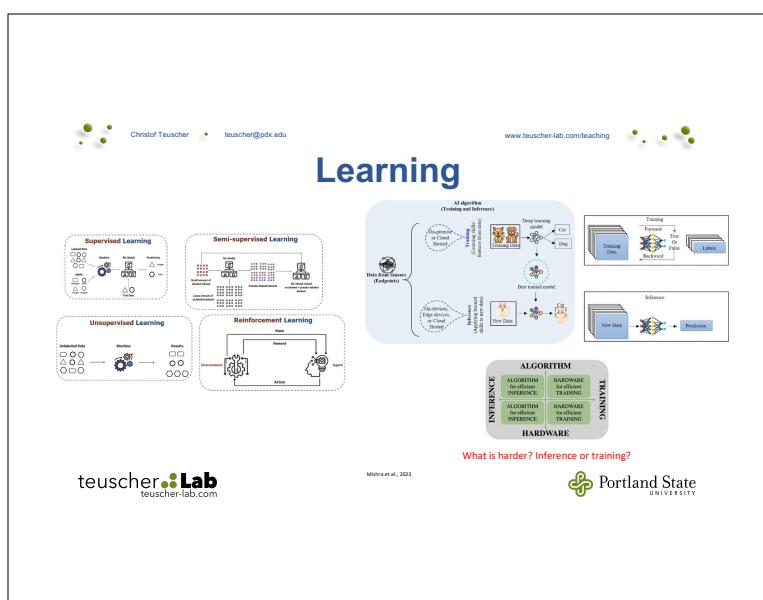
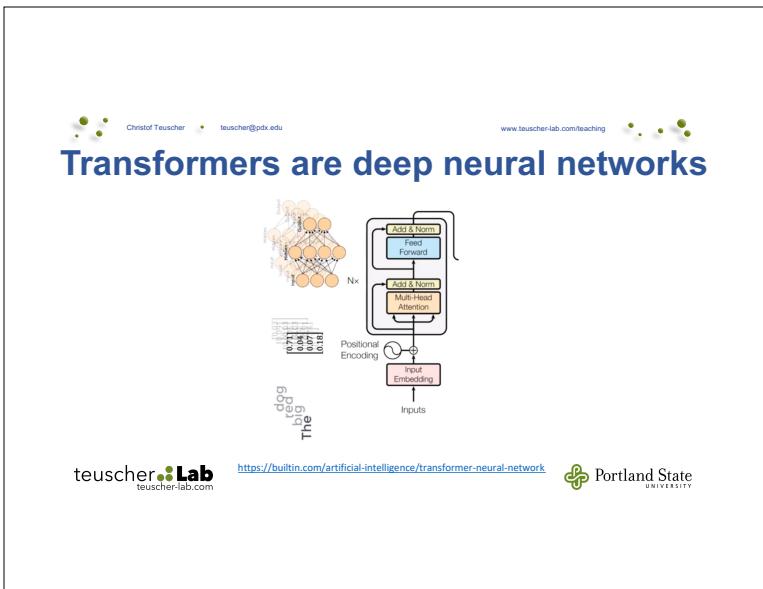
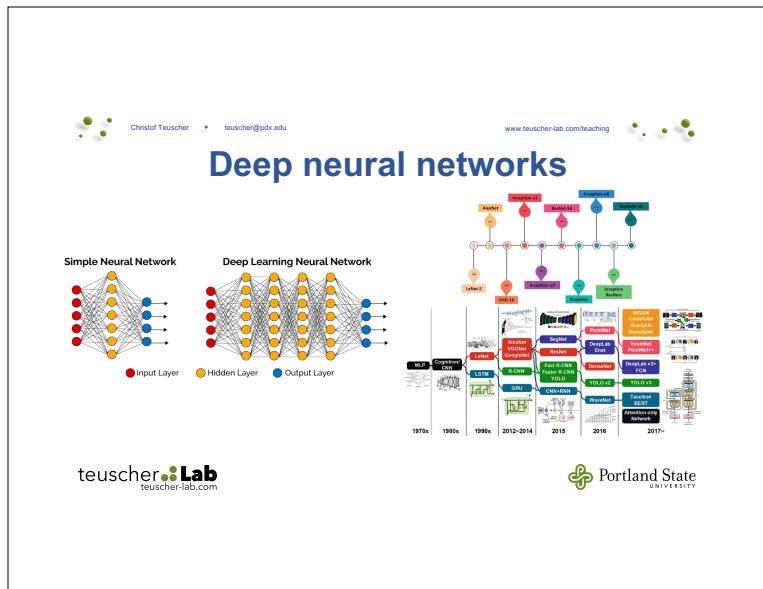


A neural network



Types of neural networks





Hardware for AI/ML

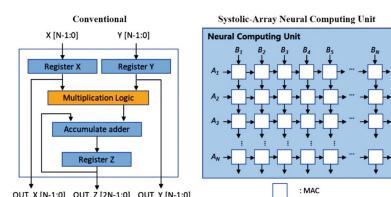
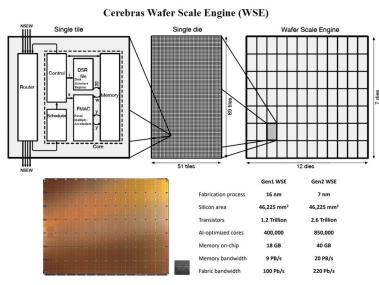


Fig. 19 A comparison of MAC operations performed on conventional and neural computing units.

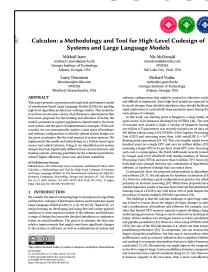
Hardware for AI/ML

Processor	Power Consumption	Strengths	Limitations
CPU	High	<ul style="list-style-type: none"> Flexible General-purpose processing Complex instructions and tasks System management 	<ul style="list-style-type: none"> Possible memory access bottlenecks Few cores (4-16)
GPU	High	<ul style="list-style-type: none"> Parallel cores (~1000s of cores) High Performance AI processing 	<ul style="list-style-type: none"> Power consumption Large footprint
FPGA	Medium	<ul style="list-style-type: none"> Configurable logic gates Flexible In-field re-programmability 	<ul style="list-style-type: none"> Programming complexity
ASIC	Low	<ul style="list-style-type: none"> Custom logic designed with libraries Optimized for computing Small footprint 	<ul style="list-style-type: none"> Fixed function Expensive custom design
Vision Processing Unit (VPU)	Ultra-low	<ul style="list-style-type: none"> Dedicated image and vision co-processor Small footprint 	<ul style="list-style-type: none"> Limited dataset and batch size Limited network support
Tensor Processing Unit (TPU)	Low to medium	<ul style="list-style-type: none"> Specialized tool support Optimized for TensorFlow 	<ul style="list-style-type: none"> Proprietary design Limited framework support

Hardware for AI/ML



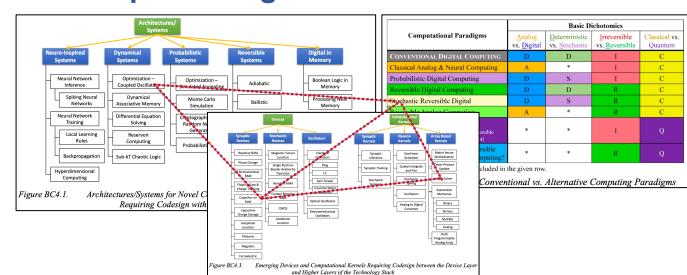
Co-design & LLM



(HW/SW) Co-design

- Classical approach:** in college, different technologies are taught in different classes.
- Hardware-software co-design:**
 - Started in the 1990s.
 - Its core idea is the concurrent designs of hardware and software components of complex electronic systems.
- More broadly:**
 - The concurrent design of hardware and software across all layers of the compute stack.
- What are possible metrics and trade-offs?**

Deep co-design across the entire stack





Food for thought

- Why is executing a neural network on a general purpose computer inefficient?
- What circuitry/processor would you design to find the shortest path in large graphs?
 - What if the graph is dense?
 - What if the graph is sparse?
- What circuitry/processor would you design to recognize QR codes for an embedded camera?
 - Recognition needs to be fast and low-power.