

MACHINE LEARNING: ALGORITHMS & HARDWARE

KAUSHIK ROY

CENTER FOR BRAIN-INSPIRED COMPUTING (C-BRIC)

PURDUE UNIVERSITY

WEST LAFAYETTE, INDIANA



Adopted from lectures by Andrew Ng



We live in a data-driven world!!!

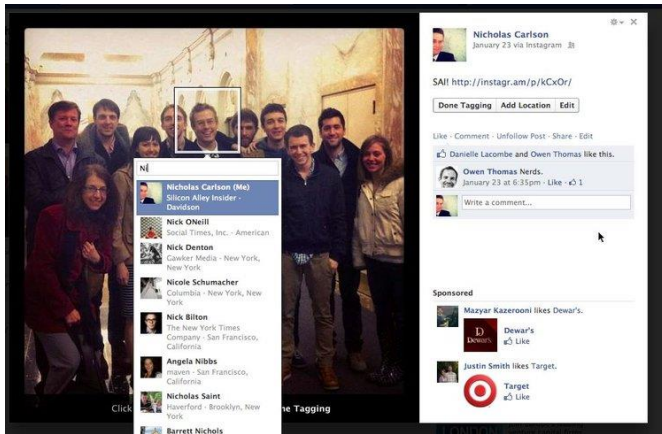


Artificial Intelligence (AI): Why now?



**Data is the new oil and AI is
the new electricity!!!!**

AI enables Cognition



Facebook photo-tagging



Alexa

Siri

Google Now

Cortana



Self-driving Car

amazon.com

Recommended for You

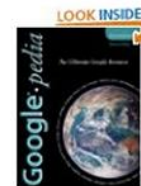
Amazon.com has new recommendations for you based on [items](#) you purchased or told us you own.



[Google Apps Deciphered: Compute in the Cloud to Streamline Your Desktop](#)



[Google Apps Administrator Guide: A Private-Label Web Workspace](#)



[Googlepedia: The Ultimate Google Resource \(3rd Edition\)](#)

AI has even defeated humans!!

1997



IBM Deep Blue vs. Kasparov

2011



IBM Watson vs. Brad Ritter & Ken Jennings

2016



Google AlphaGo vs. Lee Sedol

<https://www.youtube.com/watch?v=jGyCsVhtWOM>

A great documentary on Alphago and AI in general!!

Key driver of AI: Neural Networks/Deep learning

More recent success stories ...



Identifying snow leopards with AI



Google AlphaGo vs. Lee Sedol
(1920 CPUs, 280 GPUs)



IBM's AI loses to Harish Natarajan, but still is persuasive...



Artificial intelligence used to recognize primate faces in the wild

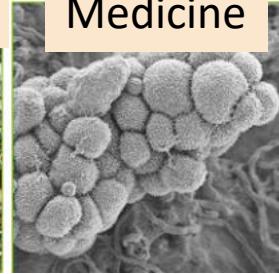
AI is scaling...

- Across applications
- Investments
- and people...

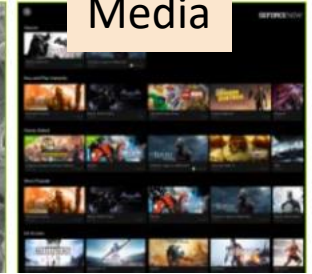
Image analytics



Medicine



Media



Security/Defense



Transportation

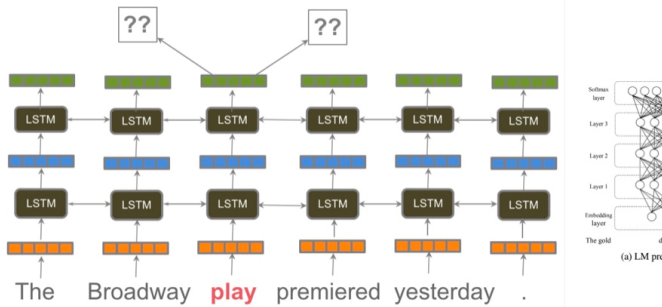


Customer Care

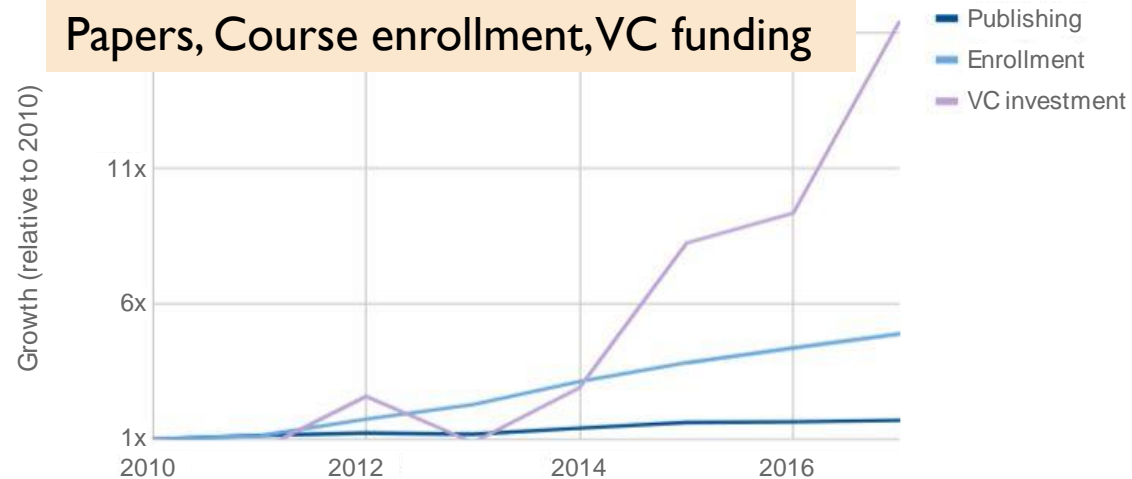


NLP's ImageNet moment has arrived

08.JUL.2018



Papers, Course enrollment, VC funding



However... Many challenges yet to be addressed

AI Stats News: Chatbots Lead To 80%

Forbes CommunityVoice Connecting expert communities to the Forbes audience. What is This?

2,400 views | Feb 22, 2019, 08:30am

1d

Exp
Ope



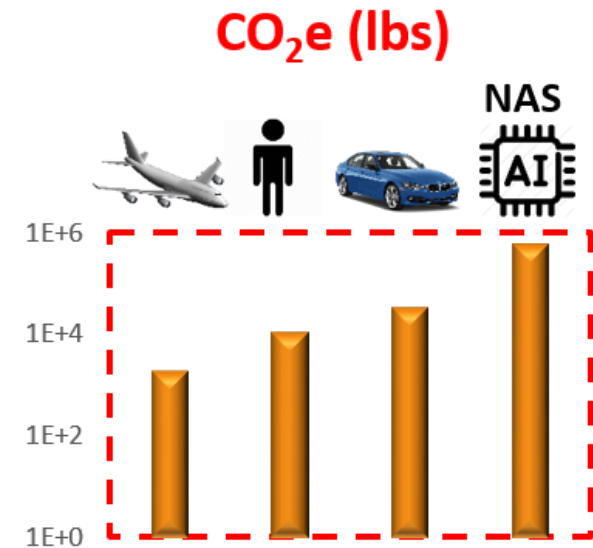
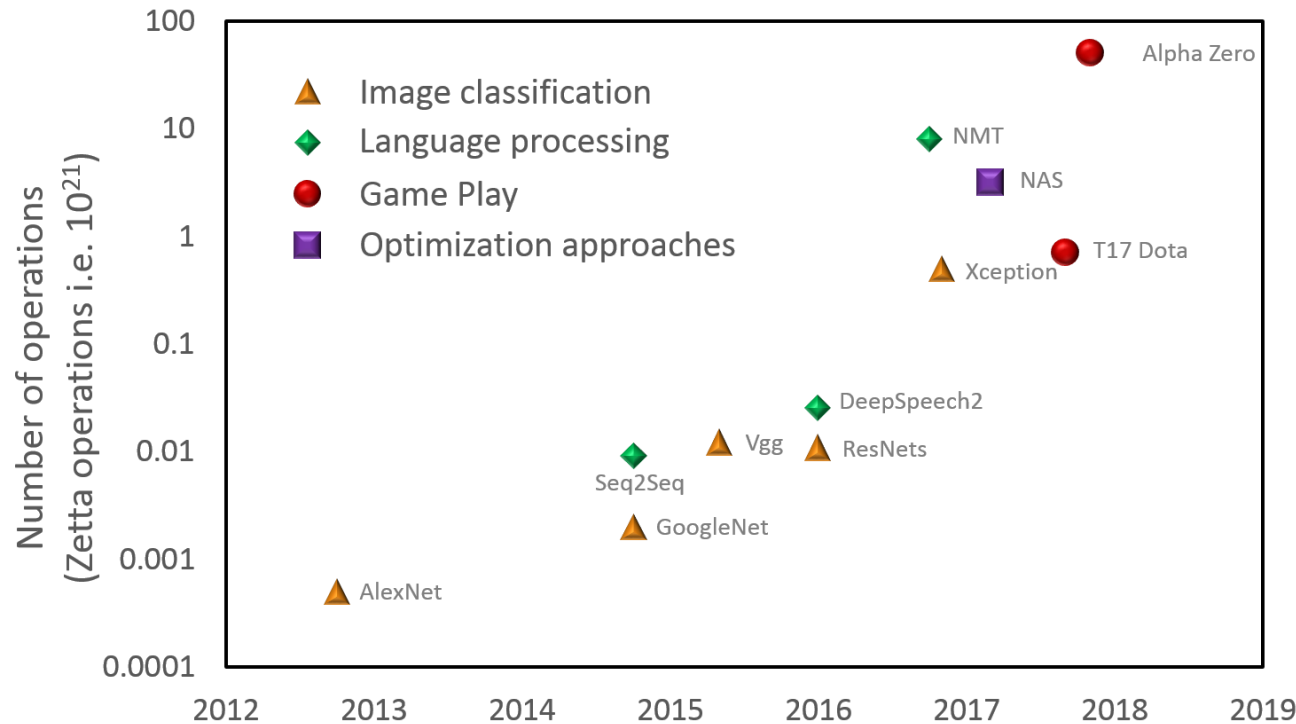
Consumer Concerns About Self-Driving Cars

% of respondents naming the following reasons for their reluctance to use self-driving cars



AI Compute Demands (Training)

ML application trends (Training)



- Estimated CO₂ emissions from **NAS on Transformer (big)** is:
 - **315x** higher than **Air Travel** from NY to SF/passenger
 - **17x** higher than **average American (1 year)**
 - **5x** higher than **Car (1 lifetime)**

Efficiency Gap in AI

- Case study: Object recognition in a smart glass with a state-of-the-art accelerator



Google Edge TPU

Retinanet DNN* on a smart glass

Performance	
Frames/sec	13.3
Battery Life	
Energy/op	0.5 pJ/op
Energy/frame	0.15 J/frame
Time-to-die (2.1WH)	64 mins

*300 GOPs/inference

Where do the in-efficiencies come from?

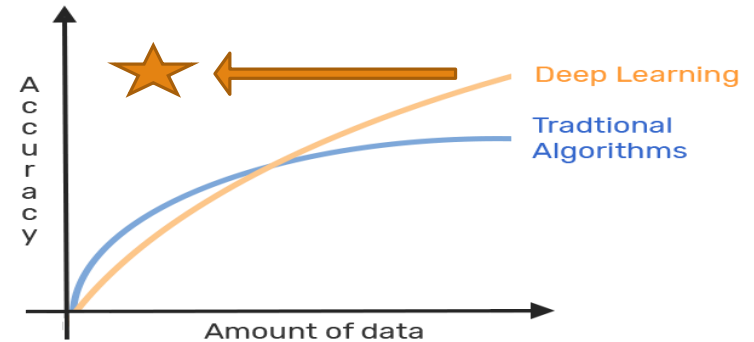
Algorithms

Hardware Architecture

Circuits and Devices

Beyond compute efficiency....

- Learning with less data
- Generalization & Robustness
- Lifelong learning

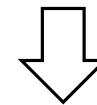
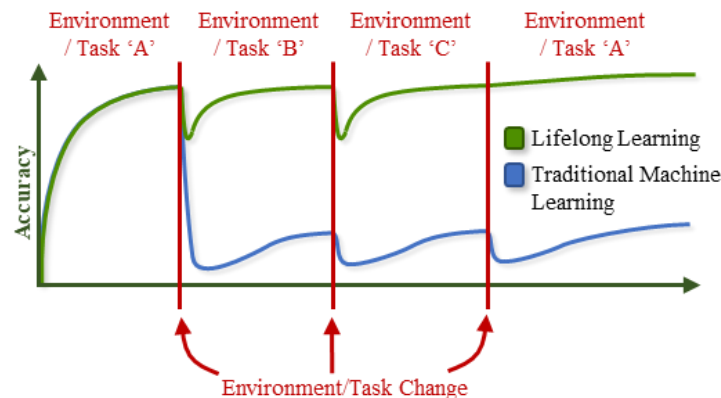


$$X_{clean} + 0.005 \times \Delta = X_{adversary}$$

97.3% confidence Macaw

Δ Adversarial Perturbation

88.9% confidence Bookcase



Neural Networks- Loosely brain-inspired*

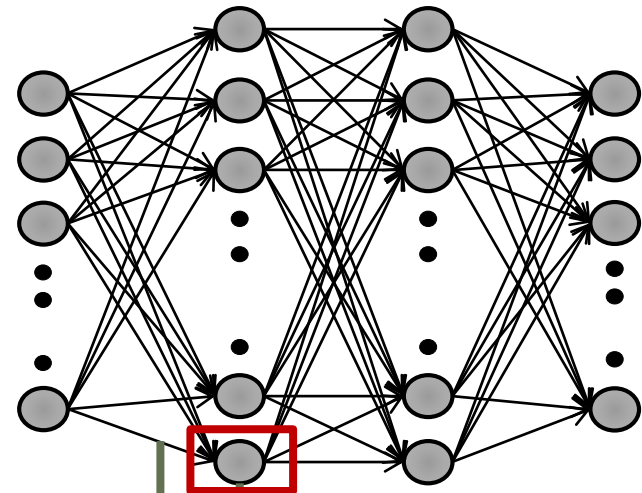
Biological Neural Network



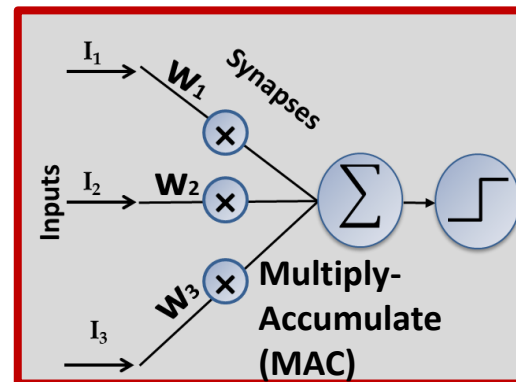
Interconnected web of neurons/synapses

- **Neurons** are the computing elements
- **Synapses/Weights** store memory and take part in learning/training
→ Intelligence

Artificial Neural Network



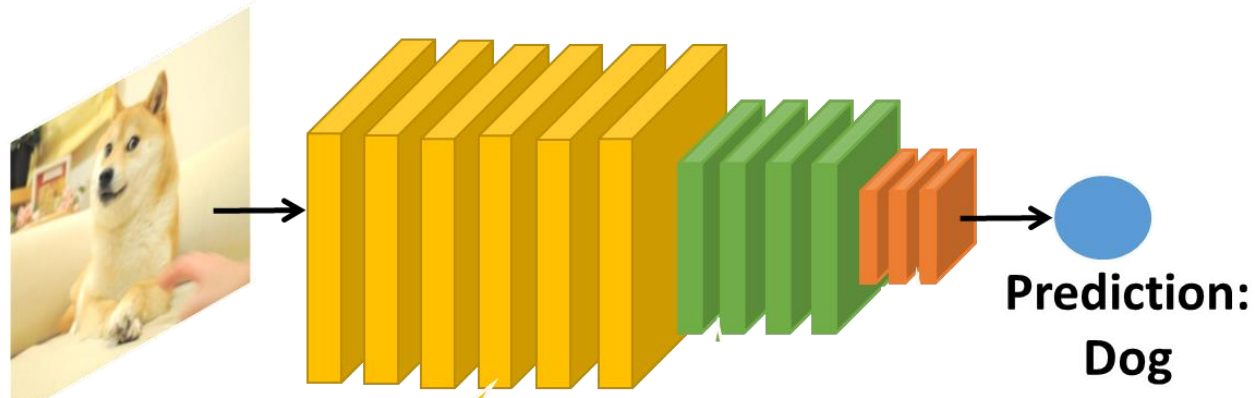
Weights ← → Neurons



Key operation:

- **MAC/ Dot-product**
(or $\sum_i w_i * I_i$)
- **Non-Linearity**
(e.g. threshold, sigmoid, etc.)

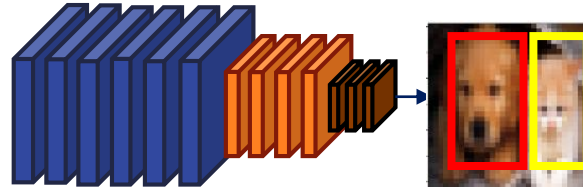
Training a Neural Network → Intelligence



**Network discovers features from these pixel values,
Pretty incredible!!!**

Neural Networks: Different Levels of Bio-fidelity

Deep Learning Revolution

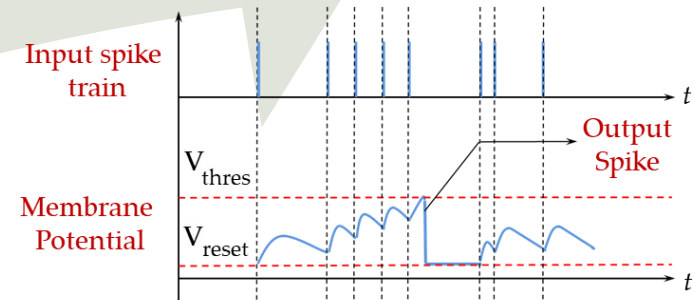
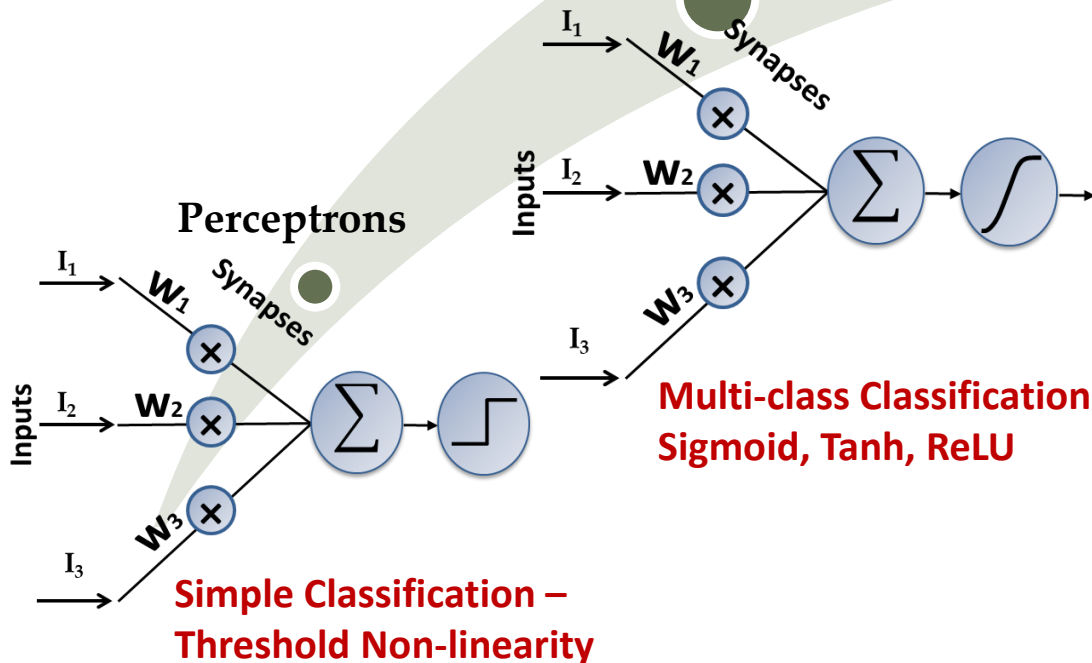


Convolutional Neural Networks for Complex Recognition

Spiking Neural Networks



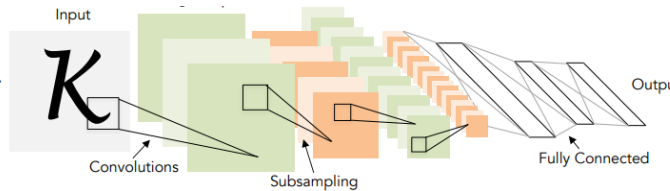
Artificial Neural Networks



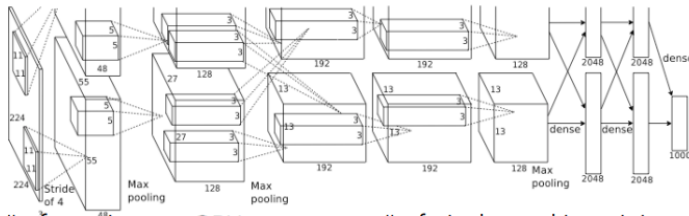
- Neural computation through spikes
- Energy savings due to computation with sparse spiking events

Key Enablers of Deep Learning

LeCun et al. 1990



Krizhevsky et al. 2012



ALGORITHMS- Improved with better training, regularization, optimization strategies

of transistors



10^6

~ 1 Million Dot Product Operations

of transistors



10^9

~ 1 Billion Dot Product Operations

GPUs



HARDWARE- Improved with architectural innovations, transistor scaling etc. for more compute power

of pixels used in training

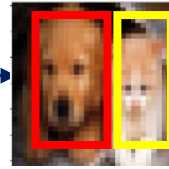
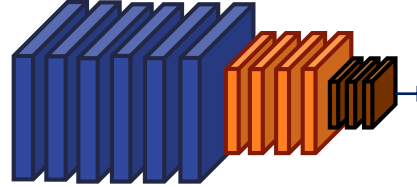
10^7 **NIST**

of pixels used in training

10^{14} **IMAGENET**

DATA- Improved training/learning with more data

Neural Networks: Different Levels of Bio-fidelity



Convolutional Neural Networks for Complex Recognition

Spiking Neural Networks



ALPHAGO
01:55:46

AlphaGo: 1920 CPUs and 280 GPUs (~1MegaWatt)

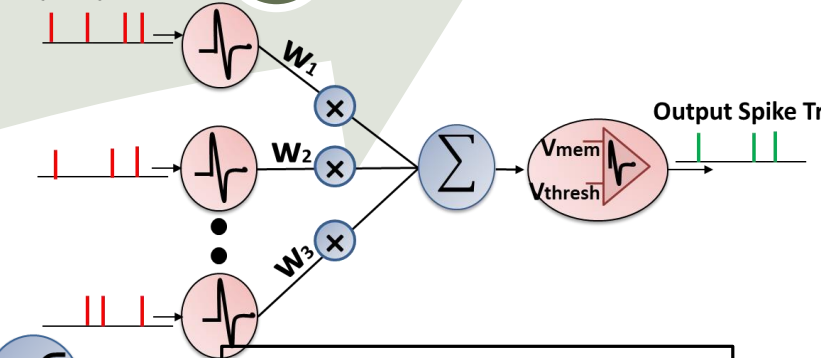
Vs.

Lee Sedol: 1 human brain (~20 Watts)

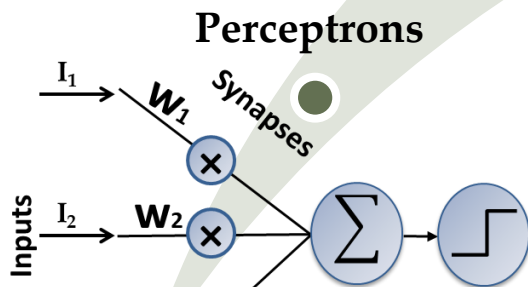
LEE SEDOL
01:55:41

Artificial Neural Networks

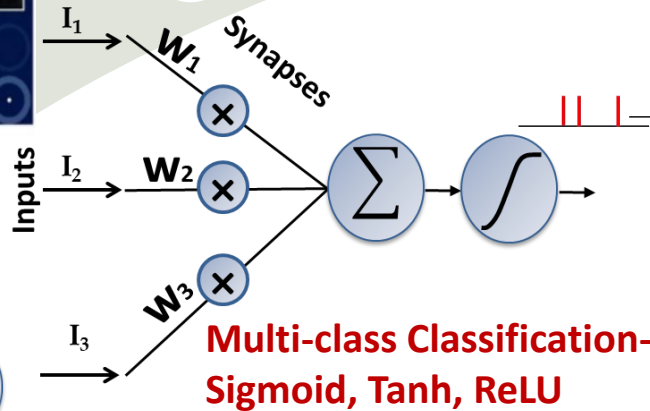
Input Spike Train



Neuromorphic Computing- Integrate-and-Fire Non-linearity



Simple Classification – Threshold Non-linearity



Multi-class Classification- Sigmoid, Tanh, ReLU

Why is Energy-Efficiency a Concern?



Case study: Object recognition in a smart glass

- Battery powered devices (smart-phones, smart-watches, drones etc.) have resource or energy constraints.
- Enabling intelligence on these platforms necessitates energy-efficient deep learning or neural network implementations



Overfeat DNN on a smart glass

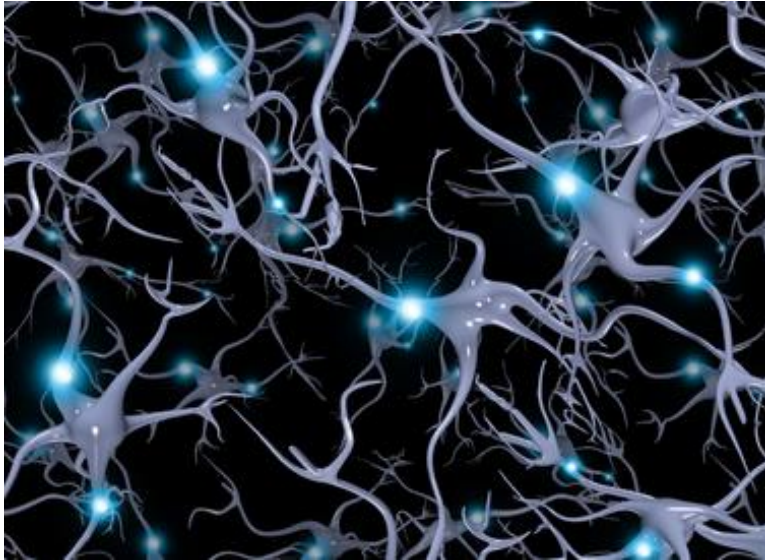
Performance*	
OMAP 4430	1.5 fps (ideal)
Battery Life	
Energy/op (mobile GPU)	5×10^{-2} nJ/op
Energy/frame	0.16 J/frame
Time-to-die (2.1WH)	25 min (ideal)

*3.2 GigaOPS/inference

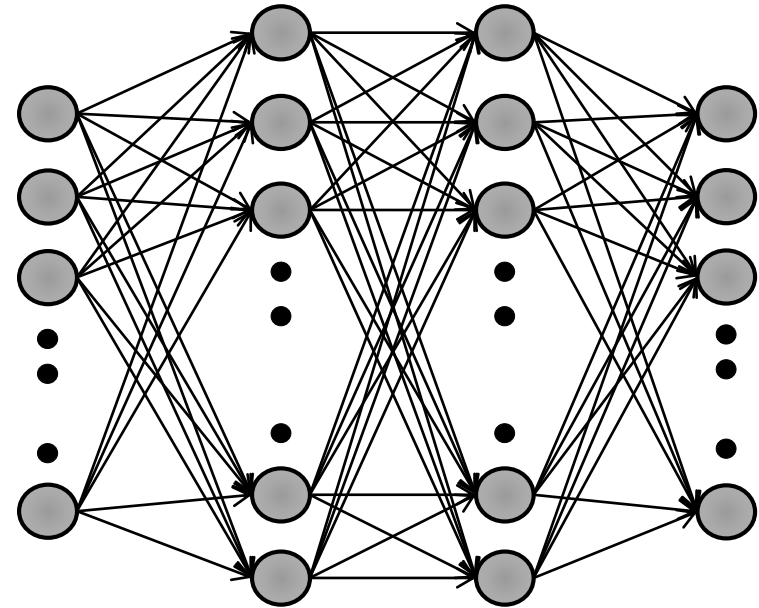
1 OPS \approx 1 dot product operation

Where do Inefficiencies Come From?

Biological Neural Network



Artificial Neural Network



Algorithms, Computing Architecture, Neurons, Synapses.....

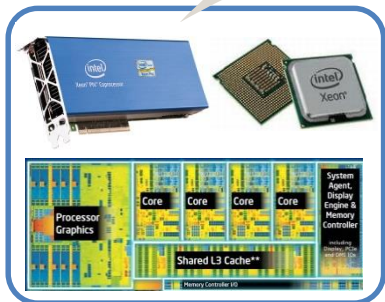
- Vastly interconnected web of neurons/synapses (1 billion neurons, 10000 synapses per neuron) → Compute and Memory are all intertwined and co-located
- Approximate and stochastic computation
- Sparse, irregular, event-driven, massively parallel networks

Hardware for Addressing Inefficiencies

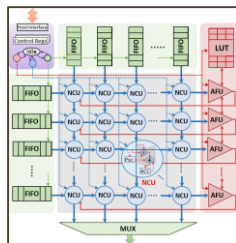
- CMOS and Post-CMOS neuro-mimetic devices and interconnects
- Compute-near-memory / Compute-in-memory
- Approximate and stochastic neuronal and synaptic hardware
- Architectures that embody computing principles from the brain (sparse, irregular, event-driven, massively parallel)
- Programming and evaluation frameworks



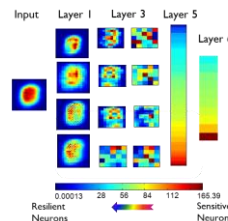
Multicores/GPUs



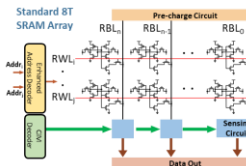
Accelerators



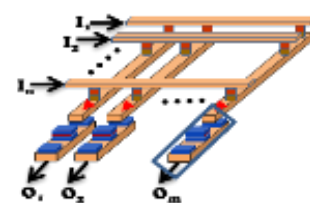
Approximate & Stochastic Hardware



In-memory computing



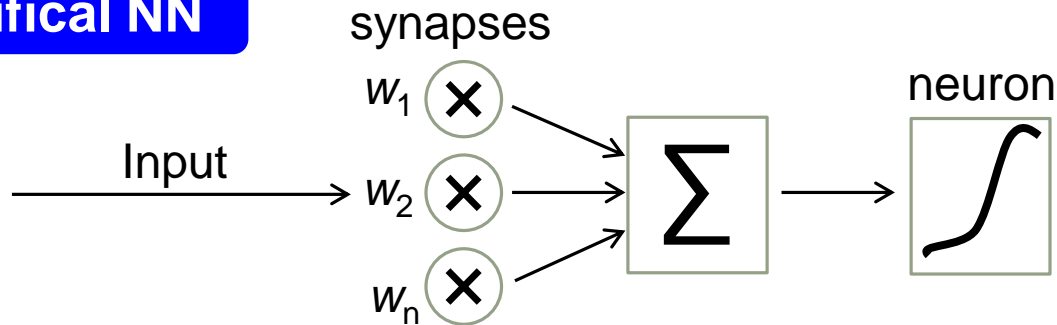
Post-CMOS Devices



~10⁴
Energy
Gap

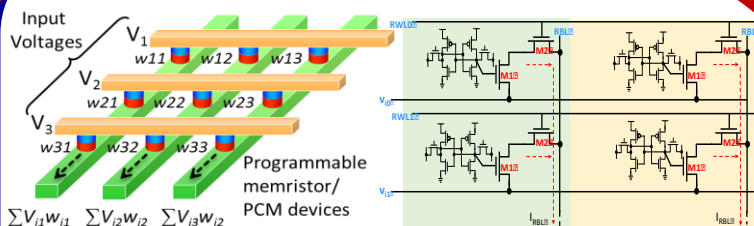
Neural Networks: Simple Hardware Model

Artificial NN

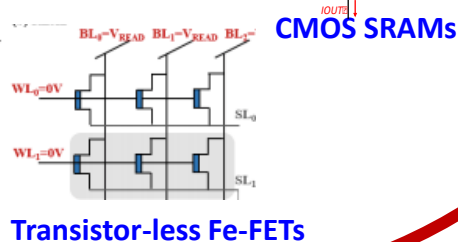


Weighted Summation

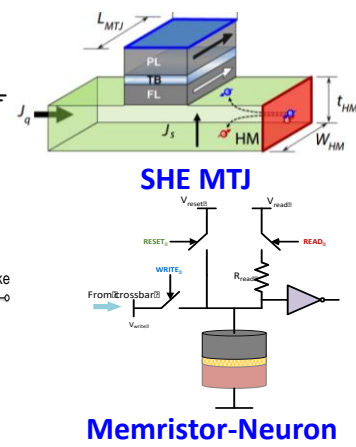
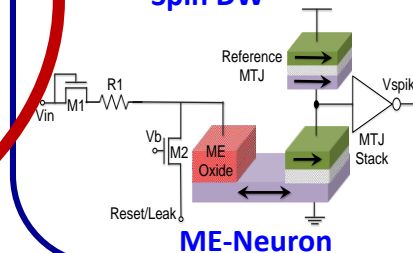
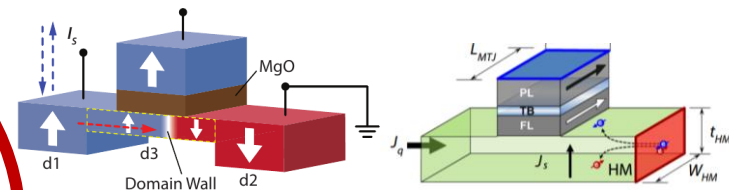
Non-linear function



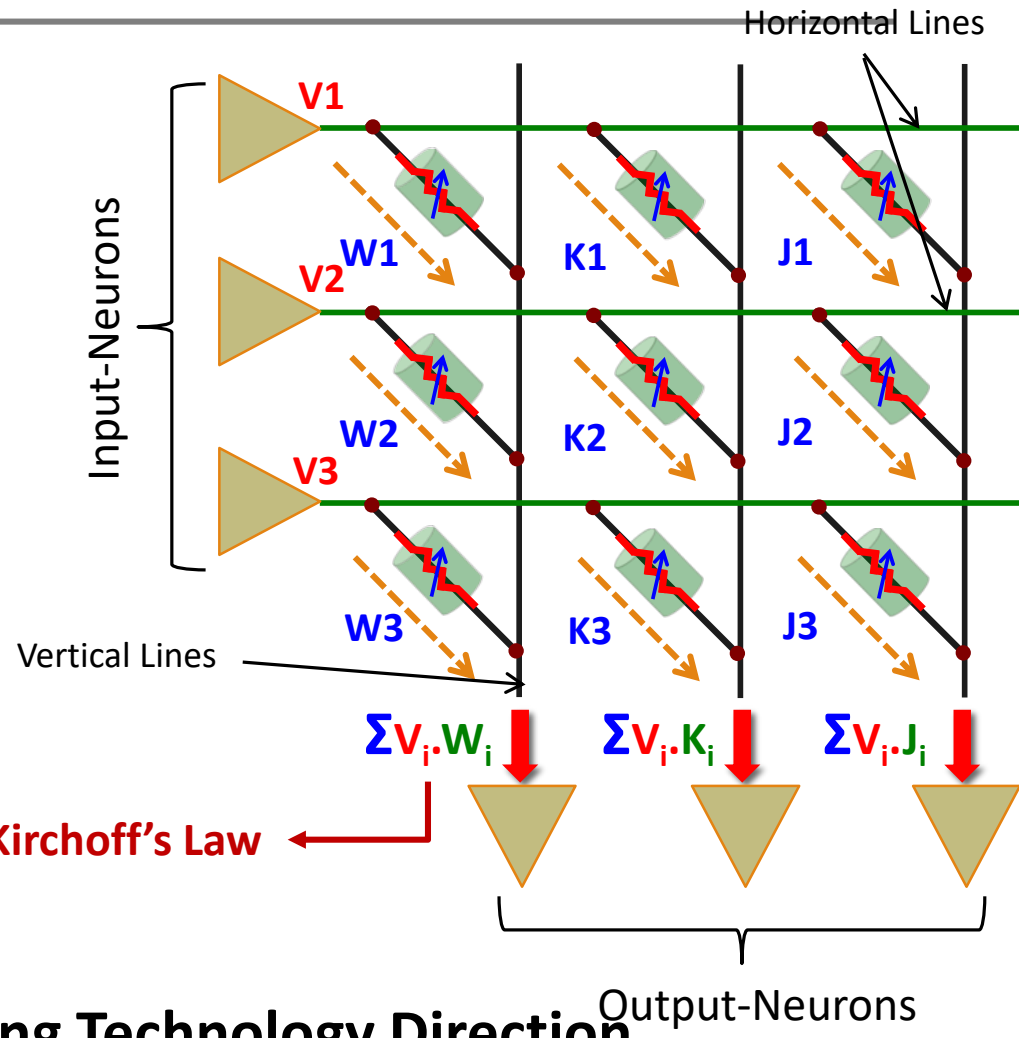
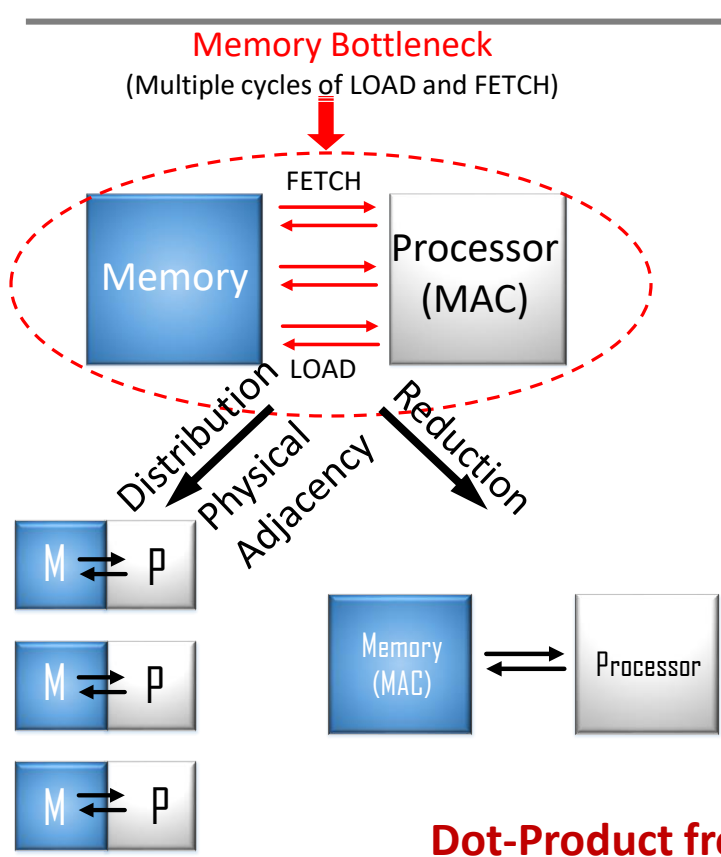
Spin/Memristors Crossbar Engine



Neurons



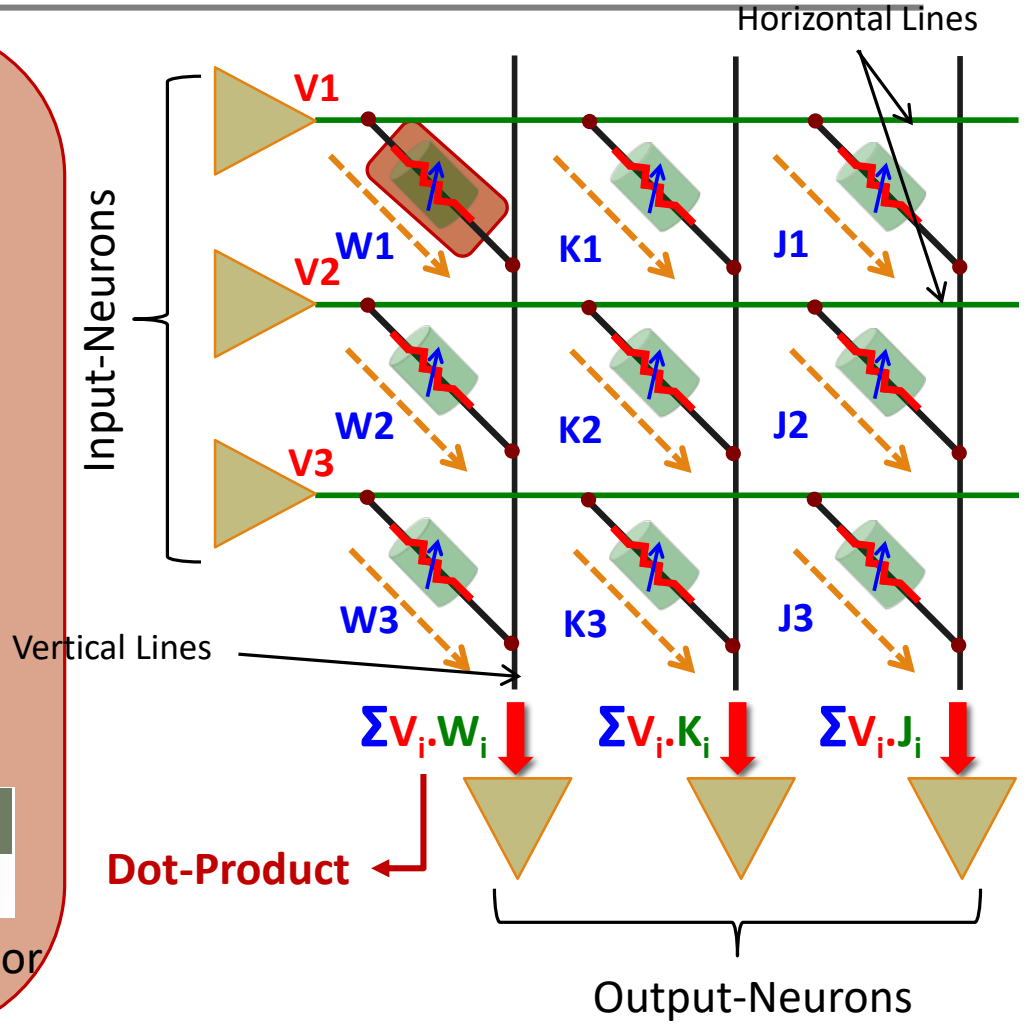
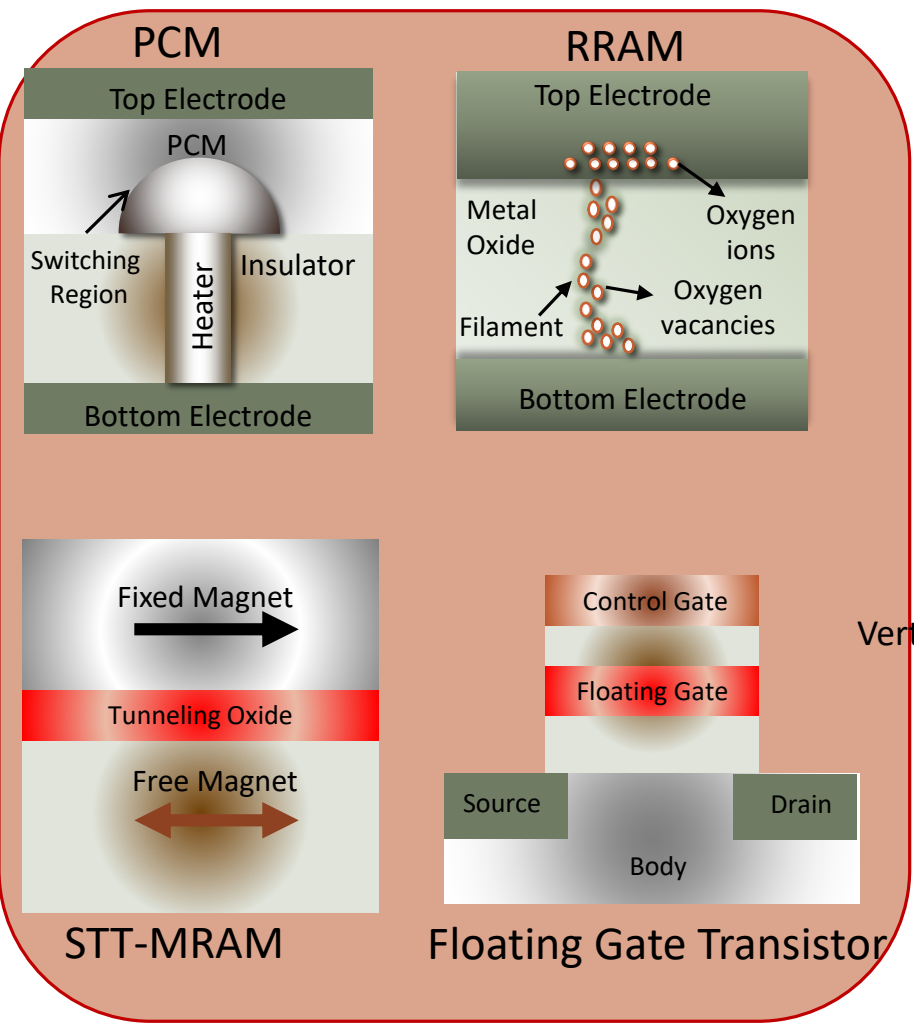
Memory Architecture/Circuits



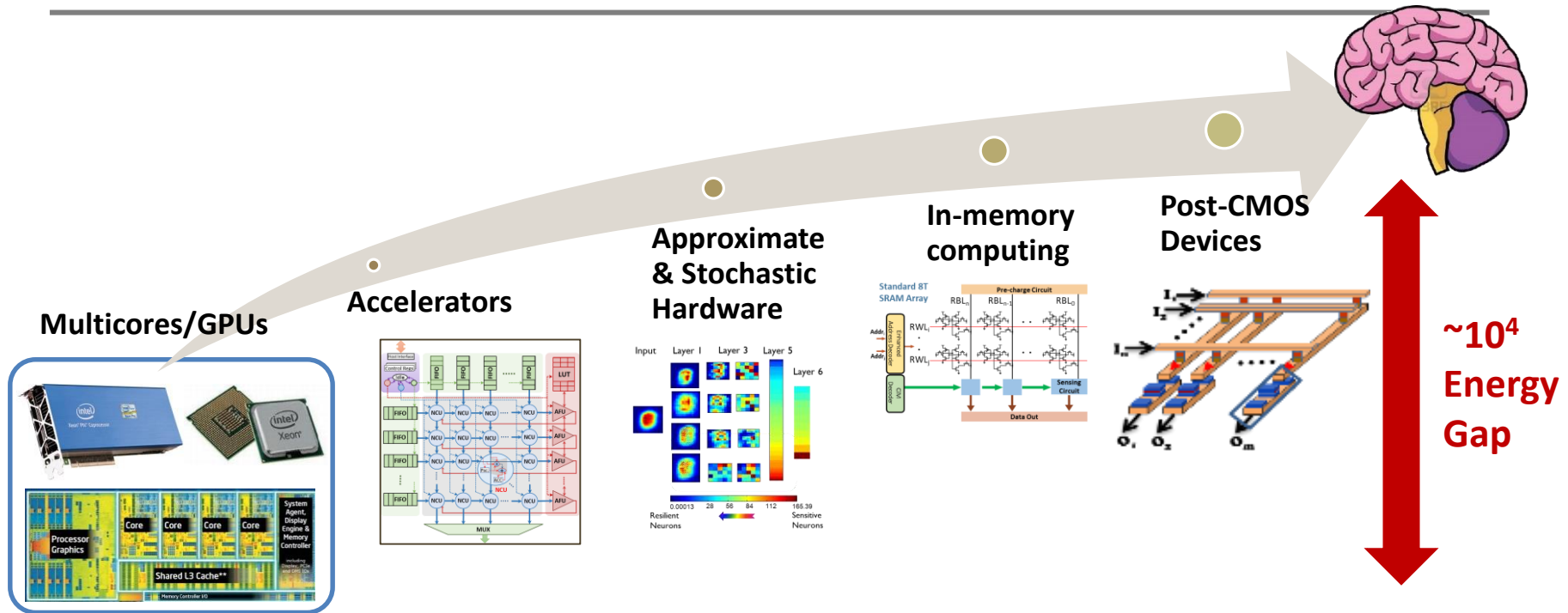
von Neumann CMOS-to-

Non-von Neumann Emerging Technology Direction

Post-CMOS Devices as Synaptic Memory Elements



Hardware for Addressing Inefficiencies

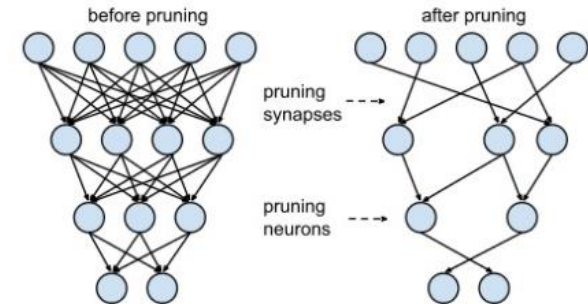


- Can we have algorithms that can yield energy-efficiency?
- Can the algorithms be hardware compatible?

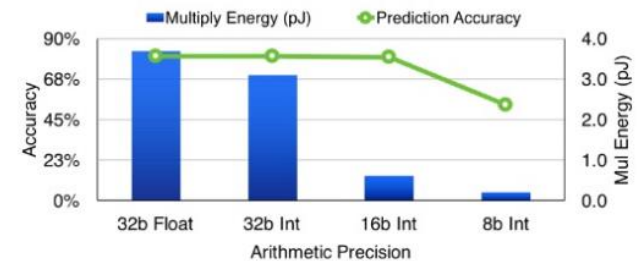
Algorithms for Addressing Inefficiencies

- Pruning/ Compression
 - Quantization
 - Binarized Networks – Dot Product to simplified XNOR
 - Early Exit
 - Spiking Networks
- Compatible with post-CMOS implementations**

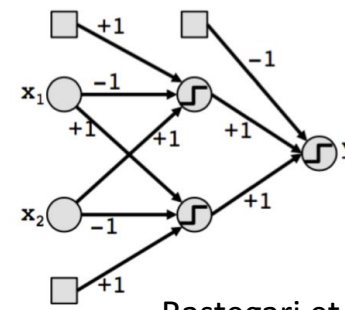
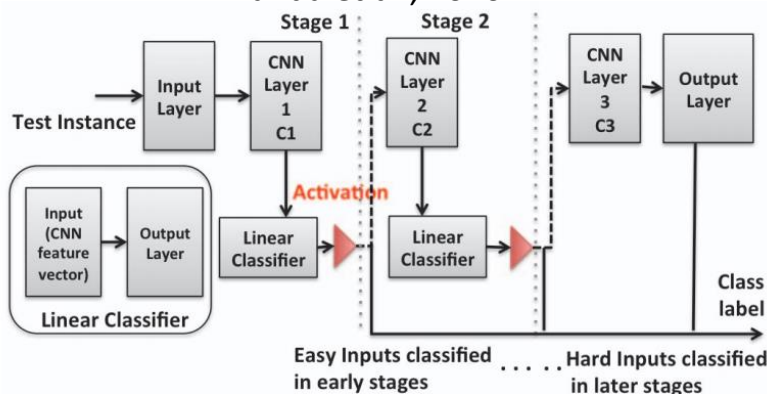
Han et al., 2016



Han et al., 2016, 2017



Panda et al., 2016



- 32x less memory with XNOR than MAC
- 23x faster than MAC

Rastegari et al., 2016

Algorithms, Systems, Circuits, & Devices

Top-Down: Device-driven Algorithms and Models

Spiking Networks Ising Networks/Boltzmann Machine Bayesian Networks

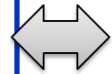
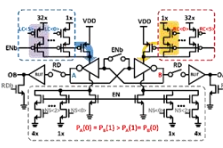
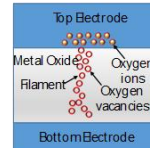
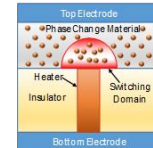
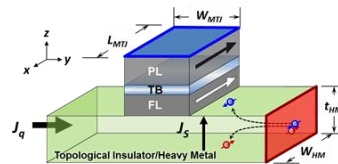
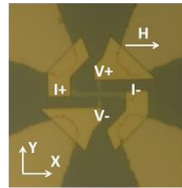
work

Algorithm-Hardware Co-Design is necessary to reap full benefits!!

Bottom-Up: Neuro-mimetic Device Design & Fabrication

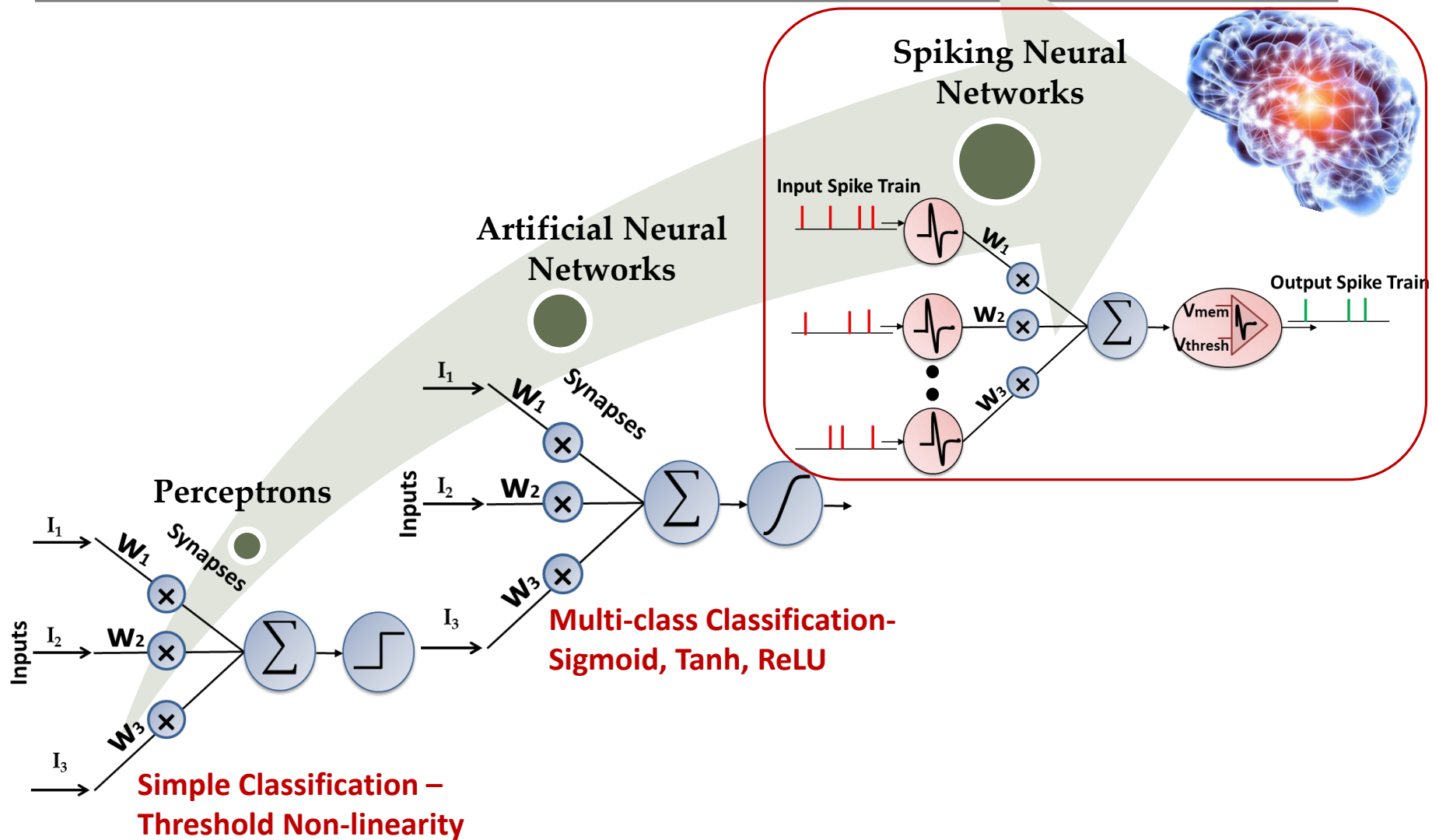
Device Fabrication TI/HM Driven MTJ Switching₁ PCM RRAM CMOS

Device physics to
mimic probabilistic
functionalities



Device-Circuit

Taking the Bio-plausible Route



Which Cues from the Brain can yield Energy-Efficiency?

Discrete Spiking (LIF)?
HH? Stochastic Neurons?

**Neuron
models**

Long-range
connections?
Feedback?

**Network
topology**



STDP? Backpropagation?

**Learning
algorithms**

**Computing
principles**

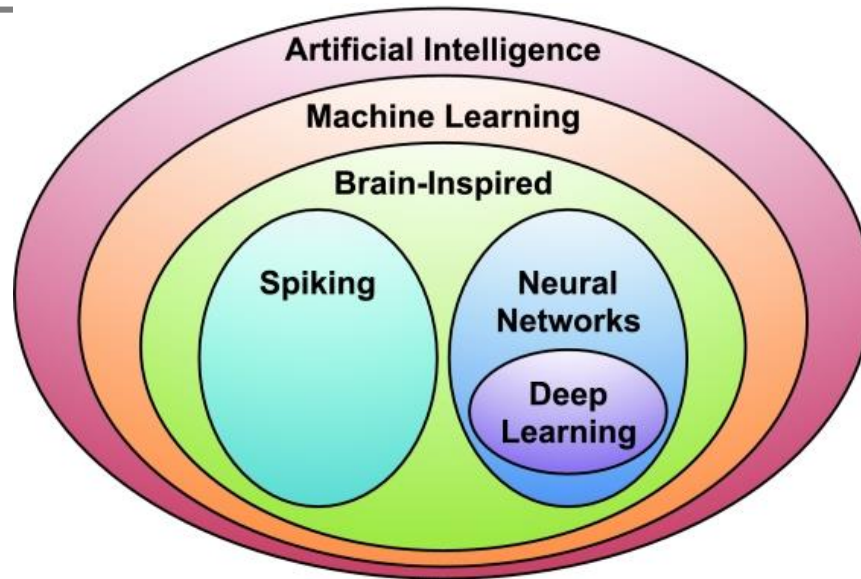
Comm. & comp. with
spikes? Tightly
integrated compute
& memory?
Asynchronous,
event-driven?
Time-based coding?

What is Machine Learning

Machine learning (ML) is the [scientific study](#) of [algorithms](#) and [statistical models](#) that [computer systems](#) use to progressively improve their performance on a specific task. Machine learning algorithms build a mathematical model of sample data, known as "[training data](#)", in order to make predictions or decisions “without being explicitly programmed” to perform the task.

"A computer program is said to learn from experience E with respect to some class of tasks T and performance measure P if its performance at tasks in T , as measured by P , improves with experience E ."

AI, ML, Deep Learning, Neural Networks, Spiking

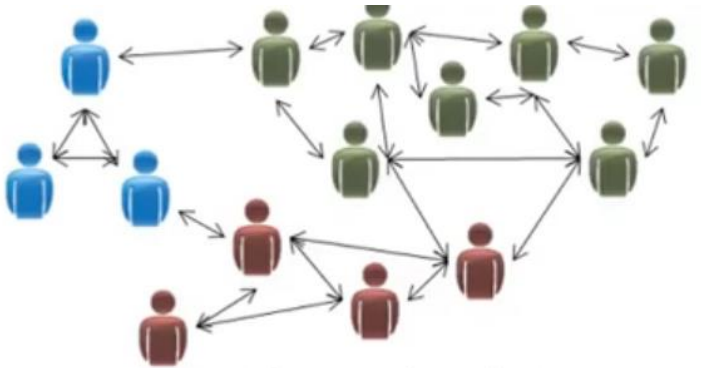


**Solving some aspect of
'intelligence'**

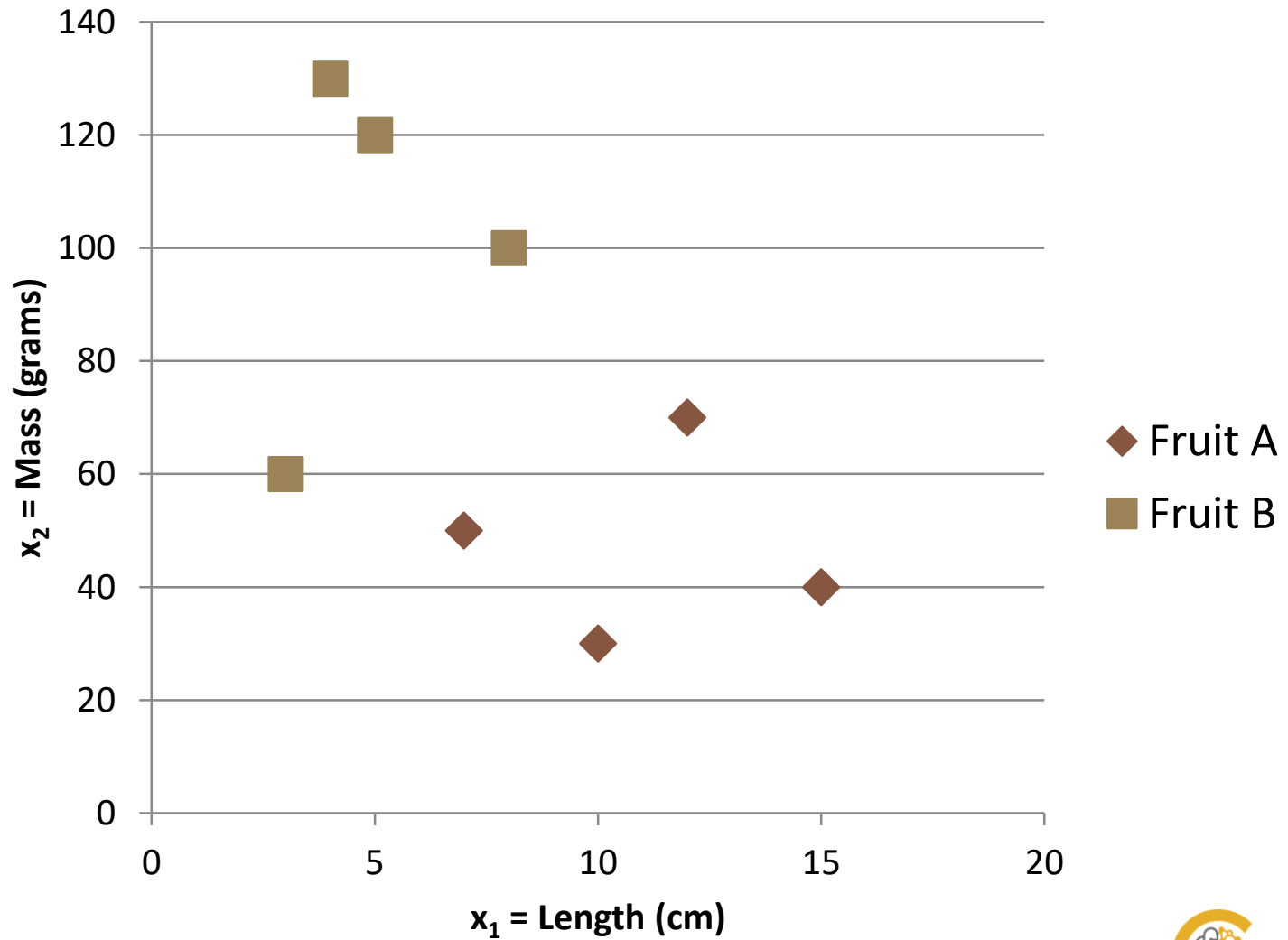
- **Machine Learning** - Field of study that gives computers ability to learn without being explicitly programmed
- **Brain-inspired computation** is a program or algorithm that takes some aspects of its basic form or functionality from the way the brain works ('Brain' is the best source of inspiration for intelligent applications)
- **Neural Networks, Deep Learning** have dot product or weighted summation of inputs, notion inspired from brain's synaptic/neuronal configuration
- **Spiking** – More deeply inspired from brain-like computations 'spikes' or 'events'

ML algorithms

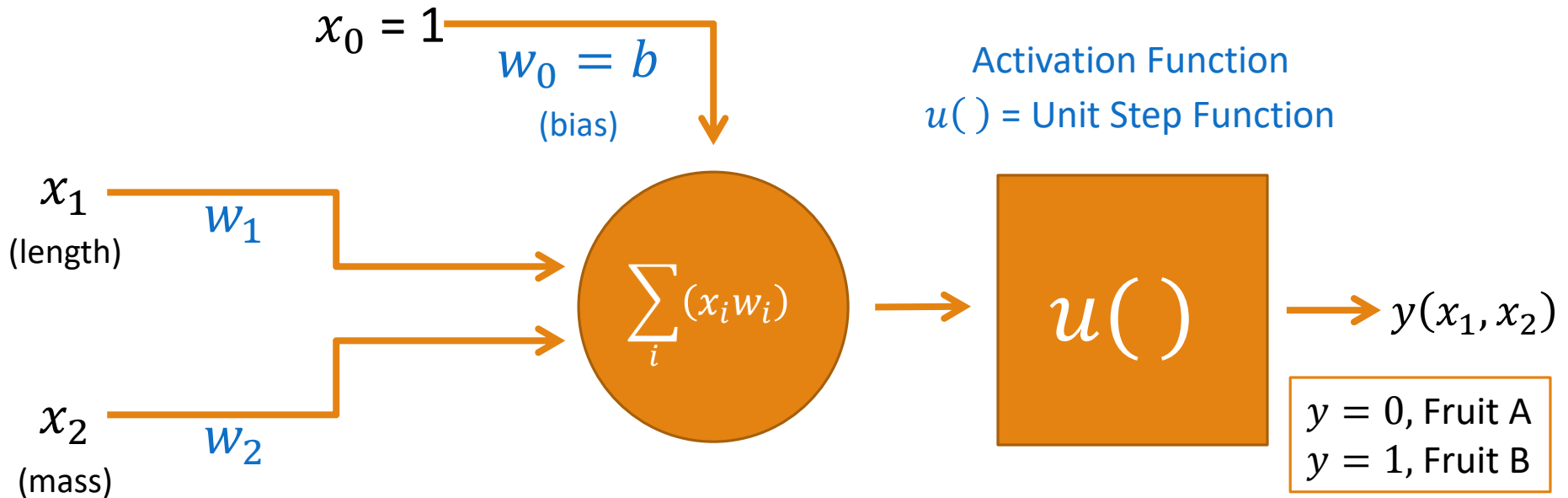
- Supervised Learning – Regression, Classification (Labels given)
- Unsupervised Learning – Clustering (No labels)
 - Application – Google News, Social Network Analysis, Market Segmentation etc.
- Others: Reinforcement Learning, Semi-supervised Learning



VERY SIMPLE NEURAL NETWORK: FRUIT CLASSIFICATION



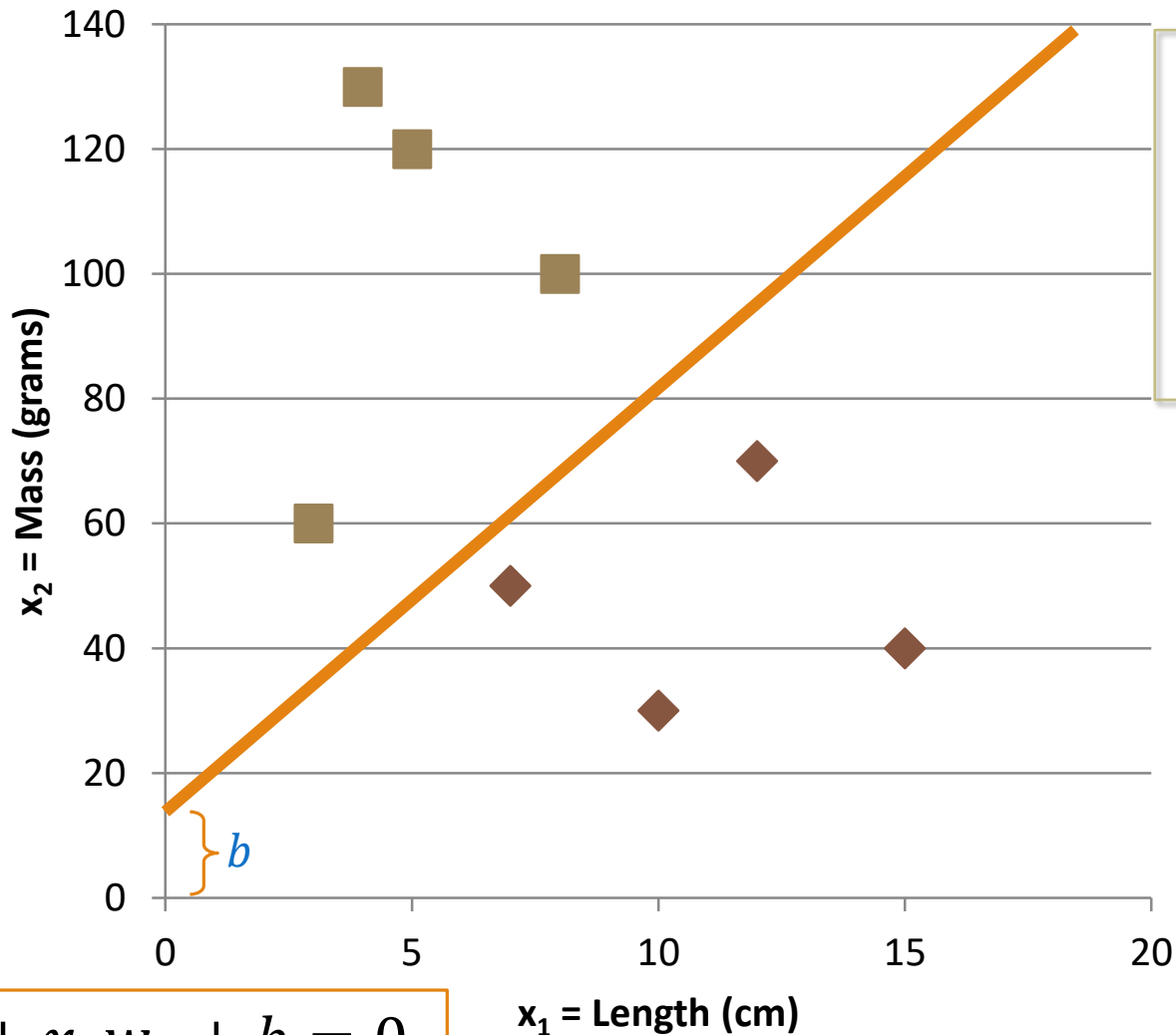
EXAMPLE: FRUIT CLASSIFICATION



Note: a single neuron ANN is called a “perceptron” and can model only linear functions.

$$x_1 w_1 + x_2 w_2 + b = 0$$

EXAMPLE: FRUIT CLASSIFICATION



Note: adding additional layers to the NN allows for non-linear function modeling.

◆ Fruit A
■ Fruit B

$$b = 15$$
$$w_1 = 6.5$$
$$w_2 = -1$$

$$x_1 w_1 + x_2 w_2 + b = 0$$