

CS 758: Advanced Topics in Computer Architecture

Lecture #11: Intro to DNNs

Professor Matthew D. Sinclair

Some of these slides were developed by Tushar Krishna at Georgia Tech
and Joel Emer and Vivienne Sze at MIT

Slides enhanced by Matt Sinclair

AI/ML Applications

“AI is the new electricity” – Andrew Ng

Object Detection

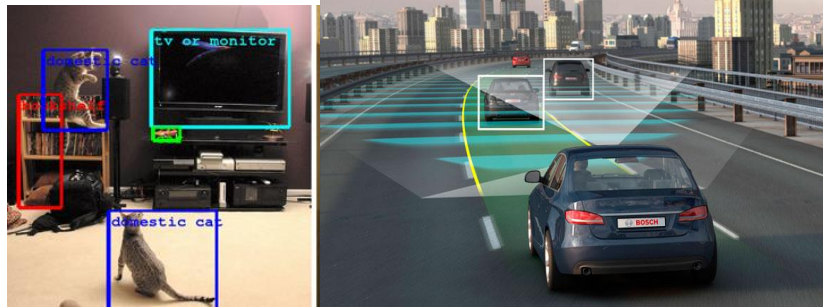
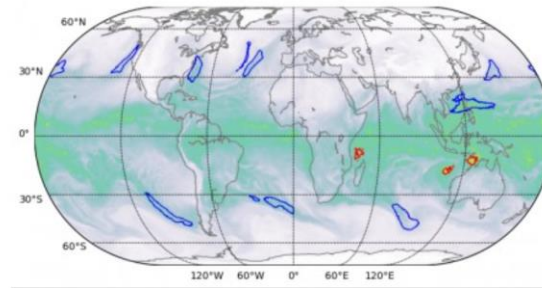
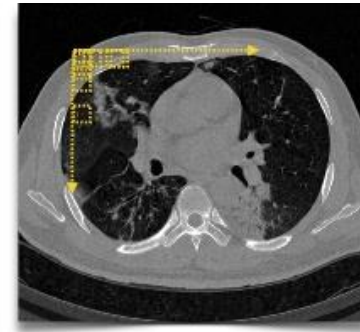


Image Segmentation



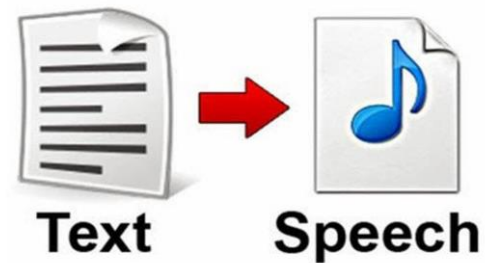
Medical Imaging



Speech Recognition



Text to Speech



Recommendations

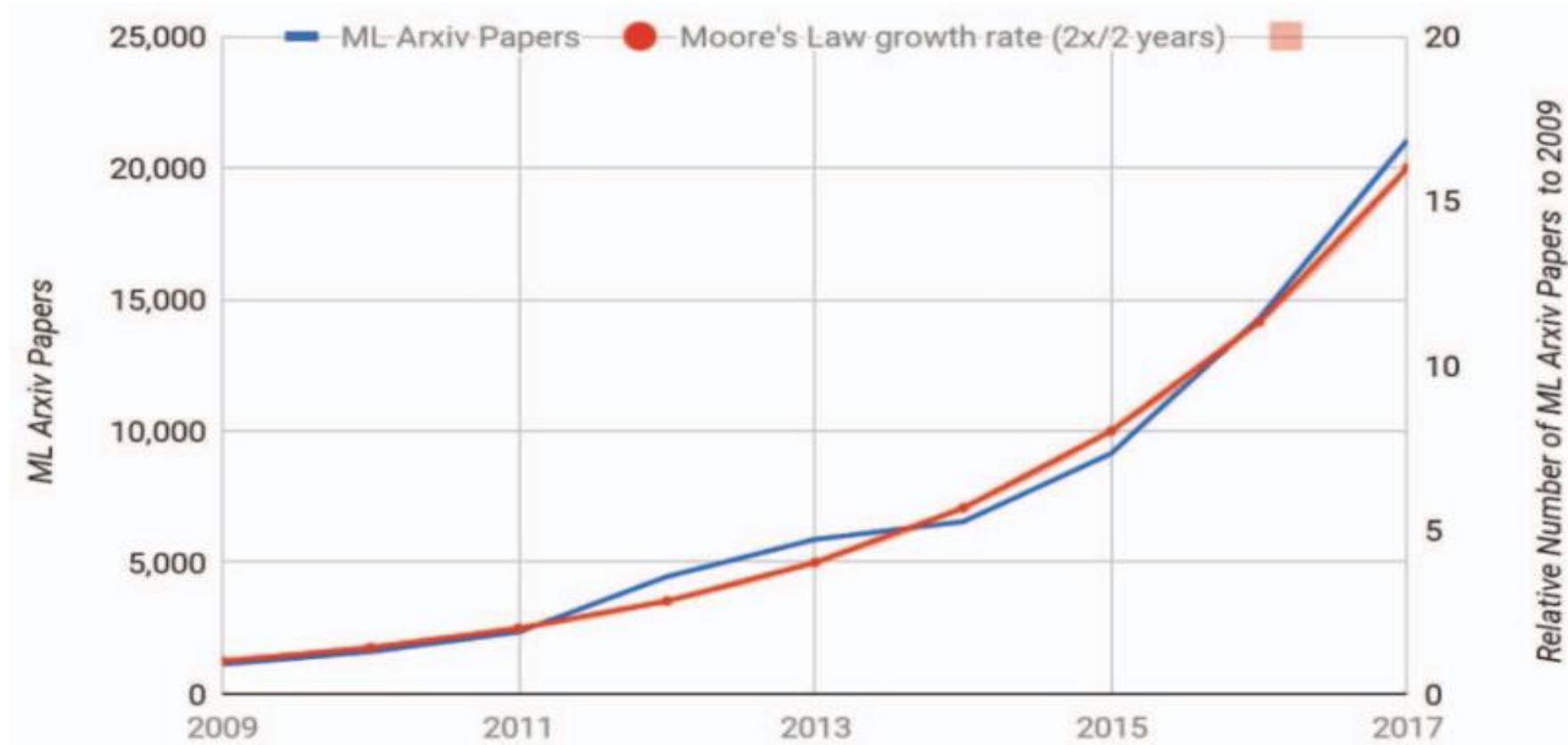


Games



Autonomous Agents, Intelligent Personal Assistants, ...

ML is HOT

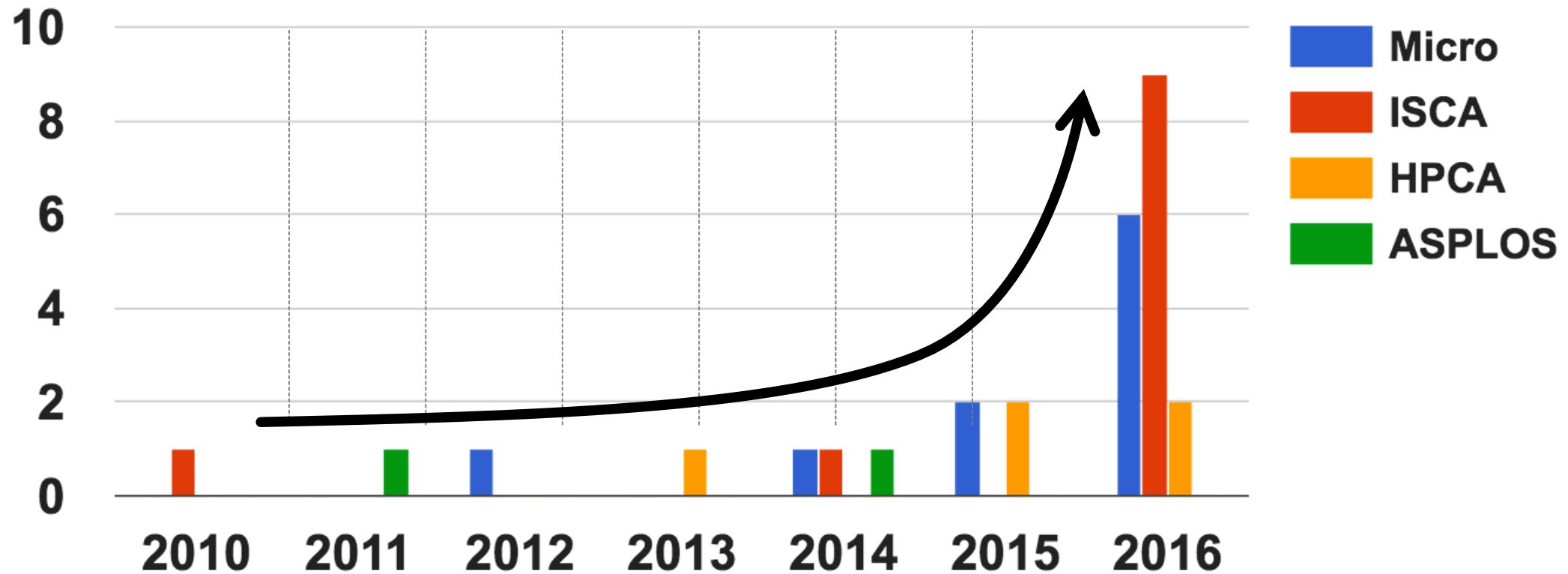


Source: Jeff Dean, Google

Over 1.5B invested in AI chip startups in 2017 alone
Not clear what the endgame will be

Publications at Architecture Conferences

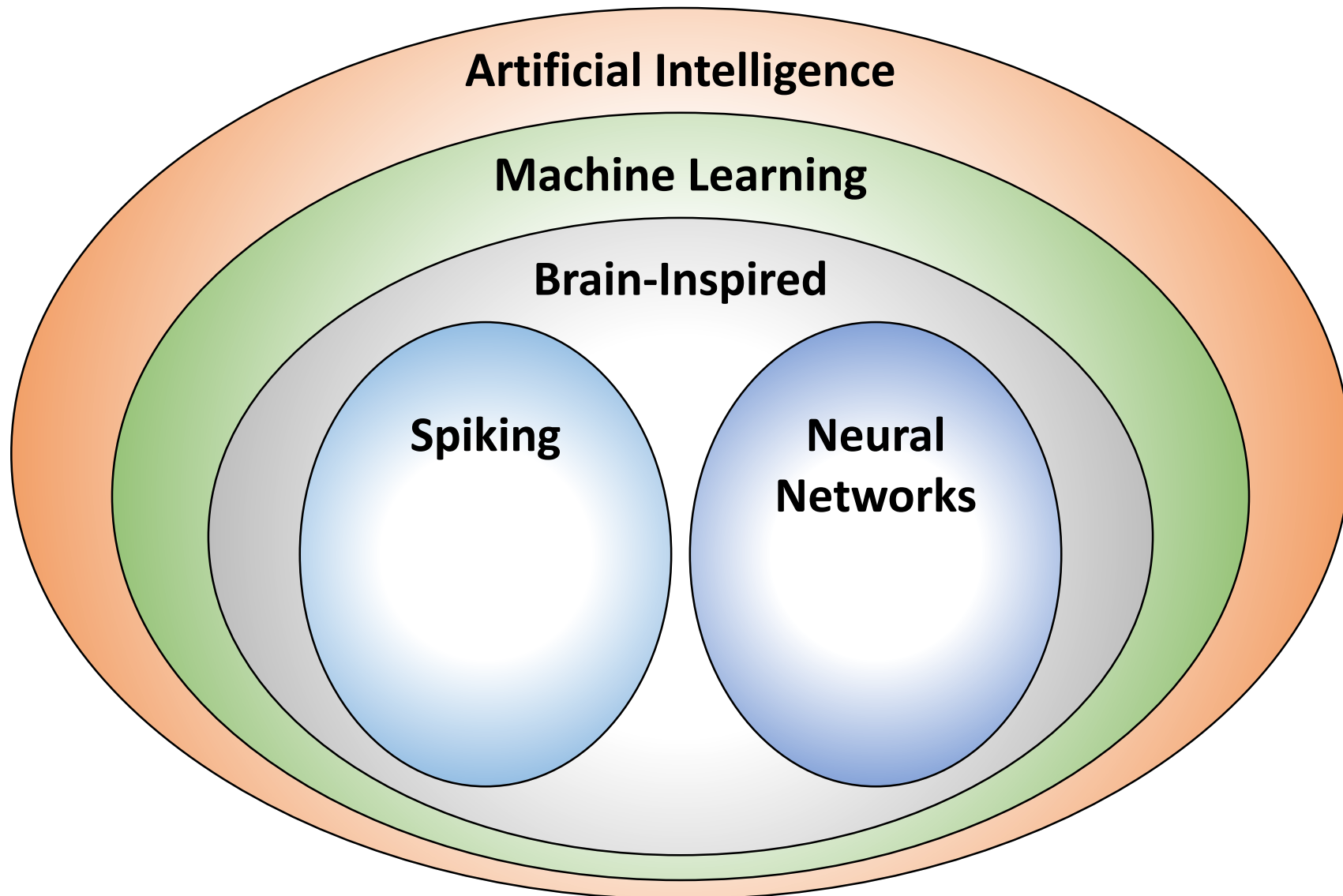
of Publications over the Years



Key Questions to Answer

- Separating the Hype from the (Remaining) Opportunity
- Computer Architecture perspective
 - The role of HW Acceleration in the AI/ML Boom
 - Understand Compute and Memory behavior of Deep Learning Workloads
 - Understand Limitations of Current Platforms (CPU and GPU)
 - What about Moore's Law?
 - Opportunities and Challenges with custom HW Accelerators
 - Techniques to reduce computation
 - Techniques to to reduce data movement
 - Techniques to to reduce memory overhead
 - Performance vs Energy vs \$\$

Machine Learning with Neural Networks



Slide Courtesy: Joel Emer and Vivienne Sze

NeurON: Weighted Sum

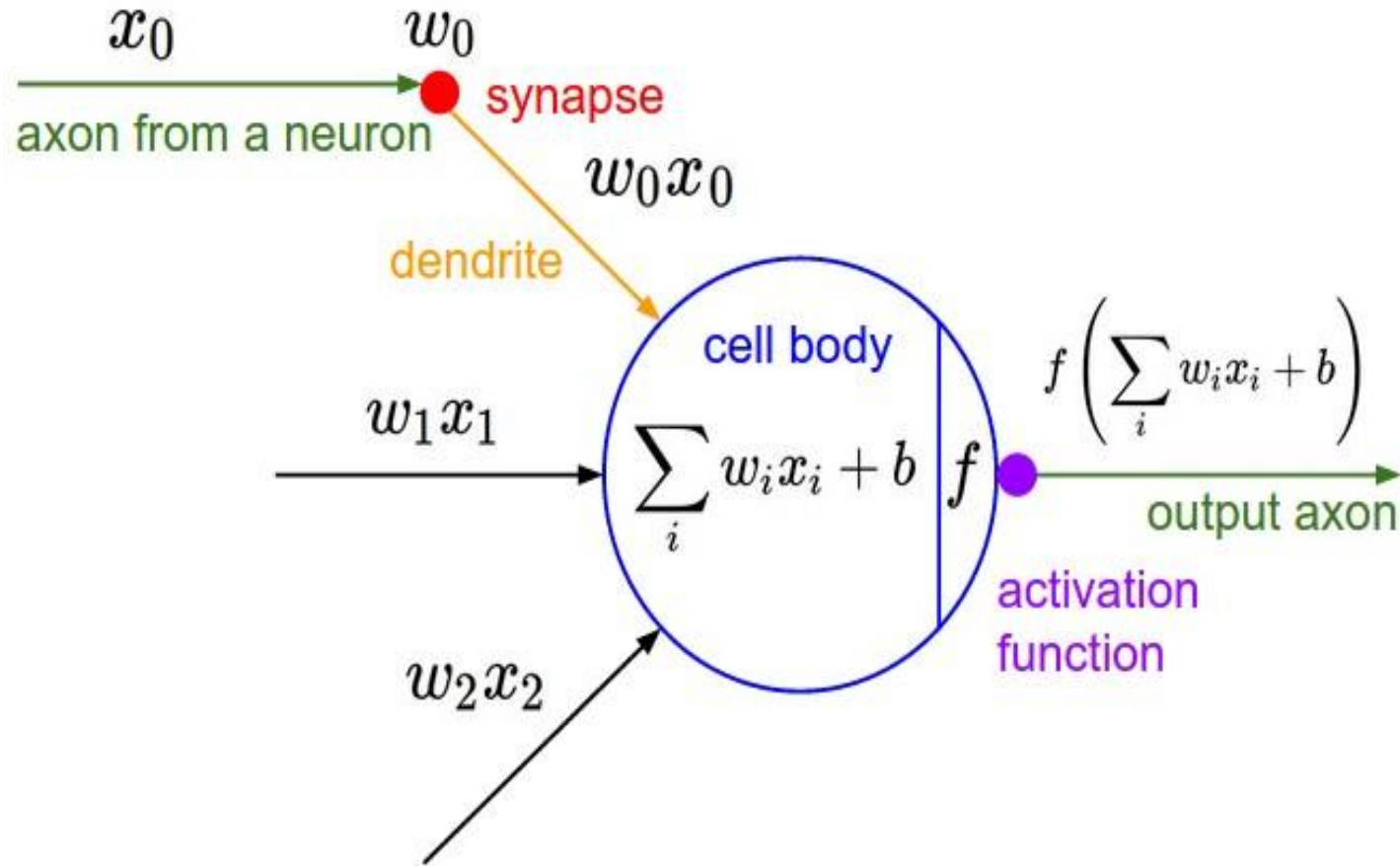
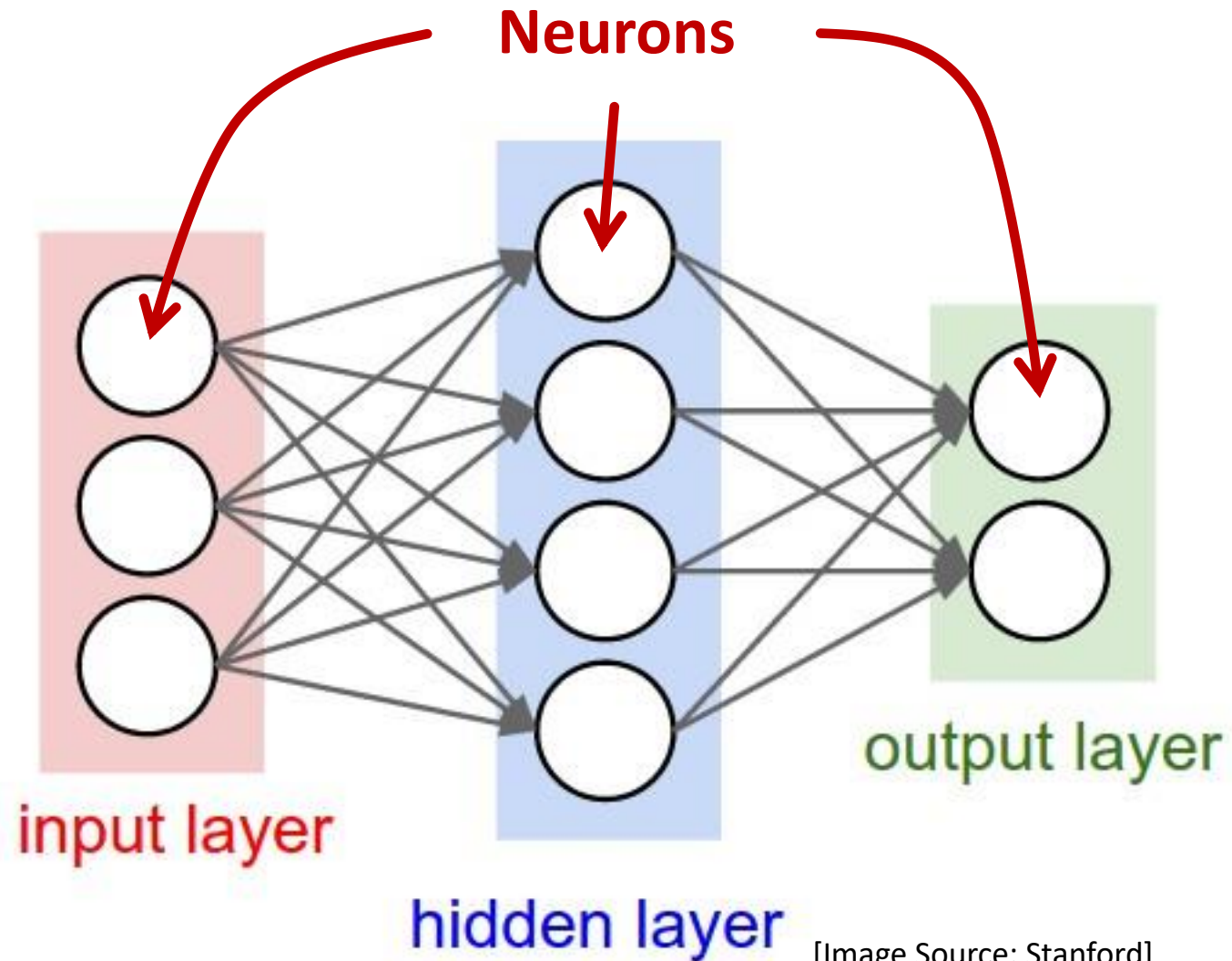


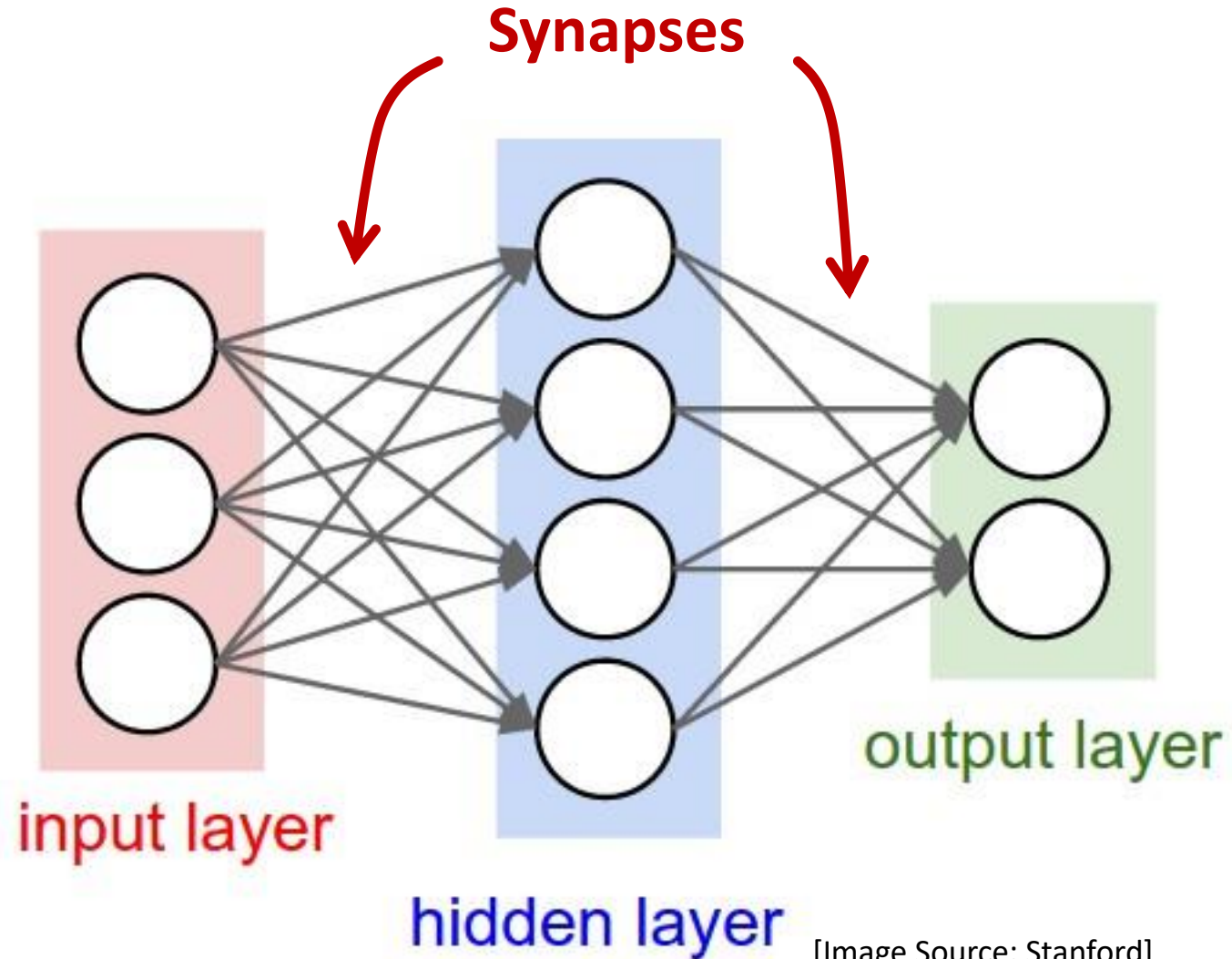
Image Source: Stanford

NEURAL NETWORK



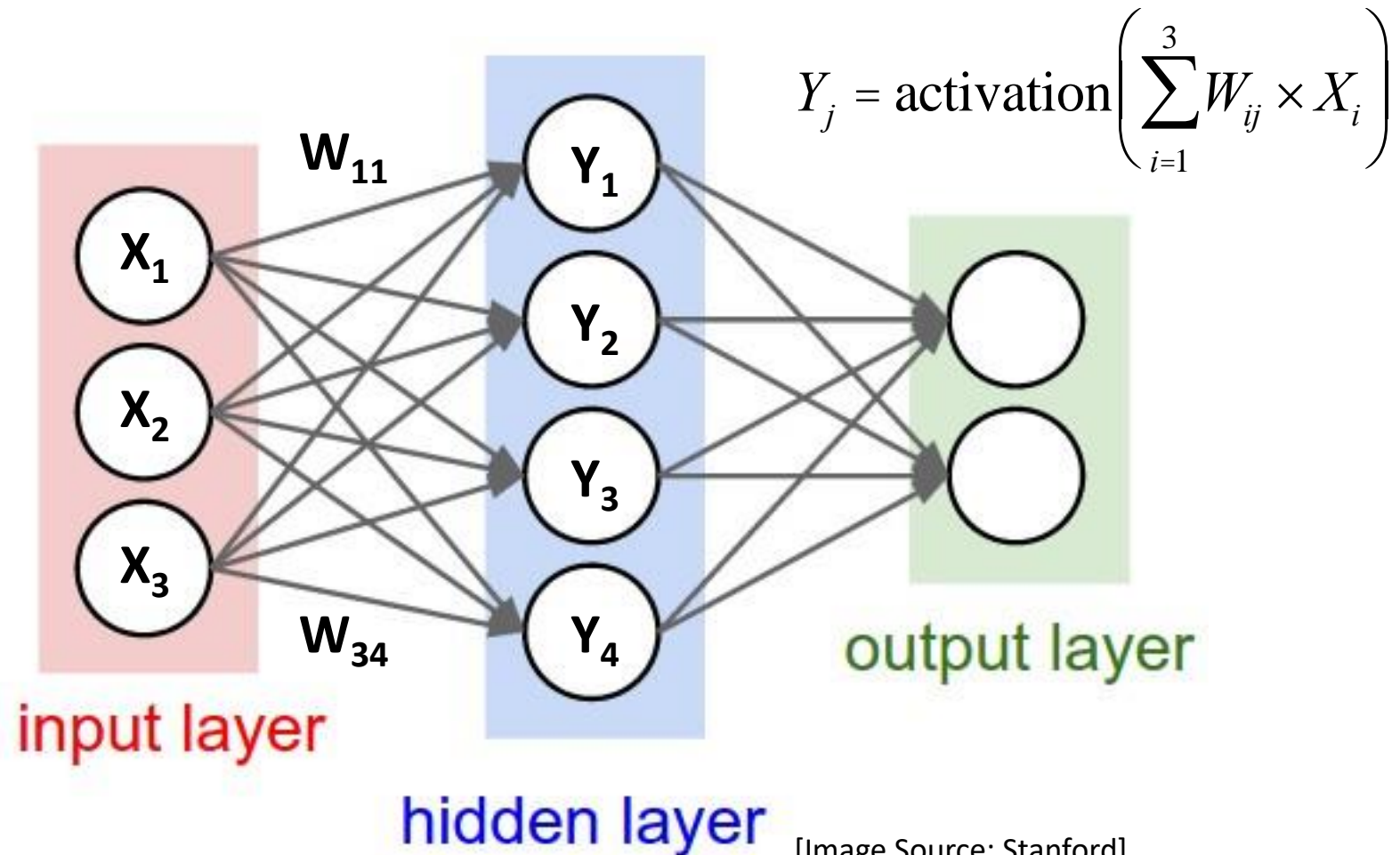
[Image Source: Stanford]

NEURAL NETWORK



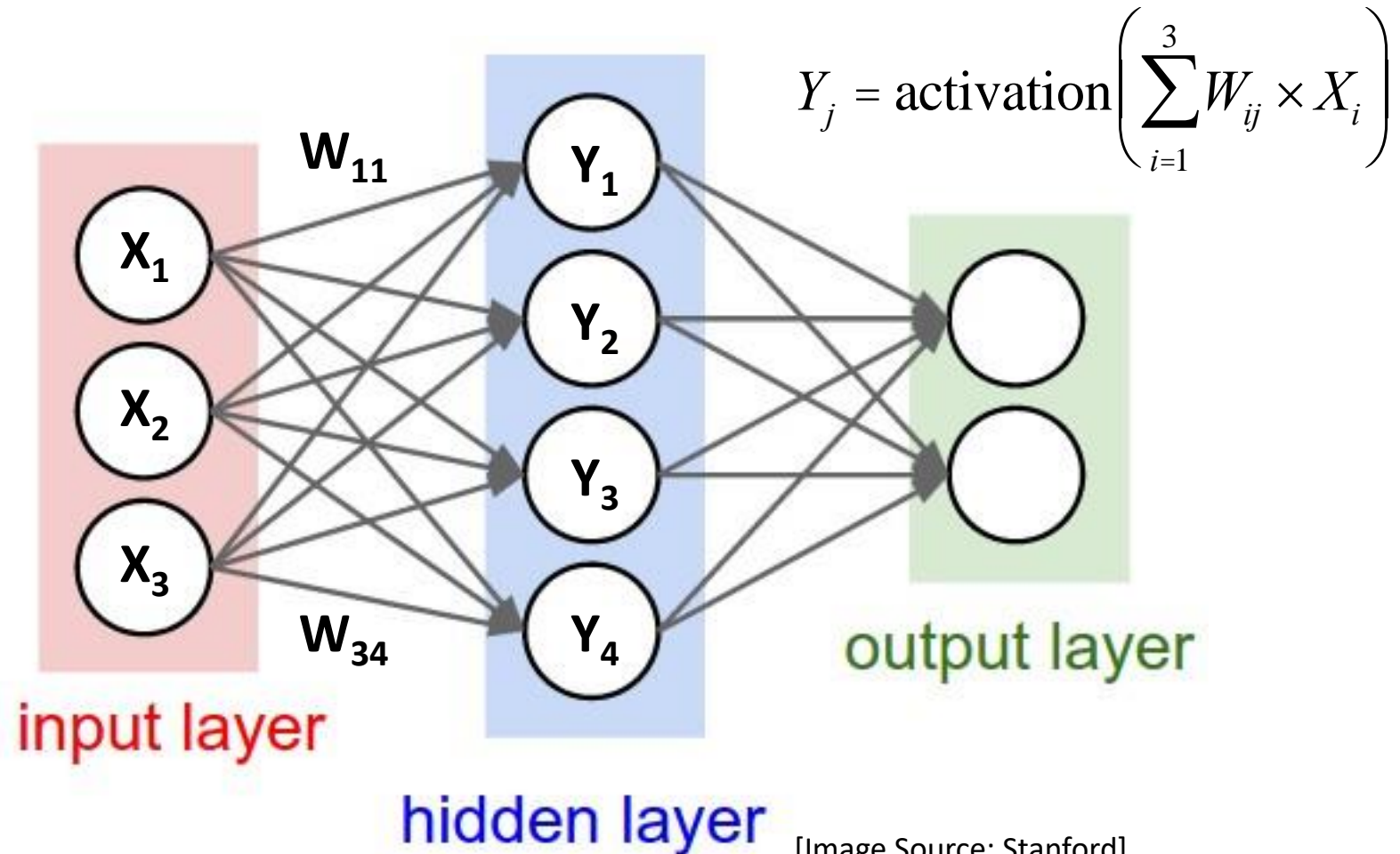
NEURAL NETWORK: Multiple weighted sums

Each **synapse** has a **weight** for neuron **activation**



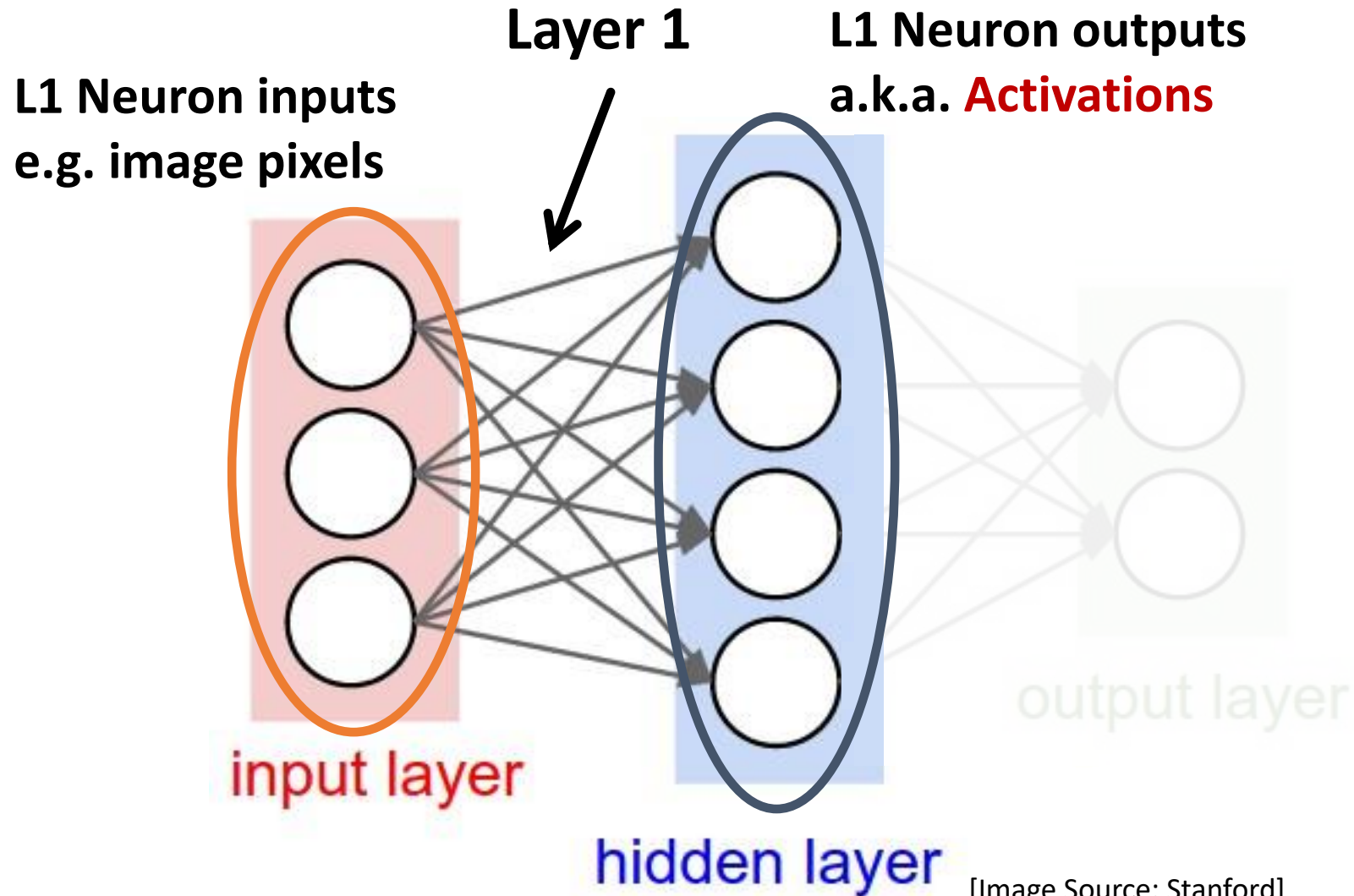
NEURAL NETWORK: Multiple weighted sums

Weight Sharing: multiple synapses use the **same weight value**



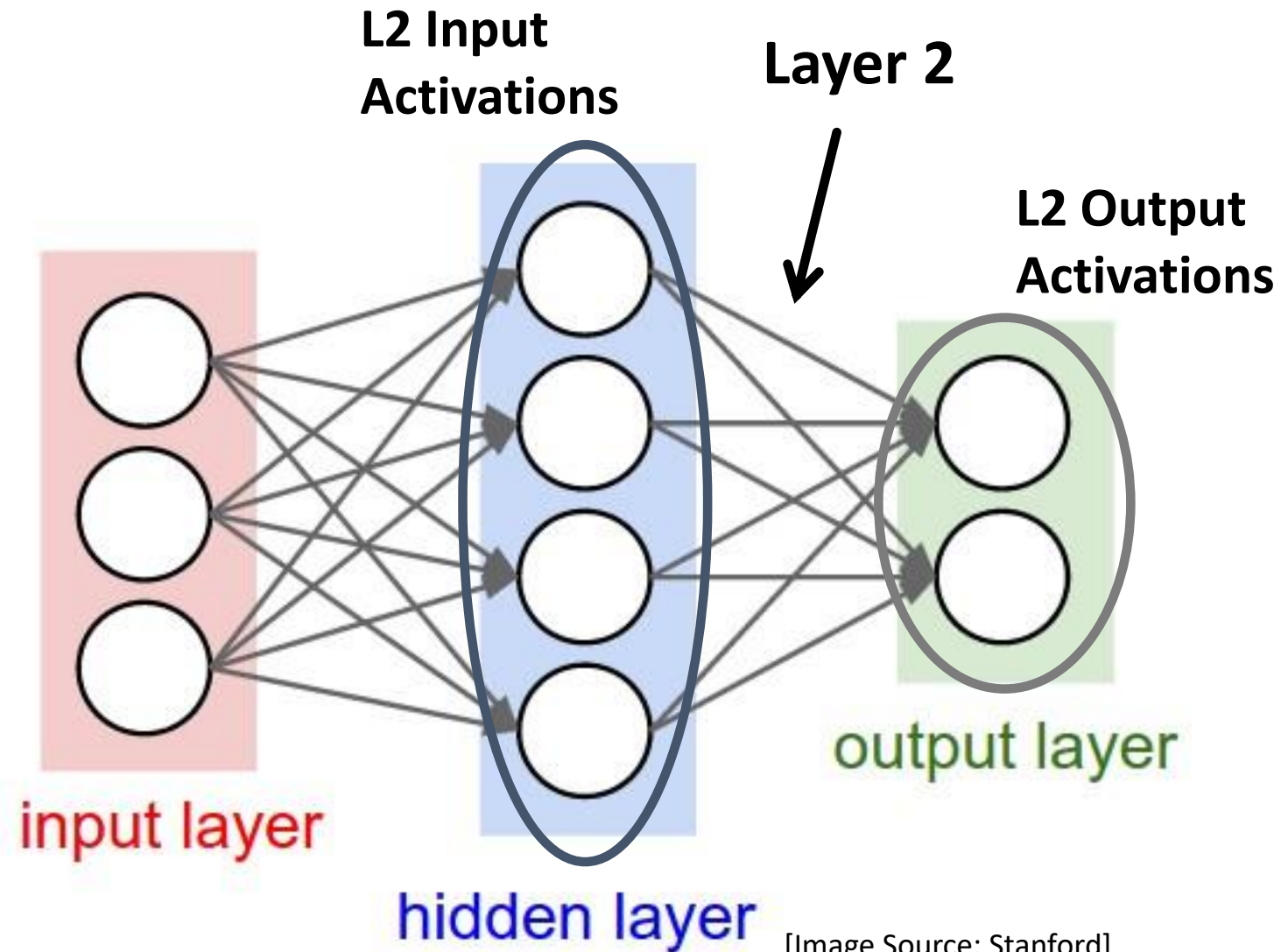
[Image Source: Stanford]

NEURAL NETWORK Terminology



[Image Source: Stanford]

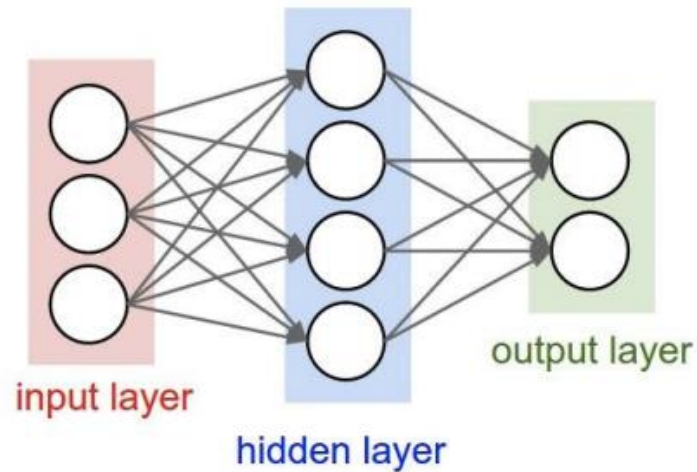
NEURAL NETWORK Terminology



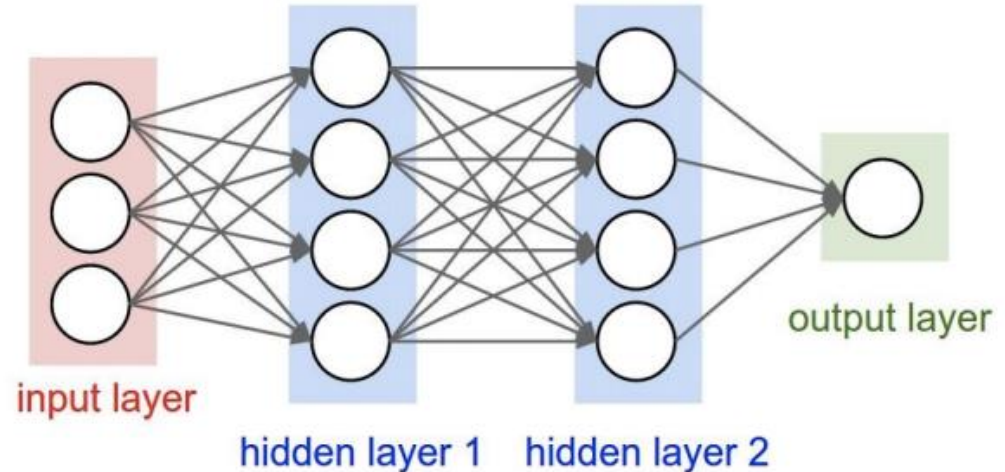
[Image Source: Stanford]

NEURAL NETWORK Terminology

A **layer** can refer to a set activations or a set of weights.
In this class, we use **layer** to refer to a set of weights.



“2-layer Neural Net”, or
“1-hidden-layer Neural Net”

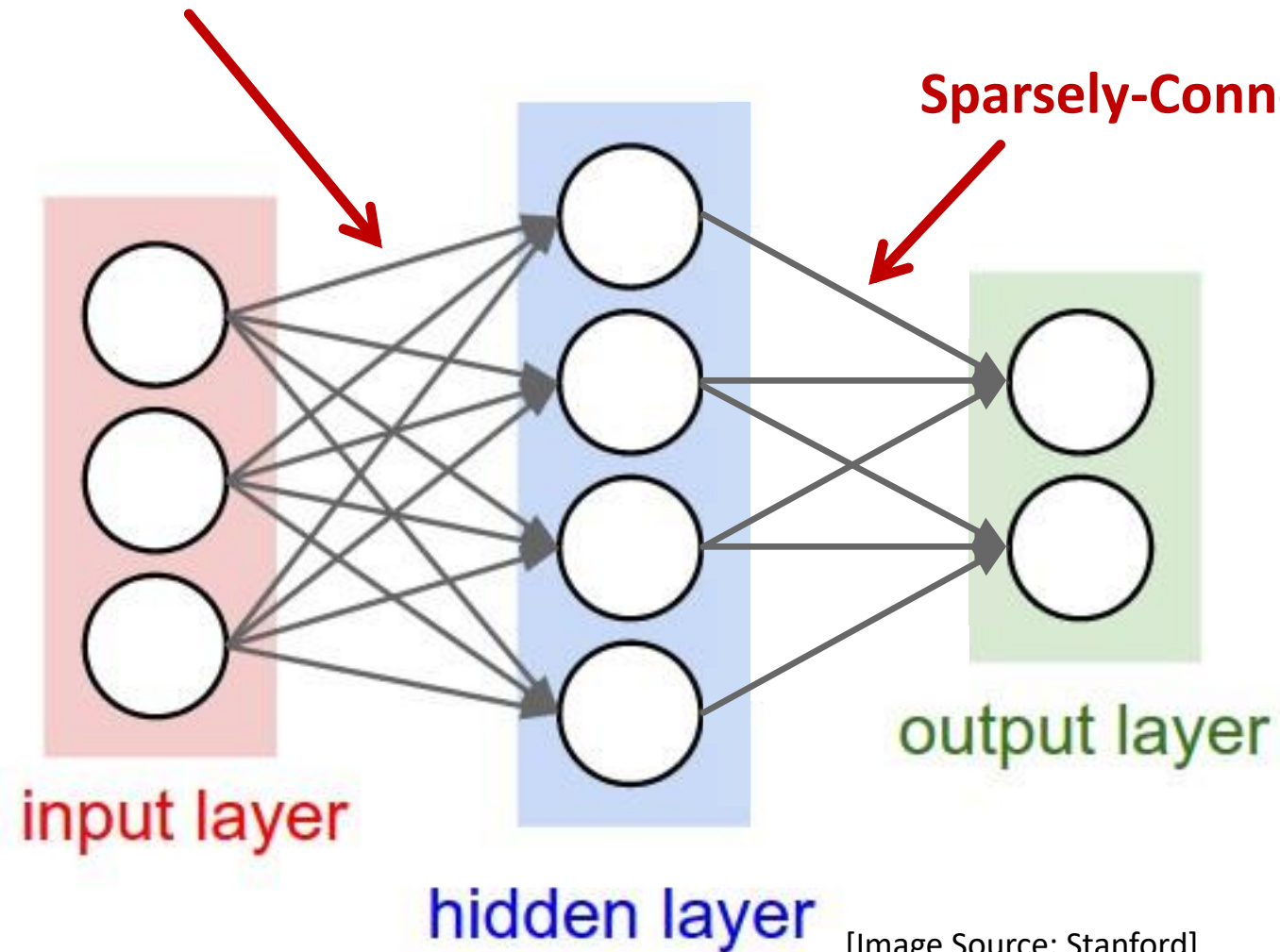


“3-layer Neural Net”, or
“2-hidden-layer Neural Net”

[Image Source: Stanford]

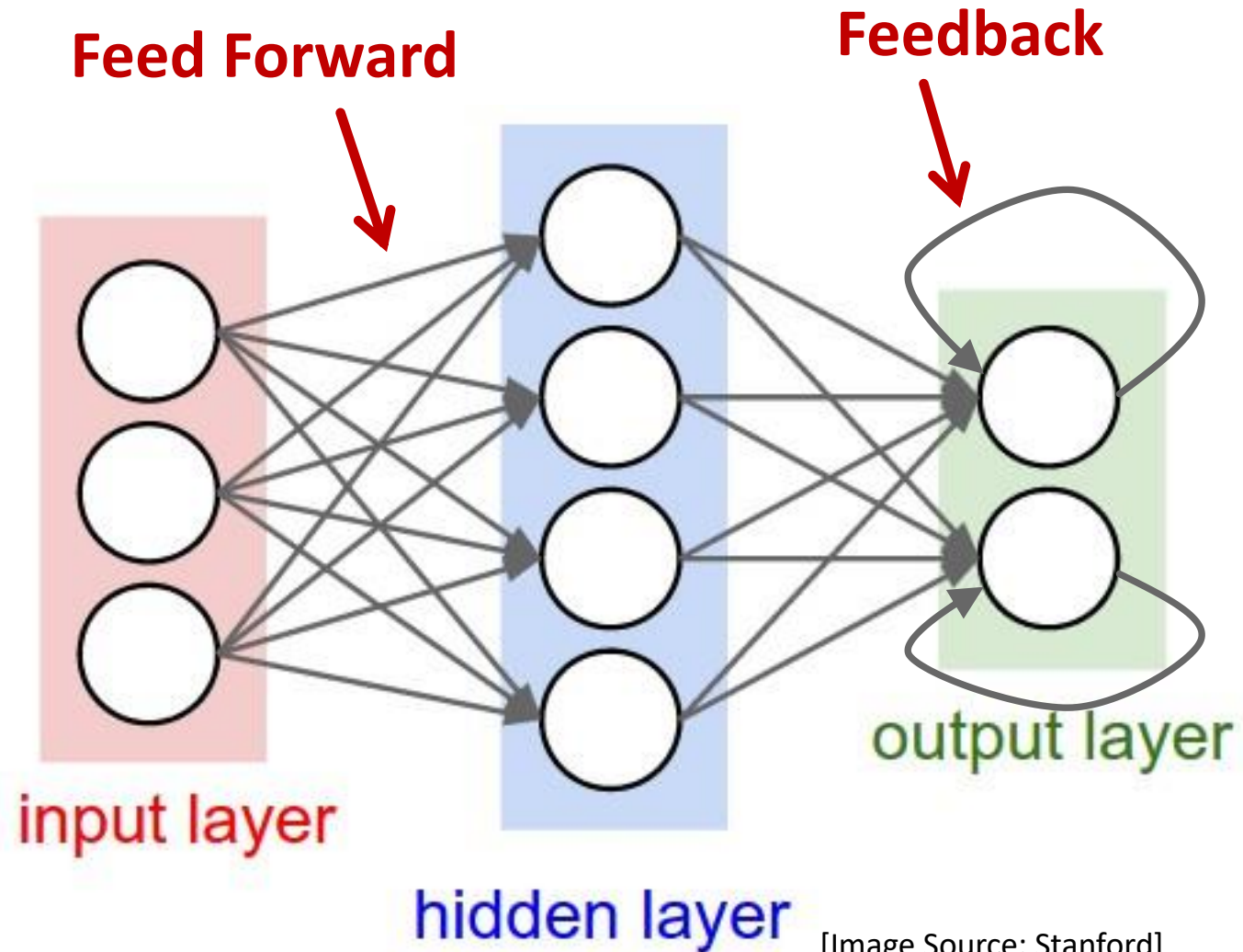
NEURAL NETWORK Terminology

Fully-Connected: all input neurons connected to all output neurons



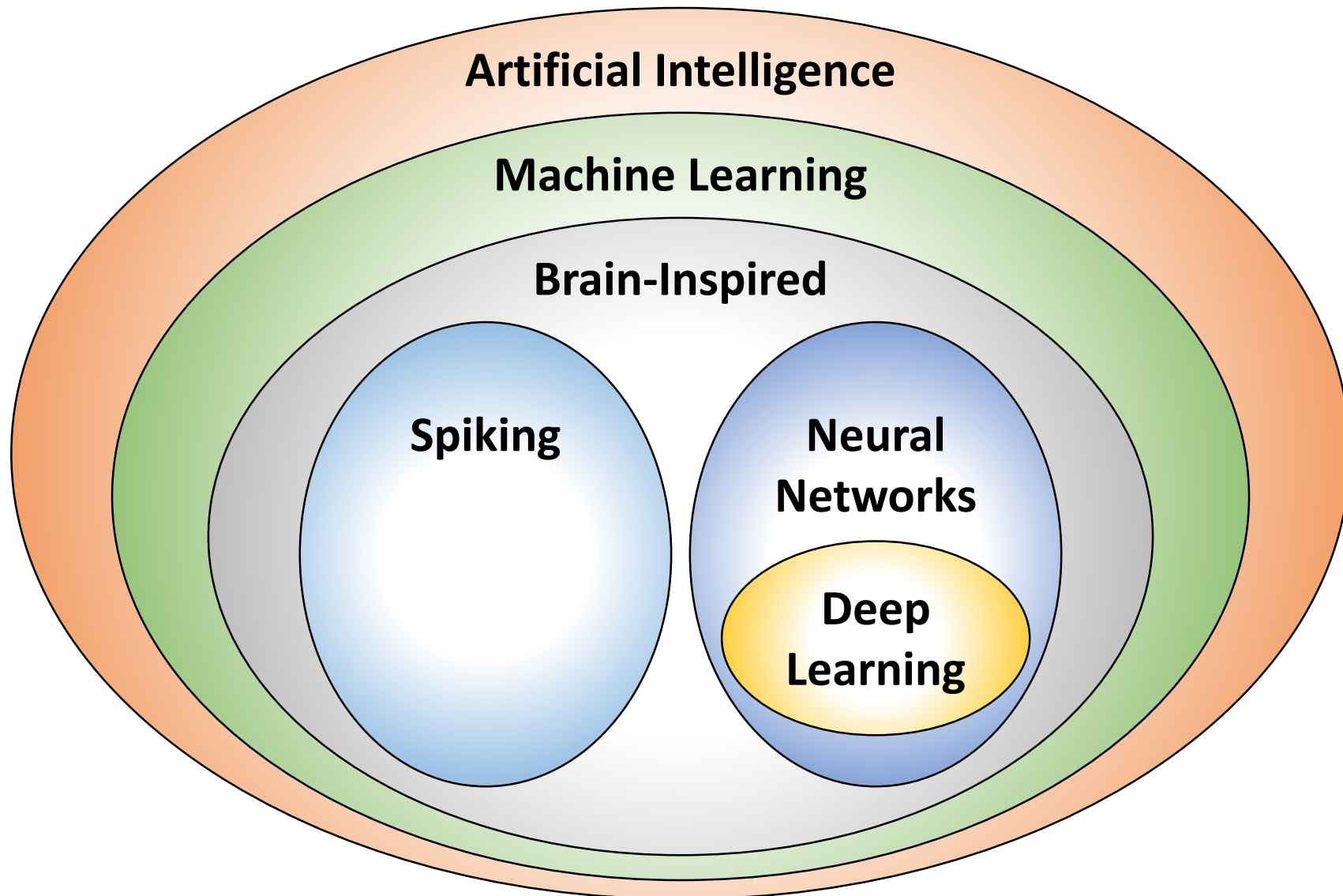
[Image Source: Stanford]

NEURAL NETWORK Terminology



[Image Source: Stanford]

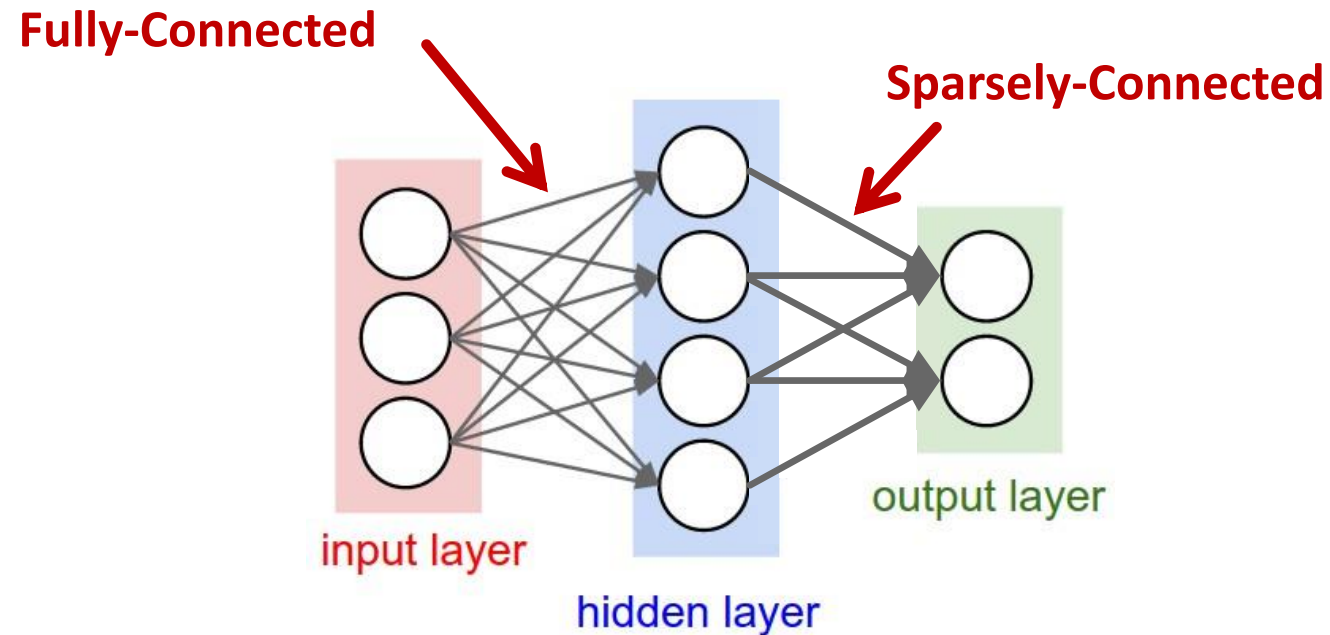
Deep Learning



Slide Courtesy: Joel Emer and Vivienne Sze

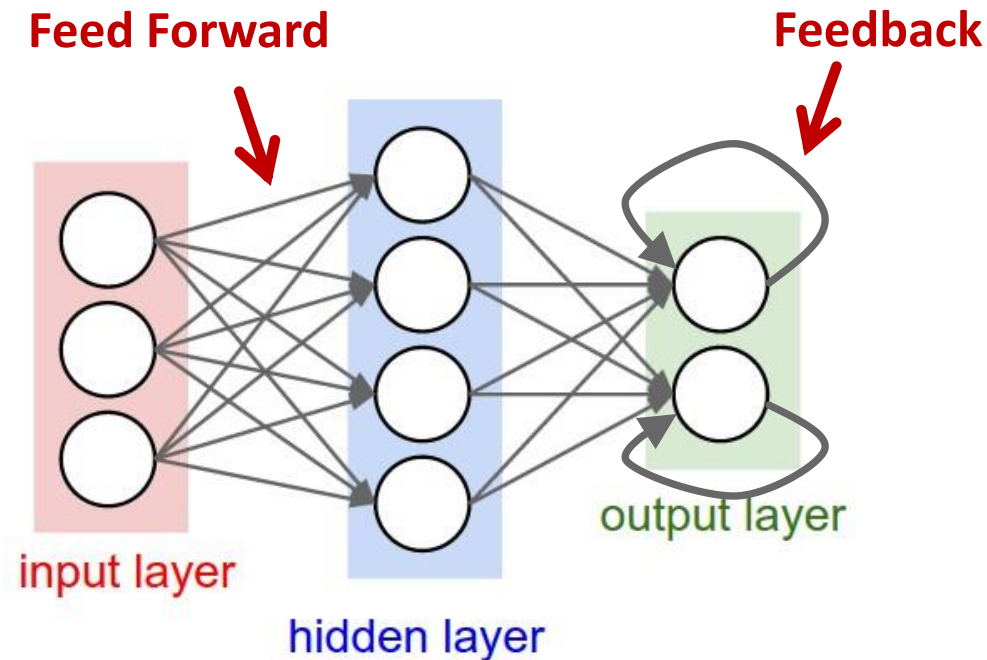
Popular Types of DEEP NEURAL NETWORKS

- Fully-Connected NN
 - feed forward, a.k.a. multilayer perceptron (MLP)
- Convolutional NN (CNN)
 - feed forward, sparsely-connected w/ weight sharing

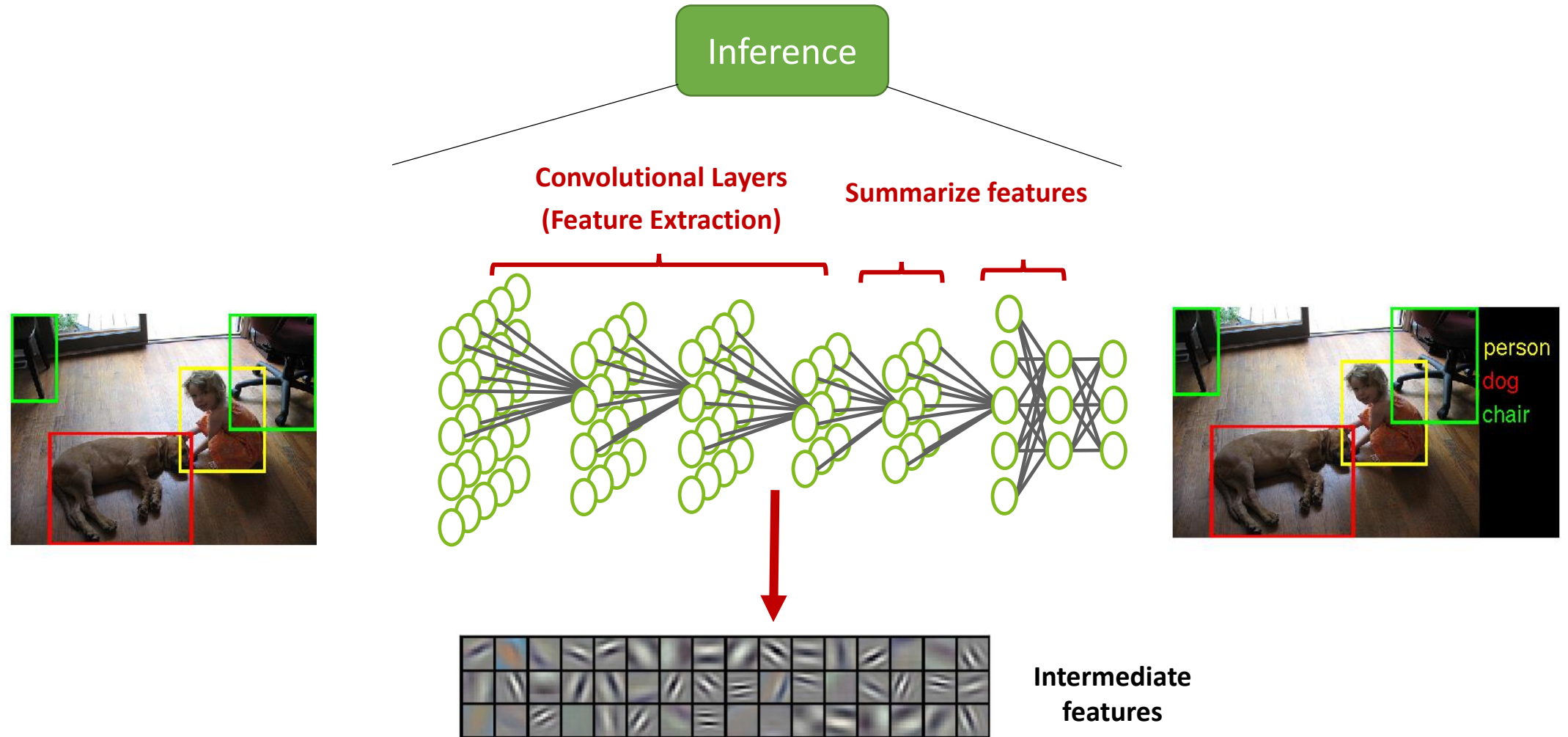


Popular Types of DEEP NEURAL NETWORKS

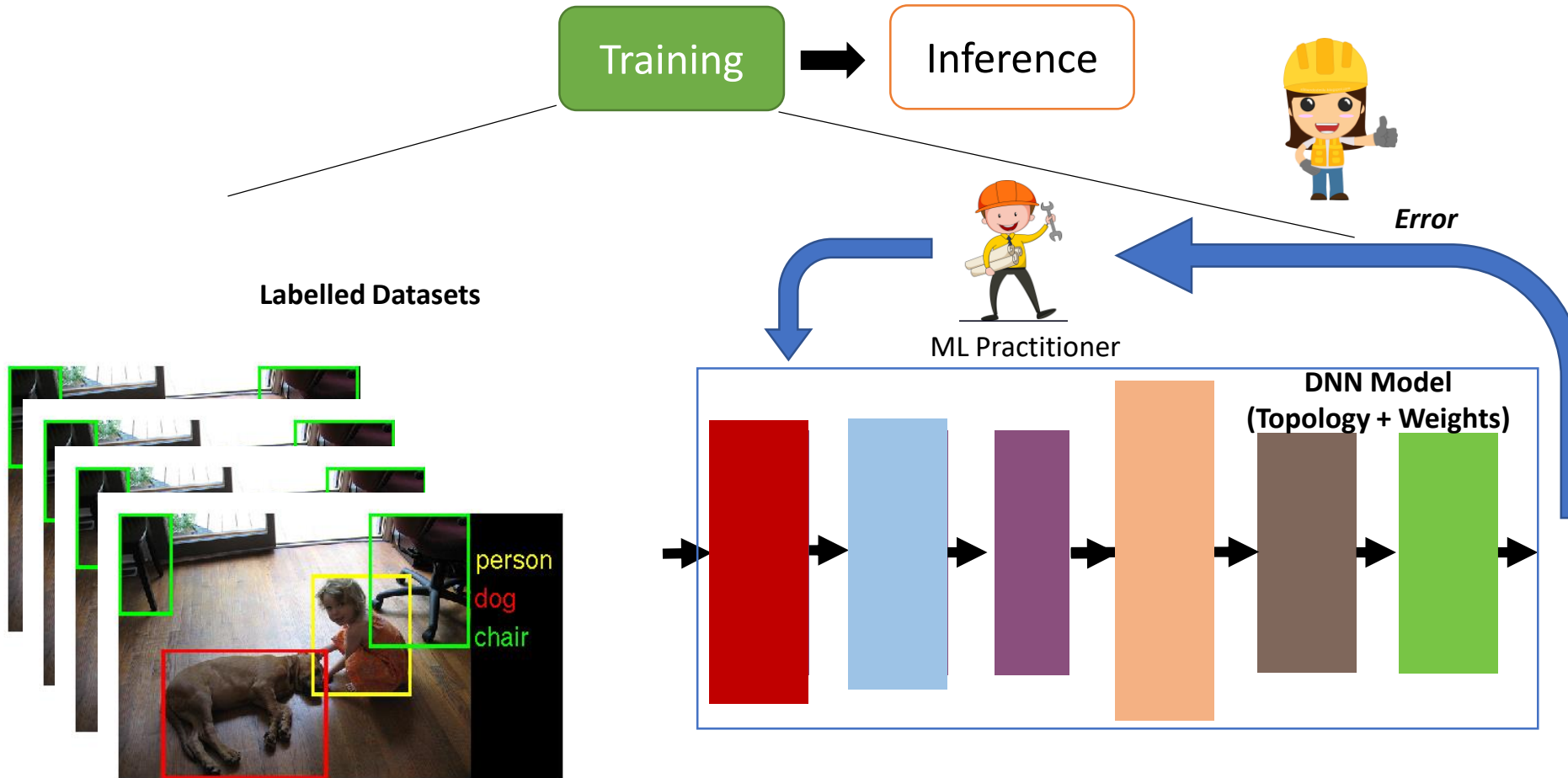
- Recurrent NN (RNN)
 - feedback
- Long Short-Term Memory (LSTM)
 - feedback + storage



Deep Learning Landscape



Deep Learning Landscape



Why is deep learning hot now?

**Big Data
Availability**

**New ML
Techniques**

**GPU
Acceleration**

facebook

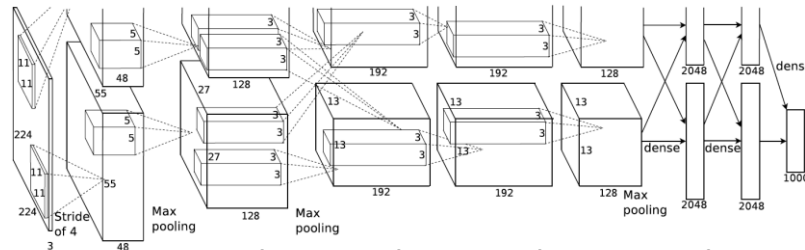
350M images
uploaded per
day

You Tube

300 hours of
video uploaded
every minute

Walmart

2.5 Petabytes of
customer data
hourly



Convolutional Neural Network

How deep learning started: ImageNet Challenge

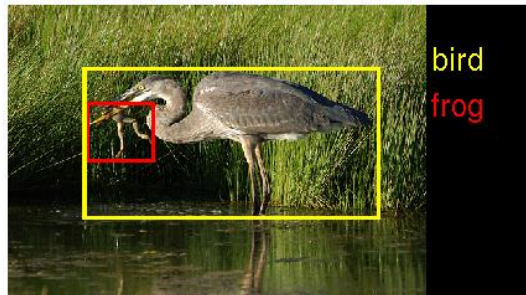
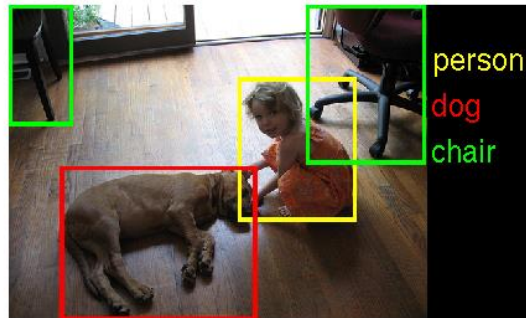


Image Classification Task:

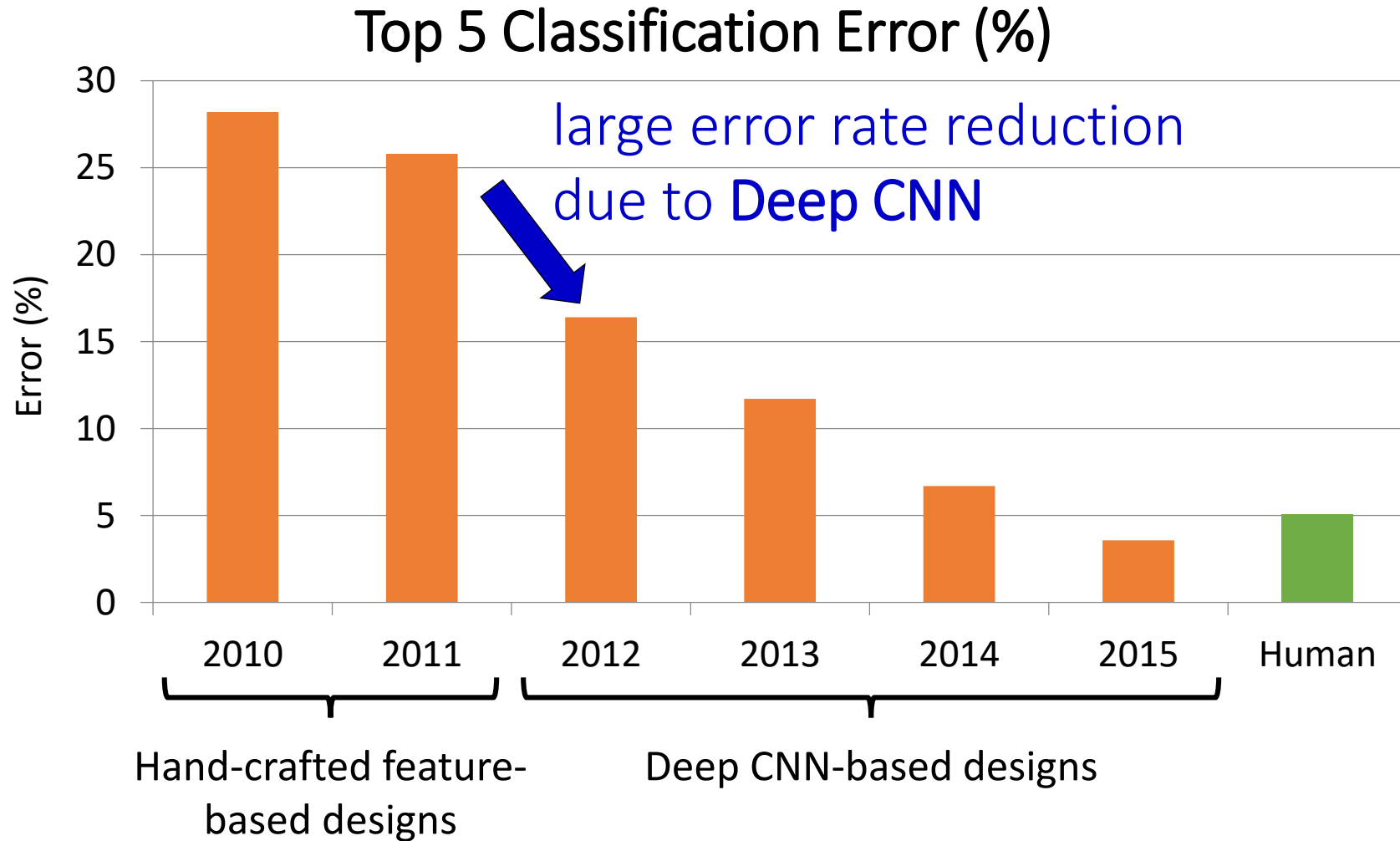
1.2M training images • 1000 object categories

Object Detection Task:

456k training images • 200 object categories

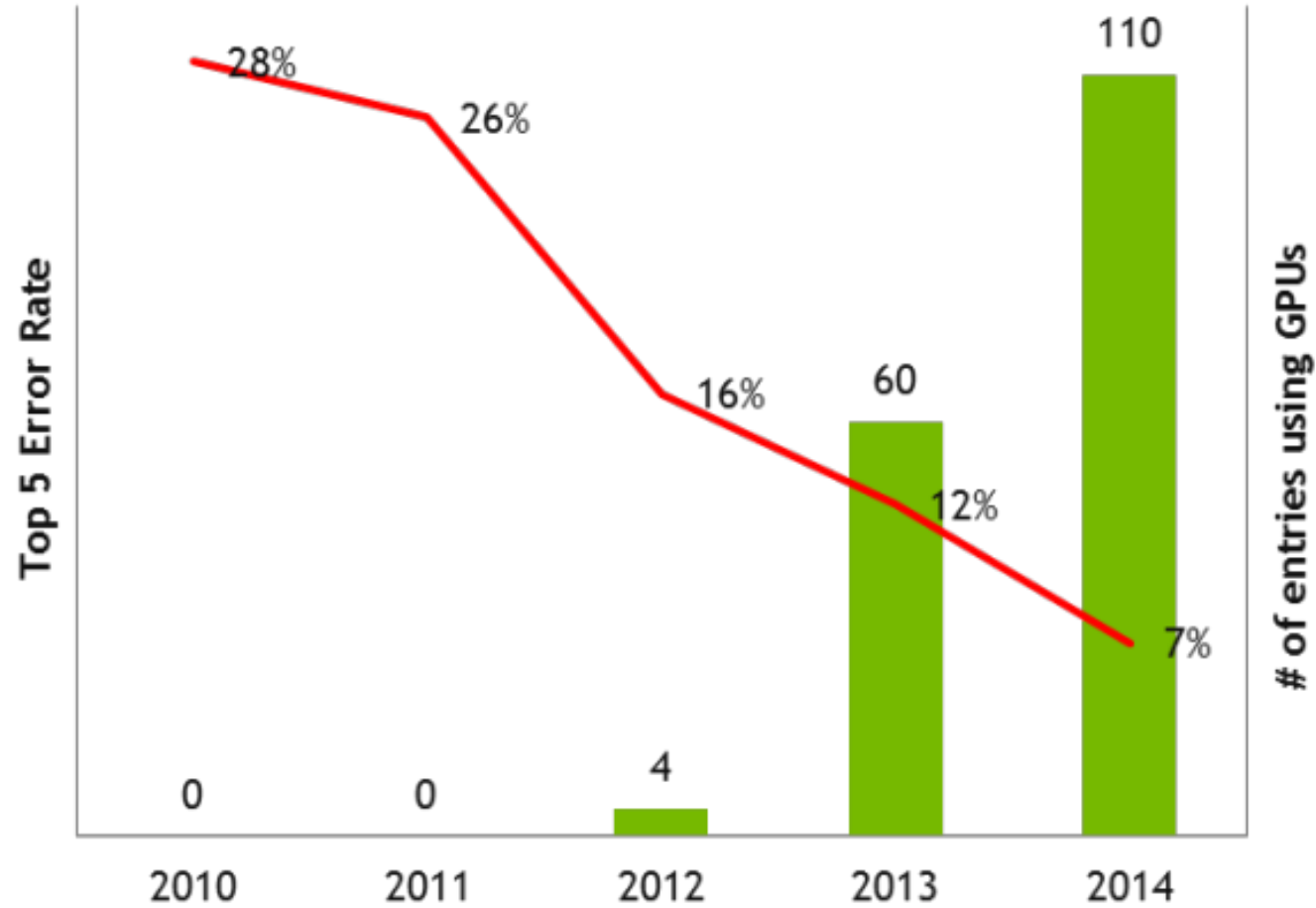


ImageNet: Image Classification Task



[Russakovsky et al., IJCV 2015]

GPU Usage for ImageNet Challenge



Mature Applications

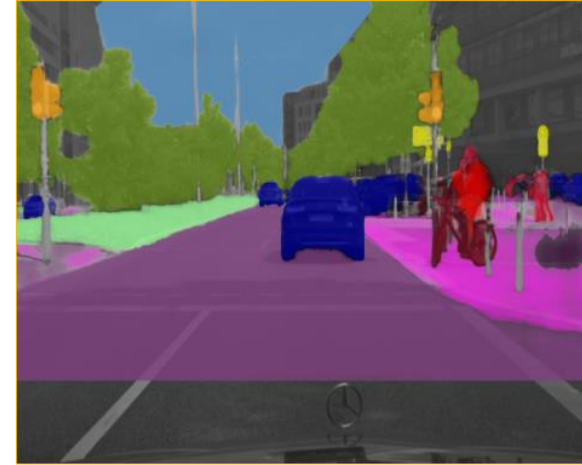
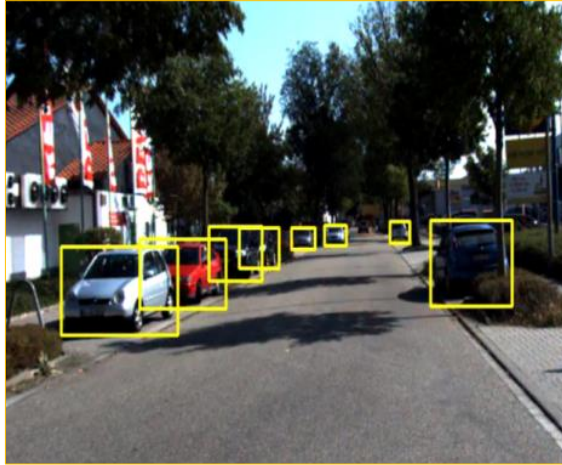
- **Image**
 - Classification: image to object class
 - Recognition: same as classification (except for faces)
 - Detection: assigning bounding boxes to objects
 - Segmentation: assigning object class to every pixel
- **Speech & Language**
 - Speech Recognition: audio to text
 - Translation
 - Natural Language Processing: text to meaning
 - Audio Generation: text to audio
- **Games**

Emerging Applications

- **Medical** (Cancer Detection, Pre-Natal)
- **Finance** (Trading, Energy Forecasting, Risk)
- **Infrastructure** (Structure Safety and Traffic)
- Weather Forecasting and Event Detection

<http://www.nextplatform.com/2016/09/14/next-wave-deep-learning-applications/>

Deep Learning for Self-driving Cars



Opportunities

\$500B Market over 10 Years!

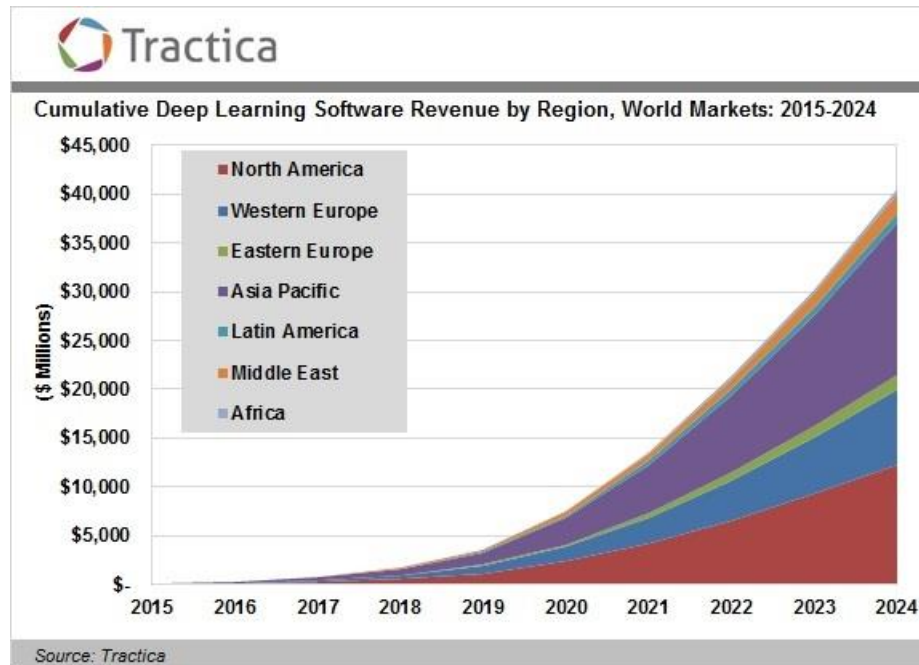
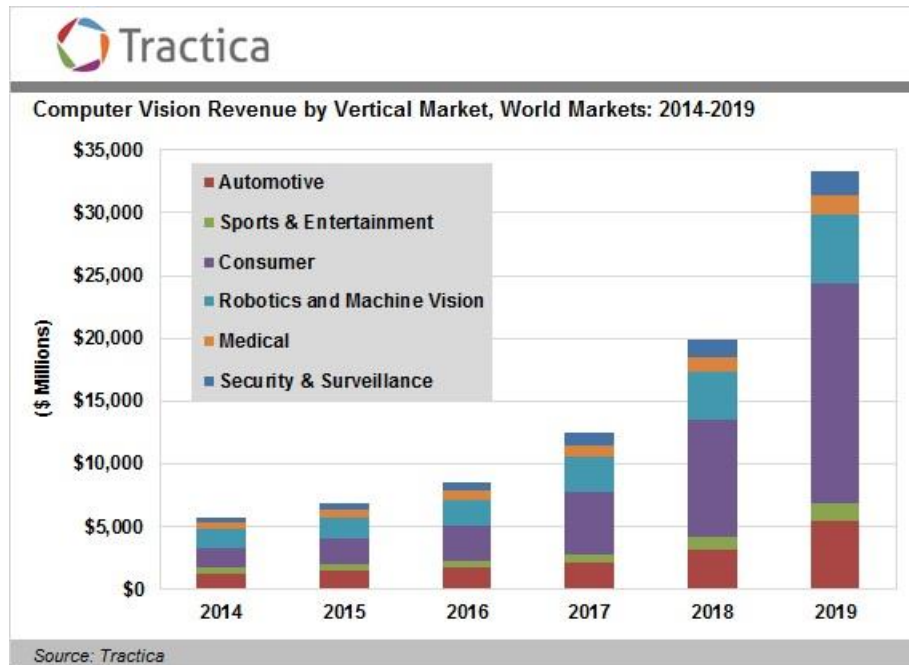


Image Source: Tractica

Opportunities

From EE Times – September 27, 2016

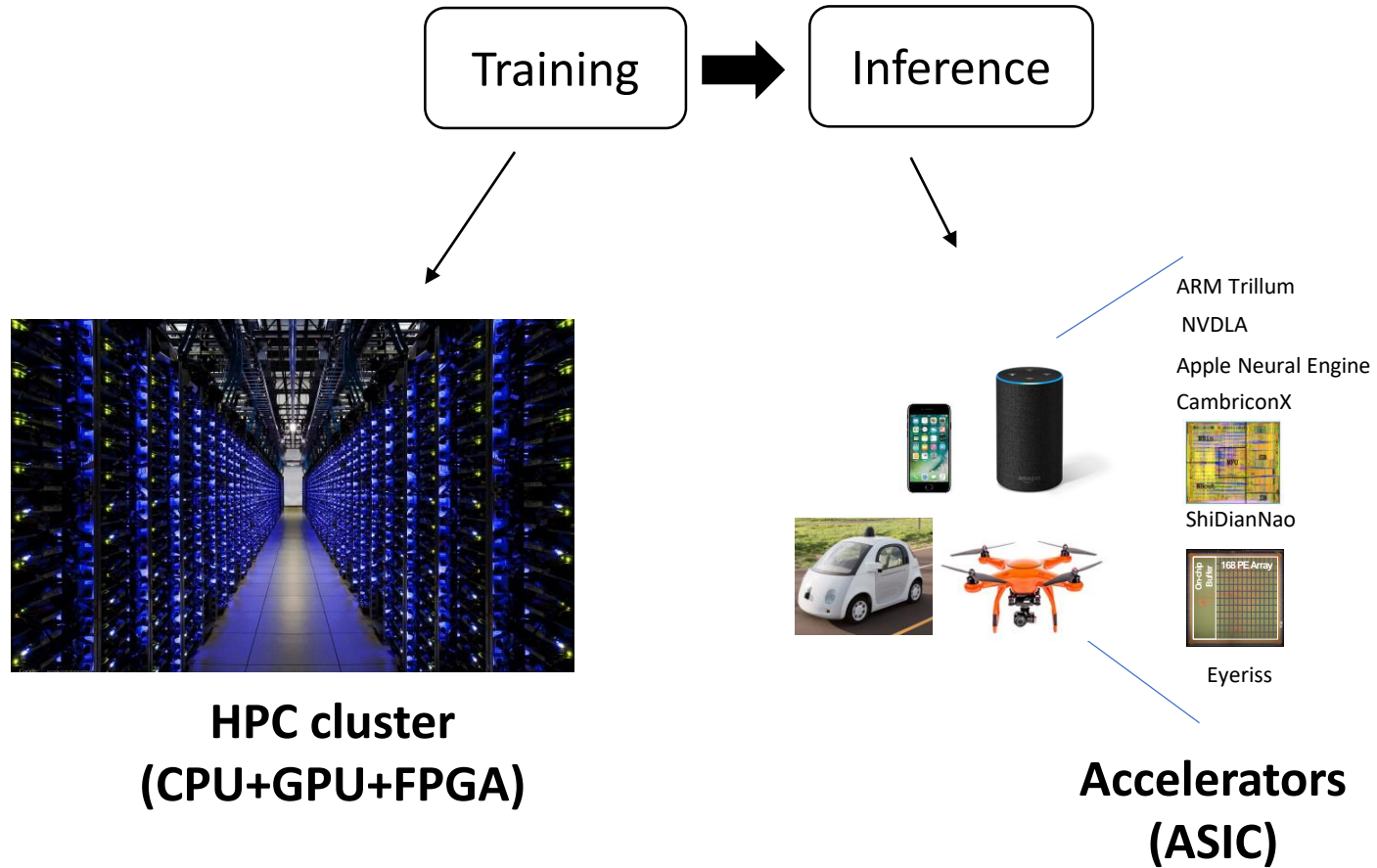
“Today the job of training machine learning models is limited by compute, if we had faster processors we’d run bigger models...in practice we train on a reasonable subset of data that can finish in a matter of months. We could use improvements of several orders of magnitude – 100x or greater.”

- Greg Diamos, Senior Researcher, (formerly) SVAIL, Baidu

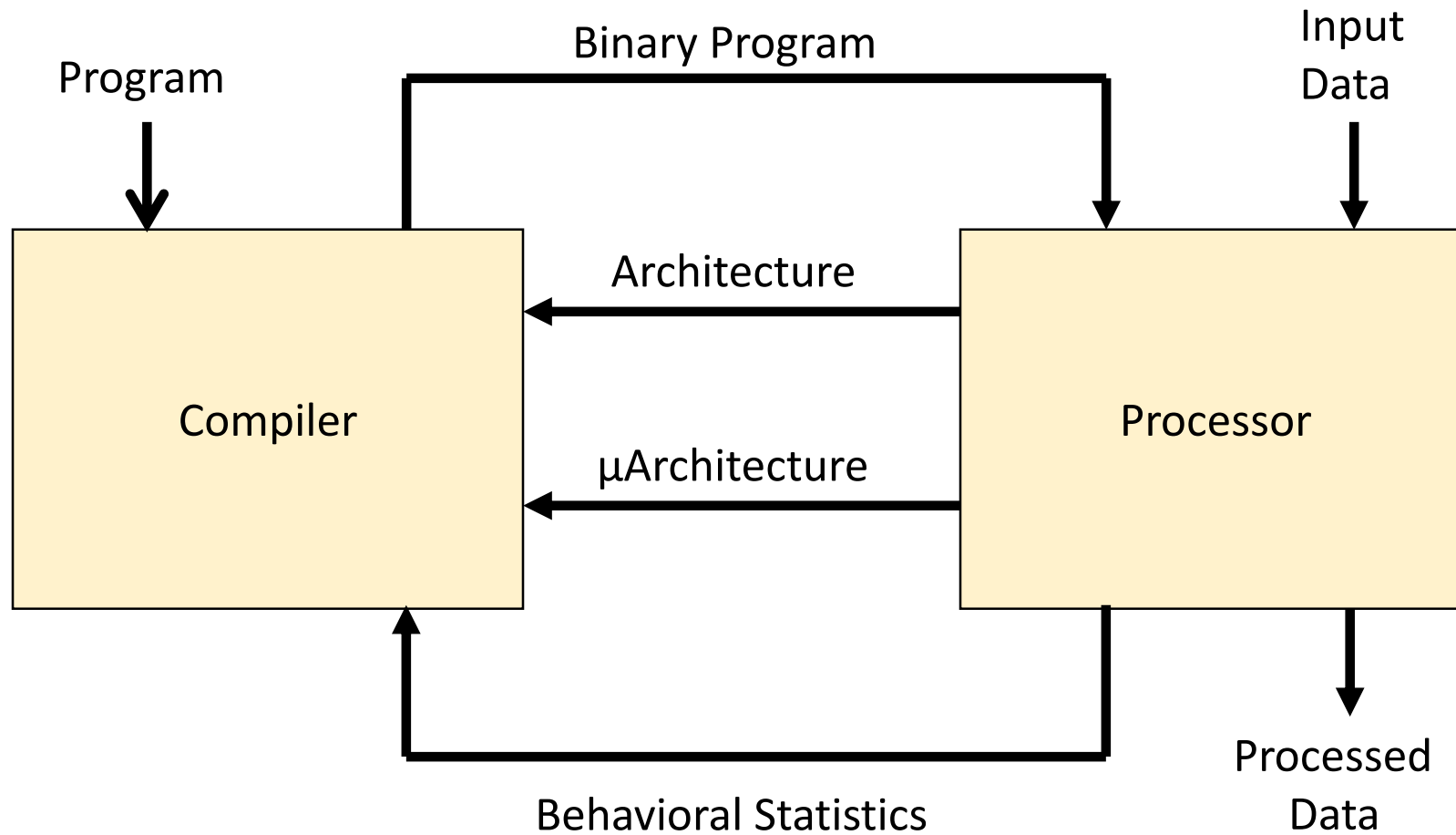
COMPUTATION Platforms

- CPU
 - Intel, ARM, AMD...
- GPU
 - NVIDIA, AMD...
- Fine Grained Reconfigurable (FPGA)
 - Microsoft BrainWave
- Coarse Grained Programmable/Reconfigurable
 - Wave Computing, Plasticine, Graphcore...
- Application Specific
 - NeufLOW, *DianNao, Eyeriss, Cnvlutin, SCNN, TPU, ...

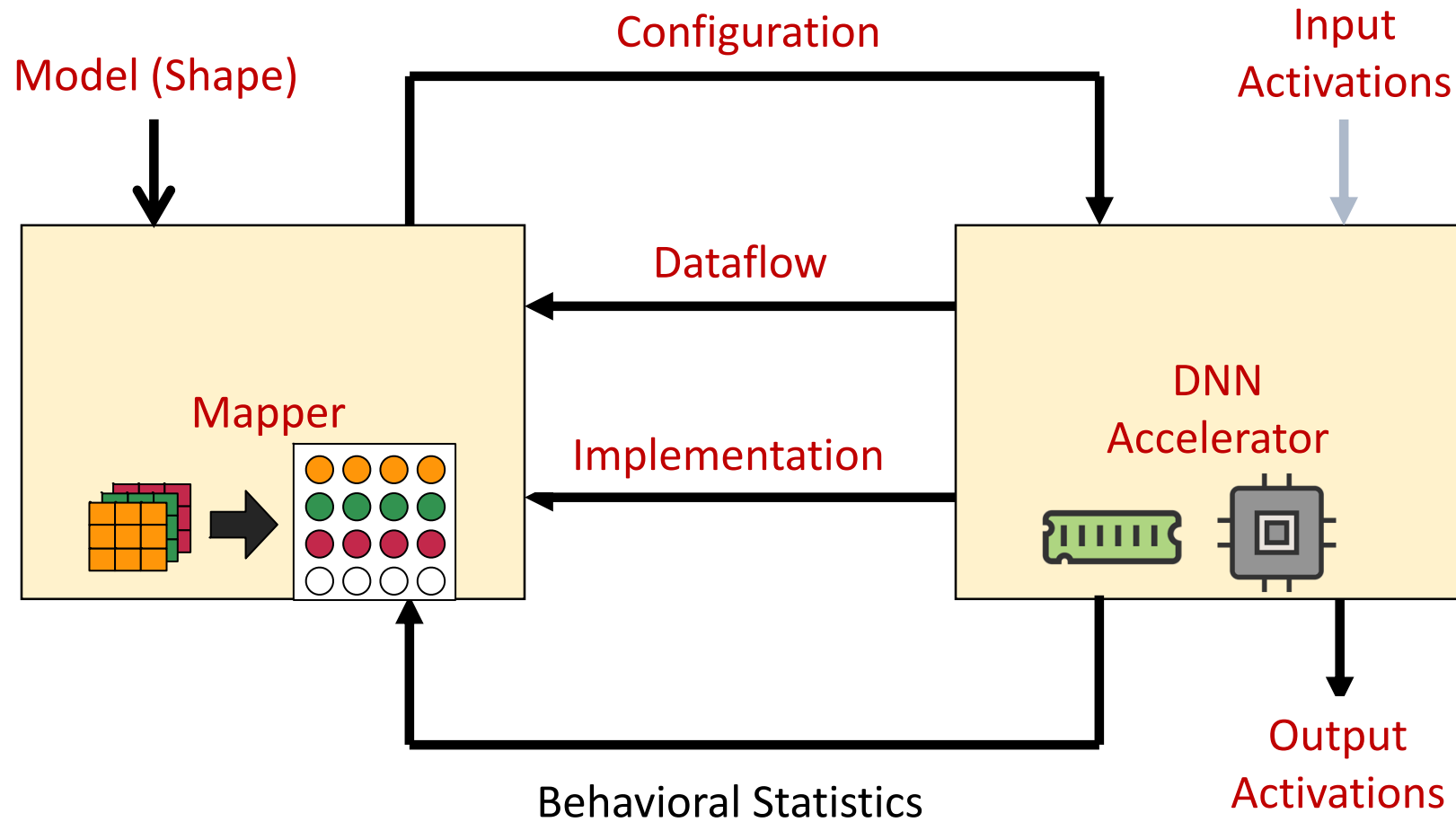
Computation Platforms



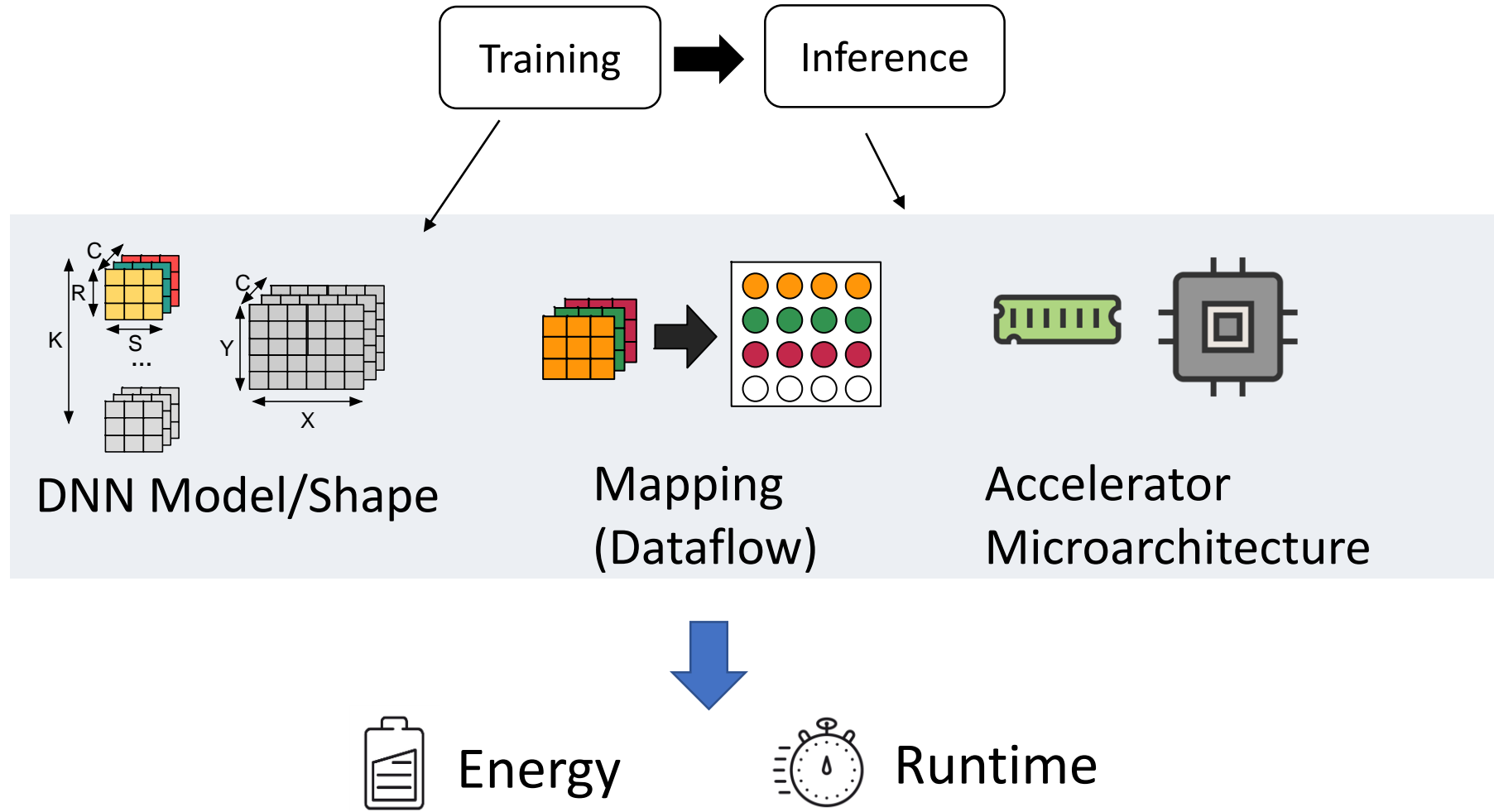
CPU Compute Model



DNN Compute Model



Challenges in Design and Deployment



DNN 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)
- 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.