# Introduction to Deep Learning

Greg Tsagkatakis

ICS - FORTH

# Machine Learning

Play Video

https://www.youtube.com/watch?v=f_uwKZIAeM0

# Agenda

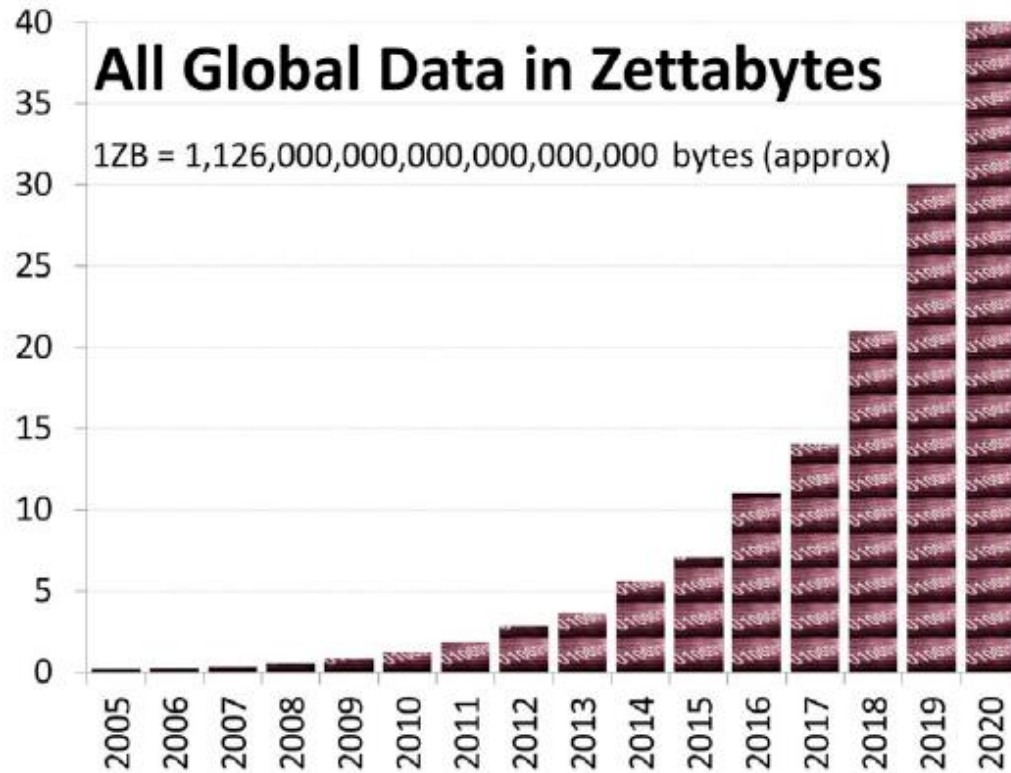Lecture #1
- ◦ Introduction
- ◦ Supervised Deep Learning

Lecture #2
- ◦ Unsupervised Deep Learning
- ◦ Deep Reinforcement learning

# Big Data

The 5Vs
➢ **Volume**



**All Global Data in Zettabytes**

1ZB = 1,126,000,000,000,000,000,000 bytes (approx)

The growth in data as seen by United Nations Economic Commission for Europe.

# Big Data

The 5Vs
- ➢Volume
- ➢**Velocity**



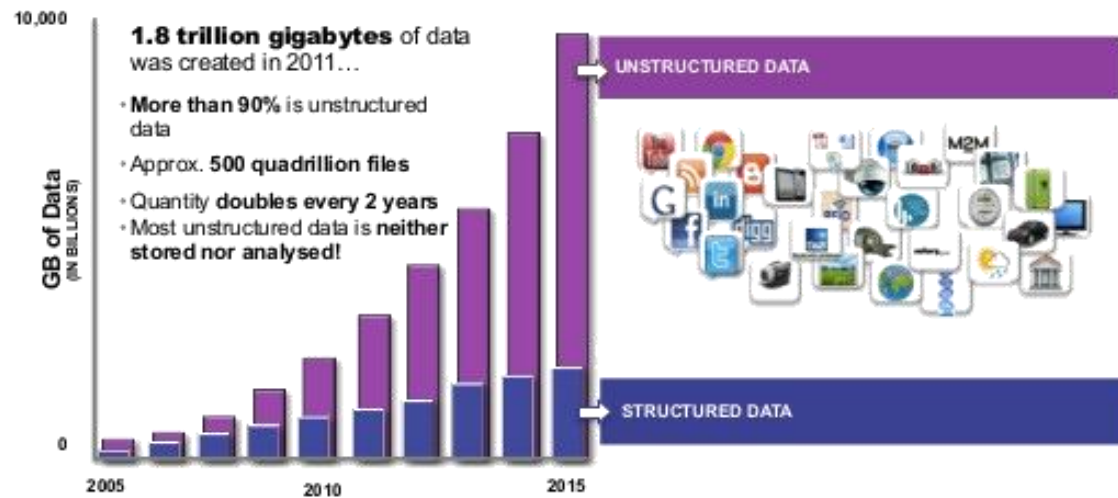2017 This Is What Happens In An Internet Minute

facebook
900,000 Logins

Google
3.5 Million Search Queries

NETFLIX
70,017 Hours Watched

$751,522 Spent Online

1.8 Million Snaps Created

15,000 GIFs Sent via Messenger

120 New Accounts Created
Linked in

50 Voice-First Devices Shipped
amazon echo

16 Million Text Messages

4.1 Million Videos Viewed
You Tube

342,000 Apps Downloaded
Google play / App Store

46,200 Posts Uploaded Instagram

452,000 Tweets Sent

990,000 Swipes
tinder

156 Million Emails Sent

40,000 Hours Listened
Spotify

Created By:
@LoriLewis
@OfficiallyChadd

60 SECONDS

# Big Data

The 5Vs
- ➢ Volume
- ➢ Velocity
- ➢ **Variety**

## The Explosion of Unstructured Data

CISCO

10,000

**1.8 trillion gigabytes** of data was created in 2011…

- More than **90%** is unstructured data
- Approx. **500 quadrillion files**
- Quantity **doubles every 2 years**
- Most unstructured data is **neither stored nor analysed!**

GB of Data (IN BILLIONS)

UNSTRUCTURED DATA

M2M

STRUCTURED DATA

0

2005          2010          2015

Source: Cloudera

# Big Data

The 5Vs
- ➤Volume
- ➤Velocity
- ➤Variety
- ➤**Veracity**

## Who's Winning? Daily track of Clinton and Trump's support

Updated daily.

More from the poll, and why it differs from others.

— Hillary Clinton    — Donald Trump

☐ Area of uncertainty

RNC  DNC    1st debate    3rd deb. **48.0%**

50

45

43.2%

40

2nd deb.

Jul. 10    Nov. 07

Note: Shaded gray area indicates the race is too close to call.

Sources: USC Dornsife/LA Times Presidential Election Daybreak Poll

# Big Data

The 5Vs
- ➤Volume
- ➤Velocity
- ➤Variety
- ➤Veracity
- ➤**Value**



Big Data Market Forecast ($US BILLIONS)

$5.1 — 2012
$10.2 — 2013
$16.8 — 2014
$32.1 — 2015
$48.0 — 2016
$53.4 — 2017

# BD & Healthcare

# Big Data & Genetics



Growth of DNA Sequencing

# Big Data & Astrophysics

Astronomy & Astrophysics

| Sky Survey Project | Volume | Velocity | Variety |
|---|---|---|---|
| Sloan Digital Sky Survey (SDSS) | 50 TB | 200 GB per day | Images, redshifts |
| Large Synoptic Survey Telescope (LSST ) | ~ 200 PB | 10 TB per day | Images, catalogs |
| Square Kilometer Array (SKA ) | ~ 4.6 EB | 150 TB per day | Images, redshifts |

Astrophysics and Big Data: Challenges, Methods, and Tools. Mauro Garofalo, Alessio Botta, and Giorgio Ventre.

# Handling Big Data

**Machine Learning + Big Data -> Data science**



Source: https://s3.amazonaws.com/aws.drewconway.com/viz/venn_diagram/data_science.html

# Big Data &the Brain

Human Visual System



© Stephen E. Palmer, 2002

# How does the Brain do it?
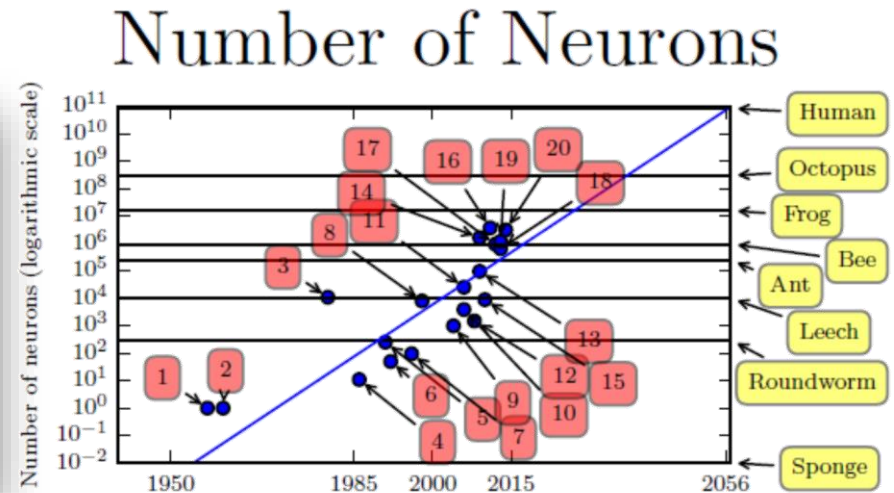
$10^{11}$ neurons

$10^{14}$-$10^{15}$ synapses
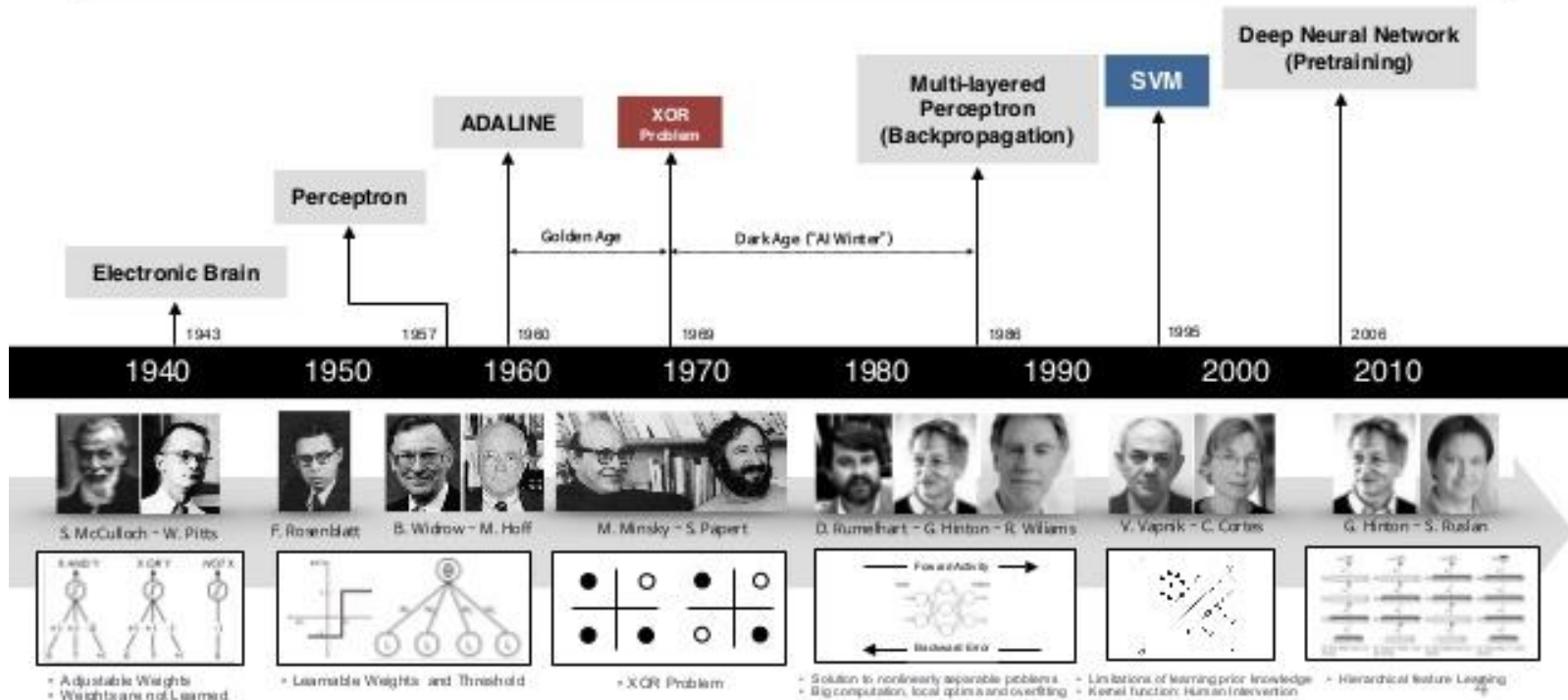


Figure 1.11

(Goodfellow 2016)

# Artificial Neural Networks



| Biological NN | Artificial NN |
|---|---|
| soma | unit |
| axon, dendrite | connection |
| synapse | weight |
| potential | weighted sum |
| threshold | bias weight |
| signal | activation |

# Brief history of DL

# Why Today?

Lots of Data

# Why Today?

Lots of Data

Deeper Learning

# Why Today?

Lots of Data

Deep Learning

More Power





https://blogs.nvidia.com/blog/2016/01/12/accelerating-ai-artificial-intelligence-gpus/
https://www.slothparadise.com/what-is-cloud-computing/

# Apps: Gaming



History of Game AI
By: Andrey Kurenkov

# Apps: Self-driving cars

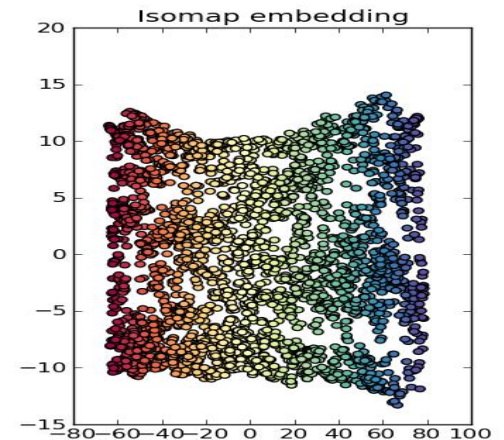https://www.youtube.com/watch?v=VG68SKoG7vE
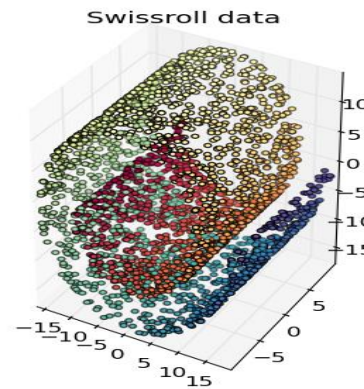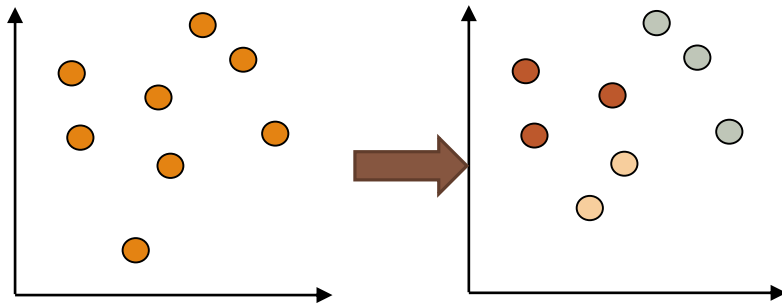
# Intro to ML

# Types of Machine Learning

**Supervised learning:** present example inputs and their desired outputs (**labels**) → learn a general rule that maps inputs to outputs.
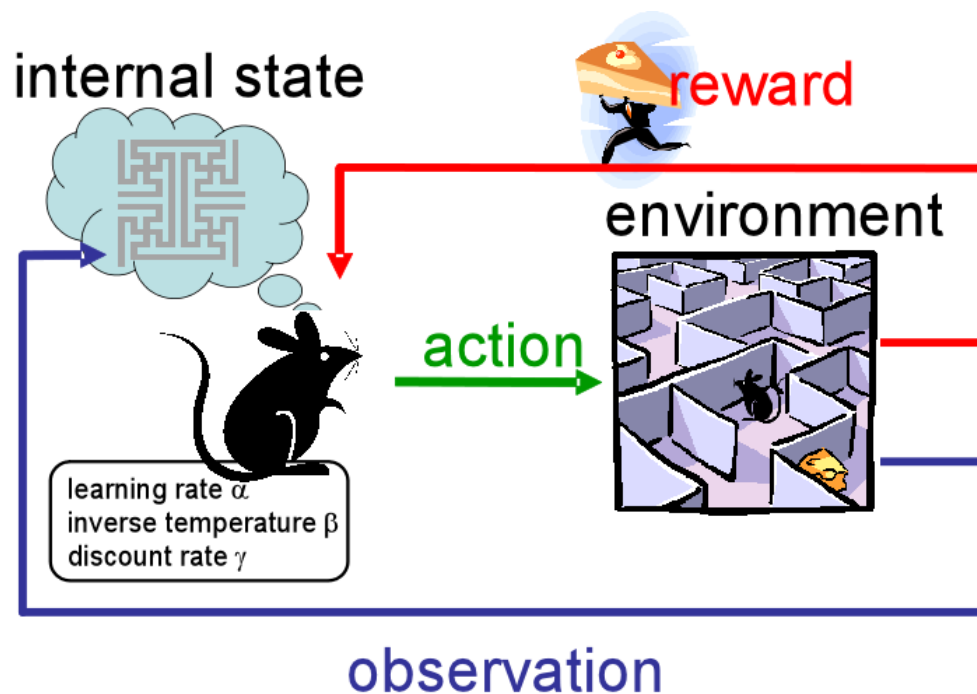
# Types of Machine Learning

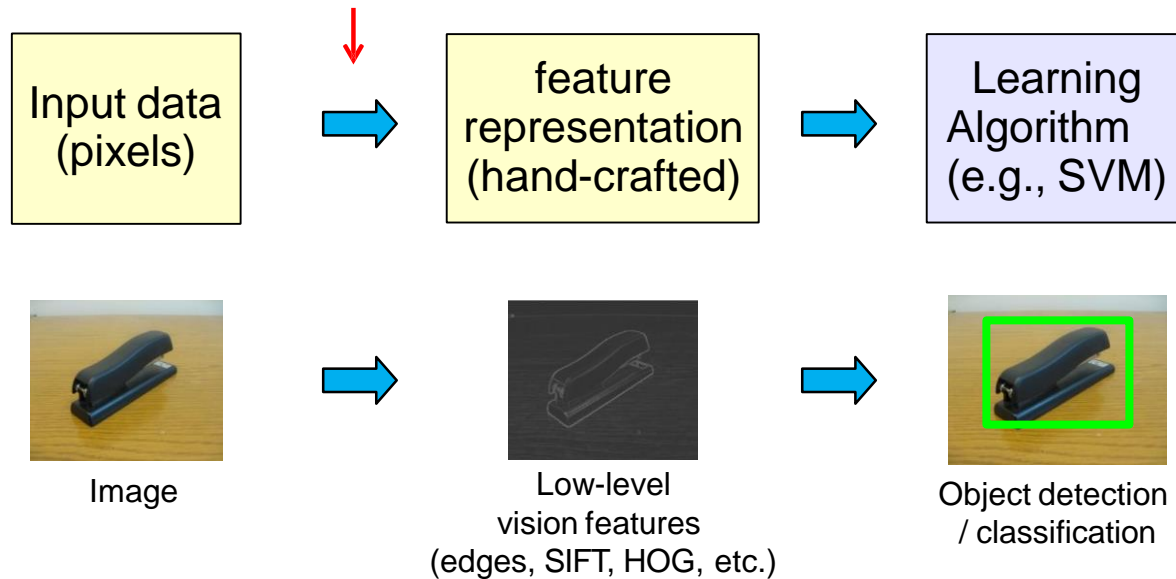**Unsupervised learning:** no labels are given → find structure in input.

# Types of Machine Learning

**Reinforcement learning:** system interacts with environment and must perform a certain goal without explicitly telling it whether it has come close to its goal or not.

# Feature extraction in ML

Features are not learned

| Input data (pixels) | → | feature representation (hand-crafted) | → | Learning Algorithm (e.g., SVM) |



Image



Low-level vision features (edges, SIFT, HOG, etc.)



Object detection / classification

# Feature learning

Pixel 1



Learning Algorithm

❌ No Car
⭐ Car

Pixel 2

Pixel 1

Pixel 2

# Feature learning

Pixel 1

Feature Representation

Learning Algorithm
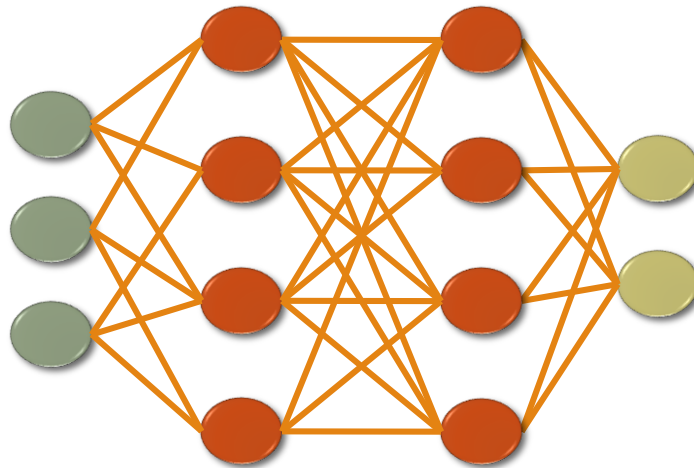
Pixel 2

❌ No Car
⭐ Car

Pixel 1

Pixel 2
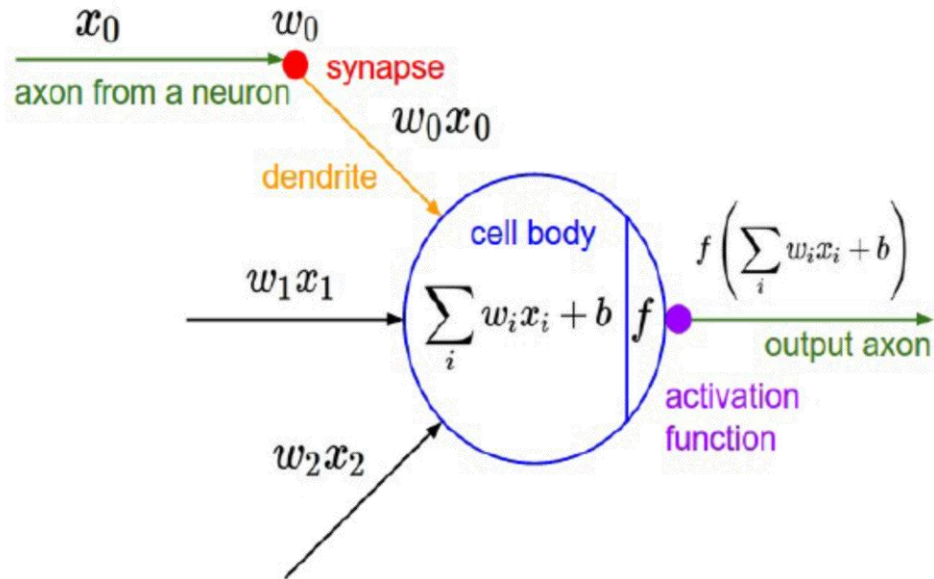
Feature 1

Feature 2

# Fundamentals of ANN

# Key components of ANN

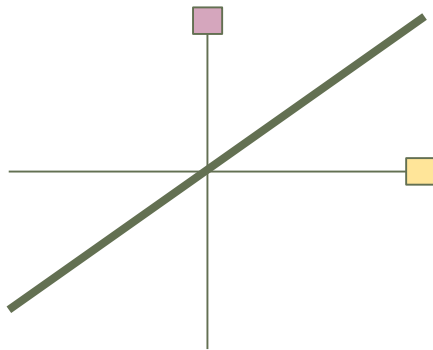➤ Architecture (input/hidden/output layers)

# Key components of ANN

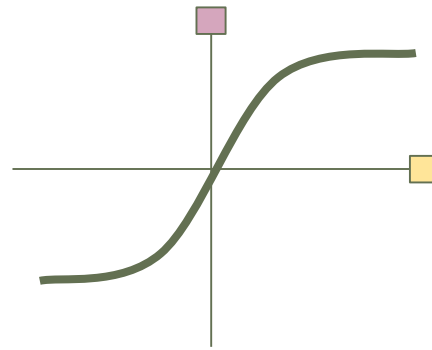➢ Architecture (input/hidden/output layers)

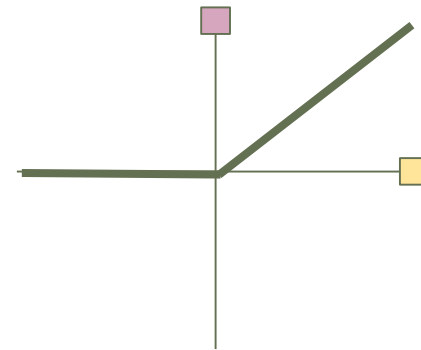➢ Weights

# Key components of ANN

➤ Architecture (input/hidden/output layers)

➤ Weights

➤ Activations

**LINEAR**

**LOGISTIC /
SIGMOIDAL / TANH**

**RECTIFIED
LINEAR (ReLU)**

# Perceptron: an early attempt

*Activation* function

$$\hat{f}(x) = \sigma(w \cdot x + b) \qquad \sigma(y) = \begin{cases} 1, & y > 0 \\ 0, & o/w \end{cases}$$

Need to tune $w$ and $b$

# Multilayer perceptron



$w_{1A}$

$w_1$

$w_{2B}$

$w_2$

$w_3$

$w_{1D}$

A

B

C

D

$w_{AE}$

$w_{DE}$

E

We just added a neuron layer!

We just introduced non-linearity!

A neuron is of the form **σ(w.x + b)** where **σ** is an *activation* function

A mostly complete chart of

# Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org
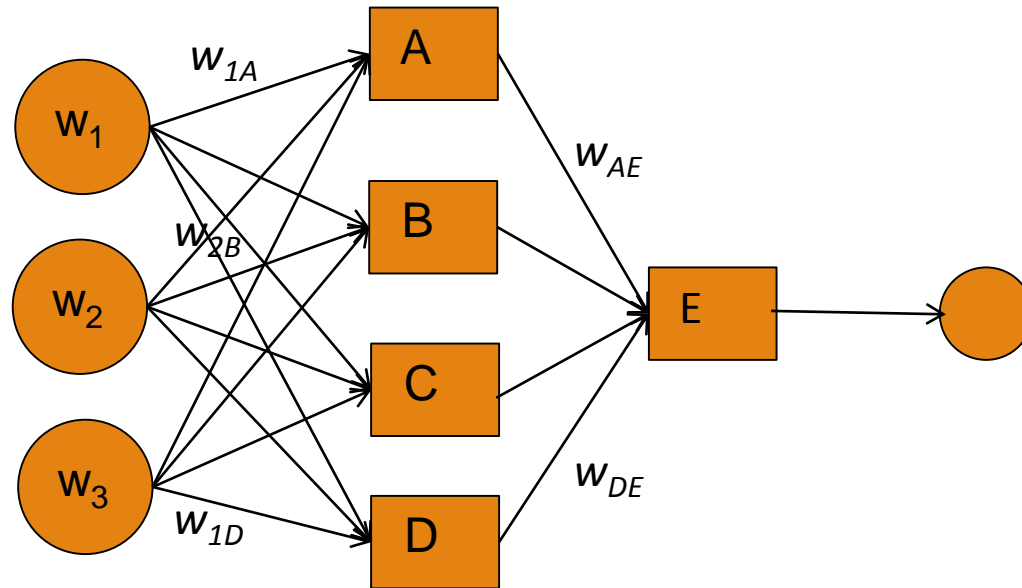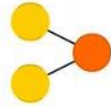
Backfed Input Cell
Input Cell
Noisy Input Cell
Hidden Cell
Probablistic Hidden Cell
Spiking Hidden Cell
Output Cell
Match Input Output Cell
Recurrent Cell
Memory Cell
Different Memory Cell
Kernel
Convolution or Pool
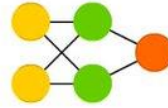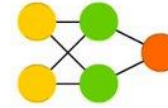
Deep Feed Forward (DFF)
Perceptron (P)
Feed Forward (FF)
Radial Basis Network (RBF)
Recurrent Neural Network (RNN)
Long / Short Term Memory (LSTM)
Gated Recurrent Unit (GRU)
Auto Encoder (AE)
Variational AE (VAE)
Denoising AE (DAE)
Sparse AE (SAE)
Markov Chain (MC)
Hopfield Network (HN)
Boltzmann Machine (BM)
Restricted BM (RBM)
Deep Belief Network (DBN)

Sasen Cain (@spectralradius)

35

# Training & Testing

Training: determine weights
- ◦ Supervised: labeled training examples
- ◦ Unsupervised: no labels available
- ◦ Reinforcement: examples associated with rewards

Testing (Inference): apply weights to new examples

# Training DNN

1. Get batch of data

2. Forward through the network -> estimate loss

3. Backpropagate error

4. Update weights based on gradient

# BackPropagation

Chain Rule in Gradient Descent: Invented in 1969 by Bryson and Ho

**Defining a loss/cost function**

Assume a function $J(x, y; \theta) = \dfrac{1}{2} \sum (y - f(x; \theta))^2$

$$f(x; \theta) = w^T x + b \quad , \quad \theta = \{w, b\}$$

Types of Loss function

- Hinge $\quad J(x, y) = max\{0, 1 - xy\}$

- Exponential $\quad J(x, y) = exp(-xy)$

- Logistic $\quad J(x, y) = log_2(1 + exp(-xy))$

# Gradient Descent

➢Minimize function J w.r.t. parameters θ

New weights $\longrightarrow$ $\theta^* = \theta - n * \nabla J(y, x; \theta)$ $\longleftarrow$ Gradient

Old weights

Learning rate

- Gradient

$$\nabla J(x) = (\frac{\partial J(x)}{\partial x_1}, \frac{\partial J(x)}{\partial x_2}, ..., \frac{\partial J(x)}{\partial x_n})$$

- Chain rule



$J$

$x$

# Visualization



Legend:
- SGD
- Momentum
- NAG
- Adagrad
- Adadelta
- Rmsprop

# Training Characteristics



Error

Testing Error

Training Error

Over-fitting

Under-fitting

**Training steps**

# References

Stephens, Zachary D., et al. "Big data: astronomical or genomical?." *PLoS biology* 13.7 (2015): e1002195.

LeCun, Yann, Yoshua Bengio, and Geoffrey Hinton. "Deep learning." *Nature* 521.7553 (2015): 436-444.

Goodfellow, Ian, Yoshua Bengio, and Aaron Courville. *Deep learning*. MIT press, 2016.

Kietzmann, Tim Christian, Patrick McClure, and Nikolaus Kriegeskorte. "Deep Neural Networks In Computational Neuroscience." *bioRxiv* (2017): 133504.

# Introduction to Deep Learning

Greg Tsagkatakis

ICS - FORTH

# Fundamentals of ANN

# Key components of ANN

➢ Architecture (input/hidden/output layers)

# Key components of ANN

➢ Architecture (input/hidden/output layers)

➢ Weights

# Key components of ANN

➢ Architecture (input/hidden/output layers)

➢ Weights

➢ Activations



**LINEAR**

**LOGISTIC / SIGMOIDAL / TANH**

**RECTIFIED LINEAR (ReLU)**

# Perceptron: an early attempt

*Activation* function

$$\hat{f}(x) = \sigma(w \cdot x + b) \quad \sigma(y) = \begin{cases} 1, & y > 0 \\ 0, & o/w \end{cases}$$

Need to tune $w$ and $b$

# Multilayer perceptron



A neuron is of the form $\sigma(\mathbf{w}.\mathbf{x} + \mathbf{b})$ where $\boldsymbol{\sigma}$ is an *activation* function

We just added a neuron layer!

We just introduced non-linearity!

A mostly complete chart of

# Neural Networks

©2016 Fjodor van Veen - asimovinstitute.org

Sasen Cain (@spectralradius)

# Training & Testing

Training: determine weights
- ◦ Supervised: labeled training examples
- ◦ Unsupervised: no labels available
- ◦ Reinforcement: examples associated with rewards

Testing (Inference): apply weights to new examples

# Training DNN

1. Get batch of data

2. Forward through the network -> estimate loss

3. Backpropagate error

4. Update weights based on gradient



$u_j$

$w_{ij}$    $w'_{jk}$

$u'_k$

Input layer    Hidden layer    Output layer

$x_i$ ⟶ $y_j$ ⟶ $z_k$ ⟵ $o_k$    Target

Errors

# BackPropagation

Chain Rule in Gradient Descent: Invented in 1969 by Bryson and Ho

**Defining a loss/cost function**

Assume a function $J(x, y; \theta) = \frac{1}{2} \sum (y - f(x; \theta))^2$

$$f(x; \theta) = w^T x + b \quad, \quad \theta = \{w, b\}$$

Types of Loss function

- Hinge $\quad J(x, y) = max\{0, 1 - xy\}$

- Exponential $\quad J(x, y) = exp(-xy)$

- Logistic $\quad J(x, y) = log_2(1 + exp(-xy))$

# Gradient Descent

➢Minimize function J w.r.t. parameters θ

New weights ⟶ $\theta^* = \theta - n * \nabla J(y, x; \theta)$ ⟵ Gradient

Old weights

Learning rate

- Gradient

$$\nabla J(x) = (\frac{\partial J(x)}{\partial x_1}, \frac{\partial J(x)}{\partial x_2}, ..., \frac{\partial J(x)}{\partial x_n})$$

- Chain rule

Loss function

Tangent line

$J$

$-\alpha \Delta_w L(X, \mathbf{w}_t, b)$ ⟵ $\Delta_w L(X, \mathbf{w}_t, b)$

$\mathbf{w}_t$

$\mathbf{w}_{t+1}$

$x$

# BackProp

**Given:** $y = g(u)$ and $u = h(x)$.

**Chain Rule:**

$$\frac{dy_i}{dx_k} = \sum_{j=1}^{J} \frac{dy_i}{du_j}\frac{du_j}{dx_k}, \quad \forall i, k$$

# BackProp

x ⟶ g ⟶ y=g(x) ⟶ f ⟶ z=f(y)=f(g(x))

**Chain rule:**

◦ **Single variable**

$$\frac{dz}{dx} = \frac{dz}{dy}\frac{dy}{dx}.$$

◦ **Multiple variables**

$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j}\frac{\partial y_j}{\partial x_i}.$$

Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4



$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$

$$\frac{\partial f}{\partial z}$$

Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



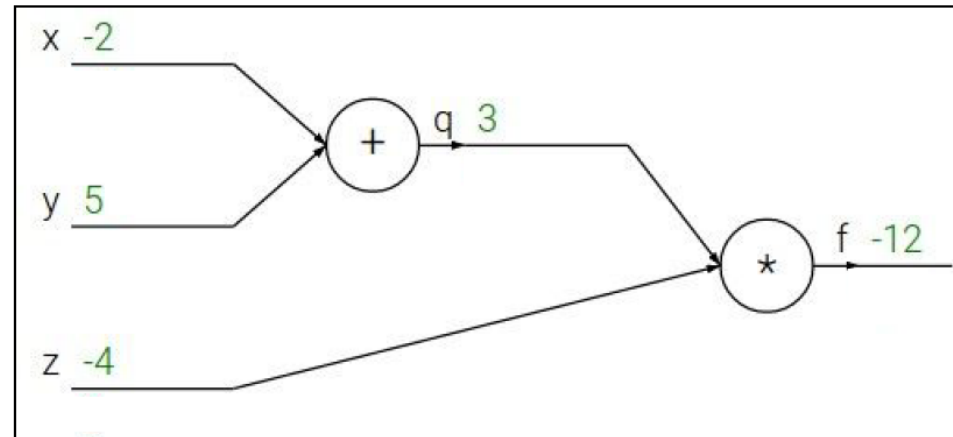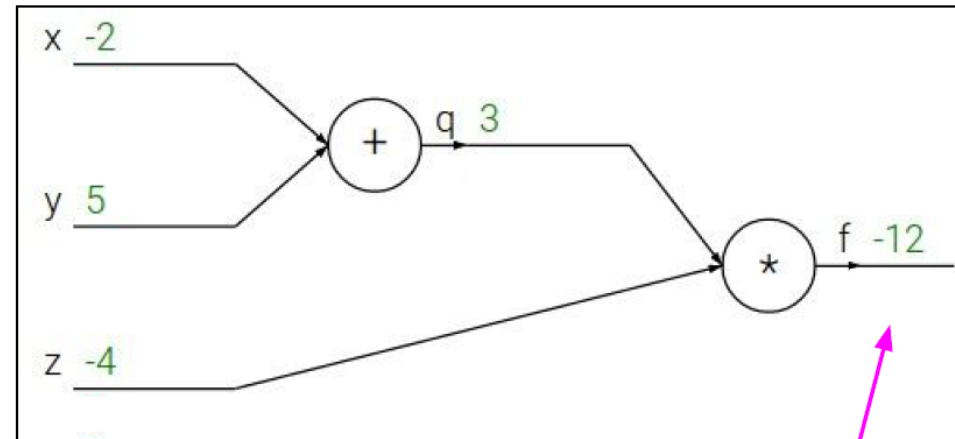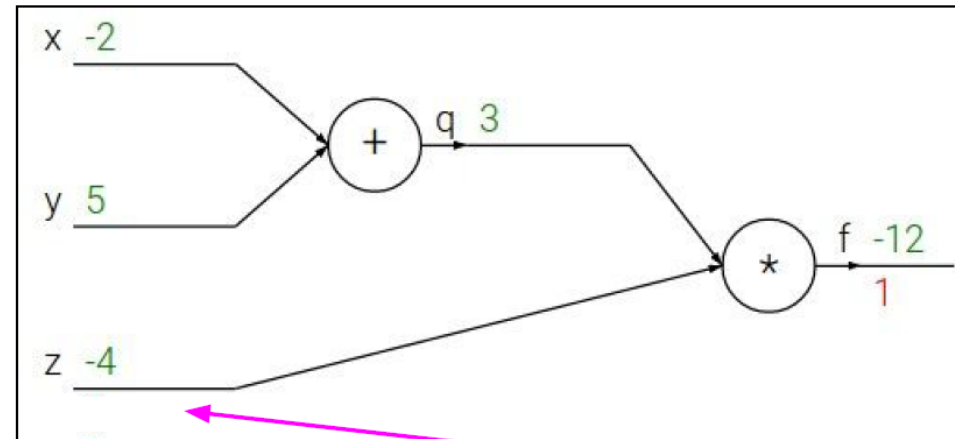$$\frac{\partial f}{\partial y}$$

# Backpropagation: a simple example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$
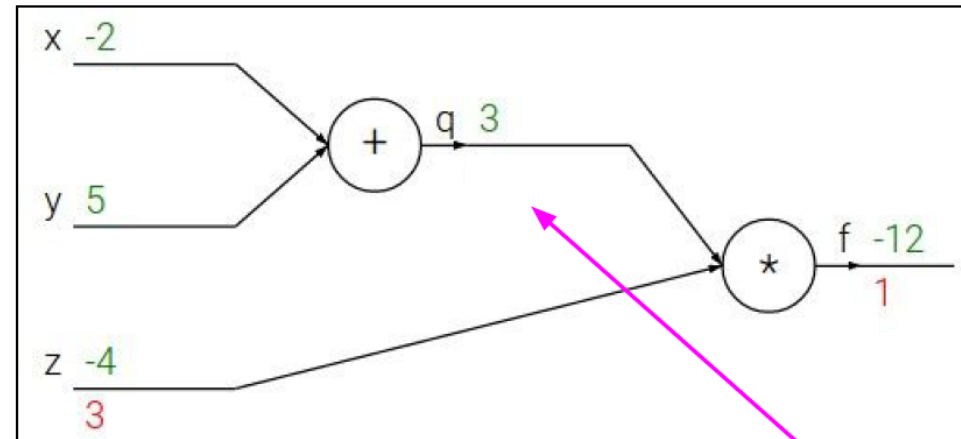


$$\frac{\partial f}{\partial y}$$

Chain rule:

$$\frac{\partial f}{\partial y} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial y}$$

# Visualization

# Training Characteristics

# Supervised Learning

# Supervised Learning

Data
Labels
Model
Prediction

**Exploiting prior knowledge**

➢ Expert users

➢ Crowdsourcing

➢ Other instruments



Spiral

Elliptical

?

# State-of-the-art (before Deep Learning)

Support Vector Machines
- Binary classification

# State-of-the-art (before Deep Learning)

Support Vector Machines
- Binary classification
- Kernels <-> non-linearities



Data projected to R^2 (hyperplane projection shown)

Data in R^3 (separable w/ hyperplane)

# State-of-the-art (before Deep Learning)

Support Vector Machines
- Binary classification
- Kernels <-> non-linearities

Random Forests
- Multi-class classification

# State-of-the-art (before Deep Learning)

Support Vector Machines
- Binary classification
- Kernels <-> non-linearities

Random Forests
- Multi-class classification

Markov Chains/Fields
- Temporal data

# State-of-the-art (since 2015)

Deep Learning (DL)

Convolutional Neural Networks (CNN) <-> Images

Recurrent Neural Networks (RNN) <-> Audio

# Convolutional Neural Networks



(Convolution  + Subsampling) + ()    …          + Fully Connected

# Convolutional Layers

32x32x1 Image

28x28xK activation map

width

height

channels

5x5x1 filter

K filters

width

height

$$(I * K)_{ij} = \sum_{m=0}^{k_1-1} \sum_{n=0}^{k_2-1} I(i-m, j-n)K(m,n)$$

$$= \sum_{m=0}^{k_1-1} \sum_{n=0}^{k_2-1} I(i+m, j+n)K(-m,-n)$$

# Convolutional Layers

Characteristics

➢ Hierarchical features

➢ Location invariance

Parameters

➢ Number of filters (32,64...)

➢ Filter size (3x3, 5x5)

➢ Stride (1)

➢ Padding (2,4)



"Machine Learning and AI for Brain Simulations" –
Andrew Ng Talk, UCLA, 2012

# Subsampling (pooling) Layers



<-> downsampling

➢ Scale invariance

Parameters

• Type

• Filter Size

• Stride

# Activation Layer

Introduction of non-linearity
- ◦ Brain: thresholding -> spike trains

Identity (Linear)

$$identity(x) = x$$

Sigmoid

$$sigmoid(x) = \frac{1}{1 + e^{-x}}$$

Tanh (Hypertangent)

$$tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

Gaussian

$$gaussian(x) = e^{-x^2/\sigma^2}$$

# Activation Layer

ReLU: x=max(0,x)

- ✓ Simplifies backprop
- ✓ Makes learning faster
- ✓ Avoids saturation issues
- ✓ ~ non-negativity constraint

(Note: The brain)

No saturated gradients

# Fully Connected Layers

Full connections to all activations in previous layer

Typically at the end

Can be replaced by conv

# LeNet [1998]



[LeCun et al., 1998]

# AlexNet [2012]



**Conv 1: Edge+Blob**        **Conv 3: Texture**        **Conv 5: Object Parts**        **Fc8: Object Classes**

Alex Krizhevsky, Ilya Sutskever and Geoff Hinton, ImageNet ILSVRC challenge in 2012
http://vision03.csail.mit.edu/cnn_art/data/single_layer.png

# VGGnet [2014]



maxpool

maxpool

maxpool

maxpool

maxpool

depth=64
3x3 conv
conv1_1
conv1_2

depth=128
3x3 conv
conv2_1
conv2_2

depth=256
3x3 conv
conv3_1
conv3_2
conv3_3
conv3_4

depth=512
3x3 conv
conv4_1
conv4_2
conv4_3
conv4_4

depth=512
3x3 conv
conv5_1
conv5_2
conv5_3
conv5_4

size=4096
FC1
FC2
size=1000
softmax

K. Simonyan, A. Zisserman Very Deep Convolutional Networks for Large-Scale Image Recognition,
arXiv technical report, 2014

# VGGnet

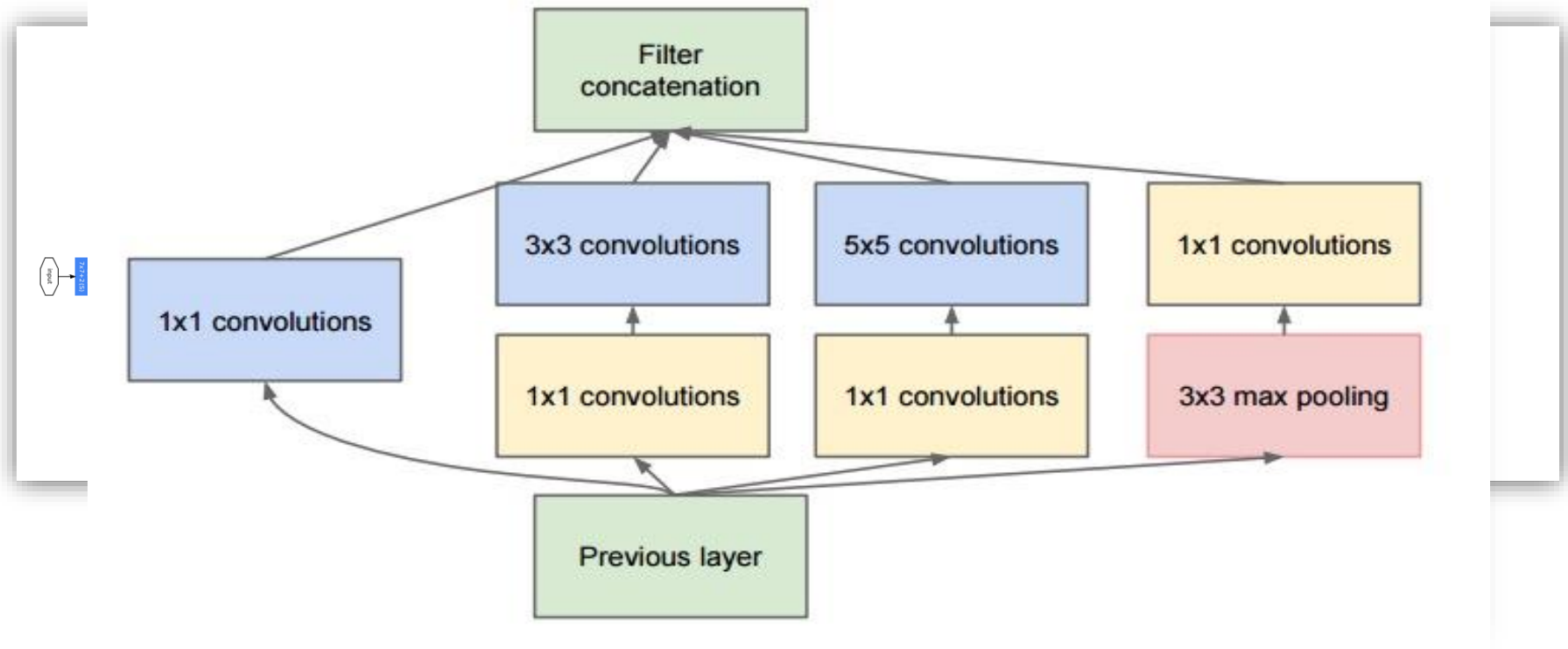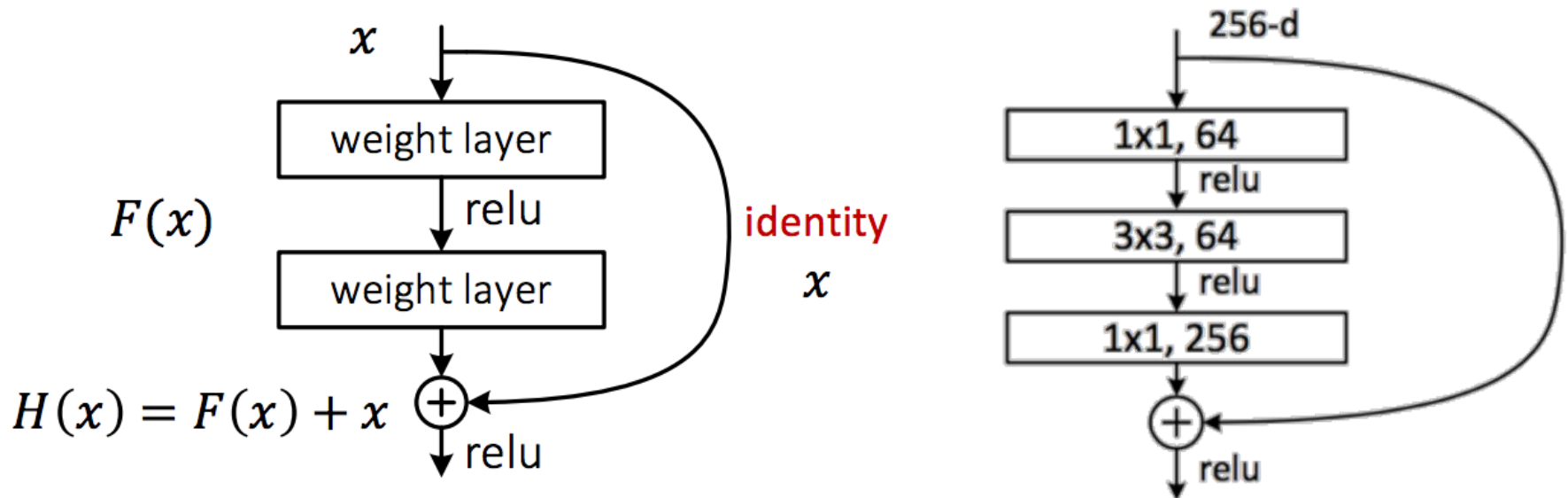| ConvNet Configuration | | | | | |
|---|---|---|---|---|---|
| A | A-LRN | B | C | D | E |
| 11 weight layers | 11 weight layers | 13 weight layers | 16 weight layers | 16 weight layers | 19 weight layers |
| input (224 × 224 RGB image) | | | | | |
| conv3-64 | conv3-64 **LRN** | conv3-64 **conv3-64** | conv3-64 conv3-64 | conv3-64 conv3-64 | conv3-64 conv3-64 |
| maxpool | | | | | |
| conv3-128 | conv3-128 | conv3-128 **conv3-128** | conv3-128 conv3-128 | conv3-128 conv3-128 | conv3-128 conv3-128 |
| maxpool | | | | | |
| conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 | conv3-256 conv3-256 **conv1-256** | conv3-256 conv3-256 **conv3-256** | conv3-256 conv3-256 conv3-256 **conv3-256** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 | conv3-512 conv3-512 **conv1-512** | conv3-512 conv3-512 **conv3-512** | conv3-512 conv3-512 conv3-512 **conv3-512** |
| maxpool | | | | | |
| FC-4096 | | | | | |
| FC-4096 | | | | | |
| FC-1000 | | | | | |
| soft-max | | | | | |

D: VGG16
E: VGG19
All filters are 3x3

More layers
smaller filters

# Inception (GoogLeNet, 2014)
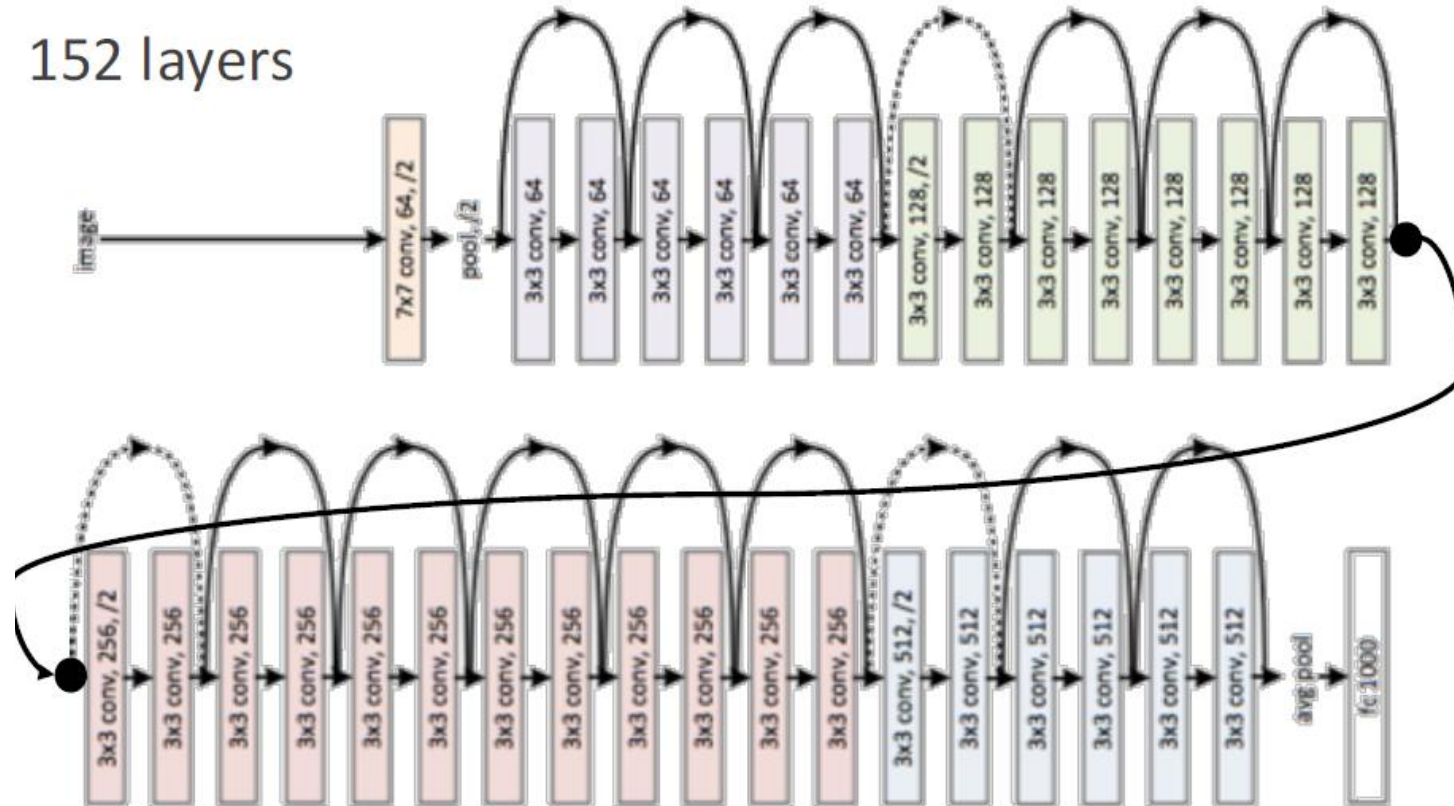


**Inception module** with dimensionality reduction

# Residuals

# ResNet, 2015



He, Kaiming, et al. "Deep residual learning for image recognition." *IEEE CVPR*. 2016.

# Training protocols

Fully Supervised
- Random initialization of weights
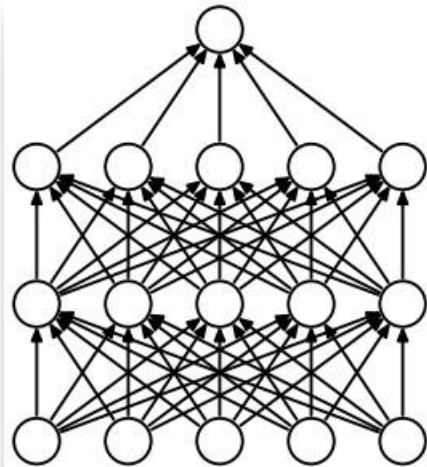- Train in supervised mode (example + label)

Unsupervised pre-training + standard classifier
- Train each layer unsupervised
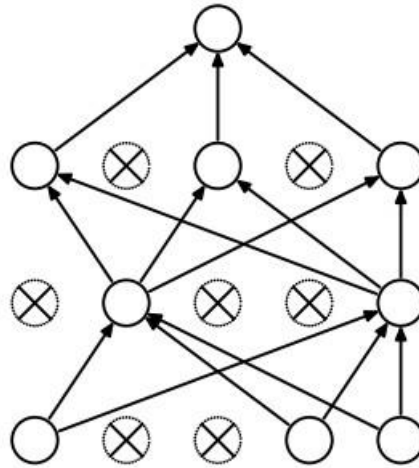- Train a supervised classifier (SVM) on top

Unsupervised pre-training + supervised fine-tuning
- Train each layer unsupervised
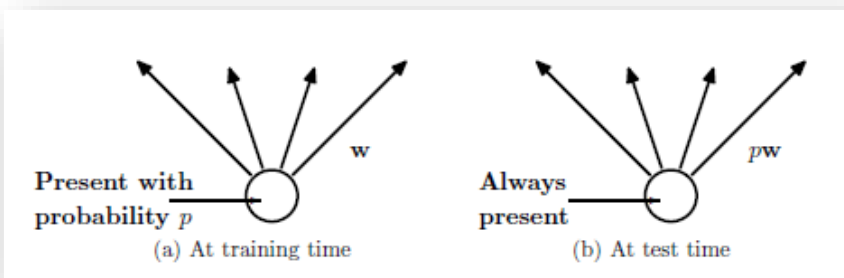- Add a supervised layer

# Dropout


(a) Standard Neural Net      (b) After applying dropout.





Present with probability $p$    w    Always present    $p\mathbf{w}$

(a) At training time      (b) At test time

Srivastava, Nitish, et al. "Dropout: a simple way to prevent neural networks from overfitting." *Journal of machine learning research* 15.1 (2014): 1929-1958.

# Batch Normalization



$$\textbf{Input:} \quad \text{Values of } x \text{ over a mini-batch: } \mathcal{B} = \{x_{1...m}\};$$
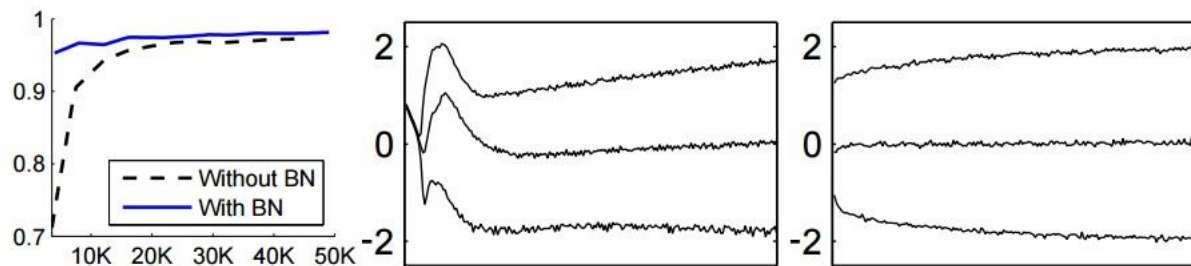$$\text{Parameters to be learned: } \gamma, \beta$$
$$\textbf{Output:} \quad \{y_i = BN_{\gamma,\beta}(x_i)\}$$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i \qquad\qquad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_{\mathcal{B}})^2 \qquad // \text{ mini-batch variance}$$

$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad\qquad // \text{ normalize}$$

$$y_i \leftarrow \gamma\widehat{x}_i + \beta \equiv BN_{\gamma,\beta}(x_i) \qquad // \text{ scale and shift}$$
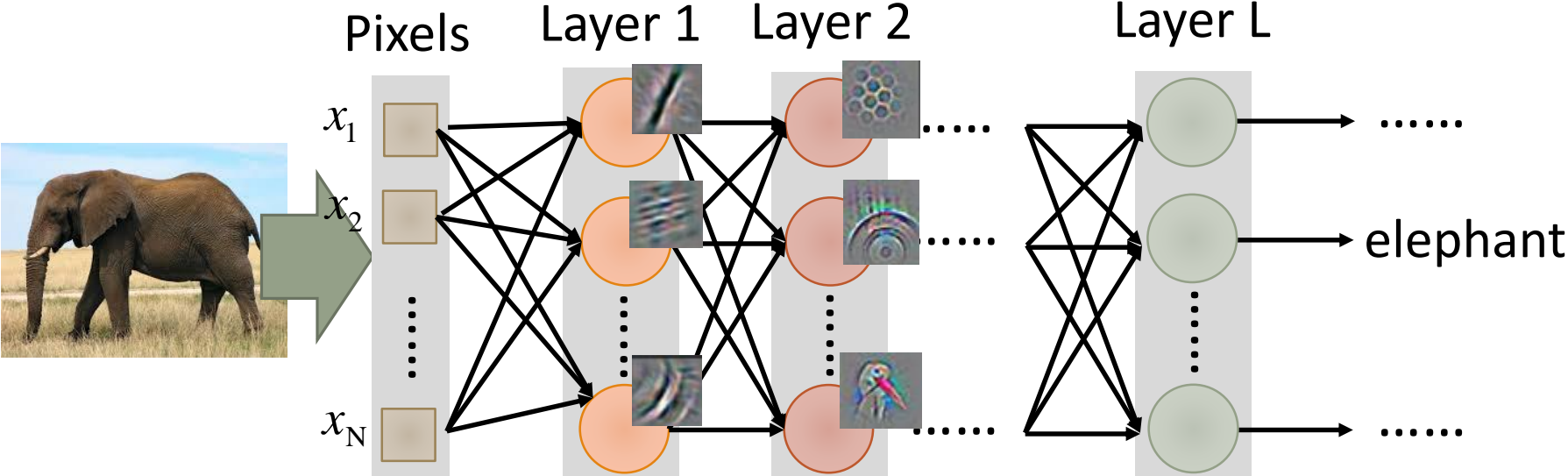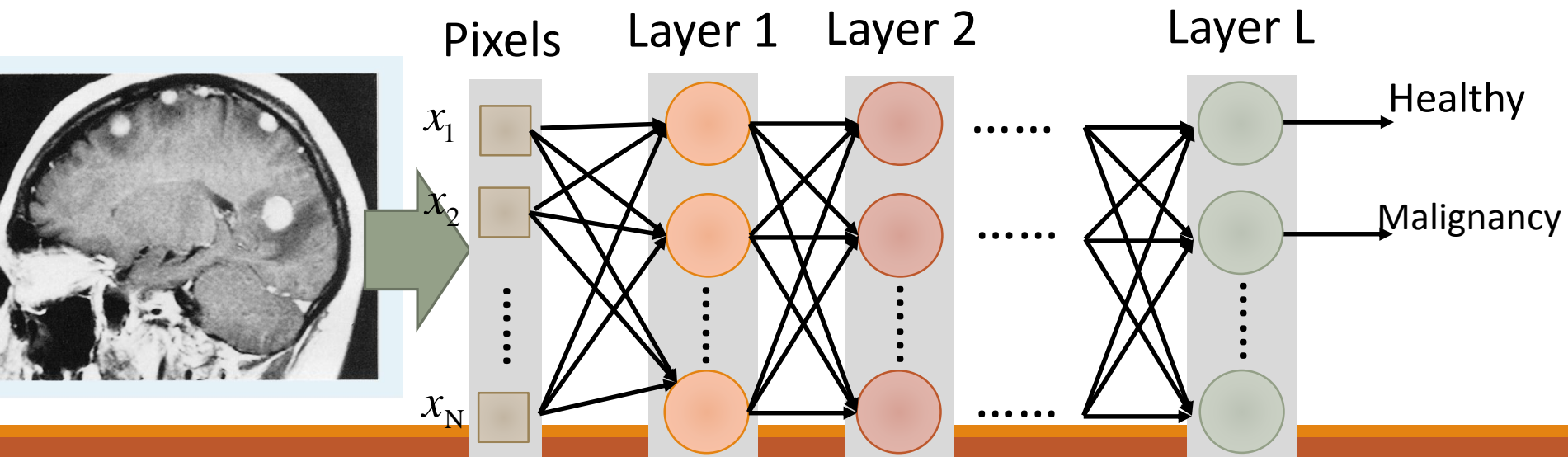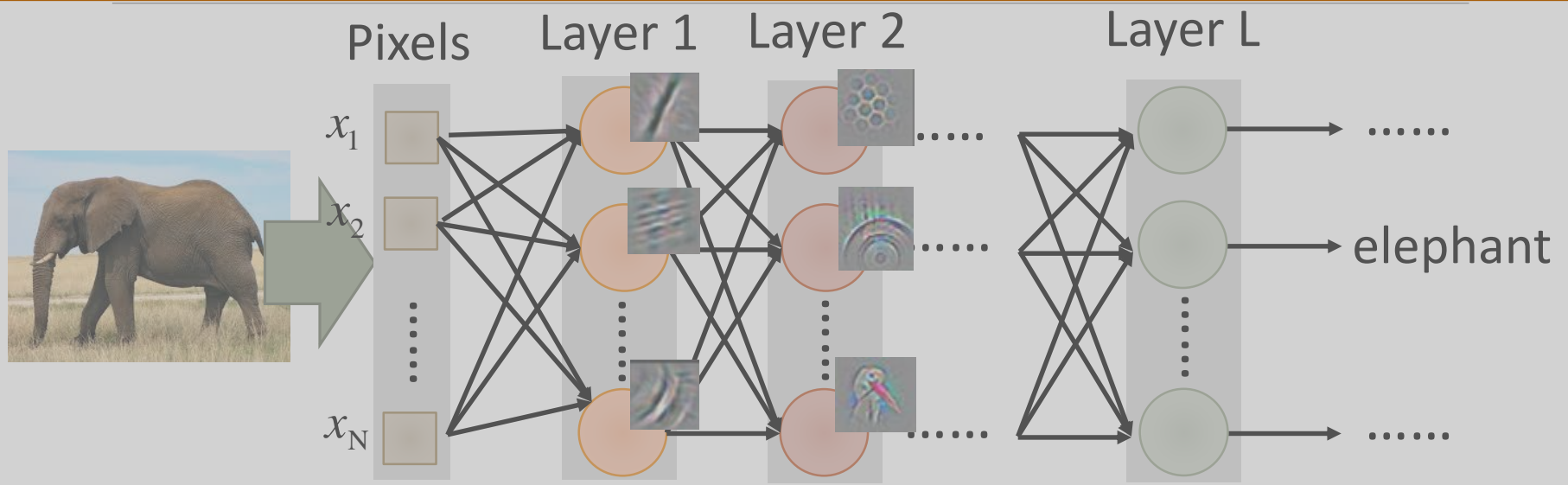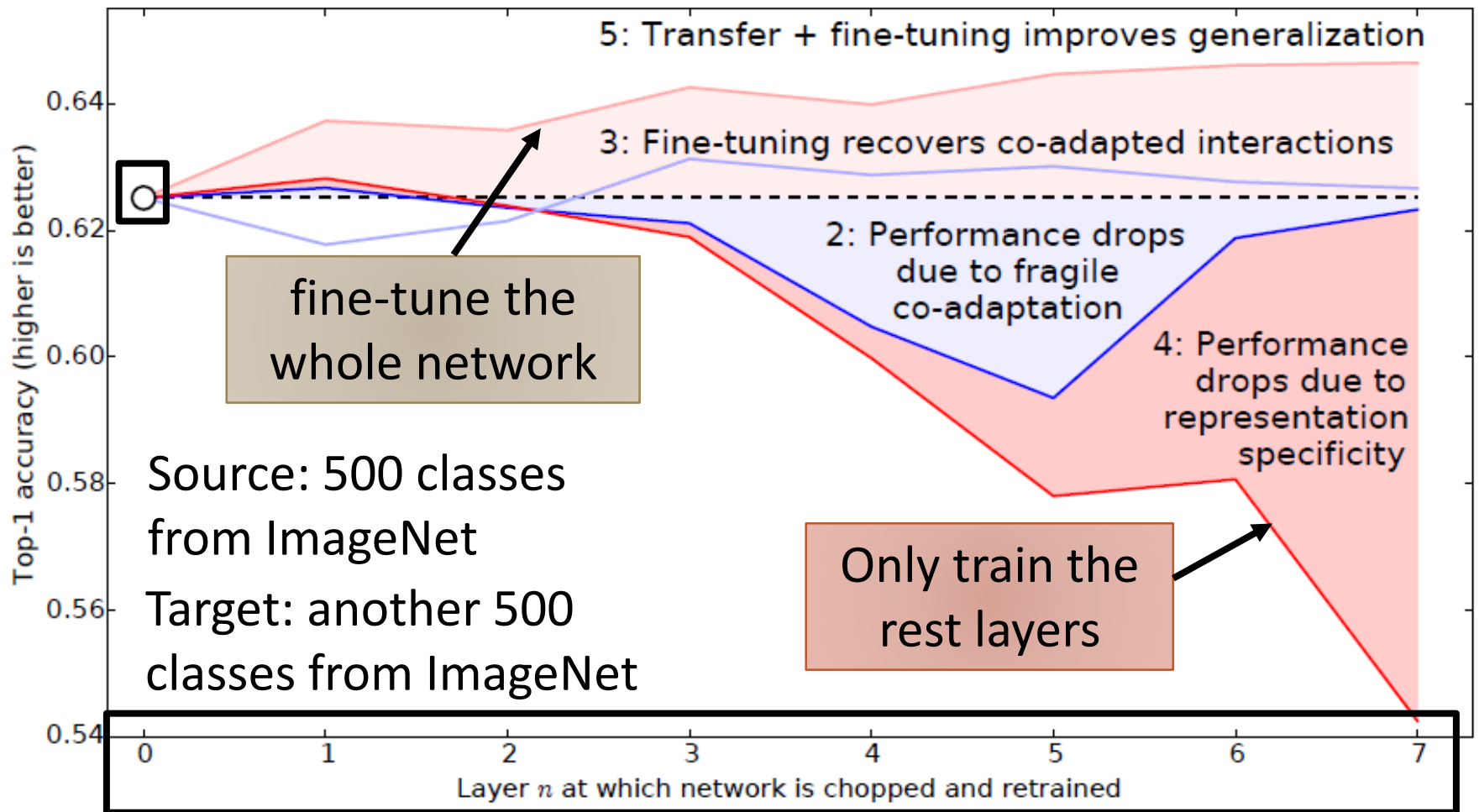
(a)    (b) Without BN    (c) With BN

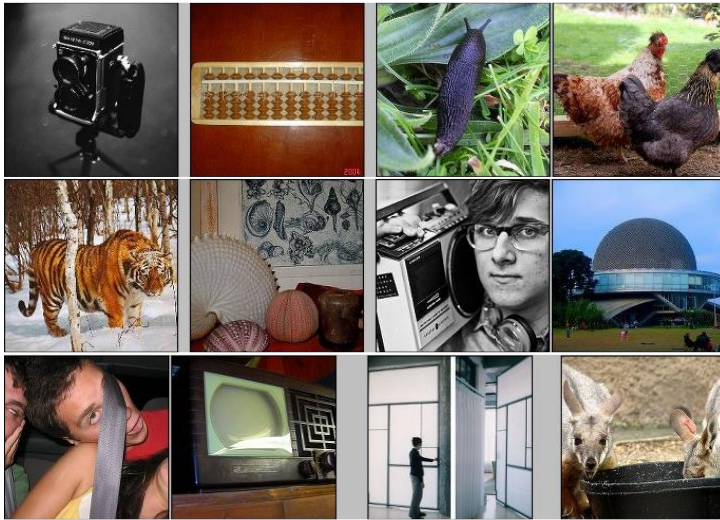# Transfer Learning

# Transfer Learning

# Layer Transfer - Image



Jason Yosinski, Jeff Clune, Yoshua Bengio, Hod Lipson, "How transferable are features in deep neural networks?", NIPS, 2014

# ImageNET



- ~14 million labeled images, 20k classes

- Images gathered from Internet

- Human labels via Amazon MTurk

- ImageNet Large-Scale Visual Recognition Challenge (ILSVRC):
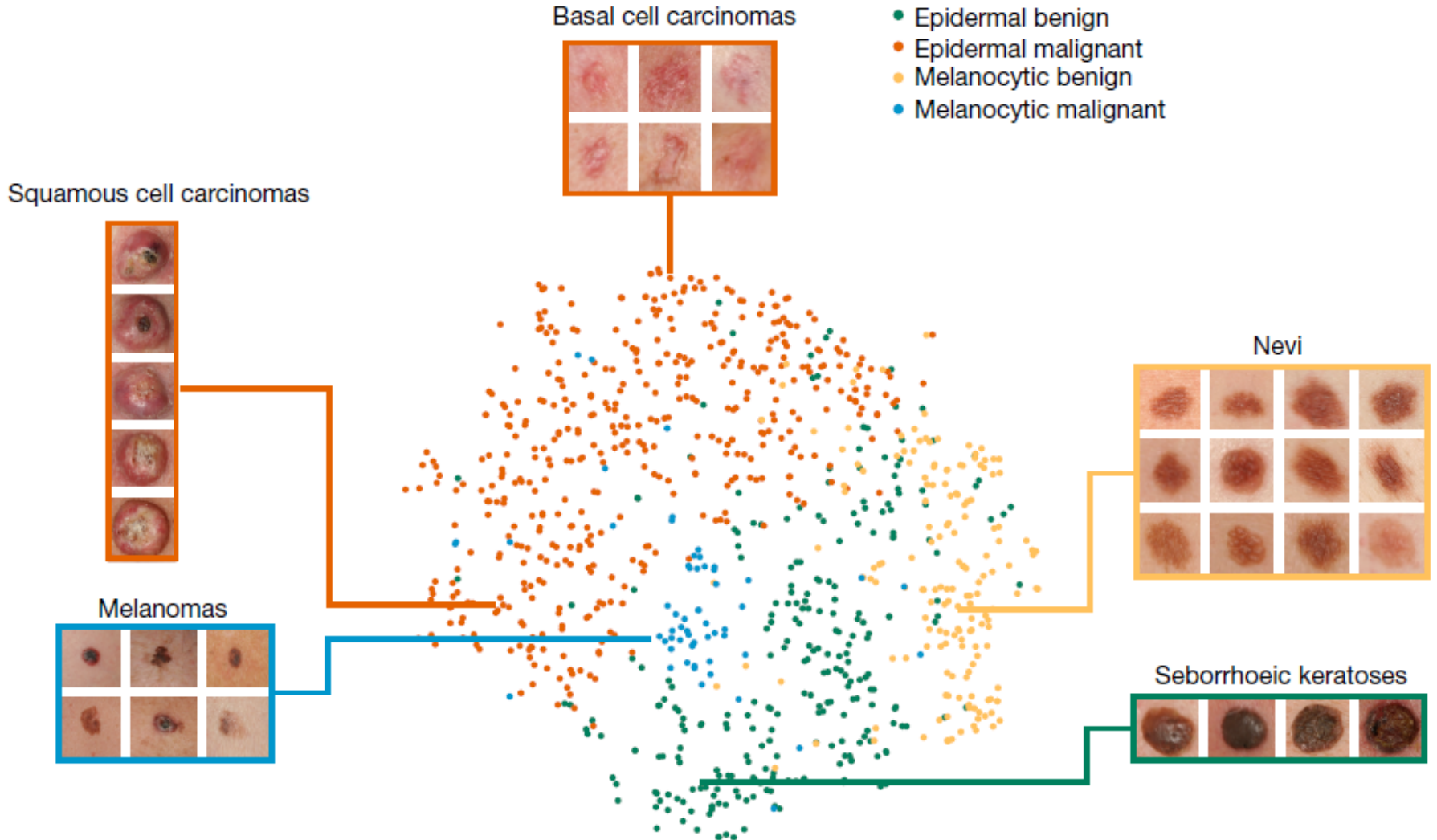1.2 million training images, 1000 classes

www.image-net.org/challenges/LSVRC/
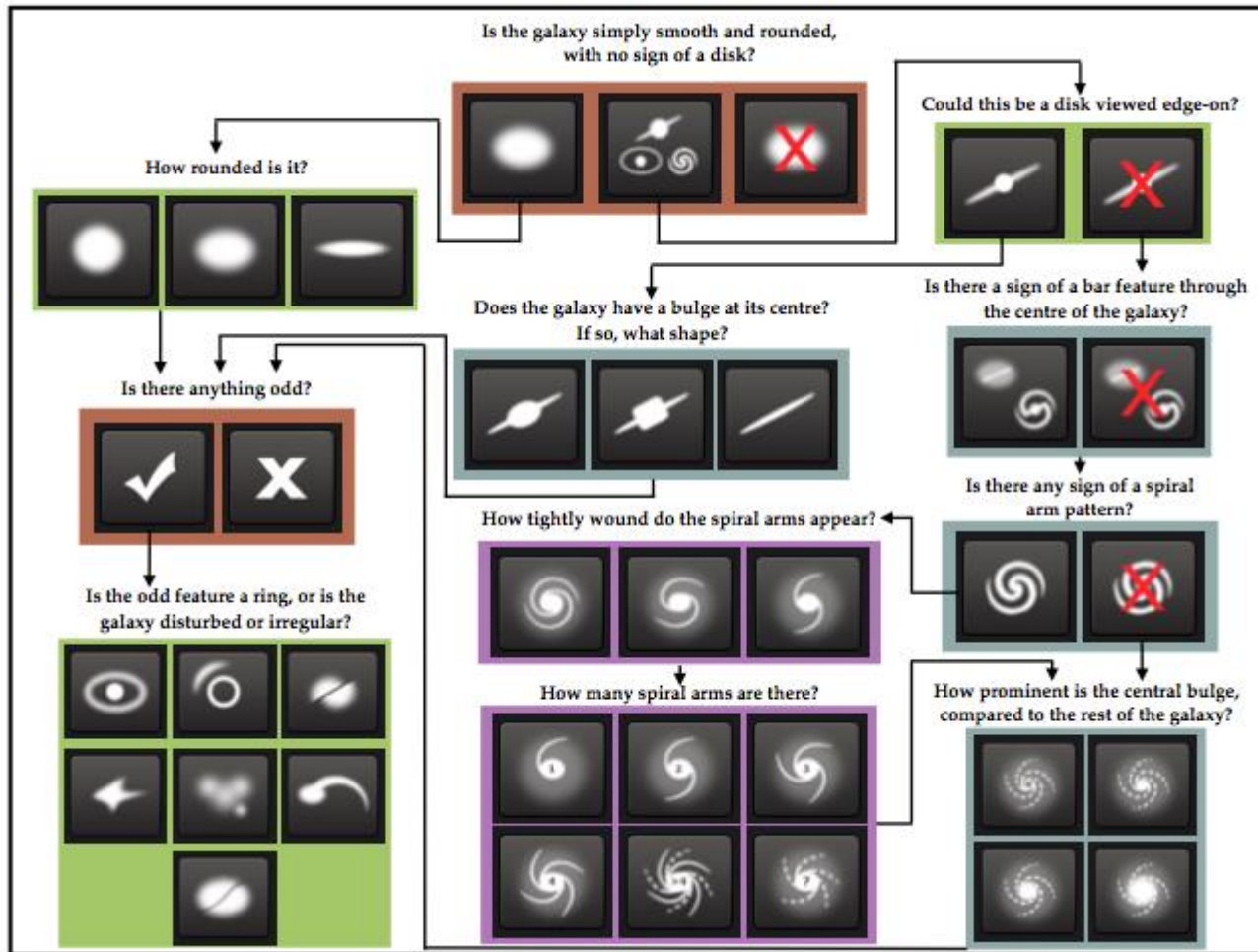
# Summary: ILSVRC 2012-2015

| Team | Year | Place | Error (top-5) | External data |
|------|------|-------|---------------|---------------|
| (AlexNet, 7 layers) | 2012 | - | 16.4% | no |
| SuperVision | 2012 | 1st | 15.3% | ImageNet 22k |
| Clarifai – NYU (7 layers) | 2013 | - | 11.7% | no |
| Clarifai | 2013 | 1st | 11.2% | ImageNet 22k |
| VGG – Oxford (16 layers) | 2014 | 2nd | 7.32% | no |
| GoogLeNet (19 layers) | 2014 | 1st | 6.67% | no |
| ResNet (152 layers) | 2015 | 1st | 3.57% | |
| Human expert* | | | 5.1% | |

http://karpathy.github.io/2014/09/02/what-i-learned-from-competing-against-a-convnet-on-imagenet/
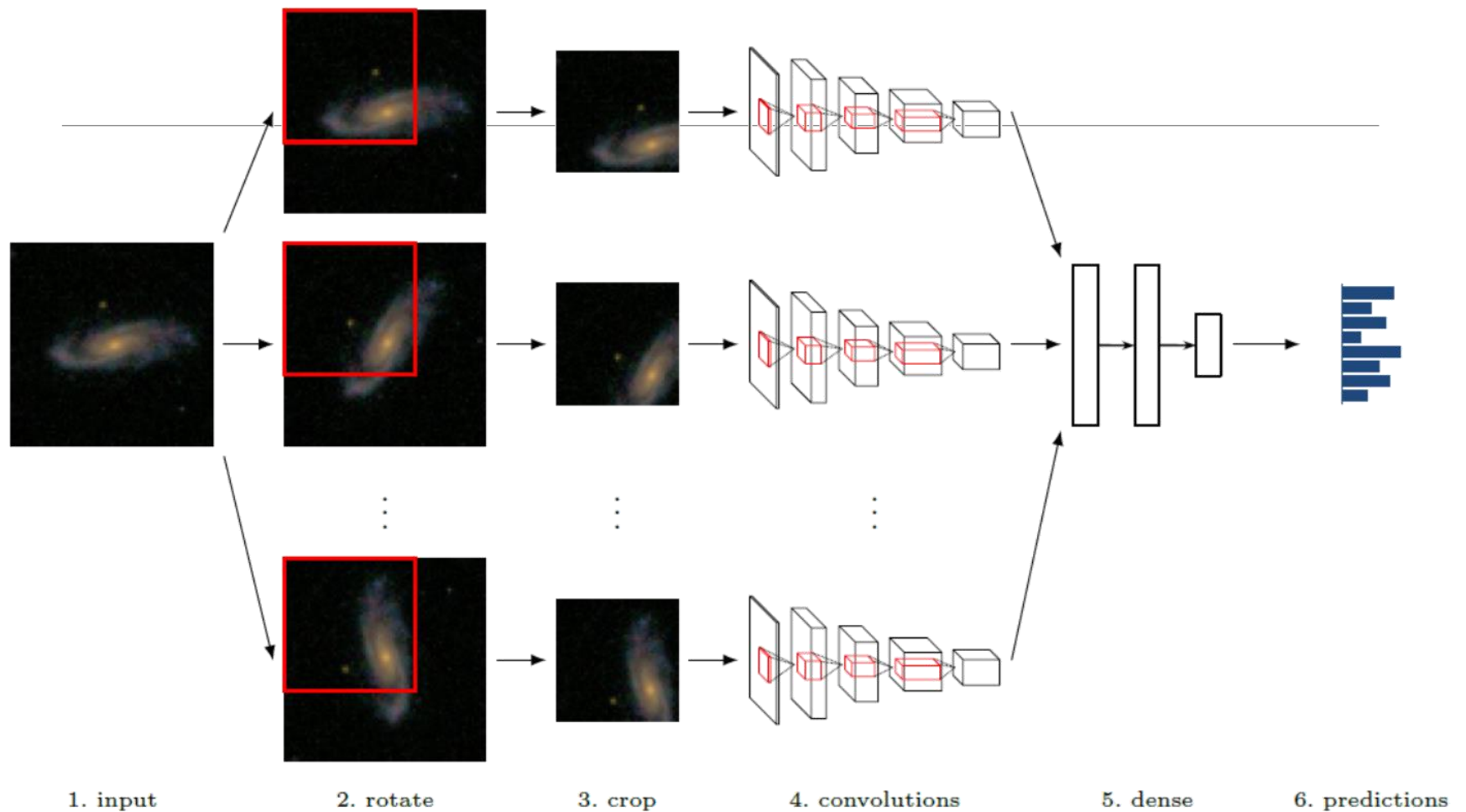
# Skin cancer detection
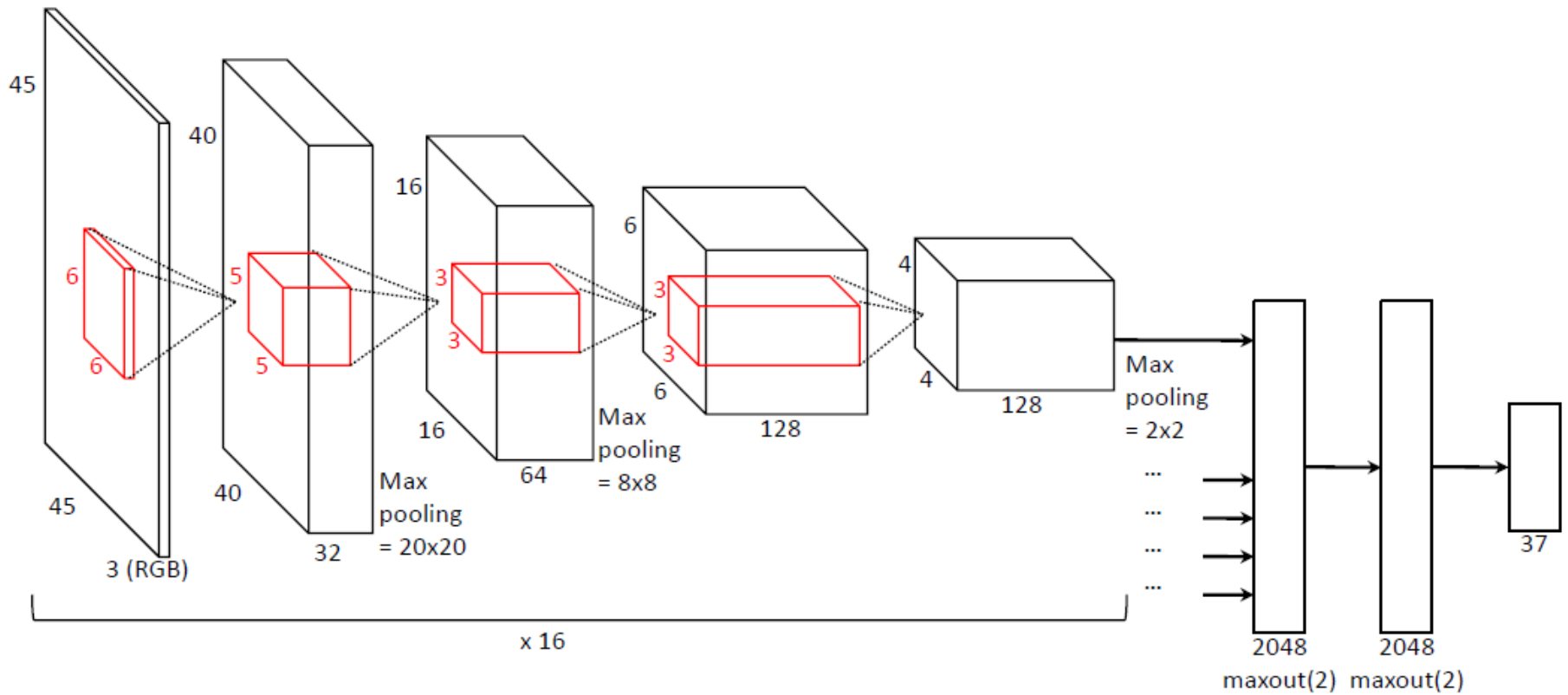
# The Galaxy zoo challenge



Online crowdsourcing project where users describe the morphology of galaxies based on color images 1 million galaxies imaged by the Sloan Digital Sky Survey (2007)

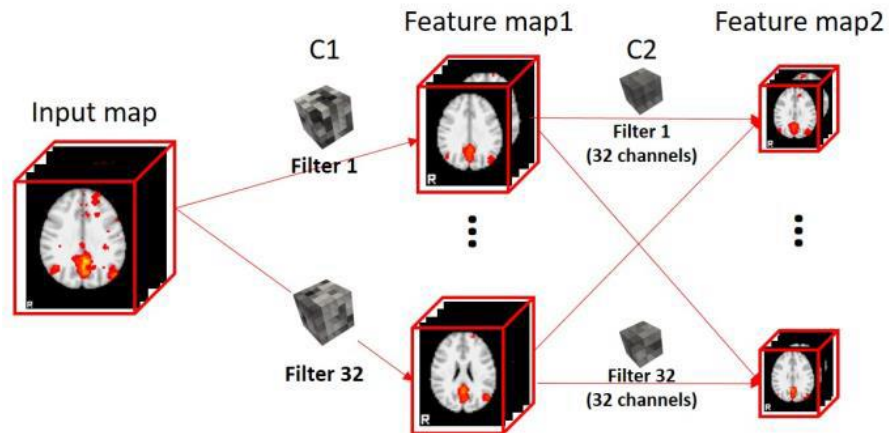1. input    2. rotate    3. crop    4. convolutions    5. dense    6. predictions
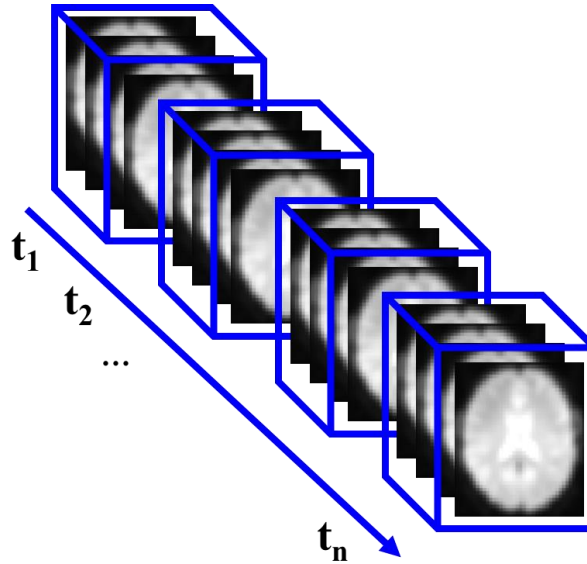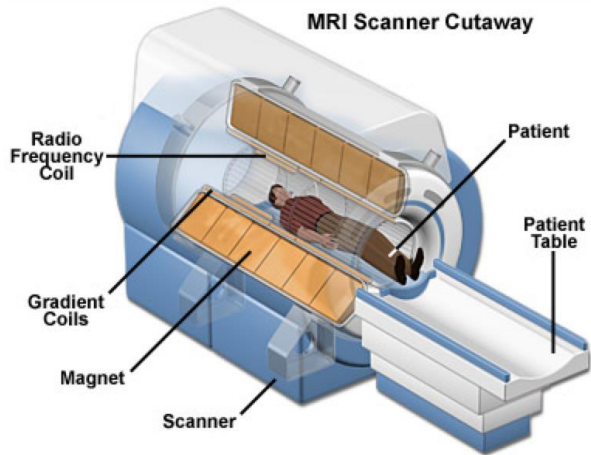
Dieleman, S., Kyle W. W., and Joni D.. "Rotation-invariant convolutional neural networks for galaxy morphology prediction." Monthly notices of the royal astronomical society, 2015
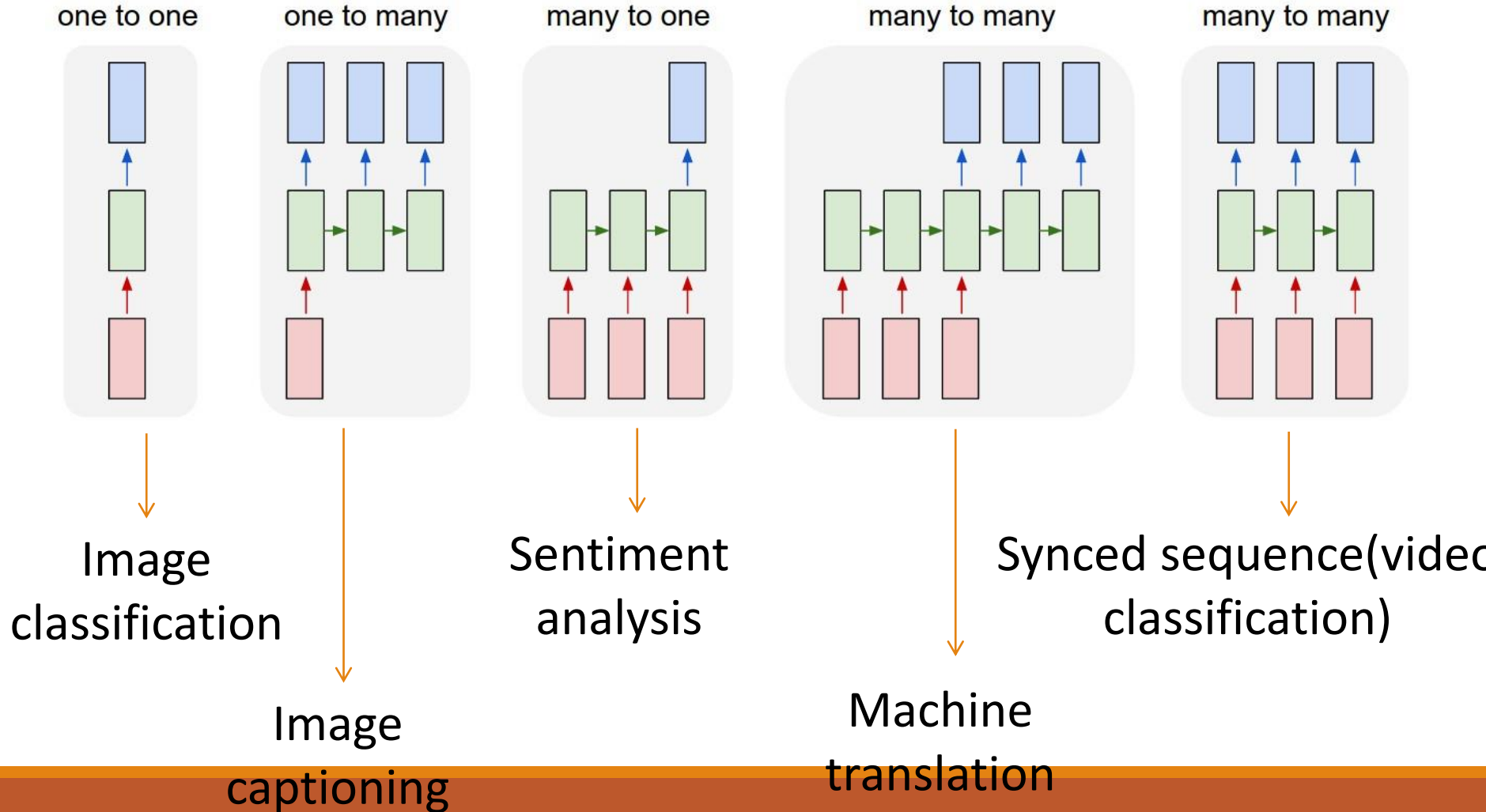
# Component

# CNN & FMRI

# Demos

https://www.clarifai.com/demo

# Different types of mapping



one to one — Image classification

one to many — Image captioning

many to one — Sentiment analysis

many to many — Machine translation

many to many — Synced sequence(video classification)
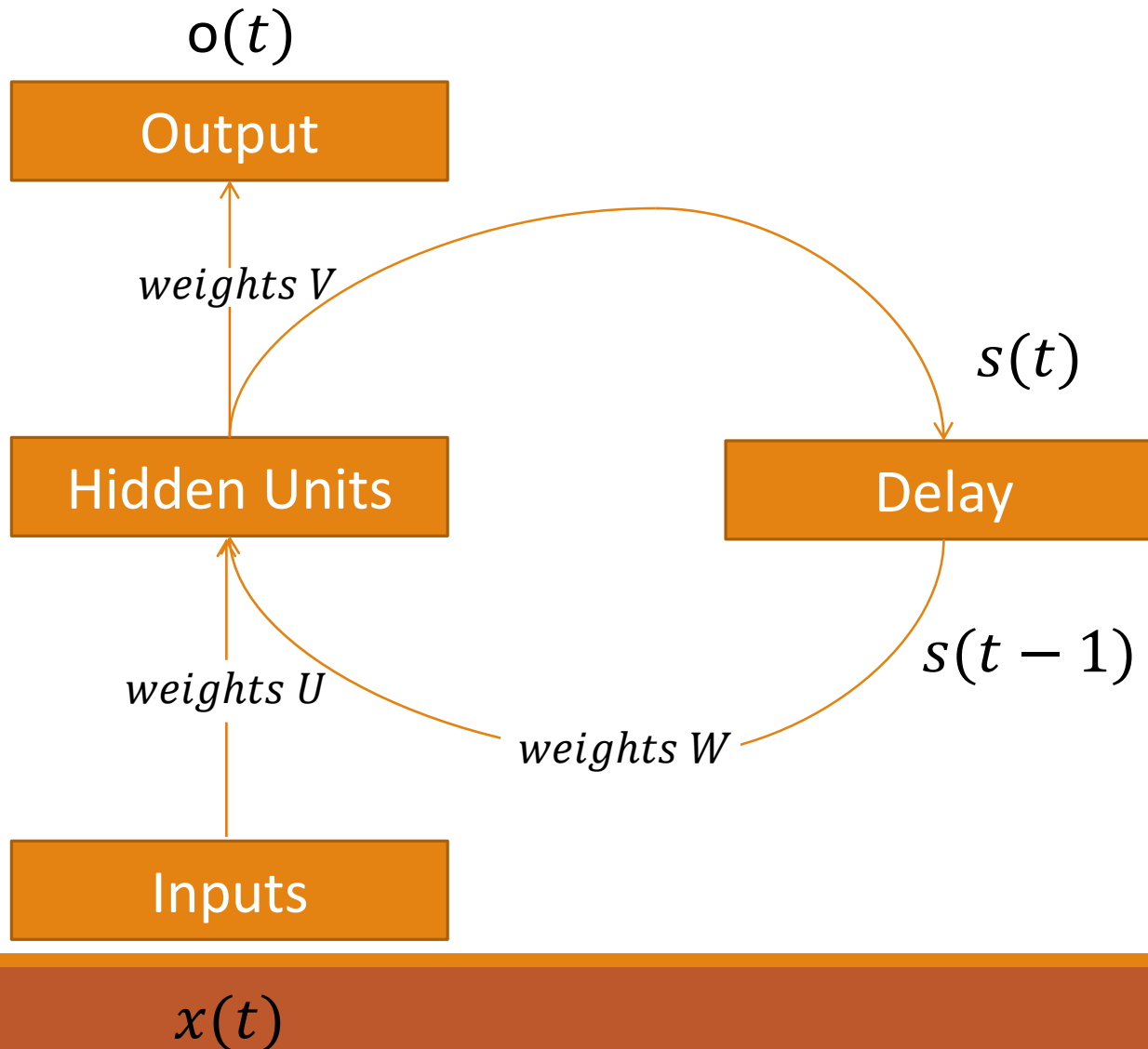
# Recurrent Neural Networks

**Motivation**

➢Feed forward networks accept a fixed-sized vector as input and produce a fixed-sized vector as output

➢fixed amount of computational steps

➢recurrent nets allow us to operate over *sequences* of vectors
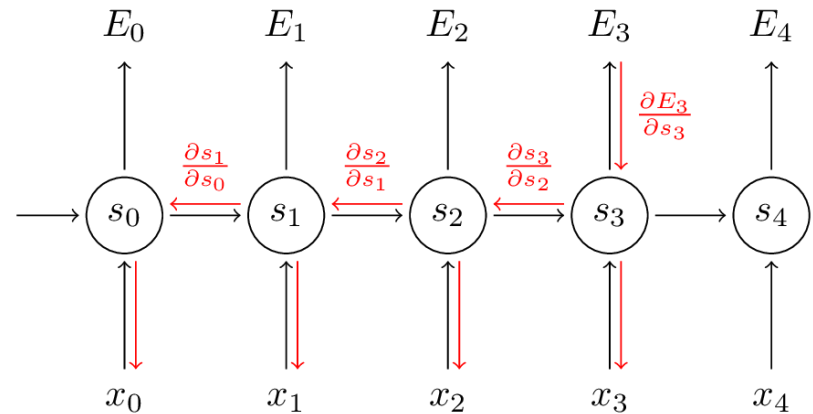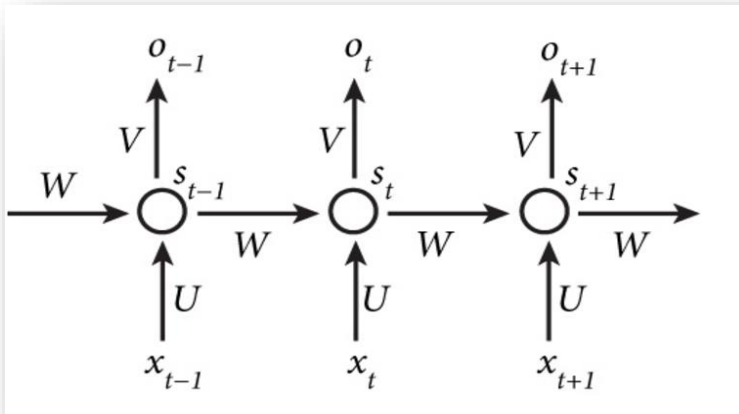
**Use cases**

➢ Video

➢ Audio

➢ Text

# RNN Architecture

# Unfolding RNNs

➢Each node represents a layer of network units at a single time step.

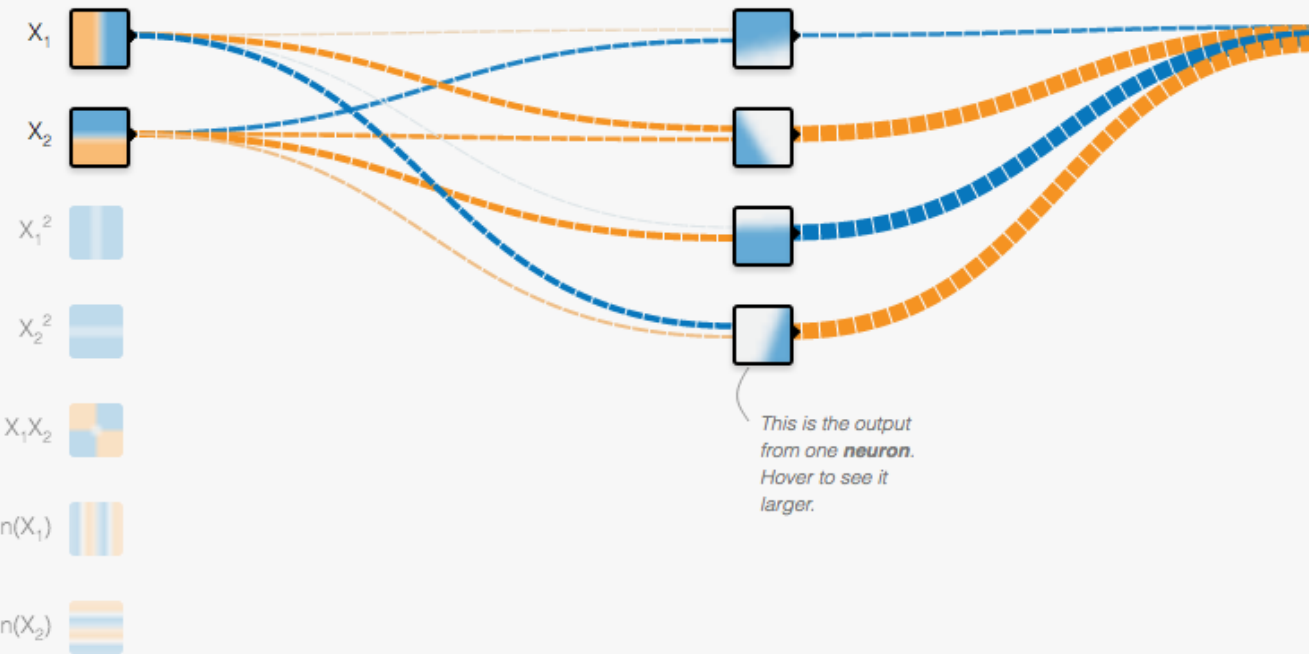➢The same weights are reused at every time step.

# Multi-Layer Network Demo

http://playground.tensorflow.org/