

# Deep Learning

## COMP 5329

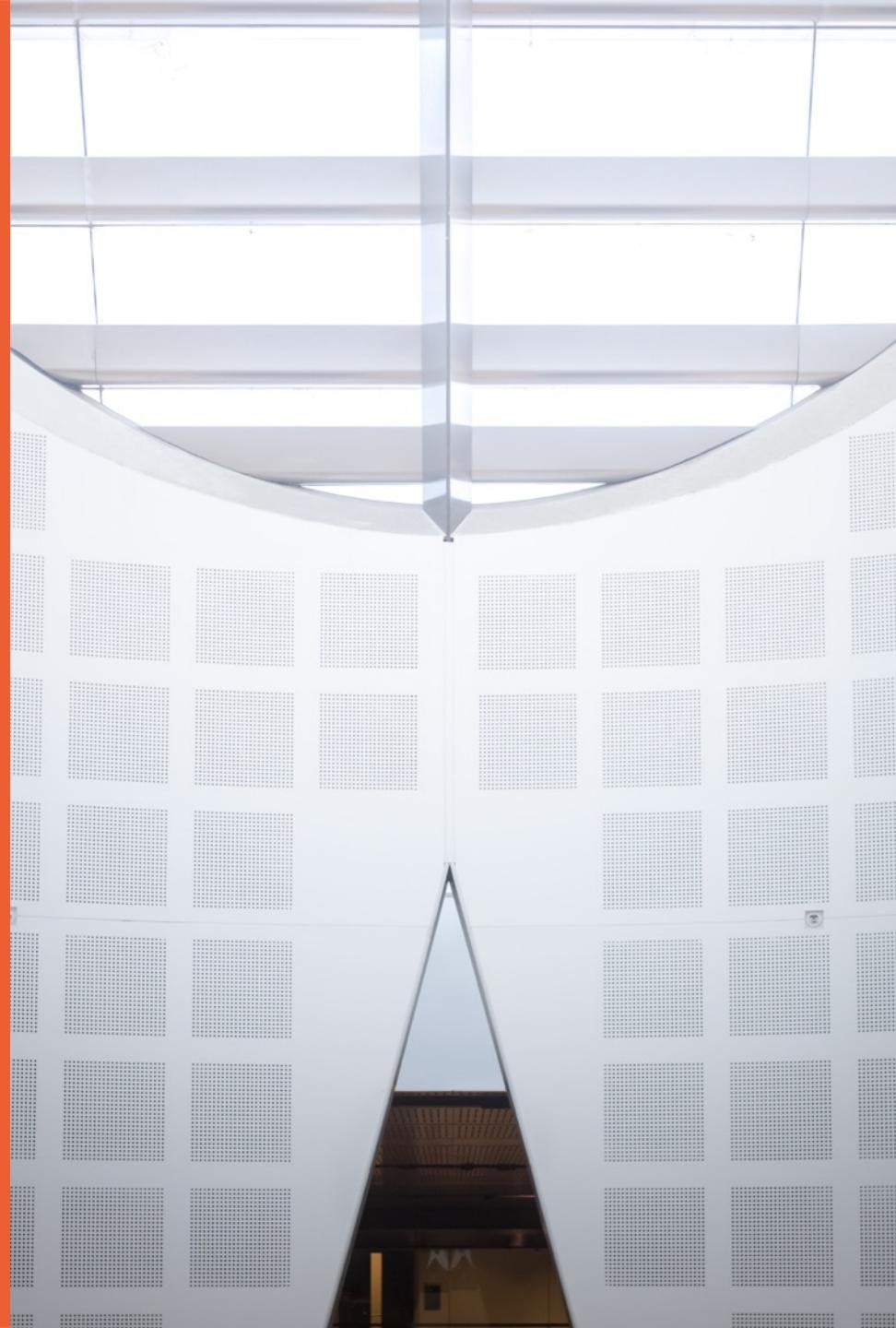
Dr Chang Xu

Prof Dacheng Tao

UBTECH Sydney AI Centre



THE UNIVERSITY OF  
SYDNEY



Daryl Scott Bella Emil  
Xinyu David Trung Lily  
Prerona Richard Handuo  
Aditya Sriram  
Daniel Sabarish Darius  
Briar Nicholas William  
Chenghao Alister  
Wenzhe Fangzhou Elija  
Michael Vincent  
James Quang Rishu Chengzhi  
Sean Shailesh Charles  
Flora Luke Bolivar Kristopher Nisal Josh  
Lucas Rui Thomas Chester Yixuan  
Ryan Rui Thomas Zegun  
Rennie Tushar Yanyu Tim Sicheng  
Jialu Nick Owen & Mike Yun  
Nick Minjung Scott Kevin Zhan  
Yixiong Amy Tom Ken Junpeng Tianfeng  
Gerhardus Youyou Grant Xia Xinyu Ujjwal  
Samuel Leslie Elaine Andrew Joseph

# Learning

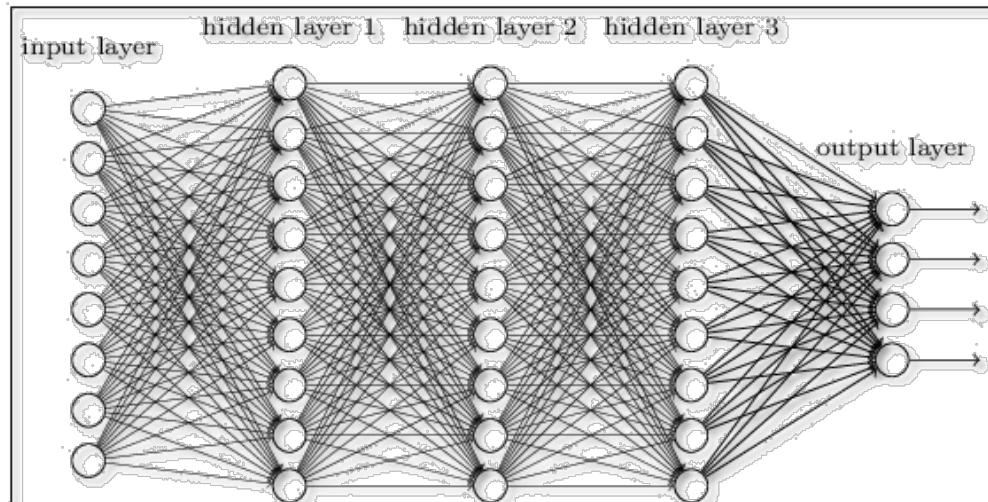
# **Unit of Study Survey (USS)**



# What is Deep Learning?

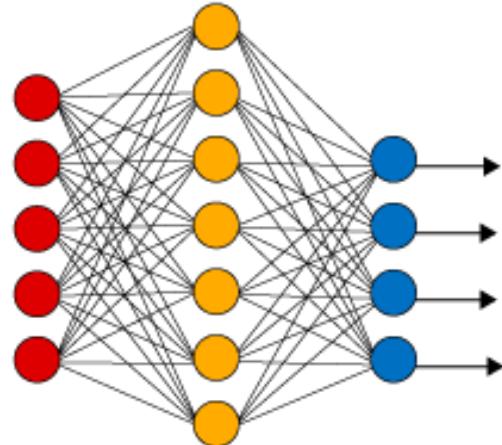
## The short answers

1. ‘Deep Learning’ means using a **neural network** with several layers of nodes between input and output
2. the series of layers between input & output to do feature identification and processing in a series of stages, just as our brains seem to.

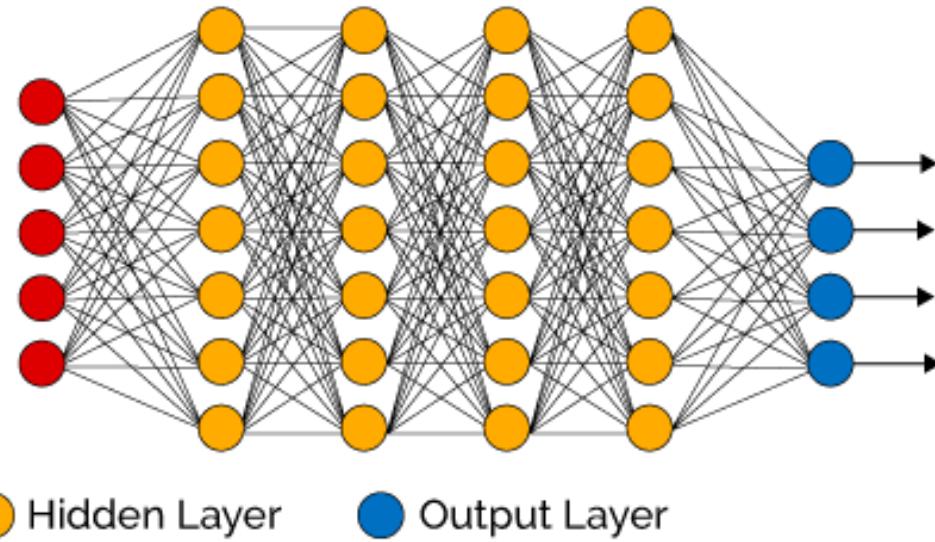


# What is Deep Learning?

Simple Neural Network



Deep Learning Neural Network

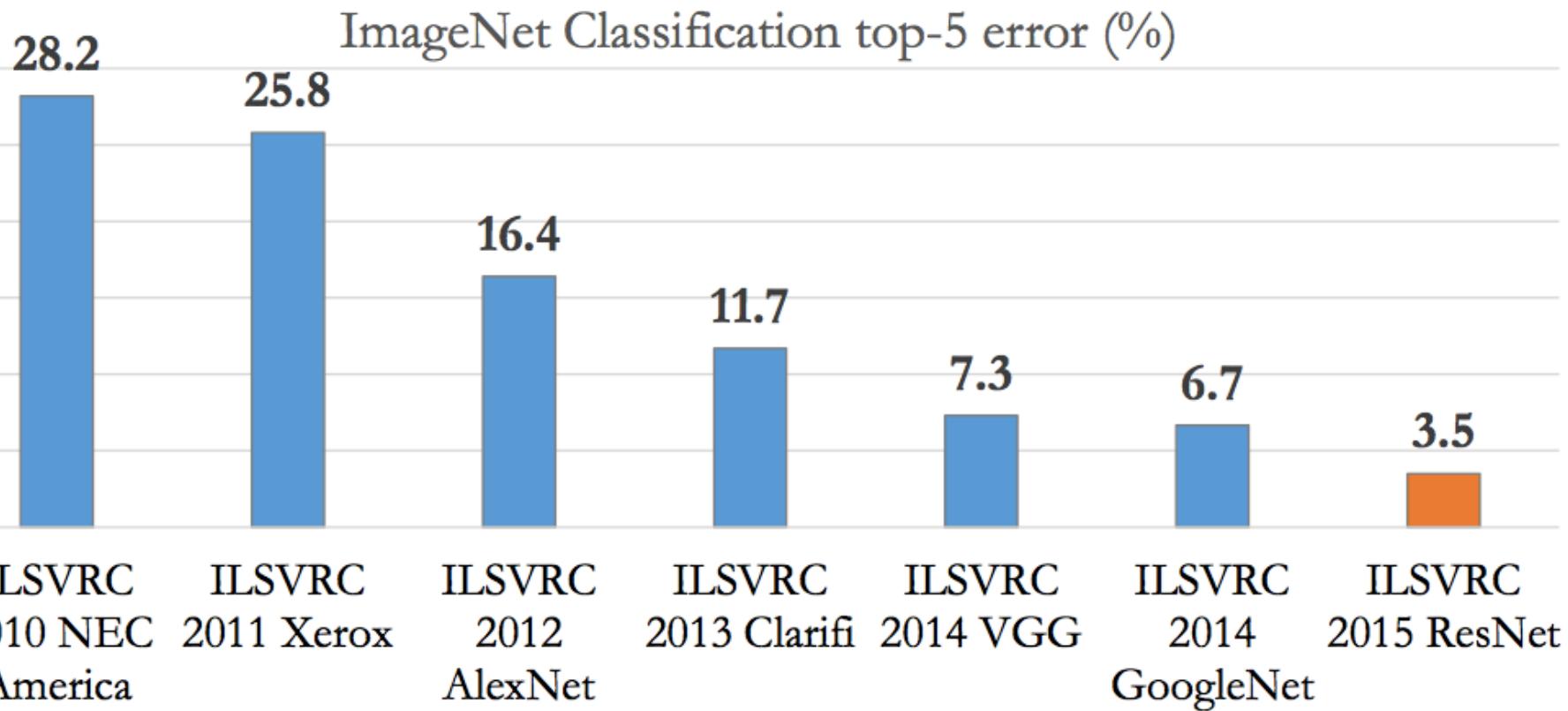


We have always had good algorithms for learning the weights in networks with 1 hidden layer;

But these algorithms are not good at learning the weights for networks with more hidden layers

Image credit to [becominghuman.ai](http://becominghuman.ai)

# Why to Study Deep Learning?



Human: 5.1% Top-5 error rate.

# Why to Study Deep Learning?



Google Deep Mind

AlphaGO defeats world champion

# Why to Study Deep Learning?



# What to Learn?

1. Introduction to Deep Learning
2. Unsupervised feature learning
3. Multilayer neural network
4. Convolutional neural networks (CNN)
5. Optimization for Training Deep Models
6. Network structures
7. Recurrent Neural Networks and LSTM
8. Deep learning applications
9. Generative Models
10. Deep reinforcement learning

# **Unsupervised Feature Learning**

# Principal Component Analysis (PCA)

**Objective:** Maximize Variance

Let  $\mathbf{X} = \{\mathbf{x}_i \in \mathbb{R}^d\}_{1 \leq i \leq n}$  be the sample, project  $\mathbf{X}$  onto a  $k$  dimensional subspace ( $k < d$ ) such that the variance of the projected data is maximized.

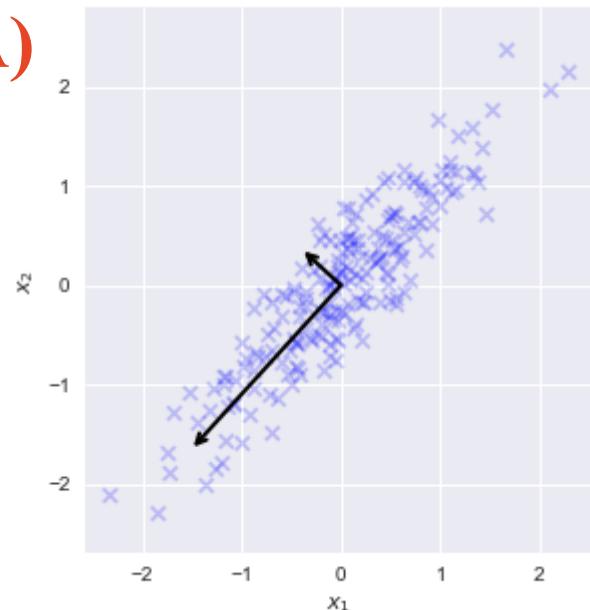
For example,  $k = 1$  :

$$\widehat{Var} = \frac{1}{n} \sum_{i=1}^n (\mathbf{u}_1^\top \mathbf{x}_i - \frac{1}{n} \sum_{i=1}^n \mathbf{u}_1^\top \mathbf{x}_i)^2 = \mathbf{u}_1^\top \mathbf{S} \mathbf{u}_1$$

**Objective:** Minimize Error

$$\min_{\mathbf{U}} \frac{1}{n} \sum_{i=1}^n \|\mathbf{x}_i - \mathbf{U}\mathbf{U}^\top \mathbf{x}_i\|^2 \quad s.t. \quad \mathbf{U}^\top \mathbf{U} = \mathbf{I}$$

**Algorithm:** Compute scatter matrix  $\mathbf{S} = \frac{1}{n} \sum_{i=1}^n \mathbf{z}_i \mathbf{z}_i^\top$  of normalized data, Compute the eigenvectors and eigenvalues by using SVD of  $\mathbf{S}$ . Then transform the data  $\mathbf{y}_i = \mathbf{U}^\top \mathbf{x}_i$ .



# Canonical Correlation Analysis (CCA)

**Objective:** Maximize correlation of  $\mathbf{u}^\top X$  and  $\mathbf{v}^\top Y$ , which are projected variables of two random vectors under projection vectors  $\mathbf{u}$  and  $\mathbf{v}$ .

For example,  $k = 1$ :

With the definition of  $\widehat{\text{cov}}$  and  $\widehat{\text{corr}}$ , model the problem as

$$\max_{\mathbf{u}_1, \mathbf{v}_1} \frac{\mathbf{u}_1^\top \mathbf{S}_{xy} \mathbf{v}_1}{\sqrt{\mathbf{u}_1^\top \mathbf{S}_{xx} \mathbf{u}_1} \sqrt{\mathbf{v}_1^\top \mathbf{S}_{yy} \mathbf{v}_1}}$$

and equals to:  $\max_{\mathbf{u}_1, \mathbf{v}_1} \mathbf{u}_1^\top \mathbf{S}_{xy} \mathbf{v}_1 \quad s.t. \quad \mathbf{u}_1^\top \mathbf{S}_{xx} \mathbf{u}_1 = 1, \mathbf{v}_1^\top \mathbf{S}_{yy} \mathbf{v}_1 = 1$

**Algorithm:** Scatter matrix computation, SVD of  $\mathbf{S}_{xx}^{-1/2} \mathbf{S}_{xy} \mathbf{S}_{yy}^{-1/2} = \mathbf{W}_u \mathbf{D} \mathbf{W}_v^\top$ ,

and projection matrix:  $\mathbf{U} = \mathbf{S}_{xx}^{-1/2} \mathbf{W}_u[:, 1:k]$ , and  $\mathbf{V} = \mathbf{S}_{yy}^{-1/2} \mathbf{W}_v[:, 1:k]$ .

# Autoencoder (AE)

**Objective:** Learn discriminative representations (through encoder and decoder)

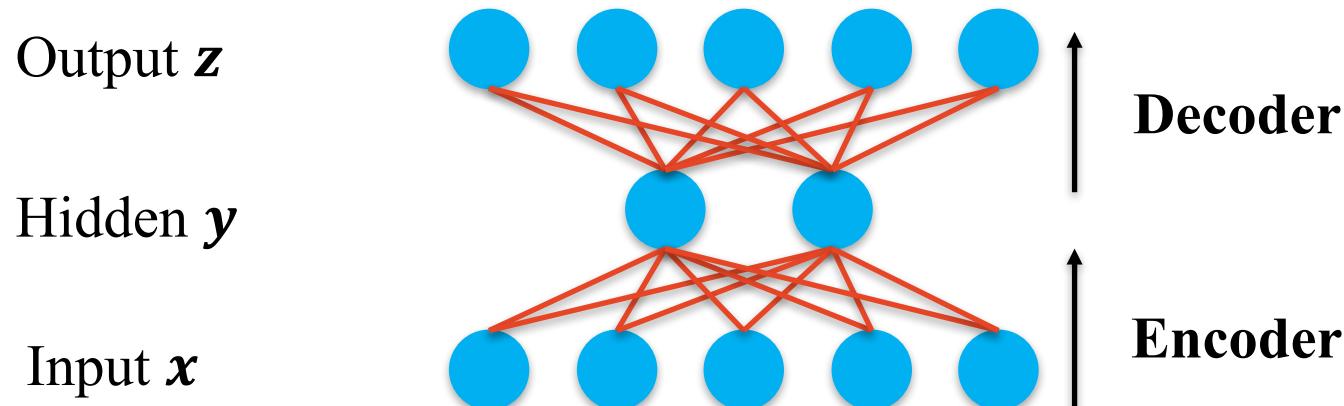
Encoder and decoder function:

$$\mathbf{y} = f(\mathbf{W}^{(1)}\mathbf{x} + \mathbf{b}^{(1)}), \mathbf{z} = f(\mathbf{W}^{(2)}\mathbf{y} + \mathbf{b}^{(2)})$$

Minimize reconstruction error (difference between  $\mathbf{z}$  and  $\mathbf{x}$ )

- Square loss:  $J(\theta) = \|\mathbf{x} - \mathbf{z}\|^2, \theta = \{\mathbf{W}^{(1)}, \mathbf{b}^{(1)}, \mathbf{W}^{(2)}, \mathbf{b}^{(2)}\}$
- Cross-entropy:  $J(\theta) = -\sum_{i=1}^d x_i \log(z_i) + (1 - x_i)\log(1 - z_i)$

**Algorithm:** Backpropagation (BP)

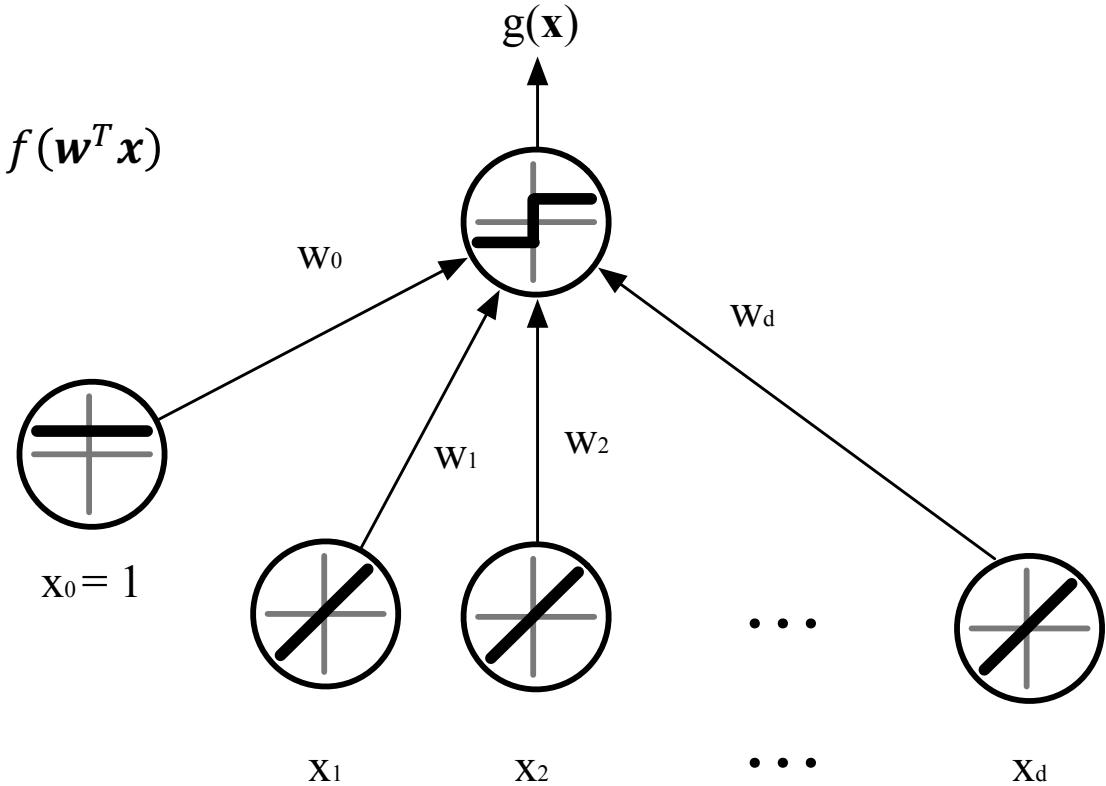


# Multilayer Neural Network

# Perceptron

$$g(\mathbf{x}) = f\left(\sum_{i=1}^d x_i w_i + w_0\right) = f(\mathbf{w}^T \mathbf{x})$$

$$f(s) = \begin{cases} 1, & \text{if } s \geq 0 \\ -1, & \text{if } s < 0 \end{cases}$$

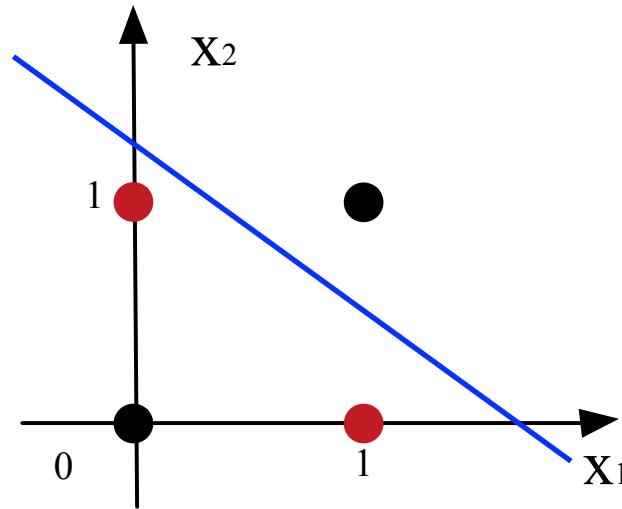


Linear classifier!

# Perceptron: Limitations

- Solving XOR (exclusive-or) with Perceptron?

$x_1$	$x_2$	$y$
0	0	0
0	1	1
1	0	1
1	1	0

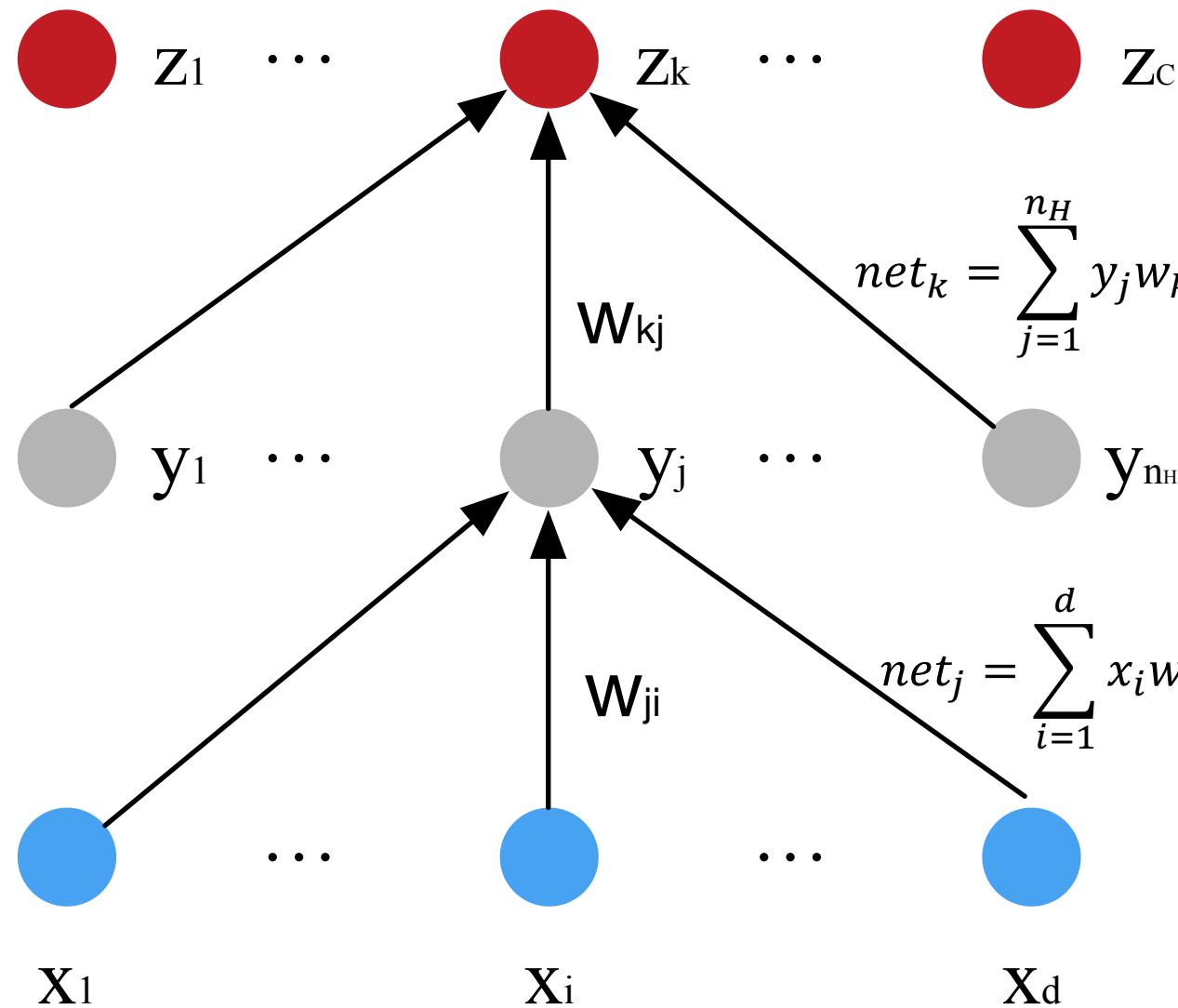


- Linear classifier cannot identify non-linear patterns.

$$w_1x_1 + w_2x_2 + b$$

This failure caused the majority of researchers to walk away.

## Feedforward



$$z_k = f(\text{net}_k)$$

$$= f\left(\sum_{j=1}^{n_H} w_{kj} f(\text{net}_j) + w_{k0}\right)$$

$Z_C$

$$\text{net}_k = \sum_{j=1}^{n_H} y_j w_{kj} + w_{k0} = \sum_{j=0}^{n_H} y_j w_{kj} = \mathbf{w}_k^T \mathbf{y}$$

$$\text{net}_j = \sum_{i=1}^d x_i w_{ji} + w_{j0} = \sum_{i=0}^d x_i w_{ji} = \mathbf{w}_j^T \mathbf{x}$$

## Backpropagation summary

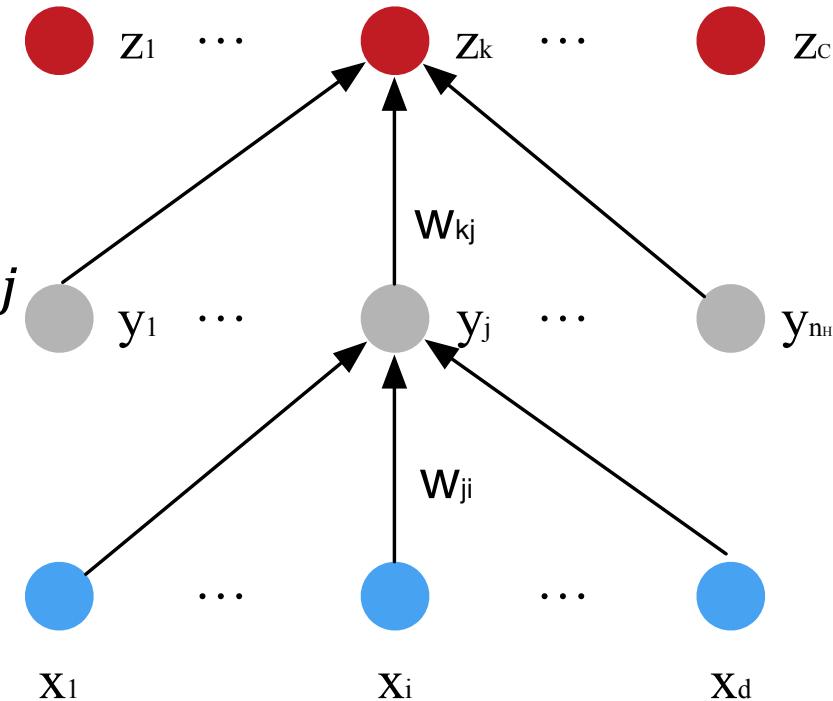
$$\delta_k = -\frac{\partial J}{\partial net_k}$$

$$\delta_j = -\frac{\partial J}{\partial net_j}$$

- Hidden-to-output weight

$$\Delta w_{kj} = \eta \delta_k y_j$$

$$= \eta(t_k - z_k)f'(net_k)y_j$$



- Input-to-hidden

$$\Delta w_{ji} = \eta \delta_j x_i$$

$$= \eta \left[ \sum_{k=1}^C \delta_k w_{kj} \right] f'(net_j) x_i$$

# Universal Approximation Theorem

Let  $f(\cdot)$  be a **nonconstant, bounded, and monotonically-increasing continuous** function. Let  $I_m$  denote the  $m$ -dimensional unit hypercube  $[0,1]^m$ . The space of continuous functions on  $I_m$  is denoted by  $C(I_m)$ . Then, given any  $\varepsilon > 0$  and any function  $f \in C(I_m)$ , **there exist an integer  $N$ , real constants  $v_i, b_i \in \mathbb{R}^m$  and real vectors  $w_i \in \mathbb{R}^m$ , where  $i = 1, \dots, N$ , such that we may define:**

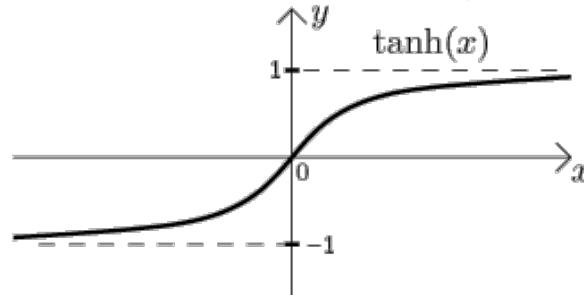
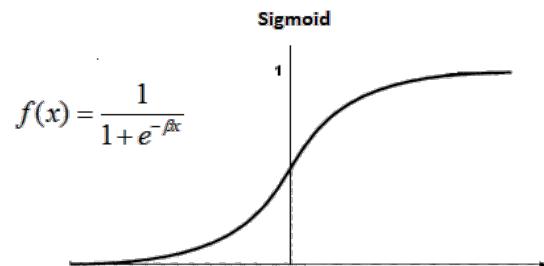
$$\hat{F}(x) = \sum_i^N v_i f(w_i^T x + b_i)$$

as an approximate realization of the function  $F$  where  $F$  is independent of  $f$ ; that is,

$$|\hat{F}(x) - F(x)| < \varepsilon$$

for all  $x \in I_m$ . In other words, functions of the form  $\hat{F}(x)$  are dense in  $C(I_m)$ .

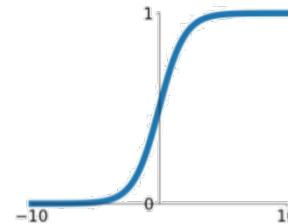
**Universality theorem (Hecht-Nielsen 1989):** “Neural networks with a single hidden layer can be used to approximate any continuous function to any desired precision.”



# Activation function $f(s)$

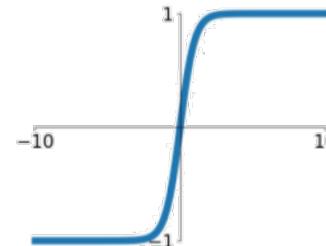
- Popular choice of activation functions (single input)
  - Sigmoid function

$$f(s) = \frac{1}{1 + e^{-s}}$$



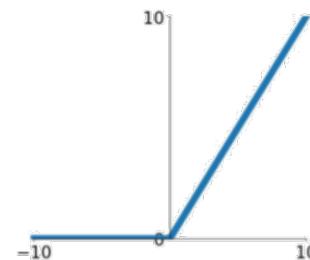
- Tanh function: shift the center of Sigmoid to the origin

$$f(s) = \frac{e^s - e^{-s}}{e^s + e^{-s}}$$



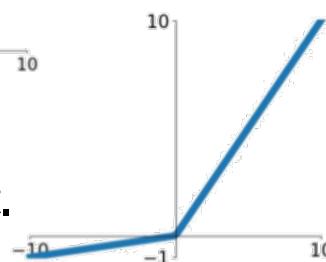
- Rectified linear unit (ReLU)

$$f(s) = \max(0, s)$$



- Leaky ReLU

$$f(s) = \begin{cases} s, & \text{if } s \geq 0 \\ \alpha s, & \text{if } s < 0 \end{cases}, \quad \alpha \text{ is a small constant.}$$



# **Optimization for Training Deep Models**

# Momentum

The momentum term increases for dimensions whose gradients point in the same directions and reduces updates for dimensions whose gradients change directions.

$$v_t = \gamma v_{t-1} + \eta \nabla_{\theta} J(\theta)$$

$$\theta_t = \theta_{t-1} - v_t$$

Momentum term  $\gamma$  is usually set to 0.9 or a similar value.

# Nesterov accelerated gradient (NAG)

First make a big jump in the direction of the previous accumulated gradient.

Then measure the gradient where you end up.

Finally accumulate the gradient and make a correction.

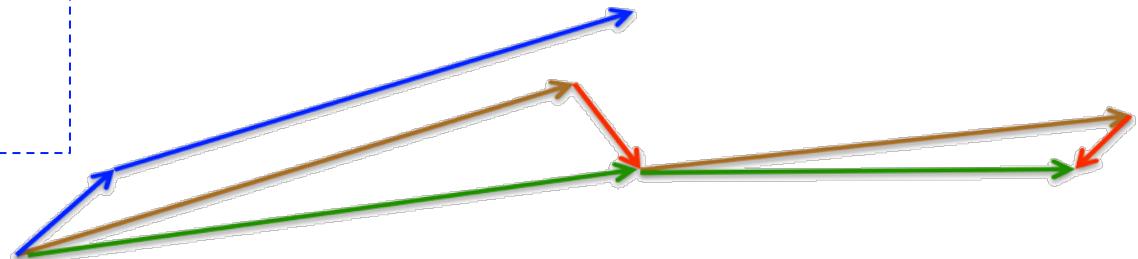
$$v_t = \gamma v_{t-1} + \theta - \theta - v_t$$

$$\eta \nabla_{\theta} J(\theta - \gamma v_{t-1})$$

Finally accumulate the gradient and make a correction

First make a big jump

Then measure the gradient where you end up



blue vectors = standard momentum

brown vector = jump, red vector = correction,  
green vector = accumulated gradient ,

- **Adagrad**  $\theta_t = [\theta_{t,1}, \theta_{t,2}, \dots, \theta_{t,i}, \dots, \theta_{t,d}]$

- **Original gradient descent:** The same for all parameters

$$\theta_{t+1,i} = \theta_{t,i} - \eta \cdot \nabla_{\theta} J(\theta_i)$$

Perform an update for all parameters using the same learning rate.

- **Adagrad:** Adapt the learning rate to the parameters based on the past gradients that have been computed for  $\theta_i$ .

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii} + \epsilon}} \cdot \nabla_{\theta} J(\theta_i)$$

$G_t \in R^{d \times d}$  is a diagonal matrix where each diagonal element is the sum of the squares of the gradients w.r.t.  $\theta_i$  up to time step t.

## Adadelta – 1 Accumulate over window

- Storing fixed previous squared gradients cannot accumulate to infinity and instead becomes a local estimate using recent gradients.
- Adadelta implements it as an exponentially **decaying average** of all past squared gradients.

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2$$

$$RMS[g]_t = \sqrt{E[g^2]_t + \epsilon}$$

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t$$

## Adadelta – 2 Correct units with Hessian approximation

$$H^{-1}g = \Delta\theta$$

$$\Delta\theta = \frac{\frac{\partial J}{\partial\theta}}{\frac{\partial^2 J}{\partial\theta^2}} \Rightarrow H^{-1} \propto \frac{1}{\frac{\partial^2 J}{\partial\theta^2}} \propto \frac{\Delta\theta}{\frac{\partial J}{\partial\theta}} \propto \frac{\Delta\theta}{g}$$

We assume the curvature is locally smooth and approximate  $\Delta\theta_t$  by computing the exponentially decaying RMS of  $\Delta\theta$ .

$$\Delta\theta_t = -\frac{\eta}{RMS[g]_t} g_t \xrightarrow{\text{red arrow}} \Delta\theta_t = -\frac{RMS[\Delta\theta]_{t-1}}{RMS[g]_t} g_t$$

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

## RMSprop

- RMSprop and Adadelta have both been developed independently around the same time.
- RMSprop in fact is identical to the first update vector of Adadelta.

$$\theta_{t+1} = \theta_t + \Delta\theta_t$$

$$\Delta\theta_t = -\frac{\eta}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2$$

Hinton suggests  $\gamma$  to be set to 0.9, while a good default value for the learning rate  $\eta$  is 0.001.

# Adam

Combine the advantages of **Adgrad** and **RMSprop**.

- Adgrad: adaptive learning rate for different parameters.

$$\theta_{t+1} = \theta_t - \frac{\eta}{\sqrt{v_t + \epsilon}} m_t$$

- RMSprop: exponentially decaying average.

- Store an exponentially decaying average of **past squared gradients**  $v_t$  .

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

Keep an exponentially decaying average of **past gradients**  $m_t$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

# Initialization

All zero initialization:

Every neuron computes the same output.



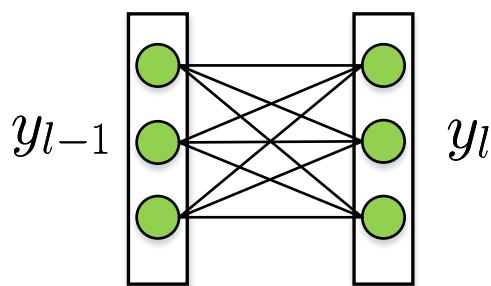
They will also compute the same gradients.



They will undergo the exact same parameter updates.



There will be no difference between all the neurons.



$$y_l = W y_{l-1} + b, \quad \frac{\partial \mathcal{L}}{\partial y_{l-1}} = W^\top \frac{\partial \mathcal{L}}{\partial y_l}$$

if  $W = 0$ , then  $\frac{\partial \mathcal{L}}{\partial y_{l-1}} = 0$

then,  $\frac{\partial \mathcal{L}}{\partial y_m} = 0, \quad m = 0, \dots, l-1$ .

## Initialization (Xavier) – [Xavier Glorot, Yoshua Bengio 2010]

The variance of the inputs and outputs are the same for all layers.

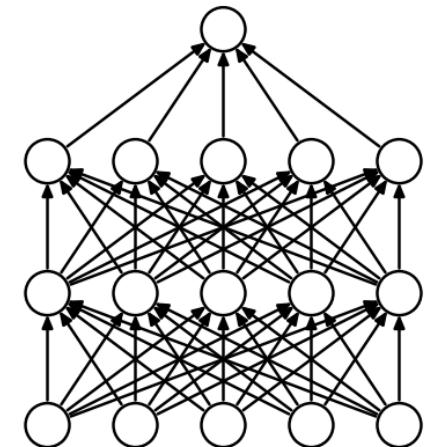
But the variance of the gradients might still vanish or explode as we consider deeper networks.

$$s^i = Z^i W^i + b^i$$

$$z^{i+1} = f(s^i)$$

$$\text{Var}\left[\frac{\partial \text{Cost}}{\partial s^i}\right] = [n\text{Var}[W]]^{d-i}\text{Var}[x]$$

$$\text{Var}\left[\frac{\partial \text{Cost}}{\partial w^i}\right] = [n\text{Var}[W]]^d\text{Var}[x]\text{Var}\left[\frac{\partial \text{Cost}}{\partial s^i}\right]$$



$$W \sim U\left[-\frac{\sqrt{6}}{\sqrt{n_j+n_{j+1}}}, \frac{\sqrt{6}}{\sqrt{n_j+n_{j+1}}}\right]$$

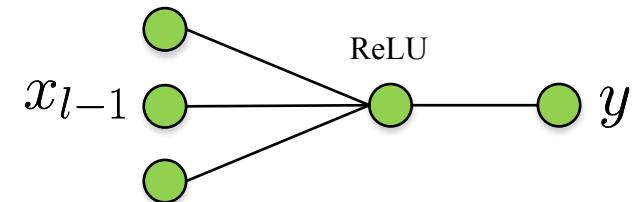
$U[-a, a]$  is the uniform distribution in the interval  $(-a, a)$ .

## Initialization (He) -- [Kaiming He, Xiangyu Zhang, et al. 2015]

If we use ReLU as the activation function,  $n$  is the number of inputs, then for layer  $l$  we have:

$$\text{Var}(y_l) = (n\text{Var}(w)) E[x_{l-1}^2]$$

$$x_{l-1} = \max(0, y_{l-1})$$



If  $y_{l-1}$  has zero mean, and has a symmetric distribution around zero:

$$E[x_{l-1}^2] = E[(\max(0, y_{l-1}))^2] = \frac{1}{2}\text{Var}[y_{l-1}]$$

So we have:

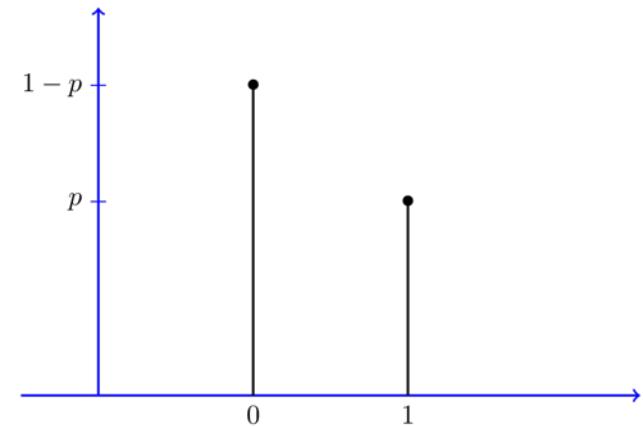
$$\boxed{\text{Var}(w) = \frac{2}{n}}$$

# Dropout

With dropout:

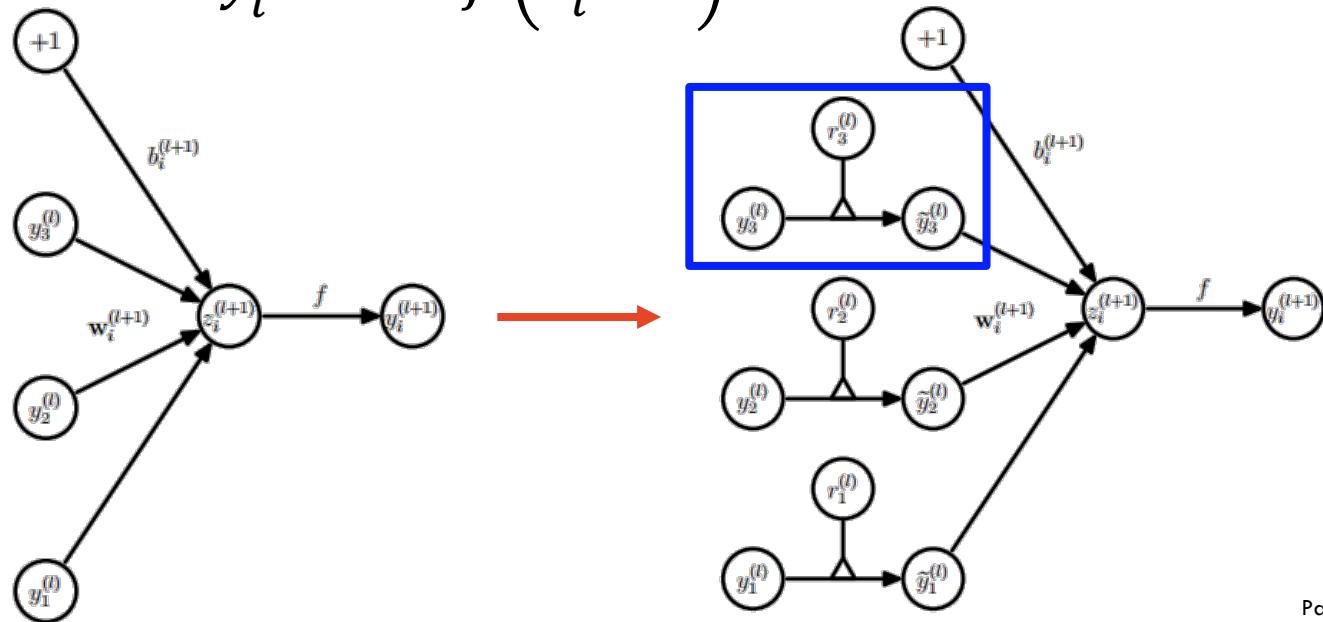
$$r_i^{(l)} \sim \text{Bernoulli}(p)$$

$$\tilde{\mathbf{y}}^{(l)} = \mathbf{r}^{(l)} * \mathbf{y}^{(l)}$$



$$z_i^{(l+1)} = \mathbf{w}_i^{(l+1)} \tilde{\mathbf{y}}^{(l)} + b_i^{(l+1)}$$

$$y_i^{(l+1)} = f(z_i^{(l+1)})$$



# Batch Normalization

**Input:** Network  $N$  with trainable parameters  $\Theta$ ;  
 subset of activations  $\{x^{(k)}\}_{k=1}^K$

**Output:** Batch-normalized network for inference,  $N_{BN}^{inf}$

- 1:  $N_{BN}^{tr} \leftarrow N$  // Training BN network
- 2: **for**  $k = 1 \dots K$  **do**
- 3:   Add transformation  $y^{(k)} = BN_{\gamma^{(k)}, \beta^{(k)}}(x^{(k)})$  to  $N_{BN}^{tr}$  (Alg. 1) ←
- 4:   Modify each layer in  $N_{BN}^{tr}$  with input  $x^{(k)}$  to take  $y^{(k)}$  instead
- 5: **end for**
- 6: Train  $N_{BN}^{tr}$  to optimize the parameters  $\Theta \cup \{\gamma^{(k)}, \beta^{(k)}\}_{k=1}^K$
- 7:  $N_{BN}^{inf} \leftarrow N_{BN}^{tr}$  // Inference BN network with frozen // parameters
- 8: **for**  $k = 1 \dots K$  **do**
- 9:   // For clarity,  $x \equiv x^{(k)}$ ,  $\gamma \equiv \gamma^{(k)}$ ,  $\mu_B \equiv \mu_B^{(k)}$ , etc.
- 10:   Process multiple training mini-batches  $\mathcal{B}$ , each of size  $m$ , and average over them:  

$$E[x] \leftarrow E_{\mathcal{B}}[\mu_B]$$

$$\text{Var}[x] \leftarrow \frac{m}{m-1} E_{\mathcal{B}}[\sigma_B^2]$$
- 11:   In  $N_{BN}^{inf}$ , replace the transform  $y = BN_{\gamma, \beta}(x)$  with  

$$y = \frac{\gamma}{\sqrt{\text{Var}[x]+\epsilon}} \cdot x + \left(\beta - \frac{\gamma E[x]}{\sqrt{\text{Var}[x]+\epsilon}}\right)$$
- 12: **end for**

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots m\}$ ;  
 Parameters to be learned:  $\gamma, \beta$

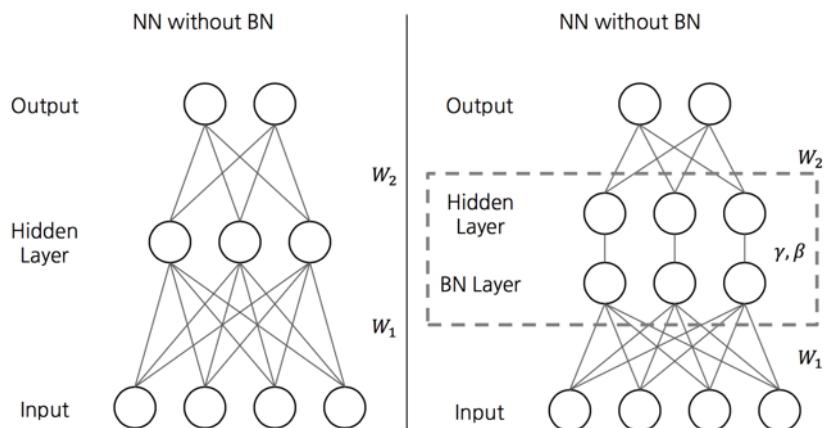
**Output:**  $\{y_i = BN_{\gamma, \beta}(x_i)\}$

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

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

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

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv BN_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$



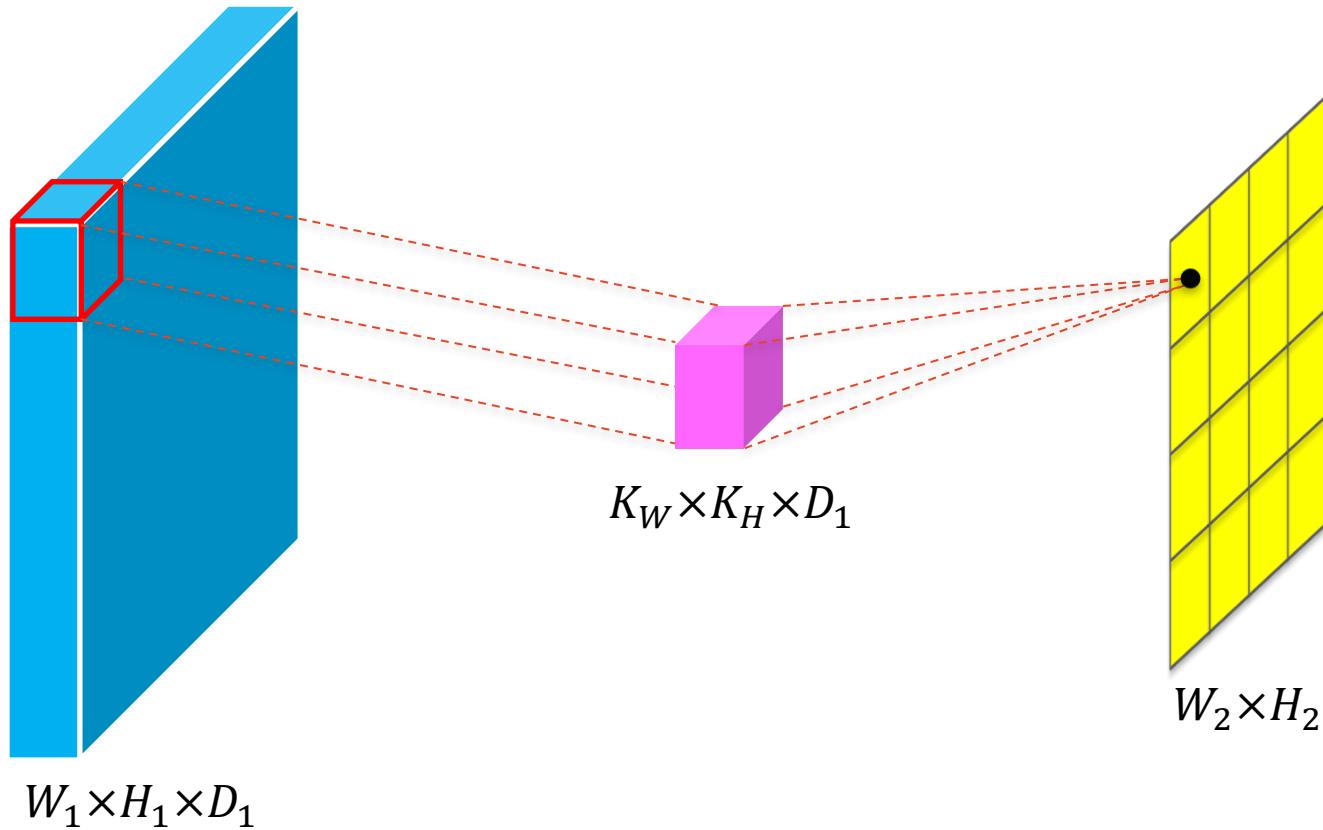
**Algorithm 2:** Training a Batch-Normalized Network

# **Convolutional Neural Networks**

# Convolutional Layer

- Suppose stride is  $(S_W, S_H)$  and pad is  $(P_W, P_H)$

$$W_2 = \frac{W_1 + 2P_W - K_W}{S_W} + 1 \quad \text{and} \quad H_2 = \frac{H_1 + 2P_H - K_H}{S_H} + 1$$



# Pooling

## Max pooling

- Filter size: (2,2)
- Stride: (2,2)
- Pooling ops:  $\max(\cdot)$

-1	2	0	0
0	1	3	-2
0	0	-1	4
3	-1	-2	-2

Feature map

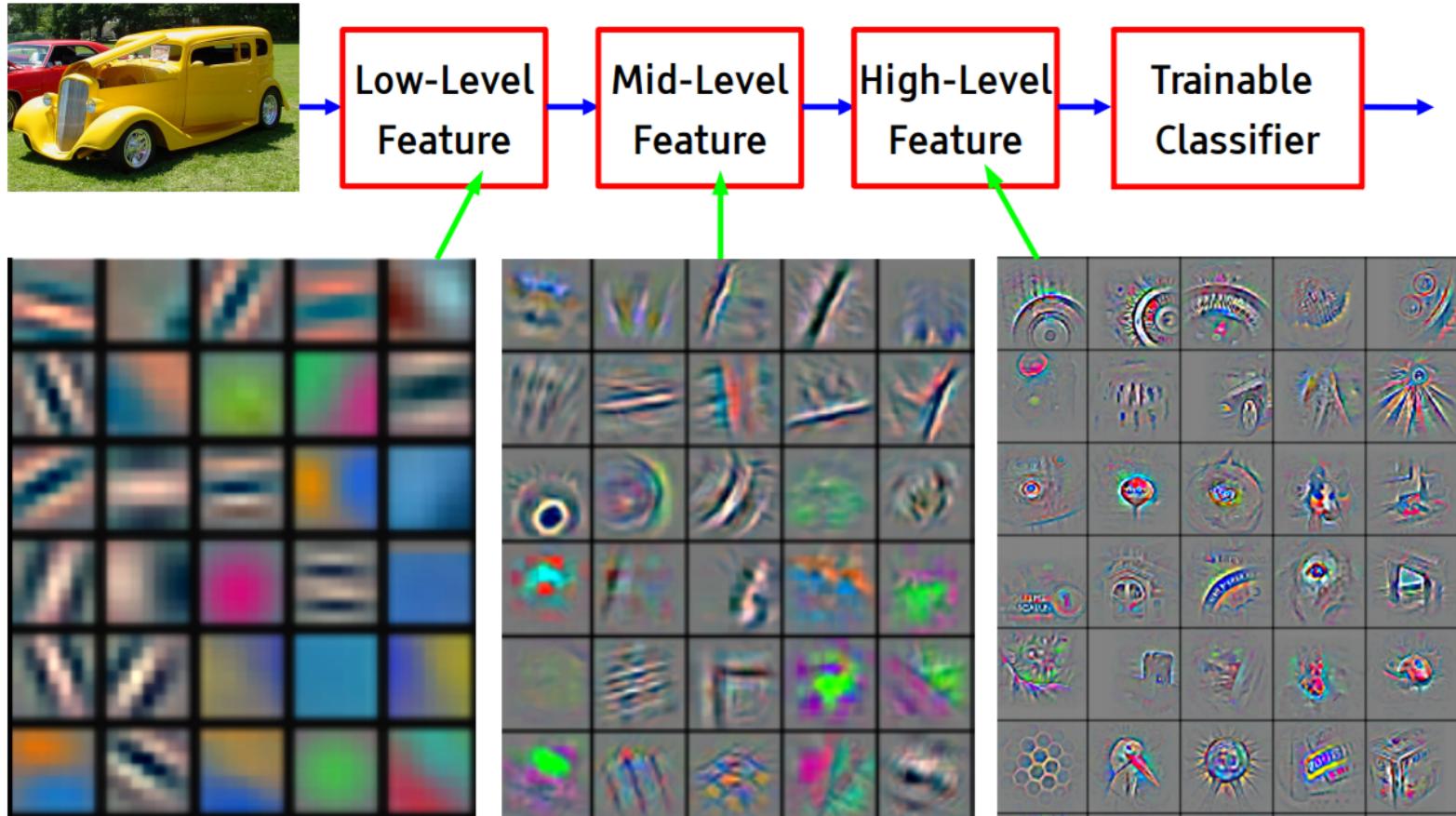


2	3
3	4

Subsample map

# Visualize features

- Give insight into the function of intermediate feature layers and the operation of the classifier



(Zeiler and Fergus, 2014)

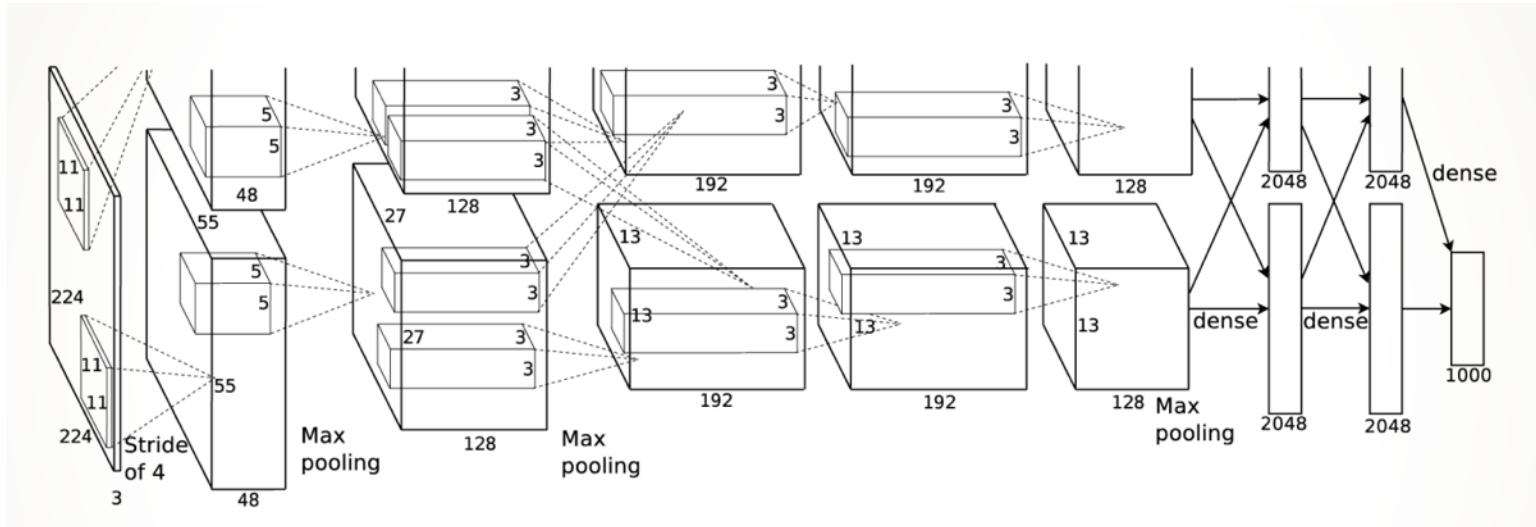
# **Network Structures**

# Winning CNN Architectures

Model	AlexNet	ZF Net	GoogLeNet	Resnet
Year	2012	2013	2014	2015
#Layer	8	8	22	152
Top 5 Acc	15.4%	11.2%	6.7%	3.57%
Data augmentation	✓	✓	✓	✓
Dropout	✓	✓		
Batch normalization				✓

# CNN Architecture: Part I

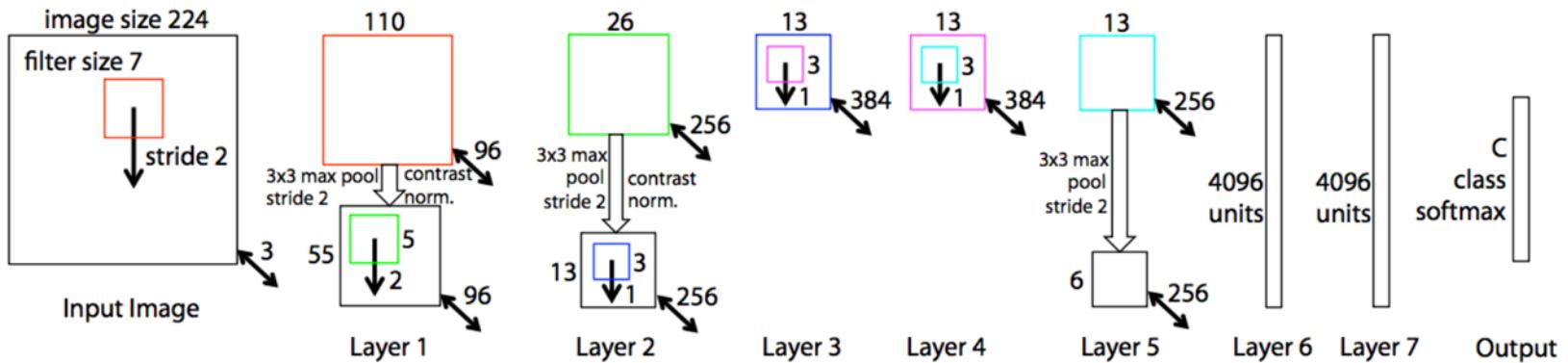
- **AlexNet** [Krizhevsky et al. 2012] - 5 conv layers, 3 fully connected layers.



- Activation function: ReLU
- Data Augmentation
- Dropout (drop rate=0.5)
- Local Response Normalization
- Overlapping Pooling

# CNN Architecture: Part I

## □ ZFNet [Zeiler and Fergus, 2013]



- An improved version of AlexNet: top-5 error from 16.4% to 11.7%
- First convolutional layer: 11x11 filter, stride=4 -> 7x7 filter, stride=2
- The number of filters increase as we go deeper.

# CNN Architecture: Part I

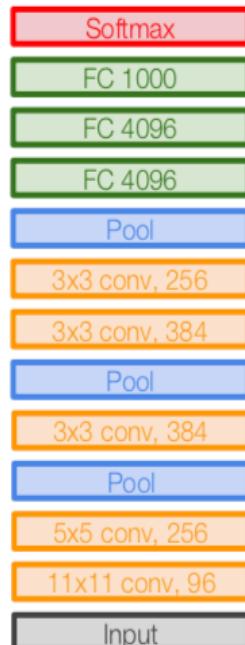
## □ VGGNet [Simonyan and Zisserman, 2014]

### Small filters: Why?

- 3x3 convolutional layers (stride 1, pad 1)
- 2x2 max-pooling, stride 2

### Deeper networks:

- AlexNet: 8 layers
- VGGNet: 16 or 19 layers



AlexNet

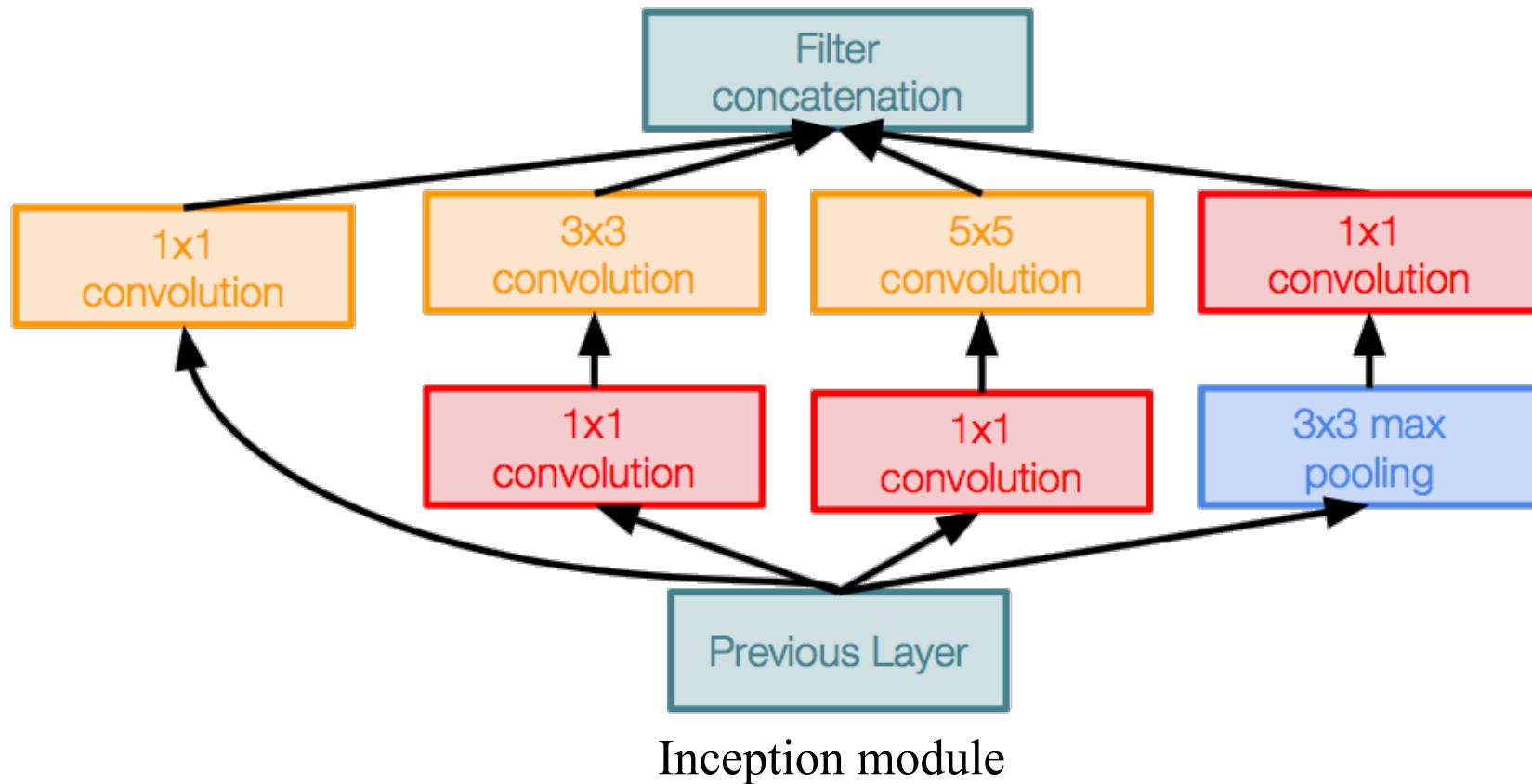


VGG16

VGG19

# CNN Architecture: Part I

## □ GoogLeNet [Szegedy et al., 2014]



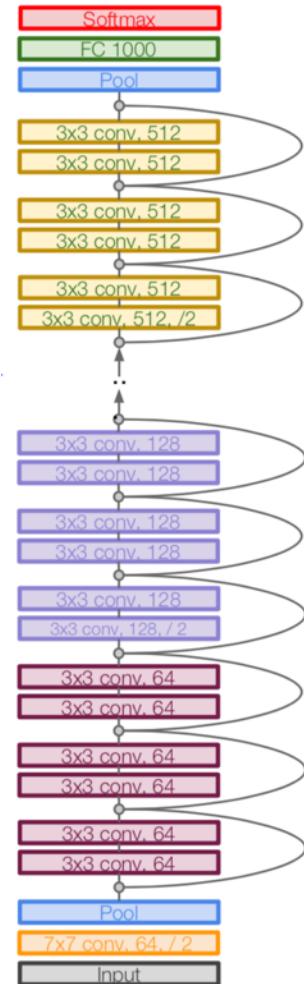
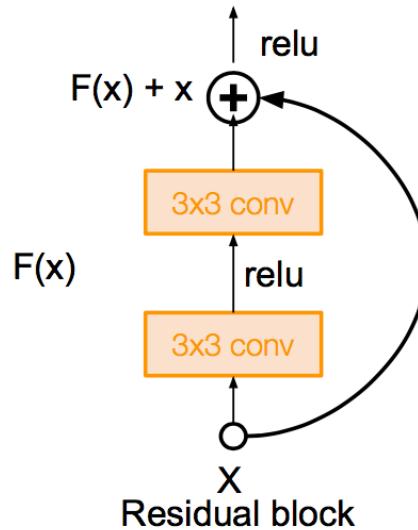
# CNN Architecture: Part I

## □ ResNet [He et al., 2015]

Very deep networks using residual connections

Basic design:

- all 3x3 conv (almost)
- spatial size /2 => #filter x2
- just deep
- average pooling (remove FC)
- no dropout

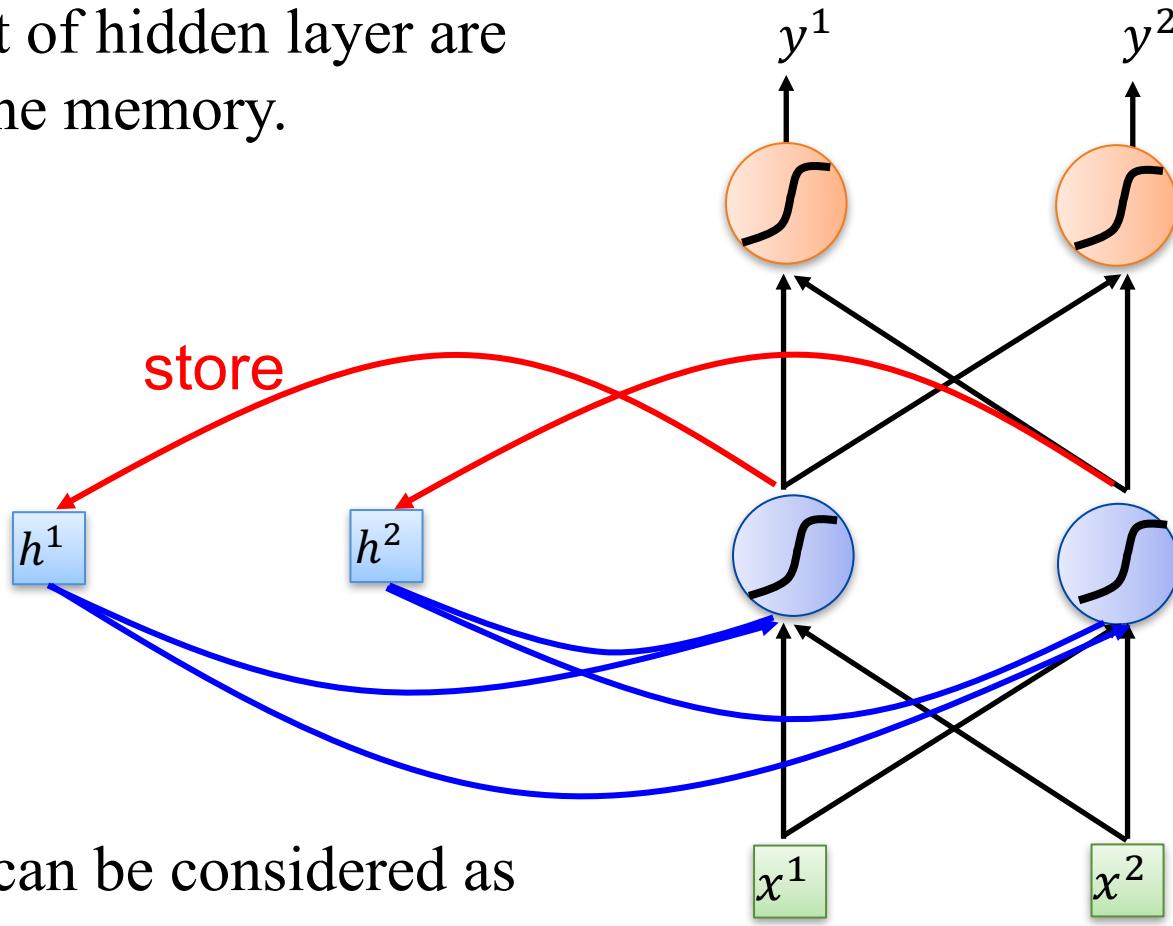


ResNet

# **Recurrent Neural Network and LSTM**

# Recurrent Neural Network

The output of hidden layer are stored in the memory.



Memory can be considered as another input.

# Recurrent Neural Network

□ Formally

□  $x_t$  is the input

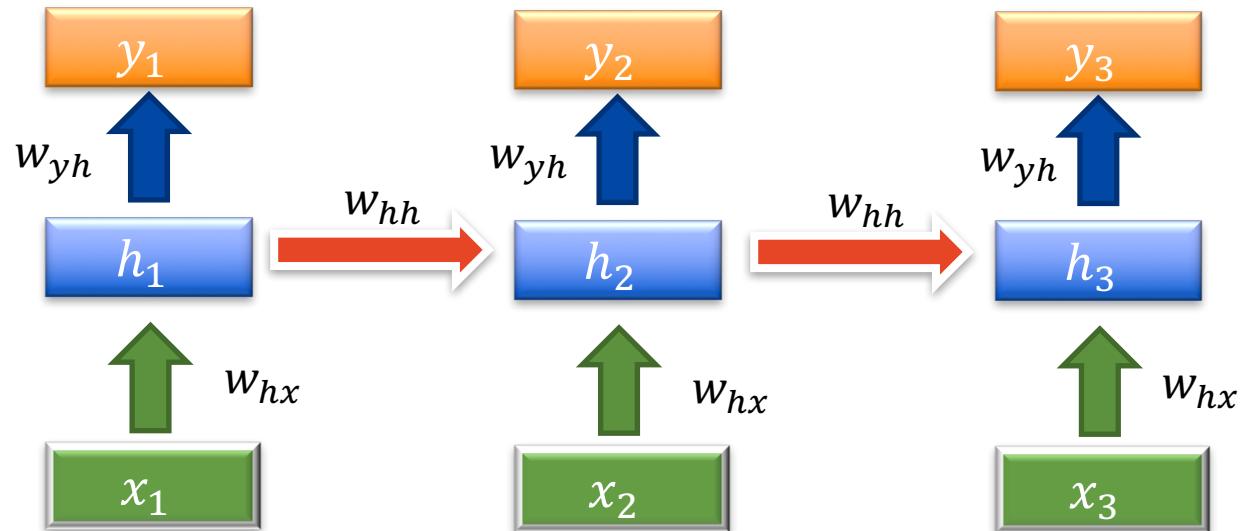
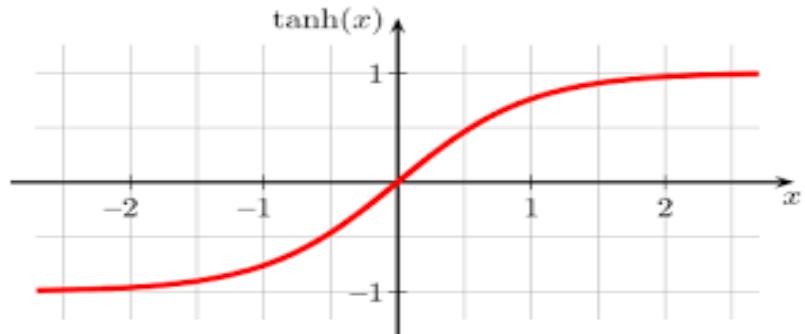
□  $h_t$  is the hidden state

□  $y_t$  is the output

□  $h_0 = \vec{0}$

□  $h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$

□  $y_t = w_{yh}h_t$



# Recurrent Neural Network

## □ Vanishing gradient problem

$$\square h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$$

$$\square y_t = w_{yh}h_t$$

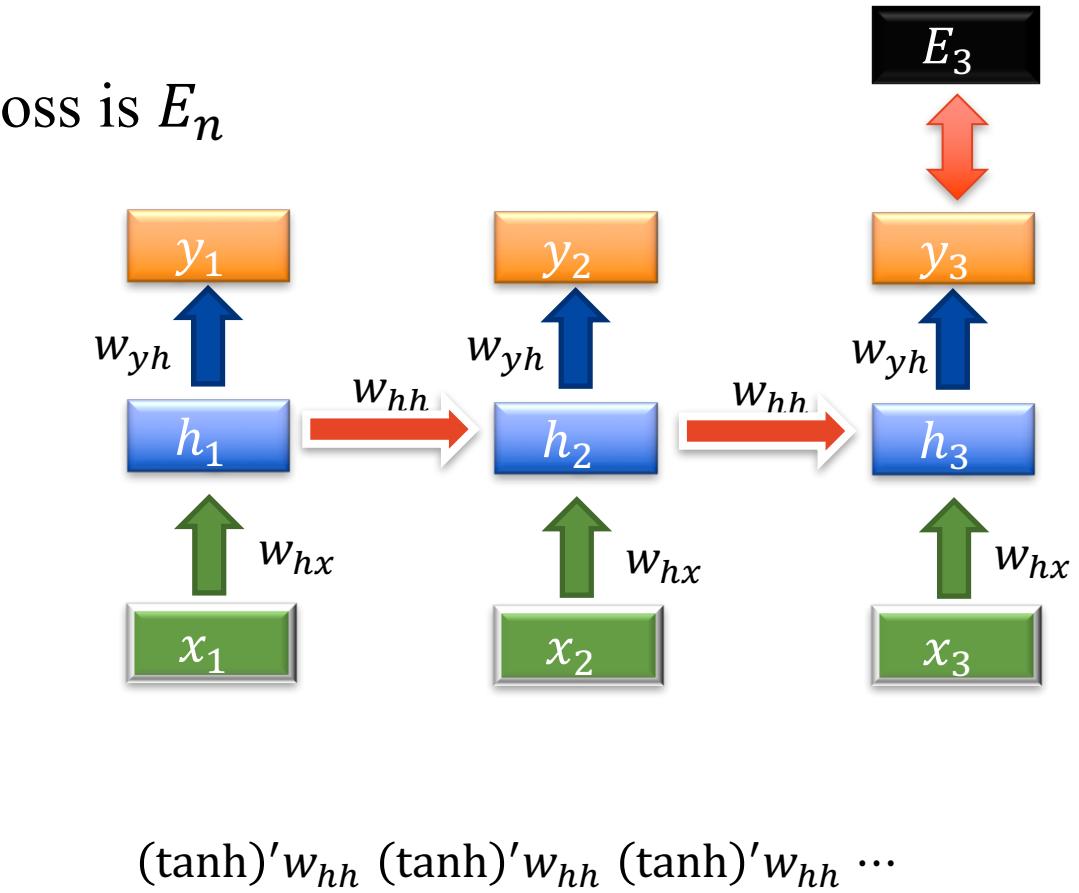
□ For every step, its loss is  $E_n$

$$\frac{\partial E_3}{\partial w_{hh}} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial w_{hh}}$$

$$h_3 = \tanh(w_{hh}h_2 + w_{hx}x_3 + b_h)$$

$$\frac{\partial h_3}{\partial w_{hh}} = (\tanh)' \left( h_2 + w_{hh} \frac{\partial h_2}{\partial w_{hh}} \right)$$

$$\frac{\partial h_2}{\partial w_{hh}} = (\tanh)' \left( h_1 + w_{hh} \frac{\partial h_1}{\partial w_{hh}} \right)$$



# Recurrent Neural Network

## □ Vanishing gradient problem

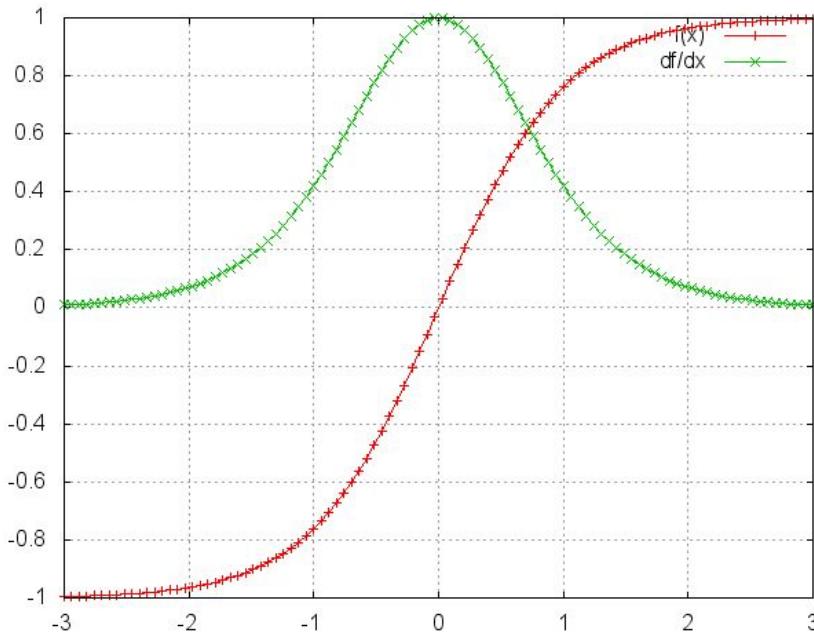
$$\square h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$$

$$\square y_t = w_{yh}h_t$$

□ For every step, its loss is  $E_n$

$$\frac{\partial E_3}{\partial w_{hh}} = \frac{\partial E_3}{\partial y_3} \frac{\partial y_3}{\partial h_3} \frac{\partial h_3}{\partial w_{hh}}$$

$$(\tanh)'w_{hh}(\tanh)'w_{hh}(\tanh)'w_{hh} \dots$$



$$(\tanh)' \leq 1$$

If  $0 < w_{hh} < 1$ , **Vanishing Gradients**

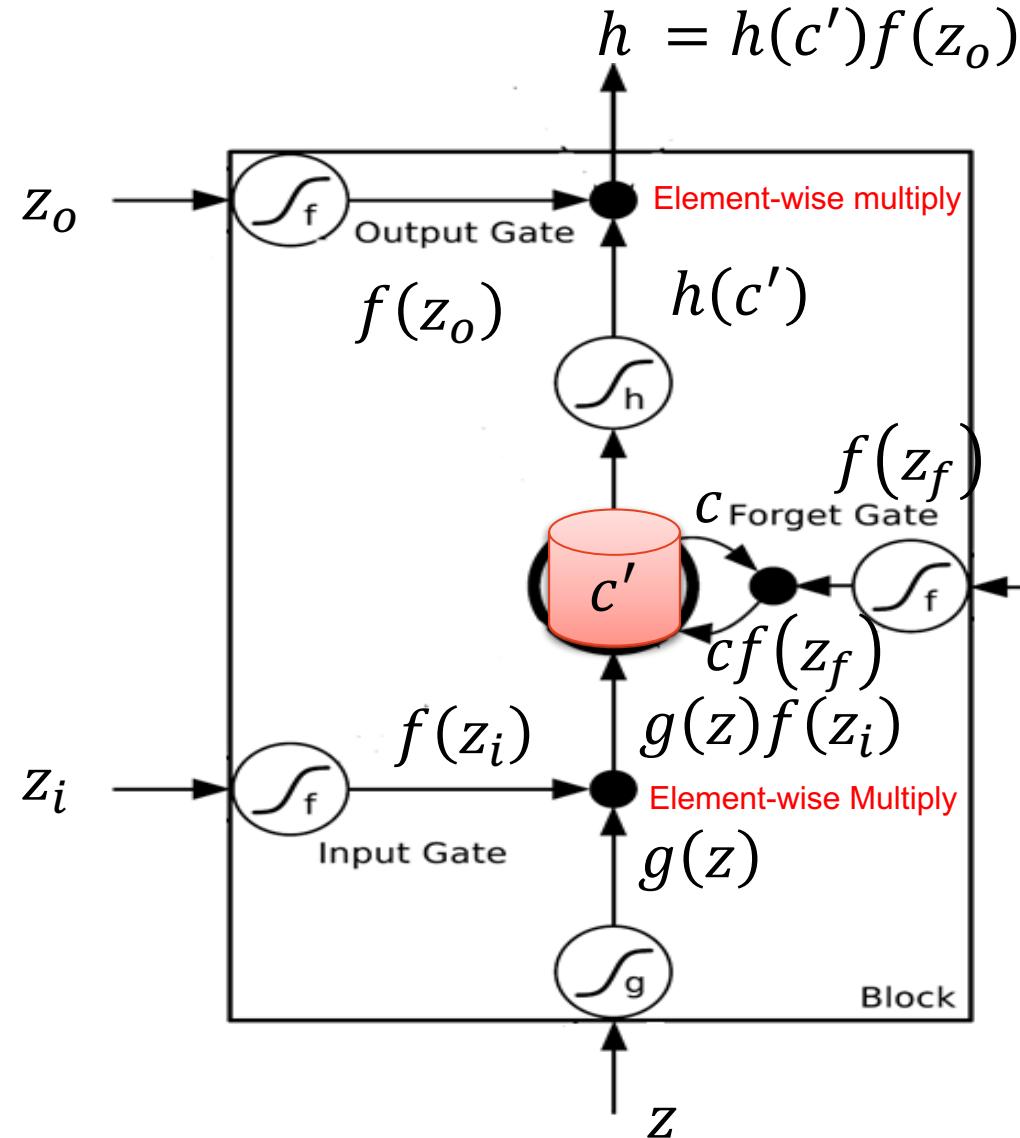
$$(\tanh)'w_{hh}(\tanh)'w_{hh}(\tanh)'w_{hh} \dots \rightarrow 0$$

If  $w_{hh}$  is larger, **Exploding Gradients**

$$(\tanh)'w_{hh}(\tanh)'w_{hh}(\tanh)'w_{hh} \dots \rightarrow \infty$$

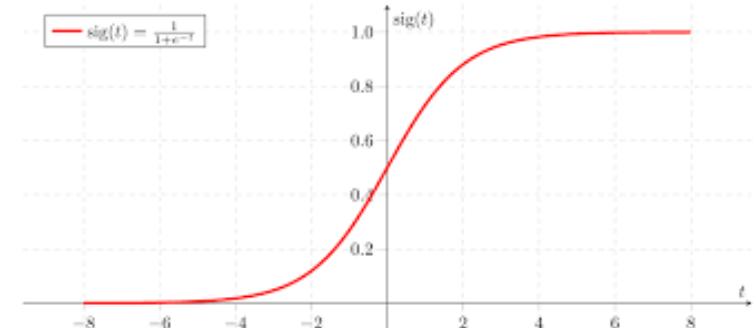
$$\begin{aligned}0.9^{1000} &\approx 0 \\1.01^{1000} &\approx 21000\end{aligned}$$

# Long Short-Term Memory



$f$  is usually a sigmoid function  
Value From 0 to 1  
Mimic open and close gate

$$c' = g(z)f(z_i) + cf(z_f)$$



# Long Short-Term Memory

## □ Forget gate

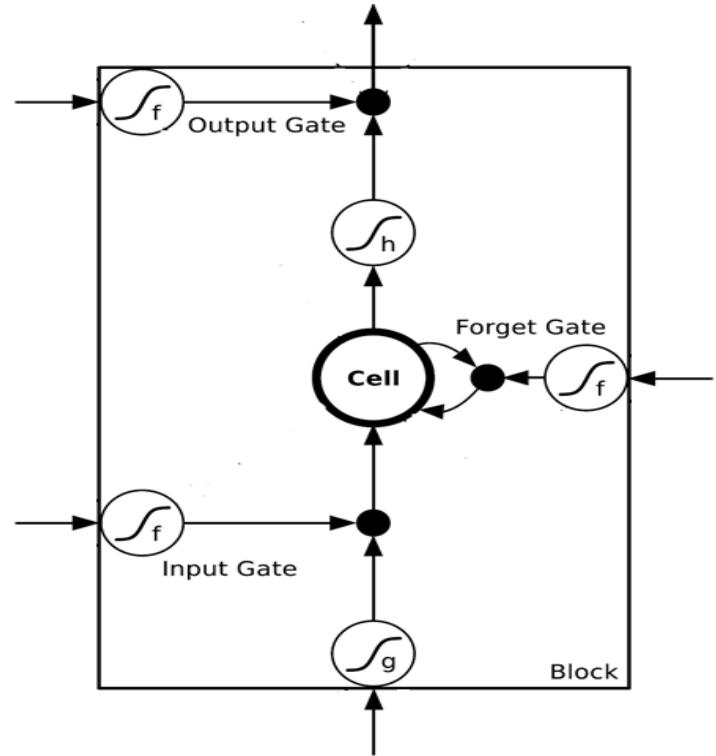
- $f_t = \sigma(w_{fh}h_{t-1} + w_{fx}x_t + b_f)$

## □ Input gate

- $i_t = \sigma(w_{ih}h_{t-1} + w_{ix}x_t + b_i)$

## □ Output gate

- $o_t = \sigma(w_{oh}h_{t-1} + w_{ox}x_t + b_o)$



# Long Short-Term Memory

## □ Candidate cell state

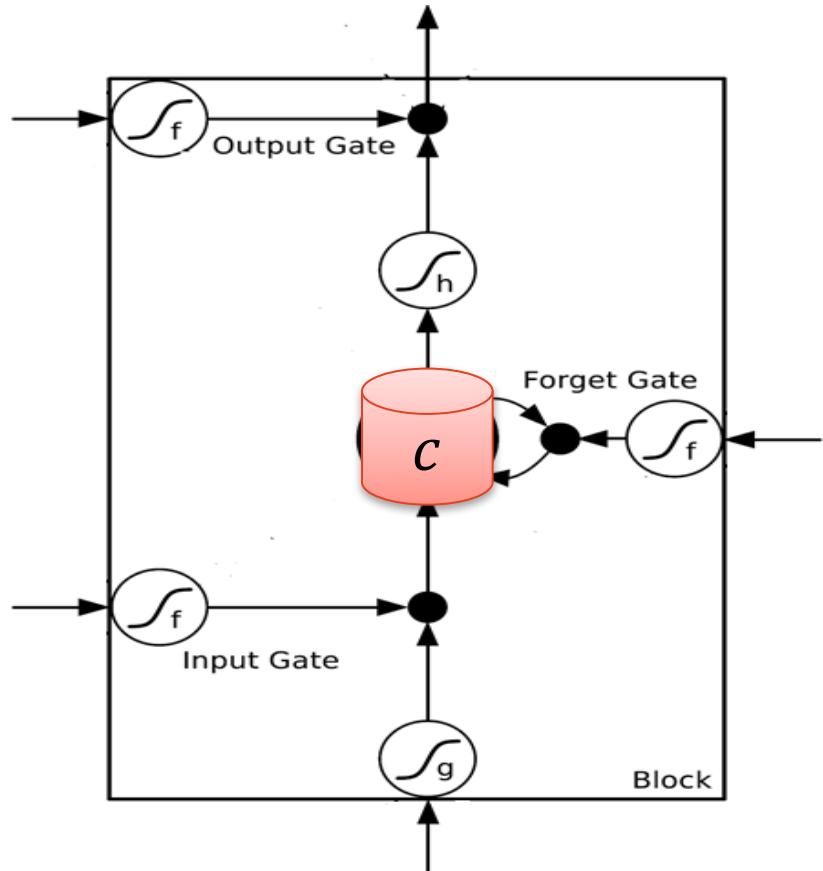
- $\tilde{c}_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$

## □ Update cell state

- $c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$

## □ Output state

- $h_t = o_t * \tanh(c_t)$



# Prevent vanishing gradient

## □ RNN

$$\bullet h_t = \tanh(w_{hh}h_{t-1} + w_{hx}x_t + b_h)$$

## □ LSTM

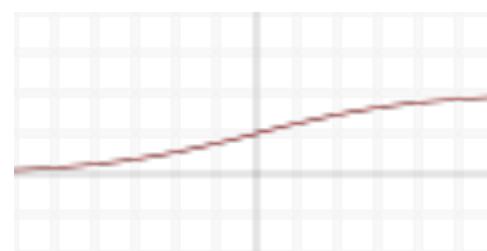
$$\bullet c_t = f_t * c_{t-1} + i_t * \tilde{c}_t$$

Back propagation from  $c_t$  to  $c_{t-1}$  only involves elementwise multiplication with  $f_t$ , no matrix multiplication with  $w_{hh}$  or tanh.



$$f(x) = \tanh(x)$$

$$f'(x) = 1 - f(x) \leq 1$$

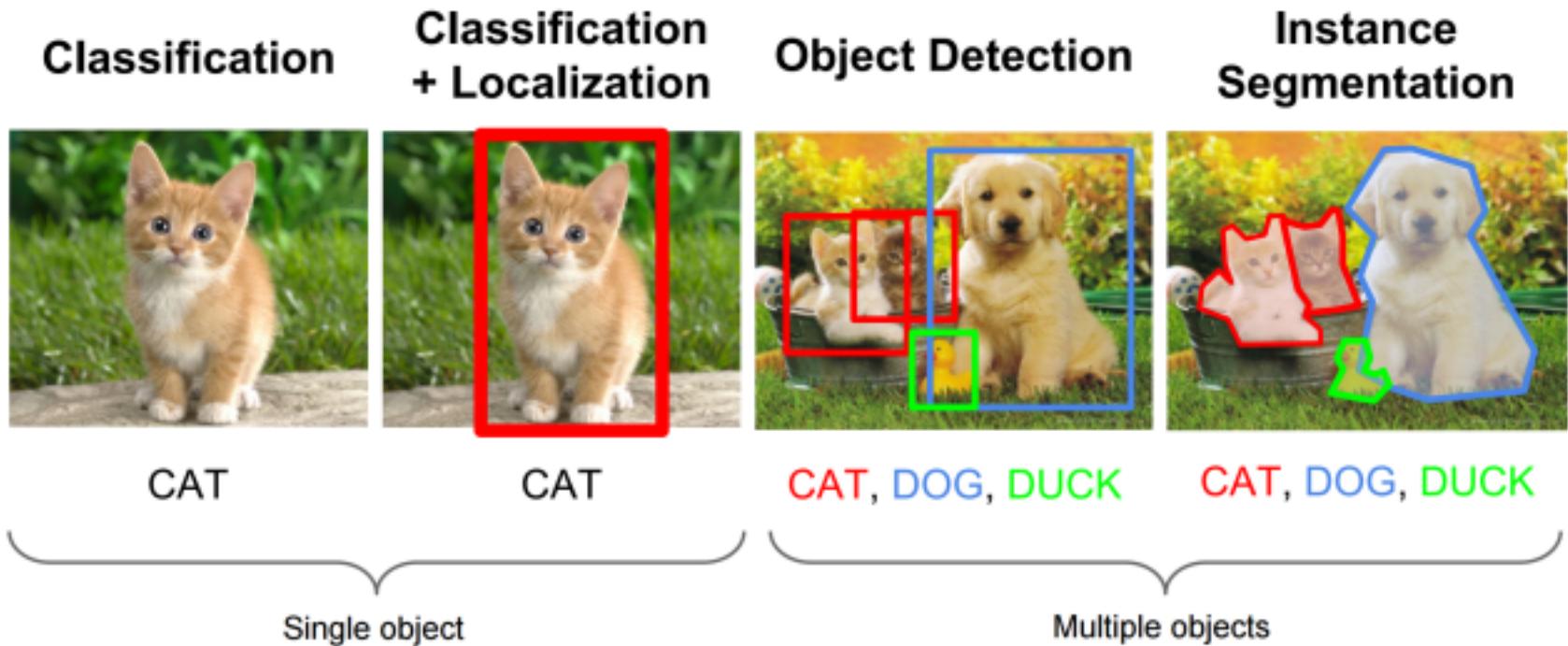


$$f(x) = \sigma(x)$$

$$\begin{aligned} f'(x) \\ = f(x)(1 - f(x)) \\ \leq 0.25 \end{aligned}$$

# **Deep Learning Applications**

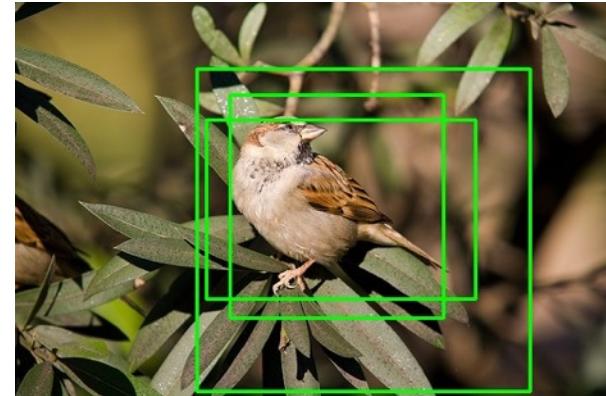
# Computer Vision Tasks



# Typical architecture



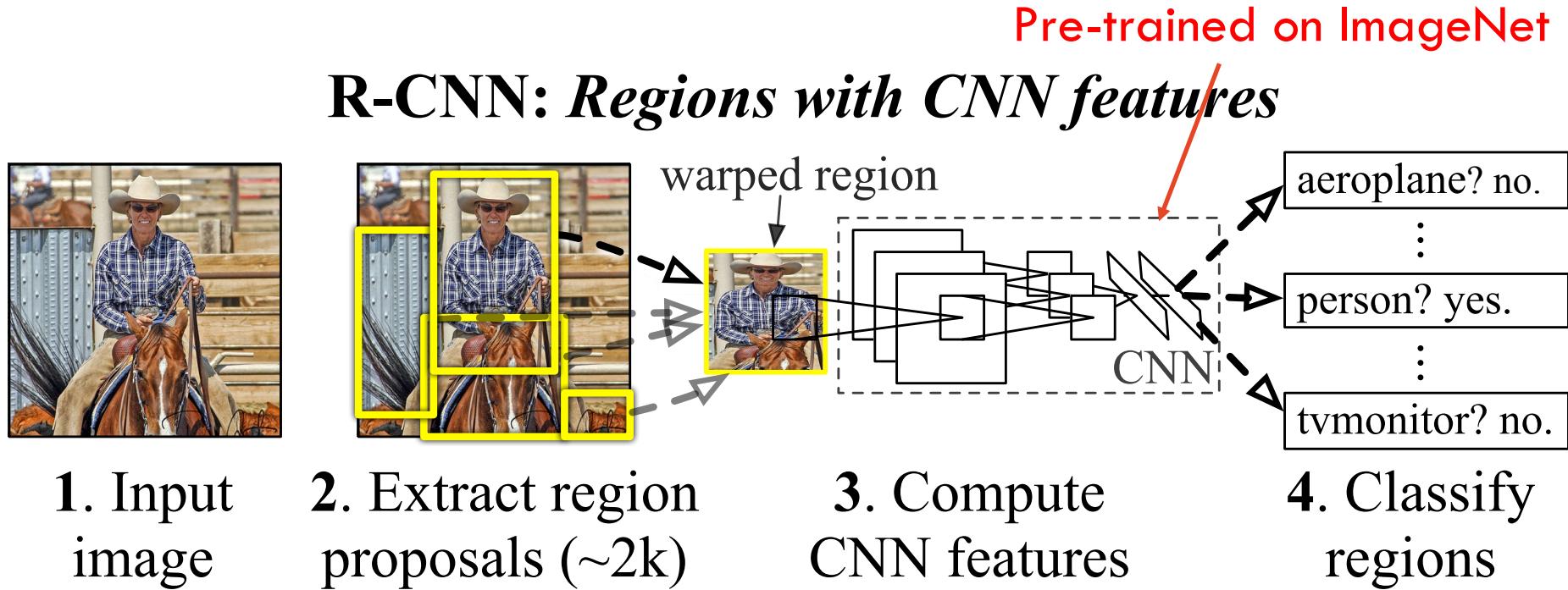
Proposals



**1. Region proposal:** Given an input image find all possible places where objects can be located. The output of this stage should be a list of bounding boxes of likely positions of objects. These are often called region proposals or regions of interest.

**2. Final classification:** for every region proposal from the previous stage, decide whether it belongs to one of the target classes or to the background. Here we could use a deep convolutional network.

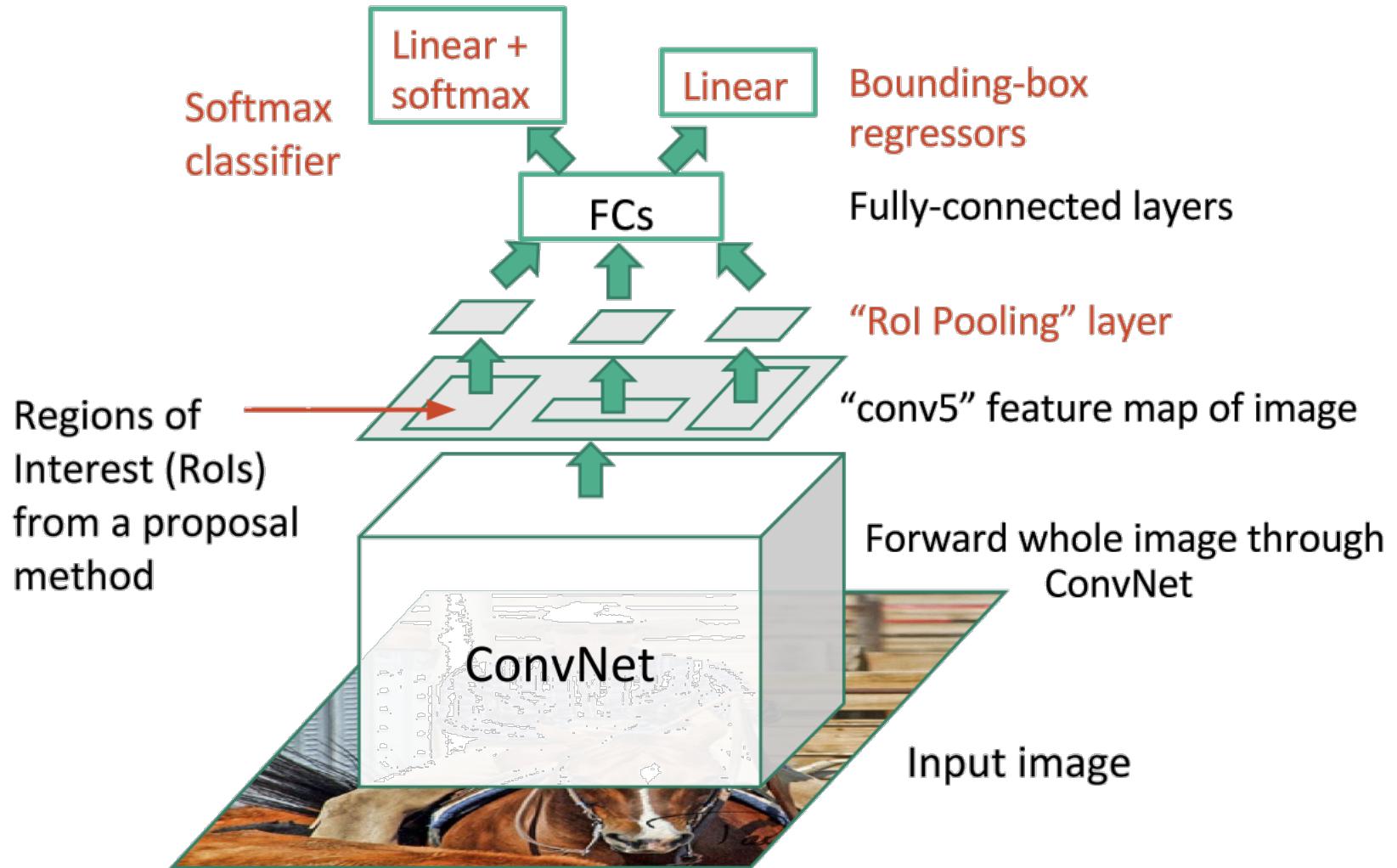
# Region CNN (R-CNN)



## Object Detection to Image Classification

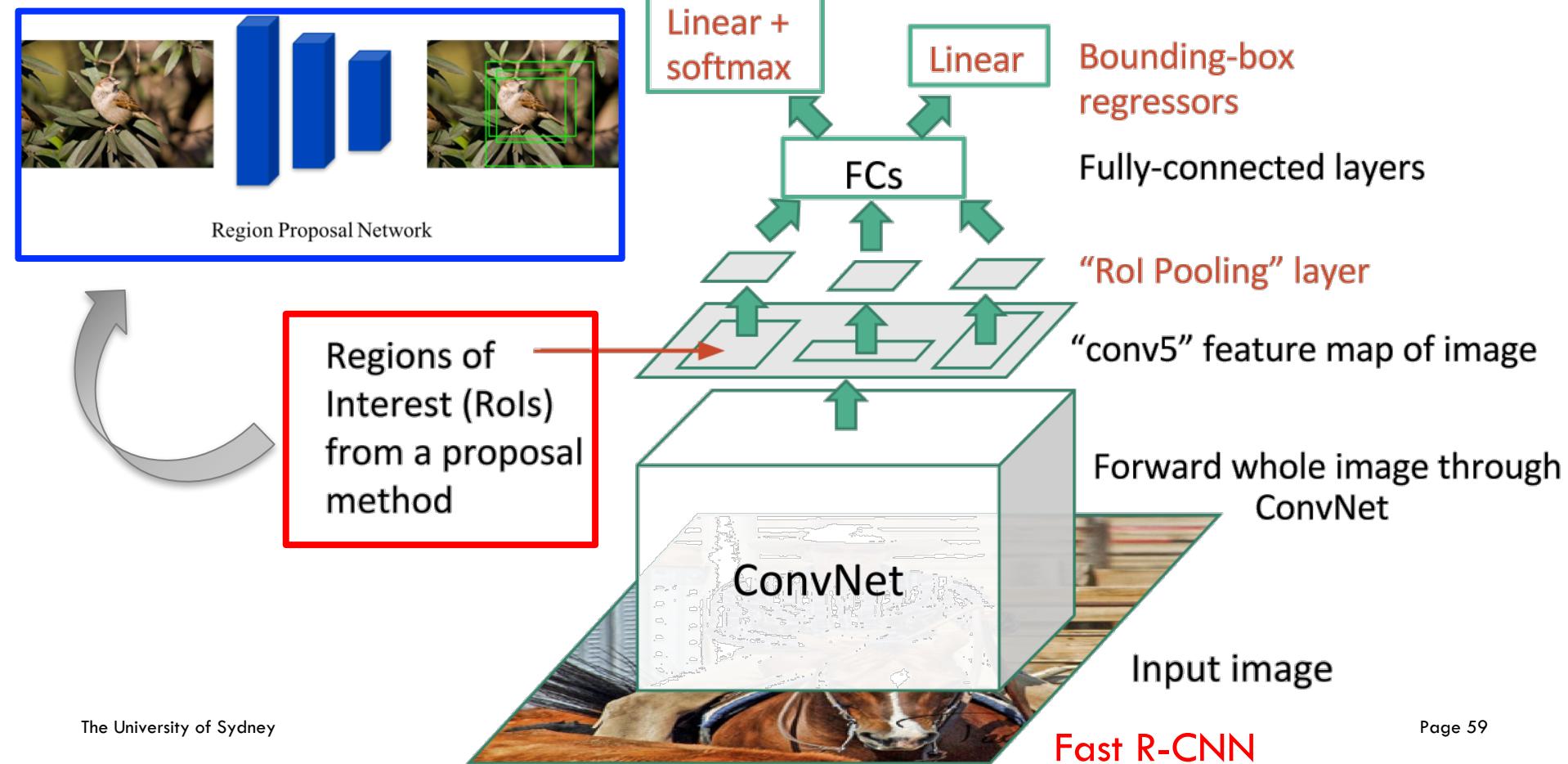
[Girshick14] R. Girshick, J. Donahue, S. Guadarrama, T. Darrell, J. Malik: **Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation**, CVPR 2014

# Fast R-CNN

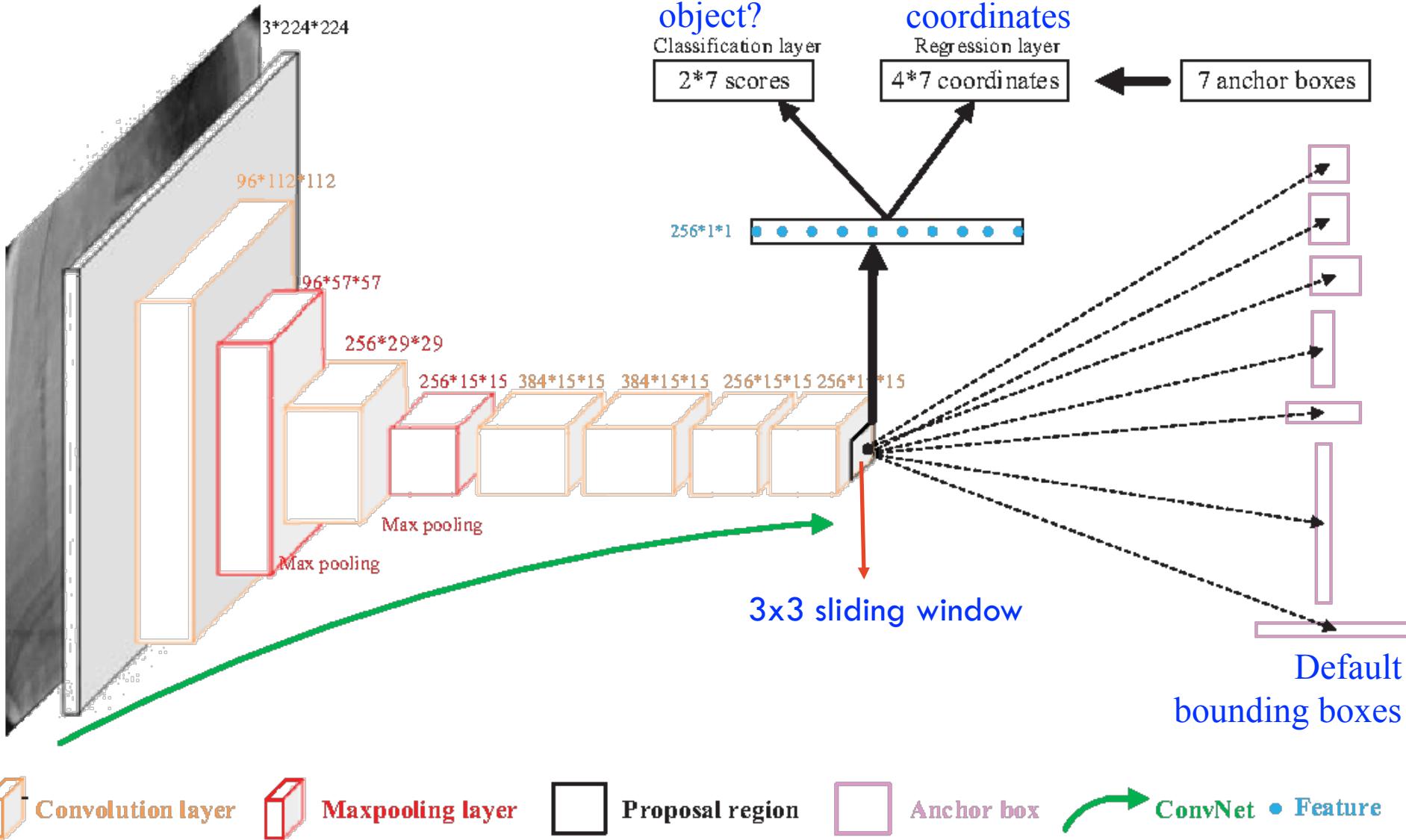


# Faster R-CNN

Replace the slow selective search algorithm with a fast neural net - region proposal network (RPN).



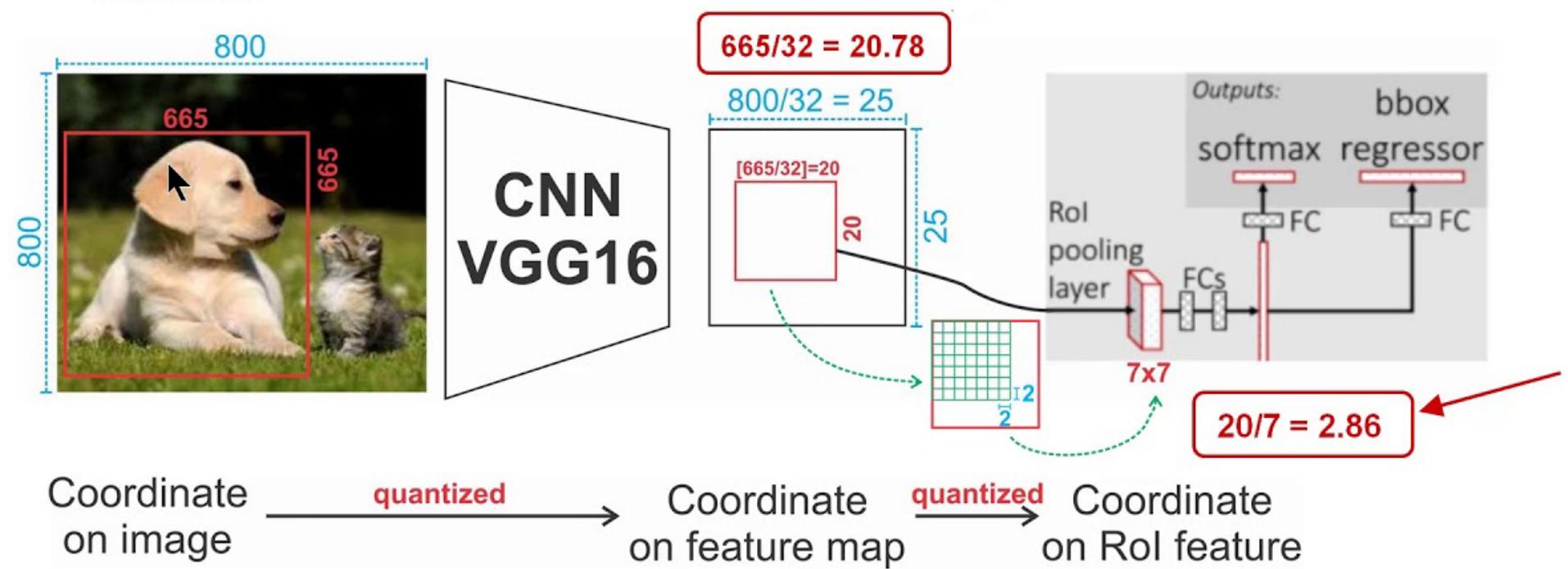
# Region Proposal Network



# RoIAlign

Realigning RoI Pooling to be More Accurate.

“RoiPool”



# Mask R-CNN



Faster R-CNN

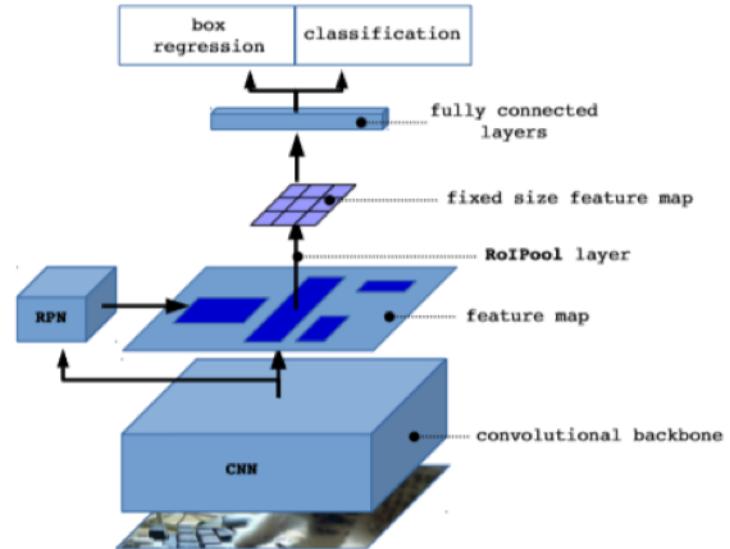


Figure 1. The Mask R-CNN framework for instance segmentation.

**Mask R-CNN**  
= Instance Segmentation+ Faster R-CNN.

Credit To: <https://www.slideshare.net/windmdk/mask-rcnn>

# **Deep Generative Models**

# Generative Modeling

- Density Estimation



- Sample Generation



# Generative Adversarial Networks

- A **counterfeiter-police game** between two components: a generator **G** and a discriminator **D**
- **G**: counterfeiter, trying to fool police with fake currency
- **D**: policy, trying to detect the counterfeit currency
- Competition drives both to improve, until counterfeits are *indistinguishable* from genuine currency



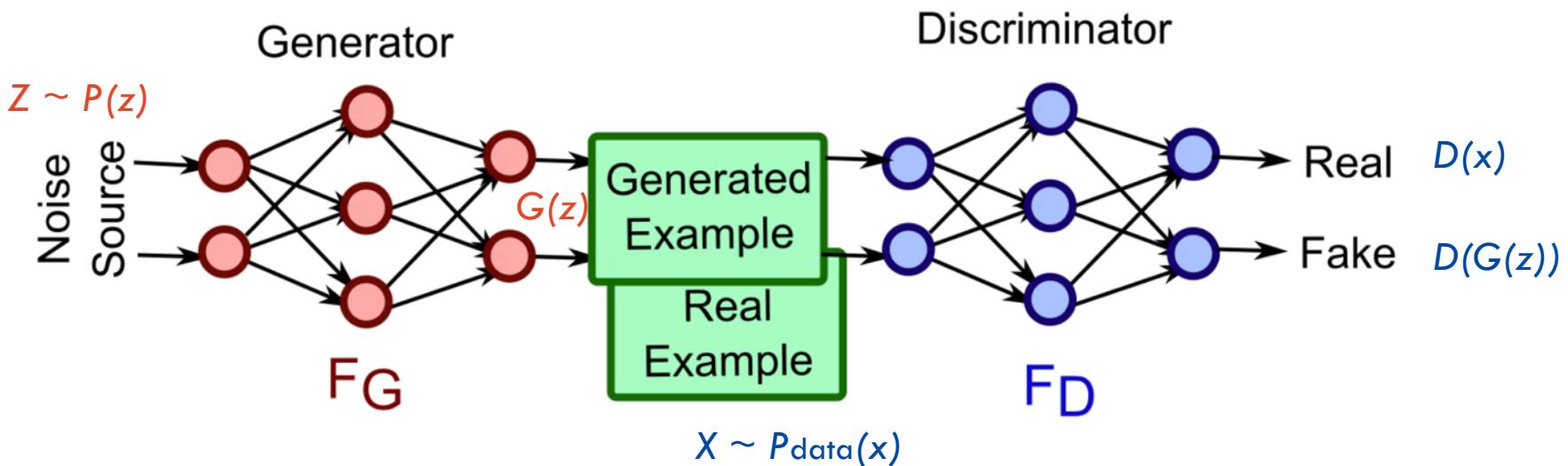
# Generative Adversarial Networks

- A **two-player game** between two components: a generator **G** and a discriminator **D**

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbf{x})} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_{\mathbf{z}}(\mathbf{z})} [\log(1 - D(G(\mathbf{z})))]$$

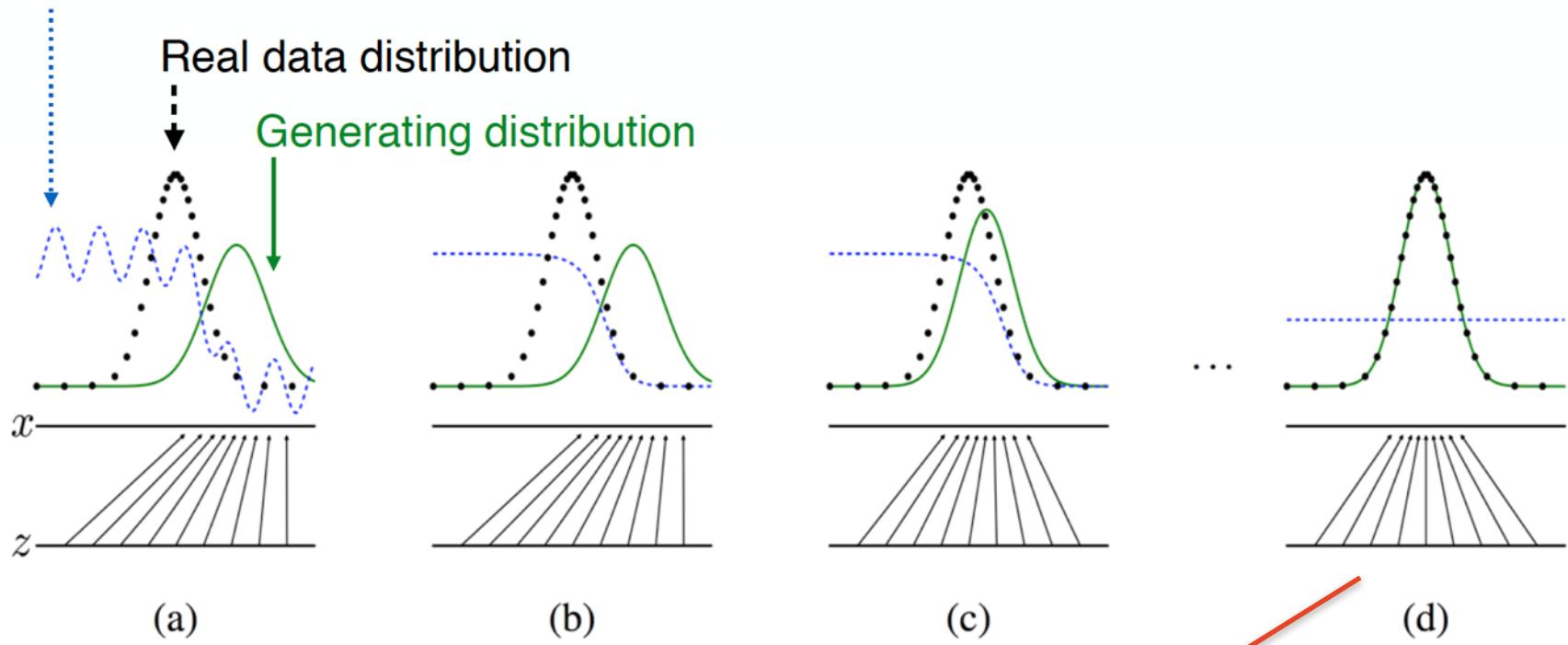
D predicting that  
real data is genuine

D predicting that  
G's generated data is fake



# A simple example

Discriminative distribution

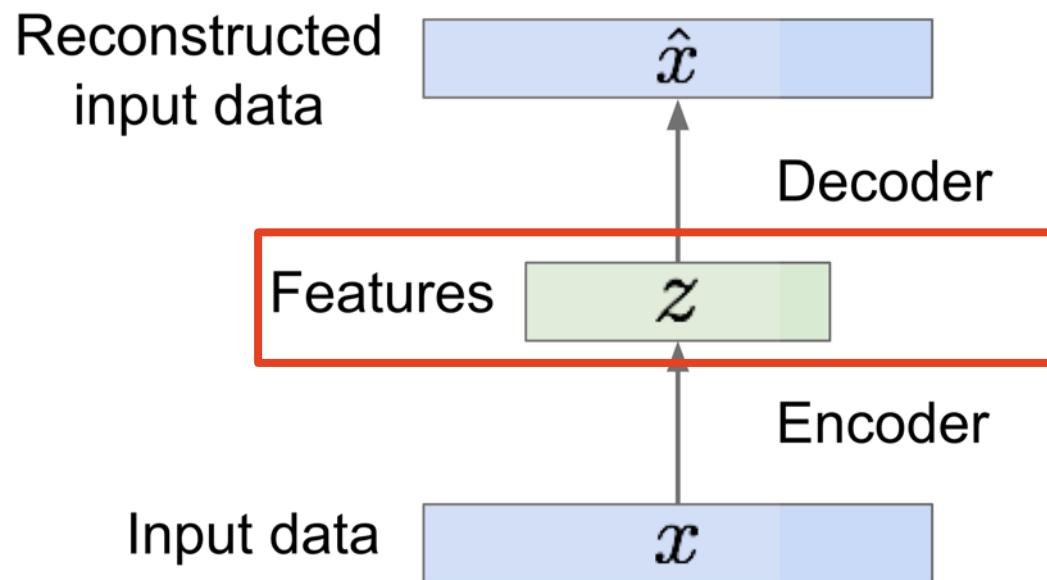


After several steps of training, if G and D have enough capacity, they will reach a point at which both cannot improve because  $p_g = p_{data}$ . The discriminator is unable to differentiate between the two distributions, i.e.  $D(x) = 1/2$ .

# Recap: Autoencoder

Autoencoders can reconstruct data, and can learn features to initialize a supervised model. Features capture factors of variation in training data.

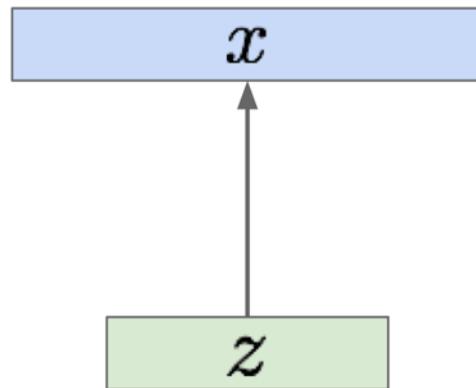
Can we generate new images from an autoencoder?



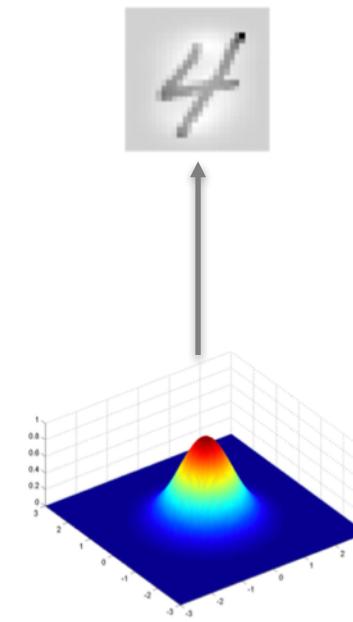
# Variational Autoencoders

- Probabilistic spin on autoencoders - will let us sample from the model to generate data!
- Assume training data  $\{x^{(i)}\}_{i=1}^N$  is generated from underlying unobserved (latent) representation  $z$

Sample from  
true conditional  
 $p_{\theta^*}(x | z^{(i)})$

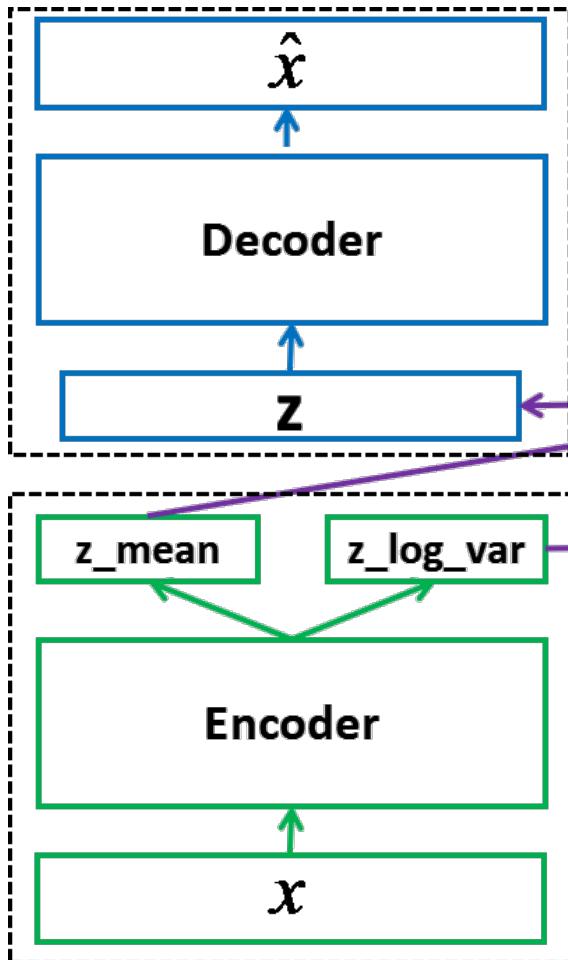


Sample from  
true prior  $p_{\theta^*}(z)$



Choose prior  $p(z)$  to  
be simple, e.g.  
Gaussian.  
Reasonable for  
latent attributes, e.g.  
pose.

# Variational Autoencoders



4

Cross Entropy

$$\sum_i^n -[x_i \log(\hat{x}_i) + (1 - x_i) \log(1 - \hat{x}_i)]$$

3

$$\text{MSE} \sum_i^n (x_i - \hat{x}_i)^2$$

KL Divergence

$$-0.5(1 + \log \sigma^2 - \mu^2 - \sigma^2)$$

1

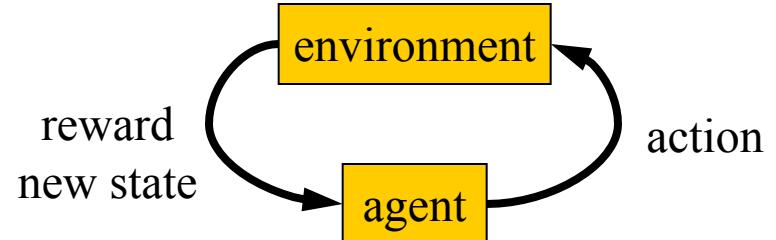
Min CrossEntropy + KLDivergence

h

# **Deep Reinforcement Learning**

# Markov Decision Process (MDP)

- set of states  $S$ , set of actions  $A$ , initial state  $S_0$
- transition model  $P(s,a,s')$ 
  - $P(\text{frame}_{(t)}, \text{right}, \text{frame}_{(t')}) = 0.8$
- reward function  $r(s)$ 
  - $r(\text{frame}_{(t)}) = +1$
- goal: maximize cumulative reward in the long run
- policy: mapping from  $S$  to  $A$ 
  - $a=\pi(s)$  or  $\pi(s, a)$  (deterministic vs. stochastic)
- reinforcement learning
  - transitions and rewards usually not available
  - how to change the policy based on experience
  - how to explore the environment



# Computing return from rewards

- episodic (vs. continuing) tasks
  - “game over” after N steps
  - optimal policy depends on N; harder to analyze
- additive rewards
  - $V(s_0, s_1, \dots) = r(s_0) + r(s_1) + r(s_2) + \dots$
  - infinite value for continuing tasks
- discounted rewards
  - $V(s_0, s_1, \dots) = r(s_0) + \gamma * r(s_1) + \gamma^2 * r(s_2) + \dots$
  - value bounded if rewards bounded

The goal of RL is to find the policy which maximizes the expected return.

# Value Function

- A value function is the prediction of the future return
- Two definitions exist for the value function
  - State value function
  - “How much reward will I get from state s?”
  - expected return when starting in s and following  $\pi$

$$V^\pi(s) = \mathbb{E} \left\{ \sum_{k=0}^{\infty} \gamma^k r_k \middle| s_0 = s, \pi \right\}$$



$$V^\pi(s) = \mathbb{E} \{ Q^\pi(s, a) | a \sim \pi(s, \cdot) \}$$

- State-action value function
- “How much reward will I get from action a in state s?”
- expected return when starting in s, performing a, and following  $\pi$

$$Q^\pi(s, a) = \mathbb{E} \left\{ \sum_{k=0}^{\infty} \gamma^k r_k \middle| s_0 = s, a_0 = a, \pi \right\}$$

# Bellman Equation and Optimality

- Value functions decompose into **Bellman equations**, i.e. the value functions can be decomposed into *immediate reward plus discounted value of successor state*

$$V^\pi(s) = \mathbb{E} \left\{ \sum_{k=0}^{\infty} \gamma^k r_k \middle| s_0 = s, \pi \right\} \quad \longrightarrow \quad V^\pi(s) = \mathbb{E} \{ R(s, a, s') + \gamma V^\pi(s') \}$$

$$Q^\pi(s, a) = \mathbb{E} \left\{ \sum_{k=0}^{\infty} \gamma^k r_k \middle| s_0 = s, a_0 = a, \pi \right\}$$

↓

$$Q^\pi(s, a) = \mathbb{E} \{ R(s, a, s') + \gamma Q^\pi(s', a') \}$$

- An **optimal** value function is the maximum achievable value.

$$V^*(s) = \max_{\pi} V^\pi(s) \quad Q^*(s, a) = \max_{\pi} Q^\pi(s, a)$$

# Bellman Optimality Equation

- Optimality for value functions is governed by the **Bellman optimality equations**.
- Two equations:

$$V^*(s) = \max_a \mathbb{E} \{ R(s, a, s') + \gamma V^*(s') \}$$

$$Q^*(s, a) = \mathbb{E} \left\{ R(s, a, s') + \max_{a'} \gamma Q^*(s', a') \right\}$$

$$\pi^*(s) = \arg \max_a Q^*(s, a)$$

# Q-Learning

$$Q^*(s, a) = \mathbb{E} \left\{ R(s, a, s') + \max_{a'} \gamma Q^*(s', a') \right\}$$

Initialize  $Q(s, a)$  arbitrarily.

Start with  $s$ .

Before taking action  $a$ , we calculate the current expected return as

$$Q(s, a)$$

After taking action  $a$ , and observing  $r$  and  $s'$ , we calculate the target expected return as

$$\overbrace{R(s, a, s')} + \max_{a'} \gamma Q(s', a')$$

$$\Delta Q(s, a) = R(s, a, s') + \max_{a'} \gamma Q(s', a') - Q(s, a)$$

$$Q(s, a) \leftarrow Q(s, a) + \alpha \Delta Q(s, a)$$

# Q-Learning

$$Q^*(s, a) = \mathbb{E} \left\{ R(s, a, s') + \max_{a'} \gamma Q^*(s', a') \right\}$$

Initialize  $Q(s, a)$  arbitrarily.

Repeat (for each episode)

    Initialize  $s$

    Repeat (for each step of the episode)

        Choose  $a$  from  $s$  using a policy

        Take action  $a$ , observe  $r, s'$

$$Q(s, a) \leftarrow Q(s, a) + \alpha[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$$

$s \leftarrow s'$

# Exploration and Exploitation

Initialize  $Q(s, a)$  arbitrarily.

Repeat (for each episode)

    Initialize  $s$

    Repeat (for each step of the episode)

        Choose  $a$  from  $s$  using a policy

        Take action  $a$ , observe  $r, s'$

        Update  $Q$  and  $s$ ....

Random policy; Exploration

Greedy policy; Exploitation

$$\pi(s) = \arg \max_a Q(s, a)$$

$\epsilon$ -greedy policy: With probability  $\epsilon$  select a random action

# Deep Q-Learning



	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
Noop	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Fire	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Right	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
Left	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0

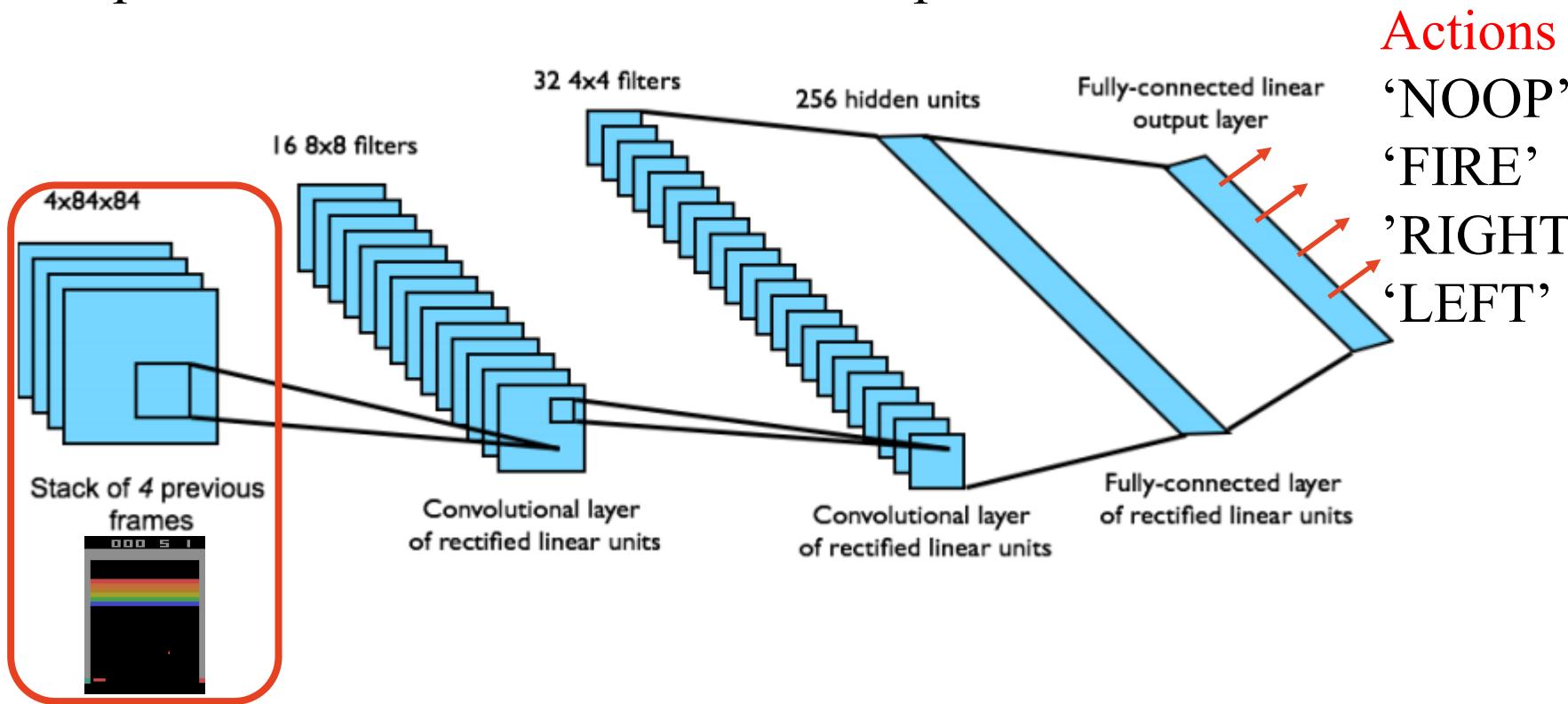
$$\# \text{states} = 256^{210 \times 160}$$

Value Function Approximation

$$Q(s, a) = f(s, a, w)$$

# Deep Q-Learning

- End-to-end learning of state-action values from raw pixels
- Input state is stack of raw pixels from last 4 frames
- Output are state-action values from all possible actions



From the Tutorial: Deep Reinforcement Learning by David Silver, Google DeepMind

# Deep Q-Networks (DQN)

- Optimal Q-values obey Bellman equation

$$Q^*(s, a) = \mathbb{E}_{s'} \left\{ r + \gamma \max_{a'} Q(s', a') | s, a \right\}$$

- Treat right-hand size as a target and minimize MSE loss by SGD

$$l = \frac{\left( r + \gamma \max_{a'} Q_w(s', a') - Q_w(s, a) \right)^2}{\text{Target}}$$

- Divergence issues using neural networks due to
  - Correlations between samples
  - Non-stationary targets

# Algorithm

## Algorithm 1 Deep Q-learning with Experience Replay

Initialize replay memory  $\mathcal{D}$  to capacity  $N$

Initialize action-value function  $Q$  with random weights

**for** episode = 1,  $M$  **do**

    Initialise sequence  $s_1 = \{x_1\}$  and preprocessed sequenced  $\phi_1 = \phi(s_1)$

**for**  $t = 1, T$  **do**

        With probability  $\epsilon$  select a random action  $a_t$        **$\epsilon$ -greedy**  
        otherwise select  $a_t = \max_a Q^*(\phi(s_t), a; \theta)$

        Execute action  $a_t$  in emulator and observe reward  $r_t$  and image  $x_{t+1}$

        Set  $s_{t+1} = s_t, a_t, x_{t+1}$  and preprocess  $\phi_{t+1} = \phi(s_{t+1})$

        Store transition  $(\phi_t, a_t, r_t, \phi_{t+1})$  in  $\mathcal{D}$

    Sample random minibatch of transitions  $(\phi_j, a_j, r_j, \phi_{j+1})$  from  $\mathcal{D}$

    Set  $y_j = \begin{cases} r_j & \text{for terminal } \phi_{j+1} \\ r_j + \gamma \max_{a'} Q(\phi_{j+1}, a'; \theta) & \text{for non-terminal } \phi_{j+1} \end{cases}$

    Perform a gradient descent step on  $(y_j - Q(\phi_j, a_j; \theta))^2$  according to

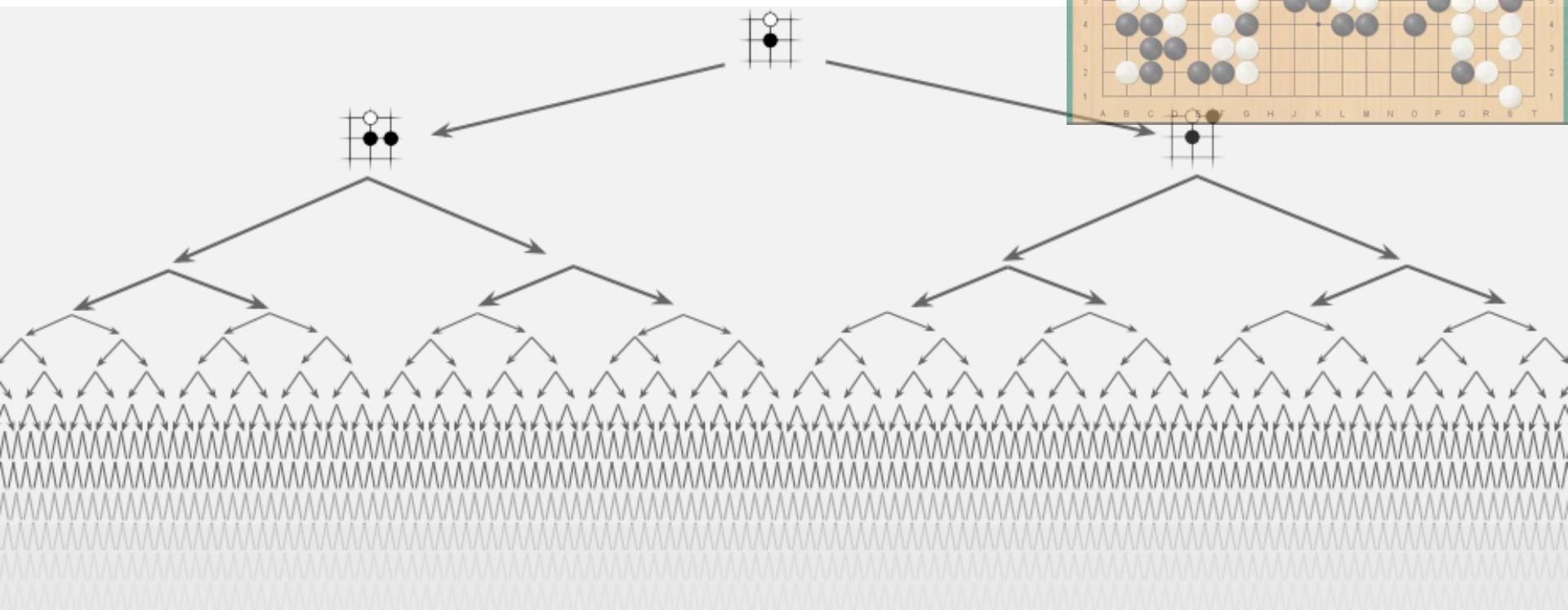
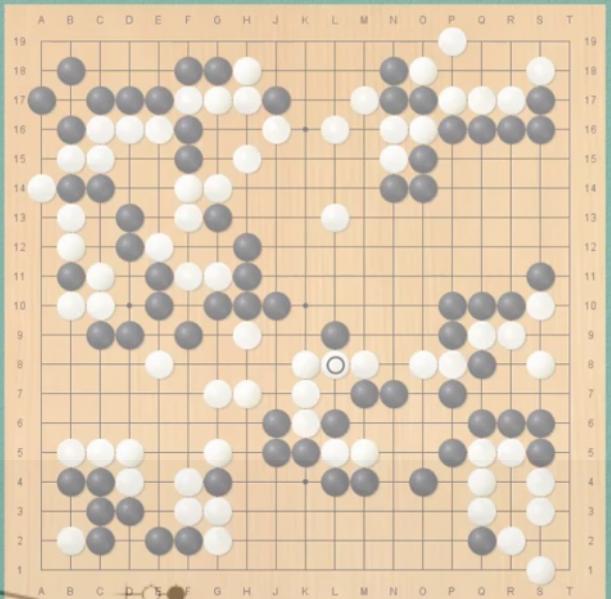
**end for**

**end for**

# Deep RL in Alpha Go

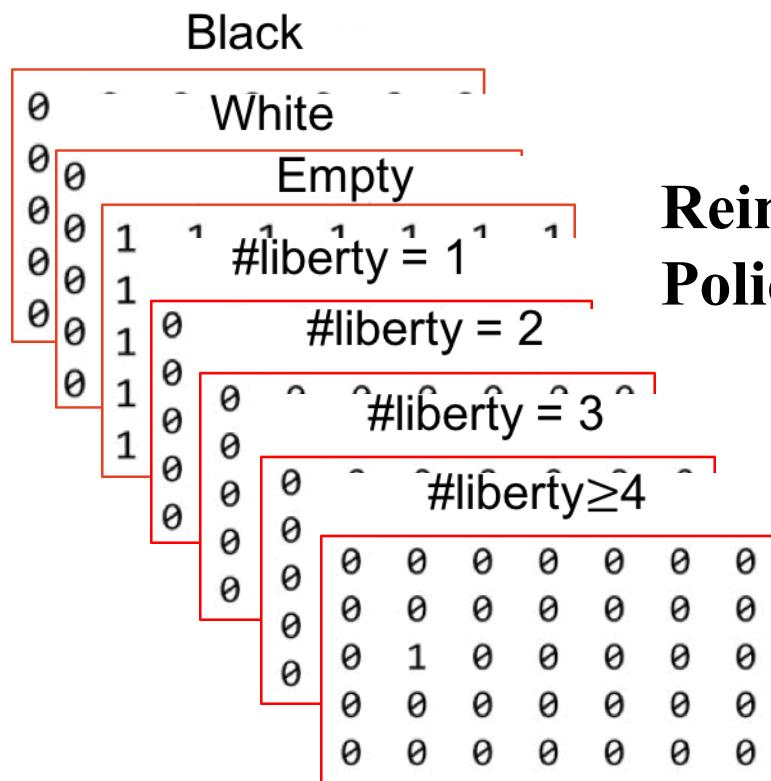
- To win a game, exhaustive search?

Monte Carlo search tree



From <https://icml.cc/2016/tutorials/AlphaGo-tutorial-slides.pdf>

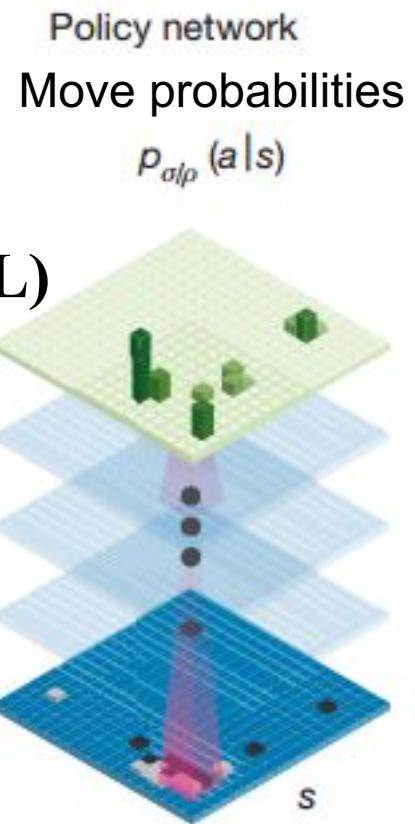
# Policy Network



## Supervised Learning (SL) Policy Network



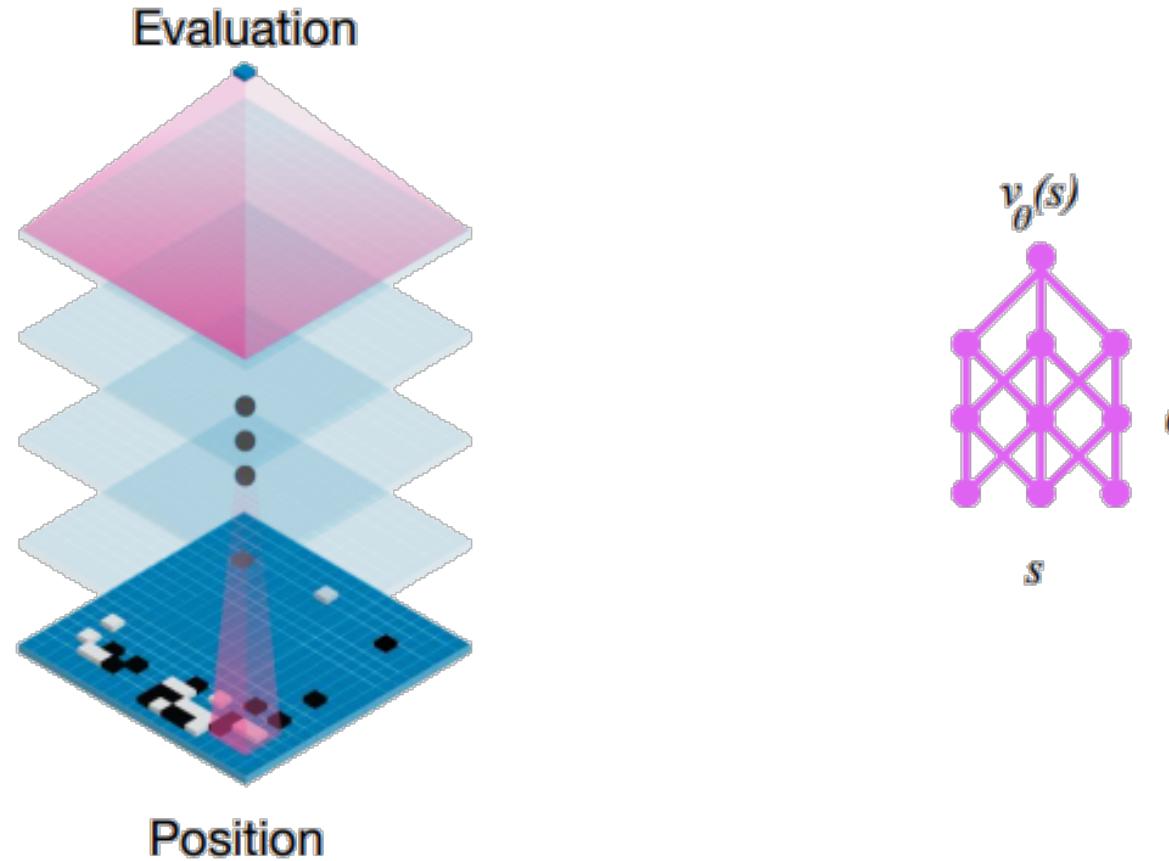
## Reinforcement Learning (RL) Policy Network



Convolutional Neural Network

# Value Network

- To estimating probability that current move leads to win



From <https://icml.cc/2016/tutorials/AlphaGo-tutorial-slides.pdf>

<b>Date</b>	<b>Time</b>	<b>Writing Time (minutes)</b>	<b>Venue</b>	<b>Building</b>	<b>Room</b>
Friday 22/06/2018	1:50pm	120	Teachers College 323A	A22	323A
Friday 22/06/2018	1:50pm	120	VSCC Seminar Room 115	B22	115
Friday 22/06/2018	1:50pm	120	Quadrangle Building Latin 2 S225	A14	S225
Friday 22/06/2018	1:50pm	120	Quadrangle Building Room S241	A14	S241
Friday 22/06/2018	1:50pm	120	Woolley N208	A20	N208

**EXAM WRITING TIME: 2 hours**

- 8 choice questions
  - 2 calculation questions
  - 4 essay questions
- 
- Calculator - non-programmable
  - One A4 sheet of handwritten and/or typed notes double sided

**Thanks!**