

# *Deep Learning: Theory and Practice*

---

**Deep Learning - Practical  
Considerations**

02-04-2020

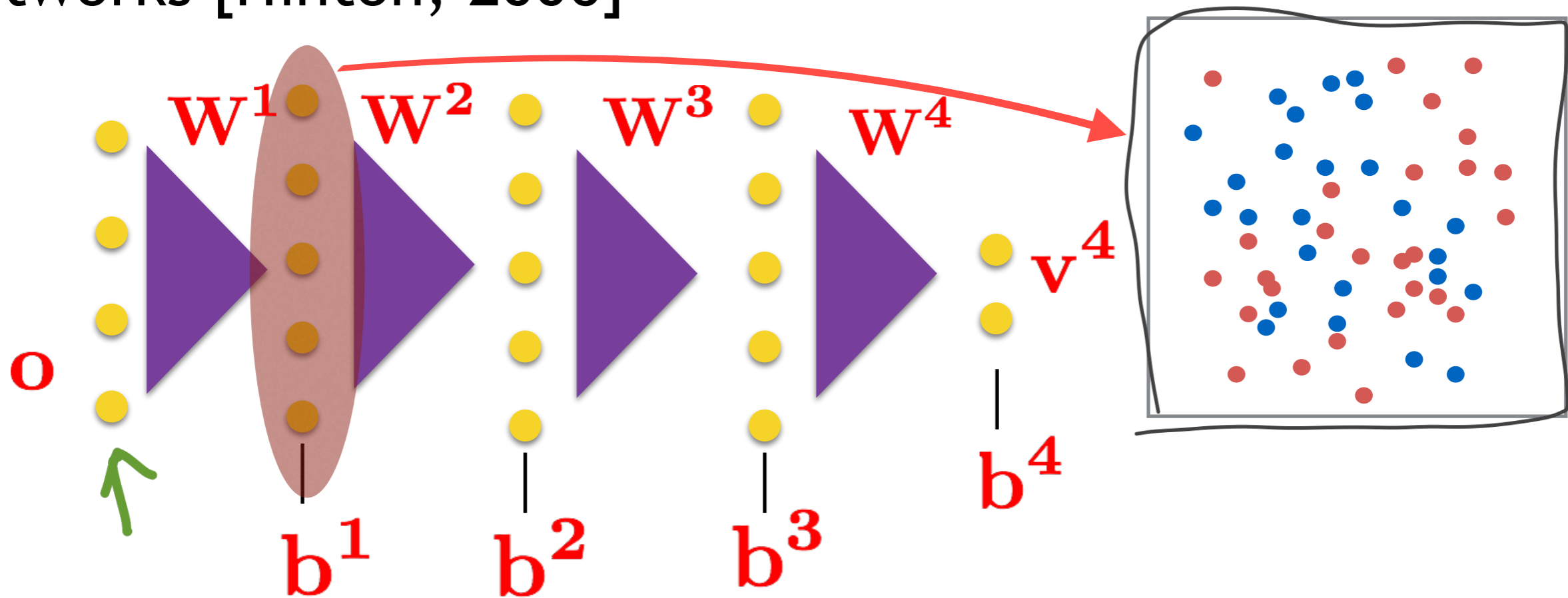
---

*deeplearning.cce2020@gmail.com*



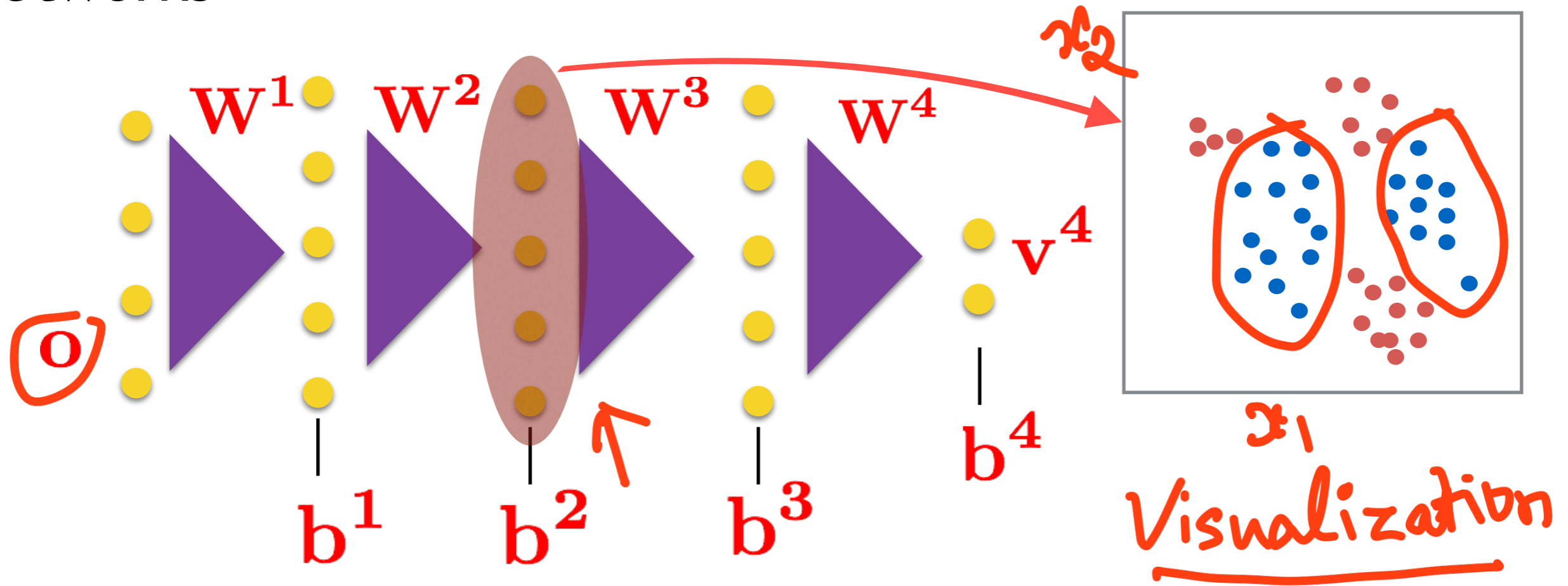
# Deep Networks Intuition

Neural networks with multiple hidden layers - Deep networks [Hinton, 2006]



# Deep Networks Intuition

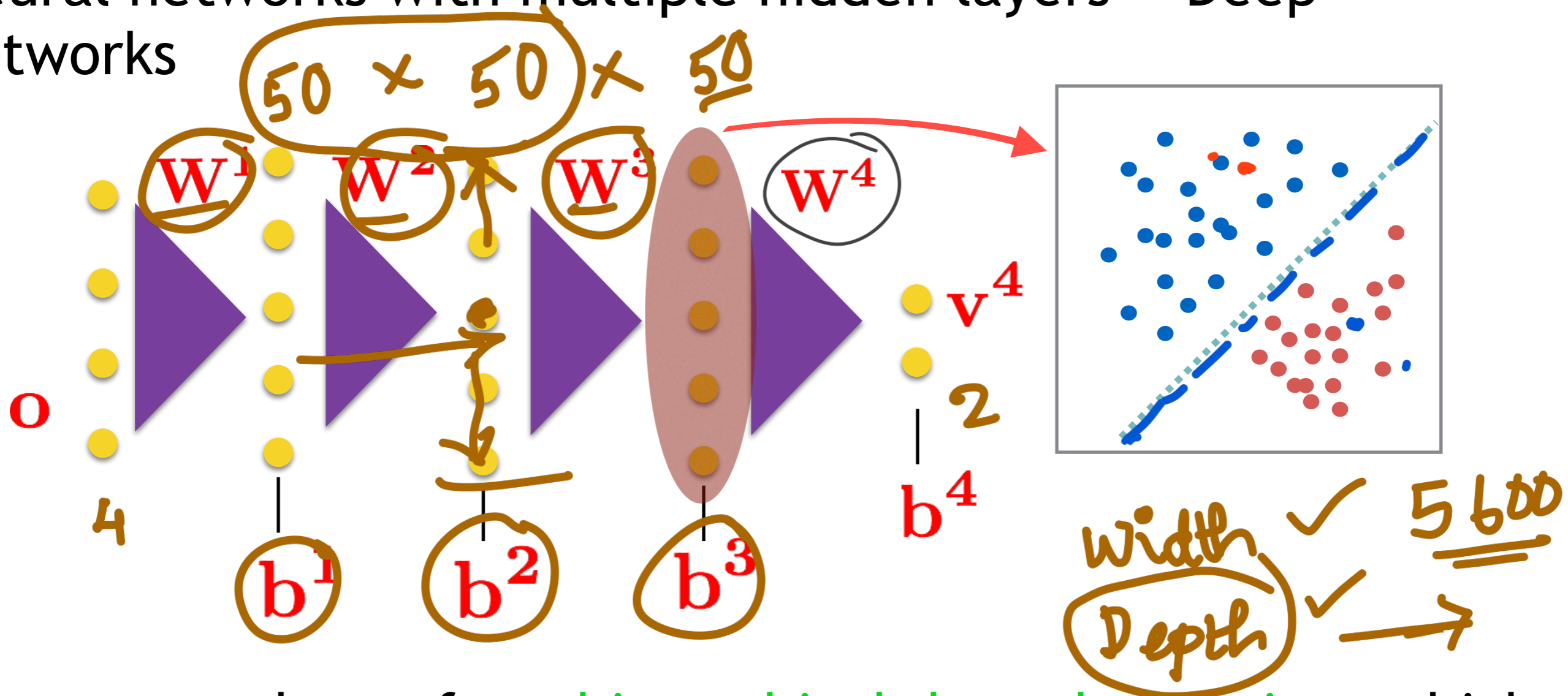
Neural networks with multiple hidden layers - Deep networks





# Deep Networks Intuition

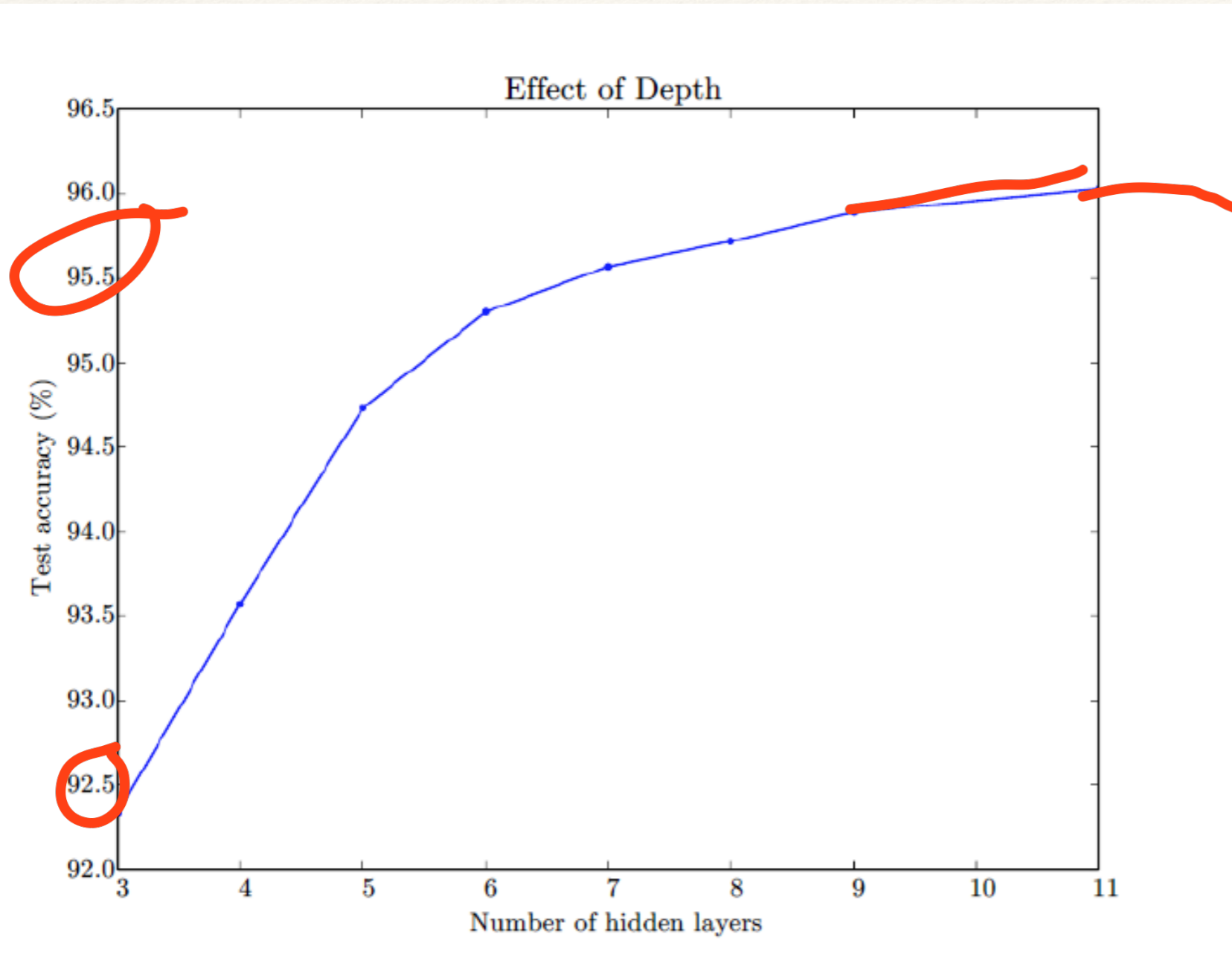
Neural networks with multiple hidden layers - Deep networks



Deep networks perform hierarchical data abstractions which enable the non-linear separation of complex data samples.



# Need for Depth

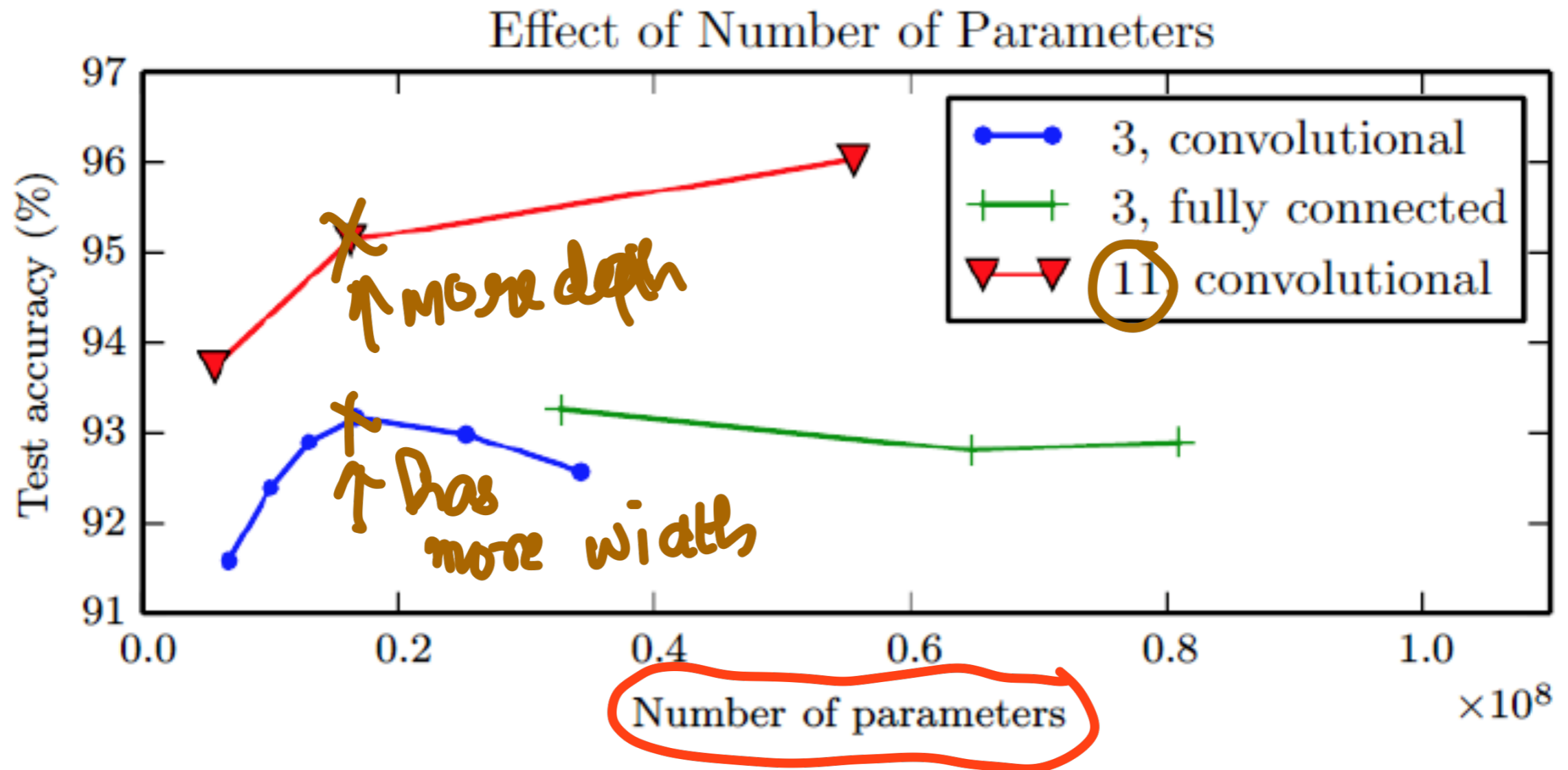


MNIST

$$\mathbf{h}^{(1)} = g^{(1)} \left( \mathbf{W}^{(1)\top} \mathbf{x} + \mathbf{b}^{(1)} \right)$$

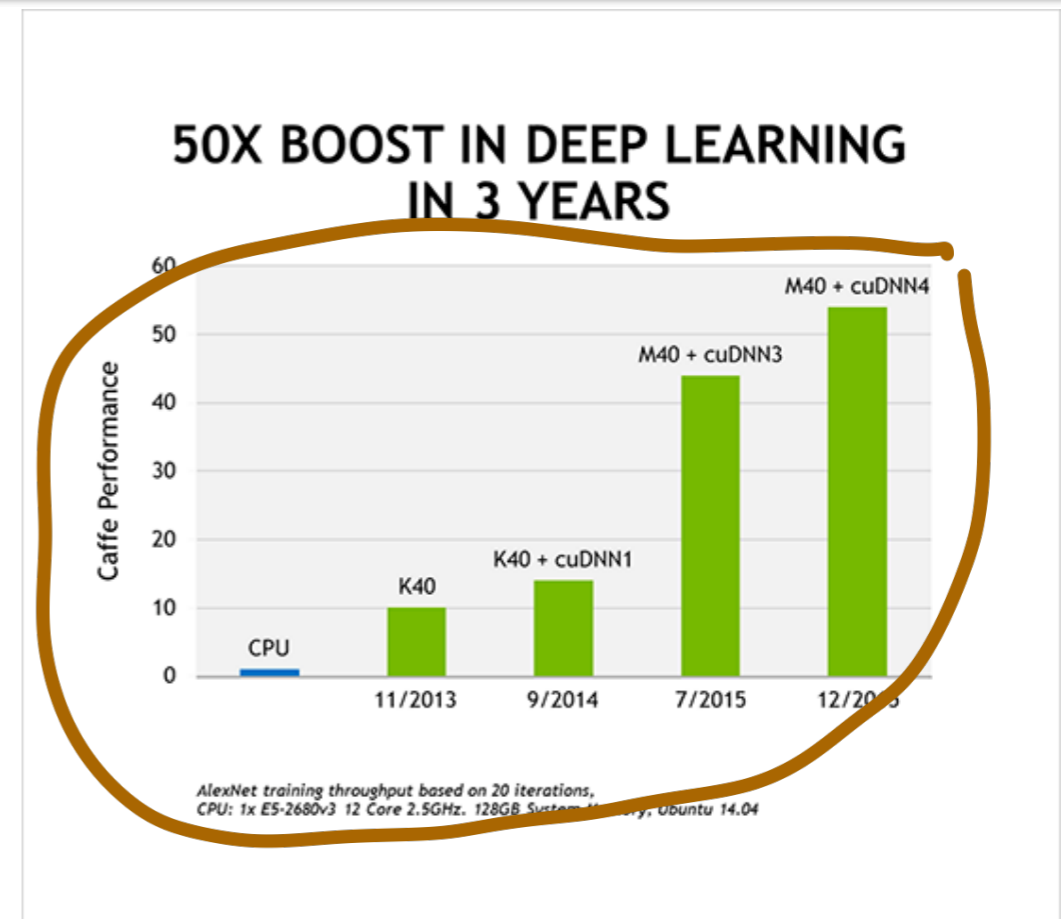
$$\mathbf{h}^{(2)} = g^{(2)} \left( \mathbf{W}^{(2)\top} \mathbf{h}^{(1)} + \mathbf{b}^{(2)} \right)$$

# Need for Depth





# Deep Networks



- Are these networks trainable ?

- Advances in computation and processing
- **Graphical processing units** (GPUs) performing multiple parallel multiply accumulate operations.
- Large amounts of supervised data sets



# Deep Networks

- Will the networks generalize with deep networks *← Deviation from training*
  - DNNs are **quite data hungry** and performance improves by increasing the data.
  - Generalization problem is tackled by providing training data from all possible conditions.
    - Many artificial data augmentation methods have been successfully deployed
  - Providing the **state-of-art performance in several real world applications.**

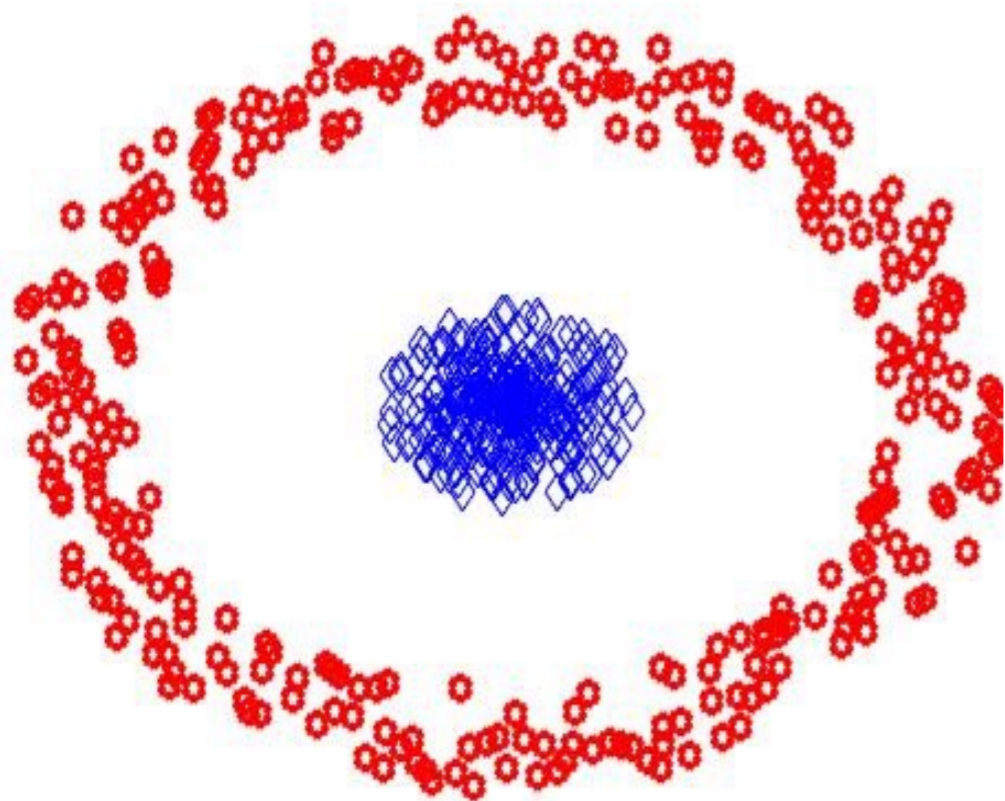
# Representation Learning in Deep Networks

- The input data representation is one of most important components of any machine learning system.

$$\|y - t\|^2$$

Cartesian Coordinates

Polar Coordinates



Non-linear



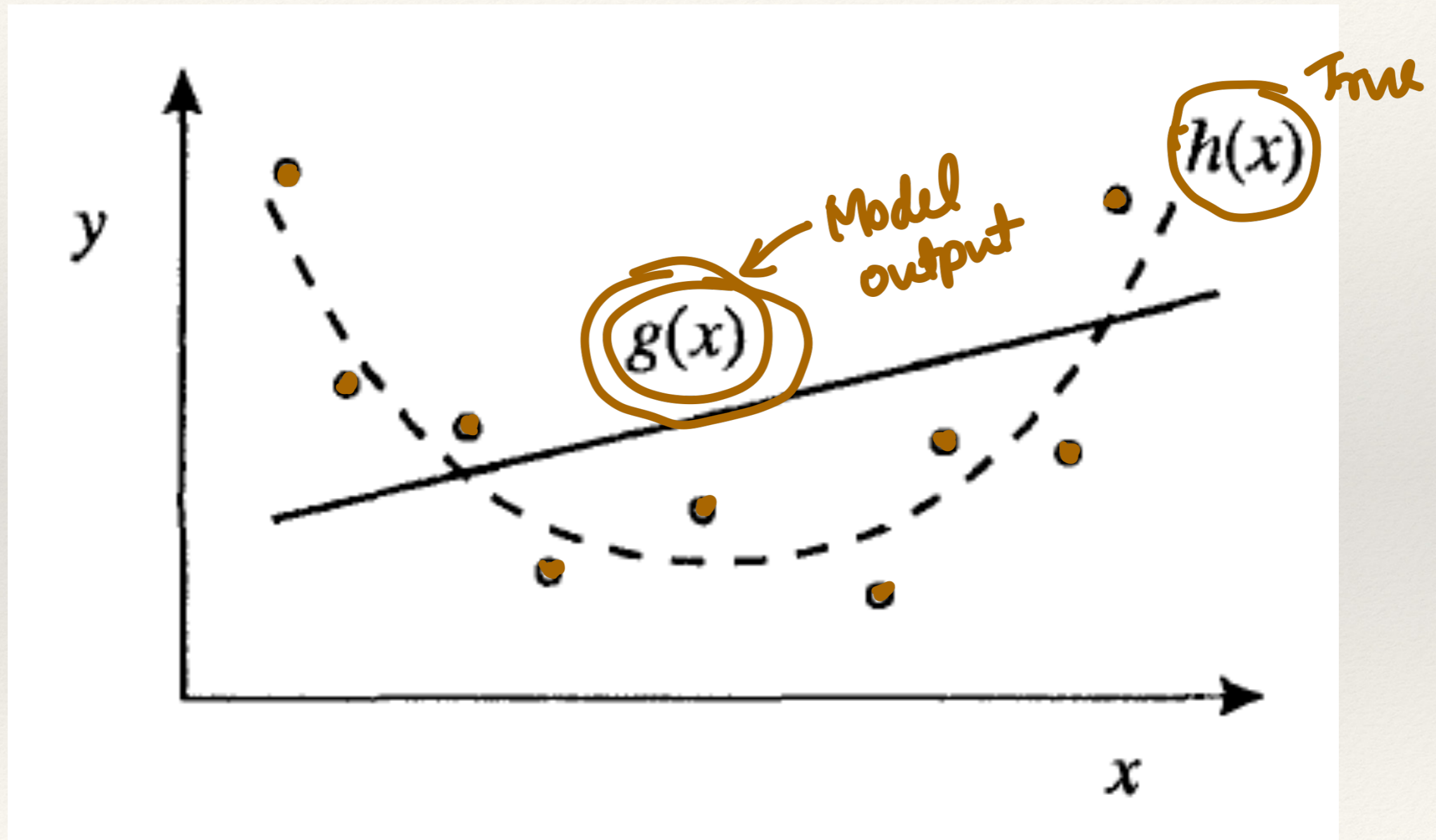
$$\|\log y - \log t\|^2$$

# Representation Learning in Deep Networks

- The input data representation is one of most important components of any machine learning system.
  - Extract factors that enable classification while suppressing factors which are susceptible to noise.
- Finding the right representation for real world applications - substantially challenging.
  - ✱ Deep learning solution - build complex representations from simpler representations.
  - ✱ The dependencies between these hierarchical representations are refined by the target.

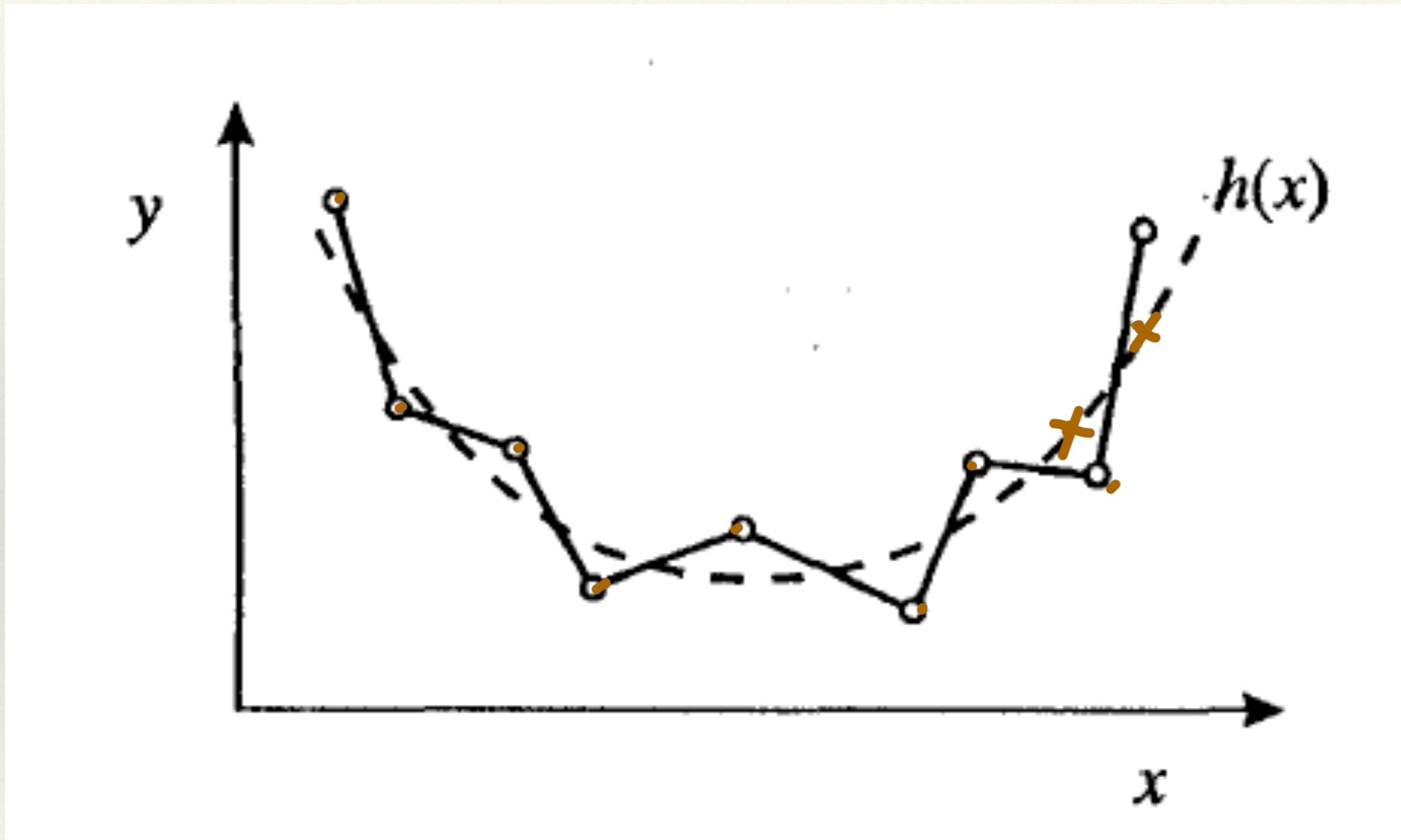


# Underfit



Overfit

→ Lack of Generalization.





---

---

# Avoiding OverFitting In Practice



# Regularization

$E_D(w)$  ← data dependent error function.

$$\|w'\|_2^2 = \sum_i \sum_j (w'_{ij})^2$$
$$\|w''\|_2^2 = \sum_p \sum_q (w''_{pq})^2$$

Overfitting ← following training data too closely.

## Learn patterns

$E_W(w)$  ← Data independent loss function.  
Model dependent

Backpropagation  
propagation →

$$E(w)$$

$$E(w) = E_D(w) + \lambda E_W(w)$$

$\lambda$  - small overfitting  
 $\lambda$  - large (underfit)

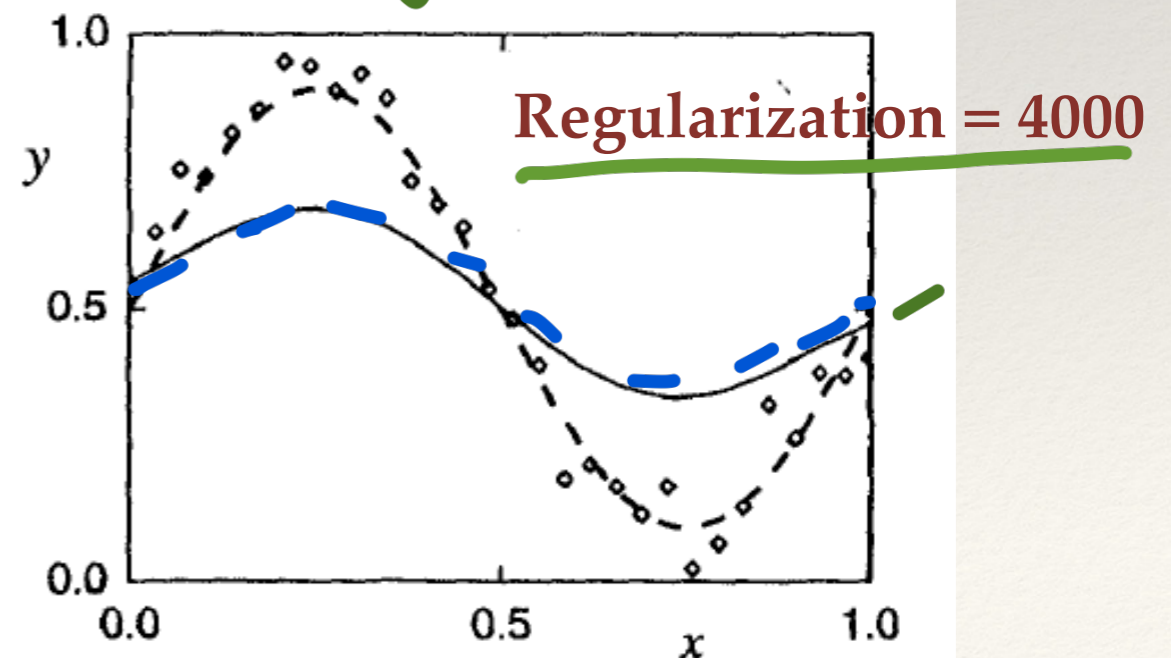
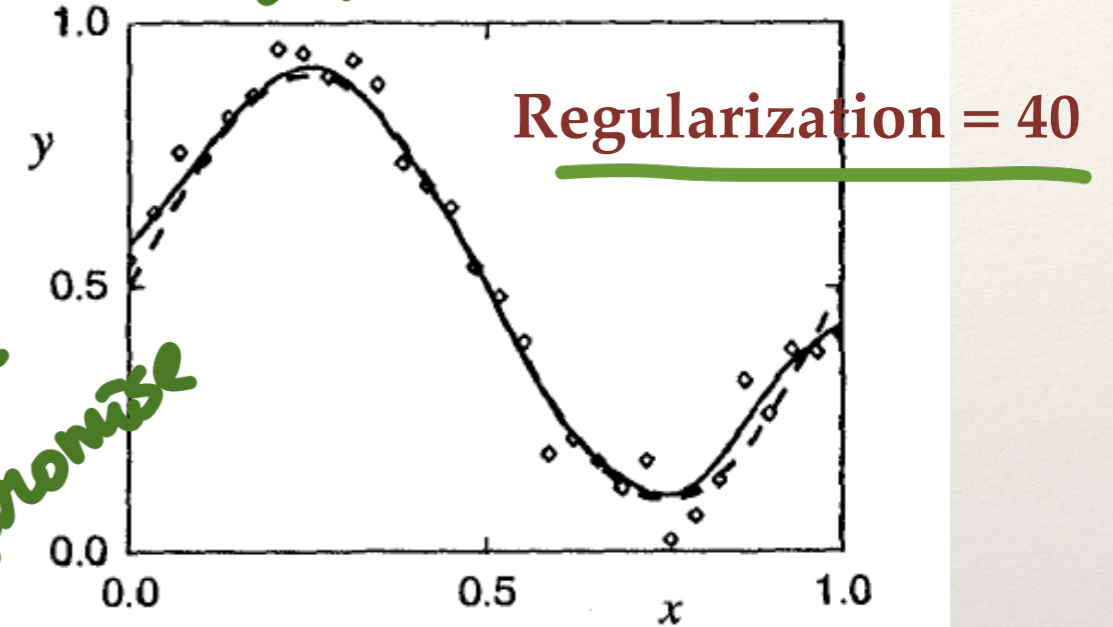
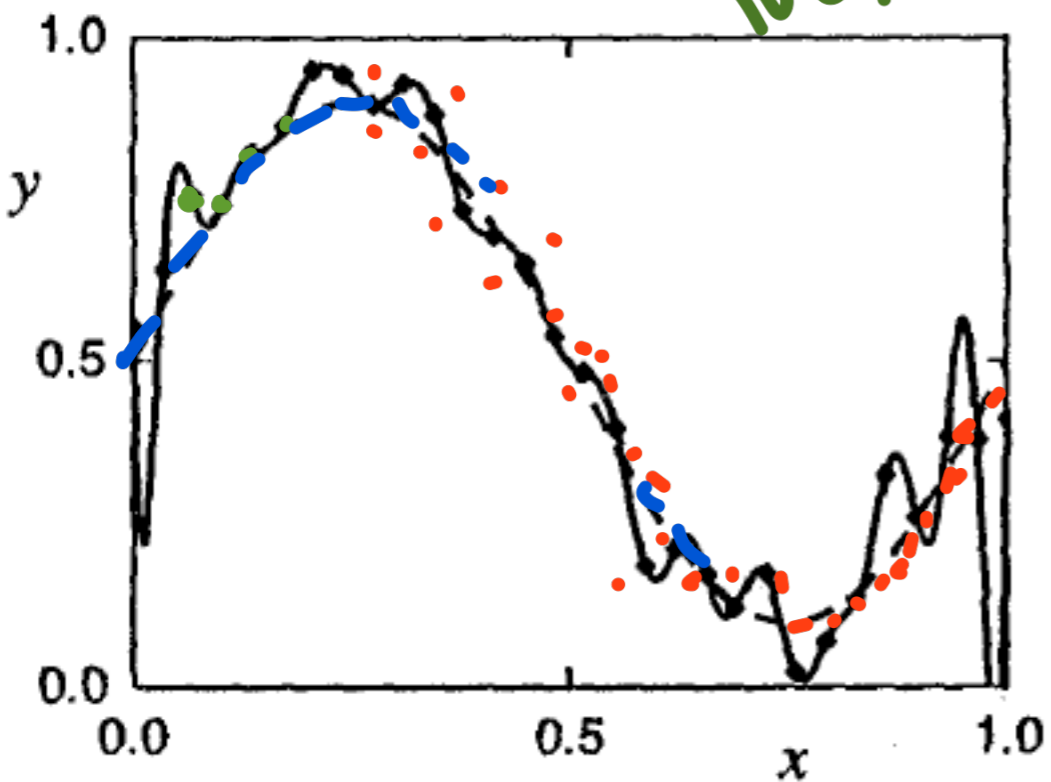
↑  
Regularization Parameter.

Weight Decay Regularization  $E_W(w) = \|w\|_2^2$



# I. Weight Decay Regularization

Regularization = 0

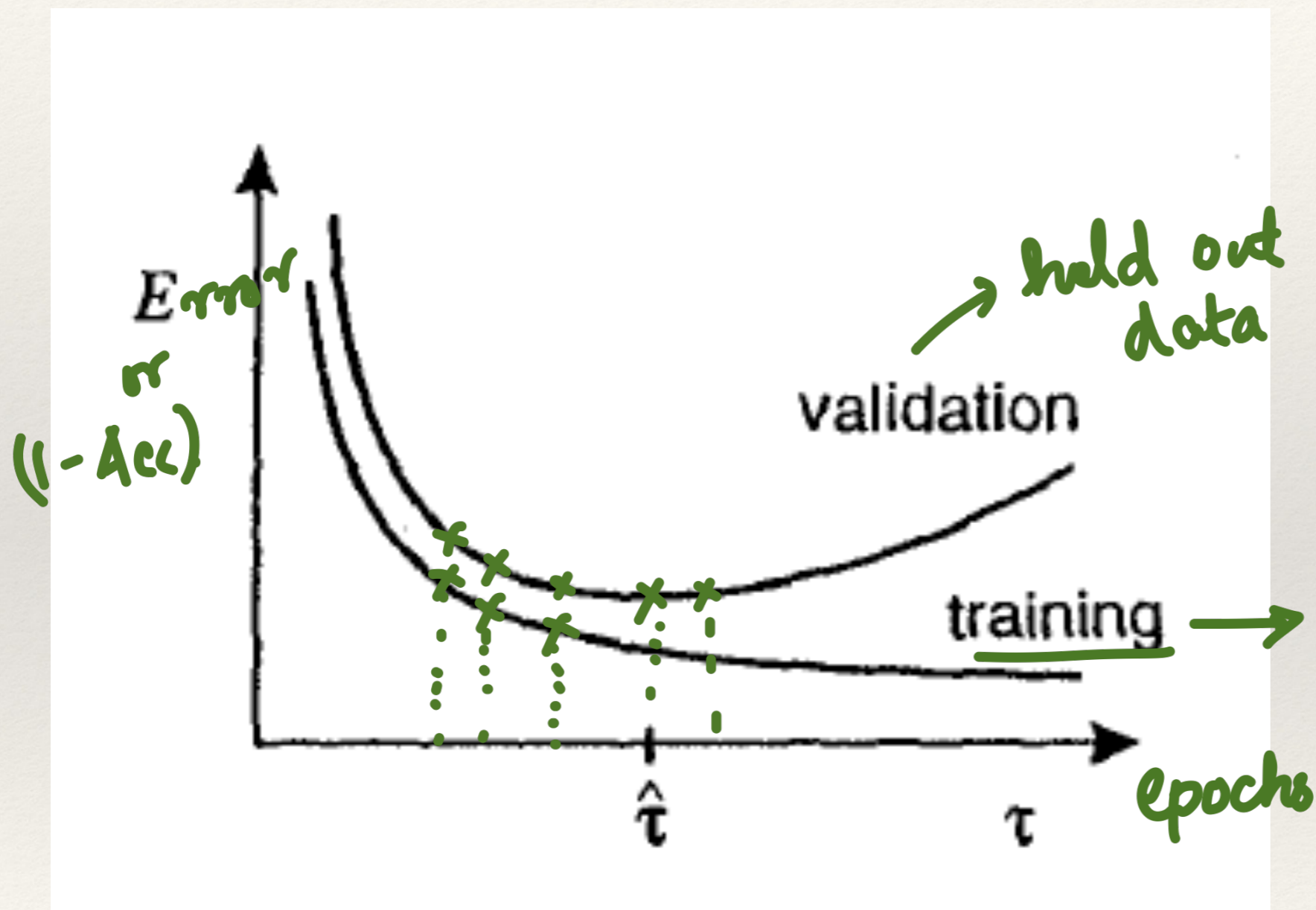


$\lambda$  - small number gives you best compromise.



II.

# Early Stopping



Most Popular in Practice

Stop training when validation performance does not improve.



---

# III. Batch Normalization

---

Batch Normalization: Accelerating Deep Network Training by  
Reducing Internal Covariate Shift

2015

Sergey Ioffe

Google Inc., [sioffe@google.com](mailto:sioffe@google.com)

Christian Szegedy

Google Inc., [szegedy@google.com](mailto:szegedy@google.com)



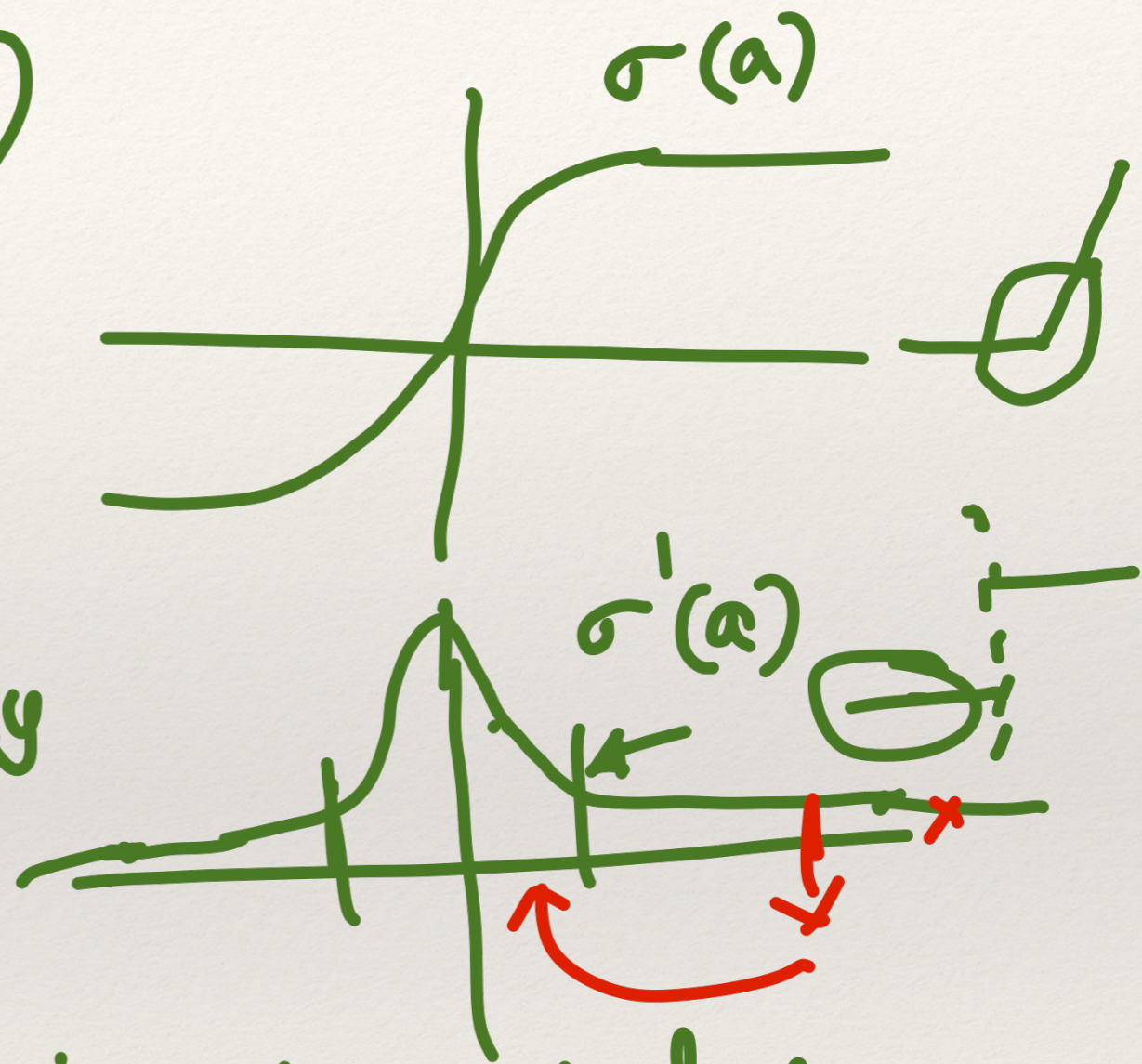
# Sigmoid as non-linearity (Two layers)

$$h^2 = (W^2 \underbrace{\sigma(W^1 \underline{x} + b^1)}_{z^1} + b^2)$$

Model saturation if weight values or inputs are very high.

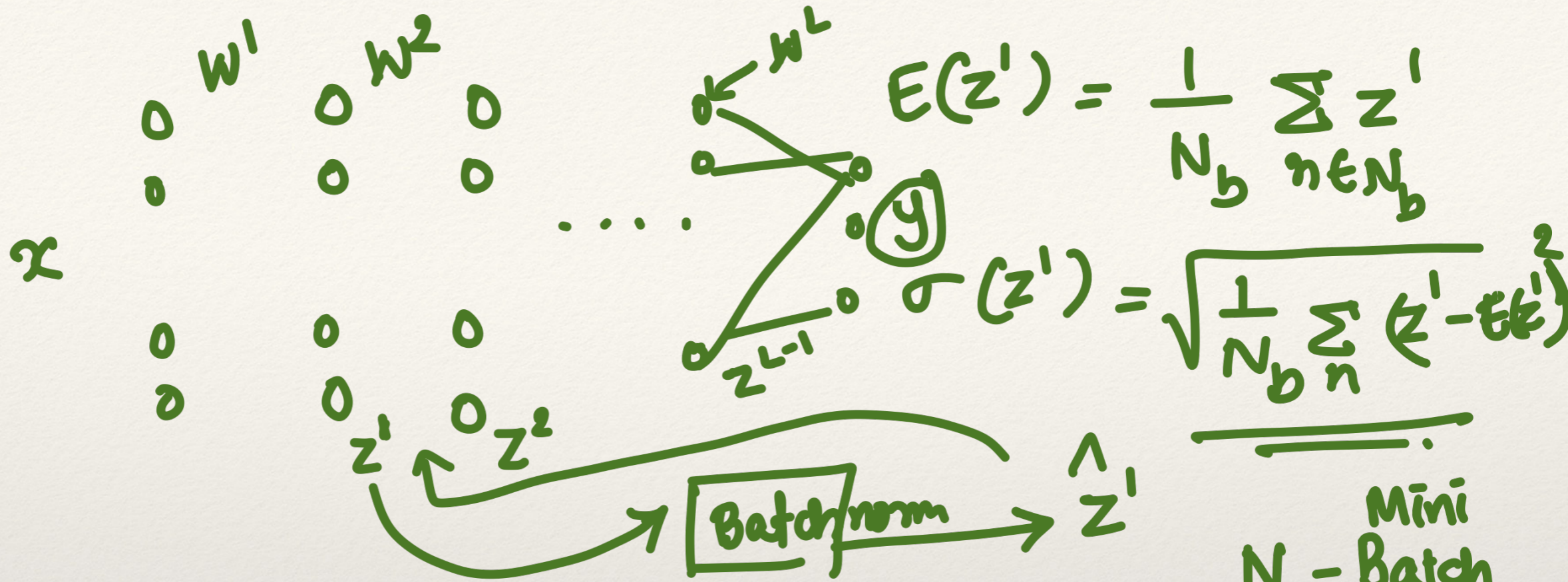
Batch normalization at every layer.

$$\hat{\underline{z}}^i = \frac{z^i - E(z^i)}{\sigma(z^i)}$$



← This makes each layer output normalized.  
(close to 0, have small variance only).





$$E(z^1) = \frac{1}{N_b} \sum_{n \in N_b} z^1$$

$$\sigma(z^1) = \sqrt{\frac{1}{N_b} \sum_n (z^1 - E(z^1))^2}$$

Mini  
 $N_b$  - Batch  
size

should be  
 $O(10^2)$

Can be done at every layer.

\* higher learning rate

\* Reduced number of epochs



# Effect of Batch Normalization

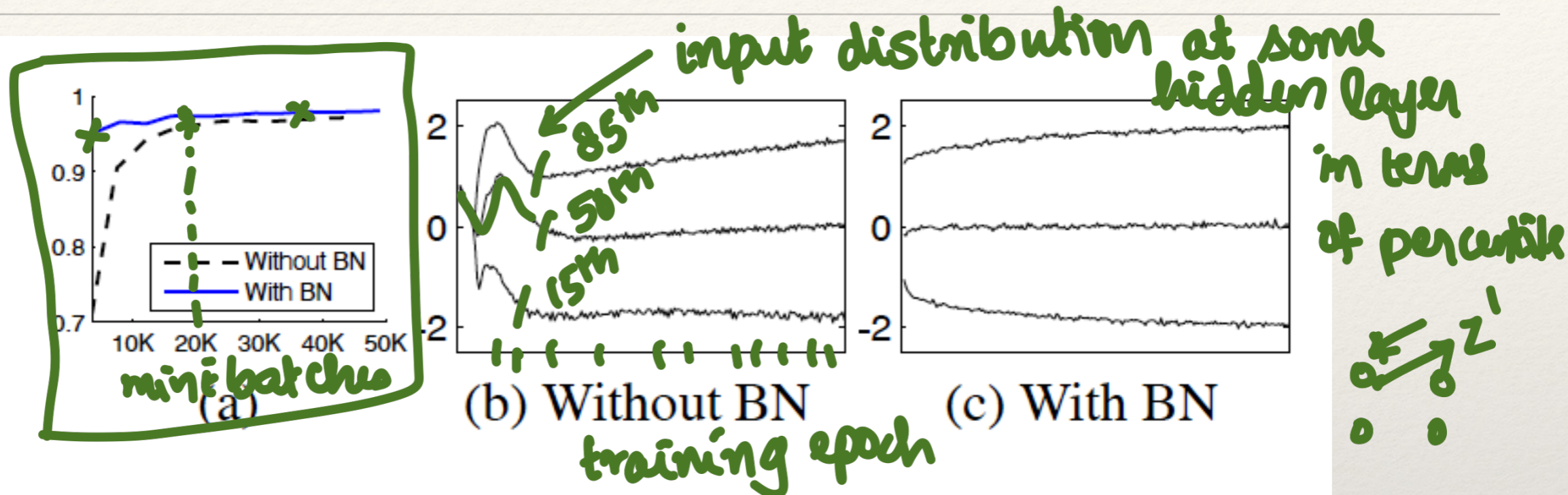


Figure 1: (a) The test accuracy of the MNIST network trained with and without Batch Normalization, vs. the number of training steps. Batch Normalization helps the network train faster and achieve higher accuracy. (b, c) The evolution of input distributions to a typical sigmoid, over the course of training, shown as {15, 50, 85}th percentiles. Batch Normalization makes the distribution more stable and reduces the internal covariate shift.

$$\{z^i\}_{i=1}^N$$



# IV. Dropout Strategy in Neural Network Training

## Dropout: A Simple Way to Prevent Neural Networks from Overfitting

Nitish Srivastava

Geoffrey Hinton

Alex Krizhevsky

Ilya Sutskever

Ruslan Salakhutdinov

*Department of Computer Science*

*University of Toronto*

*10 Kings College Road, Rm 3302*

*Toronto, Ontario, M5S 3G4, Canada.*

NITISH@CS.TORONTO.EDU

HINTON@CS.TORONTO.EDU

KRIZ@CS.TORONTO.EDU

ILYA@CS.TORONTO.EDU

RSALAKHU@CS.TORONTO.EDU

2013

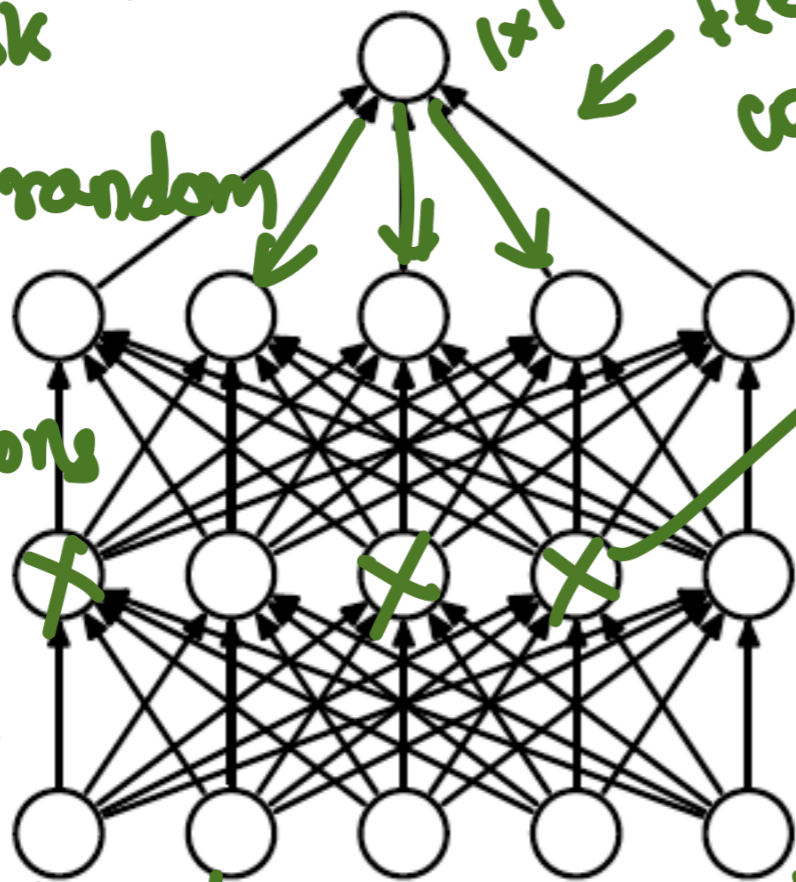
Editor: Yoshua Bengio



# Dropouts in Neural Networks

\* Nodes (neurons) are too high to perform task

\* Any random subset of neurons should still be able to predict o/p



(a) Standard Neural Net

feed forward connection

probability of switching it off

$W_2$   
100x100

$W_1$   
100x100



(b) After applying dropout.

softmax

$W_2 = [ ]$   
70x70

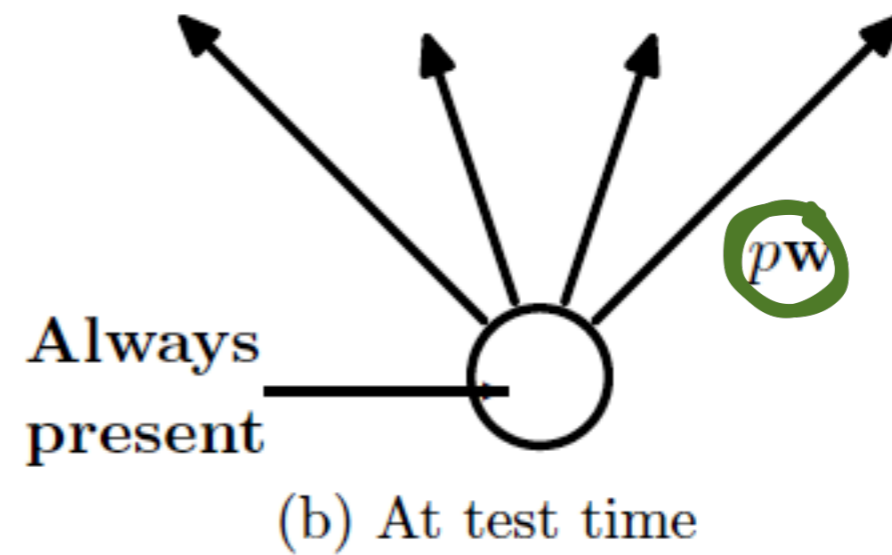
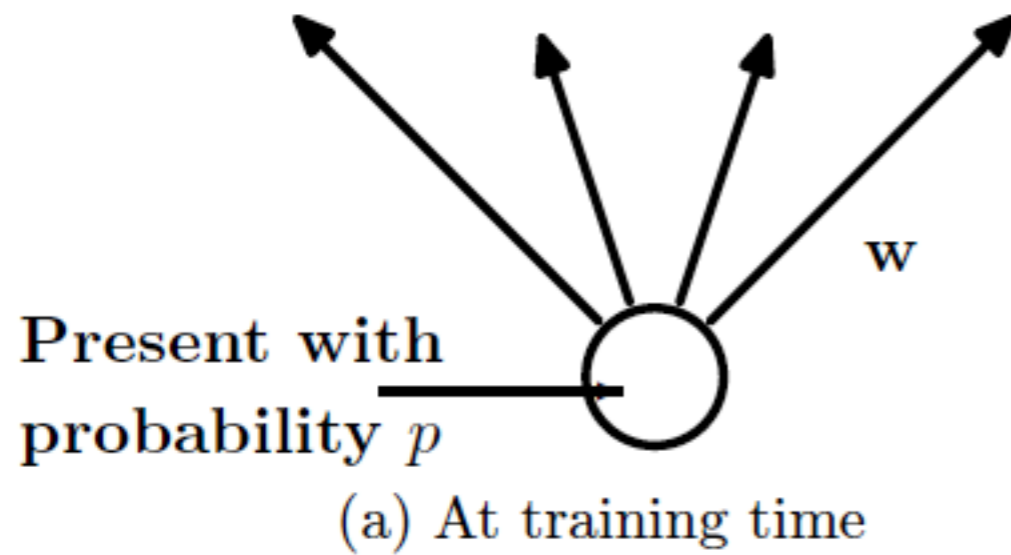
$W_1 = [ ]$   
70x70

Fixed Dropout probability

(Sparsity)



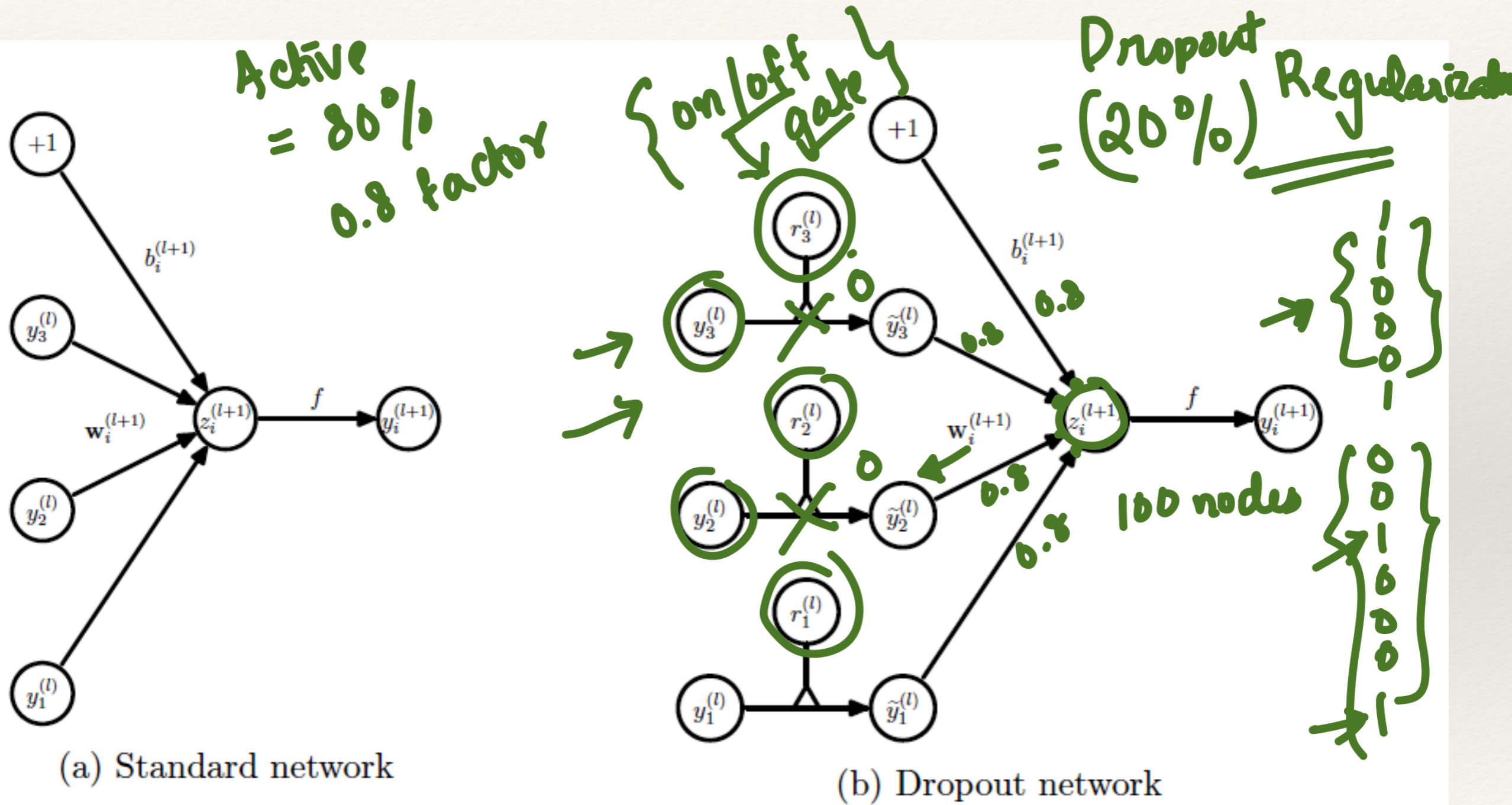
# Dropout in Training and Test



Consistent o/p  
in testing



# Dropout Application



(a) Standard network

(b) Dropout network

Figure 3: Comparison of the basic operations of a standard and dropout network.



# Effect of Dropouts

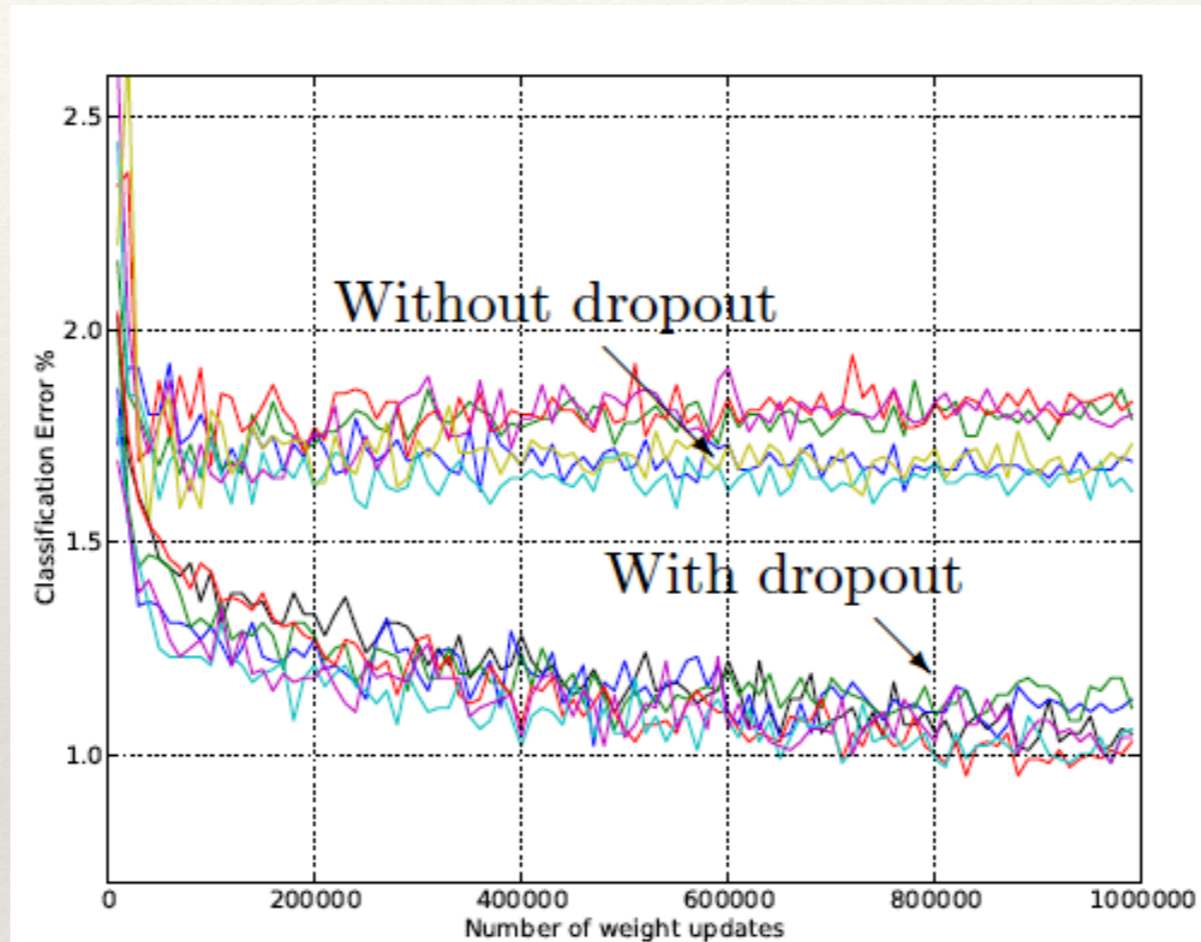


Figure 4: Test error for different architectures with and without dropout. The networks have 2 to 4 hidden layers each with 1024 to 2048 units.



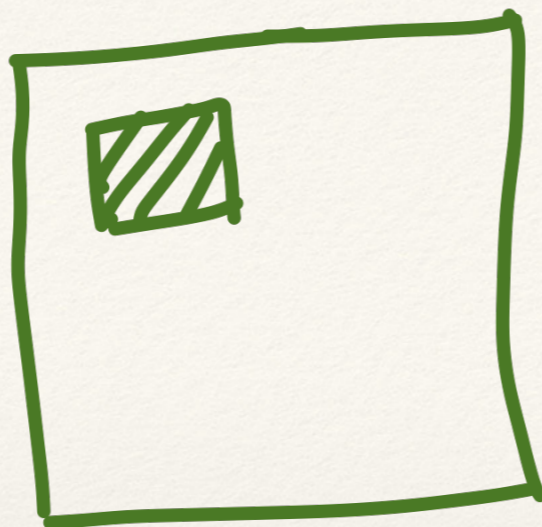
---

---

# Convolutional Neural Networks



# CNNs



2-D patterns

$v/s$



# DNN



Weights



$$(Wx + b)$$

1-D

↓ Can I also have weights which are 2-D and local. [process parts of Image].

## Convolution

$$Y(i, j)$$

## Operations

$$= \sum_{m=0}^{M-1} \sum_{n=0}^{N-1} X(i+m, j+n) K^p(m, n)$$

$$K \in \mathbb{R}^{M \times N}$$

$$K^p(m, n)$$

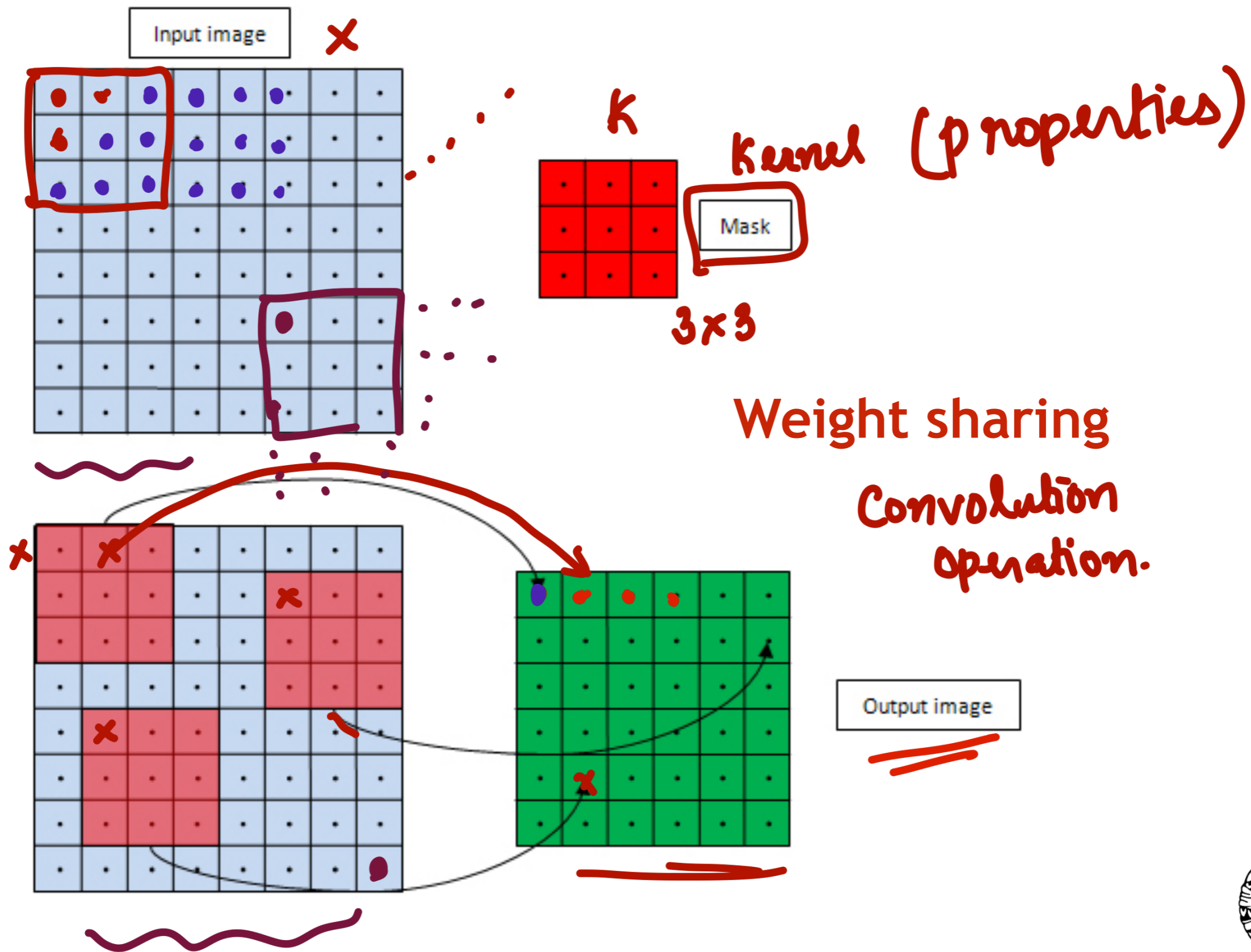
$p=1 \dots P$

Learnable parameter

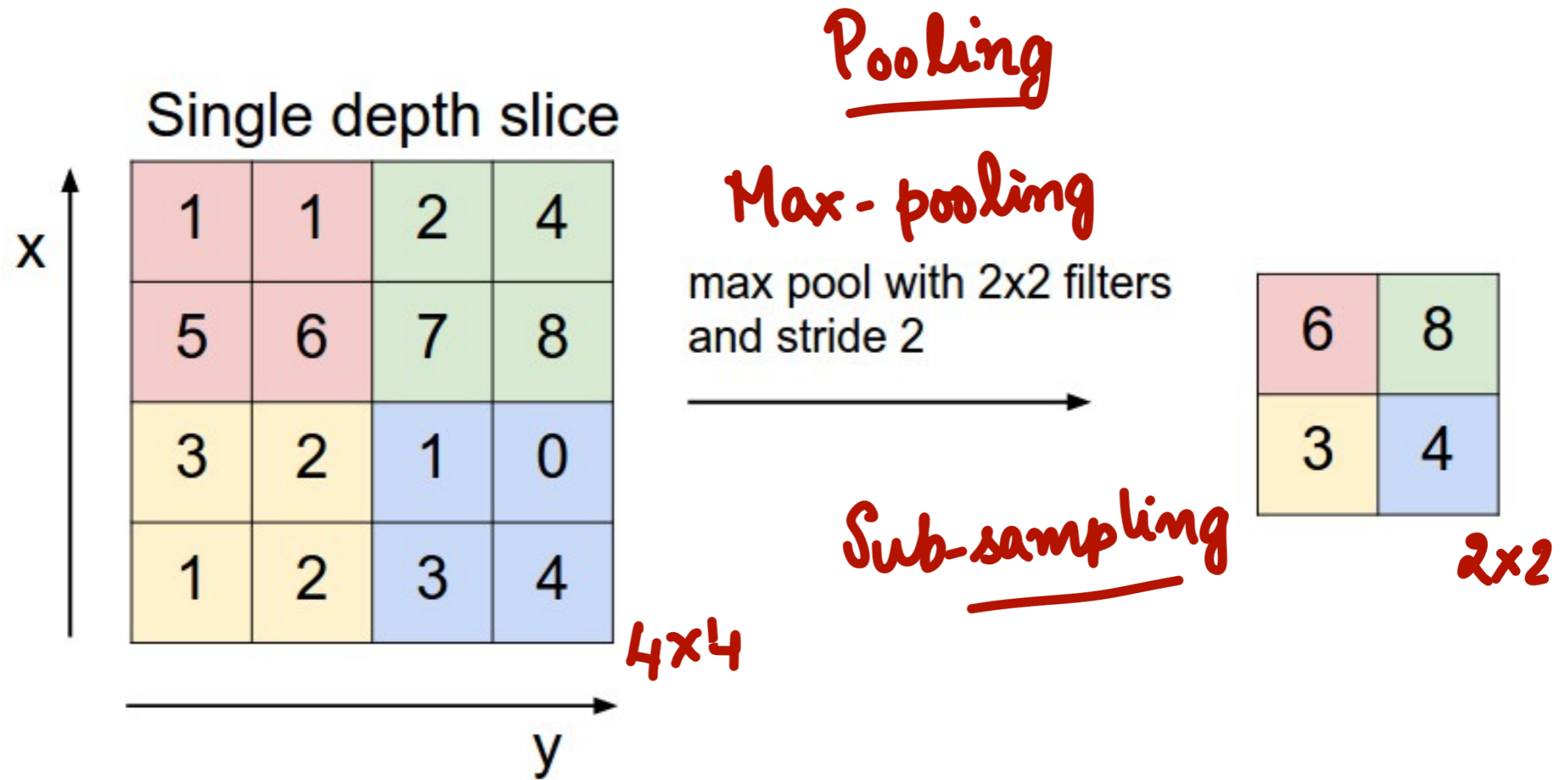
Destroys 2-D structure



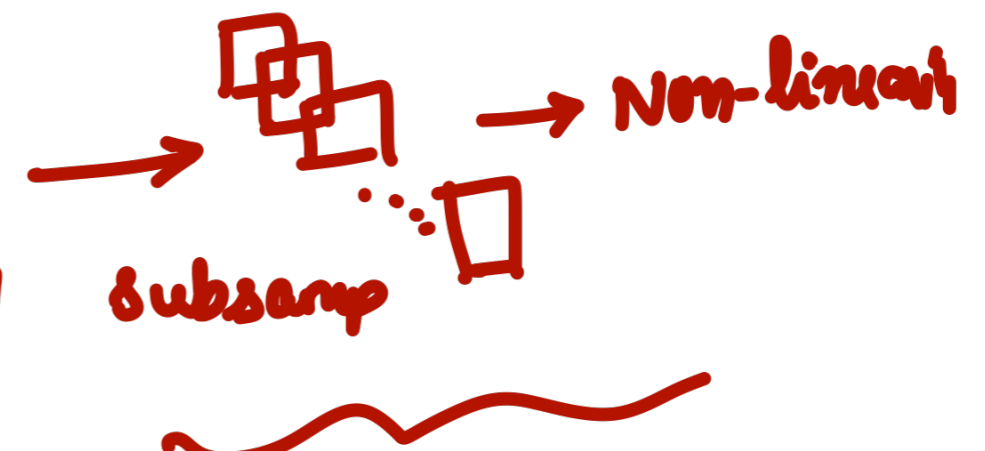
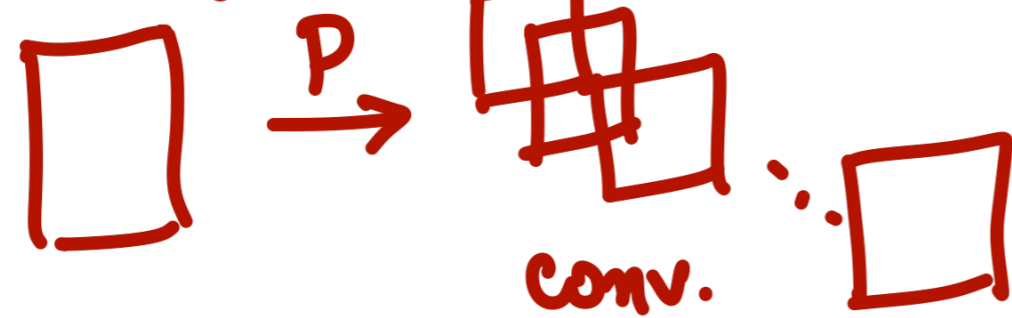
# Other Architectures - Convolution Operation



# Max Pooling Operation

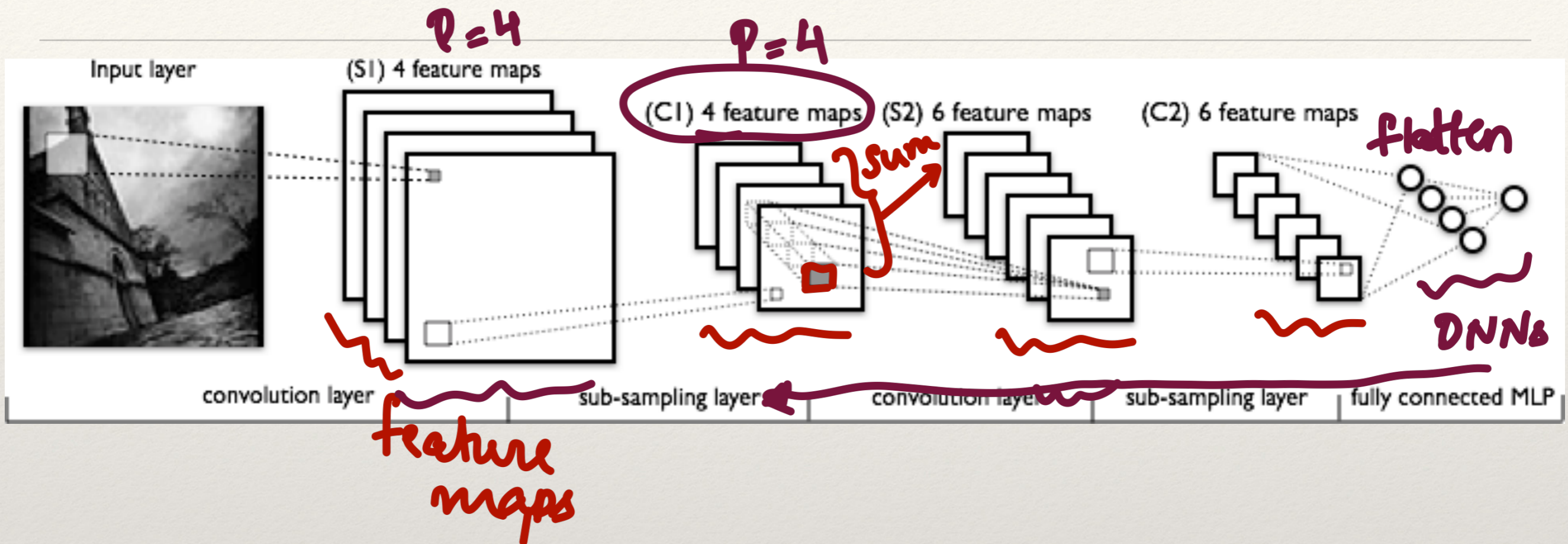


A layer of CNN





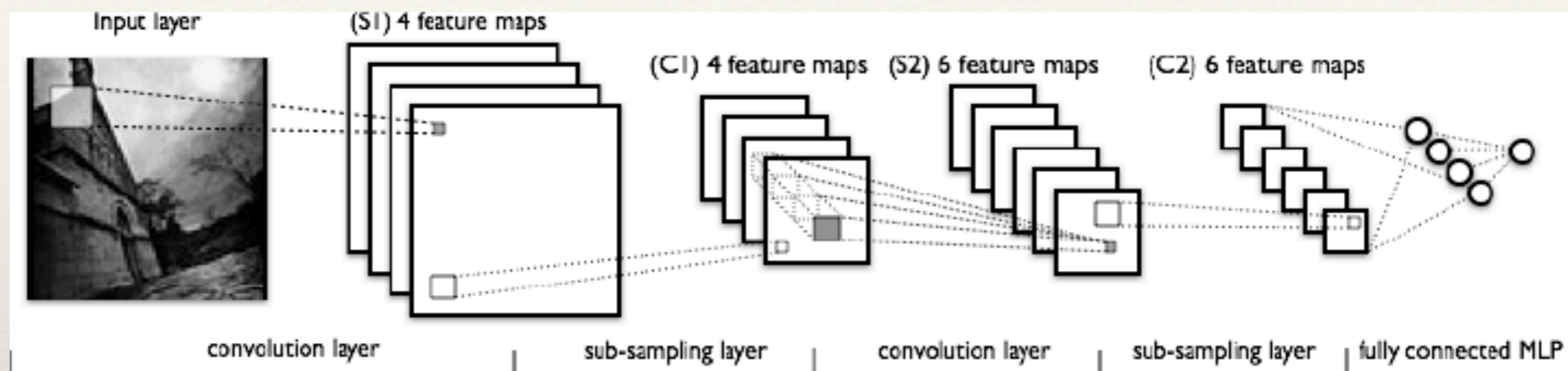
# Convolutional Neural Networks



- Multiple levels of filtering and subsampling operations.
- Feature maps are generated at every layer.



# Convolutional Neural Networks



- Multiple levels of filtering and subsampling operations.
- Feature maps are generated at every layer.



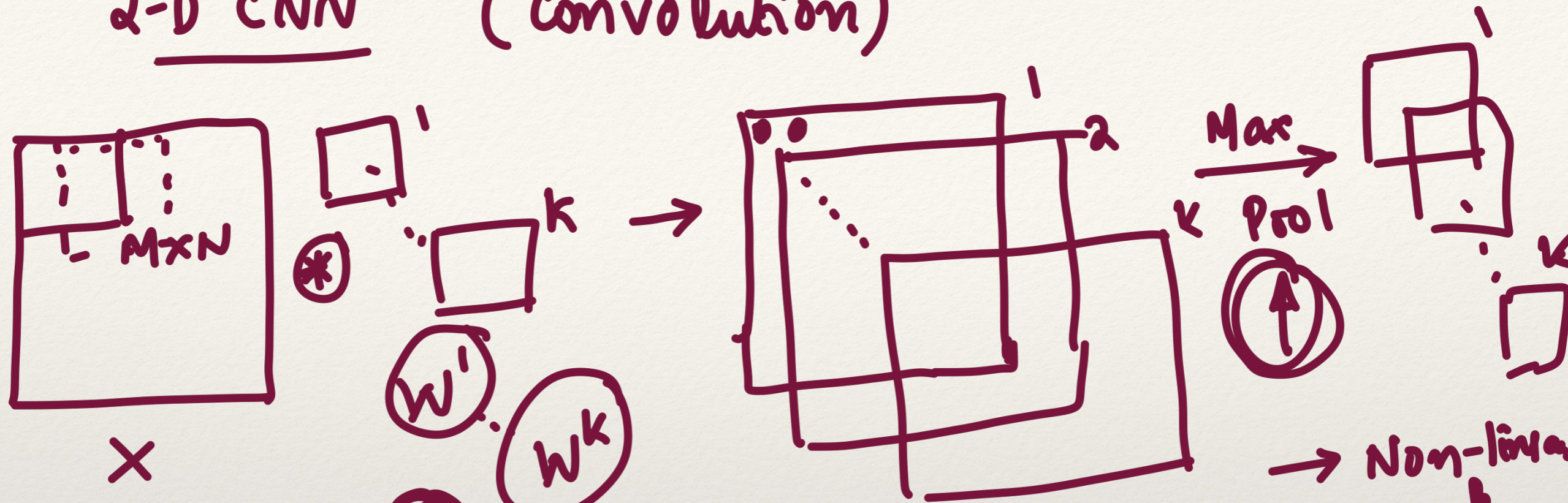
---

# Back Propagation in CNNs

---



# 2-D CNN (Convolution)



$$E = \sum_k E [t^n, y^n]$$

$$\frac{\partial E}{\partial w^k}$$

$$\frac{\partial E}{\partial y^k}$$

$$\frac{\partial E}{\partial w^k}$$

$$X \in \mathbb{R}^{D_1 \times D_2}$$

If no padding

$$y^k \in \mathbb{R}^{D_1 - M, D_2 - N}$$

$$y^k = \sum_{m,n} x(i+m, j+n) w^k(m,n)$$

Non-linear

$$[w^k]_{(M+1, N+1)}$$



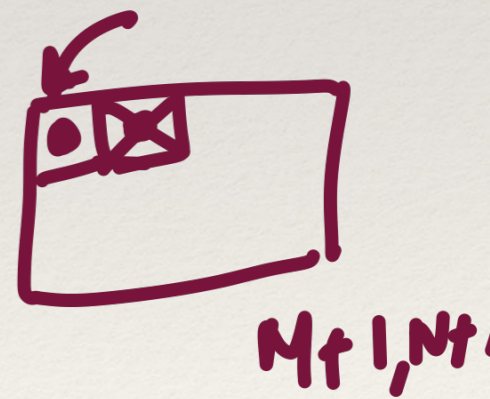
$$\underline{\underline{y^k[0,0]}} = \underbrace{x[0,0]}_{\text{convolution}} \underbrace{w^k[0,0]} + \underbrace{x[0,1]} \underbrace{w^k[0,1]} + \dots + x[0,N] w^k[0,N] + x[1,0] w^k[1,0] + \dots + x[M,N] w^k[M,N]$$

$$\underline{\underline{y^k[0,1]}} = \underbrace{x[0,1]} w^k[0,0] + \dots + x[0,N+1] w^k[0,N] + x[1,1] w^k[1,0] + \dots + x[M+1,N] w^k[M,N]$$

$$\begin{matrix} \vdots \\ \rightarrow \text{scalar} \\ \left[ \begin{array}{c} \frac{\partial E}{\partial y^k} \end{array} \right] \end{matrix} \parallel \underline{\underline{(D_1 - M, D_2 - N)}}$$

$$\left[ \begin{array}{c} \frac{\partial E}{\partial w^k} \end{array} \right] \parallel \underline{\underline{M+1, N+1}}$$

$[w^k]_{M+1, N+1}$





$$\frac{\partial \mathcal{E}}{\partial w^k [0,0]} = x[0,0] \frac{\partial \mathcal{E}}{\partial y^k [0,0]} + \dots + x[0,1] \frac{\partial \mathcal{E}}{\partial y^k [0,1]}$$

$$* \frac{\partial \mathcal{E}}{\partial w^k [0,1]} = \underline{x[0,1]} \frac{\partial \mathcal{E}}{\partial y^k [0,0]} + \dots + x[0,-M, D_2 - N + 1] \frac{\partial \mathcal{E}}{\partial y^k [0,-M, D_2 - N]}$$

$\vdots$   
 $M, N$

$$\frac{\partial \mathcal{E}}{\partial w^k [m,n]} = \sum_{p,q=0,0}^{D_1-M, D_2-N} x[m+p, n+q] \frac{\partial \mathcal{E}}{\partial y^k [p,q]}$$

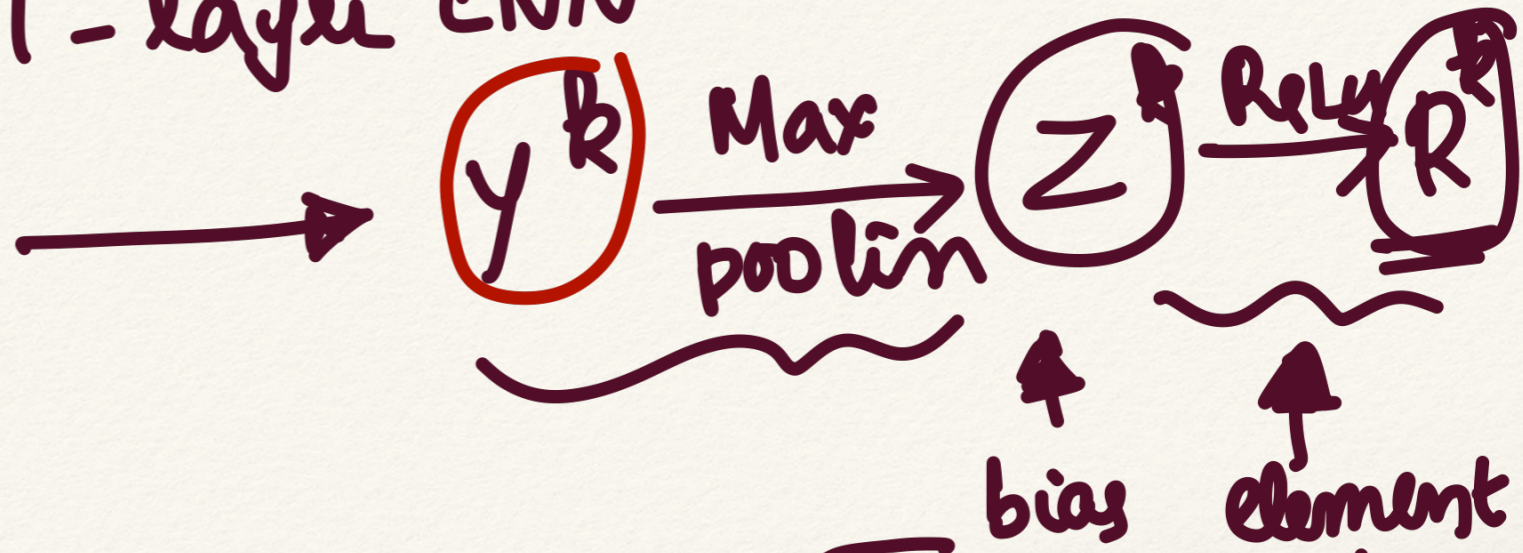
Conv.



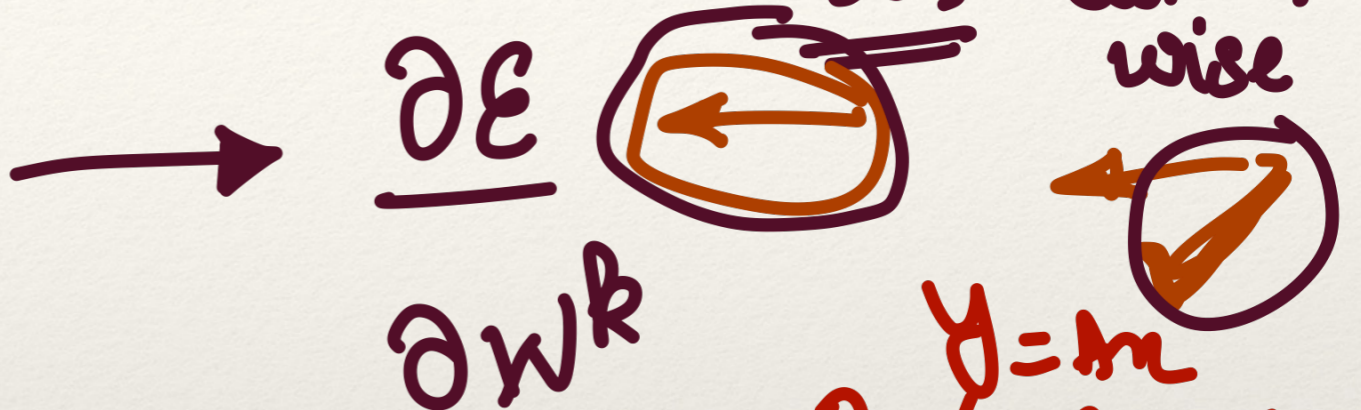
Forward



1-layer CNN

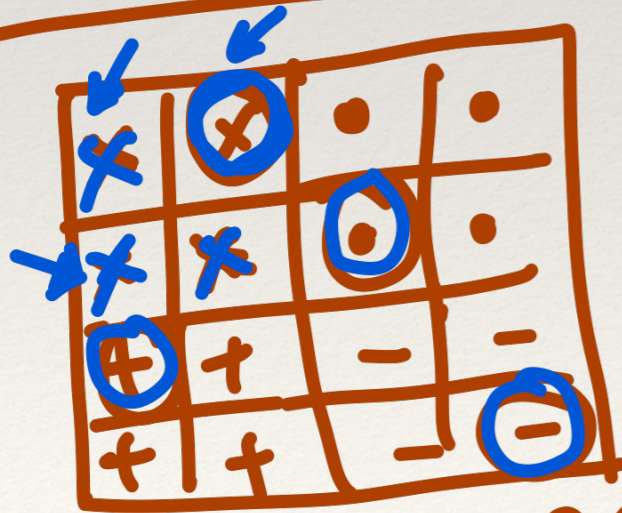


Backward

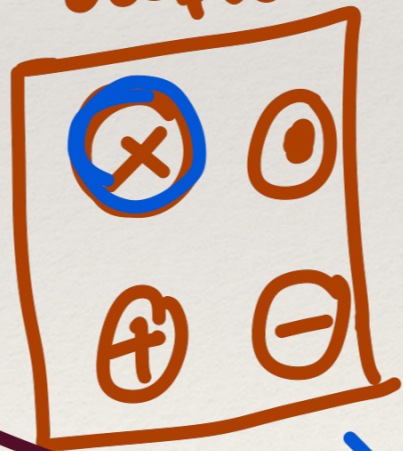


$$\frac{\partial}{\partial x} (Ax) = A^T$$

Max Pooling

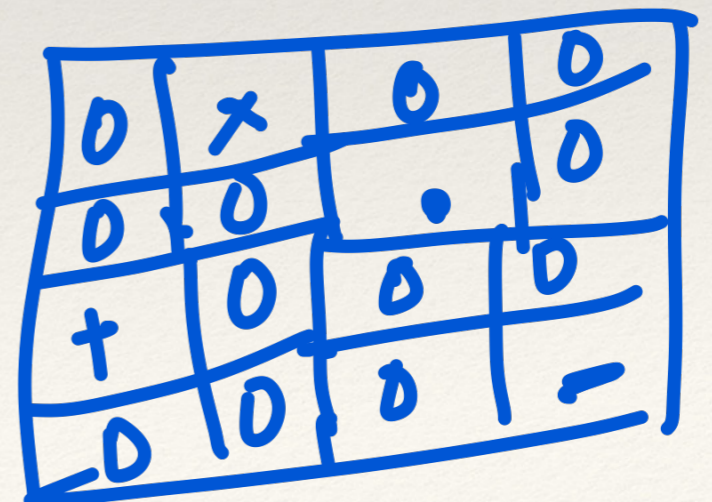


output



$$\frac{\partial \mathcal{E}}{\partial z^k}$$

Unpooling

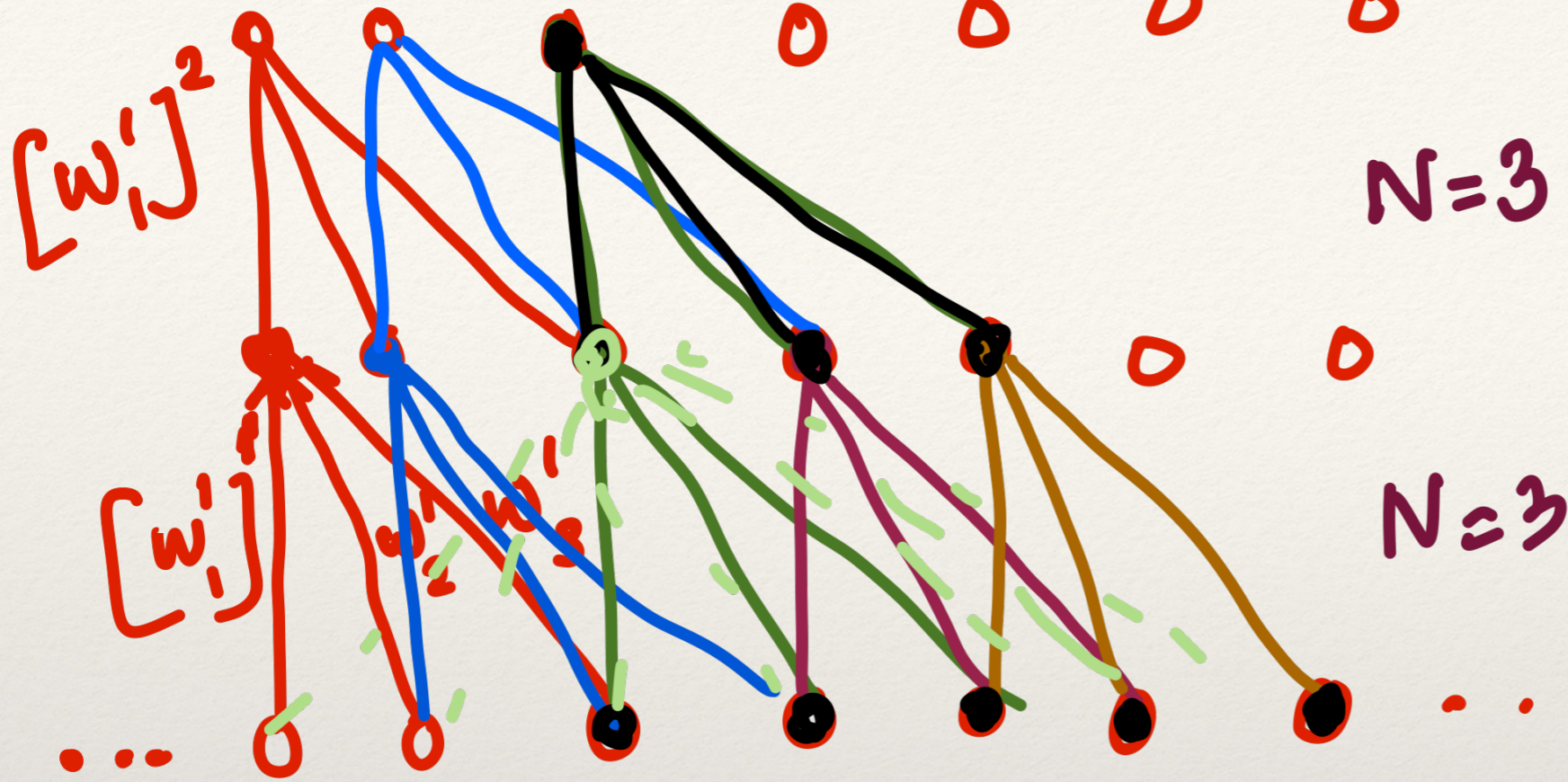


$$\frac{\partial \mathcal{E}}{\partial z^k}$$

argmax



# 1-D Convolutional Network



2<sup>nd</sup> layer

$N=3$

$N=3$

1 kernel  
output

$k$ -such outputs  $k$ .

$w_k$   $\rightarrow$  weight

$x[n]$

Layer in more depth are looking at more and more global properties of the input

for  $i$ <sup>th</sup> lag and  $k$ <sup>th</sup> kernel

$$\underline{w}^2 = \begin{bmatrix} 0 \\ 1 \\ 3 \end{bmatrix} \quad \underline{w}^1 = \begin{bmatrix} 1/2 \\ 1 \\ 2 \end{bmatrix} \quad \dots$$



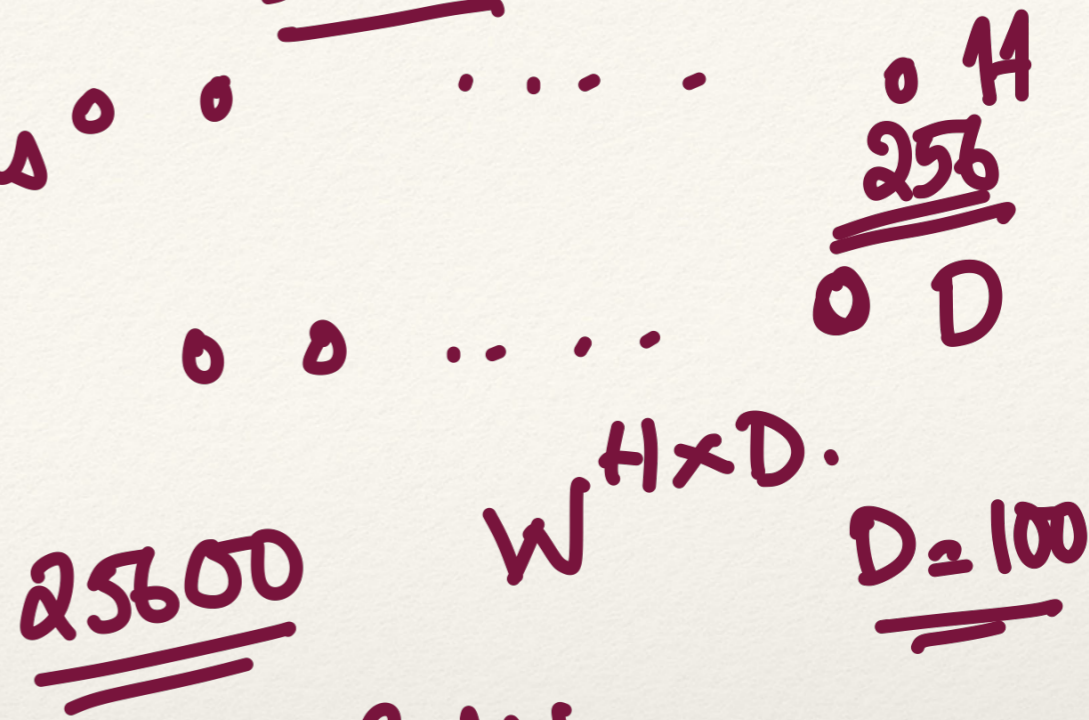
# Number of parameters

$$N = \underline{3} \text{ or } \underline{5} \text{ or } \underline{7}$$

$$K = \underline{64} \text{ or } 128 \text{ or } \underline{256}$$

Typical values

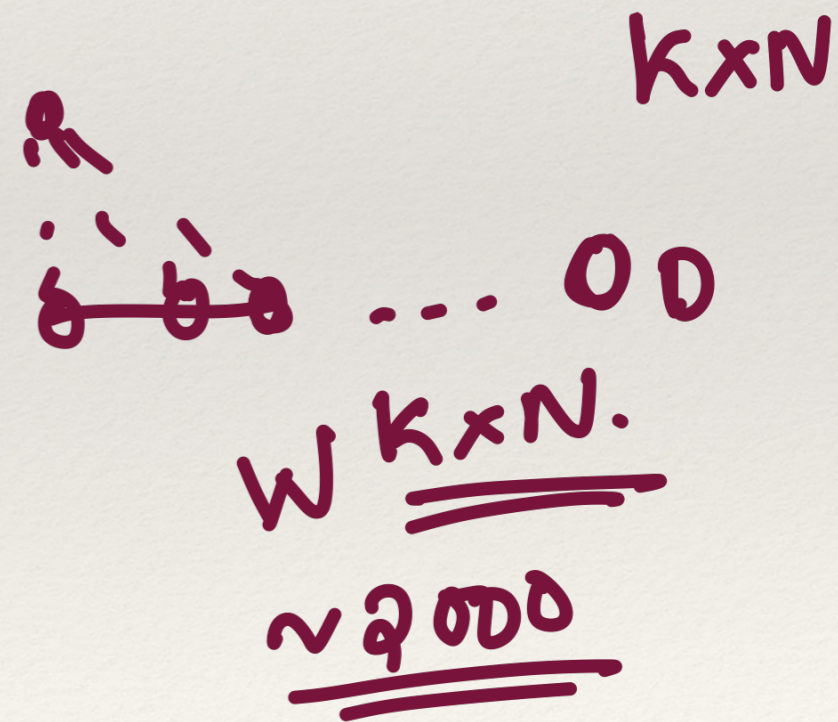
DNN



CNN

CNN layers have much smaller # of parameters

↳ weight sharing





# Expression for a 1-D CNN

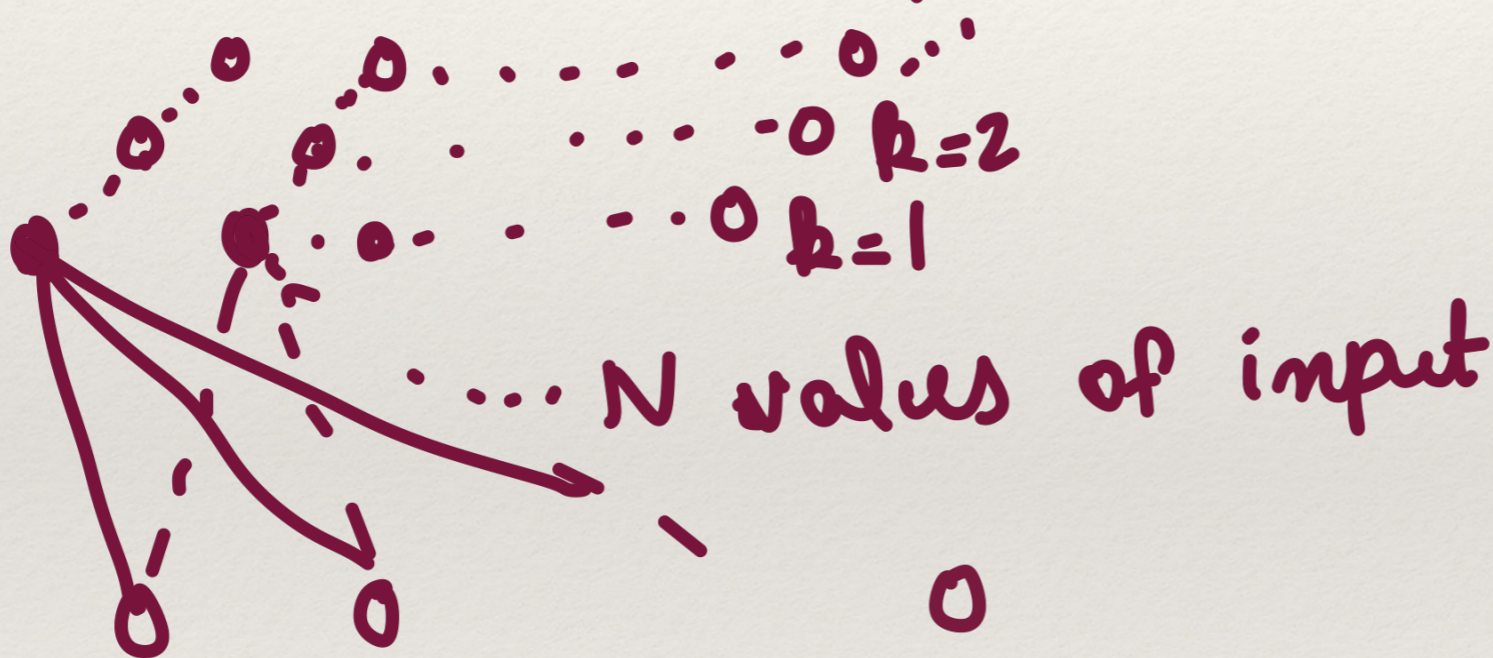
$$\underbrace{y^k[n]}_{\text{kernel}} = \sum_{i=1}^N \underline{x[n+i]} \quad \underbrace{w^k}_{\text{kernel width}}$$

$N$  - kernel width

$K$  - # kernels

$$k = \underline{1 \dots K}$$

$k \times k$



# parameters is  $\underline{k \times N}$

$$\underbrace{w^k}_{\text{kernel width}} = \begin{bmatrix} -1 \\ 0 \\ 1 \end{bmatrix} \quad \underbrace{w^2}_{\text{kernel width}} = \begin{bmatrix} 1 \\ -1 \\ 1 \end{bmatrix}$$