



# Convnets

Nando de Freitas



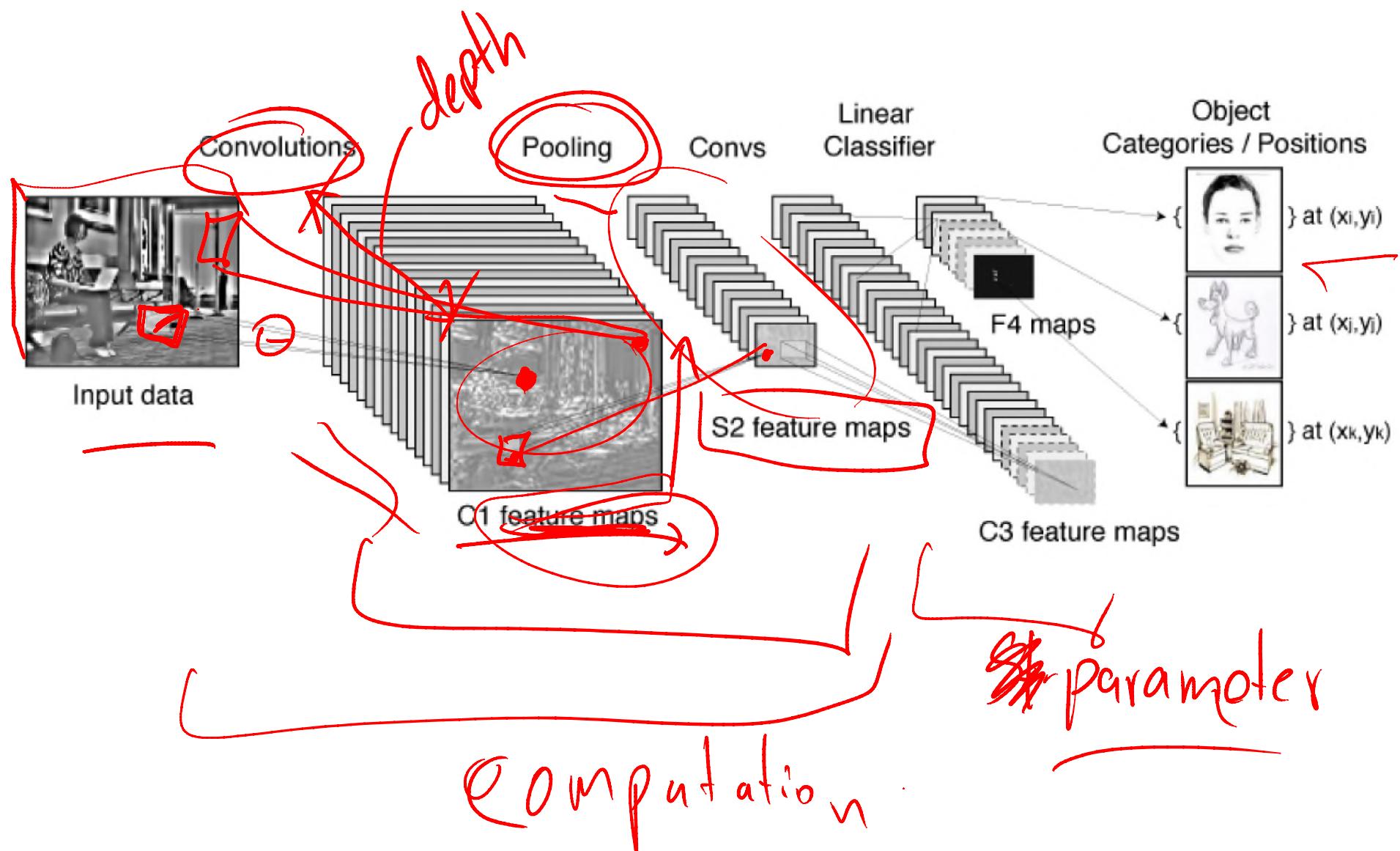
UNIVERSITY OF  
**OXFORD**

# Outline of the lecture

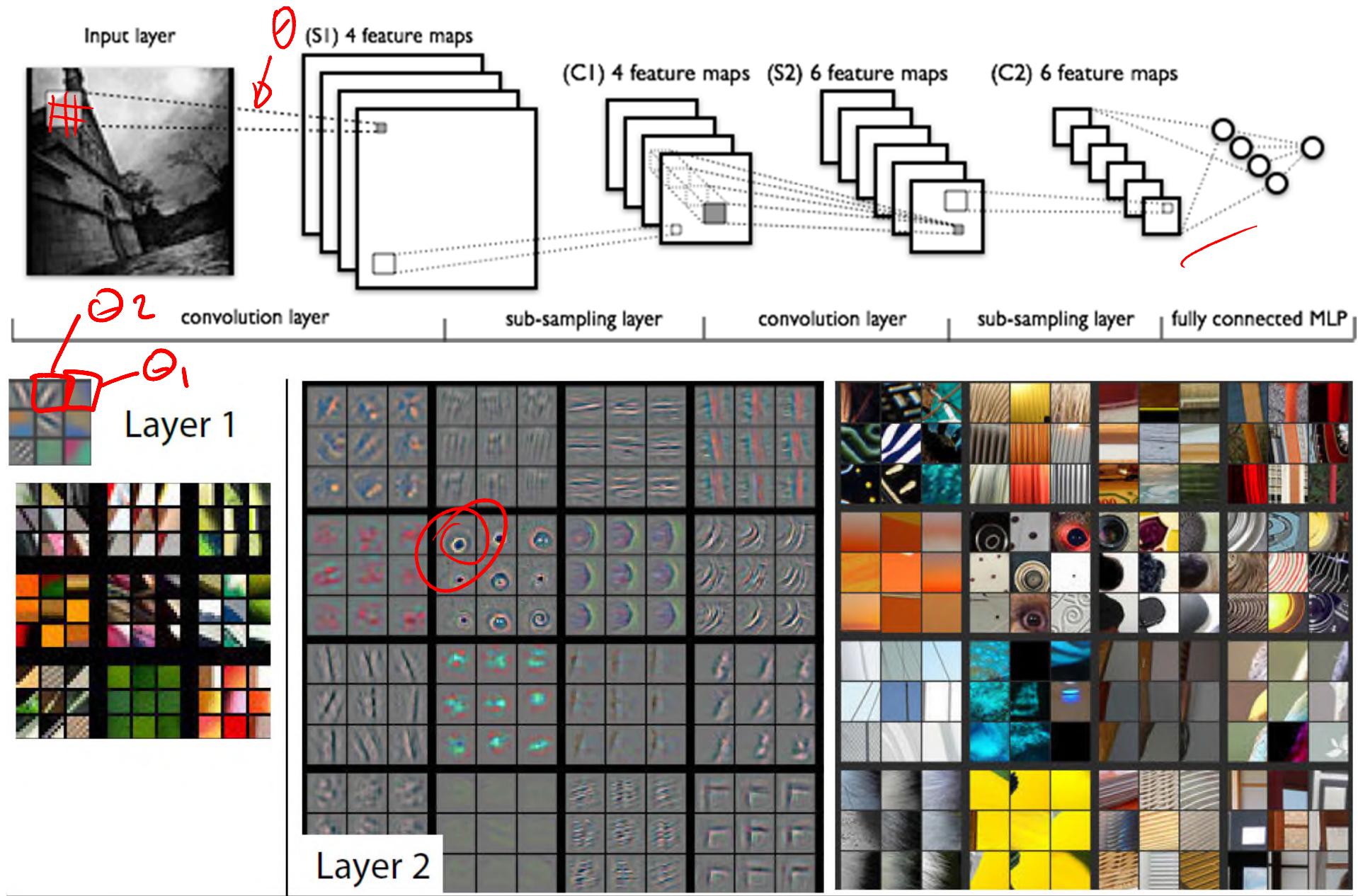
This lecture introduces you to convolutional neural networks. These models have revolutionized speech and object recognition. The goal is for you to learn

- Convnets for object recognition and language
- How to design convolutional layers
- How to design pooling layers
- How to build convnets in torch

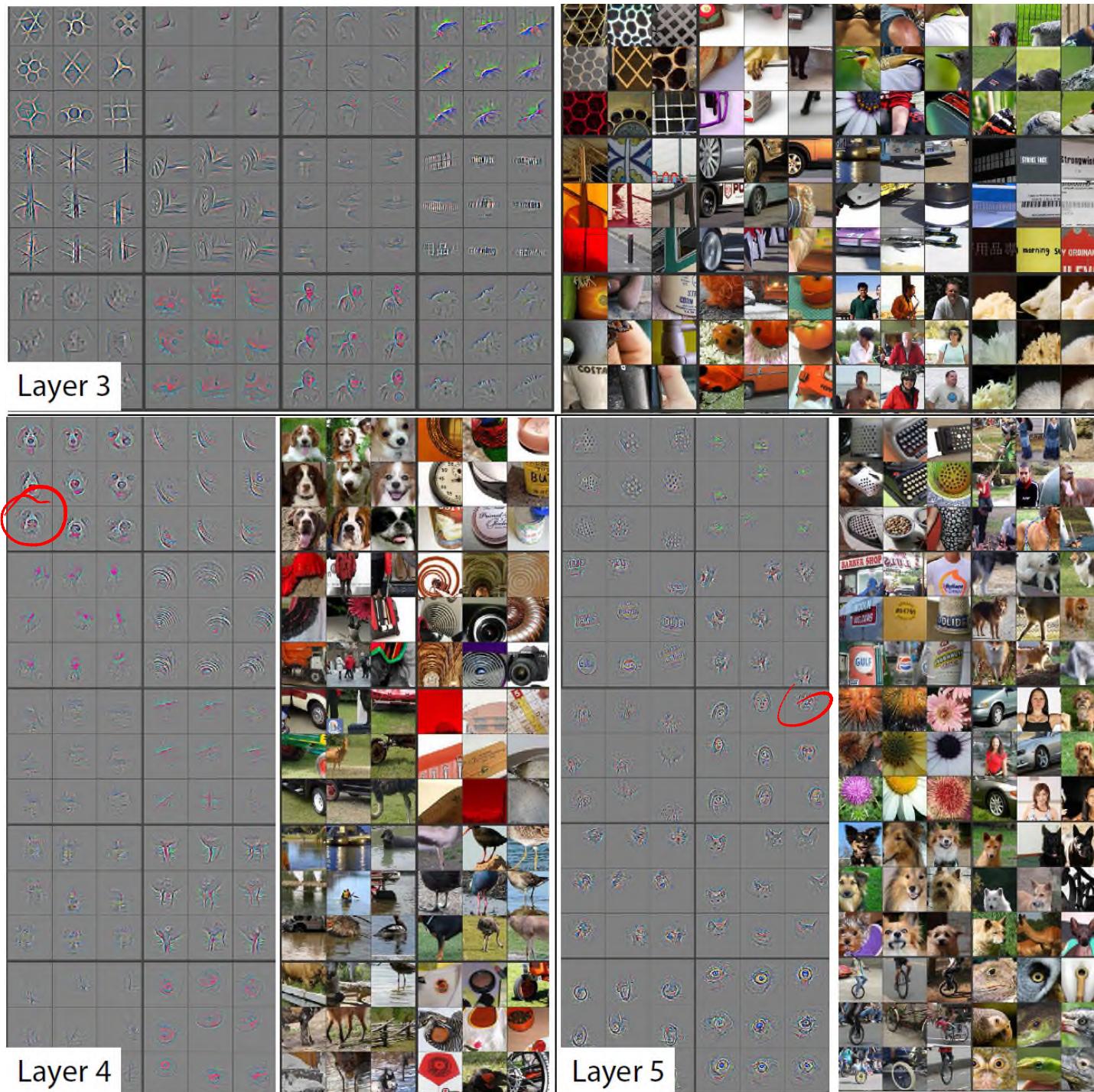
# Convnets (Fukushima, LeCun, Hinton)



# Convolutional networks



[Matthew Zeiler & Rob Fergus]



# Convolution

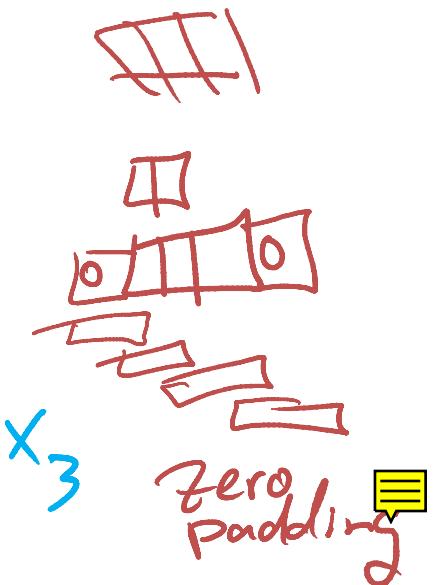
$$Z = \begin{bmatrix} Z_1 & Z_2 \end{bmatrix} \quad M_Z = 2$$

$$W = \begin{bmatrix} w_1 & w_2 \end{bmatrix} \quad \leftarrow M_W = 2$$

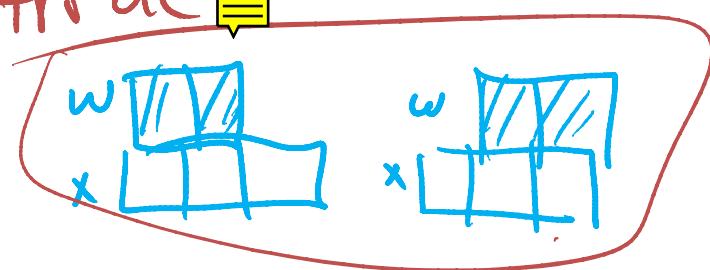
$$X = \begin{bmatrix} x_1 & x_2 & x_3 \end{bmatrix} \quad M_X = 3$$

$$Z_1 = w_1 x_1 + w_2 x_2$$

$$Z_2 = w_1 x_2 + w_2 x_3$$



Stride



$$\text{flip } \bar{W} = \begin{bmatrix} w_2 & w_1 \end{bmatrix}$$

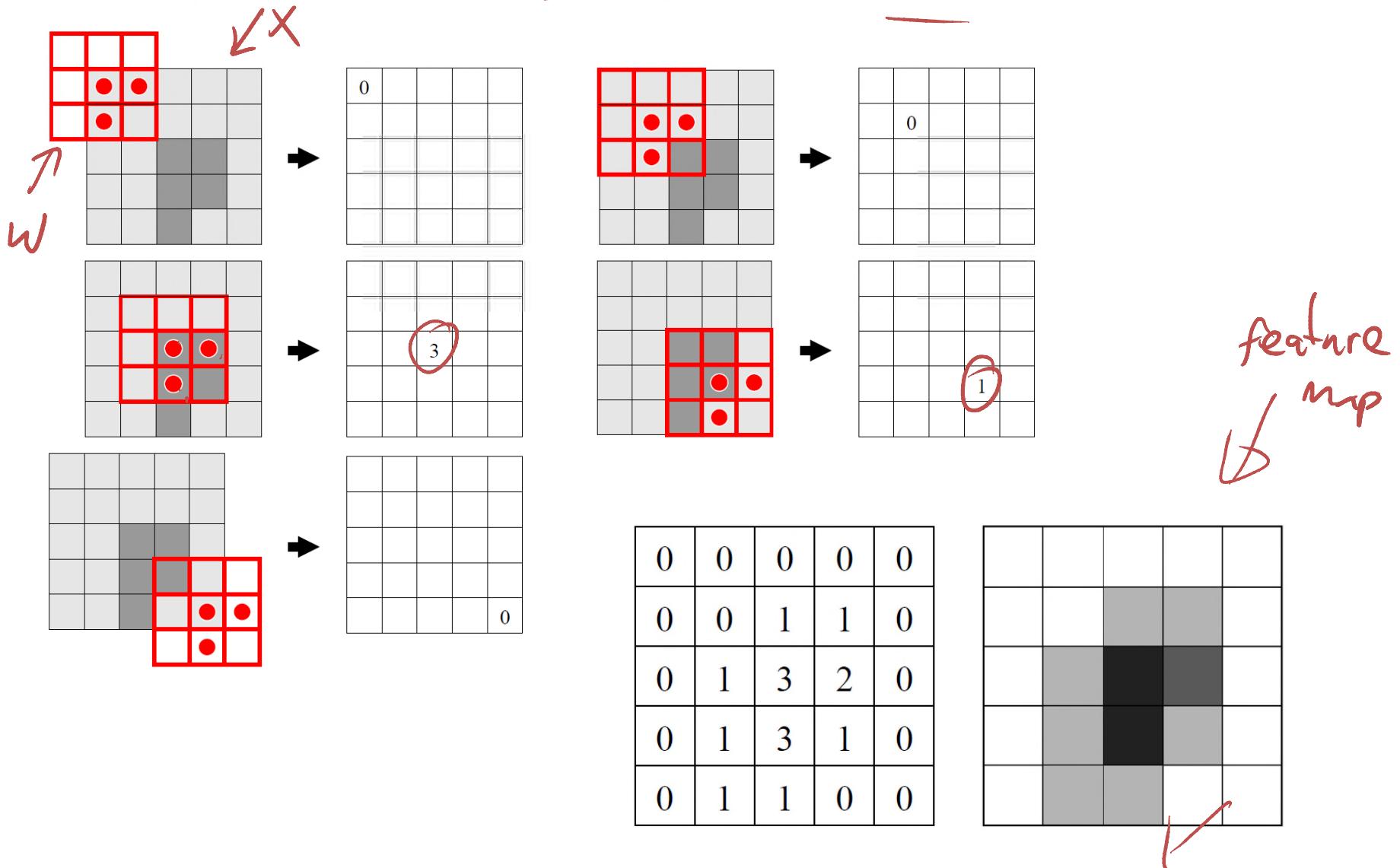
$$Z_{i'} = \sum_{i=1}^{M_F=2} \frac{w_i x_{i'+i-1}}{\parallel} \quad \text{Correlation (Similarity)}$$

$$Z_{i'} = \sum_{i=1}^{M_F=2} \frac{x_{i'+i-1} \bar{W}_{M_F-i+1}}{\parallel} \quad \text{(Convolution)}$$

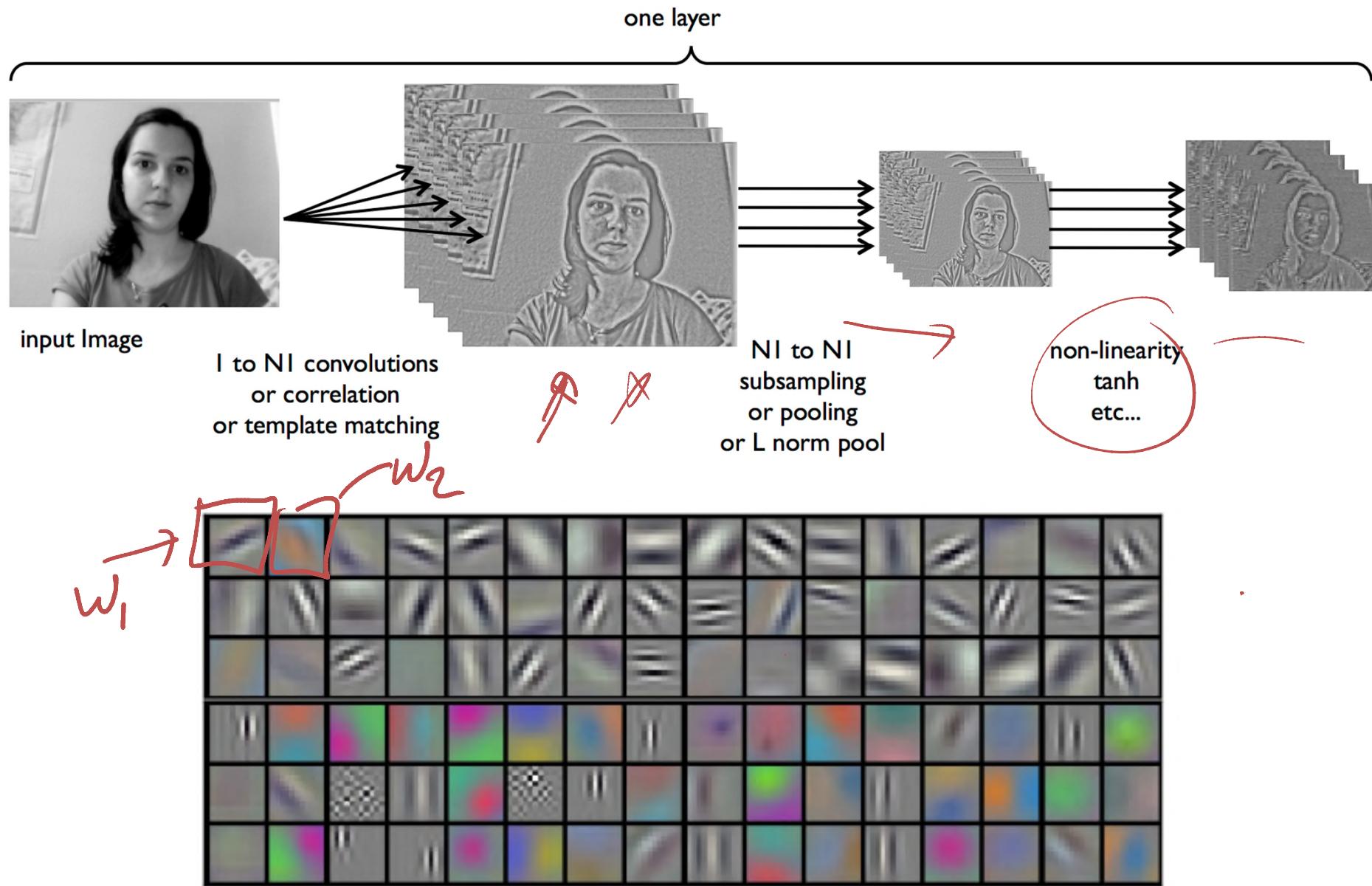
Andrej Karpathy

# Image convolution

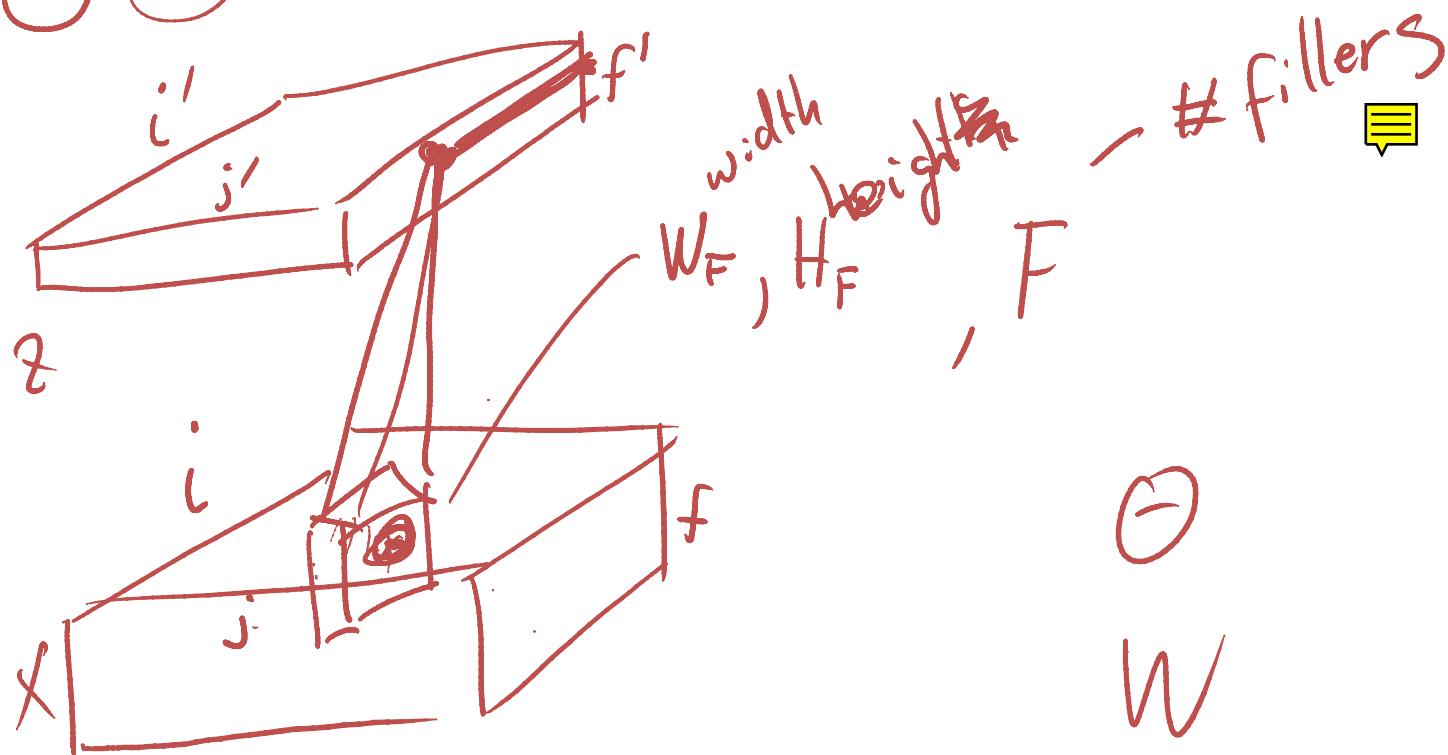
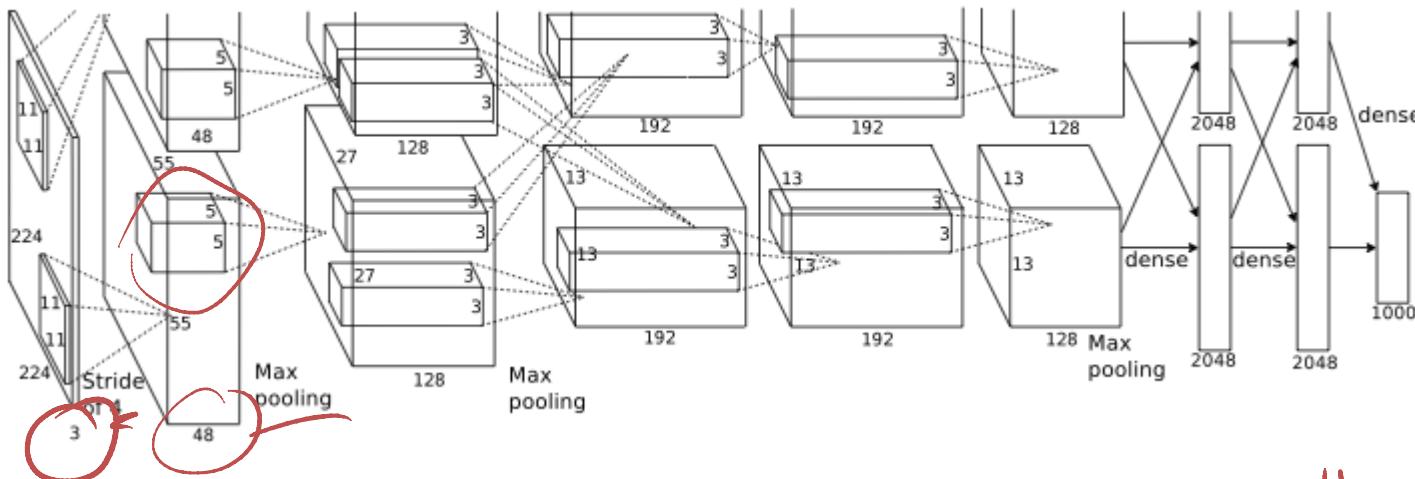
For the image, take dark pixel value = 1, light pixel value = 0.



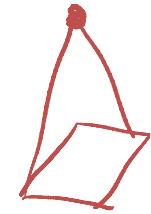
# Convnets (Fukushima, LeCun, Hinton)



# Image convolution layer



# Image convolution layer



$$f = \underbrace{\mathbf{y}_{i', j', f'}}_{\text{Output}} = \underbrace{b_{f'}}_{\text{Bias}} + \underbrace{\sum_{i=1}^{H_f} \sum_{j=1}^{W_f} \sum_{f=1}^F \mathbf{x}_{i'+i-1, j'+j-1, f} \theta_{ijff'}}_{\text{Convolutional Feature Map}}$$

$$\frac{\partial E}{\partial \theta_{ijff'}} = \sum_{i'j'f'} \delta_{i'j'f'}^{l+1} \underbrace{\frac{\partial f_{i'j'f'}(\mathbf{x}; \theta_{f'})}{\partial \theta_{ijff'}}}_{\text{Backpropagation Step}}$$

$$= \sum_{i'j'} \delta_{i'j'f'}^{l+1} \mathbf{x}_{i'+i-1, j'+j-1, f}$$

# Image convolution layer

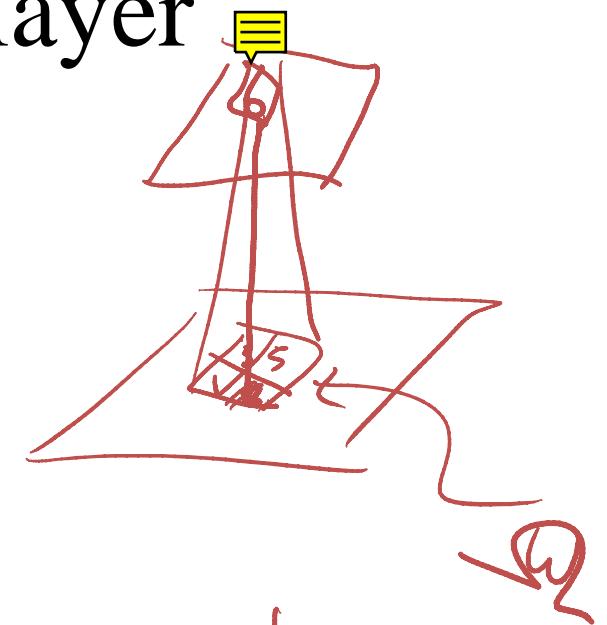
$$\mathbf{y}_{i',j',f'} = b_{f'} + \sum_{i''=1}^{H_f} \sum_{j''=1}^{W_f} \sum_{f''=1}^F \mathbf{x}_{i'+i''-1, j'+j''-1, f''} \theta_{i''j''f''f'}$$

$$\delta_{ijf}^l = \sum_{i'j'f'} \delta_{i'j'f'}^{l+1} \frac{\partial f_{i'j'f'}(\mathbf{x}; \boldsymbol{\theta}_{f'})}{\partial \mathbf{x}_{ijf}}$$
$$i = i' + i'' - 1$$
$$i'' = i - i' + 1$$

$$= \sum_{i'j'f'} \delta_{i'j'f'}^{l+1} \theta_{i-i'+1, j-j'+1, f, f'}$$

# Image max-pooling layer

$$\mathbf{y}_{i',j'} = \max_{ij \in \Omega(i',j')} \mathbf{x}_{ij}$$



$$\delta_{ij}^l = \sum_{i'j'} \delta_{i'j'}^{l+1} \frac{\partial f_{i'j'}(\mathbf{x})}{\partial \mathbf{x}_{ij}}$$

$$= \delta_{ij}^{l+1} \mathbb{I}_{ij = \arg \max_{i''j'' \in \Omega(i',j')} \mathbf{x}_{i''j''}}$$

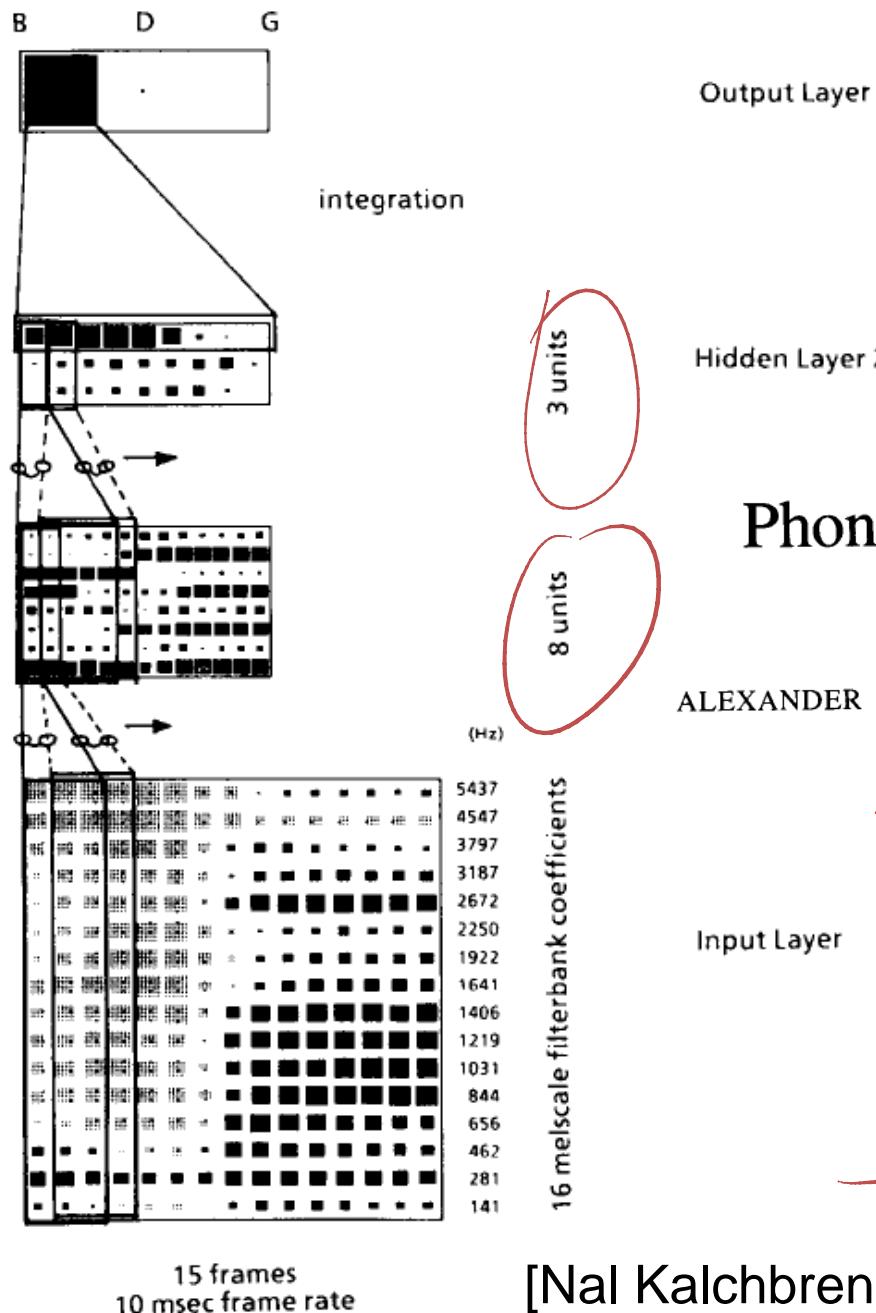
*Capsules Draw*

# Convnets in Torch

Xavier Glorot

```
1 model = nn.Sequential() ✗
2 model:add(nn.Reshape(1,32,32)) ↴
3 -- layer 1:
4 model:add(nn.SpatialConvolution(1, 16, 5, 5))
5 model:add(nn.Tanh())
6 model:add(nn.SpatialMaxPooling(2, 2, 2, 2))
7 -- layer 2:
8 model:add(nn.SpatialConvolution(16, 128, 5, 5))
9 model:add(nn.Tanh())
10 model:add(nn.SpatialMaxPooling(2, 2, 2, 2))
11 -- layer 3, a simple 2-layer neural net:
12 model:add(nn.Reshape(128*5*5))
13 model:add(nn.Linear(128*5*5, 200))
14 model:add(nn.Tanh())
15 model:add(nn.Linear(200,10)) ✗
16 model:add(nn.LogSoftMax()) ↴
```

# ConvNets for Language



## Phoneme Recognition Using Time-Delay Neural Networks

ALEXANDER WAIBEL, MEMBER, IEEE, TOSHIYUKI HANAZAWA, GEOFFREY HINTON,  
KIYOHIRO SHIKANO, MEMBER, IEEE, AND KEVIN J. LANG

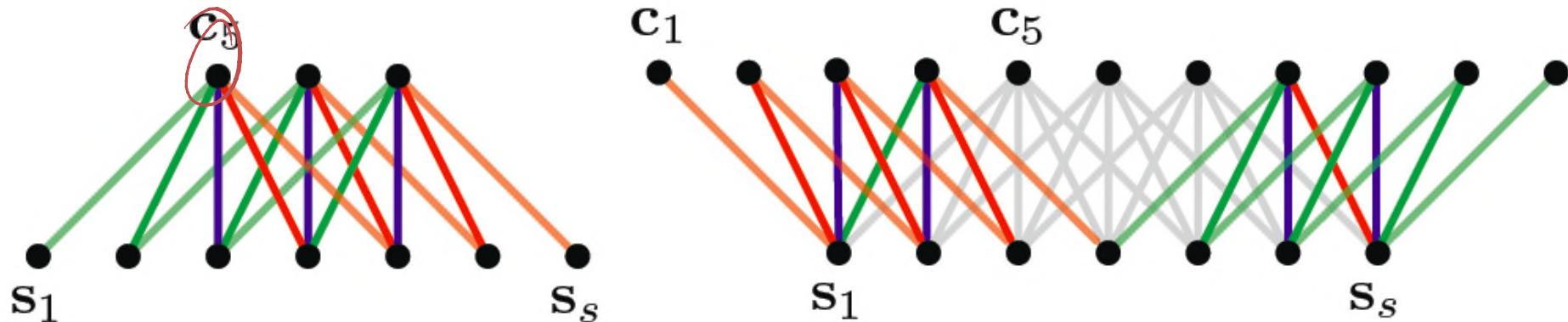
NLP almost from scratch

→ [Ronan Collobert, Jason Weston, 2008]

[Nal Kalchbrenner, Ed Grefenstette, Phil Blunsom, 2014]

Leon Bottou

# Sentence ConvNets



$$\mathbf{c}_j = \mathbf{m}^\top \mathbf{s}_{j-m+1:j}$$

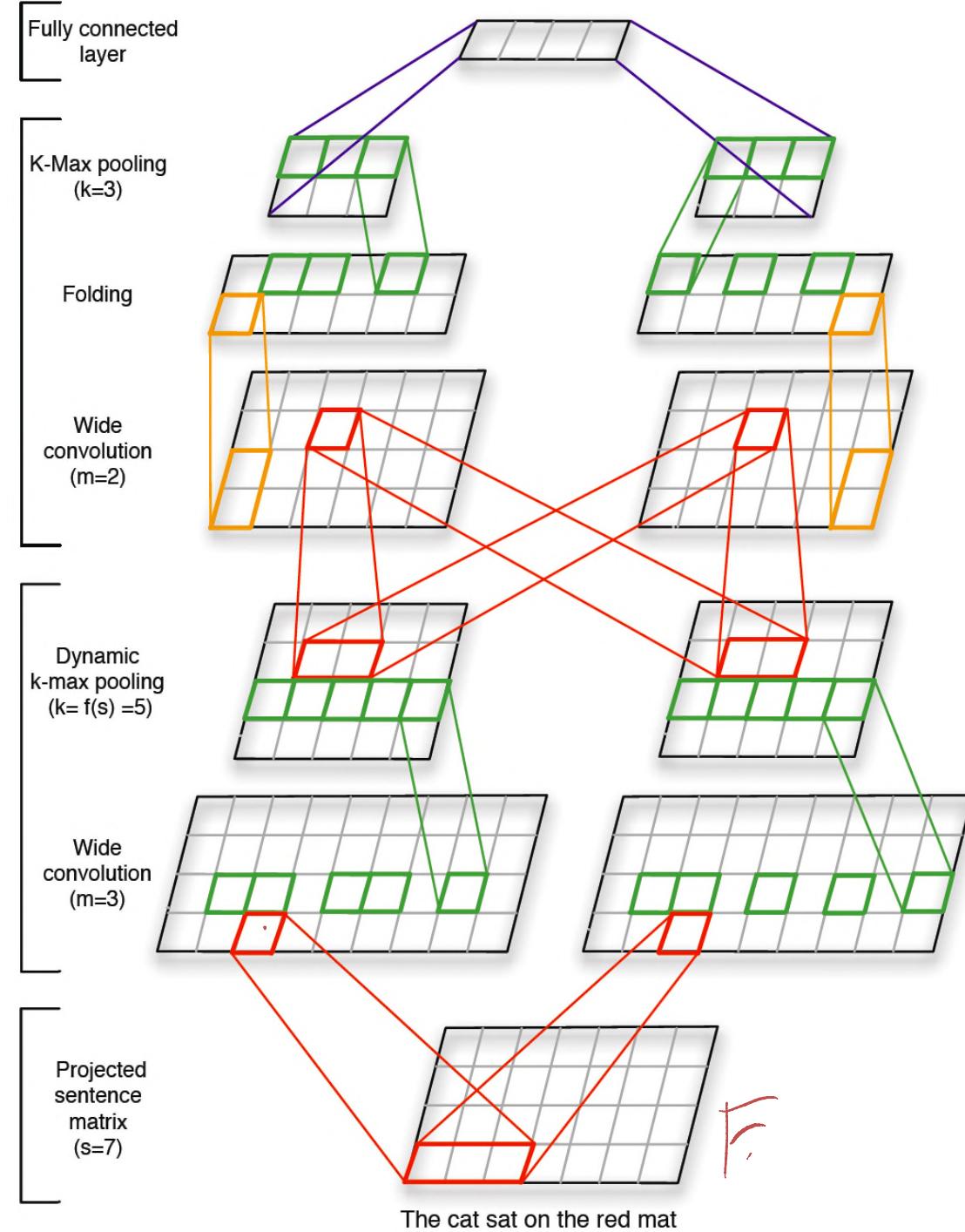
$\mathbf{m}$  is the *filter* of the convolution

$$\mathbf{s} = \begin{bmatrix} | & | & | \\ \mathbf{w}_1 & \dots & \mathbf{w}_s \\ | & | & | \end{bmatrix}$$
$$\mathbf{w}_i \in \mathbb{R}^d$$

$$N_T \begin{bmatrix} \mathbf{w}_1 & \dots & \mathbf{w}_s \end{bmatrix} \begin{bmatrix} 0 & 0 \\ 0 & 1 \\ 0 & 0 \\ \vdots & \vdots \\ 0 & 0 \end{bmatrix} \leftarrow \text{cat}$$

$\leftarrow$  onehot  
encoding

[Kalchbrenner, Grefenstette, Blunsom, 2014]



# Sentence DynConvNet

[Kalchbrenner et al, 2014]

# Document models ( Misha Denil )

