



Practical Machine Learning

Lecture 8 Convolutional Neural Networks (CNN)

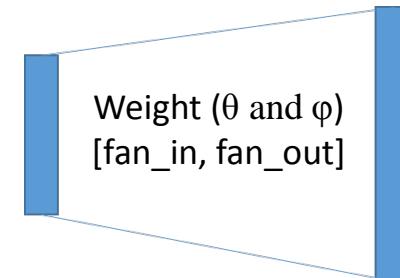
Dr. Suyong Eum



A question from the last class: initialization of parameters

- Xavier (2010): <http://proceedings.mlr.press/v9/glorot10a/glorot10a.pdf>
 - Either using an uniform distribution or a normal distribution
 - It is implemented in the current Tensorflow: `tf.contrib.layers.xavier_initializer()`

$$U\left(-\sqrt{\frac{6}{\text{fan_in} + \text{fan_out}}}, \sqrt{\frac{6}{\text{fan_in} + \text{fan_out}}}\right), N\left(0, \frac{2}{\text{fan_in} + \text{fan_out}}\right)$$



- He (2015): <https://arxiv.org/pdf/1502.01852.pdf>
 - Using a normal distribution

$$N\left(0, \frac{2}{\text{fan_in}}\right)$$

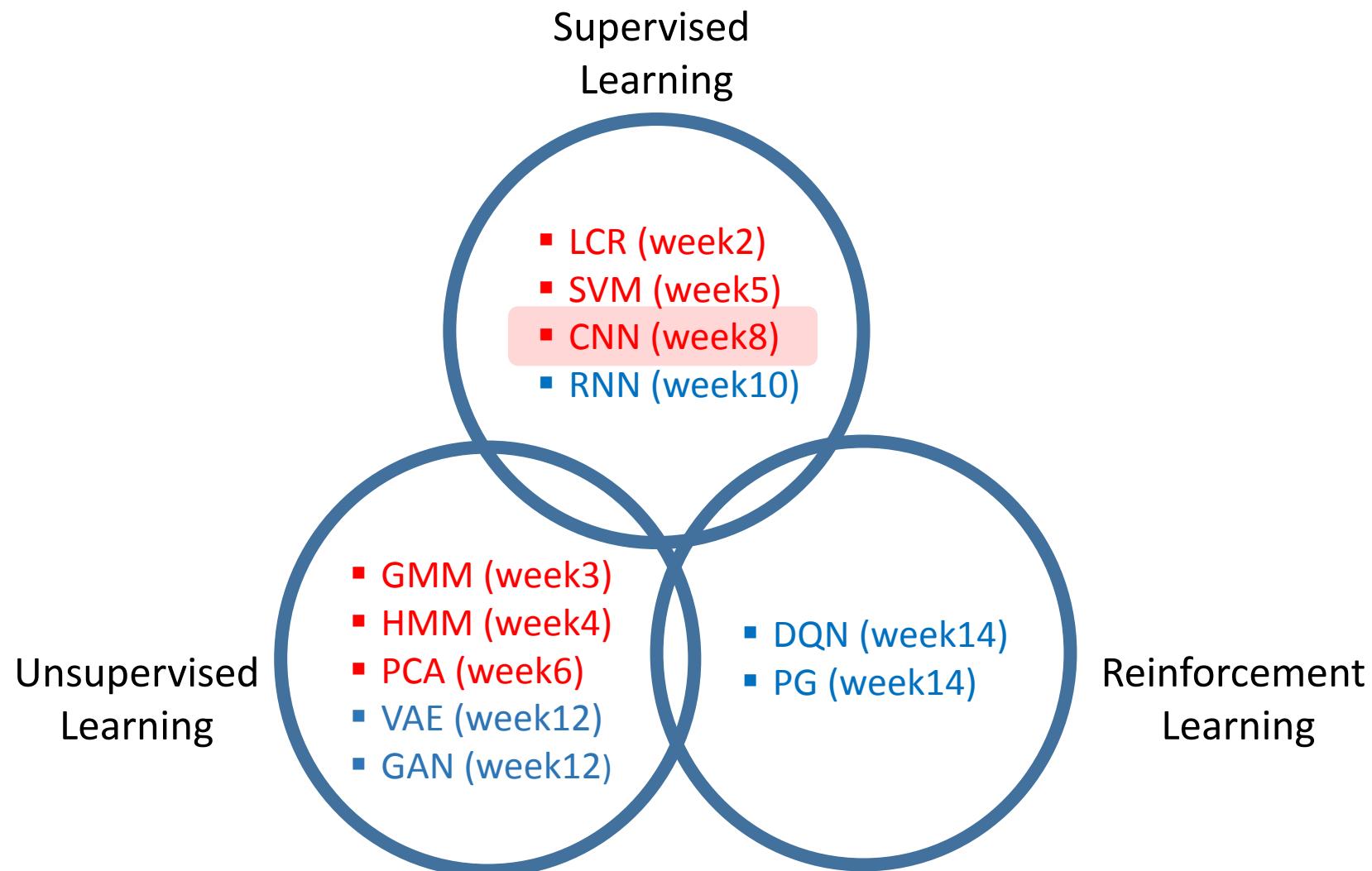
```
def xavier_init(size):
    in_dim = size[0]
    xavier_stddev = 1. / tf.sqrt(in_dim / 2.)
    return tf.random_normal(shape=size, stddev=xavier_stddev)
```

in_dimension
: fan_in

out_dimension
: fan_out

```
def xavier_init(fan_in, fan_out, constant=1):
    low = -constant*np.sqrt(6.0/(fan_in + fan_out))
    high = constant*np.sqrt(6.0/(fan_in + fan_out))
    return tf.random_uniform((fan_in, fan_out),
                           minval=low, maxval=high,
                           dtype=tf.float32)
```

Where we are

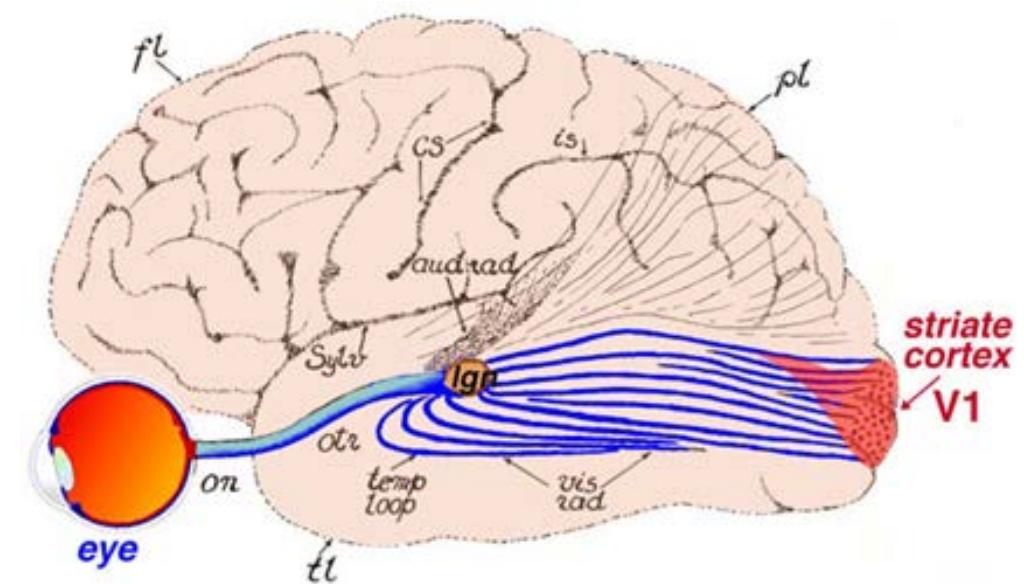
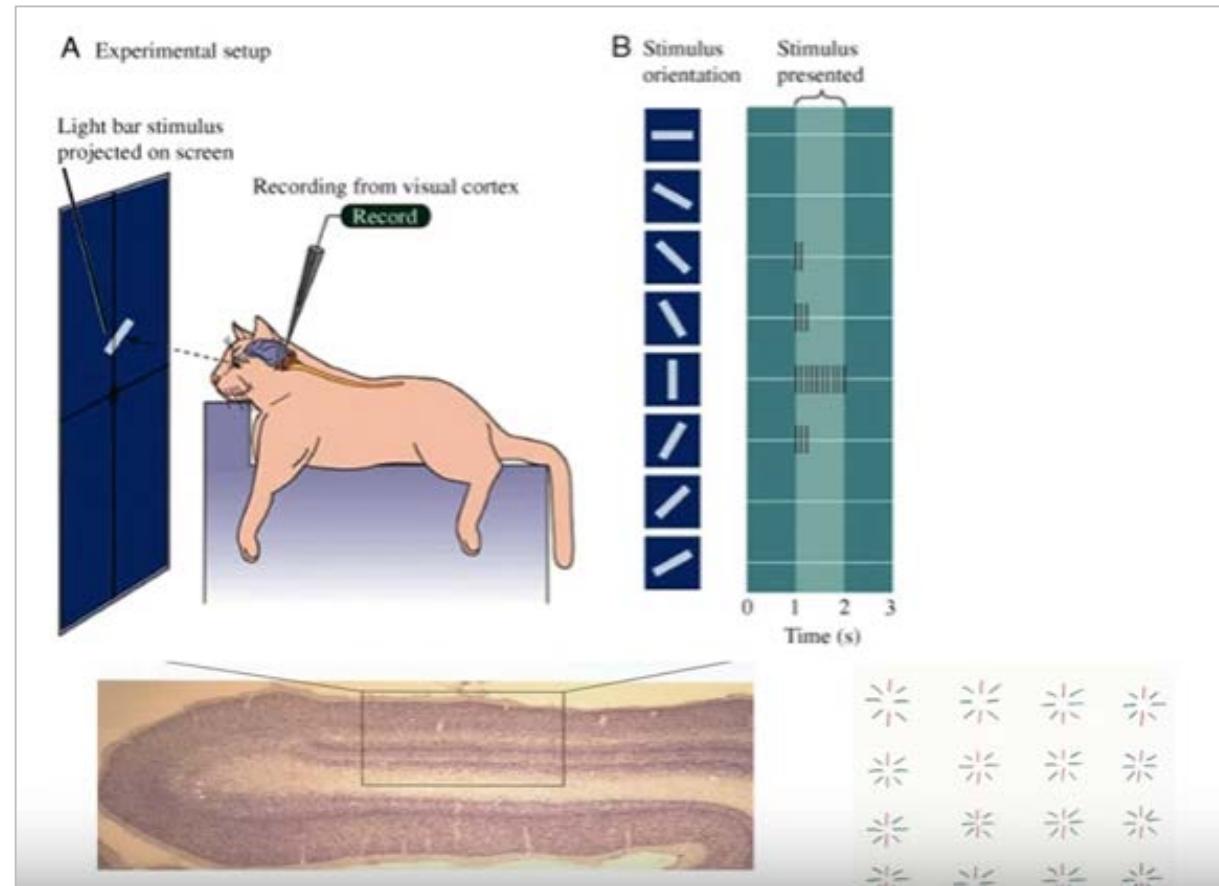


You are going to learn

- What a Convolutional Neural Network (CNN) is,
- How they work,
- How to train CNNs using backpropagation.

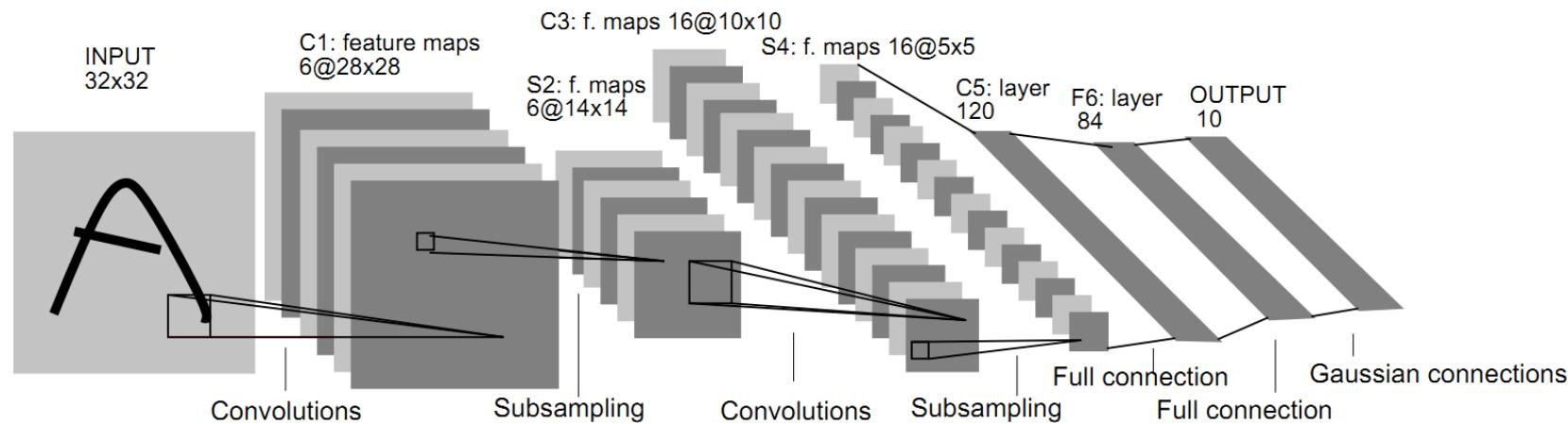
History of CNN

□ Hubel and Wiesel in 1962

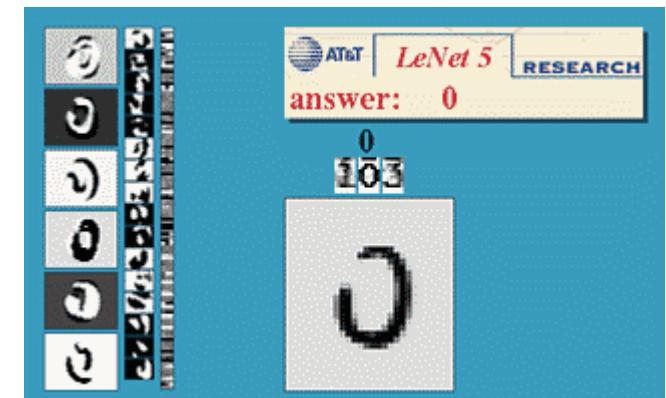


History of CNN – 1995

- LeCun et al, “Convolutional Networks for Images, Speech, and Time-Series”
 - <http://yann.lecun.com/exdb/publis/pdf/lecun-bengio-95a.pdf>



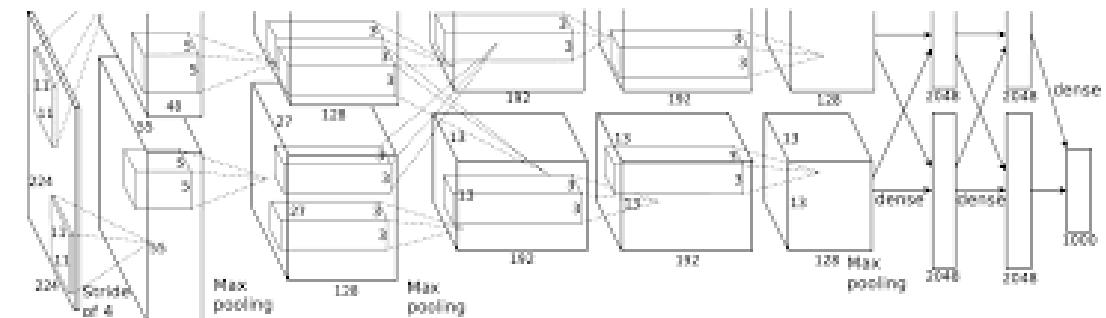
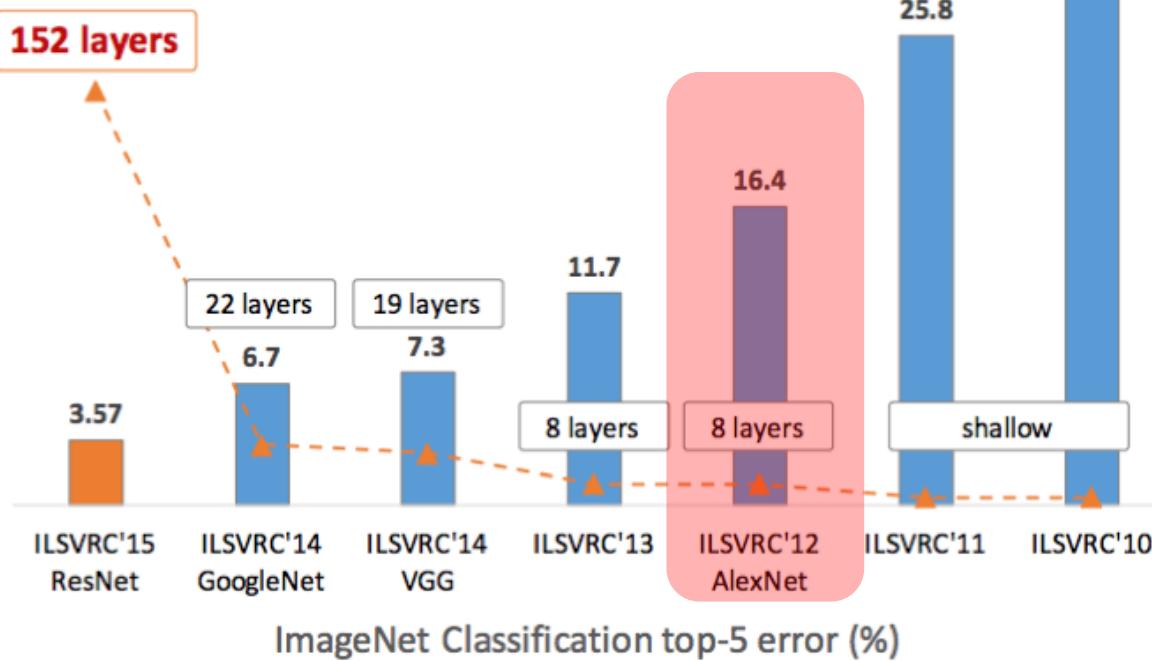
Lecun Yann
Director of AI Research, Facebook



History of CNN – 2012

- ❑ Alex Krizhevsky et al, “*ImageNet Classification with Deep Convolutional Neural Networks*” (2012)
 - Win the imageNet competition: annual Olympics of computer vision with astounding results compared to previously existing approaches (26% to 15%).

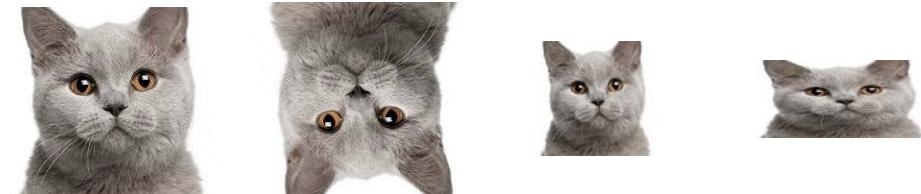
Revolution of Depth



Why Convolutional Neural Networks (CNN)?

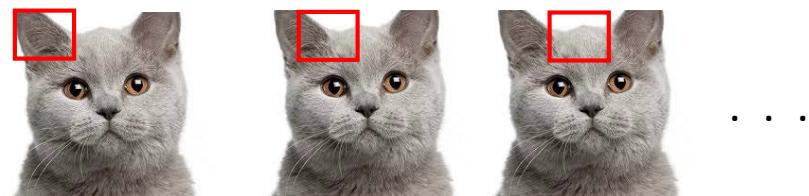
- Invariance to rotation, transformation, size, etc.

0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	9	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9



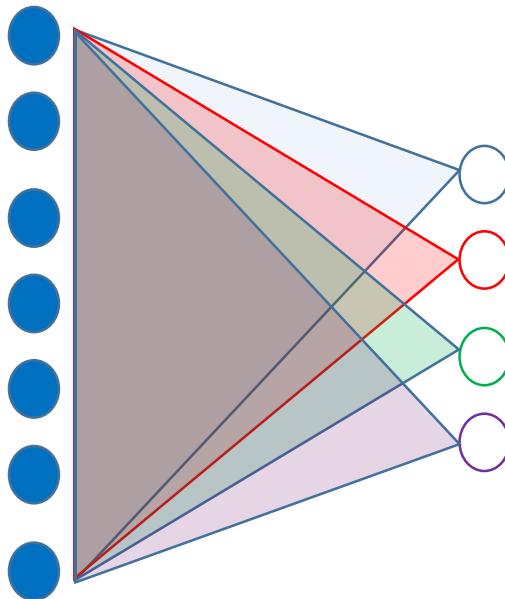
- Three architectural ideas of CNN

1. Local receptive fields
2. Weight sharing
3. Spatial or temporal subsampling



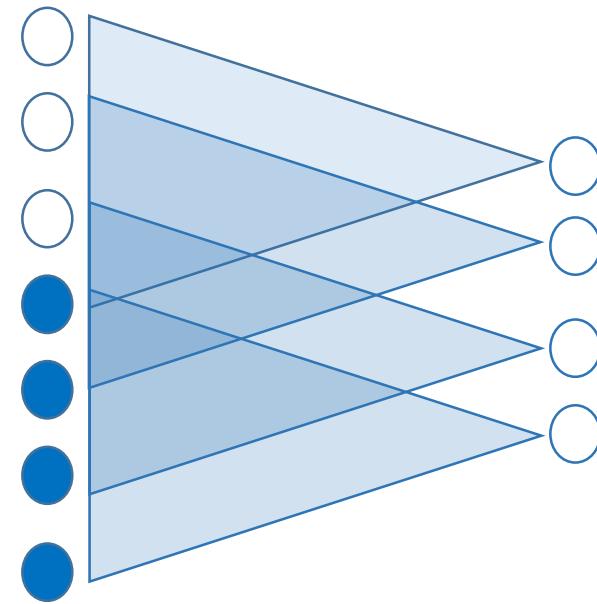
An idea of CNN

- A node takes responsibility for a portion of input data.



Fully connected Layer

- $4 \times 7 = 28$ weights

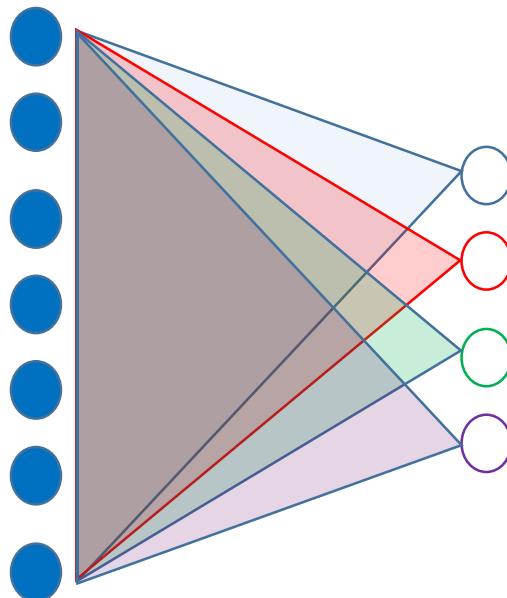


Convolutional Layer

- 4 weights

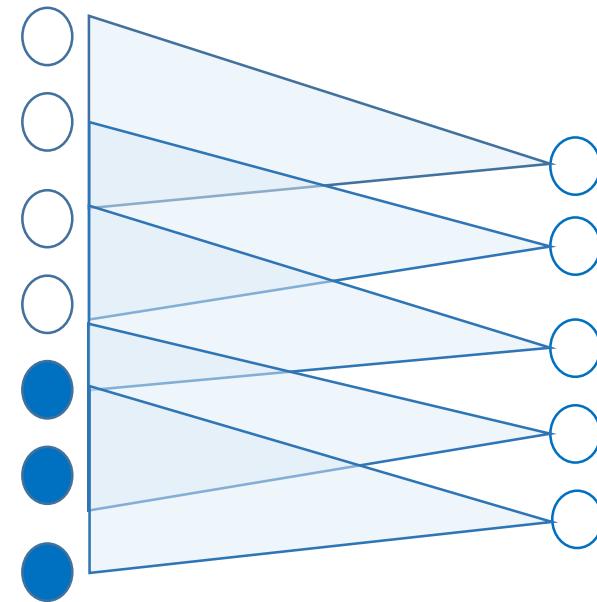
An idea of CNN

- A node takes responsibility for a portion of input data.



Fully connected Layer

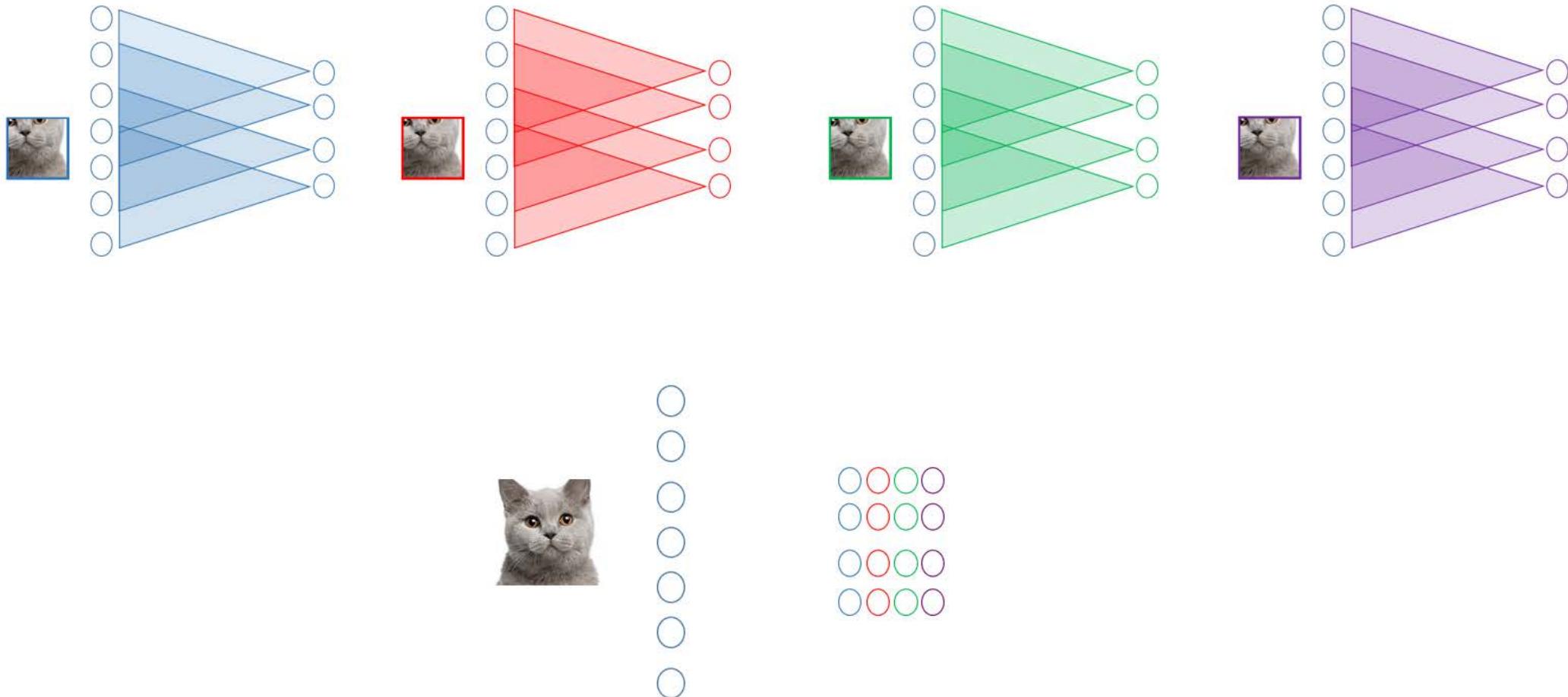
- $4 \times 7 = 28$ weights



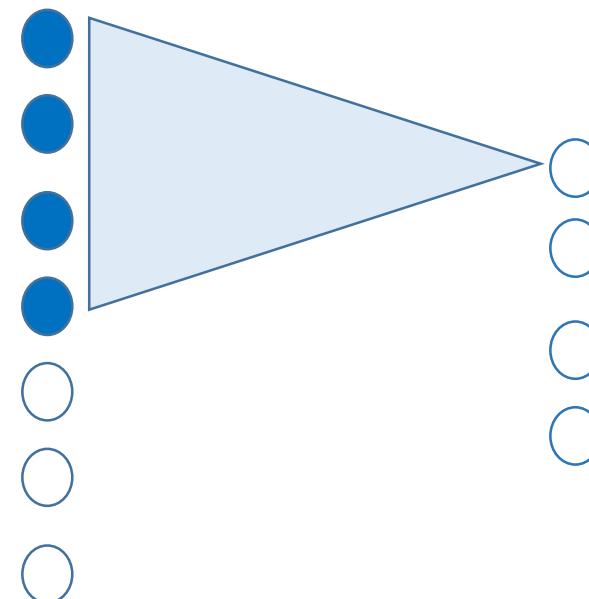
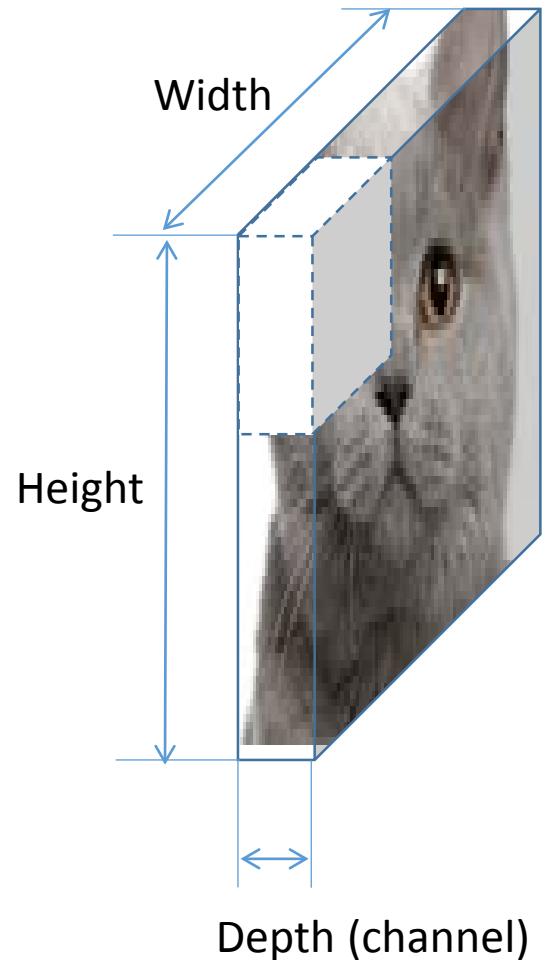
Convolutional Layer

- 3 weights

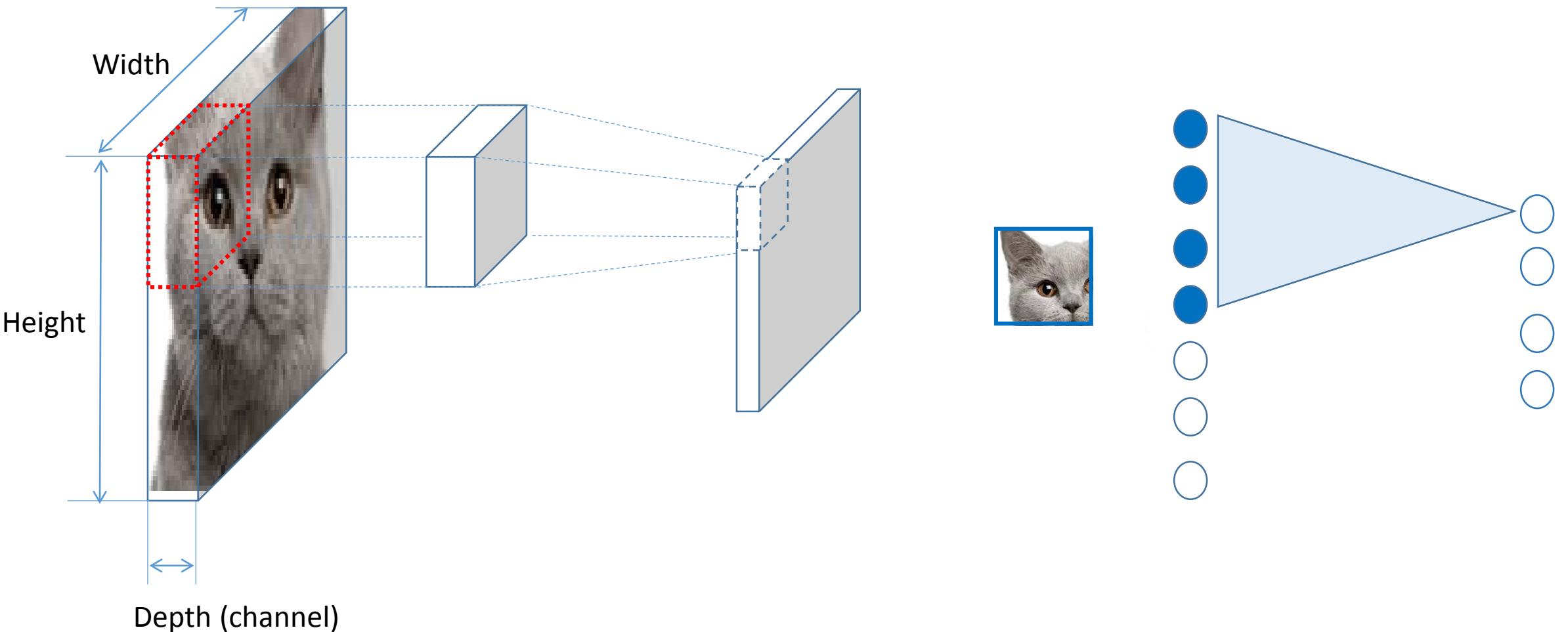
Convolution Layer tends to produce multiple results



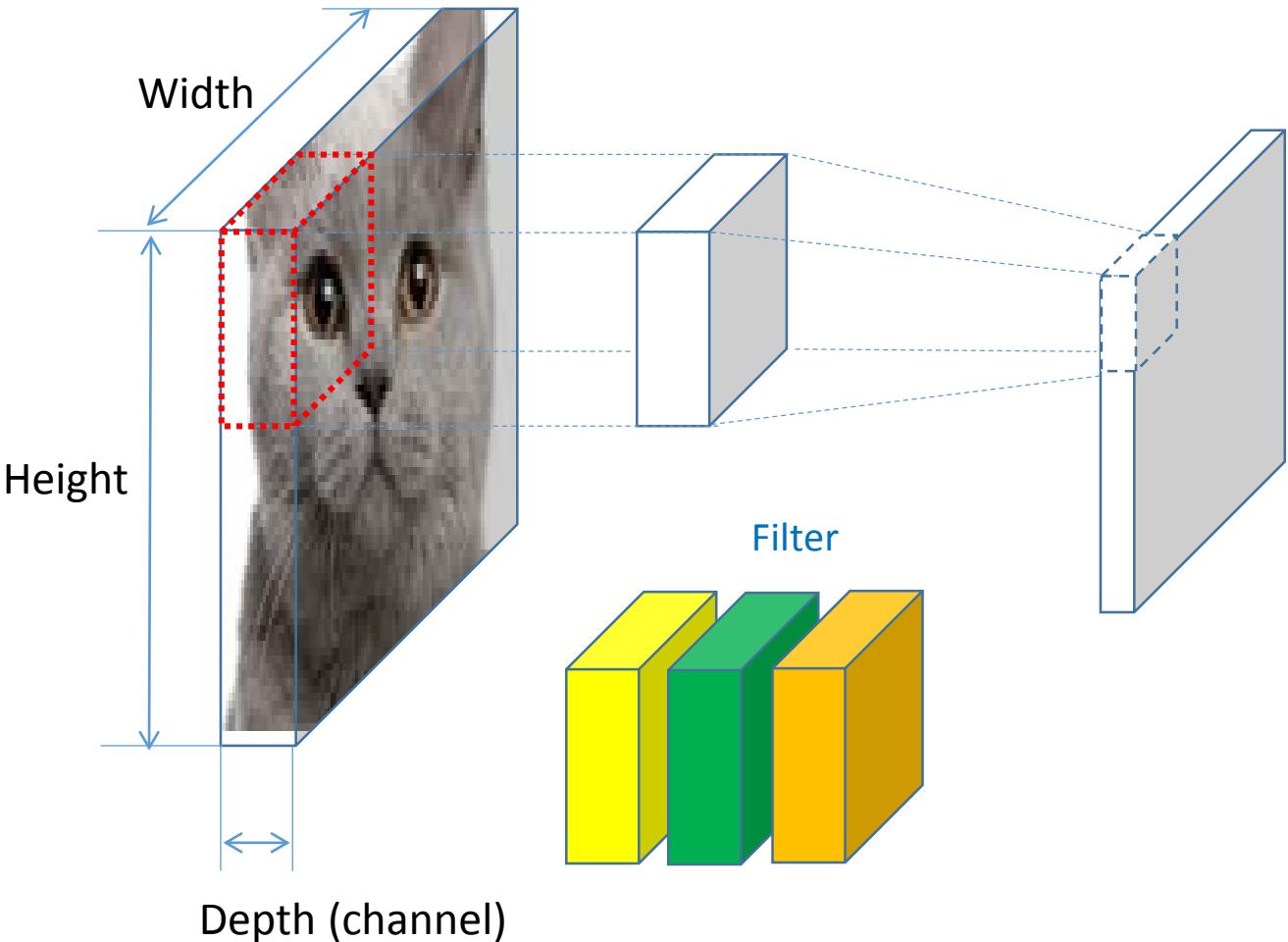
2-D representation vs 1-D representation



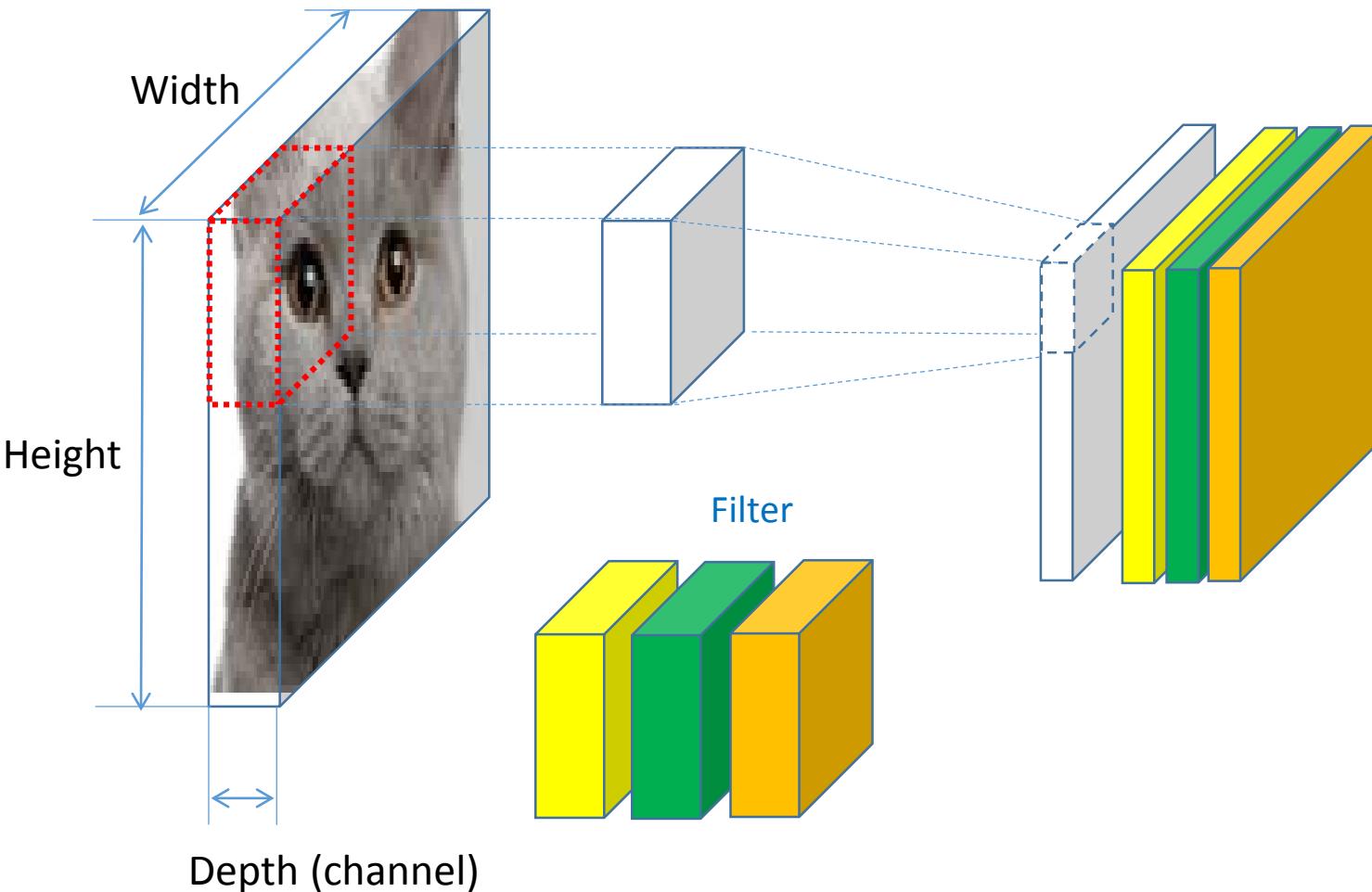
2-D representation vs 1-D representation



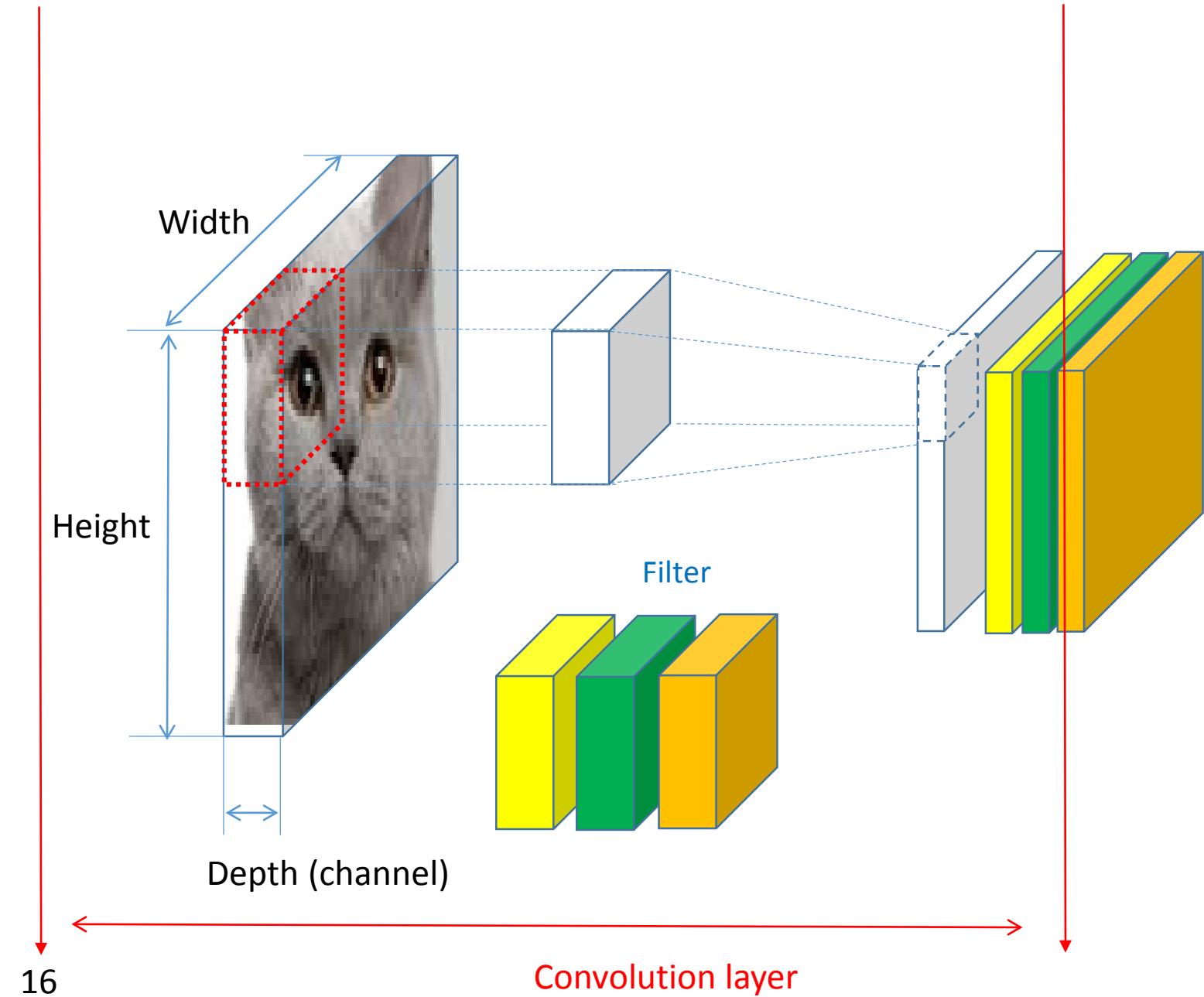
2-D representation of CNN: How does it work?



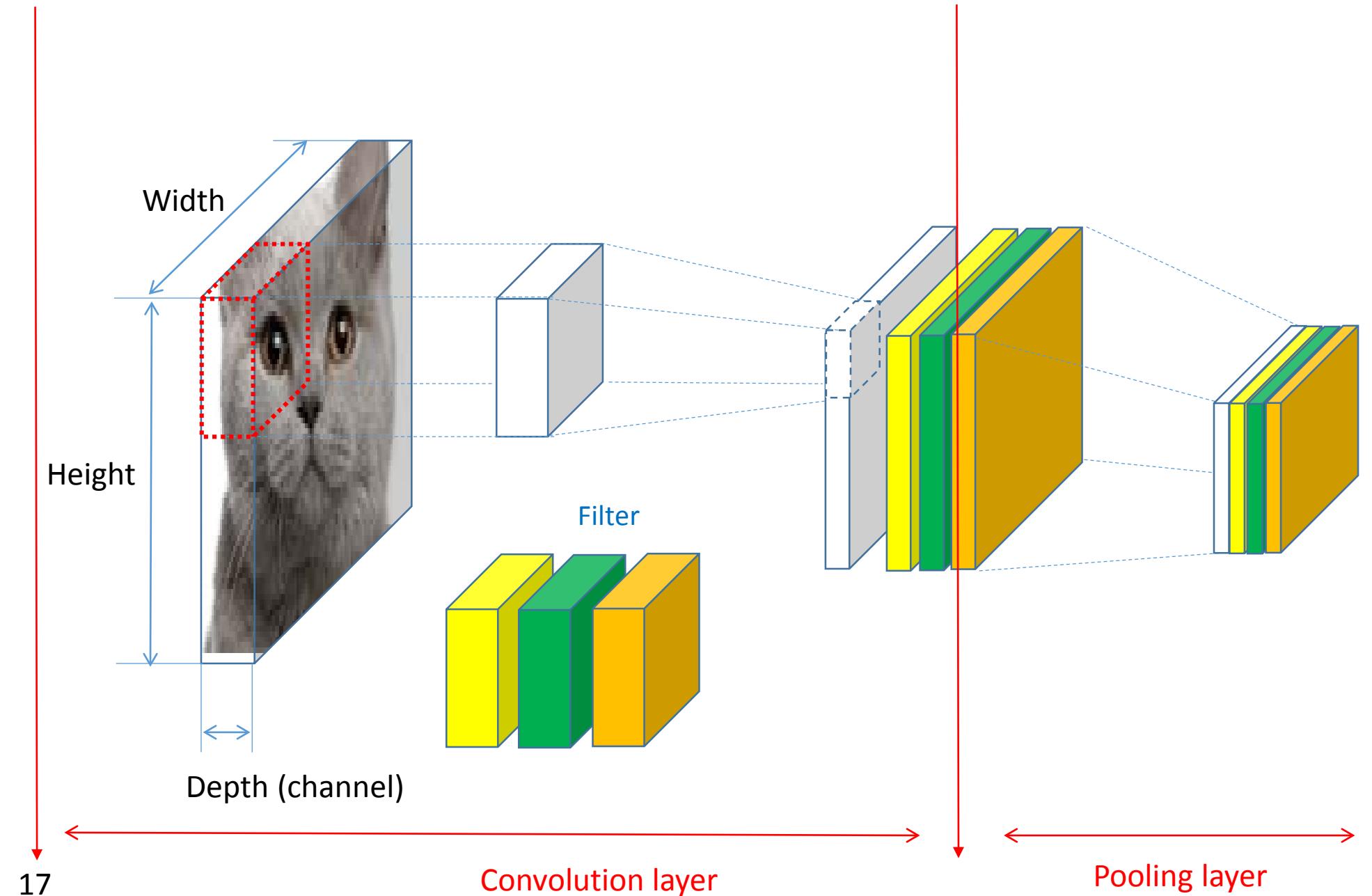
2-D representation of CNN: How does it work?



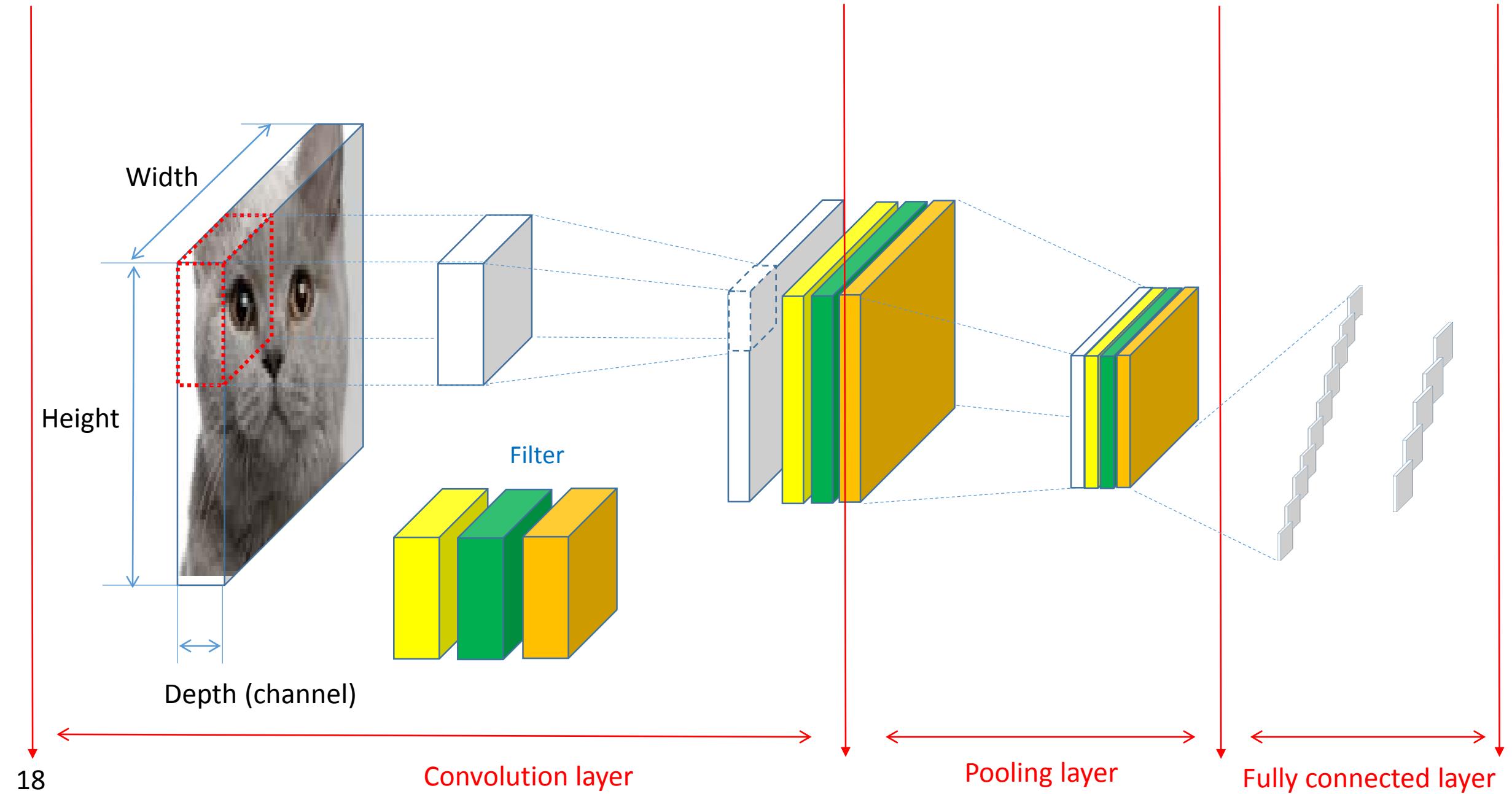
2-D representation of CNN: How does it work?



2-D representation of CNN: How does it work?

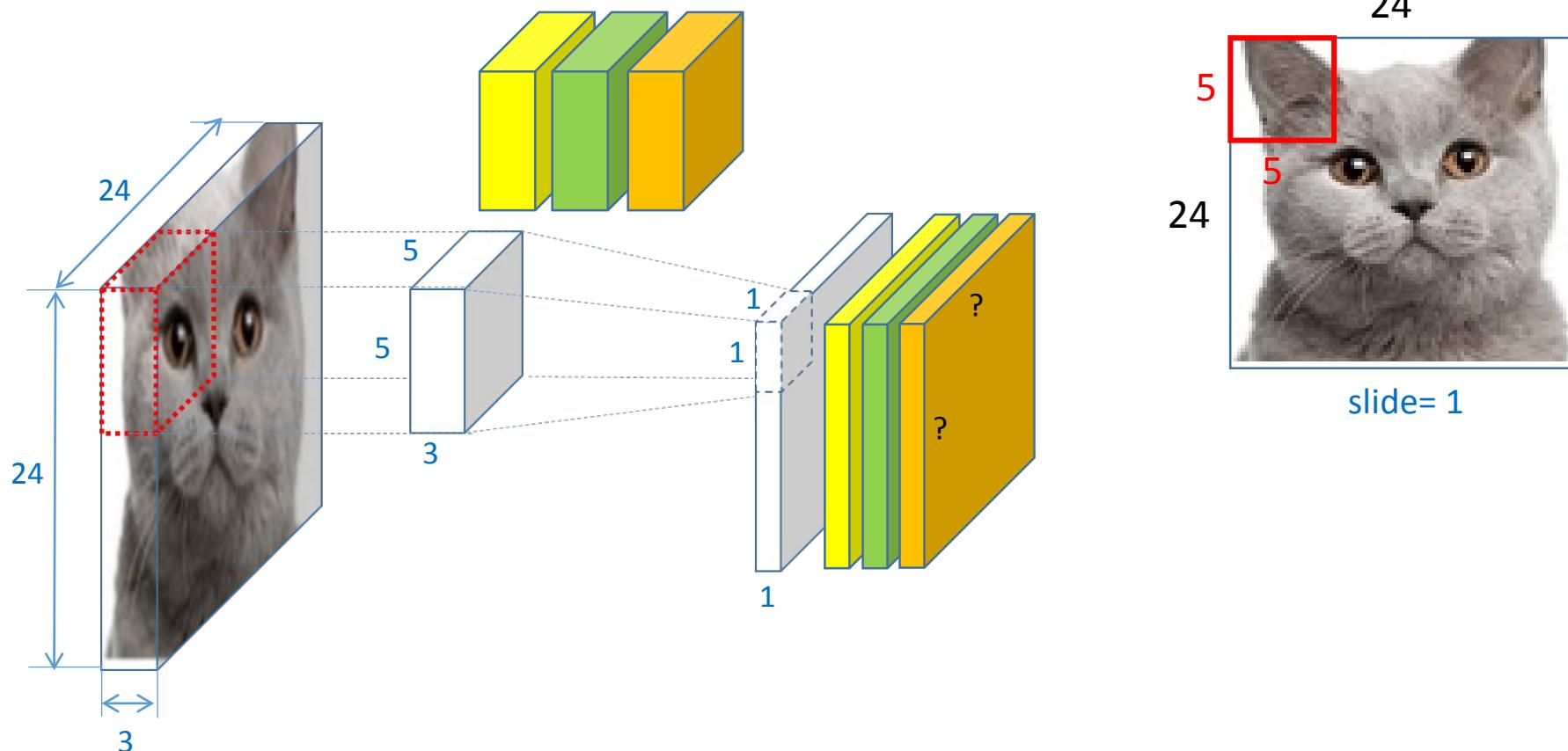


2-D representation of CNN: How does it work?



Convolution Layer

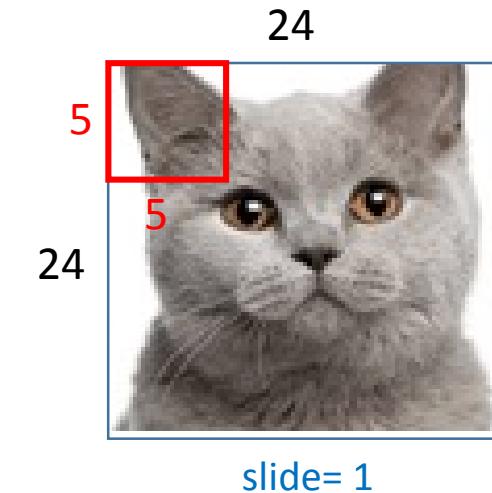
Convolutional Layer



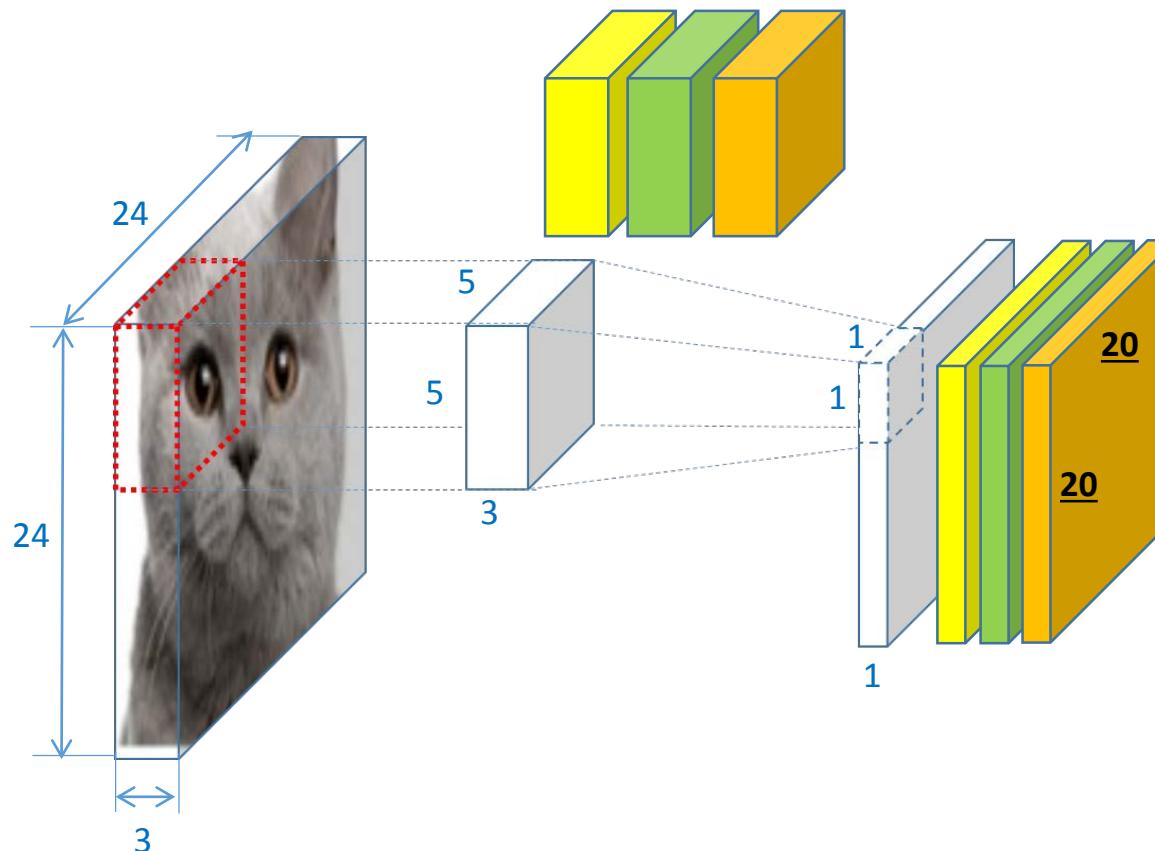
Input layer

Filters
Kernels

Activation maps
Feature maps



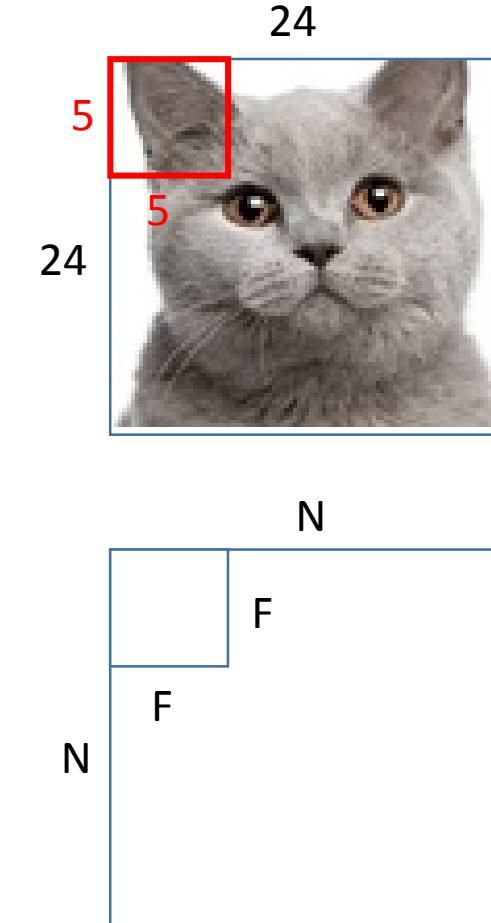
Convolutional Layer



Input layer

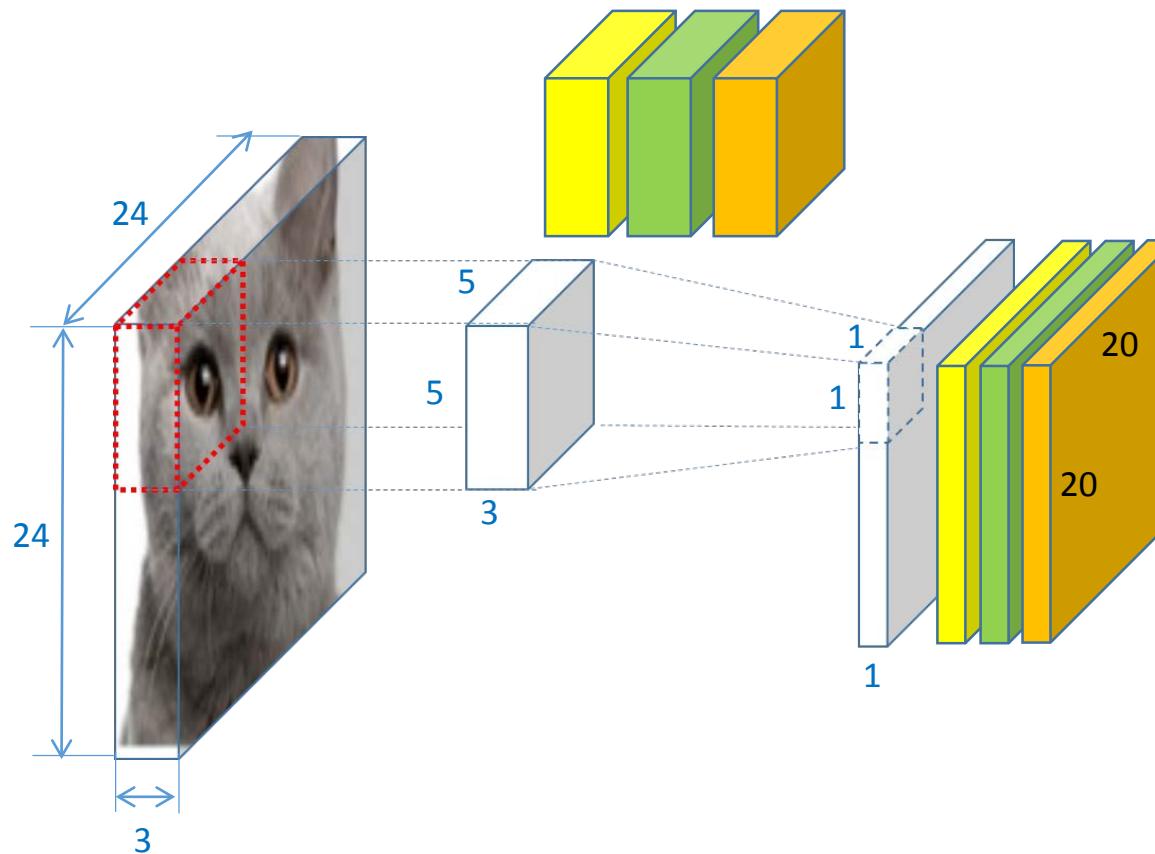
Filters
Kernels

Activation maps
Feature maps



$$\text{Size of activation map} = \frac{(N - F)}{\text{stride}} + 1$$

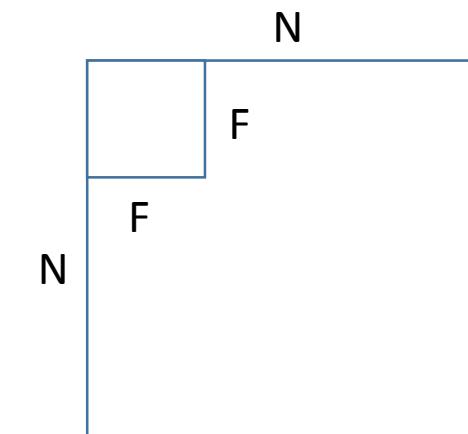
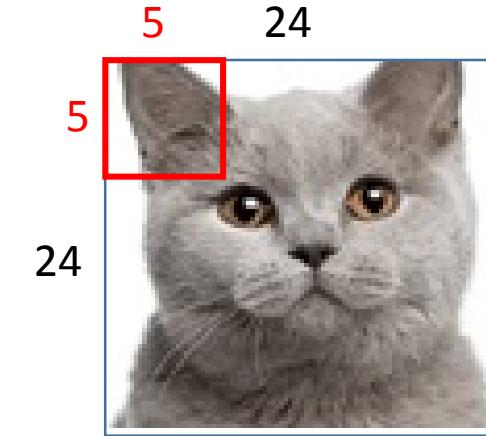
Convolutional Layer



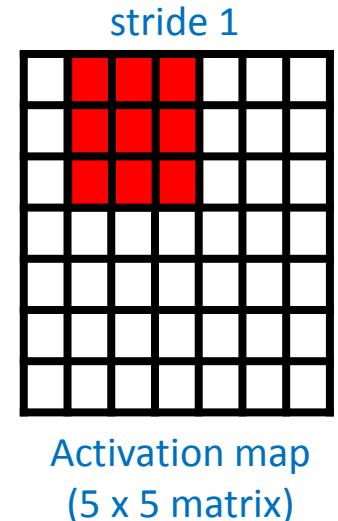
Input layer

Filters
Kernels

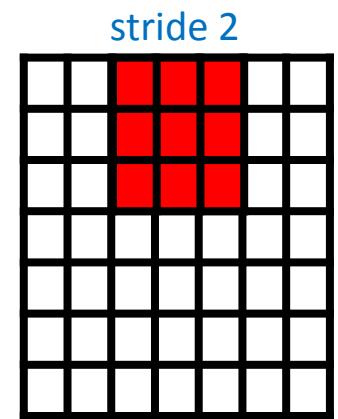
Activation maps
Feature maps



$$\text{Size of activation map} = \frac{(N - F)}{\text{stride}} + 1$$



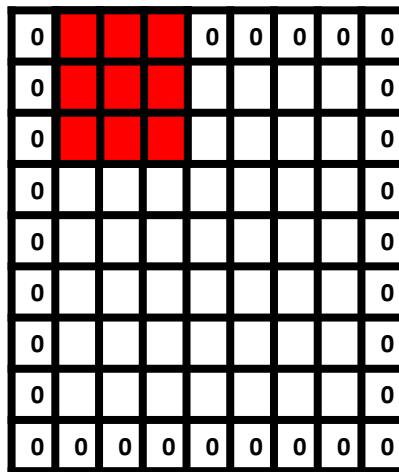
Activation map
(5 x 5 matrix)



Activation map
(3 x 3 matrix)

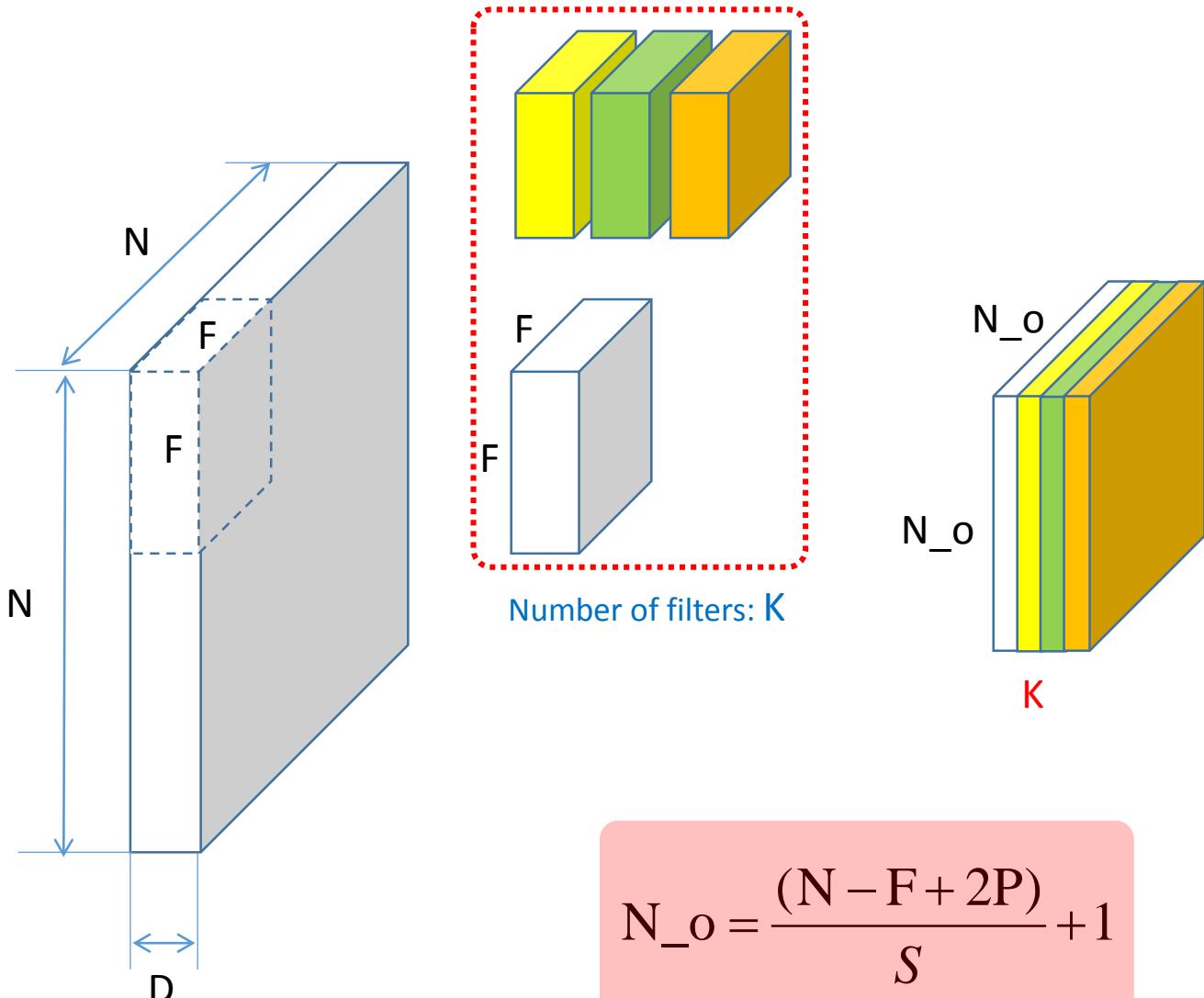
Convolutional Layer

Hyper parameters	Symbols
Number of filters	K
Size of the filter	F
Stride	S
Padding	P



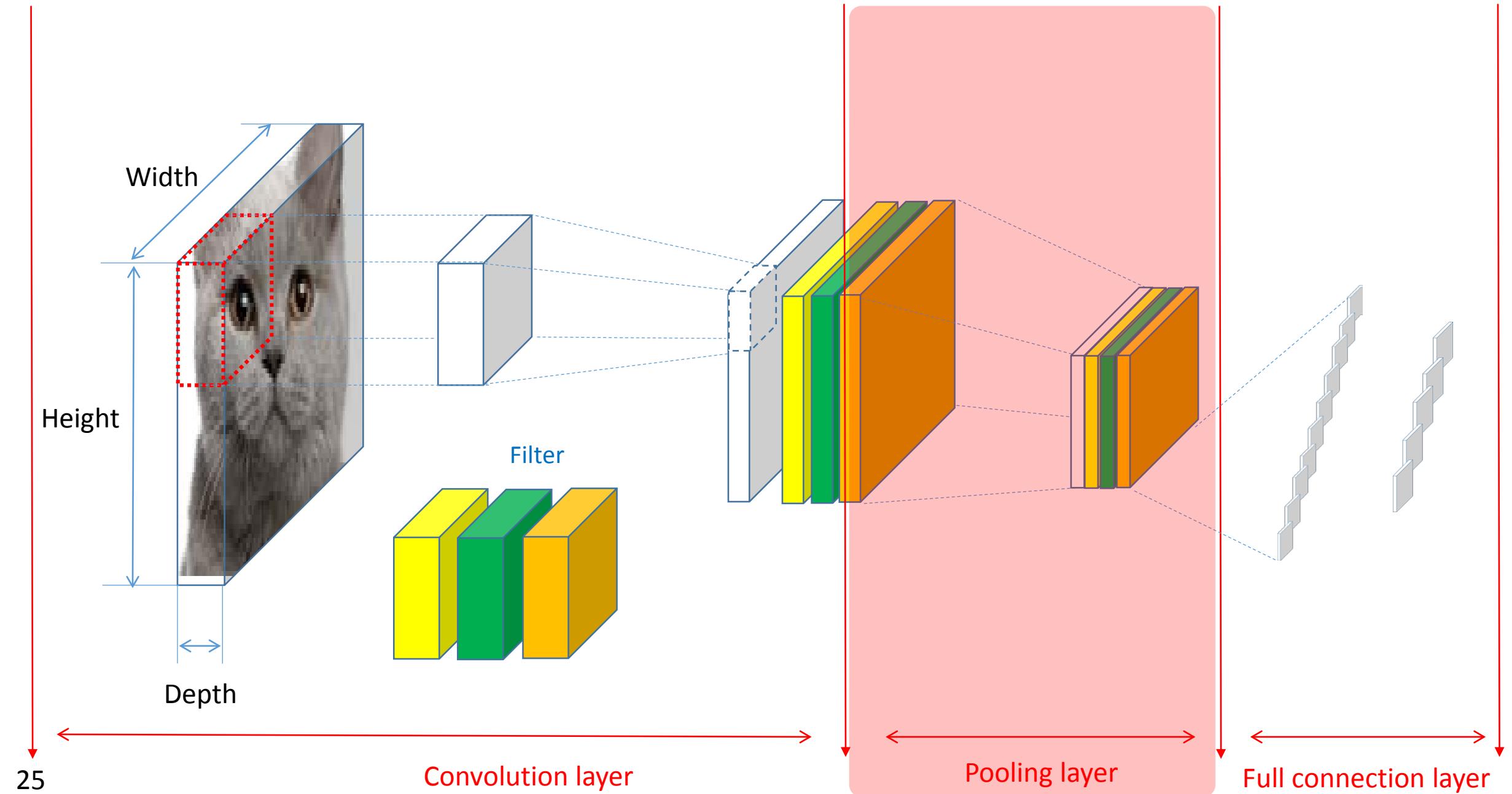
- Padding aims to maintain the original dimension of the original data.

- Padding: P=1
- Stride: S=1
- activation map becomes is 7 x 7 matrix



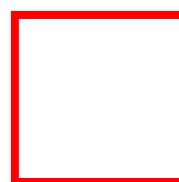
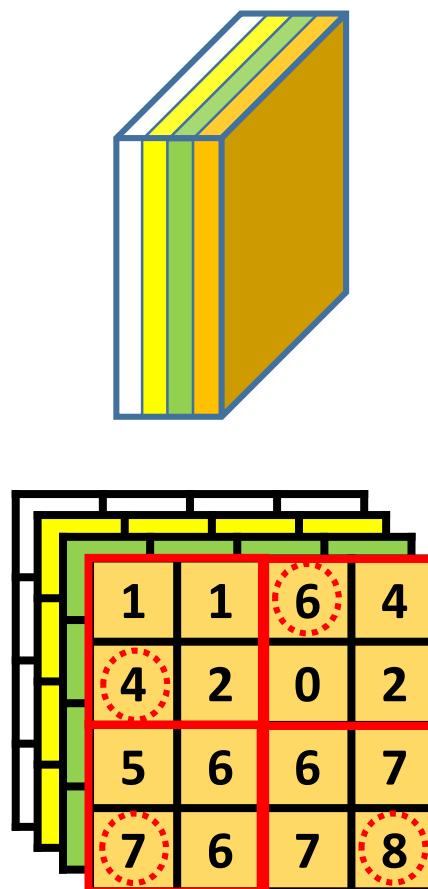
Pooling Layer

2-D representation of CNN: How does it work?



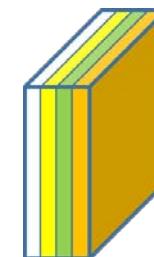
Pooling Layer

- This layer aims to reduce the dimension of each activation map.

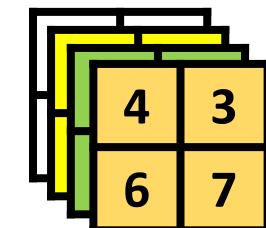
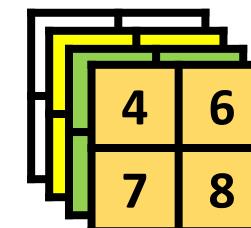


- 2x2 filters
- No overlapped

Max pooling

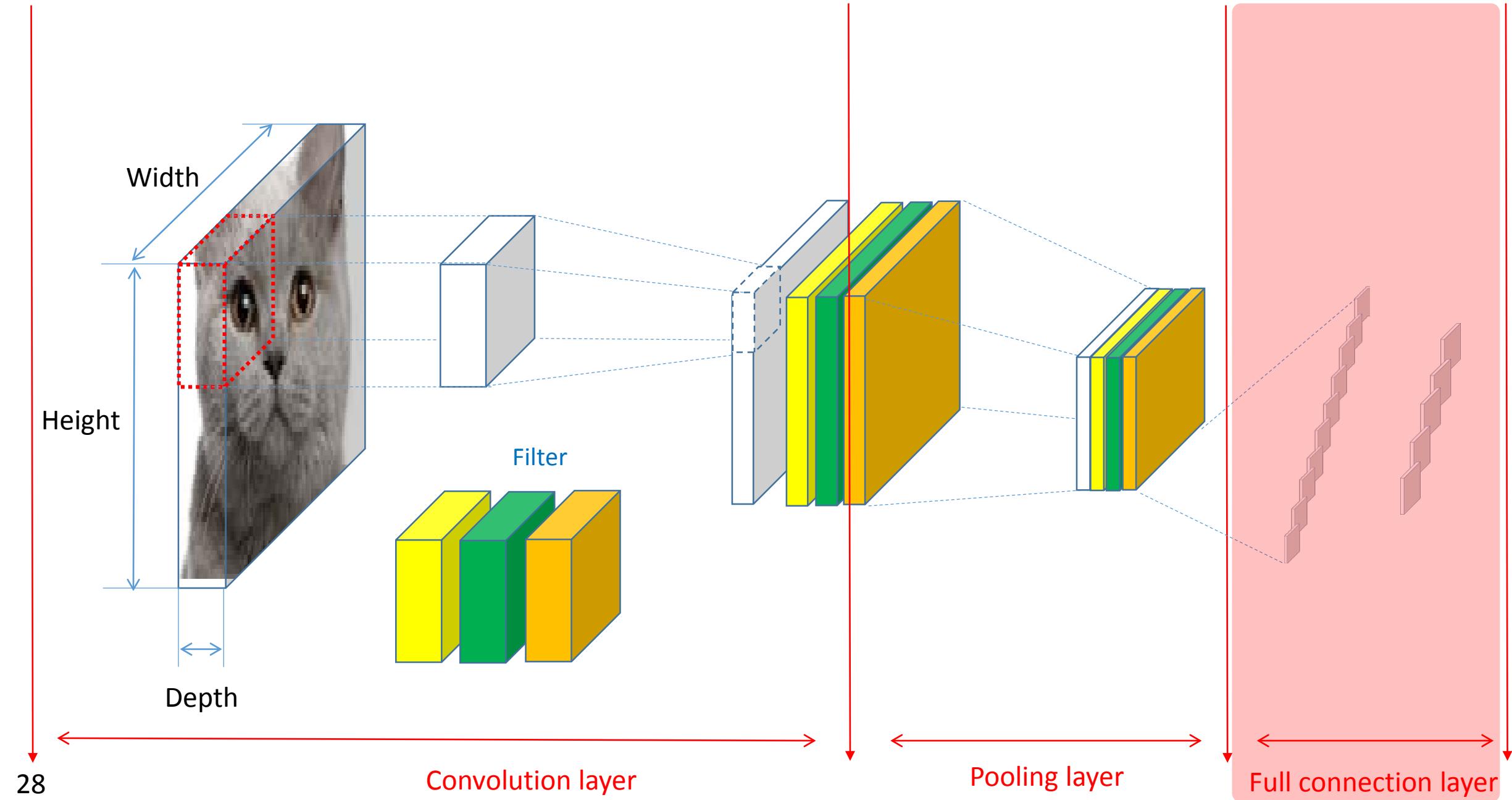


Average pooling

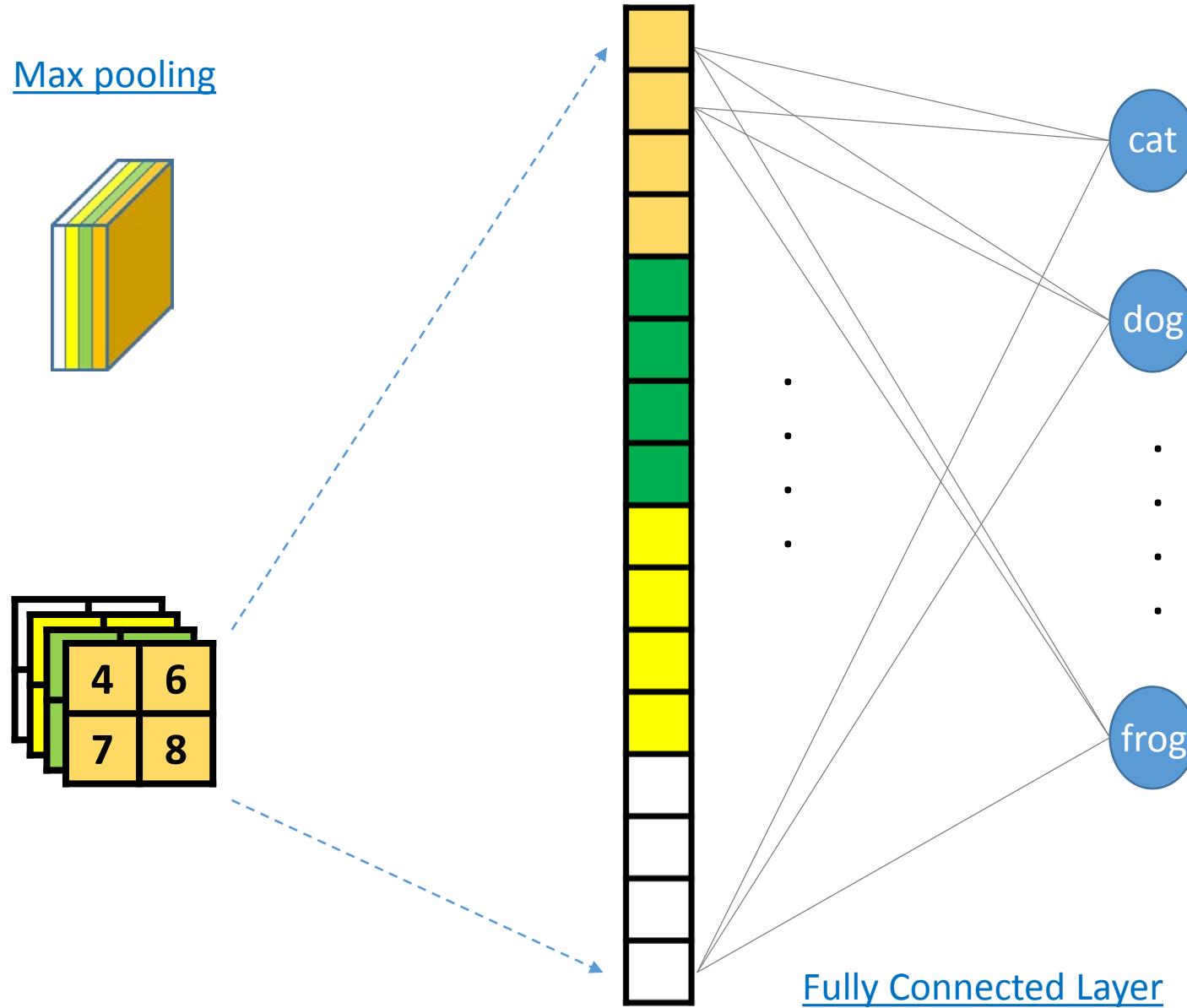


Fully Connected Layer

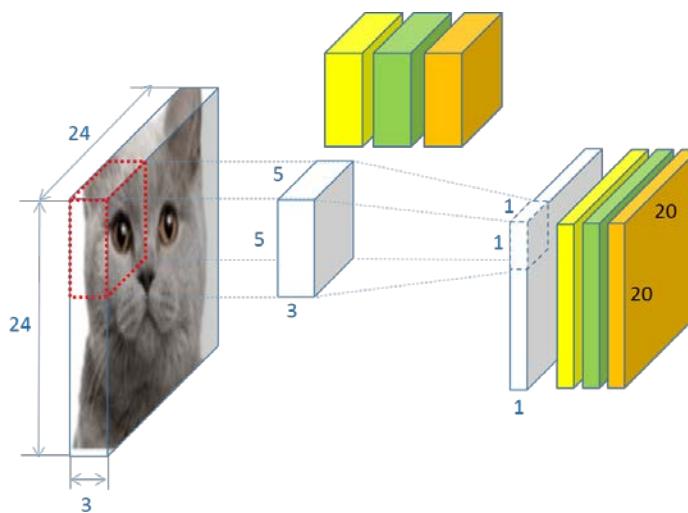
2-D representation of CNN: How does it work?



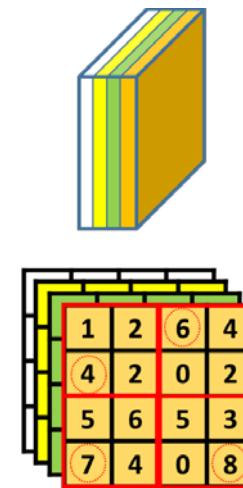
Fully Connected Layer



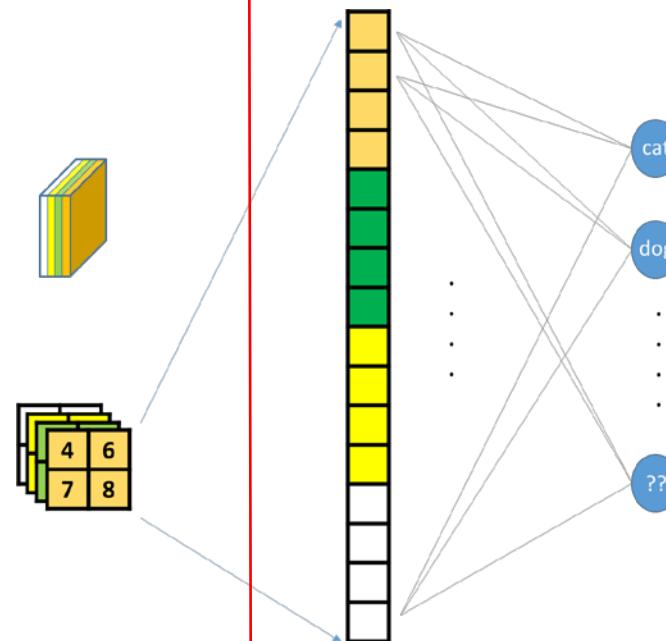
Summary



Convolution layer



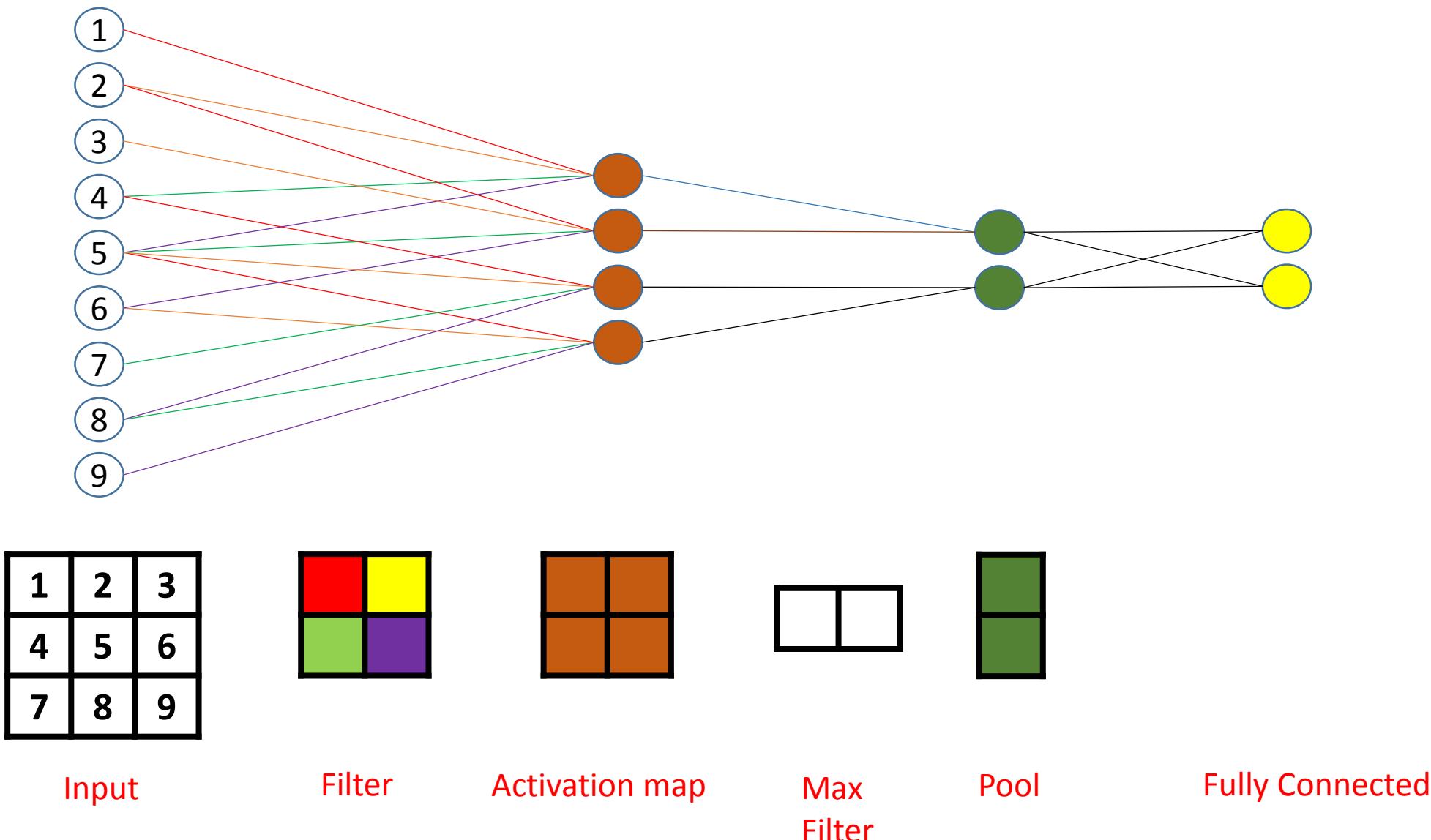
Pooling layer



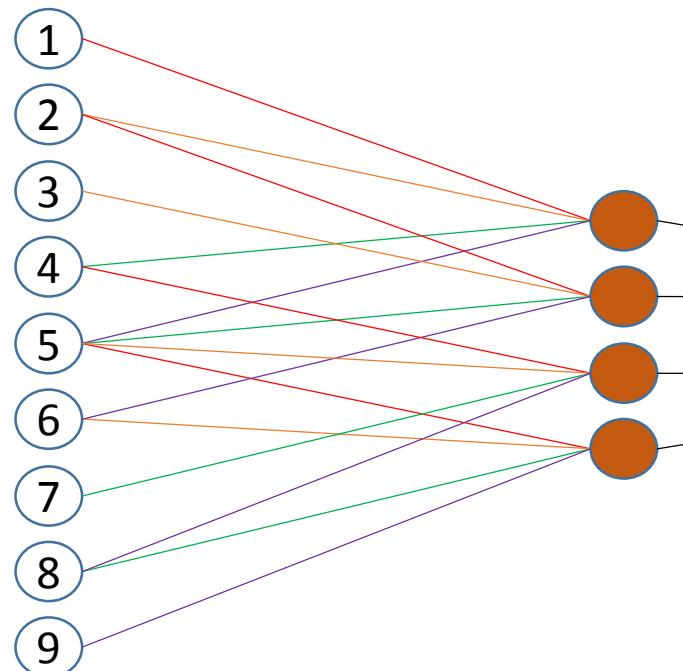
Full connection layer

Backpropagation

Backpropagation



Backpropagation: Fully Connected Layer



Softmax function and Cross Entropy Error

$$p(y_i) = \frac{e^{y_i}}{\sum_k e^{y_k}}$$

SoftMax
function



- t_i (Label)=1 or 0

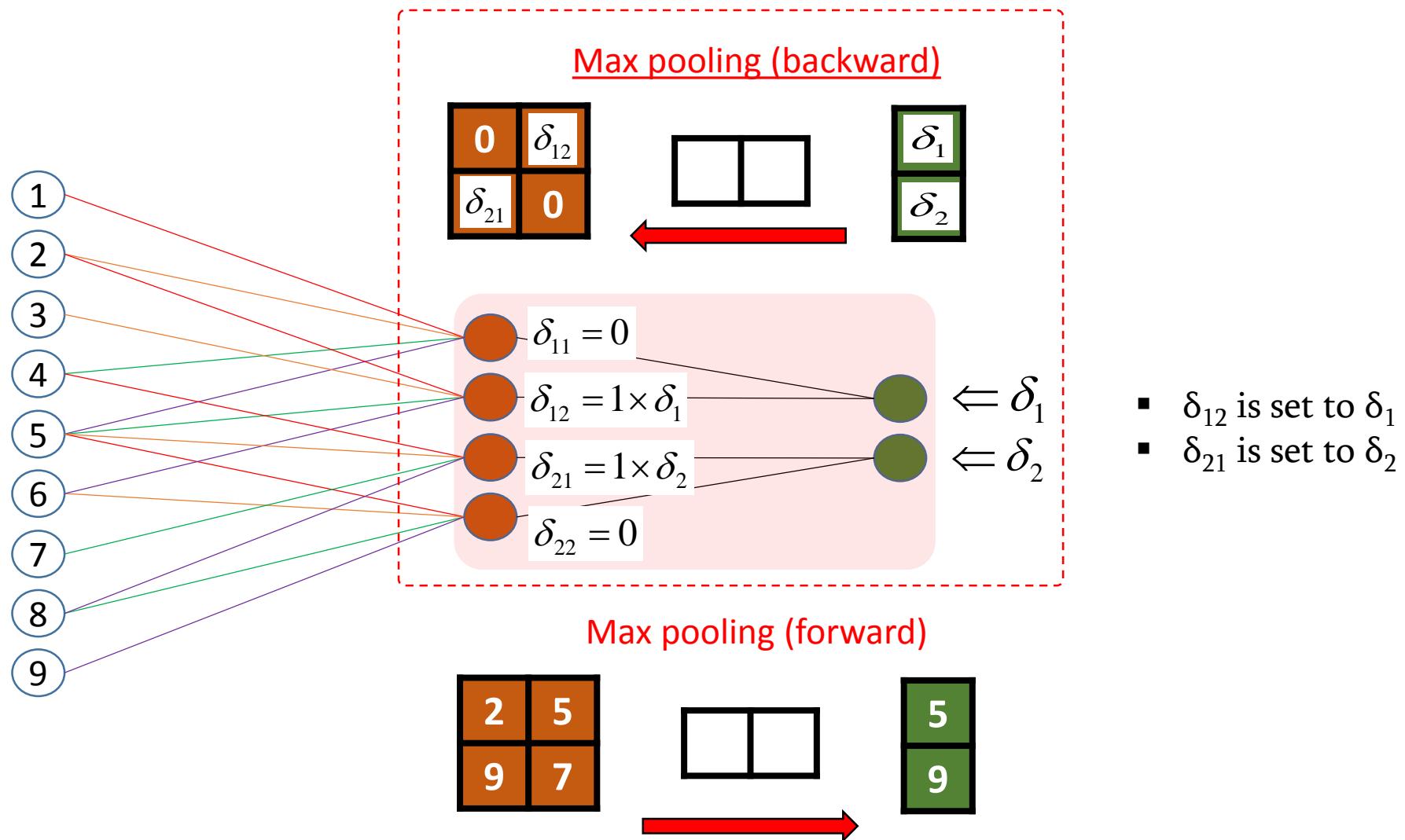
$$E(p) = - \sum_k^{#class} t_k \log p_k$$

Cross entropy
Loss function

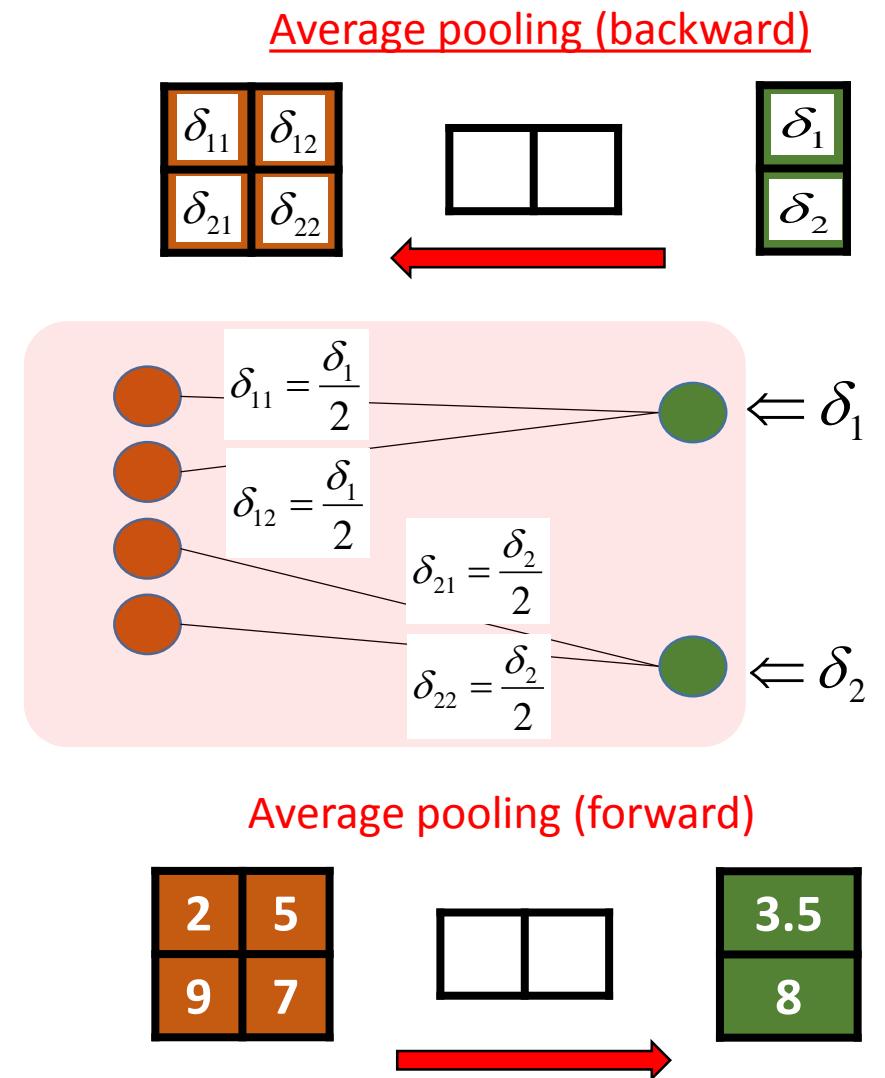
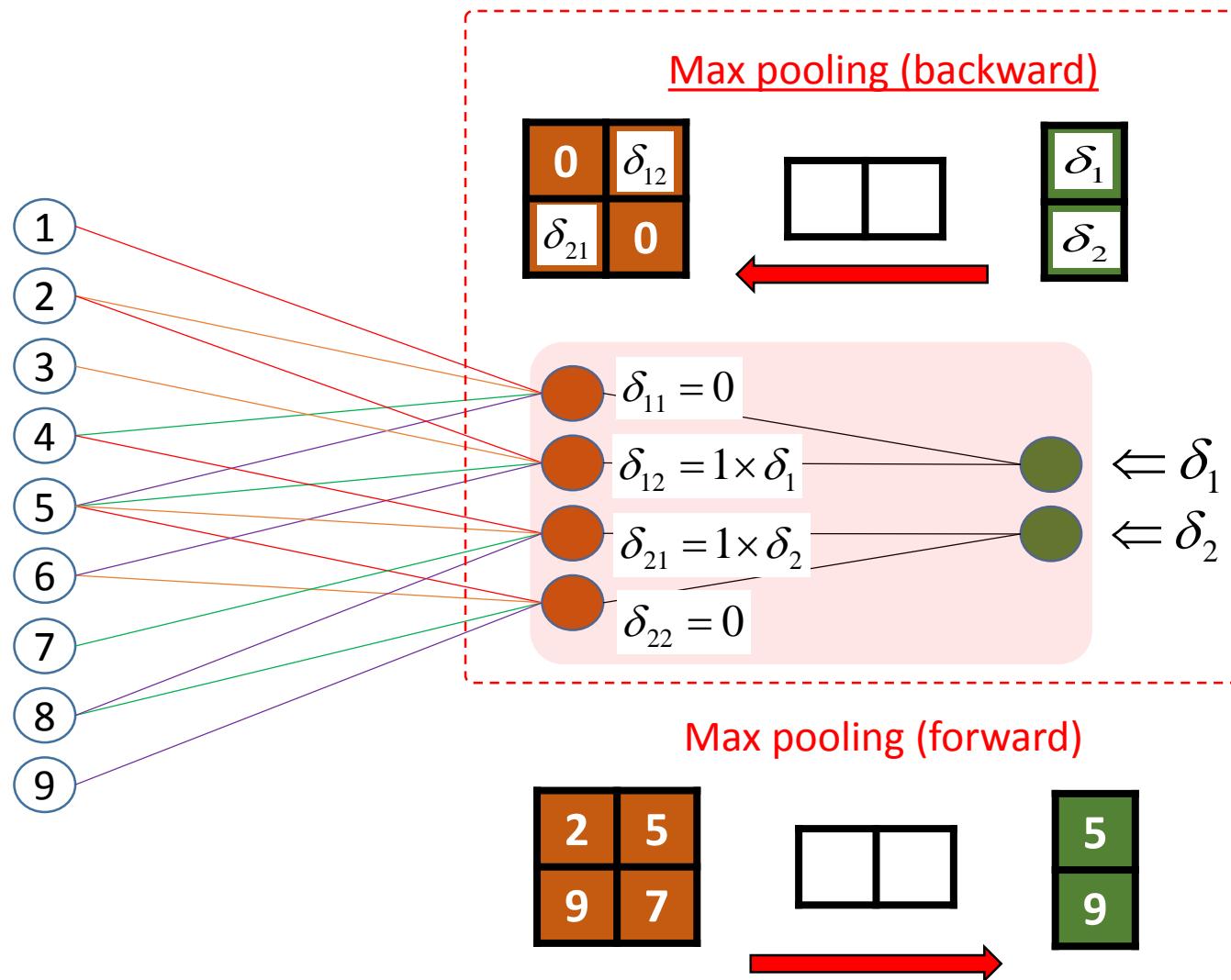
y p

$$\frac{\partial E(p)}{\partial y_i} = \frac{\partial E(p)}{\partial p} \times \frac{\partial p}{\partial y} = p_i - t_i \quad \text{Q12}$$

Backpropagation: Pooling Layer (Max pooling)

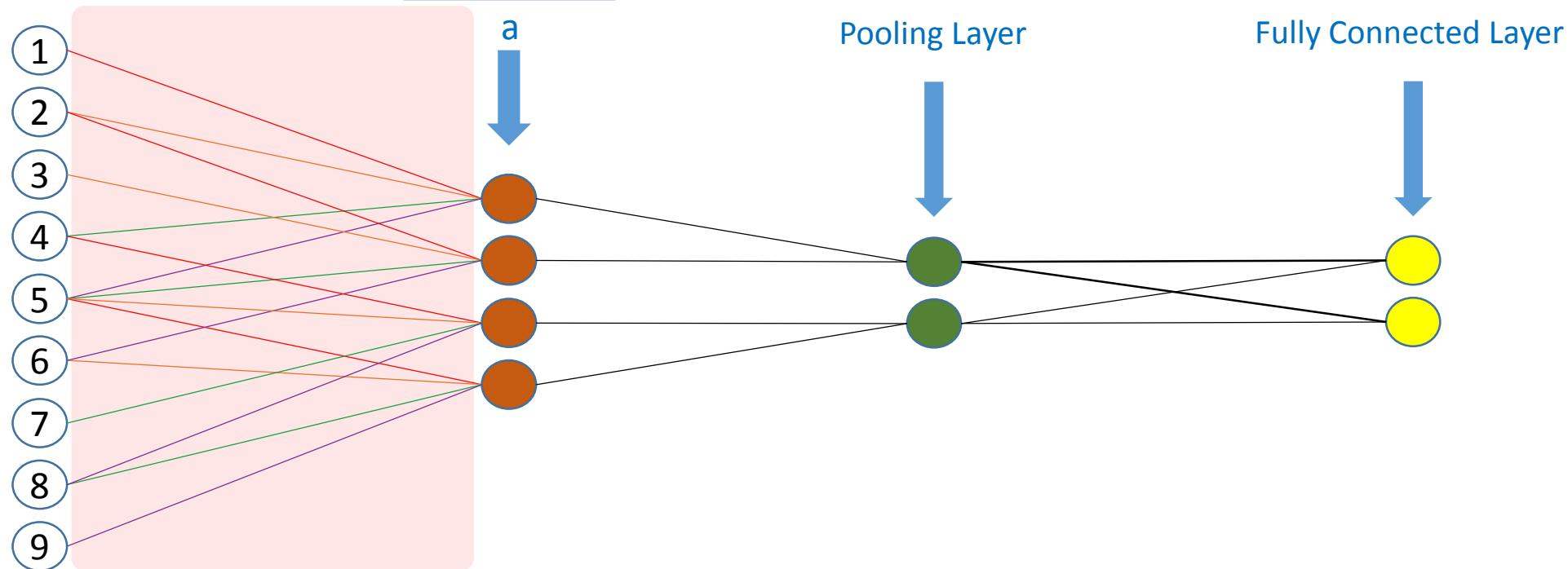
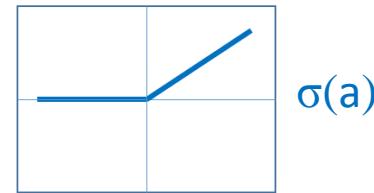


Backpropagation: Pooling Layer (Average pooling)

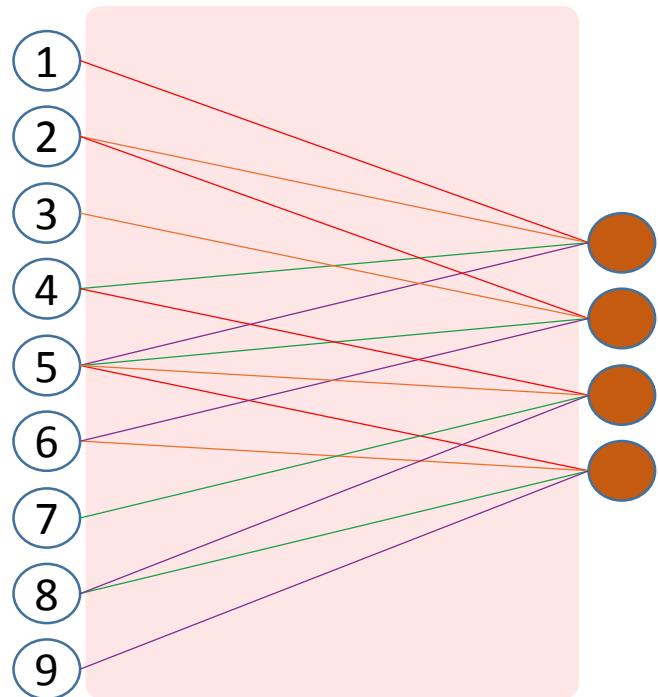


Backpropagation: Convolution Layer

$$\frac{\partial \sigma(a)}{\partial a} = \begin{cases} C & \text{If } a > 0 \\ 0 & \text{Otherwise} \end{cases}$$



Backpropagation: Convolution Layer

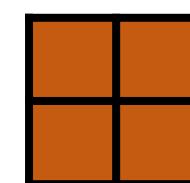


x_{11}	x_{12}	x_{13}
x_{21}	x_{22}	x_{23}
x_{31}	x_{32}	x_{33}

Input

w_{11}	w_{12}
w_{21}	w_{22}

Kernel



Output

Backpropagation: Convolution Layer

□ Convolution

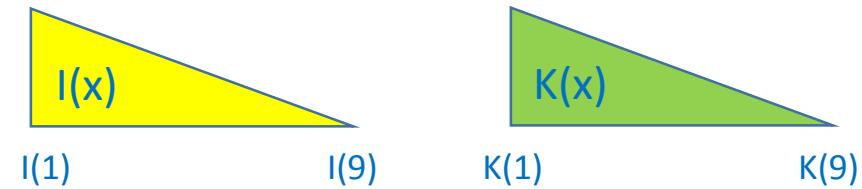
$$(I * K)_{ij} = \sum_{m=0} \sum_{n=0} I(m, n)K(i-m, j-n)$$

Individual elements in both functions, located in the **opposite index**, are multiplied and added.

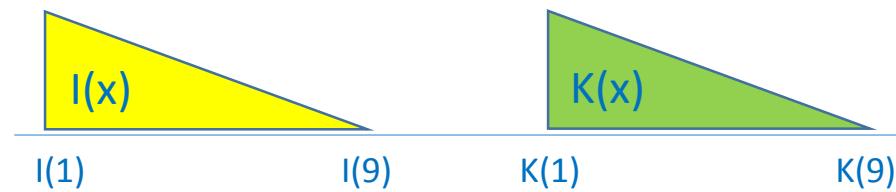
□ Cross-correlation

$$(I \circ K)_{ij} = \sum_{m=0} \sum_{n=0} I(i+m, j+n)K(m, n)$$

Individual elements in both functions, located in the **same index**, are multiplied and added.



□ Convolution



□ Cross-correlation



Convolution $(I * K)_{ij} = (I \circ \text{flip}(K))_{ij}$ Correlation

Backpropagation: Convolution Layer

□ Convolution

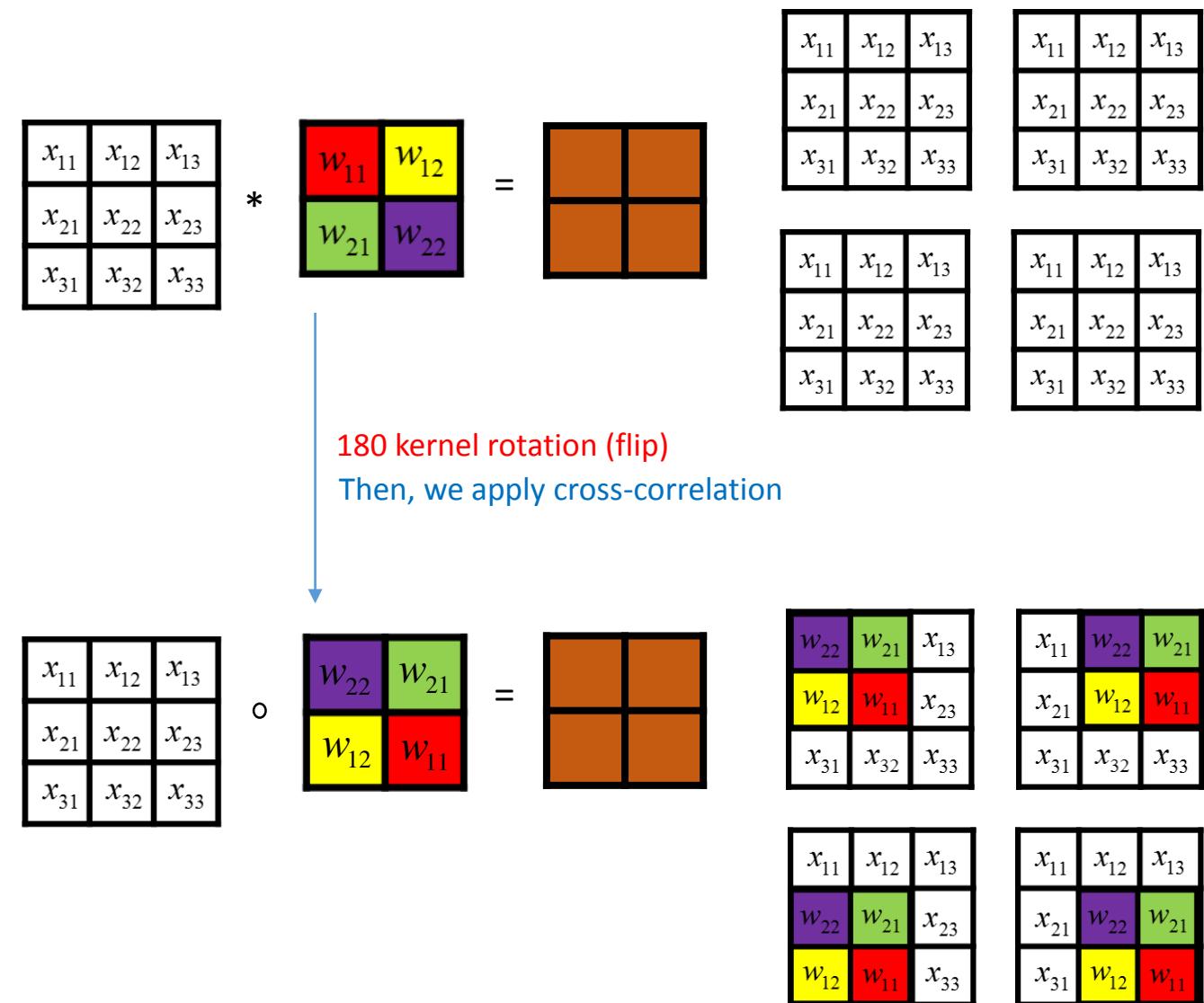
$$(I * K)_{ij} = \sum_{m=0} \sum_{n=0} I(m, n)K(i-m, j-n)$$

Individual elements in both functions, located in the **opposite index**, are multiplied and added.

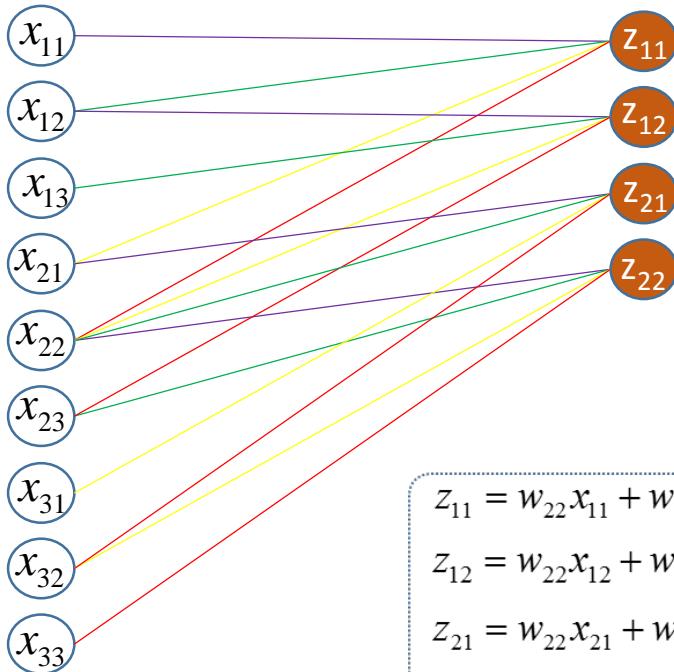
□ Cross-correlation

$$(I \circ K)_{ij} = \sum_{m=0} \sum_{n=0} I(i+m, j+n)K(m, n)$$

Individual elements in both functions, located in the **same index**, are multiplied and added.



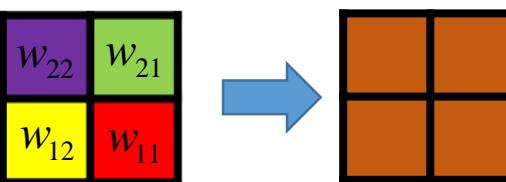
Backpropagation: Convolution Layer - forwarding



x_{11}	x_{12}	x_{13}
x_{21}	x_{22}	x_{23}
x_{31}	x_{32}	x_{33}

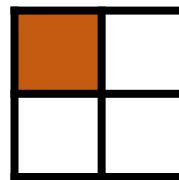
w_{22}	w_{21}
w_{12}	w_{11}

O



Activation map

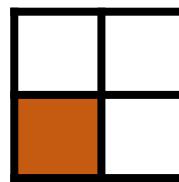
w_{22}	w_{21}	x_{13}
w_{12}	w_{11}	x_{23}
x_{31}	x_{32}	x_{33}



x_{11}	w_{22}	w_{21}	x_{13}
x_{12}	w_{21}	w_{11}	x_{23}
x_{13}			
x_{21}	w_{12}		

x_{11}	w_{22}	w_{21}	x_{13}
x_{12}	w_{21}	w_{11}	x_{23}
x_{13}			
x_{21}	w_{12}		
x_{22}			
x_{23}			
x_{31}			
x_{32}			
x_{33}			

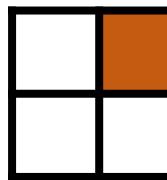
x_{11}	x_{12}	x_{13}
w_{22}	w_{21}	x_{23}
w_{12}	w_{11}	x_{33}



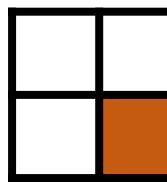
x_{11}	w_{22}	w_{21}	x_{13}
x_{12}	w_{21}	w_{11}	x_{23}
x_{13}			
x_{21}	w_{12}		

x_{11}	w_{22}	w_{21}	x_{13}
x_{12}	w_{21}	w_{11}	x_{23}
x_{13}			
x_{21}	w_{12}		
x_{22}			
x_{23}			
x_{31}			
x_{32}			
x_{33}			

x_{11}	w_{22}	w_{21}	x_{13}
x_{21}	w_{21}	w_{11}	x_{23}
x_{31}			



x_{11}	x_{12}	x_{13}
x_{12}	w_{22}	w_{21}
x_{13}	w_{21}	w_{11}



x_{11}	x_{12}	x_{13}
x_{12}	w_{22}	w_{21}
x_{13}	w_{21}	w_{11}

x_{11}	x_{12}	x_{13}
x_{12}	w_{22}	w_{21}
x_{13}	w_{21}	w_{11}

x_{11}	w_{22}	w_{21}	x_{13}
x_{12}	w_{21}	w_{11}	x_{23}
x_{13}			

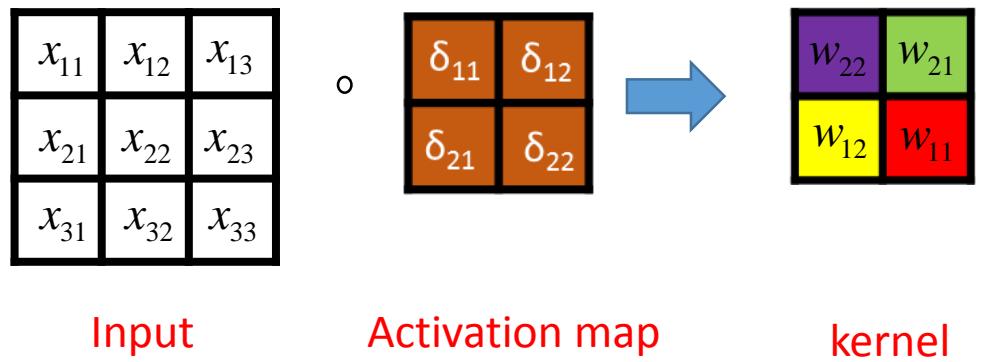
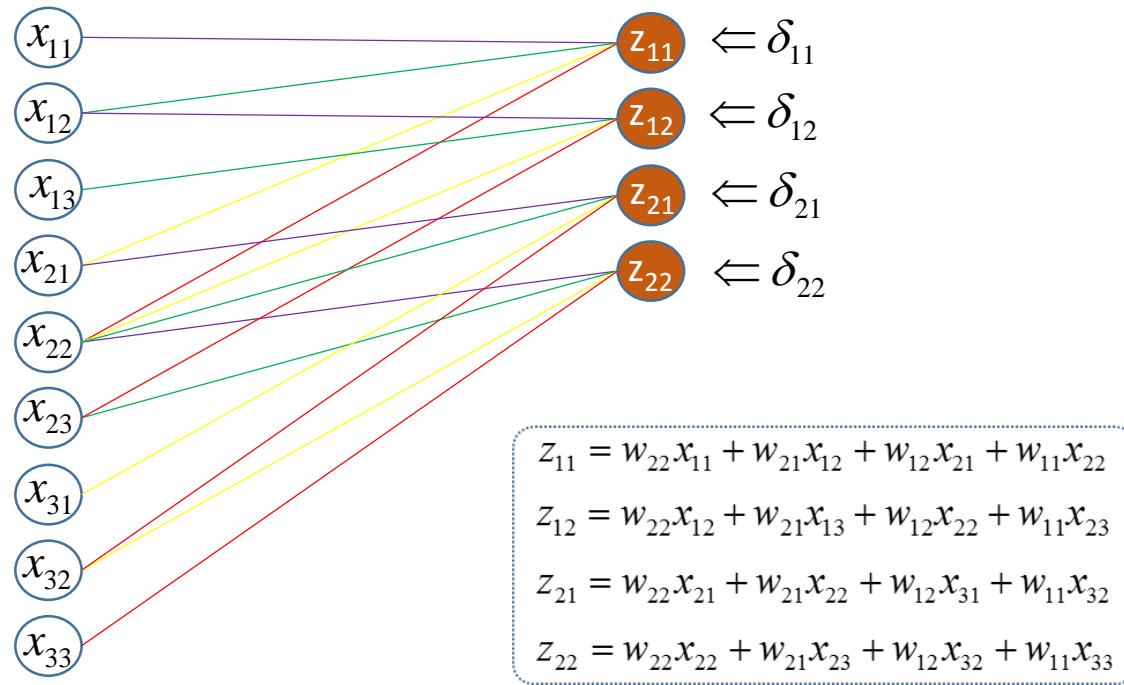
x_{11}	x_{12}	x_{13}
x_{12}	w_{22}	w_{21}
x_{13}	w_{21}	w_{11}

x_{11}	x_{12}	x_{13}
x_{12}	w_{22}	w_{21}
x_{13}	w_{21}	w_{11}

x_{11}	x_{12}	x_{13}
x_{12}	w_{22}	w_{21}
x_{13}	w_{21}	w_{11}

x_{11}	x_{12}	x_{13}
x_{12}	w_{22}	w_{21}
x_{13}	w_{21}	w_{11}

Backpropagation: Convolution Layer – backward: correlation between "w" and " δ "



$$(1) \frac{\partial E}{\partial w_{mn}} = \sum_{i=1}^{N-F+1} \sum_{j=1}^{N-F+1} \delta_{ij} \frac{\partial z_{ij}}{\partial w_{mn}}$$

$$\frac{\partial E}{\partial w_{11}} = \delta_{11}x_{22} + \delta_{12}x_{23} + \delta_{21}x_{32} + \delta_{22}x_{33}$$

$$\frac{\partial E}{\partial w_{12}} = \delta_{11}x_{21} + \delta_{12}x_{22} + \delta_{21}x_{31} + \delta_{22}x_{32}$$

$$\frac{\partial E}{\partial w_{21}} = \delta_{11}x_{12} + \delta_{12}x_{13} + \delta_{21}x_{22} + \delta_{22}x_{23}$$

$$\frac{\partial E}{\partial w_{22}} = \delta_{11}x_{11} + \delta_{12}x_{12} + \delta_{21}x_{21} + \delta_{22}x_{22}$$

δ_{11}	δ_{12}	x_{13}
δ_{21}	δ_{22}	x_{23}
x_{31}	x_{32}	x_{33}

w_{22}	w_{21}
w_{12}	w_{11}

x_{11}	x_{12}	x_{13}
δ_{11}	δ_{12}	x_{23}
δ_{21}	δ_{22}	x_{33}

x_{11}	δ_{11}	δ_{12}
x_{21}	δ_{21}	δ_{22}
x_{31}	x_{32}	x_{33}

x_{11}	x_{12}	x_{13}
x_{21}	δ_{11}	δ_{12}
x_{31}	δ_{21}	δ_{22}

Backup Slides

Backpropagation: Convolution Layer – backward: why correlation here - proof?

$$(1) \frac{\partial E}{\partial w_{mn}} = \sum_{i=1}^{N-F+1} \sum_{j=1}^{N-F+1} \delta_{ij} \frac{\partial z_{ij}}{\partial w_{mn}}$$

$$= \sum_{i=1}^{N-F+1} \sum_{j=1}^{N-F+1} \delta_{ij} \frac{\partial}{\partial w_{mn}} \left(\sum_k \sum_q w_{kq} x_{i+k, j+q} \right)$$

One in the right side was obtained using convolution previously.

$$\begin{aligned} z_{11} &= w_{22}x_{11} + w_{21}x_{12} + w_{12}x_{21} + w_{11}x_{22} \\ z_{12} &= w_{22}x_{12} + w_{21}x_{13} + w_{12}x_{22} + w_{11}x_{23} \\ z_{21} &= w_{22}x_{21} + w_{21}x_{22} + w_{12}x_{31} + w_{11}x_{32} \\ z_{22} &= w_{22}x_{22} + w_{21}x_{23} + w_{12}x_{32} + w_{11}x_{33} \end{aligned}$$

We only consider when " $w_{mn} = w_{kq}$ ", in other words " $m=k, n=q$ ".

$$= \sum_{i=1}^{N-F+1} \sum_{j=1}^{N-F+1} \delta_{ij} \frac{\partial}{\partial w_{mn}} \left(\sum_k \sum_q w_{mn} x_{i+m, j+n} \right)$$

$$= \sum_{i=1}^{N-F+1} \sum_{j=1}^{N-F+1} \delta_{ij} x_{m+i, n+j}$$

It is correlation

Finally, " ∂w " is the result of cross correlation between "x" and " δ " previously.

