

CS 4803 / 7643: Deep Learning

Topics:

- Convolutional Neural Networks
 - Stride, padding
 - Pooling layers
 - Fully-connected layers as convolutions

Dhruv Batra
Georgia Tech

Administrativa

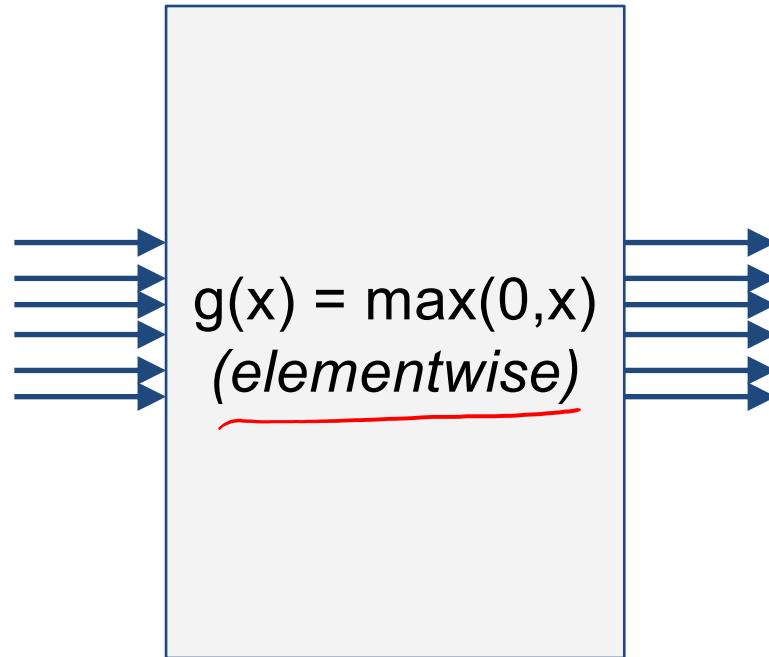
- HW1 Reminder
 - Due: 10/02, 11:55pm

Recap from last time

Jacobian of ReLU

$$\vec{h}^{l-1} \in \mathbb{R}^{4096}$$

4096-d
input vector

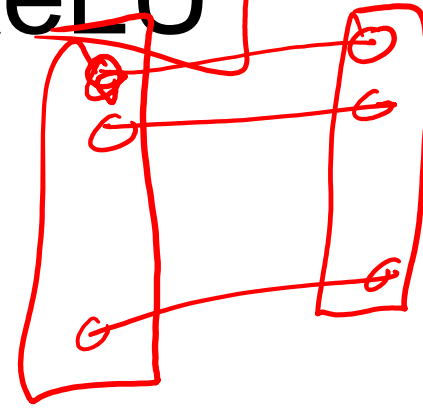


$$\vec{h}^l \in \mathbb{R}^{4096}$$

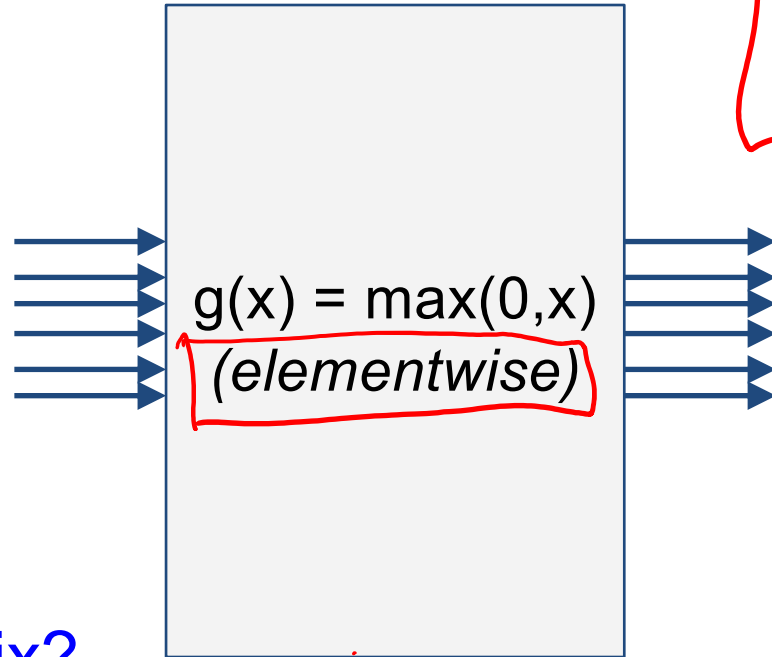
4096-d
output vector

$$\begin{bmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 0 \end{bmatrix}$$

Jacobian of ReLU



4096-d
input vector



4096-d
output vector

Q: what is the
size of the
Jacobian matrix?

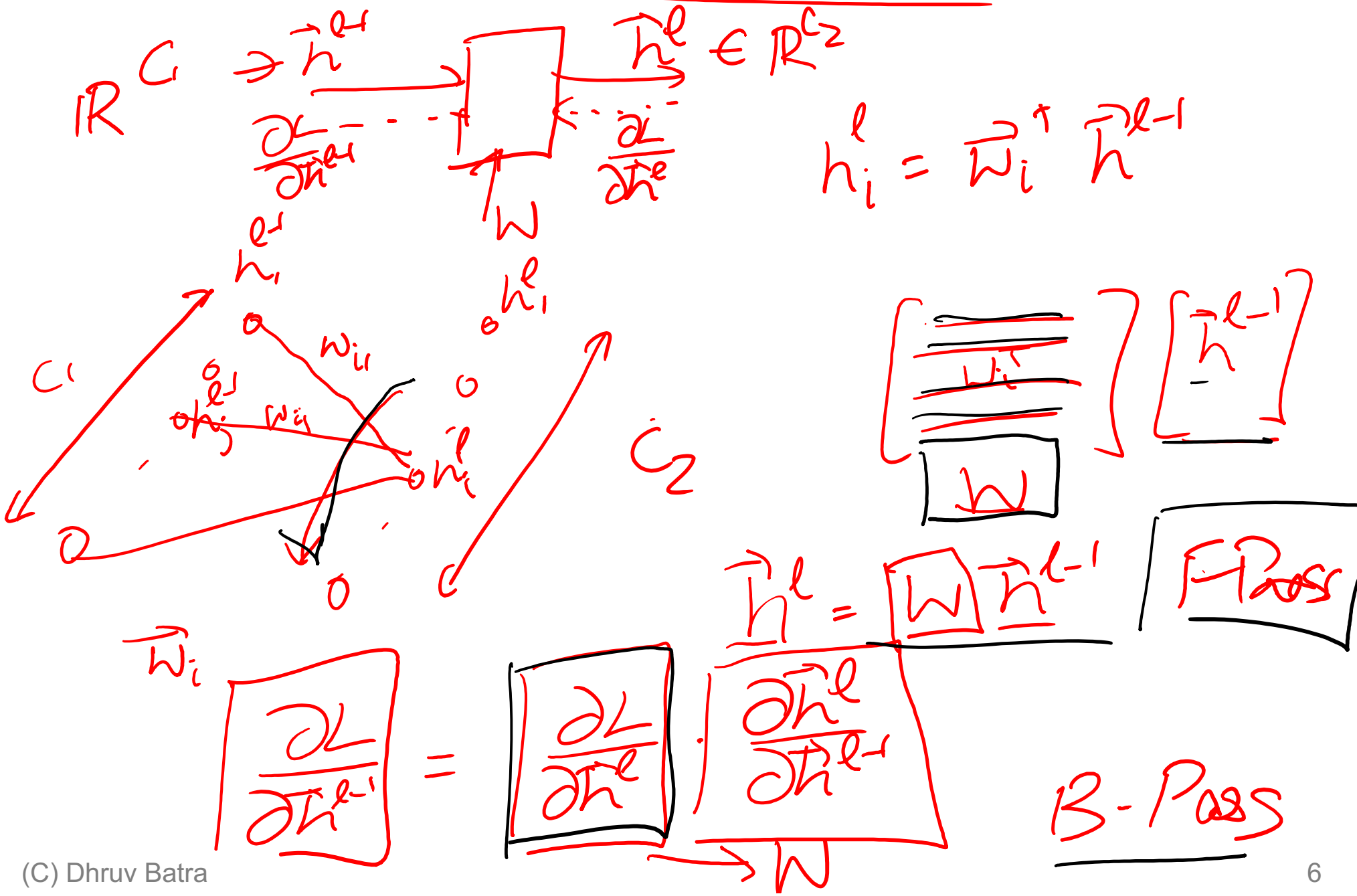
$$\frac{\partial h_i^l}{\partial h_j^{l-1}} = 0 \quad \forall i \neq j$$

$$\frac{\partial \vec{h}^l}{\partial \vec{h}^{l-1}} = \begin{bmatrix} \frac{\partial h_1^l}{\partial h_1^{l-1}} & & \\ & \ddots & \\ & & \frac{\partial h_n^l}{\partial h_n^{l-1}} \end{bmatrix} \rightarrow \begin{bmatrix} \circ & & \\ & \circ & \\ & & \circ \end{bmatrix}$$

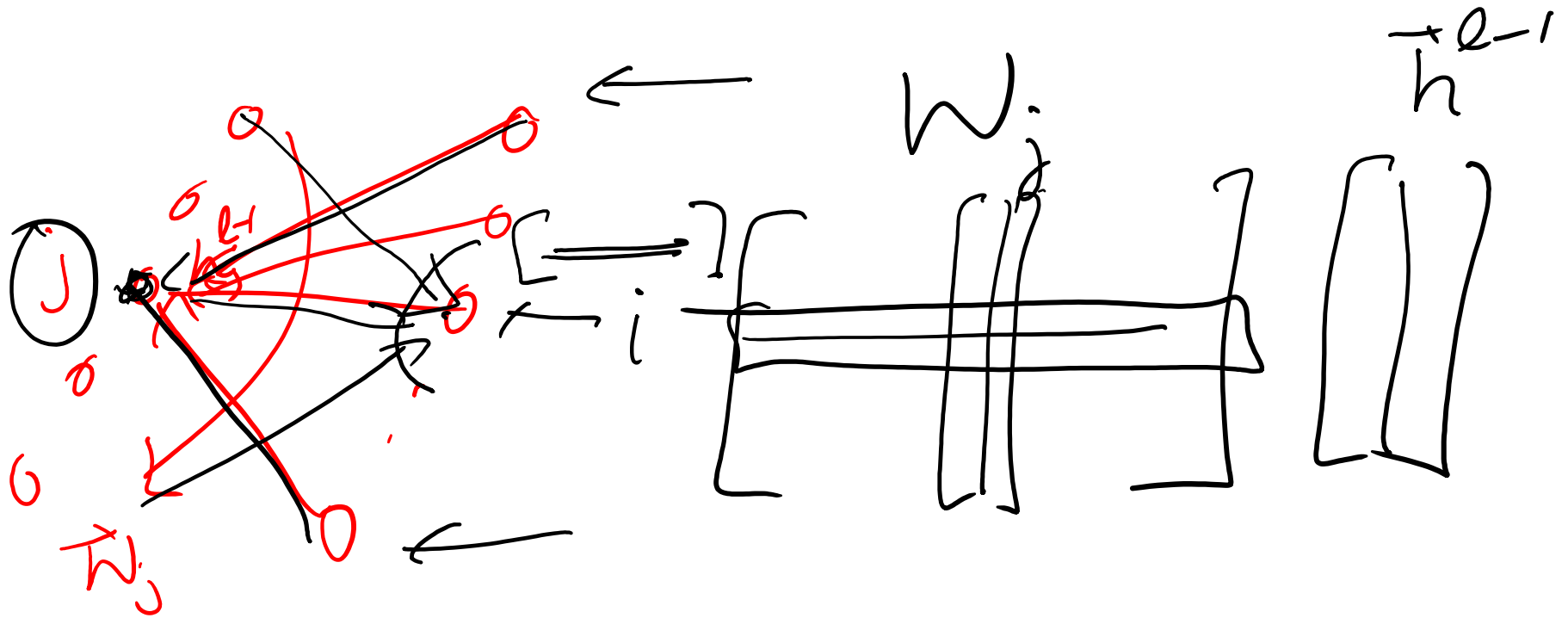
4096 x 4096

Jacobians of FC-Layer

↗



Jacobians of FC-Layer



Jacobians of FC-Layer



Convolutional Neural Networks

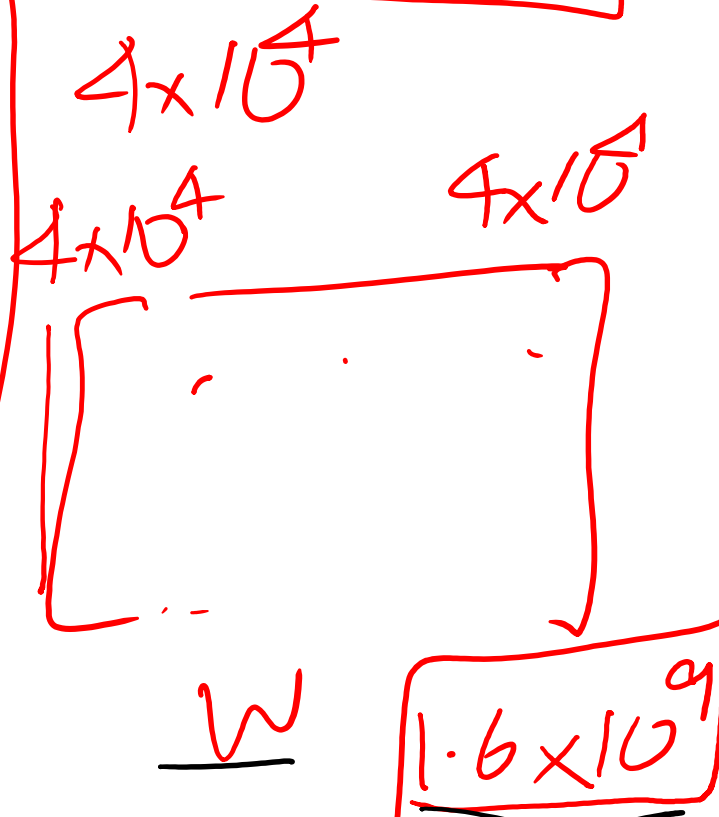
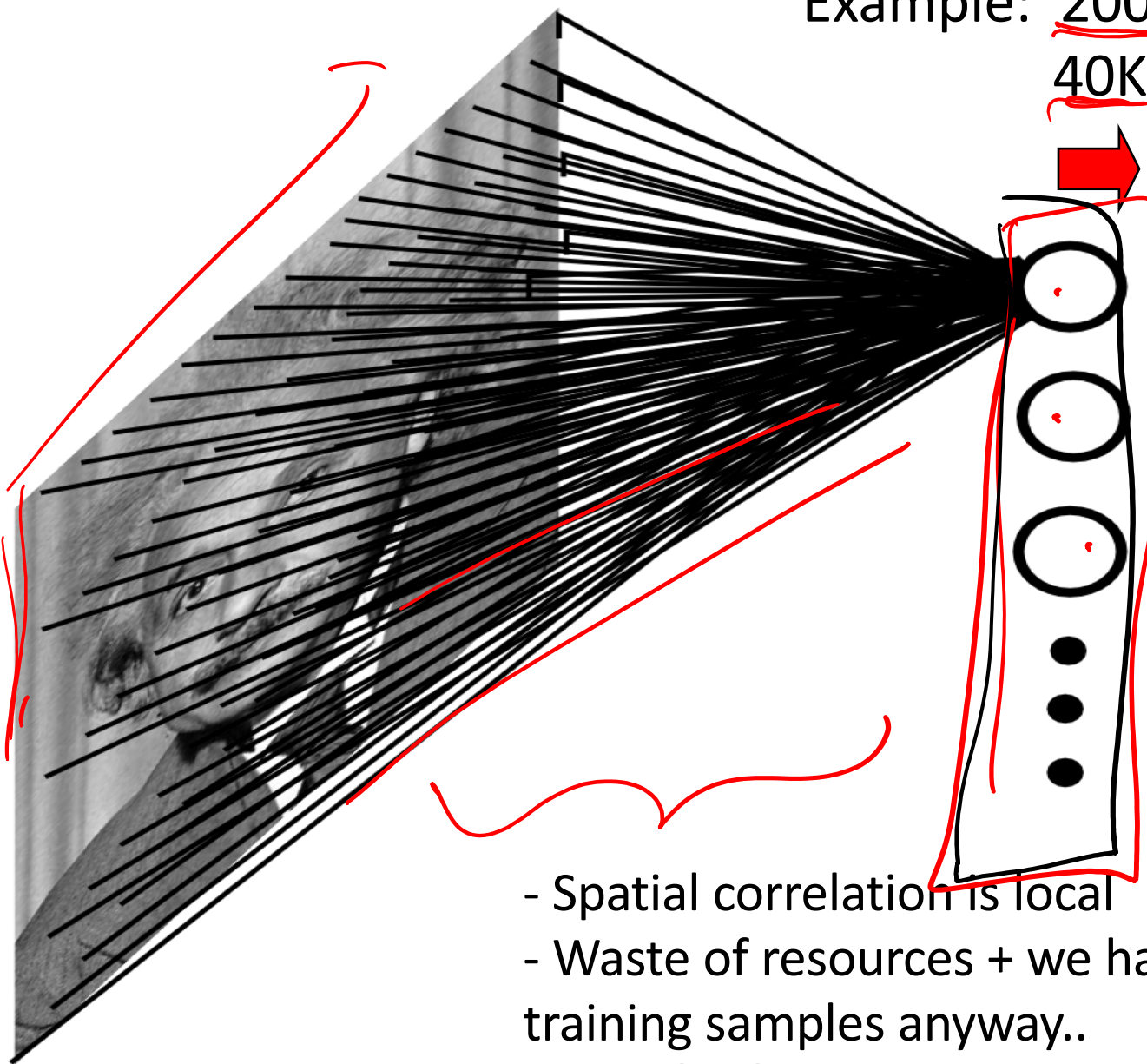
(without the brain stuff)

Fully Connected Layer

Example: 200x200 image
40K hidden units

$$4 \times 10^4$$

~2B parameters!!!

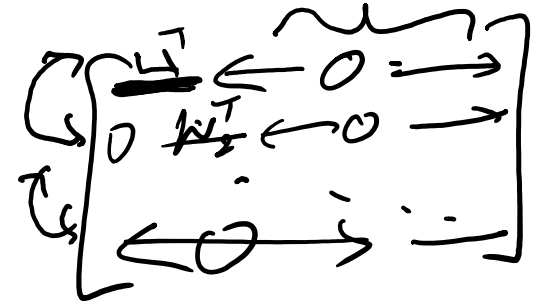
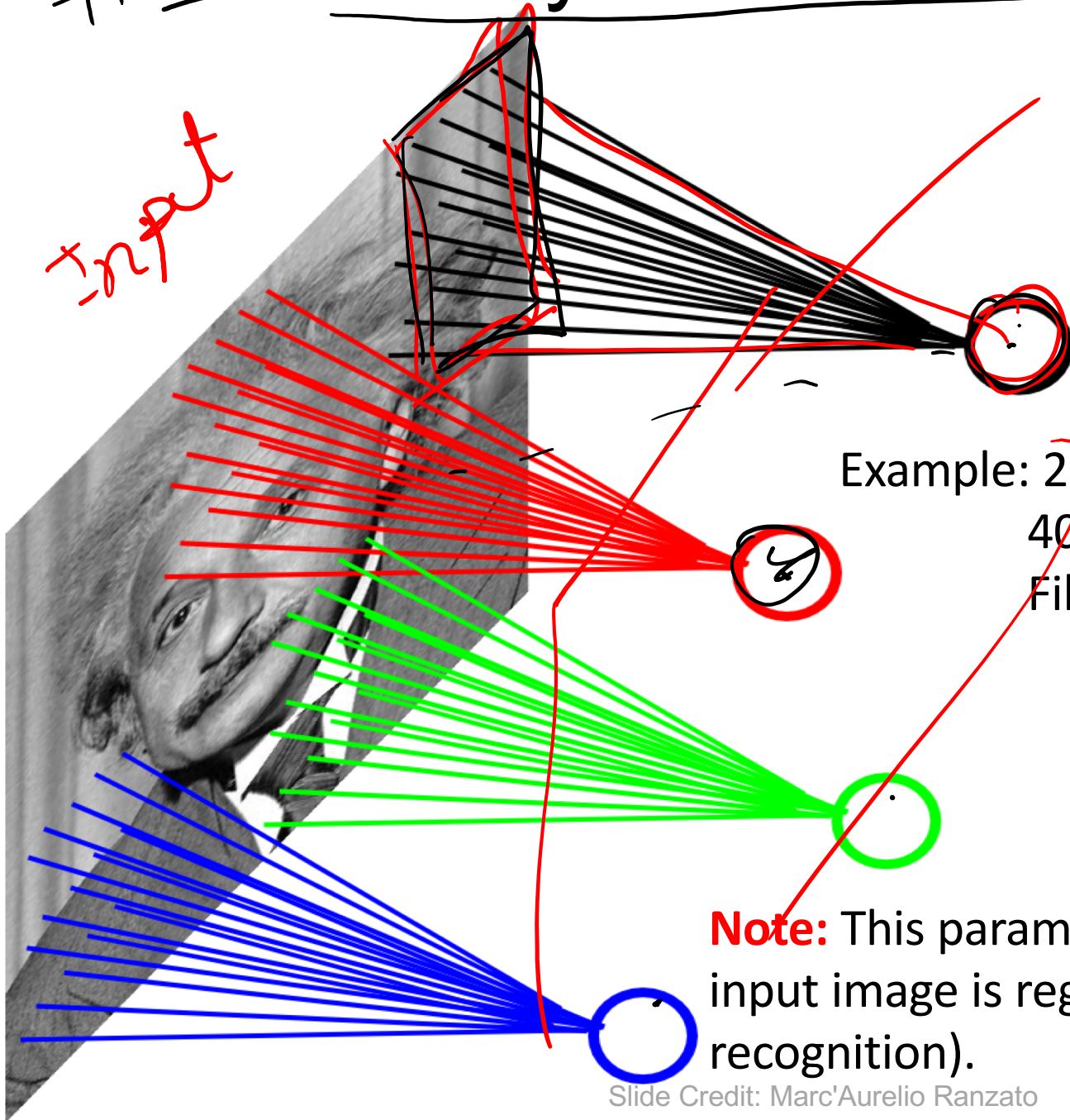


- Spatial correlation is local
- Waste of resources + we have not enough training samples anyway..

#1

Locally Connected Layer W

input



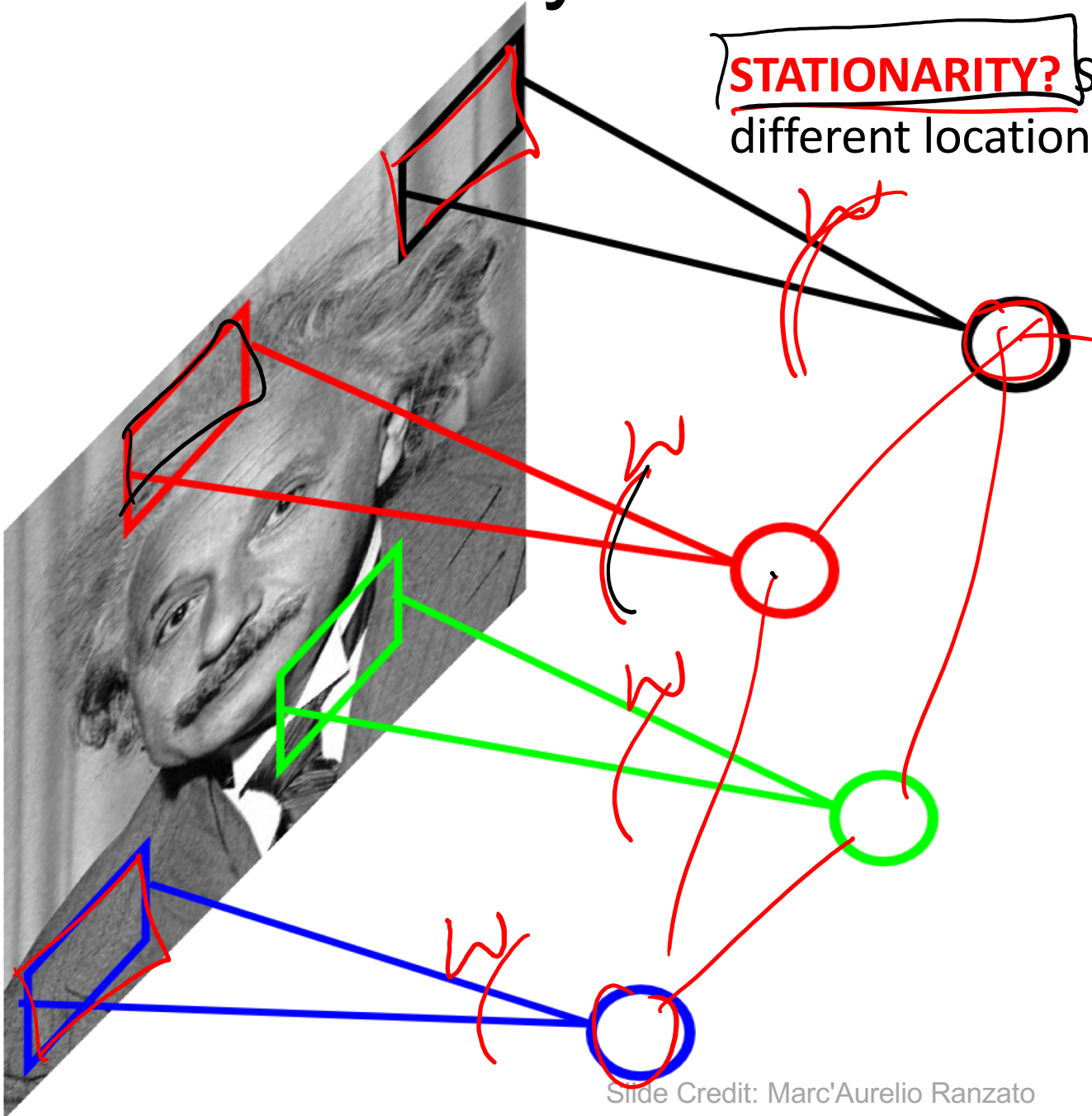
Example: 200x200 image
40K hidden units
Filter size: 10x10
4M parameters

Note: This parameterization is good when input image is registered (e.g., face recognition).

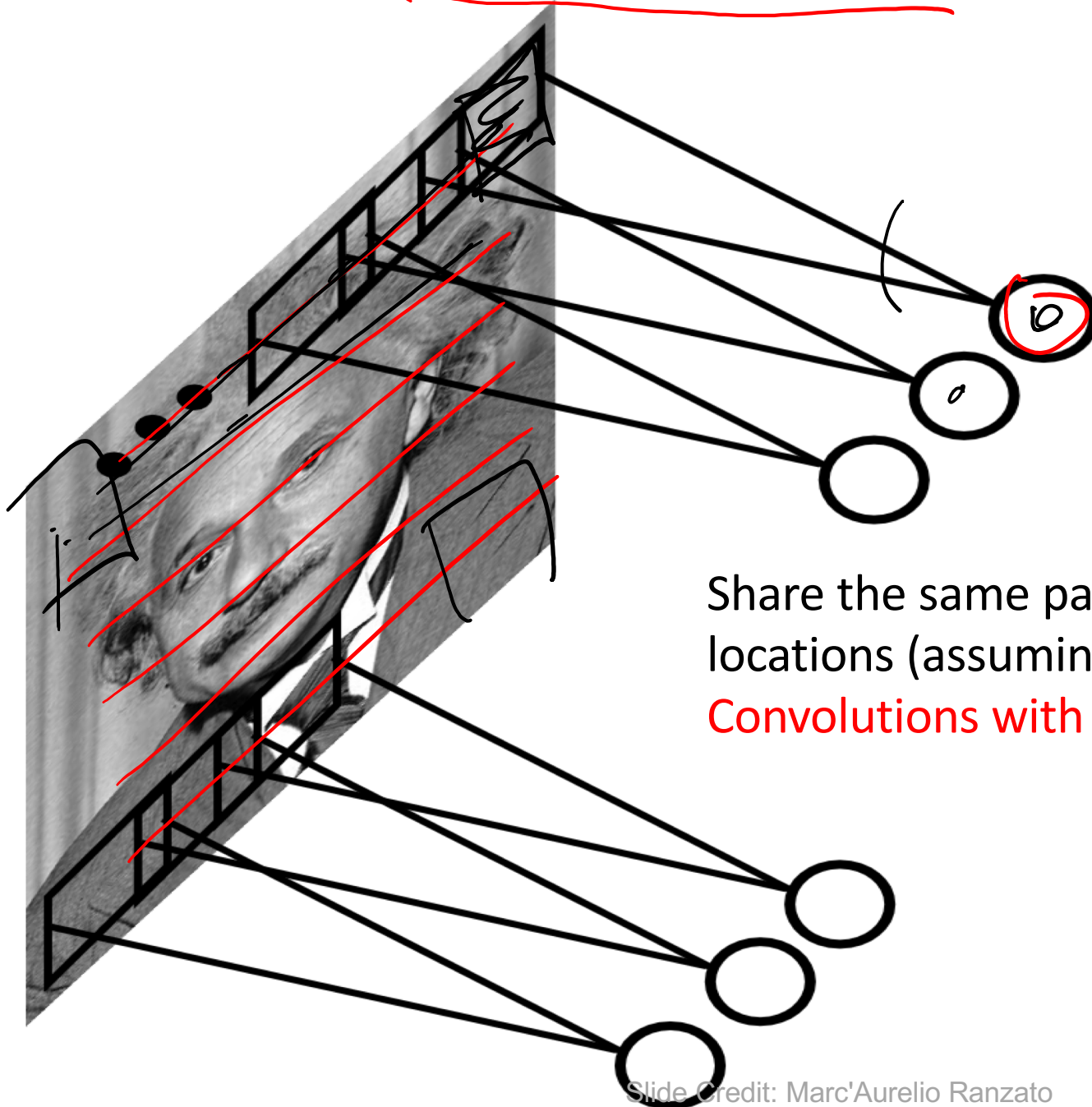
Locally Connected Layer

STATIONARITY?

Statistics is similar at different locations



Convolutional Layer



Share the same parameters across different locations (assuming input is stationary):

Convolutions with learned kernels

Convolutions!

math \rightarrow CS \rightarrow programming

Convolutions for mathematicians



$$\begin{aligned} \boxed{y(t)} &= (x * w)(t) \\ &= (w * x)(t) \\ &= \int_{-\infty}^{\infty} \underline{x(t-a)} \boxed{w(a)} da \end{aligned}$$

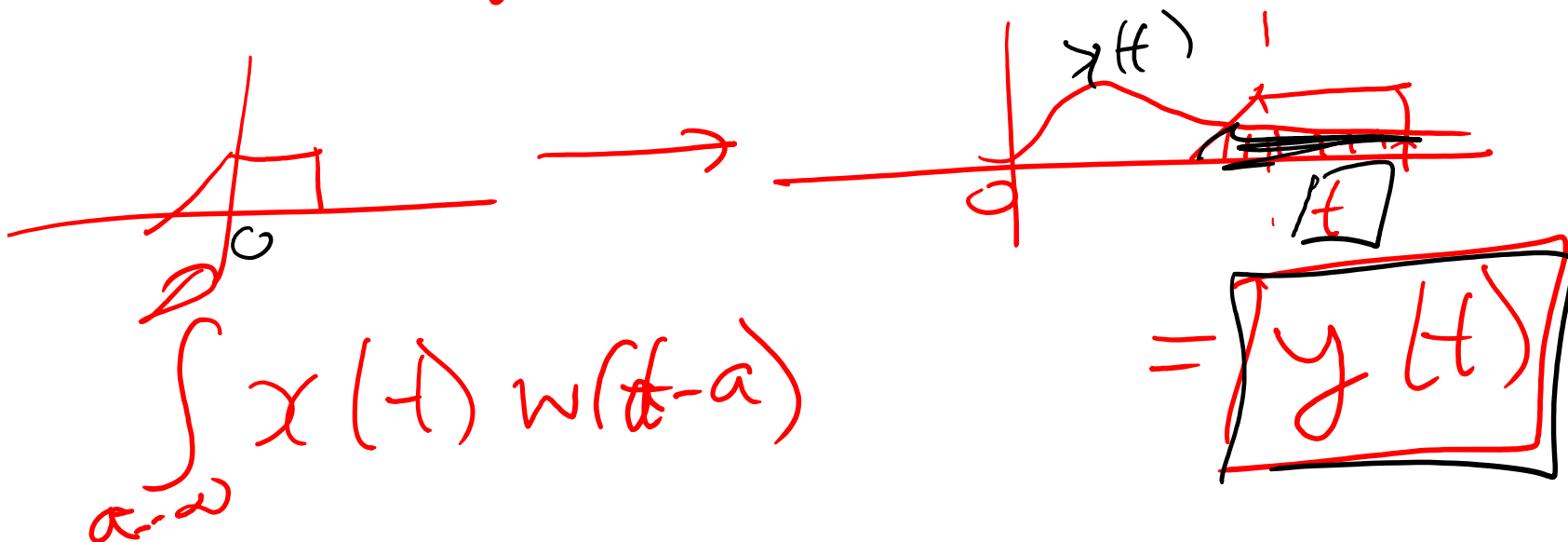
$$= \int_{-\infty}^{\infty} \underline{x(a)} \boxed{w(t-a)} da$$

Convolutions for mathematicians

$$w(a) \rightarrow w(\underline{-a})$$



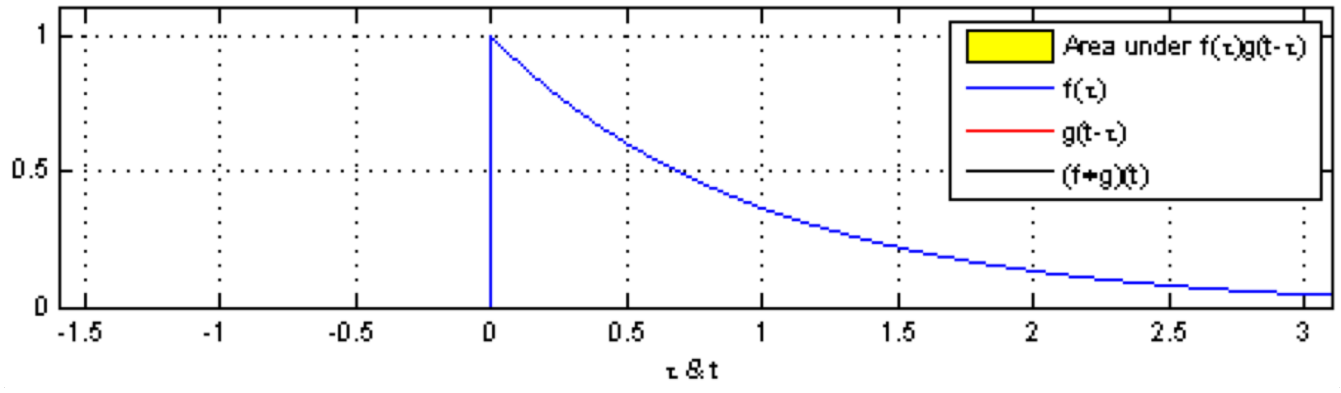
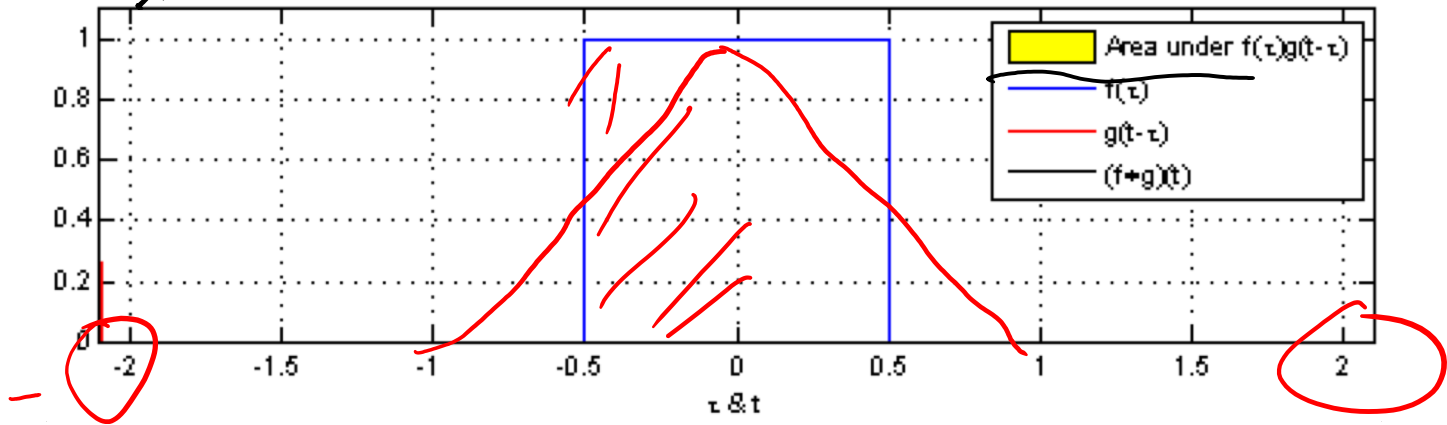
$$w(\underline{-a}) \xrightarrow{\text{-kernel}} w(\underline{-(a-t)})$$



$$\int_{a-\infty}^{\infty} x(t) w(t-a)$$

$$= y(t)$$

$$y(t_1, t_2) = \int_{a=-\infty}^{\infty} \int_{b=-\infty}^{\infty} x(t_1 - a, t_2 - b) w(a, b) da db$$



"Convolution of box signal with itself2" by Convolution_of_box_signal_with_itself.gif: Brian Amberg derivative work: Tinos (talk) - Convolution_of_box_signal_with_itself.gif. Licensed under CC BY-SA 3.0 via Commons - https://commons.wikimedia.org/wiki/File:Convolution_of_box_signal_with_itself2.gif#/media/File:Convolution_of_box_signal_with_itself2.gif

Convolutions for computer scientists

$$y[x, c] = \sum_{a=-\infty}^{\infty} \sum_{b=-\infty}^{\infty} x[x-a, c-b] w[a, b]$$

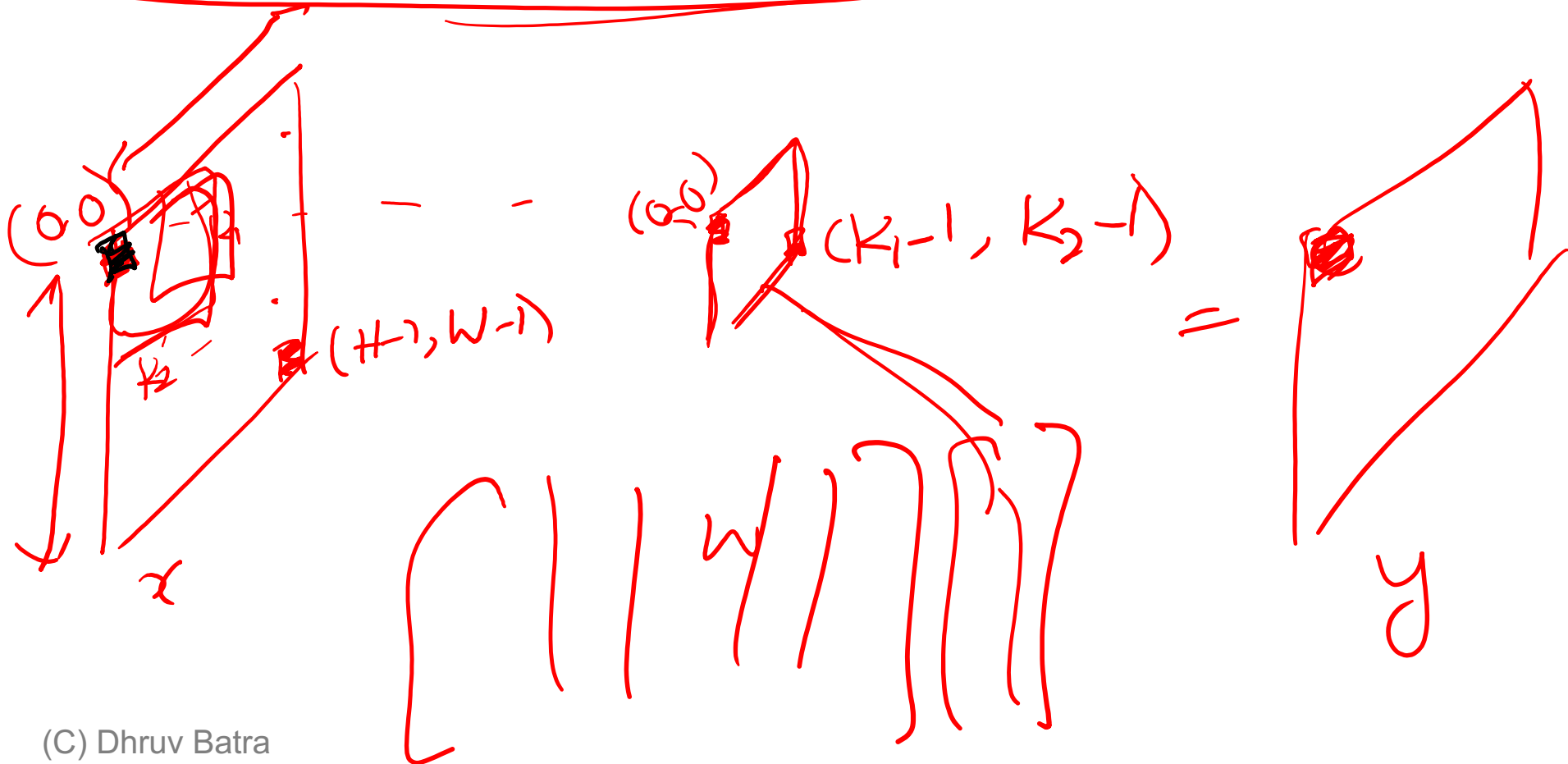


$$y[x, c] = \sum_{a=-\frac{k_1-1}{2}}^{\frac{k_1-1}{2}} \sum_{b=-\frac{k_2-1}{2}}^{\frac{k_2-1}{2}} x[x-a, c-b] w[a, b]$$

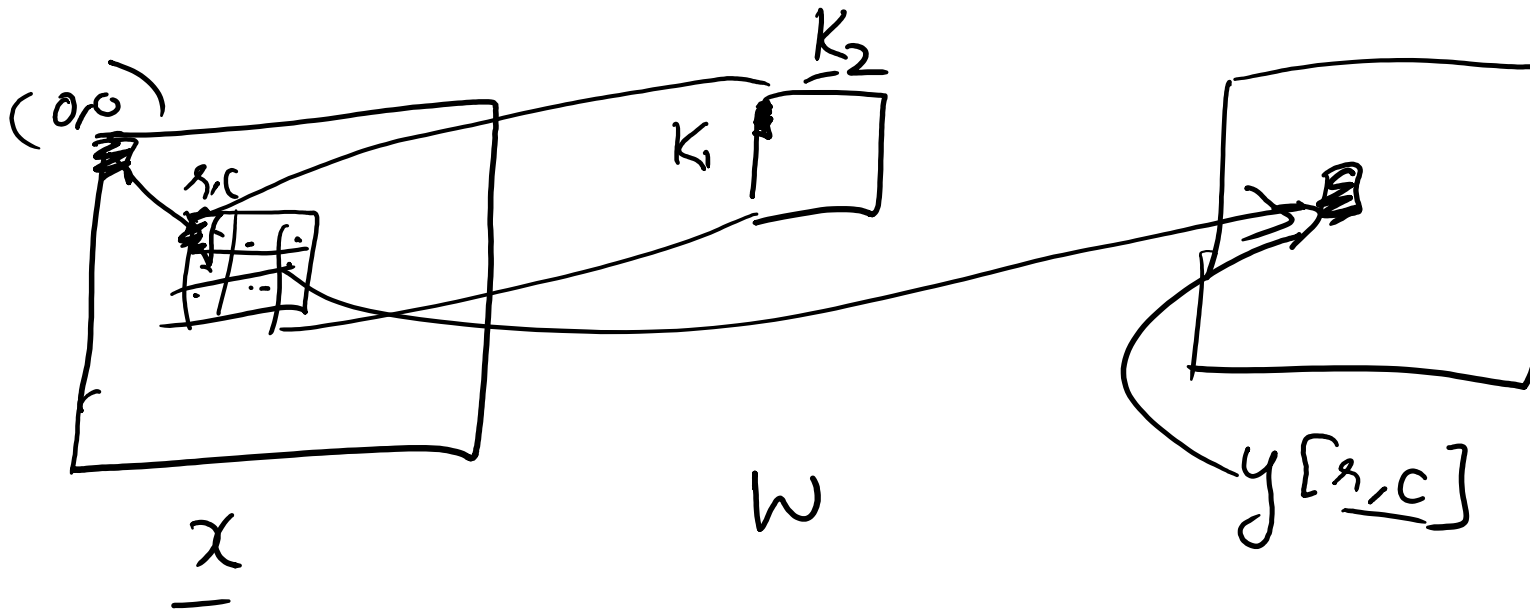
Convolutions for computer scientists

Convolutions for programmers

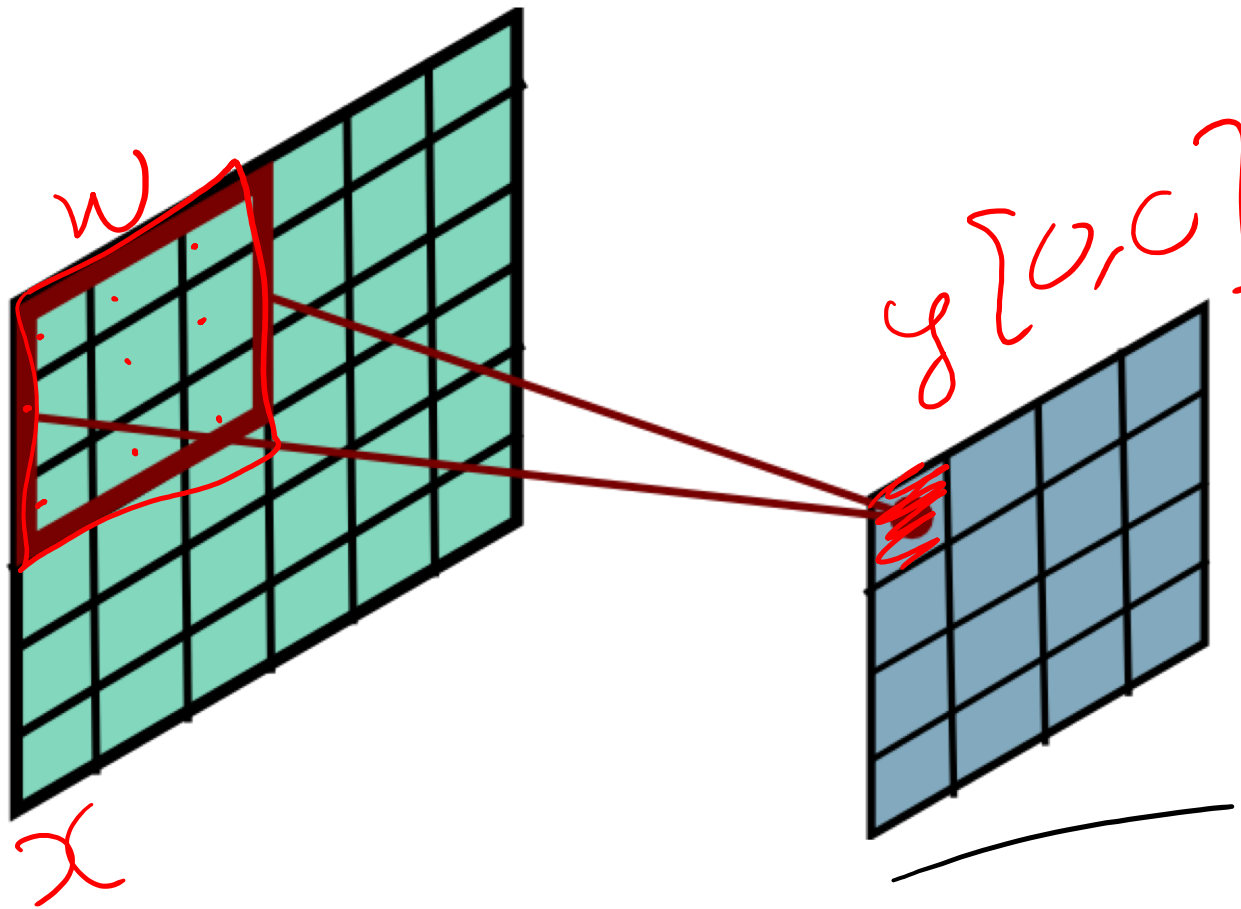
$$y[x, c] = \sum_{a=0}^{K_1-1} \sum_{b=0}^{K_2-1} x[(x+a), (y+b)] w[a, b]$$



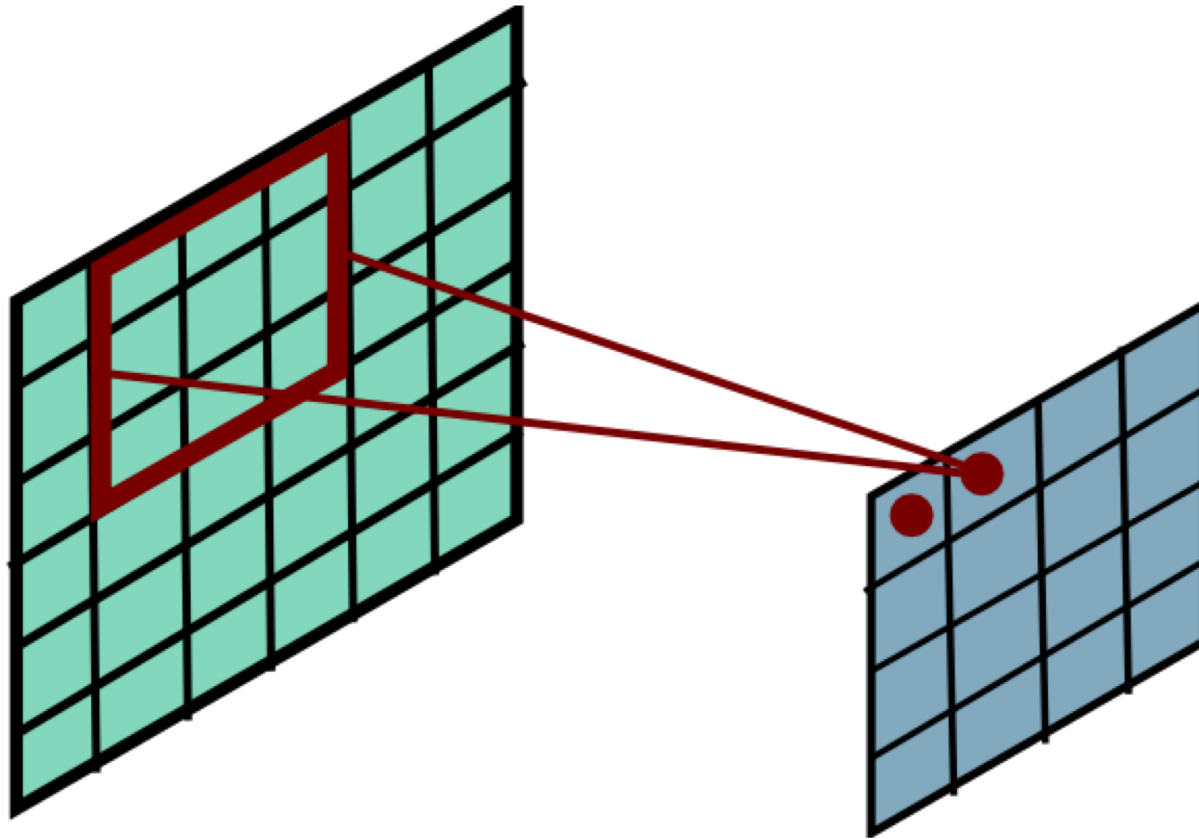
Convolutions for programmers



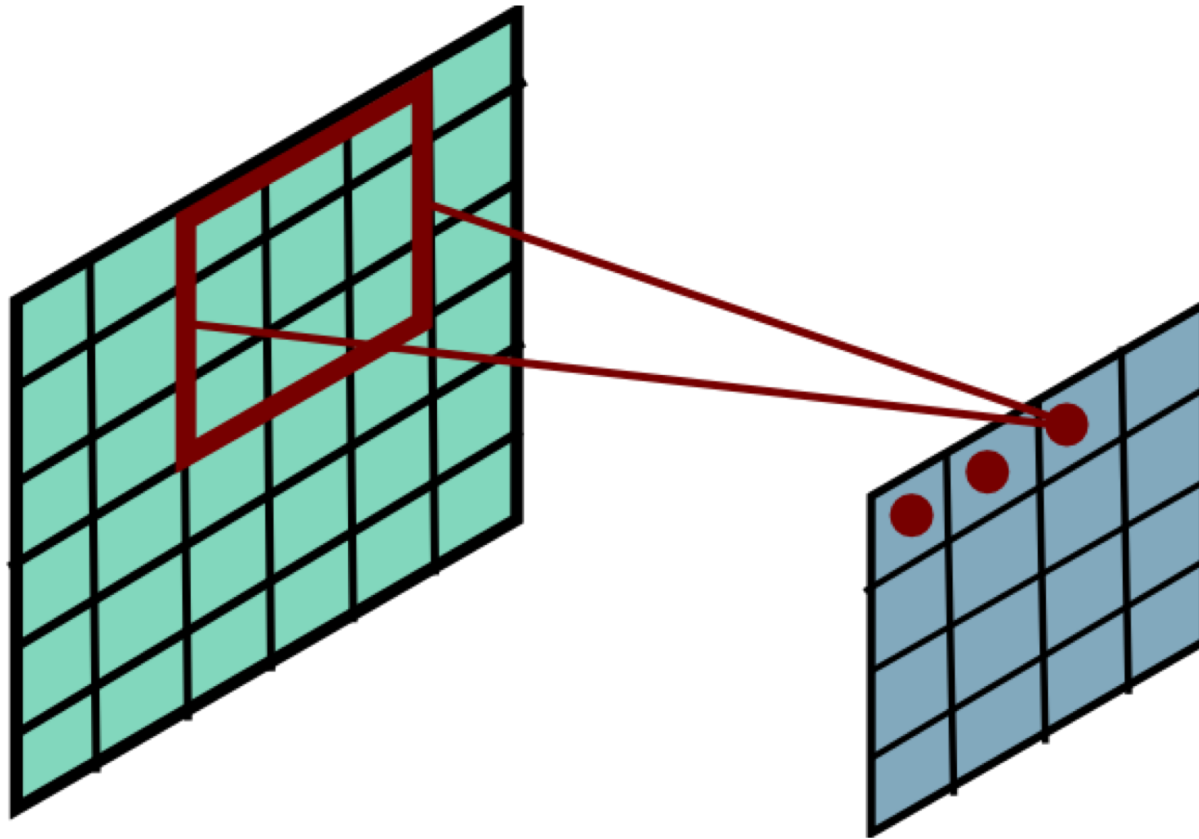
Convolutional Layer



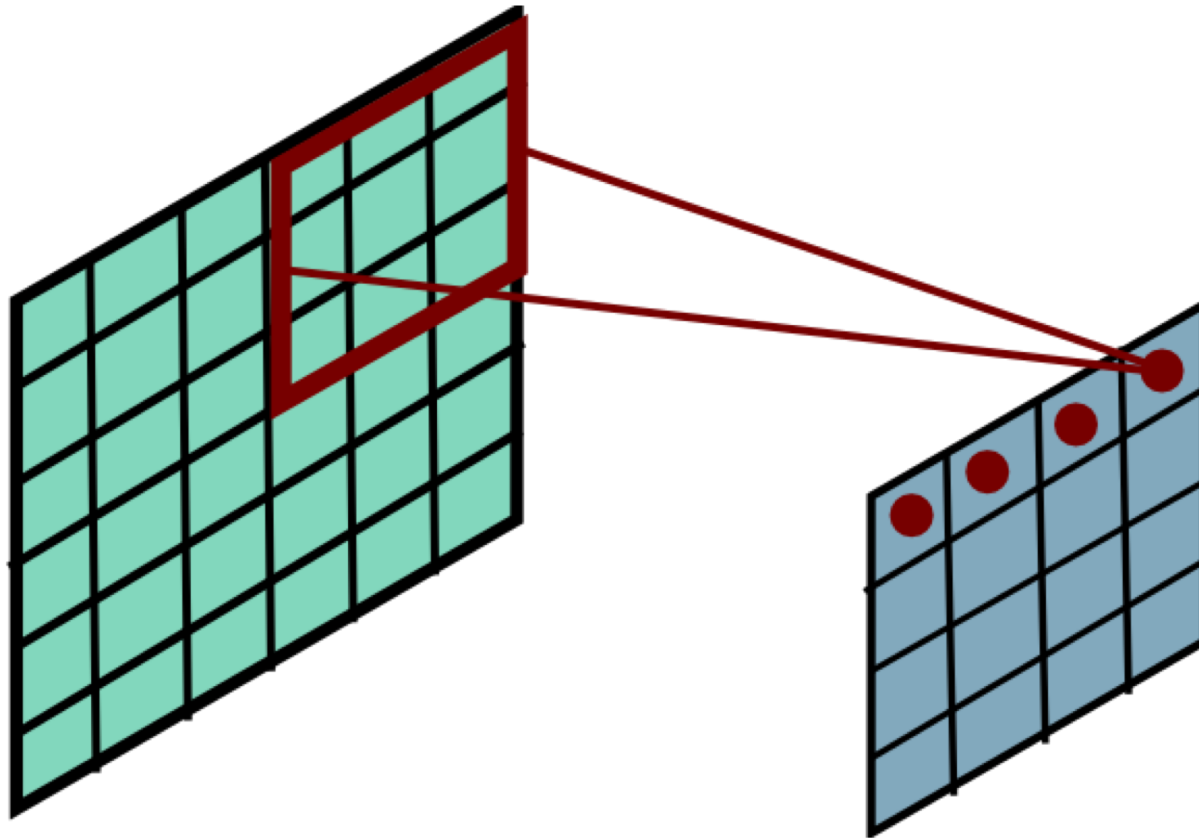
Convolutional Layer



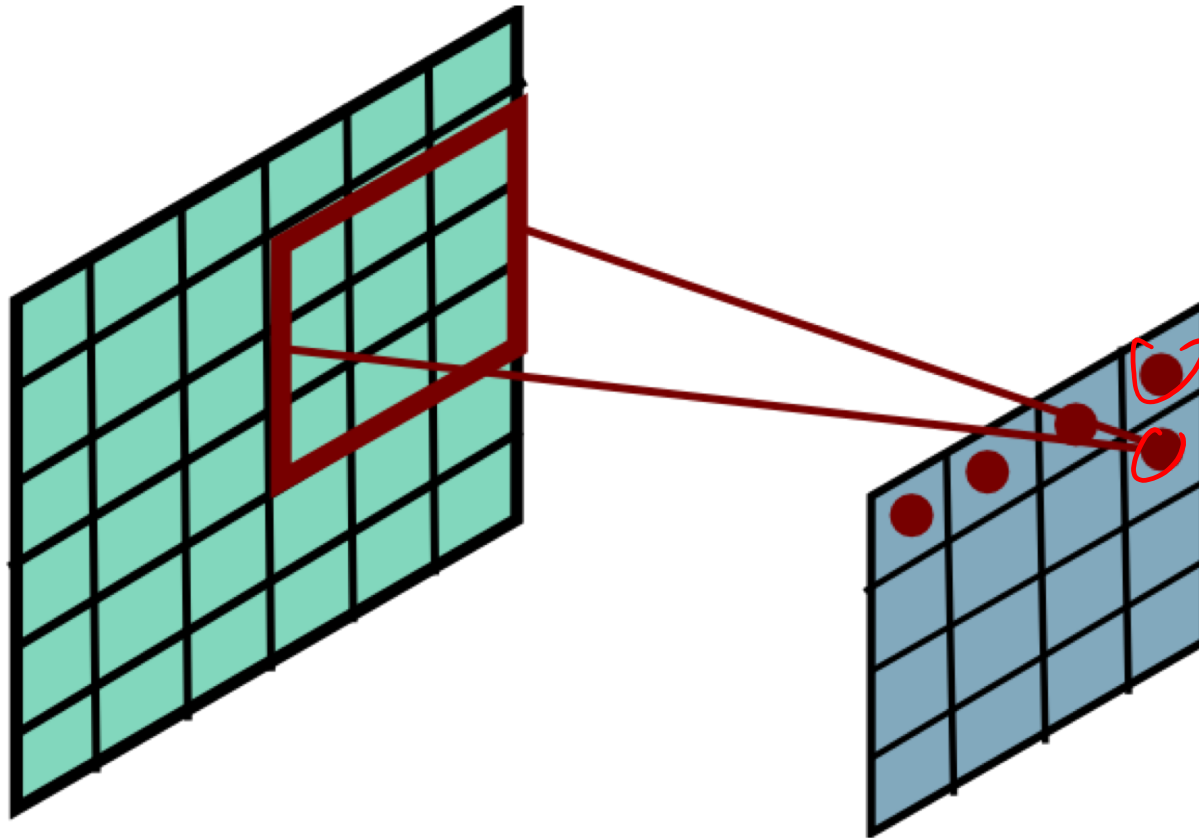
Convolutional Layer



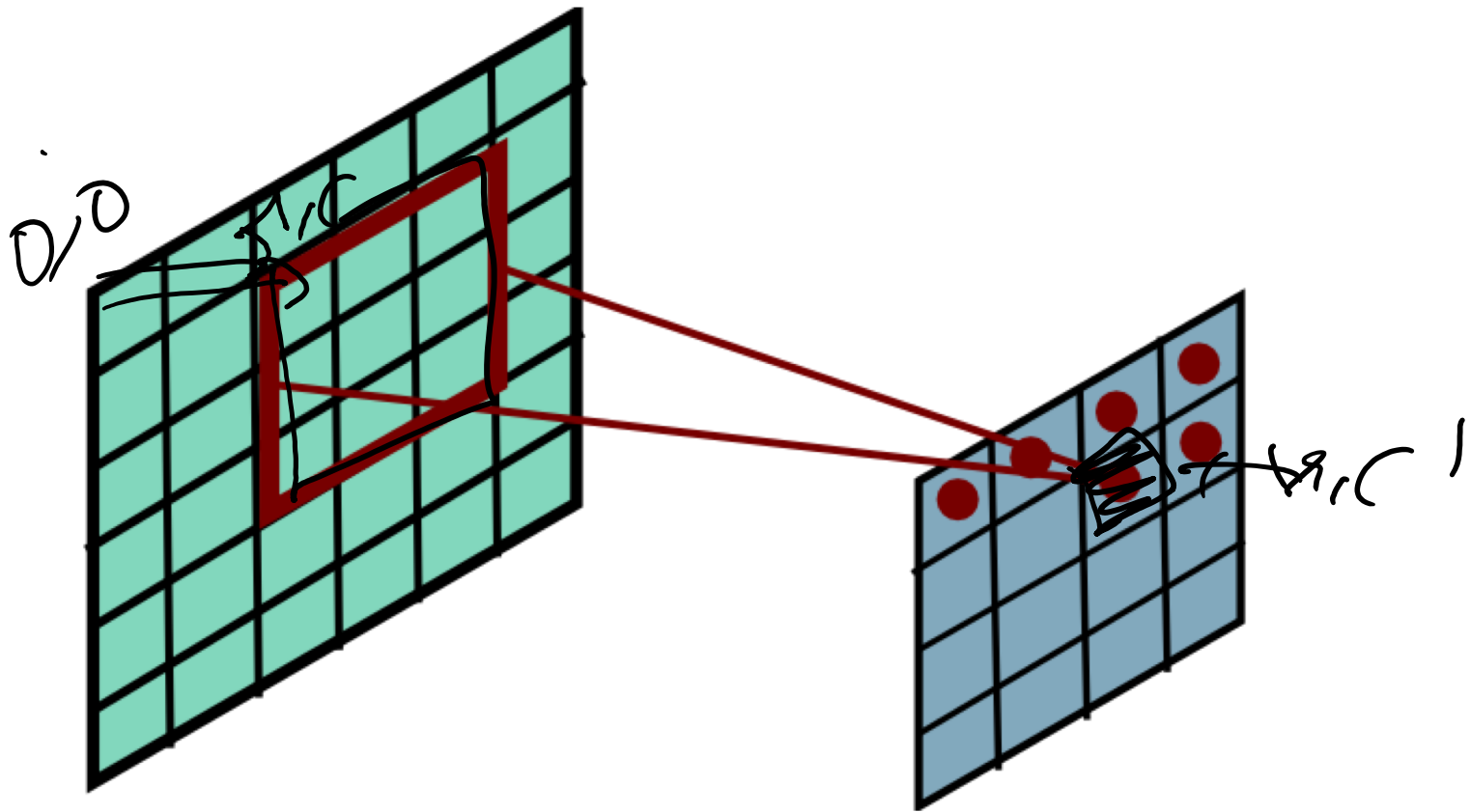
Convolutional Layer



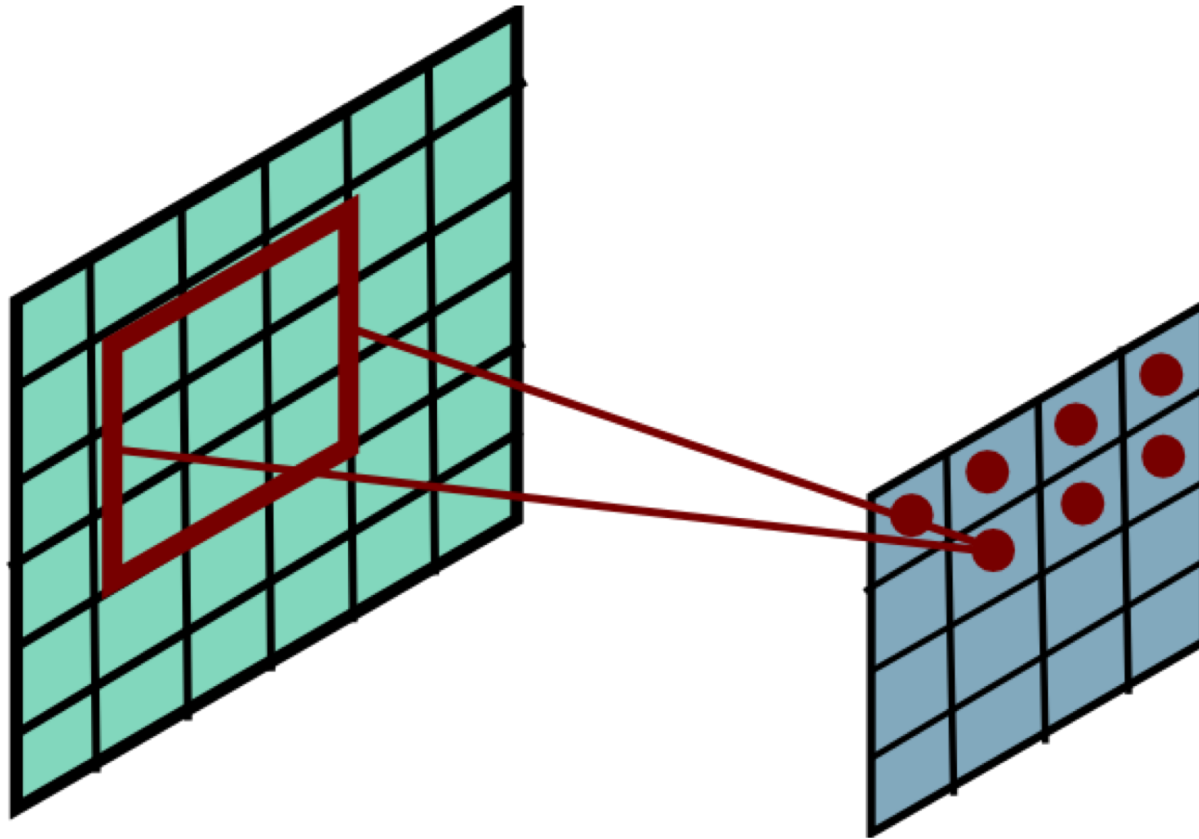
Convolutional Layer



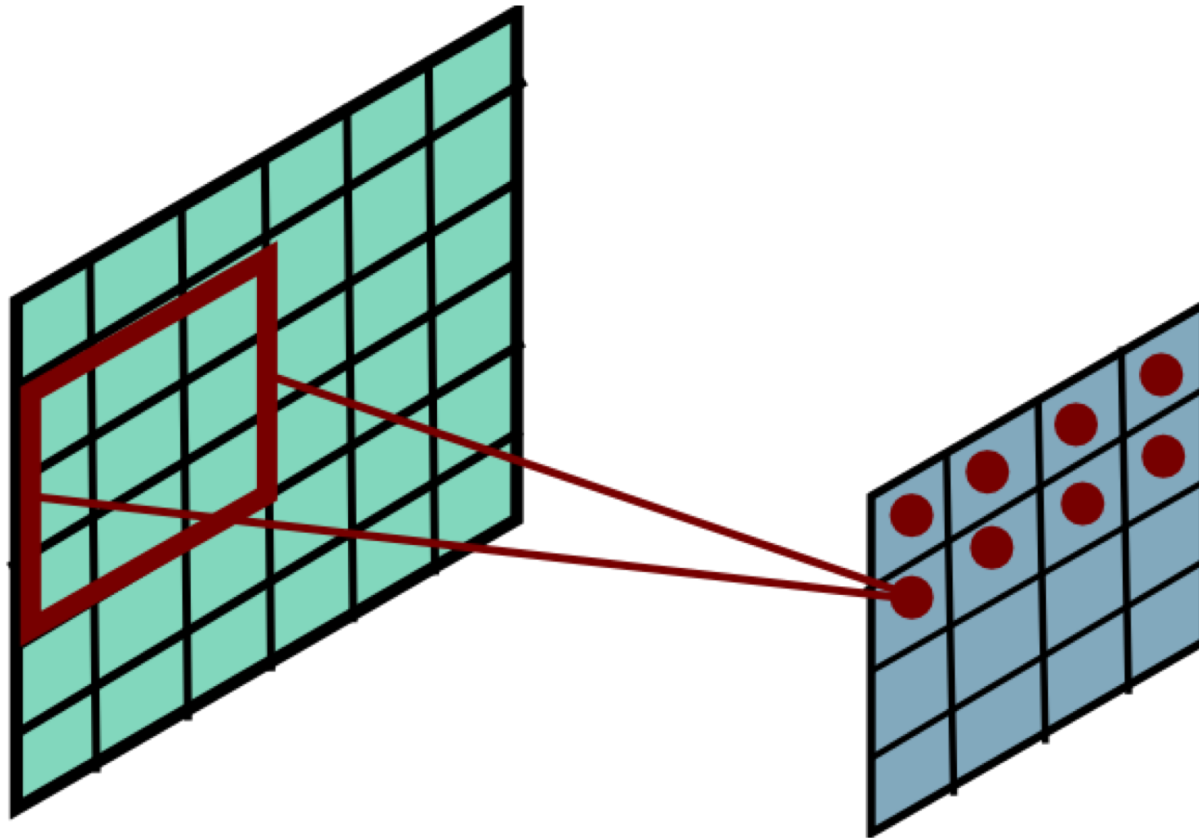
Convolutional Layer



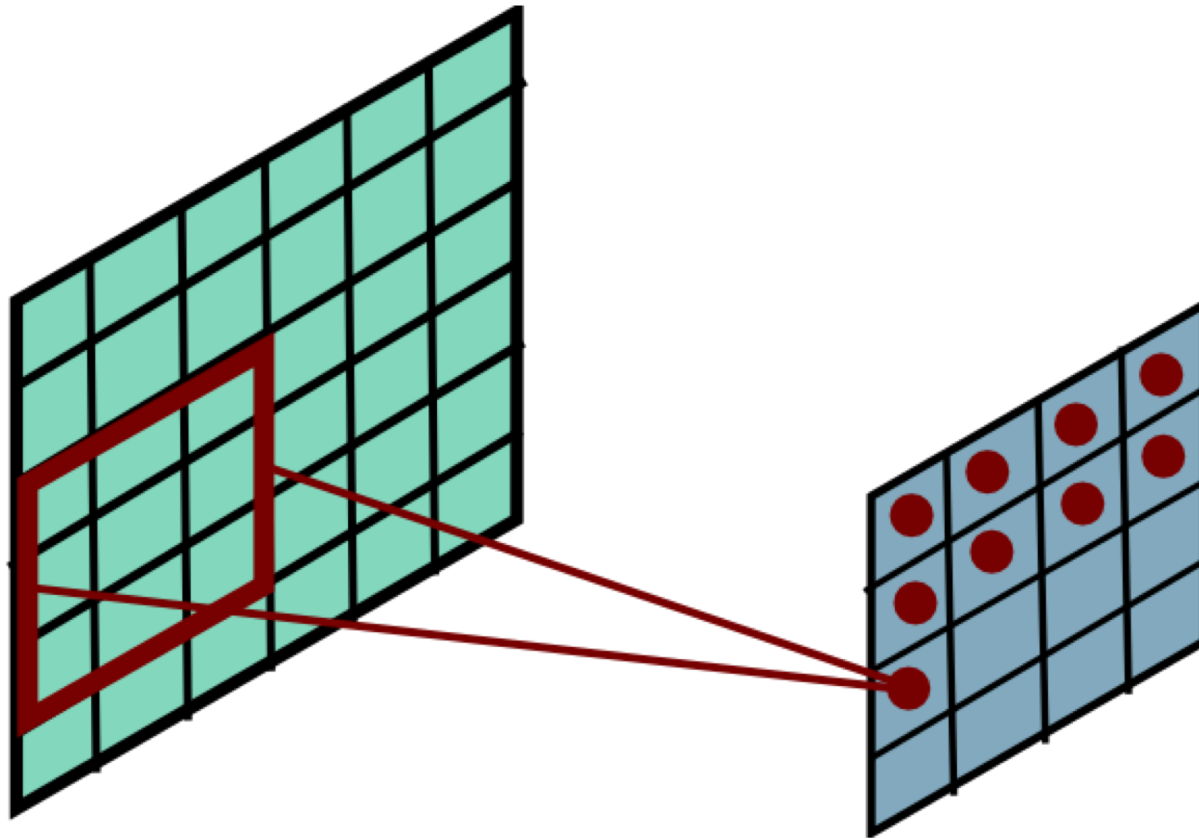
Convolutional Layer



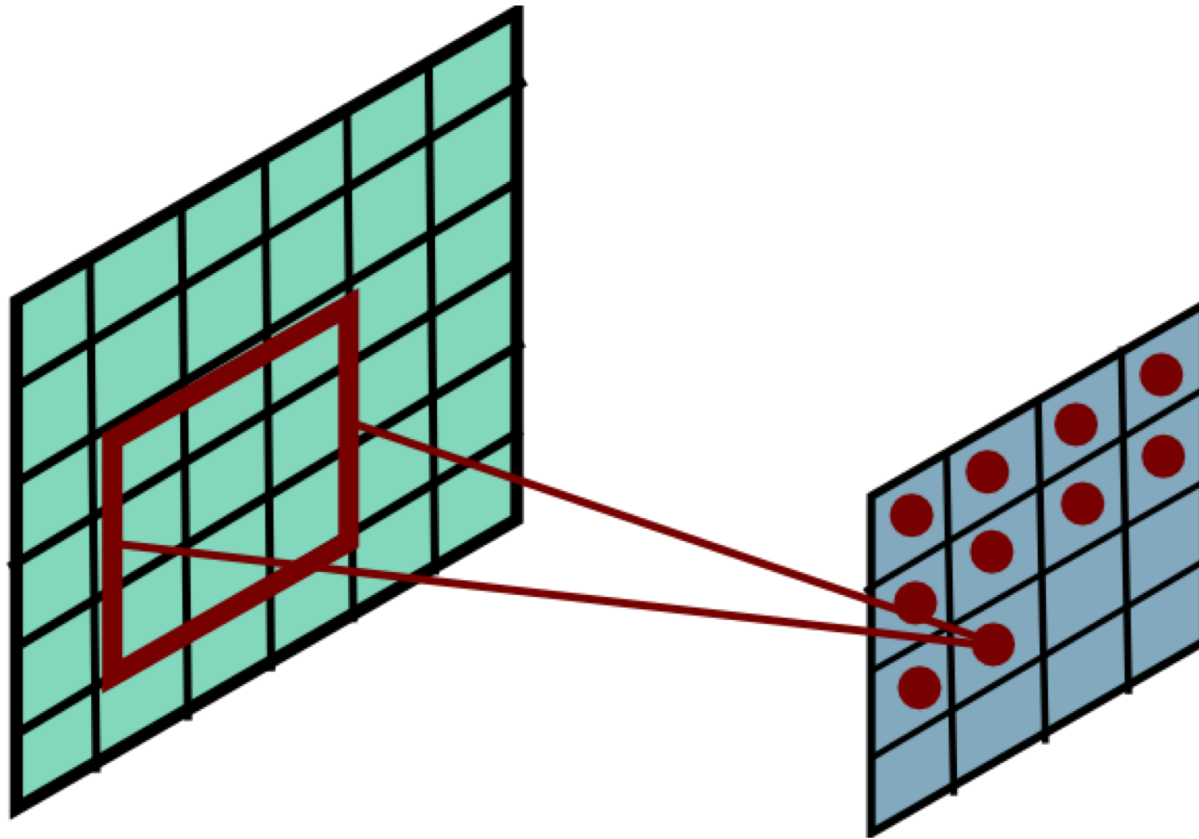
Convolutional Layer



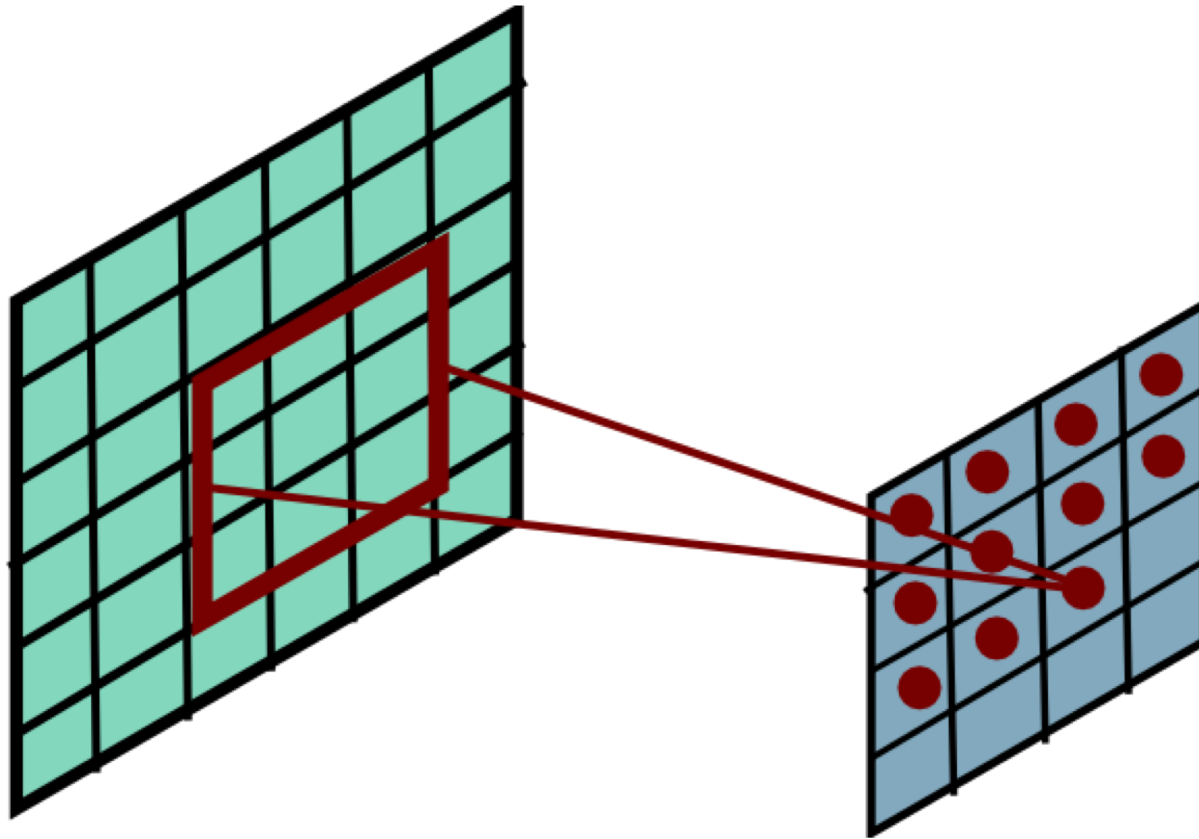
Convolutional Layer



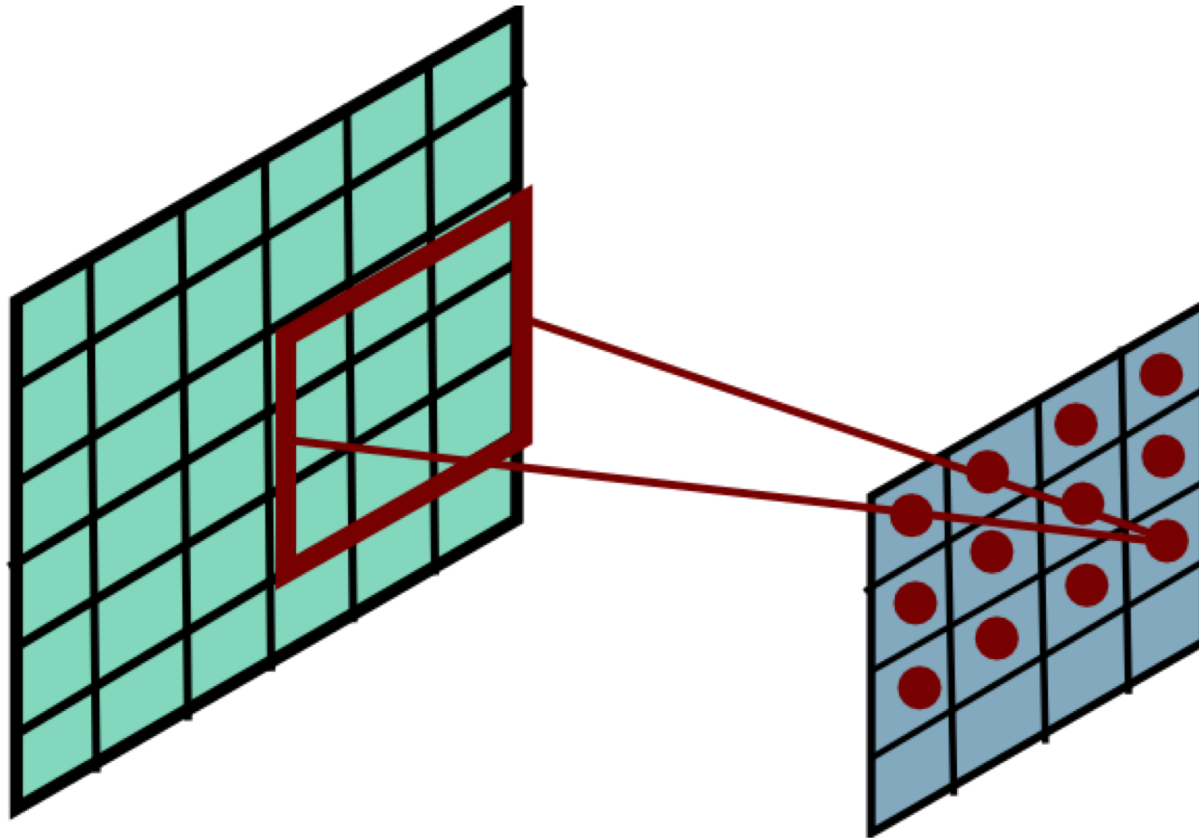
Convolutional Layer



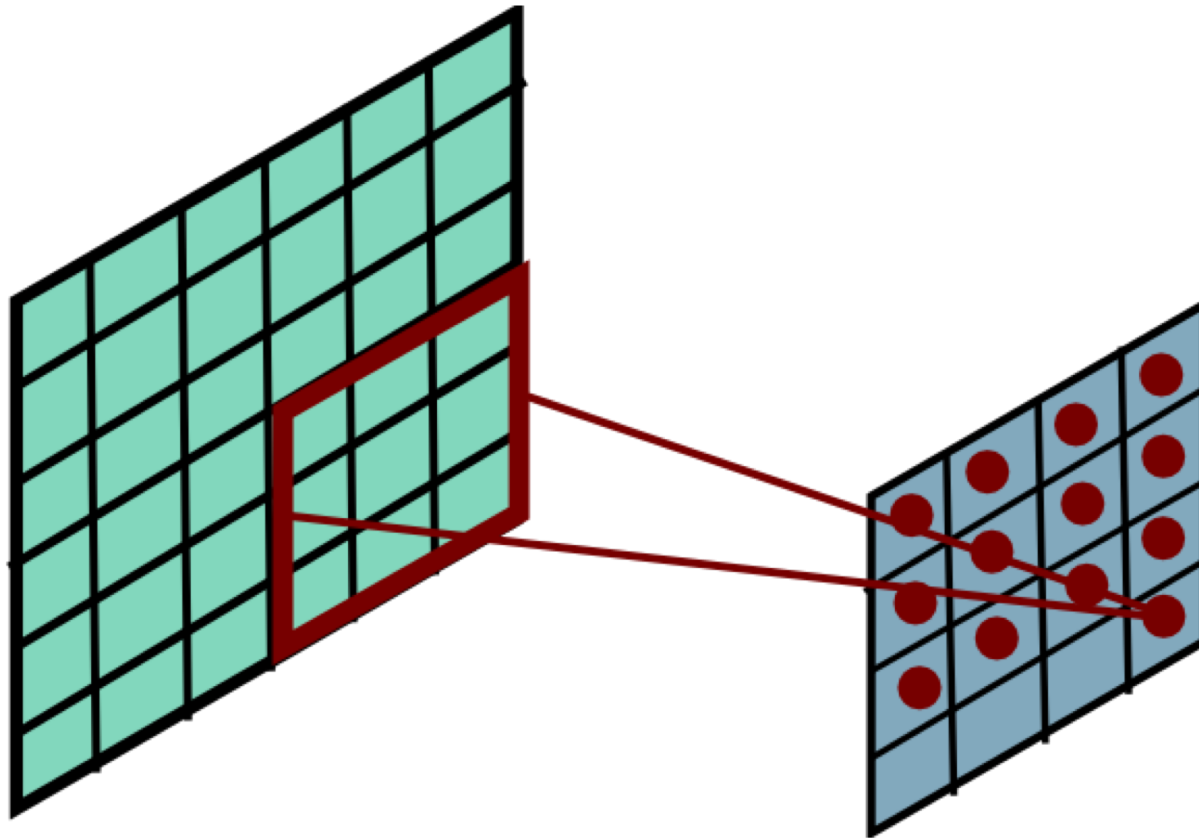
Convolutional Layer



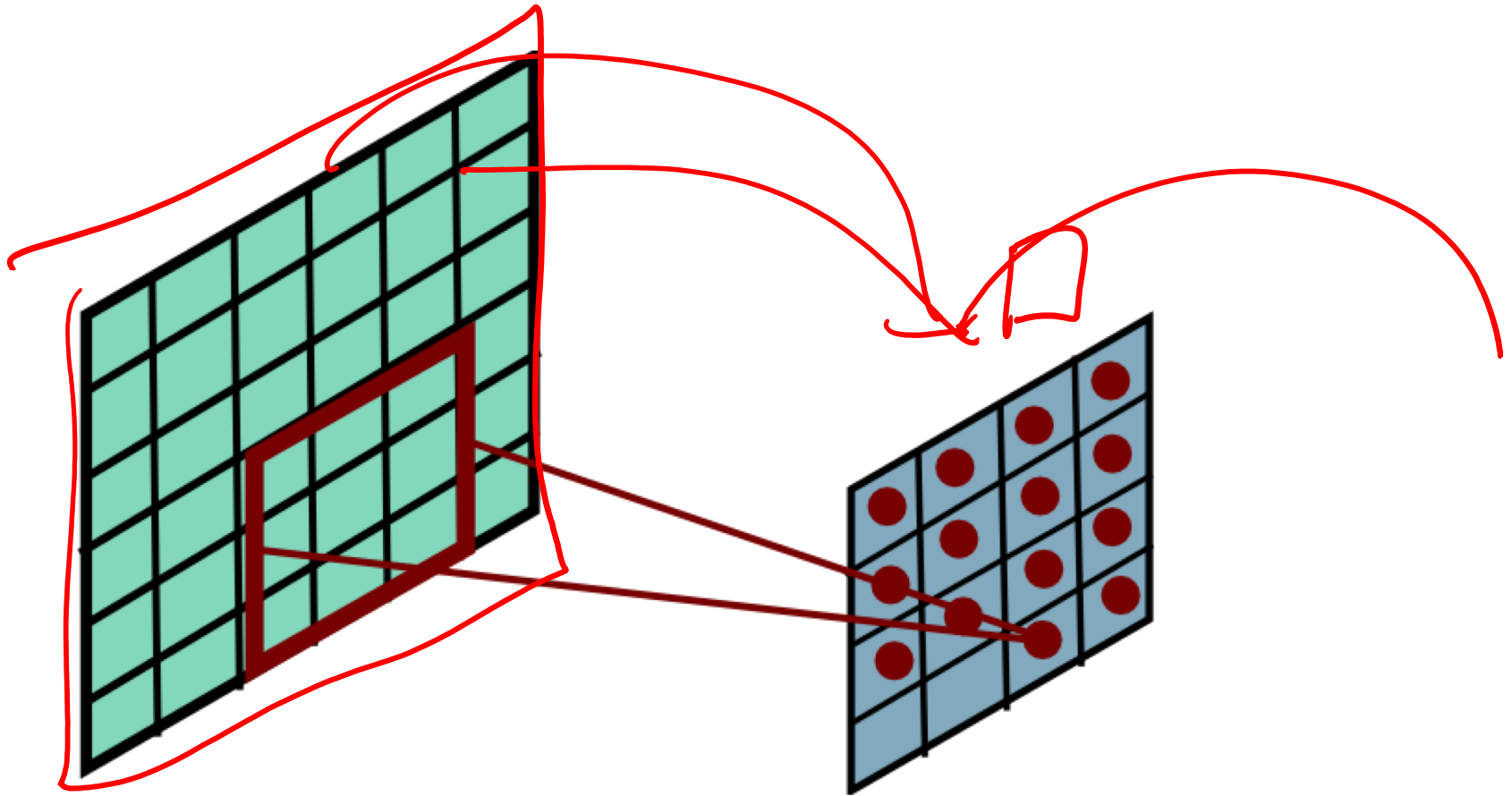
Convolutional Layer



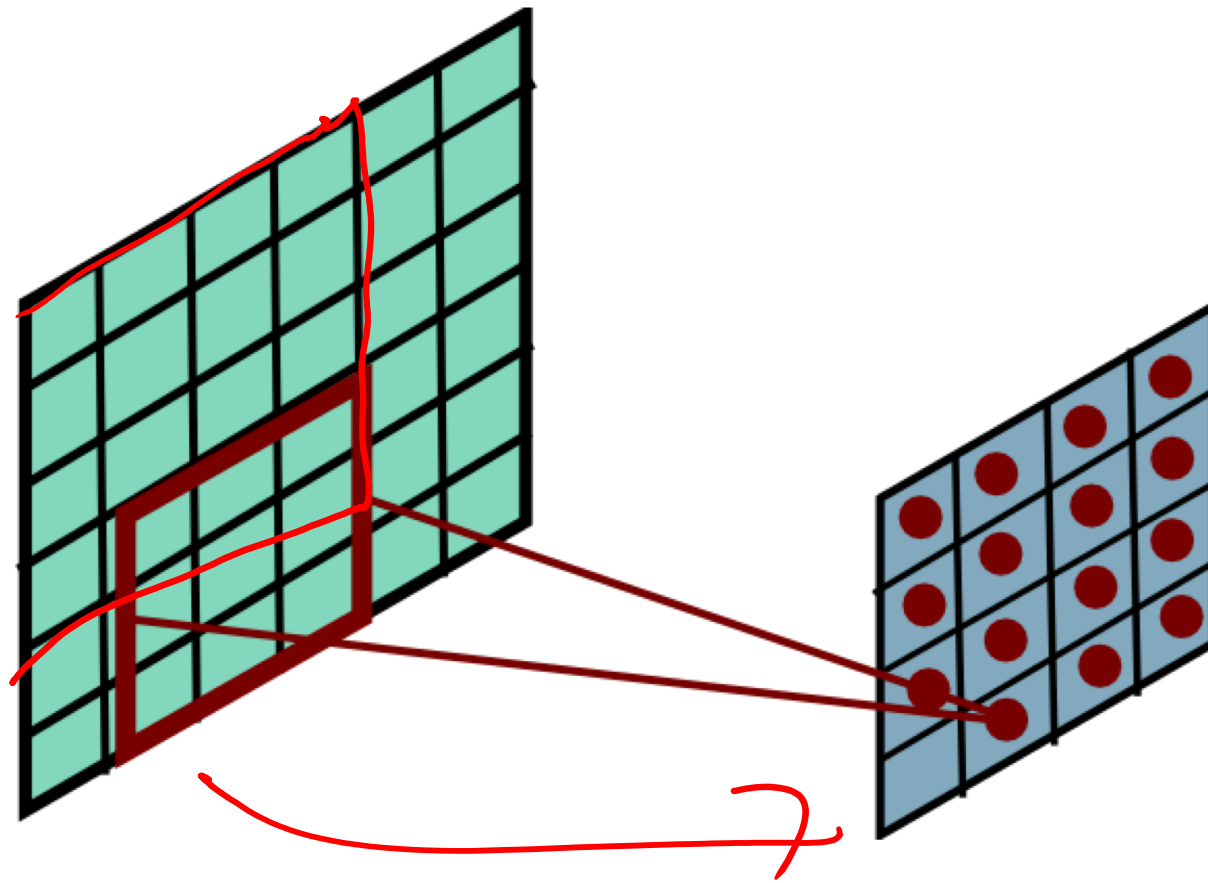
Convolutional Layer



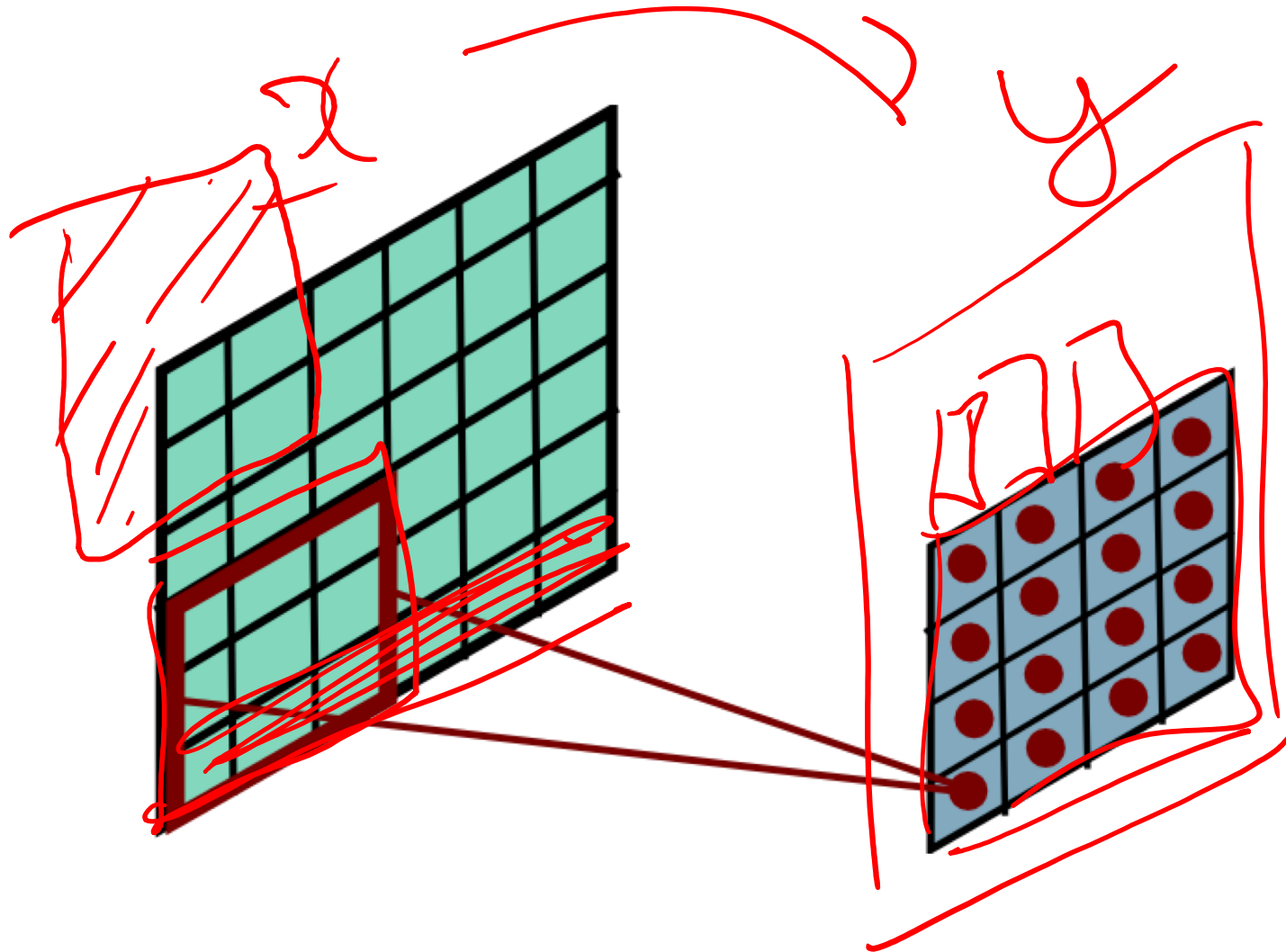
Convolutional Layer



Convolutional Layer



Convolutional Layer



Mathieu et al. "Fast training of CNNs through FFTs" ICLR 2014

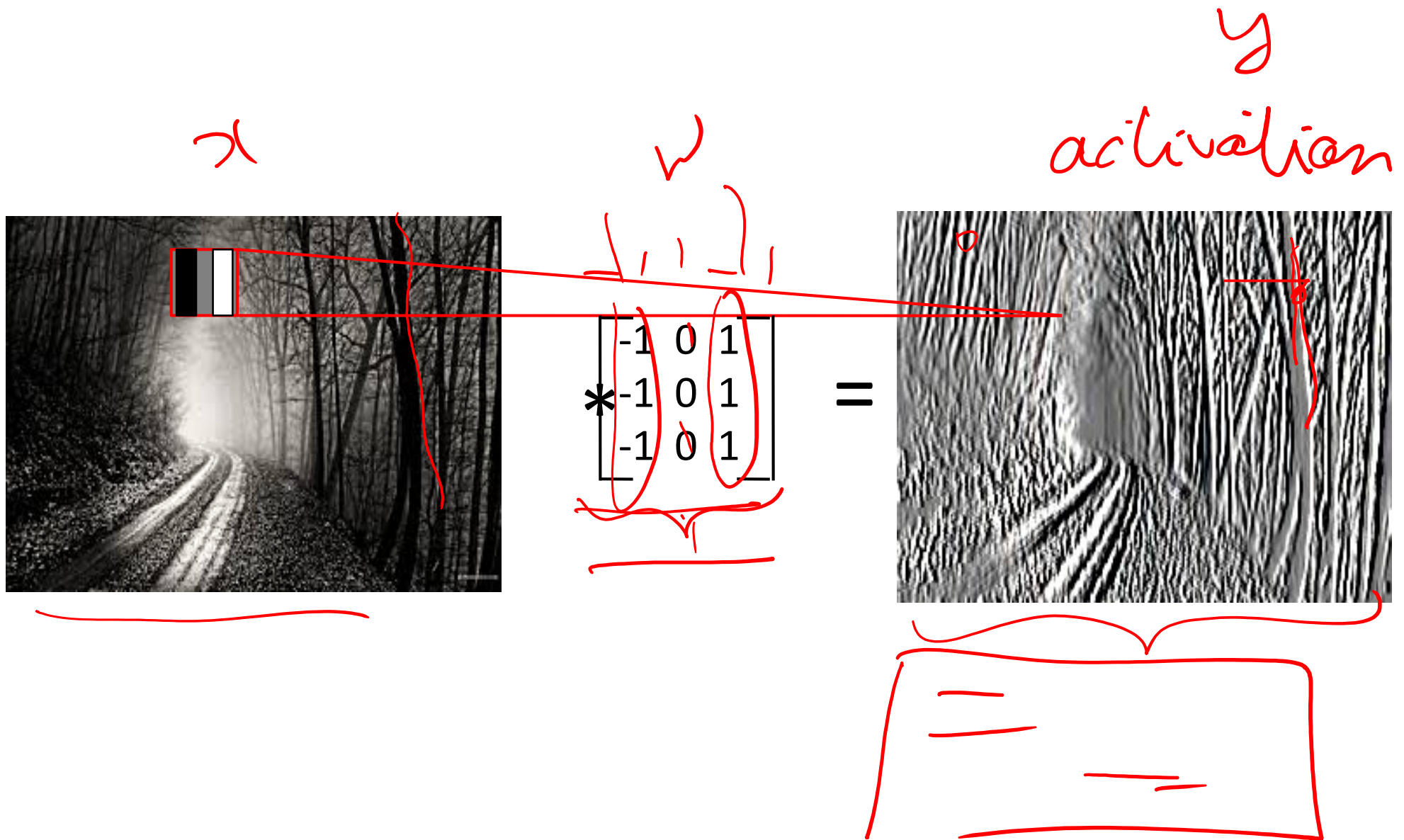
Plan for Today

- Convolutional Neural Networks
 - Stride, padding
 - Pooling layers
 - Fully-connected layers as convolutions

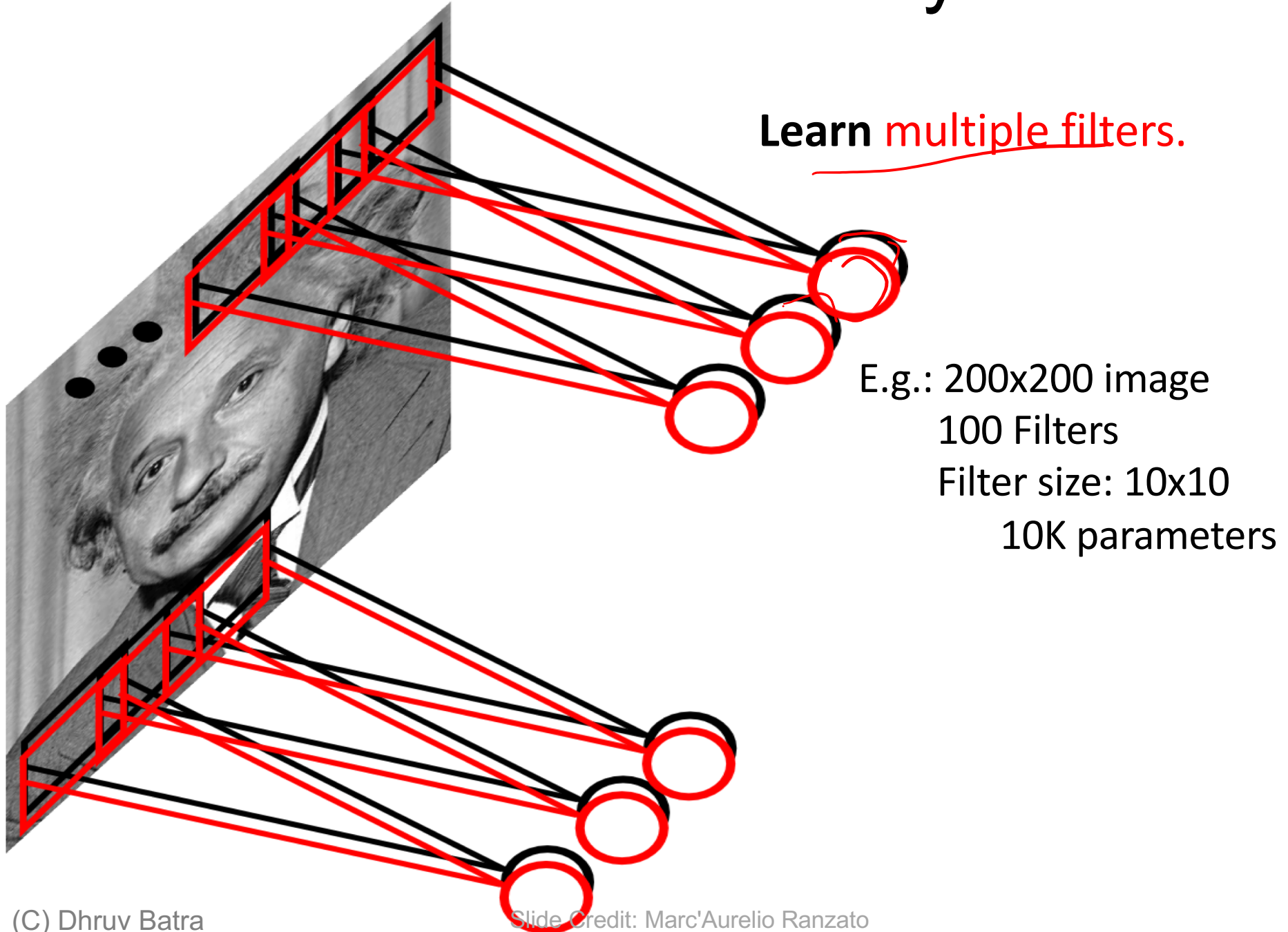
Convolution Explained

- <http://setosa.io/ev/image-kernels/>
- <https://github.com/bruckner/deepViz>

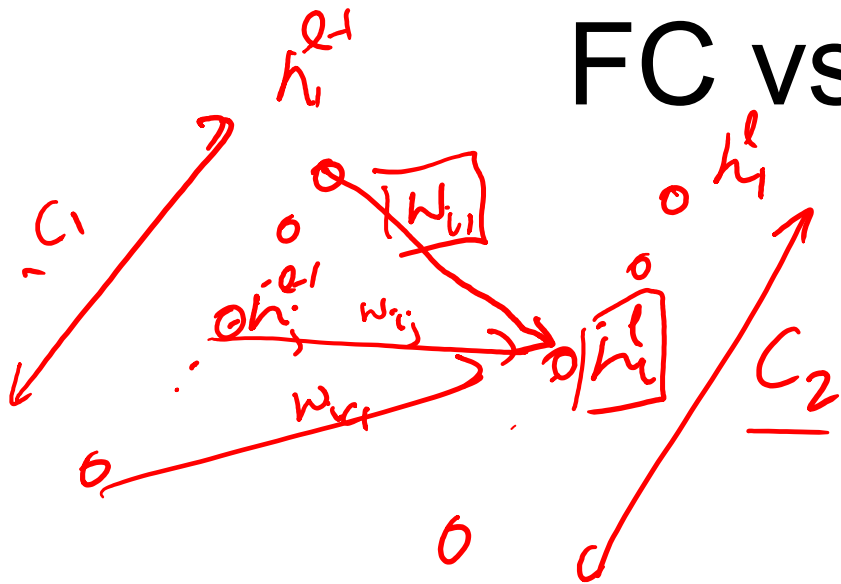
Convolutional Layer



Convolutional Layer



FC vs Conv Layer

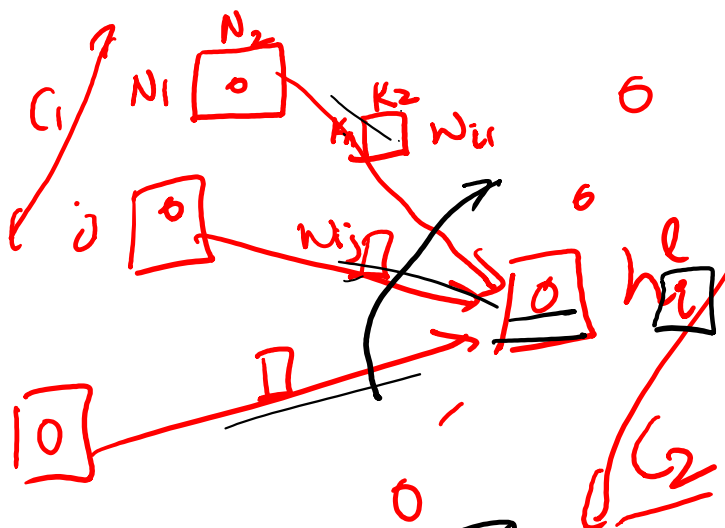


$$h_i^l = \sum_{j=1}^{C_1} h_j^{l-1} w_{ij} + b_i$$

Scalar product

~~conv~~

$$h_i^l = \sum_{j=1}^{C_1} h_j^{l-1} w_{ij} + b_i$$



$$h_i^l [r, c] = \sum_{j=1}^{C_1} \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} h_j^{l-1} [r+a, c+b] w[a, b]$$

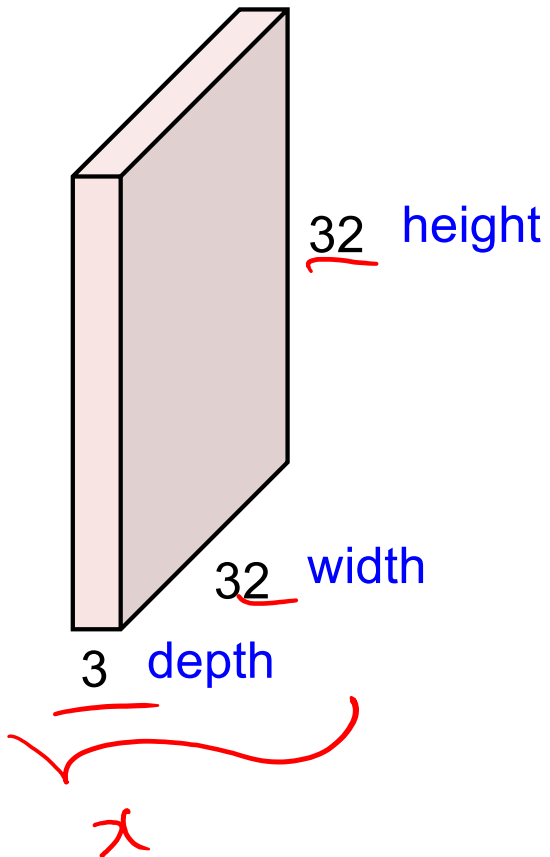


$$W \quad (k_1 \times k_2 \times C_1 \times C_2)$$

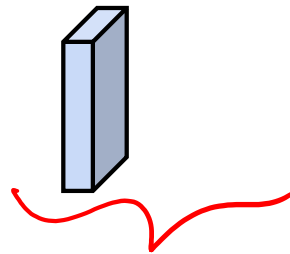
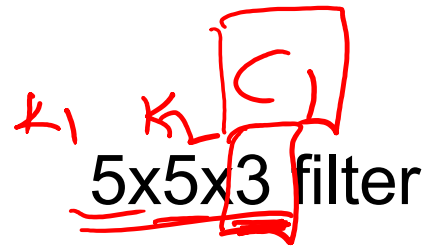
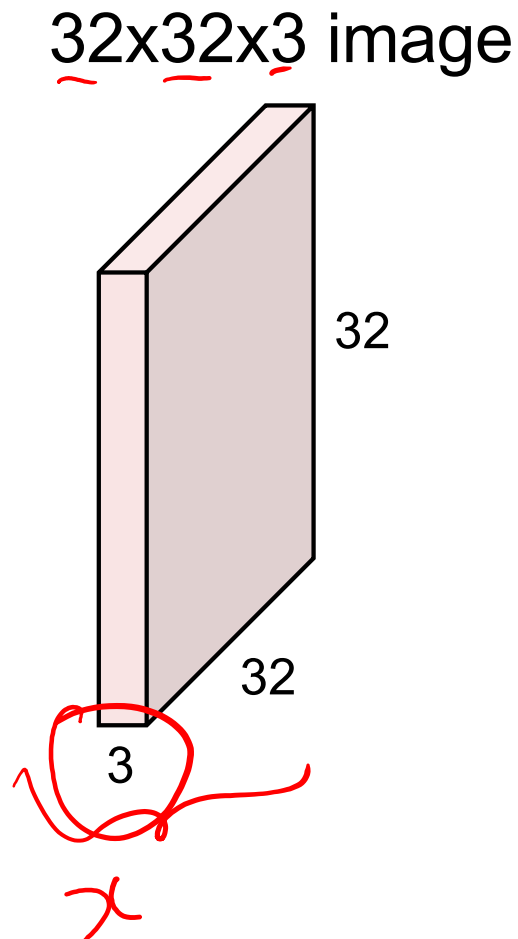
FC vs Conv Layer

Convolution Layer

32x32x3 image \rightarrow preserve spatial structure



Convolution Layer



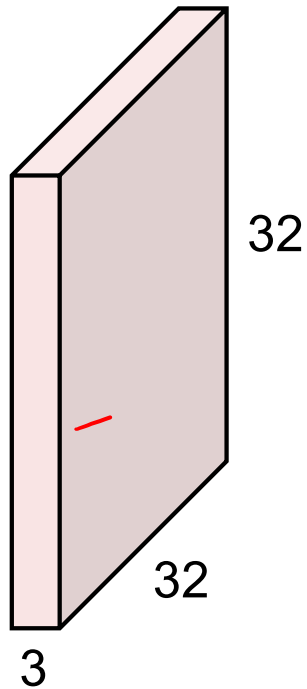
Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

w

b

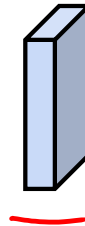
Convolution Layer

32x32x3 image



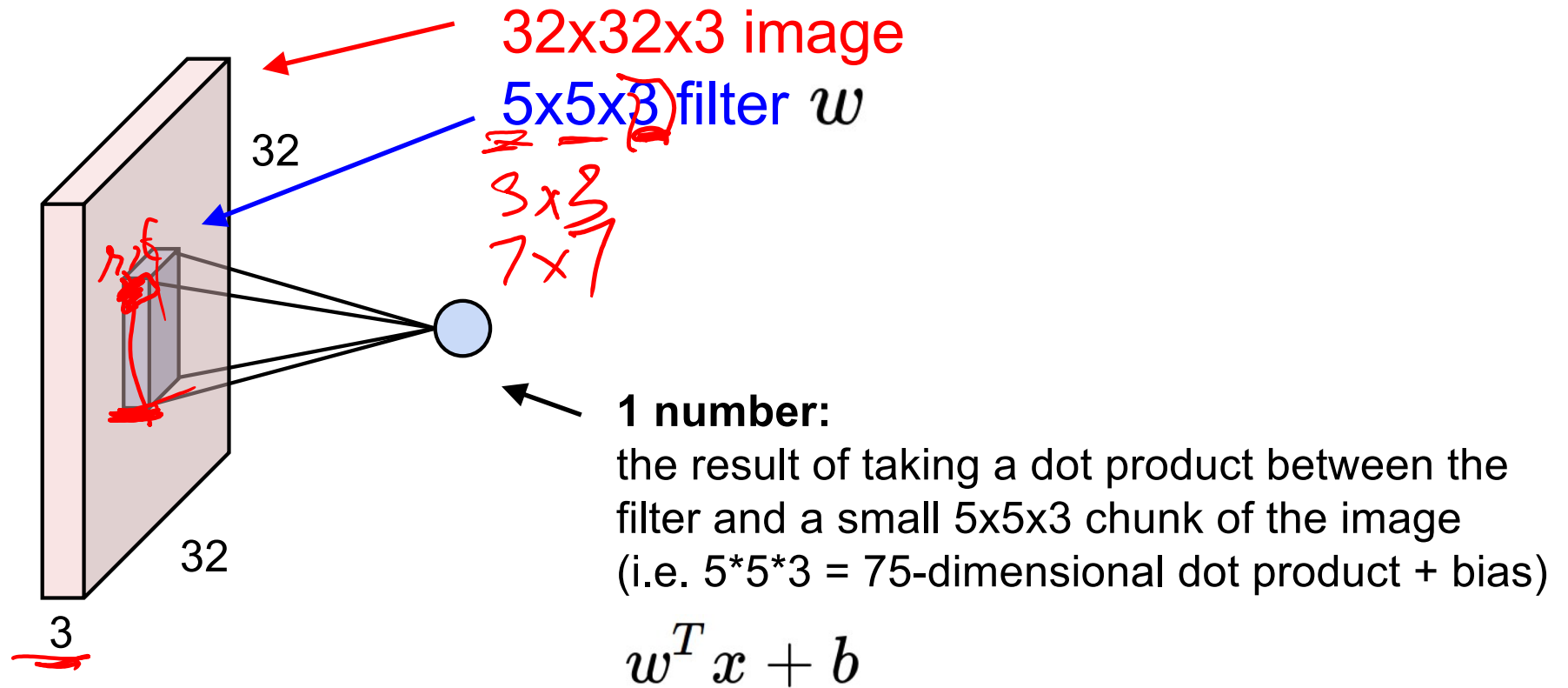
Filters always extend the full depth of the input volume

5x5x3 filter

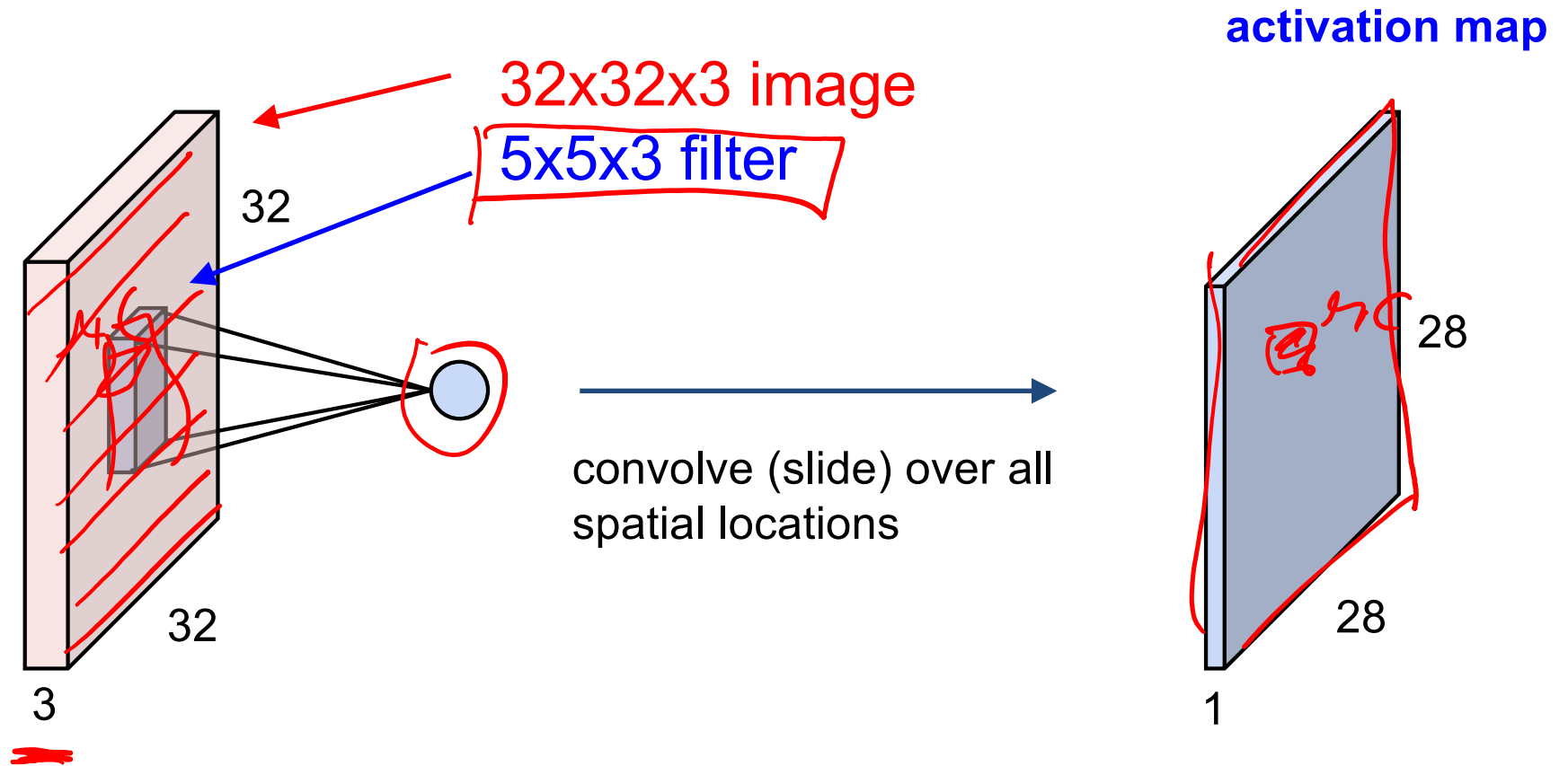


Convolve the filter with the image
i.e. “slide over the image spatially,
computing dot products”

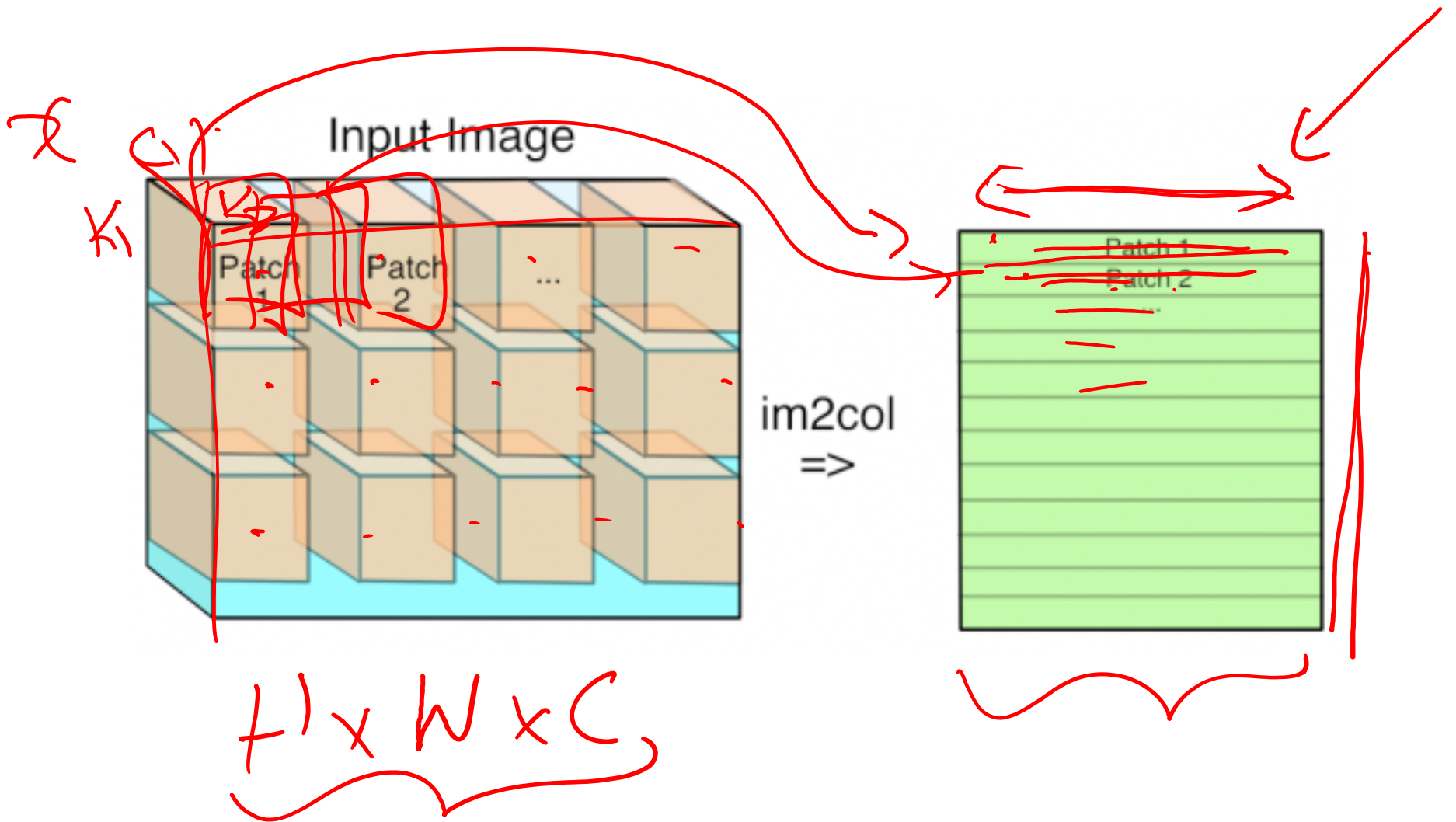
Convolution Layer



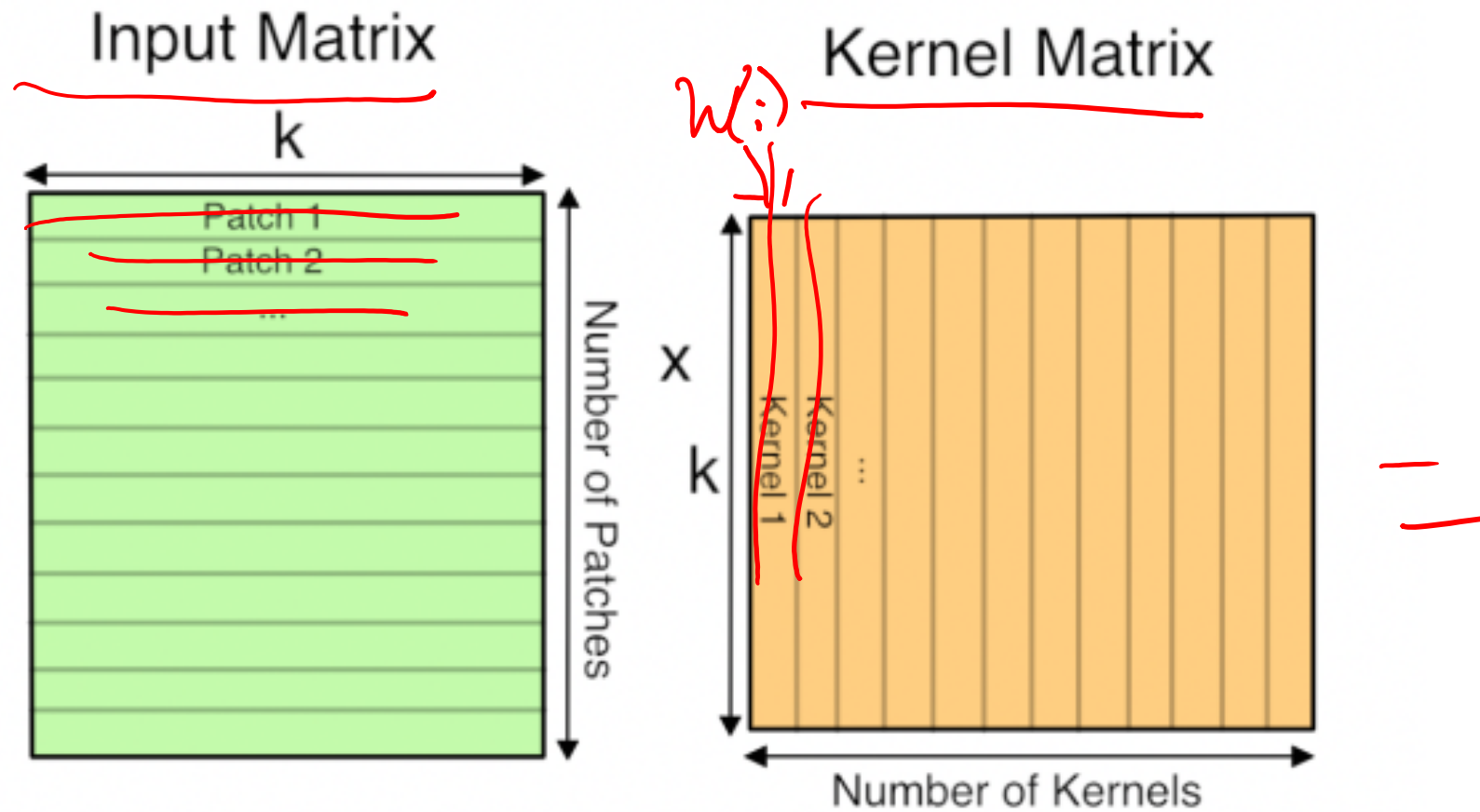
Convolution Layer



Im2Col

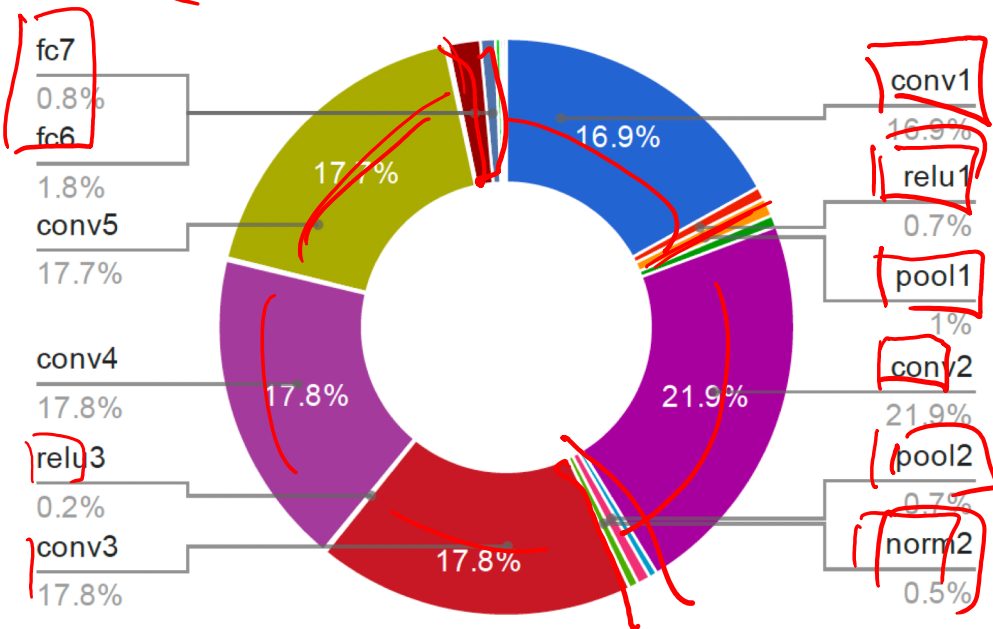


GEMM

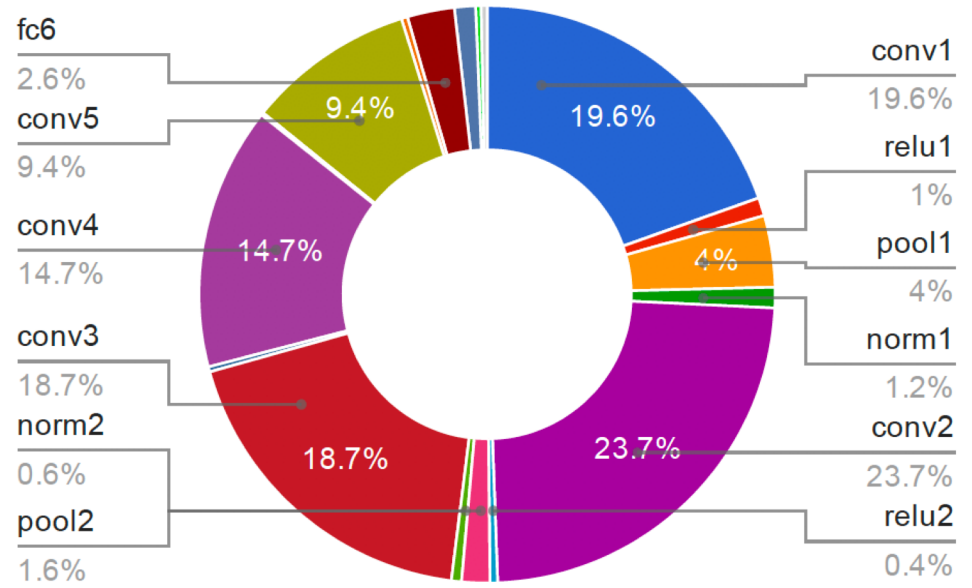


Time Distribution of AlexNet

GPU Forward Time Distribution

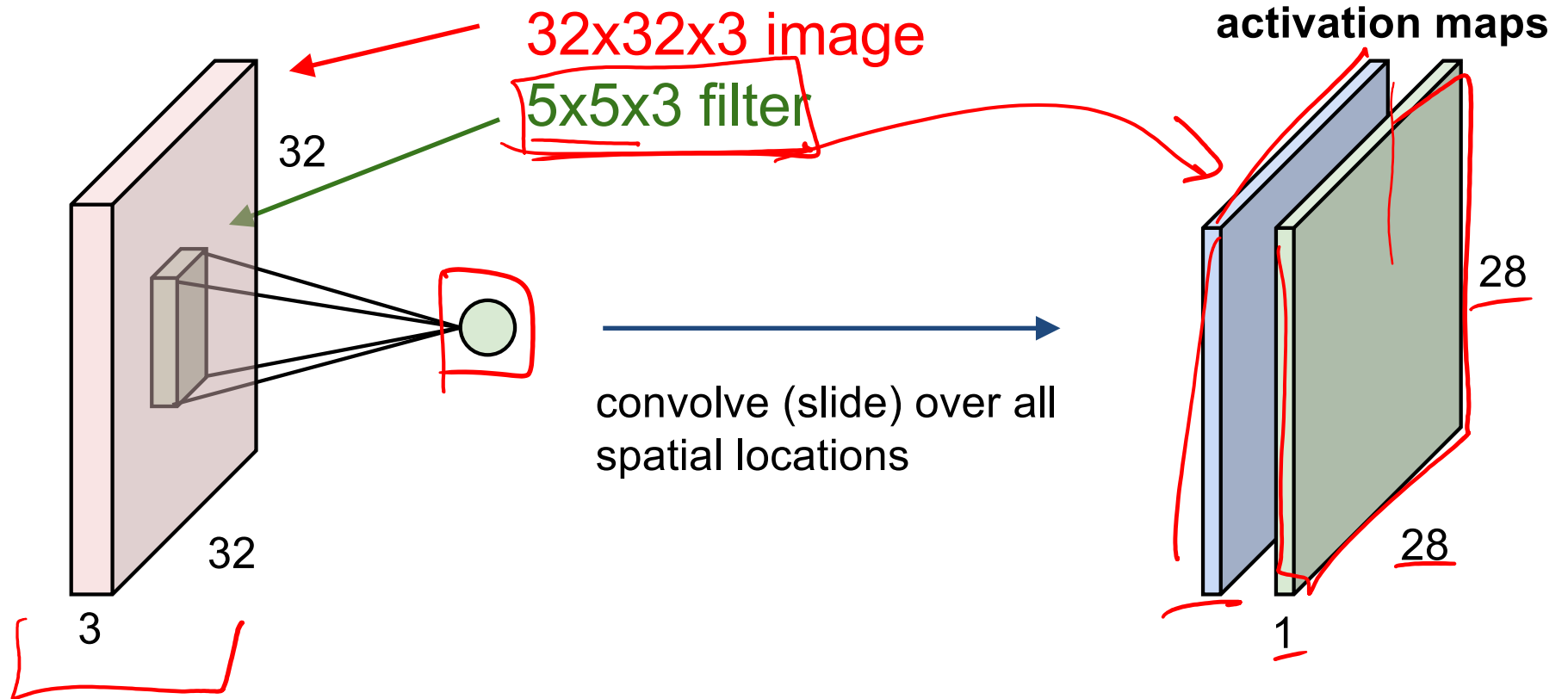


CPU Forward Time Distribution

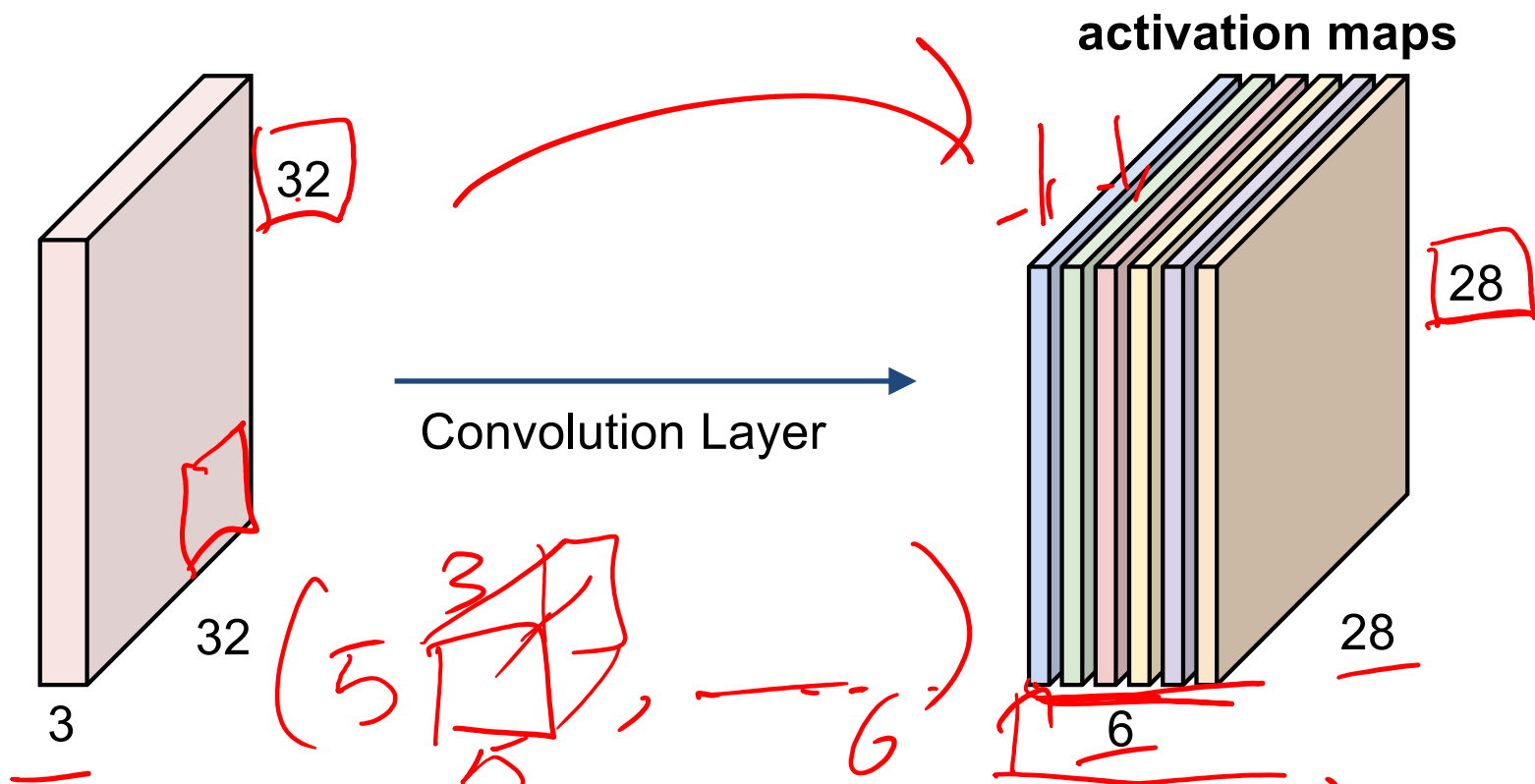


Convolution Layer

consider a second, **green** filter

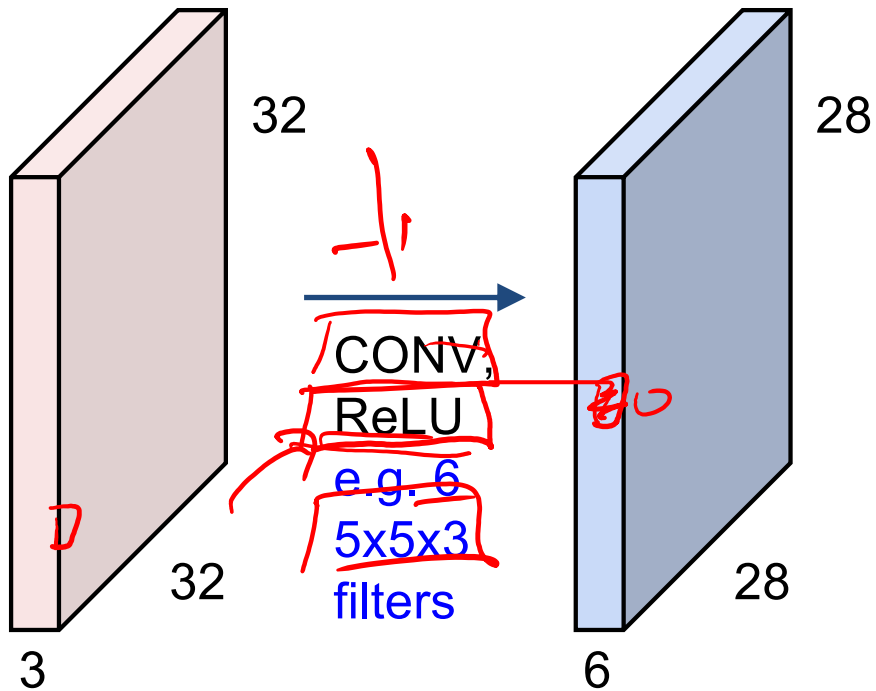


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:



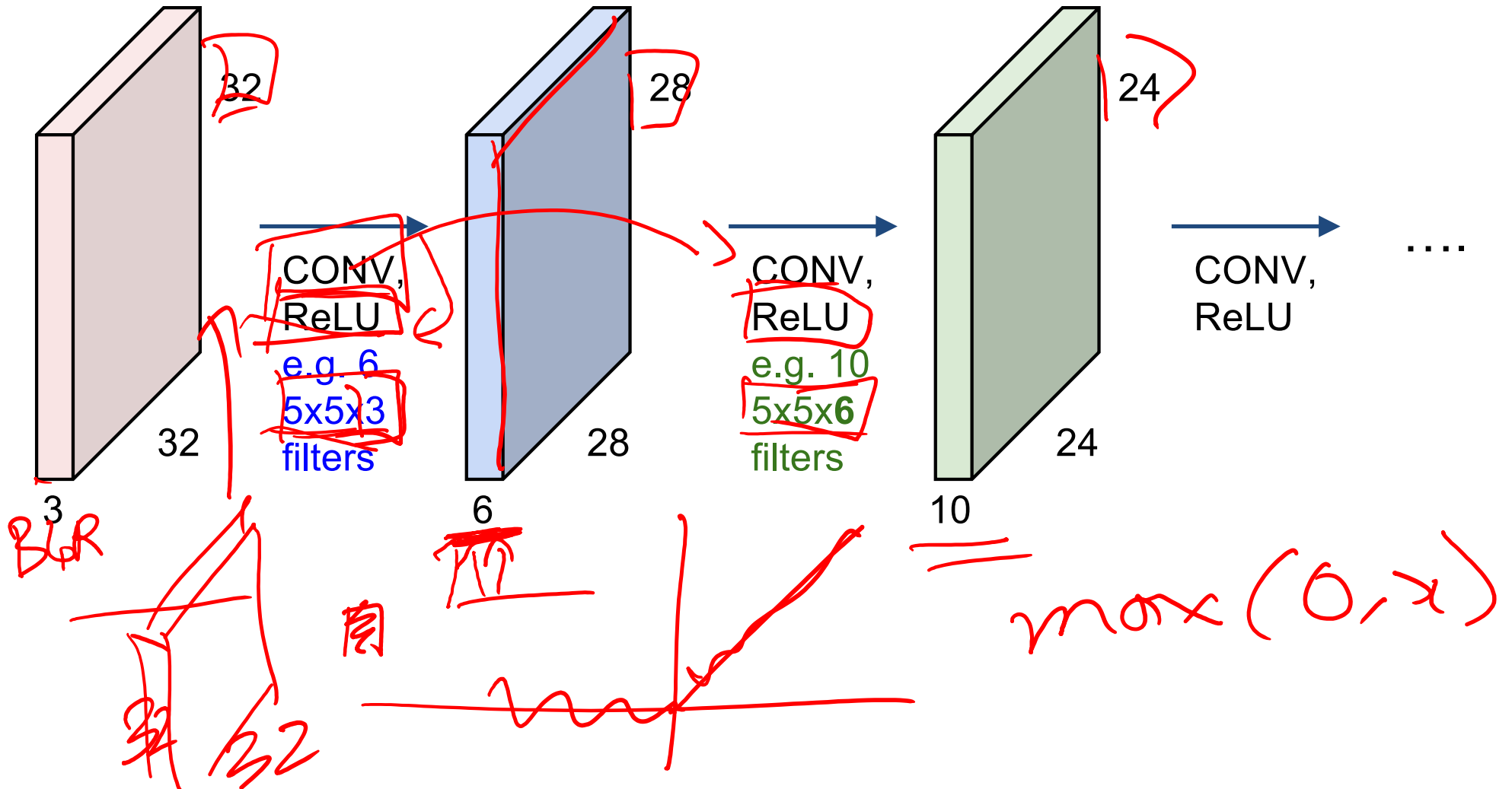
We stack these up to get a "new image" of size 28x28x6!

Preview: ConvNet is a sequence of Convolution Layers, interspersed with activation functions

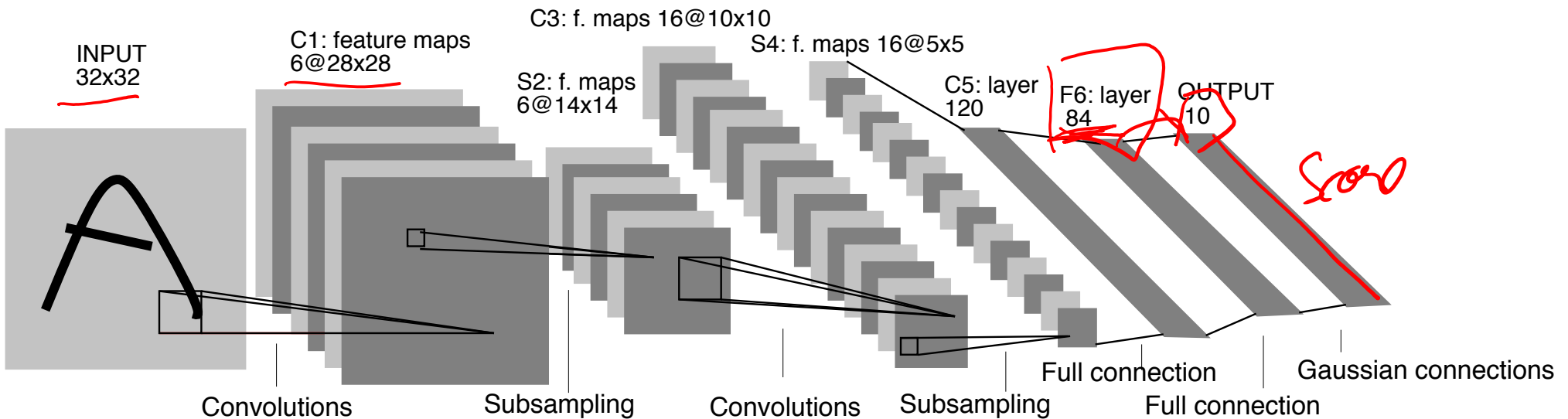
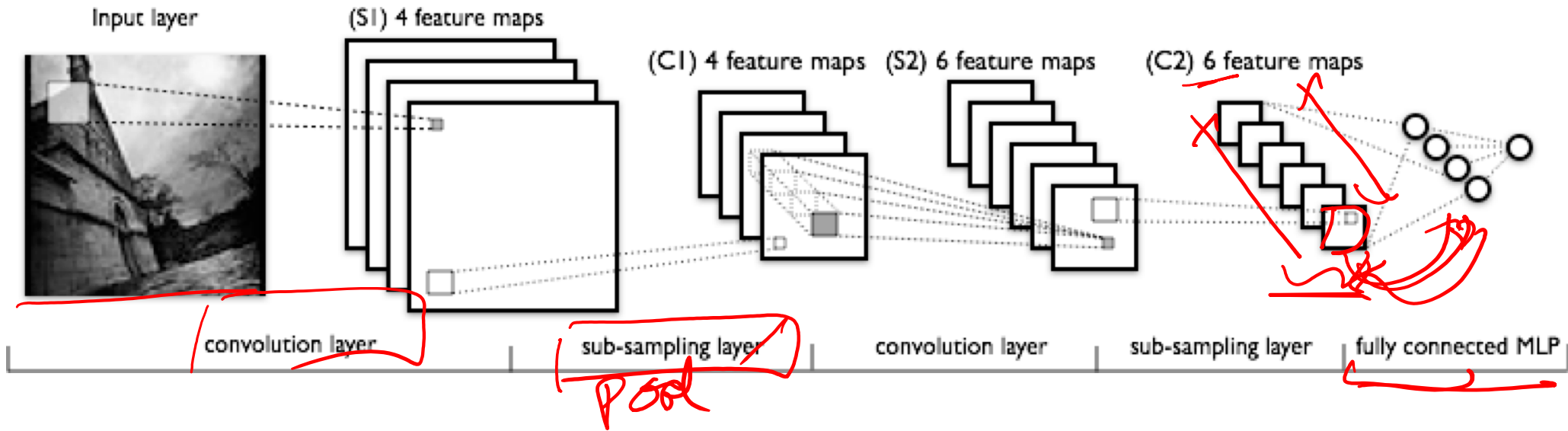


$$\max(0, x)$$

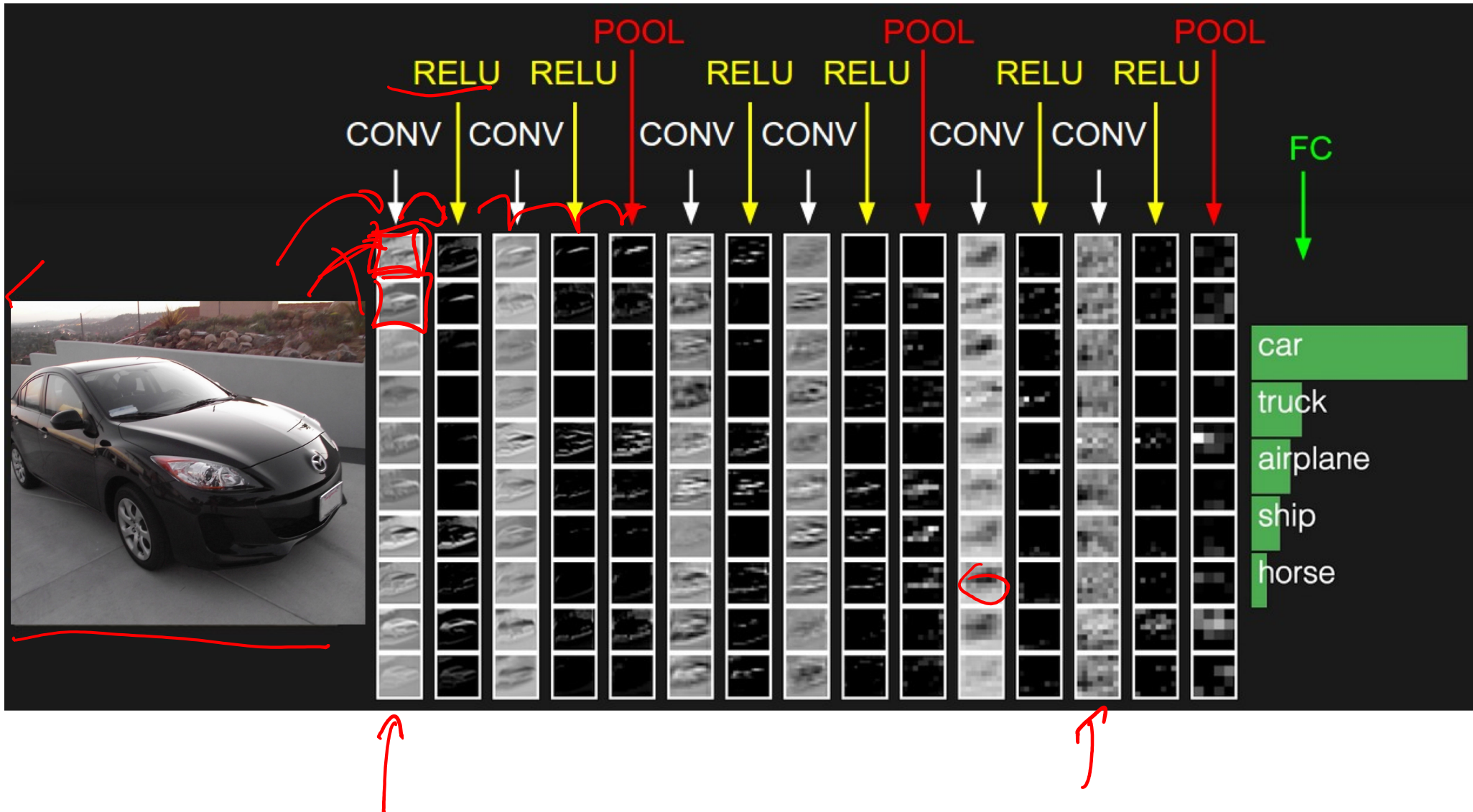
Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions



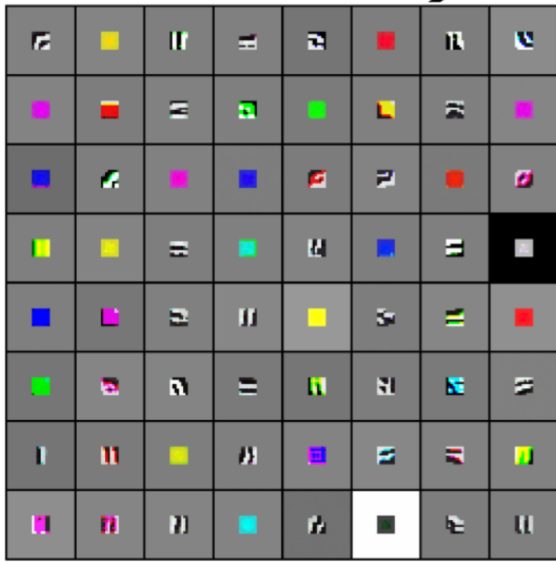
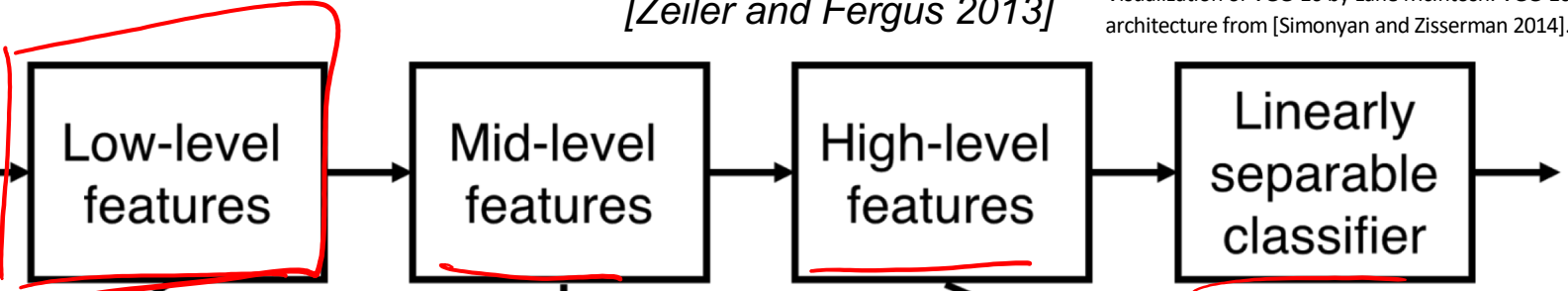
Convolutional Neural Networks



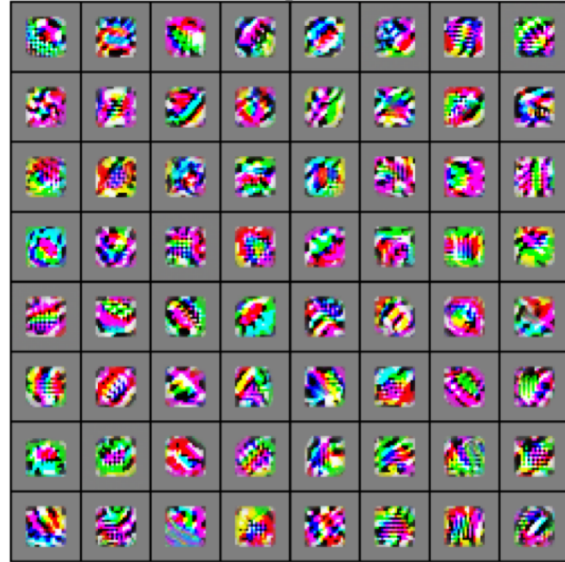
preview:



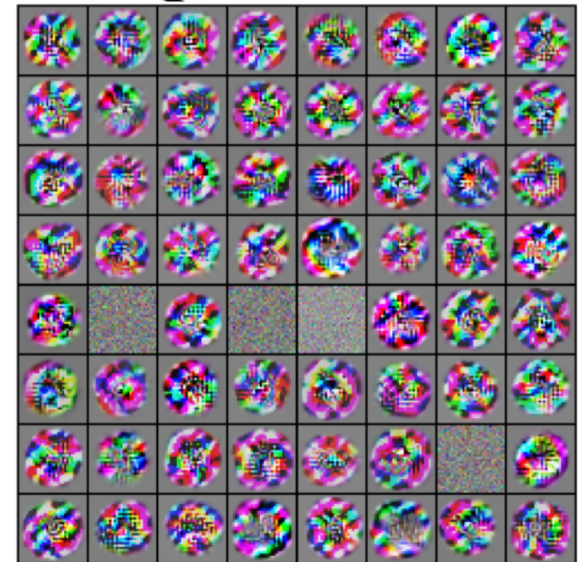
Preview



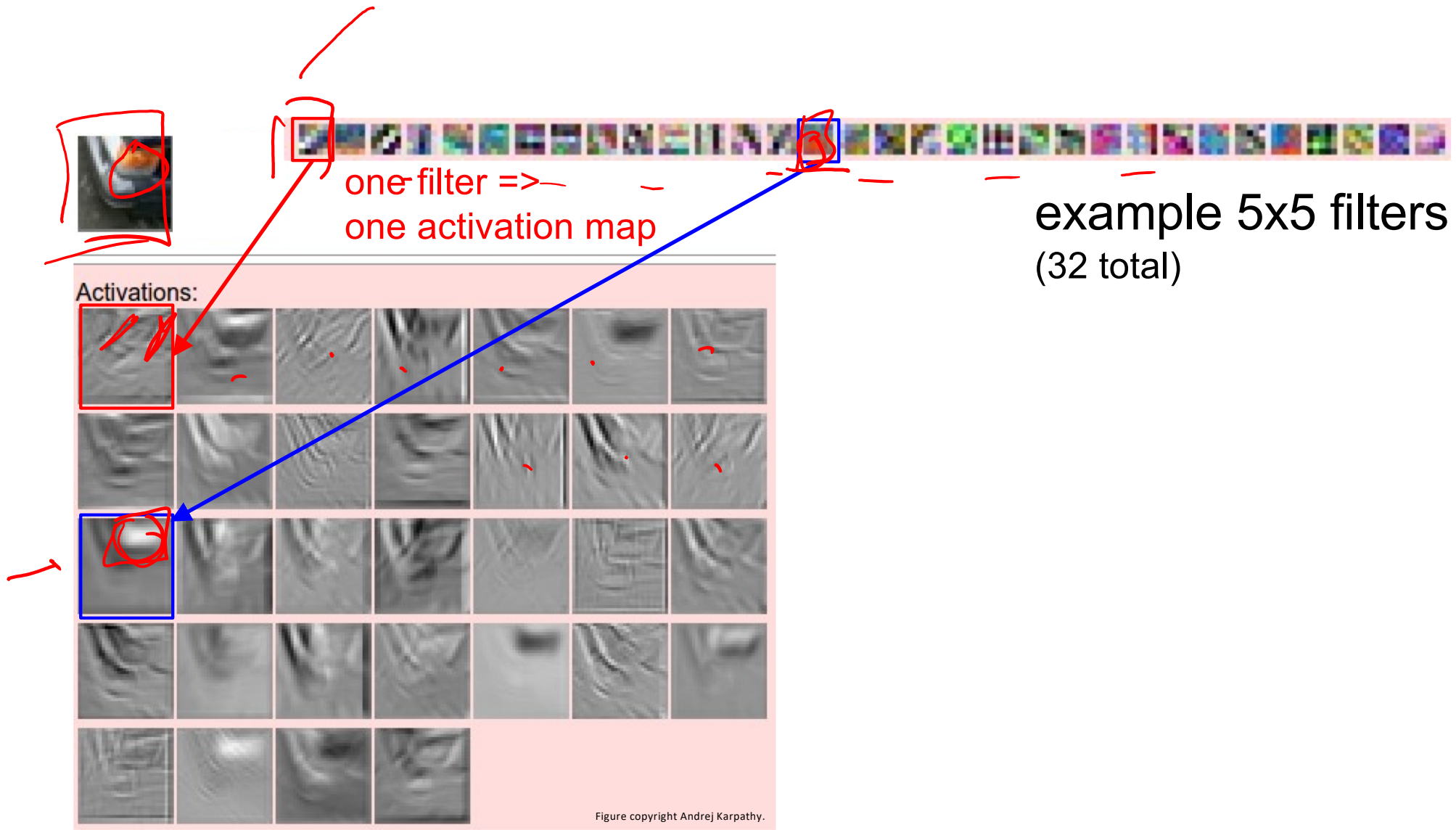
VGG-16 Conv1_1



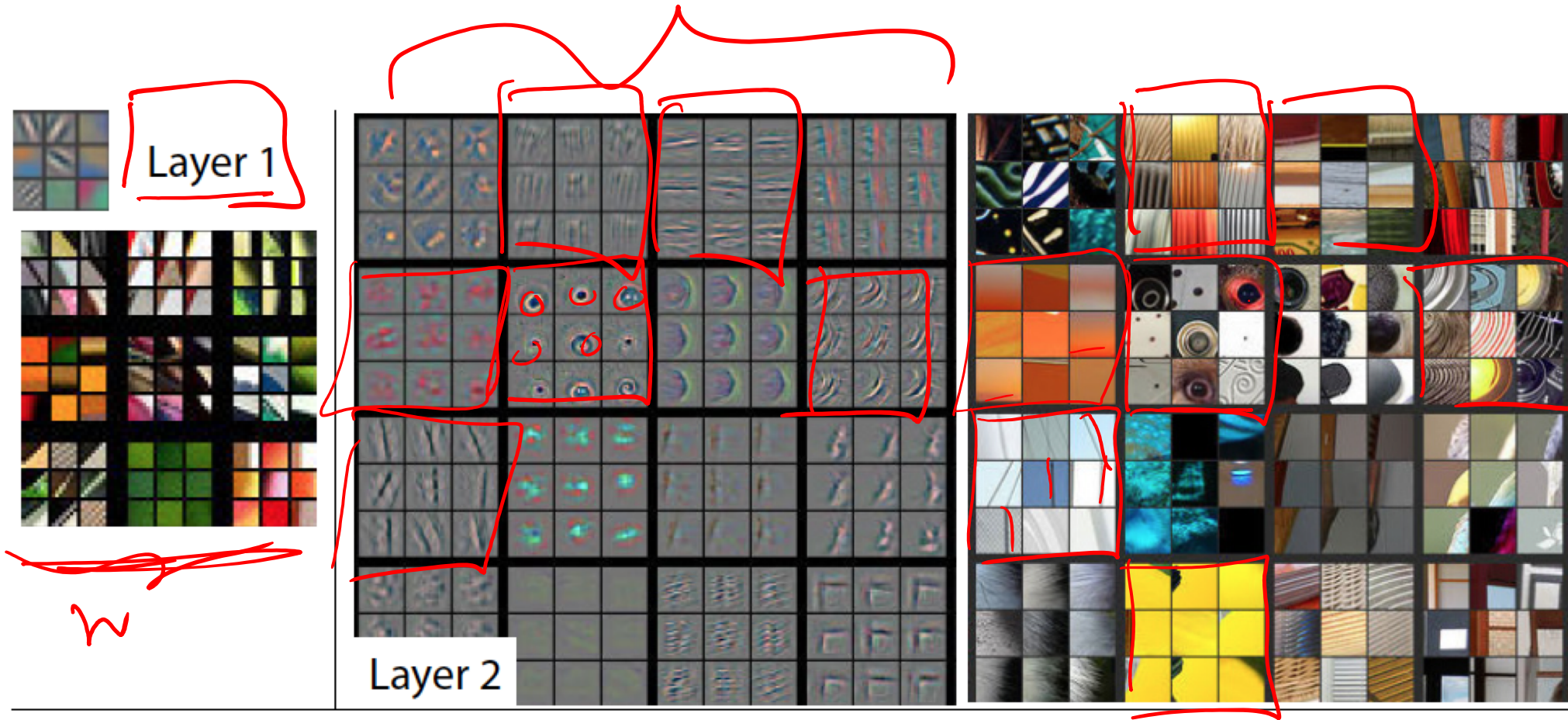
VGG-16 Conv3_2



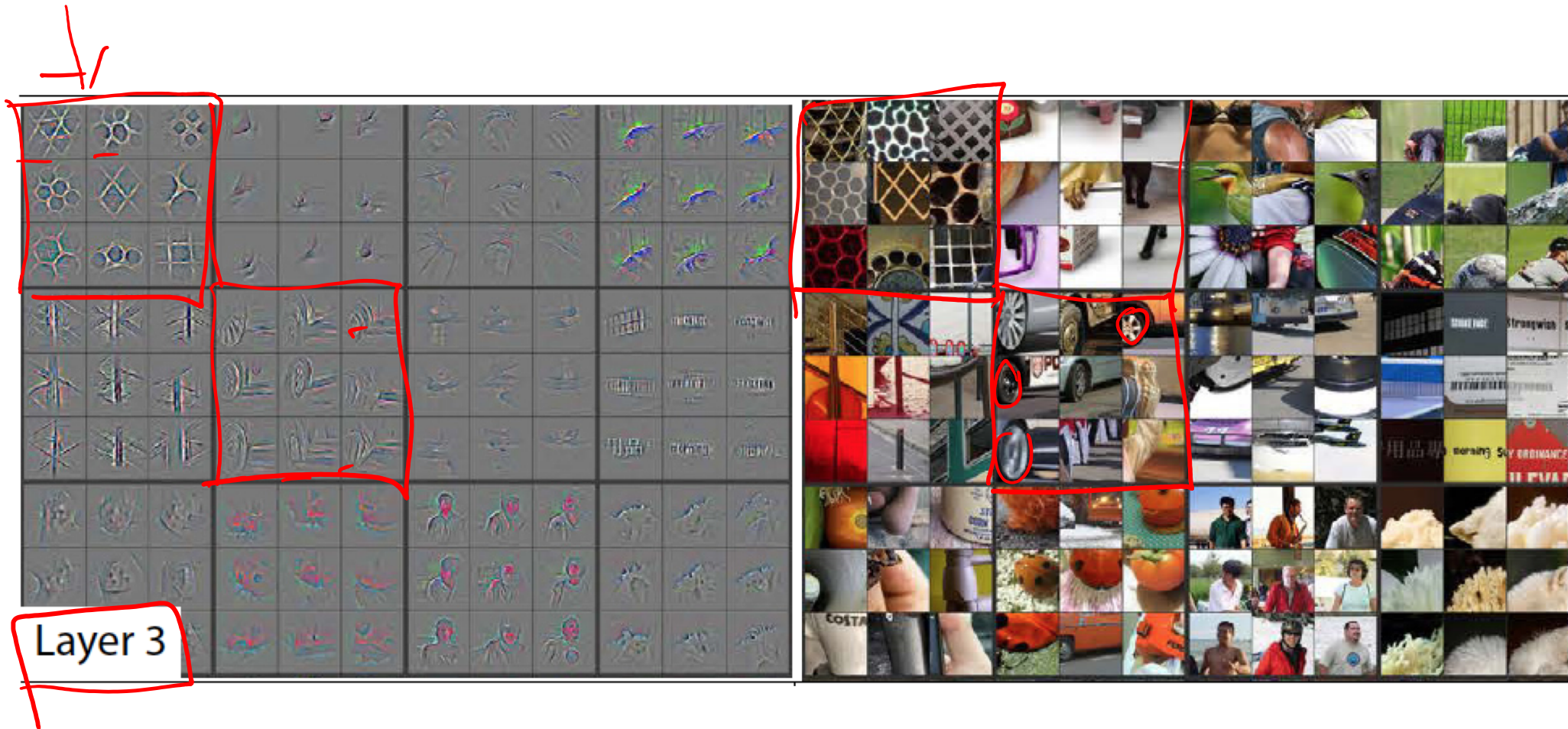
VGG-16 Conv5_3



Visualizing Learned Filters



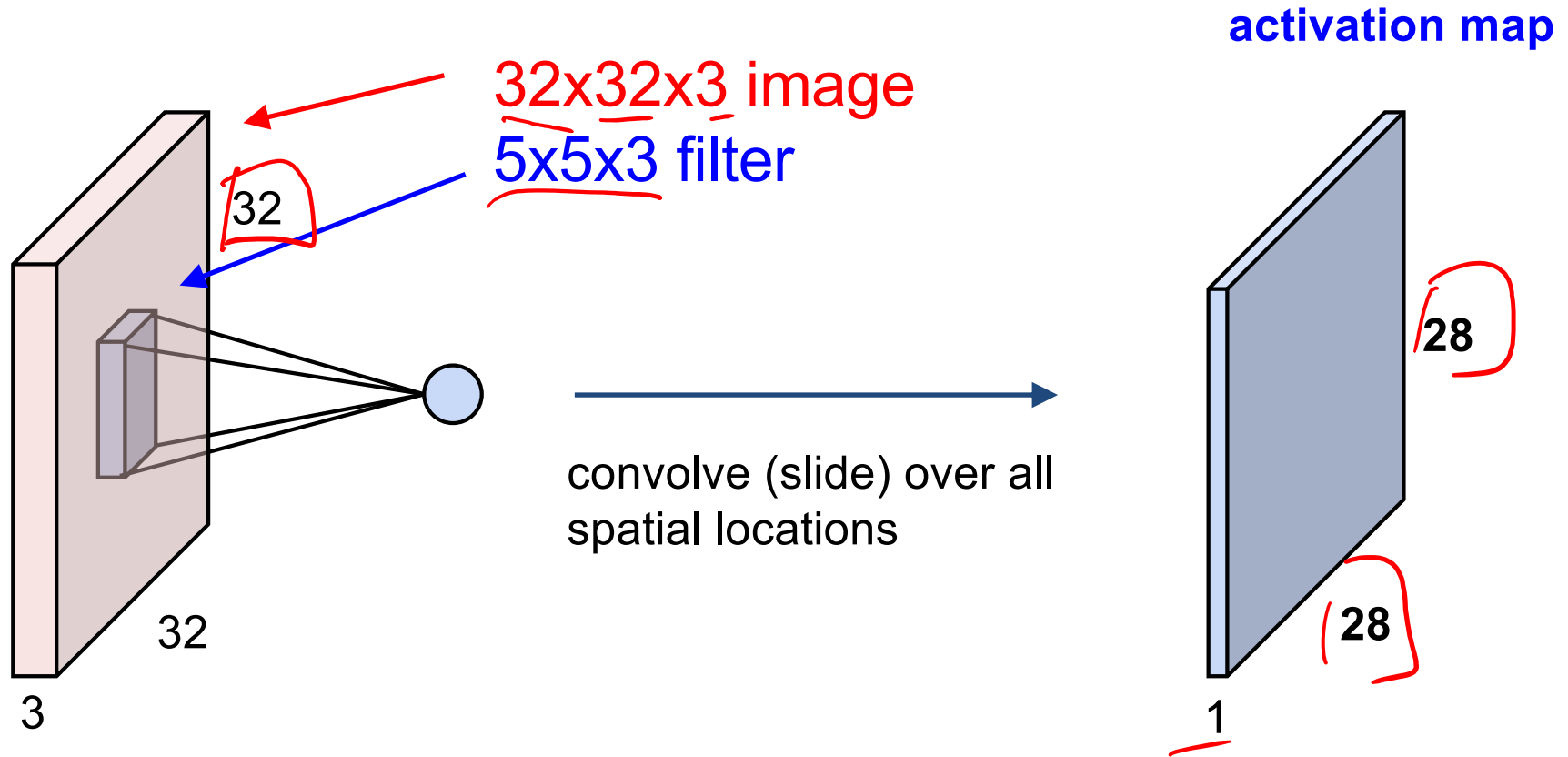
Visualizing Learned Filters



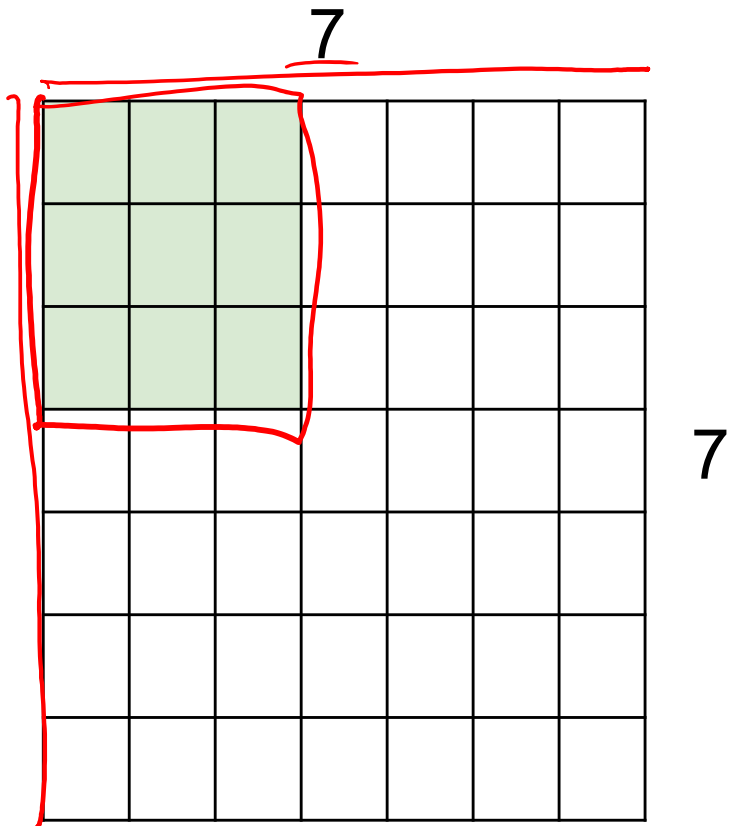
Visualizing Learned Filters $f_{\theta} \rightarrow \phi(x)$



A closer look at spatial dimensions:



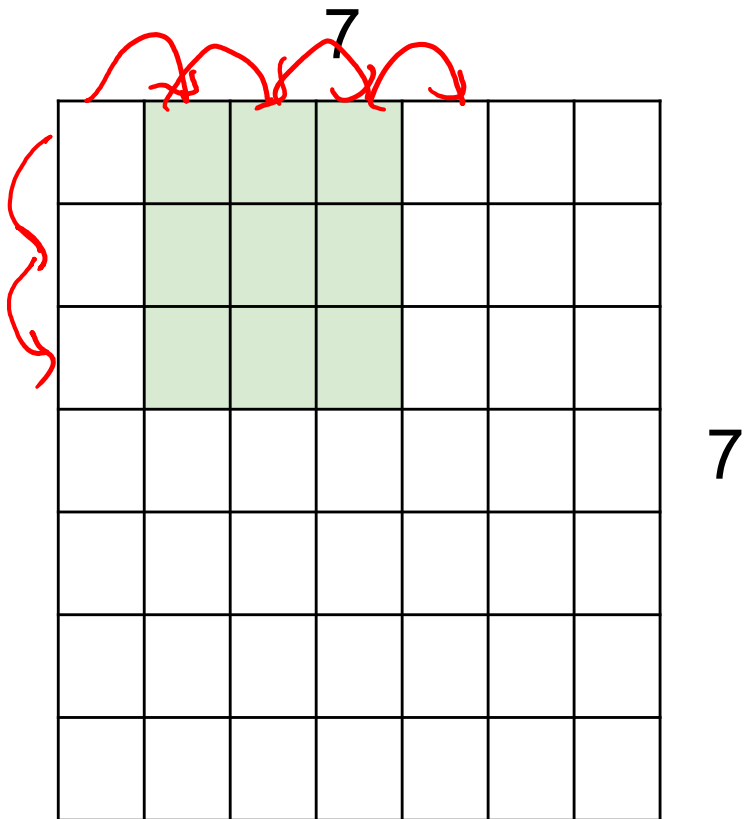
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

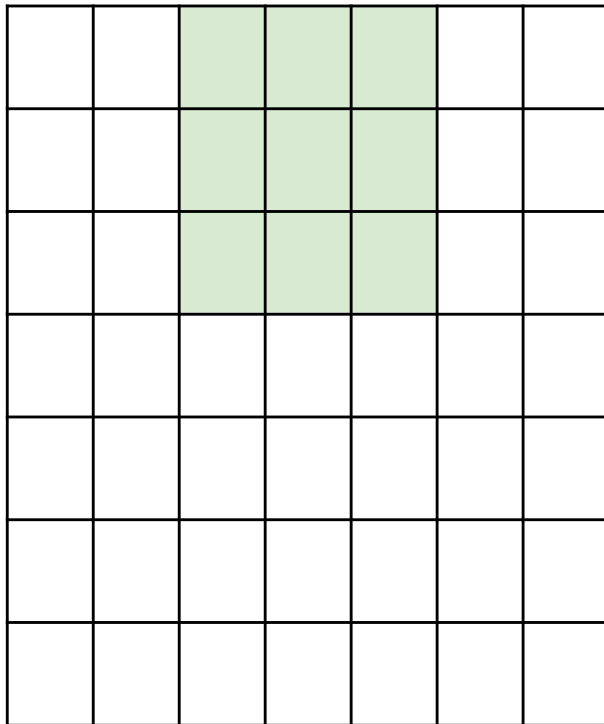
stride = 1



7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:

7

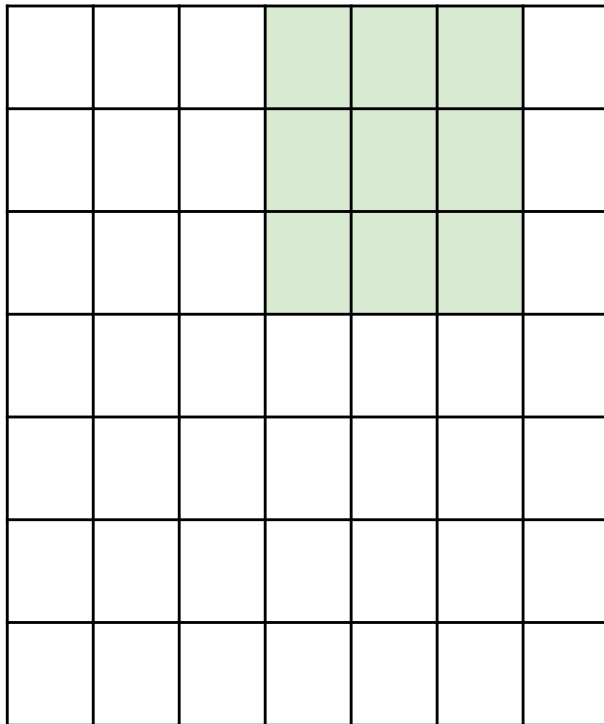


7x7 input (spatially)
assume 3x3 filter

7

A closer look at spatial dimensions:

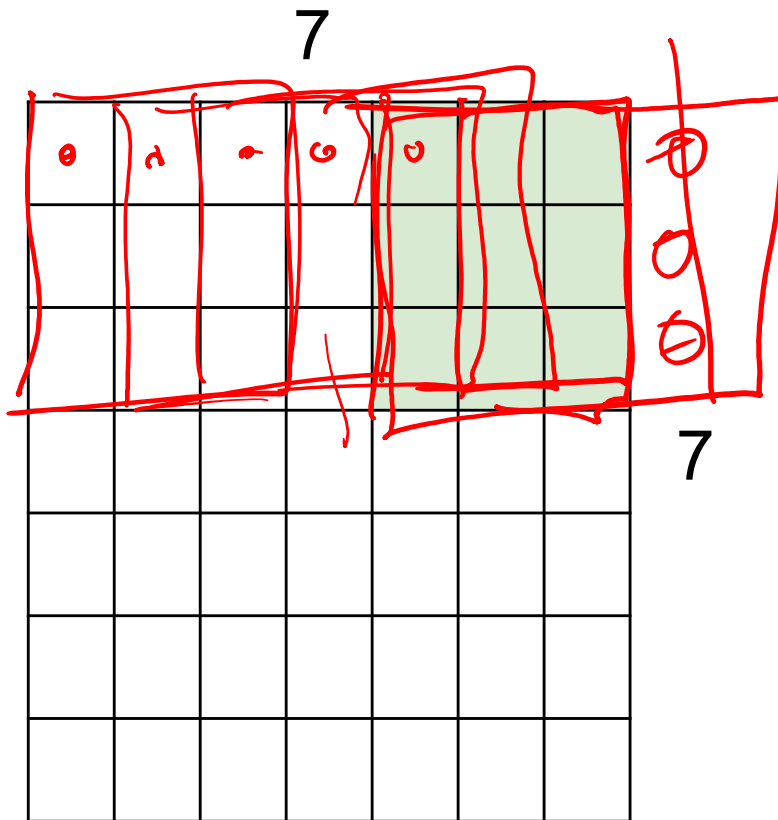
7



7x7 input (spatially)
assume 3x3 filter

7

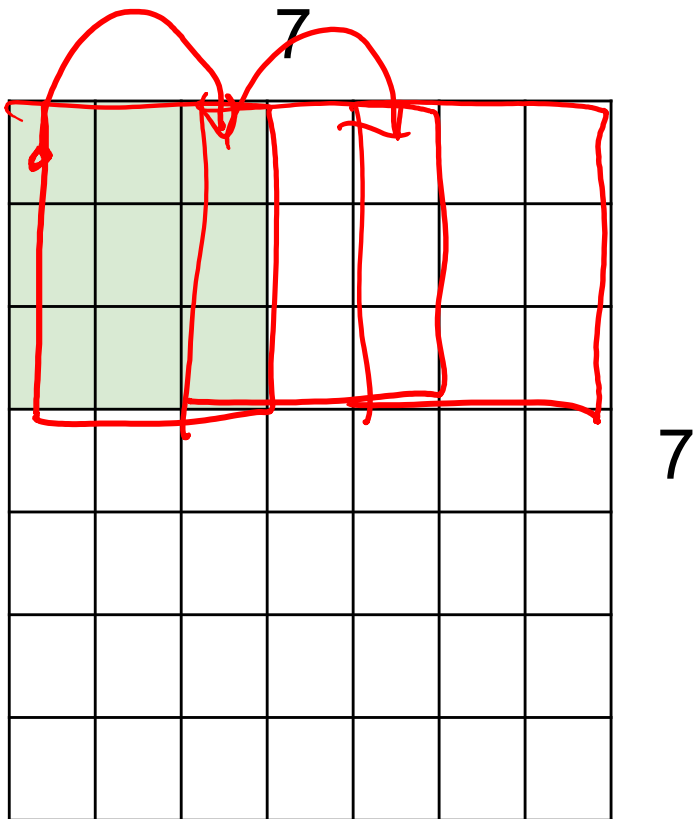
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

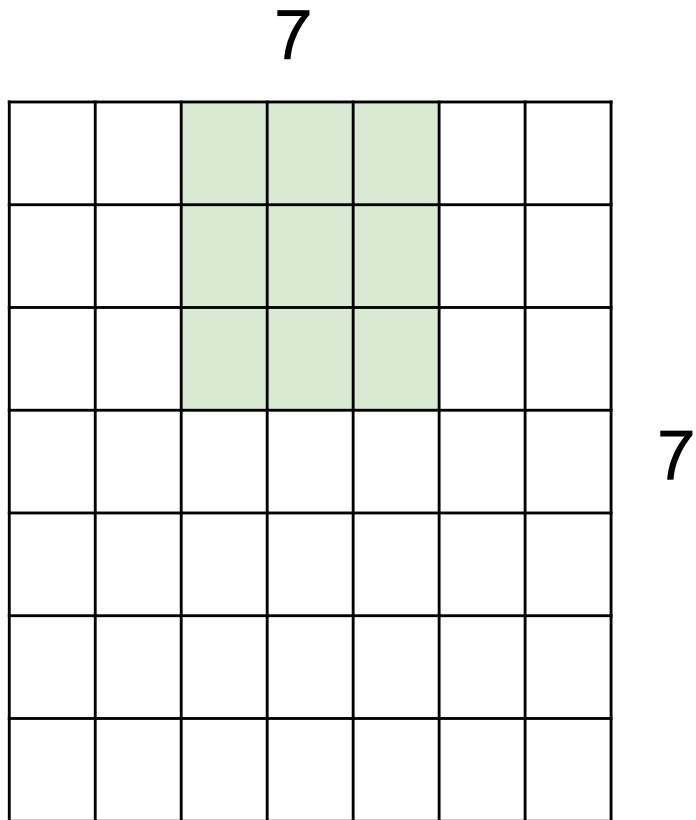
=> **5x5** output

A closer look at spatial dimensions:



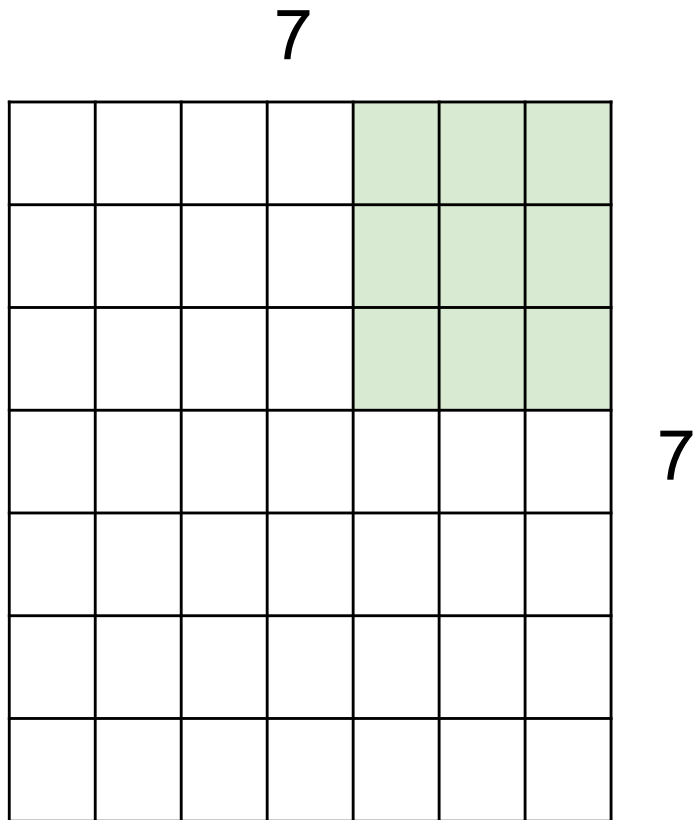
7x7 input (spatially)
assume 3x3 filter
applied with stride 2

A closer look at spatial dimensions:



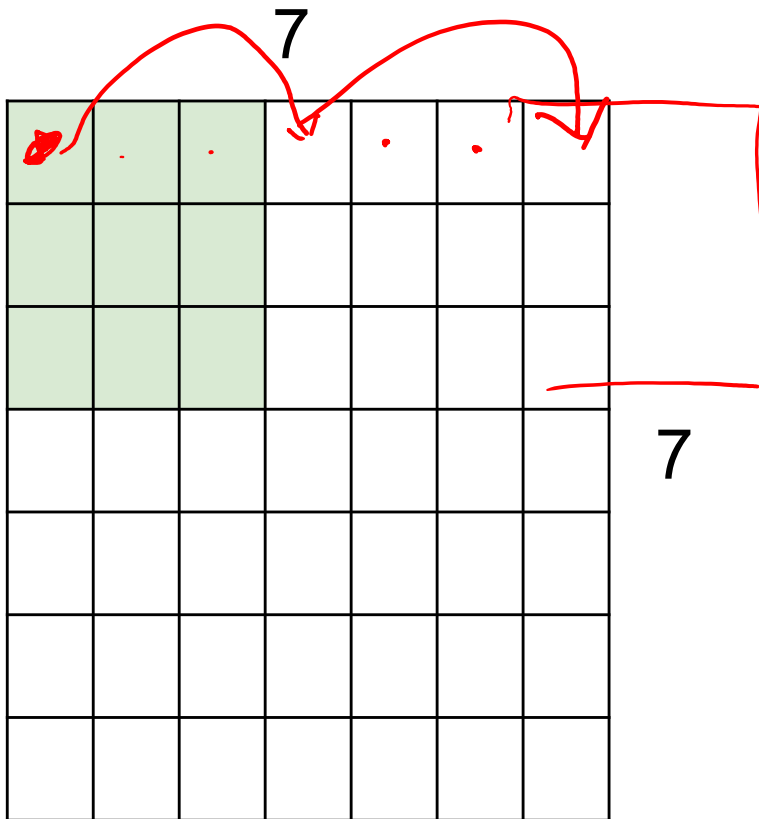
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**

A closer look at spatial dimensions:



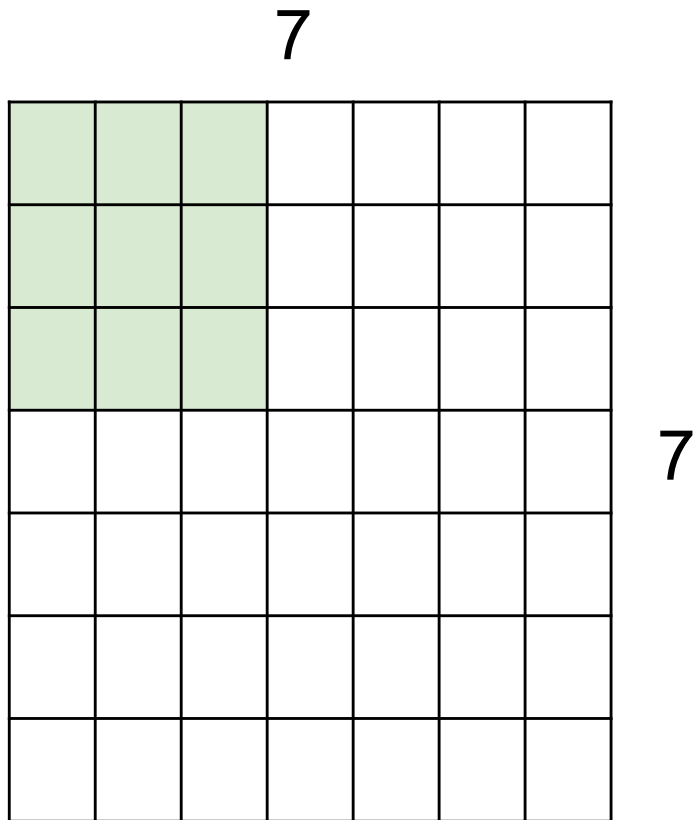
7x7 input (spatially)
assume 3x3 filter
applied **with stride 2**
=> **3x3** output!

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied with stride 3?

A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter
applied **with stride 3?**

doesn't fit!
cannot apply 3x3 filter on
7x7 input with stride 3.

$(N-F) \propto \text{stride}$

$$\frac{N-F}{\text{stride}} + 1$$

Output size:

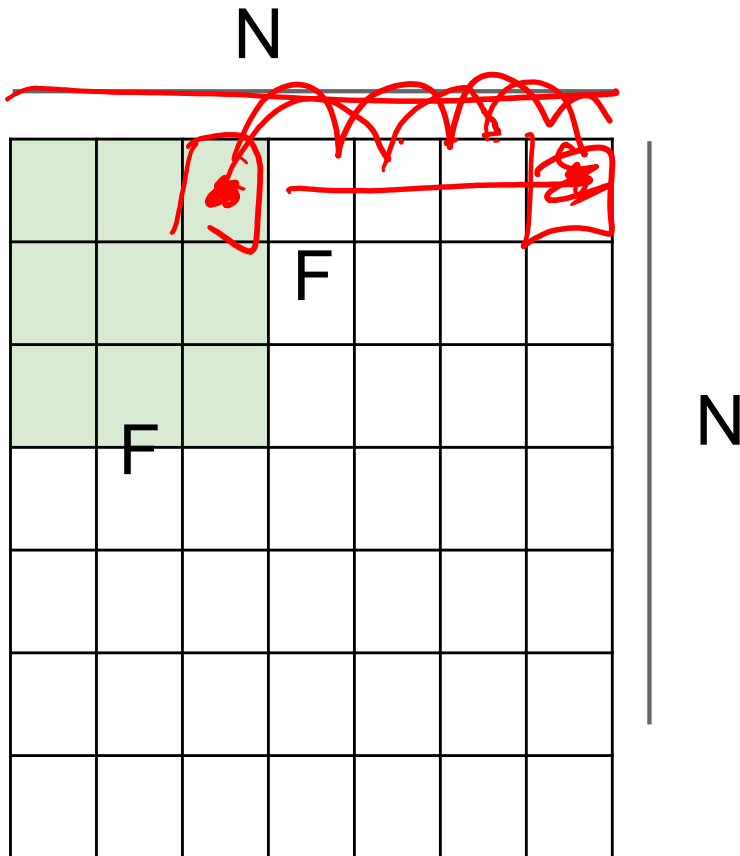
$$\boxed{(N - F) / \text{stride} + 1}$$

e.g. $N = 7$, $F = 3$:

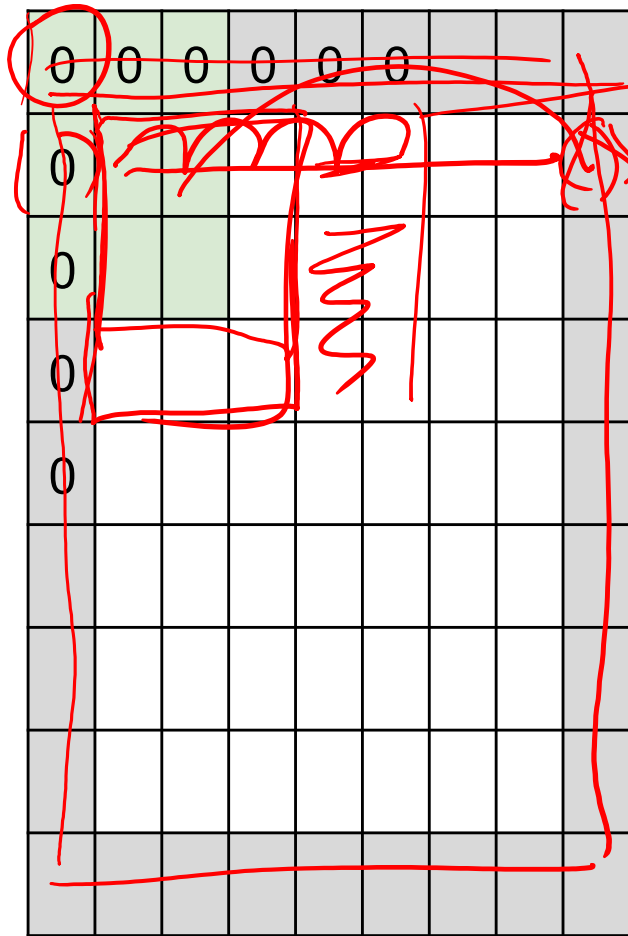
$$\text{stride } 1 \Rightarrow (7 - 3) / 1 + 1 = 5$$

$$\text{stride } 2 \Rightarrow (7 - 3) / 2 + 1 = 3$$

$$\text{stride } 3 \Rightarrow (7 - 3) / 3 + 1 = 2.33 \text{ :}\backslash$$



In practice: Common to zero pad the border



e.g. input 7x7

3x3 filter, applied with stride 1

pad with 1 pixel border => what is the output?

$N \Rightarrow \hat{N}$
 $7 \times 7 \quad 9 \times 9$

$$\frac{9 - 3}{1} + 1$$

$$\boxed{7 \times 7}$$

(recall:)

$$\underline{(N - F) / \text{stride} + 1}$$

22

$$\frac{N - F - 2 \text{ pad} + 1}{\text{Stride}} \left. \vphantom{\frac{N - F - 2 \text{ pad} + 1}{\text{Stride}}} \right\} \boxed{\frac{F-1}{2}}$$

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

7x7 output!

In practice: Common to zero pad the border

0	0	0	0	0	0			
0								
0								
0								
0								

e.g. input 7x7

3x3 filter, applied with **stride 1**

pad with 1 pixel border => what is the output?

7x7 output!

in general, common to see CONV layers with stride 1, filters of size $F \times F$, and zero-padding with $(F-1)/2$. (will preserve size spatially)

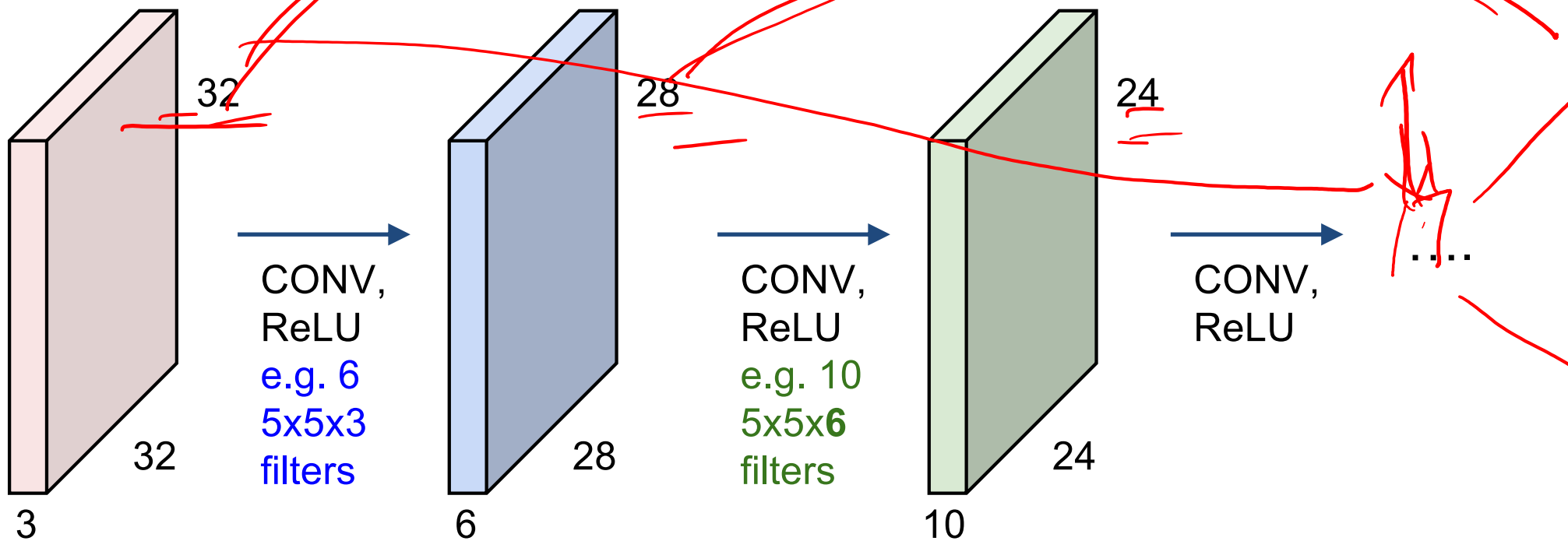
e.g. $F = 3$ => zero pad with 1

$F = 5$ => zero pad with 2

$F = 7$ => zero pad with 3

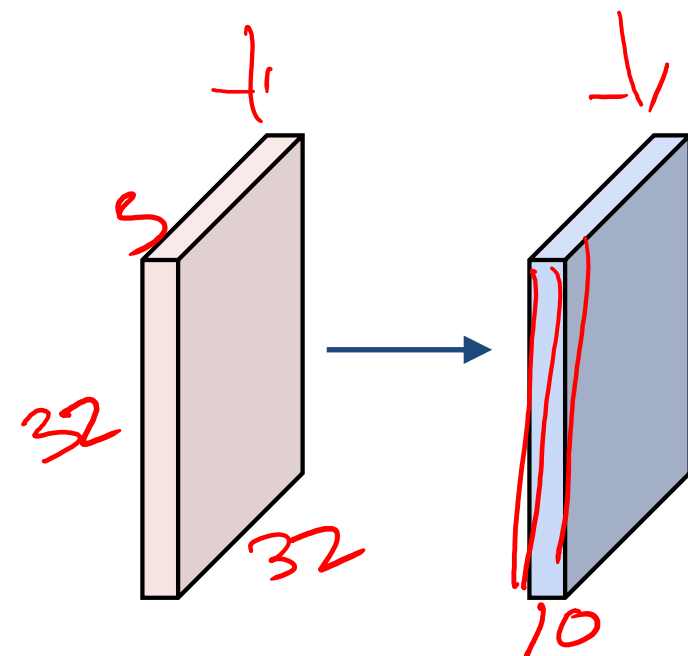
Remember back to...

E.g. 32x32 input convolved repeatedly with 5x5 filters shrinks volumes spatially! (32 -> 28 -> 24 ...). Shrinking too fast is not good, doesn't work well.



Examples time:

Input volume: $32 \times 32 \times 3$
 10 5×5 filters with stride 1, pad 2



Output volume size: ?

$$\frac{(32)}{N_1} \times \frac{(32)}{N_2} \times \frac{(10)}{C_2}$$

Examples time:

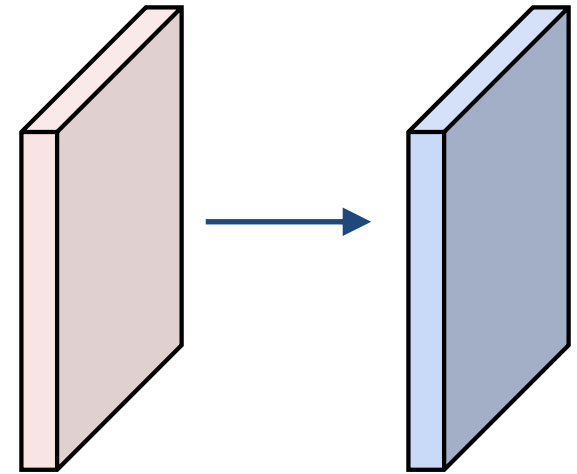
Input volume: **32x32x3**

10 **5x5** filters with stride **1**, pad **2**

Output volume size:

$(32+2*2-5)/1+1 = 32$ spatially, so

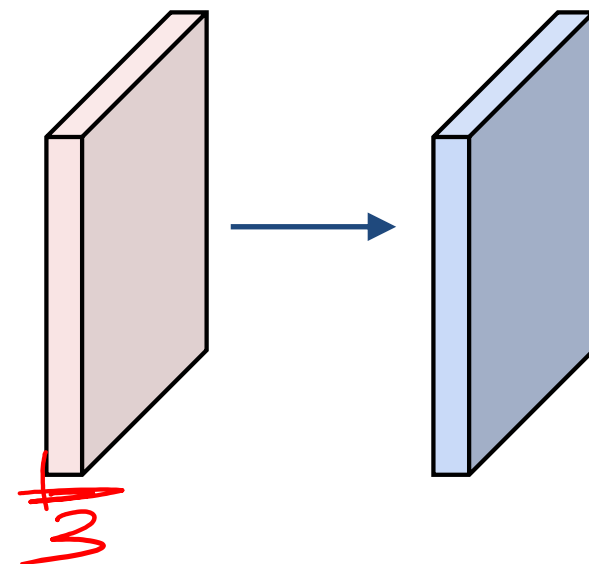
32x32x10



Examples time:

Input volume: $32 \times 32 \times 3$

10, 5×5 filters with stride 1, pad 2



Number of parameters in this layer?

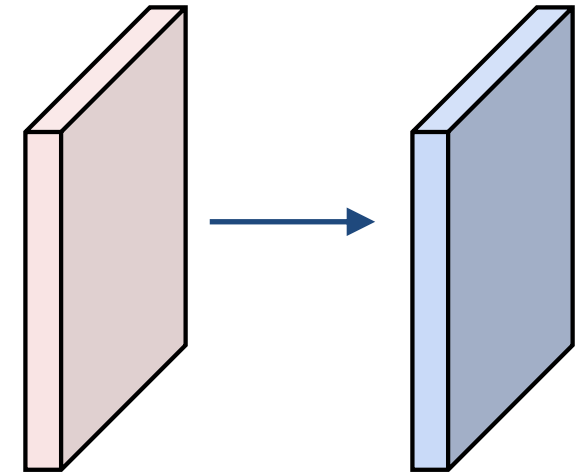
$$(5 \times 5 \times 3 + 1) \times 10 = 760$$

~~() x () x () x ()~~

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

each filter has $5*5*3 + 1 = 76$ params (+1 for bias)

$\Rightarrow 76*10 = 760$

Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Summary. To summarize, the Conv Layer:

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires four hyperparameters:
 - Number of filters K ,
 - their spatial extent F ,
 - the stride S ,
 - the amount of zero padding P .
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F + 2P)/S + 1$
 - $H_2 = (H_1 - F + 2P)/S + 1$ (i.e. width and height are computed equally by symmetry)
 - $D_2 = K$
- With parameter sharing, it introduces $F \cdot F \cdot D_1$ weights per filter, for a total of $(F \cdot F \cdot D_1) \cdot K$ weights and K biases.
- In the output volume, the d -th depth slice (of size $W_2 \times H_2$) is the result of performing a valid convolution of the d -th filter over the input volume with a stride of S , and then offset by d -th bias.

Common settings:

$K =$ (powers of 2, e.g. 32, 64, 128, 512)

- $F = 3, S = 1, P = 1$
- $F = 5, S = 1, P = 2$
- $F = 5, S = 2, P = ?$ (whatever fits)
- $F = 1, S = 1, P = 0$