




CS 4803 / 7643: Deep Learning

Topics:

- ~~Convolutional Neural Networks~~
 - Pooling layers
 - Fully-connected layers as convolutions
 - Toeplitz matrices and convolutions = matrix-mult
 - Backprop in conv layers

Dhruv Batra
Georgia Tech

Administrativa

- HW1 Reminder
 - Due: 10/02, 11:55pm
- Project Idea: ICLR19 Reproducibility Challenge
 - https://reproducibility-challenge.github.io/iclr_2019/
 - https://docs.google.com/spreadsheets/d/1BipWLvvWb7Fu6OSDd-uOCF1Lr_4drKOCRvdhxm_eSHc/edit#gid=0



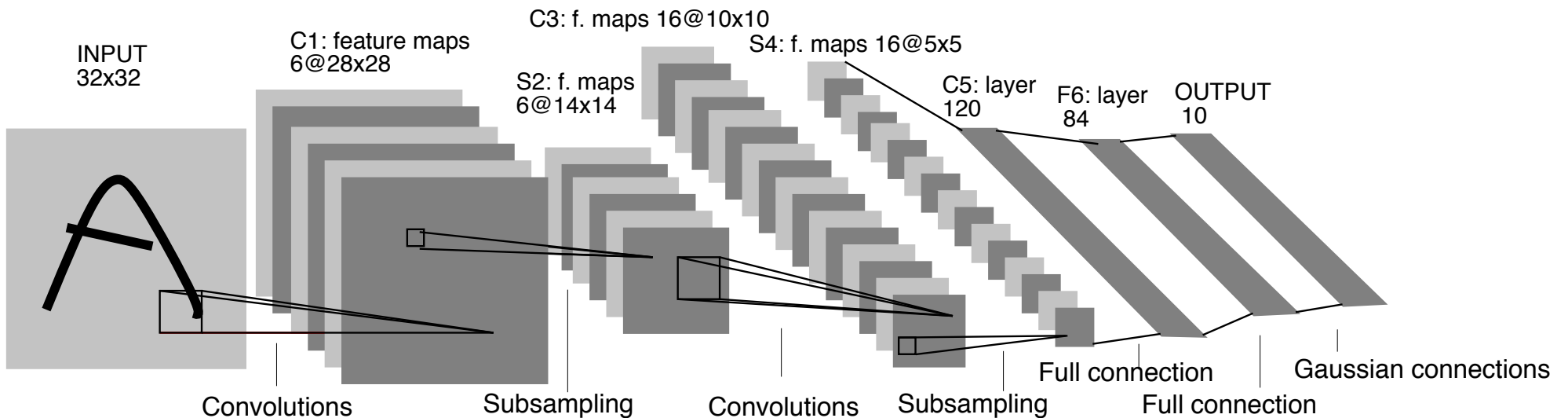
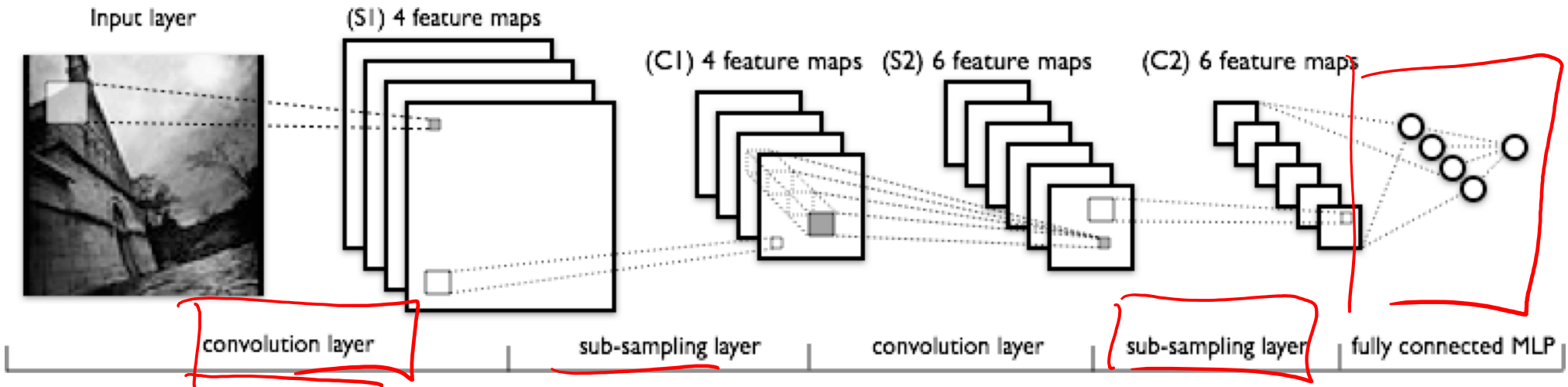
Recap from last time



Convolutional Neural Networks

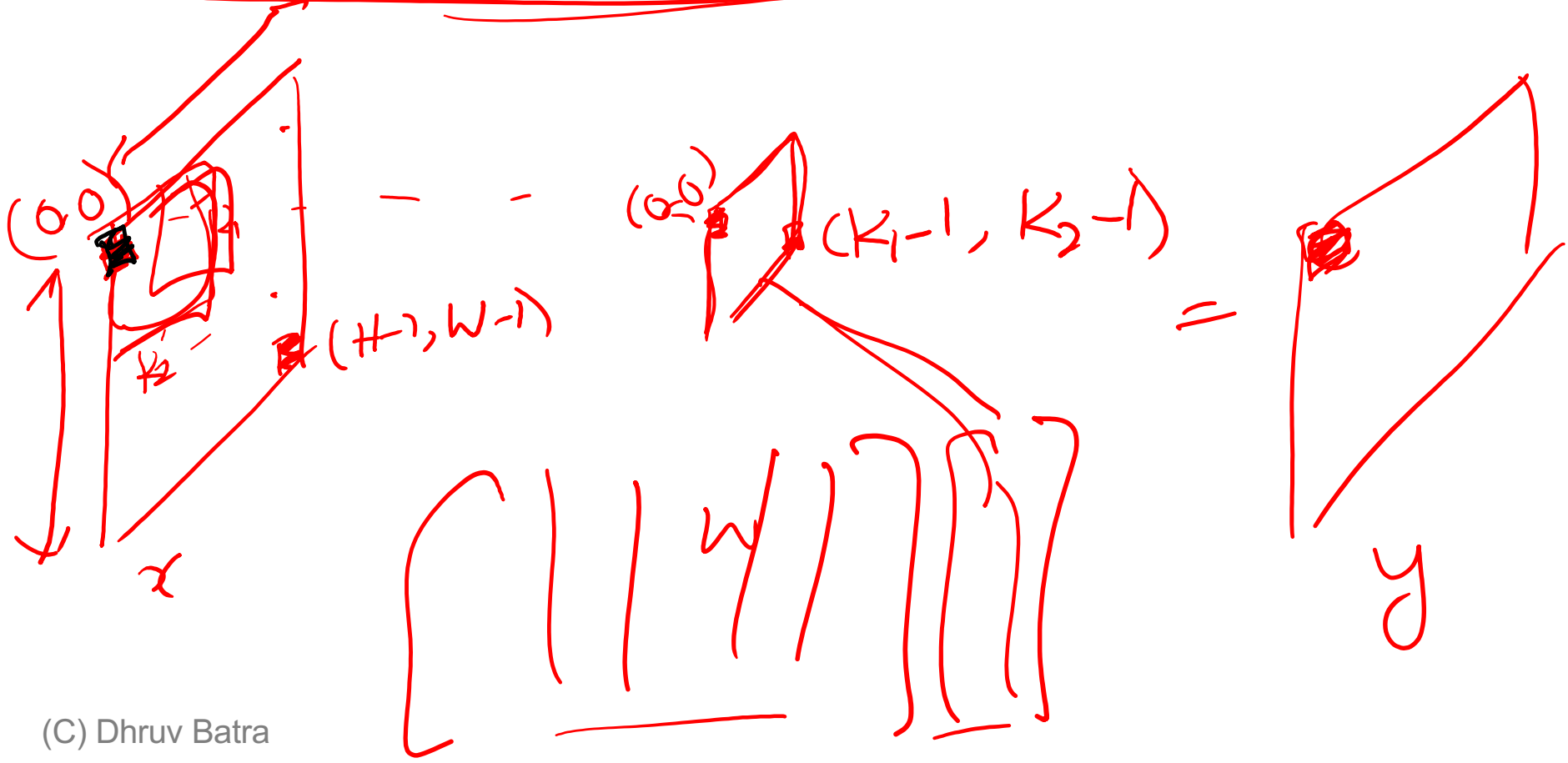
(without the brain stuff)

Convolutional Neural Networks

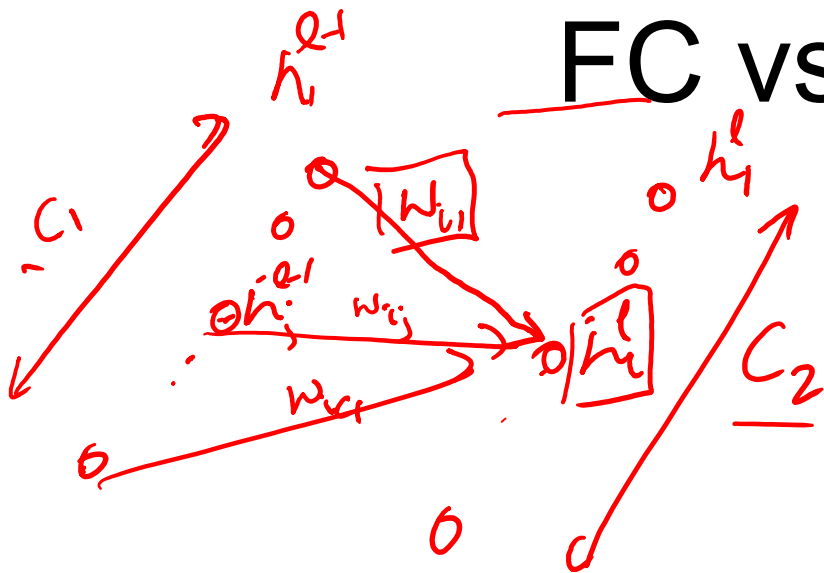


Convolutions for programmers

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

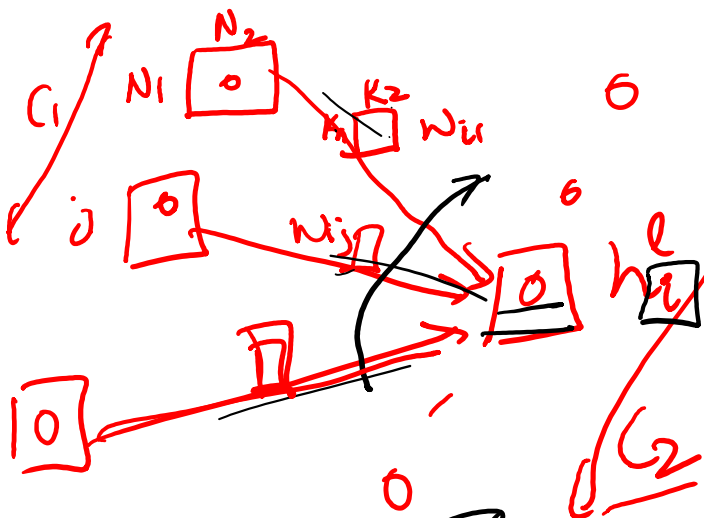


FC vs Conv Layer



$$h_i^l = \sum_{d=1}^{C_1} \underbrace{h_j^{l-1} \cdot w_{ij}}_{\text{Scalar product}} + b_i$$

$$h_i^l = \sum_{j=1}^{C_1} \underbrace{h_j^{l-1} \cdot w_{ij}}_{\text{conv}} + b_i$$

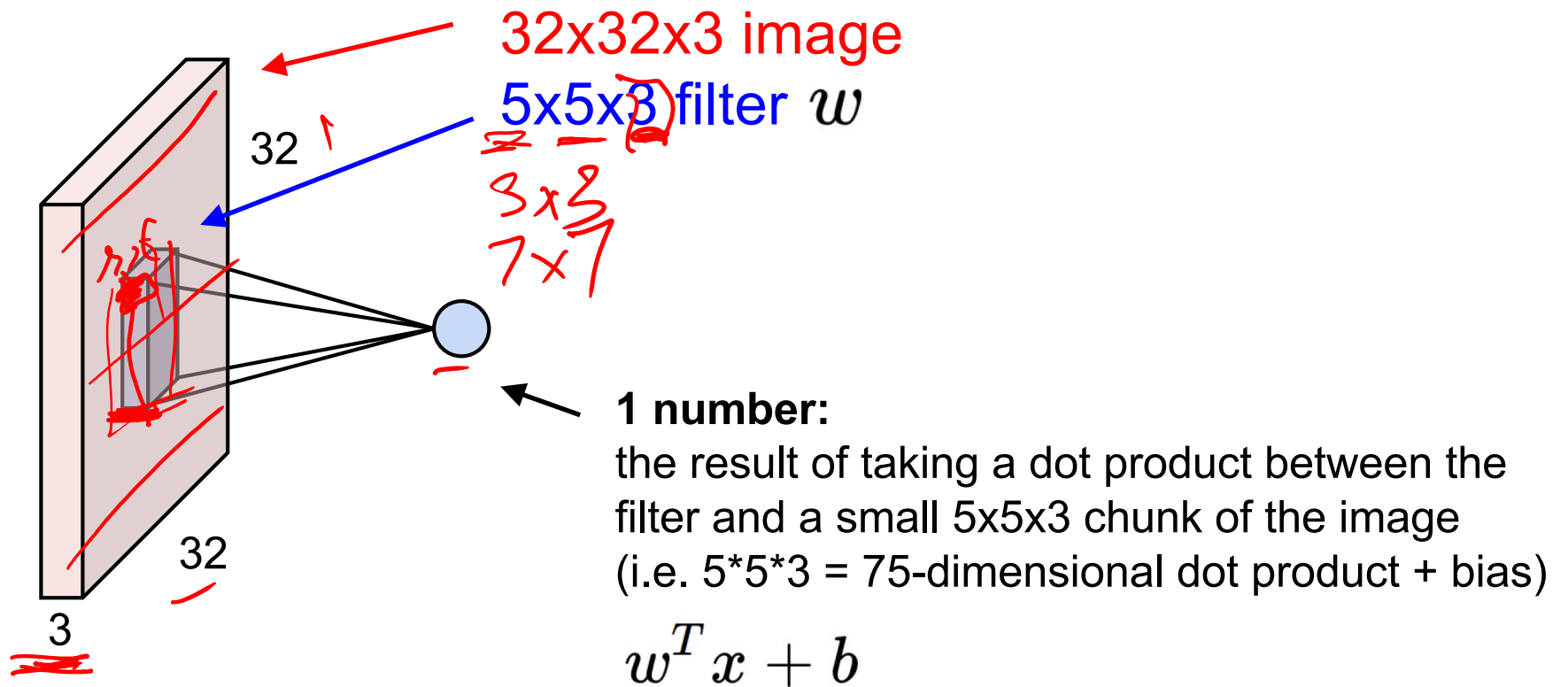


$$h_i^l [r, c] = \sum_{j=1}^{C_1} \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} h_j^{l-1} [r+a, c+b] \cdot w[a, b]$$

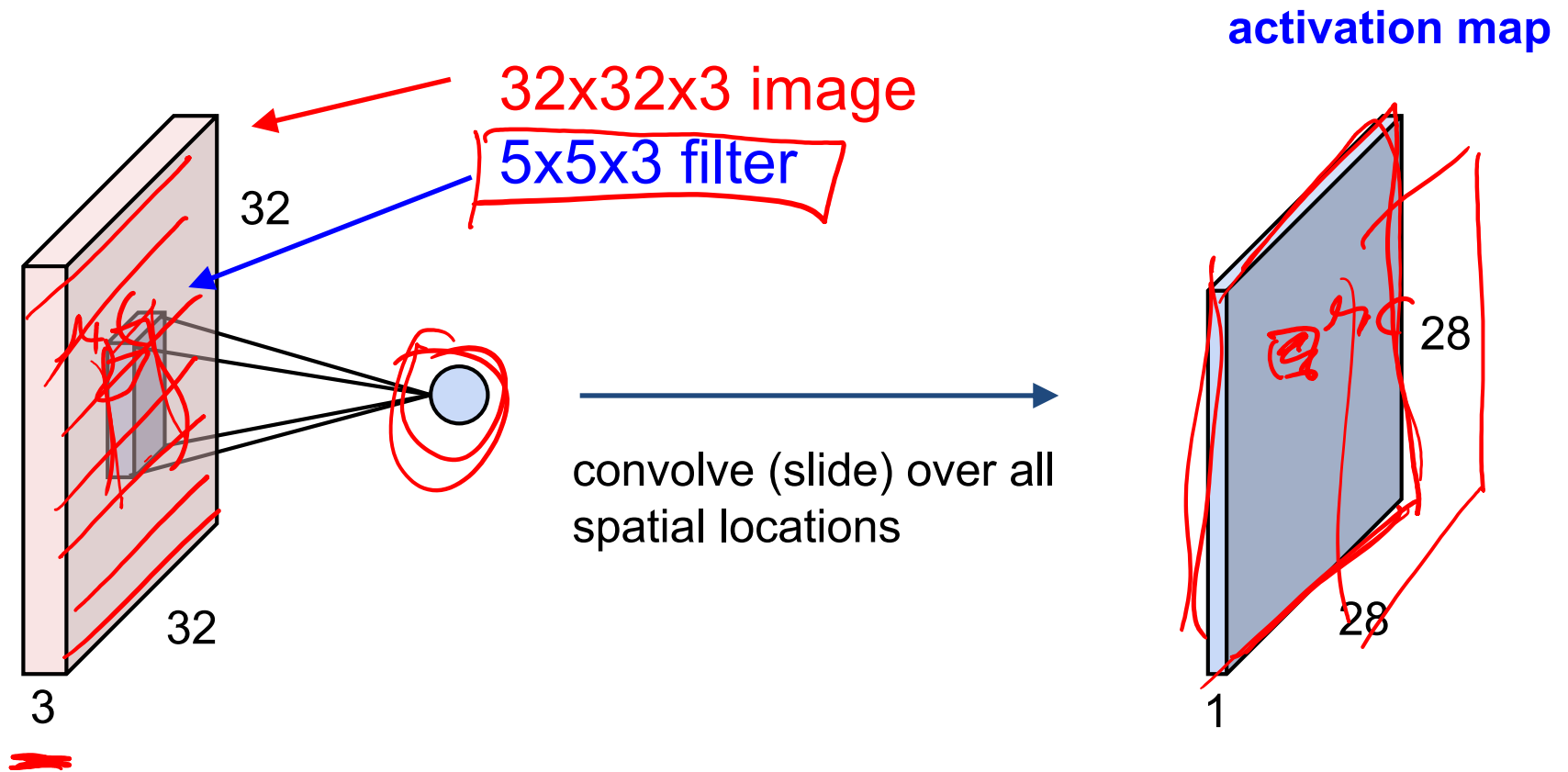


$$W \quad (K_1 \times K_2 \times C_1 \times C_2)$$

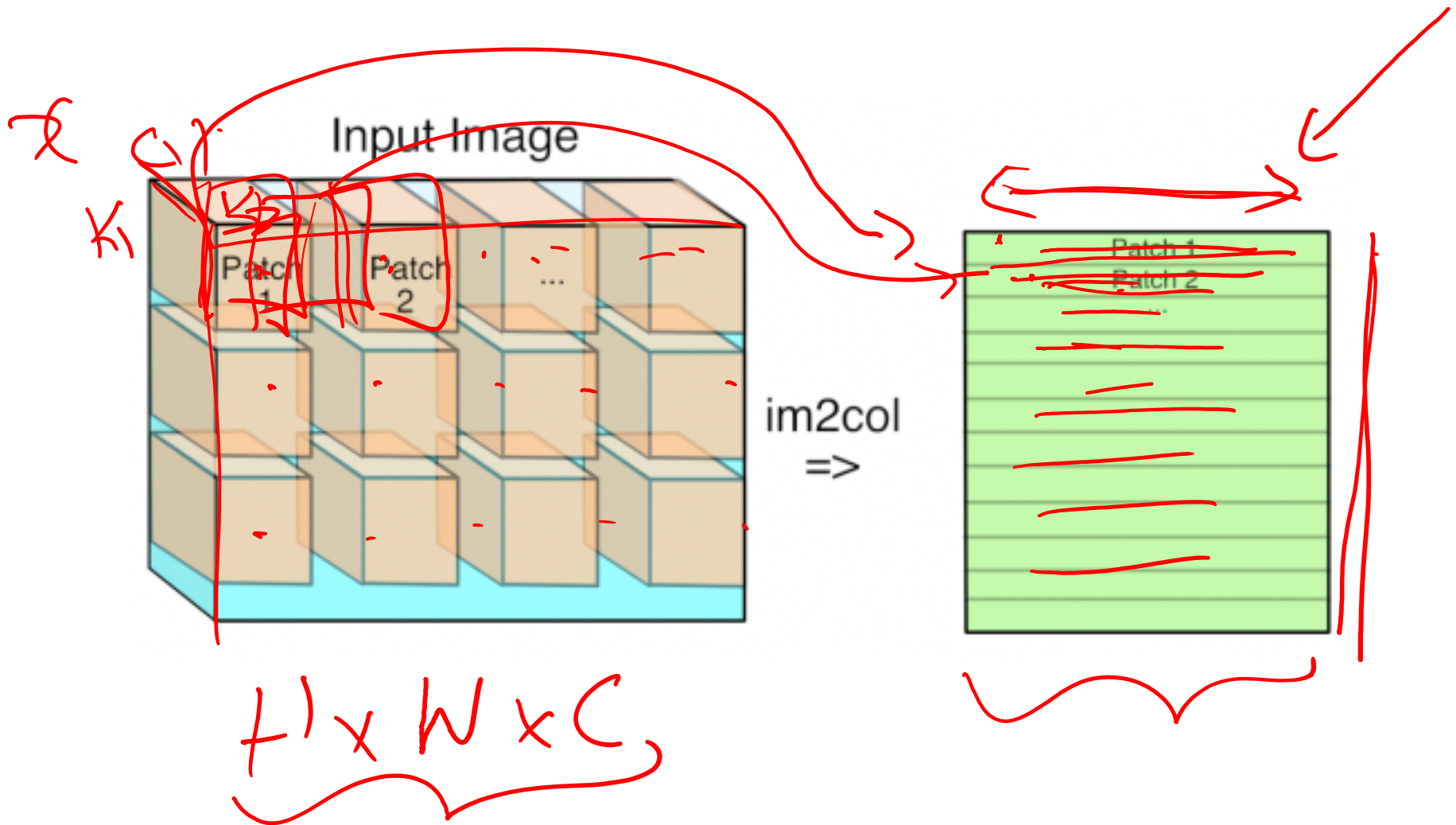
Convolution Layer



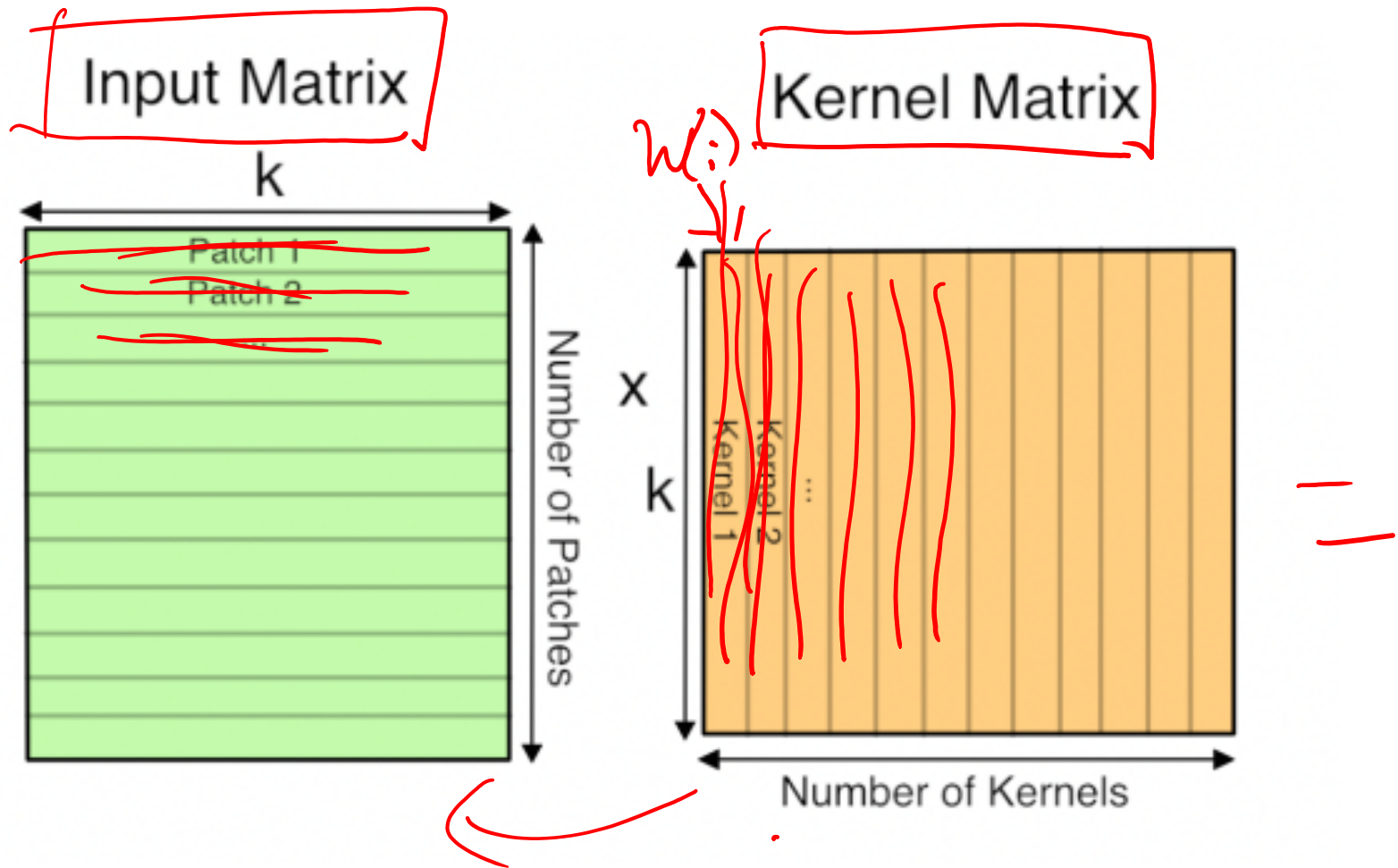
Convolution Layer



Im2Col

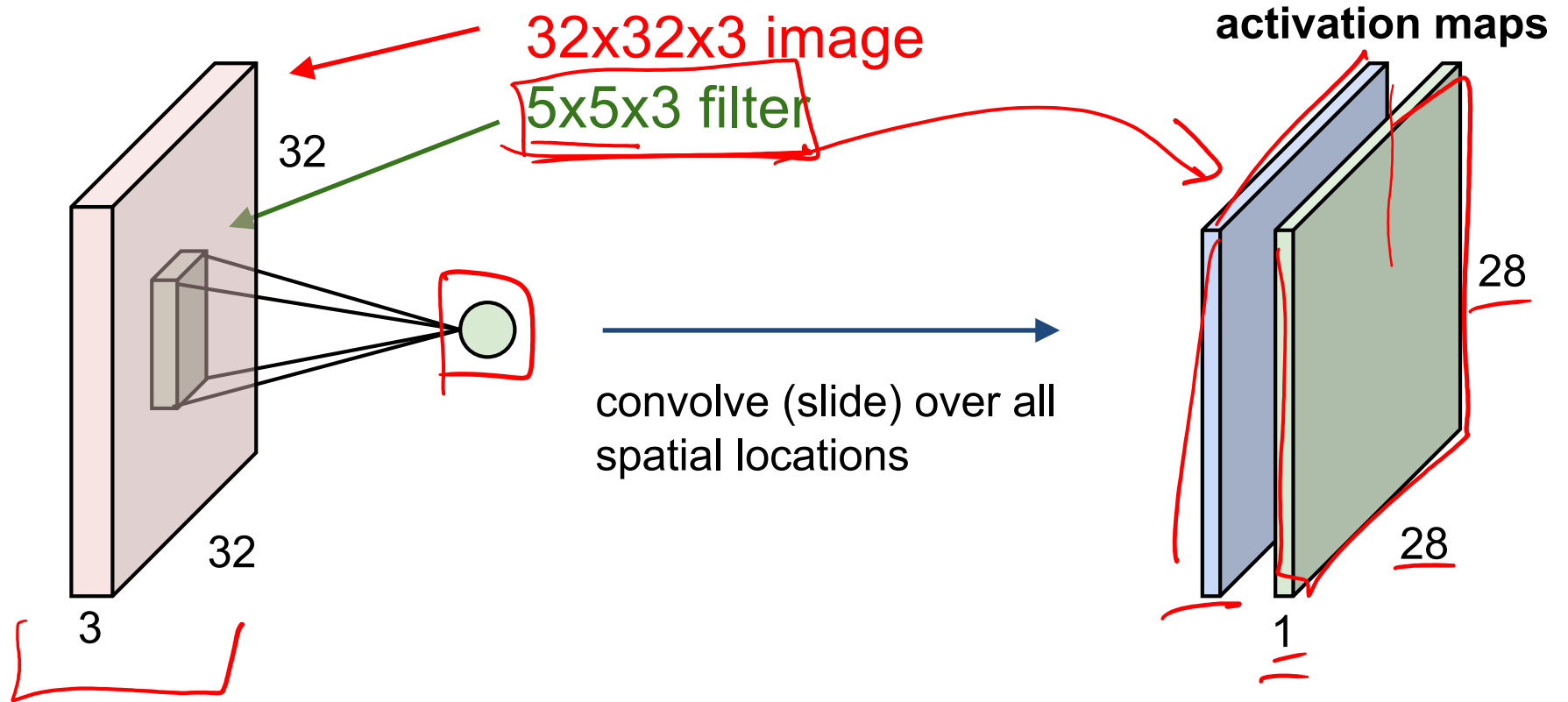


GEMM

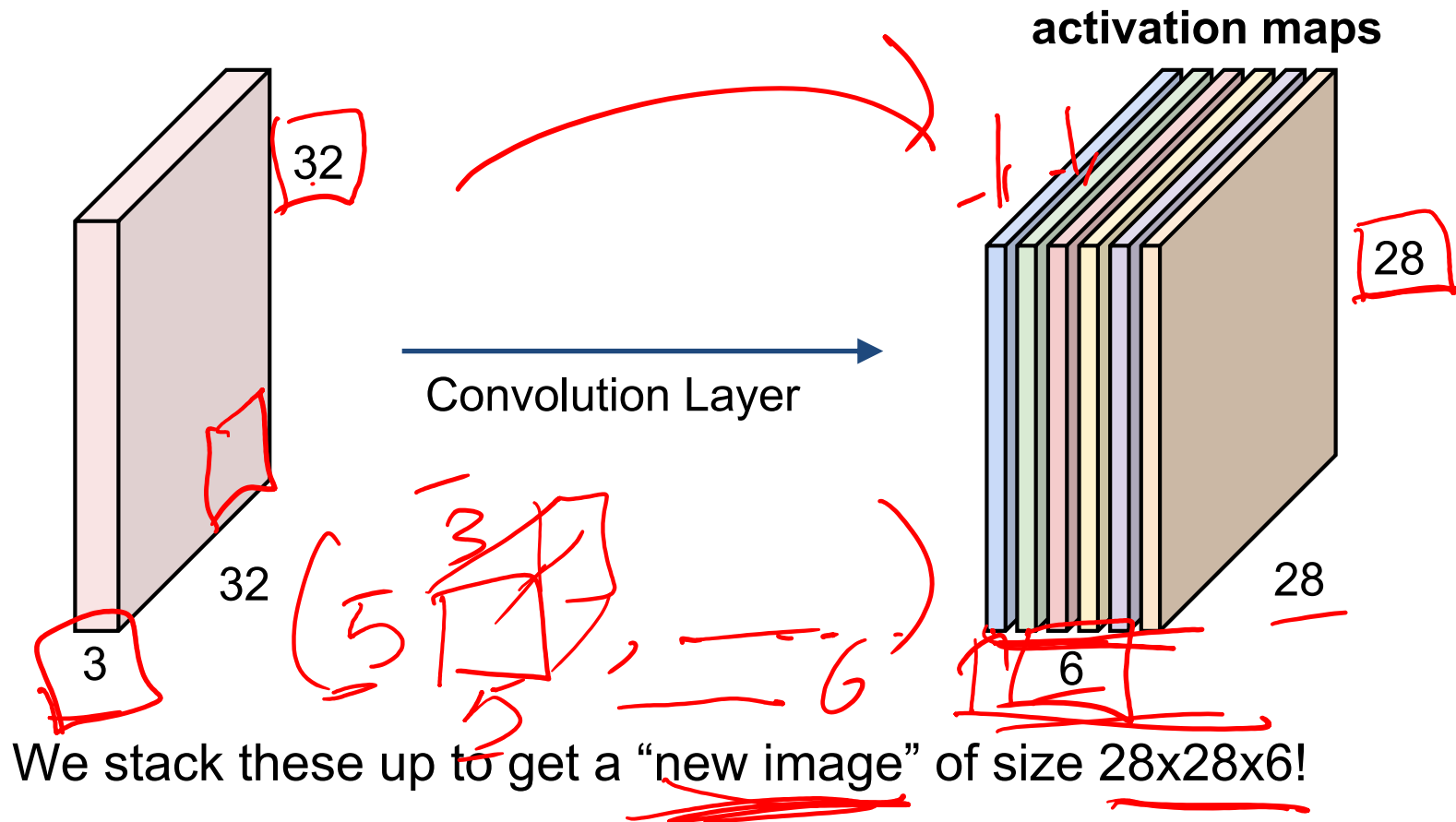


Convolution Layer

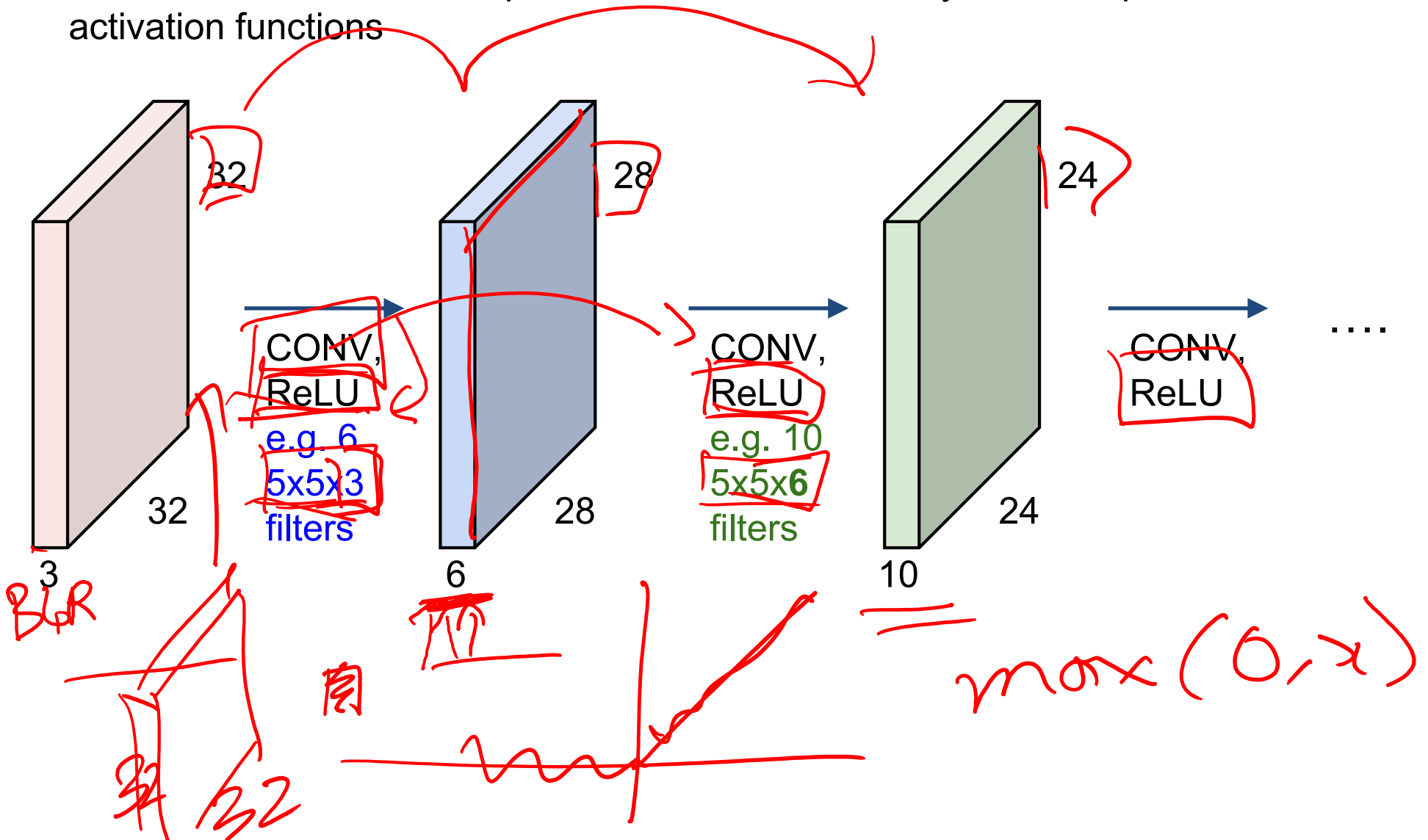
consider a second, **green** filter

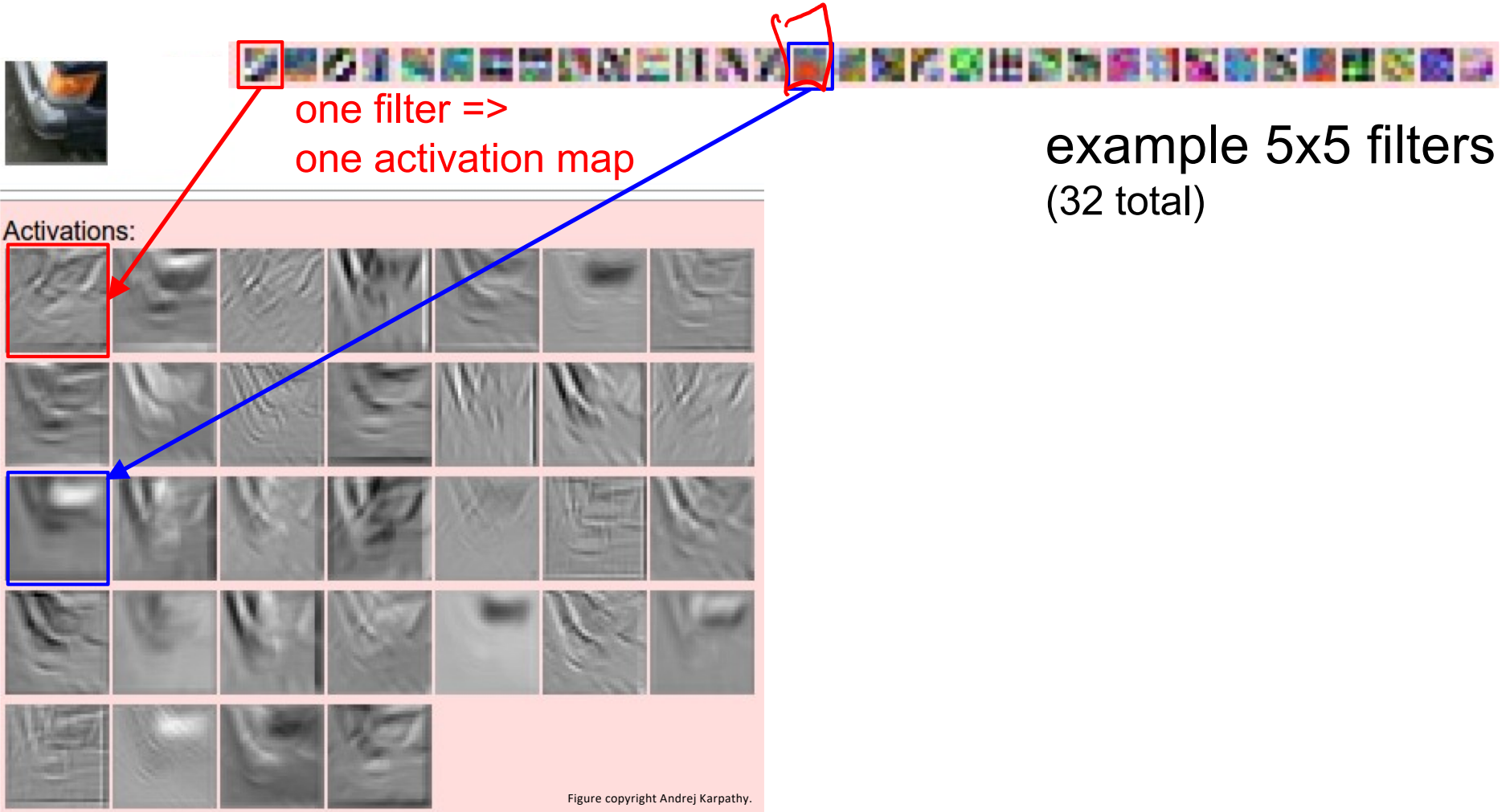


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

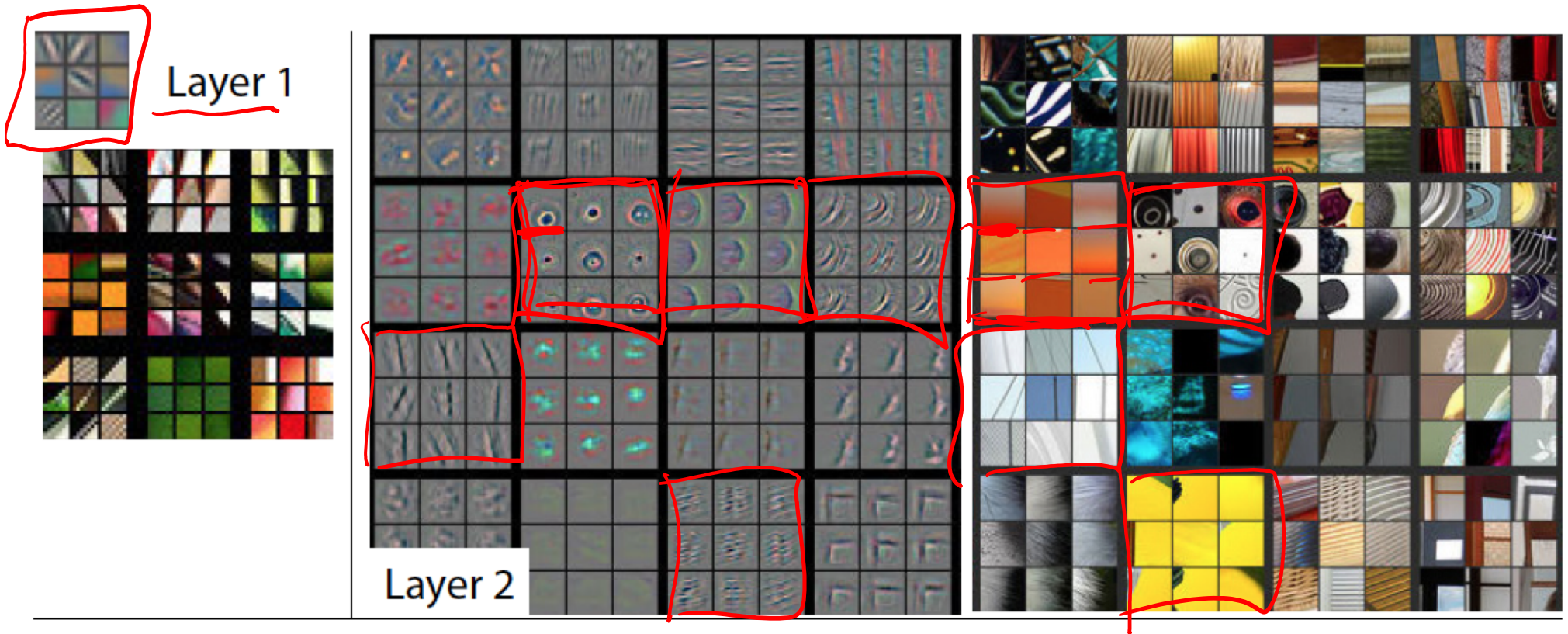


Preview: ConvNet is a sequence of Convolutional Layers, interspersed with activation functions

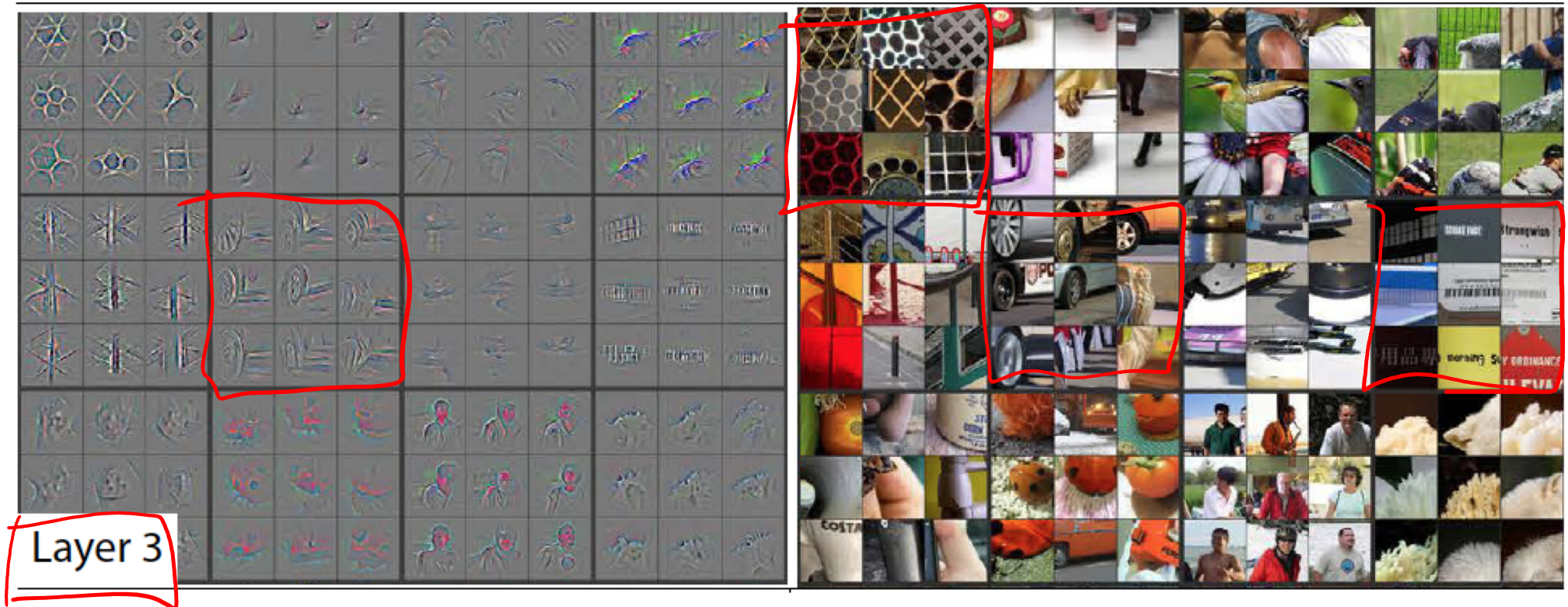




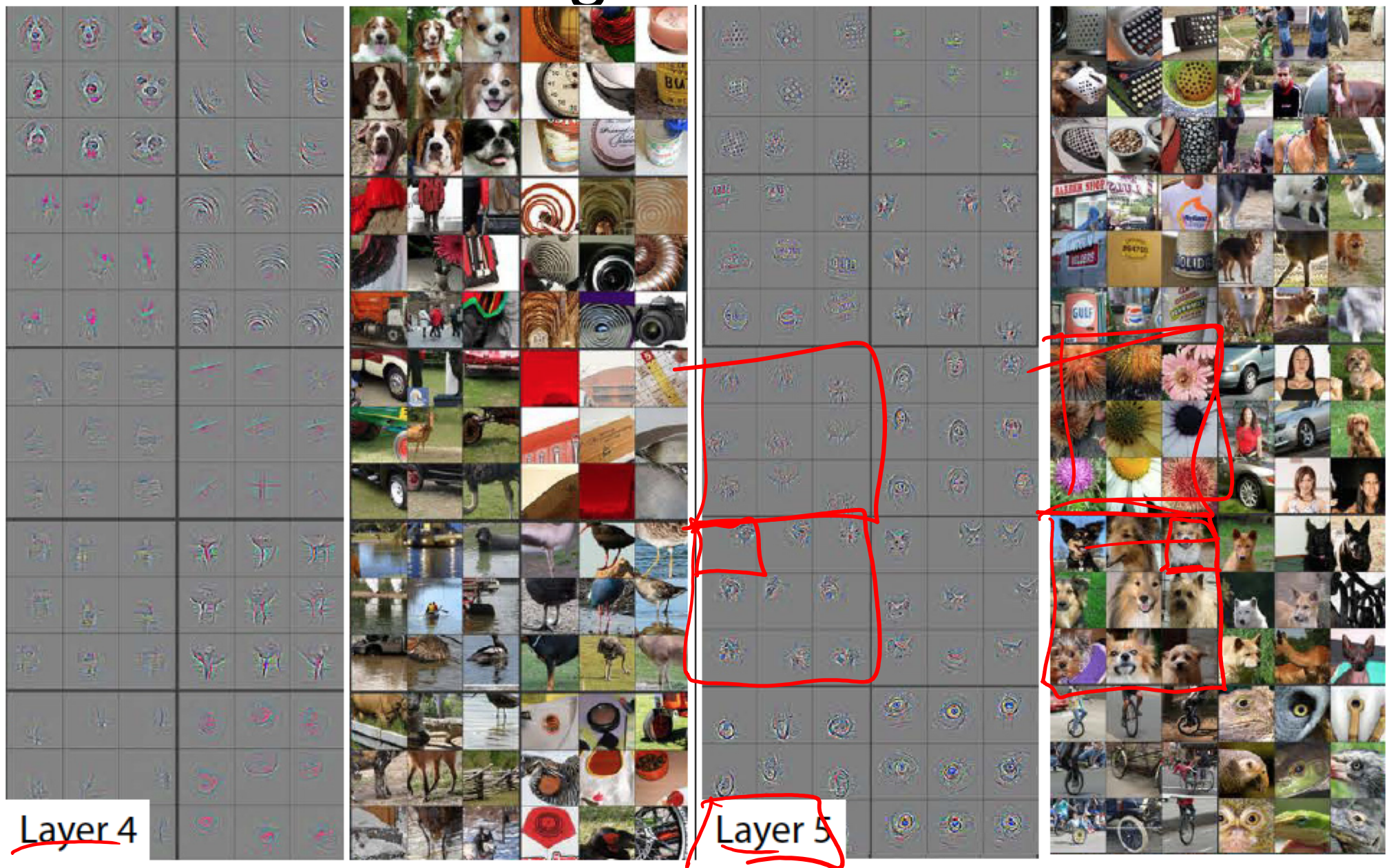
Visualizing Learned Filters



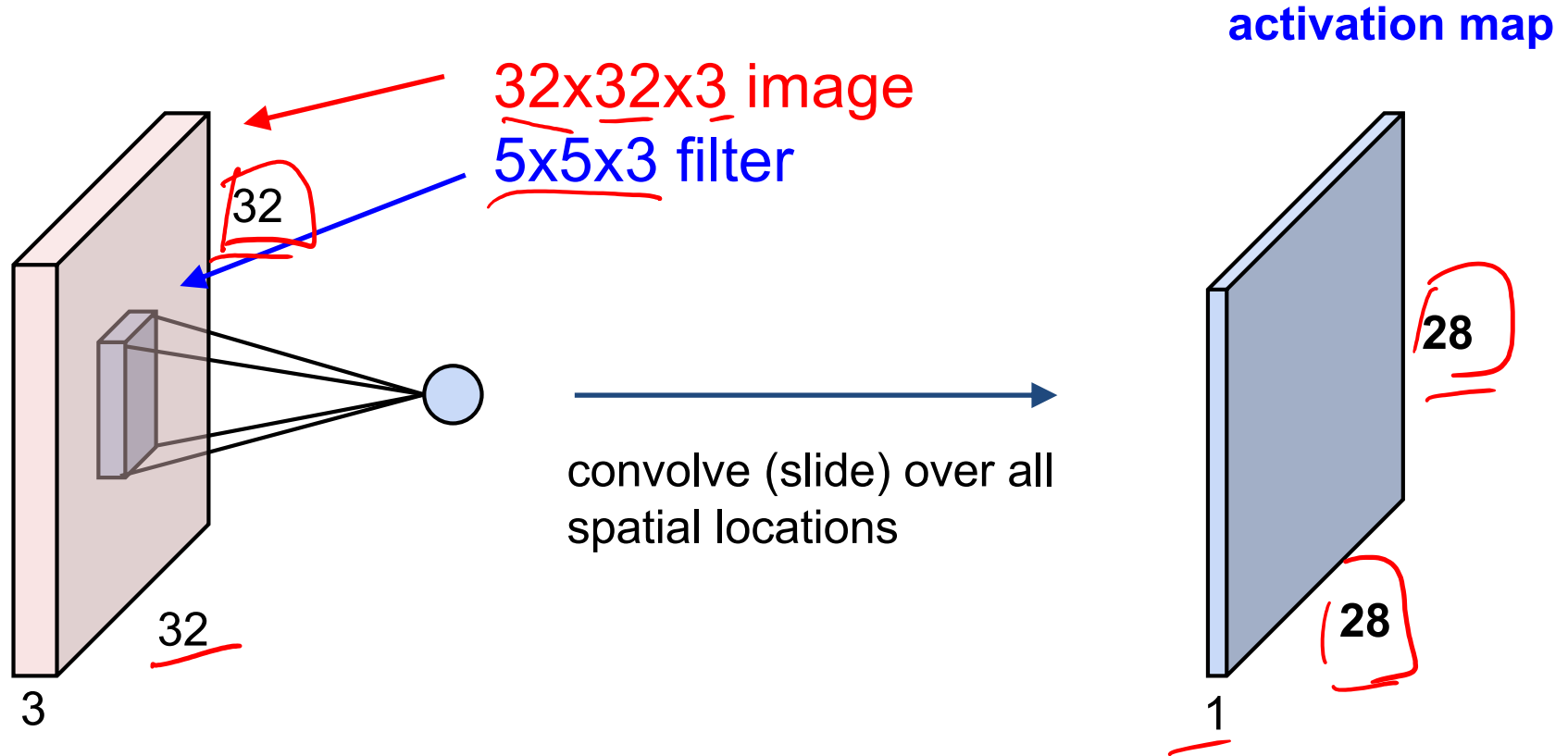
Visualizing Learned Filters



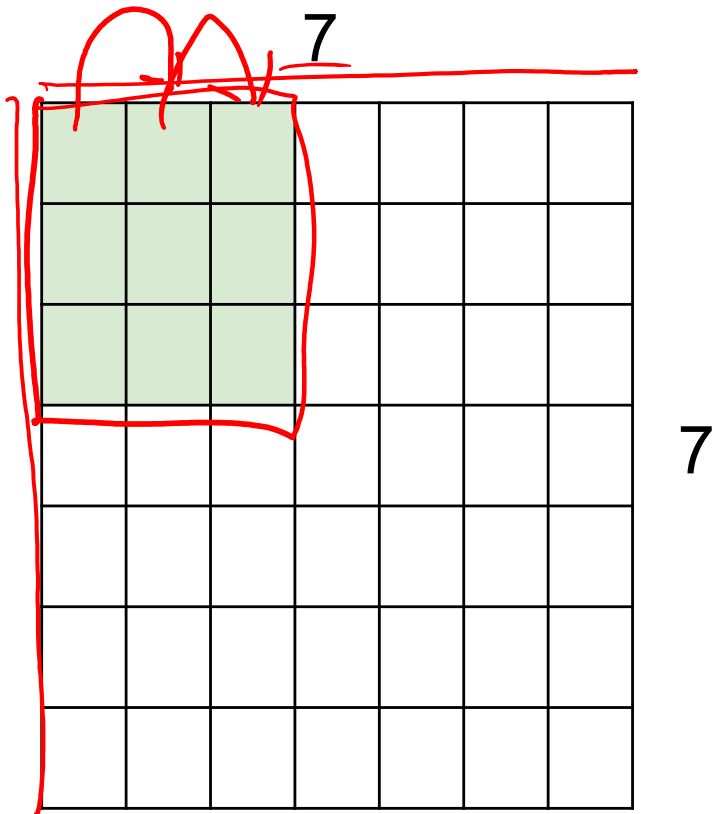
Visualizing Learned ²⁵⁶ Filters



A closer look at spatial dimensions:



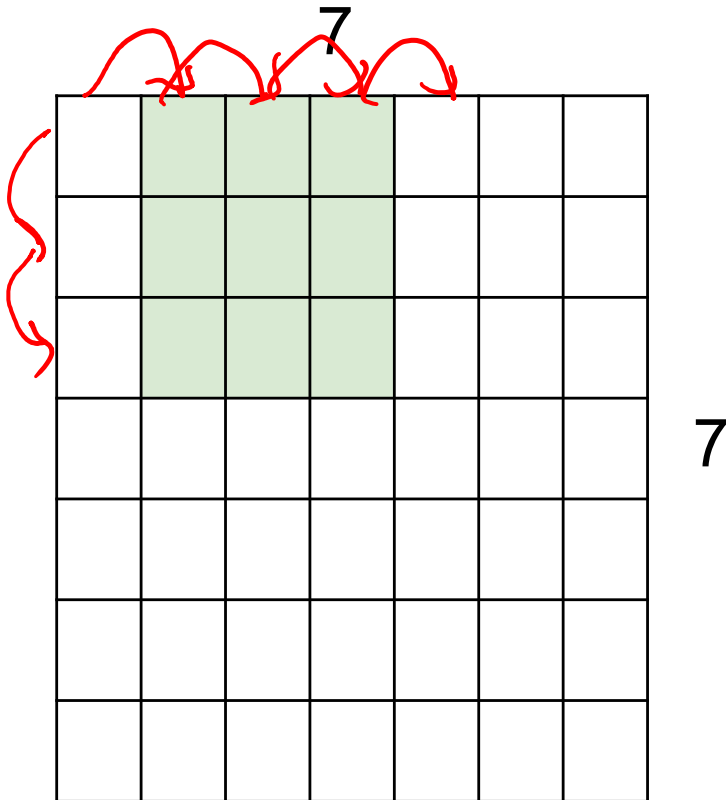
A closer look at spatial dimensions:



7x7 input (spatially)
assume 3x3 filter

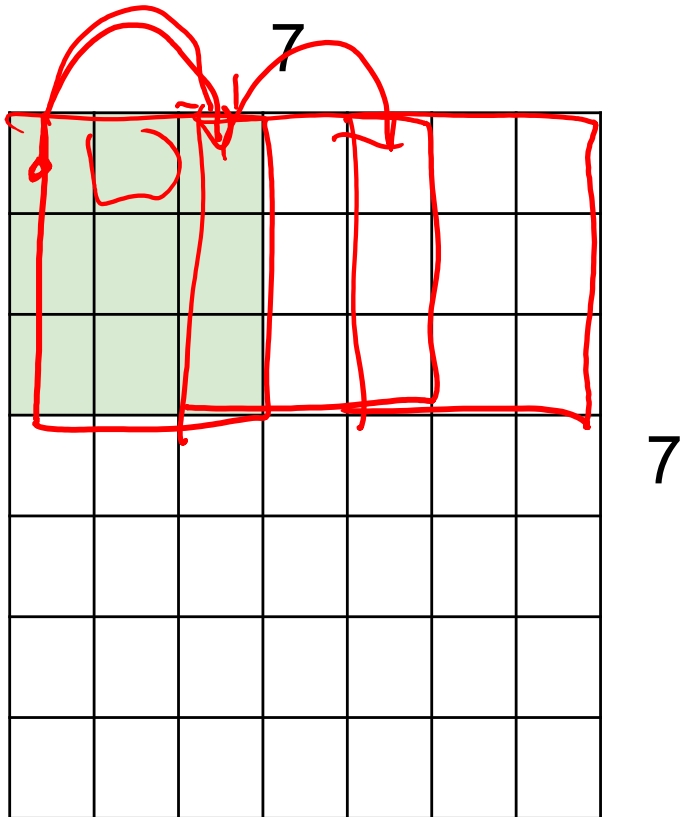
A closer look at spatial dimensions:

stride = 1



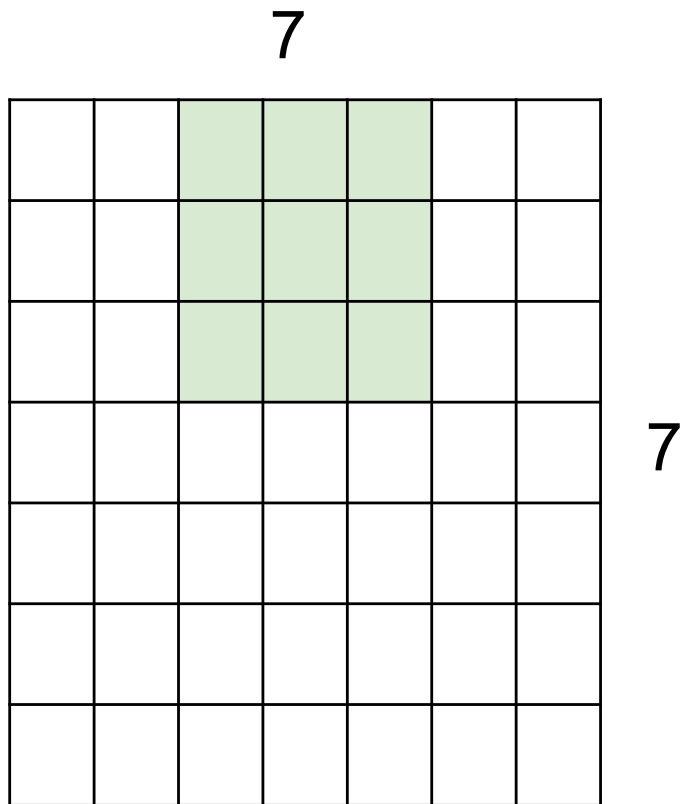
7x7 input (spatially)
assume 3x3 filter

A closer look at spatial dimensions:



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

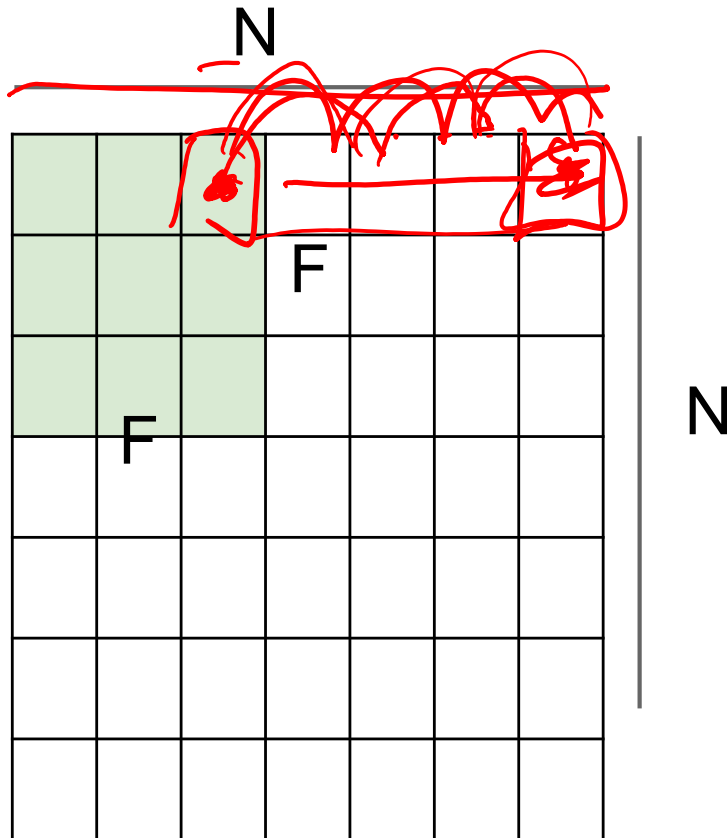
A closer look at spatial dimensions:



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

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

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



Output size:

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

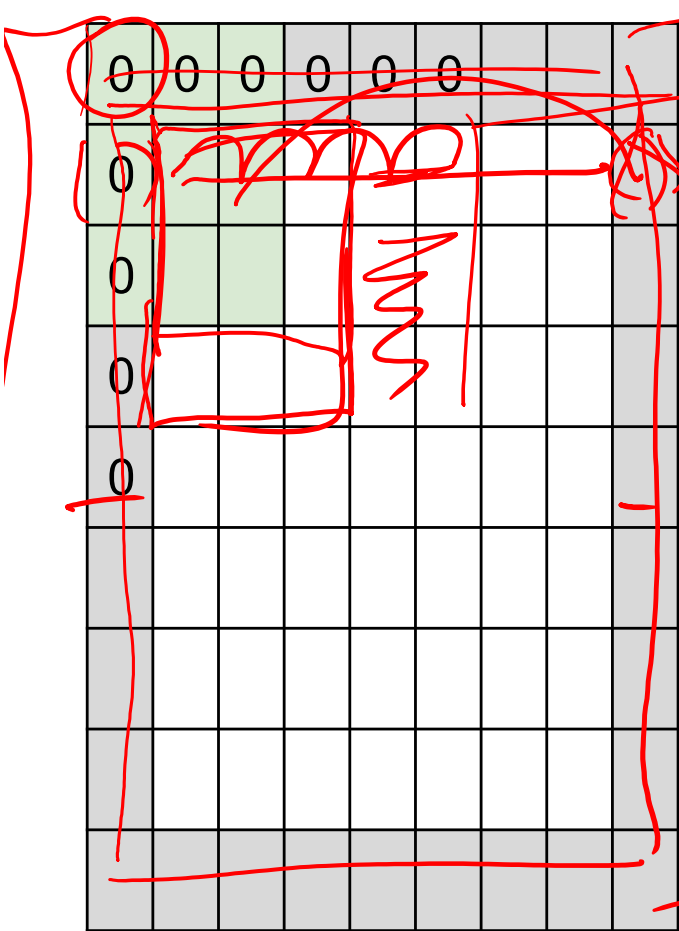
e.g. $N = 7$, $F = 3$:

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

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

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

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?

$N \Rightarrow \tilde{N}$
 $7 \times 7 \rightarrow \underline{\underline{9 \times 9}}$

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

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

(recall:)

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

$$\underline{\underline{N - F - 1}} \text{ pad} + 1 \Bigg\} \begin{matrix} F-1 \\ \hline 2 \end{matrix}$$

s-Stride

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. $F = 3 \Rightarrow$ zero pad with 1

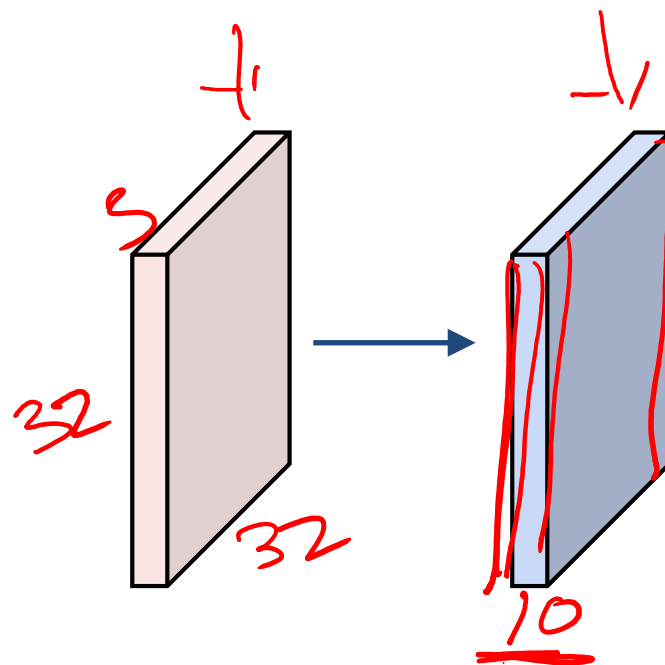
$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

Examples time:

Input volume: **32x32x3**

10 5x5 filters with stride 1, pad 2



Output volume size: ?

$$\frac{(32)}{N_f} \times \frac{(32)}{N_s} \times \frac{(10)}{C_2}$$

Examples time:

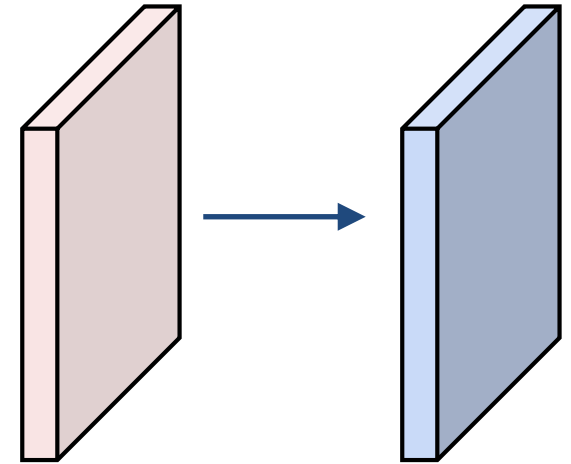
Input volume: **32x32x3**

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

Output volume size:

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

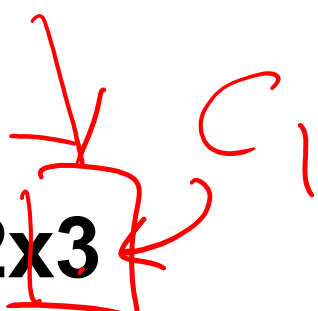
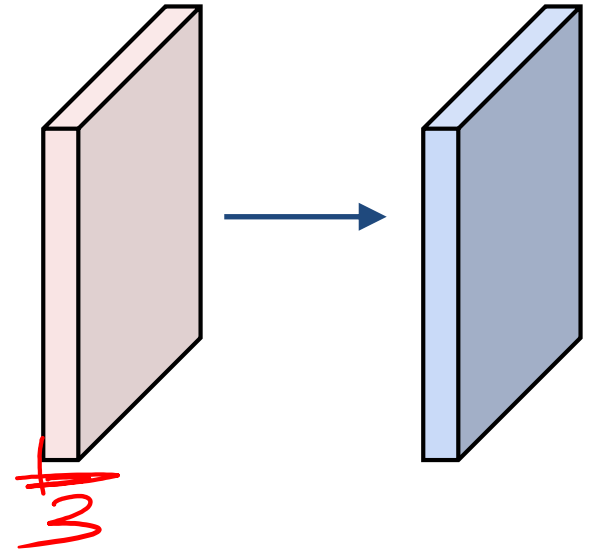
32x32x10



Examples time:

Input volume: **32x32x3**

10, 5x5 filters with stride 1, pad 2



Number of parameters in this layer?

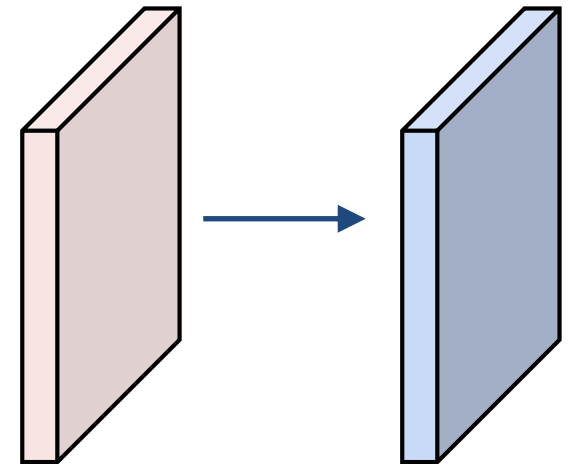
$$(5 \times 5 \times 3 + 1) \times 10 = 760$$

~~() 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)
- $F = 1, S = 1, P = 0$

Example: CONV layer in Torch

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 .

SpatialConvolution

```
module = nn.SpatialConvolution(nInputPlane, nOutputPlane, kW, kH, [dW], [dH], [padW], [padH])
```

Applies a 2D convolution over an input image composed of several input planes. The `input` tensor in `forward(input)` is expected to be a 3D tensor (`nInputPlane x height x width`).

The parameters are the following:

- `nInputPlane` : The number of expected input planes in the image given into `forward()` .
- `nOutputPlane` : The number of output planes the convolution layer will produce.
- `kW` : The kernel width of the convolution
- `kH` : The kernel height of the convolution
- `dW` : The step of the convolution in the width dimension. Default is `1` .
- `dH` : The step of the convolution in the height dimension. Default is `1` .
- `padW` : The additional zeros added per width to the input planes. Default is `0` , a good number is `(kW-1)/2` .
- `padH` : The additional zeros added per height to the input planes. Default is `padW` , a good number is `(kH-1)/2` .

Note that depending of the size of your kernel, several (of the last) columns or rows of the input image might be lost. It is up to the user to add proper padding in images.

If the input image is a 3D tensor `nInputPlane x height x width` , the output image size will be `nOutputPlane x oheight x owidth` where

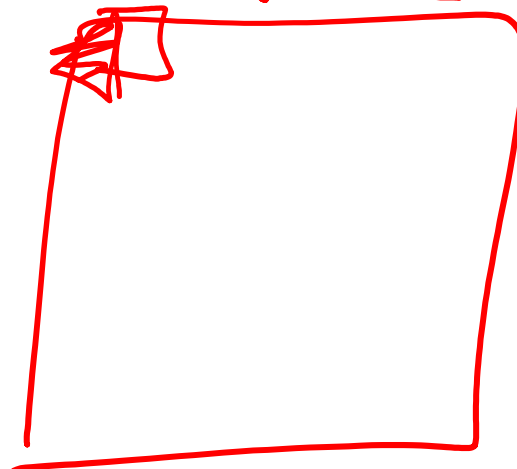
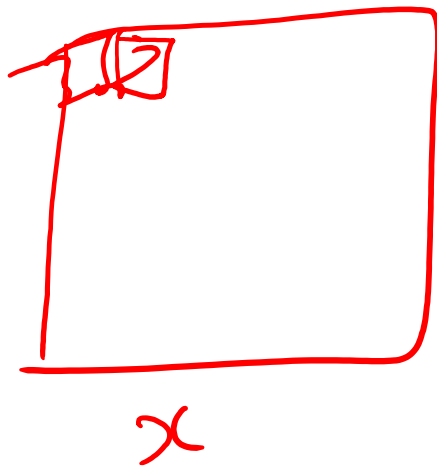
```
owidth = floor((width + 2*padW - kW) / dW + 1)  
oheight = floor((height + 2*padH - kH) / dH + 1)
```

[Torch](#) is licensed under [BSD 3-clause](#).

Plan for Today

- Convolutional Neural Networks
 - 1x1 convolutions
 - Pooling layers
 - Fully-connected layers as convolutions
 - Backprop in conv layers
 - Toeplitz matrices and convolutions = matrix-mult
 - Dilated/a-trous convolutions

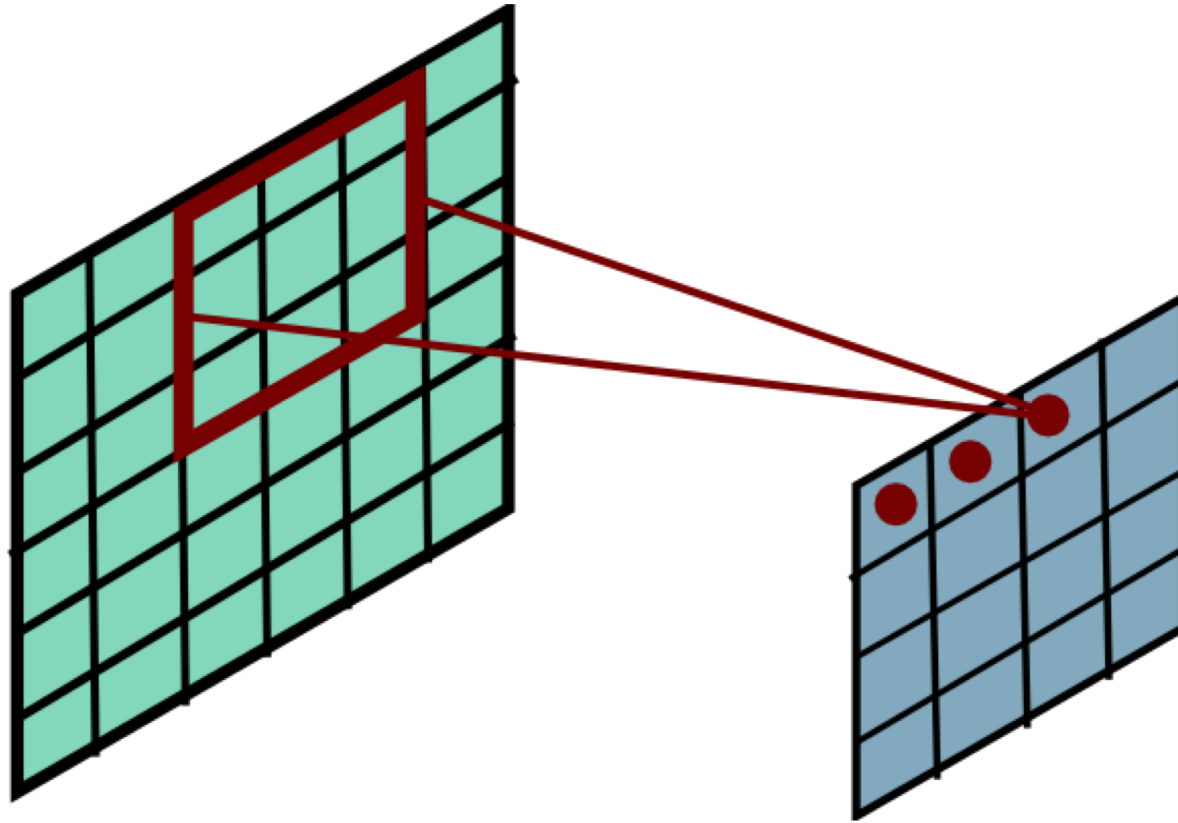
Can we have 1x1 filters?



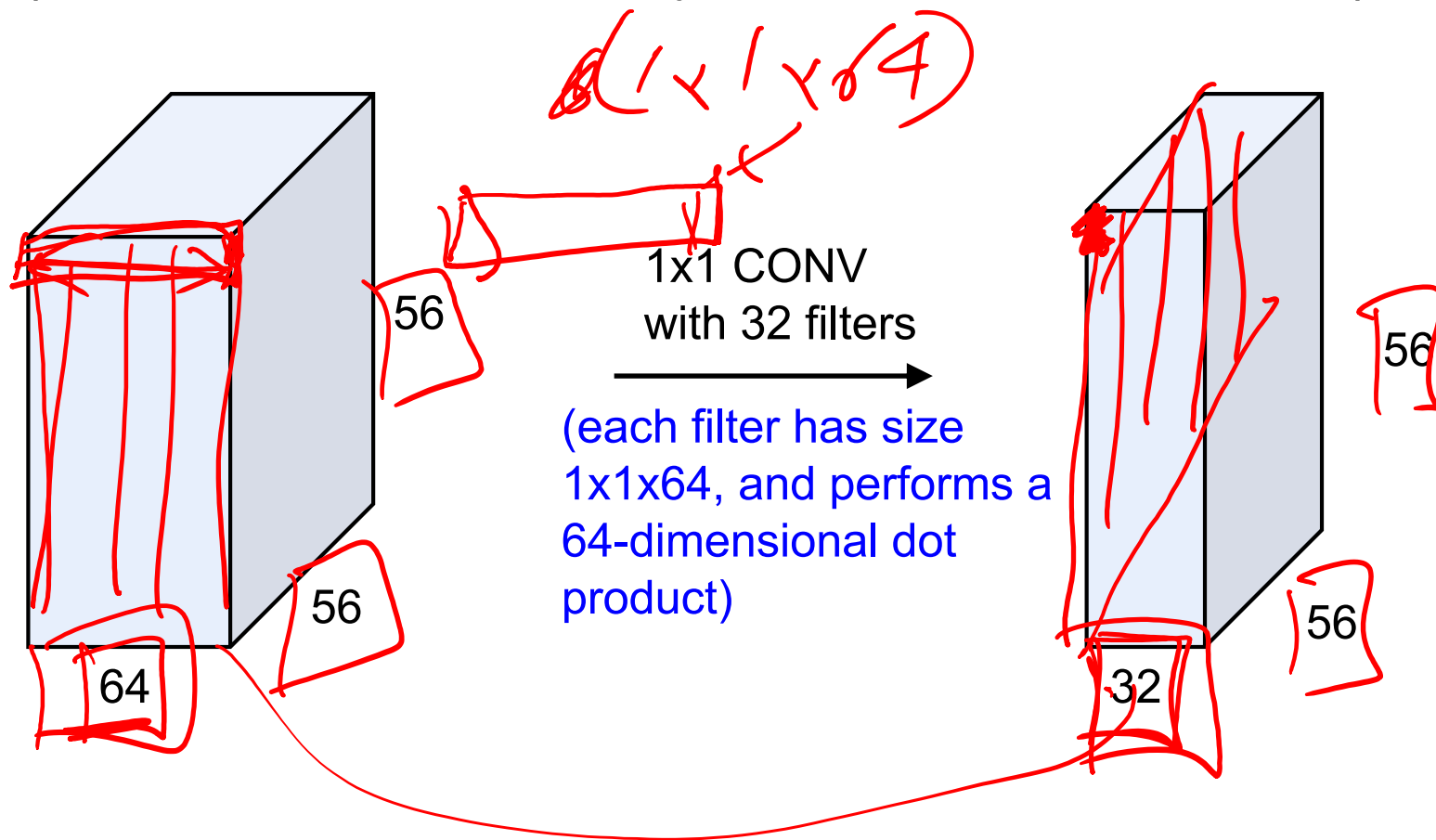
$$y[r, c] = \sum_{a=0}^{K_r-1} \sum_{b=0}^{K_c-1} x[r+a, c+b] w[a, b]$$

$$y[r, c] = x[r, c] \cdot w[0, 0]$$

Convolutional Layer



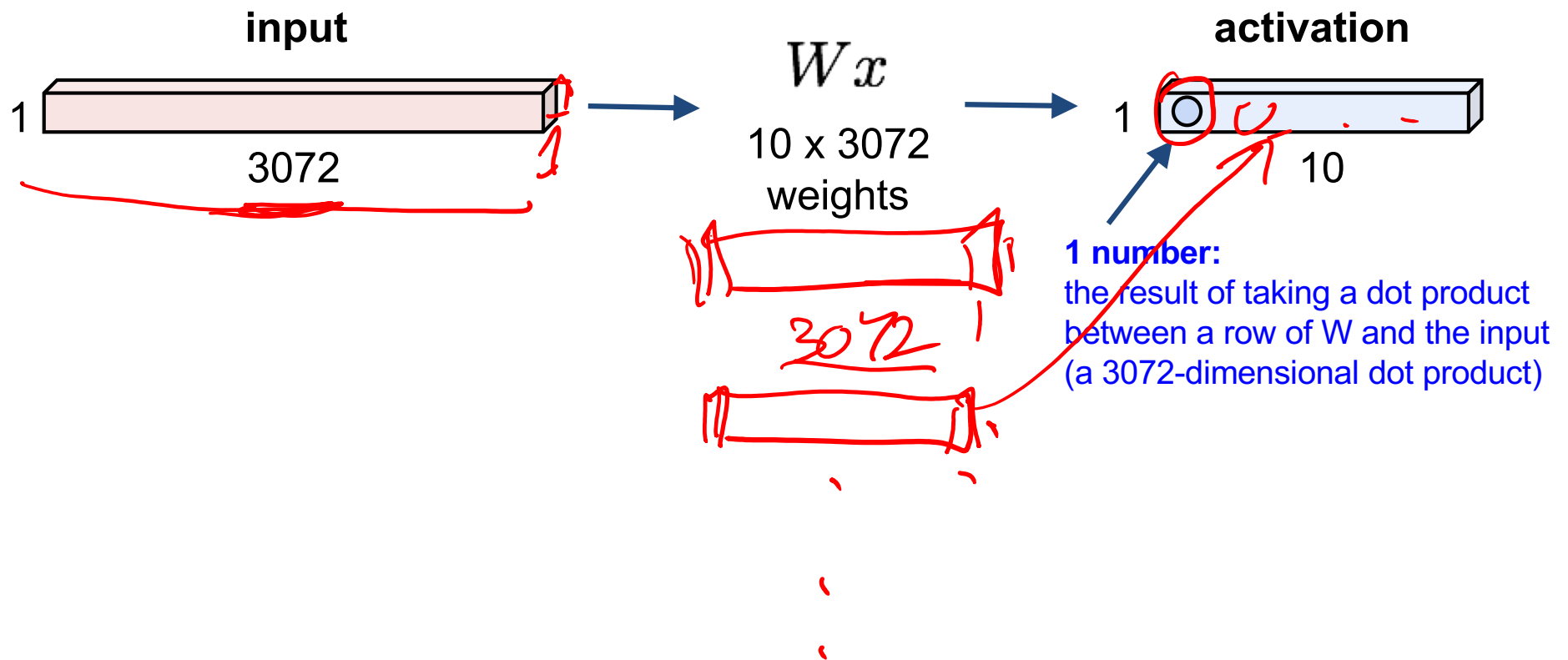
(btw, 1x1 convolution layers make perfect sense)



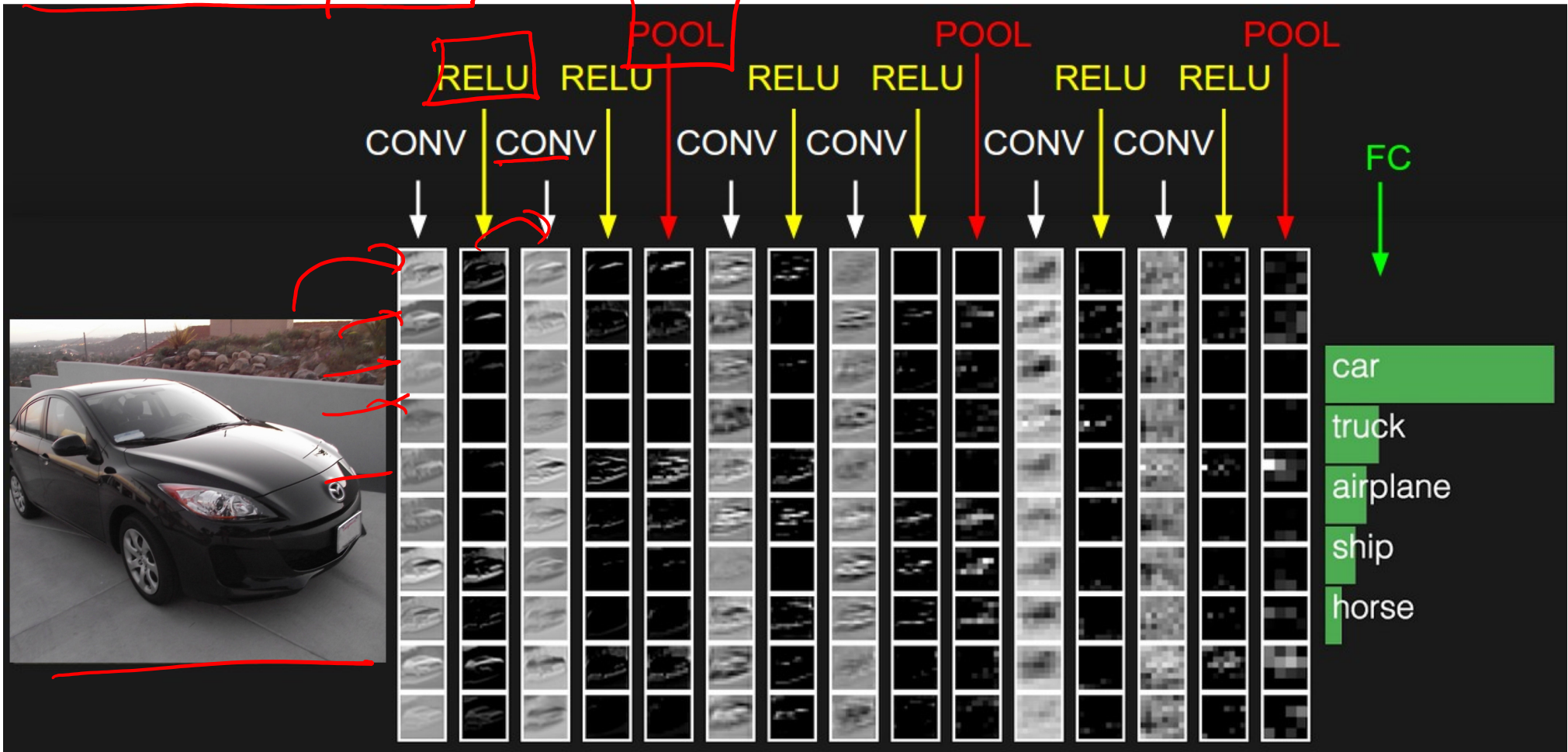
Reminder: Fully Connected Layer

32x32x3 image \rightarrow stretch to 3072 x 1

Each neuron looks at the full input volume



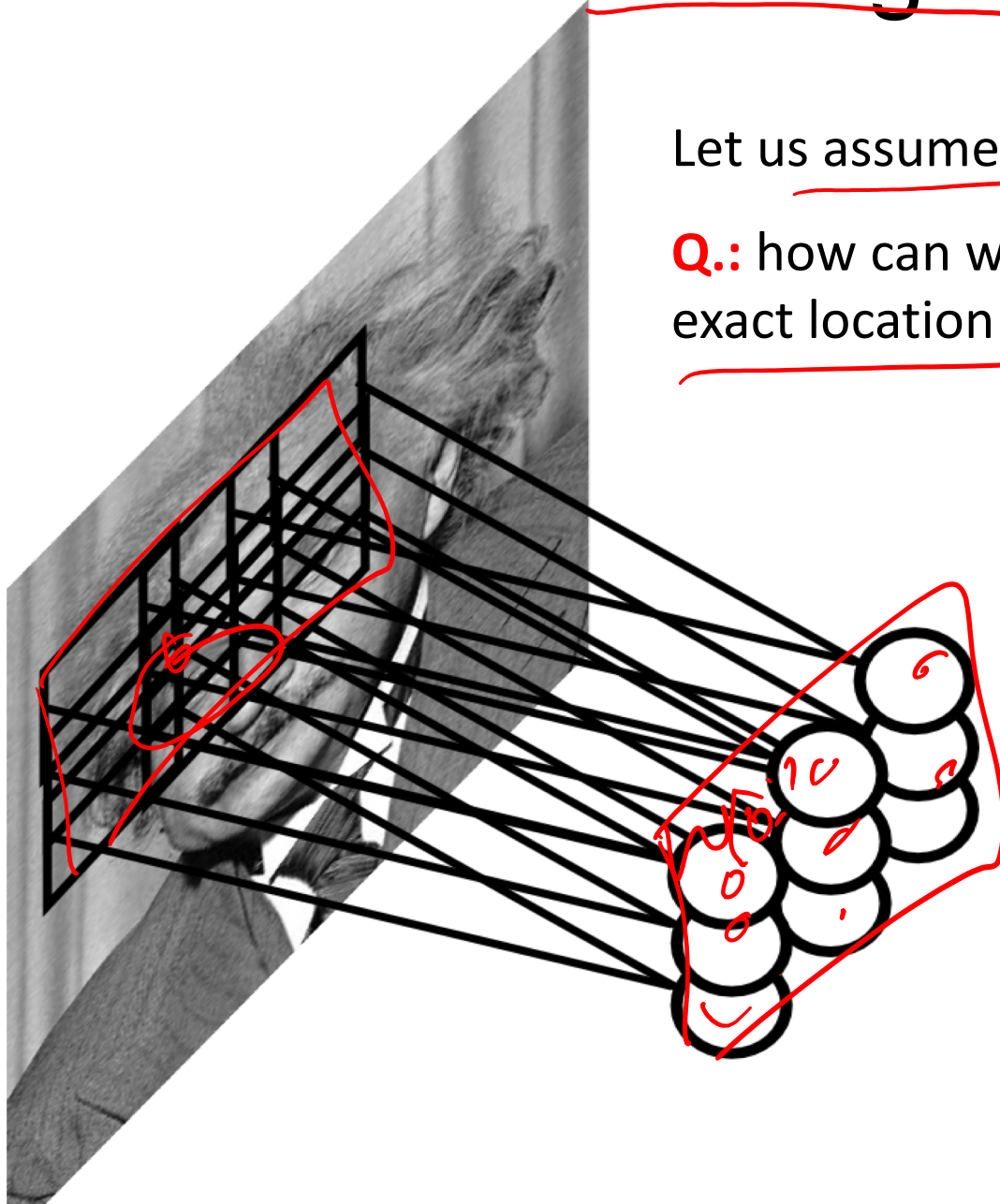
two more layers to go: POOL/FC



Pooling Layer

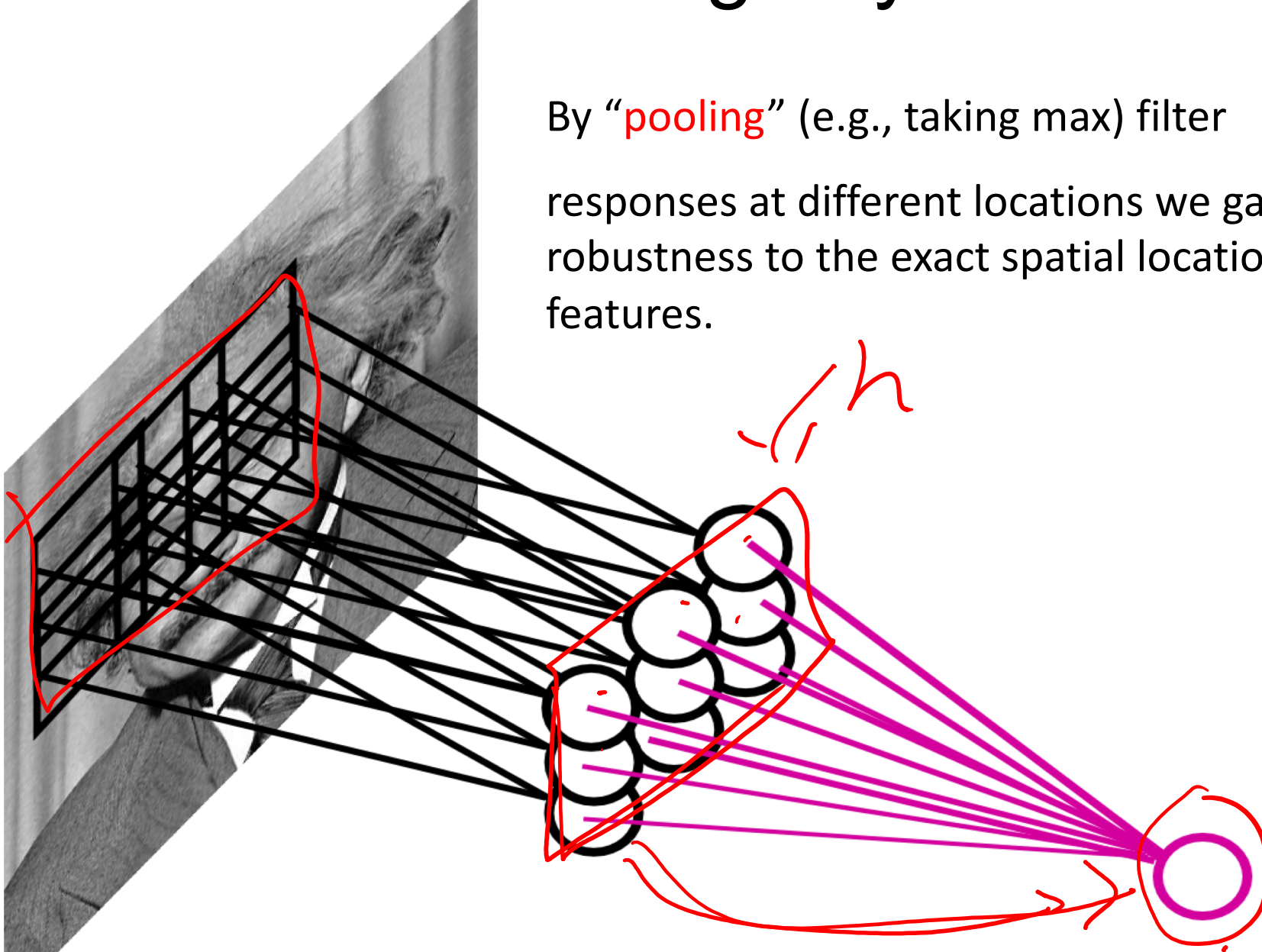
Let us assume filter is an “eye” detector.

Q.: how can we make the detection robust to the exact location of the eye?



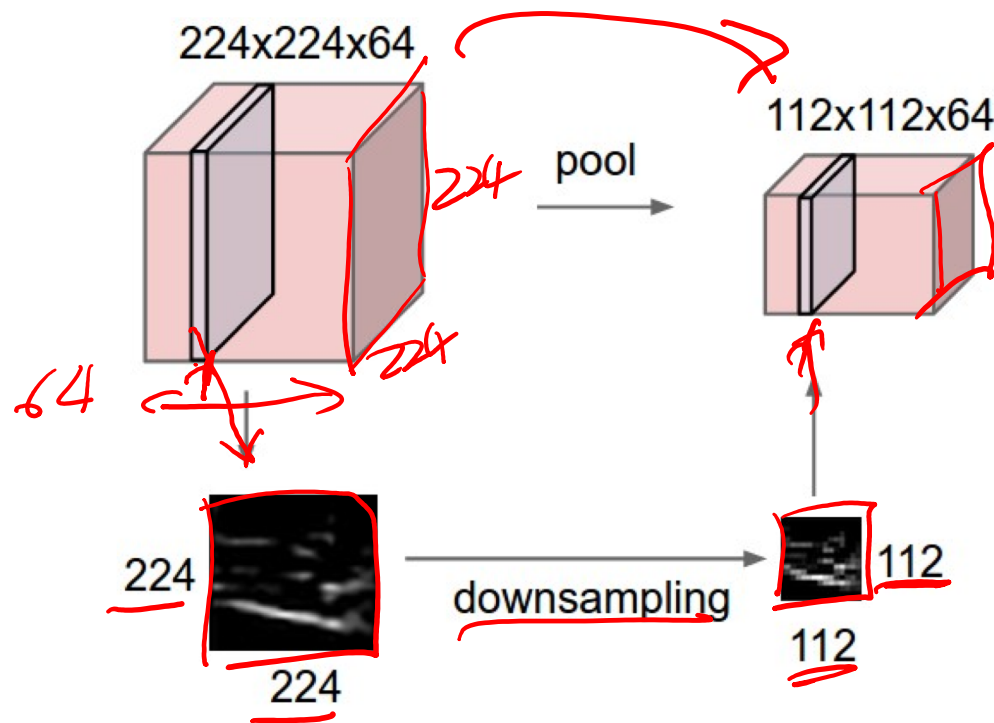
Pooling Layer

By “pooling” (e.g., taking max) filter responses at different locations we gain robustness to the exact spatial location of features.



Pooling layer

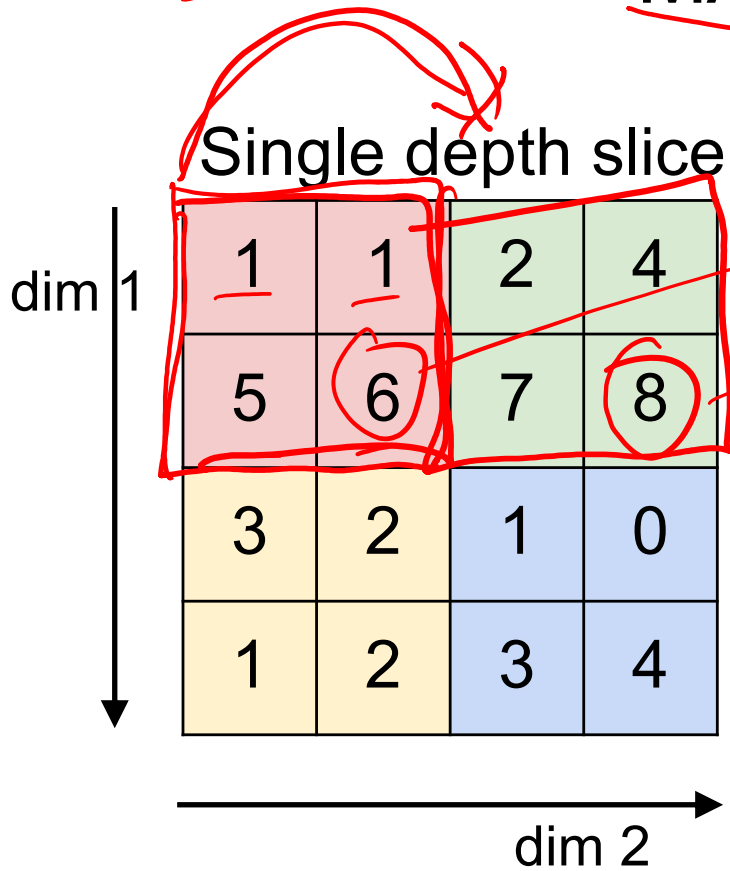
- makes the representations smaller and more manageable
- operates over each activation map independently:



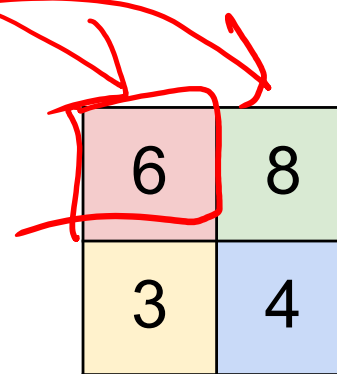
stride 2

MAX POOLING

2x2



max pool with 2x2 filters and stride 2



$$y[r, c] = \max_a \max_b x[r+a, c+b]$$

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

-	-
-	-

Pooling Layer: Examples

Max-pooling:

$$h_i^n(r, c) = \max_{\bar{r} \in N(r), \bar{c} \in N(c)} h_i^{n-1}(\bar{r}, \bar{c})$$

Average-pooling:

$$h_i^n(r, c) = \text{mean}_{\bar{r} \in N(r), \bar{c} \in N(c)} h_i^{n-1}(\bar{r}, \bar{c})$$

L2-pooling:

$$h_i^n(r, c) = \sqrt{\sum_{\bar{r} \in N(r), \bar{c} \in N(c)} h_i^{n-1}(\bar{r}, \bar{c})^2}$$

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F) / S + 1$
 - $H_2 = (H_1 - F) / S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

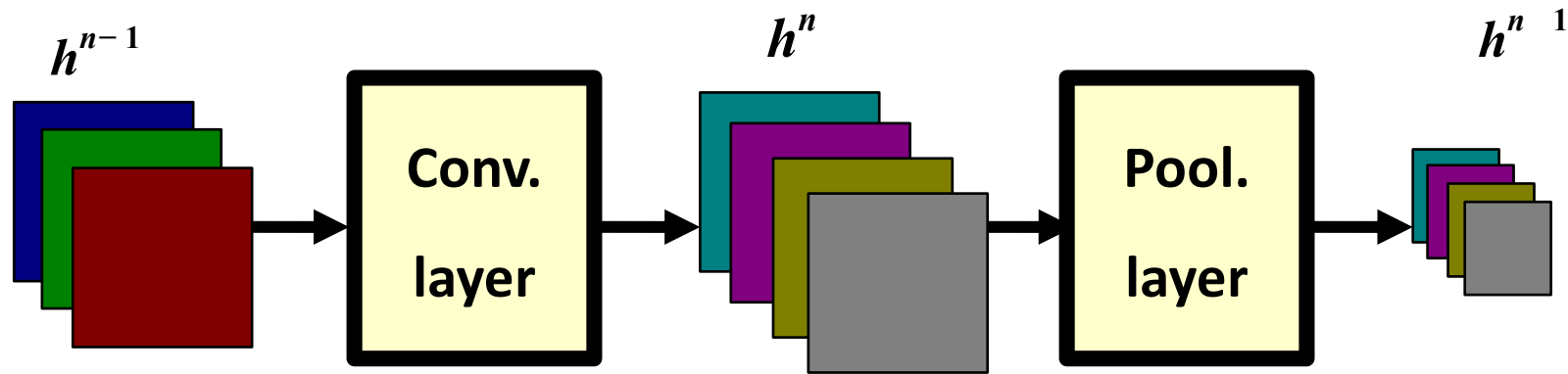
Common settings:

$$F = 2, S = 2$$

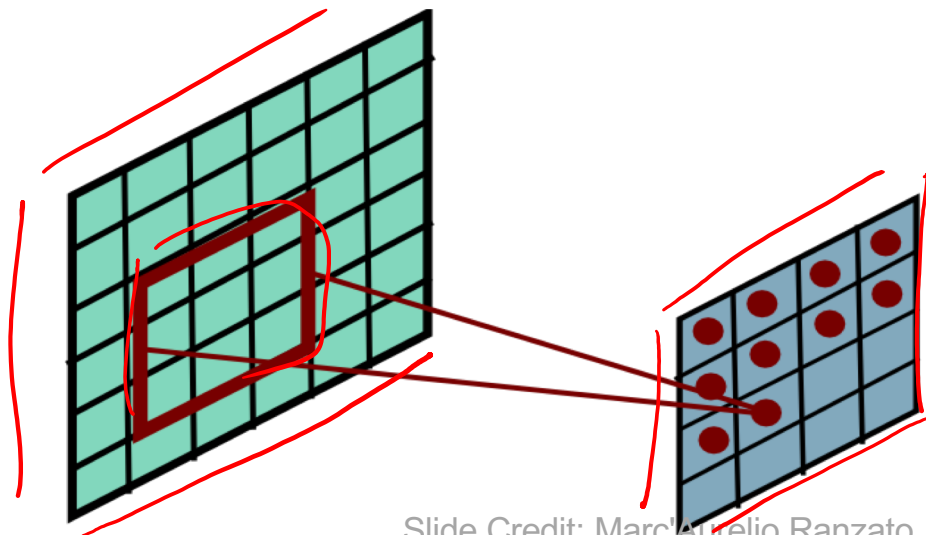
$$F = 3, S = 2$$

- Accepts a volume of size $W_1 \times H_1 \times D_1$
- Requires three hyperparameters:
 - their spatial extent F ,
 - the stride S ,
- Produces a volume of size $W_2 \times H_2 \times D_2$ where:
 - $W_2 = (W_1 - F)/S + 1$
 - $H_2 = (H_1 - F)/S + 1$
 - $D_2 = D_1$
- Introduces zero parameters since it computes a fixed function of the input
- Note that it is not common to use zero-padding for Pooling layers

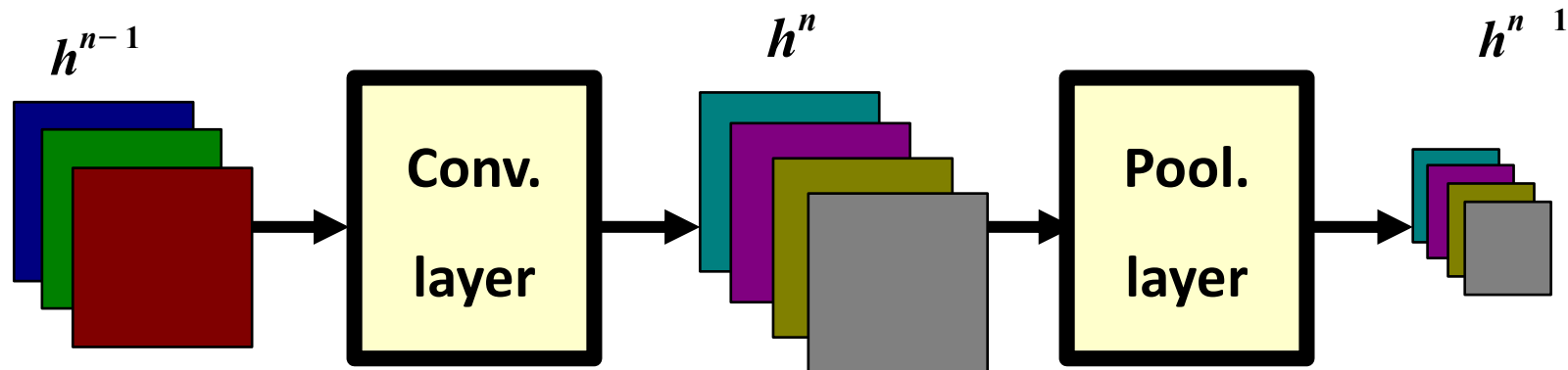
Pooling Layer: Receptive Field Size



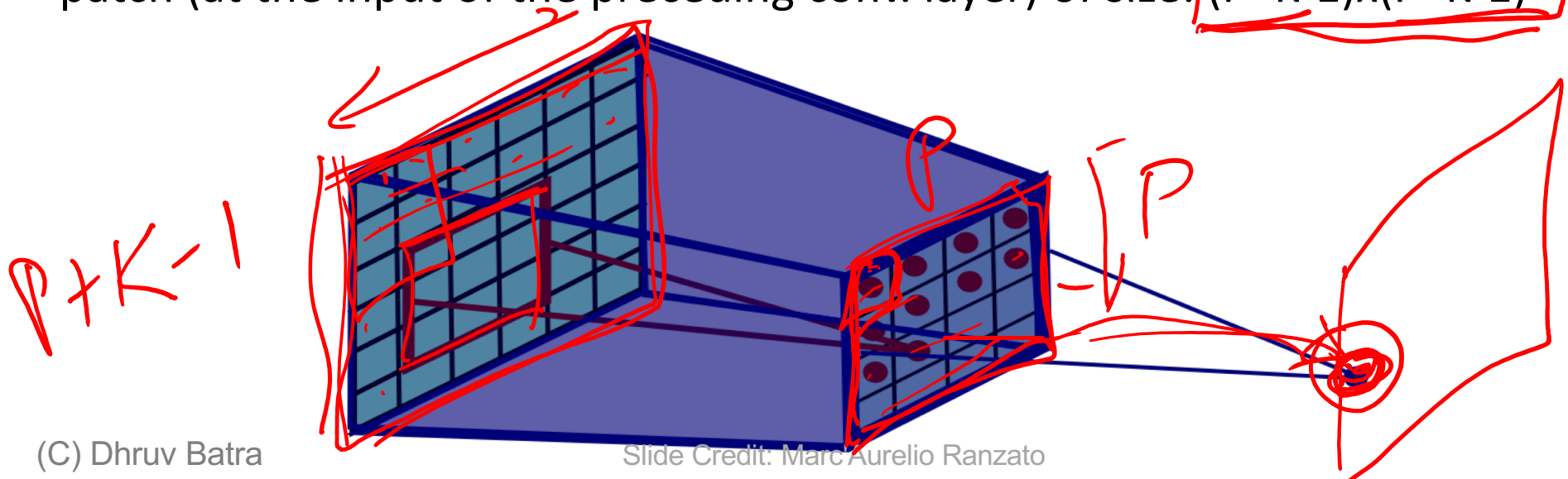
If convolutional filters have size $K \times K$ and stride 1, and pooling layer has pools of size $P \times P$, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size: $(P+K-1) \times (P+K-1)$



Pooling Layer: Receptive Field Size



If convolutional filters have size $K \times K$ and stride 1, and pooling layer has pools of size $P \times P$, then each unit in the pooling layer depends upon a patch (at the input of the preceding conv. layer) of size: $(P+K-1) \times (P+K-1)$

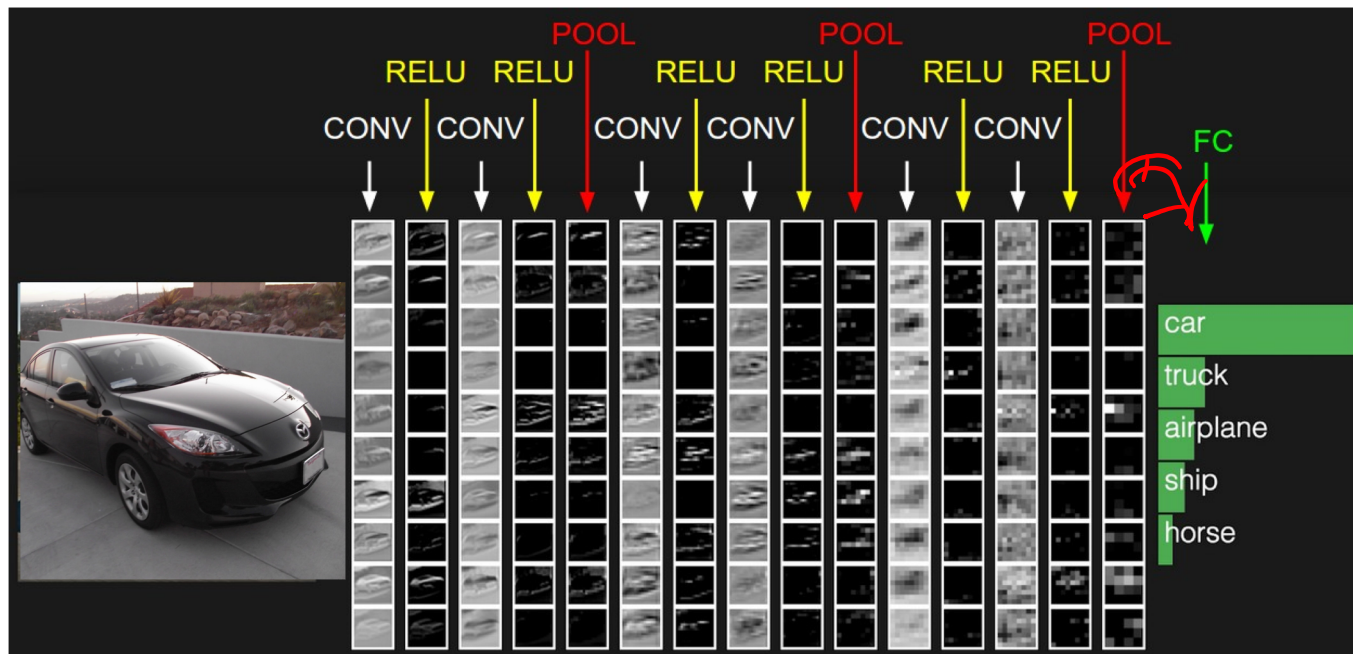


Plan for Today

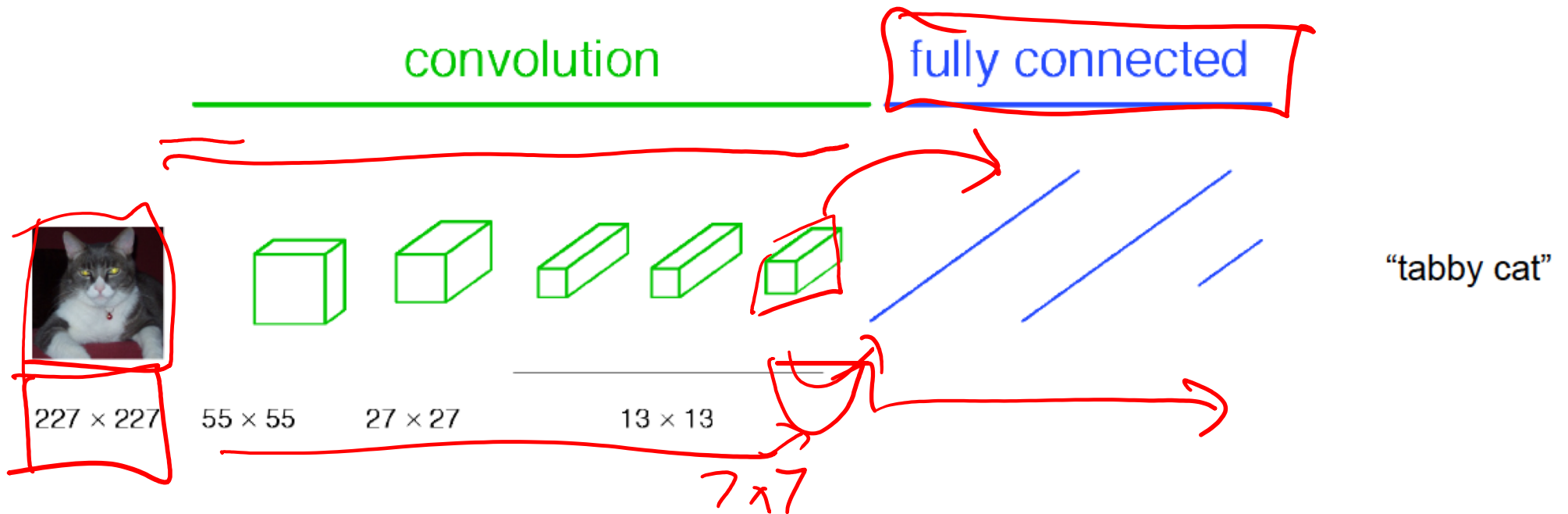
- Convolutional Neural Networks
 - 1x1 convolutions
 - Pooling layers
 - Fully-connected layers as convolutions
 - Backprop in conv layers

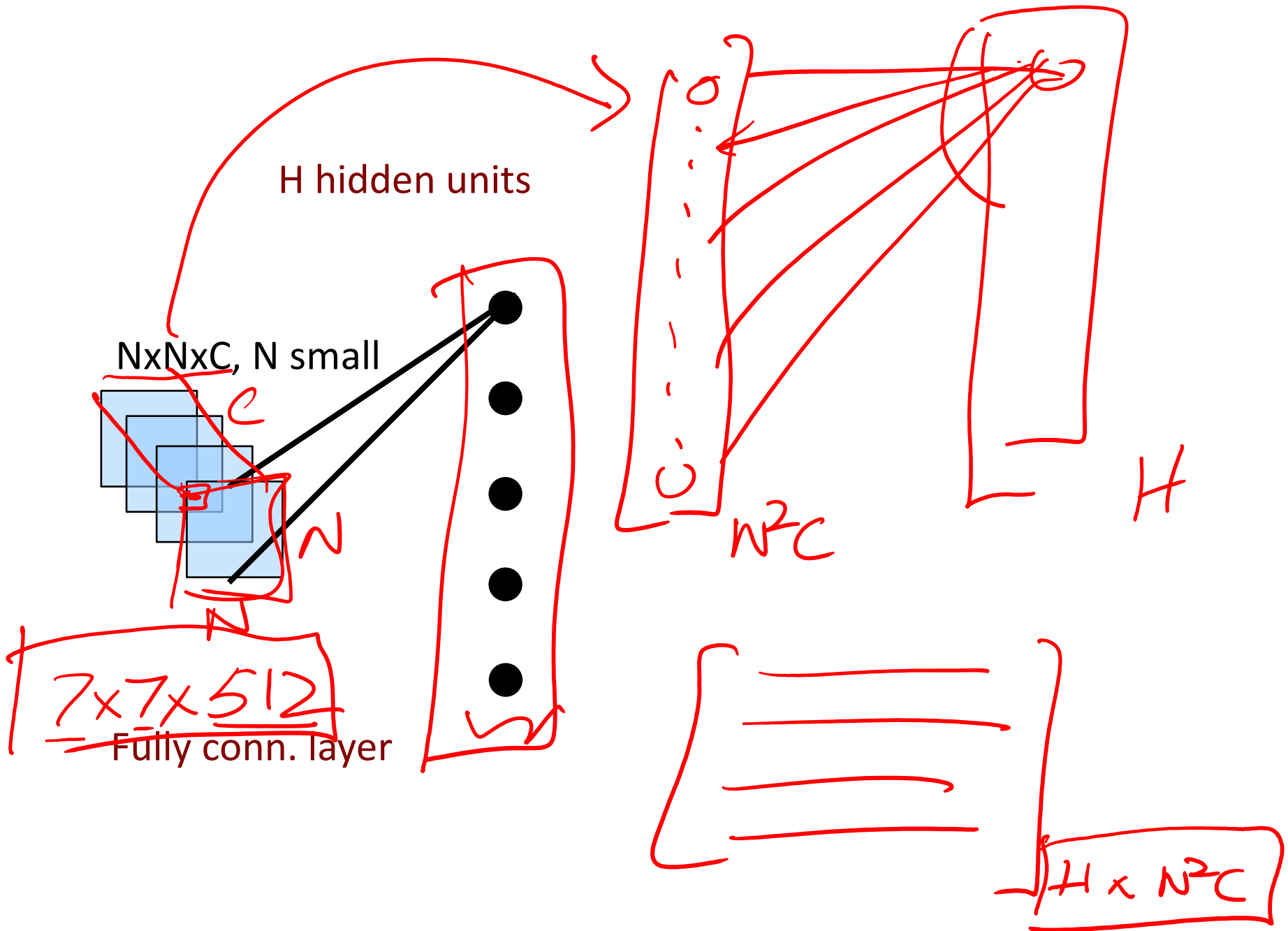
Fully Connected Layer (FC layer)

- Contains neurons that connect to the entire input volume, as in ordinary Neural Networks

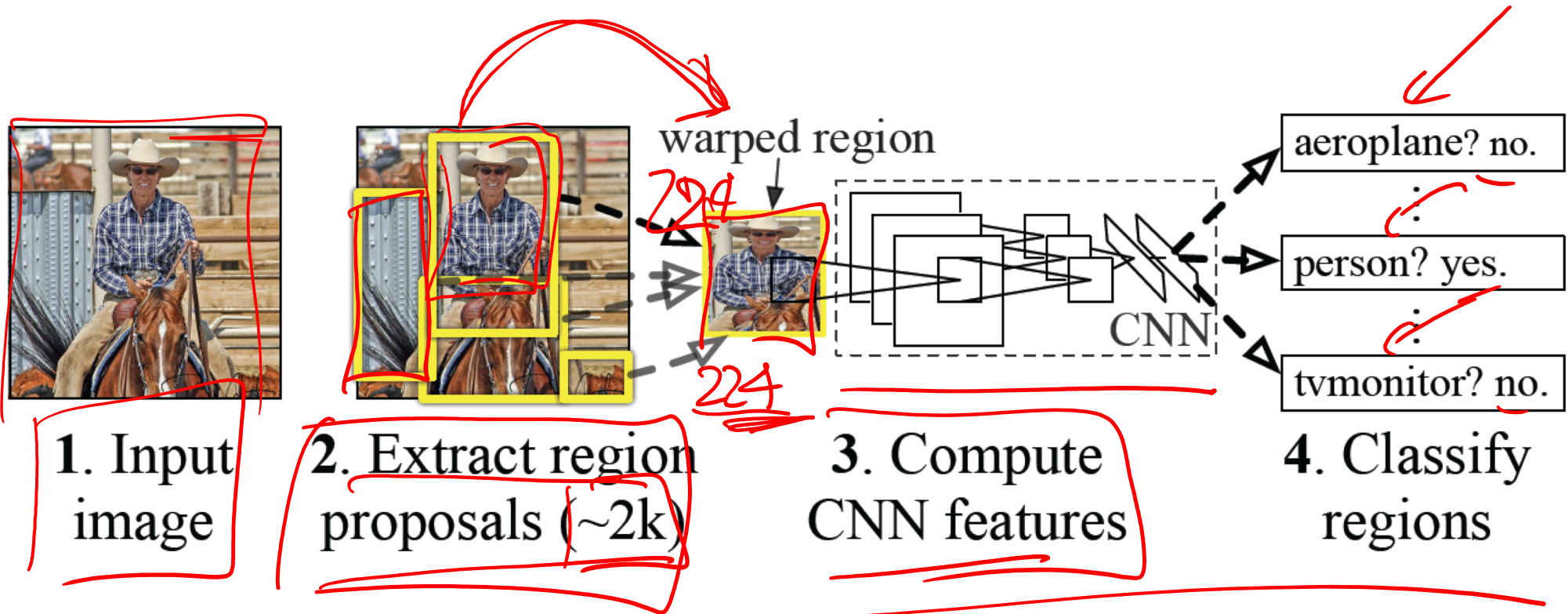


Classical View

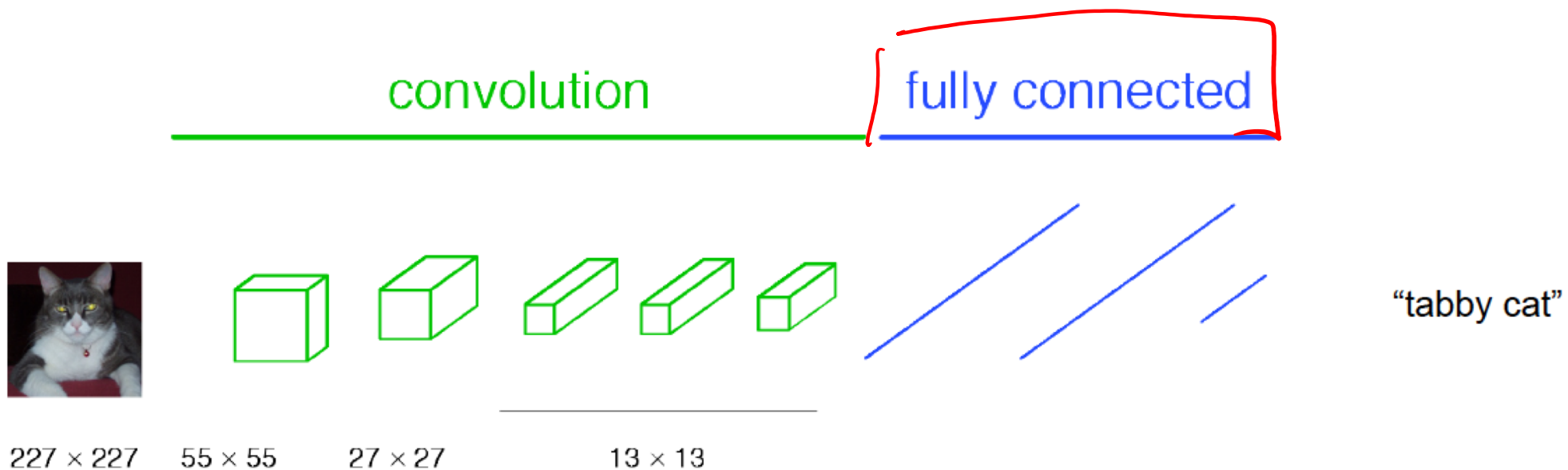




Classical View = Inefficient

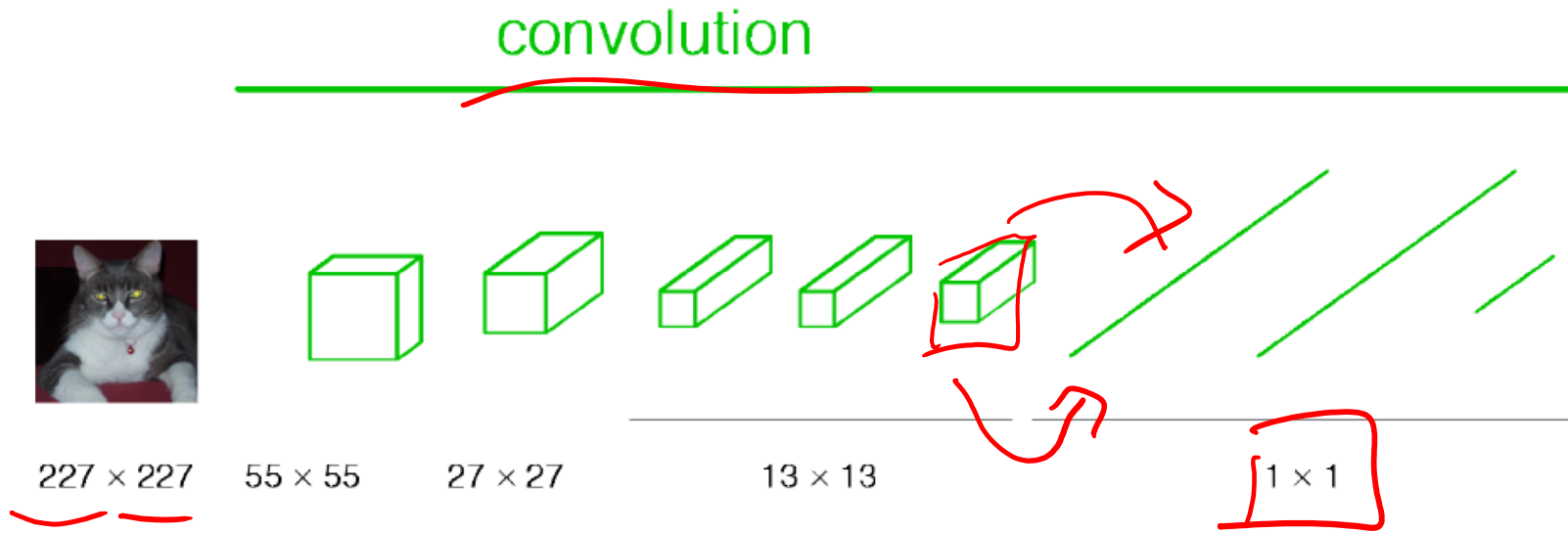


Classical View



Re-interpretation

- Just squint a little!



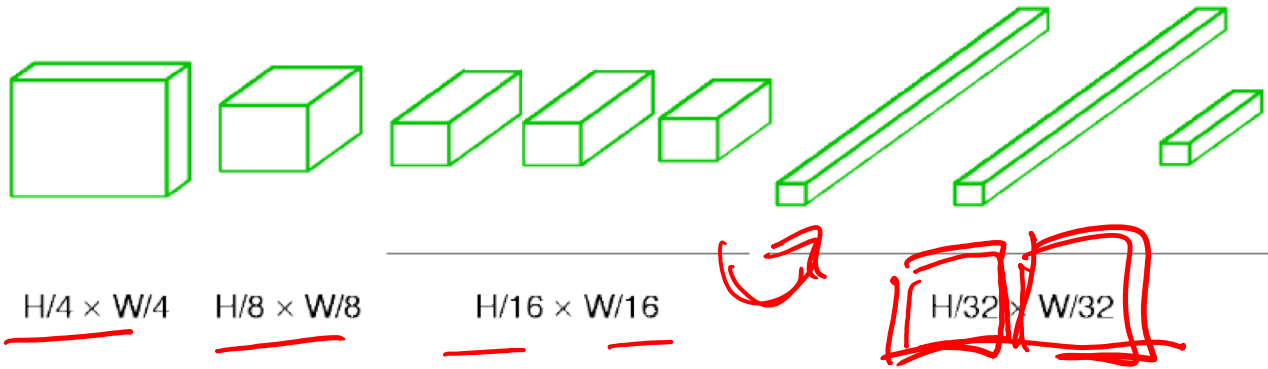
“Fully Convolutional” Networks

- Can run on an image of any size!



$H \times W$

convolution

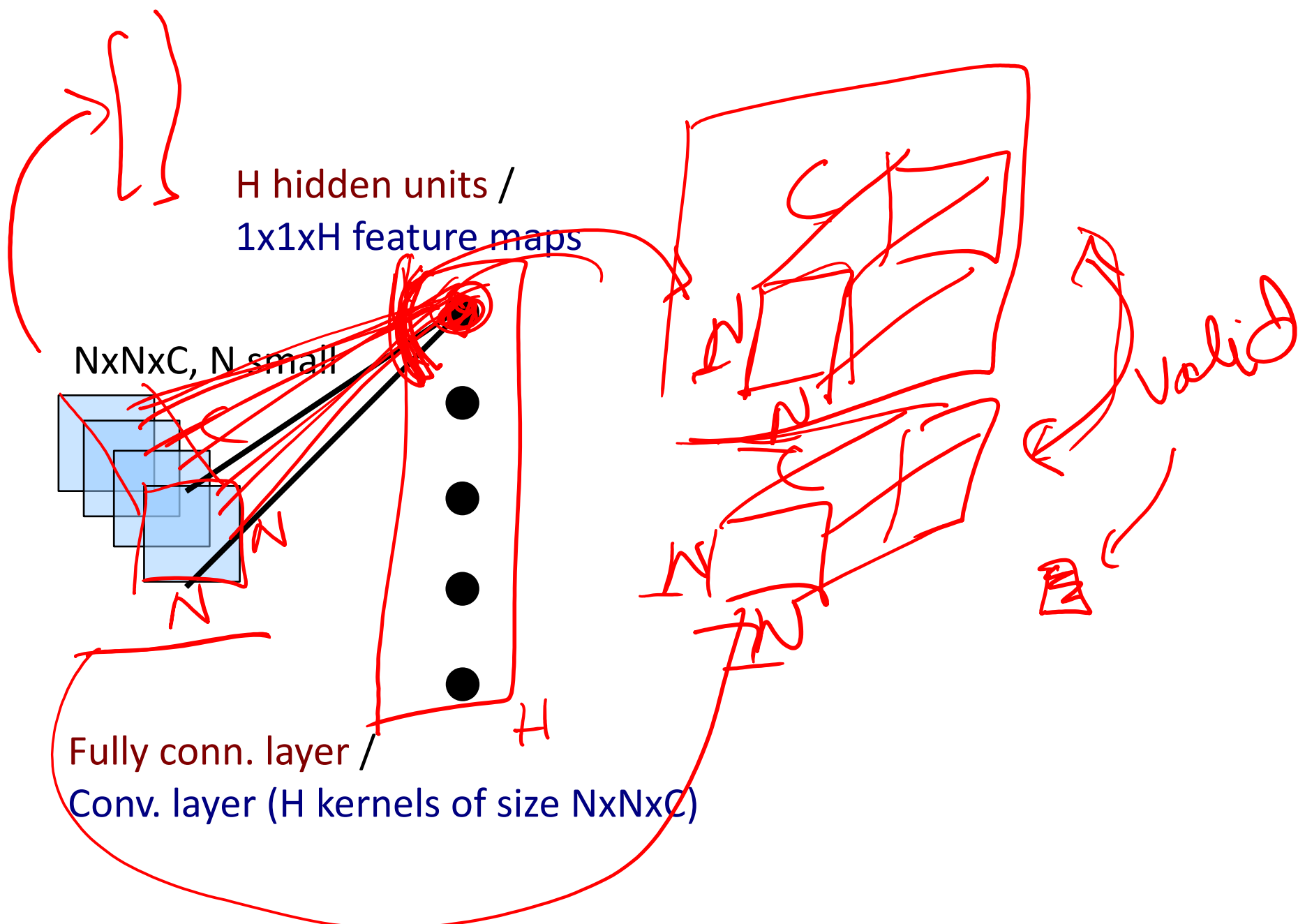


$H/4 \times W/4$

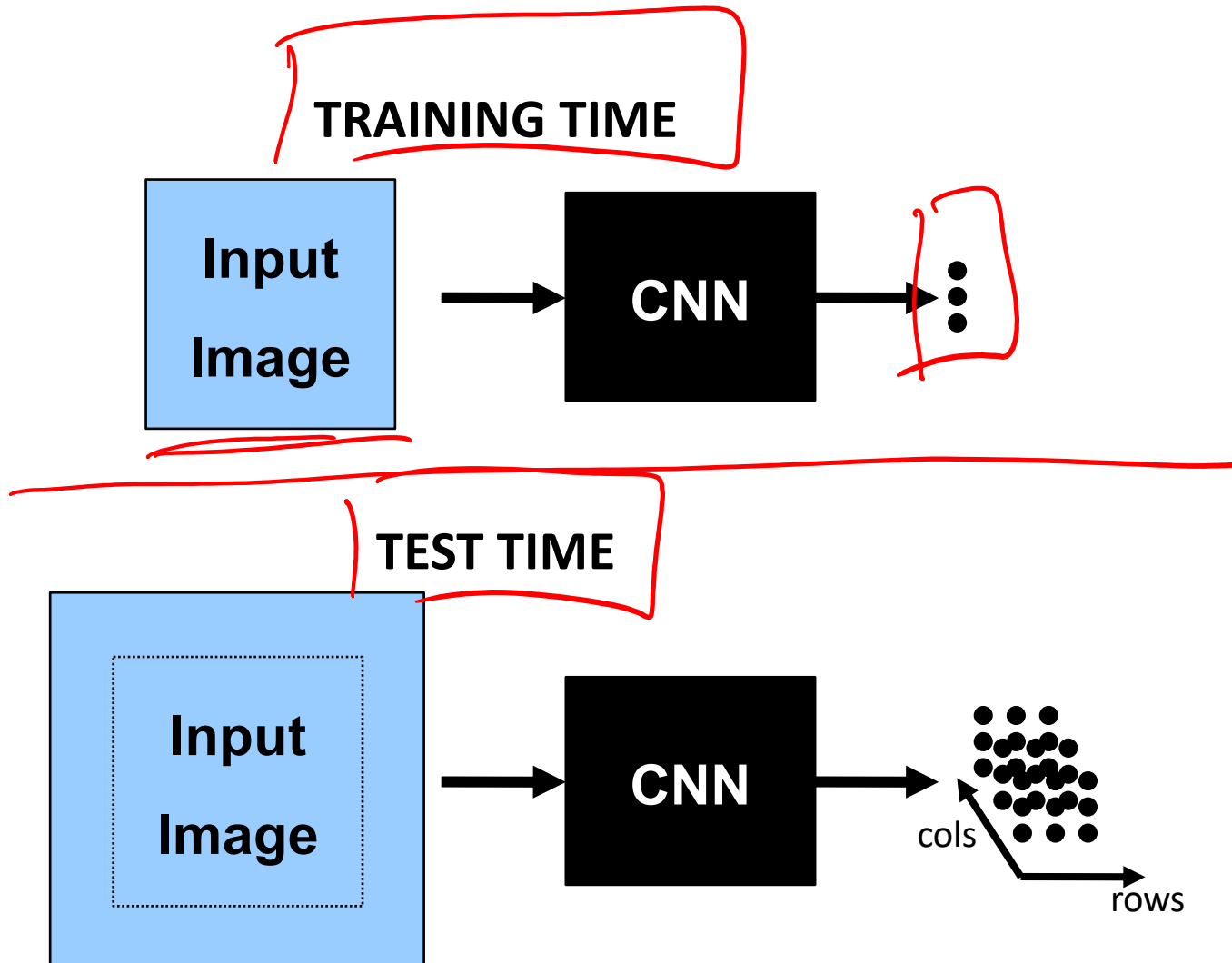
$H/8 \times W/8$

$H/16 \times W/16$

$H/32 \times W/32$

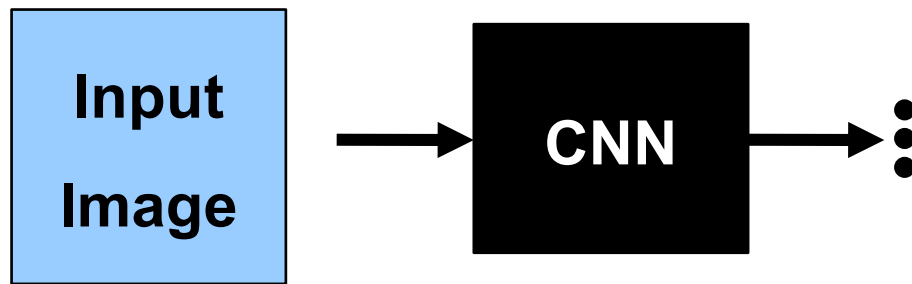


Viewing fully connected layers as convolutional layers enables efficient use of convnets on bigger images (no need to slide windows but unroll network over space as needed to re-use computation).

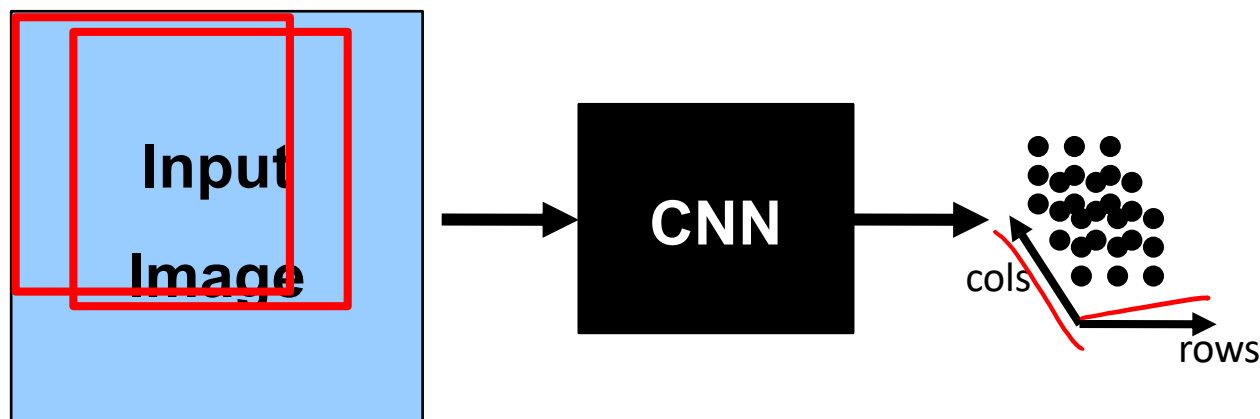


Viewing fully connected layers as convolutional layers enables efficient use of convnets on bigger images (no need to slide windows but unroll network over space as needed to re-use computation).

TRAINING TIME



TEST TIME



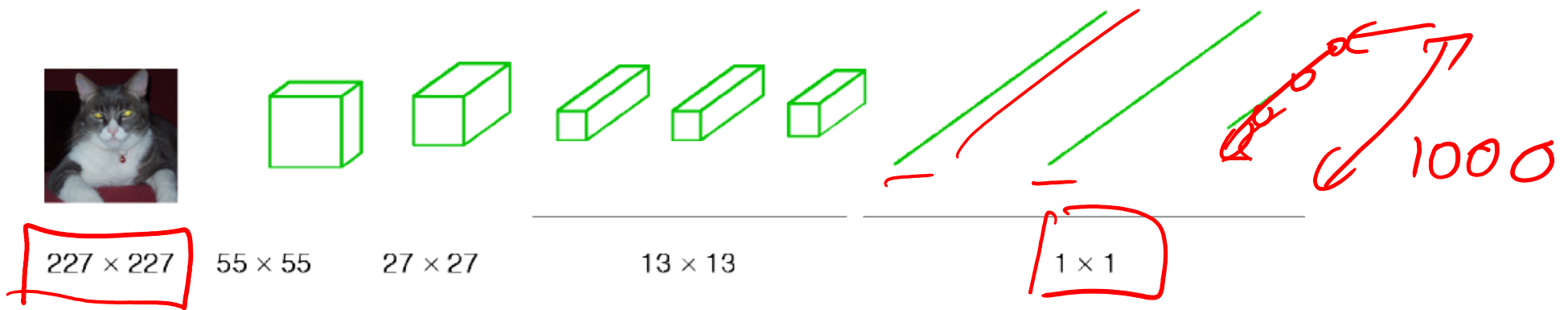
CNNs work on any image size!

Unrolling is order of magnitudes more efficient than sliding windows!

Training time

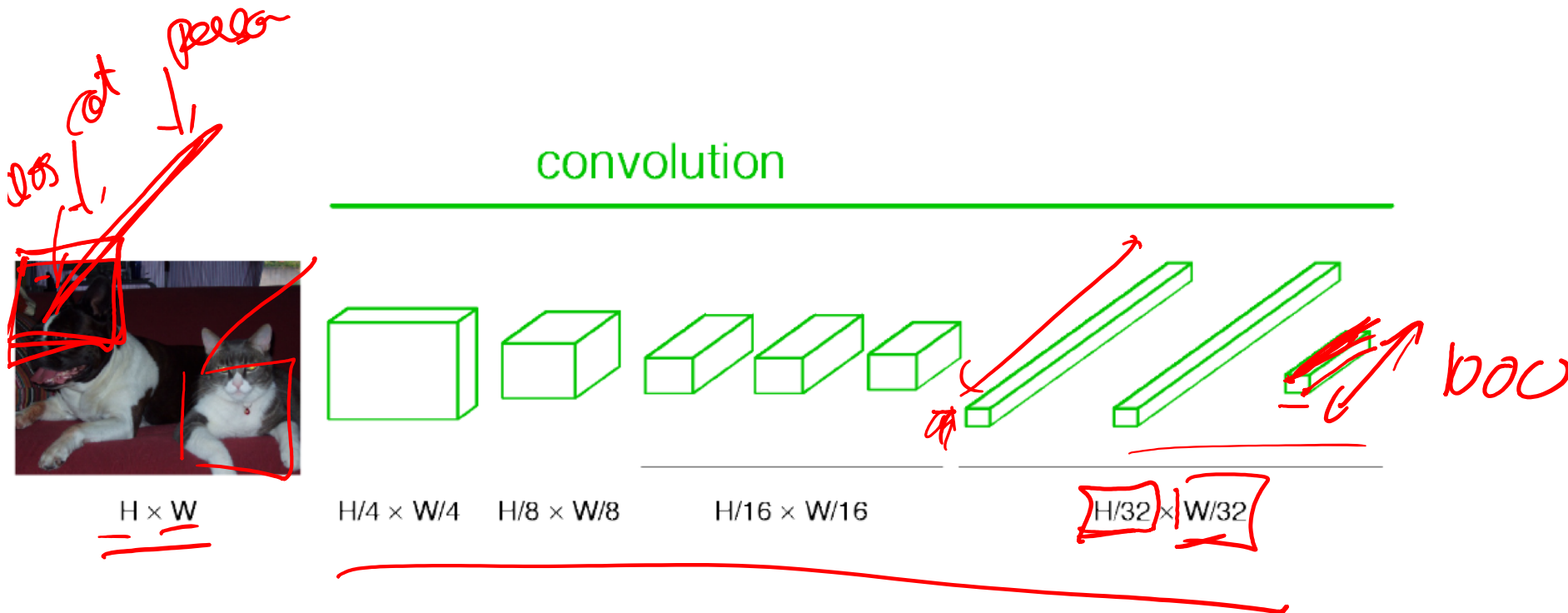
- Fixed-size images

convolution



Testing time

- Can run on an image of any size!

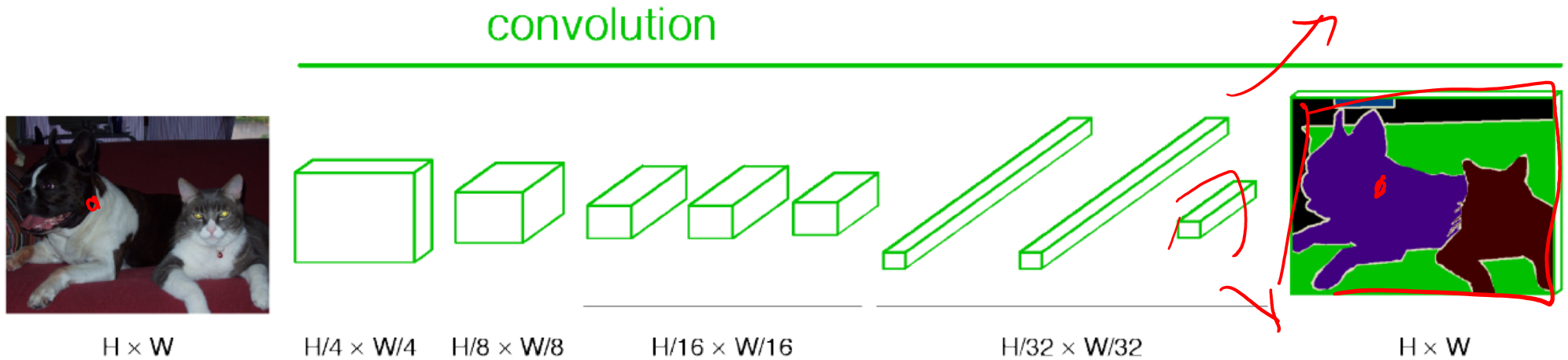


Benefit of this thinking

- Mathematically elegant
- Efficiency
 - Can run network on arbitrary image
 - Without multiple crops

“Fully Convolutional” Networks

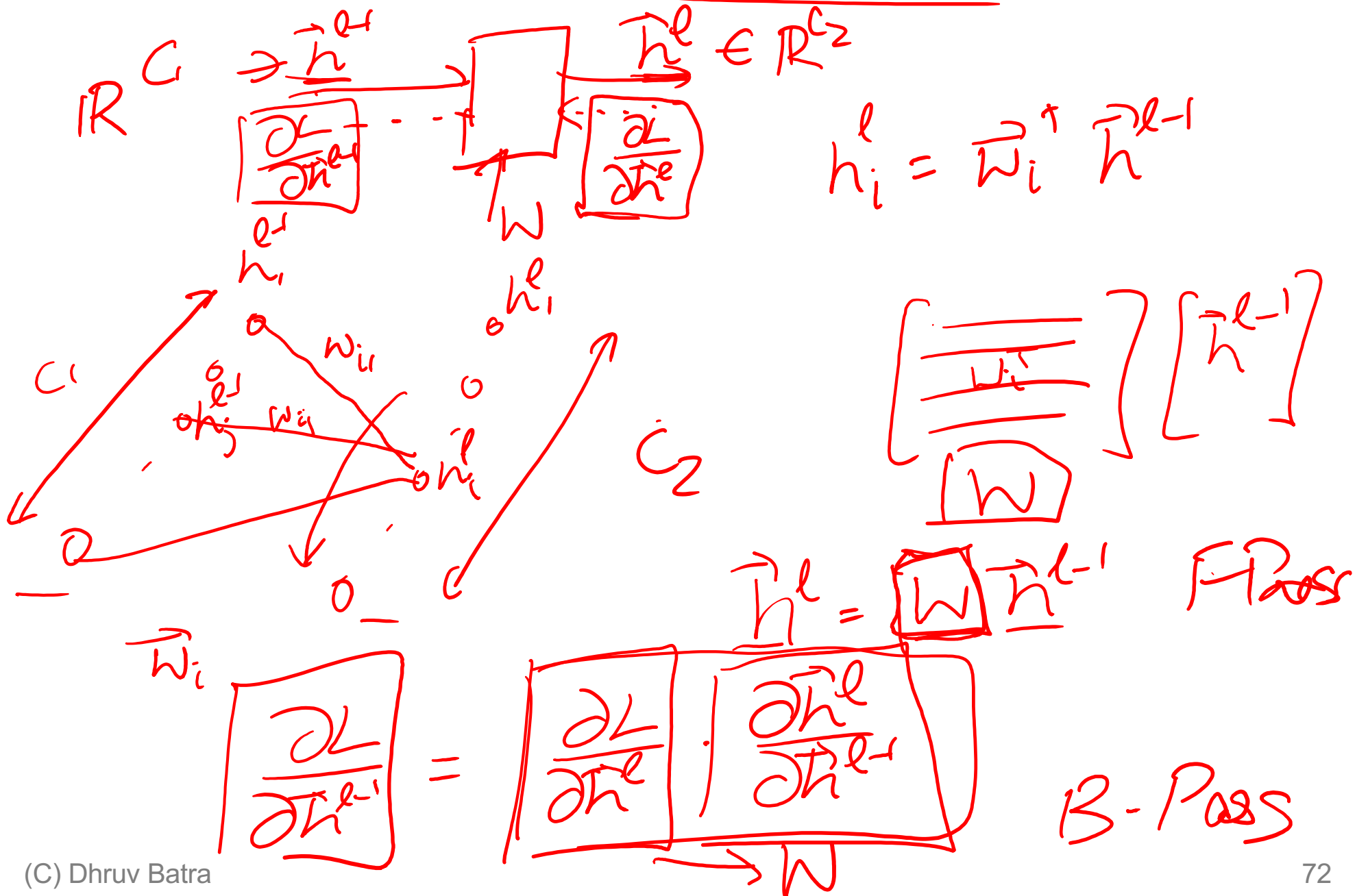
- Up-sample to get segmentation maps



Plan for Today

- Convolutional Neural Networks
 - 1x1 convolutions
 - Pooling layers
 - Fully-connected layers as convolutions
 - Backprop in conv layers

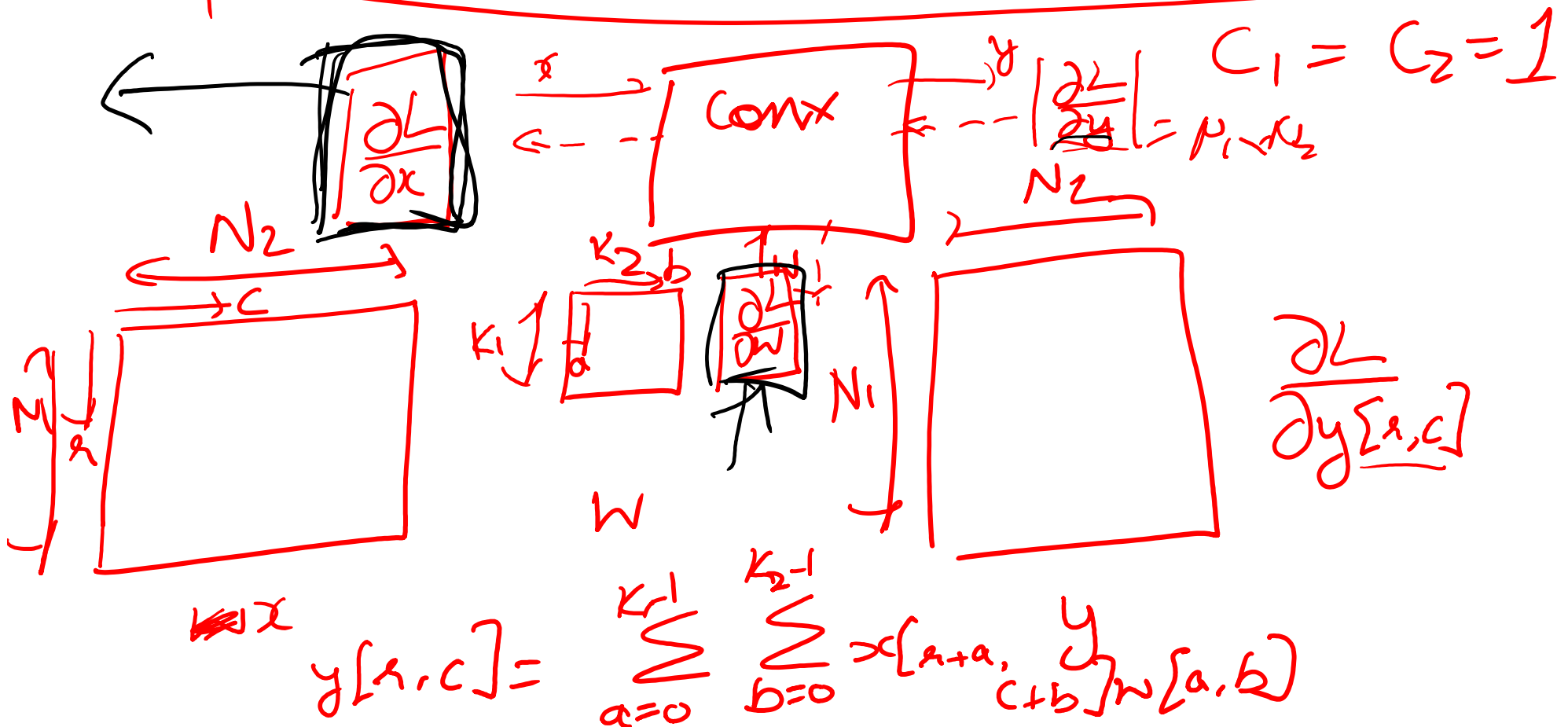
Jacobians of FC-Layer



Backprop in Convolutional Layers

- Notes

– https://www.cc.gatech.edu/classes/AY2018/cs7643_fall/slides/L6_cnns_backprop_notes.pdf



Backprop in Convolutional Layers

$$\frac{\partial L}{\partial w[a', b']} = \sum_{\text{pixels } p} \left[\frac{\partial L}{\partial y[p]} \frac{\partial y[p]}{\partial w[a', b']} \right]$$
 (input Jacobian)

$$y[r, c] = \sum_a \sum_b x[r+a, c+b] w[a, b]$$

$$= x[r+0, c+0] w[0, 0] + x[r+0, c+1] w[0, 1] + x[r+a, c+b] w[a, b]$$

$$\frac{\partial y[r, c]}{\partial w[a', b']} = x[r+a', c+b']$$

$$\frac{\partial L}{\partial w[a', b']} = \sum_{r=0}^{N_1} \sum_{c=0}^{N_2} \frac{\partial L}{\partial y[r, c]} x[r+a', c+b'] = x * \frac{\partial L}{\partial y}$$