

CS 7643: Deep Learning

Topics:

- Toeplitz matrices and convolutions = matrix-mult
- Dilated/a-trous convolutions
- Backprop in conv layers
- Transposed convolutions

Dhruv Batra
Georgia Tech

Administrativa

- HW1 extension
 - ~~09/22~~ 09/25
- HW2 + PS2 both coming out on ~~09/22~~ 09/25
- Note on class schedule coming up
 - Switching to paper reading starting next week.
 - <https://docs.google.com/spreadsheets/d/1uN31YcWAG6nhjvYPUVKMy3vHwW-h9MZCe8yKCqw0RsU/edit#gid=0>
- First review due: ~~Tue~~ ^{Mon} 09/26
- First student presentation due: Thr 09/28

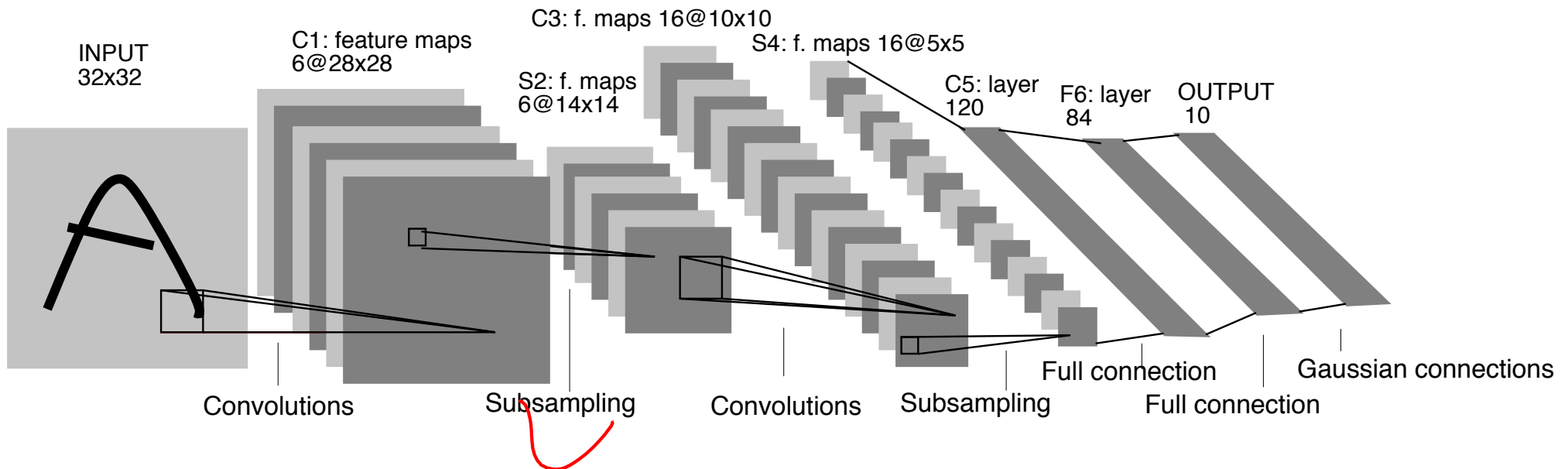
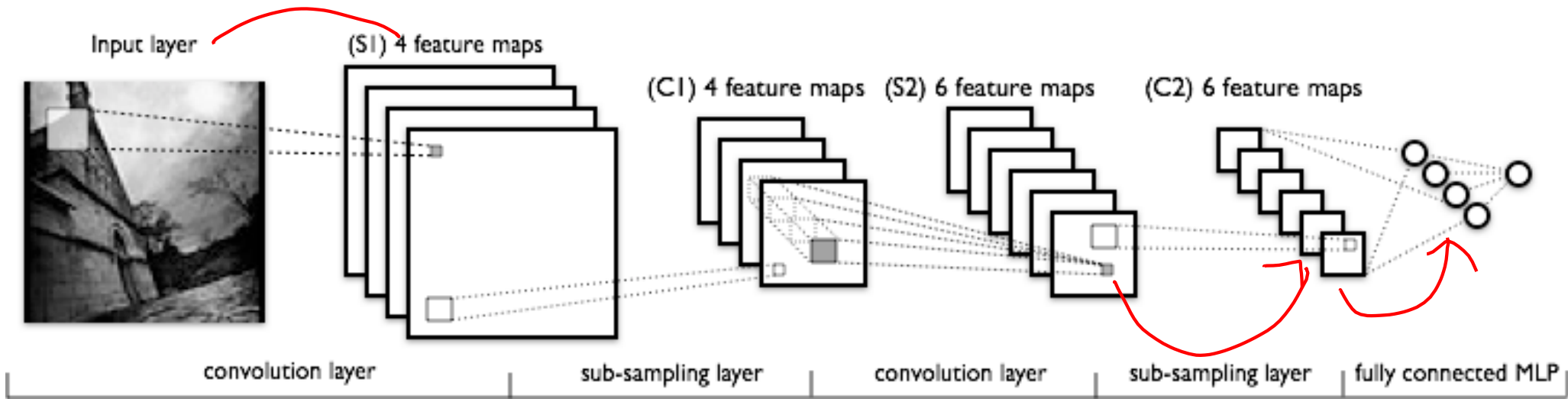
Recap of last time



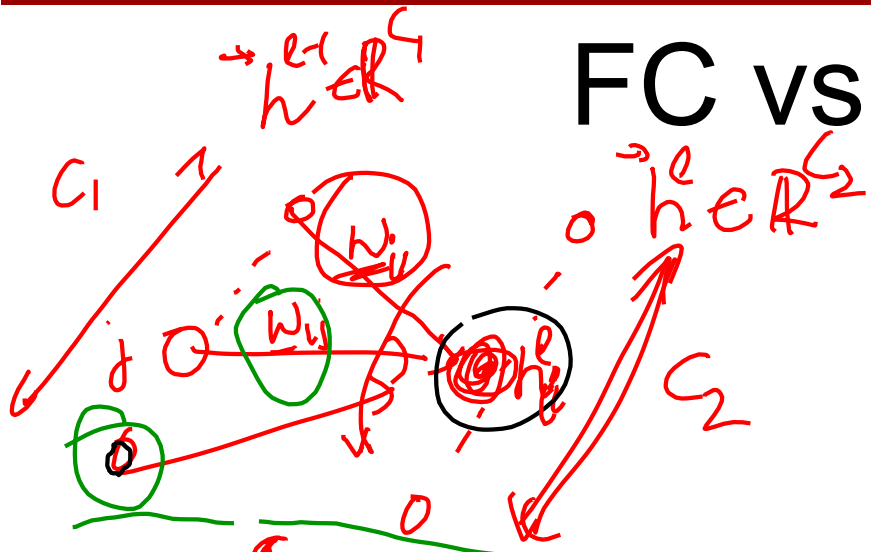
Convolutional Neural Networks

(without the brain stuff)

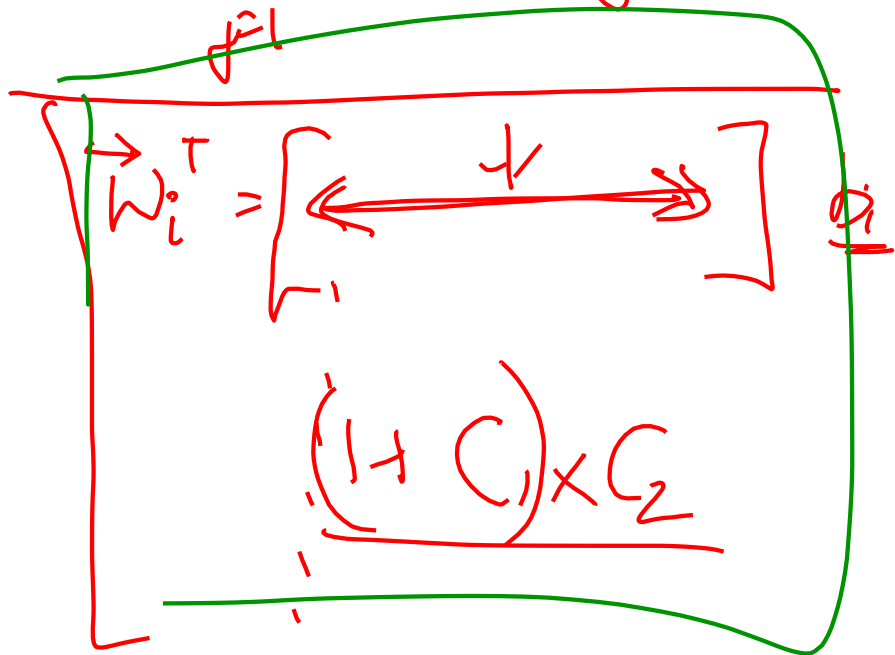
Convolutional Neural Networks



FC vs Conv Layer

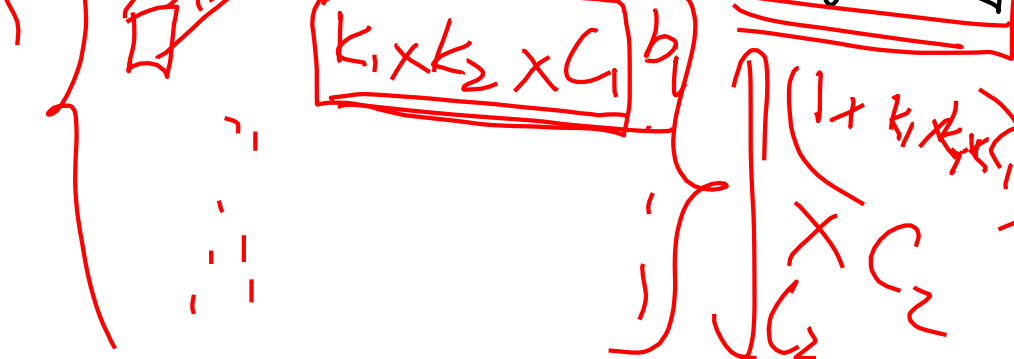


$$h_i^l = \sum_{j=1}^{C_1} h_j^{l-1} w_{ij} + b_i$$

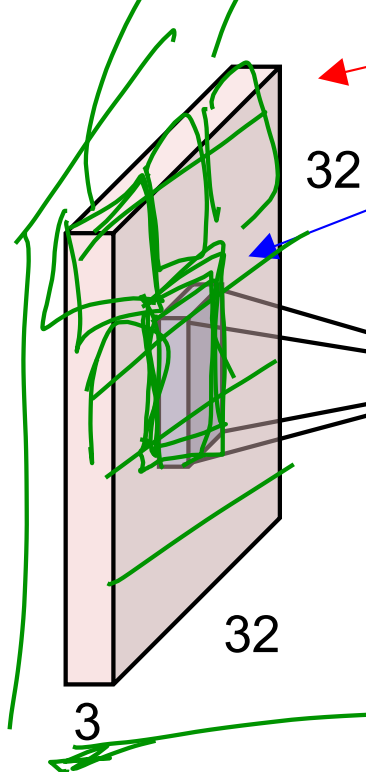


$$h_i^l = \sum_{j=1}^{C_1} h_j^{l-1} w_{ij} + b_i$$

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



Convolution Layer

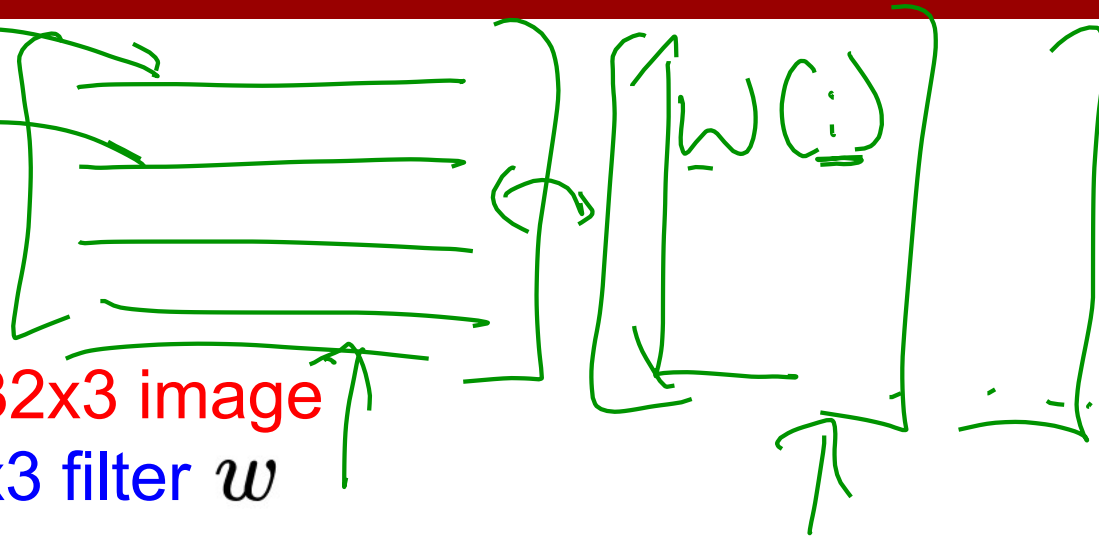


32x32x3 image
5x5x3 filter w

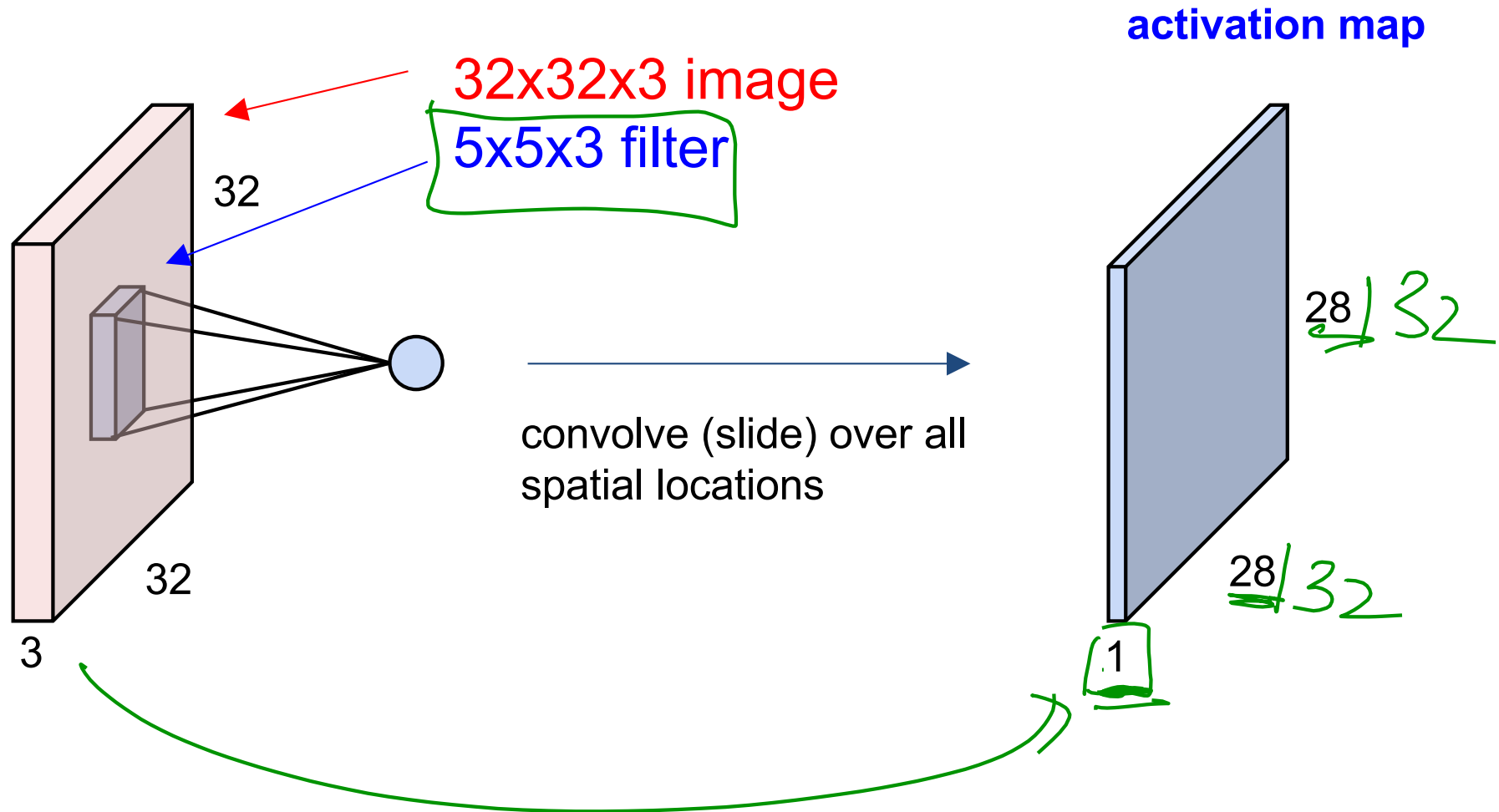
1 number:

the result of taking a dot product between the filter and a small 5x5x3 chunk of the image (i.e. $5*5*3 = 75$ -dimensional dot product + bias)

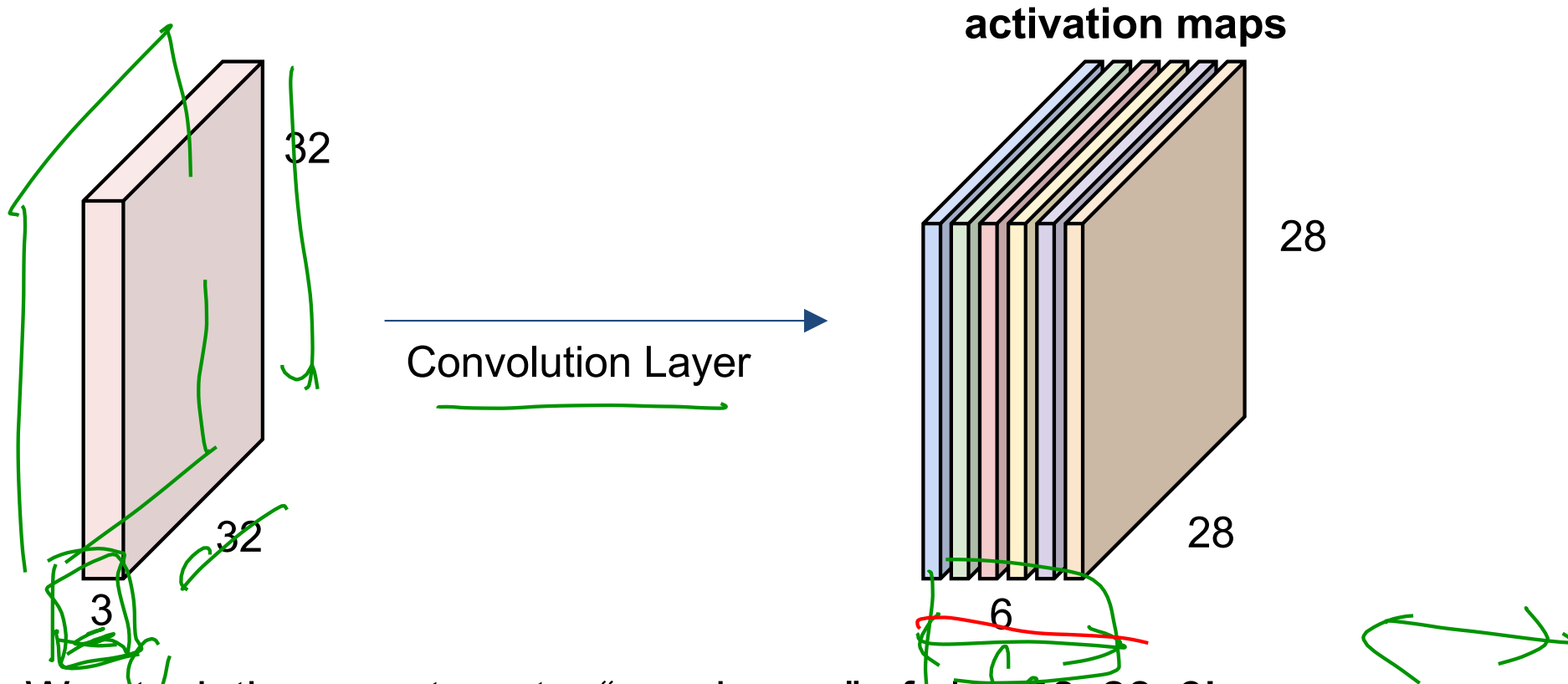
$$\underline{w^T x} + \underline{b}$$



Convolution Layer

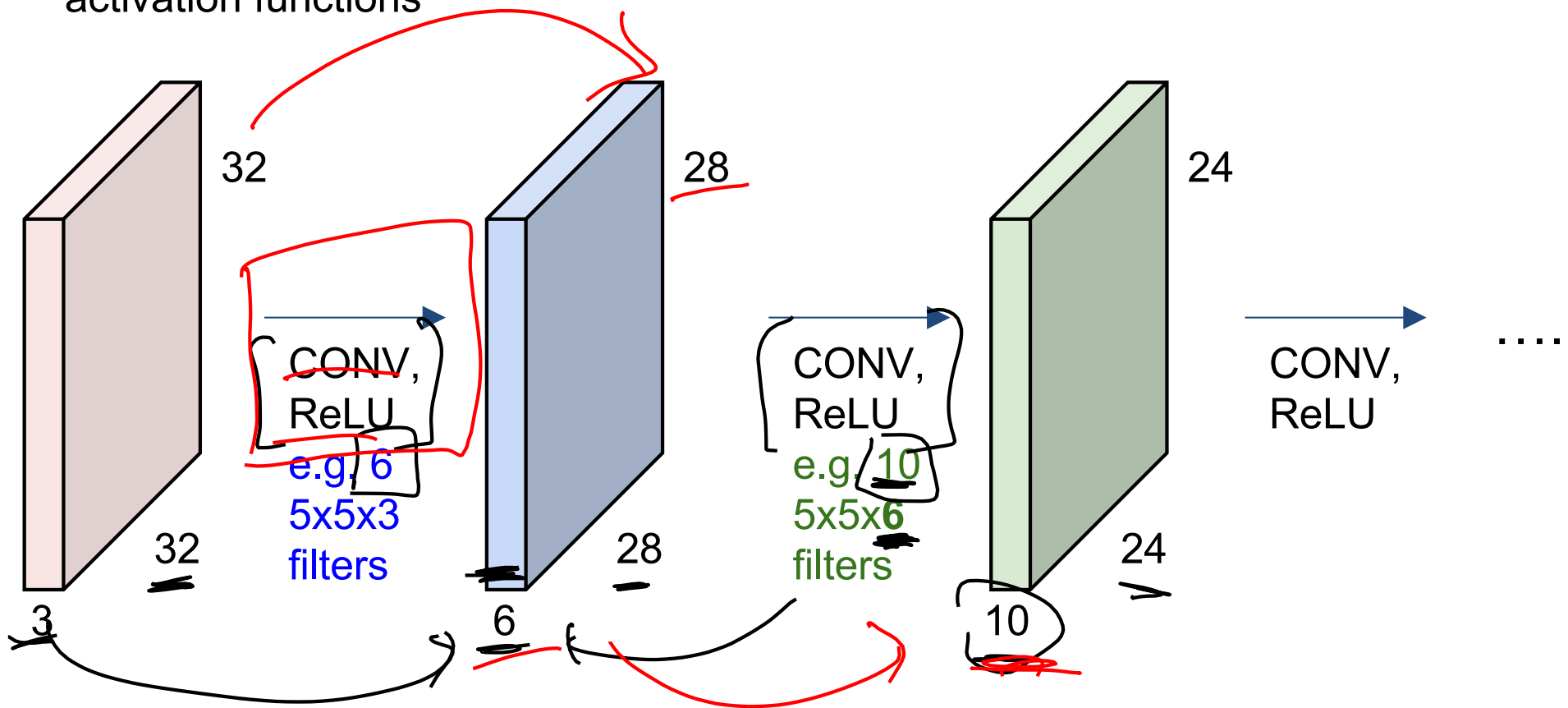


For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

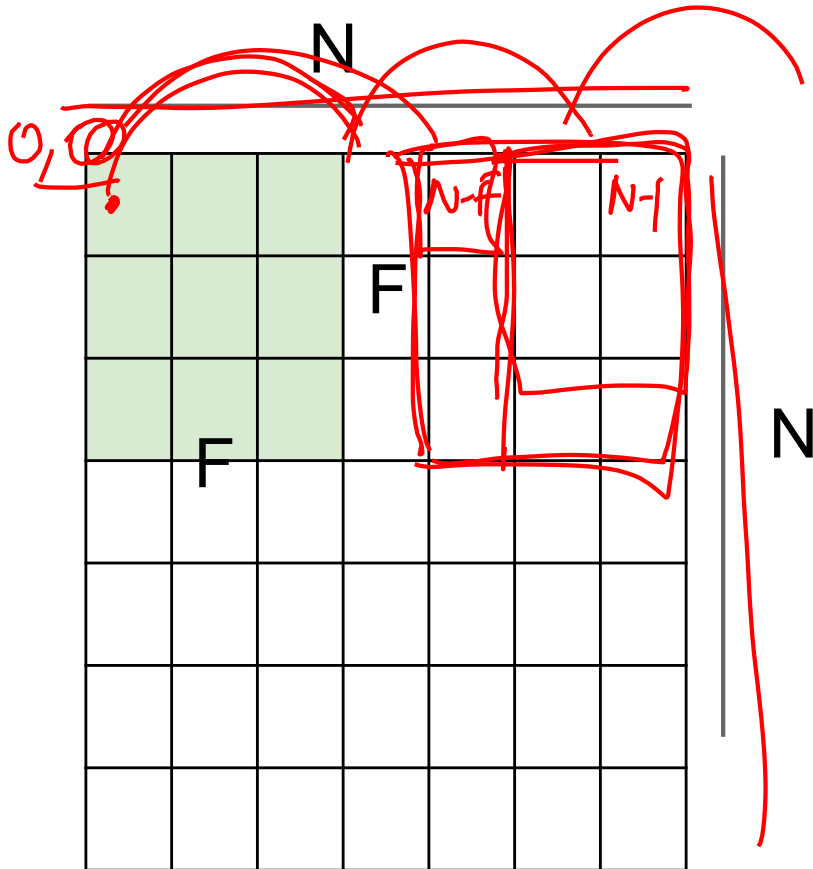


We stack these up to get a “new image” of size 28x28x6!

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



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



Output size:

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

+ 2p

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

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

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

stride 3 $\Rightarrow (7 - 3) / 3 + 1 = 2.33 \therefore \backslash$

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)

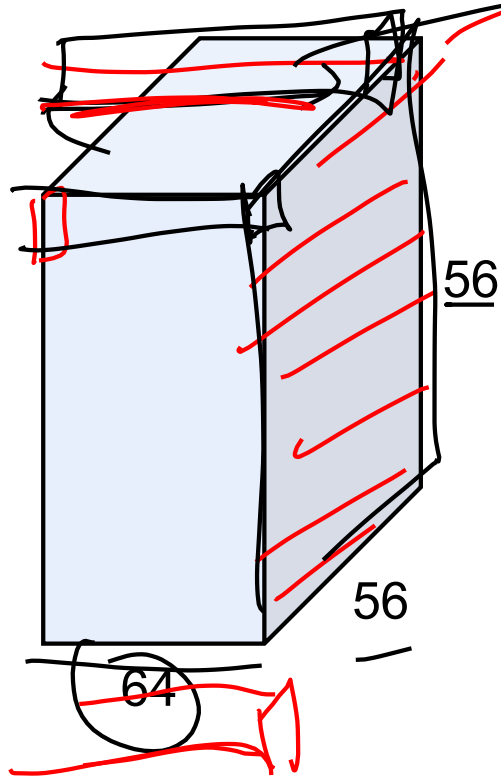
e.g. $F = 3 \Rightarrow$ zero pad with 1

$F = 5 \Rightarrow$ zero pad with 2

$F = 7 \Rightarrow$ zero pad with 3

(btw, 1x1 convolution layers make perfect sense)

$$\left[1 \times 1 \times C_1 \right] \times C_2$$



1x1 CONV
with 32 filters

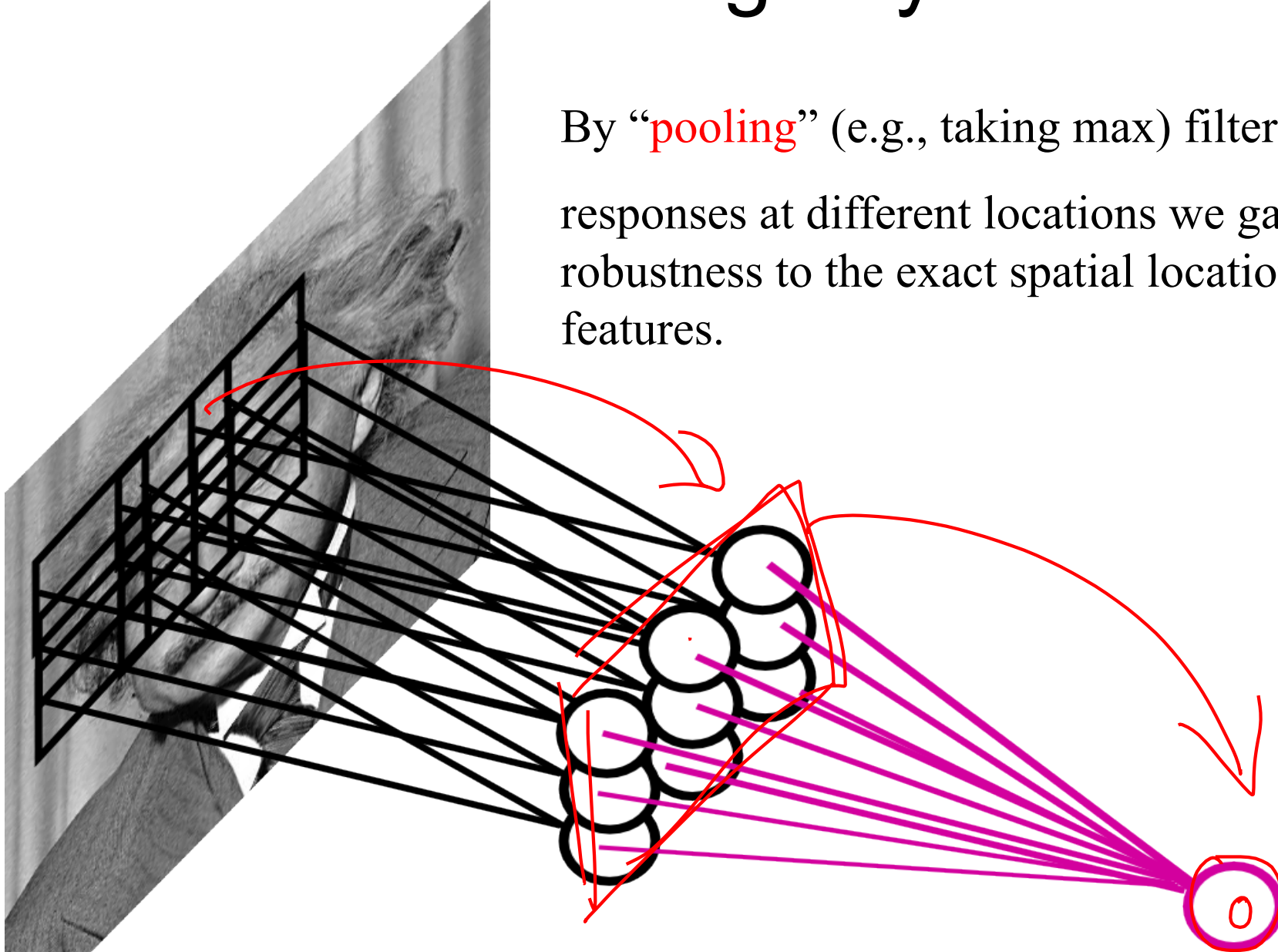
(each filter has size
1x1x64, and performs a
64-dimensional dot
product)



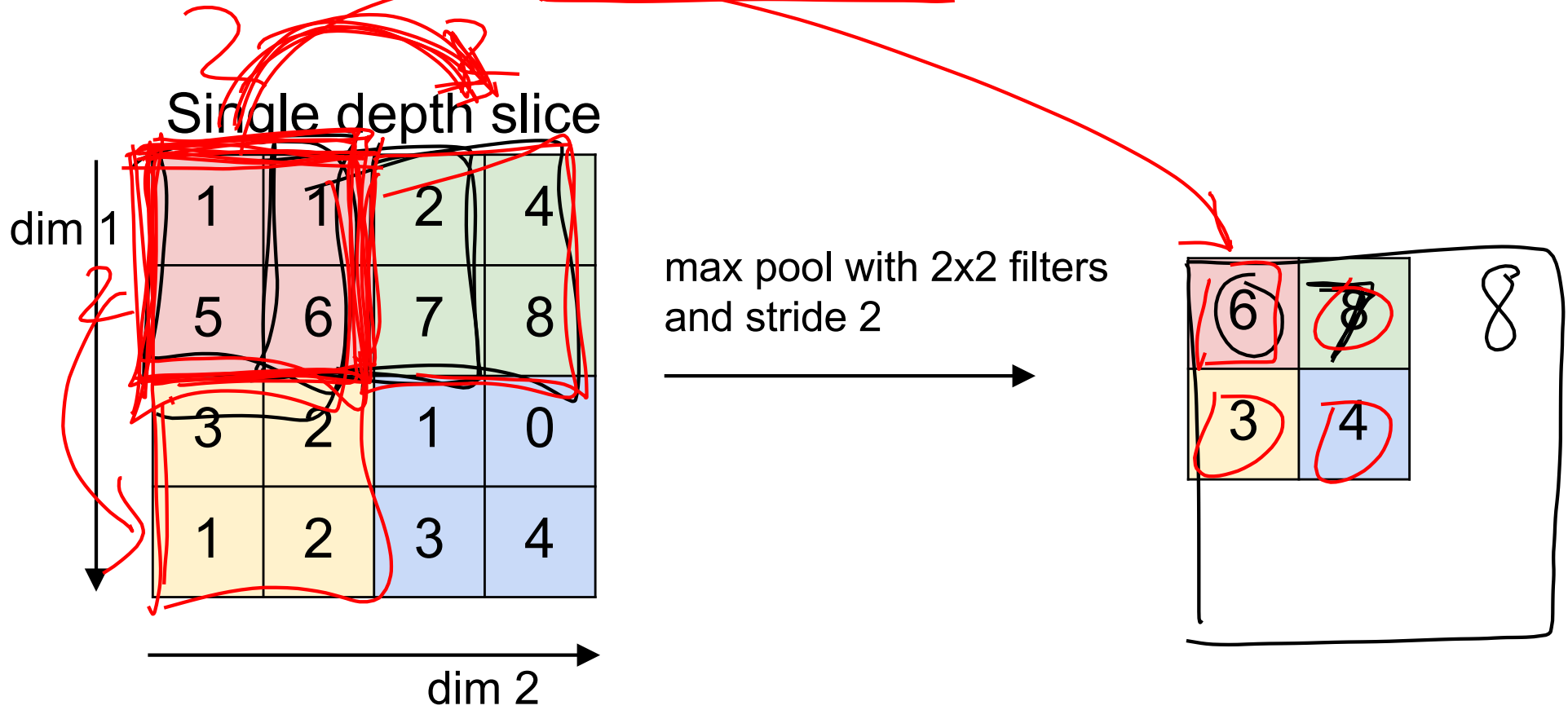
ReLU

Pooling Layer

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



MAX POOLING



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:

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

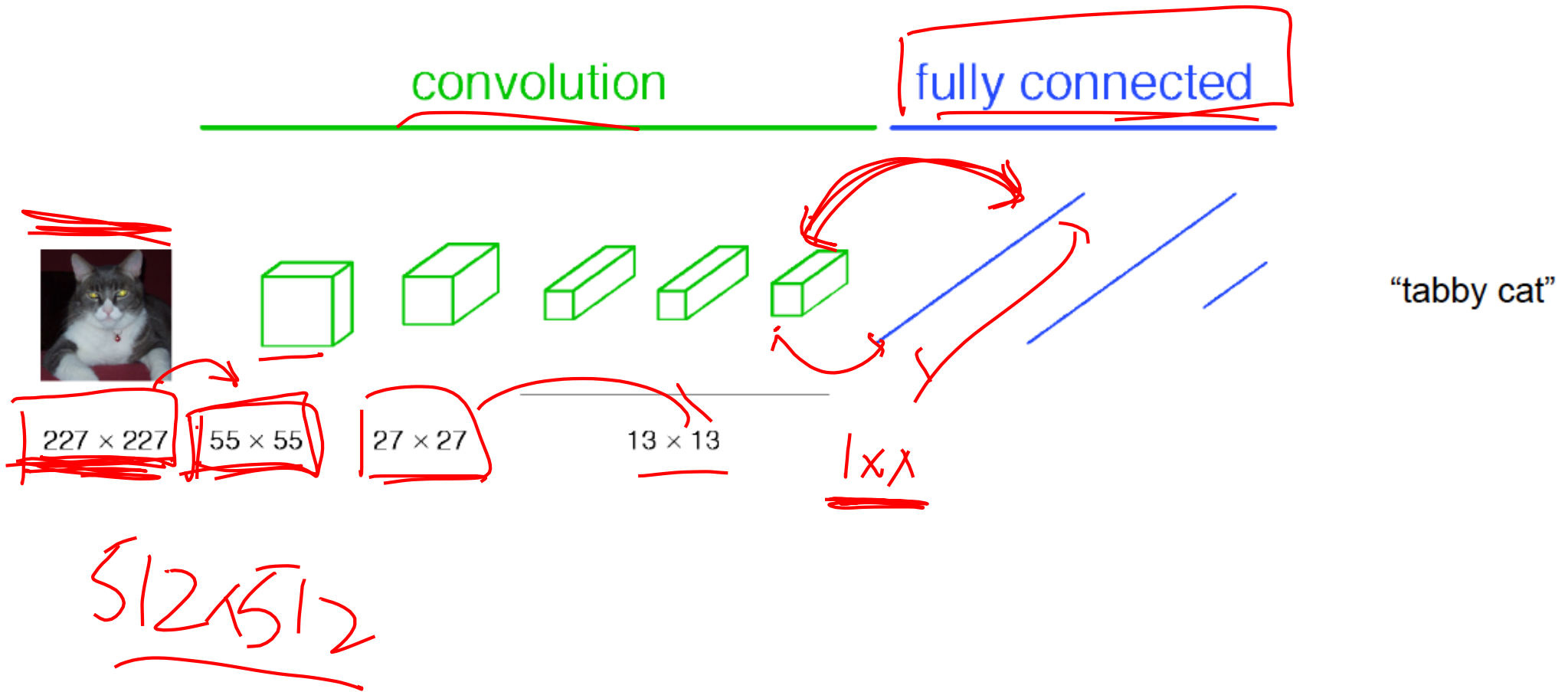
L2-pooling:

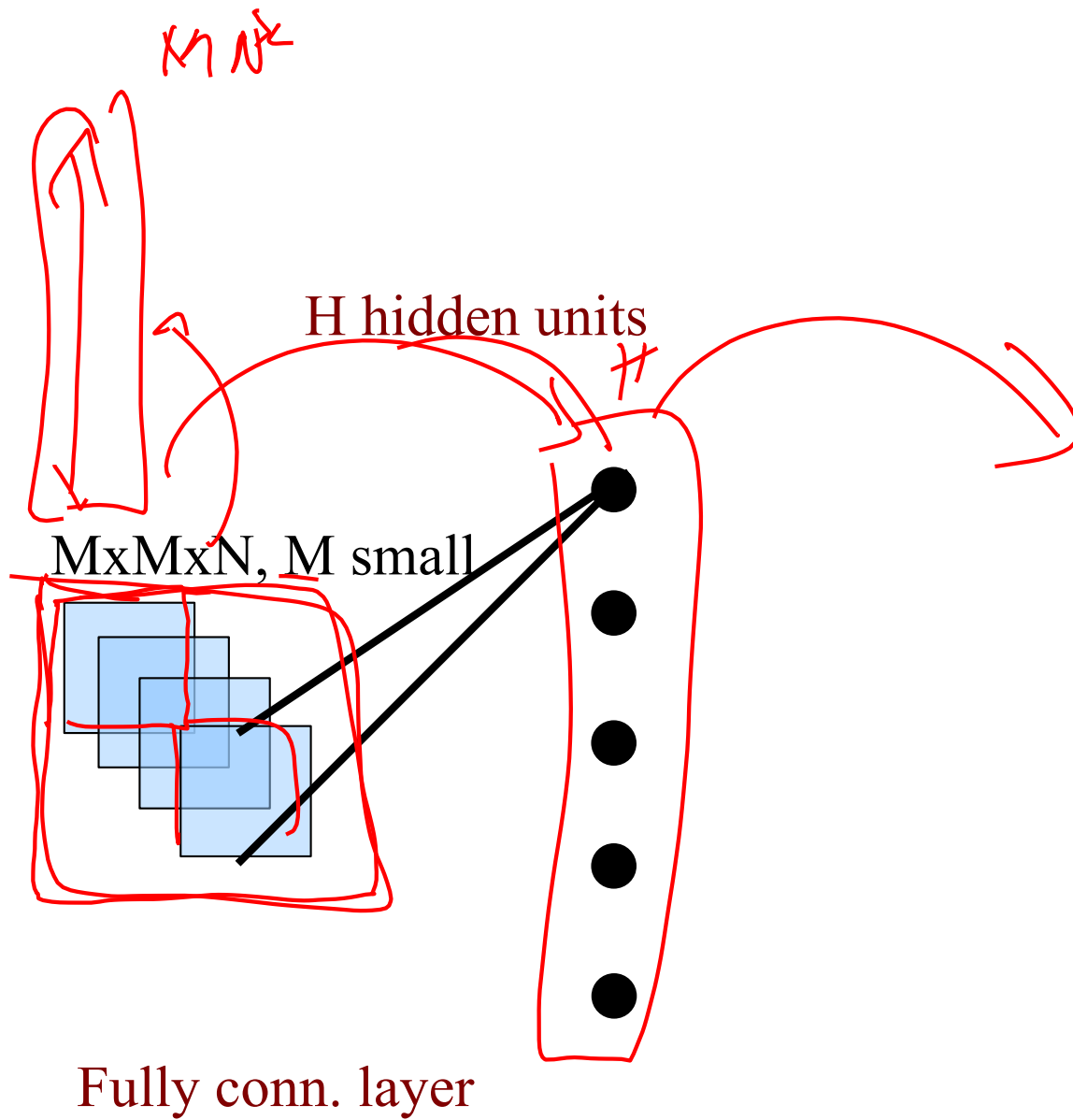
$$h_i^n(r, c) = \sqrt{\sum_{\bar{r} \in N(r), \bar{c} \in N(c)} h_i^{n-1}(\bar{r}, \bar{c})^2}$$

L2-pooling over features:

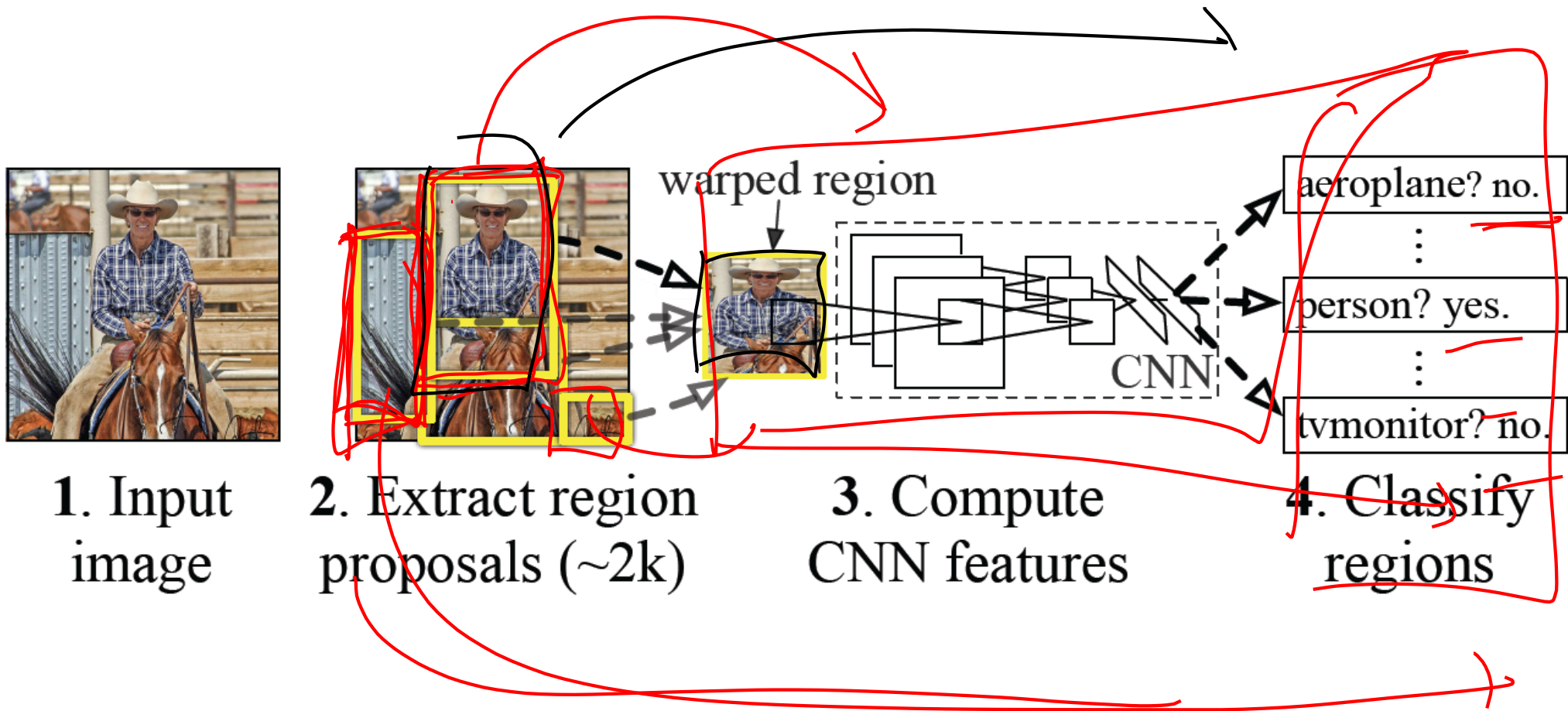
$$h_i^n(r, c) = \sqrt{\sum_{j \in N(i)} h_i^{n-1}(r, c)^2}$$

Classical View





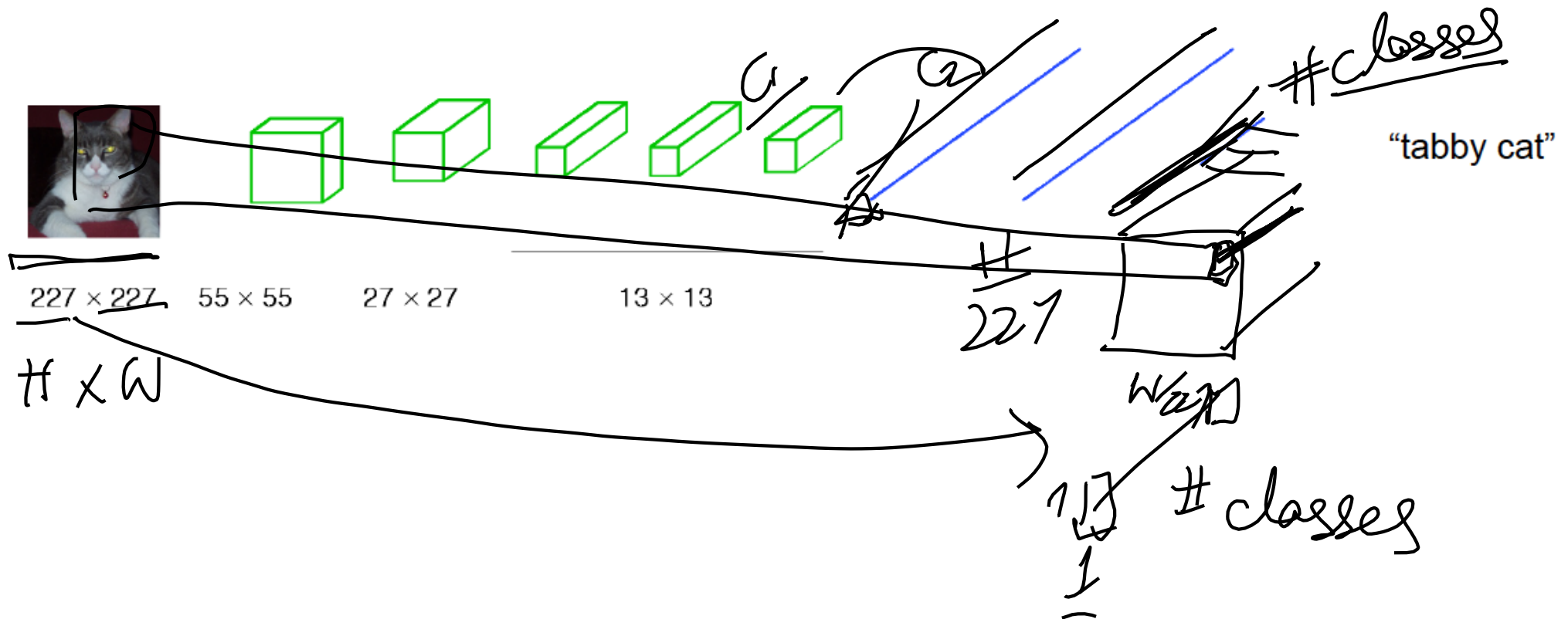
Classical View = Inefficient



Classical View

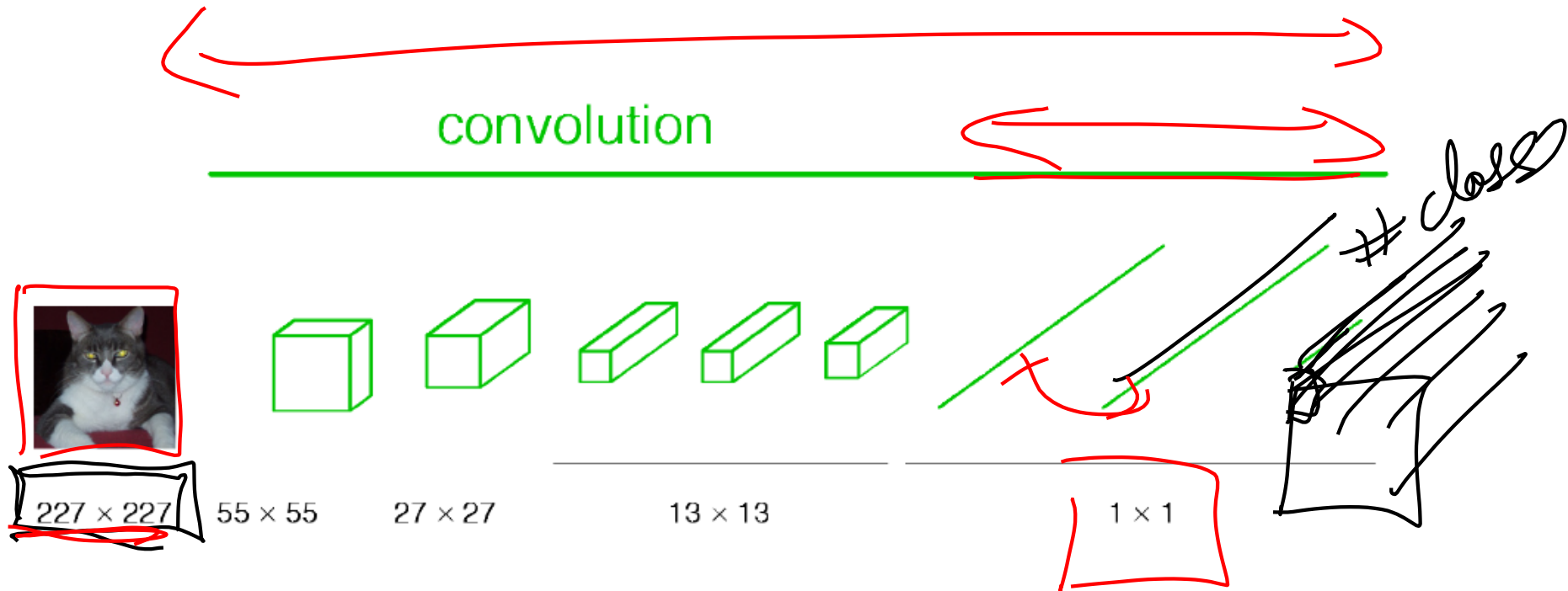
convolution

fully connected



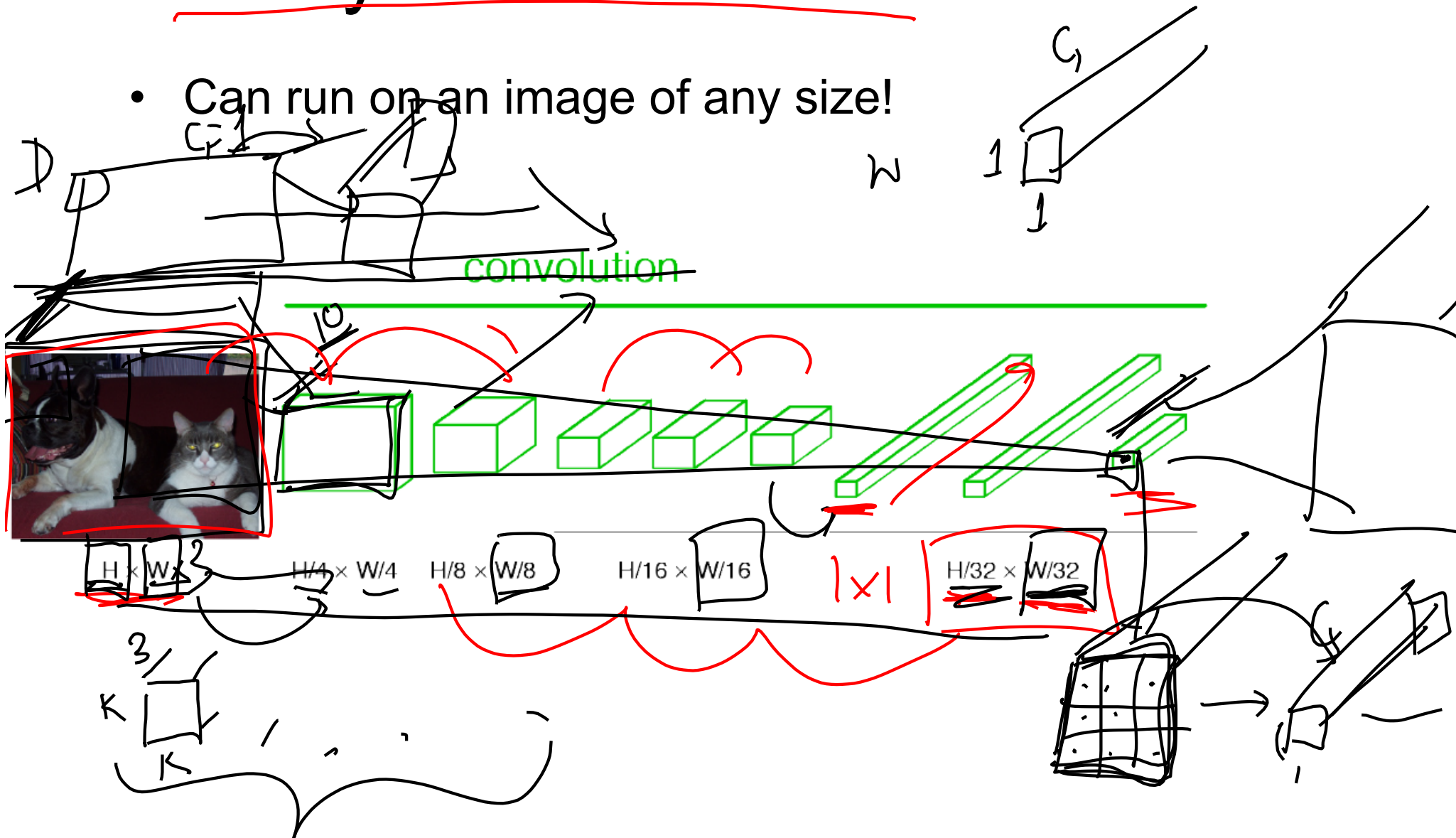
Re-interpretation

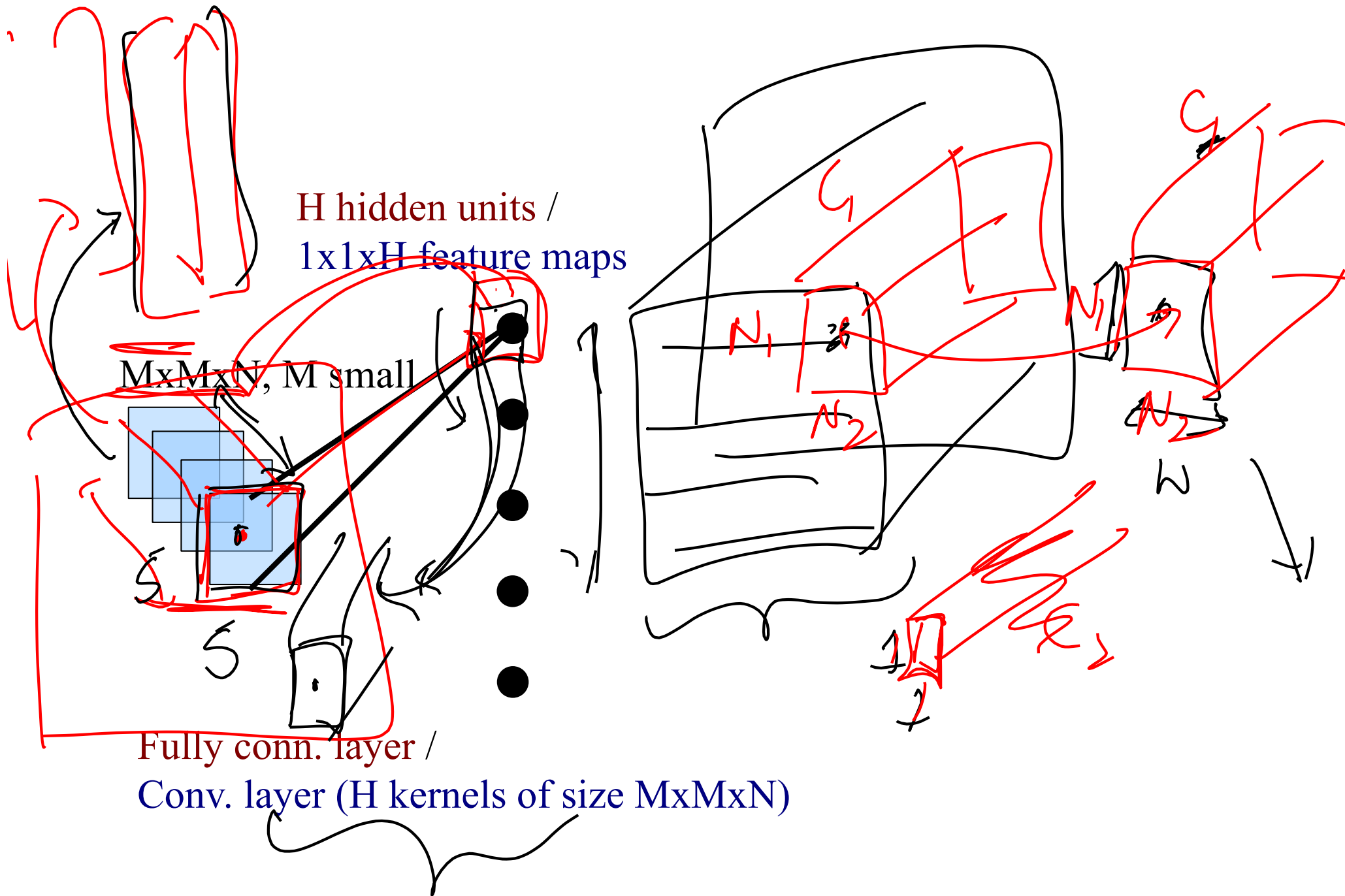
- Just squint a little!

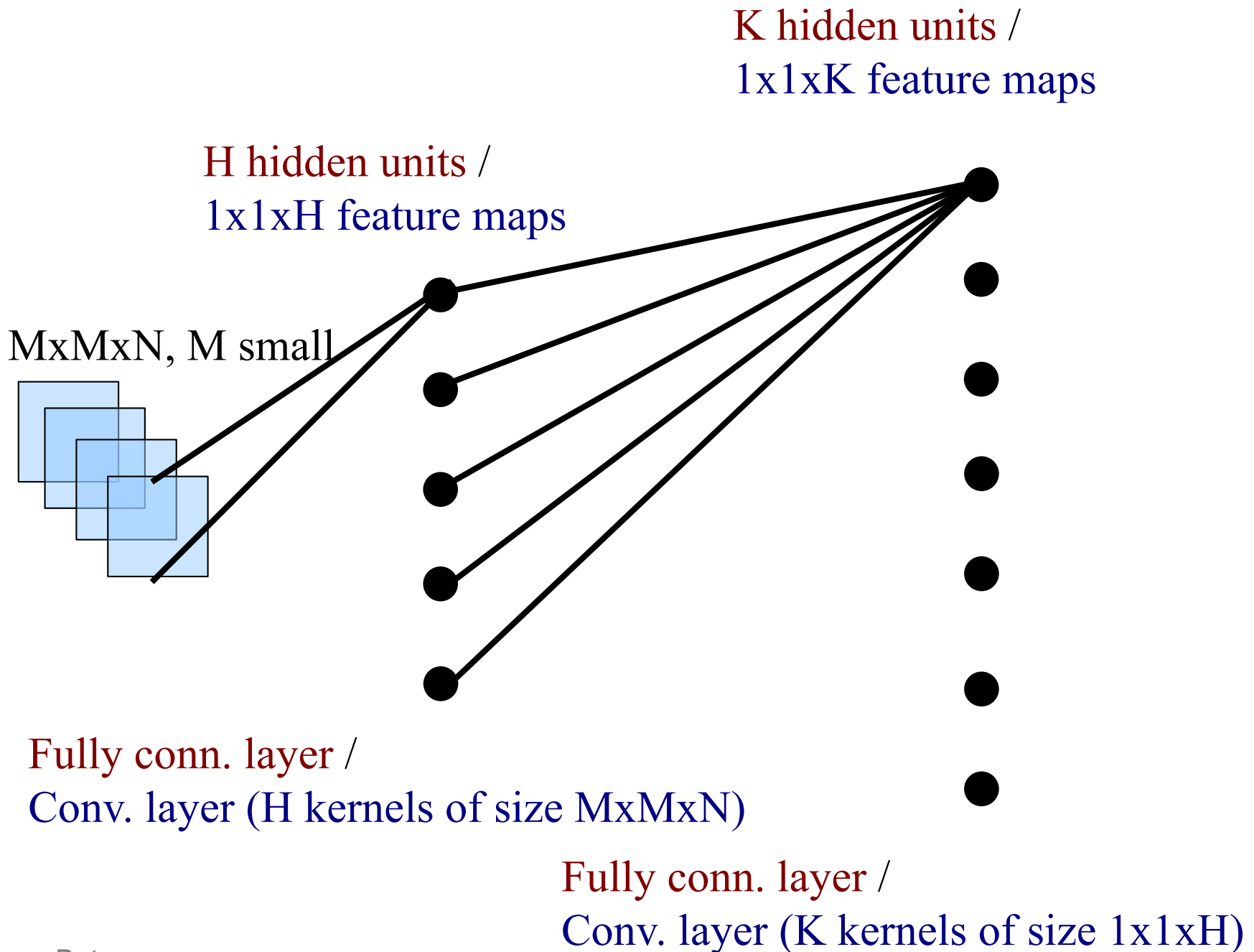


“Fully Convolutional” Networks

- Can run on an image of any size!

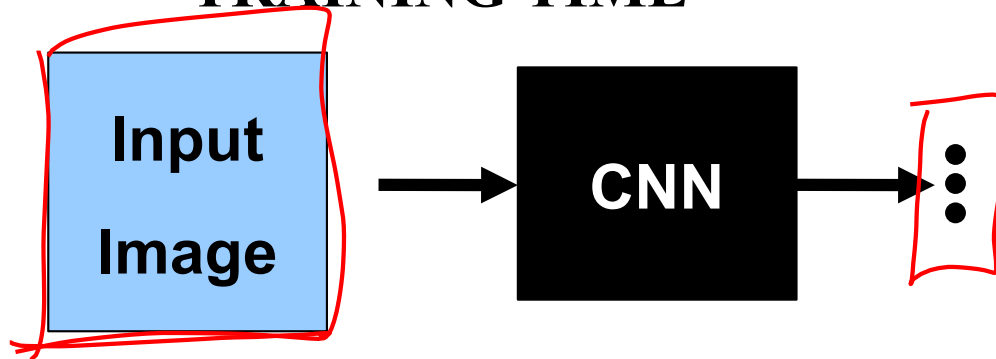




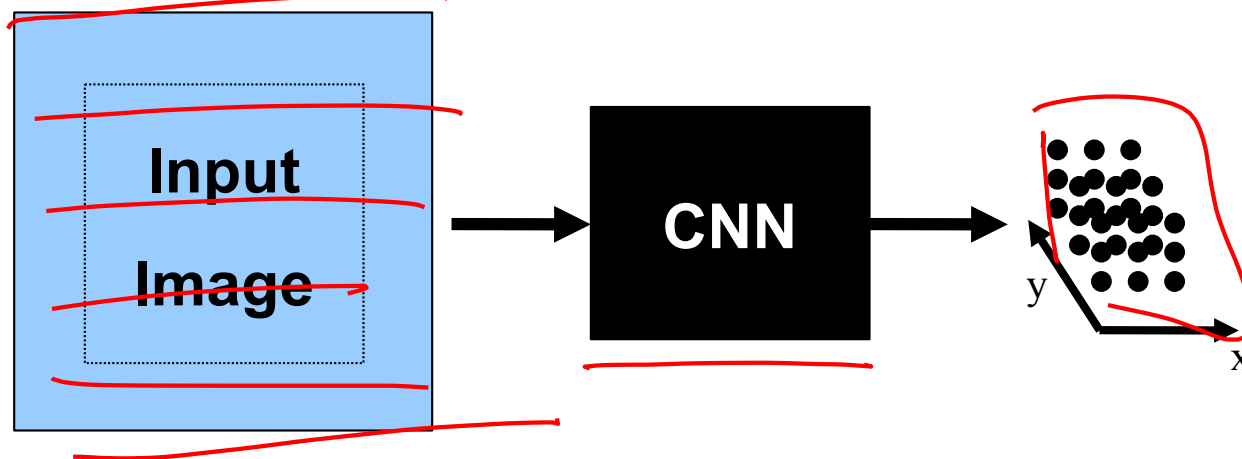


Viewing fully connected layers as convolutional layers enables efficient use of convnets on bigger images (no need to slide windows but unroll network over space as needed to re-use computation).

TRAINING TIME

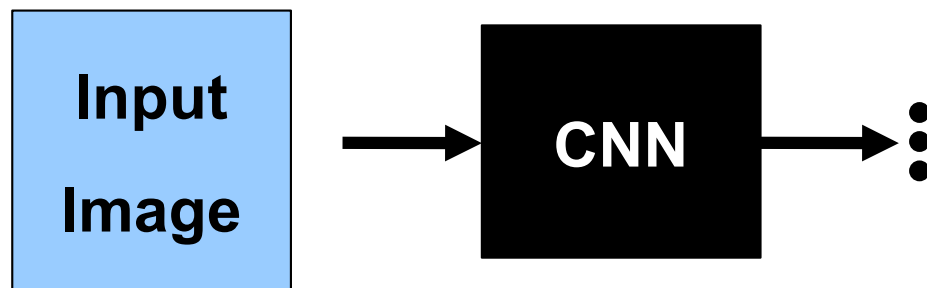


TEST TIME

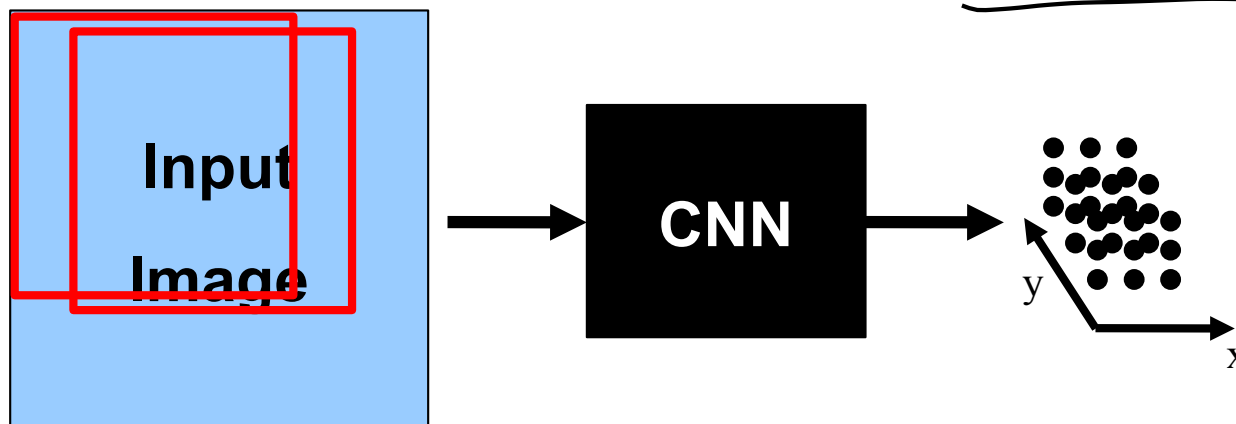


Viewing fully connected layers as convolutional layers enables efficient use of convnets on bigger images (no need to slide windows but unroll network over space as needed to re-use computation).

TRAINING TIME



TEST TIME



CNNs work on any image size!

Unrolling is order of magnitudes more efficient than sliding windows!

Benefit of this thinking

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

Summary

- ConvNets stack CONV, POOL, FC layers
- Trend towards smaller filters and deeper architectures
- Trend towards getting rid of POOL/FC layers (just CONV)
- Typical architectures look like

[CONV-RELU]*N-POOL?]*M-(FC-RELU)*K SOFTMAX

where N is usually up to ~5, M is large, $0 \leq K \leq 2$.

- but recent advances such as **ResNet/GoogLeNet** challenge this paradigm

Plan for Today

- Convolutional Neural Networks
 - Toeplitz matrices and convolutions = matrix-mult
 - Dilated/a-trous convolutions
 - Backprop in conv layers
 - Transposed convolutions

Toeplitz Matrix

- Diagonals are constants

$$\begin{bmatrix} a & b & c & d & e \\ f & a & b & c & d \\ g & f & a & b & c \\ h & g & f & a & b \\ i & h & g & f & a \end{bmatrix}.$$

- $A_{ij} = a_{i-j}$

$$A = \begin{bmatrix} a_0 & a_{-1} & a_{-2} & \dots & \dots & a_{-n+1} \\ a_1 & a_0 & a_{-1} & \ddots & & \vdots \\ a_2 & a_1 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & a_{-1} & a_{-2} \\ \vdots & & \ddots & a_1 & a_0 & a_{-1} \\ a_{n-1} & \dots & \dots & a_2 & a_1 & a_0 \end{bmatrix}$$

Why do we care?

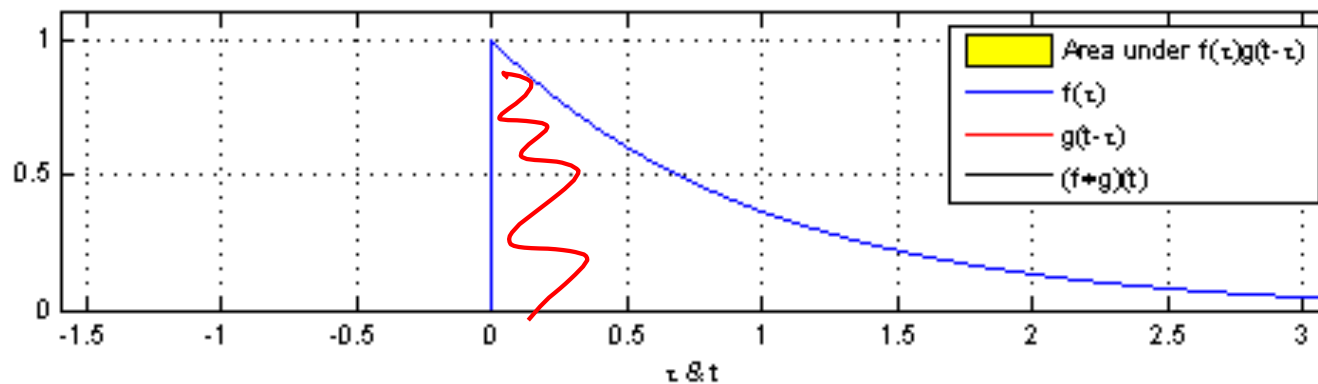
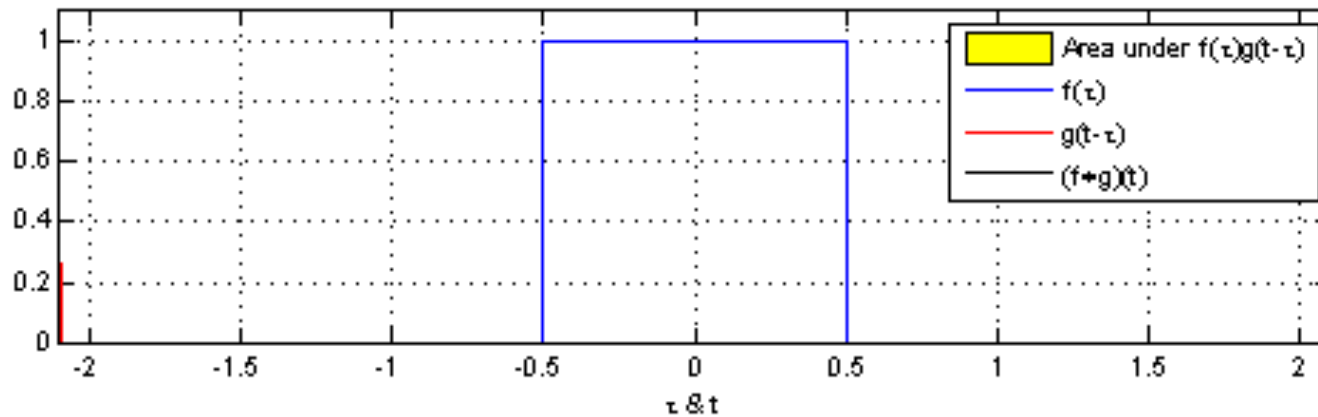
- (Discrete) Convolution = Matrix Multiplication
 - with Toeplitz Matrices

$[w_1 \dots w_k]$

$$\begin{bmatrix}
 \cancel{w_k} & 0 & \dots & 0 & 0 \\
 \cancel{w_{k-1}} & \cancel{w_k} & \dots & 0 & 0 \\
 w_{k-2} & w_{k-1} & w_k & \dots & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots \\
 w_1 & w_{k-2} & \dots & w_k & 0 \\
 \vdots & \vdots & \vdots & \vdots & \vdots \\
 0 & w_1 & \dots & w_{k-1} & w_k \\
 \vdots & \vdots & \vdots & \vdots & \vdots \\
 0 & 0 & \vdots & w_1 & w_2 \\
 0 & 0 & \vdots & 0 & w_1
 \end{bmatrix}$$

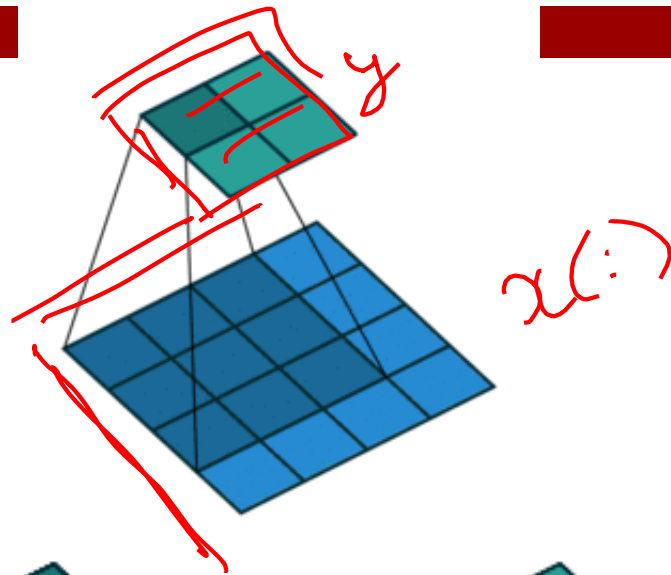
$\begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n \end{bmatrix}$

$y = w * x$

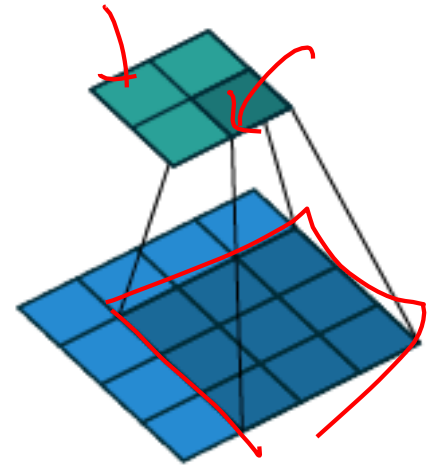
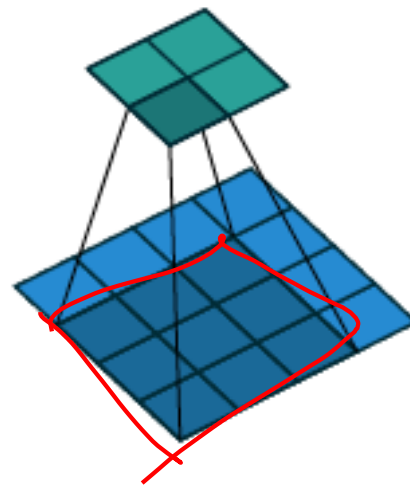
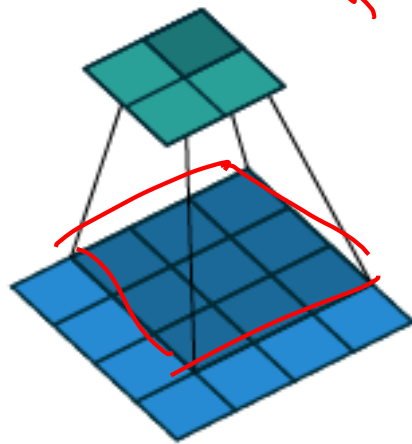
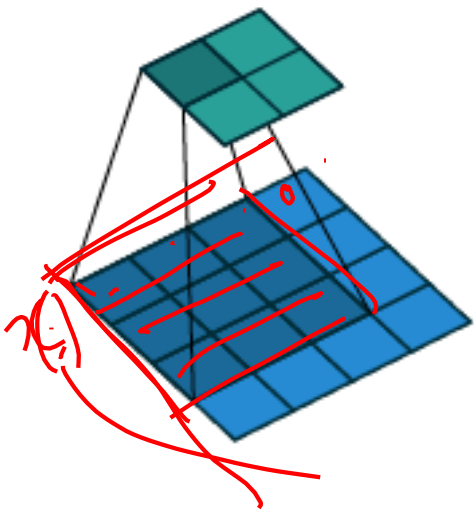


"Convolution of box signal with itself2" by Convolution_of_box_signal_with_itself.gif: Brian Ambergderivative work: Tinos (talk) - Convolution_of_box_signal_with_itself.gif. Licensed under CC BY-SA 3.0 via Commons - https://commons.wikimedia.org/wiki/File:Convolution_of_box_signal_with_itself2.gif#/media/File:Convolution_of_box_signal_wi th_itself2.gif

$$y[r,c] = \sum_w x(r) \cdot w(c)$$



W^T



4×1

$w_{0,0}$	$w_{0,1}$	$w_{0,2}$	0	$w_{1,0}$	$w_{1,1}$	$w_{1,2}$	0	$w_{2,0}$	$w_{2,1}$	$w_{2,2}$	0	0	0	0	0
0	$w_{0,0}$	$w_{0,1}$	$w_{0,2}$	0	$w_{1,0}$	$w_{1,1}$	$w_{1,2}$	0	$w_{2,0}$	$w_{2,1}$	$w_{2,2}$	0	0	0	0
0	0	0	0	$w_{0,0}$	$w_{0,1}$	$w_{0,2}$	0	$w_{1,0}$	$w_{1,1}$	$w_{1,2}$	0	$w_{2,0}$	$w_{2,1}$	$w_{2,2}$	0
0	0	0	0	0	$w_{0,0}$	$w_{0,1}$	$w_{0,2}$	0	$w_{1,0}$	$w_{1,1}$	$w_{1,2}$	0	$w_{2,0}$	$w_{2,1}$	$w_{2,2}$

$$y = Wx$$

4×16

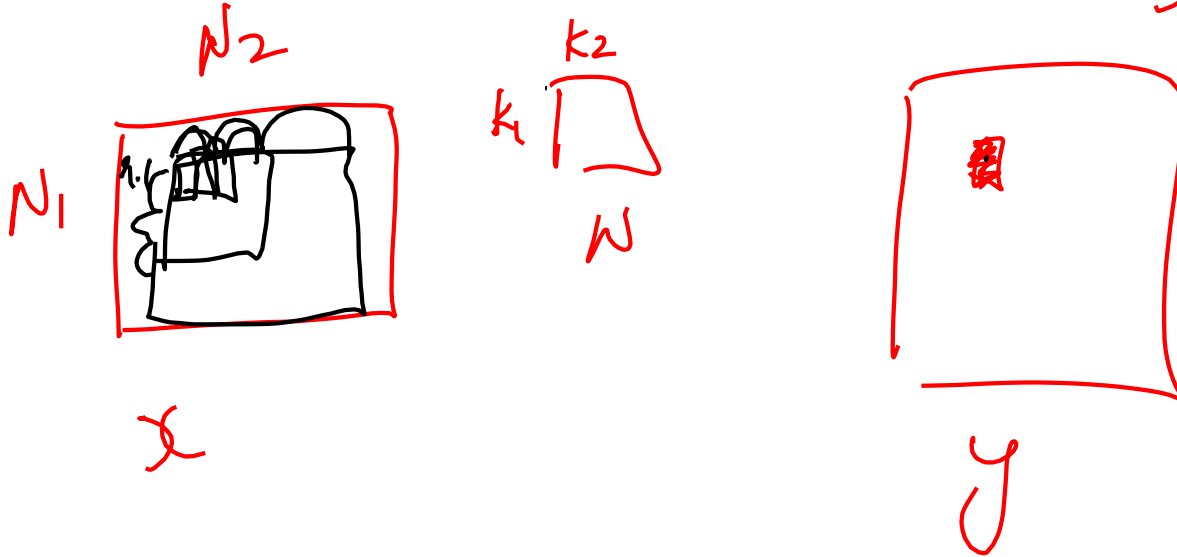
16×1

Plan for Today

- Convolutional Neural Networks
 - Toeplitz matrices and convolutions = matrix-mult
 - Dilated/a-trous convolutions
 - Backprop in conv layers
 - Transposed convolutions

Dilated Convolutions

$$y[r, c] = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x[r + \underbrace{a}_{5a}, c + \underbrace{b}_{5b}] w[a, b]$$



Dilated Convolutions

5x5

1	1	1
1	1	1
1	1	1



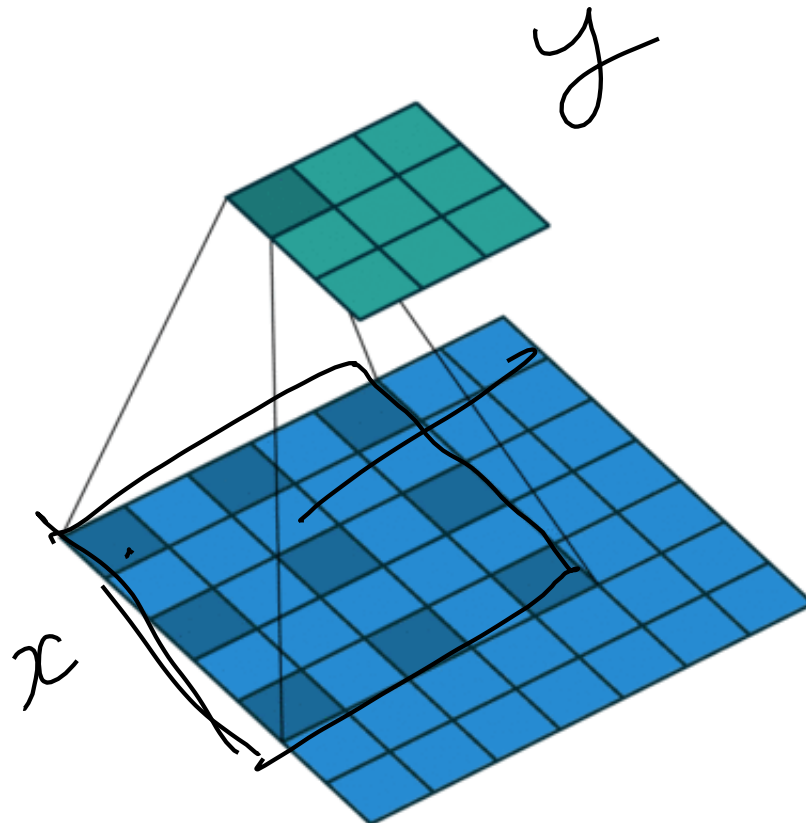
1	0	1	0	1

5x5

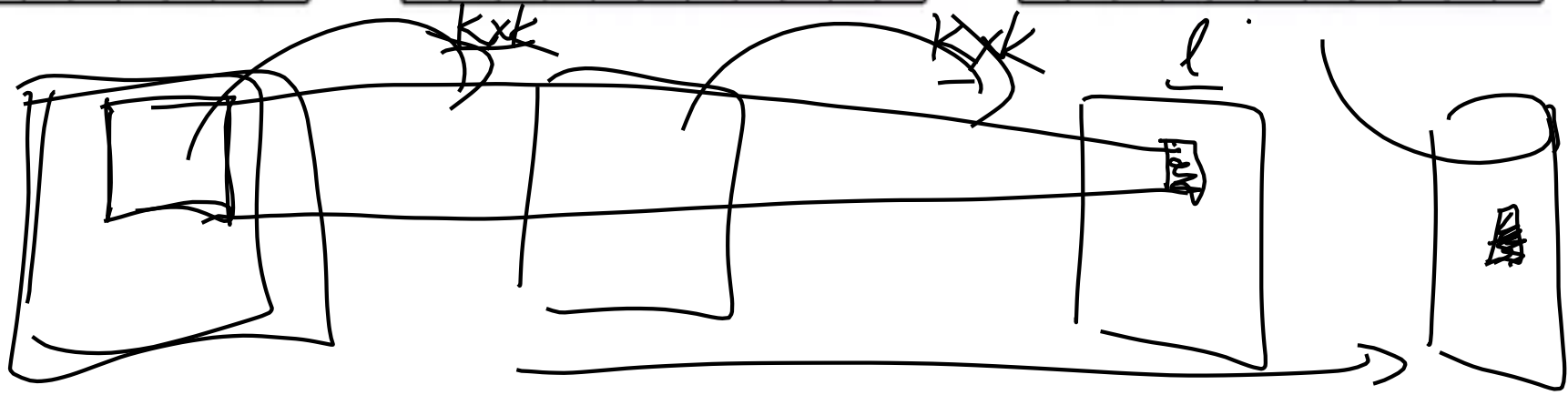
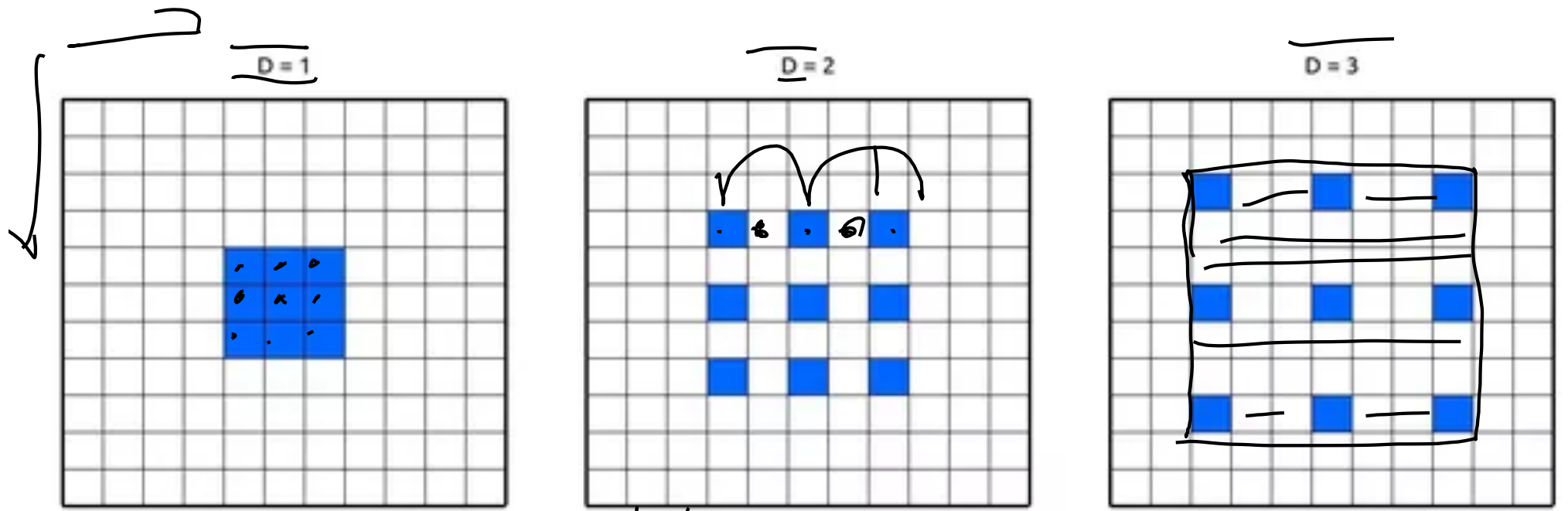
$d=2$

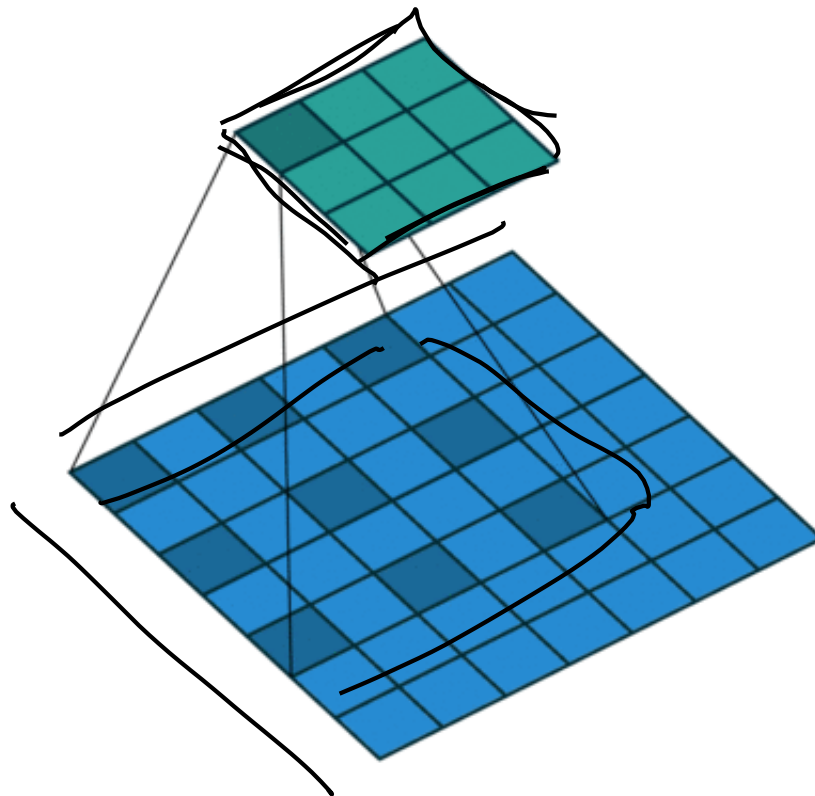
$(d-1)$

3x3



$$k \times k \rightarrow \left[k + \underline{\underline{(k-1)(d-1)}} \right] \times \left[\underline{\underline{3 \times 3}} \right]$$





(recall:)
 $(N - k) / \text{stride} + 1$

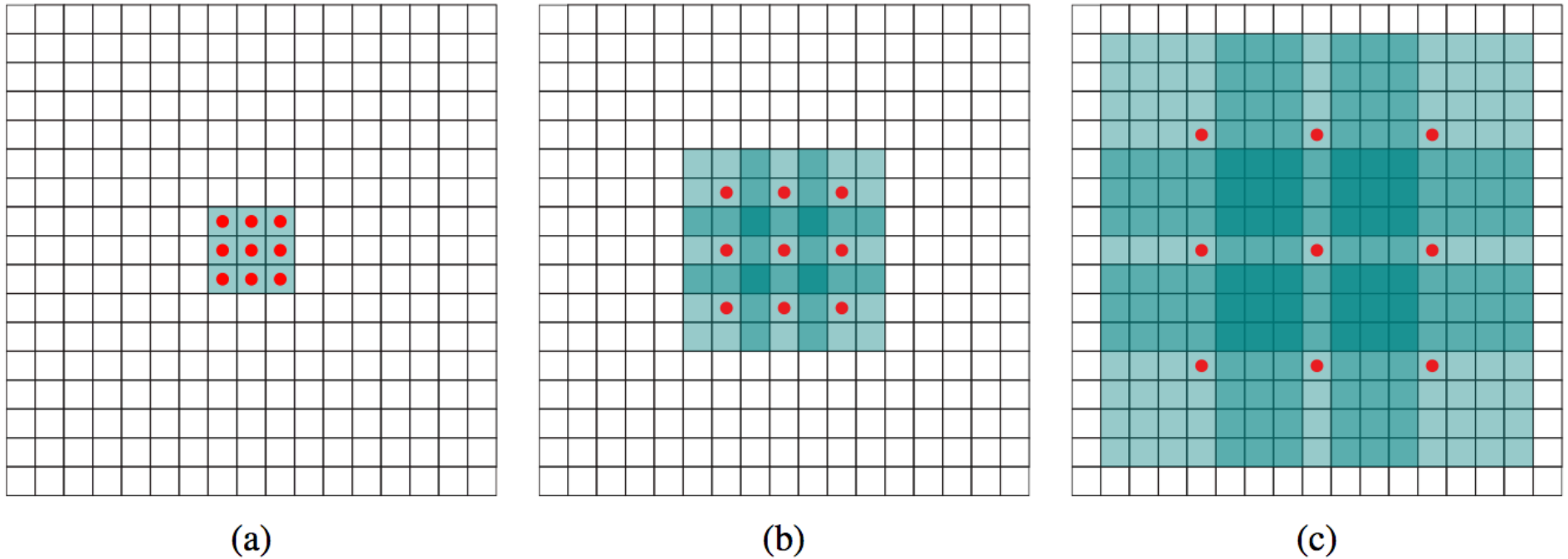
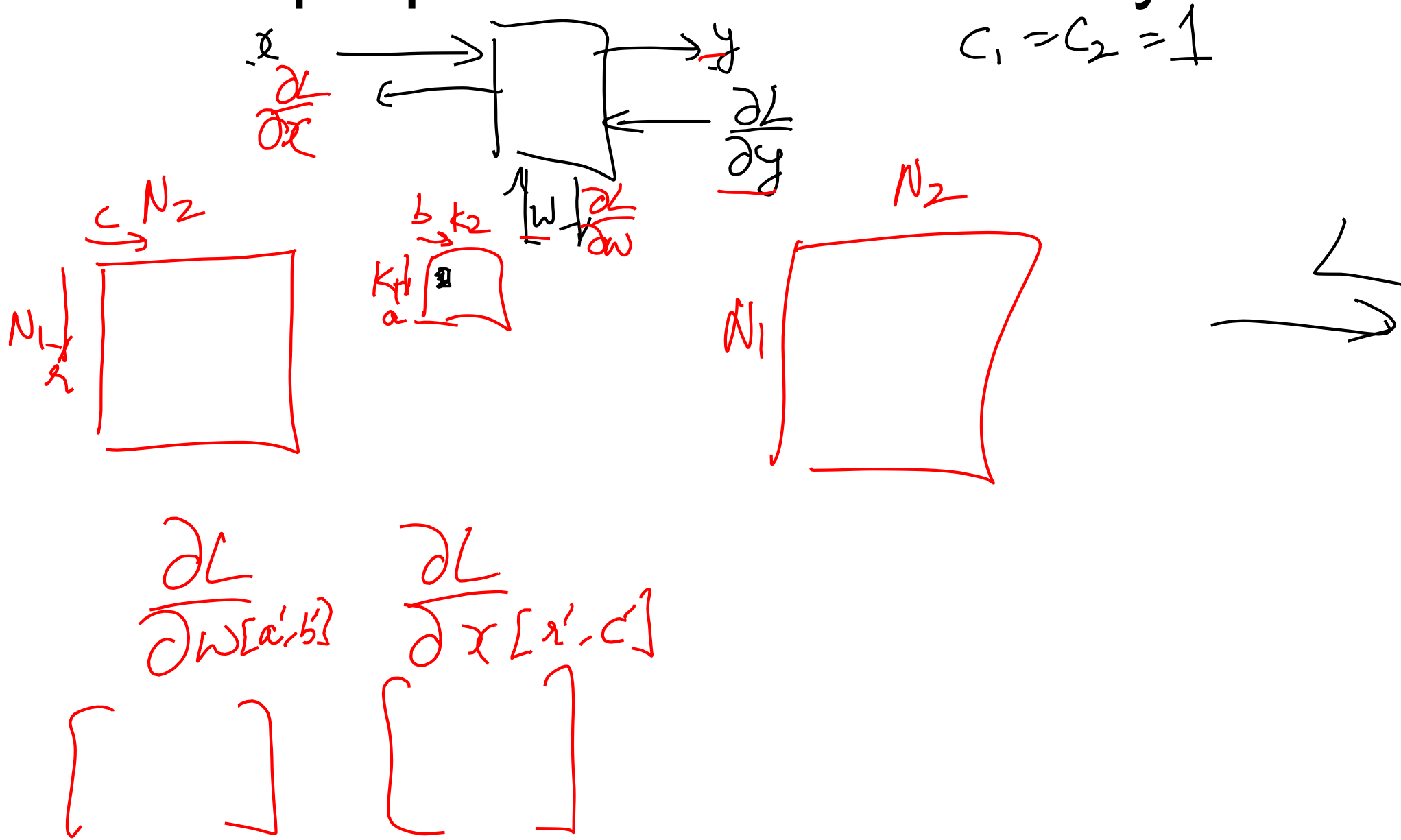


Figure 1: Systematic dilation supports exponential expansion of the receptive field without loss of resolution or coverage. (a) F_1 is produced from F_0 by a 1-dilated convolution; each element in F_1 has a receptive field of 3×3 . (b) F_2 is produced from F_1 by a 2-dilated convolution; each element in F_2 has a receptive field of 7×7 . (c) F_3 is produced from F_2 by a 4-dilated convolution; each element in F_3 has a receptive field of 15×15 . The number of parameters associated with each layer is identical. The receptive field grows exponentially while the number of parameters grows linearly.

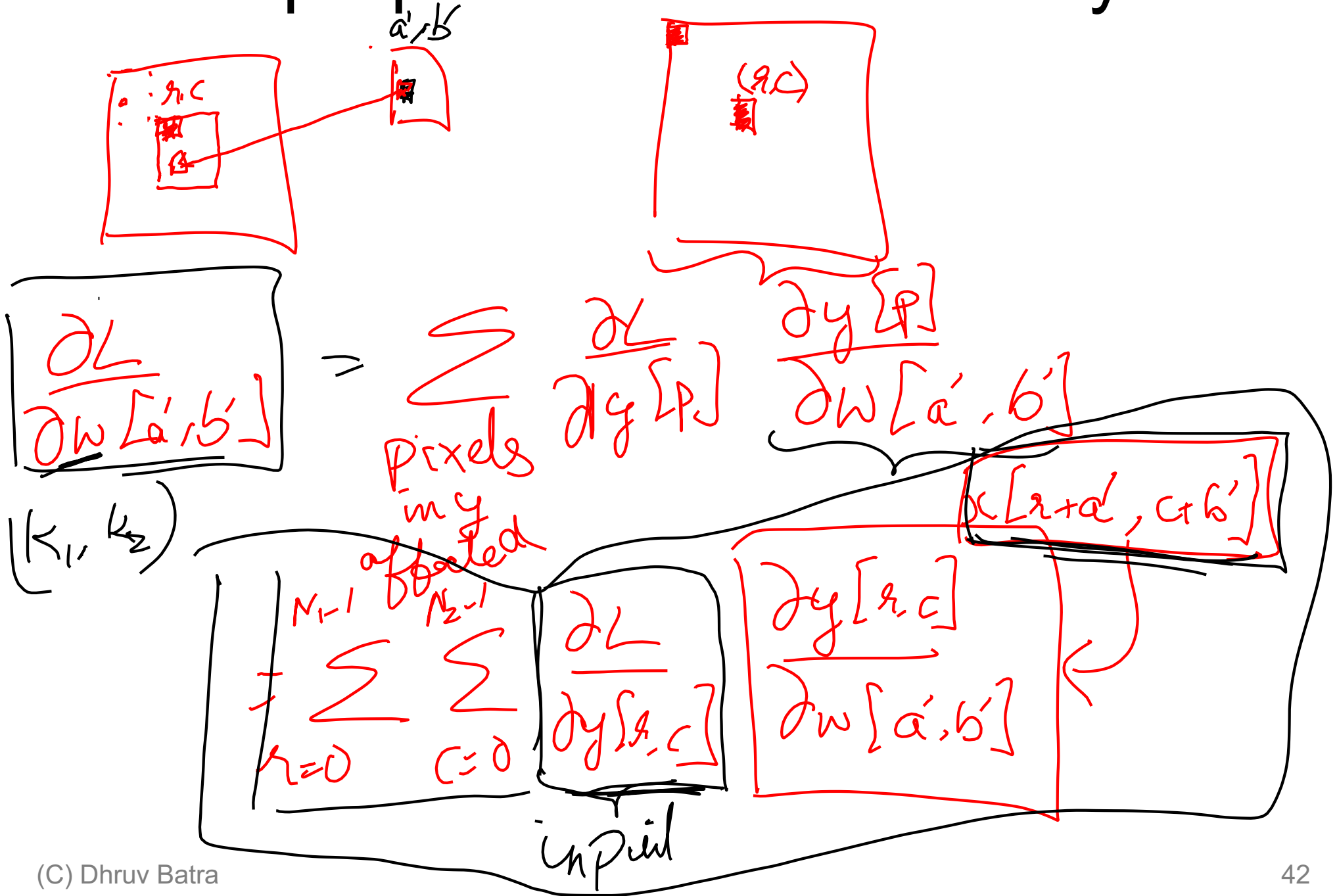
Plan for Today

- Convolutional Neural Networks
 - Toeplitz matrices and convolutions = matrix-mult
 - Dilated/a-trous convolutions
 - Backprop in conv layers
 - Transposed convolutions

Backprop in Convolutional Layers



Backprop in Convolutional Layers



$\frac{\partial L}{\partial x[r', c']}$ Backprop in Convolutional Layers



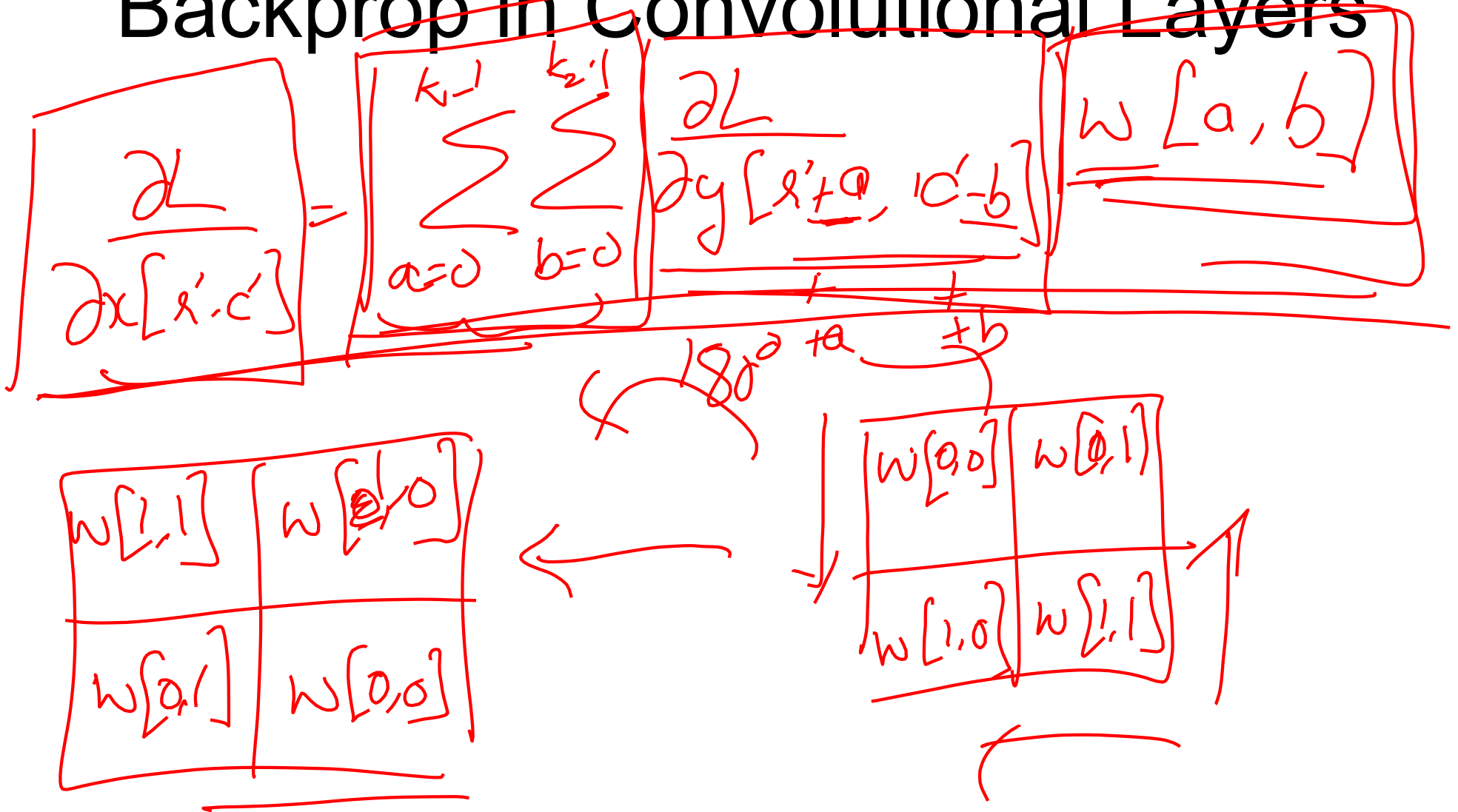
$$\frac{\partial L}{\partial x[r', c]} = \sum_{\text{pixels affected}} \frac{\partial L}{\partial y[r']} \frac{\partial y[r']} {\partial x[r', c]}$$

$$= \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} \left| \frac{\partial L}{\partial y[r'-a, c'-b]} \right| \left| \frac{\partial y[r'-a, c'-b]} {\partial x[r', c]} \right|$$

$$y[r'-a, c'-b] = \sum_{a'=0}^{k_1-1} \sum_{b'=0}^{k_2-1} x[r'+a', c'+b'] w[a', b']$$

$w[a, b]$

Backprop in Convolutional Layers



Plan for Today

- Convolutional Neural Networks
 - Toeplitz matrices and convolutions = matrix-mult
 - Dilated/a-trous convolutions
 - Backprop in conv layers
 - Transposed convolutions

Transposed Convolutions

- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

So far: Image Classification



[This image](#) is [CC0 public domain](#)

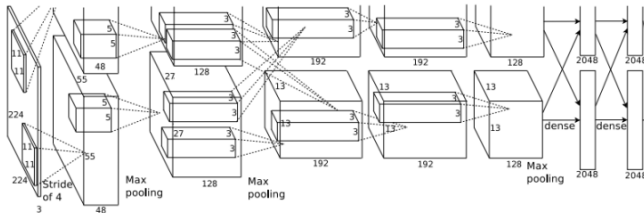


Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

Vector:
4096

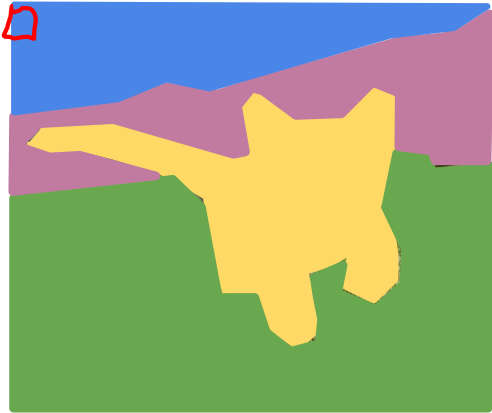
→
Fully-Connected:
4096 to 1000

Class Scores

Cat: 0.9
Dog: 0.05
Car: 0.01
...

Other Computer Vision Tasks

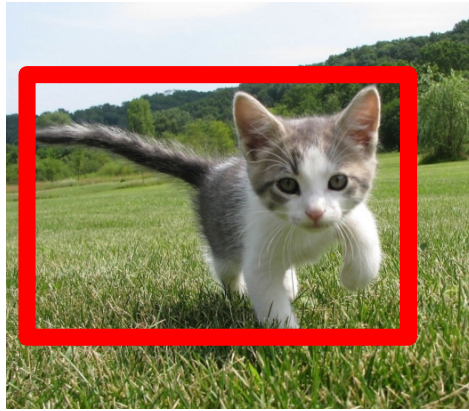
Semantic Segmentation



GRASS, CAT,
TREE, SKY

No objects, just pixels

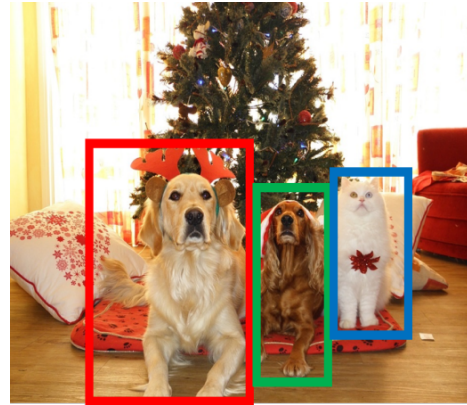
Classification + Localization



CAT

Single Object

Object Detection



DOG, DOG, CAT

Multiple Object

Instance Segmentation



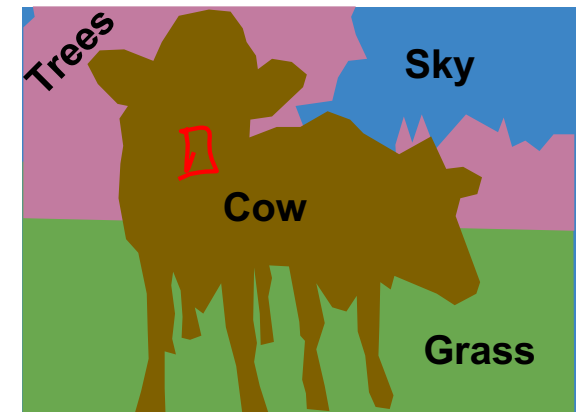
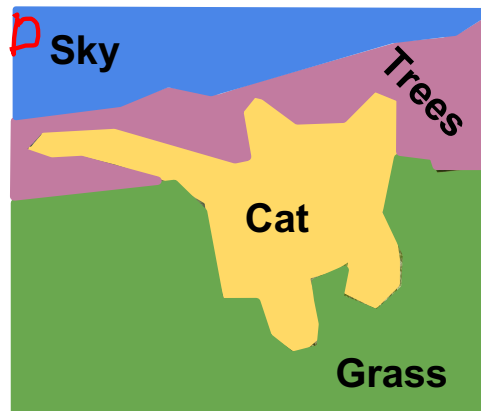
DOG, DOG, CAT

[This image is CC0 public domain](#)

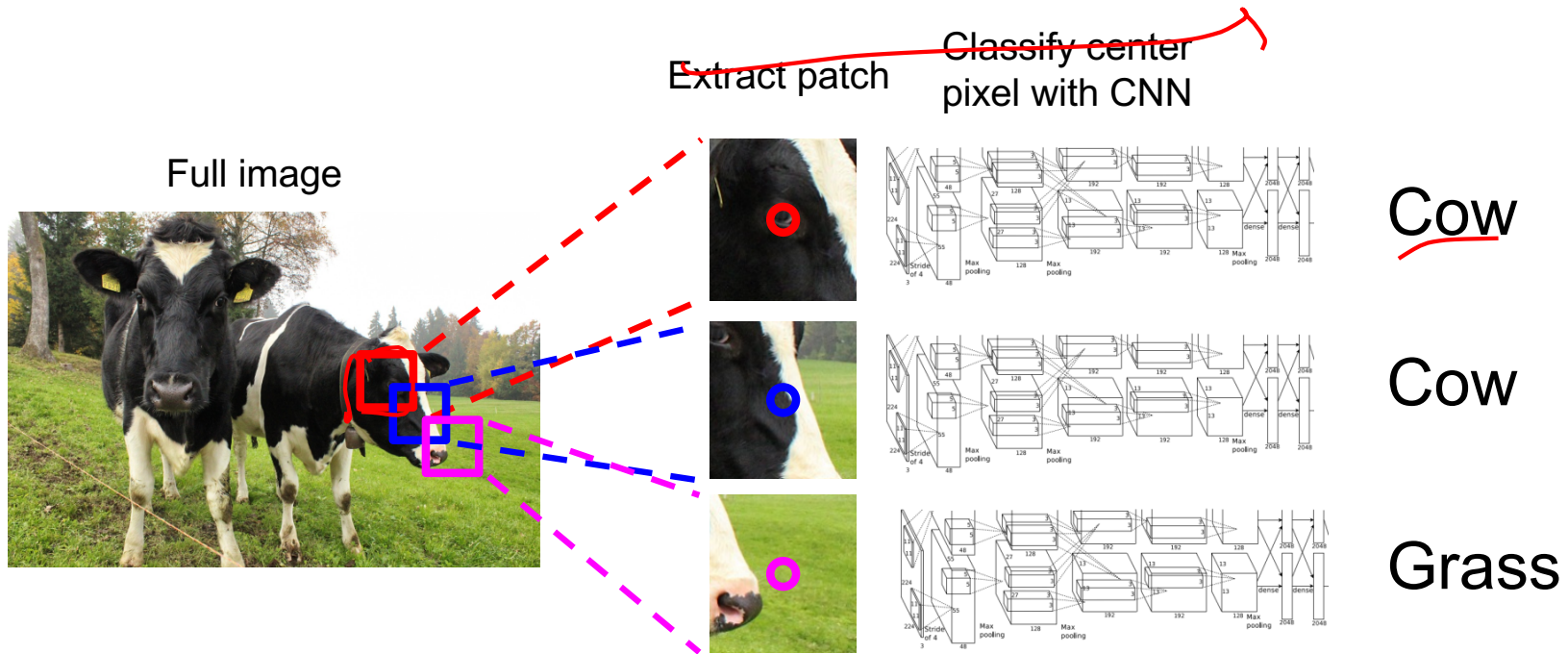
Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels



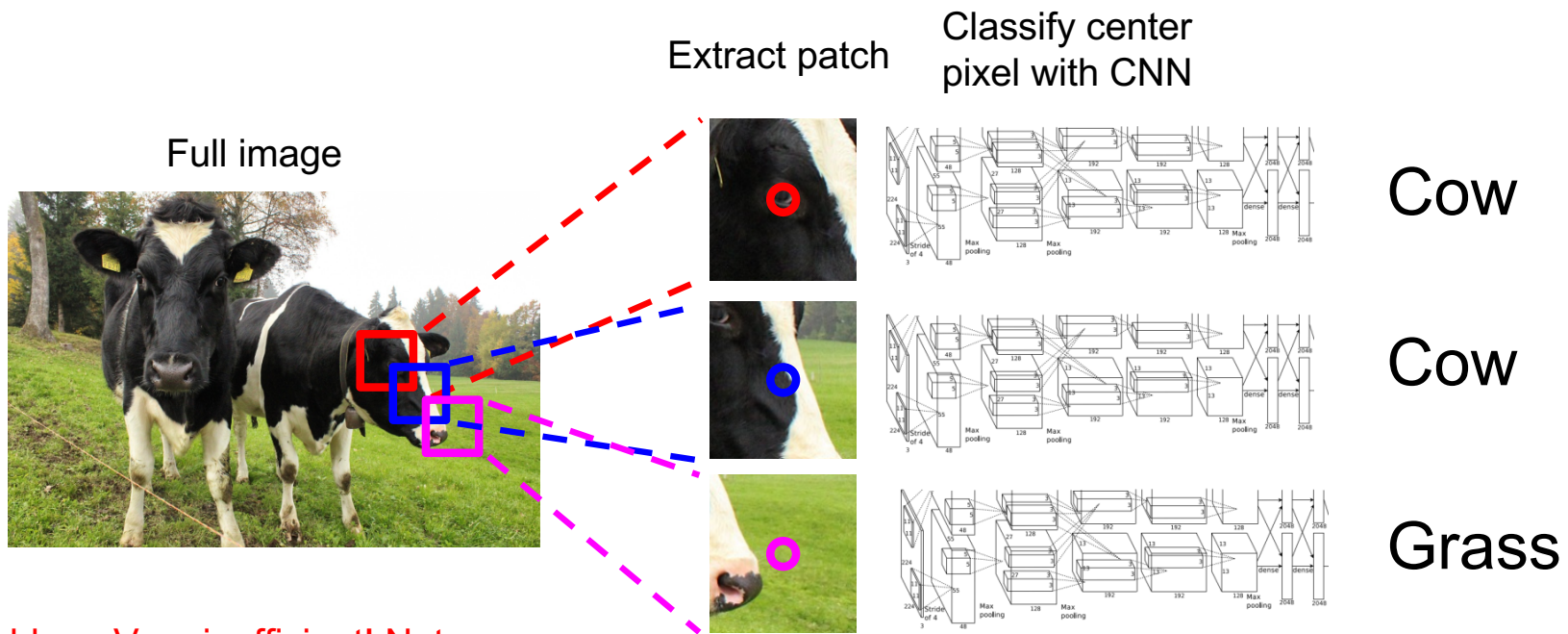
Semantic Segmentation Idea: Sliding Window



Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013

Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

Semantic Segmentation Idea: Sliding Window

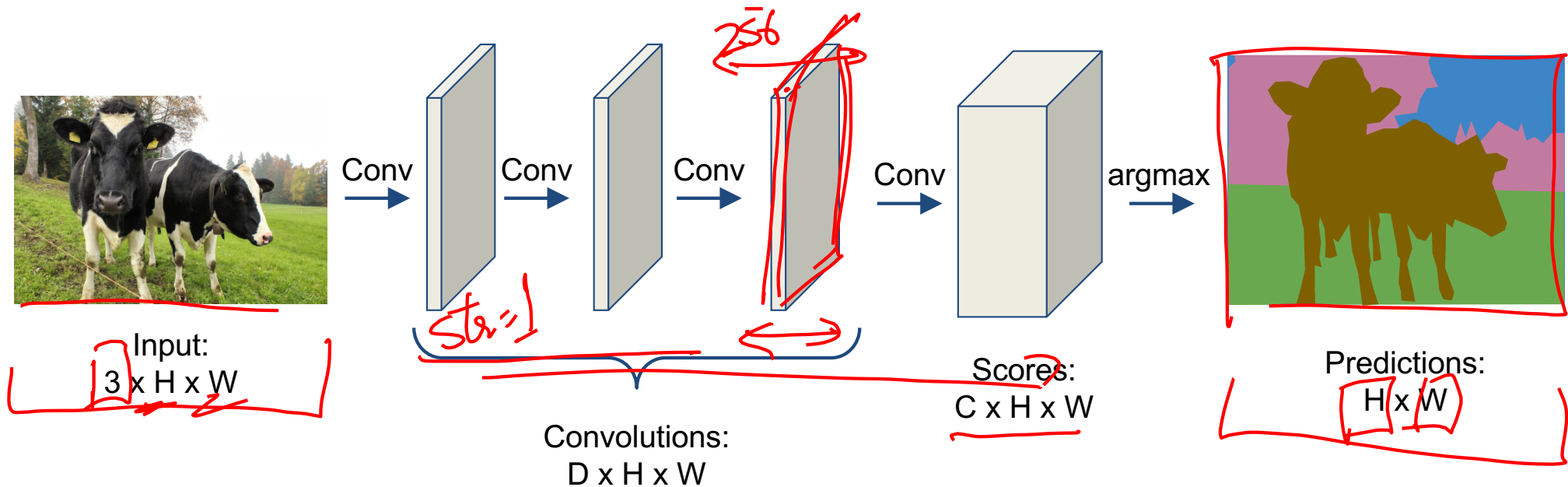


Problem: Very inefficient! Not reusing shared features between overlapping patches

Farabet et al, "Learning Hierarchical Features for Scene Labeling," TPAMI 2013
Pinheiro and Collobert, "Recurrent Convolutional Neural Networks for Scene Labeling", ICML 2014

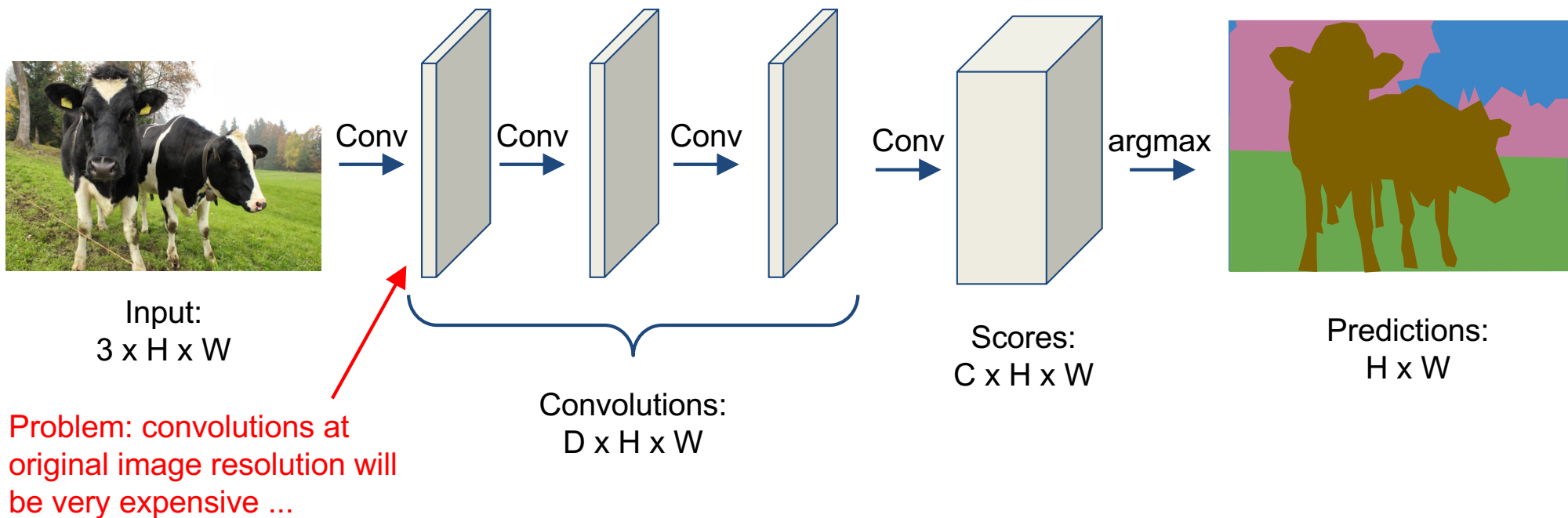
Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!



Semantic Segmentation Idea: Fully Convolutional

Design a network as a bunch of convolutional layers to make predictions for pixels all at once!

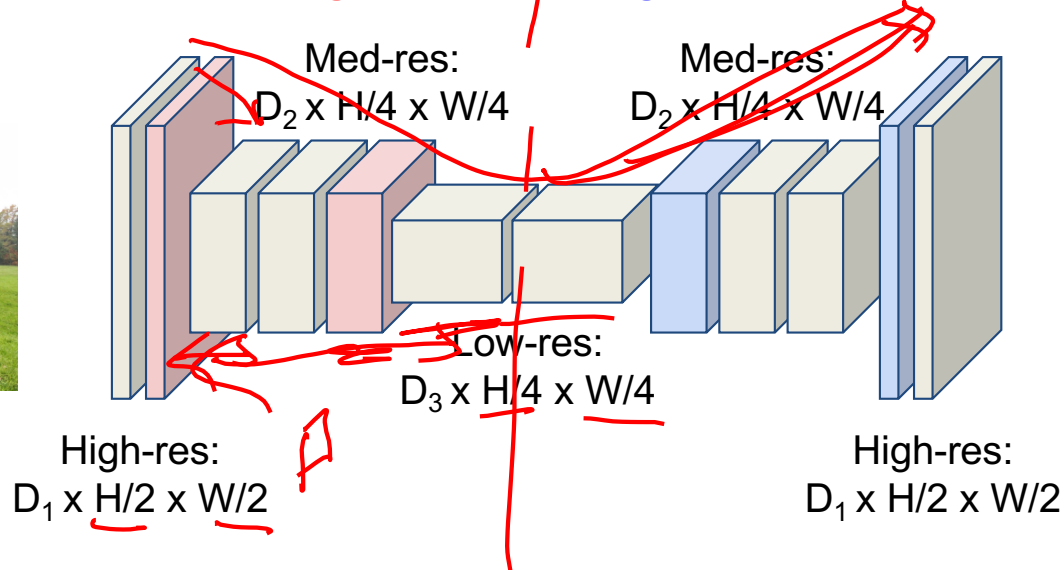


Semantic Segmentation Idea: Fully Convolutional

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



Input:
 $3 \times H \times W$



Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

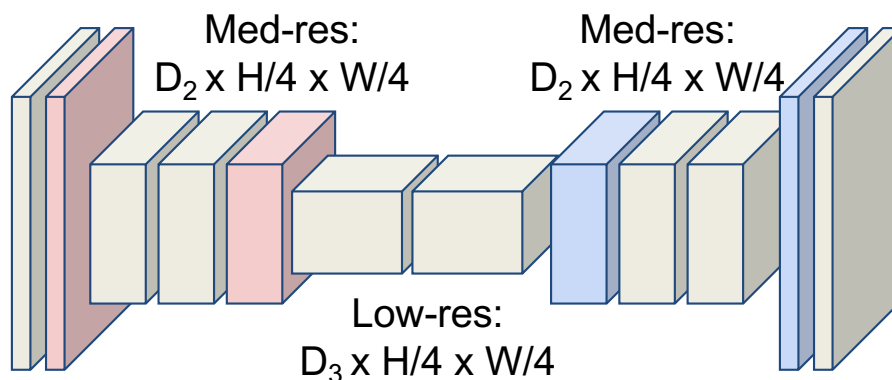
Semantic Segmentation Idea: Fully Convolutional

Downsampling:
Pooling, strided convolution



Input:
 $3 \times H \times W$

Design network as a bunch of convolutional layers, with **downsampling** and **upsampling** inside the network!



High-res:
 $D_1 \times H/2 \times W/2$

High-res:
 $D_1 \times H/2 \times W/2$

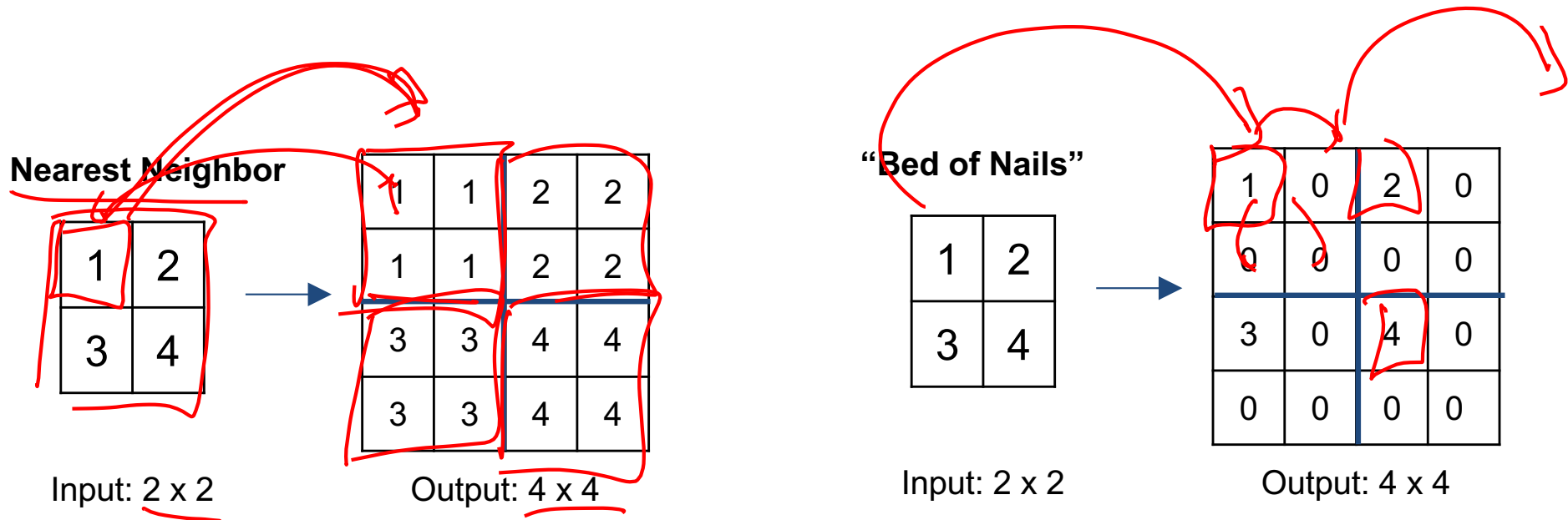
Upsampling:
???



Predictions:
 $H \times W$

Long, Shelhamer, and Darrell, "Fully Convolutional Networks for Semantic Segmentation", CVPR 2015
Noh et al, "Learning Deconvolution Network for Semantic Segmentation", ICCV 2015

In-Network upsampling: “Unpooling”



In-Network upsampling: "Max Unpooling"

Max Pooling

Remember which element was max!

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

Input: 4 x 4

argmax

5	6
7	8

Output: 2 x 2

Rest of the network

Max Unpooling

Use positions from pooling layer

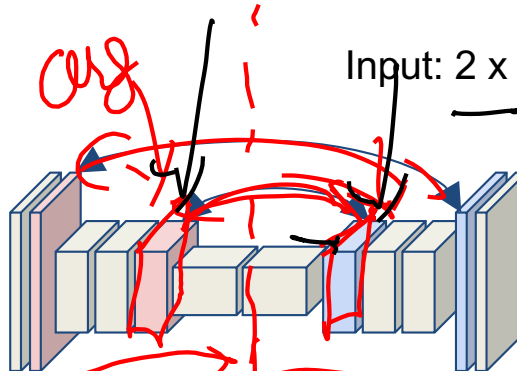
1	2
3	4

Input: 2 x 2

0	0	2	0
0	1	0	0
0	0	0	0
3	0	0	4

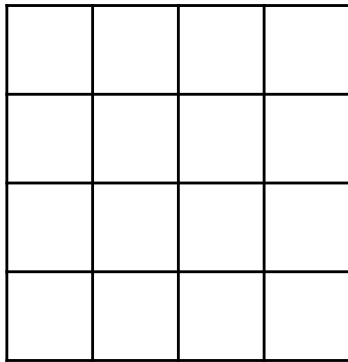
Output: 4 x 4

Corresponding pairs of downsampling and upsampling layers

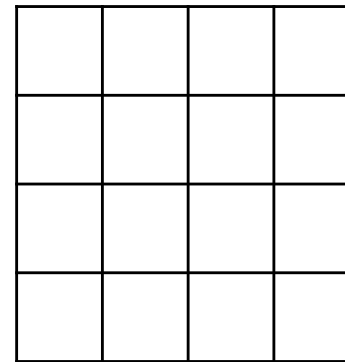


Learnable Upsampling: Transpose Convolution

Recall: Typical 3 x 3 convolution, stride 1 pad 1



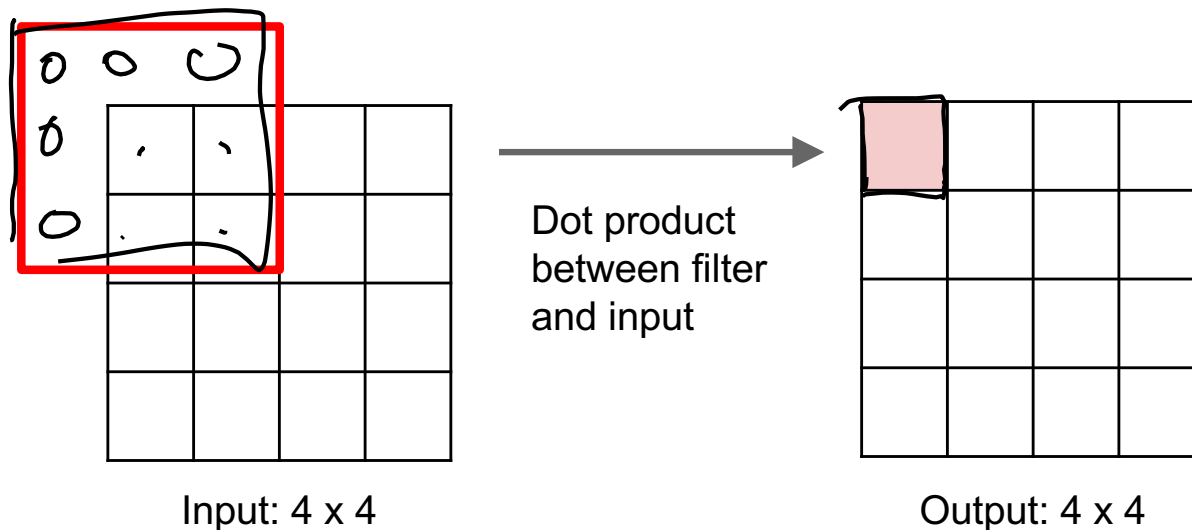
Input: 4 x 4



Output: 4 x 4

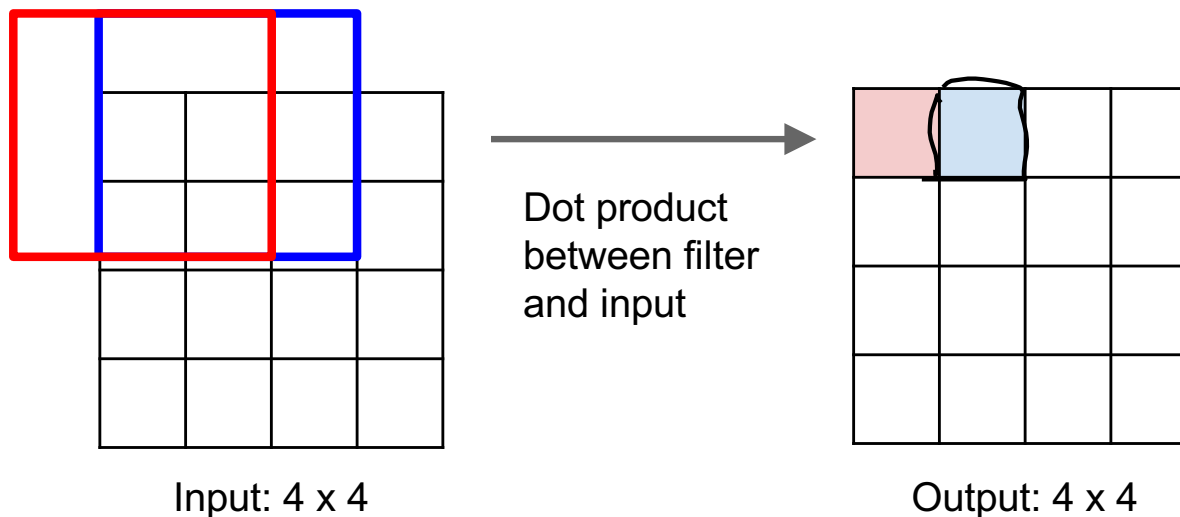
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1



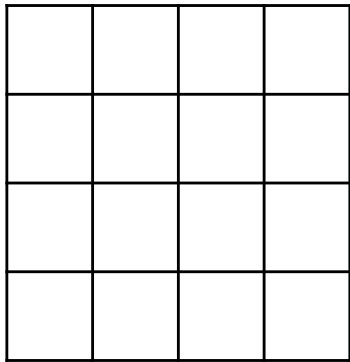
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 1 pad 1

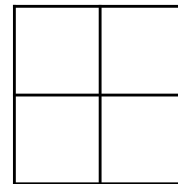


Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



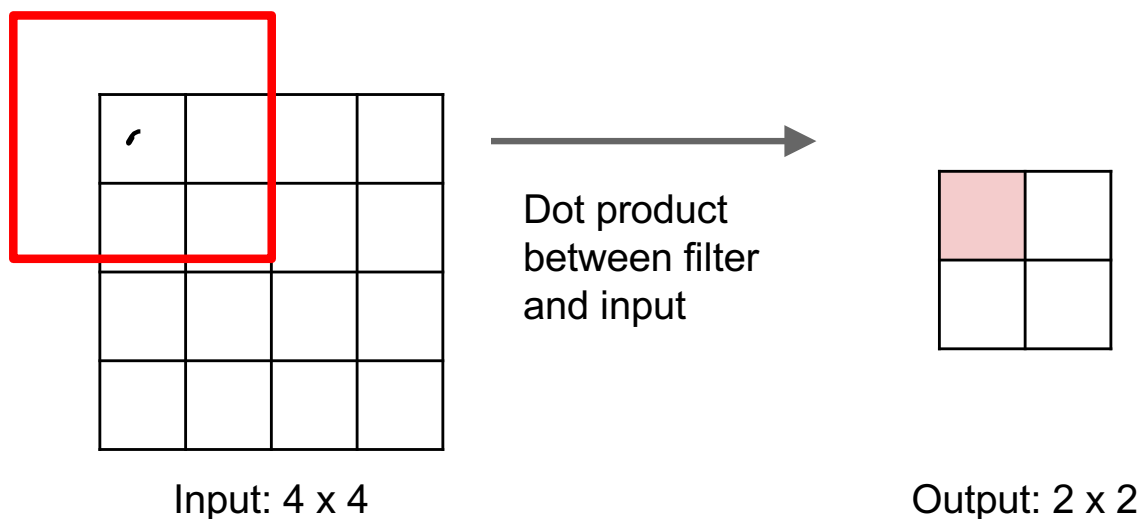
Input: 4 x 4



Output: 2 x 2

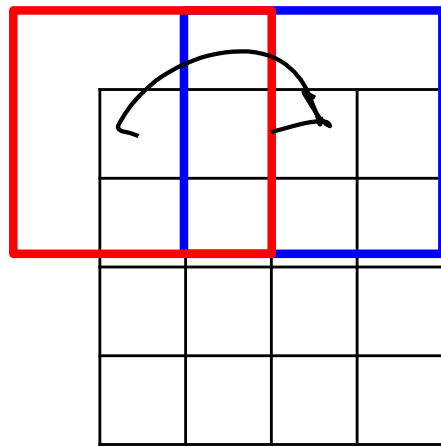
Learnable Upsampling: Transpose Convolution

Recall: Normal 3 x 3 convolution, stride 2 pad 1



Learnable Upsampling: Transpose Convolution

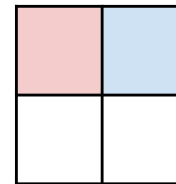
Recall: Normal 3 x 3 convolution, stride 2 pad 1



Input: 4 x 4



Dot product
between filter
and input



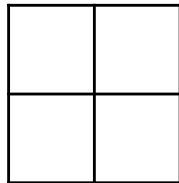
Output: 2 x 2

Filter moves 2 pixels in
the input for every one
pixel in the output

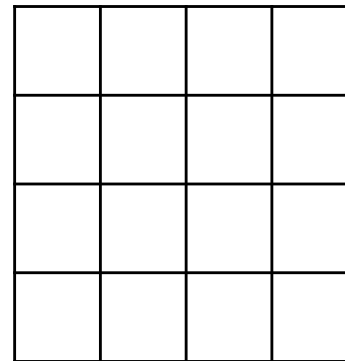
Stride gives ratio between
movement in input and
output

Learnable Upsampling: Transpose Convolution

3 x 3 **transpose** convolution, stride 2 pad 1

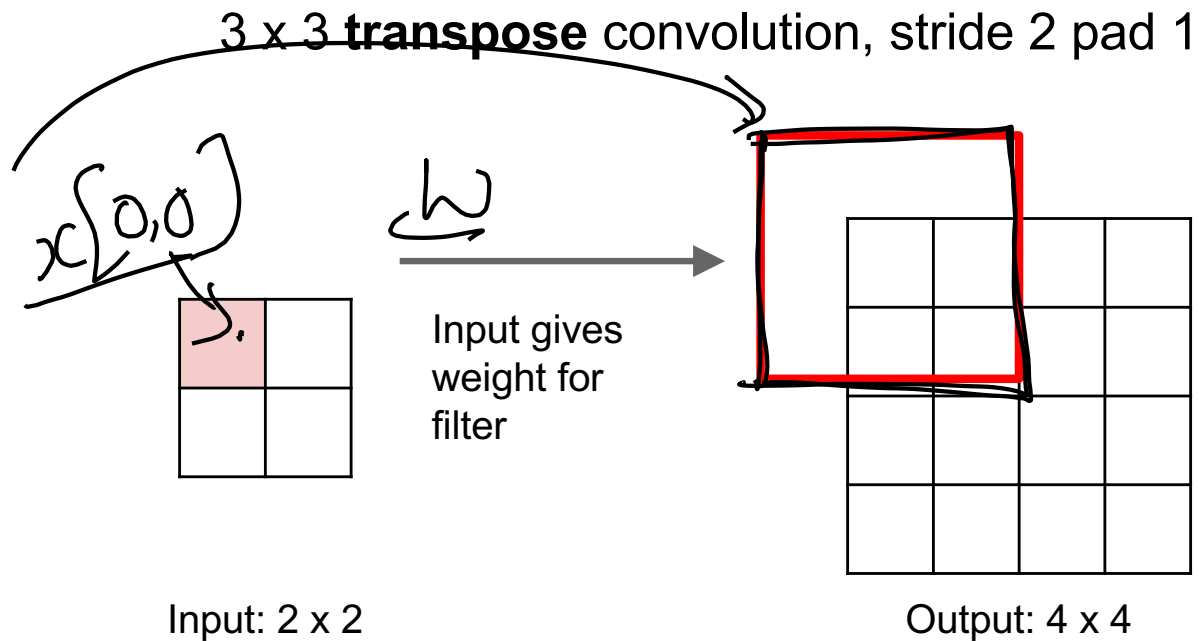


Input: 2 x 2

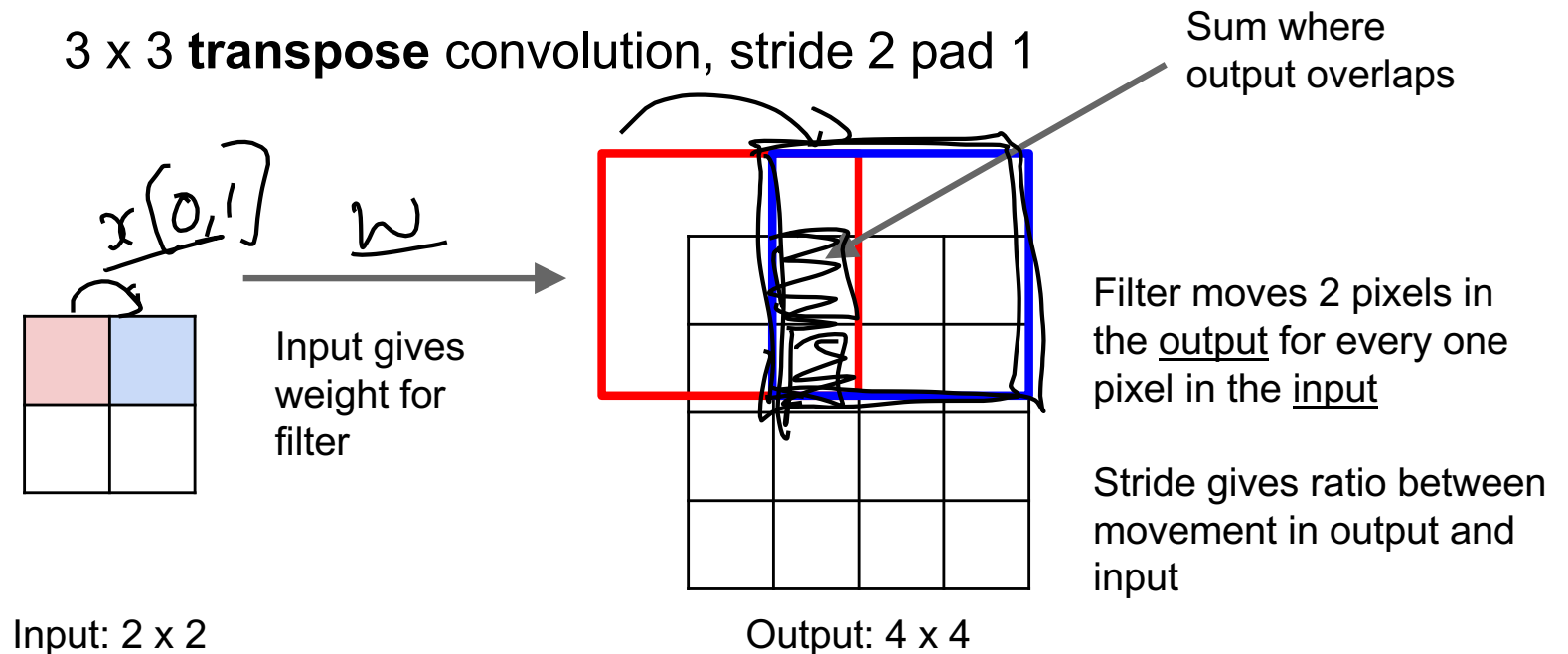


Output: 4 x 4

Learnable Upsampling: Transpose Convolution



Learnable Upsampling: Transpose Convolution

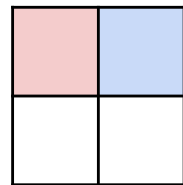


Learnable Upsampling: Transpose Convolution

Other names:

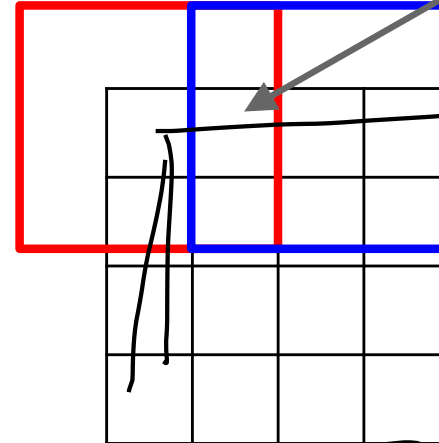
- Deconvolution (bad)
- Upconvolution
- Fractionally strided convolution
- Backward strided convolution

3 x 3 **transpose** convolution, stride 2 pad 1



Input: 2 x 2

Input gives weight for filter



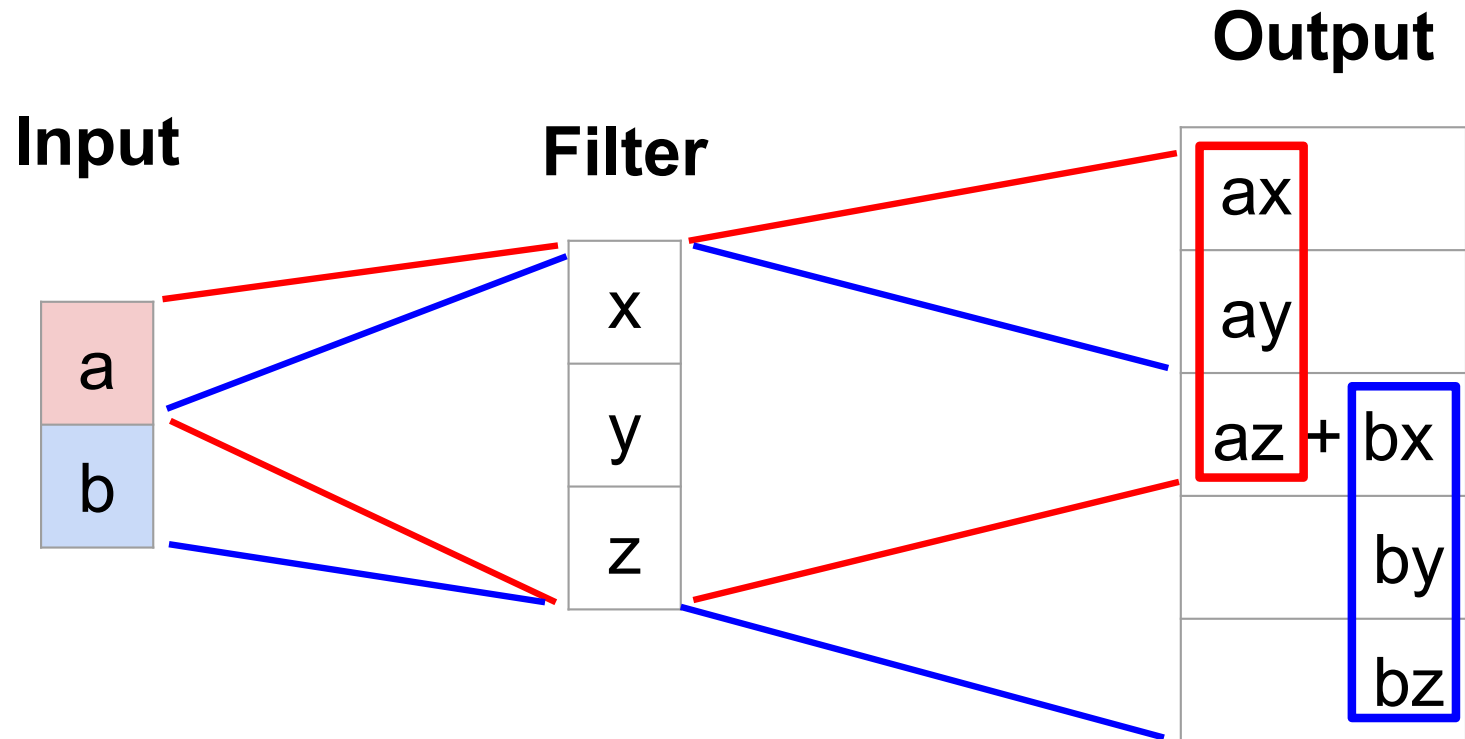
Output: 4 x 4

Sum where output overlaps

Filter moves 2 pixels in the output for every one pixel in the input

Stride gives ratio between movement in output and input

Transpose Convolution: 1D Example



Output contains copies of the filter weighted by the input, summing at where it overlaps in the output

Need to crop one pixel from output to make output exactly 2x input

Transposed Convolution

- <https://distill.pub/2016/deconv-checkerboard/>