# CS 4803 / 7643: Deep Learning

Topics:
- Convolutional Neural Networks
    - Pooling layers
    - Fully-connected layers as convolutions
    - Backprop in conv layers [Derived in notes]
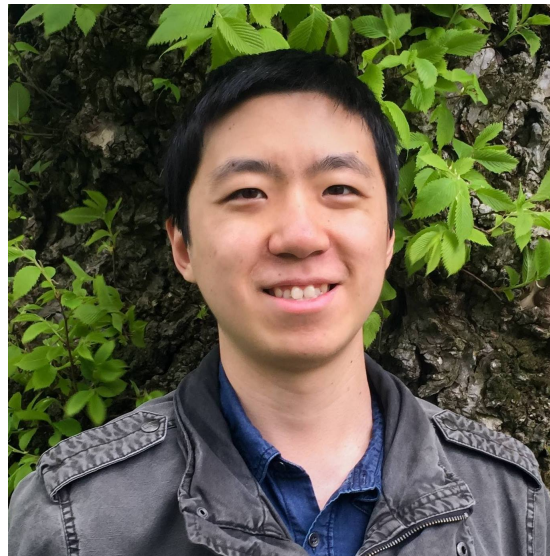    - Toeplitz matrices and convolutions = matrix-mult

Dhruv Batra

Georgia Tech

# Administrativia

- HW1 Reminder
  - Due: 09/26, 11:55pm
  - https://evalai.cloudcv.org/web/challenges/challenge-page/431/leaderboard/1200


- Project Teams Google Doc
  - https://docs.google.com/spreadsheets/d/1ouD6ctaemV_3nb2MQHs7rUOAaW9DFLu8I5Zd3yOFs7E/edit?usp=sharing
  - Project Title
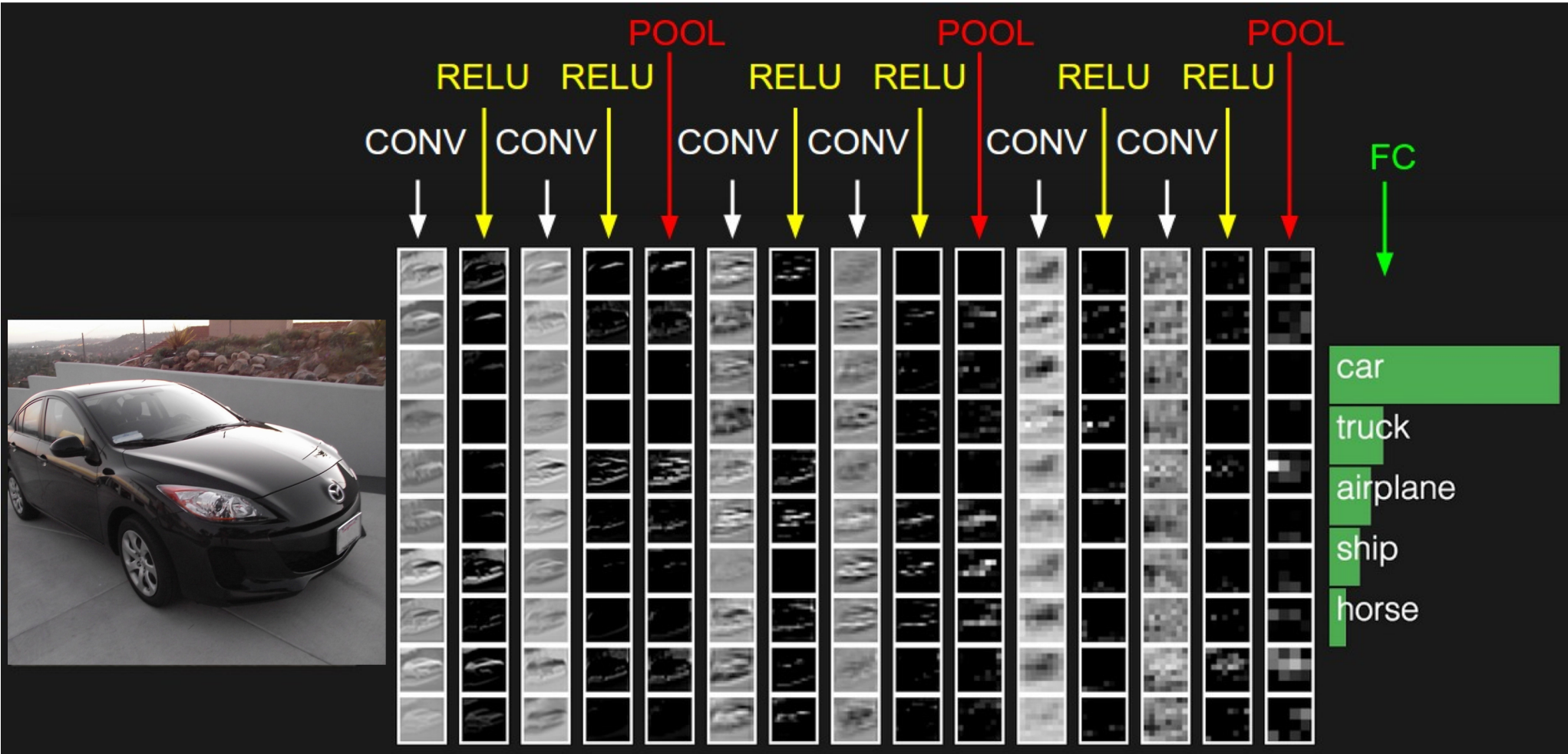  - 1-3 sentence project summary TL;DR
  - Team member names

# Administrativia

- Guest Lecture: Dr. Zhile Ren
  - Next class (09/26)
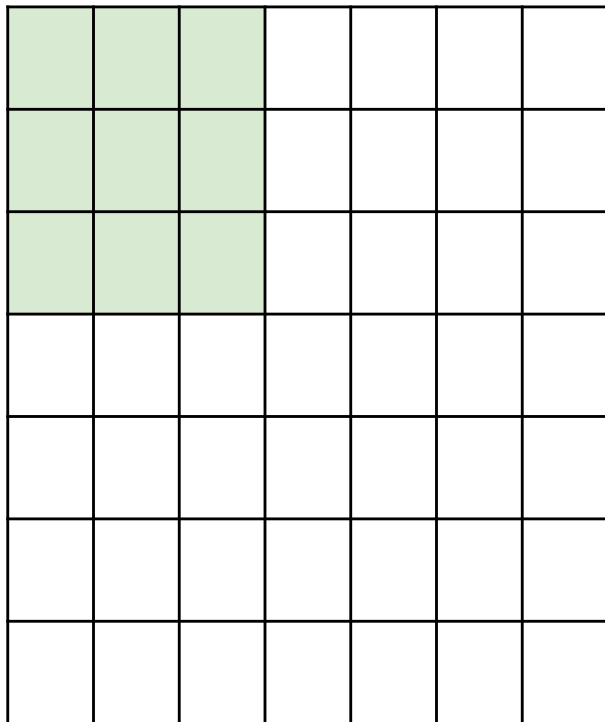  - CNN architectures for 2D & 3D Detection & Segmentation

# Recap from last time

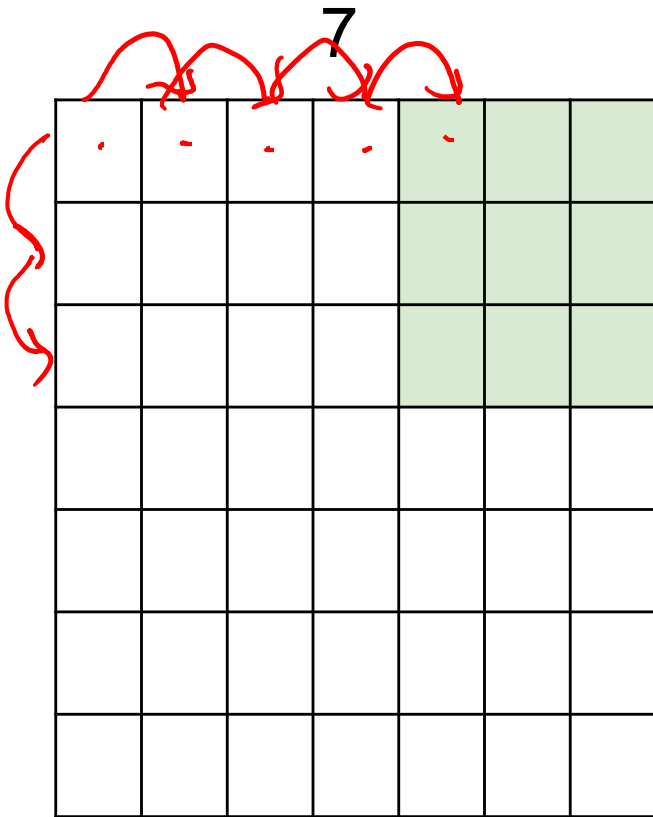preview:

# A closer look at spatial dimensions:

7



7

7x7 input (spatially)
assume 3x3 filter

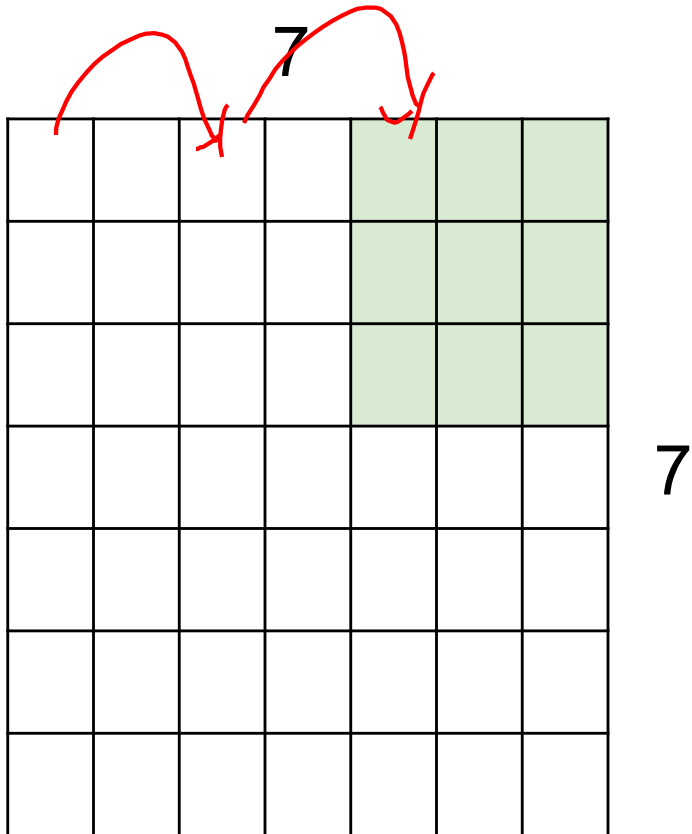A closer look at spatial dimensions:

stride = 1

['valid']

'same' )

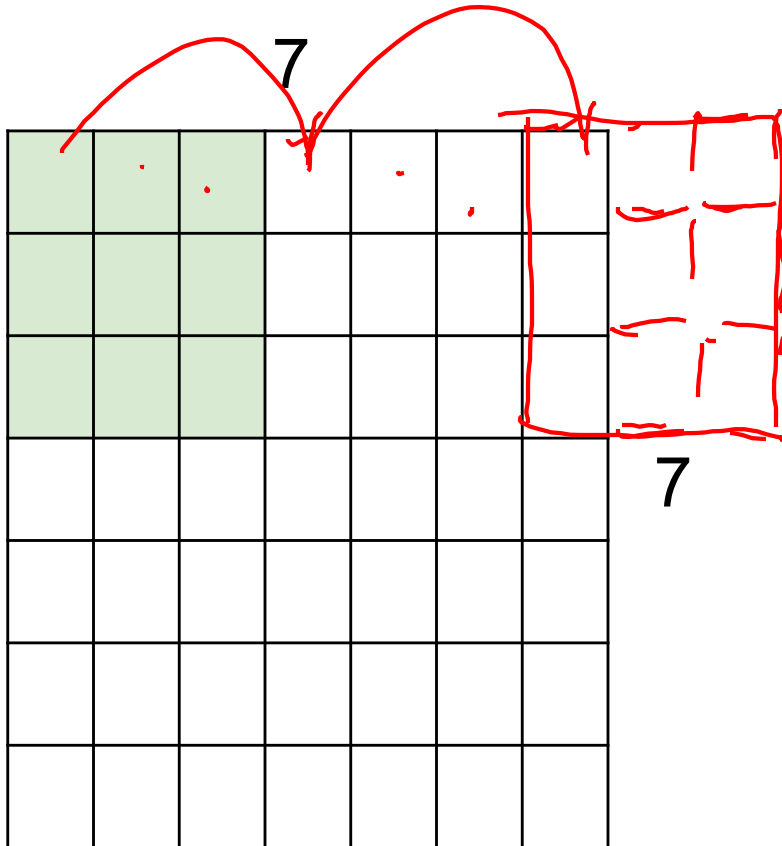7x7 input (spatially)
assume 3x3 filter

=> **5x5 output**

7

7

A closer look at spatial dimensions:



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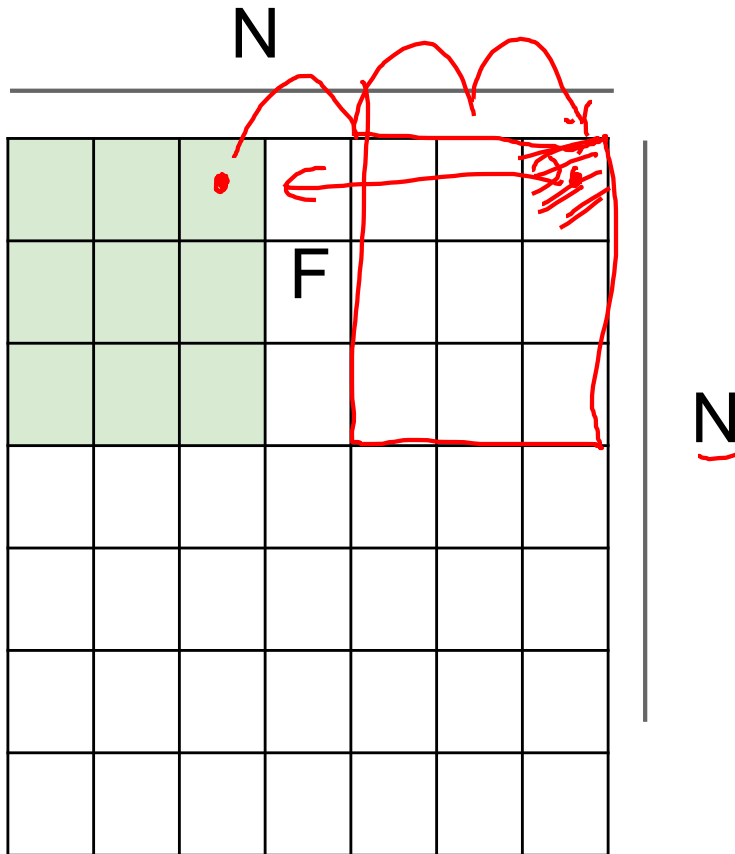assume 3x3 filter
applied **with stride 2**
**=> 3x3 output!**

A closer look at spatial dimensions:



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

'valid'

N-F & stride

N

F

N

Output size: = $\dfrac{N-F}{stride}$

**(N - F) / stride + 1**

e.g. N = 7, F = 3:
stride 1 => (7 - 3)/1 + 1 = 5
stride 2 => (7 - 3)/2 + 1 = 3
stride 3 => (7 - 3)/3 + 1 = 2.33 :\

$$\dfrac{7-3}{4} + 1$$

# 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?

$$N = \left[ \frac{N + 2pad - F}{stride} + 1 \right.$$

pad $= \frac{F-1}{2}$

(recall:)

(N - F) / stride + 1

$\tilde{N}$

$\frac{9-3}{1} + 1 = 6 + 1 = 7$

# 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?

**7x7 output!**
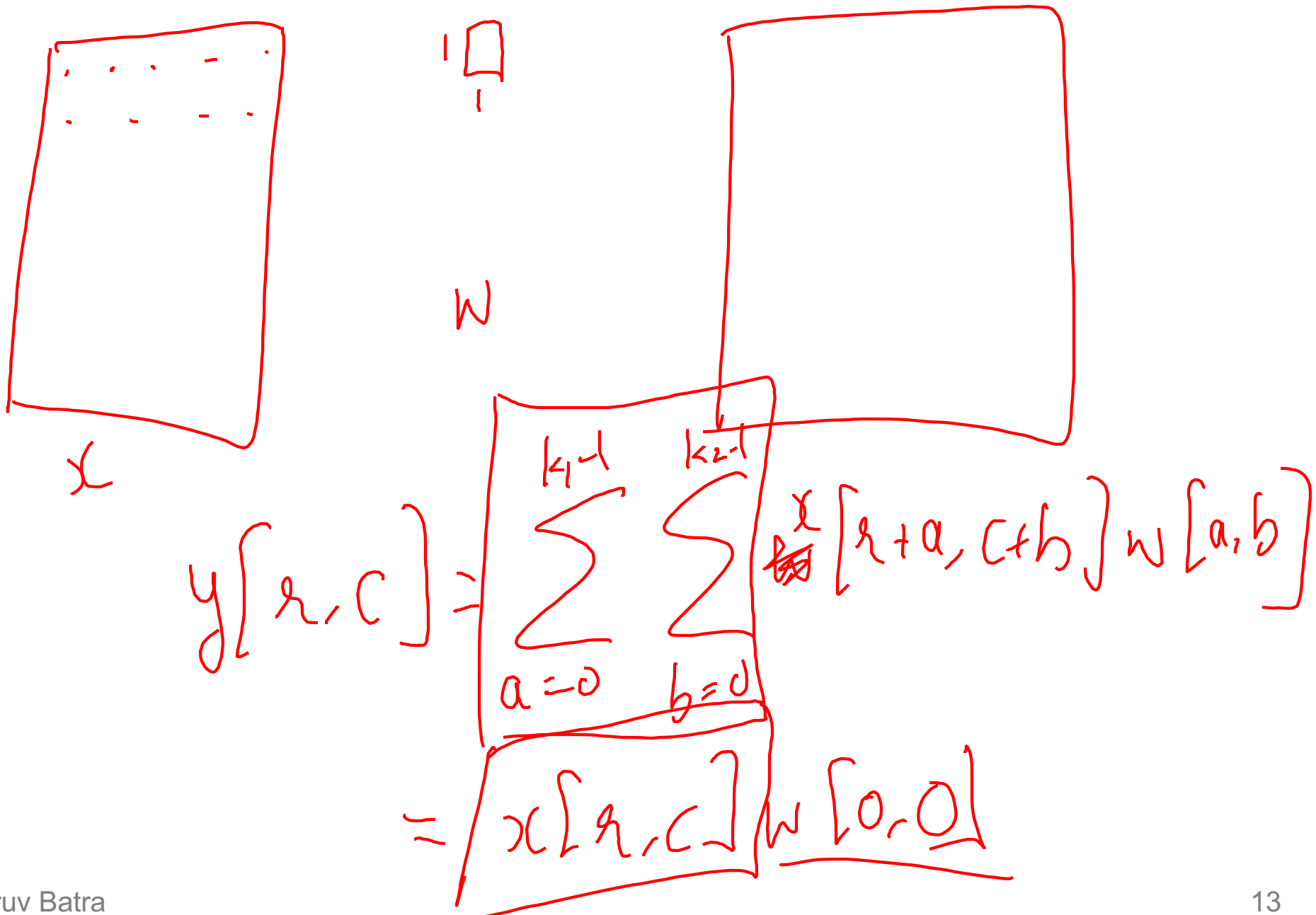in general, common to see CONV layers with stride 1, filters of size FxF, and zero-padding with $(F-1)/2$ (will preserve size spatially)

e.g. F = 3 => zero pad with 1
F = 5 => zero pad with 2
F = 7 => zero pad with 3

# Can we have 1x1 filters?



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

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

# 1x1 convolution layers make perfect sense



1x1 CONV
with 32 filters

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

# Plan for Today

- **Convolutional Neural Networks**
  - Pooling layers
  - Fully-connected layers as convolutions
  - Backprop in conv layers [Derived in notes]
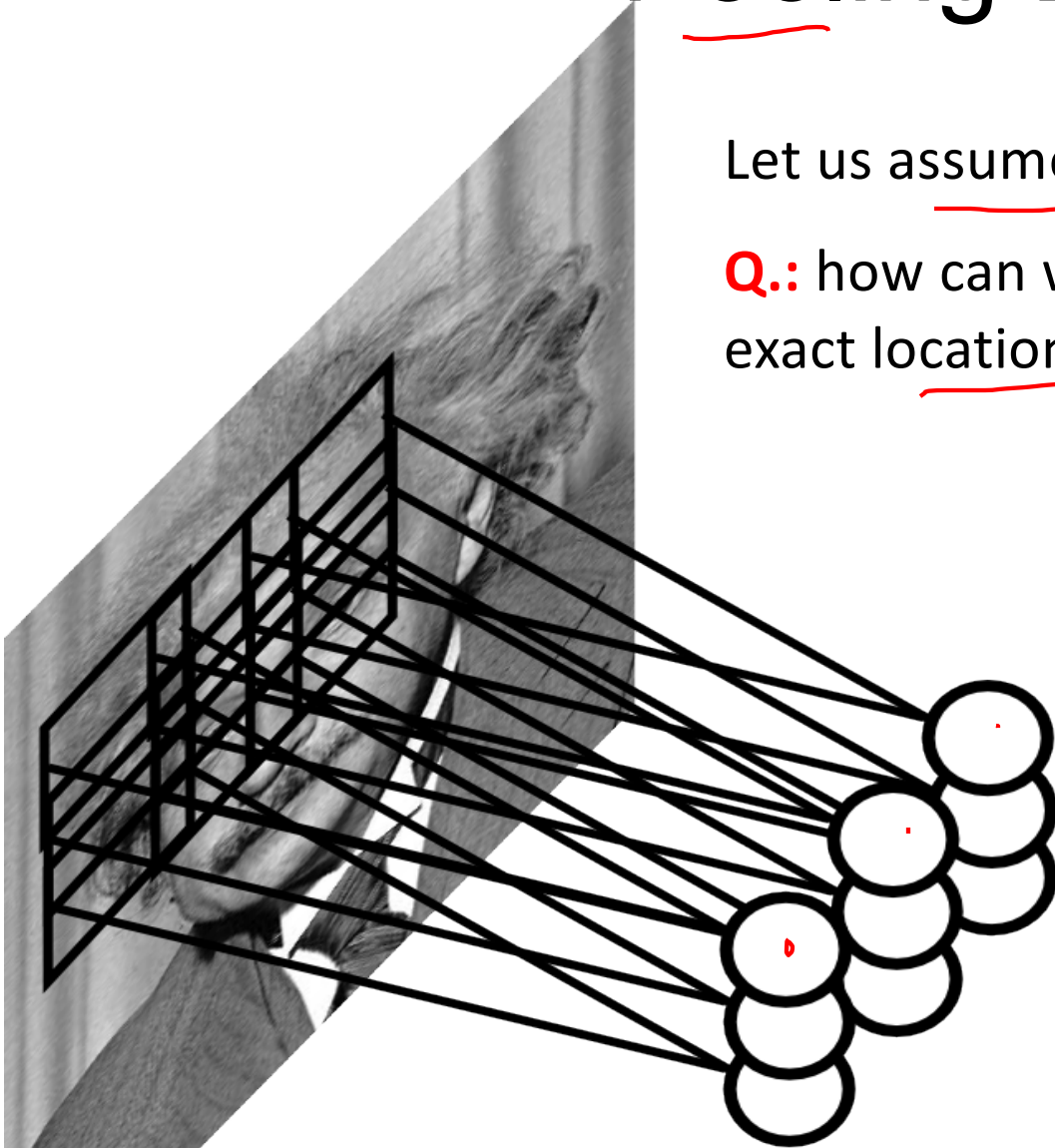  - Toeplitz matrices and convolutions = matrix-mult

two more layers to go: POOL/FC

# Pooling Layer
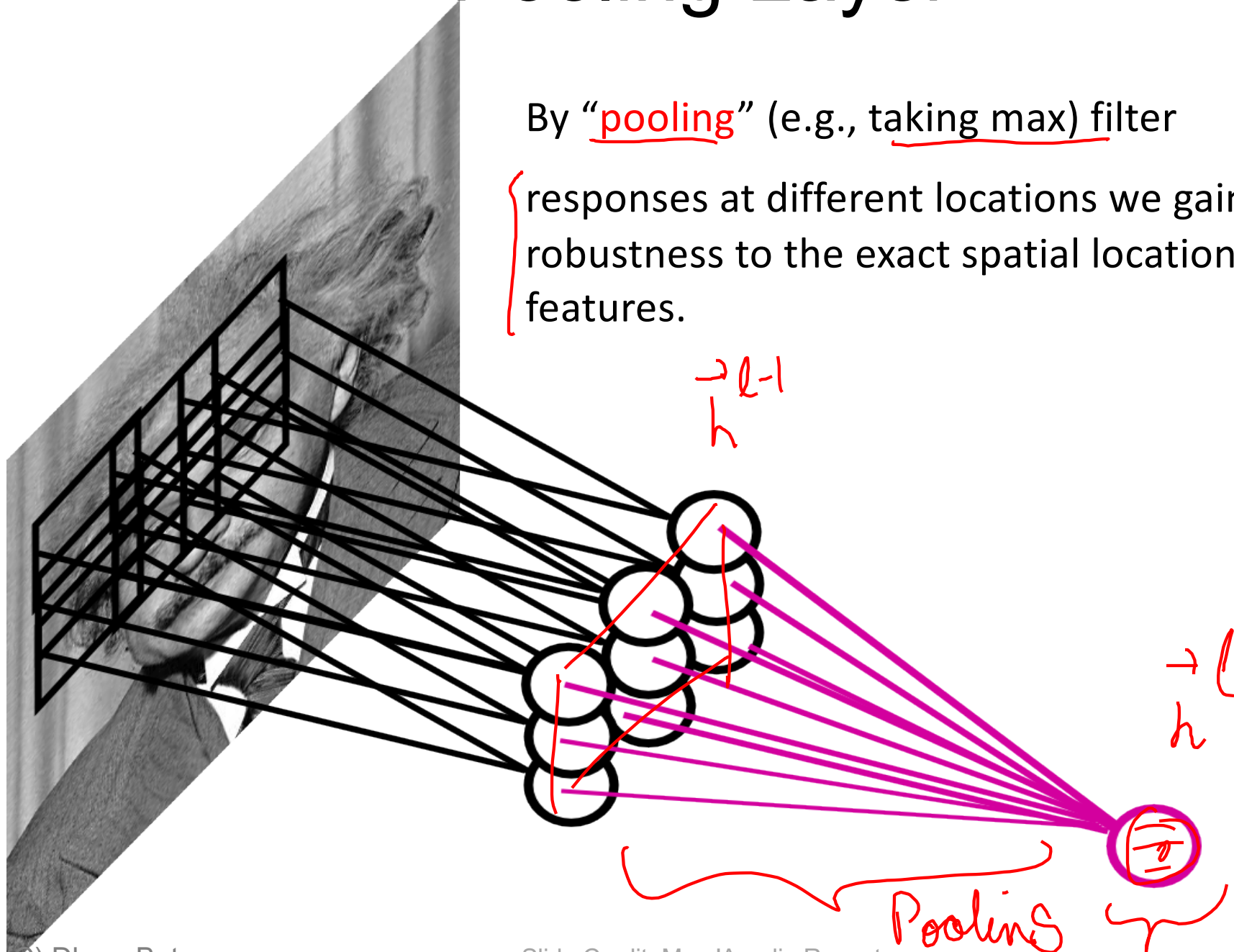
Let us assume filter is an "eye" detector.

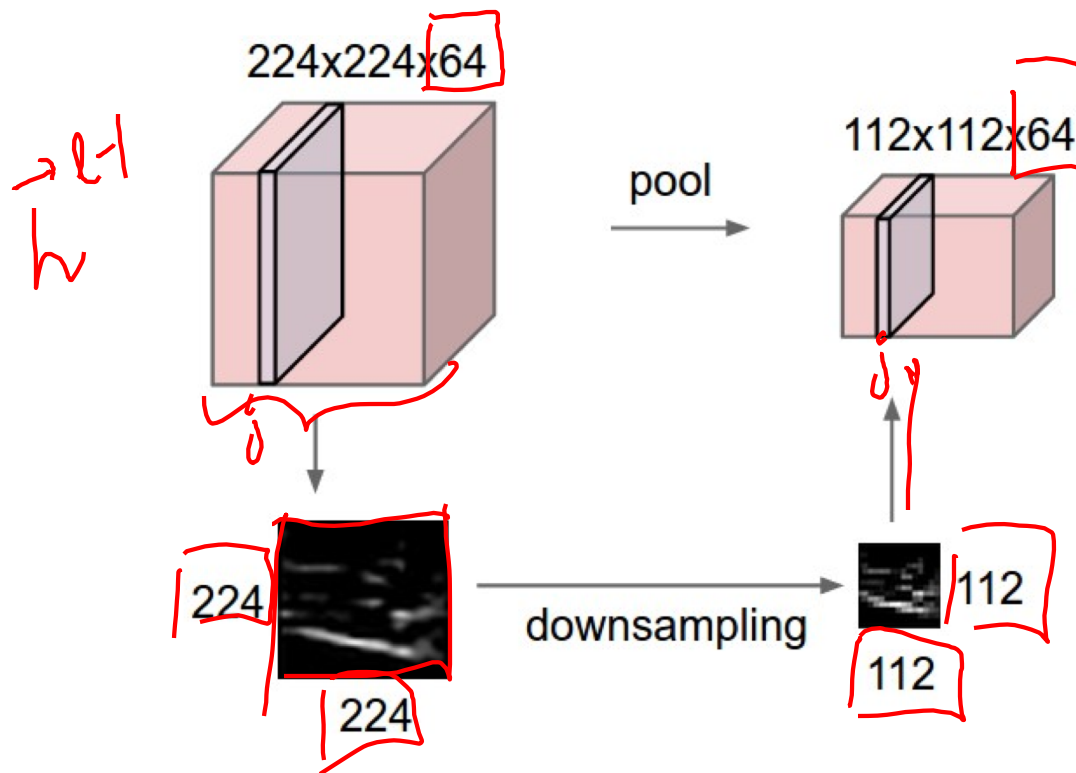**Q.:** how can we make the detection robust to the exact location of the eye?

# Pooling Layer

By "pooling" (e.g., taking max) filter responses at different locations we gain robustness to the exact spatial location of features.
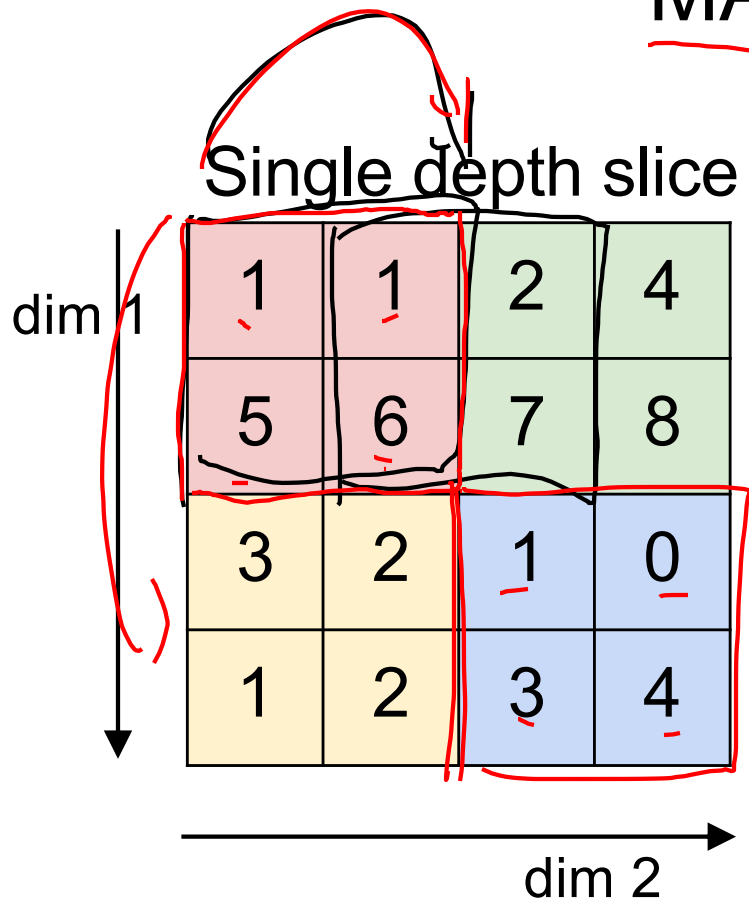
Slide Credit: Marc'Aurelio Ranzato

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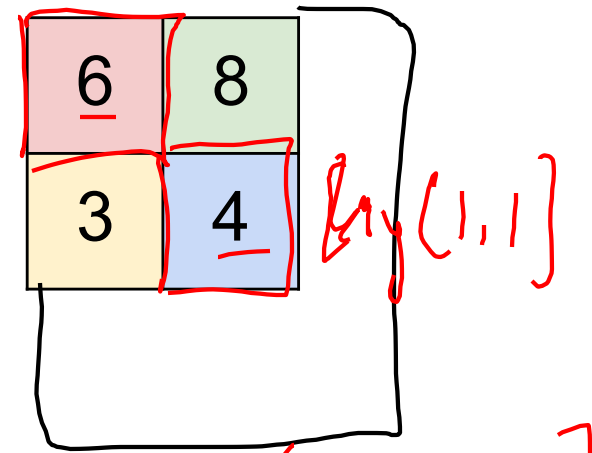
- makes the representations smaller and more manageable
- operates over each activation map independently:

# MAX POOLING

Single depth slice



dim 1

| 1 | 1 | 2 | 4 |
|---|---|---|---|
| 5 | 6 | 7 | 8 |
| 3 | 2 | 1 | 0 |
| 1 | 2 | 3 | 4 |

dim 2

max pool with 2x2 filters
and stride 2

| 6 | 8 |
|---|---|
| 3 | 4 |

$y[1,1]$

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

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

| | |
|---|---|
| | |
| | |

# Pooling Layer: Examples

Max-pooling:

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

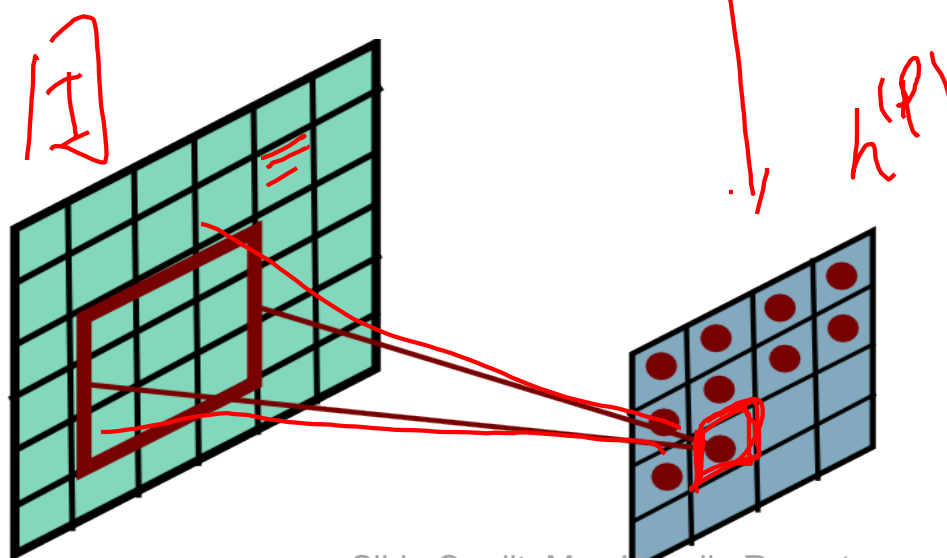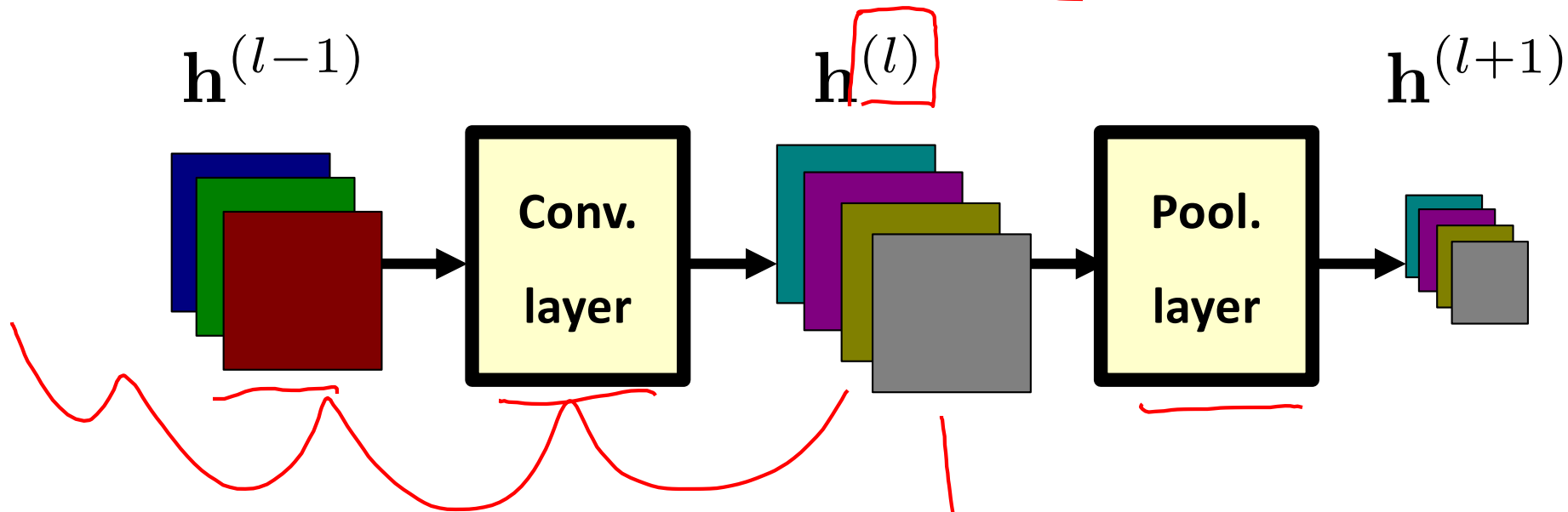Average-pooling:

$$h_i^n(r,c) = \operatorname*{mean}_{\bar{r} \in N(r),\ \bar{c} \in N(c)} h_i^{n-1}(\bar{r}, \bar{c})$$
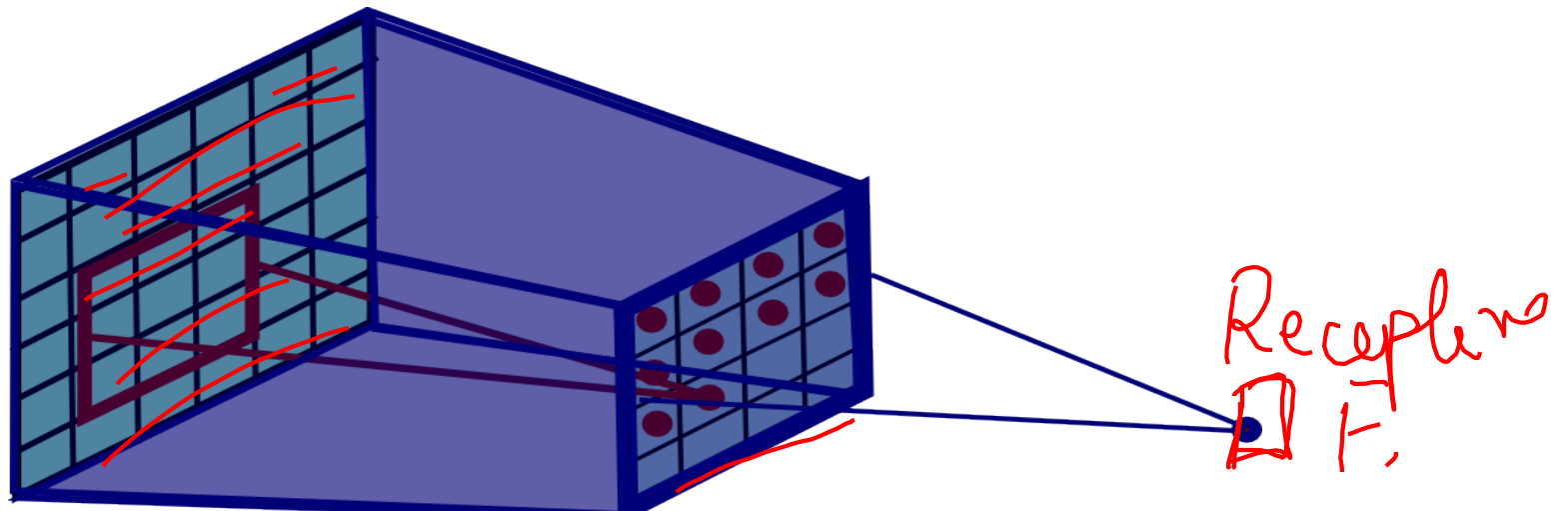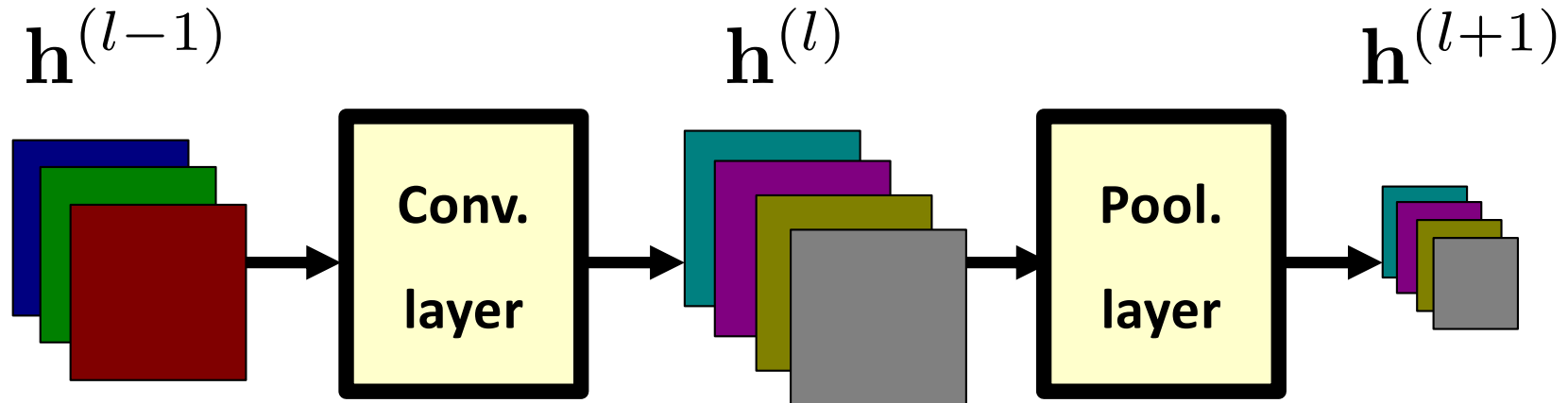
L2-pooling:

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

# Receptive Field



$\mathbf{h}^{(l-1)}$

Conv. layer

$\mathbf{h}^{(l)}$

Pool. layer

$\mathbf{h}^{(l+1)}$

# Pooling Layer: Receptive Field Size

$$\mathbf{h}^{(l-1)} \qquad \mathbf{h}^{(l)} \qquad \mathbf{h}^{(l+1)}$$

# Pooling Layer: Receptive Field Size

$$\mathbf{h}^{(l-1)} \qquad\qquad \mathbf{h}^{(l)} \qquad\qquad \mathbf{h}^{(l+1)}$$



If convolutional filters are FxF and stride 1, and
pooling layer has pools of size PxP,
then each unit in the pooling layer depends upon a patch in $\mathbf{h}^{(l-1)}$ of size:
(P+F-1)x(P+F-1)

- 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

- 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
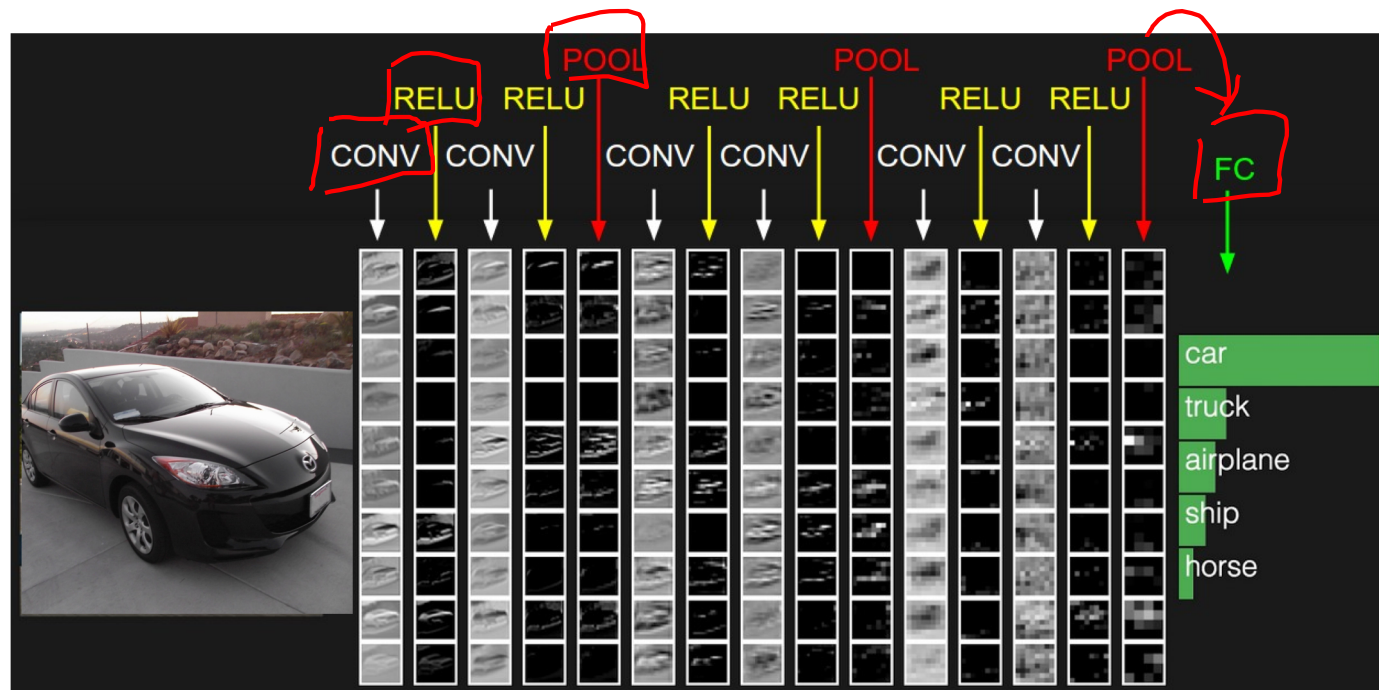- Note that it is not common to use zero-padding for Pooling layers

Common settings:

F = 2, S = 2
F = 3, S = 2

# Fully Connected Layer (FC layer)

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

# Convolutional Neural Networks



$5 \times 5 \times 6 12$

Input layer     (S1) 4 feature maps     (C1) 4 feature maps   (S2) 6 feature maps     (C2) 6 feature maps

convolution layer     sub-sampling layer     convolution layer     sub-sampling layer   fully connected MLP

# Classical View



convolution      fully connected

"tabby cat"

227 × 227    55 × 55    27 × 27      13 × 13

$C_2 \times N^2 C_1$

$C_2$ hidden units

NxNxC$_1$, N small

$N^2 C_1$

(:)

Fully conn. layer

# Classical View



convolution         fully connected

227 × 227    55 × 55    27 × 27    13 × 13    "tabby cat"

# Classical View = Inefficient



warped region

aeroplane? no.

person? yes.

tvmonitor? no.

CNN

1. Input image

2. Extract region proposals (~2k)

3. Compute CNN features

4. Classify regions

# Classical View



convolution      fully connected

227 × 227    55 × 55    27 × 27    13 × 13

"tabby cat"

# Re-interpretation

- Just squint a little!



convolution

227 × 227    55 × 55    27 × 27    13 × 13    1 × 1

Fully conn. layer /
Conv. layer ($C_2$ kernels of size $NxNxC_1$)

# Re-interpretation

- Just squint a little!



convolution

227 × 227   55 × 55   27 × 27   13 × 13   1 × 1

Figure Credit: [Long, Shelhamer, Darrell CVPR15]

# "Fully Convolutional" Networks

- Can run on an image of any size!



convolution

$H \times W$      H/4 × W/4      H/8 × W/8      H/16 × W/16      H/32 × W/32

$2 \times 2 \times C_2$

# Benefit of this thinking

- Mathematically elegant

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

# Plan for Today

- **Convolutional Neural Networks**
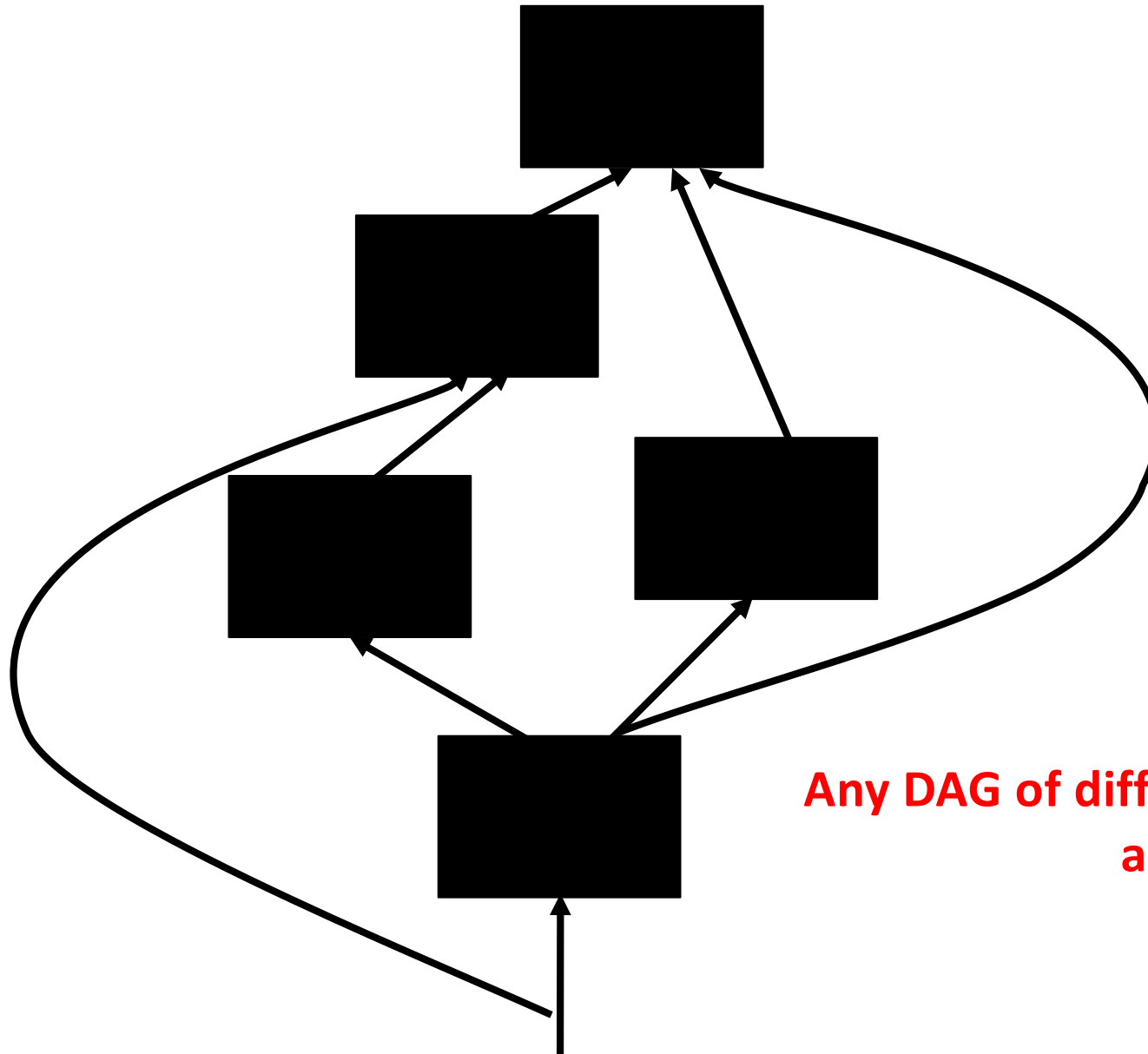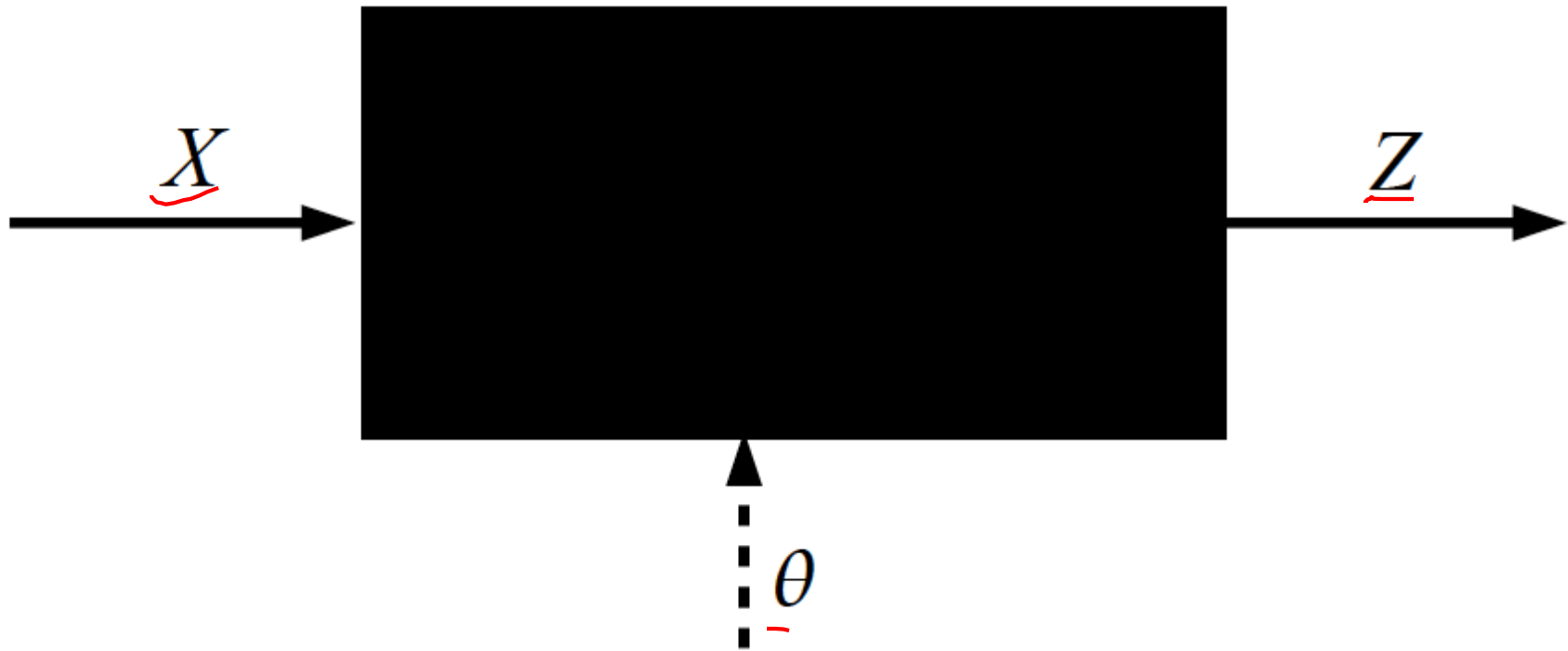    - Pooling layers
    - Fully-connected layers as convolutions
    - Backprop in conv layers [Derived in notes]
    - Toeplitz matrices and convolutions = matrix-mult

# Computational Graph



**Any DAG of differentiable modules is allowed!**
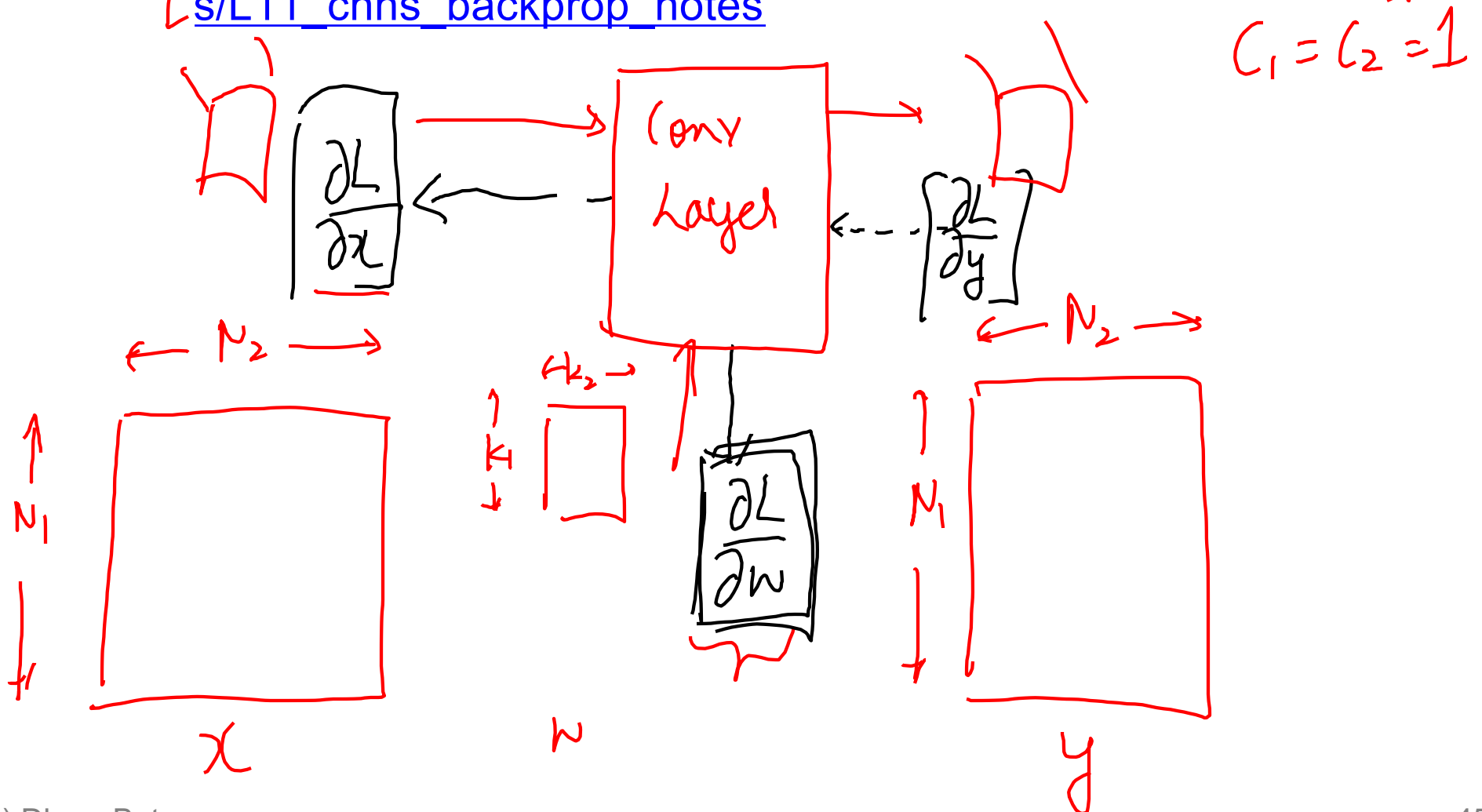
# Key Computation: Forward-Prop

# Key Computation: Back-Prop

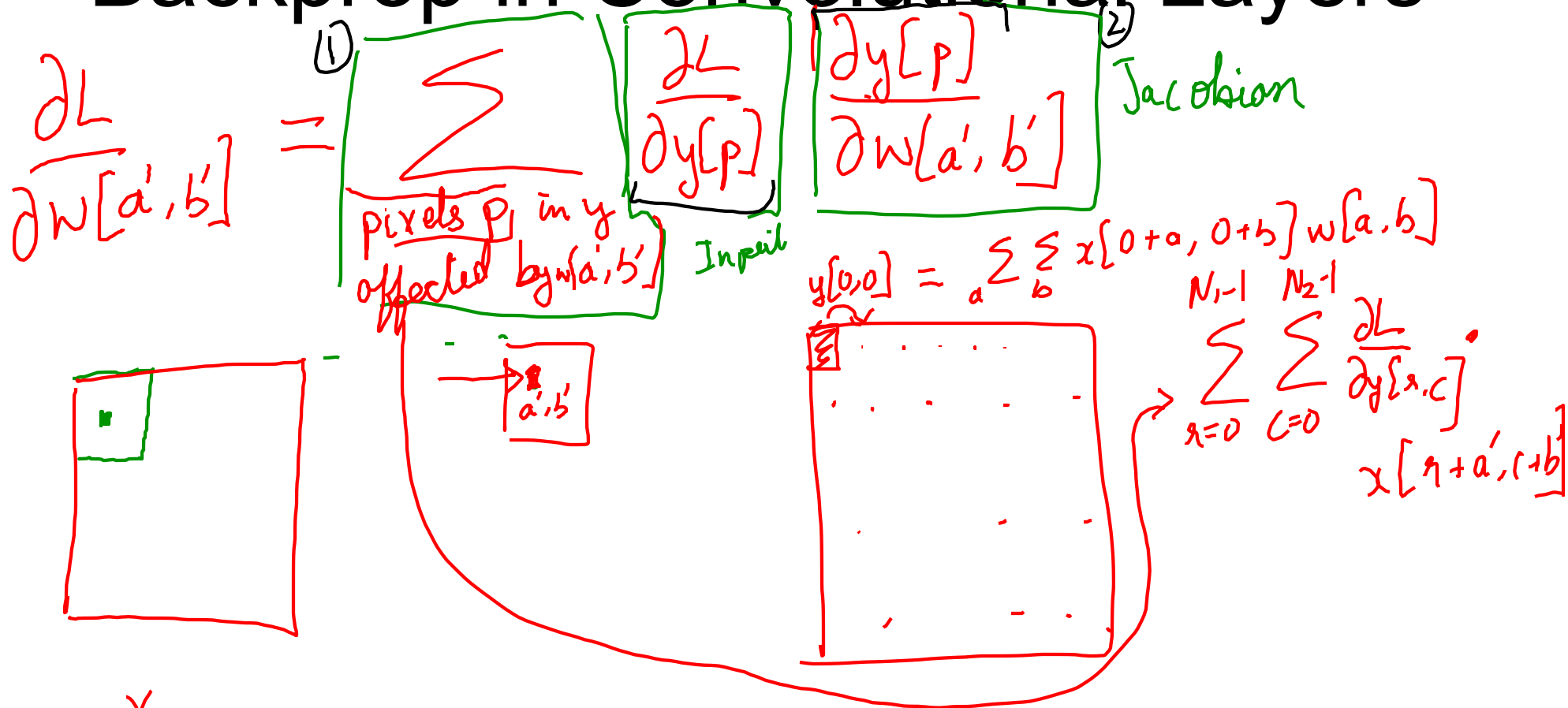Slide Credit: Marc'Aurelio Ranzato, Yann LeCun

# Backprop in Convolutional Layers

- Notes
  - https://www.cc.gatech.edu/classes/AY2020/cs7643_fall/slides/L11_cnns_backprop_notes

# Backprop in Convolutional Layers

$$\frac{\partial L}{\partial W[a', b']} = \underbrace{\sum}_{\substack{\text{pixels } p \text{ in } y \\ \text{affected by } w[a', b']}} \underbrace{\frac{\partial L}{\partial y[p]}}_{\text{Inpul}} \cdot \underbrace{\frac{\partial y[p]}{\partial W[a', b']}}_{\text{Jacobian}}$$

$$y[0,0] = \sum_a \sum_b x[0+a, 0+b]\, w[a,b]$$

$$\sum_{R=0}^{N_1-1} \sum_{C=0}^{N_2-1} \frac{\partial L}{\partial y[R,C]} \cdot x[R+a', C+b']$$

$x$

$W$

$y$

$$\frac{\partial}{\partial W[a',b']} y[R,C] = \sum_{a=0}^{k_1-1} \sum_{b=0}^{k_2-1} x[R+a, C+b]\, w[a,b]$$

$$= x[R+a', C+b']$$

# Backprop in Convolutional Layers

$$\frac{\partial L}{\partial w[a',b']} = \sum_{r=0}^{N_1-1} \sum_{c=0}^{N_2-1} \boxed{\frac{\partial L}{\partial y[r,c]}} x[r+a', c+b']$$

$$\frac{\partial L}{\partial w} = x * \frac{\partial L}{\partial y}$$

# Backprop in Convolutional Layers

# Plan for Today

- **Convolutional Neural Networks**
  - Pooling layers
  - Fully-connected layers as convolutions
  - Backprop in conv layers [Derived in notes]
  - Toeplitz matrices and convolutions = matrix-mult

# Toeplitz Matrix

- Diagonals are constants

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

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

$$A = \begin{bmatrix} a_0 & a_{-1} & a_{-2} & \dots & \dots & a_{-n+1} \\ a_1 & a_0 & a_{-1} & \ddots & & \vdots \\ a_2 & a_1 & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & a_{-1} & a_{-2} \\ \vdots & & \ddots & a_1 & a_0 & a_{-1} \\ a_{n-1} & \dots & \dots & a_2 & a_1 & a_0 \end{bmatrix}$$
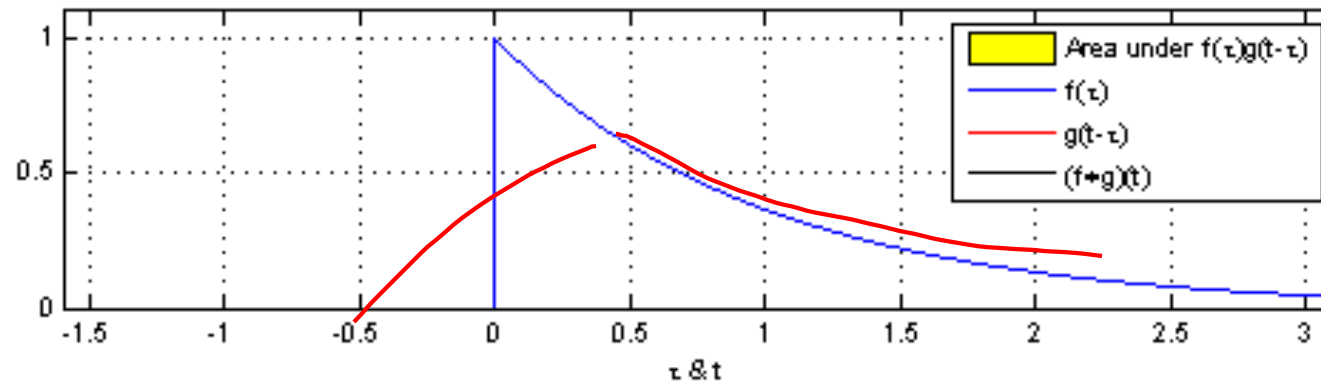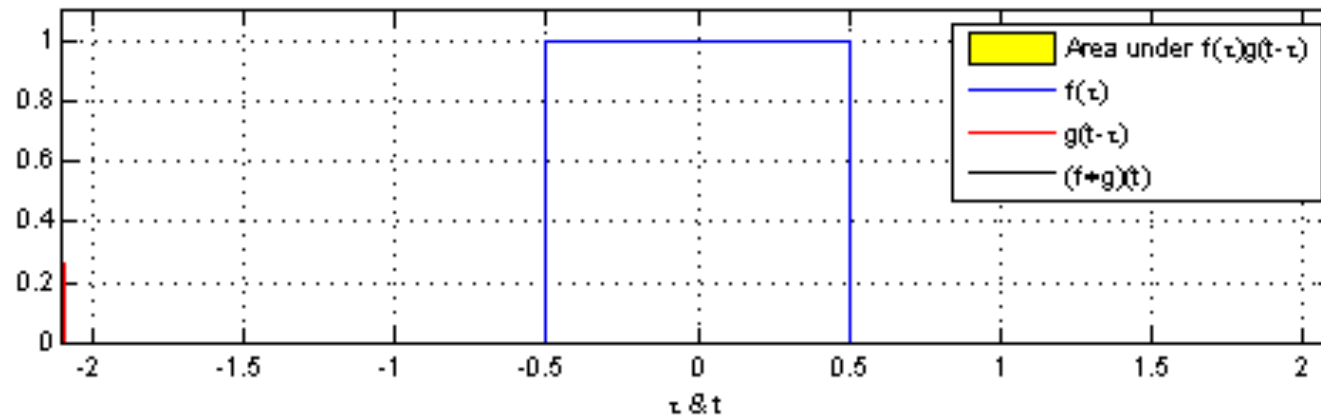
# Why do we care?

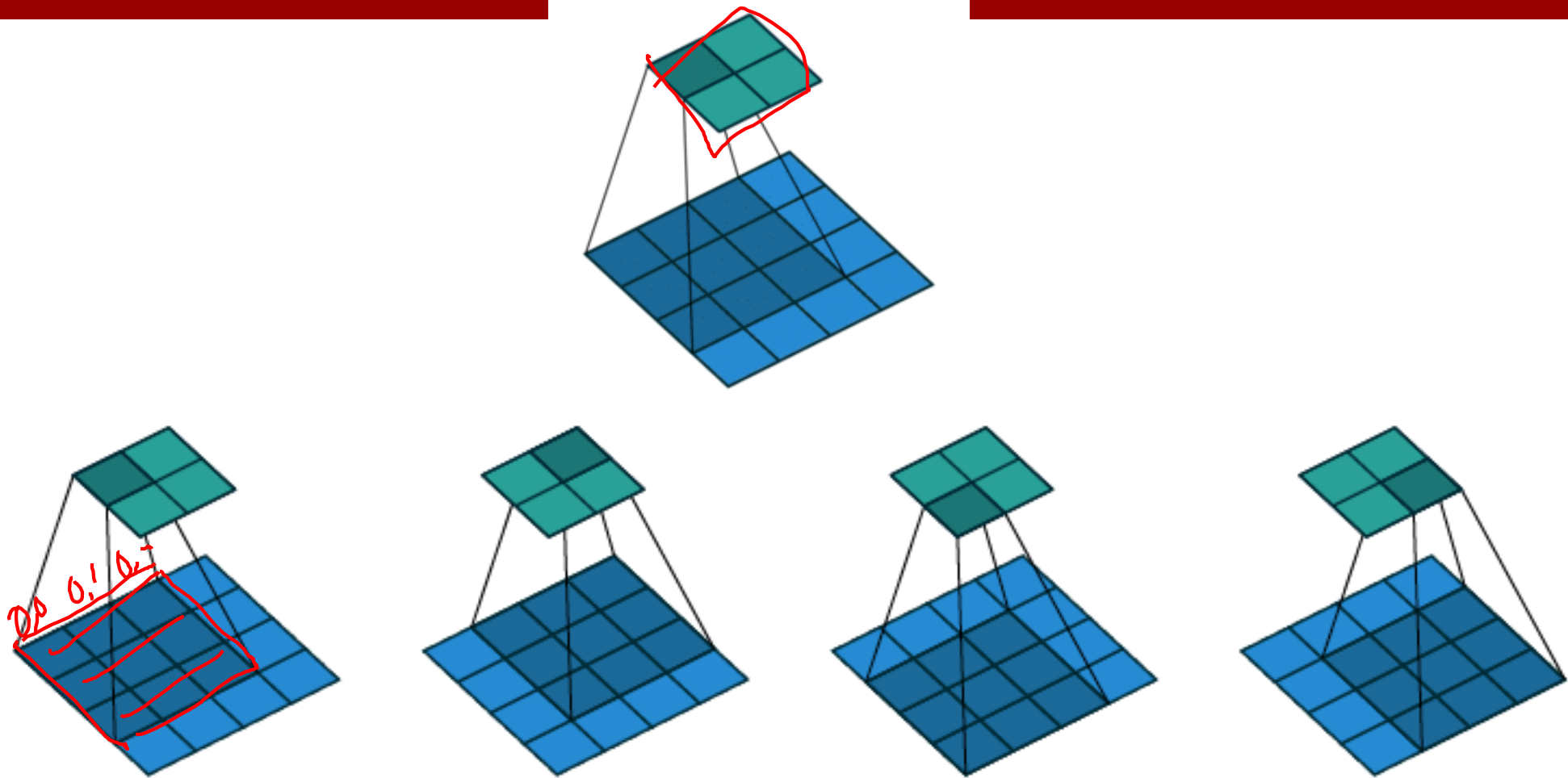- (Discrete) Convolution = Matrix Multiplication
  – with Toeplitz Matrices

$$y = w * x$$

$$
\begin{bmatrix}
w_k & 0 & \ldots & 0 & 0 \\
w_{k-1} & w_k & \ldots & 0 & 0 \\
w_{k-2} & w_{k-1} & \ldots & 0 & 0 \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
w_1 & \ldots & \ldots & w_k & 0 \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
0 & w_1 & \ldots & w_{k-1} & w_k \\
\vdots & \vdots & \vdots & \vdots & \vdots \\
0 & 0 & \vdots & w_1 & w_2 \\
0 & 0 & \vdots & 0 & w_1
\end{bmatrix}
\begin{bmatrix}
x_1 \\ x_2 \\ x_3 \\ \vdots \\ x_n
\end{bmatrix}
$$

"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_with_itself2.gif

$$y\text{-vec} = \begin{pmatrix} w_{0,0} & w_{0,1} & w_{0,2} & 0 & w_{1,0} & w_{1,1} & w_{1,2} & 0 & w_{2,0} & w_{2,1} & w_{2,2} & 0 & 0 & 0 & 0 & 0 \\ 0 & w_{0,0} & w_{0,1} & w_{0,2} & 0 & w_{1,0} & w_{1,1} & w_{1,2} & 0 & w_{2,0} & w_{2,1} & w_{2,2} & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & w_{0,0} & w_{0,1} & w_{0,2} & 0 & w_{1,0} & w_{1,1} & w_{1,2} & 0 & w_{2,0} & w_{2,1} & w_{2,2} & 0 \\ 0 & 0 & 0 & 0 & 0 & w_{0,0} & w_{0,1} & w_{0,2} & 0 & w_{1,0} & w_{1,1} & w_{1,2} & 0 & w_{2,0} & w_{2,1} & w_{2,2} \end{pmatrix}$$

# So far: Image Classification
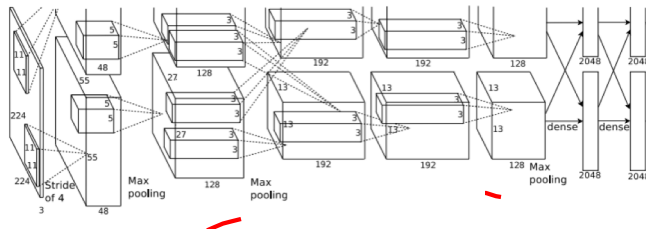


This image is CC0 public domain

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**
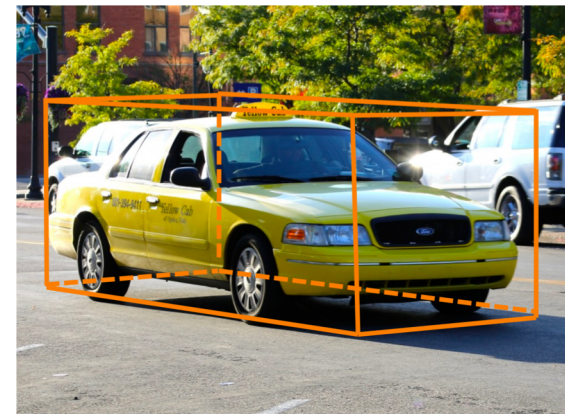
No objects, just pixels

## 2D Object Detection



**DOG**, **DOG**, **CAT**

Object categories +
2D bounding boxes

## 3D Object Detection



**Car**

Object categories +
3D bounding boxes

Slide Credit: Fei-Fei Li, Justin Johnson, Serena Yeung, CS 231n