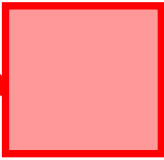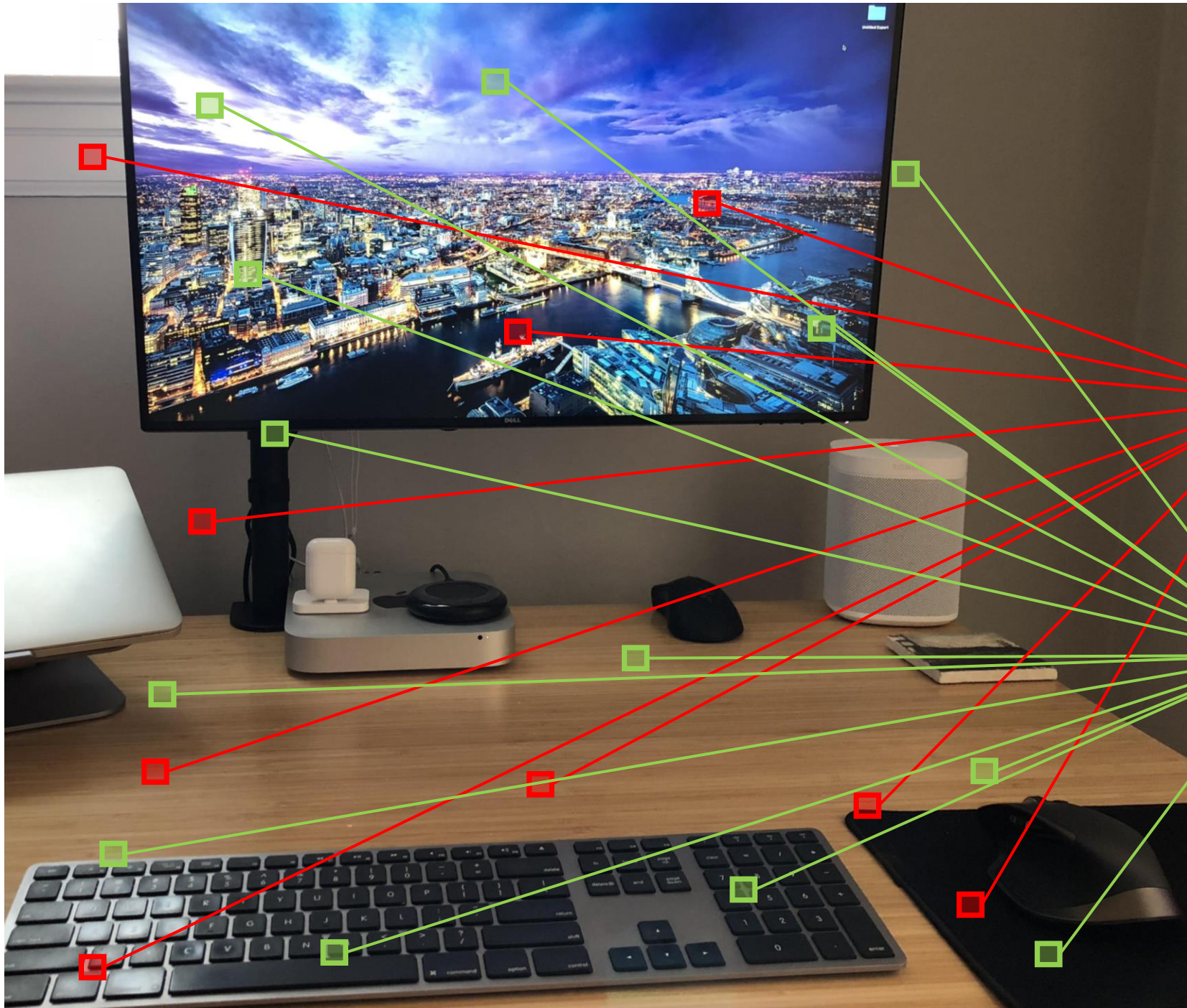# Lecture 3: Convolutional Neural Networks

# Today:

- Expectations of visual recognition systems     (15%)

- A Conv layer, stride, padding etc.            (25%)

- Conv2D                                        (20%)

- Transposed Conv                               (5%)

- Implementation and backprop                   (15%)

- What is encoded in feature maps?              (10%)

- Max-pooling, convnet and conv variants        (10%)

Monitor

Desk

Mouse

vertical line

keyboard

key key key key key key key key
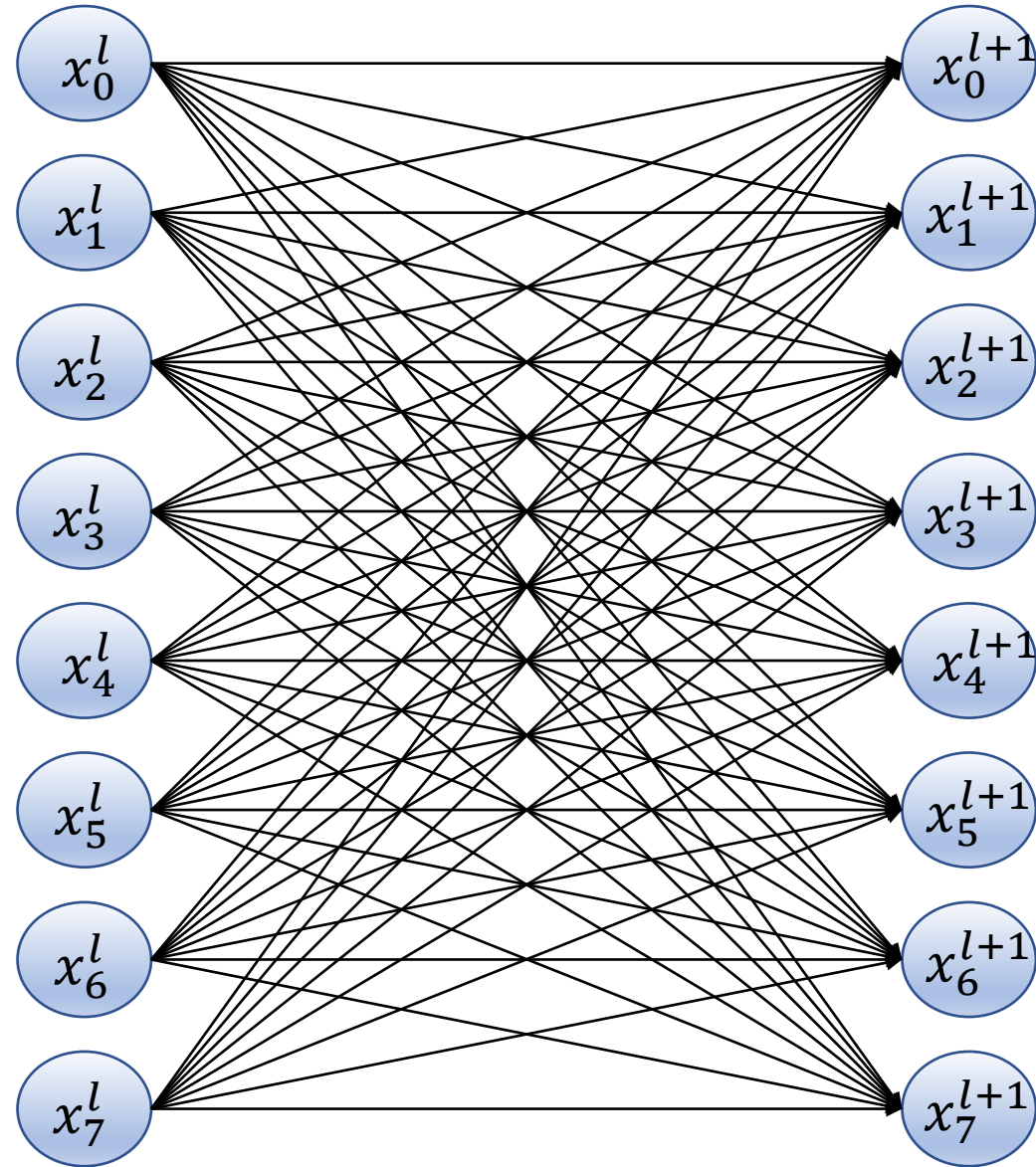
horizontal line

key key

# Expectations of visual recognition network

1. Maintain 2D structure logic

2. Shift invariant (actually, equivariant)

3. Consider only local correlations

4. Hierarchically growing field of view

5. Hierarchically progressing complexity
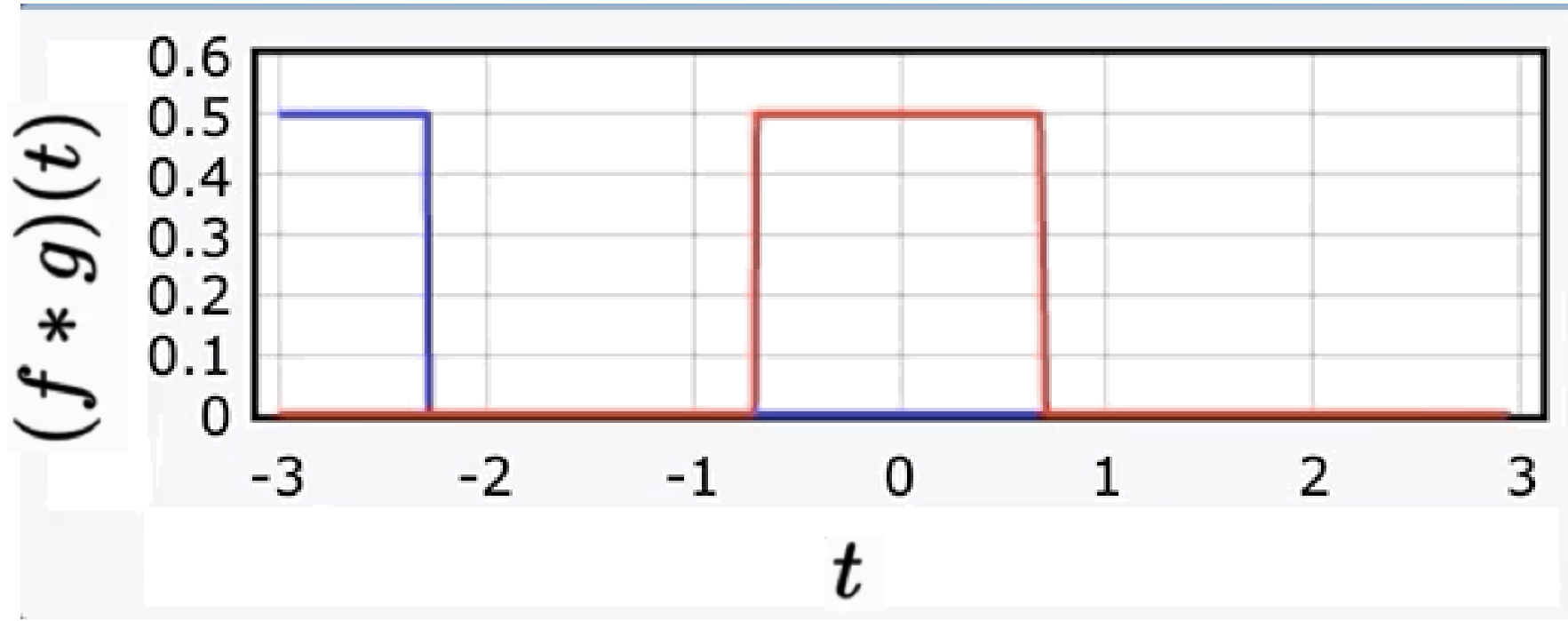
6. Reasonable amount of params

# Fully connected layer

$$
\begin{bmatrix}
w_{00}^l & w_{01}^l \\
w_{10}^l & w_{11}^l
\end{bmatrix}
\cdots
$$
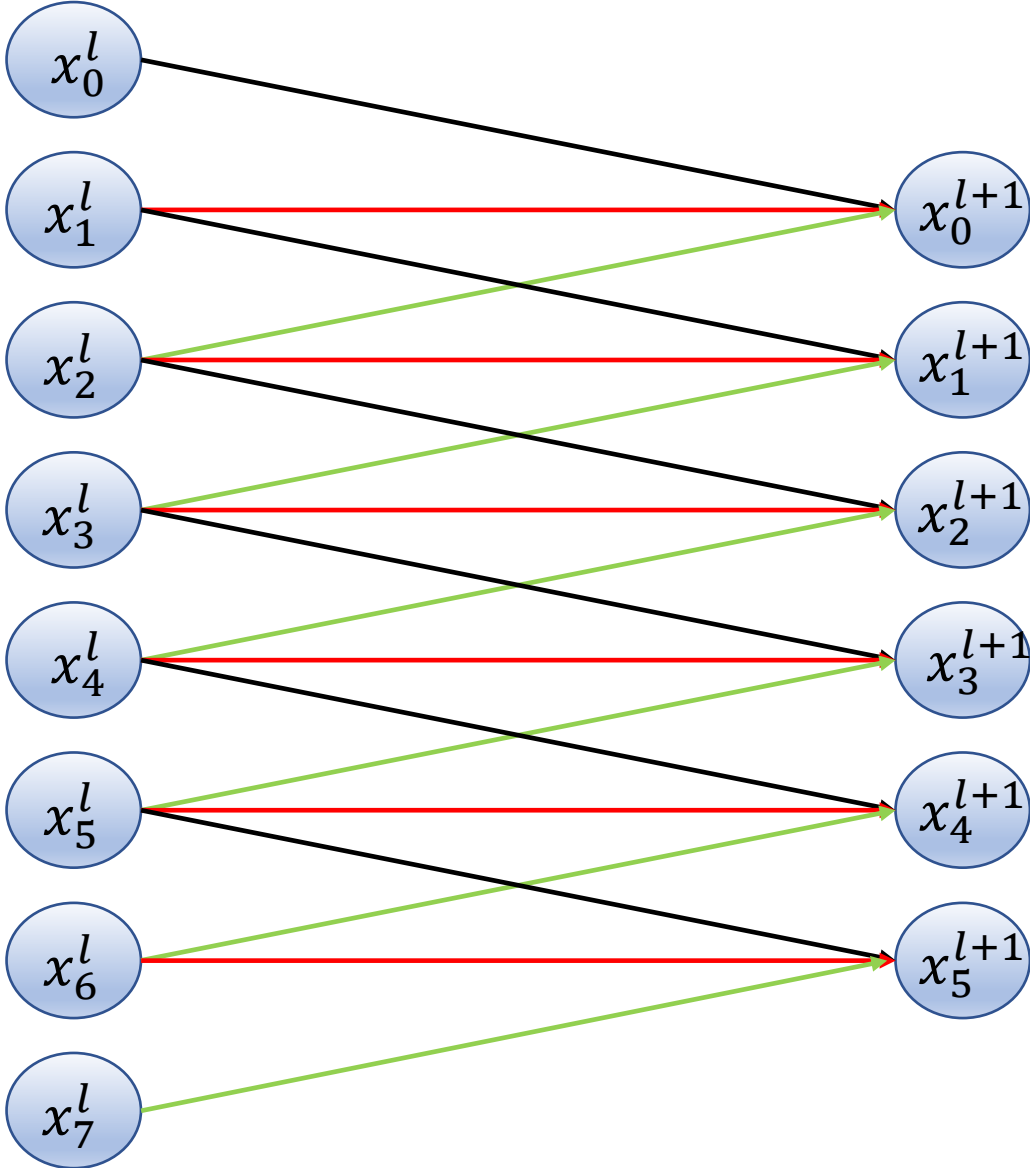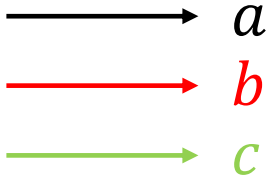
❌ 1. Maintain 2D structure logic
❌ 2. Shift invariant (actually, equivariant)
❌ 3. Consider only local correlations
❌ 4. Hierarchically growing field of view
✅ 5. Hierarchically progressing complexity
❌ 6. Reasonable amount of params
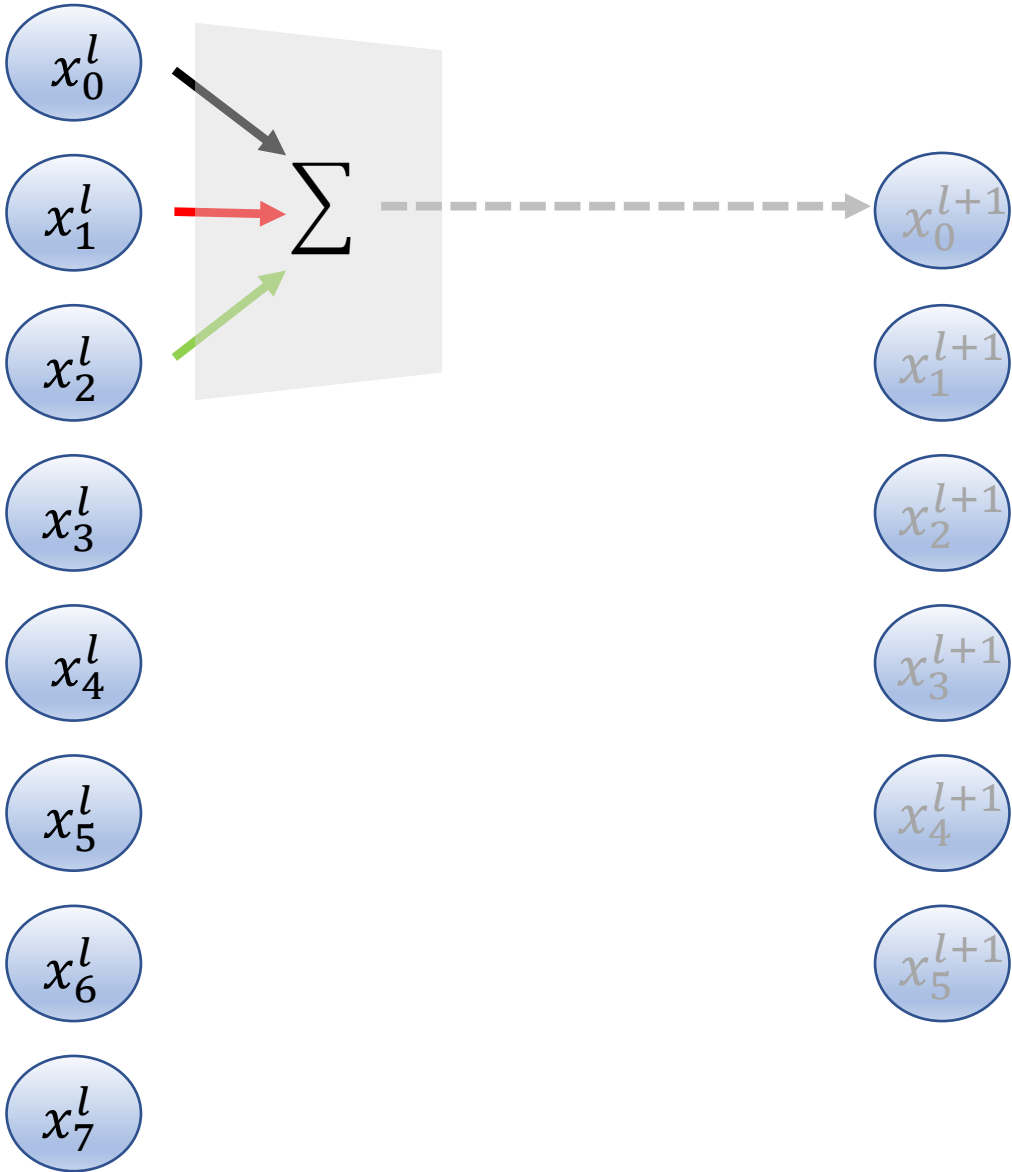
# Convolution



$$(f * g)(t) \stackrel{\text{def}}{=} \int_{-\infty}^{\infty} f(\tau)\, g(t - \tau)\, d\tau$$

# Convolution layer



Toeplitz matrix

# Convolution Filter

$x_0^l$

$x_1^l$

$x_2^l$

$x_3^l$

$x_4^l$

$x_5^l$

$x_6^l$

$x_7^l$

$\sum$

$x_0^{l+1}$

$x_1^{l+1}$

$x_2^{l+1}$

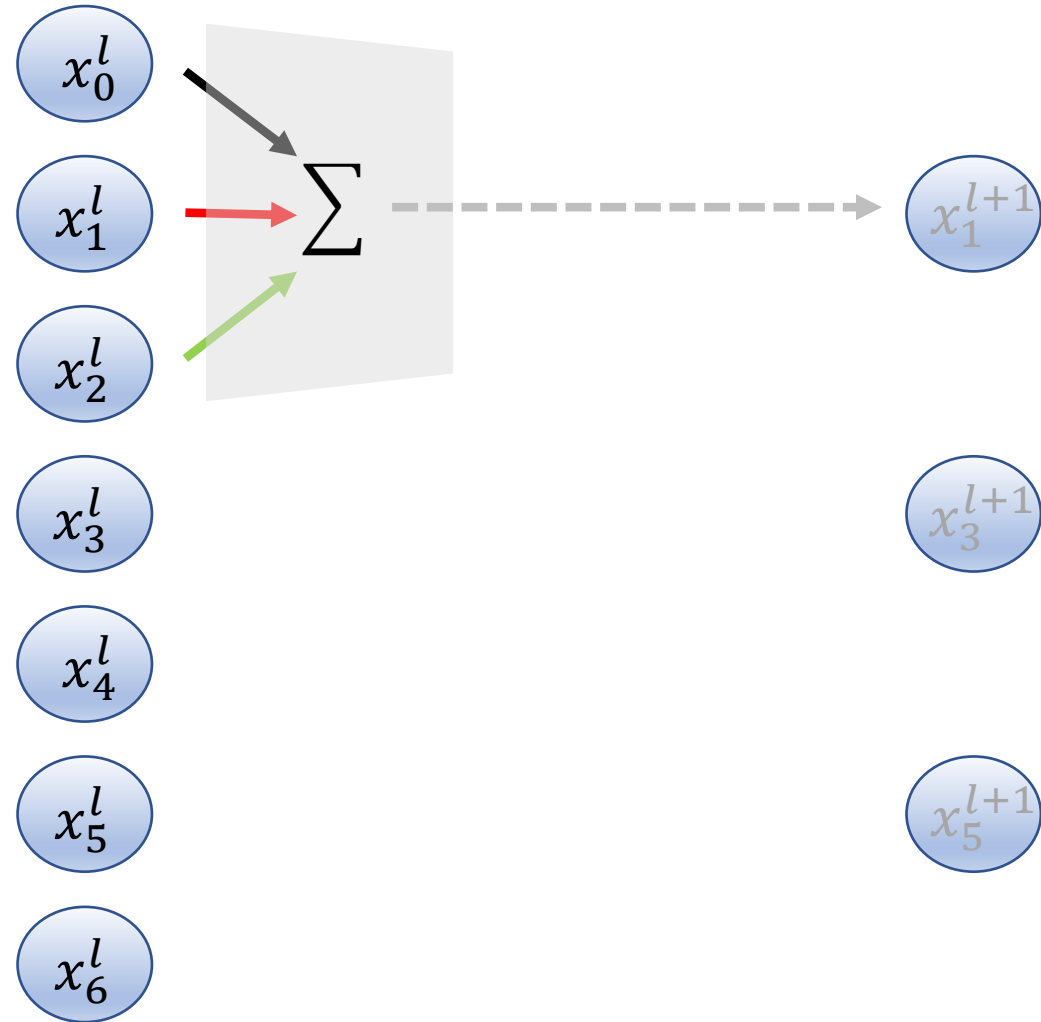$x_3^{l+1}$

$x_4^{l+1}$

$x_5^{l+1}$

Q: Given input size and filter size, find output size.

Q: is this a convolution?

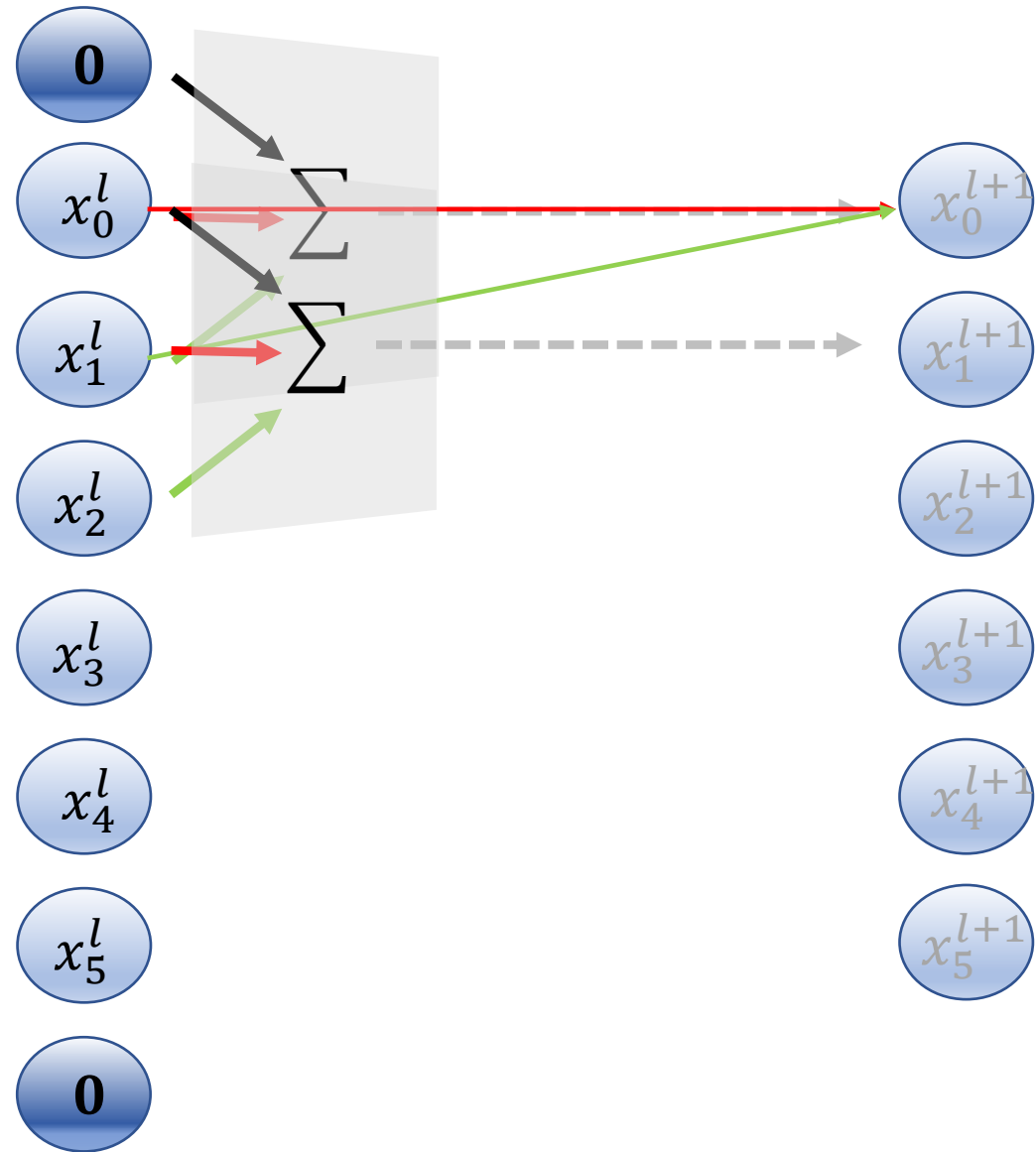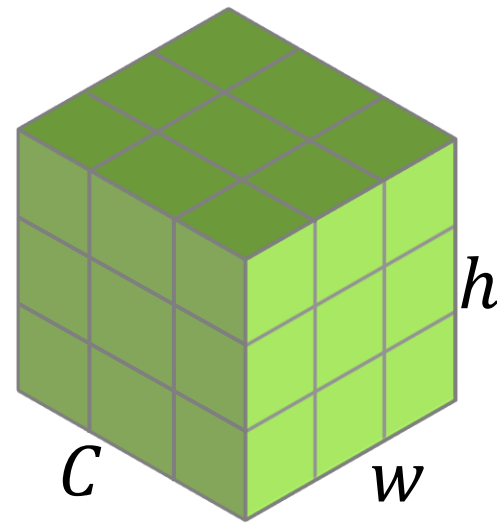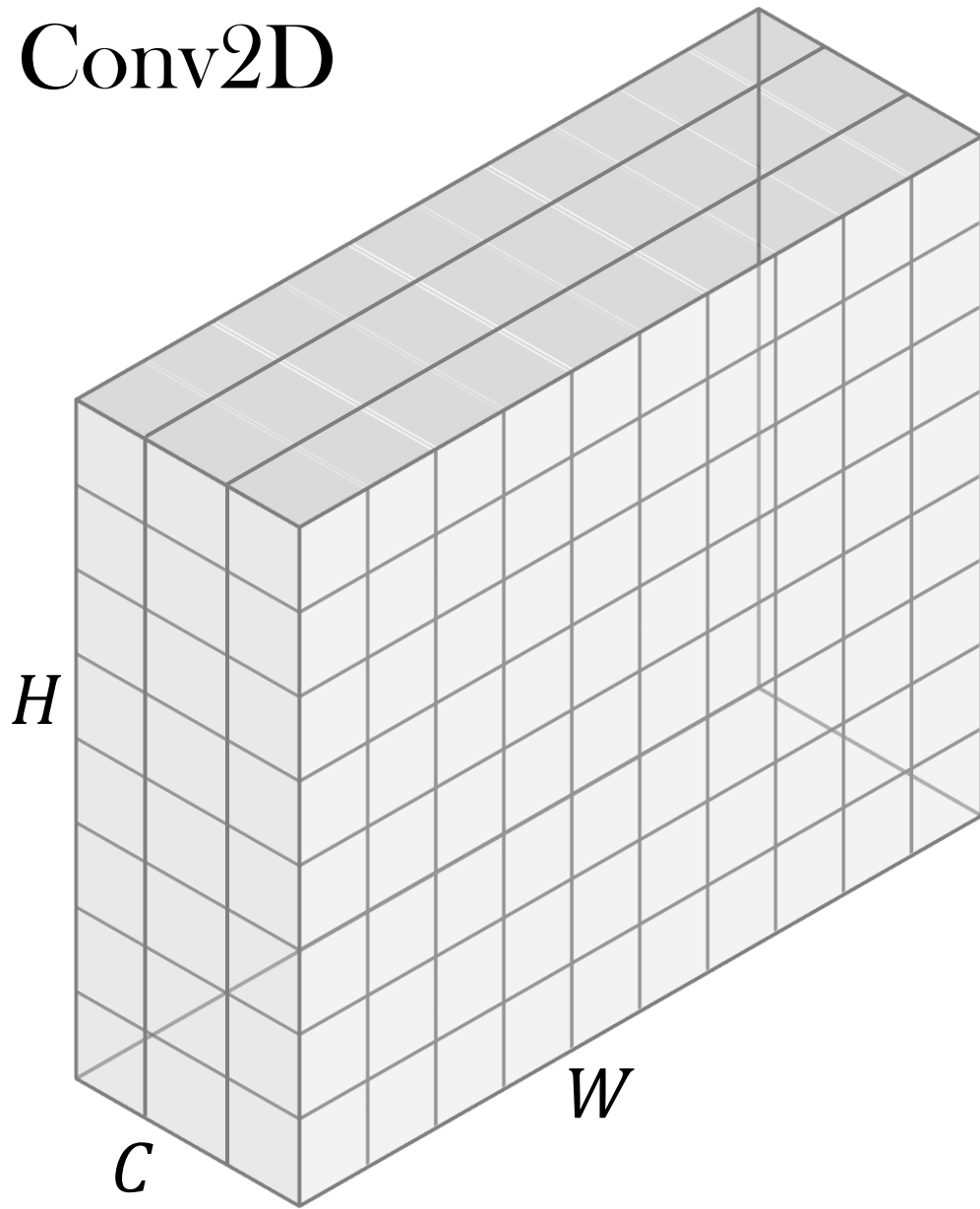A: Yes, but with the flipped filter. This is cross-correlation.

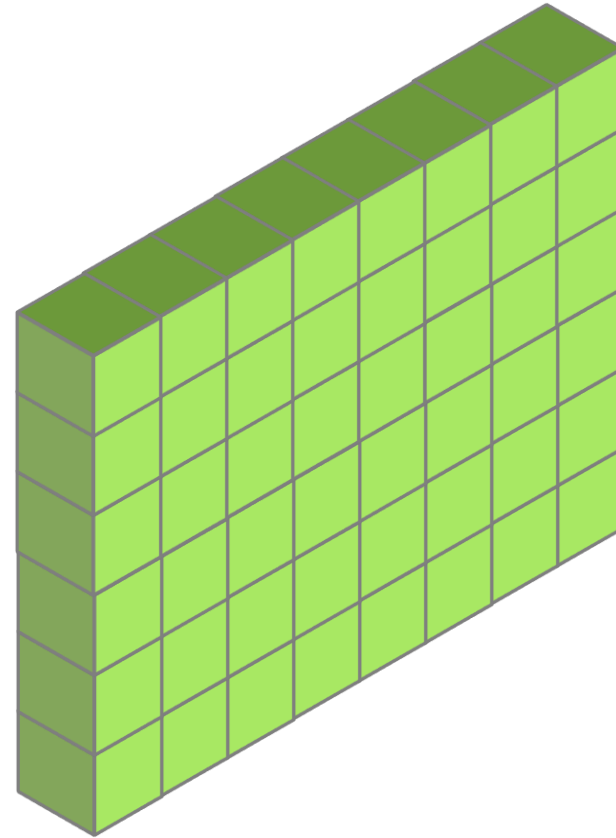# Stride

# Padding
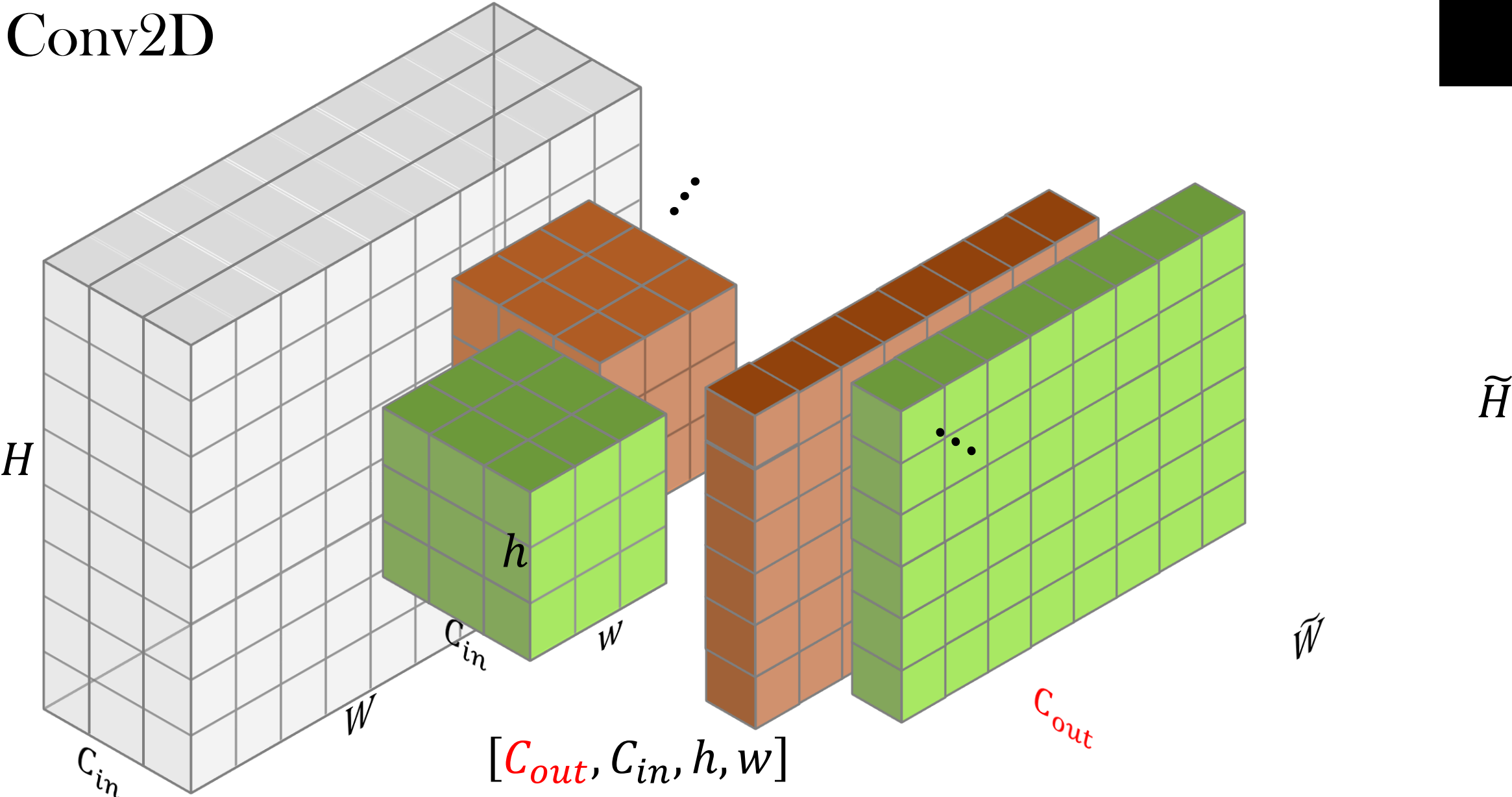
# Conv2D



$N, C, H, W$

$C, h, w$

# Conv2D



$$N, C, H, W$$

$$[N, 1, \widetilde{H}, \widetilde{W}]$$

# Conv2D



$[N, C_{in}, H, W]$

$[C_{out}, C_{in}, h, w]$

$[N, C_{out}, \widetilde{H}, \widetilde{W}]$

DL4CV@Weizmann
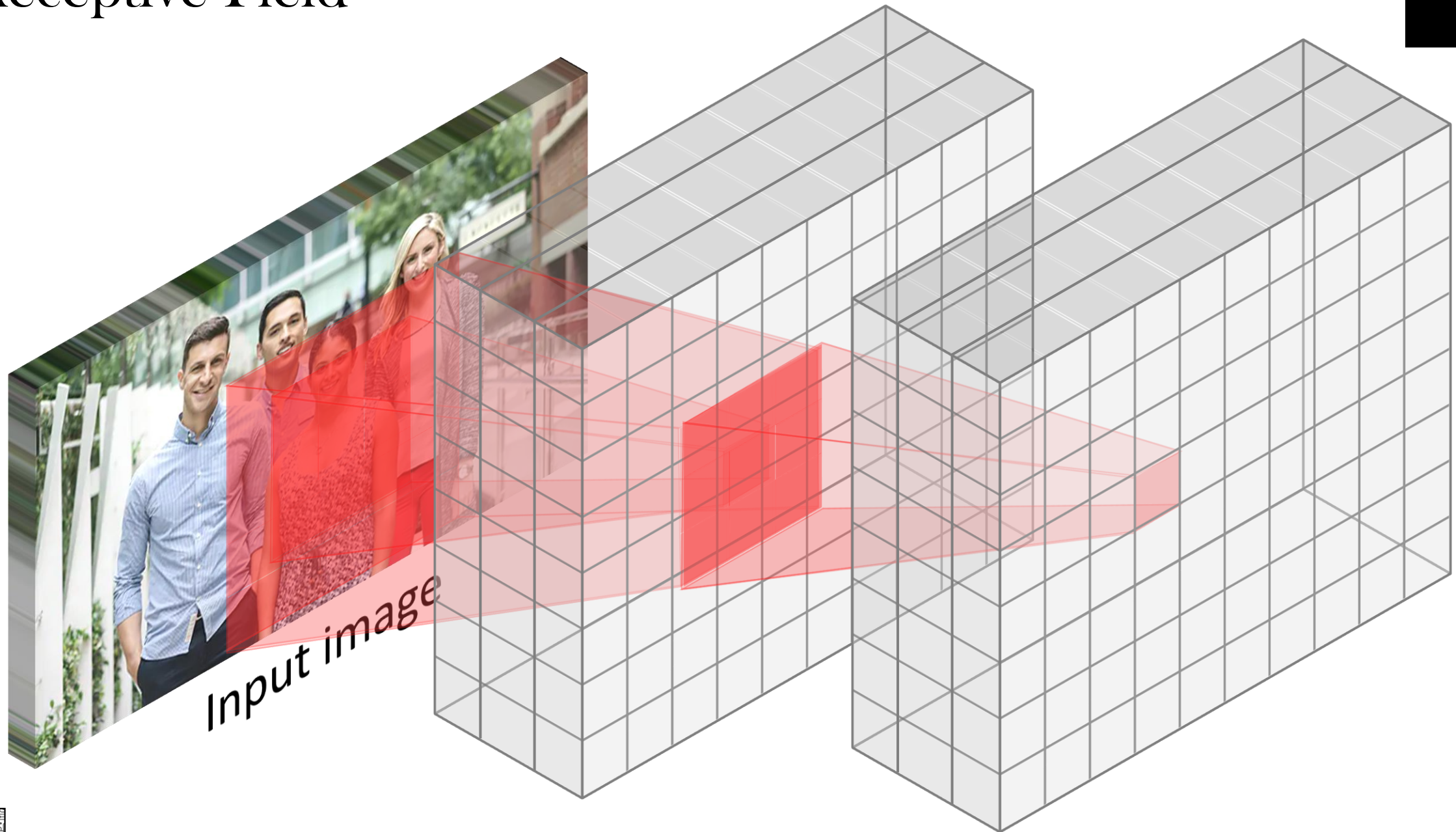
Receptive Field $[Filter\ size: 3]$

Input image

# Convs rock!

✓ 1. Maintain 2D structure logic

✓ 2. Shift invariant (actually, equivariant)

✓ 3. Consider only local correlations

✓ 4. Hierarchically growing field of view

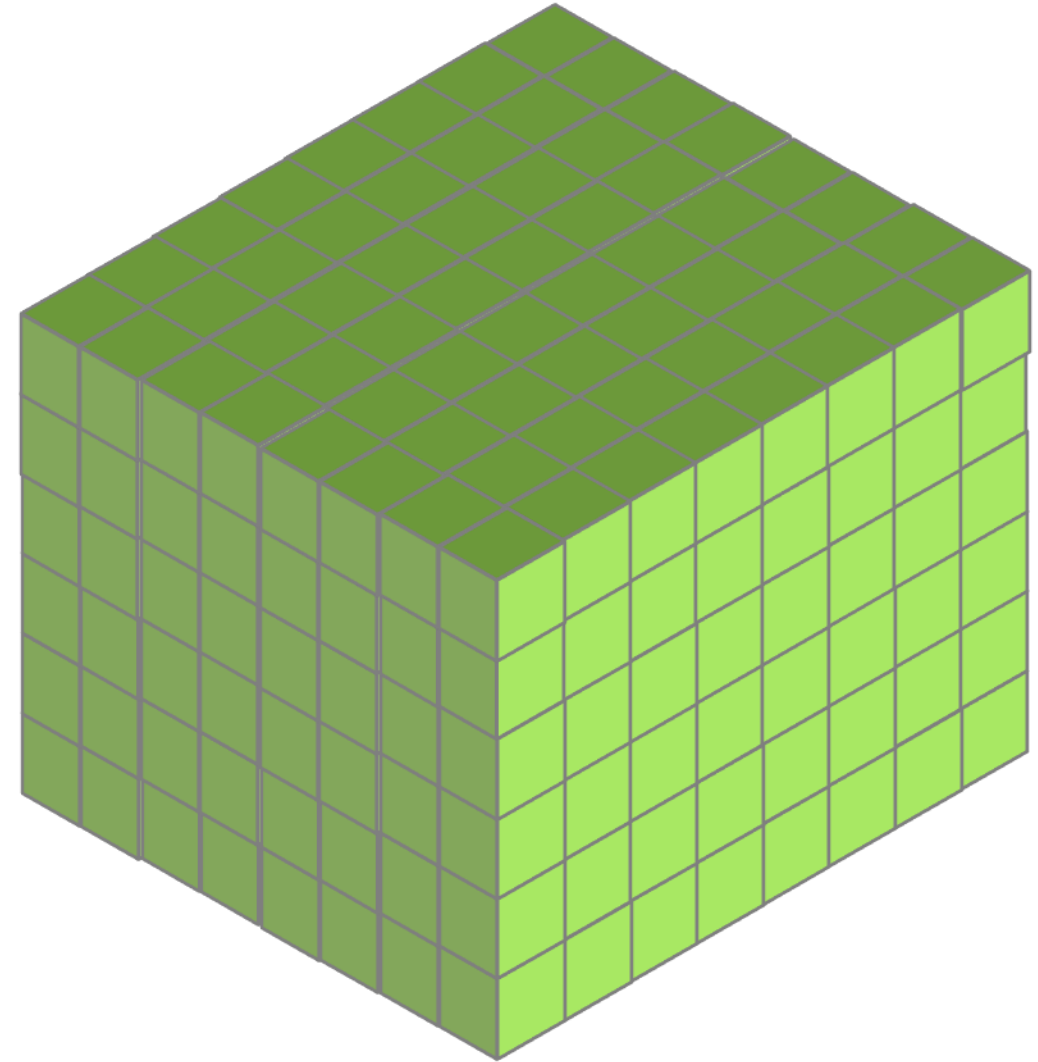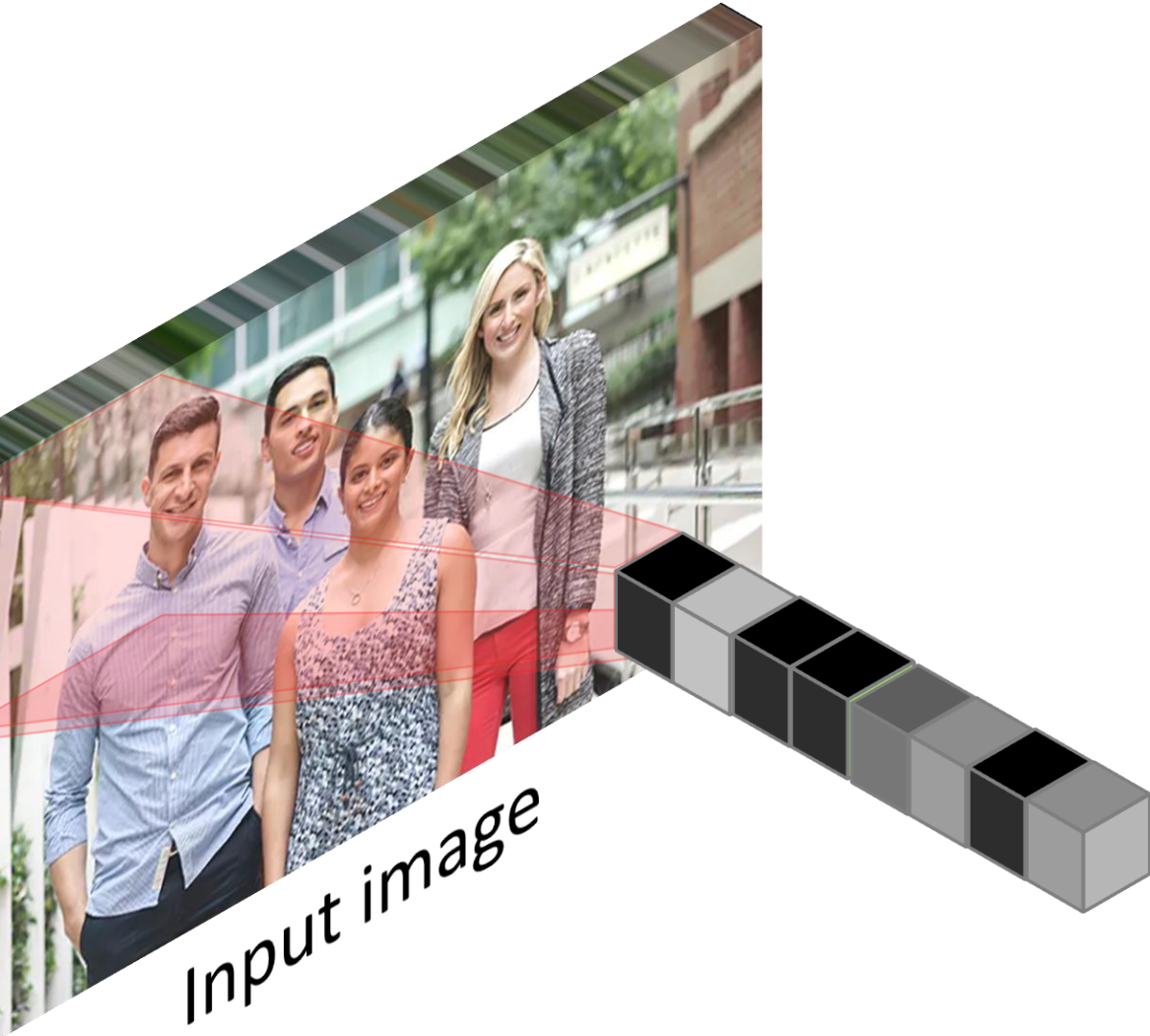✓ 5. Hierarchically progressing complexity

✓ 6. Reasonable amount of params
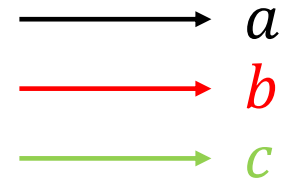
# Two important intuitions about feature maps



Input image

$H$

$C$

$W$

Deep layer feature-map

# Two important intuitions about feature maps



Input image

Deep layer feature-map

# Transposed Convolution



$$\begin{bmatrix} \boldsymbol{b} & \boldsymbol{a} & 0 & 0 & 0 & 0 & 0 \\ \boldsymbol{c} & \boldsymbol{b} & \boldsymbol{a} & 0 & 0 & 0 & 0 \\ 0 & \boldsymbol{c} & \boldsymbol{b} & \boldsymbol{a} & 0 & 0 & 0 \\ 0 & \boldsymbol{a} & \boldsymbol{c} & \boldsymbol{b} & \boldsymbol{a} & 0 & 0 \\ 0 & 0 & \boldsymbol{a} & \boldsymbol{c} & \boldsymbol{b} & \boldsymbol{a} & 0 \\ 0 & 0 & 0 & 0 & \boldsymbol{c} & \boldsymbol{b} & \boldsymbol{a} \\ 0 & 0 & 0 & 0 & 0 & \boldsymbol{c} & \boldsymbol{b} \end{bmatrix}$$

# Transposed Convolution with stride



**Transposed Conv**

$\xrightarrow{\hspace{2cm}} a$

$\xrightarrow{\hspace{2cm}} b$ (red)

$\xrightarrow{\hspace{2cm}} c$ (green)

Recall- stride 2 conv:

$$\begin{bmatrix} a & b & c & 0 & 0 & 0 & 0 \\ 0 & 0 & a & b & c & 0 & 0 \\ 0 & 0 & 0 & 0 & a & b & c \end{bmatrix}$$

# Transposed Convolution by dilation & flip

# A note about the implementation of conv



$$[C_{out}, C_{in}, h, w]$$

*Unfold*

*BMM*

*Reshape*

*Fold* (1x1)

$$[N, C_{in}, H, W]$$

$$[C_{out}, (C_{in} \cdot h, \cdot w)]$$

$$[N, (\widetilde{H} \cdot \widetilde{W}), C_{out}]$$

$$[N, (\widetilde{H} \cdot \widetilde{W}), (C_{in} \cdot h \cdot w)]$$

$$[N, C_{out}, \widetilde{H}, \widetilde{W}]$$

# Q: How do you backprop a Conv? (1D)

For derivative of loss w.r.t all neurons:
Transposed Conv with the same filter!

For derivative of loss w.r.t filter weights:
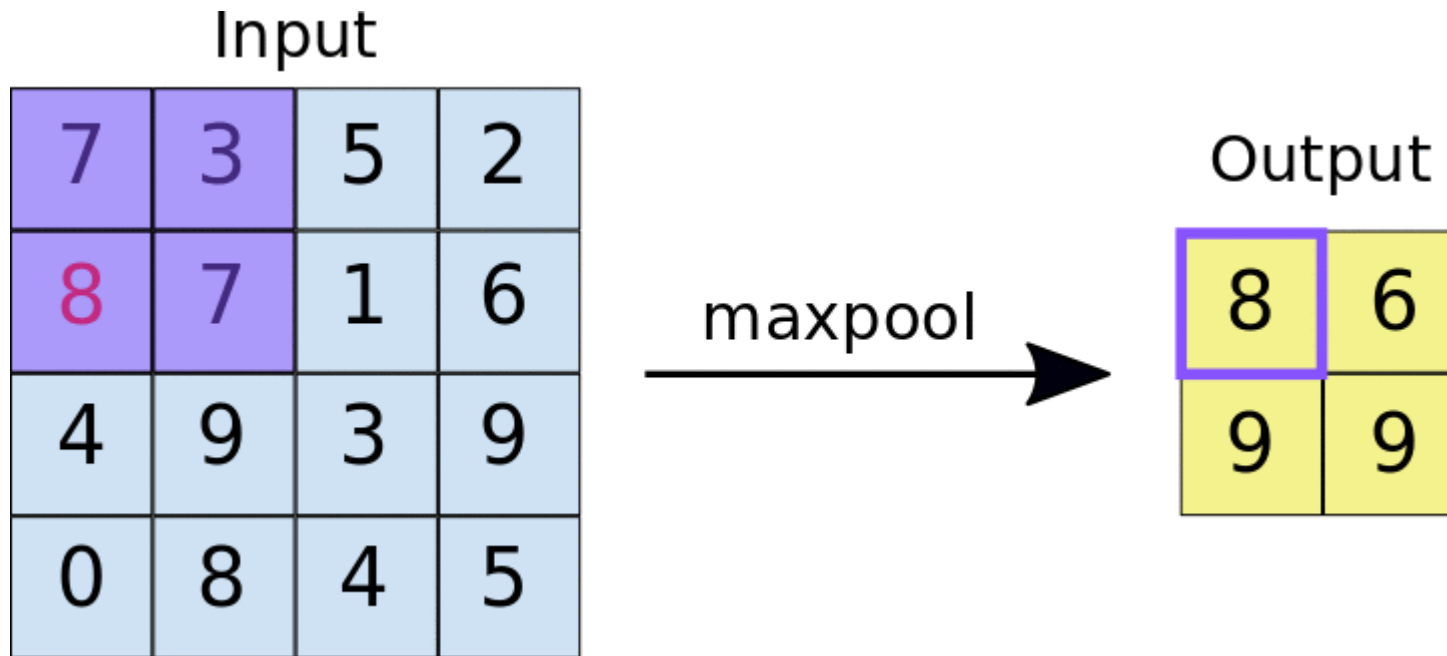Multiply conv-output-grad with the neurons it "sees", sum over all locations.

# Q: How about Conv2D?

For derivative of loss w.r.t all neurons:
1. Transpose dims of filter $c_{in}$ , $c_{out}$ (as in FC)
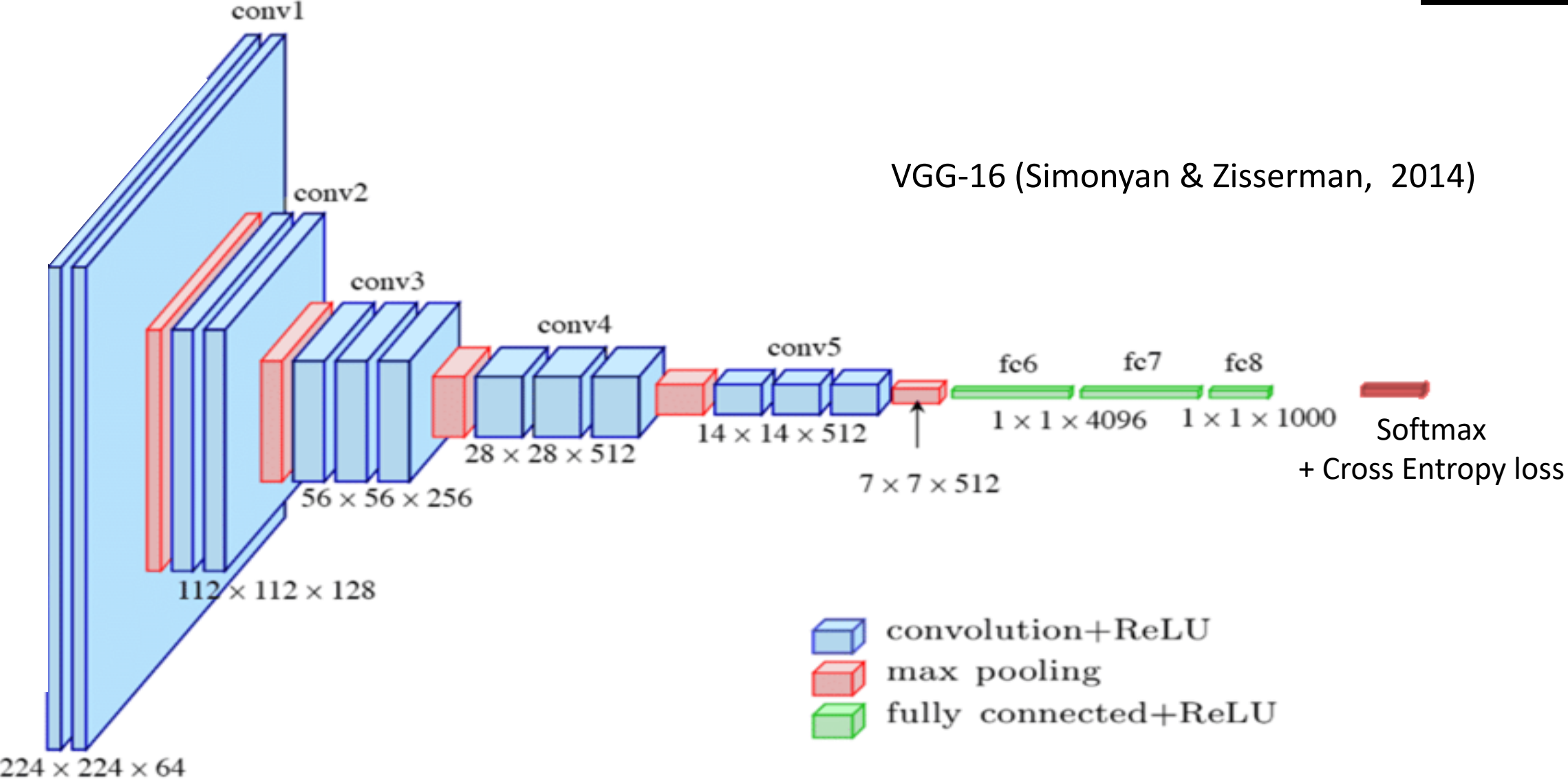2. Transposed Conv2D with the modified filter!

For derivative of loss w.r.t filter weights:
Outer-prod of conv-output-grad with the neurons it "sees", sum over all locations.
(You need to unfold conv-input like in forward, and conv-out-grad with 1x1)

# Max Pooling



- Usually stride=win-size, but not always.
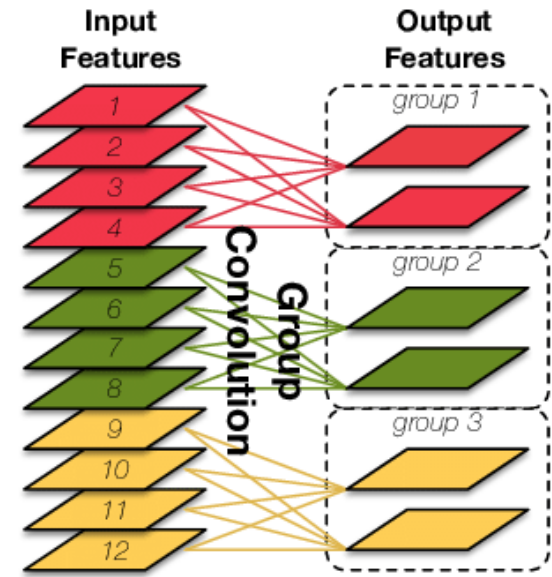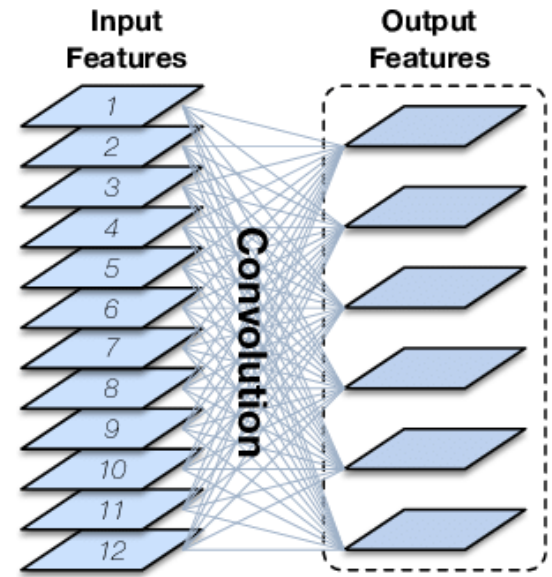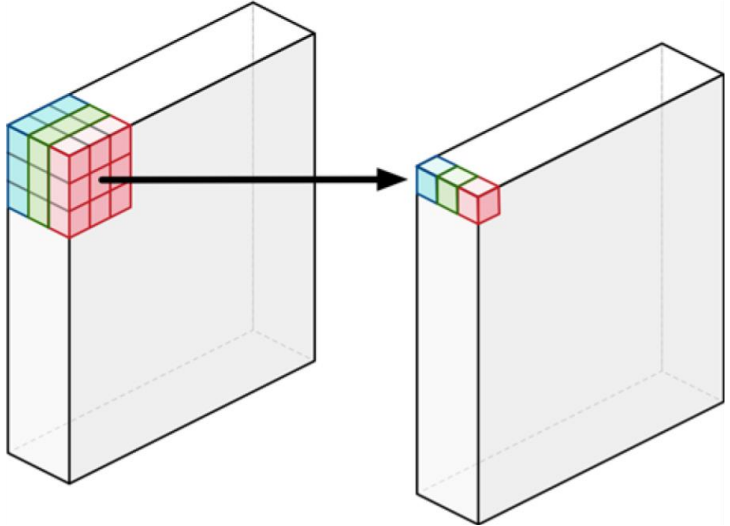- Each channel separately.

# ConvNet Example



VGG-16 (Simonyan & Zisserman, 2014)

conv1

conv2

conv3

conv4

conv5

fc6  fc7  fc8

$224 \times 224 \times 64$

$112 \times 112 \times 128$

$56 \times 56 \times 256$

$28 \times 28 \times 512$

$14 \times 14 \times 512$

$7 \times 7 \times 512$

$1 \times 1 \times 4096$

$1 \times 1 \times 1000$

Softmax
+ Cross Entropy loss

convolution+ReLU

max pooling

fully connected+ReLU

WAIC

# More special Convs!

## Group Conv (Krihzevsky 2012)



## Depthwise Conv



## Dilated Conv (Yu&Koltun 2016)



dilation=1          dilation=2          dilation=3
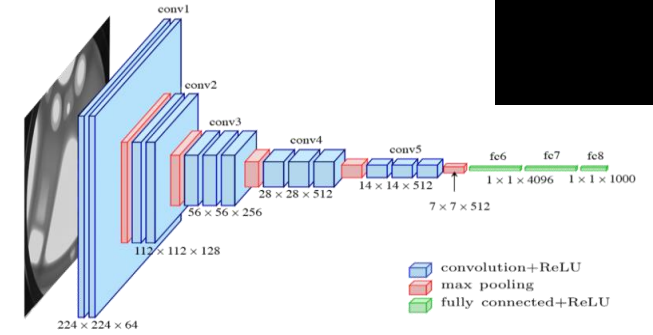
This week's tutorial:
# CNN Architectures



Dror Moran

Next week's lecture:

# Practical Training

Shai Bagon

My Neural Net

Initialization
Augmentation
Regularization
DropOut
BatchNorm

Hyperparam tuning