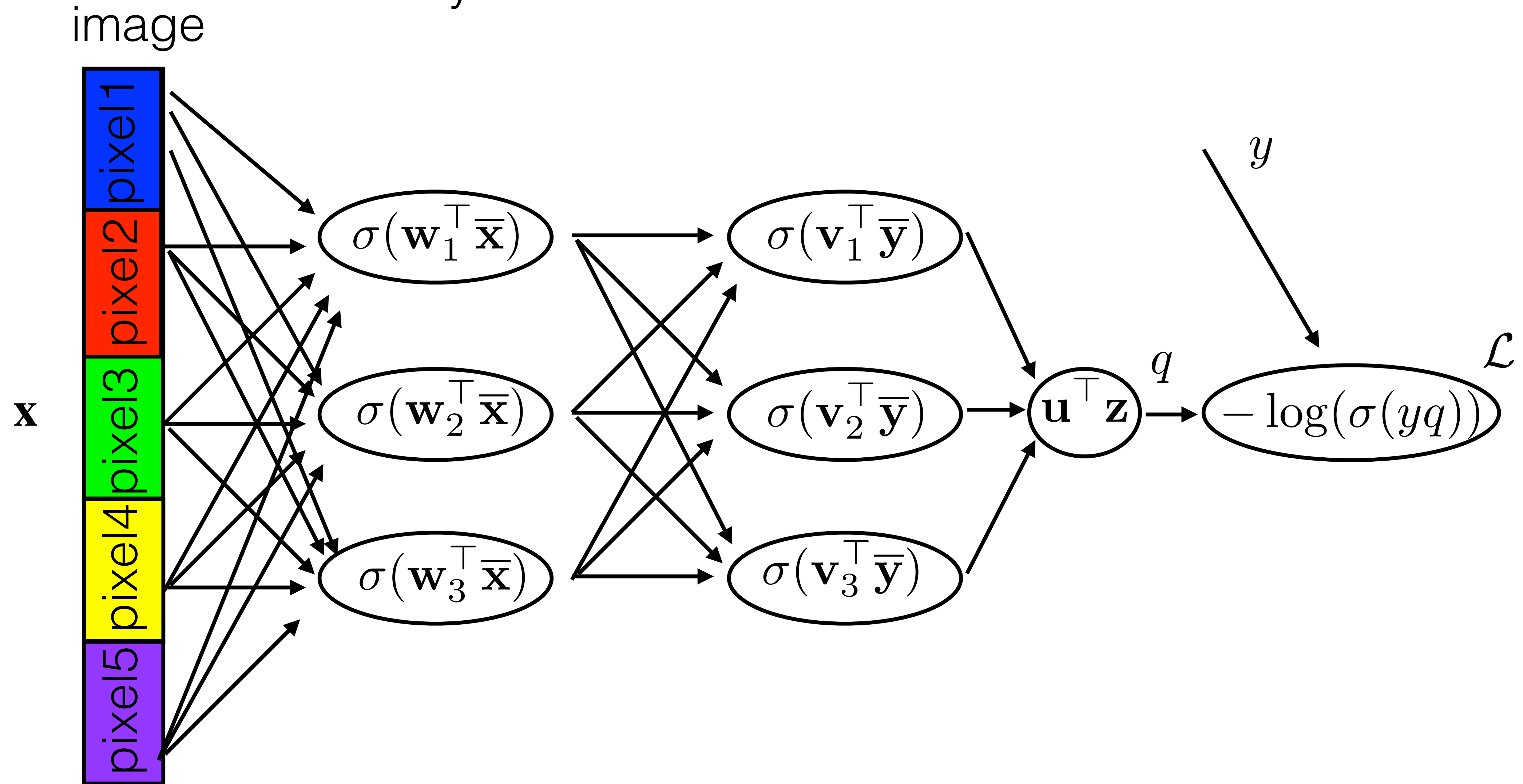# The story of the cat's brain surgery

**cortex, convolution layer, its vector-Jacobian product, feature maps, low-dimensional encoding, and fun with pre-trained convnet**

**Karel Zimmermann**
**Czech Technical University in Prague**
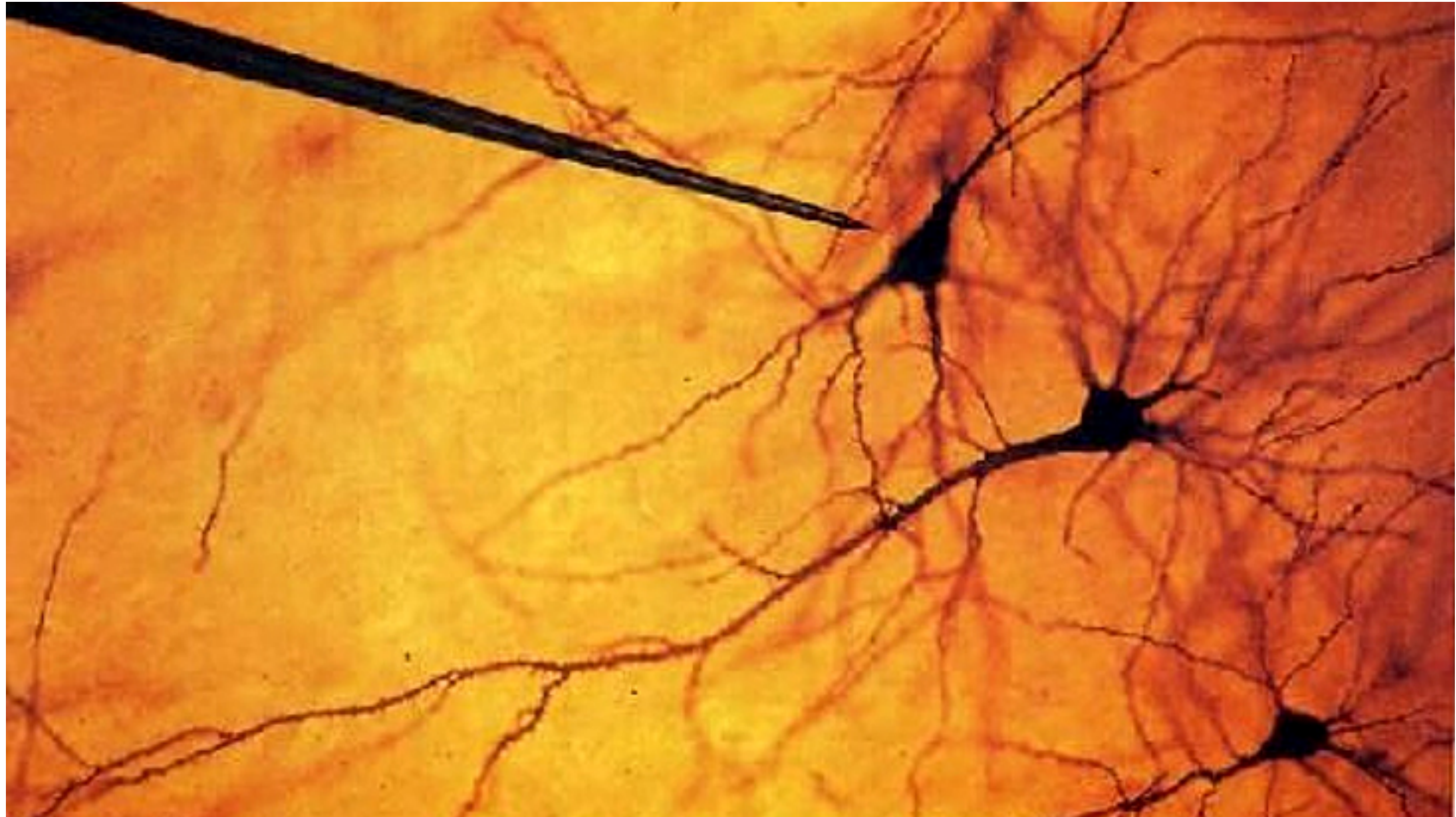**Faculty of Electrical Engineering, Department of Cybernetics**

Fully connected neural network

Learning prone to overfitting, the structure is too general, the resulting function is wild

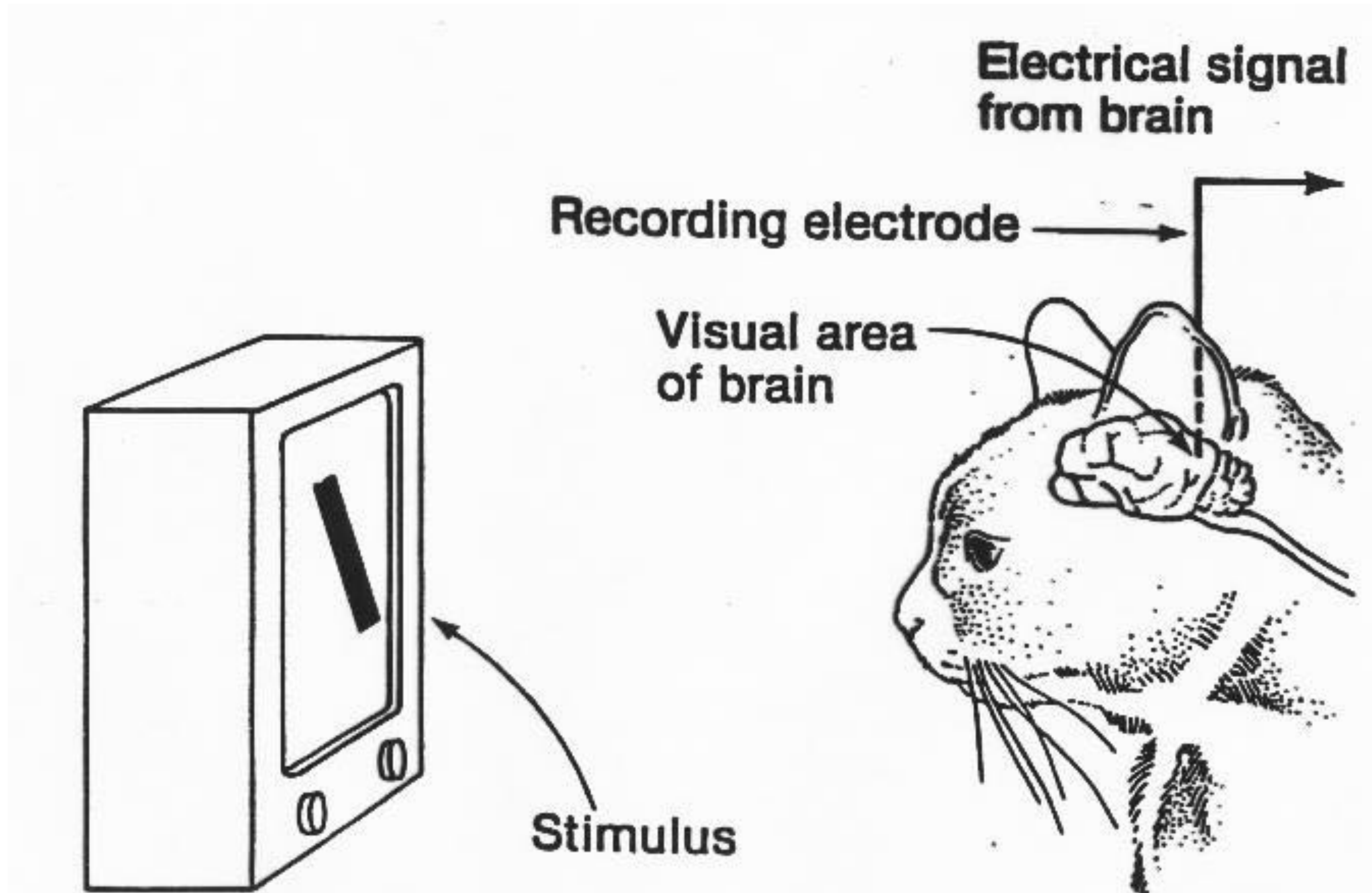# The Tungsten Electrode [Hubel-Science-1957]



http://braintour.harvard.edu/archives/portfolio-items/hubel-and-wiesel

- Device capable to record signal from a single neuron

Electrical signal from brain

Recording electrode

Visual area of brain

Stimulus

- Experiment with anaesthetised paralysed cat
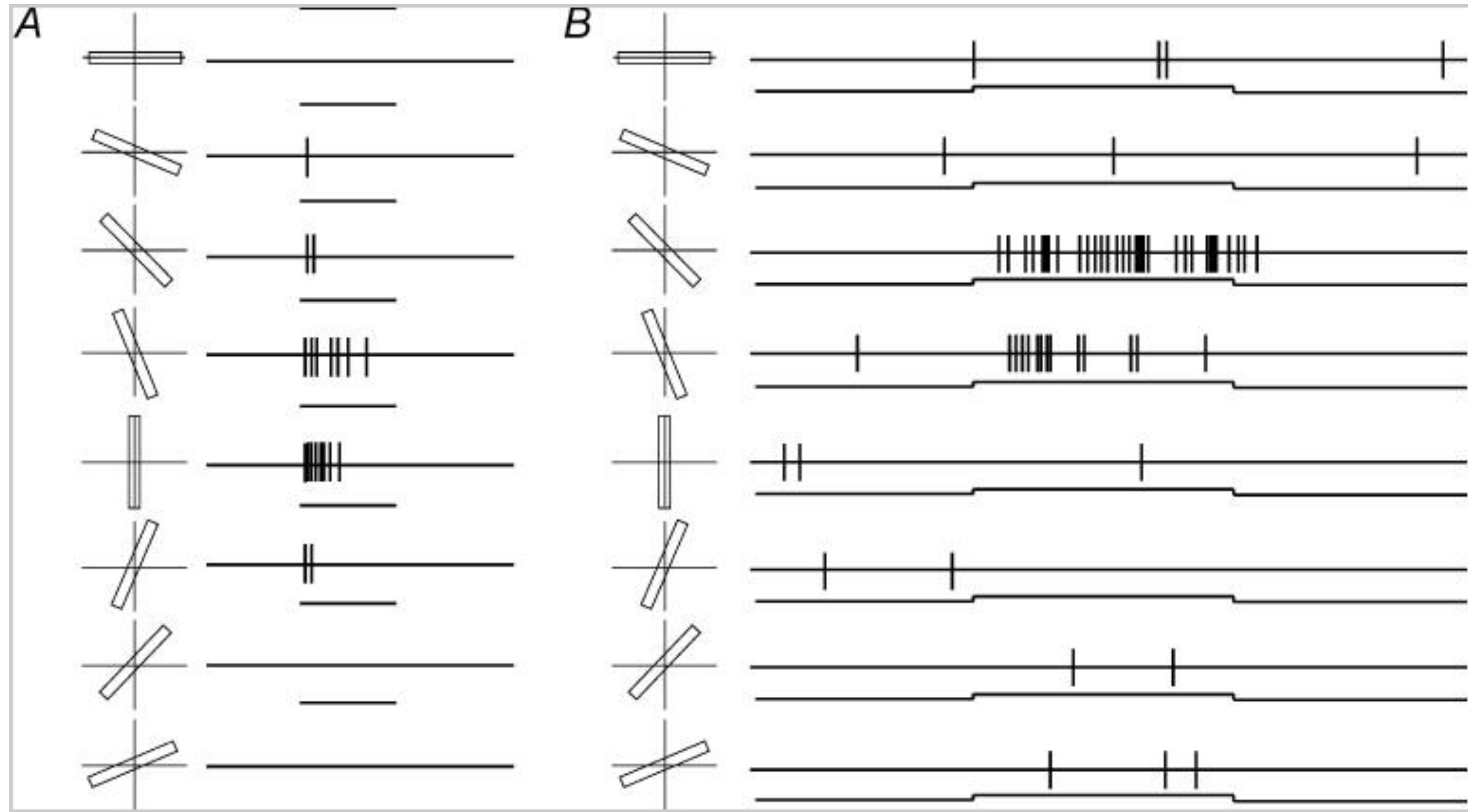
[Hubel and Wiesel 1960]



- Edge sensitivity
- Topographical mapping (nearby neurons process information from nearby visual fields)
- Translation invariance (the same edge is detected at all positions)

https://knowingneurons.com/2014/10/29/hubel-and-wiesel-the-neural-basis-of-visual-perception/

# [Hubel and Wiesel 1960]

paralysed cat

awake monkey



https://knowingneurons.com/2014/10/29/hubel-and-wiesel-the-neural-basis-of-visual-perception/
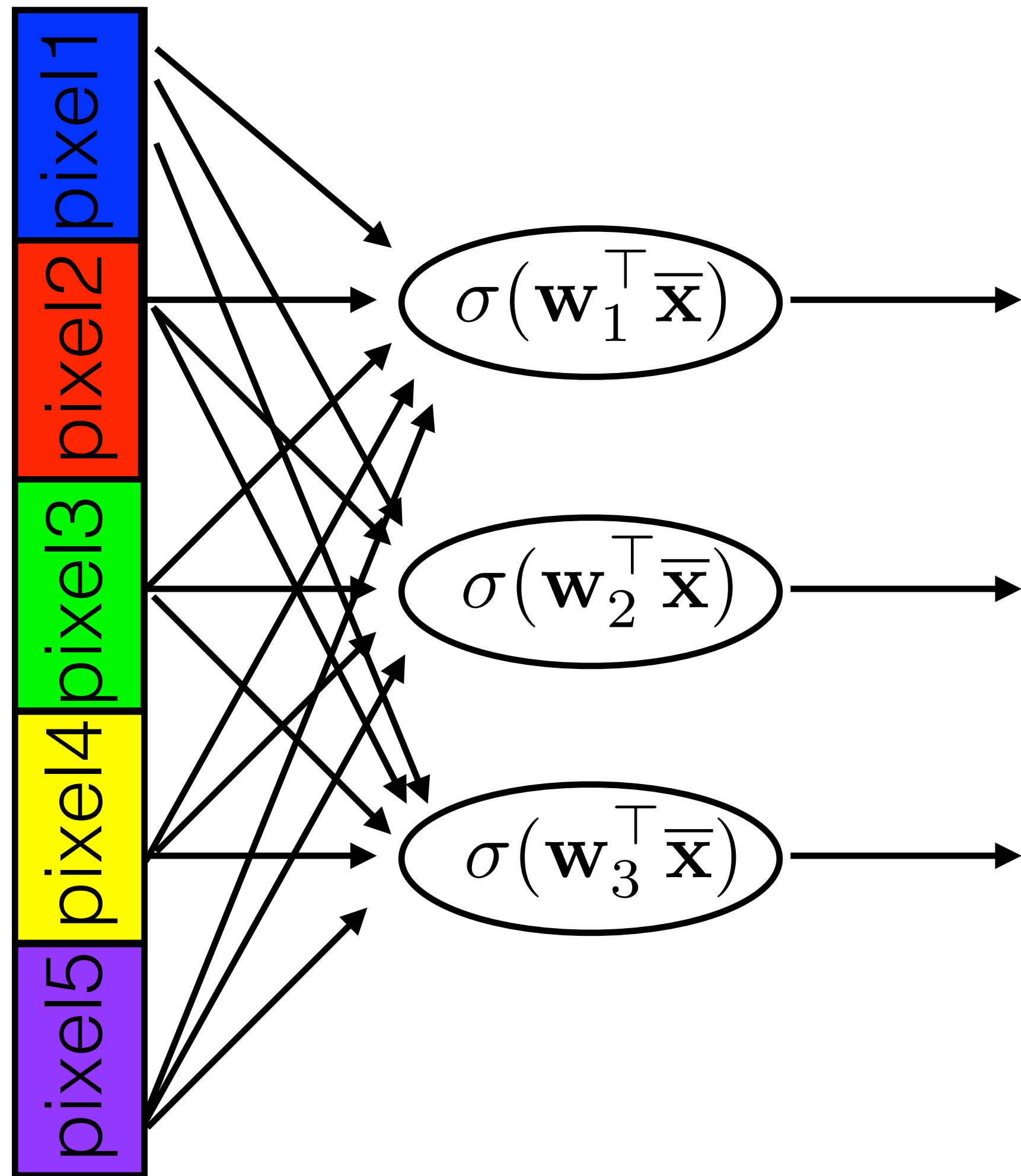
# Hubel and Wiesel experiments in 1950s and 1960s



- Nobel Prize in Physiology and Medicine in 1981
- Dr. Hubel: "There has been a myth that the brain cannot understand itself.  It is compared to a man trying to lift himself by his own bootstraps.  We feel that is nonsense. The brain can be studied just as the kidney can."

1. **Topographical map:** nearby neurons process information from nearby visual fields
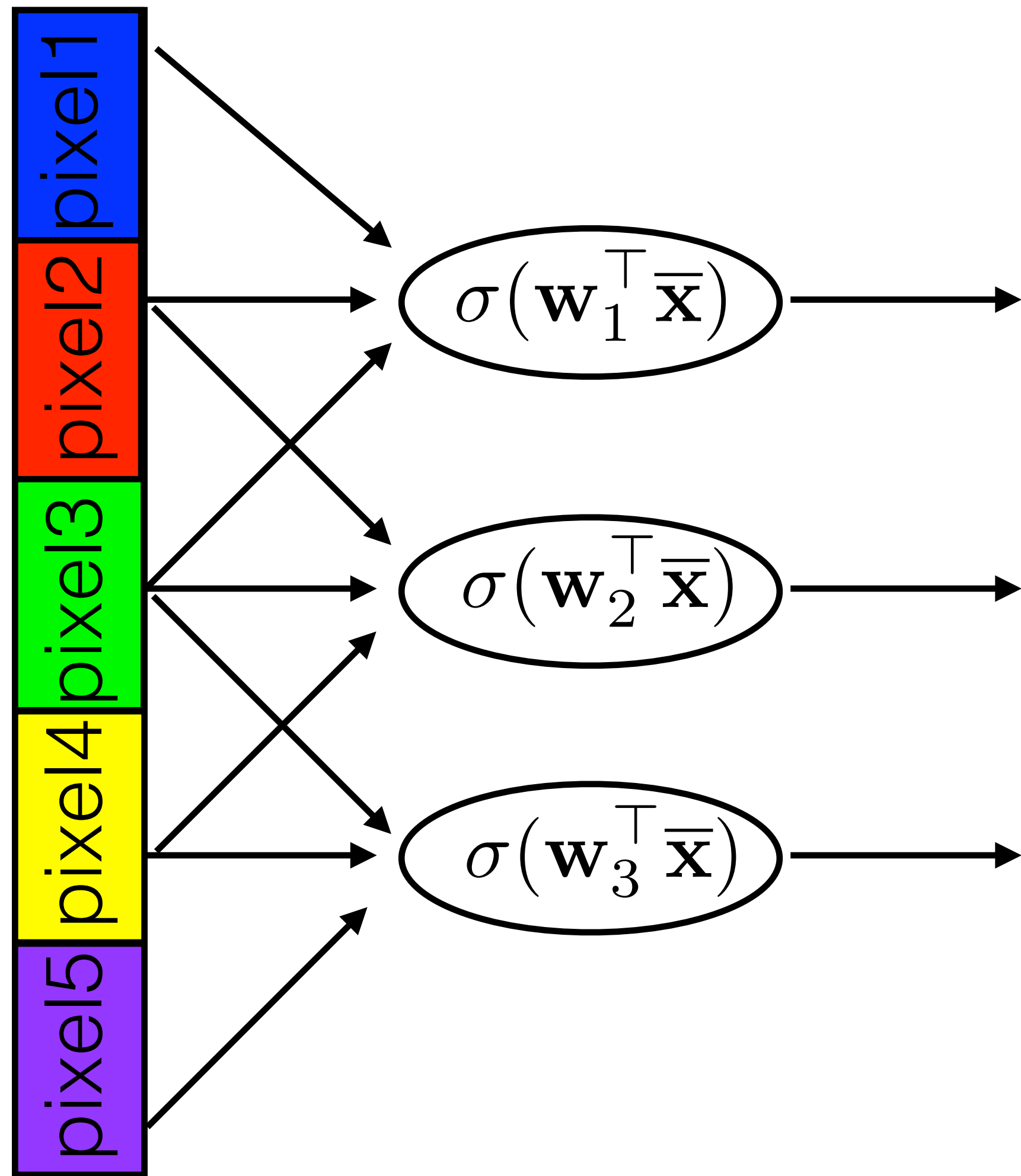
image



- Processing of visual information in cortex is not fully connected.

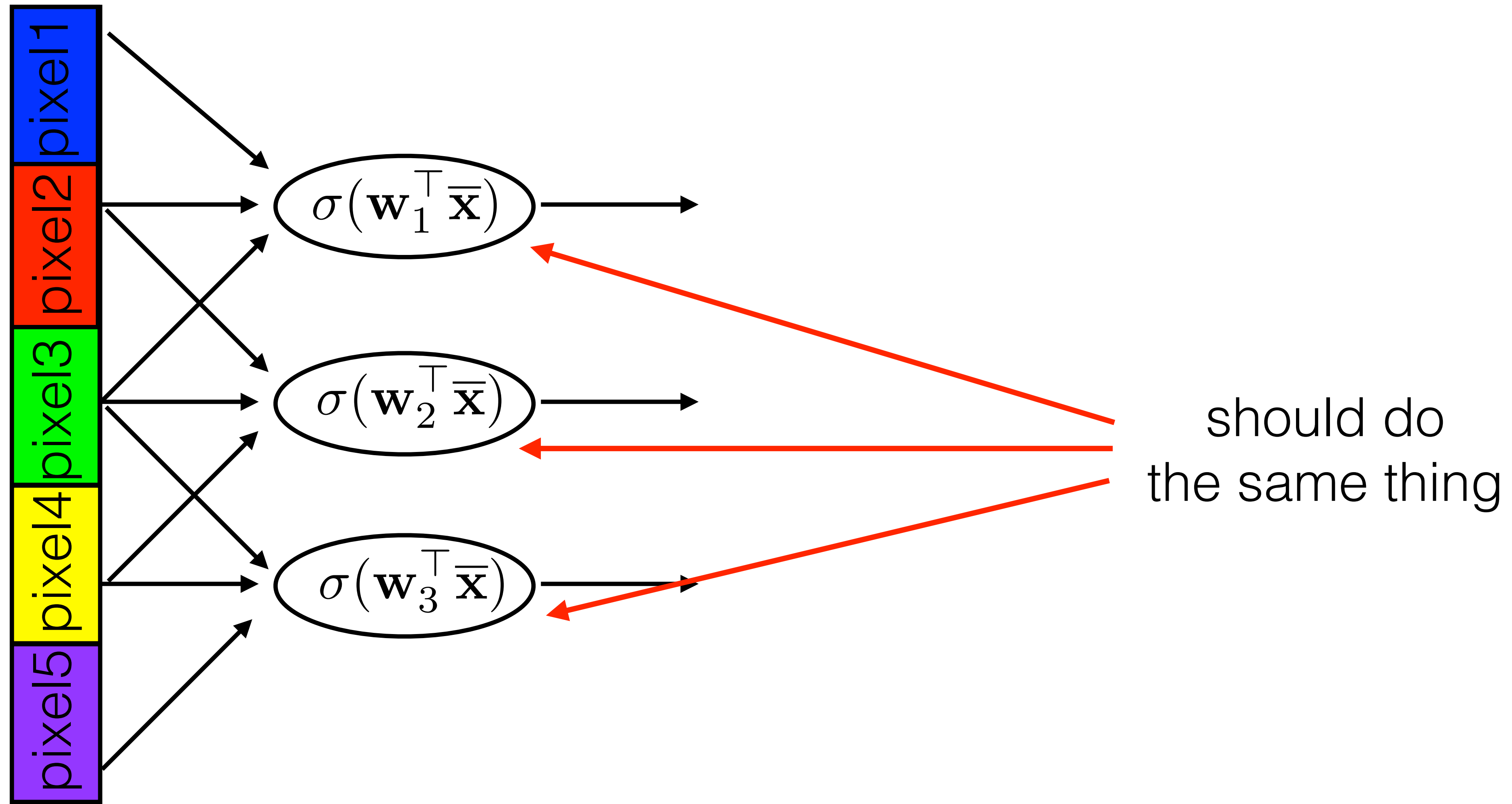1. **Topographical map:** nearby neurons process information from nearby visual fields

image



- Processing of visual information in cortex is not fully connected.

2. **Translation invariance:** the same edge is detected at all positions

image



$\sigma(\mathbf{w}_1^\top \overline{\mathbf{x}})$

$\sigma(\mathbf{w}_2^\top \overline{\mathbf{x}})$

$\sigma(\mathbf{w}_3^\top \overline{\mathbf{x}})$

should do
the same thing

2. **Translation invariance:** the same edge is detected at all positions
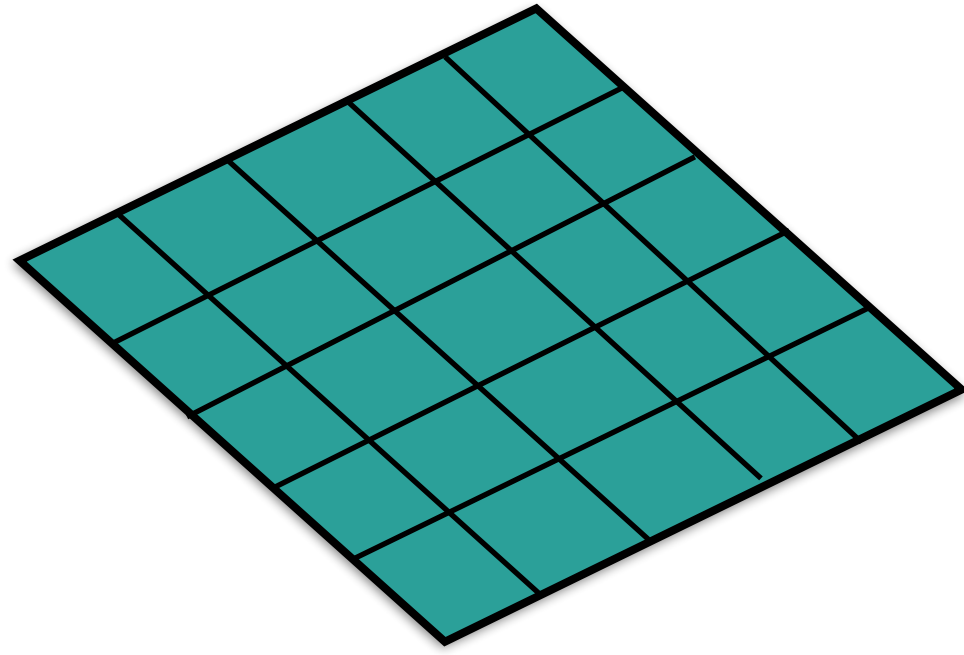
image



should do
the same thing

2. **Translation invariance:** the same edge is detected at all positions

image

# Fully Connected layer on images

output
(fcnn layer with
5x5 neurons)

input
(7x7 image)

Fully-connected layer

Fully Connected layer on images

$$\sigma(\sum \mathbf{W}_1 \cdot \mathbf{x})$$

output
(fcnn layer with
5x5 neurons)

input
(7x7 image)

Fully-connected layer

Fully Connected layer on images

$\sigma(\sum \mathbf{W}_1 \cdot \mathbf{X})$ $\sigma(\sum \mathbf{W}_2 \cdot \mathbf{X})$
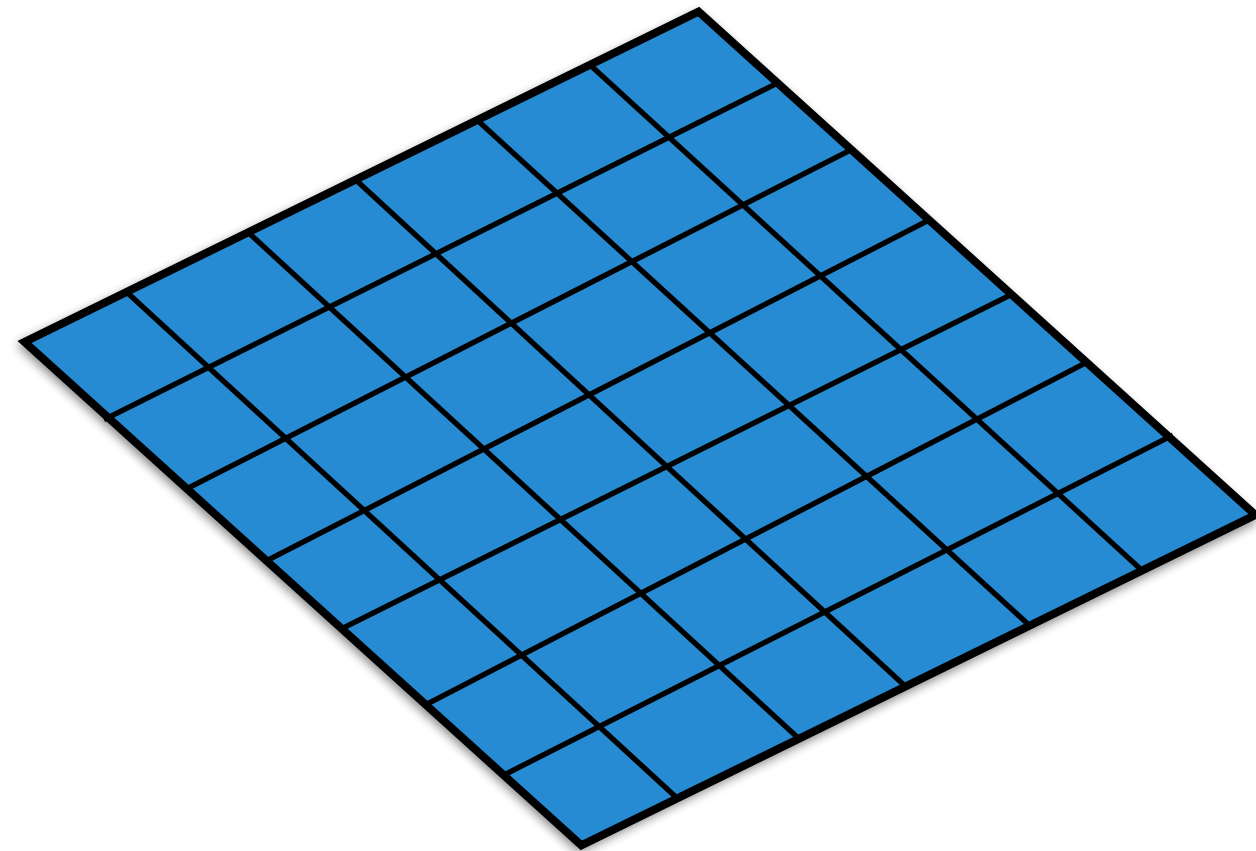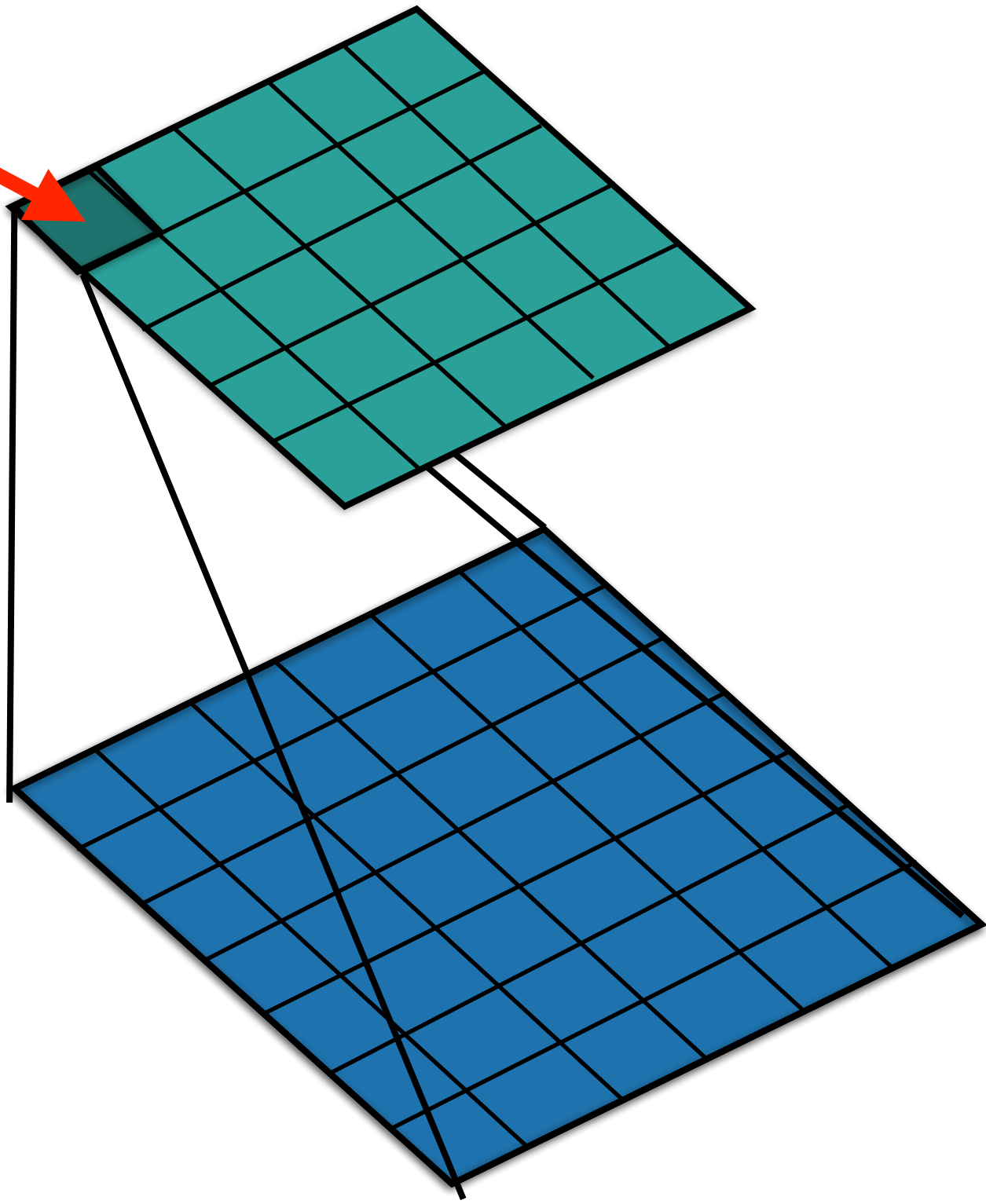
output
(fcnn layer with
5x5 neurons)

input
(7x7 image)

Fully-connected layer

Fully Connected layer on images

$$\sigma\left(\sum \mathbf{W}_1 \cdot \mathbf{X}\right) \quad \sigma\left(\sum \mathbf{W}_2 \cdot \mathbf{X}\right) \quad \sigma\left(\sum \mathbf{W}_3 \cdot \mathbf{X}\right)$$

output
(fcnn layer with
5x5 neurons)

input
(7x7 image)

Fully-connected layer

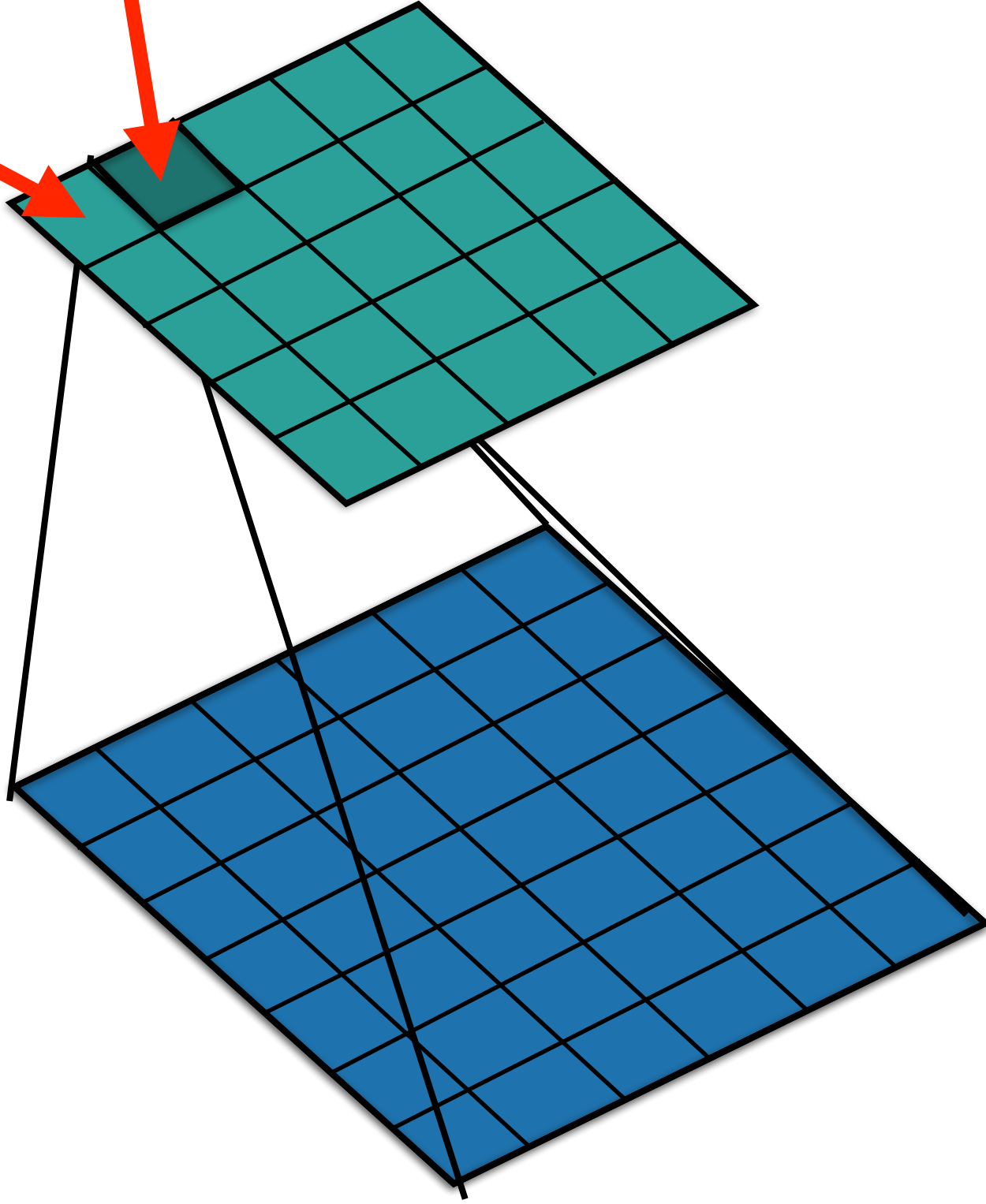**How many weights do I have in total?**

25x(7x7+1)=1250

Fully Connected layer on images

$\sigma(\sum \mathbf{W}_1 \cdot \mathbf{X})$  $\sigma(\sum \mathbf{W}_2 \cdot \mathbf{X})$  $\sigma(\sum \mathbf{W}_3 \cdot \mathbf{X})$

$\sigma(\mathrm{conv}(\mathbf{X}, \mathbf{w}))$

output
(fcnn layer with
5x5 neurons)

input
(7x7 image)
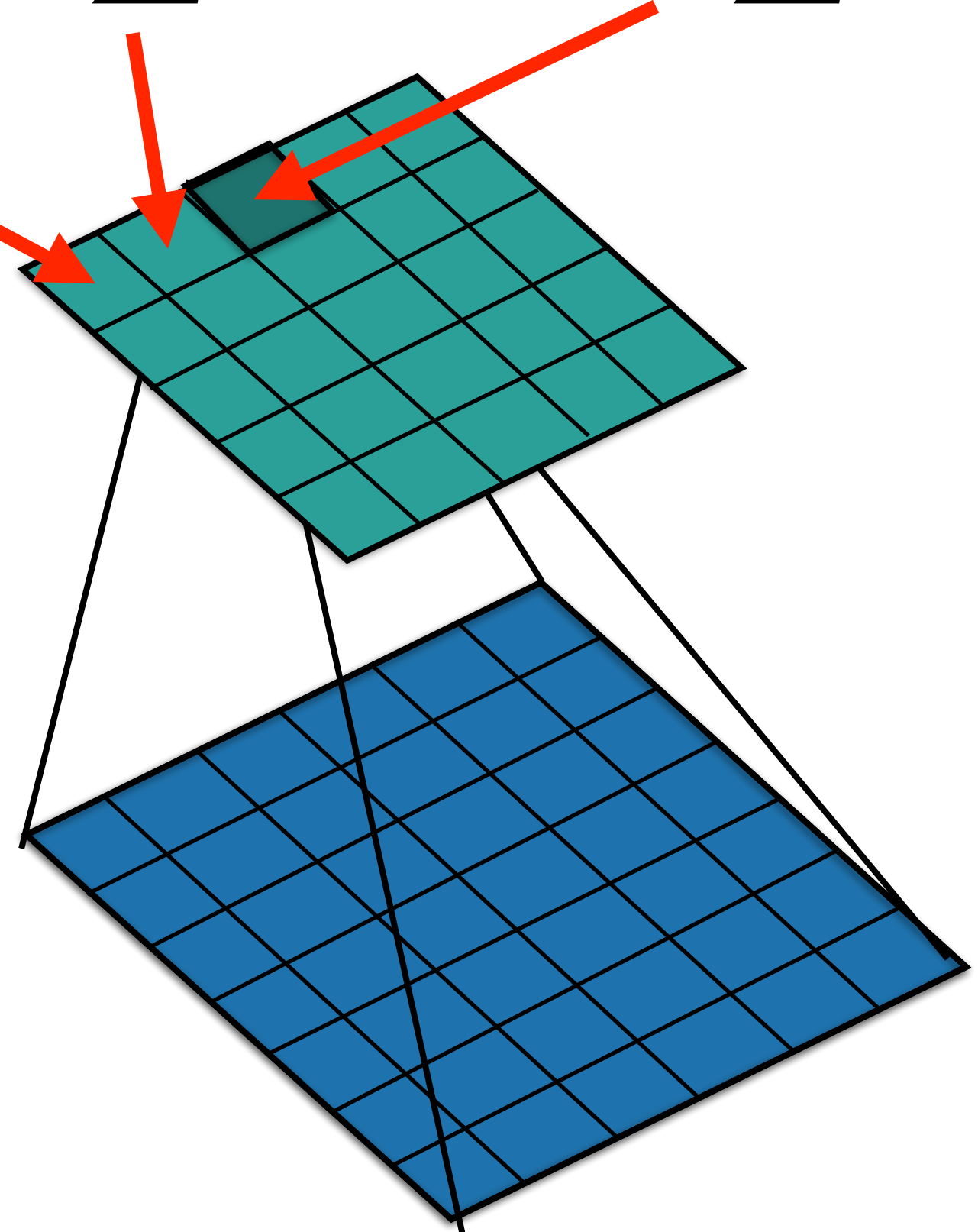


Fully-connected layer

Convolutional layer

**How many weights do I have in total?**

**How many weights do I have in total?**

25x(7x7+1)=1250

3x3+1=10

# Convolution forward pass $\mathbf{y} = \mathrm{conv}(\mathbf{x}, \mathbf{w})$

kernel/filter

| $w_{11}$ | $w_{12}$ |
|----------|----------|
| $w_{21}$ | $w_{22}$ |

filter response/
output map

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|----------|----------|----------|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|----------|----------|
| $y_{21}$ | $y_{22}$ |

image

Convolution forward pass $\mathbf{y} = \mathrm{conv}(\mathbf{x}, \mathbf{w})$

Local linear classifier run in double-for-loop over rows and columns

$$
\begin{array}{|c|c|}
\hline
y_{11} & y_{12} \\
\hline
y_{21} & y_{22} \\
\hline
\end{array}
= \mathrm{conv} \left(
\begin{array}{|c|c|c|}
\hline
x_{11} & x_{12} & x_{13} \\
\hline
x_{21} & x_{22} & x_{23} \\
\hline
x_{31} & x_{32} & x_{33} \\
\hline
\end{array}
,
\begin{array}{|c|c|}
\hline
w_{11} & w_{12} \\
\hline
w_{21} & w_{22} \\
\hline
\end{array}
\right)
$$

$$y_{11} = w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22}$$

$$y_{12} = w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23}$$

$$y_{21} = w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32}$$

$$y_{22} = w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33}$$

Convolution forward pass $\mathbf{y} = \mathrm{conv}(\mathbf{x}, \mathbf{w})$

Local linear classifier run in double-for-loop over rows and columns



$$y_{11} = w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22}$$

$$y_{12} = w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23}$$

$$y_{21} = w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32}$$

$$y_{22} = w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33}$$

# Convolution forward pass $\mathbf{y} = \mathrm{conv}(\mathbf{x}, \mathbf{w})$

Local linear classifier run in double-for-loop over rows and columns

$$
\begin{array}{|c|c|}
\hline
y_{11} & y_{12} \\
\hline
y_{21} & y_{22} \\
\hline
\end{array}
= \mathrm{conv} \left(
\begin{array}{|c|c|c|}
\hline
x_{11} & x_{12} & x_{13} \\
\hline
x_{21} & x_{22} & x_{23} \\
\hline
x_{31} & x_{32} & x_{33} \\
\hline
\end{array}
,
\begin{array}{|c|c|}
\hline
w_{11} & w_{12} \\
\hline
w_{21} & w_{22} \\
\hline
\end{array}
\right)
$$

$$y_{11} = w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22}$$

$$y_{12} = w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23}$$

$$y_{21} = w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32}$$

$$y_{22} = w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33}$$

Convolution forward pass $\quad \mathbf{y} = \mathrm{conv}(\mathbf{x}, \mathbf{w})$

Local linear classifier run in double-for-loop over rows and columns



$$y_{11} = w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22}$$

$$y_{12} = w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23}$$

$$y_{21} = w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32}$$

$$y_{22} = w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33}$$

Convolution forward pass $\mathbf{y} = \mathrm{conv}(\mathbf{x}, \mathbf{w})$

Local linear classifier run in double-for-loop over rows and columns



$$y_{11} = w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22}$$

$$y_{12} = w_{11}x_{12} + w_{12}x_{13} + w_{21}x_{22} + w_{22}x_{23}$$

$$y_{21} = w_{11}x_{21} + w_{12}x_{22} + w_{21}x_{31} + w_{22}x_{32}$$

$$y_{22} = w_{11}x_{22} + w_{12}x_{23} + w_{21}x_{32} + w_{22}x_{33}$$

# Convolution layer properties - output size

$$\mathrm{conv} \left( \text{image (5x5)} , \text{kernel (2x2)} \right) = \text{output}$$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

$$\mathrm{conv}\left(\text{image}, \text{kernel}\right) = \square$$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

# Convolution layer properties - output size

$$\mathrm{conv}\left( \; \text{image} \; , \; \text{kernel} \; \right) = \text{output}$$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

$$\mathrm{conv}\left(\quad,\quad\right)=$$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

# Convolution layer properties - output size

$$\mathrm{conv}\left( \begin{array}{|c|c|c|c|c|} \hline & & & & \\ \hline & & & & \\ \hline & & & & \\ \hline & & & & \\ \hline & & & & \\ \hline \end{array} , \begin{array}{|c|c|} \hline & \\ \hline & \\ \hline \end{array} \right) = \begin{array}{|c|c|c|c|} \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline & & & \\ \hline \end{array}$$

image          kernel          output
(5x5)          (2x2)           (4x4)

# Convolution layer properties - output size

$$M = N - K + 1$$

$$\mathrm{conv}\left(\quad,\quad\right) = $$

image
(NxN)

kernel
(KxK)

output
(MxM)

# Convolution layer properties - stride

## stride = 1

kernel moves by 1 pixel

$$\text{conv}\left( \text{image}, \text{kernel} \right) = \text{output}$$

image
(5x5)

kernel
(2x2)

output
(4x4)

# Convolution layer properties - stride

stride = 3

kernel moves by 3 pixels

$$\mathrm{conv}\left( \quad , \quad \right) = \square$$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

# Convolution layer properties - stride

stride = 3

kernel moves by 3 pixels

$$\text{conv} \left( \text{image}, \text{kernel} \right) = \text{output}$$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

# Convolution layer properties - stride

stride = 3

kernel moves by 3 pixels

→

$$\text{conv} \left( \text{image} , \text{kernel} \right) = \text{output}$$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

# Convolution layer properties - stride

stride = 3

kernel moves by 3 pixels



$$\mathrm{conv}\left(\ \text{image},\ \text{kernel}\ \right)\ =\ \text{output}$$

image
(5x5)

kernel
(2x2)

output
(2x2)

# Convolution layer properties - stride

$$M = \texttt{floor}( (N\text{-}K) / \text{stride} + 1)$$

stride

$$\text{conv} \left( \begin{array}{} \\ \end{array} , \begin{array}{} \\ \end{array} \right) = \begin{array}{} \\ \end{array}$$

image
(NxN)

kernel
(KxK)

output
(MxM)

e.g. M = (5-2) / 3 +1 = 2

# Convolution layer properties - pad

pad = 1

added border of size 1

$$\mathrm{conv}\left( \begin{array}{ccccccc} 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & & & & & & 0 \\ 0 & & & & & & 0 \\ 0 & & & & & & 0 \\ 0 & & & & & & 0 \\ 0 & & & & & & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{array}, \quad \begin{array}{cc} & \\ & \end{array} \right) = $$

image
(5x5)

kernel
(2x2)

output
(6x6)

# Convolution layer properties - pad

$$M = \texttt{floor}( (N+2*pad-K) / stride + 1)$$



conv ( image , kernel ) = output

added border of size 1

image
(N+2*pad)X(N+2*pad)

kernel
(KxK)

output
(MxM)

# Convolution + padding + stride

convolution                padding                padding+stride



■ input

■ output

https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md

# Convolution layer

Dilatation rate = 1

$$\mathrm{conv}\left( \quad , \quad \right) = $$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

# Dilated convolution layer

Dilatation rate = 2

$$\mathrm{conv}\left( \quad , \quad \right) = $$

image
(5x5)

kernel
(2x2)

output
(**? x ?**)

# Dilated convolution layer

dilatation rate = 1          dilatation rate = 2



- input
- output

https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md

# Transposed convolution
## Upsampling / smart interpolation

no padding
no stride

input

output

https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md

# Transposed convolution
## Upsampling / smart interpolation

no padding
no stride

no padding
stride

■ input

■ output

https://github.com/vdumoulin/conv_arithmetic/blob/master/README.md

# Transposed convolution
## Upsampling / smart interpolation

| no padding<br>no stride | no padding<br>stride | padding<br>stride | full padding<br>no stride |
|---|---|---|---|



■ input
■ output

## Input image



## Input kernel



| | 0 | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|---|
| 0 | 1, | 1, | 1, | 1, | 1, | 1 |
| 1 | 1, | 1, | 1, | 1, | 1, | 1 |
| 2 | 1, | 1, | 1, | 1, | 1, | 1 |
| 3 | 1, | 1, | 1, | 1, | 1, | 1 |
| 4 | 1, | 1, | 1, | 1, | 1, | 1 |
| 5 | 1, | 1, | 1, | 1, | 1, | 1 |

## Output



```
nn.Conv2d(1,1,3)
```



```
nn.Conv2d(1,1,6)
```



```
nn.Conv2d(1,1,12)
```

## Input image

## Input kernel

## Output



`nn.Conv2d(1,1,6)`

`nn.Conv2d(1,1,12)`

## Input image



## Input kernel



## Output



nn.Conv2d(1,1,6)

nn.Conv2d(1,1,12)

## Input image



## Input kernel



## Output



nn.Conv2d(1,1,6)

nn.Conv2d(1,1,12)

# Input image



# Input kernel



# Output



`nn.Conv2d(1,1,6)`



`nn.Conv2d(1,1,12)`

# Meaning of padding and stride



`nn.Conv2d(1,1,6,`
`stride=(3,3))`



`nn.Conv2d(1,1,6,`
`padding=(15,15))`

48

# Input image



# Input kernel



# Output



`nn.Conv2d(1,1,6)`



`nn.Conv2d(1,1,12)`

compare!

# Meaning of dilation



`nn.Conv2d(1,1,6, dilation=2)`

Input image　　　Input kernel　　　　　　　　Output



```
nn.Conv2d(1,1,32,
padding=(15,15))
```

Meaning of kernel

Convolution is locally applied linear classifier that computes correlation between kernel and image

Let's build a convolutional layer !

# Multi-channel convolution



$$\mathrm{conv} \left( \quad , \quad \right) =$$

RGB image
(5x5x3)

kernel
(2x2x3)

output
(4x4x1)

# Multi-channel convolution



$$\mathrm{conv} \left( \quad , \quad \right) =$$

RGB image
(5x5x3)

kernel
(2x2x3)

output
(4x4x1)

# Multi-channel convolution



$$\mathrm{conv}\left( \begin{array}{c} \text{RGB image} \end{array}, \begin{array}{c} \text{kernel} \end{array} \right) = \begin{array}{c} \text{output} \end{array}$$

RGB image
(5x5x3)

kernel
(2x2x3)

output
(4x4x1)

# Multi-channel convolution

$$\text{conv}\left(\quad,\quad\right) =$$

RGB image
(5x5x3)

kernel
(2x2x3)

output
(4x4x1)

# Multi-channel convolution



$$\mathrm{conv}\left( \text{RGB image}, \text{kernel} \right) = \text{output}$$

RGB image
(5x5x3)

kernel
(2x2x3)

output
(4x4x1)

$\mathbf{x}$

$\mathbf{w}_1$ Convolutional layer

$\mathbf{y}_1$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_1)$

$$\mathbf{x} \qquad \mathbf{w}_1 \qquad \text{Convolutional layer}$$

$$\mathbf{y}_1$$

$$\text{conv}(\mathbf{x}, \mathbf{w}_1)$$

$$\mathbf{w}_2$$

$$\text{conv}(\mathbf{x}, \mathbf{w}_2)$$

$$\mathbf{y}_2$$

**x**

**w**$_1$    Convolutional layer

**w**$_2$

**w**$_3$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_1)$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_2)$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_3)$

**y**$_1$

**y**$_2$

**y**$_3$

**x**

**w₁**  Convolutional layer

$\mathbf{w}_1$

$\mathbf{y}_1$

$\text{conv}(\mathbf{x}, \mathbf{w}_1)$

5x5x3

$\mathbf{w}_2$

$\text{conv}(\mathbf{x}, \mathbf{w}_2)$

$\mathbf{y}_2$

```python
# initialise
import torch.nn as nn
# define 2D convolutional layer
first_layer = nn.Conv2d(in_channels=3, out_channels=2, kernel_size=2
                        stride=1, padding=1)
```

**x**

**w₁**   Convolutional layer

$\mathbf{w}_1$

$\mathbf{y}_1$

$\text{conv}(\mathbf{x}, \mathbf{w}_1)$

5x5x3

$\mathbf{w}_2$

$\mathbf{y}_2$

$\text{conv}(\mathbf{x}, \mathbf{w}_2)$

also number
of kernels

4x4x2

```python
# initialise
import torch.nn as nn
# define 2D convolutional layer
first_layer = nn.Conv2d(in_channels=3, out_channels=2, kernel_size=2
                        stride=1, padding=1)
```

**x**

**w₁**  Convolutional layer

$\mathbf{w}_1$

2x2x3

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_1)$

$\mathbf{y}_1$

5x5x3

$\mathbf{w}_2$

2x2x3

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_2)$

$\mathbf{y}_2$

4x4x2

```python
# initialise
import torch.nn as nn
# define 2D convolutional layer
first_layer = nn.Conv2d(in_channels=3, out_channels=2, kernel_size=2
                        stride=1, padding=1)
```

Let's train the convolutional network !

Example: 1D convolution backward pass

$$\mathbf{w} = \begin{bmatrix} +2 \\ -2 \end{bmatrix}$$

$$\mathbf{z} = \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 2 \\ 1 \end{bmatrix}$$



$\mathrm{conv}(\mathbf{x}, \mathbf{w})$     $p(\mathbf{y})$     $p$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}} = \quad \textbf{???}$$

# Example: 1D convolution backward pass

$$\mathbf{w} = \begin{bmatrix} +2 \\ -2 \end{bmatrix}$$

$$\mathbf{z} = \begin{bmatrix} -2 \\ 0 \\ 2 \end{bmatrix}$$

$$\mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 2 \\ 1 \end{bmatrix} \longrightarrow \text{conv}(\mathbf{x}, \mathbf{w}) \longrightarrow p(\mathbf{y}) \longrightarrow p$$

$$\frac{\partial \mathbf{z}}{\partial \mathbf{w}} = \frac{\partial [w_0 x_0 + w_1 x_1, \ \ w_0 x_1 + w_1 x_2, \ \ w_0 x_2 + w_1 x_3]}{\partial \mathbf{w}} = \begin{bmatrix} x_0, x_1 \\ x_1, x_2 \\ x_2, x_3 \end{bmatrix}$$

```
def vjp_conv_w(v, x):
```

$$\texttt{return} \ \ \mathbf{v}^\top \cdot \frac{\partial \mathbf{z}}{\partial \mathbf{w}} \ = [v_0, v_1, v_2] \cdot \begin{bmatrix} x_0, x_1 \\ x_1, x_2 \\ x_2, x_3 \end{bmatrix} = [v_0 x_0 + v_1 x_1 + v_2 x_2, v_0 x_1 + v_1 x_2 + v_2 x_3]$$

$$= \text{conv}(\mathbf{x}, \mathbf{v})$$

# 2D Convolution backward pass

| $w_{11}$ | $w_{12}$ |
|----------|----------|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|----------|----------|----------|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|----------|----------|
| $y_{21}$ | $y_{22}$ |

$\mathrm{p}(\mathbf{y})$

| $p$ |
|-----|

# 2D Convolution backward pass

$$\begin{array}{|c|c|} \hline \dfrac{\partial p}{\partial w_{11}} & \dfrac{\partial p}{\partial w_{12}} \\ \hline \dfrac{\partial p}{\partial w_{21}} & \dfrac{\partial p}{\partial w_{22}} \\ \hline \end{array} \;\; = ?$$

| $w_{11}$ | $w_{12}$ |
|----------|----------|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|----------|----------|----------|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

$\mathrm{p}(\mathbf{y})$

| $y_{11}$ | $y_{12}$ |
|----------|----------|
| $y_{21}$ | $y_{22}$ |

$p$

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \quad ?$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$p$

$\mathrm{p}(\mathbf{y})$

| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream
gradient

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}}\frac{\partial y_{11}}{\partial w_{11}} + \frac{\partial p}{\partial y_{12}}\frac{\partial y_{12}}{\partial w_{11}} + \frac{\partial p}{\partial y_{21}}\frac{\partial y_{21}}{\partial w_{11}} + \frac{\partial p}{\partial y_{22}}\frac{\partial y_{22}}{\partial w_{11}}$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\text{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$\text{p}(\mathbf{y})$

$p$

| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream
gradient

69

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}} \frac{\partial y_{11}}{\partial w_{11}} + \frac{\partial p}{\partial y_{12}} \frac{\partial y_{12}}{\partial w_{11}} + \frac{\partial p}{\partial y_{21}} \frac{\partial y_{21}}{\partial w_{11}} + \frac{\partial p}{\partial y_{22}} \frac{\partial y_{22}}{\partial w_{11}}$$

$$\frac{\partial y_{11}}{\partial w_{11}} = \quad ?$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$\mathrm{p}(\mathbf{y})$

$p$

| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream gradient

## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}} \frac{\partial y_{11}}{\partial w_{11}} + \frac{\partial p}{\partial y_{12}} \frac{\partial y_{12}}{\partial w_{11}} + \frac{\partial p}{\partial y_{21}} \frac{\partial y_{21}}{\partial w_{11}} + \frac{\partial p}{\partial y_{22}} \frac{\partial y_{22}}{\partial w_{11}}$$

$$\frac{\partial y_{11}}{\partial w_{11}} = \frac{\partial(w_{11}x_{11} + w_{12}x_{12} + w_{21}x_{21} + w_{22}x_{22})}{\partial w_{11}} = x_{11}$$

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\text{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

$\text{p}(\mathbf{y})$

$p$

upstream gradient

$$\begin{array}{|c|c|}\hline \frac{\partial p}{\partial w_{11}} & \frac{\partial p}{\partial w_{12}} \\\hline \frac{\partial p}{\partial w_{21}} & \frac{\partial p}{\partial w_{22}} \\\hline\end{array}$$

## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}}x_{11} + \frac{\partial p}{\partial y_{12}}x_{12} + \frac{\partial p}{\partial y_{21}}x_{21} + \frac{\partial p}{\partial y_{22}}x_{22}$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$p$

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

$\mathrm{p}(\mathbf{y})$

| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream
gradient

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

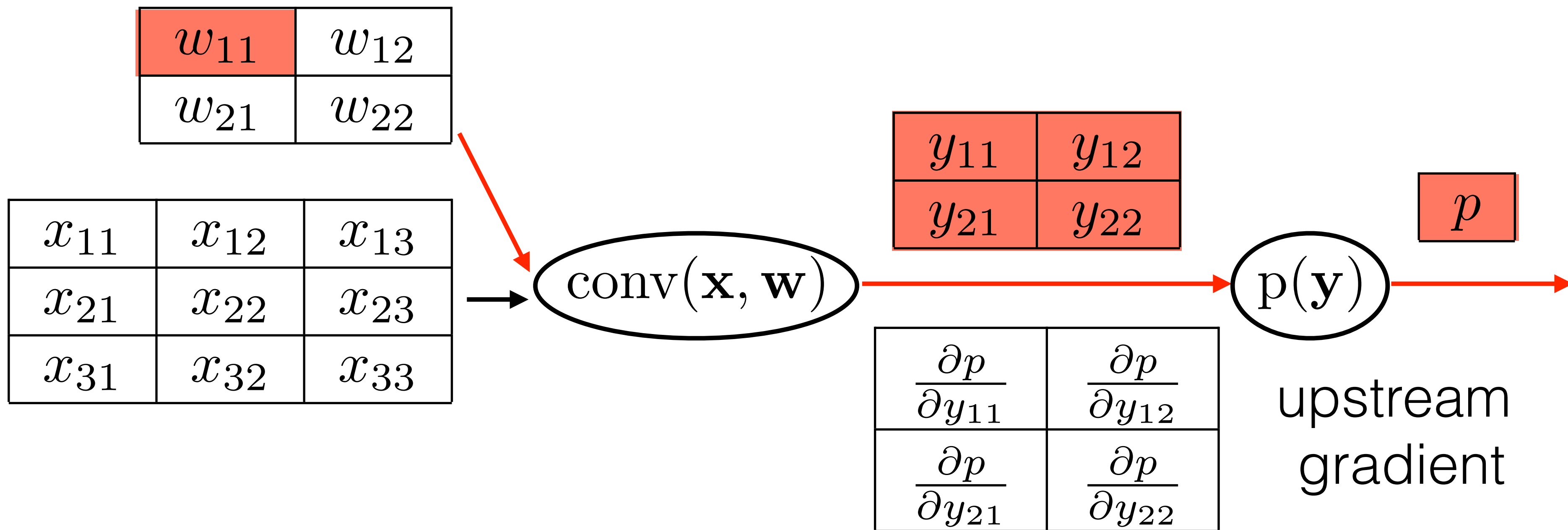## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}}x_{11} + \frac{\partial p}{\partial y_{12}}x_{12} + \frac{\partial p}{\partial y_{21}}x_{21} + \frac{\partial p}{\partial y_{22}}x_{22}$$

$$\frac{\partial p}{\partial w_{12}} = \quad ?$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\text{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$p$

$\text{p}(\mathbf{y})$

upstream gradient

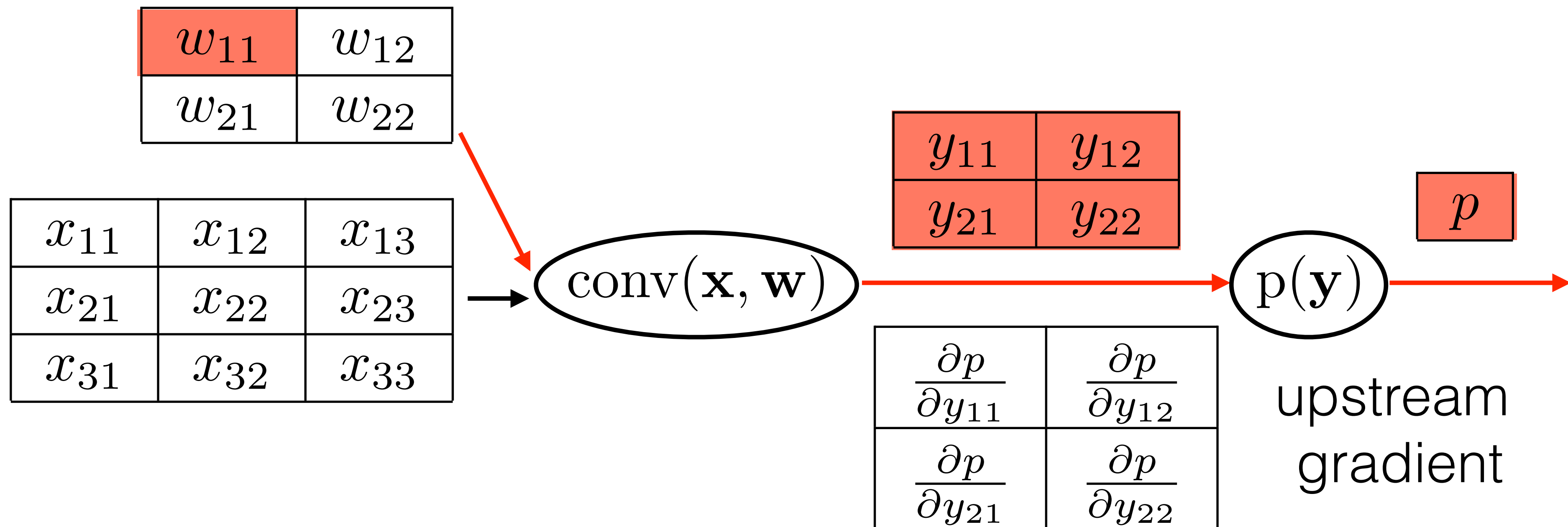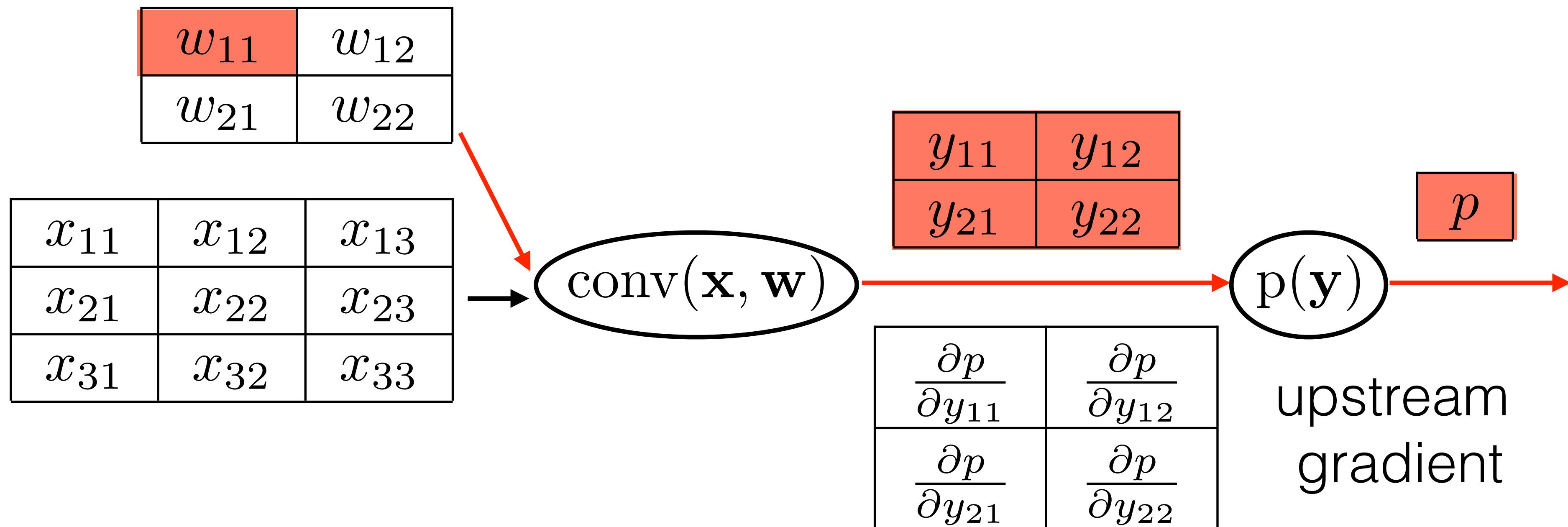| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

73

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}}x_{11} + \frac{\partial p}{\partial y_{12}}x_{12} + \frac{\partial p}{\partial y_{21}}x_{21} + \frac{\partial p}{\partial y_{22}}x_{22}$$

$$\frac{\partial p}{\partial w_{12}} = \frac{\partial p}{\partial y_{11}}x_{12} + \frac{\partial p}{\partial y_{12}}x_{13} + \frac{\partial p}{\partial y_{21}}x_{22} + \frac{\partial p}{\partial y_{22}}x_{23}$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$\mathrm{p}(\mathbf{y})$

$p$

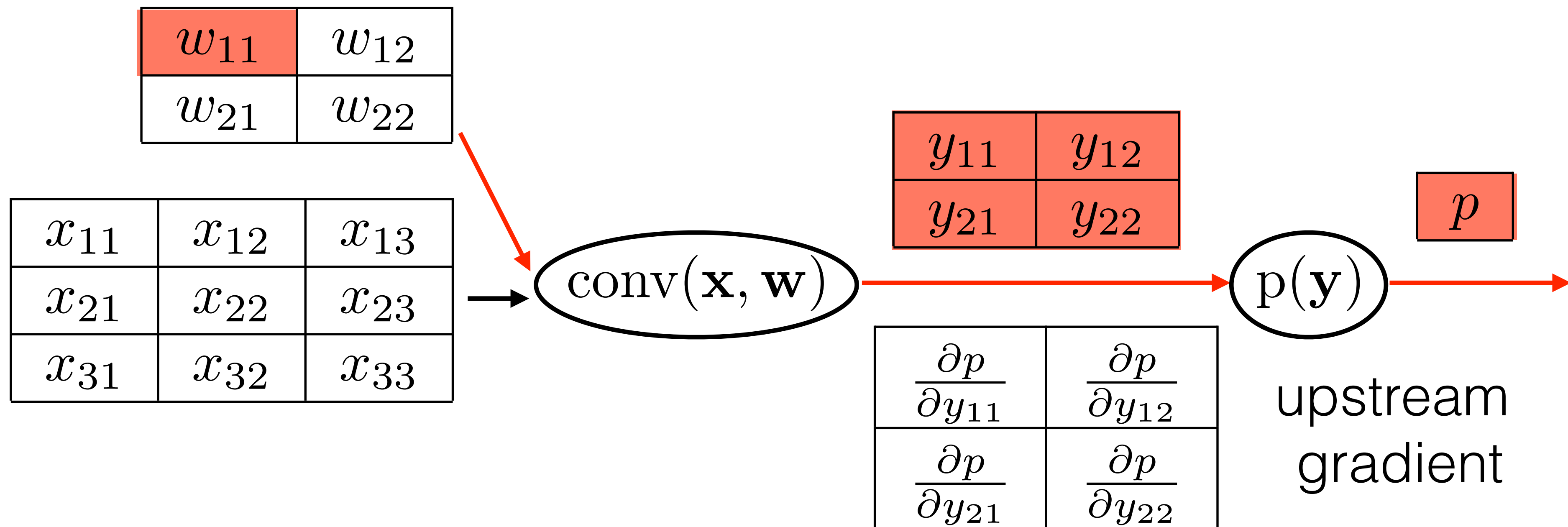| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream gradient

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}}x_{11} + \frac{\partial p}{\partial y_{12}}x_{12} + \frac{\partial p}{\partial y_{21}}x_{21} + \frac{\partial p}{\partial y_{22}}x_{22}$$

$$\frac{\partial p}{\partial w_{12}} = \frac{\partial p}{\partial y_{11}}x_{12} + \frac{\partial p}{\partial y_{12}}x_{13} + \frac{\partial p}{\partial y_{21}}x_{22} + \frac{\partial p}{\partial y_{22}}x_{23}$$

$$\frac{\partial p}{\partial w_{21}} = \frac{\partial p}{\partial y_{11}}x_{21} + \frac{\partial p}{\partial y_{12}}x_{22} + \frac{\partial p}{\partial y_{21}}x_{31} + \frac{\partial p}{\partial y_{22}}x_{32}$$
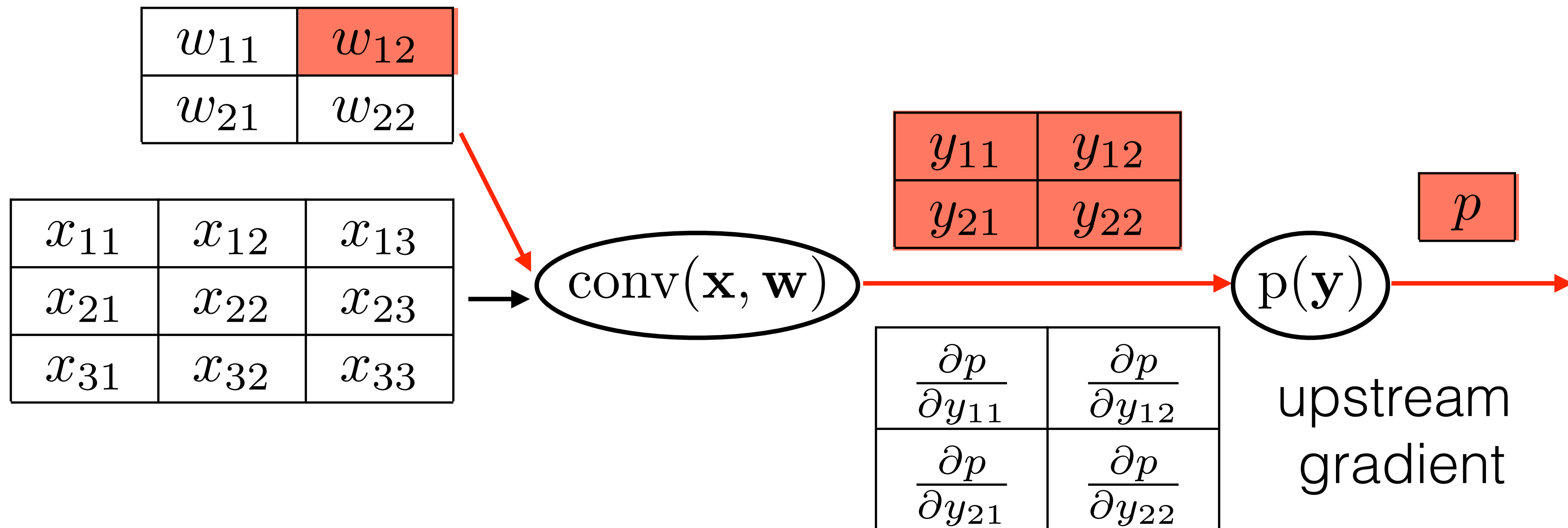
| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$\mathrm{p}(\mathbf{y})$

$p$

| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream gradient

## 2D Convolution backward pass

| $\frac{\partial p}{\partial w_{11}}$ | $\frac{\partial p}{\partial w_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial w_{21}}$ | $\frac{\partial p}{\partial w_{22}}$ |

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}}x_{11} + \frac{\partial p}{\partial y_{12}}x_{12} + \frac{\partial p}{\partial y_{21}}x_{21} + \frac{\partial p}{\partial y_{22}}x_{22}$$
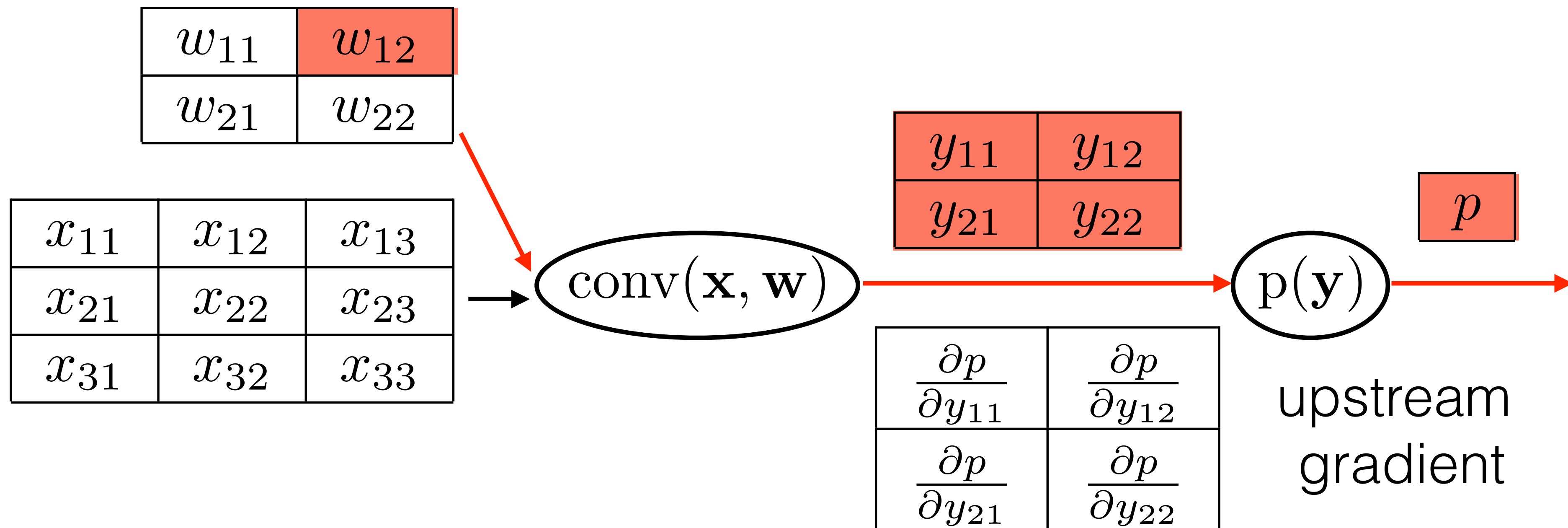
$$\frac{\partial p}{\partial w_{12}} = \frac{\partial p}{\partial y_{11}}x_{12} + \frac{\partial p}{\partial y_{12}}x_{13} + \frac{\partial p}{\partial y_{21}}x_{22} + \frac{\partial p}{\partial y_{22}}x_{23}$$
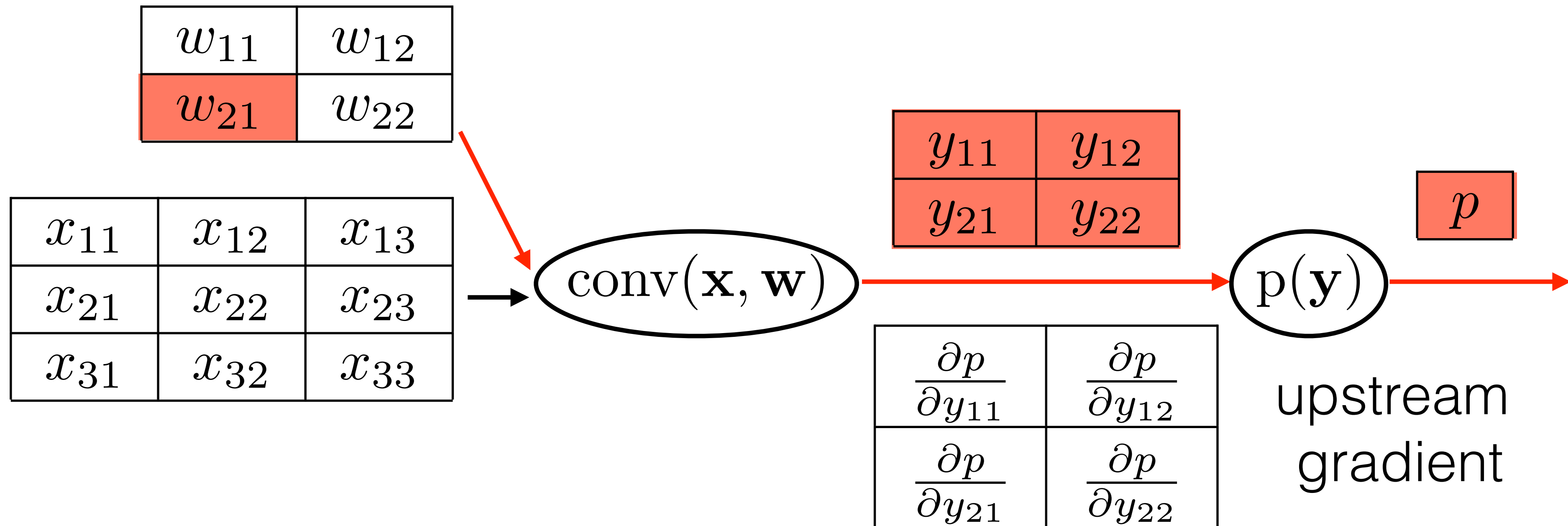
$$\frac{\partial p}{\partial w_{21}} = \frac{\partial p}{\partial y_{11}}x_{21} + \frac{\partial p}{\partial y_{12}}x_{22} + \frac{\partial p}{\partial y_{21}}x_{31} + \frac{\partial p}{\partial y_{22}}x_{32}$$

$$\frac{\partial p}{\partial w_{22}} = \frac{\partial p}{\partial y_{11}}x_{22} + \frac{\partial p}{\partial y_{12}}x_{23} + \frac{\partial p}{\partial y_{21}}x_{32} + \frac{\partial p}{\partial y_{22}}x_{33}$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$p$

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

$\mathrm{p}(\mathbf{y})$

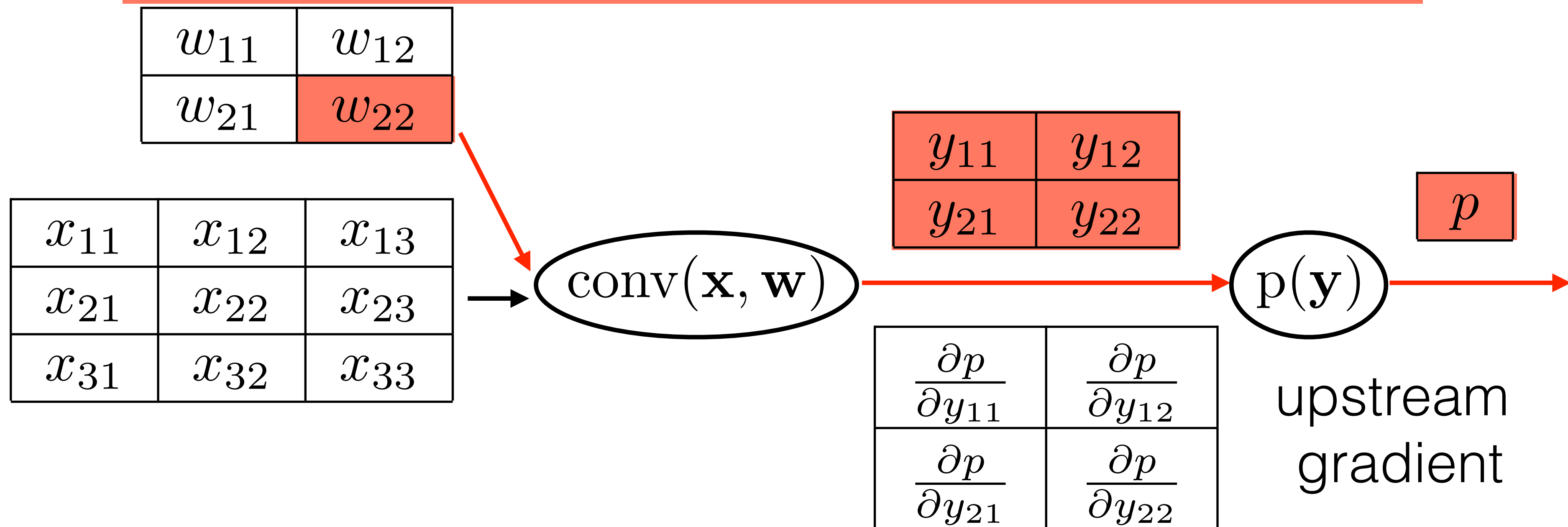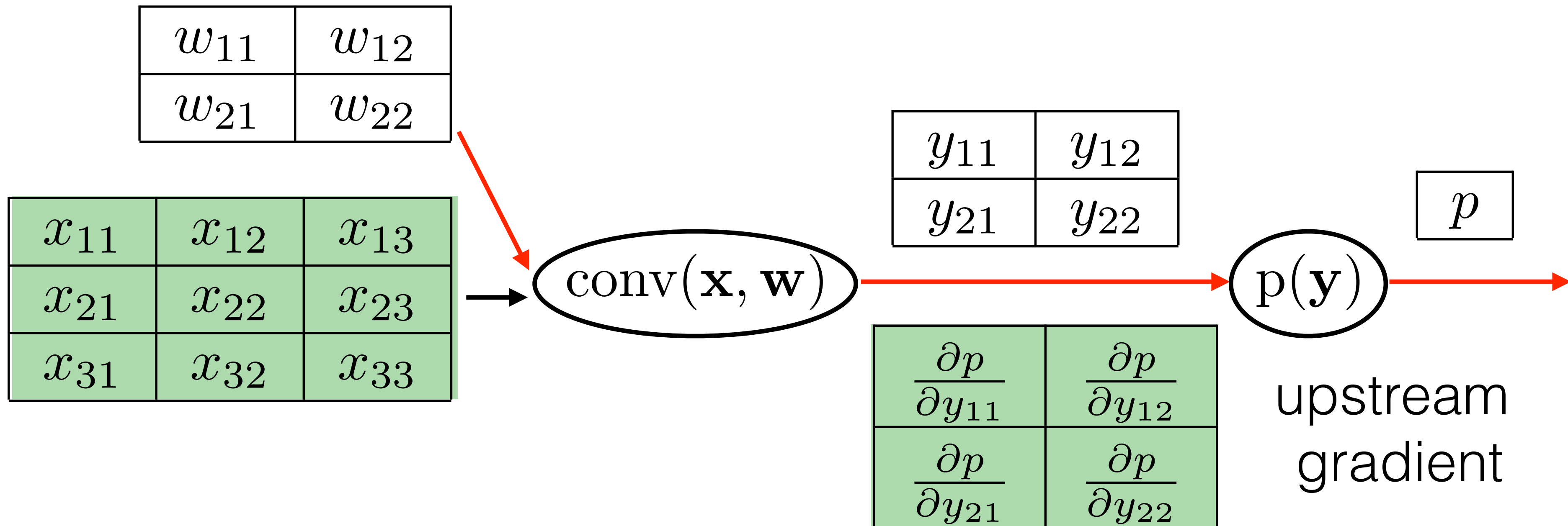| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|---|---|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream
gradient

## 2D Convolution backward pass

$$\frac{\partial p}{\partial w_{11}} = \frac{\partial p}{\partial y_{11}} x_{11} + \frac{\partial p}{\partial y_{12}} x_{12} + \frac{\partial p}{\partial y_{21}} x_{21} + \frac{\partial p}{\partial y_{22}} x_{22}$$

$$\frac{\partial p}{\partial w_{12}} = \frac{\partial p}{\partial y_{11}} x_{12} + \frac{\partial p}{\partial y_{12}} x_{13} + \frac{\partial p}{\partial y_{21}} x_{22} + \frac{\partial p}{\partial y_{22}} x_{23}$$

$$\frac{\partial p}{\partial w_{21}} = \frac{\partial p}{\partial y_{11}} x_{21} + \frac{\partial p}{\partial y_{12}} x_{22} + \frac{\partial p}{\partial y_{21}} x_{31} + \frac{\partial p}{\partial y_{22}} x_{32}$$
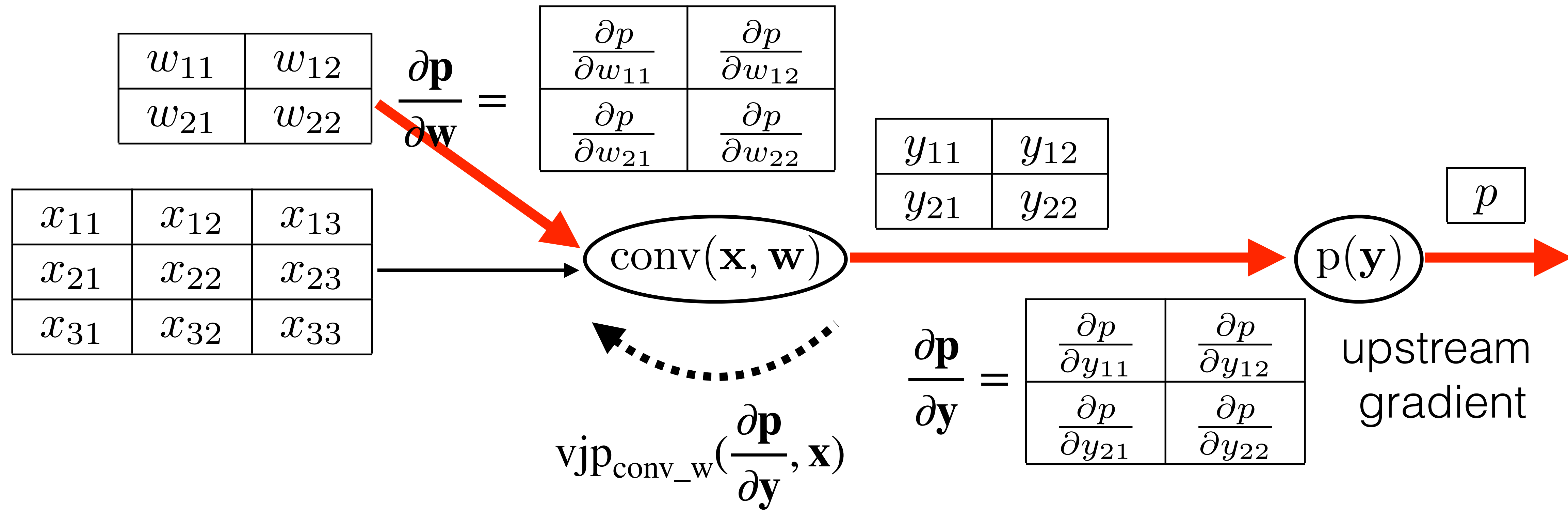
$$\frac{\partial p}{\partial w_{22}} = \frac{\partial p}{\partial y_{11}} x_{22} + \frac{\partial p}{\partial y_{12}} x_{23} + \frac{\partial p}{\partial y_{21}} x_{32} + \frac{\partial p}{\partial y_{22}} x_{33}$$

| $w_{11}$ | $w_{12}$ |
|----------|----------|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|----------|----------|----------|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$\mathrm{conv}(\mathbf{x}, \mathbf{w})$

| $y_{11}$ | $y_{12}$ |
|----------|----------|
| $y_{21}$ | $y_{22}$ |

$\mathrm{p}(\mathbf{y})$

$p$

| $\frac{\partial p}{\partial y_{11}}$ | $\frac{\partial p}{\partial y_{12}}$ |
|----------|----------|
| $\frac{\partial p}{\partial y_{21}}$ | $\frac{\partial p}{\partial y_{22}}$ |

upstream gradient

# Convolution backward pass wrt weights

$$\frac{\partial \mathbf{p}}{\partial \mathbf{w}} = \begin{array}{|c|c|} \hline \dfrac{\partial p}{\partial w_{11}} & \dfrac{\partial p}{\partial w_{12}} \\ \hline \dfrac{\partial p}{\partial w_{21}} & \dfrac{\partial p}{\partial w_{22}} \\ \hline \end{array}$$

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$$\text{conv}(\mathbf{x}, \mathbf{w})$$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$$\text{p}(\mathbf{y})$$

$$p$$

upstream gradient

$$\text{vjp}_{\text{conv\_w}}(\frac{\partial \mathbf{p}}{\partial \mathbf{y}}, \mathbf{x})$$

$$\frac{\partial \mathbf{p}}{\partial \mathbf{y}} = \begin{array}{|c|c|} \hline \dfrac{\partial p}{\partial y_{11}} & \dfrac{\partial p}{\partial y_{12}} \\ \hline \dfrac{\partial p}{\partial y_{21}} & \dfrac{\partial p}{\partial y_{22}} \\ \hline \end{array}$$

$$\text{vjp}_{\text{conv\_w}}(\frac{\partial \mathbf{p}}{\partial \mathbf{y}}, \mathbf{x}) = \begin{array}{|c|c|} \hline \dfrac{\partial p}{\partial w_{11}} & \dfrac{\partial p}{\partial w_{12}} \\ \hline \dfrac{\partial p}{\partial w_{21}} & \dfrac{\partial p}{\partial w_{22}} \\ \hline \end{array} = \text{conv} \left( \begin{array}{|c|c|c|} \hline x_{11} & x_{12} & x_{13} \\ \hline x_{21} & x_{22} & x_{23} \\ \hline x_{31} & x_{32} & x_{33} \\ \hline \end{array}, \begin{array}{|c|c|} \hline \dfrac{\partial p}{\partial y_{11}} & \dfrac{\partial p}{\partial y_{12}} \\ \hline \dfrac{\partial p}{\partial y_{21}} & \dfrac{\partial p}{\partial y_{22}} \\ \hline \end{array} \right)$$

- Backpropagation of convolutional layer wrt weights is defined as: "**convolution of input feature map with upstream gradient**"

# Convolution backward pass wrt feature map

| $w_{11}$ | $w_{12}$ |
|---|---|
| $w_{21}$ | $w_{22}$ |

| $x_{11}$ | $x_{12}$ | $x_{13}$ |
|---|---|---|
| $x_{21}$ | $x_{22}$ | $x_{23}$ |
| $x_{31}$ | $x_{32}$ | $x_{33}$ |

$$\text{conv}(\mathbf{x}, \mathbf{w})$$

| $y_{11}$ | $y_{12}$ |
|---|---|
| $y_{21}$ | $y_{22}$ |

$$\text{p}(\mathbf{y})$$

$$p$$

$$\frac{\partial \mathbf{p}}{\partial \mathbf{x}} = \begin{array}{|c|c|c|} \hline \frac{\partial p}{\partial x_{11}} & \frac{\partial p}{\partial x_{12}} & \frac{\partial p}{\partial x_{13}} \\ \hline \frac{\partial p}{\partial x_{21}} & \frac{\partial p}{\partial x_{22}} & \frac{\partial p}{\partial x_{23}} \\ \hline \frac{\partial p}{\partial x_{31}} & \frac{\partial p}{\partial x_{32}} & \frac{\partial p}{\partial x_{33}} \\ \hline \end{array}$$

$$\text{vjp}_{\text{conv\_x}}(\frac{\partial \mathbf{p}}{\partial \mathbf{y}}, \mathbf{w})$$

$$\frac{\partial \mathbf{p}}{\partial \mathbf{y}} = \begin{array}{|c|c|} \hline \frac{\partial p}{\partial y_{11}} & \frac{\partial p}{\partial y_{12}} \\ \hline \frac{\partial p}{\partial y_{21}} & \frac{\partial p}{\partial y_{22}} \\ \hline \end{array}$$

upstream
gradient

$$\text{vjp}_{\text{conv\_x}}(\frac{\partial \mathbf{p}}{\partial \mathbf{y}}) = \begin{array}{|c|c|c|} \hline \frac{\partial p}{\partial x_{11}} & \frac{\partial p}{\partial x_{12}} & \frac{\partial p}{\partial x_{13}} \\ \hline \frac{\partial p}{\partial x_{21}} & \frac{\partial p}{\partial x_{22}} & \frac{\partial p}{\partial x_{23}} \\ \hline \frac{\partial p}{\partial x_{31}} & \frac{\partial p}{\partial x_{32}} & \frac{\partial p}{\partial x_{33}} \\ \hline \end{array} = \text{conv}\left( \begin{array}{|c|c|c|c|} \hline 0 & 0 & 0 & 0 \\ \hline 0 & \frac{\partial p}{\partial y_{11}} & \frac{\partial p}{\partial y_{12}} & 0 \\ \hline 0 & \frac{\partial p}{\partial y_{21}} & \frac{\partial p}{\partial y_{22}} & 0 \\ \hline 0 & 0 & 0 & 0 \\ \hline \end{array}, \begin{array}{|c|c|} \hline w_{22} & w_{21} \\ \hline w_{12} & w_{11} \\ \hline \end{array} \right)$$

- Backpropagation of convolutional layer wrt input feature map is defined as:
  "**convolution of padded upstream gradient with mirrored weights**"

Convolution backward pass wrt input feature map

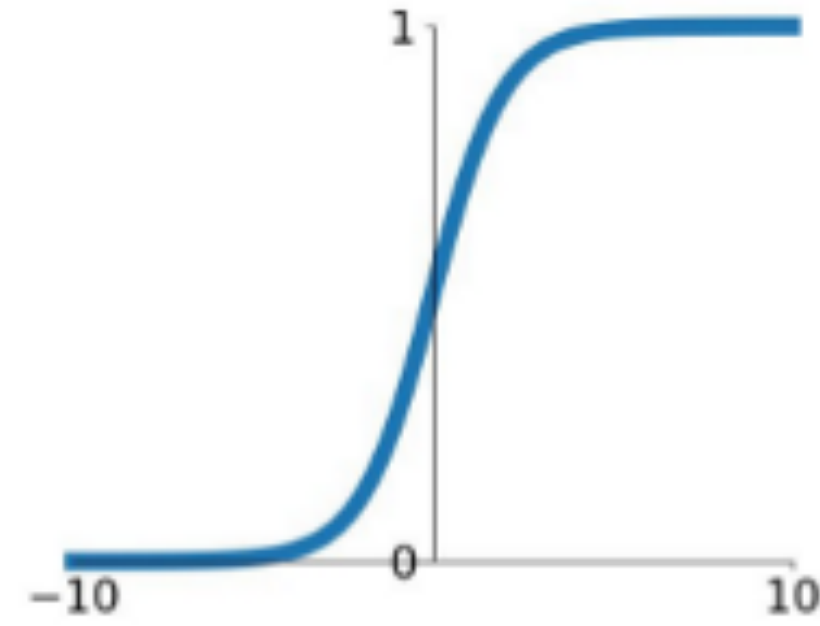Very important property of convolutional layer is:
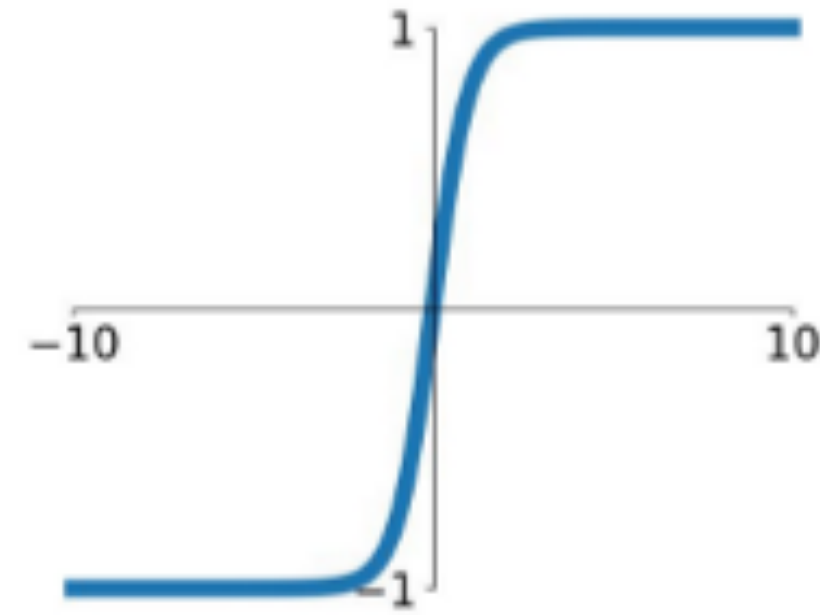
**Backpropagation is also convolution !!!**

**Sigmoid**

$$\sigma(x) = \frac{1}{1+e^{-x}}$$

**tanh**

$$\tanh(x)$$
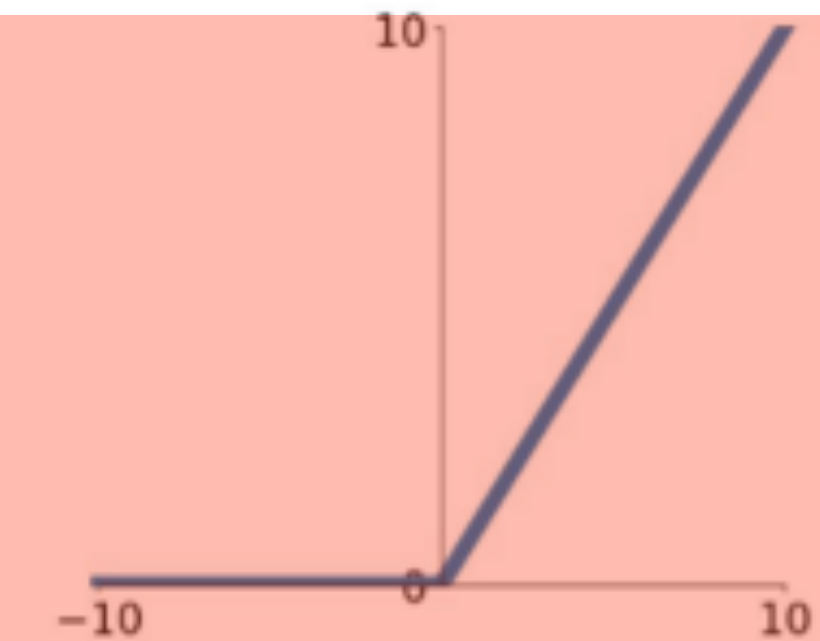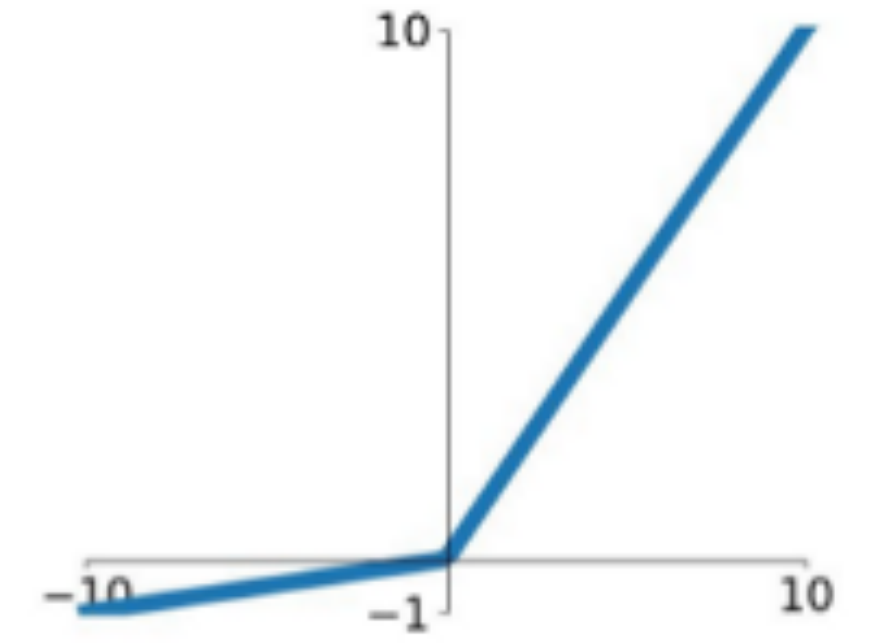
**ReLU**

$$\max(0, x)$$
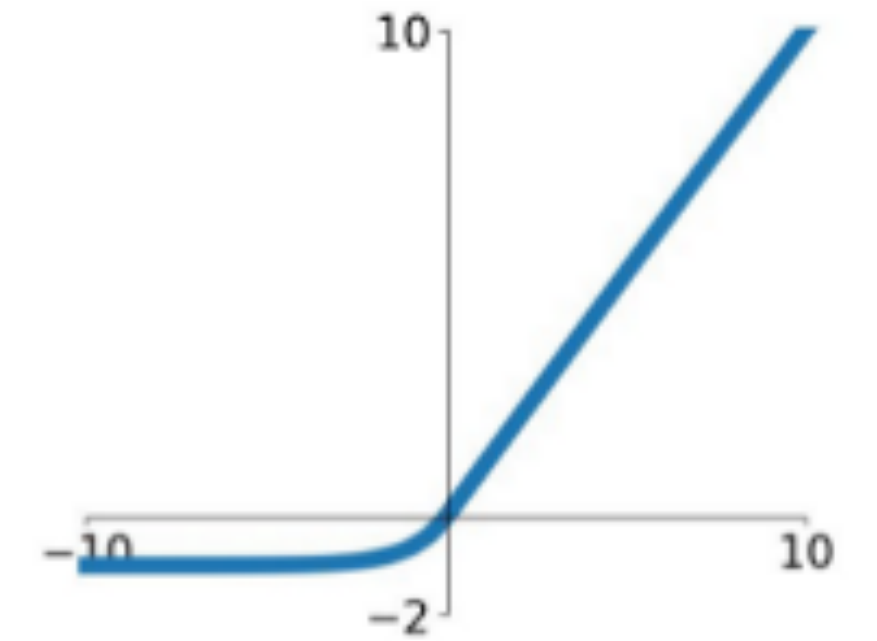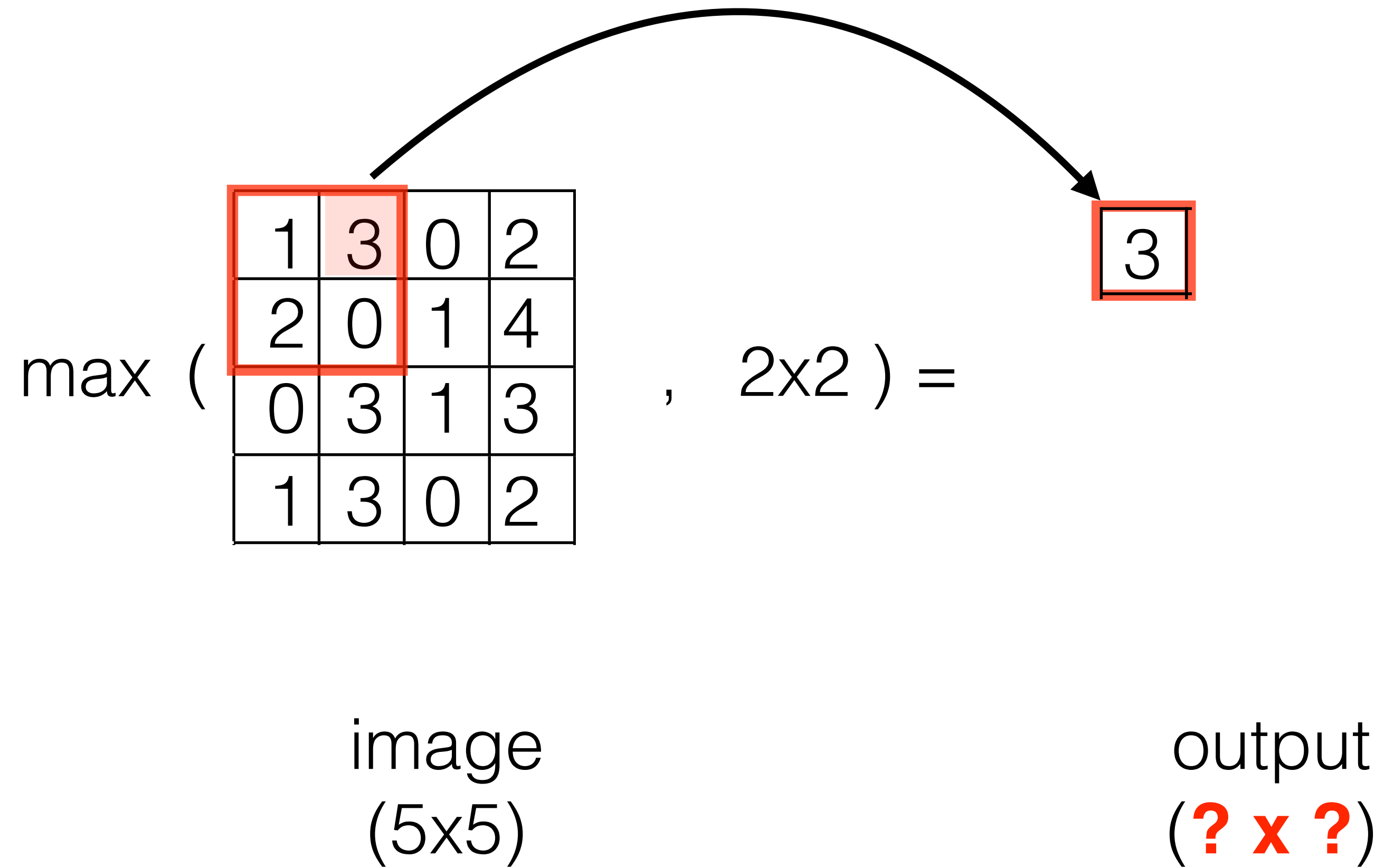
**Leaky ReLU**

$$\max(0.1x, x)$$

**Maxout**

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

**ELU**

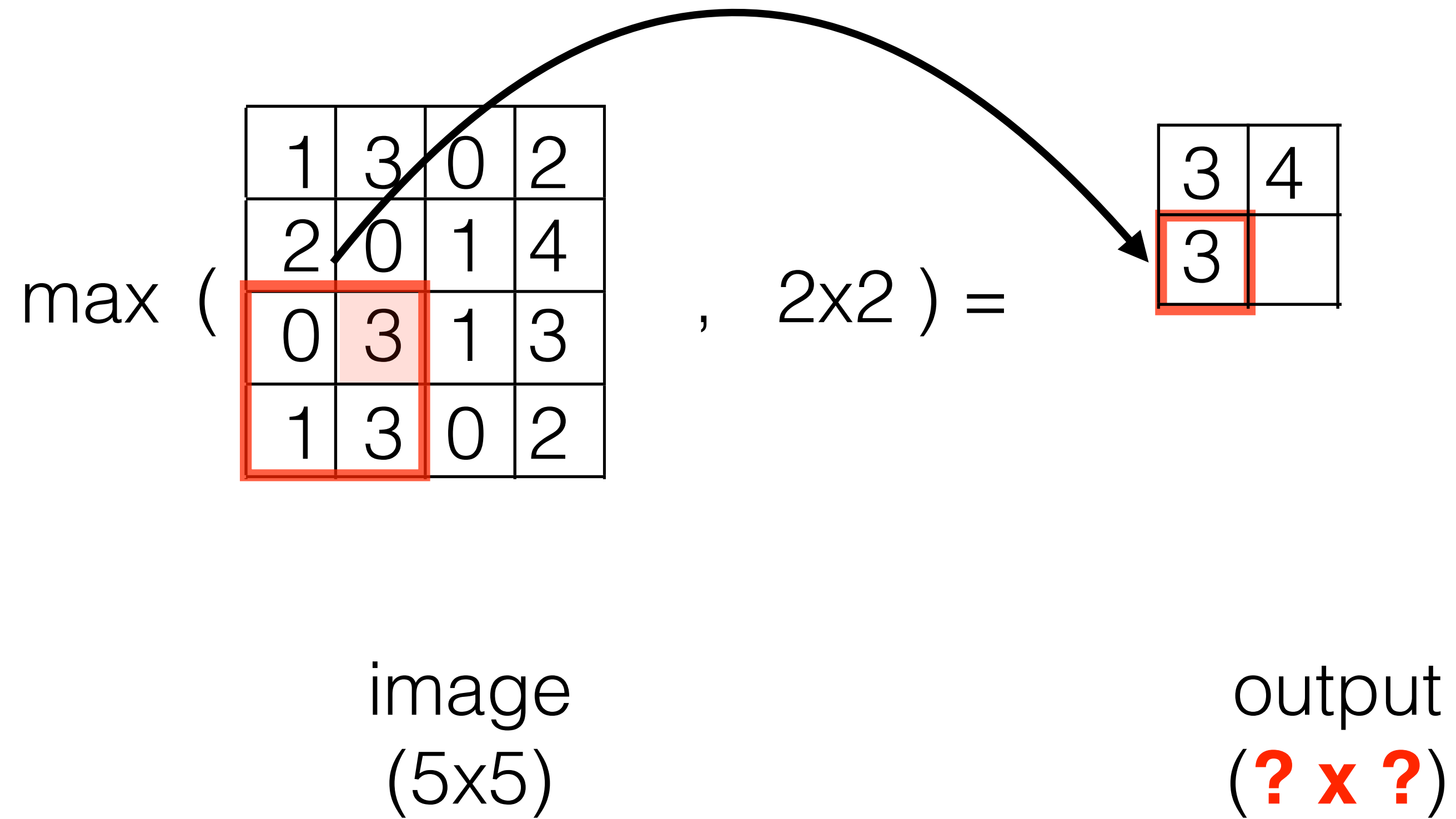$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$

# Max-pooling

max (

|   |   |   |   |
|---|---|---|---|
| 1 | 3 | 0 | 2 |
| 2 | 0 | 1 | 4 |
| 0 | 3 | 1 | 3 |
| 1 | 3 | 0 | 2 |

, 2x2 ) =

| 3 |
|---|

image
(5x5)

output
(**? x ?**)

# Max-pooling

max (
| 1 | 3 | 0 | 2 |
|---|---|---|---|
| 2 | 0 | 1 | 4 |
| 0 | 3 | 1 | 3 |
| 1 | 3 | 0 | 2 |
, 2x2 ) =

| 3 | 4 |
|---|---|

image
(5x5)

output
(**? x ?**)

# Max-pooling



$$\max \left( \begin{array}{cccc} 1 & 3 & 0 & 2 \\ 2 & 0 & 1 & 4 \\ 0 & 3 & 1 & 3 \\ 1 & 3 & 0 & 2 \end{array} , \quad 2x2 \right) = \begin{array}{cc} 3 & 4 \\ 3 & \end{array}$$

image
(5x5)

output
(**? x ?**)

# Max-pooling

$$\max \left( \begin{array}{cccc} 1 & 3 & 0 & 2 \\ 2 & 0 & 1 & 4 \\ 0 & 3 & 1 & 3 \\ 1 & 3 & 0 & 2 \end{array} \right., \ 2x2 \ ) = \begin{array}{cc} 3 & 4 \\ 3 & 3 \end{array}$$

image
(5x5)

output
(**? x ?**)

# Learning of a simple convolutional network
## input 3D tensor: channels x height x width



$\mathbf{x}$

$\mathbf{w}_1$

$\mathbf{w}_2$

$\mathbf{w}_3$

$\mathbf{y}$

$\text{conv}(\mathbf{x}, \mathbf{w}_1)$ — $\mathbf{y}_1$

$\text{conv}(\mathbf{x}, \mathbf{w}_2)$ — $\mathbf{y}_2$

$\text{conv}(\mathbf{x}, \mathbf{w}_3)$ — $\mathbf{y}_3$

3x5x5

feature
map

conv1
layer

3x4x4

feature
map

$\mathbf{x}$

$\mathbf{w}_1$

$\mathbf{w}_2$

$\mathbf{w}_3$

$\text{conv}(\mathbf{x}, \mathbf{w}_1)$ $\mathbf{y}_1$

$\text{conv}(\mathbf{x}, \mathbf{w}_2)$ $\mathbf{y}_2$

$\text{conv}(\mathbf{x}, \mathbf{w}_3)$ $\mathbf{y}_3$

$\mathbf{y}$

$\sigma(\mathbf{y})$

$\mathbf{u}$

3x5x5

3x4x4

3x4x4

feature map

feature map

feature map

conv1 layer

activ. function

87

Learning of a simple convolutional network
**input 3D tensor: channels x height x width**

$\mathbf{x}$

$\mathbf{w}_1$

$\mathbf{w}_2$

$\mathbf{w}_3$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_1)$ $\mathbf{y}_1$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_2)$ $\mathbf{y}_2$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_3)$ $\mathbf{y}_3$

$\mathbf{y}$

$\sigma(\mathbf{y})$

$\mathbf{u}$

$\mathrm{maxp}(u)$

$\mathbf{z}$

3x5x5

3x4x4

3x4x4

3x2x2

feature map

feature map

feature map

feature map

conv1 layer

activ. function

maxp layer

# Learning of a simple convolutional network
## input 3D tensor: channels x height x width



$\mathbf{x}$

$\mathbf{w}_1$

$\text{conv}(\mathbf{x}, \mathbf{w}_1)$ $\quad \mathbf{y}_1$

$\mathbf{w}_2$

$\text{conv}(\mathbf{x}, \mathbf{w}_2)$ $\quad \mathbf{y}_2$

$\mathbf{w}_3$

$\text{conv}(\mathbf{x}, \mathbf{w}_3)$ $\quad \mathbf{y}_3$

$\mathbf{y}$

$\sigma(\mathbf{y})$

$\mathbf{u}$

$\text{maxp(u)}$

$\mathbf{z}$

$\text{conv}(\mathbf{z}, \mathbf{v_1})$

$\text{conv}(\mathbf{z}, \mathbf{v_2})$

$\text{conv}(\mathbf{z}, \mathbf{v_3})$

$\text{conv}(\mathbf{z}, \mathbf{v_4})$

3x5x5

feature map

3x4x4

feature map

3x4x4

feature map

3x2x2

feature map

4x1x1

feature map

conv1 layer

activ. function

maxp layer

layer: conv2

$\mathbf{x}$

$\mathbf{w}_1$

$\mathbf{w}_2$

$\mathbf{w}_3$

$\mathbf{y}$

$\mathbf{u}$

$\mathbf{z}$

$\text{conv}(\mathbf{x}, \mathbf{w}_1)$   $\mathbf{y}_1$

$\text{conv}(\mathbf{x}, \mathbf{w}_2)$   $\mathbf{y}_2$

$\text{conv}(\mathbf{x}, \mathbf{w}_3)$   $\mathbf{y}_3$

$\sigma(\mathbf{y})$

$\text{maxp(u)}$

$\text{conv}(\mathbf{z}, \mathbf{v_1})$

$\text{conv}(\mathbf{z}, \mathbf{v_2})$

$\text{conv}(\mathbf{z}, \mathbf{v_3})$

$\text{conv}(\mathbf{z}, \mathbf{v_4})$

3x5x5

3x4x4

3x4x4

3x2x2

4x1x1

feature map

feature map

feature map

feature map

feature map

conv1 layer

activ. function

maxp layer

layer: conv2

# input 4D tensor: batch_size x channels x height x width



$\mathbf{x}$

$\mathbf{w}_1$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_1)$  $\mathbf{y}_1$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_2)$  $\mathbf{y}_2$

$\mathbf{w}_2$

$\mathbf{w}_3$

$\mathrm{conv}(\mathbf{x}, \mathbf{w}_3)$  $\mathbf{y}_3$

$\mathbf{y}$

$\sigma(\mathbf{y})$

$\mathbf{u}$

$\mathrm{maxp(u)}$

$\mathbf{z}$

$\mathrm{conv}(\mathbf{z}, \mathbf{v_1})$

$\mathrm{conv}(\mathbf{z}, \mathbf{v_2})$

$\mathrm{conv}(\mathbf{z}, \mathbf{v_3})$

$\mathrm{conv}(\mathbf{z}, \mathbf{v_4})$

Bx3x5x5

Bx3x4x4

Bx3x4x4

Bx3x2x2

Bx4x1x1

Batch
of
images

conv1
layer

Batch
of
feature
maps

activ.
function

Batch
of
feature
maps

maxp
layer

Batch
of
feature
maps

layer: conv2

Batch
of
feature
maps

91

# Convolutional net

- **Convolutional network** (ConvNet) is concatenation of convolutional layers, activation function and optionally max-pooling functions.
- **Backprop in convolutional layer** is convolution of feature maps or kernels or feature-maps with the upstream gradient.
- Feed-forward and backprop are convolutions => efficient implementation on GPU

# Kunihiko Fukushima 1980

## Neocognitron: A Self-organizing Neural Network Model for a Mechanism of Pattern Recognition Unaffected by Shift in Position

Kunihiko Fukushima

NHK Broadcasting Science Research Laboratories, Kinuta, Setagaya, Tokyo, Japan

**Abstract.** A neural network model for a mechanism of visual pattern recognition is proposed in this paper. The network is self-organized by "learning without a teacher", and acquires an ability to recognize stimulus patterns based on the geometrical similarity (Gestalt) of their shapes without affected by their positions. This network is given a nickname "neocognitron". After completion of self-organization, the network has a structure similar to the hierarchy model of the visual reveal it only by conventional physiological experiments. So, we take a slightly different approach to this problem. If we could make a neural network model which has the same capability for pattern recognition as a human being, it would give us a powerful clue to the understanding of the neural mechanism in the brain. In this paper, we discuss how to synthesize a neural network model in order to endow it an ability of pattern recognition like a human being.



https://en.wikipedia.org/wiki/Kunihiko_Fukushima

# LeCun's letter recognition 1998 (over 13k citations !!!)

## backpropagation formulated



LeCun et al, Gradient based learning applied to document recognition, IEEE, 1998
http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf

# AlexNet on ImageNet 2012 (**over 27k citations !!!**)



Alex Krizhevsky et al, Imagenet classification with deep convolutional neural networks, NIPS, 2012
https://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf

# IM𝄏GENET

1.2M images with (227x227x3) resolution

http://image-net.org/challenges/LSVRC/2017/index

**Steel drum**



| Output: | | Output: |
| --- | --- | --- |
| Scale | | Scale |
| T-shirt | | T-shirt |
| Steel drum | ✔ | Giant panda |
| Drumstick | | Drumstick |
| Mud turtle | ✘ | Mud turtle |

$$\text{Error} = \frac{1}{100,000} \sum_{100,000 \text{ images}} 1[\text{incorrect on image i}]$$

Pascal VOC object detection challenge

# Under the hood of ConvNets

# Convolution as spatial attention

FCNN = "global hard attention"        Convolution = "local hard attention"



Fully-connected layer                    Convolutional layer

Do you see any other suitable attentions?

Trained kernels

$\mathbf{x}$

ConvNet

$\mathbf{y}$

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

# Trained kernels

$\mathbf{x}$

ConvNet

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

$\mathbf{y}$

kernels/filters

Kernel for Conv 0 (Channelwise Slices)

….

Direct visualization of trained kernels does not say too much

Feature maps = low dimensional encoding



x

ConvNet

$f_1$ $f_2$ $f_3$ $f_4$ $f_5$ $f_6$ $f_7$ $f_8$ $f_9$

feedforward pass

y 0.2 dog
0.1 ship
**0.5** cat
0.0 car
0.1 airplane

Feature maps:
intermediate results
of feedforward pass

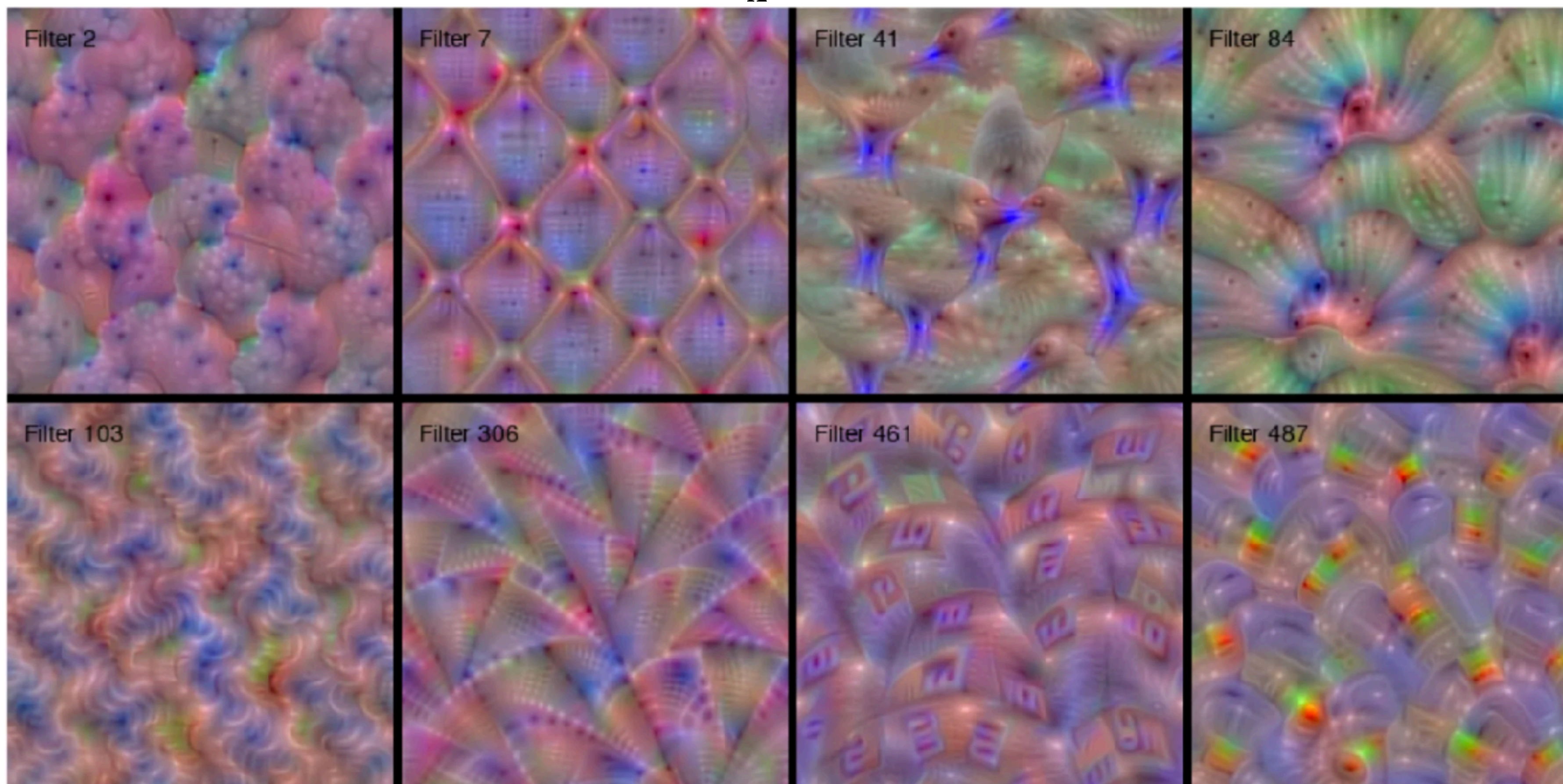# Feature maps = low dimensional encoding

Features maps that maximize filter output
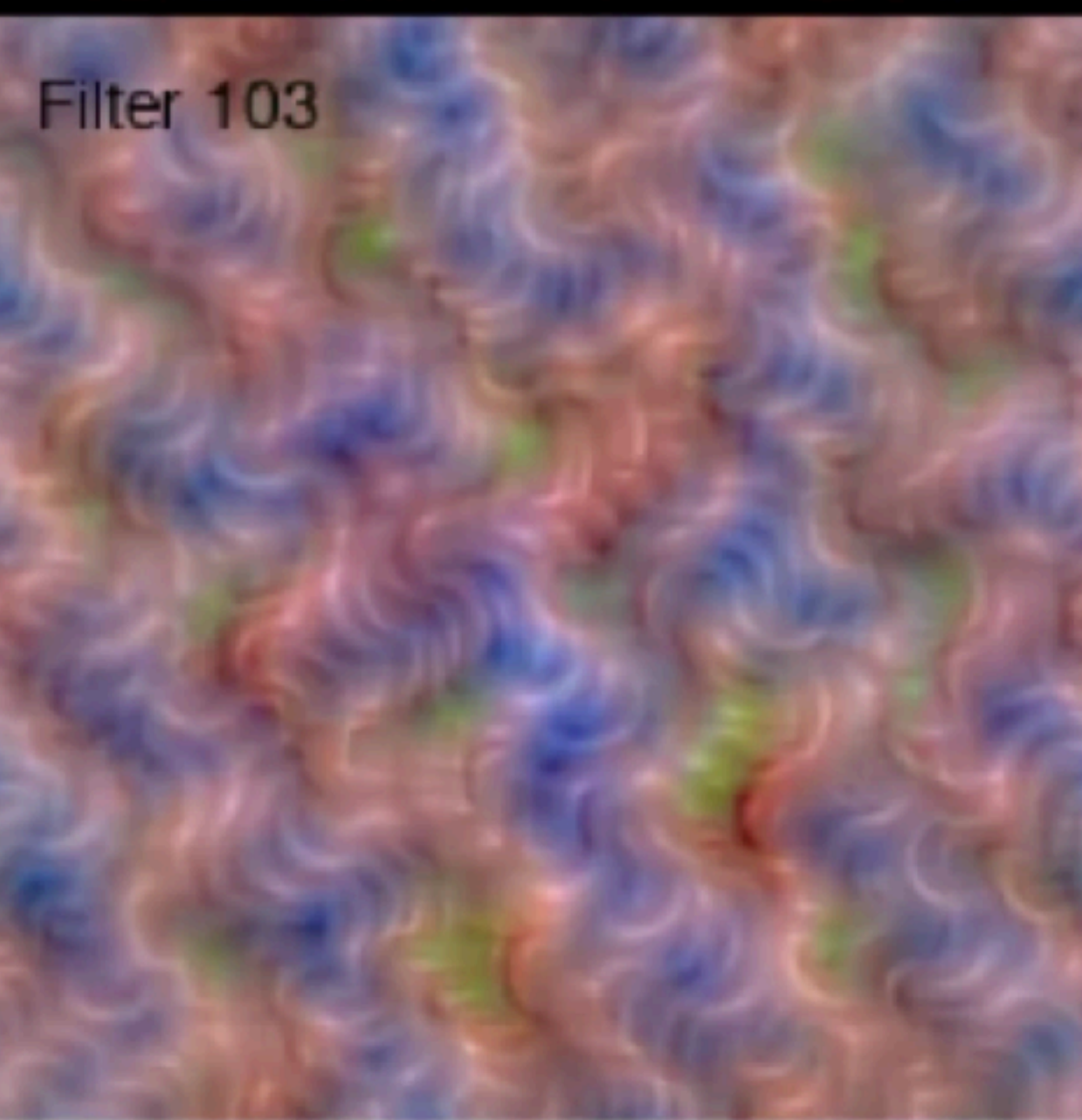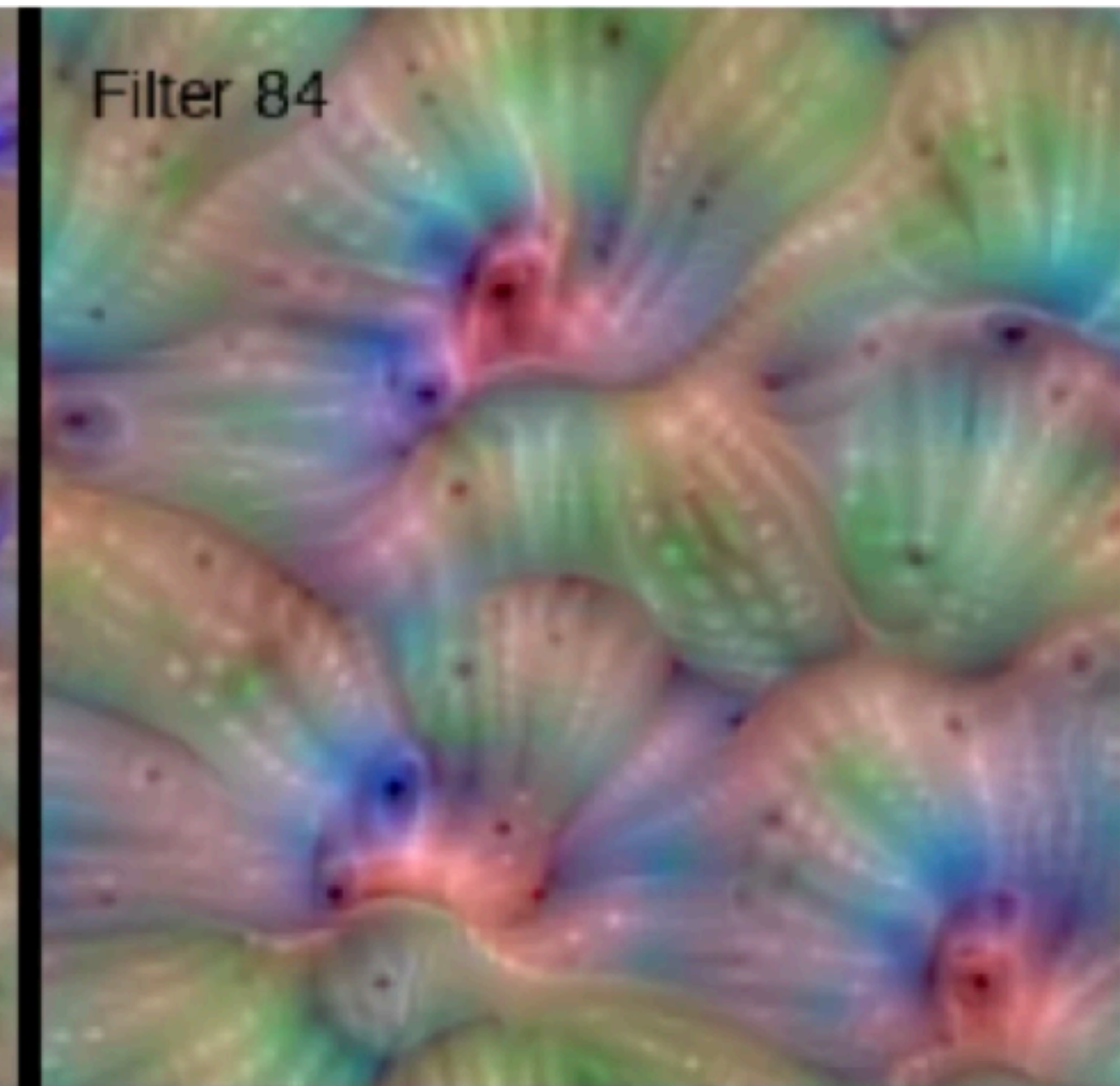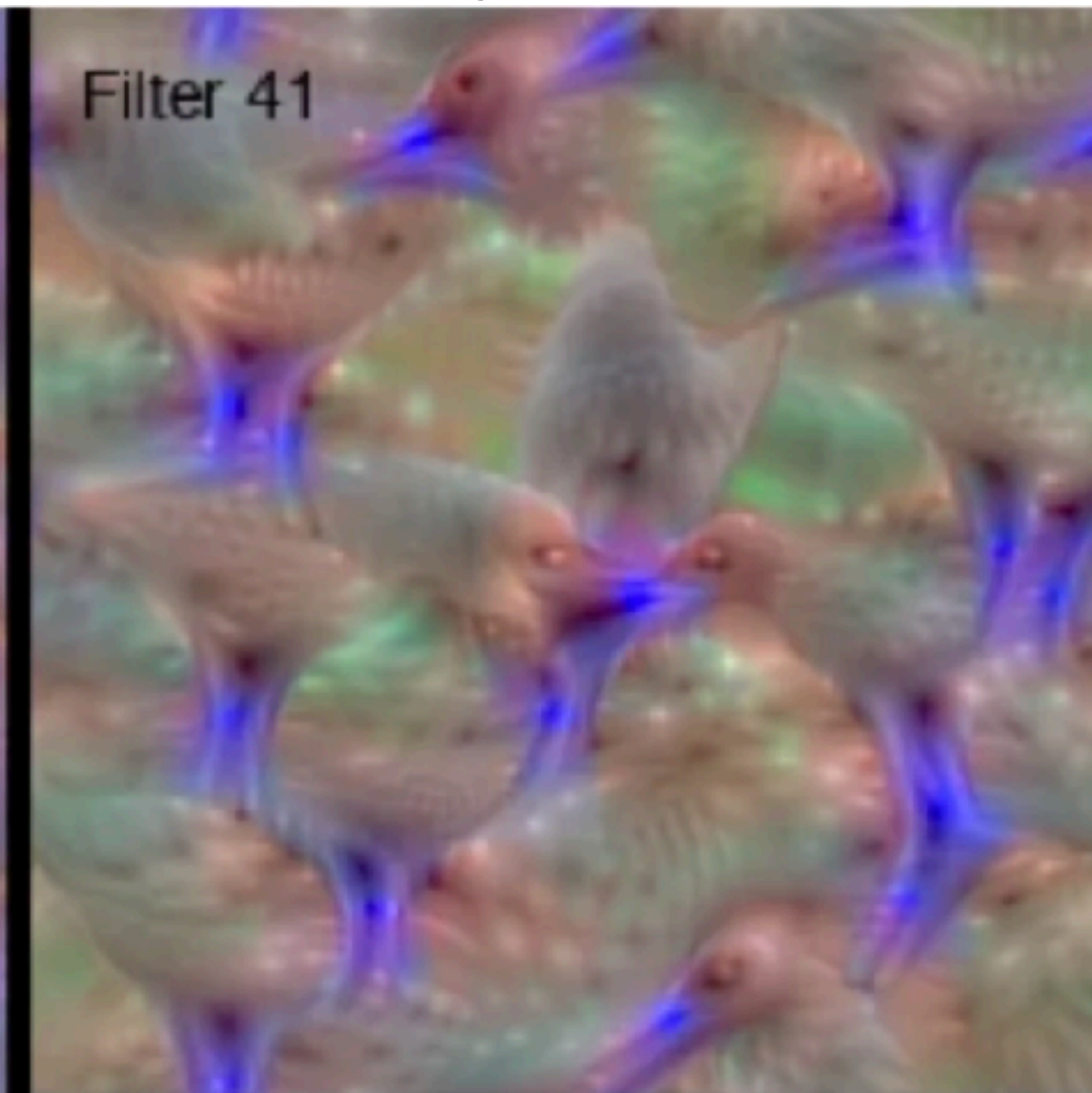
ConvNet

What is the input that maximizes a filter output?

$$\mathbf{x} \longrightarrow$$

$$\mathbf{f}_1 \; \mathbf{f}_2 \; \mathbf{f}_3 \; \mathbf{f}_4 \; \mathbf{f}_5 \; \mathbf{f}_6 \; \mathbf{f}_7 \; \mathbf{f}_8 \; \mathbf{f}_9 \longrightarrow \mathbf{y}$$

$$\frac{\partial \mathbf{f}_1}{\partial \mathbf{x}}$$

$$\arg \max_{\mathbf{x}} \|\mathbf{f}_1(\mathbf{x}, \mathbf{w})\|$$



Filter 0    Filter 3    Filter 11    Filter 25

Filter 33    Filter 42    Filter 44    Filter 62

Features maps that maximize filter output

ConvNet

**x**

What is the input that
maximizes a filter output?

**y**

$$\frac{\partial \mathbf{f}_2}{\partial \mathbf{x}}$$

$$\arg\max_{\mathbf{x}} \|\mathbf{f}_2(\mathbf{x}, \mathbf{w})\|$$



Filter 8   Filter 21   Filter 23   Filter 38

Filter 39   Filter 45   Filter 50   Filter 79

Features maps that maximize filter output

ConvNet

What is the input that
maximizes a filter output?

$\mathbf{x}$

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

$\mathbf{y}$

$\dfrac{\partial \mathbf{f}_3}{\partial \mathbf{x}}$

$$\arg\max_{\mathbf{x}} \|\mathbf{f}_3(\mathbf{x}, \mathbf{w})\|$$



Filter 40

Filter 48

Filter 52

Filter 54

Filter 81

Filter 107

Filter 224

Filter 226

Features maps that maximize filter output

ConvNet

$\mathbf{x}$

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

$\mathbf{y}$

$\dfrac{\partial \mathbf{f}_4}{\partial \mathbf{x}}$

What is the input that maximizes a filter output?

$$\arg\max_{\mathbf{x}} \|\mathbf{f}_4(\mathbf{x}, \mathbf{w})\|$$

Features maps that maximize filter output

ConvNet

$\mathbf{x}$

$\mathbf{f}_1 \ \mathbf{f}_2 \ \mathbf{f}_3 \ \mathbf{f}_4 \ \mathbf{f}_5 \ \mathbf{f}_6 \ \mathbf{f}_7 \ \mathbf{f}_8 \ \mathbf{f}_9$

$\mathbf{y}$

$\dfrac{\partial \mathbf{f}_5}{\partial \mathbf{x}}$

$\arg \max_{\mathbf{x}} \| \mathbf{f}_5(\mathbf{x}, \mathbf{w}) \|$

What is the input that maximizes a filter output?



Filter 2    Filter 7    Filter 41    Filter 84

Filter 103    Filter 306    Filter 461    Filter 487

Features maps that maximize filter output
Can you interpret functionality of a filter?

Filter 2    Filter 7    Filter 41    Filter 84

Filter 103    Filter 306    Filter 461    Filter 487
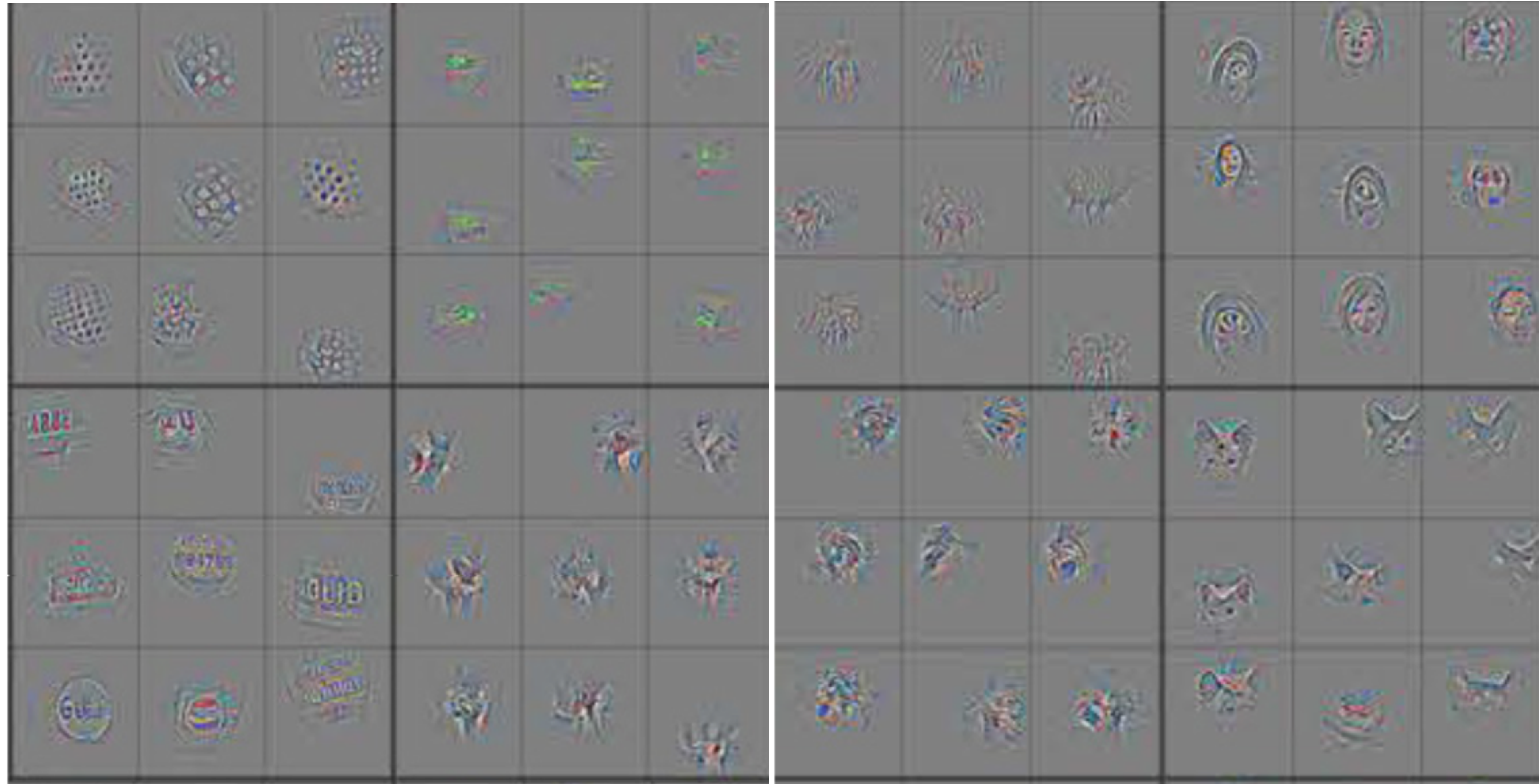
3. Neurons are sensitive to edges and its orientation

Inputs which maximized output of **layer 1**



[Zeiler and Fergus, ECCV, 2014]

3. Neurons are sensitive to edges and its orientation

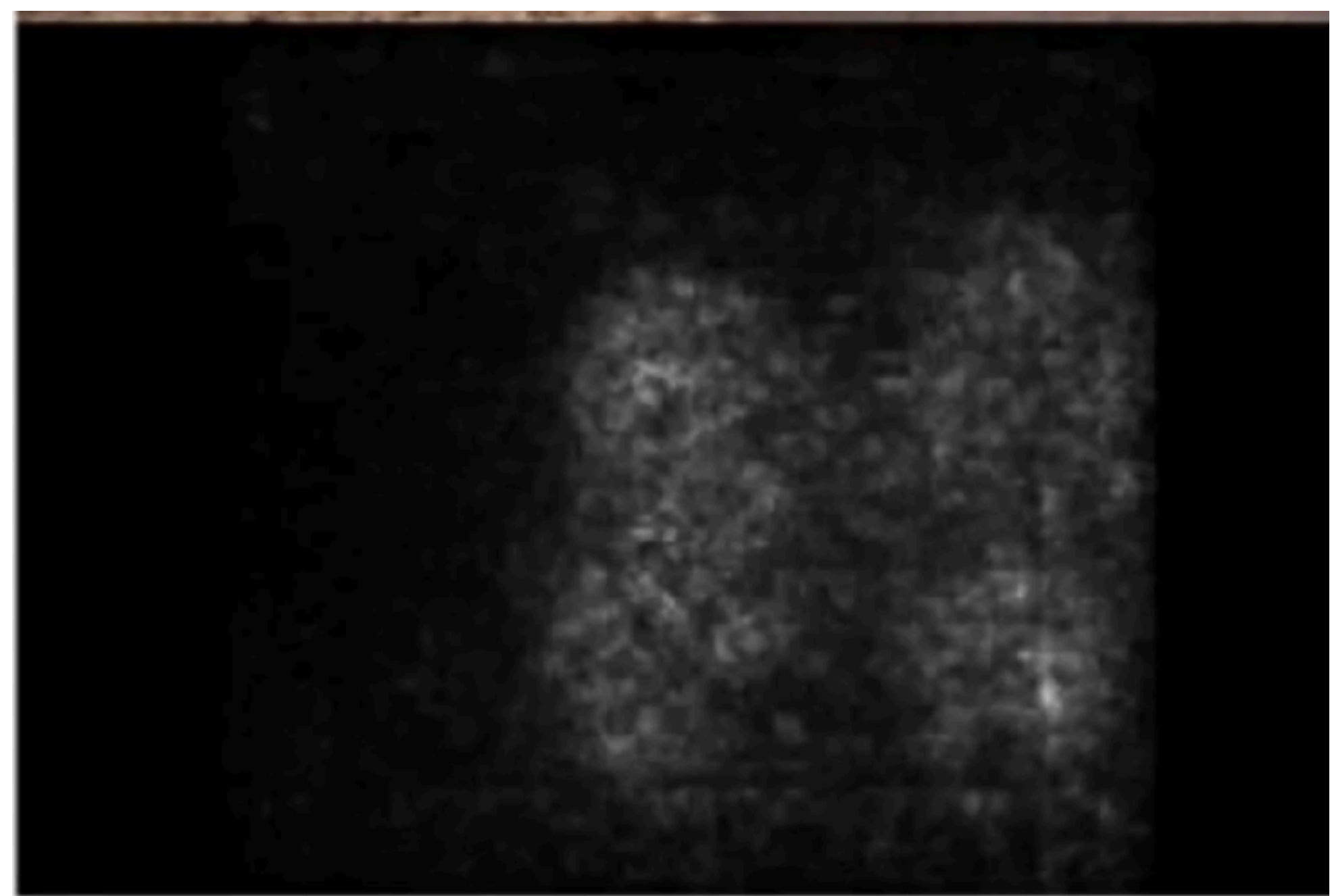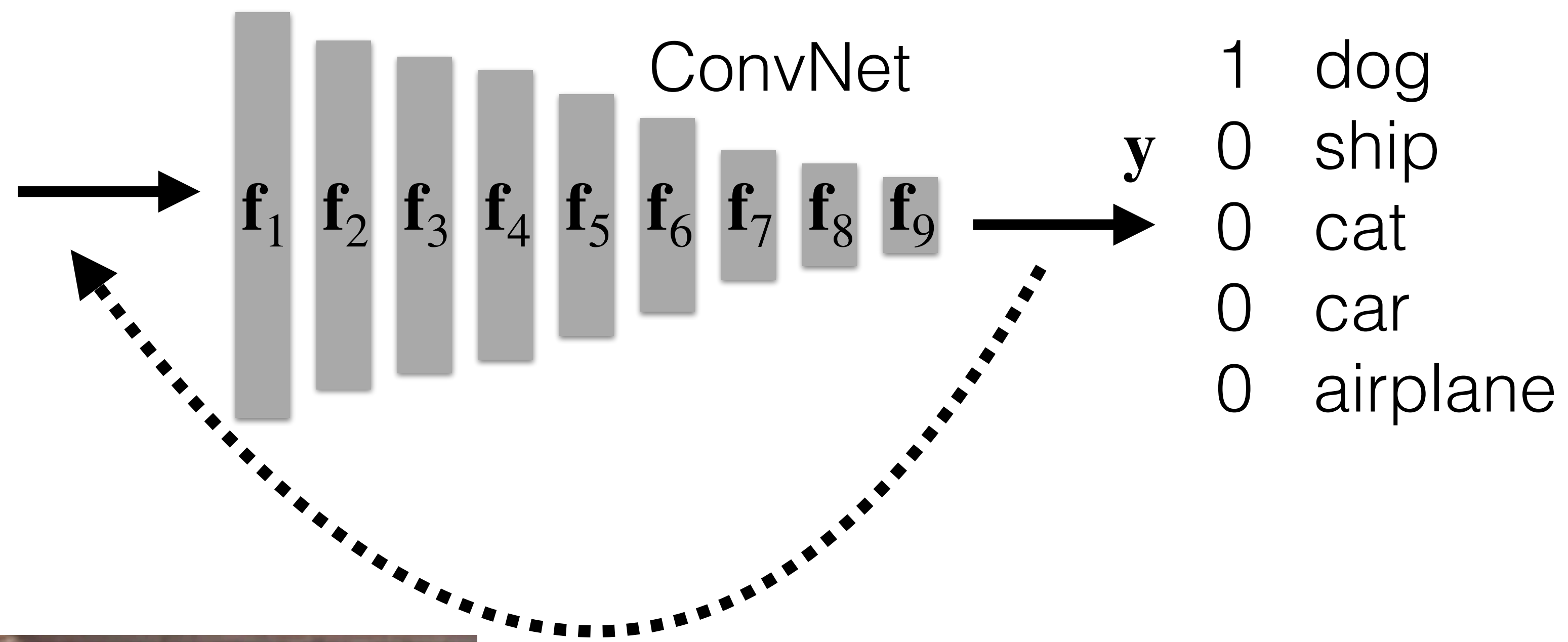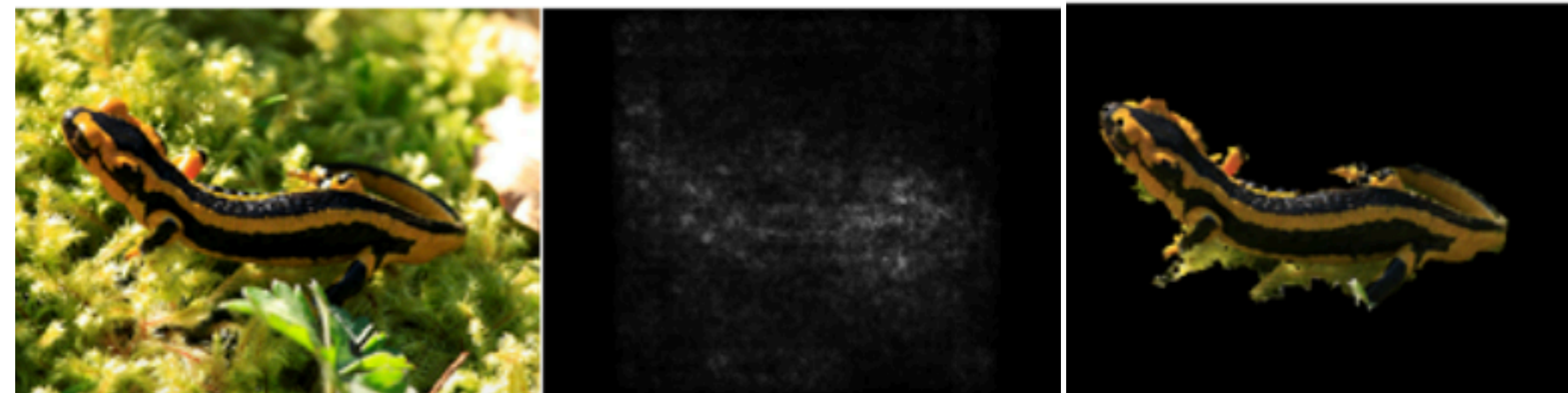Inputs which maximized output of **layer 2**



[Zeiler and Fergus, ECCV, 2014]

3. Neurons are sensitive to edges and its orientation

Inputs which maximized output of **layer 3**



[Zeiler and Fergus, ECCV, 2014]

3. Neurons are sensitive to edges and its orientation

Inputs which maximized output of **layer 4**



[Zeiler and Fergus, ECCV, 2014]

3. Neurons are sensitive to edges and its orientation

Inputs which maximized output of **layer 5**
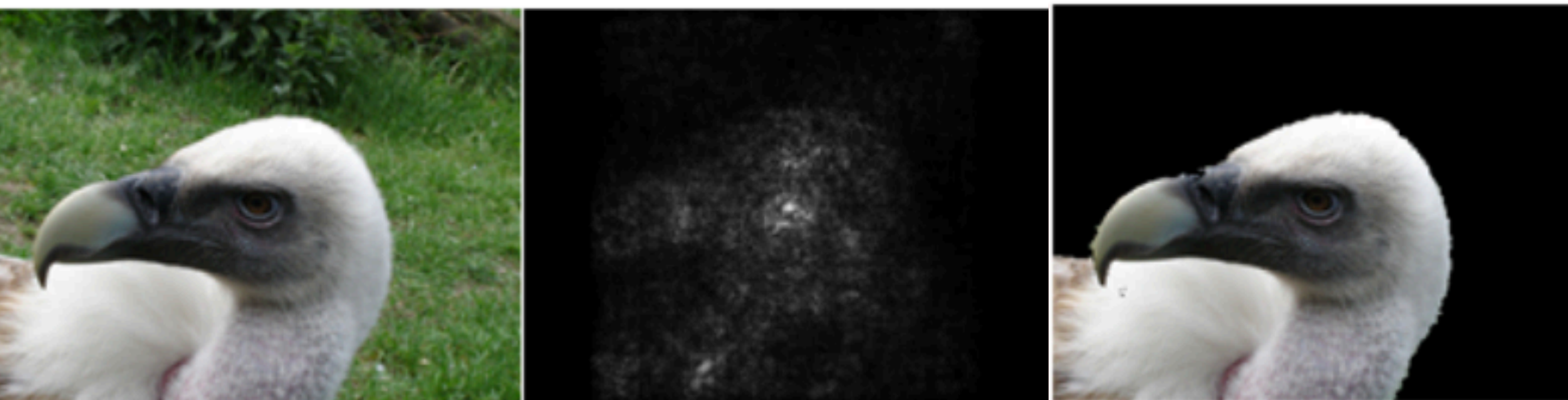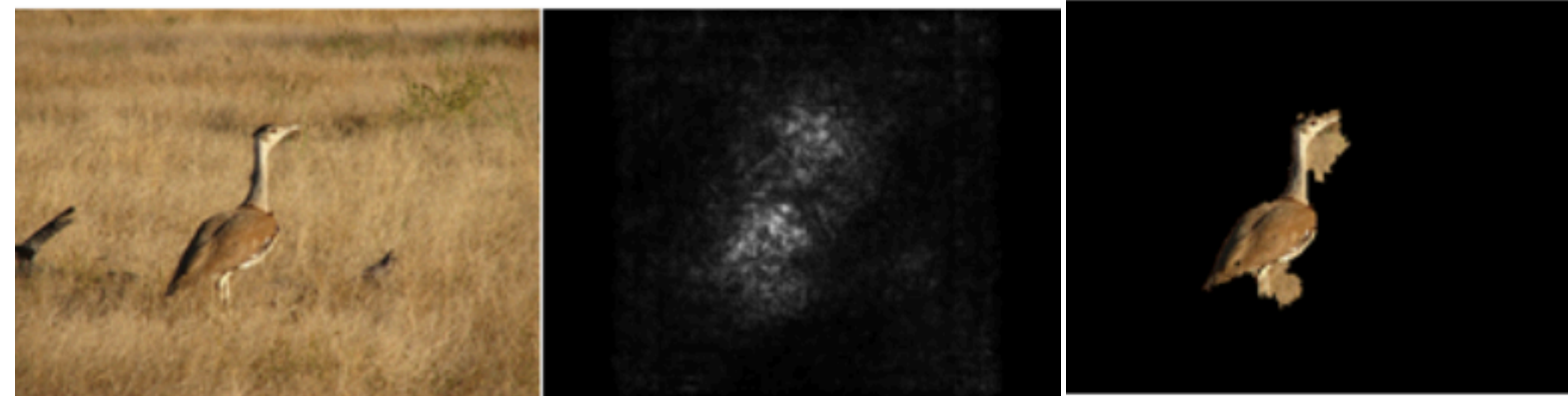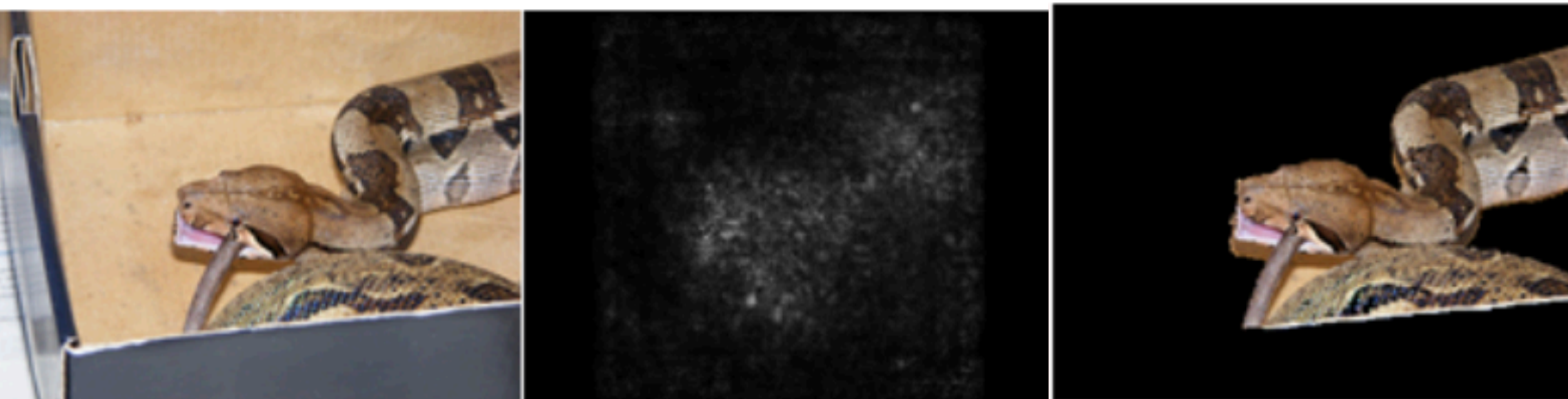


[Zeiler and Fergus, ECCV, 2014]

Saliency maps

ConvNet

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$
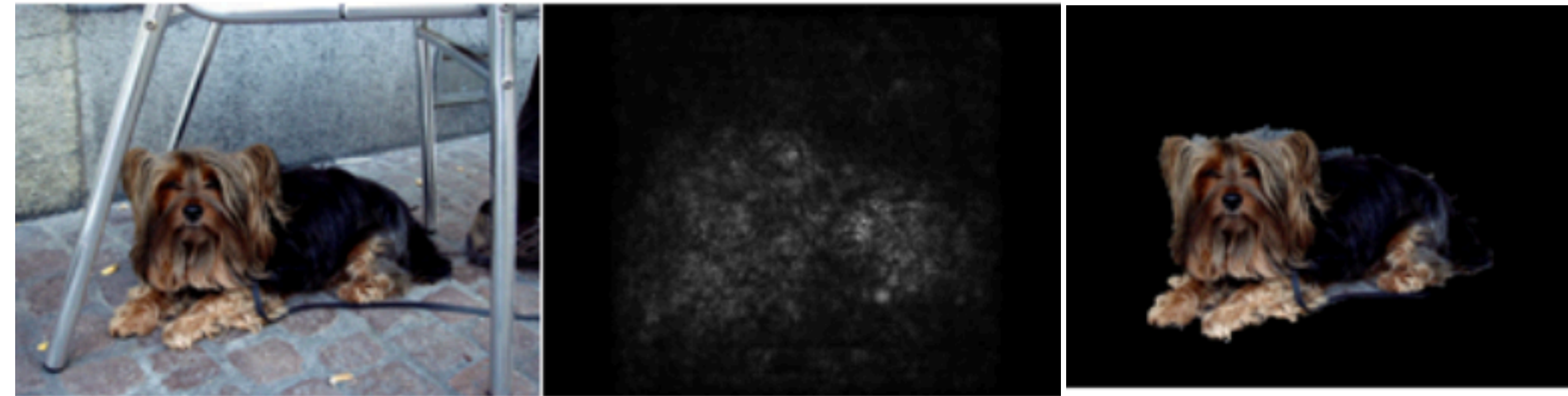
$\mathbf{y}$

1   dog
0   ship
0   cat
0   car
0   airplane

$\dfrac{\partial \mathbf{y}_1}{\partial \mathbf{x}}$

What pixels contributed the most for dog category on the output
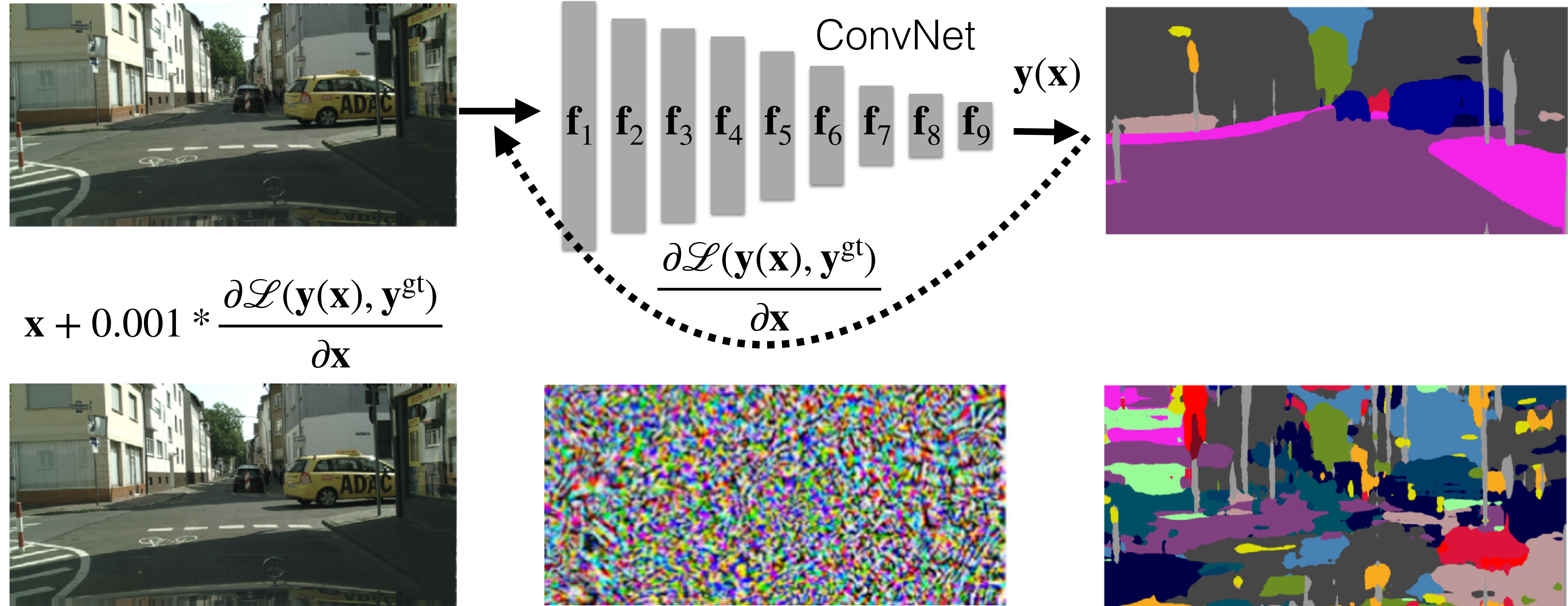
# Saliency maps
https://arxiv.org/pdf/1312.6034.pdf

# Adversarial attacks [Arnab CVPR 2018]

- Given trained network and an image, find the closest image on which the net fails

$\mathbf{x}$



ConvNet
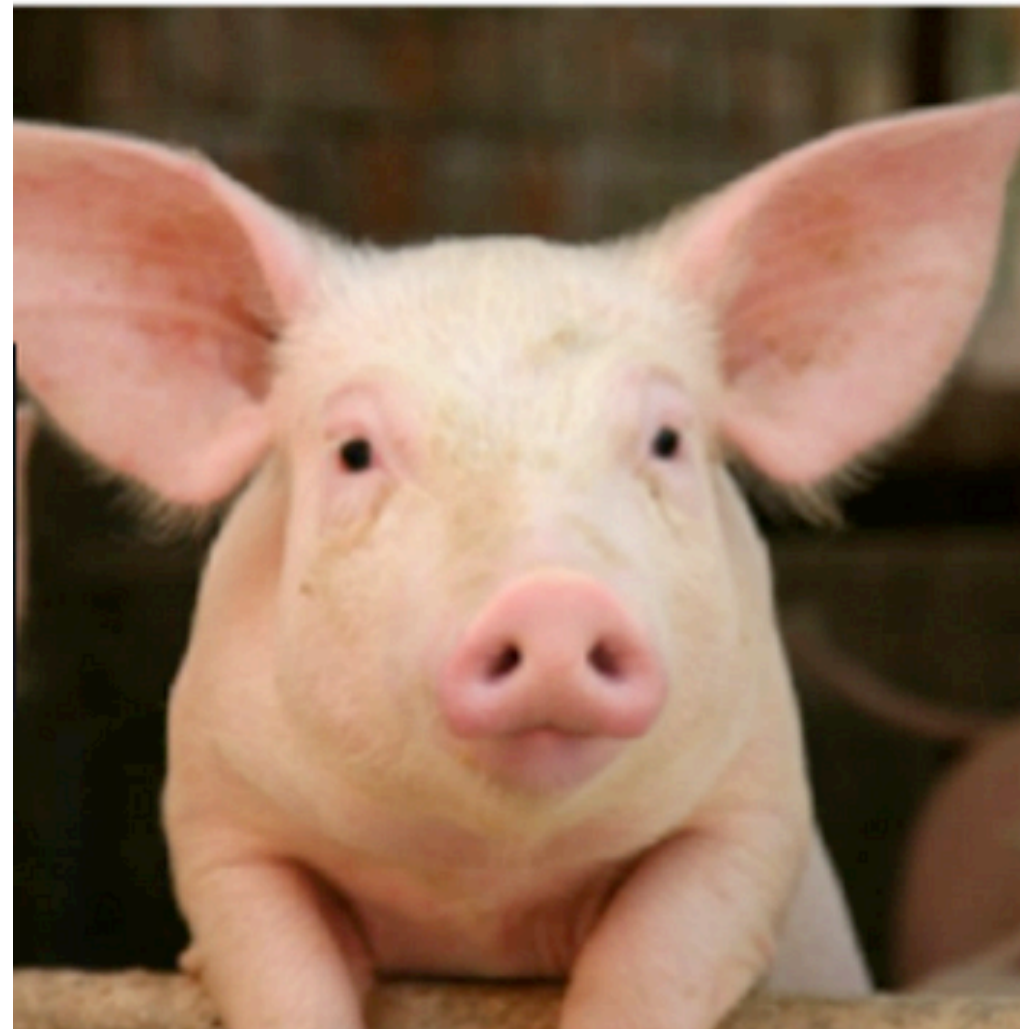
$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

$\mathbf{y}(\mathbf{x})$

$$\frac{\partial \mathscr{L}(\mathbf{y}(\mathbf{x}), \mathbf{y}^{\mathrm{gt}})}{\partial \mathbf{x}}$$

$$\mathbf{x} + 0.001 * \frac{\partial \mathscr{L}(\mathbf{y}(\mathbf{x}), \mathbf{y}^{\mathrm{gt}})}{\partial \mathbf{x}}$$

# Adversarial attacks [Arnab CVPR 2018]
## https://github.com/hmph/adversarial-attacks



"pig"

+ 0.005 x

=

"airliner"

$$\mathbf{x}$$

$$\frac{\partial \mathscr{L}(\mathbf{y}(\mathbf{x}), \mathbf{y}^{\mathrm{gt}})}{\partial \mathbf{x}}$$

$$\mathbf{x} + 0.005 * \frac{\partial \mathscr{L}(\mathbf{y}(\mathbf{x}), \mathbf{y}^{\mathrm{gt}})}{\partial \mathbf{x}}$$

Adversarial noise = direction in image domain that increases the loss the most

High frequency noise

The space is high-dimensional (10+ dimensions)

Access to network architecture + input image => physical attacks unsuccessful
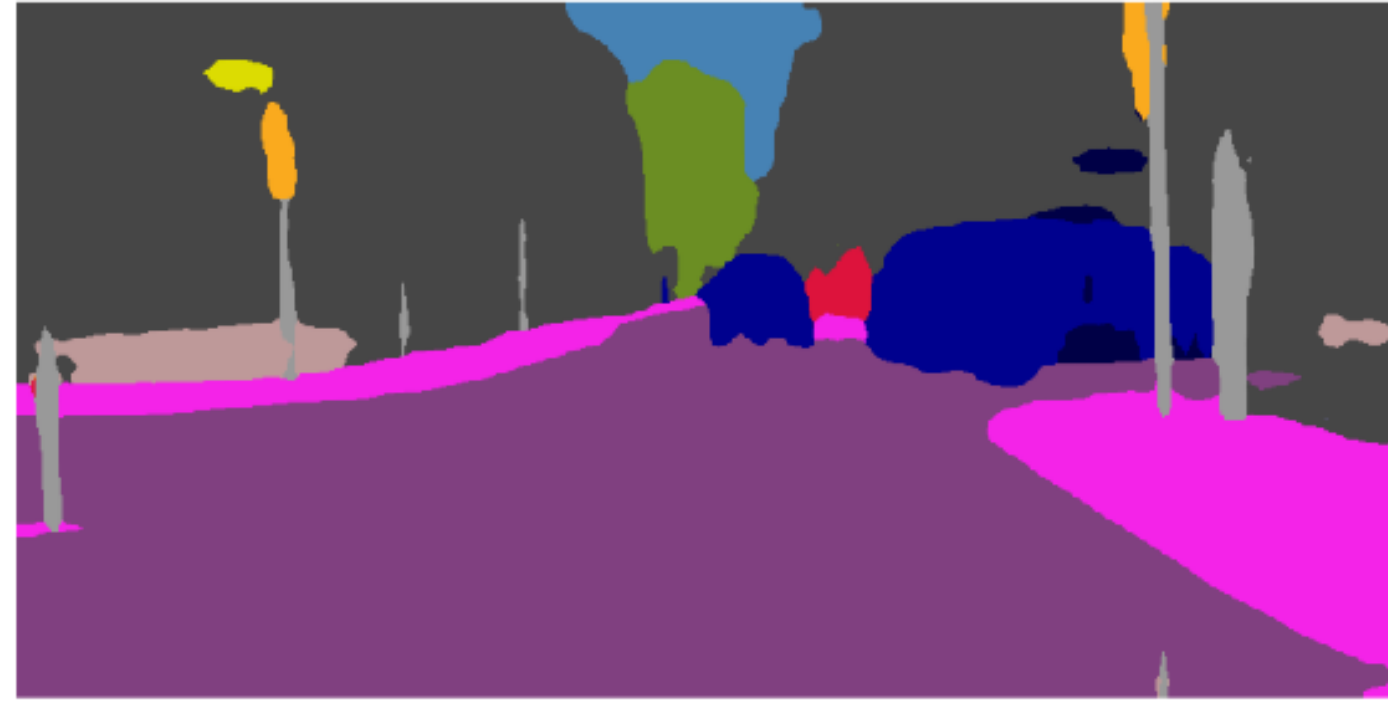
# Adversarial attacks [Arnab CVPR 2018]
https://github.com/hmph/adversarial-attacks

- Given trained network and an image, find the closest image on which the net fails



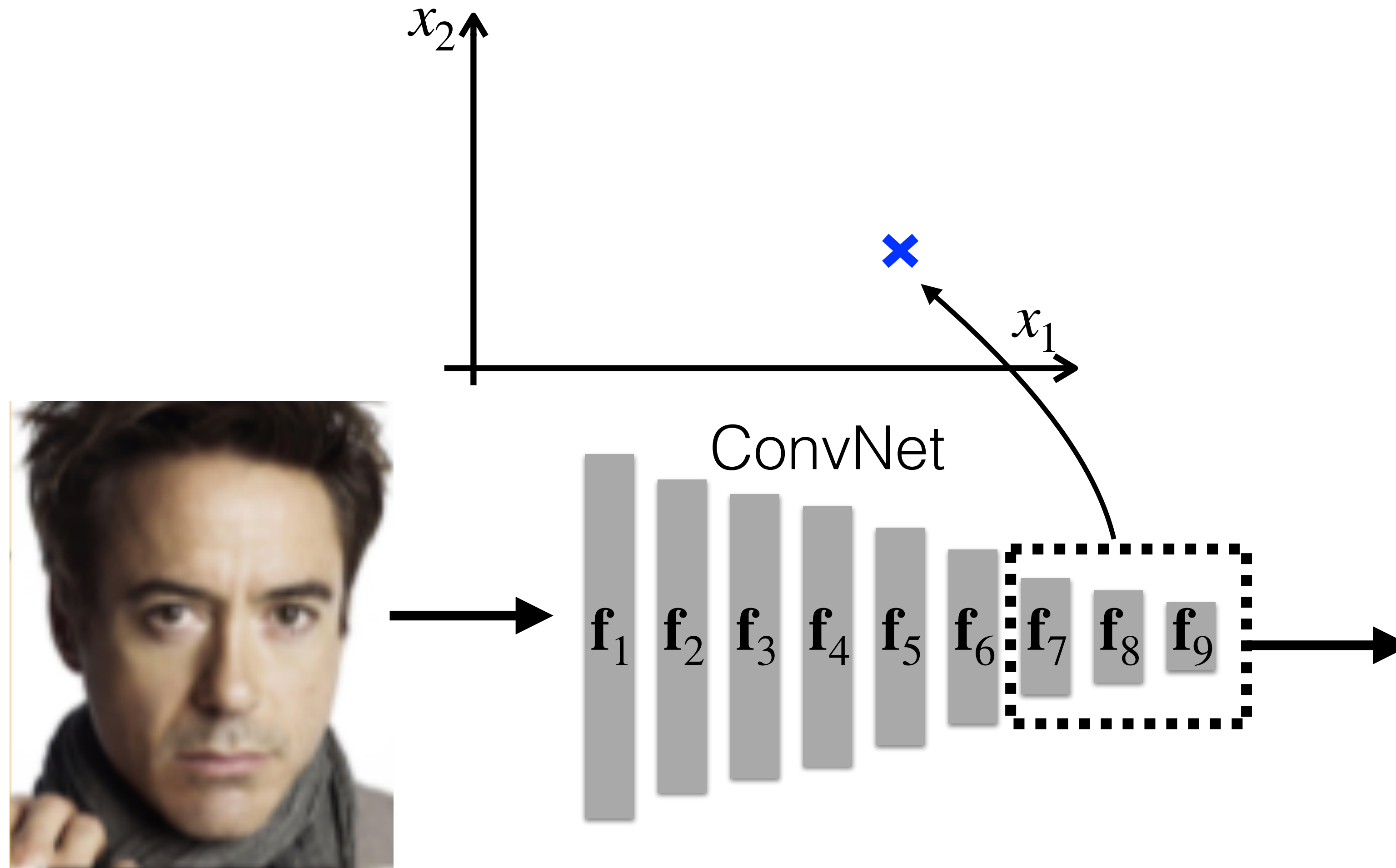| Input image | Original prediction | Adversarial prediction |



| Input image | Original prediction | Adversarial prediction |

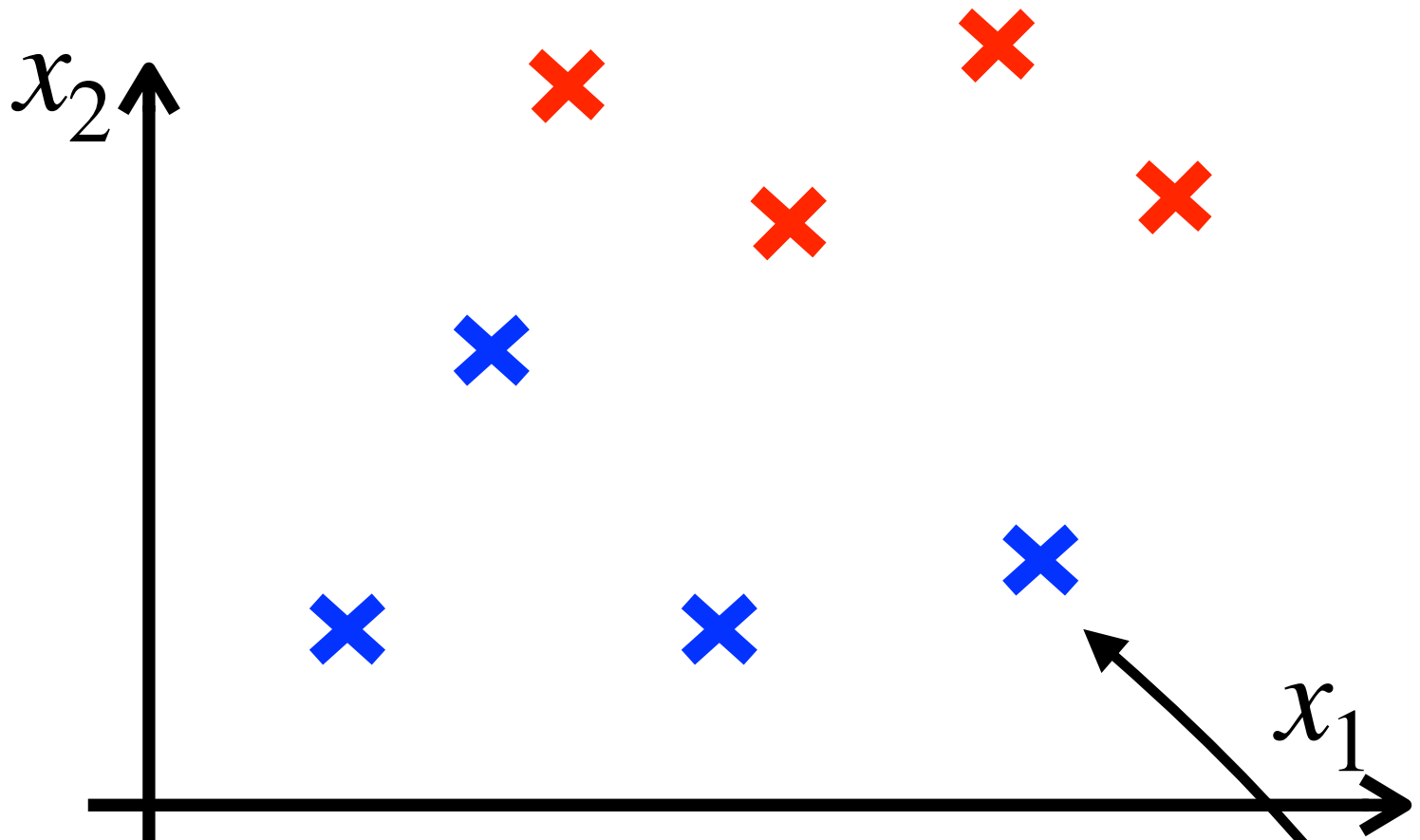# Deep Feature interpolations [Upchurch CVPR 2017]
https://arxiv.org/pdf/1611.05507.pdf



Deep ConvNet project images into a meaningful low-dimensional representation

# Deep Feature interpolations [Upchurch CVPR 2017]
## https://arxiv.org/pdf/1611.05507.pdf



faces **with** facial hair

faces **without** facial hair

$x_2$

$x_1$

ConvNet

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

Deep ConvNet project images into a meaningful low-dimensional representation
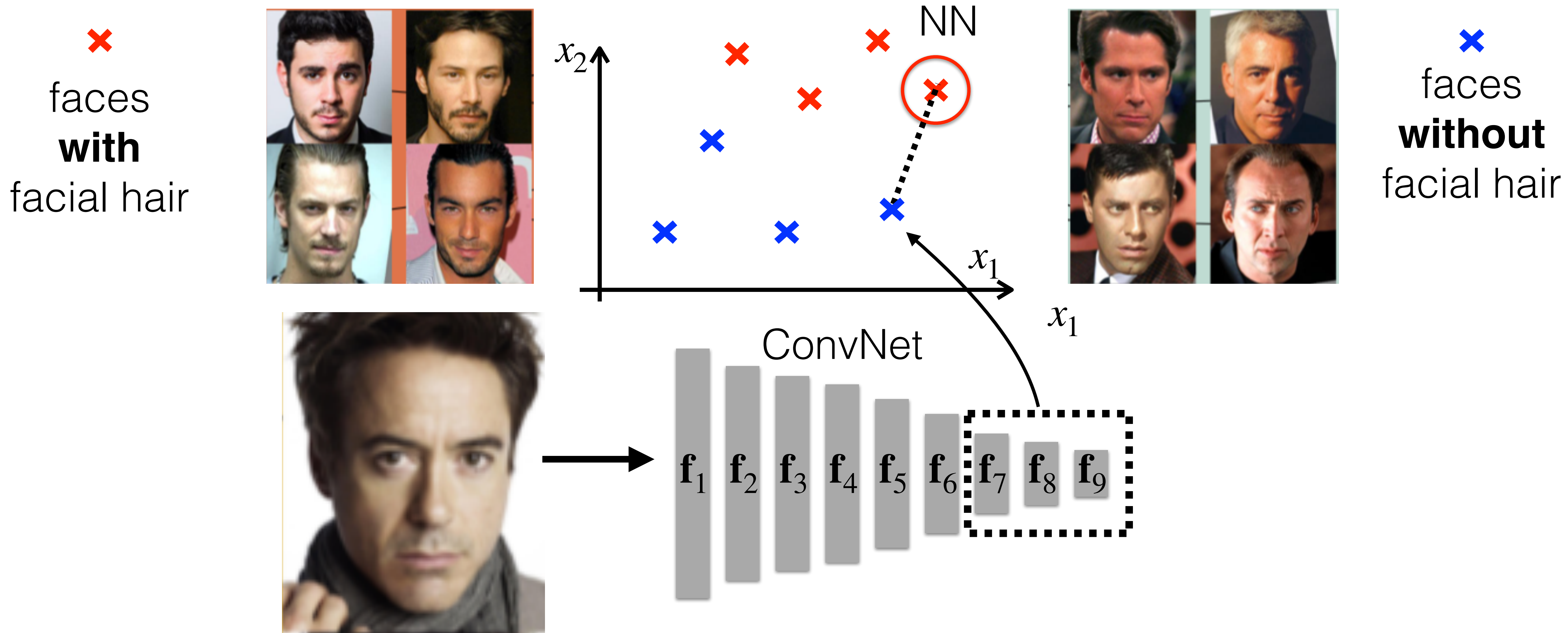
# Deep Feature interpolations [Upchurch CVPR 2017]
## https://arxiv.org/pdf/1611.05507.pdf



faces **with** facial hair

NN

$x_2$

$x_1$

$x_1$

faces **without** facial hair

ConvNet

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$
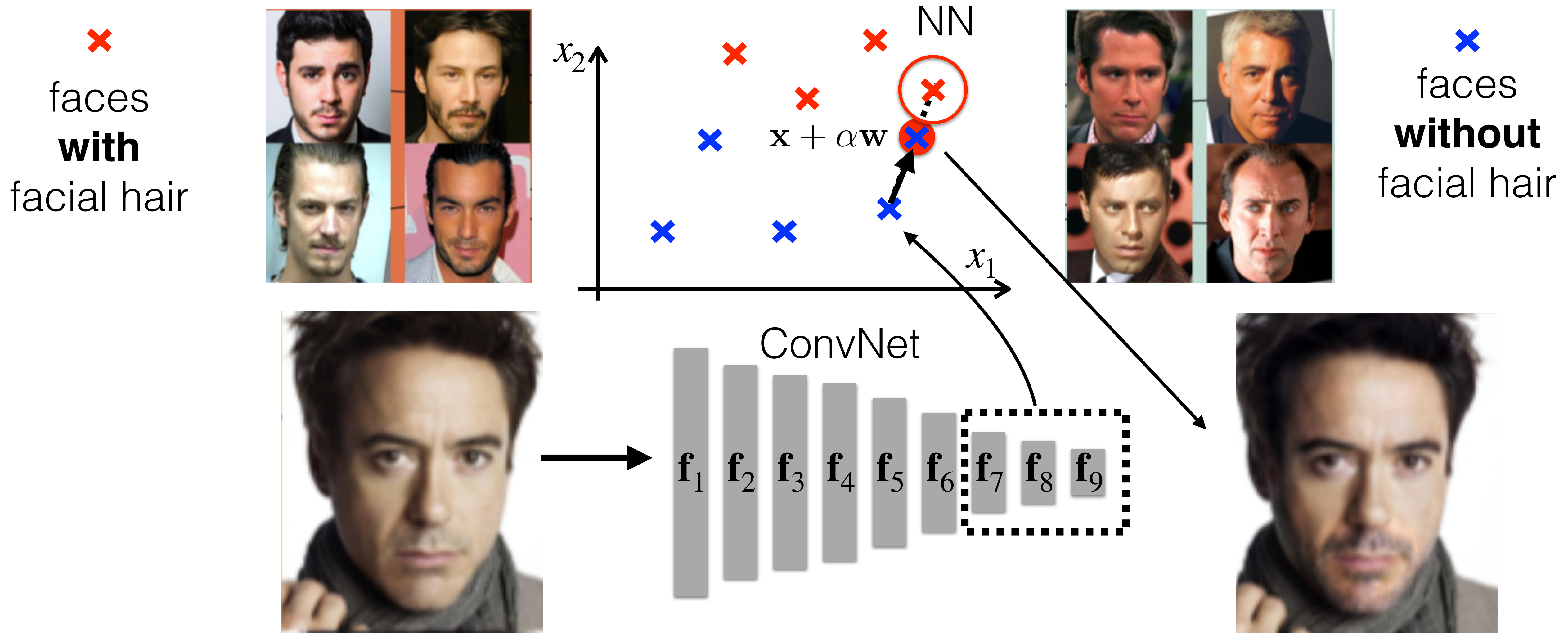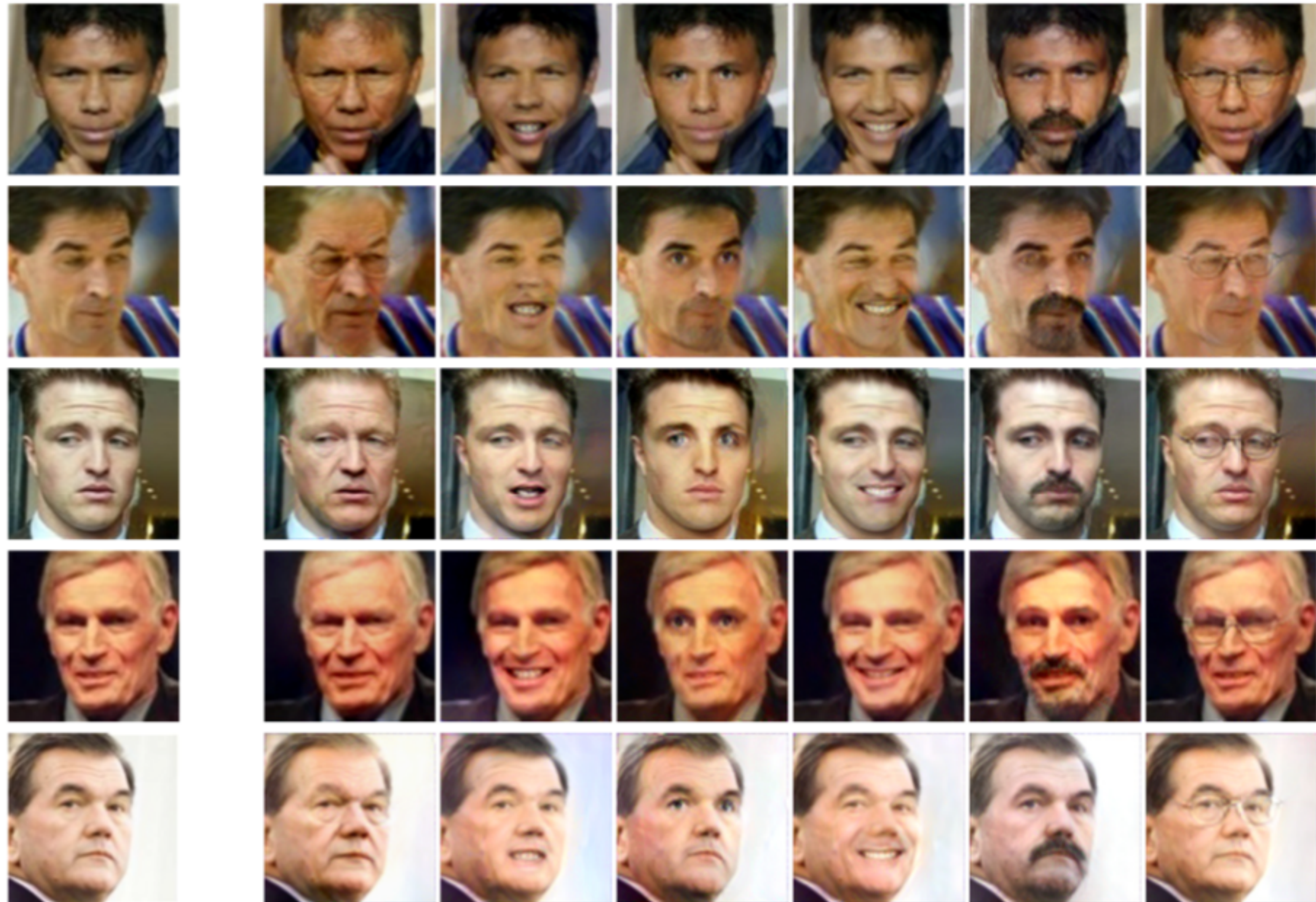
Deep ConvNet project images into a meaningful low-dimensional representation

# Deep Feature interpolations [Upchurch CVPR 2017]
https://arxiv.org/pdf/1611.05507.pdf



faces **with** facial hair

NN

$x_2$

$\mathbf{x} + \alpha\mathbf{w}$

$x_1$

faces **without** facial hair

ConvNet

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

Deep ConvNet project images into a meaningful low-dimensional representation

|  | Older | Mouth Open | Eyes Open | Smiling | Moustache | Glasses |

# Deep Feature interpolations [Upchurch CVPR 2017]
https://arxiv.org/pdf/1611.05507.pdf

# Shallow architecture
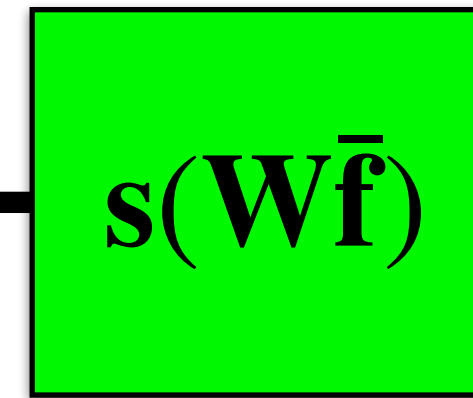
input/measurements     features     linear classifier     output



$\mathbf{x}$

$\mathbf{f}$

$s(\mathbf{W}\bar{\mathbf{f}})$

$\mathbf{y}$

0.2 dog
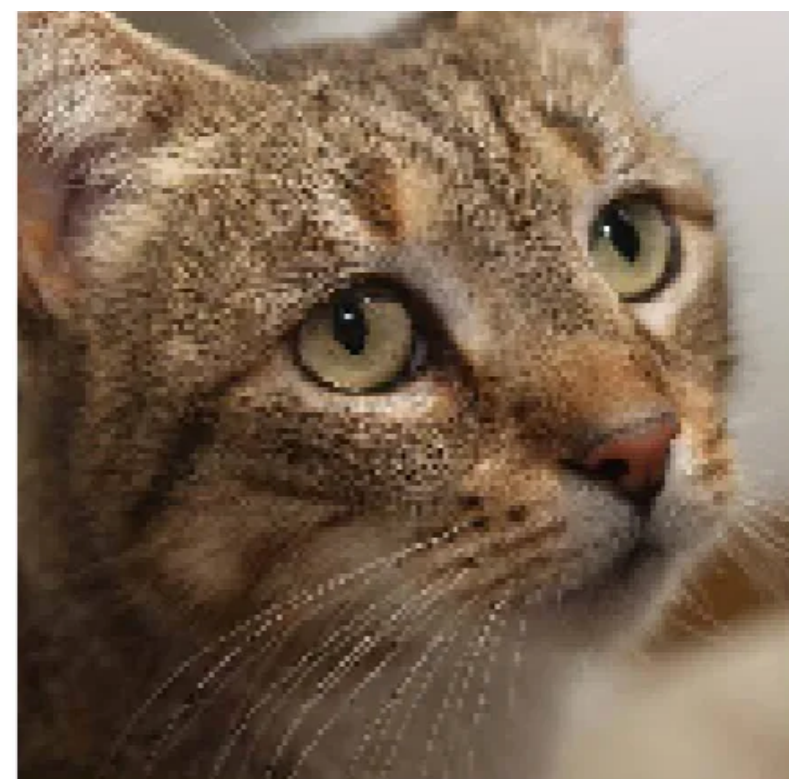0.1 ship
**0.5** cat
0.0 car
0.1 airplane

Manually estimated adhoc features
(pixels, color histograms, gradients, …)

🟥 manually estimated
🟩 trained

# Deep architecture

input/measurements     convolutional network     output



$\mathbf{x}$

$\mathbf{f}_1$ $\mathbf{f}_2$ $\mathbf{f}_3$ $\mathbf{f}_4$ $\mathbf{f}_5$ $\mathbf{f}_6$ $\mathbf{f}_7$ $\mathbf{f}_8$ $\mathbf{f}_9$

$\mathbf{y}$

0.2 dog
0.1 ship
**0.5** cat
0.0 car
0.1 airplane

Trained features

# Test T1 competencies

- ConvNet/Layer feed-forward pass
- ConvNet/Layer backpropagation
- Meaning of kernels, feature maps, stride, dilation, padding, …