# Computer Vision

## Lecture 02: Review on deep learning

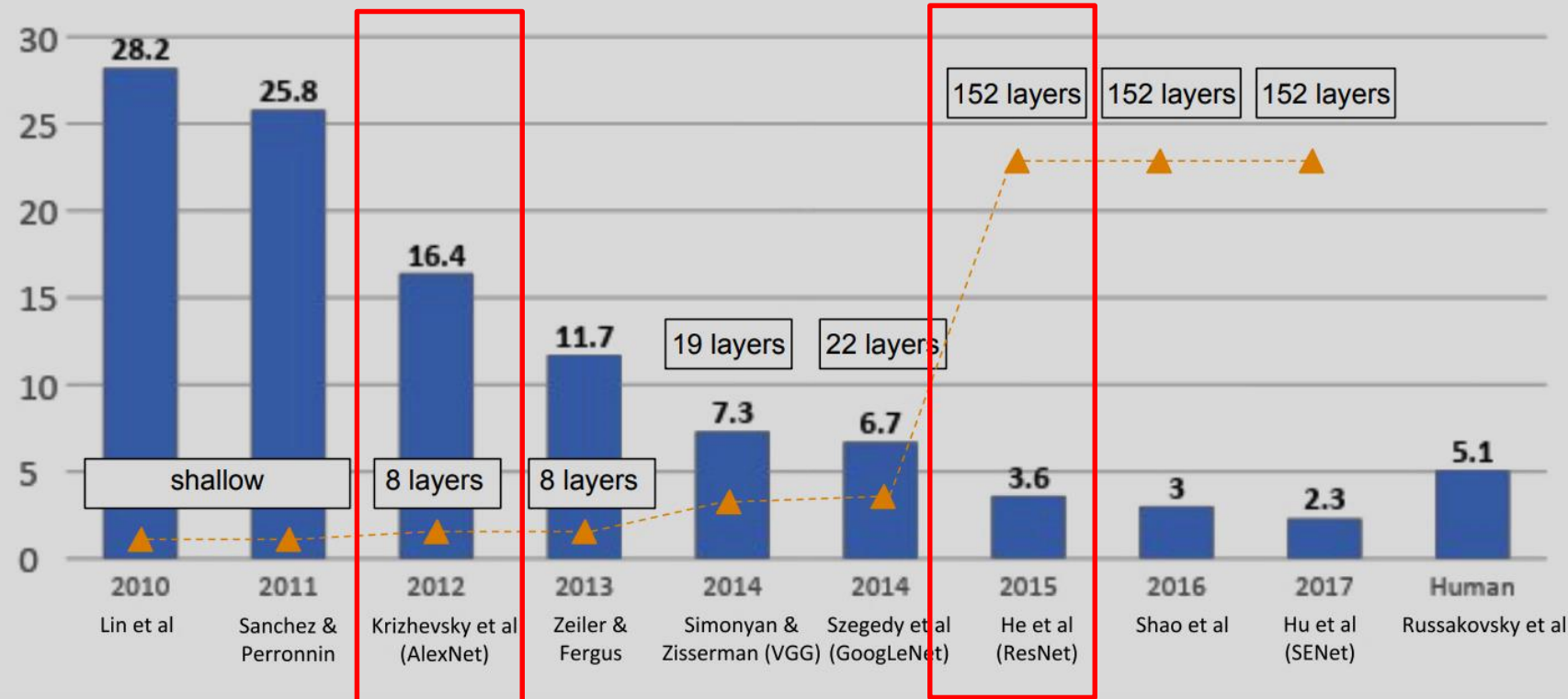# ImageNet challenge

## ImageNet Challenge



IM*A*GENET

- 1,000 object classes (categories).
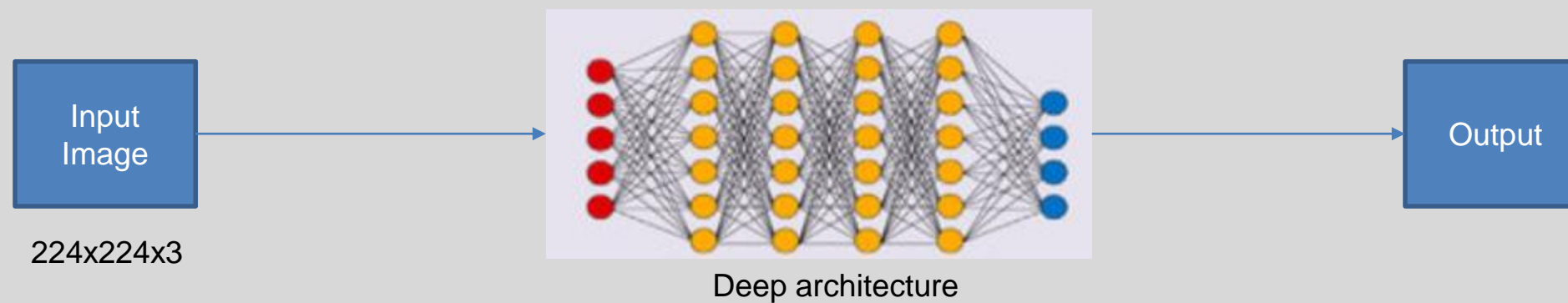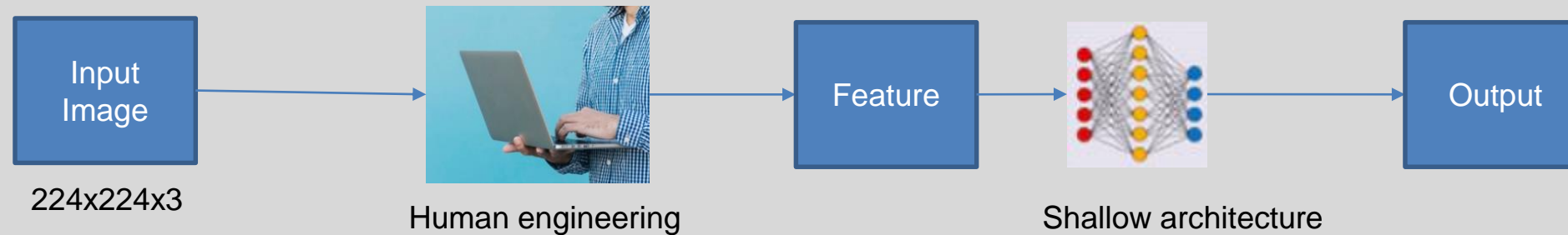- Images:
  - 1.2 M train
  - 100k test.

# Deep learning



Year 2012: Deep learning achieved the best performance on image classification task.
Year 2015: Surpasses the human performance.
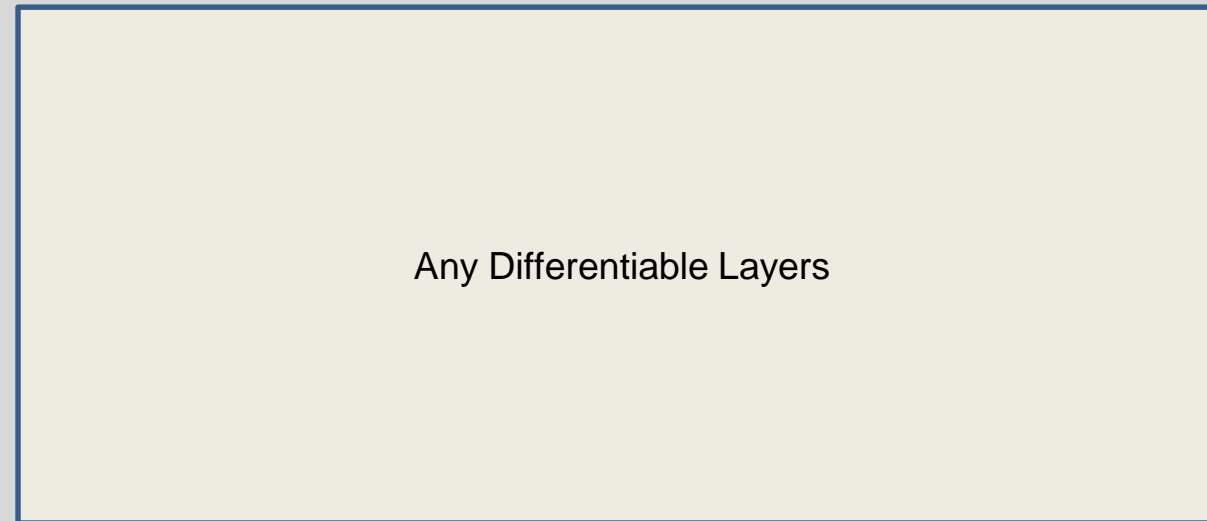
# Deep vs. Machine learning

# Deep vs. Machine learning

- ## Extract optimal representation via End-to-End learning
  - For the machine learning, we need to design dedicated representation by ourselves for each task.
  - Deep learning learns intermediate representation automatically for different tasks.

- ## Non-linearity
  - Deep learning delivers the capability to achieve the non-linear mappings.
  - Most computer vision applications involve data which requires non-linear mappings.
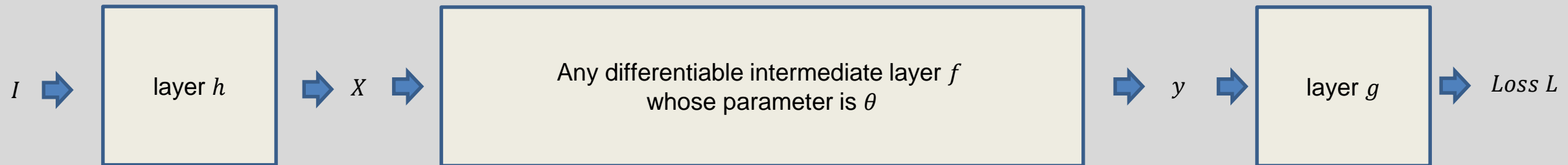
# CNNs

RGB image → Any Differentiable Layers → cat

RGB image

Semantic labels

# Differentiable layers

$I$ → layer $h$ → $X$ → Any differentiable intermediate layer $f$ whose parameter is $\theta$ → $y$ → layer $g$ → $Loss\ L$

We need to implement three things for an intermediate layer $f$ :

forward rule: $\qquad y = f(\mathrm{h}(\mathrm{I}); \theta)$  for  $g\big(f(\mathrm{h}(\mathrm{I}); \theta)\big) = \mathrm{L}$

backward rule: $\qquad \dfrac{dy}{dX}$  for  $\dfrac{dL}{d\mathrm{I}} = \dfrac{d\mathrm{X}}{d\mathrm{I}} \times \dfrac{dy}{dX} \times \dfrac{dL}{dy}$

parameter update rule: $\dfrac{dy}{d\theta}$  for  $\theta^{new} = \theta - \varepsilon \dfrac{dy}{d\theta} \times \dfrac{dL}{dy}$

Via chain rule, the entire architecture becomes differentiable, if each layer becomes differentiable.
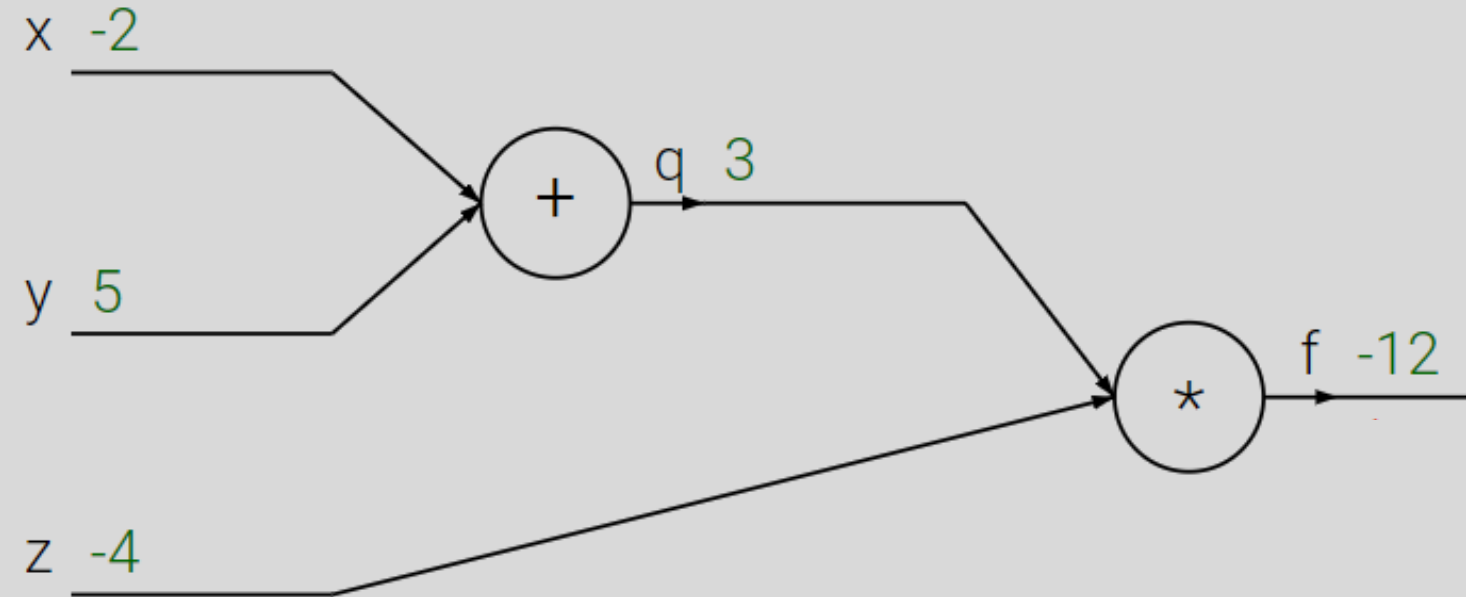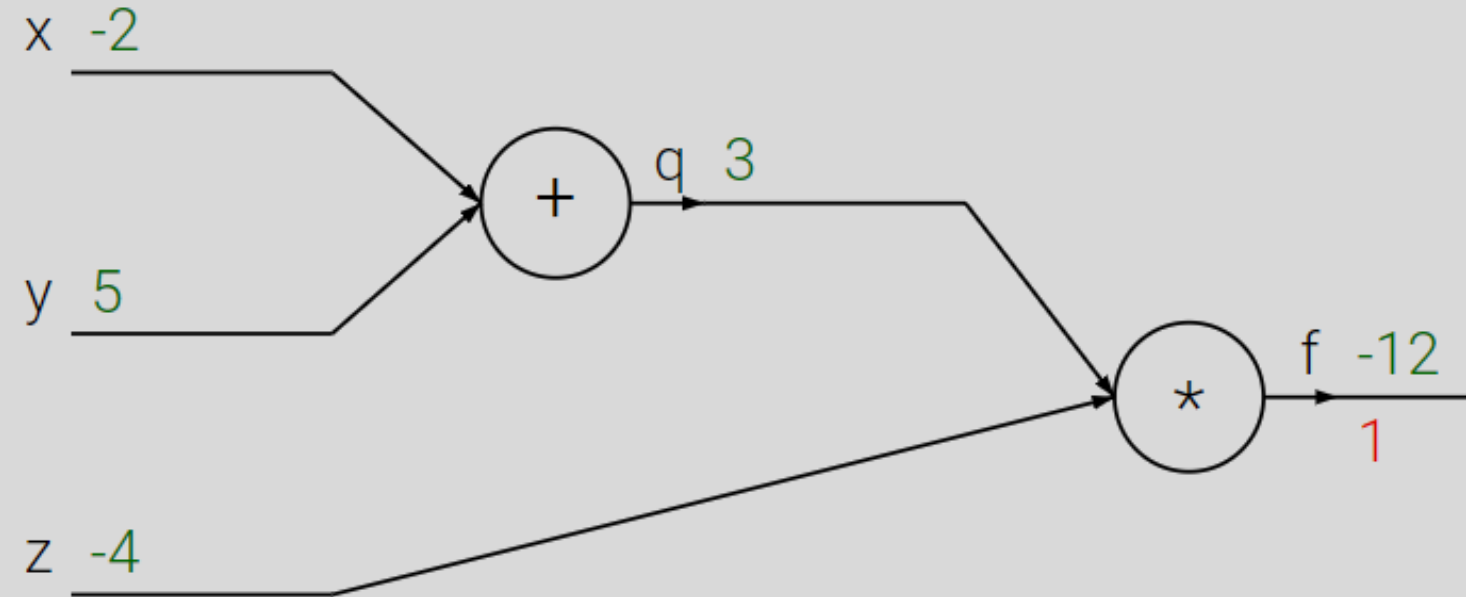
# Gradient descent

L(θ) = g(f(h(I); θ))

$$\theta^{\text{new}} = \theta - \epsilon \frac{\partial}{\partial \theta} L(\theta)$$

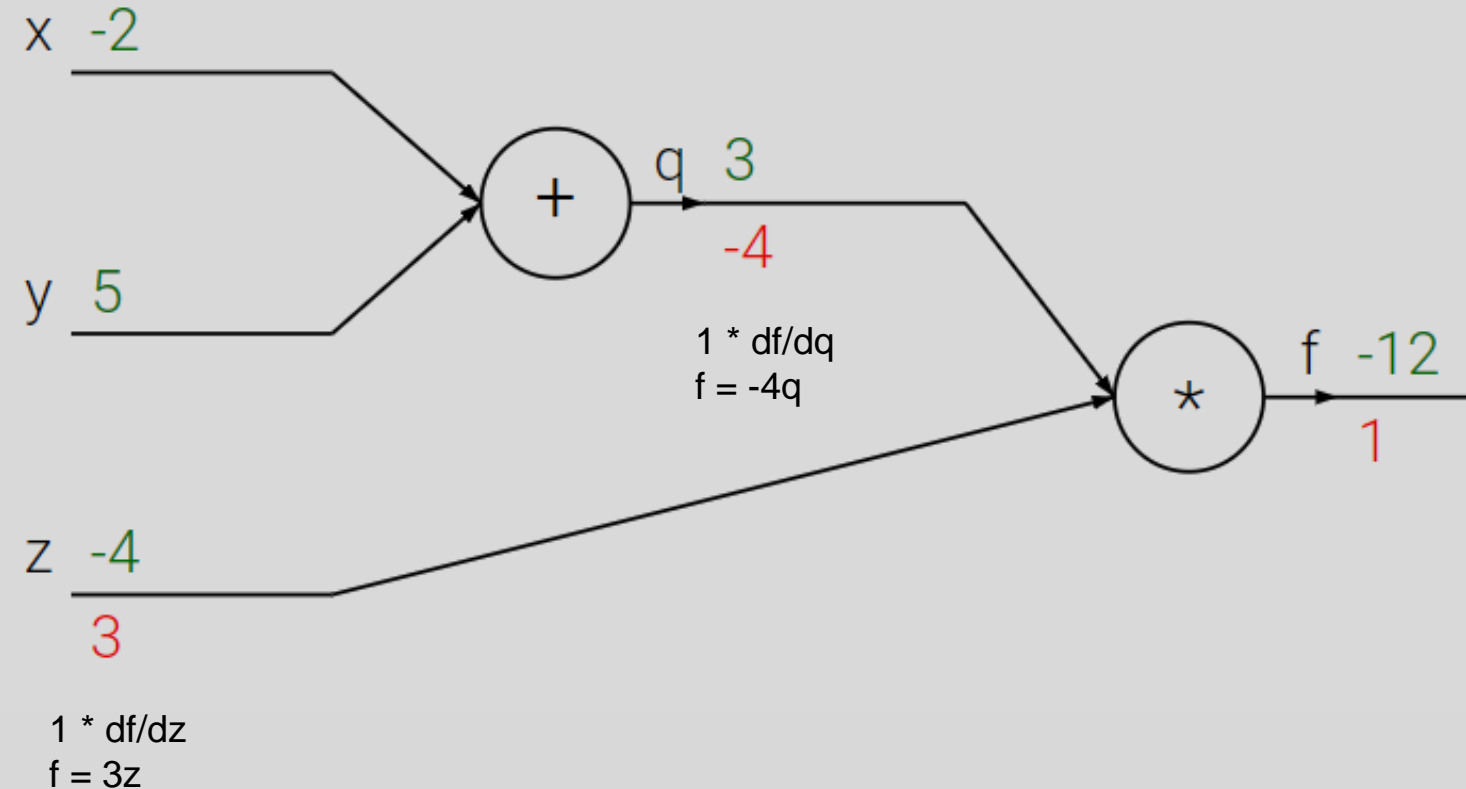$\epsilon$ : Learning rate (small value e.g. 0.1)
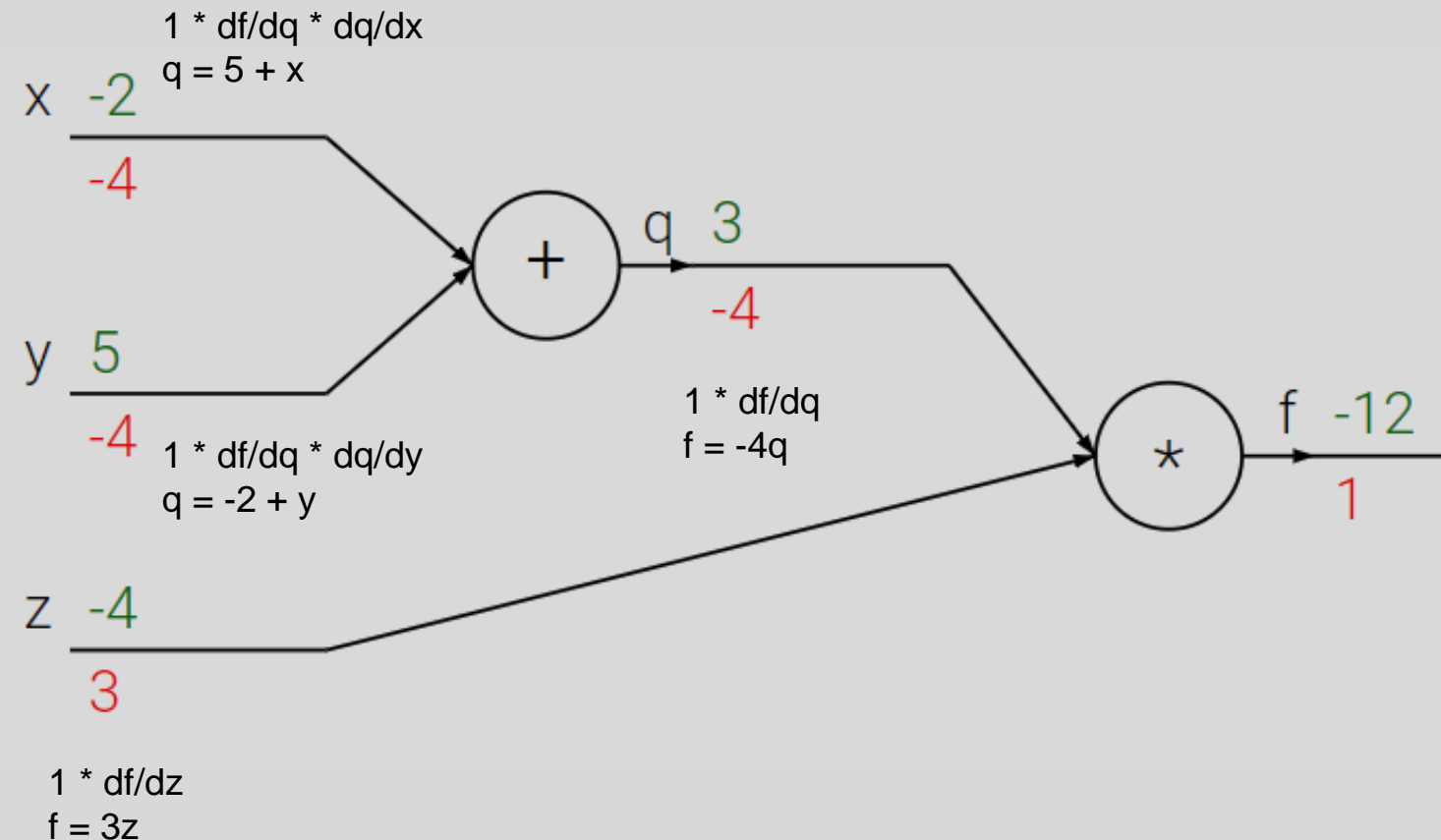
# Calculating gradients in deep learning

# Calculating gradients in deep learning

x  -2

y  5

q  3

z  -4

+

*

f  -12

1

# Calculating gradients in deep learning

# Calculating gradients in deep learning
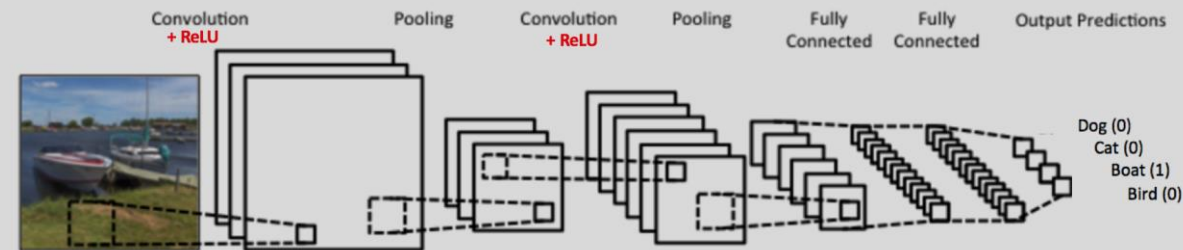
# CNNs



RGB image

Any Differentiable Layers

cat

Semantic labels

# CNNs



RGB image

Semantic labels

cat

# 2D Convolution



Image (5x5)

Filter kernel (3x3)

# 2D Convolution

$$1*0+1*0+0*-1$$
$$+ 2*1+0*-1+0*0$$
$$+ 0*0+0*1+2*-1$$

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 0 | 1 | 2 |
| 2 | 0 | 0 | 1 | 0 |
| 0 | 0 | 2 | 2 | 1 |
| 0 | 0 | 0 | 2 | 0 |
| 0 | 0 | 2 | 1 | 1 |

| | | |
|---|---|---|
| 0 | 0 | -1 |
| 1 | -1 | 0 |
| 0 | 1 | -1 |

0

Image (5x5)                    Filter kernel (3x3)

# 2D Convolution



Image (5x5)                    Filter kernel (3x3)

$$1*0+0*0+1*-1$$
$$+ 0*1+0*-1+1*0$$
$$+ 0*0+2*1+2*-1$$

# 2D Convolution

Image (5x5)

Filter kernel (3x3)

$0*0+1*0+2*-1$
$+ 0*1+1*-1+0*0$
$+ 2*0+2*1+1*-1$

# 2D Convolution

$2*0+0*0+0*-1$
$+ 0*1+0*-1+2*0$
$+ 0*0+0*1+0*-1$

| 1 | 1 | 0 | 1 | 2 |
|---|---|---|---|---|
| 2 | 0 | 0 | 1 | 0 |
| 0 | 0 | 2 | 2 | 1 |
| 0 | 0 | 0 | 2 | 0 |
| 0 | 0 | 2 | 1 | 1 |

Image (5x5)

| 0 | 0 | -1 |
|---|---|----|
| 1 | -1 | 0 |
| 0 | 1 | -1 |

Filter kernel (3x3)

| 0 | -1 | -2 |
|---|----|----|
| 0 | | |

19

# 2D Convolution



$2*0+2*0+1*-1$
$+ 0*1+2*-1+2*0$
$+ 2*0+1*1+1*-1$

Image (5x5)                    Filter kernel (3x3)                    Output feature (3x3)

# Ex. Image blurring

| 1 | 1 | 1 |
|---|---|---|
| 1 | 1 | 1 |
| 1 | 1 | 1 |

# 2D Convolution w/ zero padding

$0*0+0*0+0*-1$
$+ 0*1+1*-1+1*0$
$+ 0*0+2*1+0*-1$

Image (5x5)

Filter kernel (3x3)

# 2D Convolution w/ zero padding

$$0*0+0*0+0*-1$$
$$+ 1*1+1*-1+0*0$$
$$+ 2*0+0*1+0*-1$$

| | | | | |
|---|---|---|---|---|
| | 1 | 1 | 0 | 1 | 2 |
| | 2 | 0 | 0 | 1 | 0 |
| | 0 | 0 | 2 | 2 | 1 |
| | 0 | 0 | 0 | 2 | 0 |
| | 0 | 0 | 2 | 1 | 1 |

Image (5x5)

Filter kernel (3x3)

| 0 | 0 | -1 |
|---|---|---|
| 1 | -1 | 0 |
| 0 | 1 | -1 |

| 1 | 0 |
|---|---|

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

# 2D Convolution w/ zero padding



2*0+0*0+0*-1
+ 1*1+1*-1+0*0
+ 0*0+0*1+0*-1

Image (5x5)                          Filter kernel (3x3)                  Output feature (5x5)
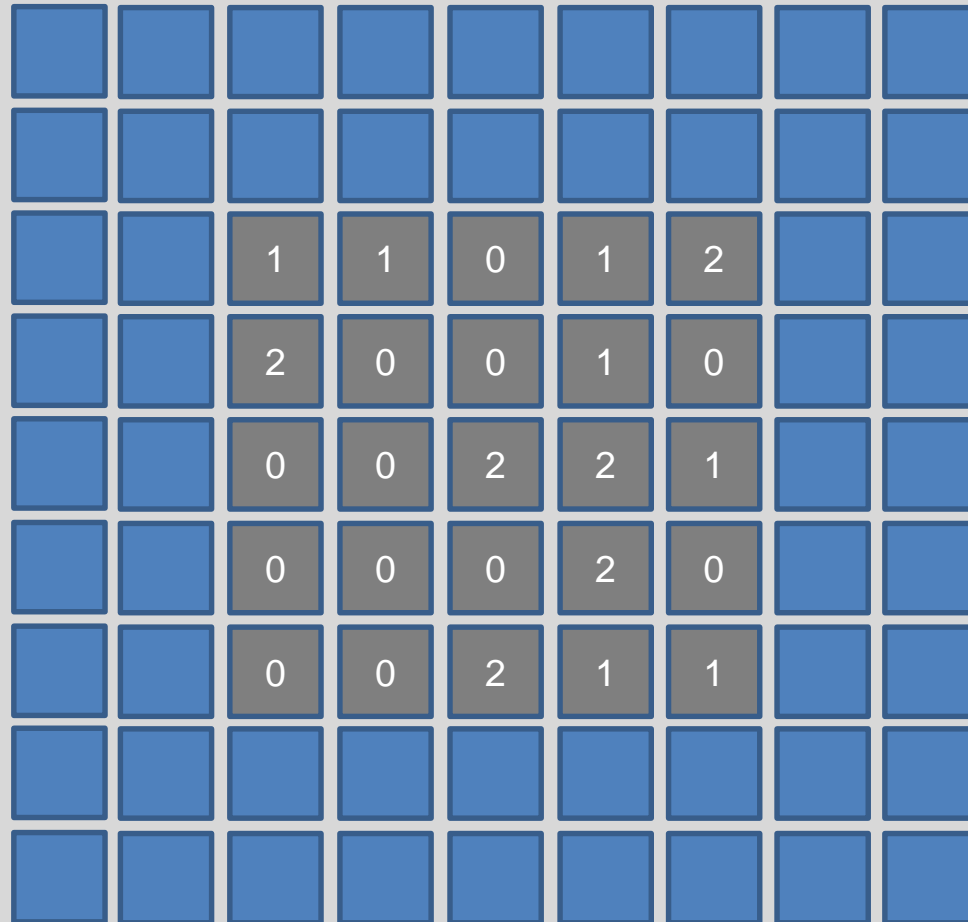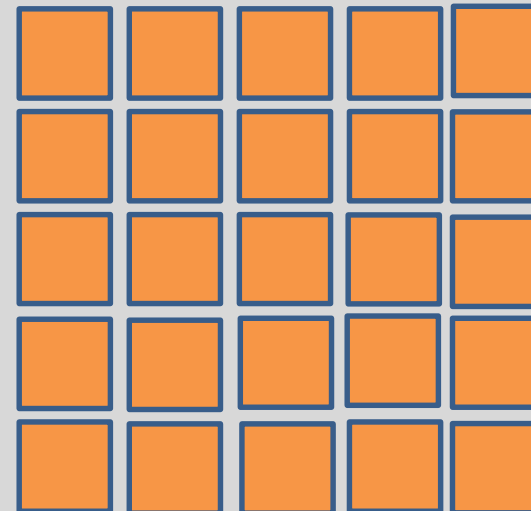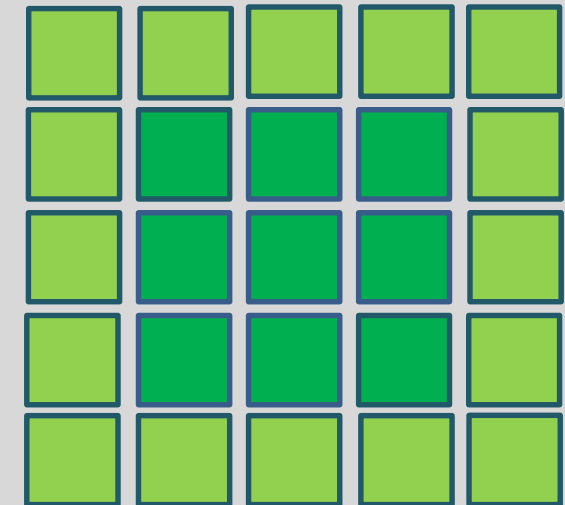
24

# 2D Convolution w/ zero padding

→ (K-1) / 2 padding is required to obtain the original size.



Image (5x5)
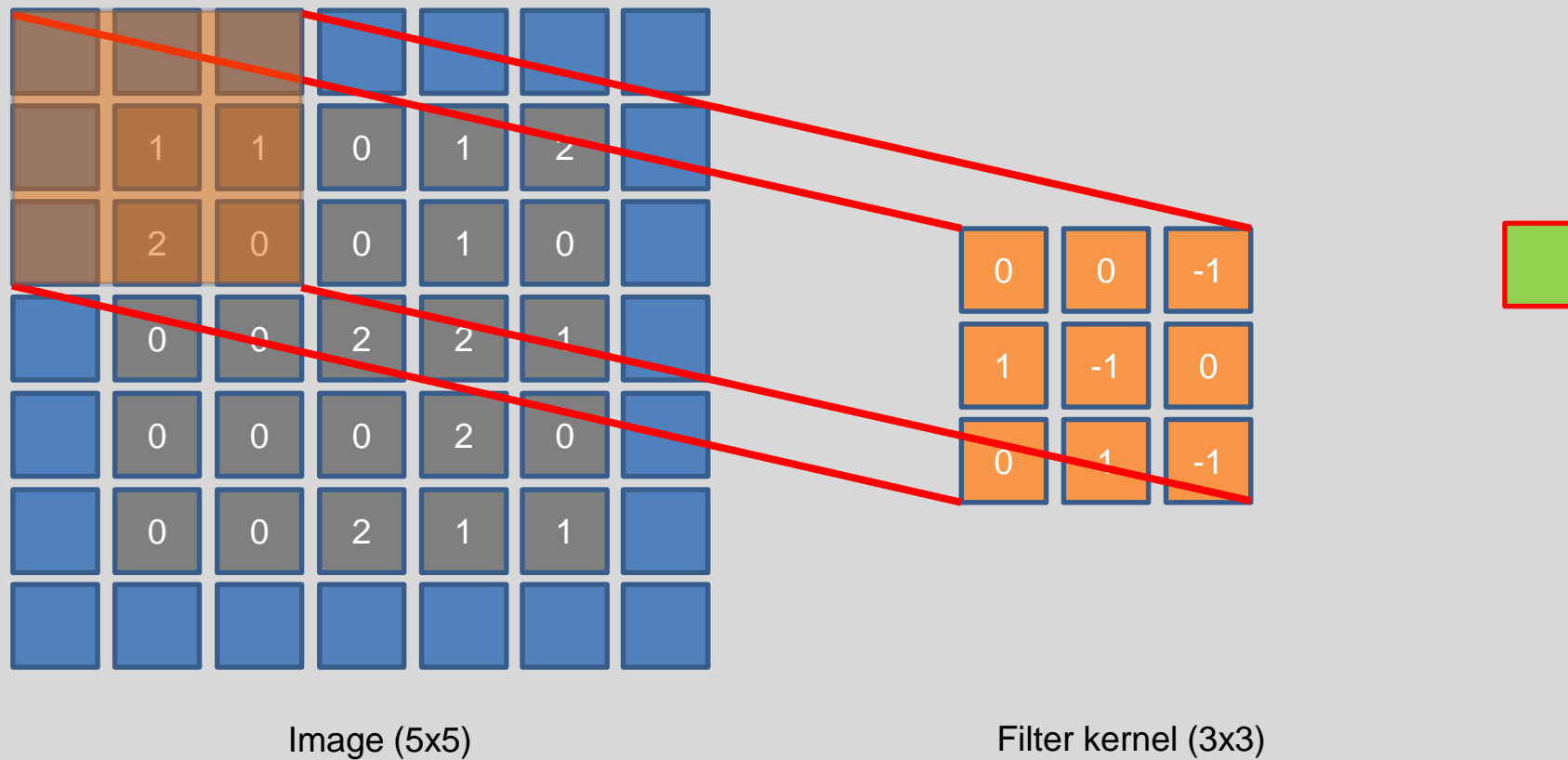
Filter kernel (5x5)

Output feature (5x5)

# 2D Convolution w/ stride

Image (5x5)

Filter kernel (3x3)

Stride=2

# 2D Convolution w/ stride



Image (5x5)                              Filter kernel (3x3)

Stride=2

27

# 2D Convolution w/ stride

Image (5x5)

Filter kernel (3x3)

Stride=2

# 2D Convolution w/ stride



Image (5x5)

Filter kernel (3x3)

Stride=2

# 2D Convolution w/ stride



Image (5x5)

Filter kernel (3x3)

Stride=2

# 2D Convolution w/ stride



Image (5x5)

Filter kernel (3x3)

Stride=2

31

# Convolutional layer

# Convolutional layer

| | | | |
|---|---|---|---|
| I1 5x5 | f11 3x3 | o11 5x5 | O1 5x5 |
| I2 5x5 | f21 3x3 | o12 5x5 | |
| I3 5x5 | f31 3x3 | o13 5x5 | |

# Convolutional layer

| | | |
|---|---|---|
| I1 5x5 | f12 3x3 | o21 5x5 |
| I2 5x5 | f22 3x3 | o22 5x5 |
| I3 5x5 | f32 3x3 | o23 5x5 |

# Convolutional layer



I1
5x5

f12
3x3

o21
5x5

I2
5x5

f22
3x3

o22
5x5

O2
5x5

I3
5x5

f32
3x3

o23
5x5

# Convolutional layer

| I1 5x5 | f13 3x3 | o31 5x5 |
| I2 5x5 | f23 3x3 | o32 5x5 |
| I3 5x5 | f33 3x3 | o33 5x5 |

# Convolutional layer

I1
5x5

f13
3x3

o31
5x5

I2
5x5

f23
3x3

o32
5x5

I3
5x5

f33
3x3

o33
5x5

O3
5x5

# Convolutional layer

| I1<br>5x5 | f1N<br>3x3 | oN1<br>5x5 |

| I2<br>5x5 | f2N<br>3x3 | oN2<br>5x5 |

| I3<br>5x5 | f3N<br>3x3 | oN3<br>5x5 |

# Convolutional layer

# Convolutional layer

# Pooling

# Pooling

max pooling



torch.nn.MaxPool2d

average pooling

torch.nn.AvgPool2d

# Activation layer

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

$$f(x) = tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$

$$f(x) = max(0, x)$$

torch.nn.Sigmoid()

torch.nn.Tanh()

torch.nn.ReLU()

43

# Overall CNN architecture



**Combination of differentiable layers → Differentiable architecture!**

# Overall CNN architecture

Convolutional layers.

Fully-connected layers.



**Combination of differentiable layers → Differentiable architecture!**

# Semantic Segmentation

## Semantic Segmentation

Label each pixel in the image with a category label

Don't differentiate instances, only care about pixels

# Semantic Segmentation

# Semantic Segmentation



Input:
3 x H x W

High-res:
$D_1$ x H/2 x W/2

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

High-res:
$D_1$ x H/2 x W/2

Predictions:
H x W

# Transposed Convolution



Convolution operation                    Transposed convolution operation

# Transposed Convolution

Transposed Convolution with 0 padding, stride 1, 2x2 kernel:

Output_size = (input_size-1)*stride – 2*padding + kernel_size + output_padding

2x2 kernel,
stride=1

Intermediate grid

Output

Input

5x5

Stride=1

6x6

# Transposed Convolution

Transposed Convolution with 0 padding, stride 2, 2x2 kernel:

Output_size = (input_size-1)*stride – 2*padding + kernel_size + output_padding

# Transposed Convolution

Transposed Convolution with 1 padding, stride 2, 2x2 kernel:

Output_size = (input_size-1)*stride – 2*padding + kernel_size + output_padding

# Conv. operation



4x4 Input                    3x3 kernel                    2x2 Output

# Conv. operation

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

4x4 Input

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

4x16 Conv. Kernel matrix

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

16x1 Input

# Conv. operation

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

4x16 Conv. Kernel matrix

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

16x1 Input

| 122 |
|-----|
| 148 |
| 126 |
| 134 |

4x1 Output

# Transposed Conv. operation

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

16x4
Conv. Kernel Matrix

| |
|---|
| 2 |
| 1 |
| 4 |
| 4 |

4x1 Input

| |
|---|
| 2 |
| 9 |
| 6 |
| 1 |
| 6 |
| 29 |
| 30 |
| 7 |
| 10 |
| 29 |
| 33 |
| 13 |
| 12 |
| 24 |
| 16 |
| 4 |

16x1 Output

| | | | |
|---|---|---|---|
| 2 | 9 | 6 | 1 |
| 6 | 29 | 30 | 7 |
| 10 | 29 | 33 | 13 |
| 12 | 24 | 16 | 4 |

4x4 Output

# Achievable by differentiable layers



Input:
3 x H x W

Med-res:
$D_2$ x H/4 x W/4

Med-res:
$D_2$ x H/4 x W/4

Low-res:
$D_3$ x H/4 x W/4

High-res:
$D_1$ x H/2 x W/2

High-res:
$D_1$ x H/2 x W/2

Predictions:
H x W

16x4
Conv. Kernel Matrix

# Achievable by differentiable layers



(a) Input context　　(b) Human artist　　(c) Context Encoder ($L2$ loss)　　(d) Context Encoder ($L2$ + Adversarial loss)

16x4
Conv. Kernel Matrix

# Achievable by differentiable layers



16x4
Conv. Kernel Matrix

# Next class…

– Review on the PyTorch.