# Reminder: Scalars, Vectors, Matrices & Tensors
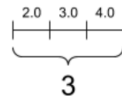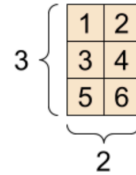
A scalar, shape: [ ]

4
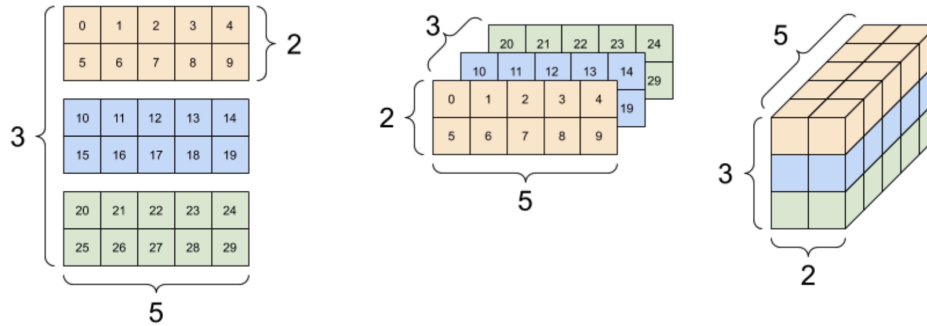
A vector, shape: [3]
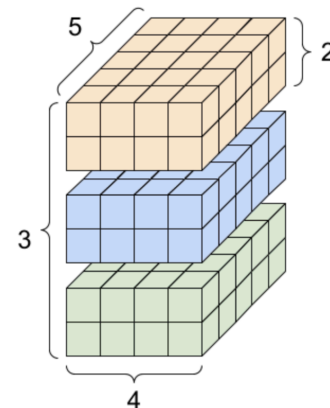
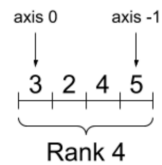2.0   3.0   4.0

3

A matrix, shape: [3, 2]

| 1 | 2 |
| 3 | 4 |
| 5 | 6 |

3

2

A 3-axis tensor, shape: [3, 2, 5]

| 0 | 1 | 2 | 3 | 4 |
| 5 | 6 | 7 | 8 | 9 |

2

| 10 | 11 | 12 | 13 | 14 |
| 15 | 16 | 17 | 18 | 19 |

| 20 | 21 | 22 | 23 | 24 |
| 25 | 26 | 27 | 28 | 29 |

3

5

3

| 20 | 21 | 22 | 23 | 24 |
| 10 | 11 | 12 | 13 | 14 | 29 |
| 0 | 1 | 2 | 3 | 4 | 19 |
| 5 | 6 | 7 | 8 | 9 |

2

5

5

3

2

A rank-4 tensor, shape: [3, 2, 4, 5]

axis 0          axis -1

3   2   4   5

Rank 4

5

2

3

4

# Representing Images



H

W

| 255 | 255 | 127 |
|-----|-----|-----|
| 255 | 127 | 0 |
| 192 | 0 | 62 |

— Images composed of pixel features

— Typically 8-bit values 0 — 255 (int)

0 → Black    ≈127    255 → White

Single obs x is a matrix: $H \times W$

Batch X̶ is a 3-D Tensor: $N \times H \times W$

# Representing Images



H

W

For Neural networks want to convert to float and rescale

255 255 127
255 127 0
192 0 62

## 0 to 1
good

| 1. | 1. | 0.5 |
| 1. | 0.5 | 0. |
| 0.7 | 0. | 0.2 |

## -1 to 1
better

| 1. | 1. | 0. |
| 1. | 0. | -1. |
| 1.2 | -1. | -.5 |

## Norm
best!

[But data set specific :-]

| 3.1 | 3.1 | 0.1 |
| 3.1 | 0.1 | -3. |
| 1.4 | -3. | -1.1 |

# RGB Images

Represent color pixels w/ 3-Channels
[3-Values per pixel]

Purple → 240 High red
→ 10 Low green
→ 200 High blue

White → 255 Max red
→ 255 Max green
→ 255 Max blue

black → 0 No red
→ 0 No green
→ 0 No blue

Single obs. $x$ is 3-D tensor: $H \times W \times 3$

Batch $X$ is 4-D tensor: $N \times H \times W \times 3$

# RGB Images – PyTorch

## Normal

Single obs. x is 3-D tensor: $H \times W \times 3$

Batch X is 4-D tensor: $N \times H \times W \times 3$

## PyTorch

Single obs. x is 3-D tensor: $3 \times H \times W$

Batch X is 4-D tensor: $N \times 3 \times H \times W$

$$X.\text{permute}(0, 3, 1, 2)$$

# Representing Images (So far)

For standard neural network Reshape into vector



Single Obs.

$$x : HW$$

Pixels

Batch

$$\underline{X} : N \times HW$$

Pixels

Image 1
Image 2
Image 3

N

Pixels

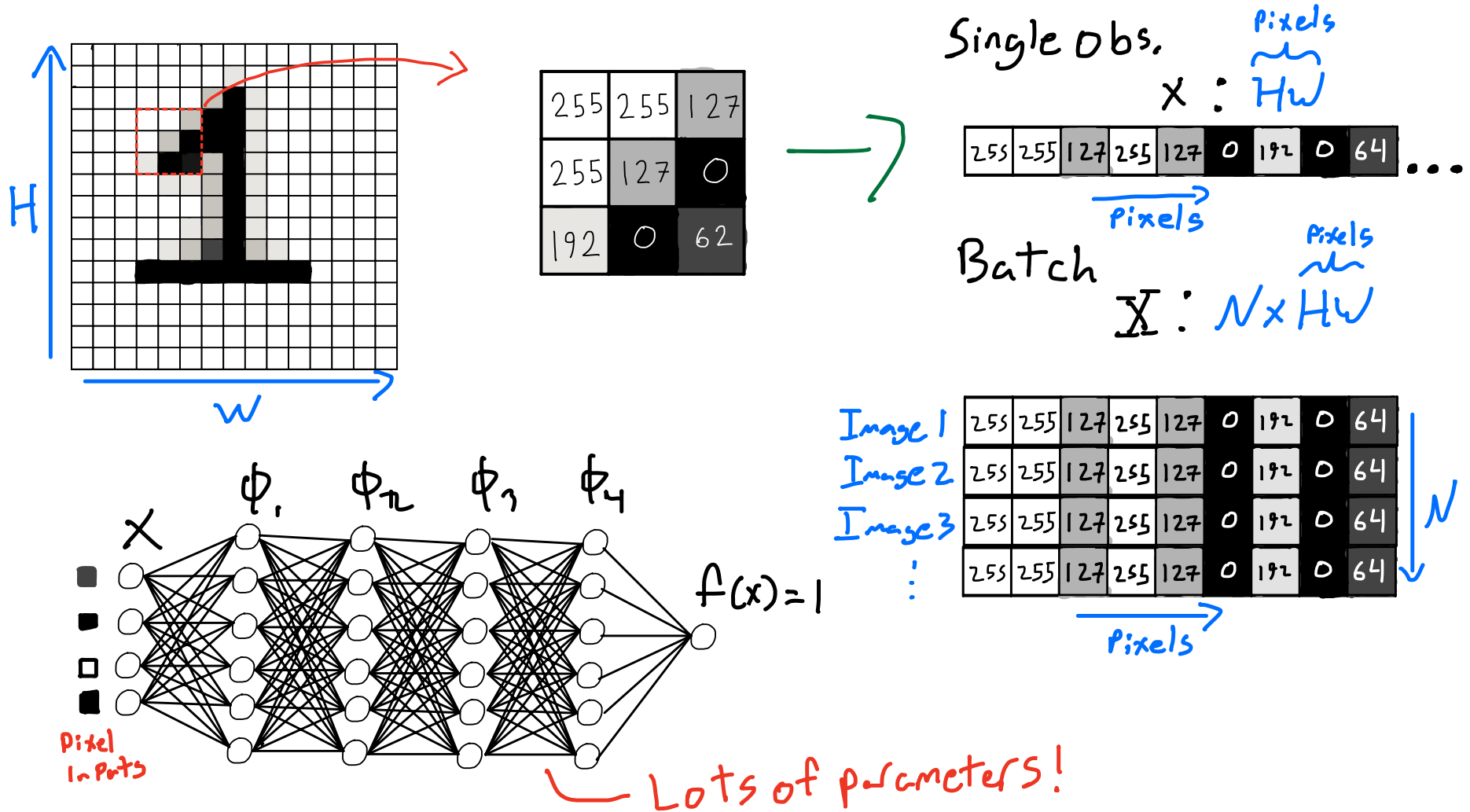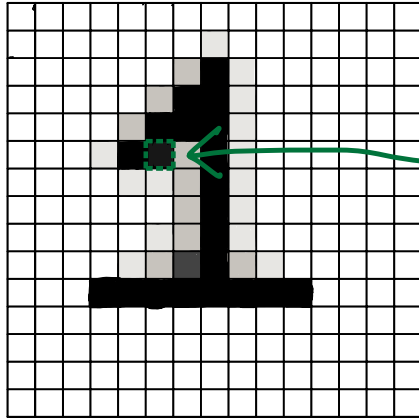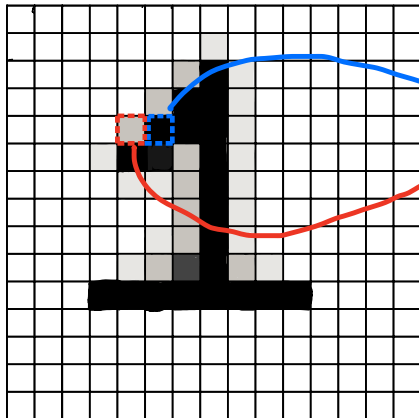$$f(x) = 1$$

Pixel Inputs

Lots of parameters!

# Image Structure
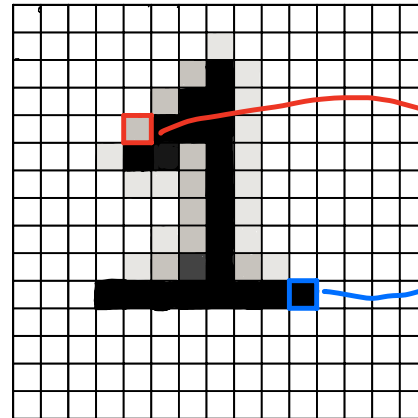
- Class determined by relationships between features

One pixel has little info by itself!

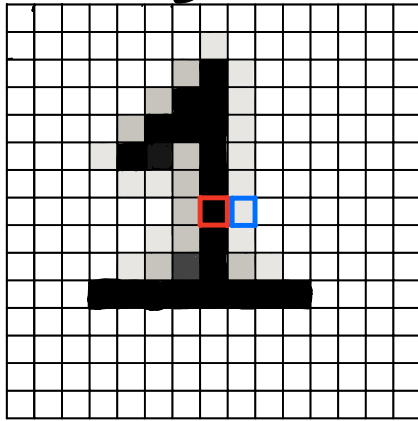- Relationships between nearby pixels are more important than far pixels
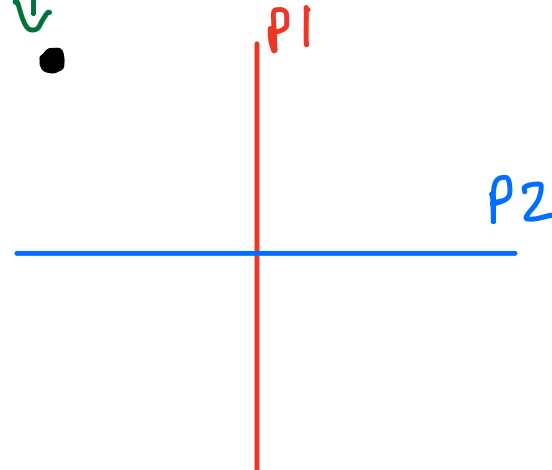
There is an edge here

?

# Image Structure
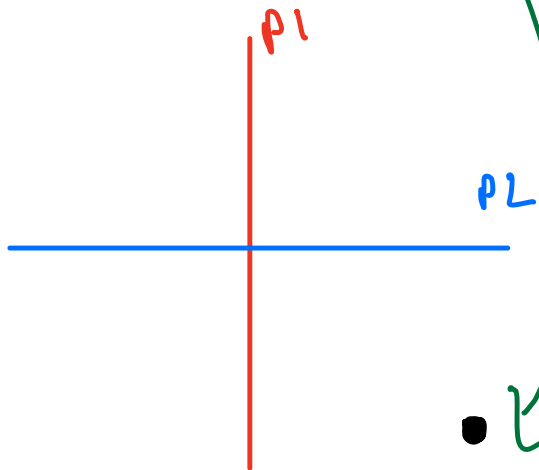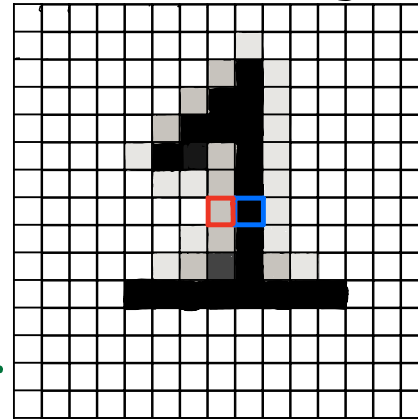
Original



Shifted right



Translation shouldn't change prediction

P1

P2

P1

P2

But features can change a lot!

# Logistic Regression

| Input | $W_c$ | $x \odot w$ | Pred | Input (shifted) | $W_c$ | $x \odot w$ | Pred |
|---|---|---|---|---|---|---|---|



— Regression weights from HW 3

— Shifting Image Changes predicted class!

Remember prediction function for multiclass Classification

→ Class w/ max output

Pred. $y = f(x) = \underset{c}{\mathrm{argmax}} \; W_c^T x$

# ConVolutional Networks (CNNs)

1) Maintain Image Structure
   (Don't flatten)

2) Shift weights to find best alignment

3) Make network sparse
   (Remove weights)

# Shifting Weights

Input (Shifted) W    C  ————  $X^T W$   At _every_ location



i.e. $C_{ij} \rightarrow$ Shift W to be centered at $i, j$ and compute $X^T W$

Predict 3 again!

$\rightarrow$ Predict with $\max(C)$

$\rightarrow$ Location where X and W most closely match

Convolution!

# Convolution in 1-dimension

1-Row



Input

0.7  0.5  0.3  1.0  0.0  0.7

×    ×    ×    ×    ×    ×     ⊕ —>

0.2  0.2  0.8  0.0  0.8  0.2

Weights

$(0.7)(0.2)+(0.5)(0.2)+(0.3)(0.8)+(1.0)(0.0)+(0.0)(0.8)+(0.7)(0.2)$

$= 0.77$ {

# Convolution in 1-dimension

1-Row

Input

0.7  0.5  0.3  1.0  0.0  0.7
×    ×    ×    ×    ×    ×    ×
0.1  0.1  0.8  0.0  0.8  0.2  —    —

⊕ —>

└ Assume 0

Weights

Try different alignment

$(0.7)(0.8)+(0.5)(0.0)+(0.3)(0.8)+(1.0)(0.2)$

$= 1.0$ {

Better than before!

# Convolution in 1-dimension

1-Row

Input

0.7  0.5  0.3  1.0  0.0  0.7

× × × × × × ×

0.2  0.2  0.8  0.0  0.8  0.2  —  —

$\oplus \rightarrow$

⌞ Assume 0

Weights

Try different alignment

$(0.7)(0.8)+(0.5)(0.0)+(0.3)(0.8)+(1.0)(0.2)$

$= 1.0$ {

Better than before!

# Convolution operator (1-D)

## Inputs:

$x$ : Array of length $d$

Kernel: Array of length $s$
(weights)

Typically: $s < d$

$x$

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

$\times \quad \times \quad \times$

kernel

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

$\oplus$

output

| | | 16 | | |
|---|---|---|---|---|

Only compute alignments
where kernel fully overlaps $x$

In torch: Padding = 'Valid'

$16 = 2 \cdot 1 + 4 \cdot 0 + 7 \cdot 2$

$$\text{Conv}(x, k)_i = \sum_{j=1}^{s} x_{i+j} k_j$$

Ouput length $= d - (s-1)$

# Convolution Animated!

## Convolution (kernel size: 3)

Output[ 0]=( 5)(-1)+( 7)( 3)+( 0)(-1)= 16

Output

| 16 | | | | | | | | | | | | | |

Kernel

| -1 | 3 | -1 |

Input

| 5 | 7 | 0 | -4 | -3 | 8 | 8 | 8 | -8 | 2 | -7 | -7 | 2 | -8 | -9 |

## Convolution (kernel size: 3)

Output
-6

Kernel
-2  0  2

Input
8  5  8  7  8  8
0  -4  -3  -8  -6  -8  -5  -5  -7

# Padding

- If we want to try <u>every</u> possible alignment we need to <u>pad</u> the input w/ 0s

- In torch: padding = 'full'

$$\text{Output length} = d + 2(s-1)$$

# Padding

**Convolution (size: 3, padding: 1)**

— often want padding
in-between

— e.g. If we want an
output size of $d$

— In torch: padding = 'same'

Output

Kernel

Input  0 | 6 | 6 | -6 | -9 | -1 | -3 | -6 | 1 | -3 | -9 | 5 | 3 | 3 | -9 | 1 | 0

Output length $= d$

# Convolution as a Layer

| Standard ('Dense') Layer | Locally-Connected Layer | Convolutional Layer |
|---|---|---|

$$\phi_1 = \sigma(x^T W + b)$$

$$\phi_1 = \sigma(x^T W^* + b)$$

$$\phi_1 = \sigma(\text{Conv}(x, k) + b)$$

Every output depends on _every_ input

Every output depends only on _Local_ inputs

_And_ weights are _Shared_ for each output!

# Derivatives of convolutions

In general:

$$\frac{dL}{dx_i} = \sum_{j=1}^{d} \frac{dL}{d\phi_j} \cdot \frac{d\phi_j}{dx_i}$$

but only $\phi_{i-1}, \phi_i, \phi_{i+1}$ depend on $x_i$ so:

$$\frac{dL}{dx_i} = \frac{dL}{d\phi_{i-1}} \frac{d\phi_{i-1}}{dx_i} + \frac{dL}{d\phi_i} \frac{d\phi_i}{dx_i} + \frac{dL}{d\phi_{i+1}} \frac{d\phi_{i+1}}{dx_i}$$

$$\parallel \qquad\qquad \parallel \qquad\qquad \parallel$$
$$K_1 \qquad\qquad K_2 \qquad\qquad K_3$$

$$\phi = conv(x, [K_1, K_2, K_3])$$

$$\frac{dL}{dx} = conv(\frac{dL}{d\phi}, [K_3, K_2, K_1])$$

X    kernel    $\phi_1$

— : $K_1$
— : $K_2$
— : $K_3$

$x_i$

$K_1$
$K_2$
$K_3$

$\phi_{i-1}$

$\phi_i$

$\phi_{i+1}$

# Derivatives Animated!

**Convolution (kernel size: 3)**

# Convolutions in 2-D

Align kernel in every 2-d location

$$\text{Conv2d}(x, k)_{ij} = \sum_{a=1}^{s} \sum_{b=1}^{s} x_{i+a, j+b} \cdot k_{ab}$$

[For padding = 'Valid']

| 5 | -1 | 3 | 8 | 4 | 6 | -2 |
|---|----|---|---|---|---|----|
| -3 | 4 | 9 | 1 | 1 | 7 | -4 |
| 5 | -6 | 3 | 2 | 1 | -2 | 0 |
| 0 | -8 | 5 | 5 | 4 | -2 | 3 |
| 4 | 7 | 2 | -7 | 3 | 1 | 4 |
| 4 | -3 | 1 | 2 | 6 | 1 | -1 |
| 0 | 0 | 3 | -2 | 3 | -1 | 0 |

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

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

$$1 = 1 \cdot (-1) + 1 \cdot (-2) + 7(-1) + 2 \cdot 0 + 1 \cdot 0 +$$
$$(-2) \cdot 0 + 5 \cdot 1 + 4 \cdot 2 + (-2) \cdot 1$$

# Convolutions in 2-D (Blur)

## Small Blur

| 0.12 | 0.12 | 0.12 |
|------|------|------|
| 0.12 | 0.12 | 0.12 |
| 0.12 | 0.12 | 0.12 |

## Large Blur

| 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |
|------|------|------|------|------|
| 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |
| 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |
| 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |
| 0.04 | 0.04 | 0.04 | 0.04 | 0.04 |

# Convolutions in 2-D (Edge detect)

## Vertical edges

pos.

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

— neg.



most outputs near 0

Low values at neg. of kernel

High values where pattern matches kernel

## Horizontal edges

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

# Convolutions in 2-D (Edge detect)

## Diagonal edges

# Convolutions for color images



Image(x)

Kernel (k)

Red | green | Blue

result

Kernel has 3 channels like Image!

$$\text{Conv}(x, k) \overset{H \times W \times 3}{=} \overset{5 \times 5 \times 3}{}$$

$$\text{Conv}(x_r, k_r) +$$

$$\text{Conv}(x_g, k_g) +$$

$$\text{Conv}(x_b, k_b)$$

# Multi-Channel Convolution (1-D)



Convolution (size: 3, channels: 3, output channels: 1)

Not limited to
3-channels!

$$\text{Conv}(\ddot{X}, K) = \sum_{c=1}^{C} \text{Conv}(X_c, K_c)$$

Sum over
Channels!

$H \times W \times C$     $S \times S \times C$          $H \times W$     $S \times S$

# Multiple output channels

Kernel can also produce multiple channels (>1 value at each location)

## 1-D

**Convolution (size: 3, channels: 3, output channels: 3)**



## In 2-D

Combine 3 filters into a multi-channel Image

# General 2-D Convolutional Layer

Input ($\underline{X}$): $N \times H \times W \times C$

<span style="color:red">\# Images (Batch)</span>    <span style="color:green">Height Width</span>    <span style="color:blue">\# Channels</span>

Kernel ($\bar{\bar{K}}$): $S \times S \times C \times O$

<span style="color:purple">Kernel size</span>   <span style="color:blue">\#channels</span>   <span style="color:orange">\# output channels</span>

Output ($\phi$): $N \times \underline{H \times W} \times O$

<span style="color:green">May change for padding $\neq$ 'same'</span>

Diagram:



Input    W    Kernel    — Output

H    S   S

C

— Single 'neuron'

$= \Sigma x_w \odot K_o$

(Still a sum-product!)

<span style="color:magenta">Follow w/ Activation</span>

O

# PyTorch 2-D Convolutional Layer

Input $(\underline{X})$: $N \times C \times H \times W$

$\qquad$ #Images (Batch) $\quad$ #Channels $\quad$ Height $\quad$ Width

Channels before Image size

Kernel $(\underline{\overline{K}})$: $C \times O \times S \times S$

$\qquad$ #channels $\quad$ #output channels $\quad$ Kernel size

Output $(\phi)$: $N \times O \times H \times S$

## CONV2D

CLASS  torch.nn.Conv2d(*in_channels*, *out_channels*, *kernel_size*, *stride=1*, *padding=0*,
$\qquad$ *dilation=1*, *groups=1*, *bias=True*, *padding_mode='zeros'*, *device=None*,
$\qquad$ *dtype=None*)  [SOURCE]

Applies a 2D convolution over an input signal composed of several input planes.

In the simplest case, the output value of the layer with input size $(N, C_{\text{in}}, H, W)$ and output $(N, C_{\text{out}}, H_{\text{out}}, W_{\text{out}})$ can be precisely described as:

$$\text{out}(N_i, C_{\text{out}_j}) = \text{bias}(C_{\text{out}_j}) + \sum_{k=0}^{C_{\text{in}}-1} \text{weight}(C_{\text{out}_j}, k) \star \text{input}(N_i, k)$$

where $\star$ is the valid 2D cross-correlation operator, $N$ is a batch size, $C$ denotes a number of channels, $H$ is a height of input planes in pixels, and $W$ is width in pixels.

Technically what we call Convolution is cross correlation

# Convolutions as feature detectors



Images

Kernel channels

Outputs

Detected eye
nose
mouth

ReLU as threshold

$-b$

$\text{Max}(\text{conv}(x,k)+b, 0)$

# Multi-Layer CNNs

Detect low-level features

Detect high-level features

Detected face

Exact location of features doesn't matter!

# Down Sampling

Reduce resolution → Easier and less expensive
to find high-level features

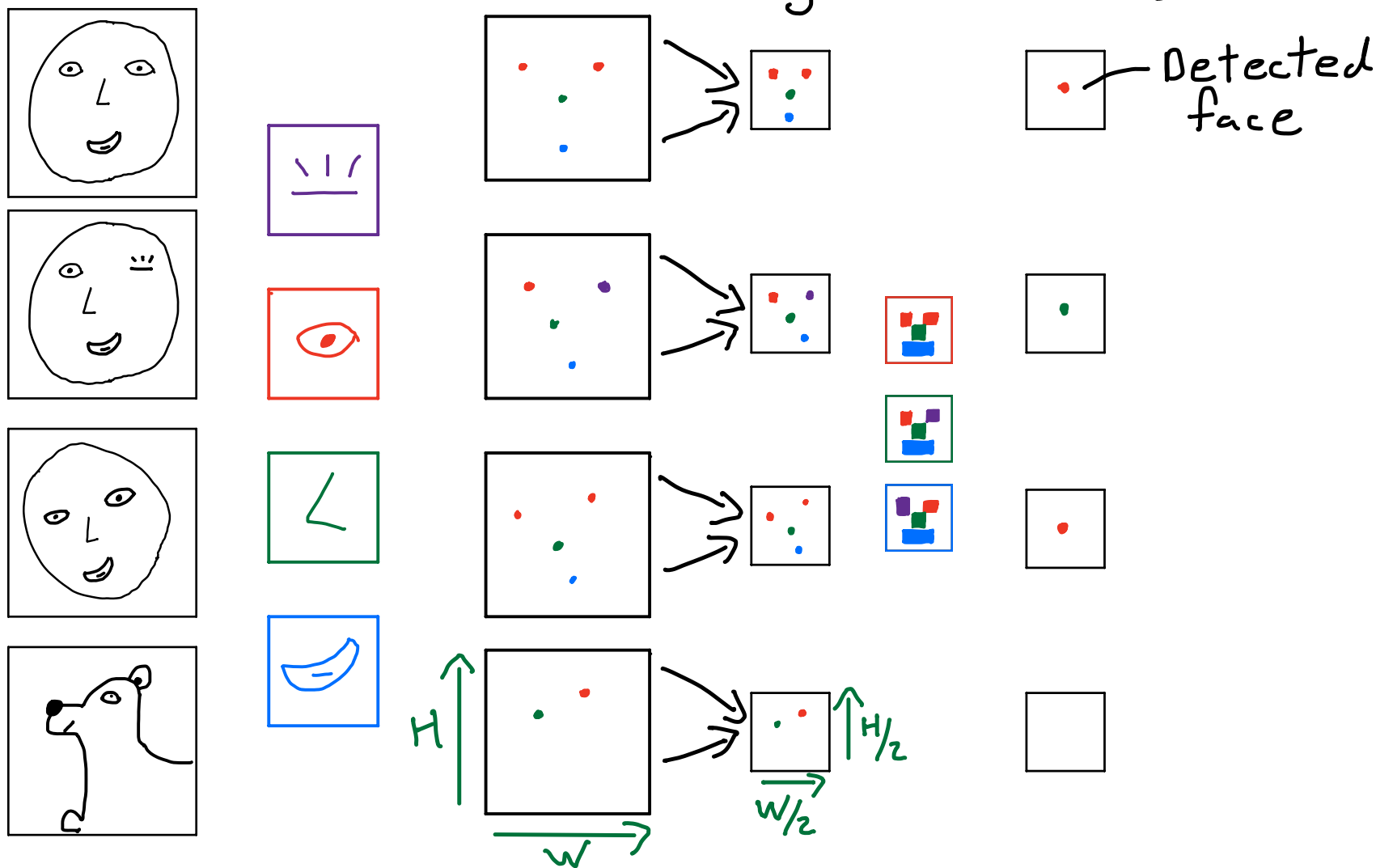Detected face

$H$ $W$ $H/2$ $W/2$

# Strided Convolution (1-D)

~ Take only $\frac{1}{stride}$ × Outputs

x  kernel $\phi_1$   x  kernel $\phi_1$   x  kernel $\phi_1$

Stride = 1      Stride = 2      Stride = 3

**Convolution (size: 3, stride: 2)**

| Output | 30 | | | | | | |

| Kernel | -1 | 3 | -1 |

| Input | -9 | 5 | -6 | -5 | 1 | -6 | 4 | 3 | 4 | -8 | -5 | 8 | 6 | -9 | 1 |

What if the optimal activation is here?

# Strided Convolutions (2-D)

## — Apply stride in each dimension



Skip Locations

| 5 | -1 | 3 | 8 | 4 | 6 | -2 |
|---|----|---|---|---|---|----|
| -3 | 4 | 9 | 1 | 1 | 7 | -4 |
| 5 | -6 | 3 | 2 | 1 | -2 | 0 |
| 0 | -8 | 5 | 5 | 4 | -2 | 3 |
| 4 | 7 | 2 | -7 | 3 | 1 | 4 |
| 4 | -3 | 1 | 2 | 6 | 1 | -1 |
| 0 | 0 | 3 | -2 | 3 | -1 | 0 |

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

Stride = 2

# Pooling operator (1-D)

Inputs:  $x$ : Array of length $d$

window size $s$

$x$

| 5 | 3 | 1 | 4 | 7 | 1 | 2 |

window

output

max

| | | 7 | | |

$\cdot 7 = \max(1, 4, 7)$

- Take max output around each location before downsampling

- Make sure we don't miss any high activations

- Can also take Avg. (Average Pooling)
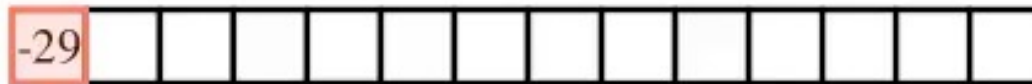
# Convolution + MaxPooling Animated!

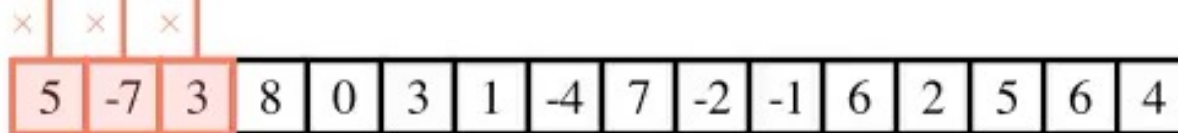## Convolution (size: 3) + Max Pooling (size: 3, stride: 2)

Output — ☐ ☐ ☐ ☐ ☐ ☐
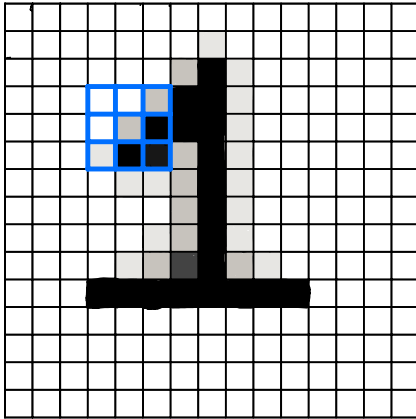
Pooling

Conv. output — | -29 | | | | | | | | | | | | | | |

Kernel — | -1 | 3 | -1 |

$\times$ $\times$ $\times$

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

# Pooling in 2-D

Align window in every 2-d location

$$\text{Max-Pool}(x)_{ij} = \overset{3}{\underset{a=1}{\text{Max}}}\left(\overset{5}{\underset{b=1}{\text{Max}}}\left(x_{i+a, j+b}\right)\right)$$

$$\text{Avg.-Pool}(x)_{ij} = \frac{1}{s^2}\sum_{a=1}^{5}\sum_{b=1}^{5}x_{i+a, j+b}$$

| 5 | -1 | 3 | 8 | 4 | 6 | -2 |
|---|----|---|---|---|---|----|
| -3 | 4 | 9 | 1 | 1 | 7 | -4 |
| 5 | -6 | 3 | 2 | 1 | -2 | 0 |
| 0 | -8 | 5 | 5 | 4 | -2 | 3 |
| 4 | 7 | 2 | -7 | 3 | 1 | 4 |
| 4 | -3 | 1 | 2 | 6 | 1 | -1 |
| 0 | 0 | 3 | -2 | 3 | -1 | 0 |

max

7

$$7 = \text{max}(1, 1, 7, 2, 1, -2, 5, 4, -2)$$

# DownSampling Comparison



|  | Image | Stride=2 | Maxpool size=2, stride=1 | Maxpool size=2, stride=2 | Avgpool size=2, stride=1 | Avgpool size=2, stride=2 |

|  | Image | Stride=4 | Maxpool size=4, stride=1 | Maxpool size=4, stride=4 | Avgpool size=4, stride=1 | Avgpool size=4, stride=4 |

|  | Image | Stride=4 | Maxpool size=4, stride=1 | Maxpool size=4, stride=4 | Avgpool size=4, stride=1 | Avgpool size=4, stride=4 |

# Pixel Shuffle

- Rearrange adjacent results into more channels — Expensive!

**Convolution (size: 3) + Shuffle**

| | | | | | |
|---|---|---|---|---|---|
| Output | -4 | | | | |

Kernel: | -1 | 3 | -1 |

× × ×

Input: | 8 | 2 | 2 | 0 | -2 | 0 | -2 | -2 | 6 | 7 | 5 | 2 | -2 | -6 | 1 | -8 |

# Flattening

— After convolutions, flatten as before



Can add 'dense' layers here!

$$N \times H \times W \times C \xrightarrow{\text{flatten}} N \times (H \cdot W \cdot C)$$
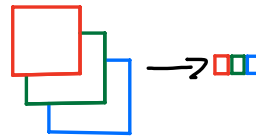
# Global Avg. Pooling

— Alt. to flattening, just Avg. over remaining Image

$$\text{Global Avg.-Pool}(x) = \frac{1}{W \cdot H} \sum_{a=1}^{W} \sum_{b=1}^{H} X_{a,b}$$

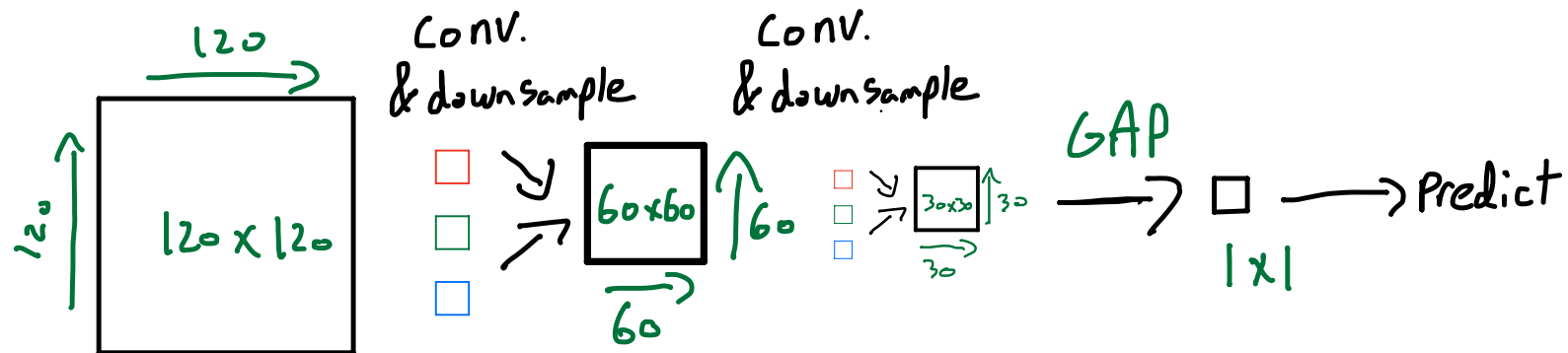| 5 | -1 | 3 | 8 | 4 | 6 | -2 |
|---|---|---|---|---|---|---|
| -3 | 4 | 9 | 1 | 1 | 7 | -4 |
| 5 | -6 | 3 | 2 | 1 | -2 | 0 |
| 0 | -8 | 5 | 5 | 4 | -2 | 3 |
| 4 | 7 | 2 | -7 | 3 | 1 | 4 |
| 4 | -3 | 1 | 2 | 6 | 1 | -1 |
| 0 | 0 | 3 | -2 | 3 | -1 | 0 |

→ 3

Remember: still multiple channels!

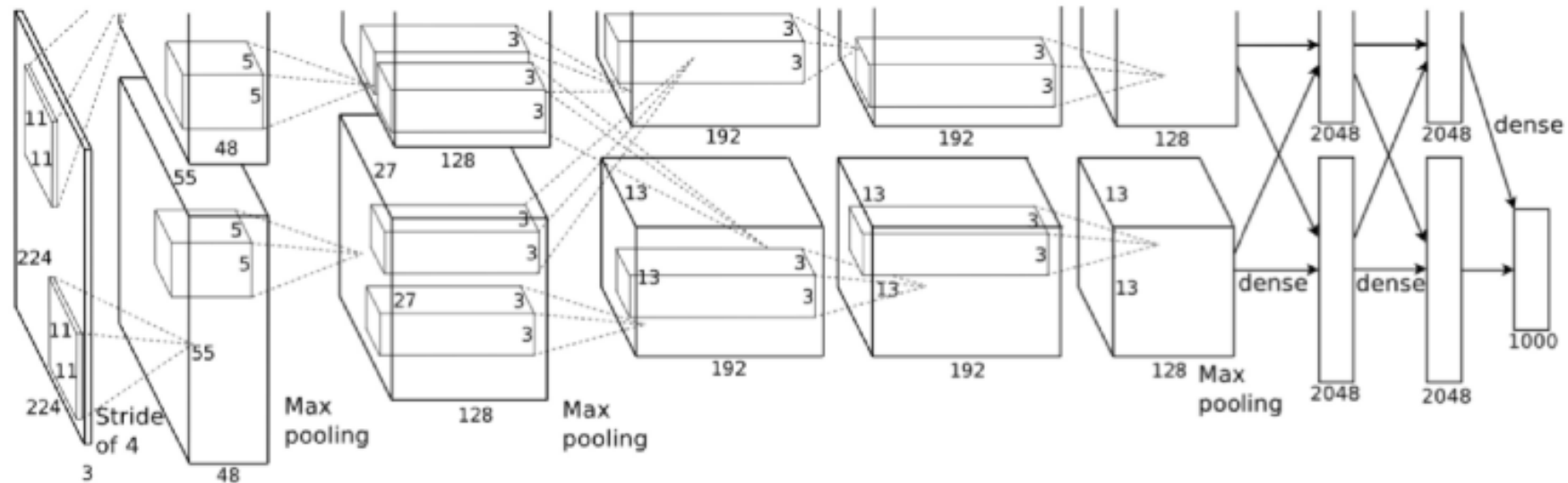$$N \times H \times W \times C \xrightarrow{\text{GAP}} N \times C$$

# Global Avg. Pooling

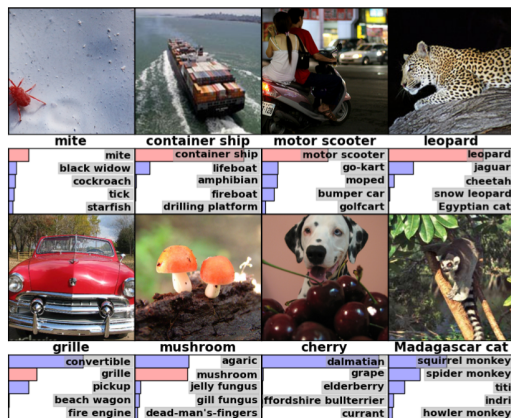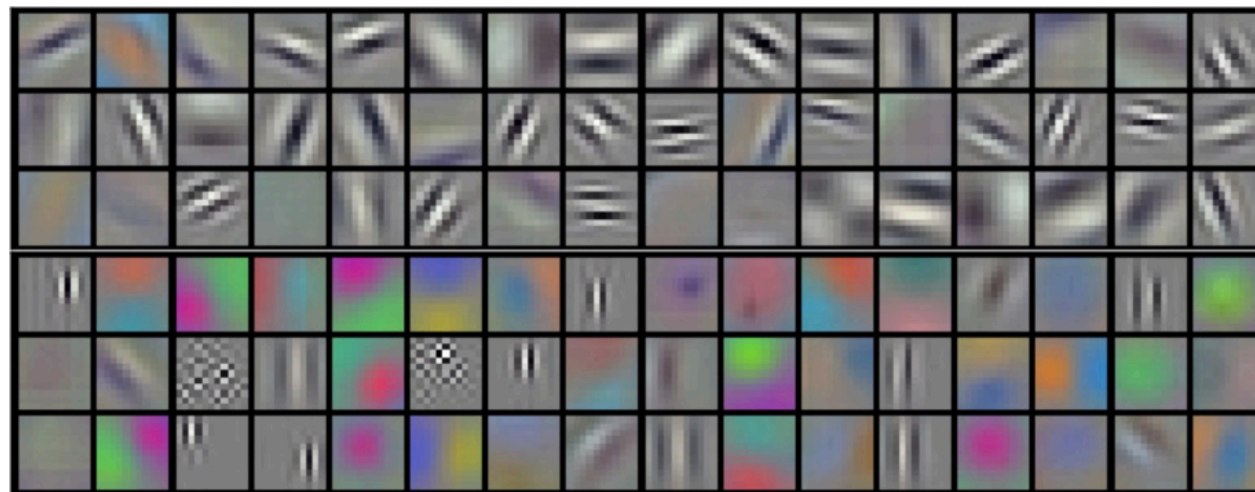— Allows for inputs of different sizes!



Conv. & downsample → Conv. & downsample → GAP → Predict

120 × 120 → 60 × 60 → 30 × 30 → 1 × 1

100 × 200 → 50 × 100 → 25 × 50 → 1 × 1

Removes H and W !

# AlexNet (2012)

## Architecture:



## 1000 Classes



## Learned kernels

# CNN Features by Layer



Layer 1

Layer 2

Layer 3

Layer 4

Channels

Images

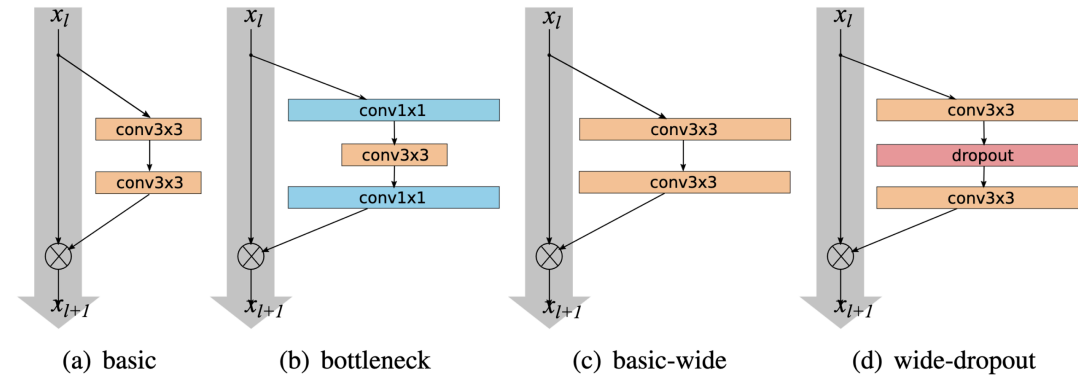# Other Architectures

← VGG, Resnet-50, etc.

Wide resnet ↓

Many, Many more!



Figure 1: Various residual blocks used in the paper. Batch normalization and ReLU precede each convolution (omitted for clarity)

# Data Augmentation

− It takes <u>a lot</u> of data to train

good Image classifiers!

~ millions to <u>Billions</u> of Images for
general object recognition (1000+ classes)

# Data Augmentation

- CNNs (mostly) invarient to translation

  – What about scale, rotation, color etc. ?



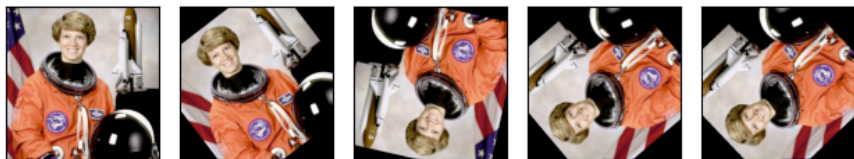Astronaut

Still Astronaut!

– Still shouldn't Change Class!

# Data Augmentation

– Augment existing data by randomly Scaling, rotating, Shifting colors etc.

– <u>Much</u> easier than collecting ~10x the data

– Can do this as we train!

# Common Augmentations
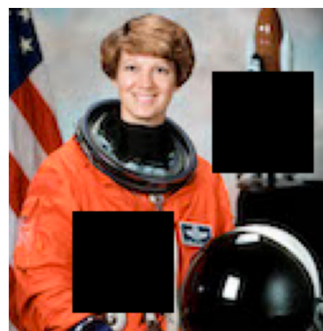
### Rotation



### Crop



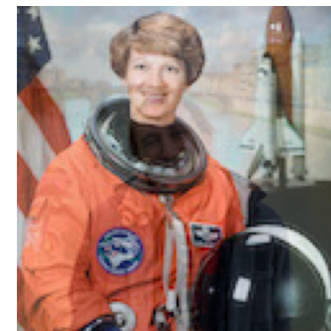### Color shift



### Shift / scale / shear



### Flip



### Cut-out



### Mix-up

# Fine Tuning

— Low-level features can often be shared between Models

**Detect low-level features**

**Old model**

• • •

→ Face?

• • •

— Exploit this by copying convolutional layers from a previously trained model

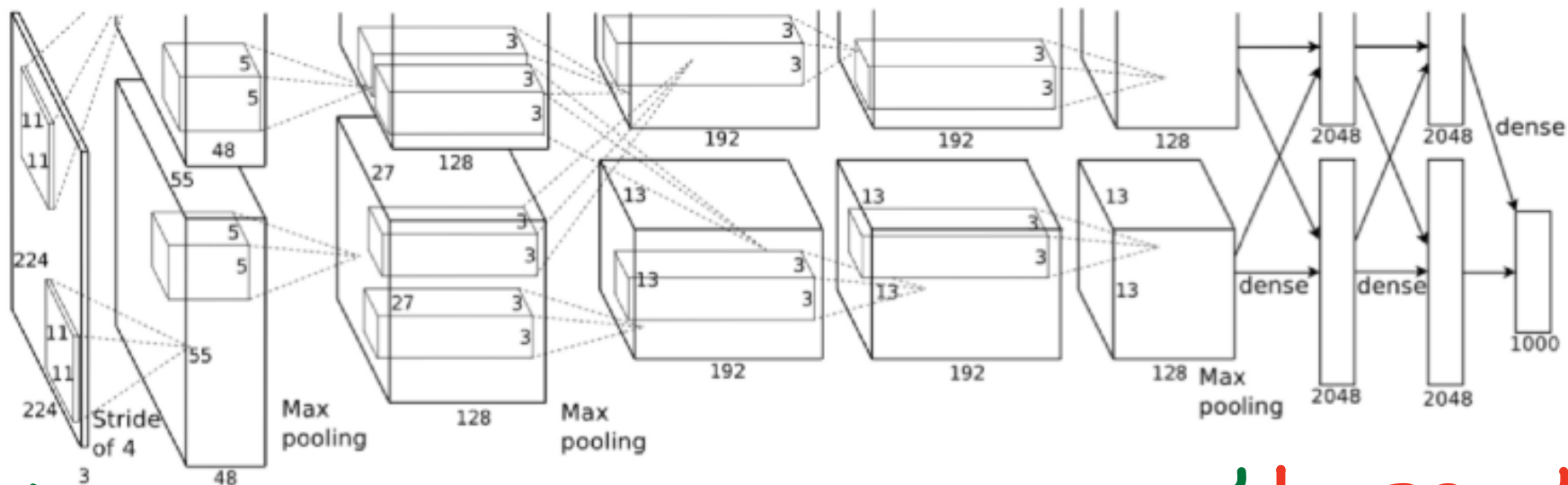**Detect low-level features**

**New model**

• • •

→ Smiling?

• • •

Actually more like this! ↓

# Fine Tuning    For new task:



keep
Convolutional Layers

Replace
Dense Layers

Typical Approach, but could keep more or less

# Fine Tuning

— Better starting point for optimization



Loss

Not Fine-tuning

Fine-tuning

— Without Fine tuning
— With Fine-tuning

Loss

Iteration