# Convolutional Neural Networks

## CS4782: Intro to Deep Learning

# Thanks to:

Varsha Kishore
Justin Lovelace
Anissa Dallmann
Stephanie Ginting
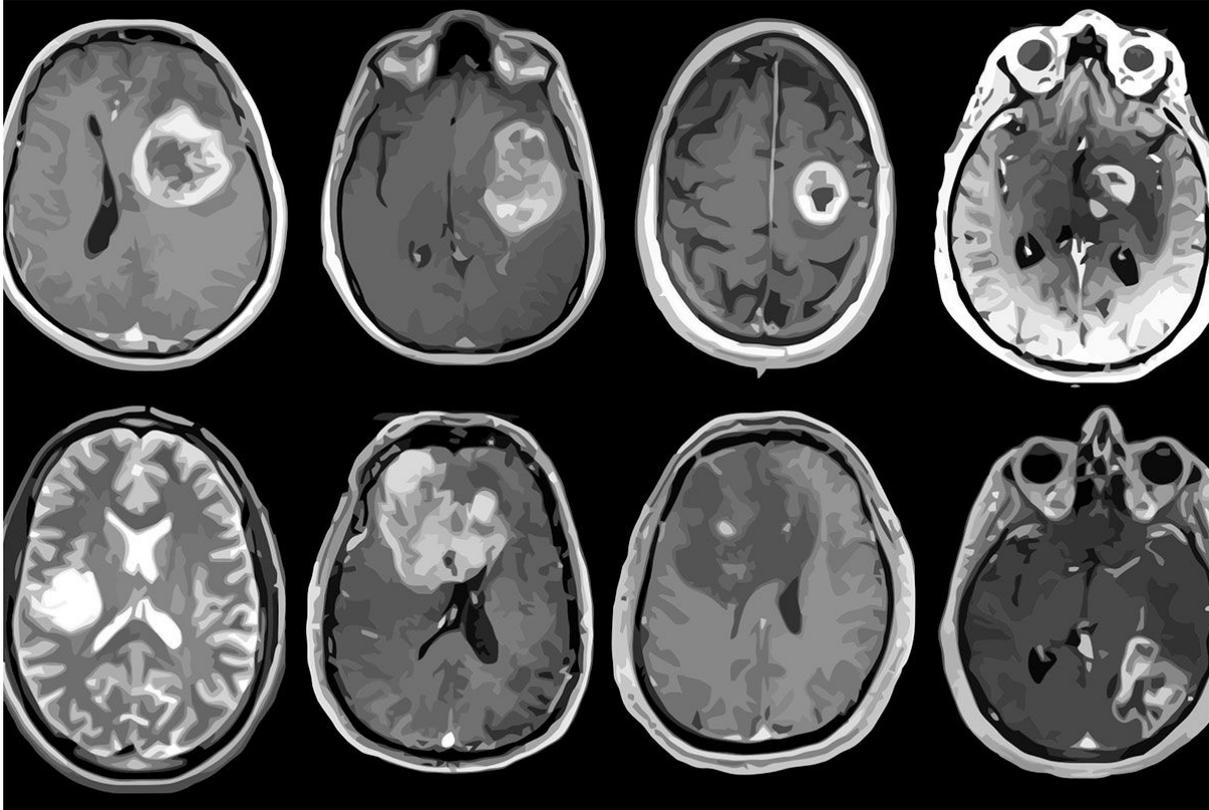Alexander Scotte

# Image Classification



input image

classification ⟶ "dog"



input image

classification ⟶ "cat"

# Applications in Medicine

# Applications in Autonomous Driving

# Image Classification



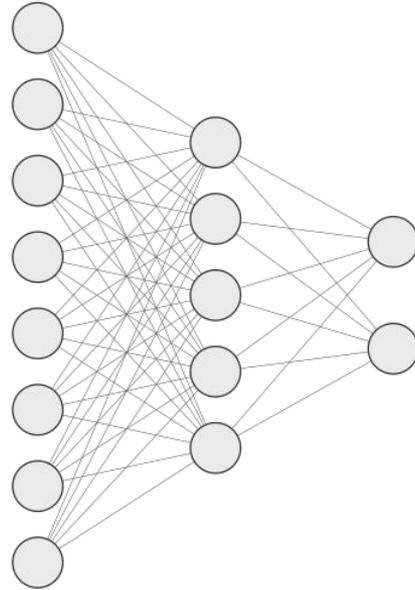input image

classification →

"dog"



input image

classification →

"cat"

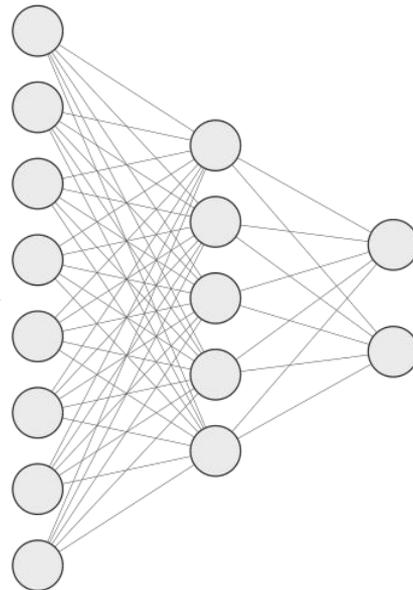# How to use Multi-Layer Perceptron for image classification?

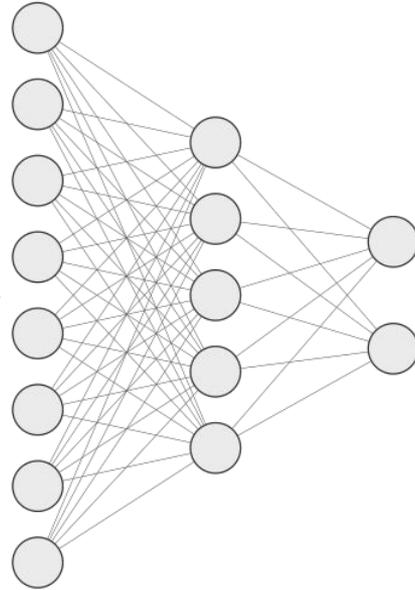# How to use Multi-Layer Perceptron for image classification?



flatten

# Why not use a Multi-Layer Perceptron?



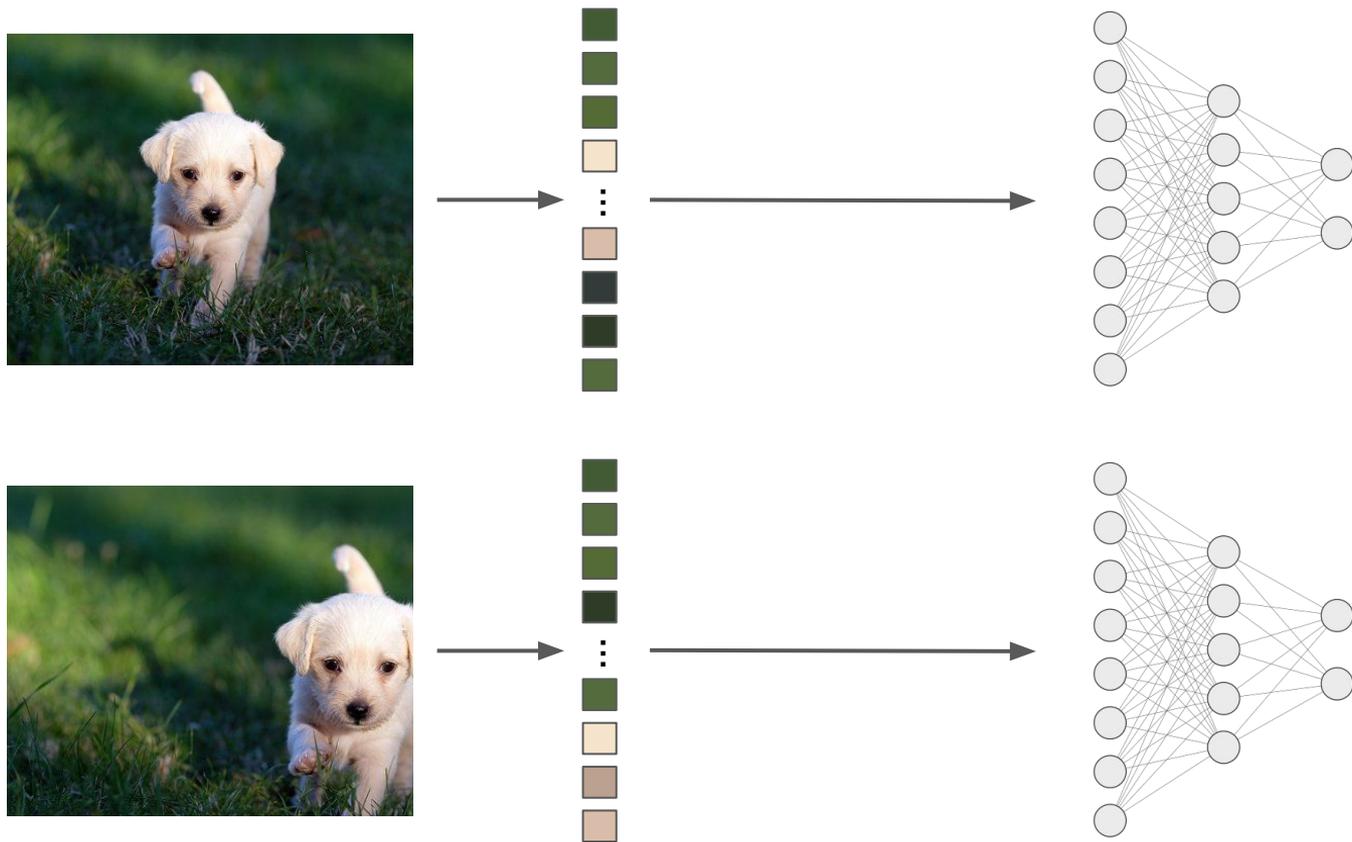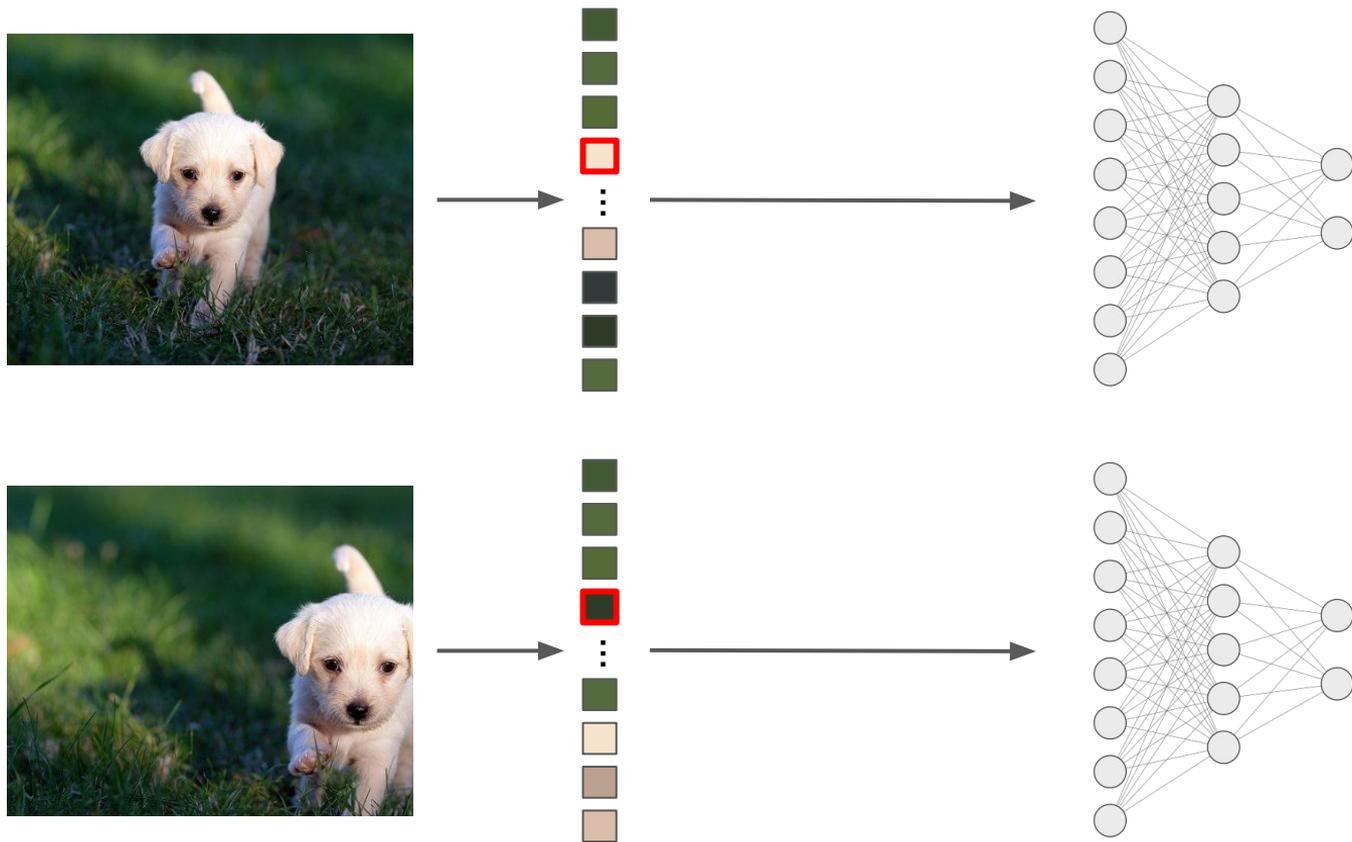flatten

Which pixels were next to each other?

# Why not use a Multi-Layer Perceptron?

# Why not use a Multi-Layer Perceptron?

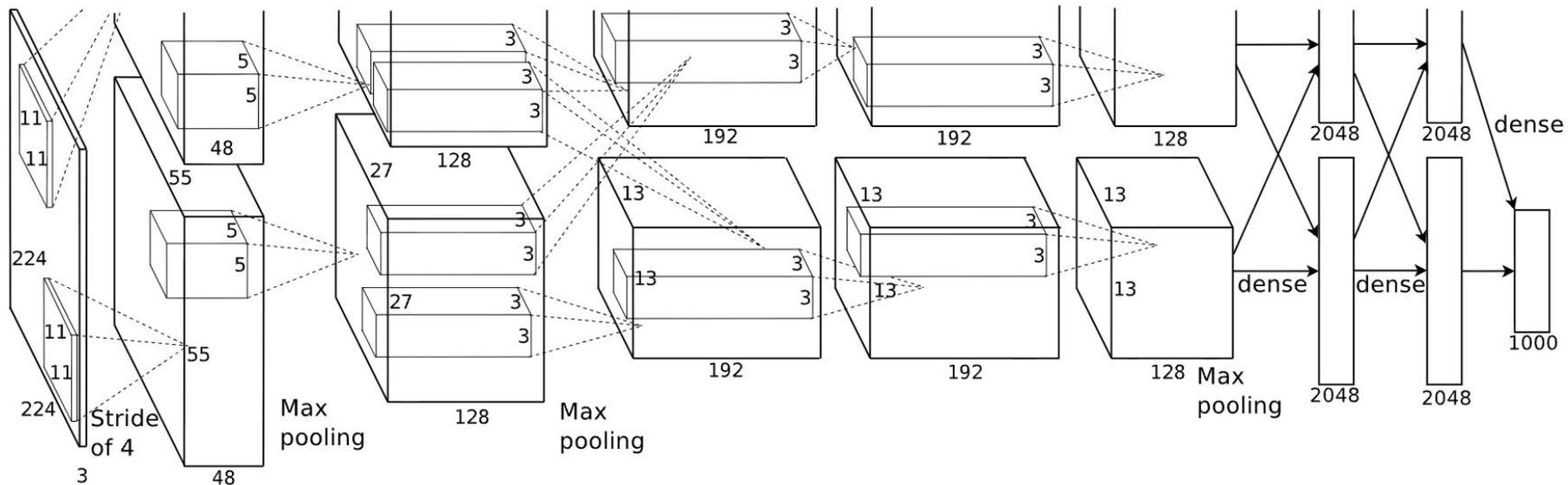# Why not use a Multi-Layer Perceptron?



💸many pixels = many parameters

# Our goal today

# Conv2d

```
class torch.nn.Conv2d(in_channels, out_channels, kernel_size, stride=1, padding=0,
dilation=1, groups=1, bias=True, padding_mode='zeros', device=None, dtype=None)
```

[source]

# Our goal today

# Our goal today

# Convolutional Filters



"image"

convolutional filter

# Convolutional Filters



"image"

convolutional filter

Convolutional Filters

# Convolutional Filters



"image" * convolutional filter =

# Convolutional Filters



"image"

convolutional filter

# Convolutional Filters



"image" * convolutional filter =

# Convolutional Filters



"image"

convolutional filter

# Convolutional Filters



"image" * convolutional filter =

# Convolutional Filters

- ❖ Aggregates information from local window around pixel

- ❖ Translational invariance

- ❖ Reduce number of parameters needed to be learned



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

"image"

\*

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

convolutional filter

=

| 3 | 2 | 2 |
|---|---|---|
| 1 | 3 | 2 |
| 2 | 1 | 2 |

# Discuss with your Neighbor!

Match the following convolutional filters with the output they produce.

input image

# Convolutional Filters

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

"image"

\*

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

convolutional filter

can learn this!

=

| 3 | 2 | 2 |
|---|---|---|
| 1 | 3 | 2 |
| 2 | 1 | 2 |

# Convolutional Filters



"image" * convolutional filter = 3

# Convolutional Filters



"image" * convolutional filter =

# CNNs - Padding

- ❖ Padding adds layers of zeros (or other number) around image border

- ❖ Prevents image shrinking and loss of information from image boundary



**Input**

**Filter**

**Result**

**Parameters:**

Size: $f = 3$

Stride: $s = 2$

**Padding: $p = 1$**

$= 0*1 + 0*0 + 0*(-1) +$
$0*1 + 4*0 + 9*(-1) +$
$0*1 + 9*0 + 6*(-1)$
$= -15$

Dimension: 6 x 6

https://indoml.com

# CNNs - Padding

‹ **mode** : *{'full', 'valid', 'same'}, optional*

‹

    **'full':**

        By default, mode is 'full'. This returns the convolution at each point of overlap, with an output shape of (N+M-1,). At the end-points of the convolution, the signals do not overlap completely, and boundary effects may be seen.

    **'same':**

        Mode 'same' returns output of length `max(M, N)`. Boundary effects are still visible.

    **'valid':**

        Mode 'valid' returns output of length `max(M, N) - min(M, N) + 1`. The convolution product is only given for points where the signals overlap completely. Values outside the signal boundary have no effect.

*Dimension: 6 x 6*

# CNNs - Padding

- **mode** : *{'full', 'valid', 'same'}, optional*

  **'full':**

  By default, mode is 'full'. This returns the convolution at each point with an output shape of (N+M-1,). At the end-points of the convolution, signals do not overlap completely, and boundary effects may be seen.

  **'same':**

  Mode 'same' returns output of length `max(M, N)`. Boundary effects visible.

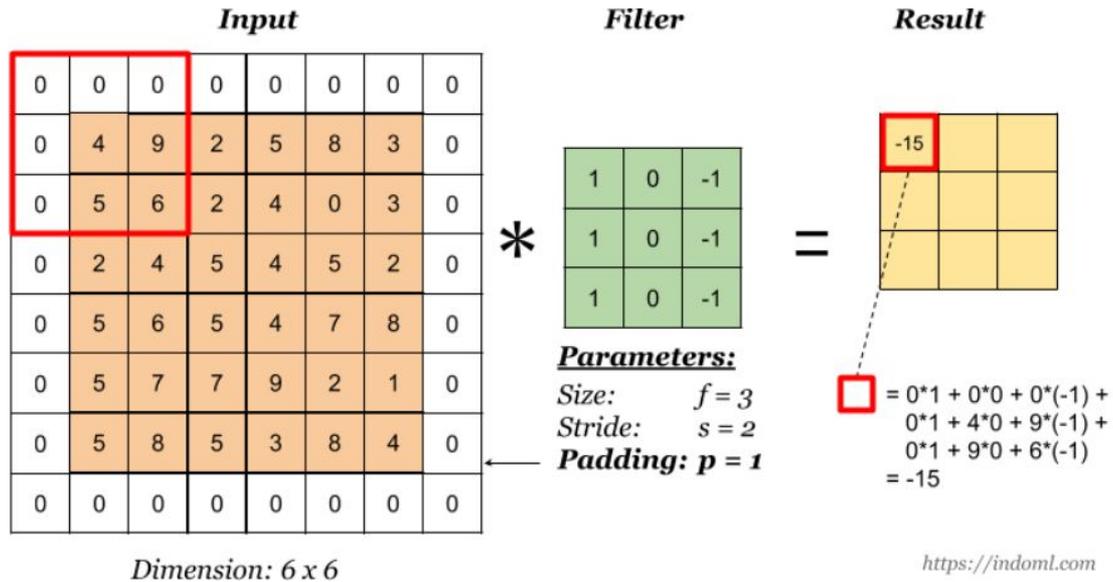  **'valid':**

  Mode 'valid' returns output of length `max(M, N) - min(M, N) + 1`. convolution product is only given for points where the signals overlap completely. Values outside the signal boundary have no effect.
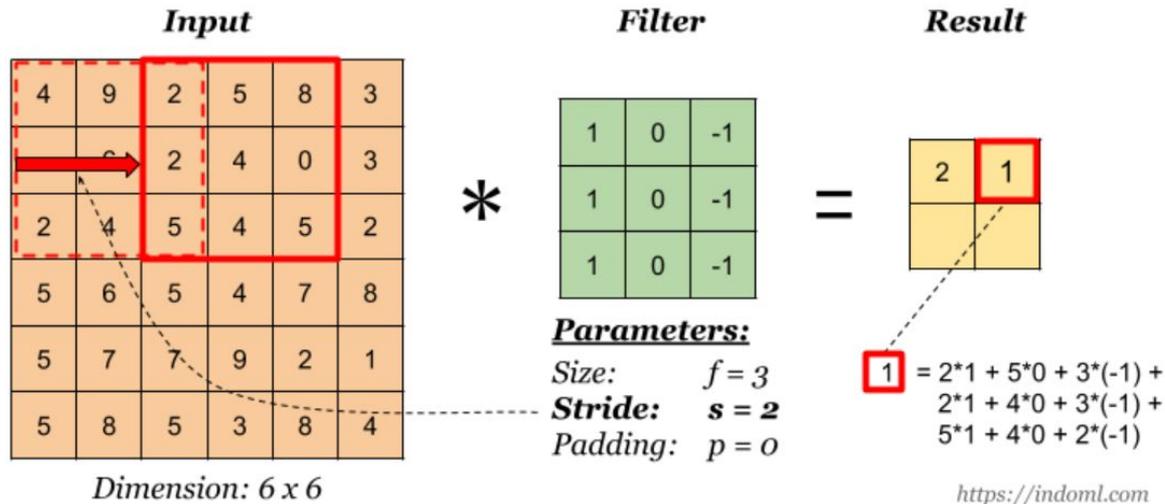
Dimension: 6 x 6

https://indoml.com

# CNNs - Padding

- `padding` controls the amount of padding applied to the input. It can be either a string {'valid', 'same'} or an int / a tuple of ints giving the amount of implicit padding applied on both sides.

**Input**

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

Dimension: 6 x 6

\*

**Filter**

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

**Parameters:**

Size: $f = 3$

Stride: $s = 2$

**Padding: $p = 1$**

=

**Result**

| -15 | | |
|-----|---|---|
| | | |
| | | |

= 0\*1 + 0\*0 + 0\*(-1) +
0\*1 + 4\*0 + 9\*(-1) +
0\*1 + 9\*0 + 6\*(-1)
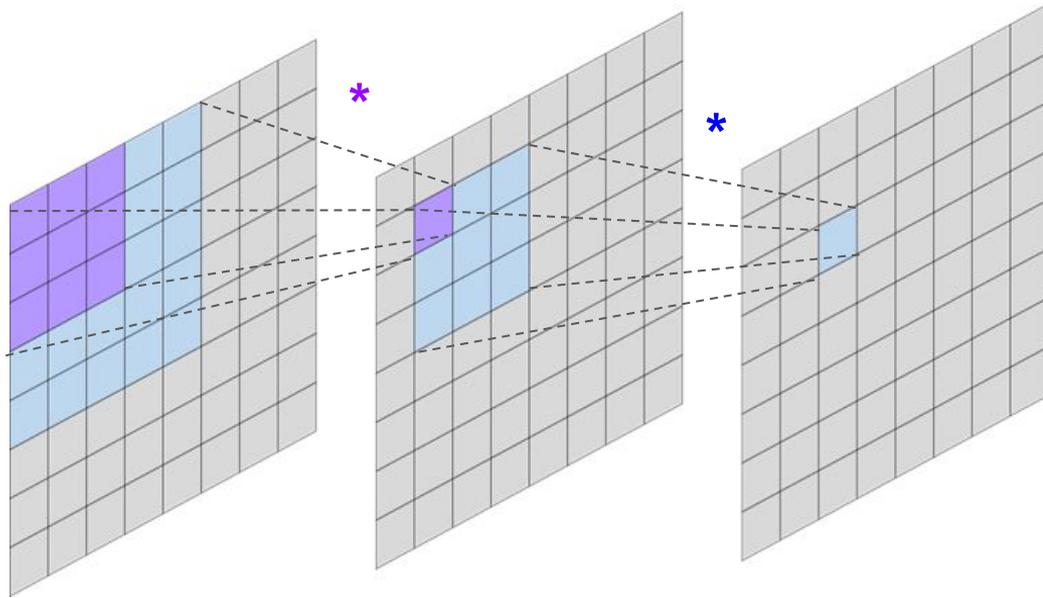= -15

https://indoml.com

# CNNs - Stride

- ❖ Stride controls how many units the filter / the receptive field shift at a time

- ❖ The size of the output image shrinks more as the stride becomes larger

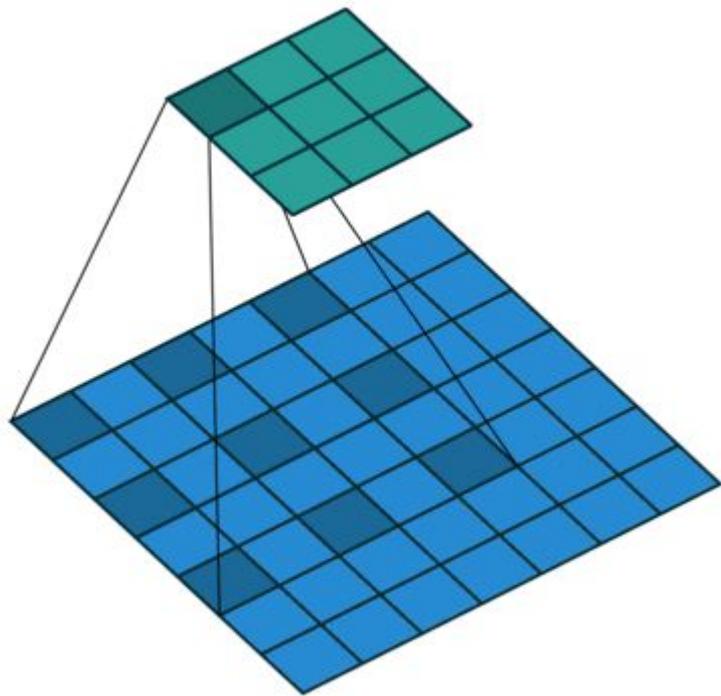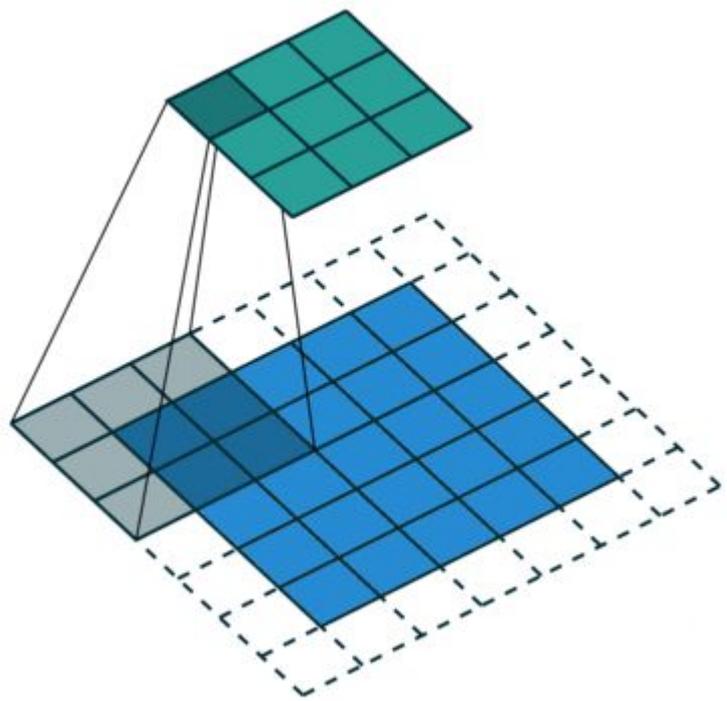- ❖ The receptive fields overlap less as the stride becomes larger



**Input**

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

Dimension: 6 x 6

**Filter**

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

\*

**Parameters:**
Size:     $f = 3$
**Stride:**    $s = 2$
Padding:   $p = 0$

=

**Result**

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

1 = 2\*1 + 5\*0 + 3\*(-1) +
2\*1 + 4\*0 + 3\*(-1) +
5\*1 + 4\*0 + 2\*(-1)

https://indoml.com

Filter with stride (s) = 2

# Stacking Convolutions



❖ Size of receptive field increases with each layer

❖ Capture more complex features

# Dilated Convolutions



https://github.com/vdumoulin/conv_arithmetic
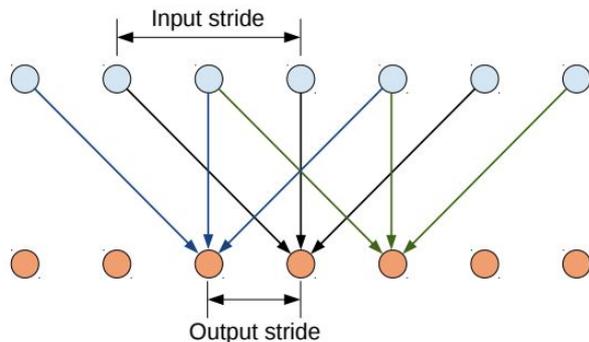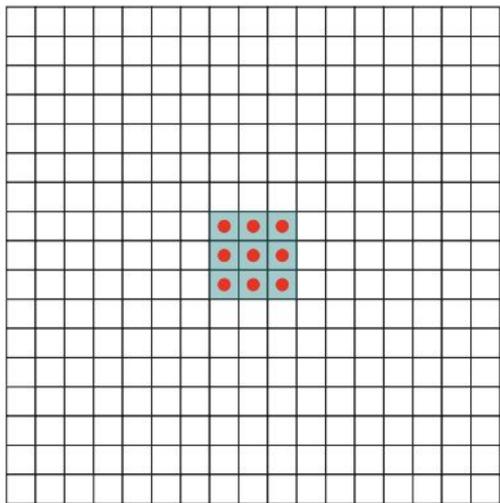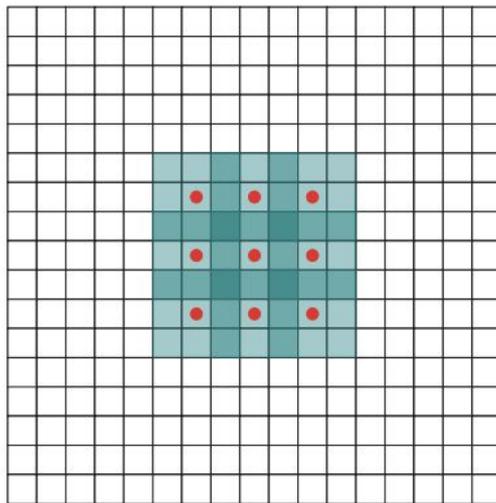
# Dilated convolution or Atrous convolution?

Figure 1: Illustration of the hole algorithm in 1-D, when *kernel_size = 3*, *input_stride = 2*, and *output_stride = 1*.
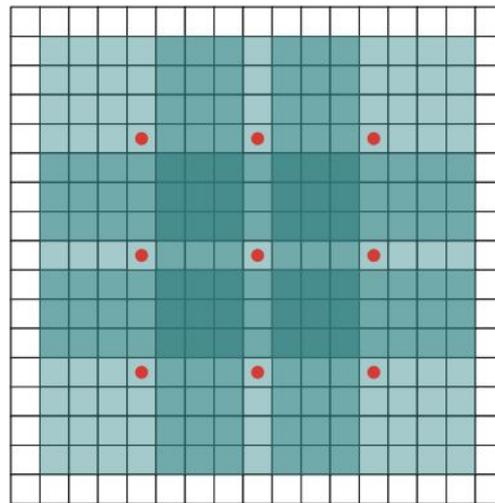
# Dilated convolution or Atrous convolution?



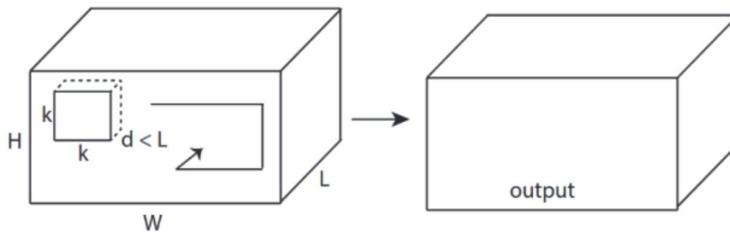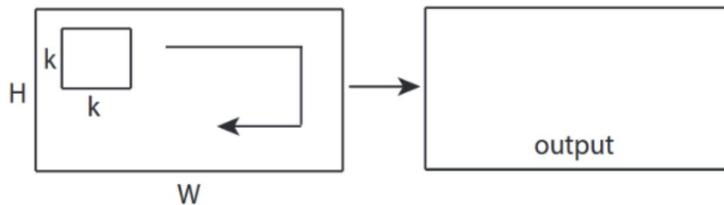MULTI-SCALE CONTEXT AGGREGATION BY DILATED CONVOLUTIONS

(a)          (b)          (c)

# 1D and 3D Convolutions

Cornell Bowers C·IS

# Convolutional Filters

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

"image"

*

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

convolutional filter

=

| 3 | 2 | 2 |
|---|---|---|
| 1 | 3 | 2 |
| 2 | 1 | 2 |

can learn this!

# Multiple Input Channels: Convolution Over Volumes
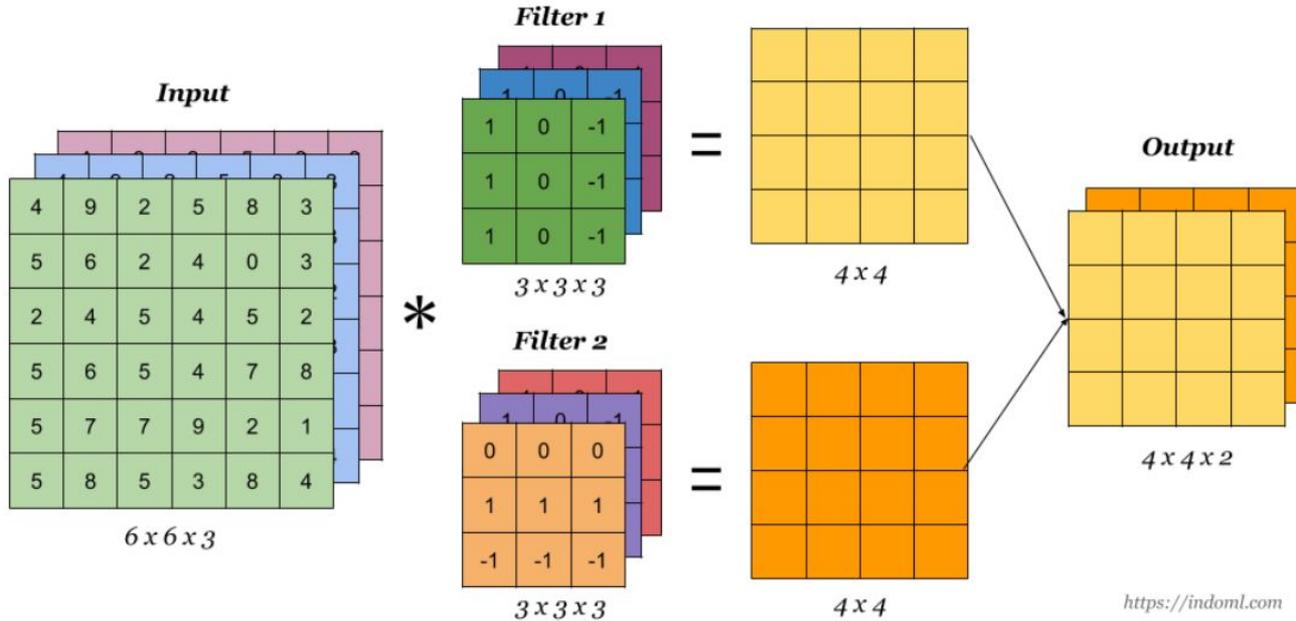
What if our input image has more than one channel?



**Input**

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

$n_H \times n_W \times n_C = 6 \times 6 \times 3$

**\***

**Filter**

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

**Parameters:**

Size: $f = 3$
#channels: $n_C = 3$
Stride: $s = 1$
Padding: $p = 0$

**=**

**Result**

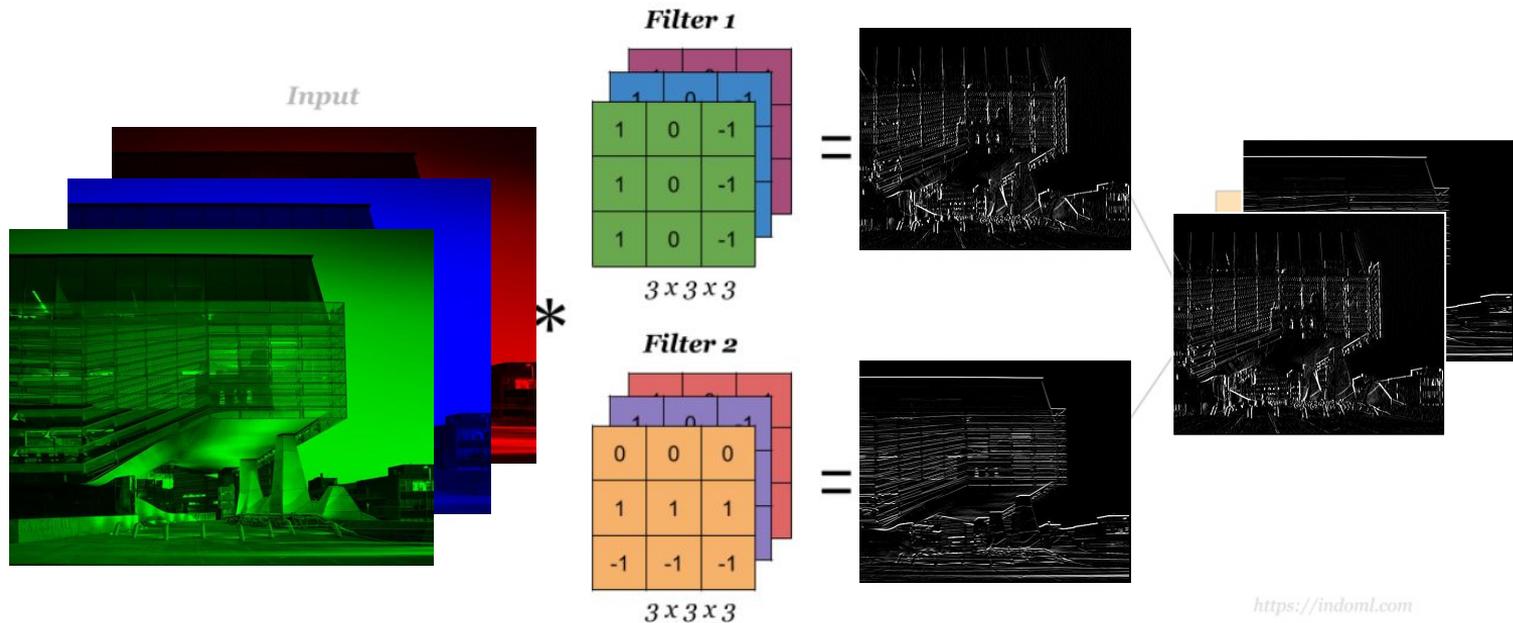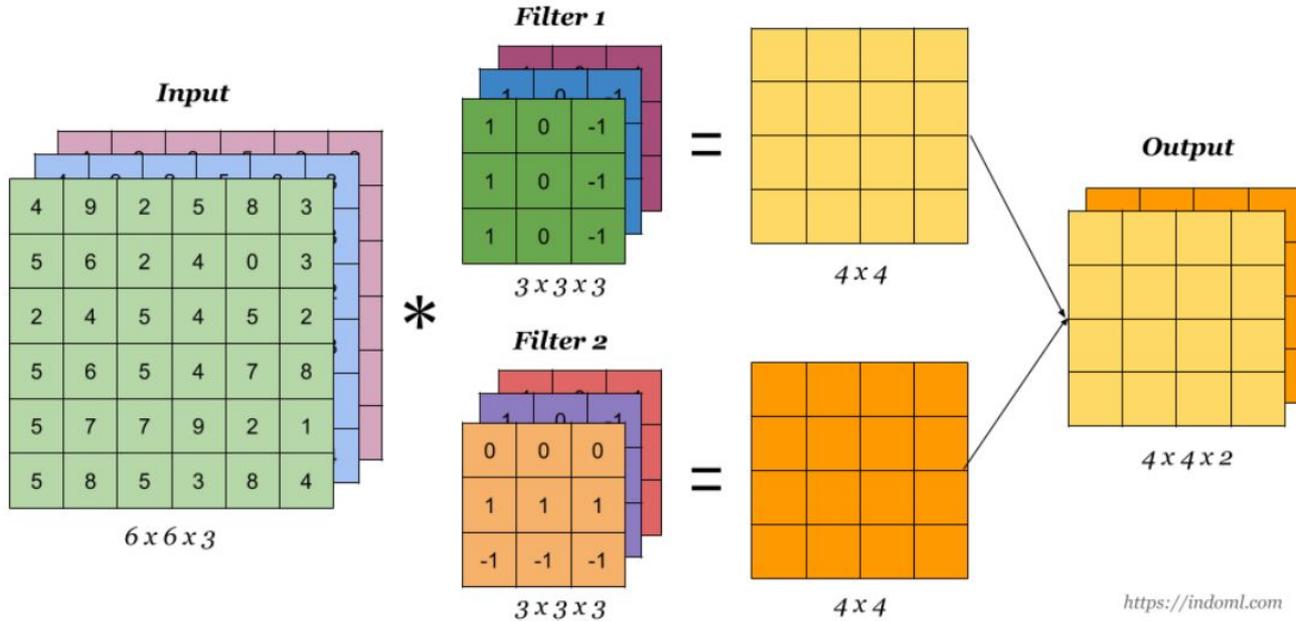| 2 |   |   |   |
|---|---|---|---|
|   |   |   |   |
|   |   |   |   |
|   |   |   |   |

https://indoml.com

# Multiple Output Channels: Multiple Filters

# Multiple Output Channels: Multiple Filters

# Multiple Output Channels: Multiple Filters

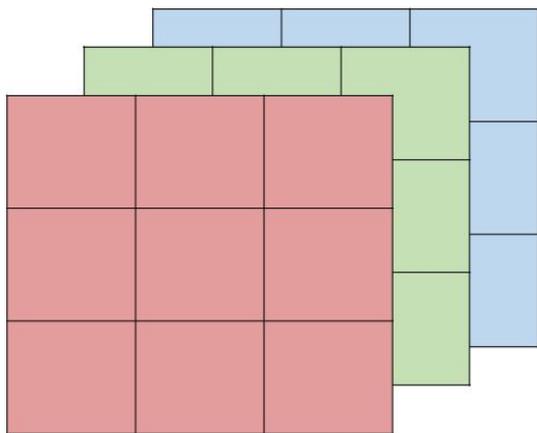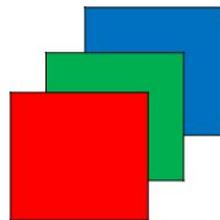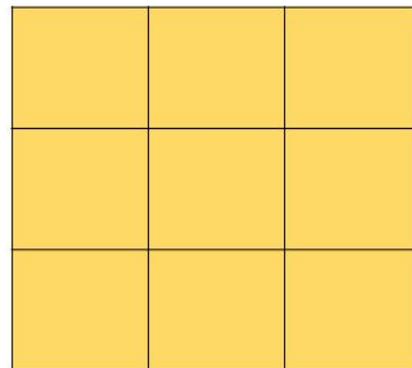# Multiple Output Channels: Multiple Filters



Input

6 x 6 x 3

*

Filter 1

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

3 x 3 x 3

=

4 x 4

Filter 2

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

3 x 3 x 3

=

4 x 4

Output

4 x 4 x 2

https://indoml.com

# Multiple Output Channels: Multiple Filters

# Multiple Output Channels: Multiple Filters

# Slight Detour: 1x1 convolutions



input      filter      output

\*    =

# Slight Detour: 1x1 convolutions

**input**        **filter**        **output**



*3x3x3*        *1x1x3*        *3x3x1*

# Slight Detour: 1x1 convolutions

# Slight Detour: 1x1 convolutions



input

filters

output

*3x3x3*

*1x1x3*

*3x3x1*

# Slight Detour: 1x1 convolutions



input

filters

output

*3x3x3*

*1x1x3*

*1x1x3*

*3x3x2*
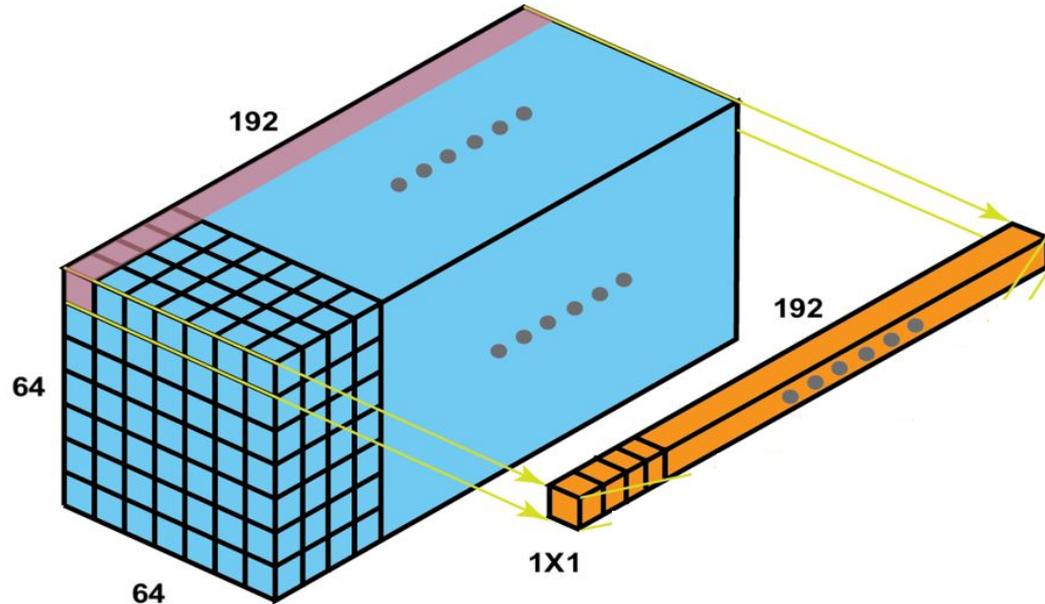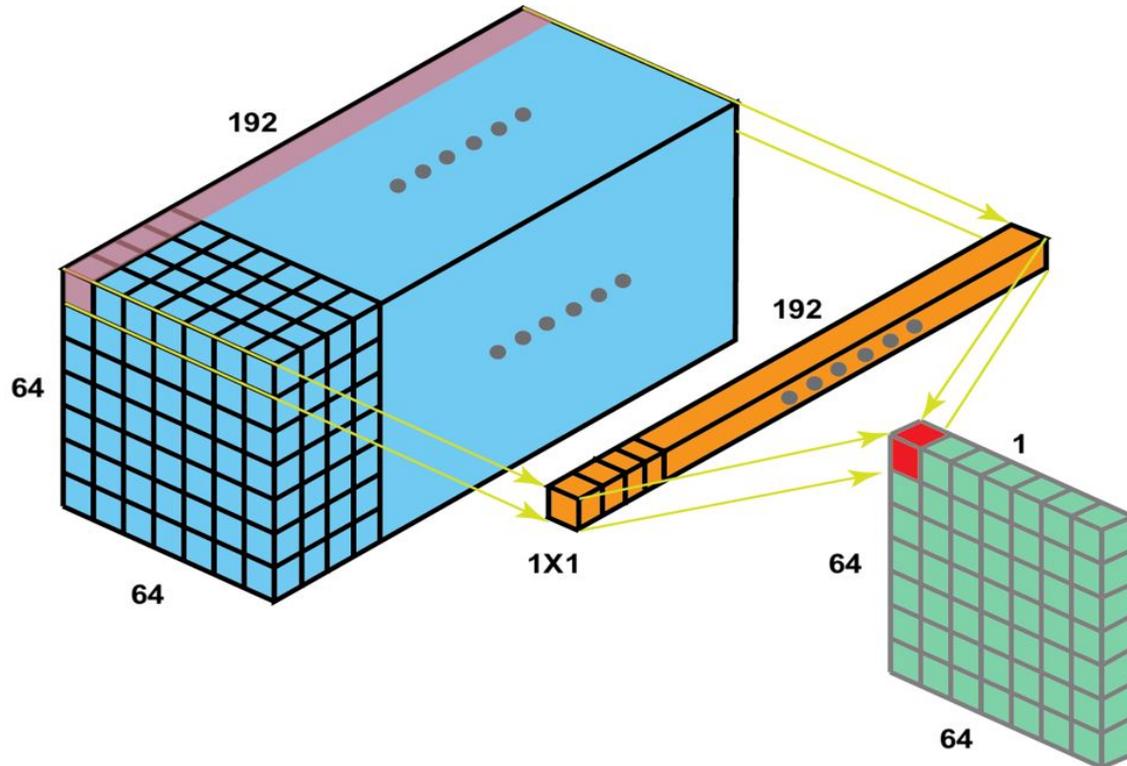
# Discuss: 1x1 Convolutions

What is the result of convolving a 64x64x192 dimensional cube with a 1x1 filter?

# 1x1 Convolutions

# Convolution Layer
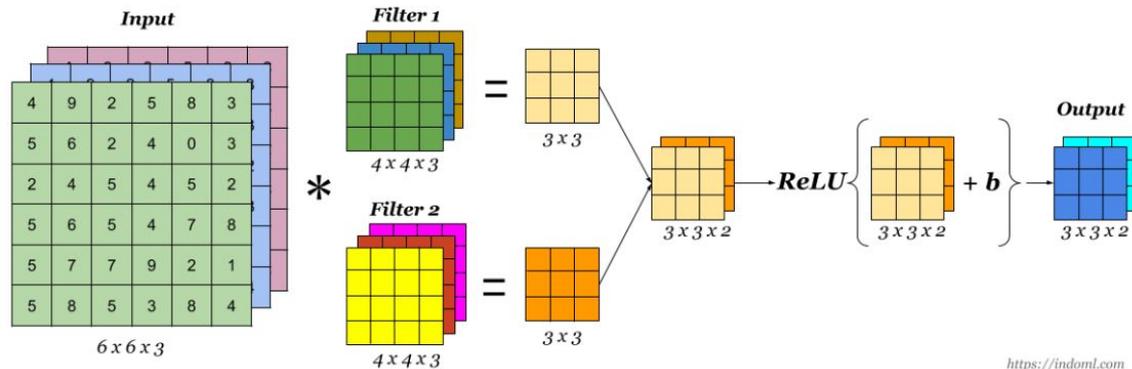
**MLP Layer**



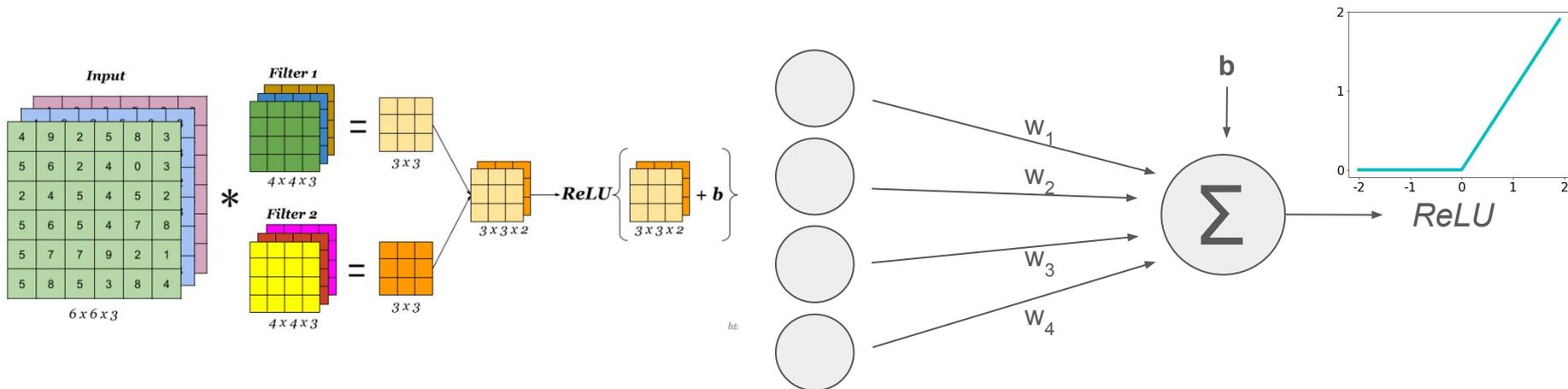**Convolution Layer**



https://indoml.com

# CNN/MLP Equivalence

Differences in a convolution layer:

- neurons are connected to a local region

- Weights are shared across multiple parameters

CONV layers can be converted to Fully connected layers and vice versa!
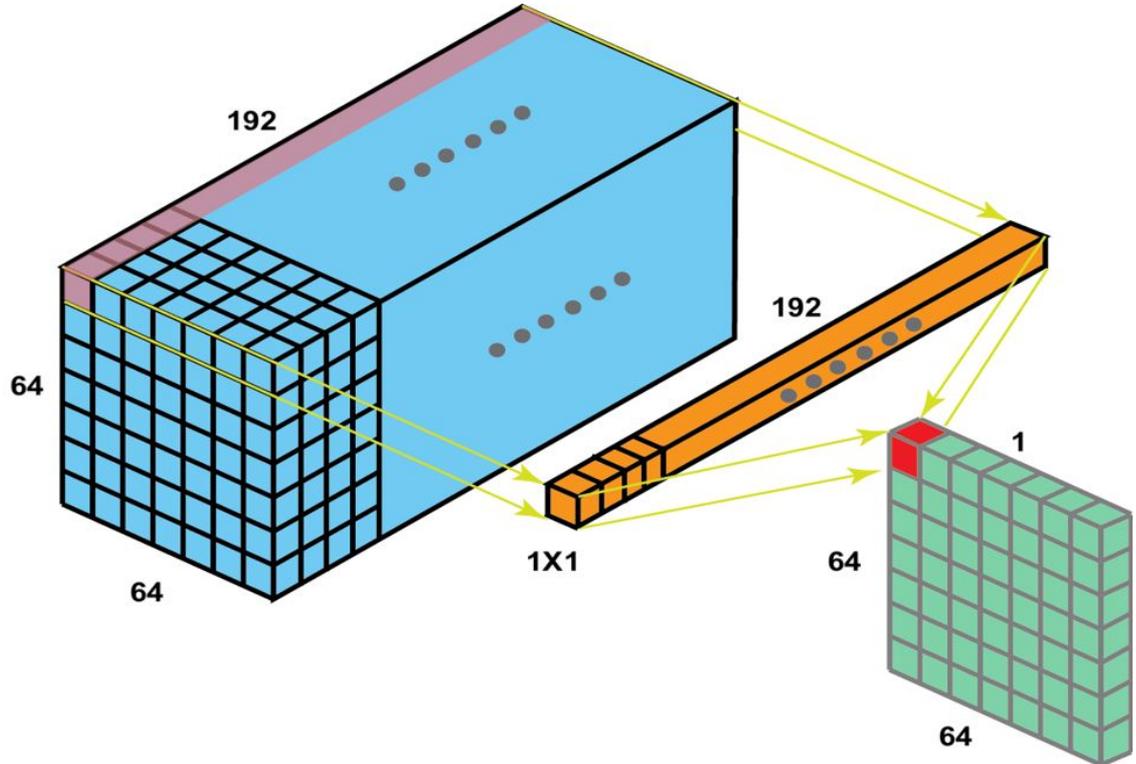
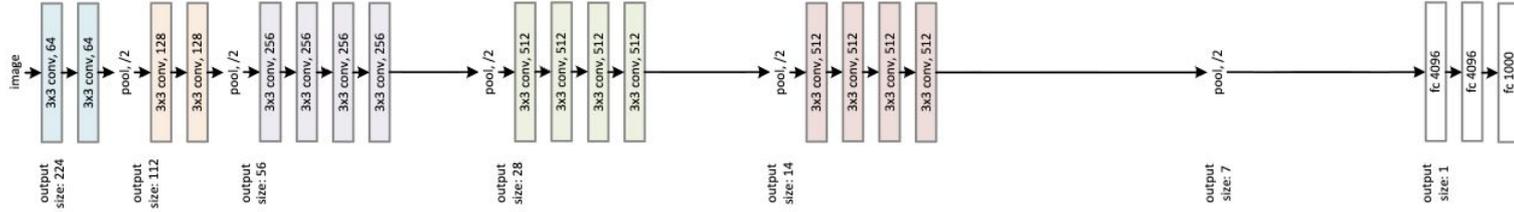# Convolutional Filters



"image"     *     convolutional filter     =

# Discuss: Trade-offs between CNNs and MLPs

How would this image change if you used an MLP instead of a 1x1 convolution filter to produce a (64x64x1) feature map? Hint: think about parameter counts and feature interactions.
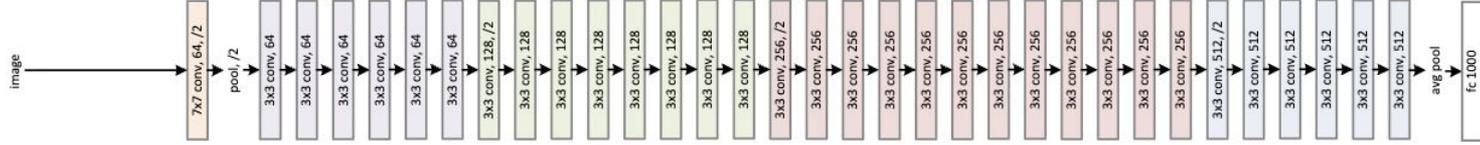
# Our goal today

# CNN Layer Output Visualization

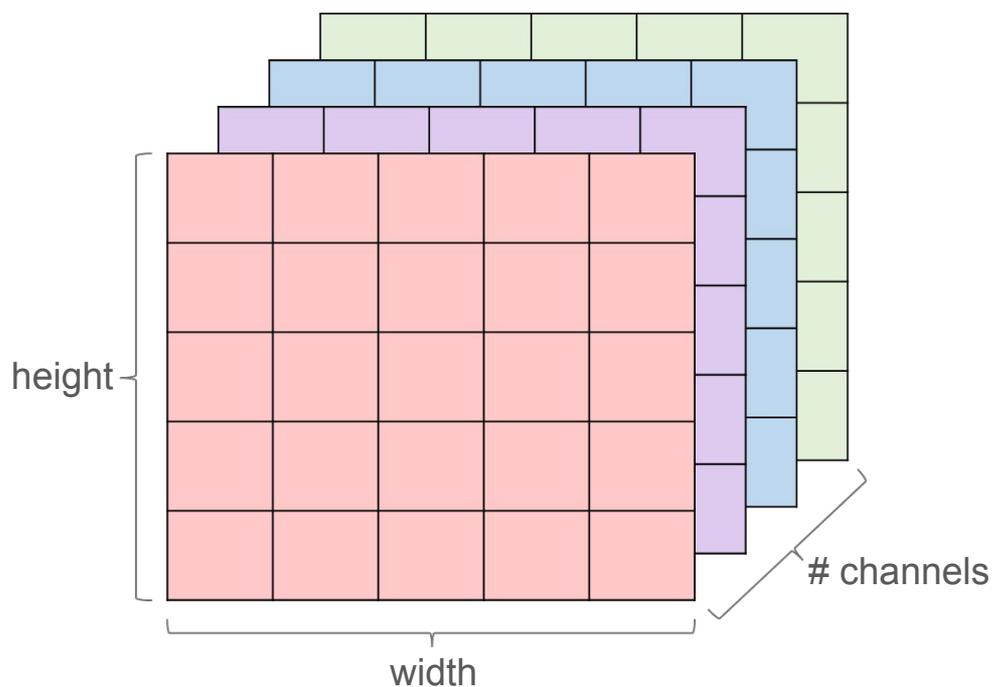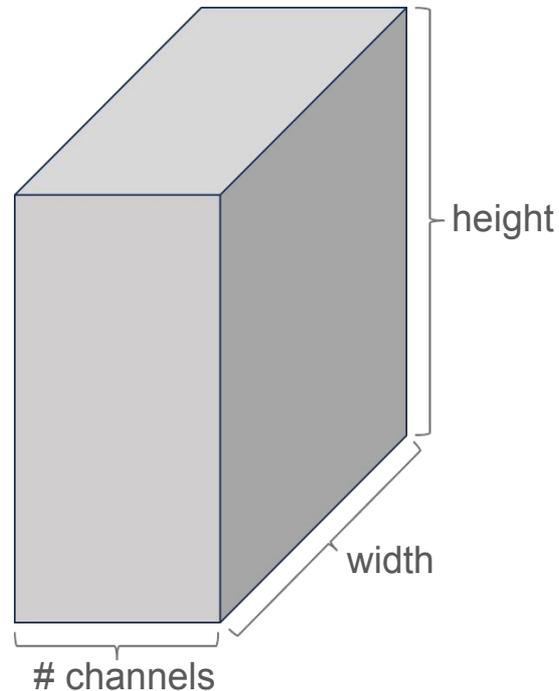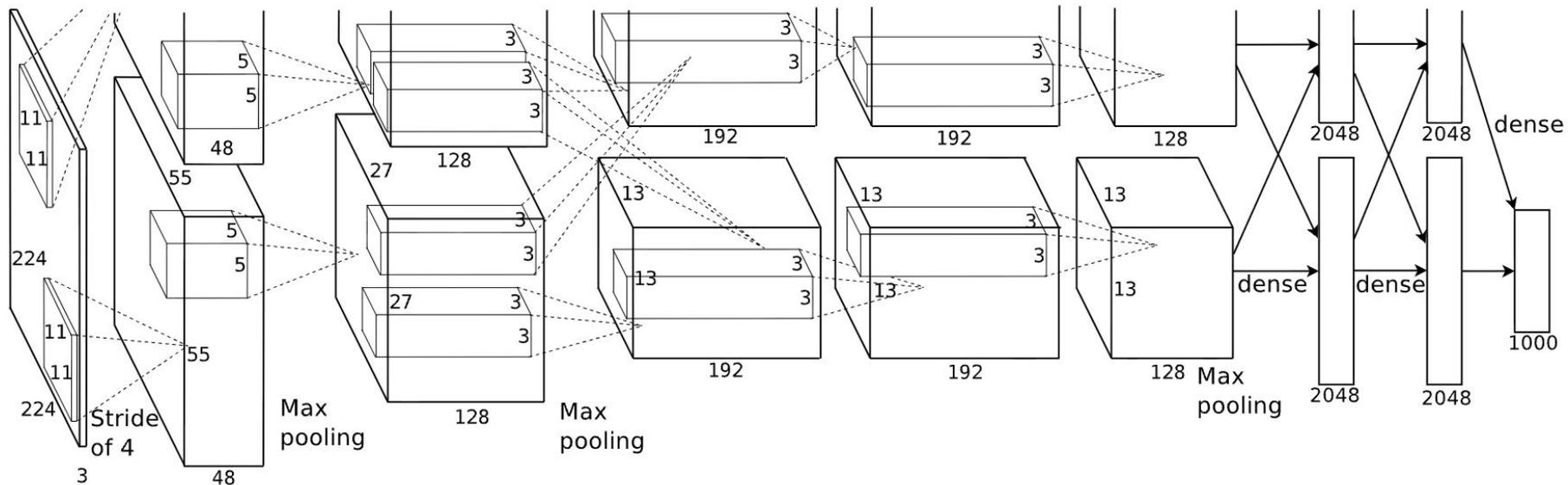# Our goal today
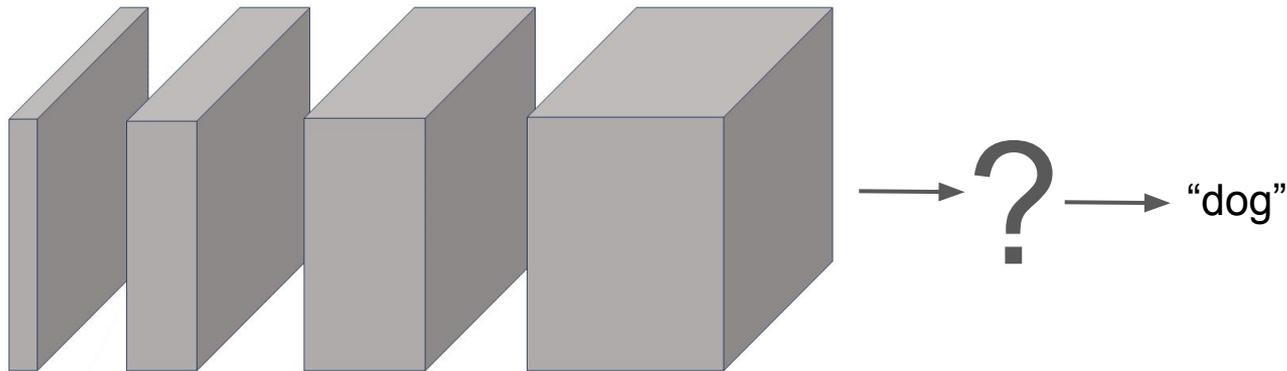
# Convolutional Neural Networks (CNNs)

✅ **Convolutions**          Maintain spatial relation between pixels

Reduce number of parameters through weight sharing



input image

"dog"

# Ensuring translational invariance

# Max Pooling

# CNNs - Pooling

# CNNs - Pooling

❖ Down sample feature maps that highlight the most present feature in the patch

❖ Improve efficiency by reducing computations with downsampling

❖ Increase receptive field size



**Max Pooling**

| 4 | 9 | 2 | 5 |
|---|---|---|---|
| 5 | 6 | 2 | 4 |
| 2 | 4 | 5 | 4 |
| 5 | 6 | 8 | 4 |

→

| 9 | 5 |
|---|---|
| 6 | 8 |

**Avg Pooling**

| 4 | 9 | 2 | 5 |
|---|---|---|---|
| 5 | 6 | 2 | 4 |
| 2 | 4 | 5 | 4 |
| 5 | 6 | 8 | 4 |

→

| 6.0 | 3.3 |
|-----|-----|
| 4.3 | 5.3 |

https://indoml.com

# Convolutional Neural Networks (CNNs)

✅ Convolutions   Maintain spatial relation between pixels
Reduce number of parameters through weight sharing
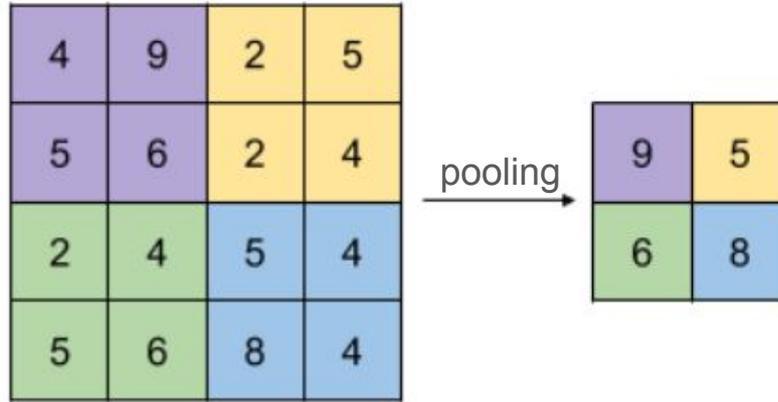
✅ **Pooling**   Captures key information from across different areas of the feature maps
Together with convolutions allows for translational invariance

pooling

input image

"dog"

# Normalization

❖ Normalize channels to mean 0 and variance 1 across each training batch

❖ Increases speed of training by enabling the use of larger learning rates
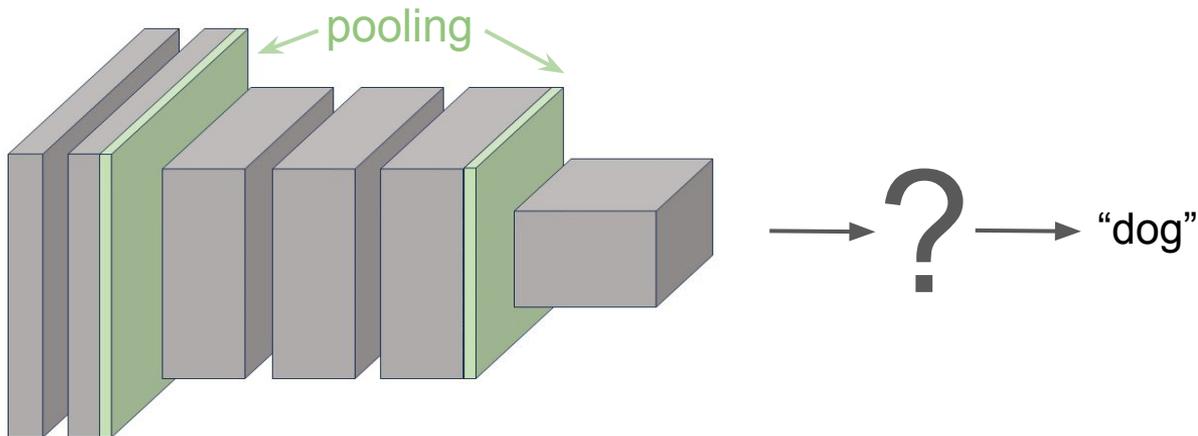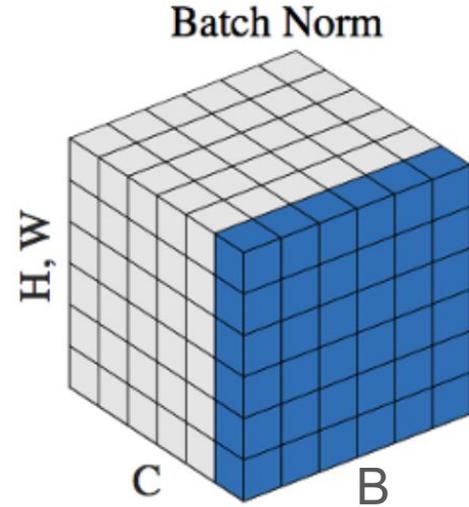
❖ Improves stability of training



Batch Norm

## The Batch Normalization Algorithm

**Input:** Values of $x$ over a mini-batch: $\mathcal{B} = \{x_{1...m}\}$;
Parameters to be learned: $\gamma, \beta$

**Output:** $\{y_i = \text{BN}_{\gamma,\beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^{m} x_i \qquad\qquad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^{m} (x_i - \mu_{\mathcal{B}})^2 \qquad\qquad // \text{ mini-batch variance}$$
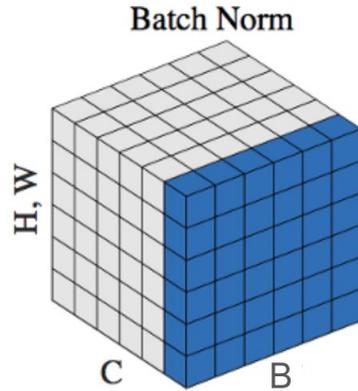
$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad\qquad // \text{ normalize}$$

$$y_i \leftarrow \gamma \widehat{x}_i + \beta \equiv \text{BN}_{\gamma,\beta}(x_i) \qquad\qquad // \text{ scale and shift}$$

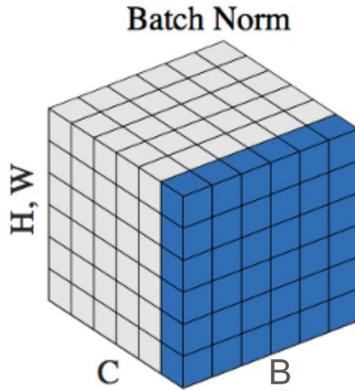**Algorithm 1:** Batch Normalizing Transform, applied to activation $x$ over a mini-batch.

# Discuss!

What is the dimension of the mean when you compute the batch norm of a volume of dimension (b x c x h x w)?



Batch Norm

# Discuss!

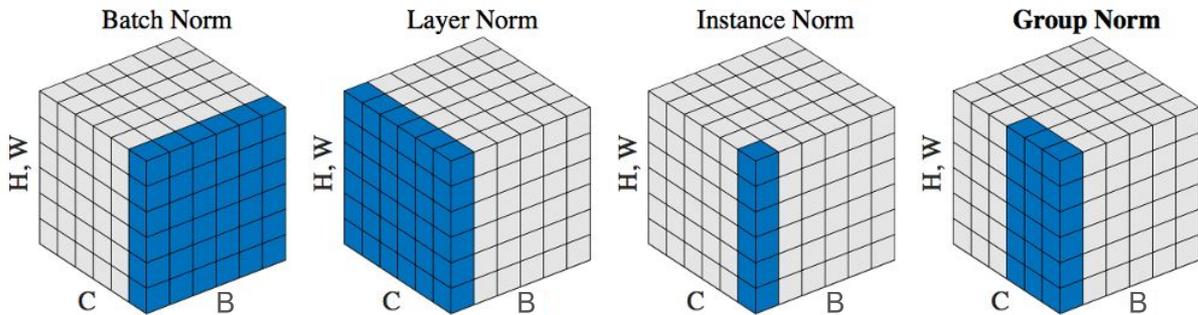What is the dimension of the mean when you compute the batch norm of a volume of dimension (b x c x h x w)?



Batch Norm

H, W

C    B

**The mean should be (1 x c x 1 x 1)!**

# Normalization Layers

- Normalization layers improve training stability
- Can train with larger learning rates
  - Faster training
- A large learning rate acts as an implicit regularizer
  - Better generalization
- Normalization can also be applied across different dimensions for different use cases

# Convolutional Neural Networks (CNNs)

✅ Convolutions        Maintain spatial relation between pixels
Reduce number of parameters through weight sharing

✅ Pooling        Captures key information from across different areas of the feature maps
Together with convolutions allows for translational invariance

✅ **BatchNorm**        Increases speed and stability of training

BatchNorm Layer



input image

→ ? → "dog"

# Convolutional Neural Networks (CNNs)

# Convolutional Neural Networks (CNNs)

# Convolutional Neural Networks (CNNs)

# Convolutional Neural Networks (CNNs)

✅ Convolutions      Maintain spatial relation between pixels
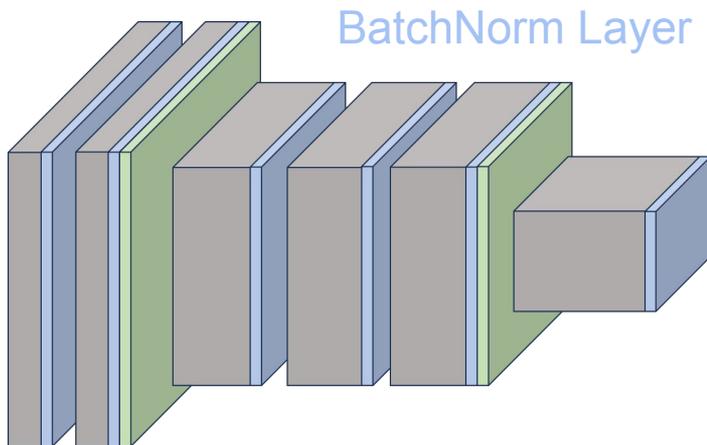Reduce number of parameters through weight sharing

✅ Pooling      Captures key information from across different areas of the feature maps
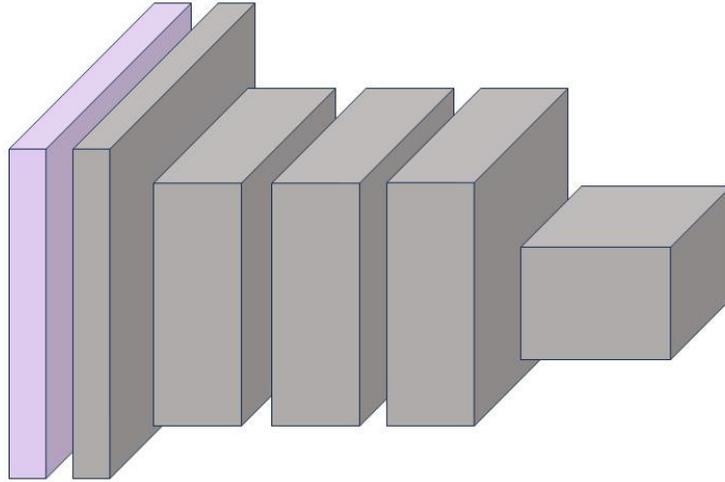Together with convolutions allows for translational invariance

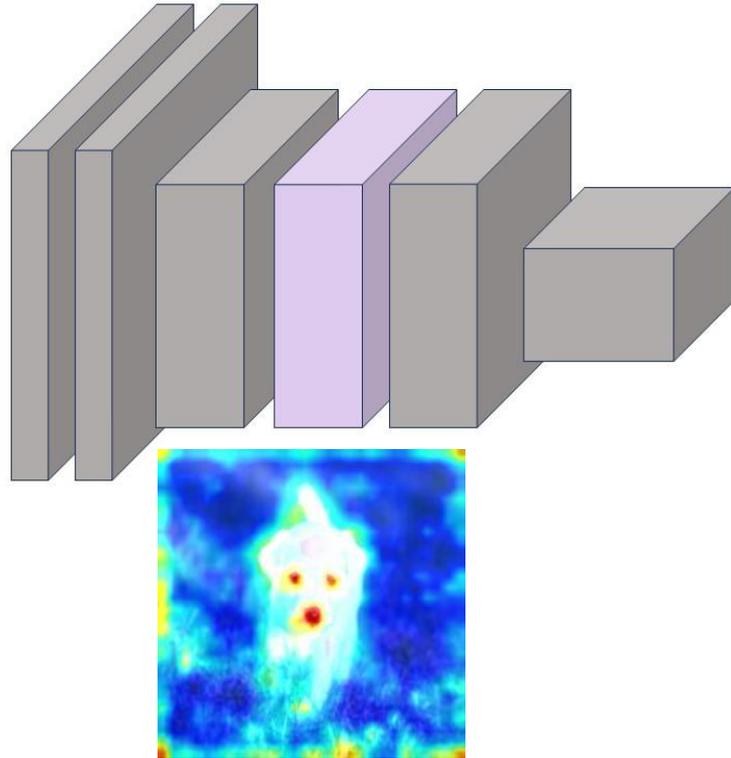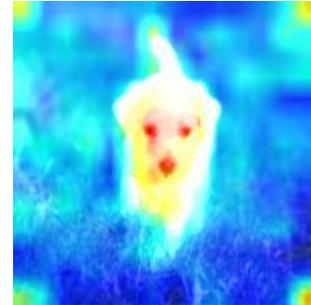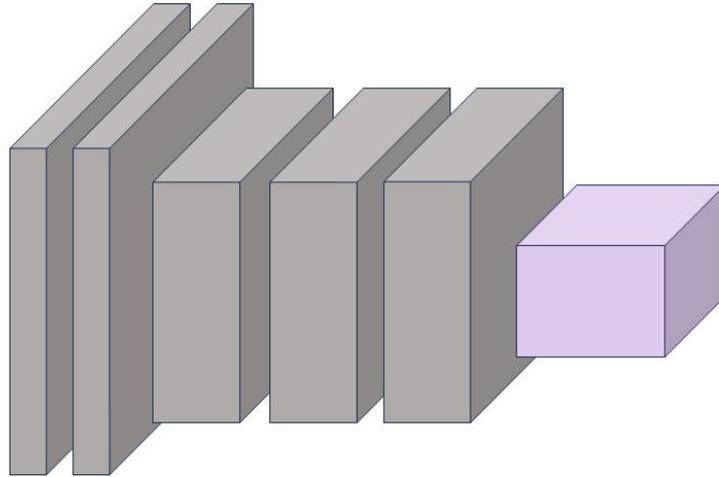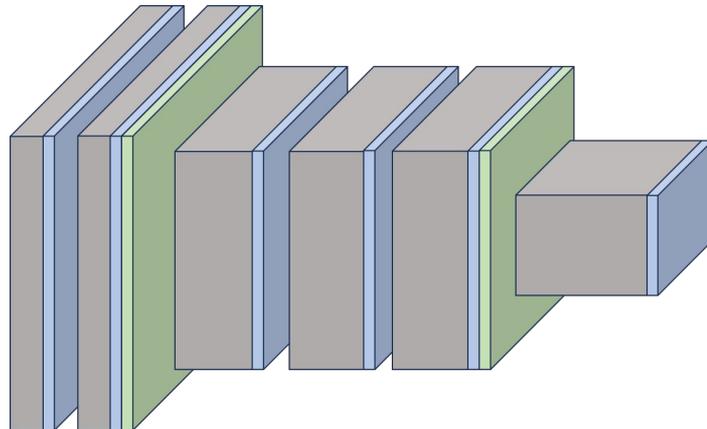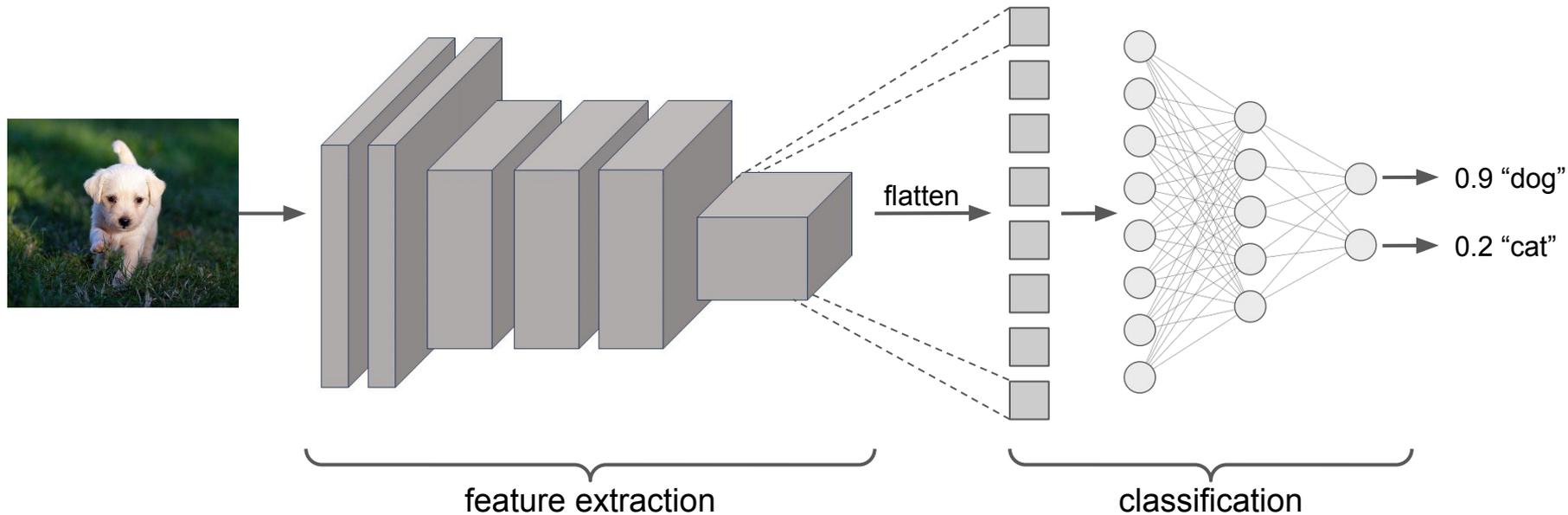✅ BatchNorm      Increases speed and stability of training



input image

→ **?** → "dog"

Image Classification

# Practical Guide

- Input image dimensions is divisible by 2

- Small conv filters (3x3 or 5x5)

- Zero padding is used to maintain spatial resolution

- Max pooling for downsampling

- Pooling layers have a receptive field of 2 and stride of 2

# Summary

- CNNs are primarily designed to process and analyze visual data, such as images and videos.

- Key components: convolution layers, pooling layers, activation functions, normalization layers

- Advantages:
  - Translational Invariance
  - Parameter sharing
  - Feature learning

- Can be trained with backprop

- Used for tasks such as segmentation, classification, object detection, etc.