**Cornell Bowers C·IS**
College of Computing and Information Science

# Convolutional Neural Networks

CS4782: Intro to Deep Learning

---

## Logistics

- HW1 has been released
  - Due next Thursday (February 15)
- Office hours are listed on the course website
- Homework clarifications are listed as pinned posts under HW1 on Ed
- Post questions on Ed

---

### The Batch Normalization Algorithm

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

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

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

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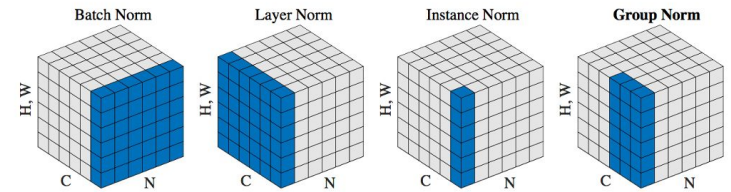
$$\widehat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \qquad \text{// normalize}$$

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

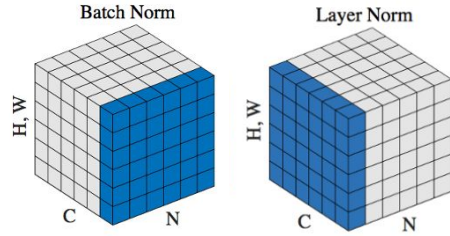**Algorithm 1:** Batch Normalizing Transform, applied to activation $x$ over a mini-batch.
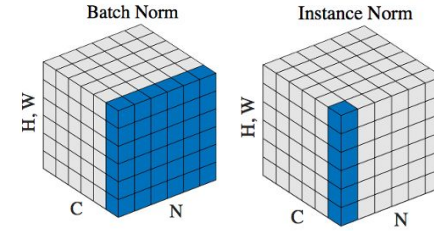
---

## Many Kinds of Normalization Layers



Normalization Methods

"Group Normalization" by Wu et al., 2018

## Layer Normalization

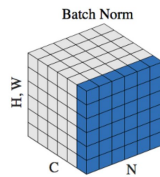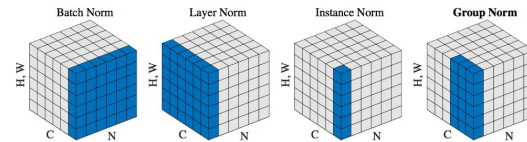## Instance Normalization

## 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)?
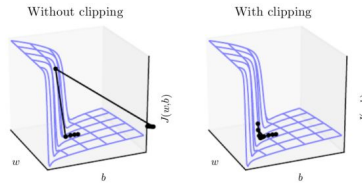
## Normalization Layers

- Normalization layers improve training stability
- Can train with larger learning rates
  - Faster training
- A large learning rate acts as an implicit regularizer
  - Better generalization

**Cornell Bowers C·IS**

## Gradient Clipping

- Exploding gradients result in unstable training
- Optimization is hard when you have very large gradients
- Fixes:
  - Clip by value
  - Clip by norm

Without clipping

With clipping

---

**Cornell Bowers C·IS**

## So far…

- MLPs learn complex decision boundaries
- Optimization algorithms use the gradient of the loss to find network parameters
- Different training strategies like regularization, early stopping and normalization can improve training and generalization
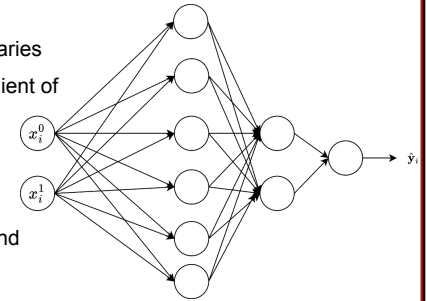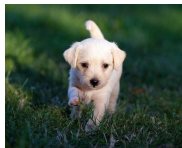
$x_i^0$

$x_i^1$

$\hat{y}_i$

---

**Cornell Bowers C·IS**

## Image Classification

input image

classification → "dog"

input image

classification → "cat"

---

**Cornell Bowers C·IS**
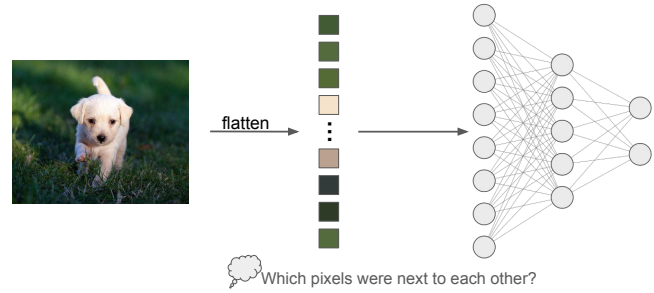
## Applications in Medicine

Applications in Autonomous Driving

Why not use a Multi-Layer Perceptron?



flatten

Which pixels were next to each other?
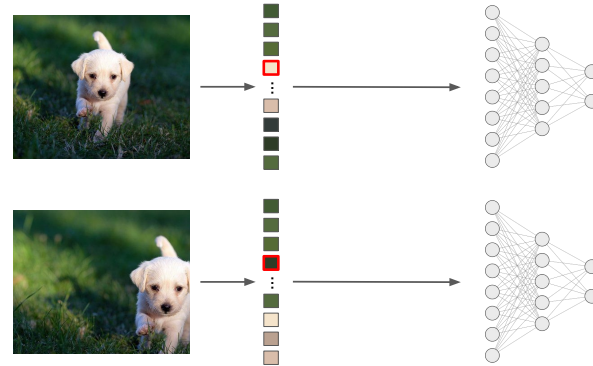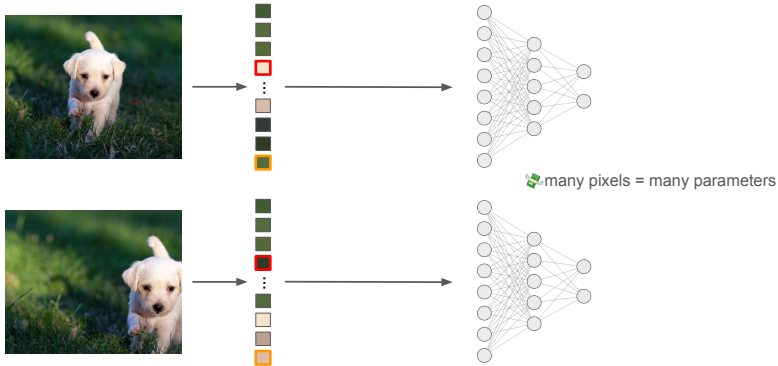
Why not use a Multi-Layer Perceptron?

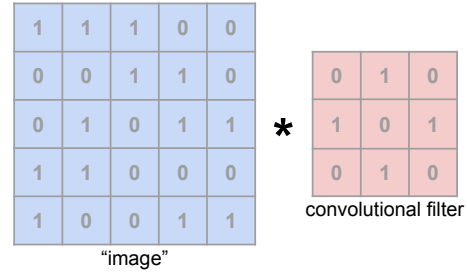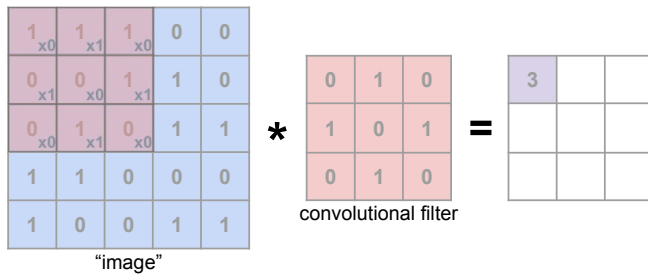Why not use a Multi-Layer Perceptron?

Why not use a Multi-Layer Perceptron?

🐌 many pixels = many parameters
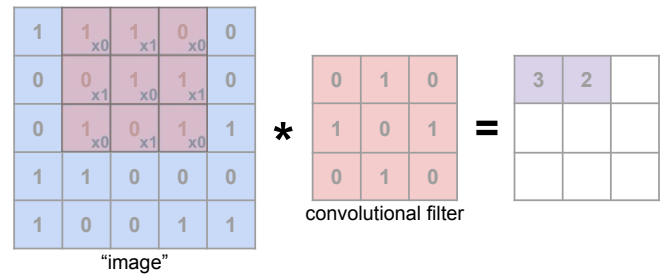
Convolutional Filters

| 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

Convolutional Filters

| $1_{x0}$ | $1_{x1}$ | $1_{x0}$ | 0 | 0 |
|---|---|---|---|---|
| $0_{x1}$ | $0_{x0}$ | $1_{x1}$ | 1 | 0 |
| $0_{x0}$ | $1_{x1}$ | $0_{x0}$ | 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 |  |  |
|---|---|---|
|  |  |  |
|  |  |  |

Convolutional Filters

| 1 | $1_{x0}$ | $1_{x1}$ | $0_{x0}$ | 0 |
|---|---|---|---|---|
| 0 | $0_{x1}$ | $1_{x0}$ | $1_{x1}$ | 0 |
| 0 | $1_{x0}$ | $0_{x1}$ | $1_{x0}$ | 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 |  |
|---|---|---|
|  |  |  |
|  |  |  |

**Cornell Bowers CIS**

## Convolutional Filters

$$
\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0_{\times 0} & 1_{\times 1} & 0_{\times 0} & 1 & 1 \\ 1_{\times 1} & 1_{\times 0} & 0_{\times 1} & 0 & 0 \\ 1_{\times 0} & 0_{\times 1} & 0_{\times 0} & 1 & 1 \end{bmatrix}
* \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}
= \begin{bmatrix} 3 & 2 & 2 \\ 1 & 3 & 2 \\ 2 & & \end{bmatrix}
$$

"image"  convolutional filter

---

**Cornell Bowers CIS**

## Convolutional Filters

$$
\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1_{\times 0} & 0_{\times 1} & 1_{\times 0} & 1 \\ 1 & 1_{\times 1} & 0_{\times 0} & 0_{\times 1} & 0 \\ 1 & 0_{\times 0} & 0_{\times 1} & 1_{\times 0} & 1 \end{bmatrix}
* \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}
= \begin{bmatrix} 3 & 2 & 2 \\ 1 & 3 & 2 \\ 2 & 1 & \end{bmatrix}
$$

"image"  convolutional filter

---

**Cornell Bowers CIS**

## Convolutional Filters

$$
\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0_{\times 0} & 1_{\times 1} & 1_{\times 0} \\ 1 & 1 & 0_{\times 1} & 0_{\times 0} & 0_{\times 1} \\ 1 & 0 & 0_{\times 0} & 1_{\times 1} & 1_{\times 0} \end{bmatrix}
* \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}
= \begin{bmatrix} 3 & 2 & 2 \\ 1 & 3 & 2 \\ 2 & 1 & 2 \end{bmatrix}
$$

"image"  convolutional filter

---

**Cornell Bowers CIS**

## Convolutional Filters

$$
\begin{bmatrix} 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 0 & 1 & 1 \end{bmatrix}
* \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 1 \\ 0 & 1 & 0 \end{bmatrix}
= \begin{bmatrix} 3 & 2 & 2 \\ 1 & 3 & 2 \\ 2 & 1 & 2 \end{bmatrix}
$$

"image"  convolutional filter

can learn this!

**Slide 1**

## 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 |

**Slide 2**

## Discuss with your Neighbor!

Match the following convolutional filters with the output they produce.



input image

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

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

| 1/9 | 1/9 | 1/9 |
|---|---|---|
| 1/9 | 1/9 | 1/9 |
| 1/9 | 1/9 | 1/9 |

**Slide 3**

## Dilated Convolutions



https://towardsdatascience.com/review-dilated-convolution-semantic-segmentation-9d5a5bd768f5

**Slide 4**

## 1D and 3D Convolutions



https://wandb.ai/ayush-thakur/dl-question-bank/reports/Intuitive-understanding-of-1D-2D-and-3D-convolutions-in-convolutional-neural-networks---VmlldzoxOTk2MDA
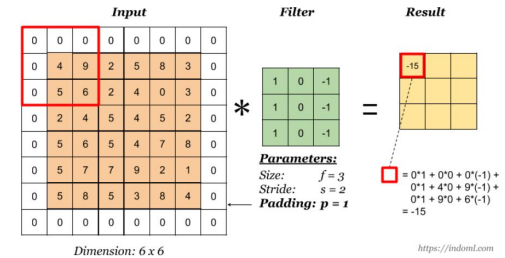
## CNNs - Stride

❖ 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

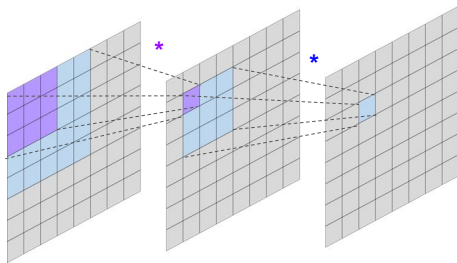❖ The receptive fields overlap less as the stride becomes larger



*Filter with stride (s) = 2*

## CNNs - Padding

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

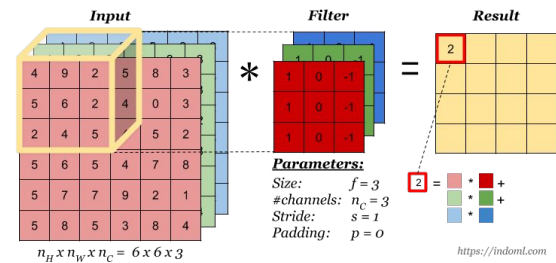❖ Prevents image shrinking and loss of information from image boundary
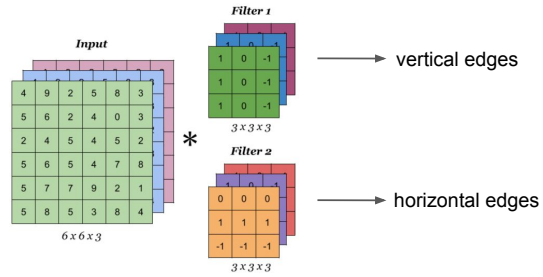
## Stacking Convolutions



❖ Size of receptive field increases with each layer

❖ Capture more complex features

## Convolution Over Volumes

What if our input image has more than one channel?

## Convolution Operation with Multiple Filters
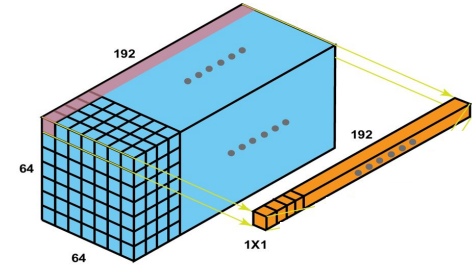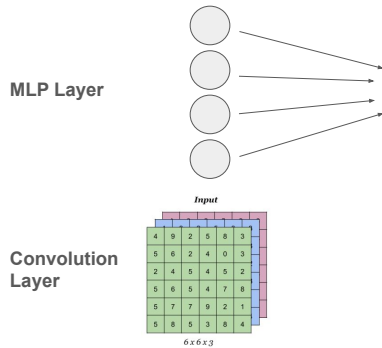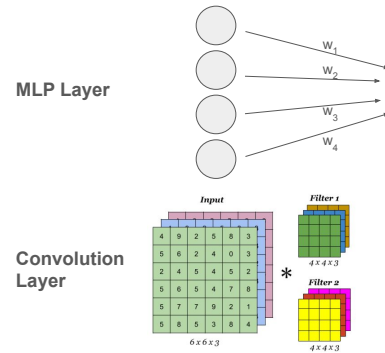
## Discuss: 1x1 Convolutions

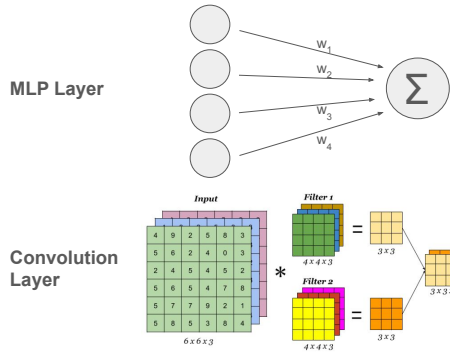What is the result of convolving a 64x64x192 dimensional cube with a 1x1 filter?

## Convolution Layer

**MLP Layer**

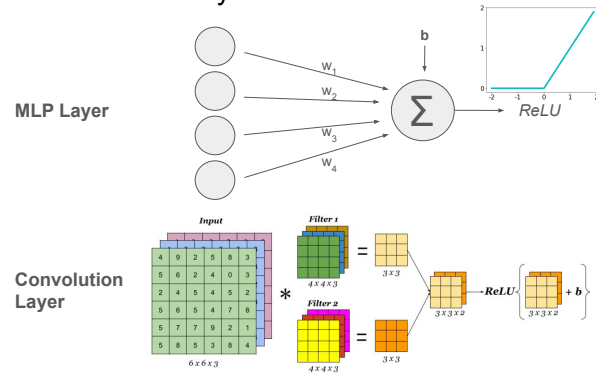**Convolution Layer**

## Convolution Layer

**MLP Layer**

**Convolution Layer**

# Slide 1

## Convolution Layer

**MLP Layer**

$w_1$ $w_2$ $w_3$ $w_4$ $\sum$

**Convolution Layer**

Input $6 \times 6 \times 3$ * Filter 1 $4 \times 4 \times 3$ = $3 \times 3$ Filter 2 $4 \times 4 \times 3$ = $3 \times 3 \times 2$



# Slide 2

## Convolution Layer

**MLP Layer**

**b**

$w_1$ $w_2$ $w_3$ $w_4$ $\sum$ ReLU

**Convolution Layer**

Input $6 \times 6 \times 3$ * Filter 1 $4 \times 4 \times 3$ = $3 \times 3$ Filter 2 $4 \times 4 \times 3$ = $3 \times 3$ $\rightarrow$ $3 \times 3 \times 2$ ReLU $+ b$ $\{ 3 \times 3 \times 2 \}$



# Slide 3

## Convolution Layer

**MLP Layer**

**b**

$w_1$ $w_2$ $w_3$ $w_4$ $\sum$ ReLU → output ℝ

**Convolution Layer**

Input $6 \times 6 \times 3$ * Filter 1 $4 \times 4 \times 3$ = $3 \times 3$ Filter 2 $4 \times 4 \times 3$ = $3 \times 3$ → $3 \times 3 \times 2$ ReLU $+ b$ $\{ 3 \times 3 \times 2 \}$ Output $3 \times 3 \times 2$
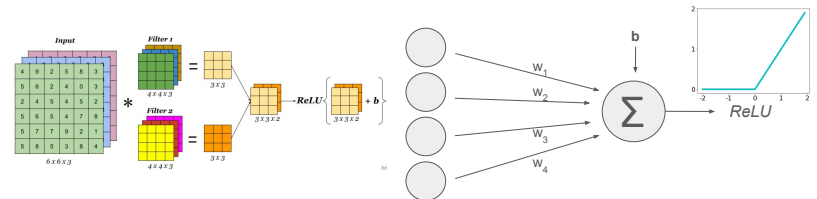
https://indoml.com



# Slide 4

## CNN/MLP Equivalence

Differences in a convolution layer:

- neurons are connected to a local region
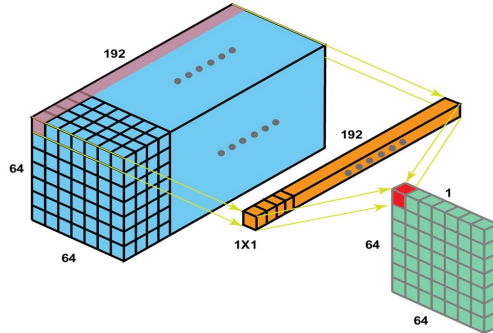- Weights are shared across multiple parameters

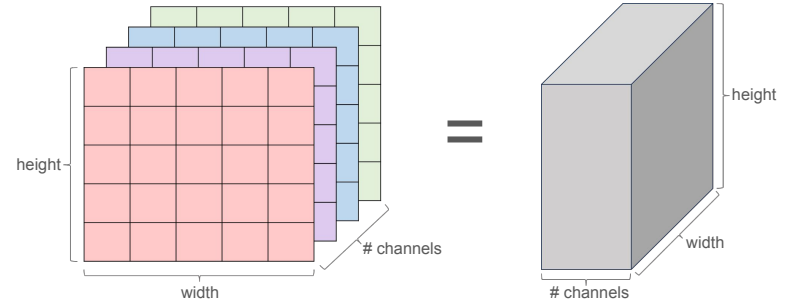CONV layers can be converted to Fully connected layers and vice versa!

## 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.
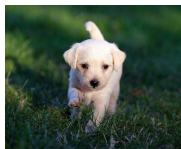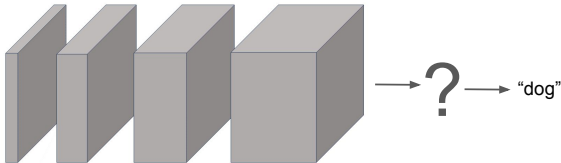
## CNN Layer Output Visualization

## Convolutional Neural Networks (CNNs)

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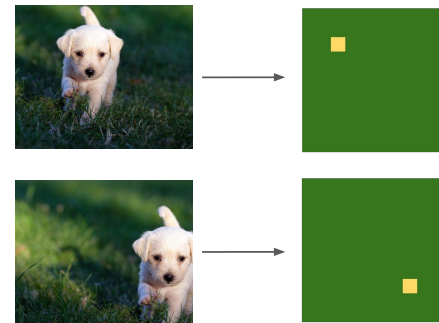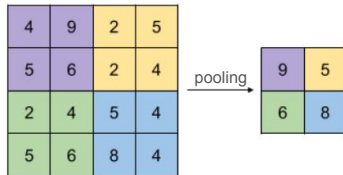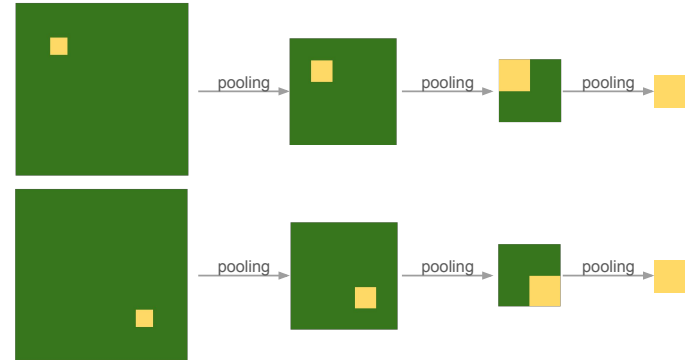Reduce number of parameters through weight sharing

? → "dog"

input image

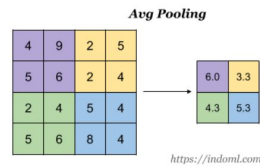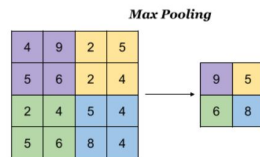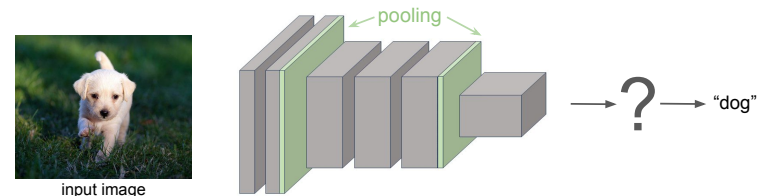## Ensuring translational invariance

## 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
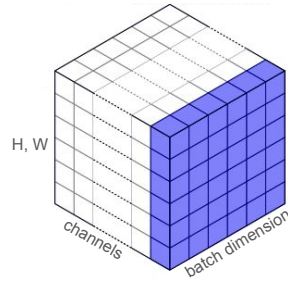
❖ Increase receptive field size

## 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
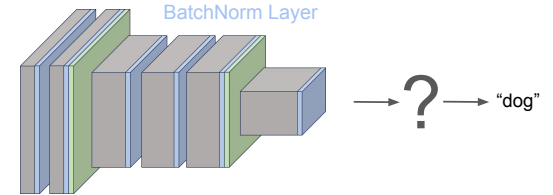


input image

## Review - Batch 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
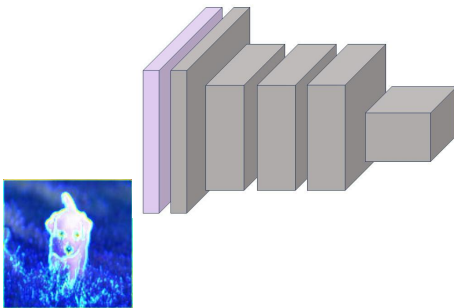
❖ Improves stability of training

H, W

channels

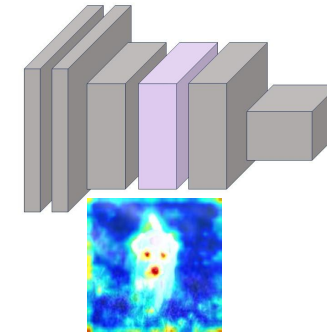batch dimension

## 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

? → "dog"

input image

## Convolutional Neural Networks (CNNs)

## Convolutional Neural Networks (CNNs)

**Cornell Bowers C·IS**

## Convolutional Neural Networks (CNNs)



---

**Cornell Bowers C·IS**

## 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
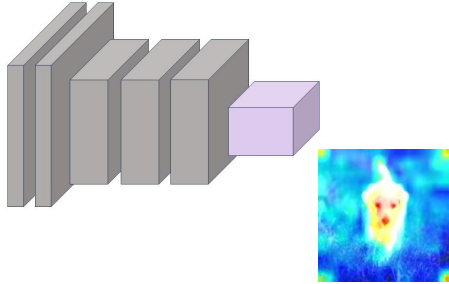
✅ BatchNorm      Increases speed and stability of training



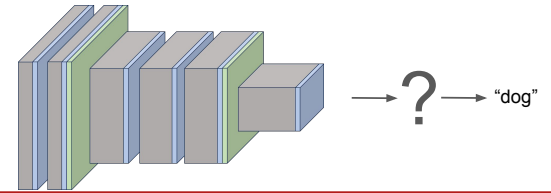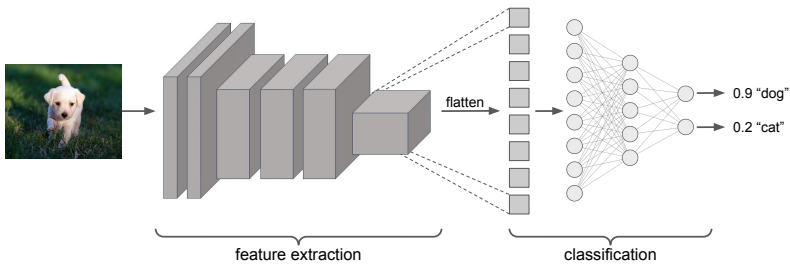input image      **?** → "dog"

---

**Cornell Bowers C·IS**

## Image Classification



flatten → 0.9 "dog" / 0.2 "cat"

feature extraction      classification

---

**Cornell Bowers C·IS**

## 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