# Convolutional Neural Networks

*PROF LIM KWAN HUI*

### 50.021 Artificial Intelligence

*The following notes are compiled from various sources such as textbooks, lecture materials, Web resources and are shared for academic purposes only, intended for use by students registered for a specific course. In the interest of brevity, every source is not cited. The compiler of these notes gratefully acknowledges all such sources.*
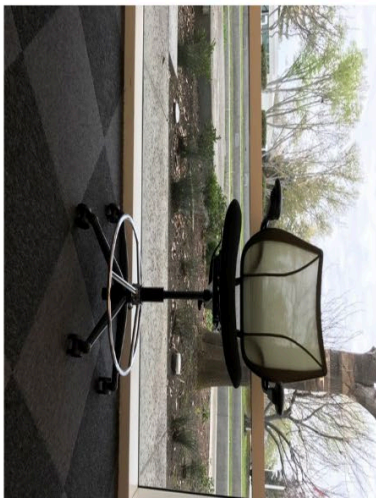
# Outline & Objectives

o Be able to use neural networks for generating word representations, e.g., Word2Vec

o Have a general understanding of how sequence models work, including RNNs and LSTMs

Last week

o Understand how convolution neural networks work

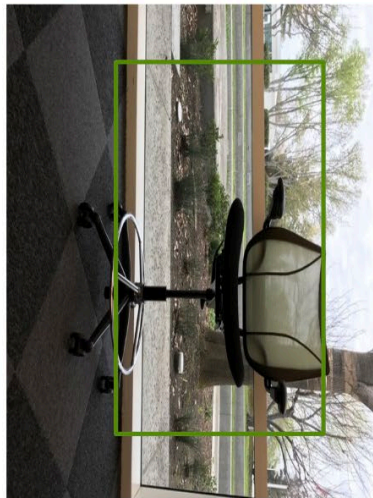o Be able to apply the various convolution-related operations in a simple example

This week

# Computer vision tasks



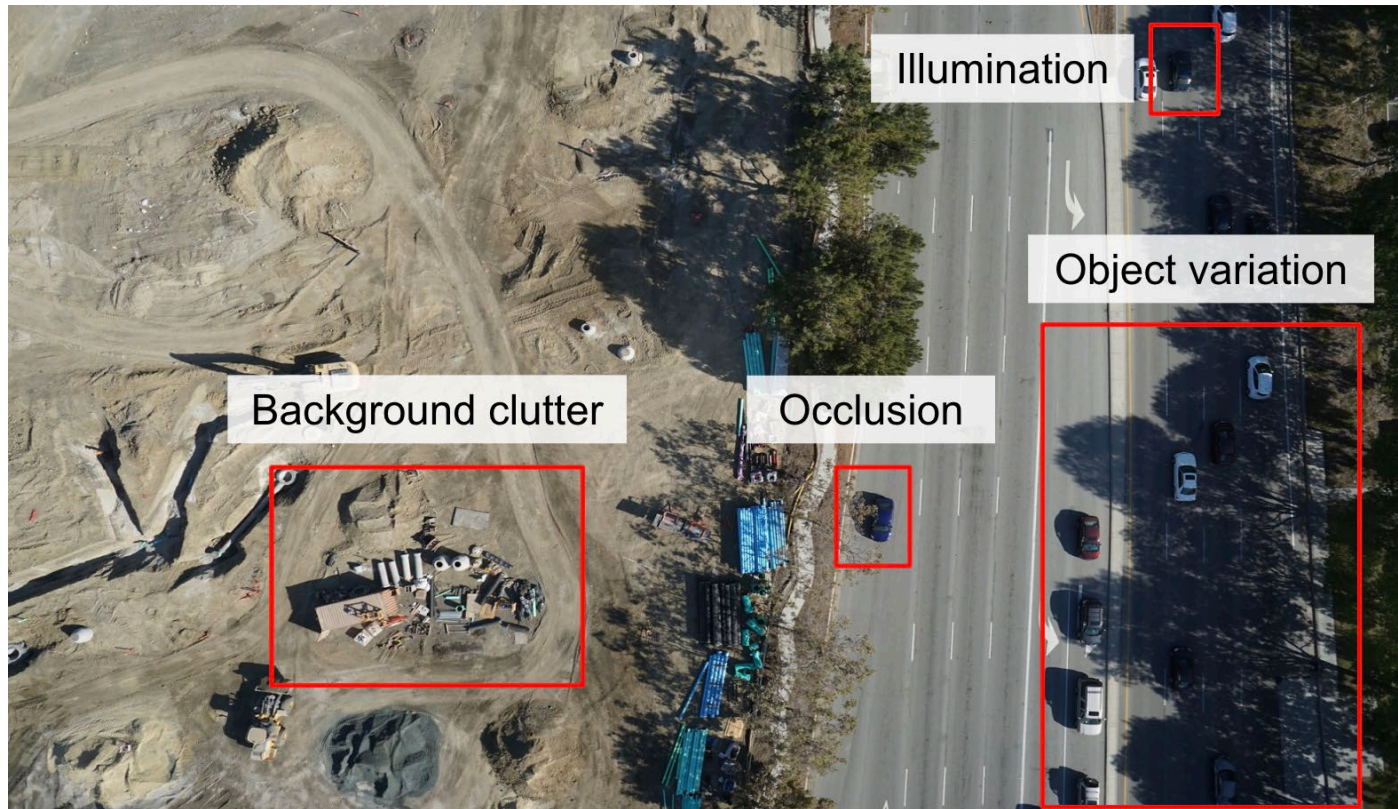**Image Classification** | **Image Classification + Localization** | **Object Detection** | **Image Segmentation**
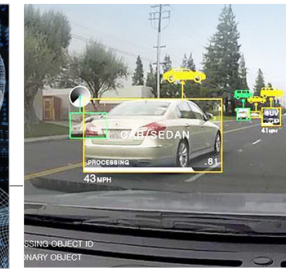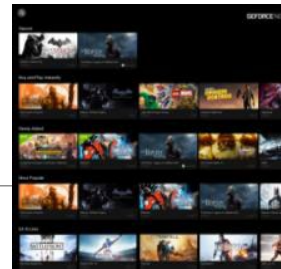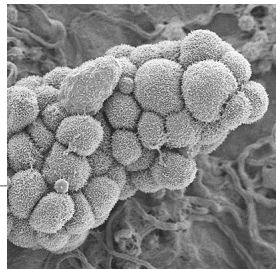
(inspired by a slide found in cs231n lecture from Stanford University)

# Challenges in images

# Deep learning with CNNs everywhere



**INTERNET & CLOUD**

Image Classification
Speech Recognition
Language Translation
Language Processing
Sentiment Analysis
Recommendation

**MEDICINE & BIOLOGY**

Cancer Cell Detection
Diabetic Grading
Drug Discovery

**MEDIA & ENTERTAINMENT**

Video Captioning
Video Search
Real Time Translation

**SECURITY & DEFENSE**

Face Detection
Video Surveillance
Satellite Imagery

**AUTONOMOUS MACHINES**

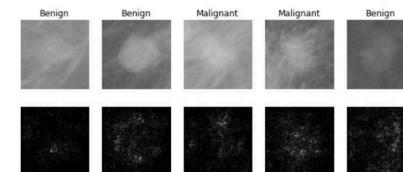Pedestrian Detection
Lane Tracking
Recognize Traffic Sign



[Toshev, Szegedy 2014]

Images are examples of pose estimation, not actually from Toshev & Szegedy 2014. Copyright Lane McIntosh.

frame: t-3   t-2   t-1   t

"submarine"

"diver"

"enemy"

"enemy+diver"

[Guo et al. 2014]

Figures copyright Xiaoxiao Guo, Satinder Singh, Honglak Lee, Richard Lewis,
and Xiaoshi Wang, 2014. Reproduced with permission.

Benign   Benign   Malignant   Malignant   Benign

[Levy et al. 2016]

Figure copyright Levy et al. 2016.
Reproduced with permission.

[Dieleman et al. 2014]

From left to right: public domain by NASA, usage permitted by
ESA/Hubble, public domain by NASA, and public domain.

[Sermanet et al. 2011]
[Ciresan et al.]

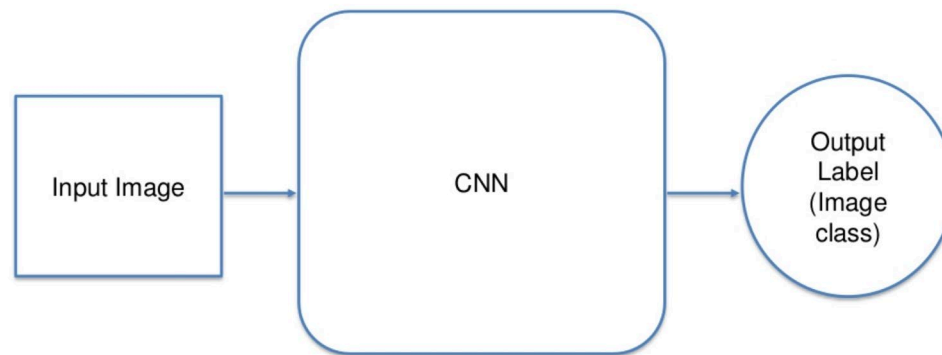Photos by Lane McIntosh.
Copyright CS231n 2017.

# Convolutional NN



o LeCun, 1989 (Chief scientist FB)

o "...are a specialized kind of neural network for processing data that has a known **grid-like** topology. *Examples* include *time-series data*, which can be thought of as a 1-D grid taking samples at regular time intervals, and *image data*, which can be thought of as a 2-D grid of pixels. Convolutional networks have been tremendously successful in practical applications. The name "convolutional neural network" indicates that the network employs a **mathematical operation called convolution**. Convolution is a specialized kind of linear operation." (Goodfellow et al., 2016)
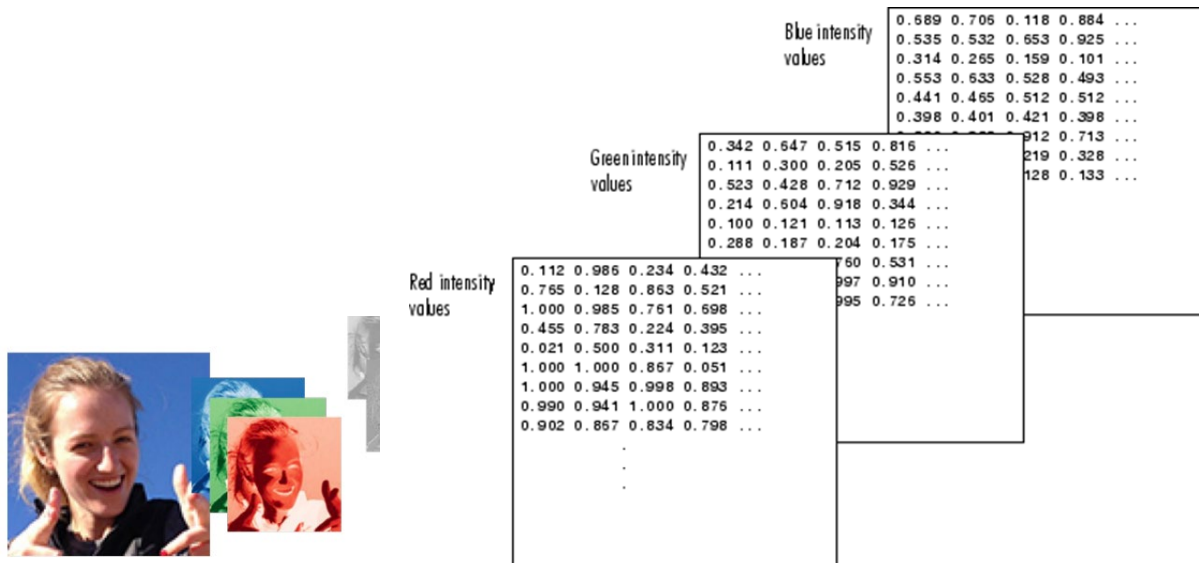
# Convolutional neural networks

o Neural networks that use convolution in place of general matrix multiplication in at least one of their layers.
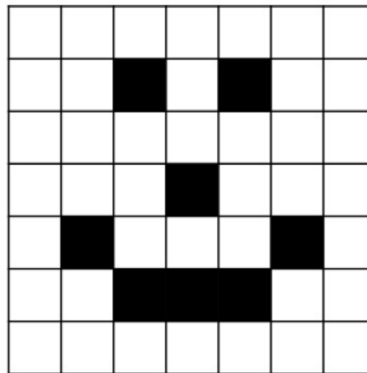
# Images are numbers



Original image    RGB channels

Blue intensity values

0.689 0.706 0.118 0.884 ...
0.535 0.532 0.653 0.925 ...
0.314 0.265 0.159 0.101 ...
0.553 0.633 0.528 0.493 ...
0.441 0.465 0.512 0.512 ...
0.398 0.401 0.421 0.398 ...

Green intensity values

0.342 0.647 0.515 0.816 ...
0.111 0.300 0.205 0.526 ...
0.523 0.428 0.712 0.929 ...
0.214 0.604 0.918 0.344 ...
0.100 0.121 0.113 0.126 ...
0.288 0.187 0.204 0.175 ...

Red intensity values

0.112 0.986 0.234 0.432 ...
0.765 0.128 0.863 0.521 ...
1.000 0.985 0.761 0.698 ...
0.455 0.783 0.224 0.395 ...
0.021 0.500 0.311 0.123 ...
1.000 1.000 0.867 0.051 ...
1.000 0.945 0.998 0.893 ...
0.990 0.941 1.000 0.876 ...
0.902 0.867 0.834 0.798 ...

https://blog.datawow.io/interns-explain-cnn-8a669d053f8b

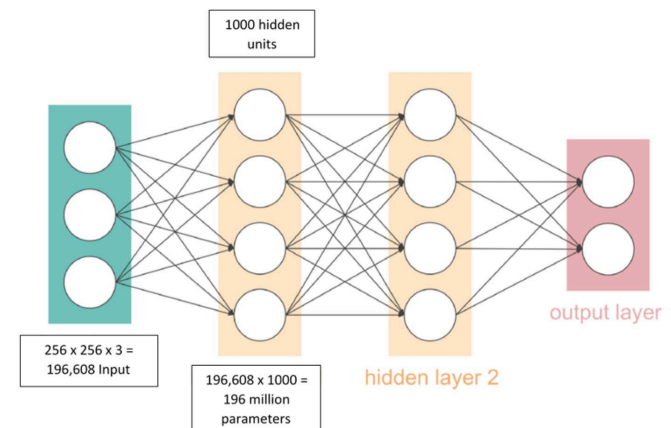# Simple black and white image

o 2D matrix, no grayscale

# Before convolution

o Original values of a 24-bit color images (True Color):

  ◦ 8-bit per color: 0 – 255.

  ◦ Total: 256 *  256 * 256 = 16,777,216 colors

  ◦ Value for red, green, and blue.

o Preprocessing color values:  normalized between 0 and 1
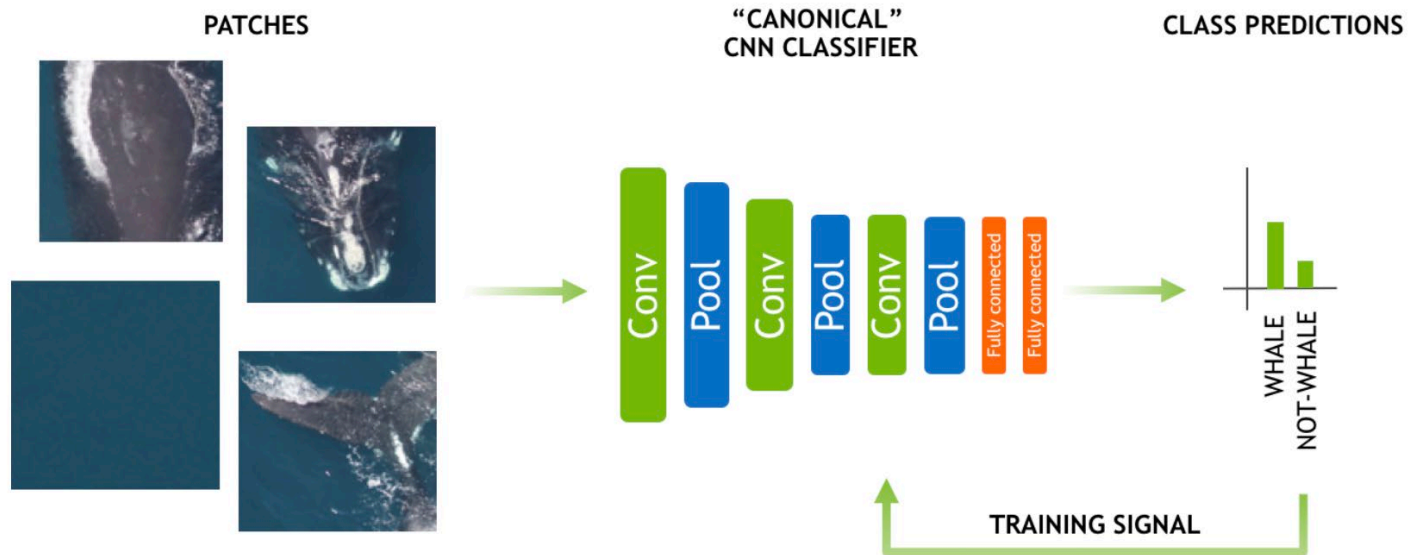  -> will increase performance

# Simple FC Network

o A color image with size 300 x 300 would have 300 x 300 x 3 input values which is equal to 270,000 inputs. If, for example, we have 1,000 hidden units in our first hidden layer, there would be approximately 270 million parameters or weights for us to train which is infeasible.

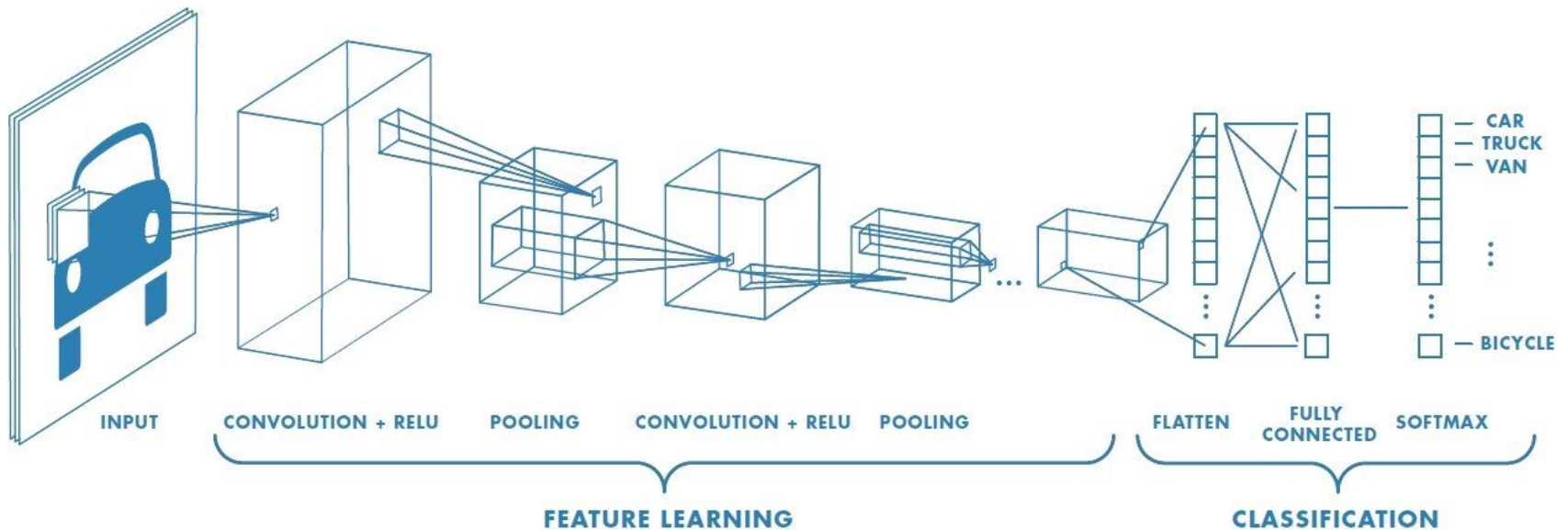o High chance of overfitting and highly complex network



1000 hidden units

256 x 256 x 3 = 196,608 Input

196,608 x 1000 = 196 million parameters

hidden layer 2

output layer

# Solution: convolution

o Reduces the number of parameters we need to learn.

o Preserves locality. We don't have to flatten the image matrix into a vector, thus the relative positions of the image pixels are preserved.
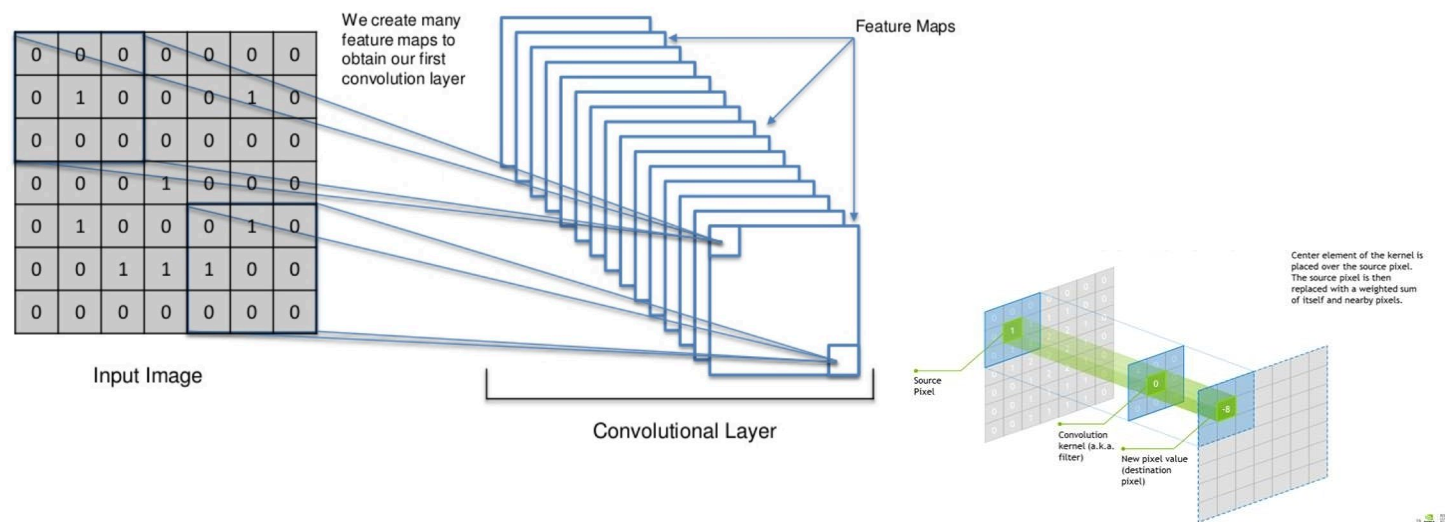
# A typical CNN

# Convolutional layer

o Many feature maps are created, using filters (also called kernels).

o Kernels (or filters) are learned to best fit the task at hand.



[Image source](#)

# Convolution - example

o Edge detection filter/kernel

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

6 x 6

\*

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

3 x 3
filter

=

# Convolution - example

o Edge detection filter/kernel

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

6 x 6

\*

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

3 x 3
filter

=

# Convolution - example

○ Edge detection filter/kernel

= 3x1 + 1x1 + 2x1 + 0x0 + 5x0 + 7x0 + 1x(-1) + 8x(-1) + 2x(-1)

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

6 x 6

\*

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

3 x 3
filter

=

-5

# Convolution - example

○ Edge detection filter/kernel, **stride size = 1**

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

6 x 6

$*$

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

3 x 3
filter

$=$

| | |
|---|---|
| -5 | -4 |

# Convolution - example

o Edge detection filter/kernel, **stride size = 1**

$$
\begin{array}{cccccc}
3 & 0 & 1 & 2 & 7 & 4 \\
1 & 5 & 8 & 9 & 3 & 1 \\
2 & 7 & 2 & 5 & 1 & 3 \\
0 & 1 & 3 & 1 & 7 & 8 \\
4 & 2 & 1 & 6 & 2 & 8 \\
2 & 4 & 5 & 2 & 3 & 9
\end{array}
\quad * \quad
\begin{array}{ccc}
1 & 0 & -1 \\
1 & 0 & -1 \\
1 & 0 & -1
\end{array}
\quad = \quad
\begin{array}{ccc}
-5 & -4 & 0
\end{array}
$$

6 x 6

3 x 3 filter

# Convolution - example

o  Edge detection filter/kernel, **stride size = 1**

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

6 x 6

\*

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

3 x 3
filter

=

| -5 | -4 | 0 | 8 |
|----|----|---|---|
| -10 | -2 | 2 | 3 |
| 0 | -2 | -4 | -7 |
| -3 | -2 | -3 | -16 |

4 x 4

# Convolution - exercise

o Edge detection filter/kernel, **stride size = 1**

| 10 | 10 | 10 | 0 | 0 | 0 |
|----|----|----|---|---|---|
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |

6 x 6

\*

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

3 x 3
filter

=

**?**

? x ?

# Convolution - exercise

o Edge detection filter/kernel, **stride size = 1**

| 10 | 10 | 10 | 0 | 0 | 0 |
|----|----|----|---|---|---|
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |
| 10 | 10 | 10 | 0 | 0 | 0 |

6 x 6

\*

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

3 x 3
filter

=

| 0 | 30 | 30 | 0 |
|---|----|----|---|
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |
| 0 | 30 | 30 | 0 |

4 x 4

# Activation maps



32x32x3 image
5x5x3 filter

activation map

32

32

3

convolve (slide) over all
spatial locations

28

28

1

# Activation maps

o Each filter creates an activation map



32x32x3 image
5x5x3 filter

32

32

3

convolve (slide) over all
spatial locations

activation maps

28

28

1

# Activation maps

For example, if we had 6 5x5 filters, we'll get 6 separate activation maps:

**activation maps**

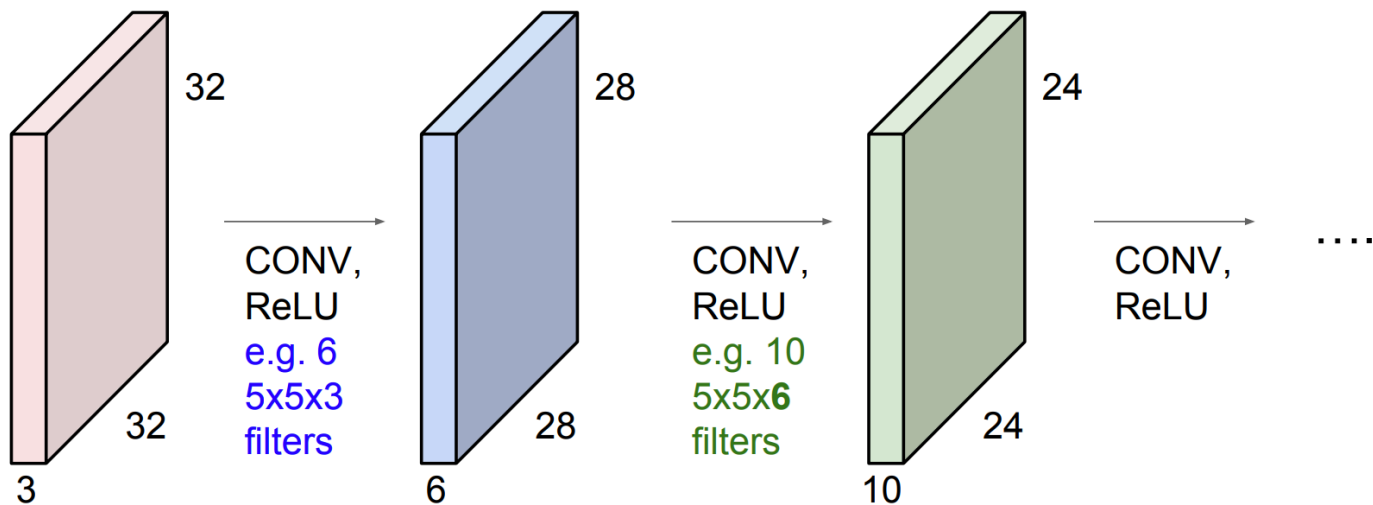32

32

3

Convolution Layer →

28

28

6

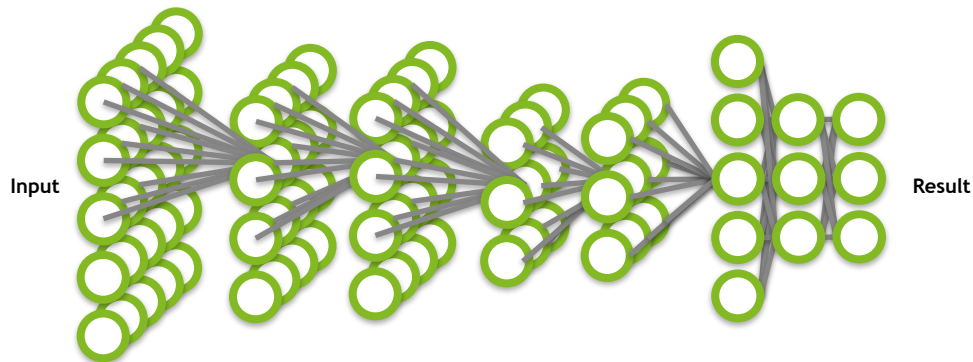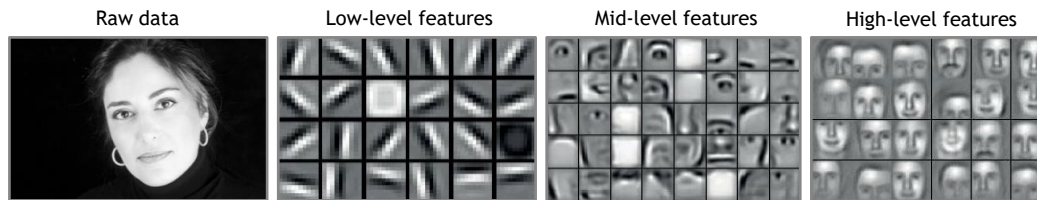We stack these up to get a "new image" of size 28x28x6!

# ConvNet

**Preview:** ConvNet is a sequence of Convolutional Layers, interspersed with activation functions

# Different level of filters



Raw data      Low-level features      Mid-level features      High-level features

Input                                     Result

http://matlabtricks.com/post-5/3x3-convolution-kernels-with-online-demo
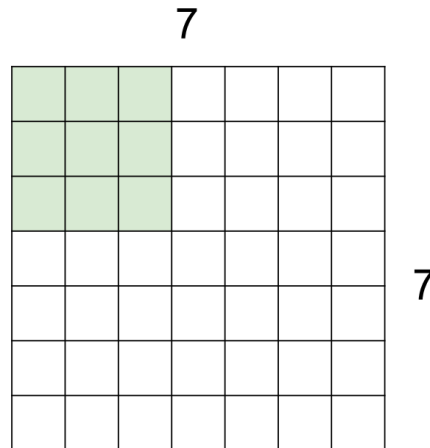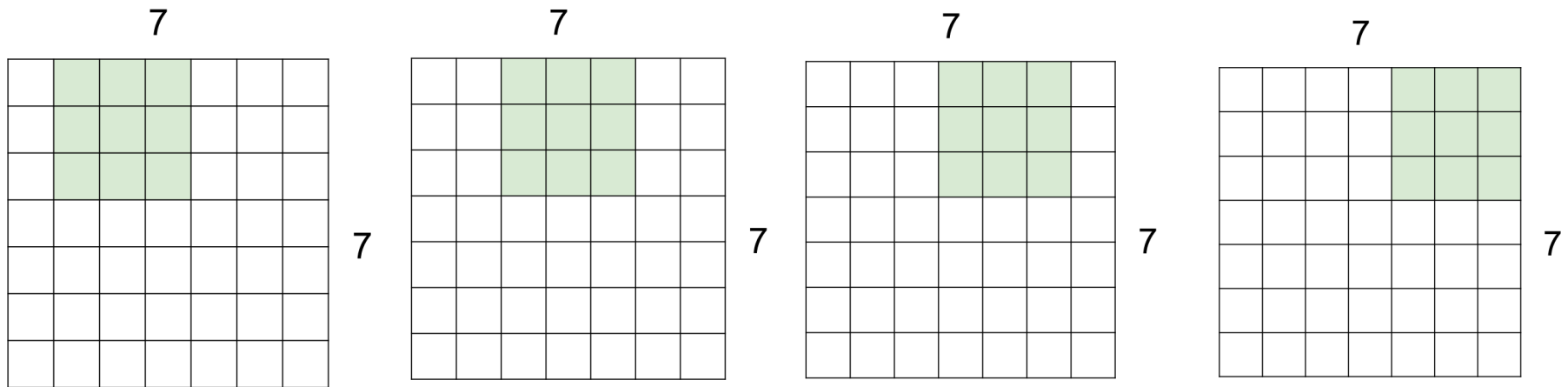
# A closer look at dimensions

○ A 7x7 input with a 3x3 filter:

7

7

# A closer look at dimensions

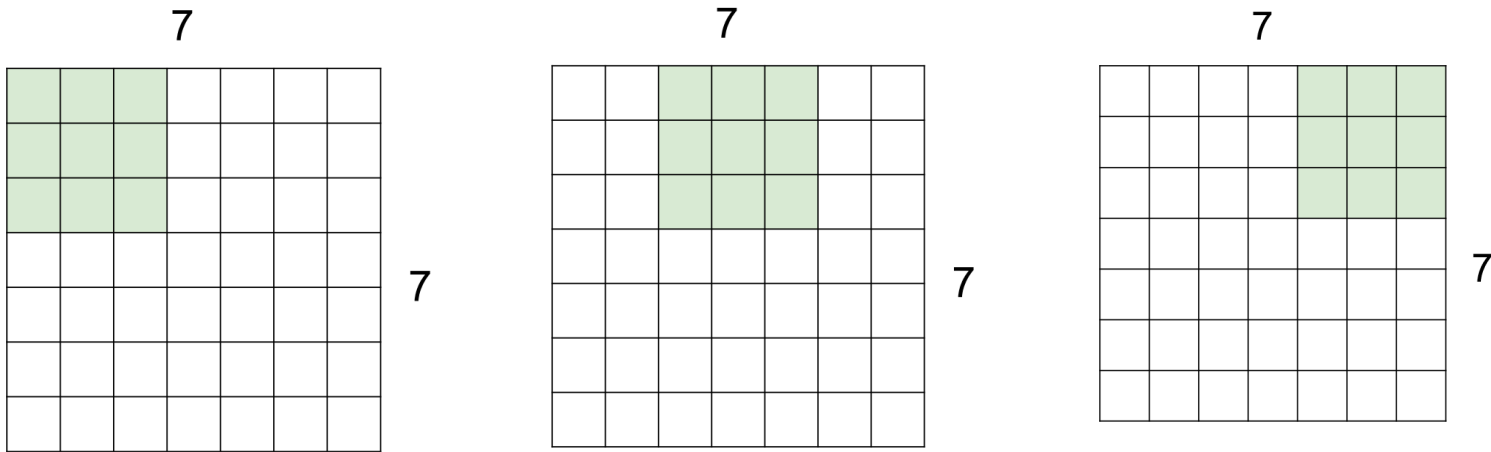o  A 7x7 input with a 3x3 filter, stride = 1

Result: 5x5 output!

# A closer look at dimensions

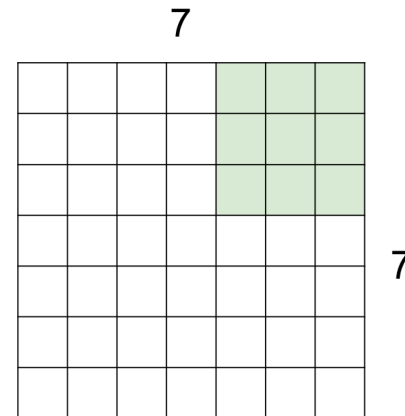○ A 7x7 input with **stride 2** and a 3x3 filter:



○ 3 x 3 output!

# Ex: A closer look at dimensions

o A 7x7 input with stride 3? What is the output size?

o [(N - F) / stride] + 1

o Exercise: N = 7, F = 3:

◦ stride 1 => ?

◦ stride 2 => ?

◦ stride 3 => ?

# Ex: A closer look at dimensions

o A 7x7 input with stride 3? What is the output size?

o [(N - F) / stride] + 1

7

o Exercise: N = 7, F = 3:

◦ stride 1 => (7 - 3)/1 + 1 = 5

◦ stride 2 => (7 - 3)/2 + 1 = 3

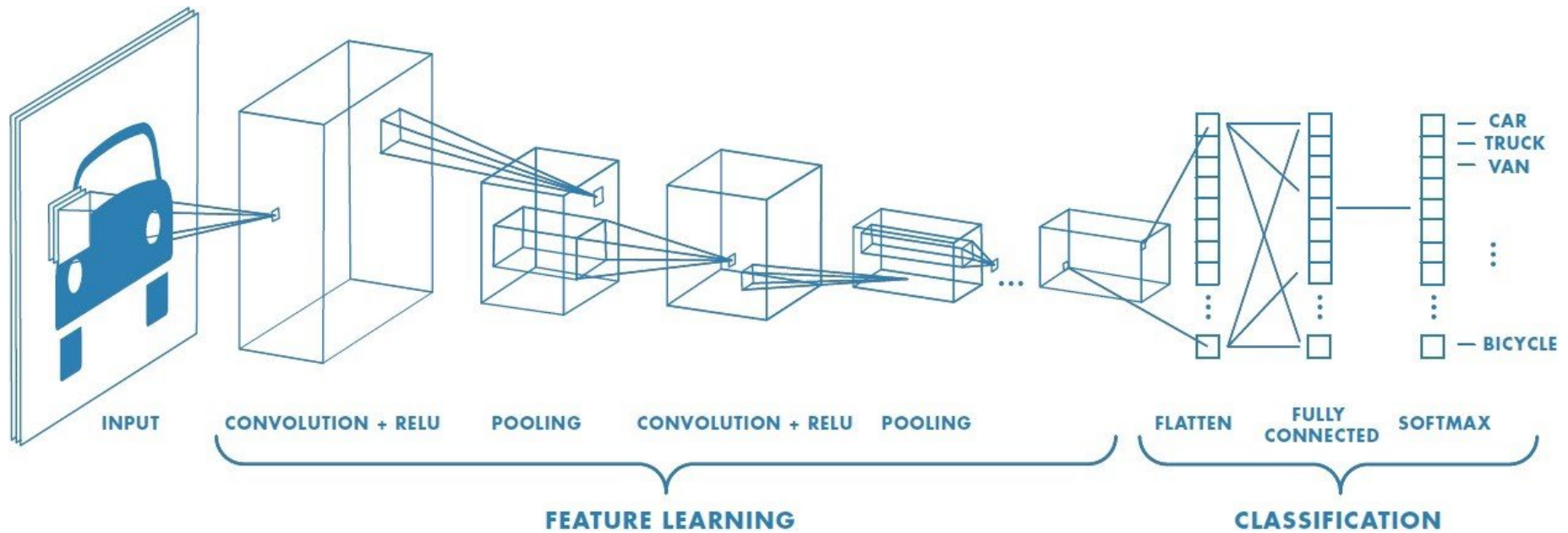◦ stride 3 => (7 - 3)/3 + 1 = 2.33

7

# Padding

o  e.g. input 7x7 image, 3x3 filter, applied with stride 1 pad with 1 pixel border

o  What is the output?
  ◦ 7x7 output!

o  In general, common to see conv. layers with:
  ◦ stride 1
  ◦ filters of size FxF
  ◦ zero-padding with (F-1)/2

o  This will preserve size spatially

# After convolution

# Non-linear activation
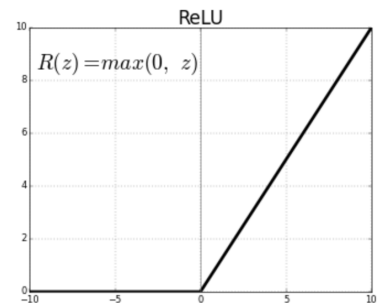
$$f(x) = \max(0, x)$$

o Tanh, sigmoid, or ReLu activation function

o Most popular (often performs best): ReLu

=> to introduce non-linearity in our ConvNet, since most of the real-world data we would want our ConvNet to learn would be non-linear.

o After each conv layer, it is convention to apply a non-linear layer (or activation layer) immediately afterwards.

ReLU

$R(z) = max(0, z)$

# Pooling

o Non-linear down-sampling to simplify the output of the convolutional layer.

o ConvNets often use pooling layers to reduce the size of the representation, speed up the computation, as well as make some of the detected features a bit more robust.

o Types of pooling:

   ◦ Max pooling (popular)

   ◦ Average pooling

o Typical shape: 2x2 or sometimes 4x4

o Too large window: dramatic loss of information

o Non-overlapping windows perform the best

# Pooling

o Hyperparameters:

◦ Stride size

◦ Pooling window size



o Pooling layers don't learn themselves, they just reduce the size of the problem

Image: https://medium.com/data-science-group-iitr/building-a-convolutional-neural-network-in-python-with-tensorflow-d251c3ca8117

# Exercise: Pooling

o Hyperparameters:

◦ Stride size: 2

◦ Pooling window size: 2 x 2

Max Pooling

**?**

? x ?

| 10 | 3  | 0  | 18 |
|----|----|----|----|
| 0  | 3  | 6  | 4  |
| 10 | 12 | 0  | 8  |
| 0  | 2  | 30 | 2  |

Average Pooling

**?**

? x ?

# Exercise: Pooling

o Hyperparameters:

◦ Stride size: 2

◦ Pooling window size: 2 x 2

| 10 | 3 | 0 | 18 |
|----|----|----|----|
| 0 | 3 | 6 | 4 |
| 10 | 12 | 0 | 8 |
| 0 | 2 | 30 | 2 |

**Max Pooling**

| 10 | 18 |
|----|----|
| 12 | 30 |

2 x 2

**Average Pooling**

**?**

? x ?

# Exercise: Pooling

o Hyperparameters:

◦ Stride size: 2

◦ Pooling window size: 2 x 2

**Max Pooling**

| 10 | 18 |
|----|----|
| 12 | 30 |

2 x 2

| 10 | 3  | 0  | 18 |
|----|----|----|----|
| 0  | 3  | 6  | 4  |
| 10 | 12 | 0  | 8  |
| 0  | 2  | 30 | 2  |

**Average Pooling**

| 4 | 7  |
|---|----|
| 6 | 10 |

2 x 2

# Flatten
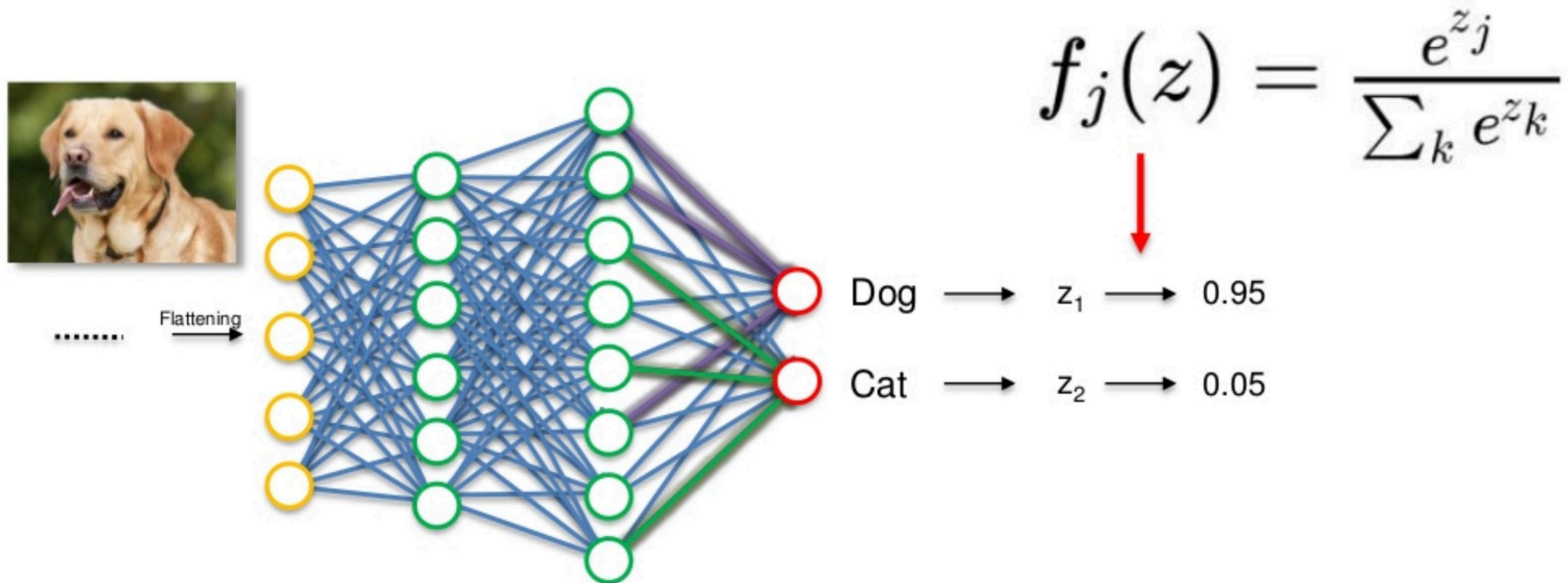


Pooled Feature Map

Flattening

# FC

o The high-level reasoning in the neural network is done via fully connected layers. Neurons in a fully connected layer have connections to all activations in the previous layer, as seen in regular neural networks. Their activations can hence be computed with a matrix multiplication followed by a bias offset

o Training Loss: how training penalizes the deviation between the predicted and true labels and is normally the final layer. Various loss functions appropriate for different tasks may be used there:

◦ Softmax loss (a Softmax activation plus a Cross-Entropy loss) is used for multiclass classification (it distributes the probability throughout each output node, meaning that the sum of all probabilities is 1)

◦ Sigmoid cross-entropy loss (Sigmoid activation plus a Cross-Entropy loss) is used for binary classification

◦ Euclidean loss is used for regressing to real-values. (could be mean squared error: mse)

# Softmax



$$f_j(z) = \frac{e^{z_j}}{\sum_k e^{z_k}}$$

Dog $\longrightarrow$ $z_1$ $\longrightarrow$ 0.95

Cat $\longrightarrow$ $z_2$ $\longrightarrow$ 0.05

```python
class myCNN(nn.Module):
    def __init__(self):
        super(myCNN, self).__init__()

        # Convolutional layers.
        self.conv1 = nn.Conv2d(1, 6, 5) #in_channels, out_channels, kernel_size
        self.conv2 = nn.Conv2d(6, 12, 5)

        # Linear layers.
        self.fc1 = nn.Linear(12*4*4, 120)

    def forward(self, x):
        # Conv1 + ReLU + MaxPooling.
        out = F.relu(self.conv1(x))
        out = F.max_pool2d(out, 2)

        # Conv2 + ReLU + MaPooling.
        out = F.relu(self.conv2(out))
        out = F.max_pool2d(out, 2)
        out = out.view(out.size(0), -1) #flatten

        # Linear layer + ReLU.
        out = F.relu(self.fc1(out))
        Return out

model = myCNN()
loss_fn = nn.CrossEntropyLoss() #includes softmax
```
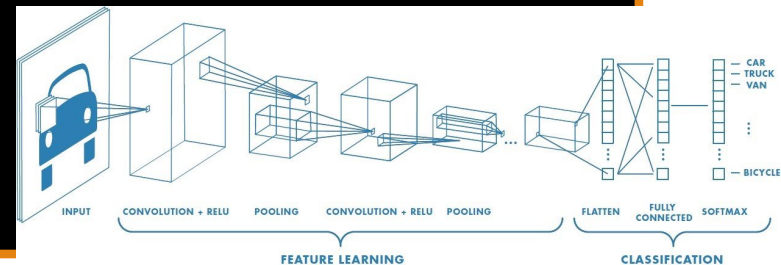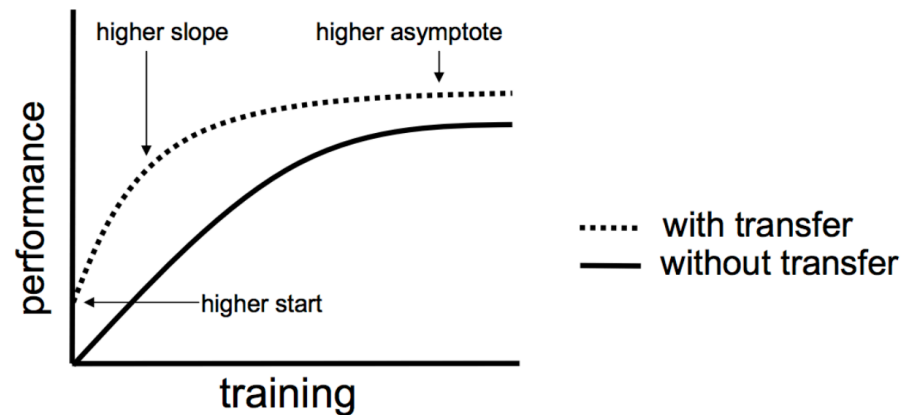
# Popular CNNs

o LeNet-5:        Developed 1998, Parameters 60k

o AlexNet:        Developed 2012, Parameters 60M

o VGG-16:        Developed 2014, Parameters 138M

o InceptionV3:    Developed 2014, Parameters 23M

o ResNet50:      Developed 2015, Parameters 25M

o ResNext50:      Developed 2016, Parameters 25M

o DenseNet201:  Developed 2017, Parameters 20M

# Transfer Learning

○ Transfer learning is an optimization, a shortcut to saving time or getting better performance.

○ Don't reinvent the wheel
  ◦ Use pertained models from a larger dataset or related task and use those to represent your input
  ◦ Then you can finetune these weights further

# PyTorch finetune pertained networks

```
cnn_model = models.inception_v3(pretrained = True)


#or


cnn_model = models.resnet50(pretrained = True)


# print(cnn_model) # Print the model to see what you can modify.

# We are modifying the last layer which is stored in the fc property
# for this model as you can see by printing out the network.

cnn_model.fc = nn.Linear(2048, len(train_dataset.classes))


# print(cnn_model)   # Verify that the last linear layer was changed.
```

# Similarly for text

o Remember we were not reinventing the wheel…

- Google - word2vec
  - Google News: 3 million 300-dimension English word vectors)
- Stanford - GLOVE
  - Wikipedia 2014 + Gigaword 5 (6B tokens, 400K vocab, uncased, 50d, 100d, 200d, & 300d vectors, 822 MB download)
- Google – BERT
  - BooksCorpus (800M words) and English Wikipedia (2,500M words). Available on https://huggingface.co/

# Credits

o Images thanks to:
  - MIT 6.S191
  - https://towardsdatascience.com/https-medium-com-piotr-skalski92-deep-dive-into-deep-networks-math-17660bc376ba
  - https://mattmazur.com/2015/03/17/a-step-by-step-backpropagation-example/
  - https://colah.github.io/posts/2015-08-Understanding-LSTMs/
  - https://colah.github.io/
  - https://www.deeplearningbook.org/contents/convnets.html
  - https://github.com/BlackBindy/MNIST-invert-color
  - https://www.slideshare.net/GauravMittal68/convolutional-neural-networks-cnn
  - http://slazebni.cs.illinois.edu/spring17/lec01_cnn_architectures.pdf
  - https://www.jeremyjordan.me/convnet-architectures/#resnext
  - https://www.superdatascience.com/ppt-the-ultimate-guide-to-convolutional-neural-networks-cnn/