



# Intro to Neural Nets

Week 5: Convolutional NNs (ConvNets)

# Today's Agenda

## **Convolutional Neural Networks (CNNs)**

- What CNNs try to accomplish
- What is a convolution?
  - Padding, strides, filters
- What is pooling?
  - Max, min, avg pooling.

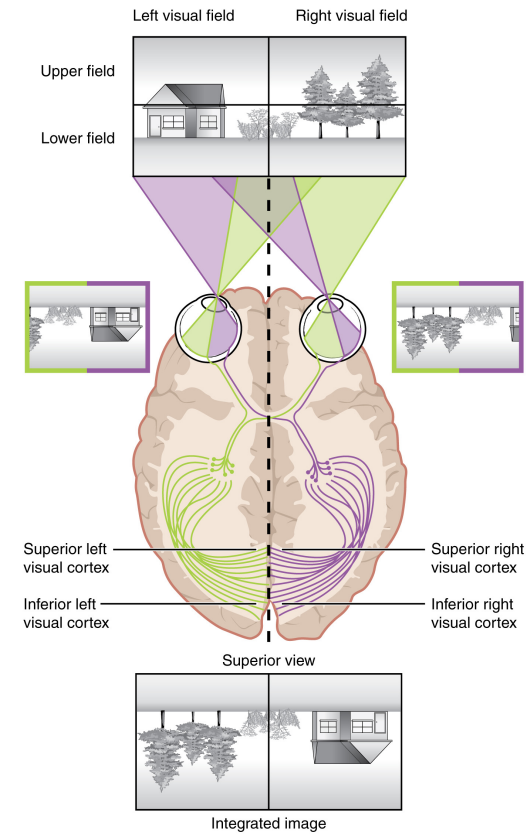
## **Other Stuff**

- CNN specific techniques to avoid overfitting (data augmentation).
- Extracting feature representations from your trained model.
- Adapting pre-trained models (transfer learning).

# Inspiration for Convnets

## Our Visual System

- Human eye is basically a 576-megapixel video camera.
- For comparison, the Pixel 6 camera is 50-megapixels.
- The human field of vision is not a square; something like a video camera that records individual image frames comprised of 24,000 x 24,000 pixels.



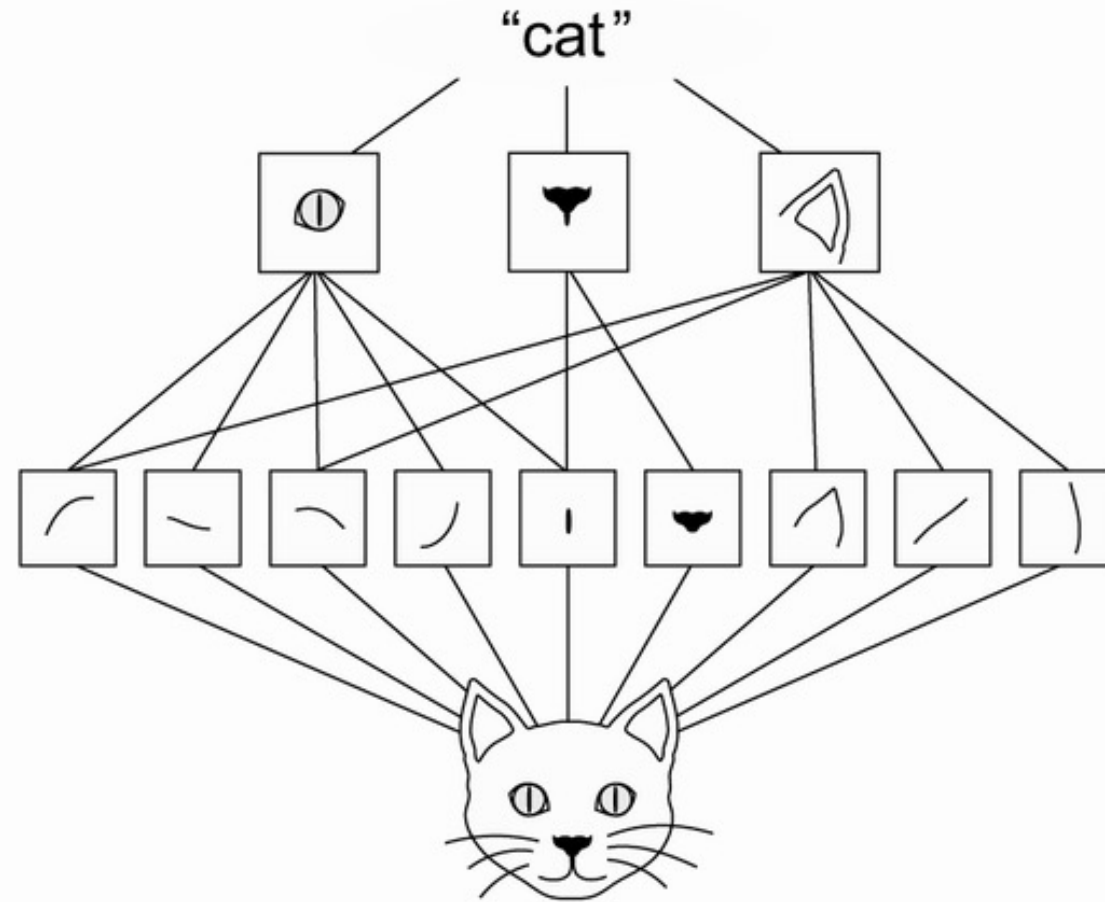
# Feature Detection

## How Does Your Visual System Work?

- We think that your brain processes individual visual receptors in groups, identifies combinations of inputs in proximity to one another that imply something like an edge (edge detection), combines that with color and so on. These low-level features are then processed together to arrive at higher level objects (e.g., a nose, a mouth, an eye).
- Those higher-level features are then processed together to yield a face (perhaps someone we know or do not know). Hence why you might have a hard time recognizing someone who has a new haircut, or who is wearing a facemask!



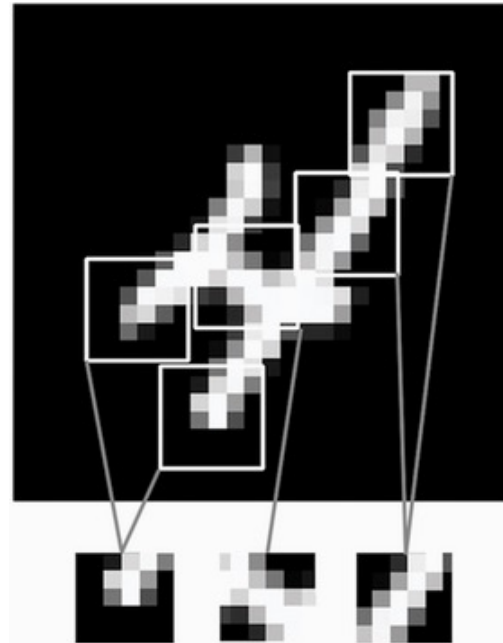
# Feature Detection / Aggregation



# What is Convolution?

## Hone-in On Sub-sections of the Image

- So, if we have a 28x28 image, we might separately consider 3x3 pixel subsection of that image. Each subsection (they can be overlapping) is represented by its own node in the first hidden layer.
- That local input matrix (subfield) is considered in tandem with a 'filter' a matrix of weights. A filter might be something like



# The Convolution Operation

## Consider in Matrix Representation

- We have the raw image data (a.k.a. input feature map), the filter, and the result of passing our filter over our image (a.k.a. output feature map). We will have one output feature map for a given image, per filter (each filter is intended to detect a different type of feature).
- The filter elements are just weights for the Conv layer; we learn the filter values as part of the backpropagation process. So, the CNN will figure out what features to look for to predict the label (probably what a baby does when its first board and first learning how to process visual information).

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

--	--

  
 $*$ 

1	0	-1
1	0	-1
1	0	-1

--	--	--

  
 $=$ 

0	30	30	0
0	30	30	0
0	30	30	0
0	30	30	0

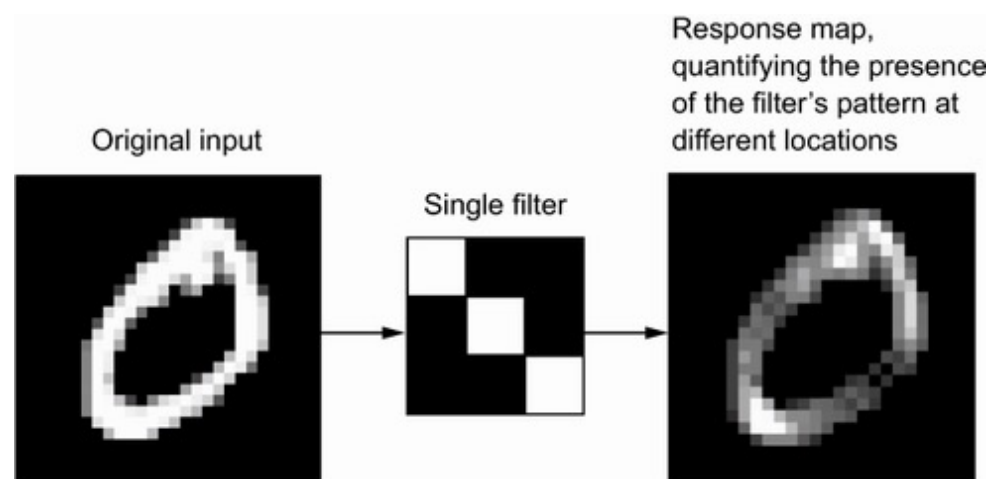
  

--	--	--

# The Convolution Operation

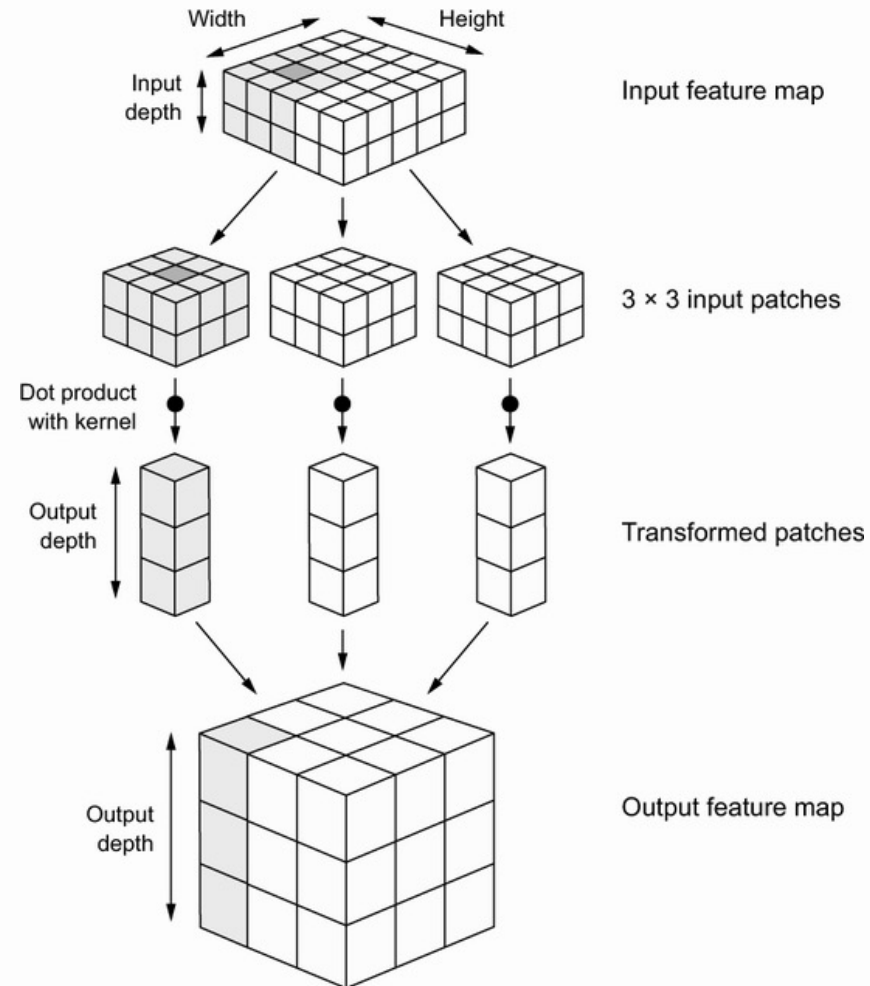
## Consider in Matrix Representation

- We have the raw image data (a.k.a. input feature map), the filter, and the result of passing our filter over our image (a.k.a. output feature map). We will have one output feature map for a given image, per filter (each filter is intended to detect a different type of feature).
- The filter elements are just weights for the Conv layer; we learn the filter values as part of the backpropagation process. So, the CNN will figure out what features to look for to predict the label (probably what a baby does when its first board and first learning how to process visual information).





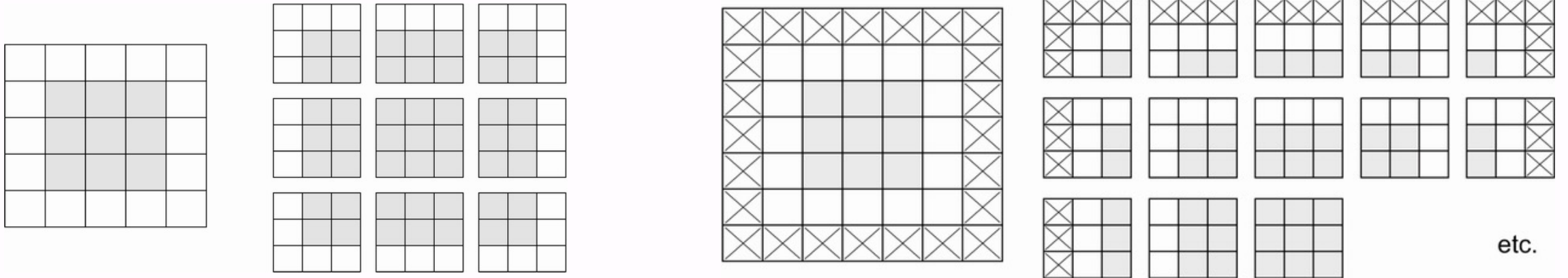
# Depth of Output = Filter #, !Color Channels



# Padding

## Padding

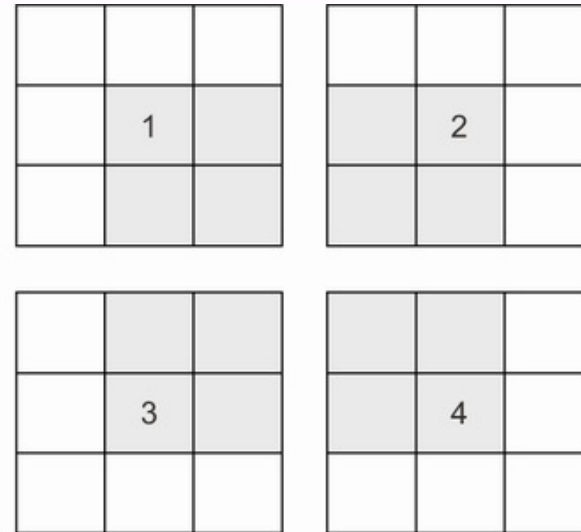
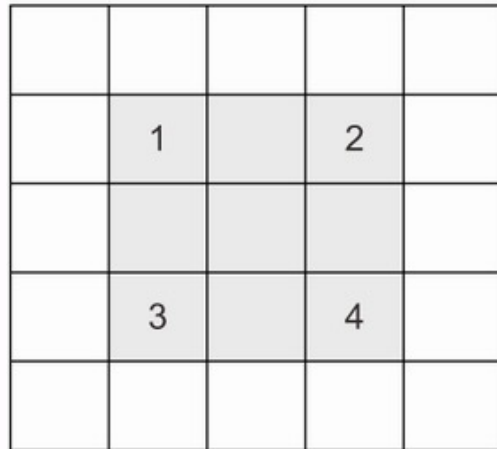
- To prevent the transformation from down-sampling (reducing the size of the matrix during convolution to output), we can pad the edges of the image with 0's.



# Strides

## Strides

- Often, we will pass the filter over every pixel cell, but we don't have to; we might pass over every other cell. This is what strides refers to (skipping).



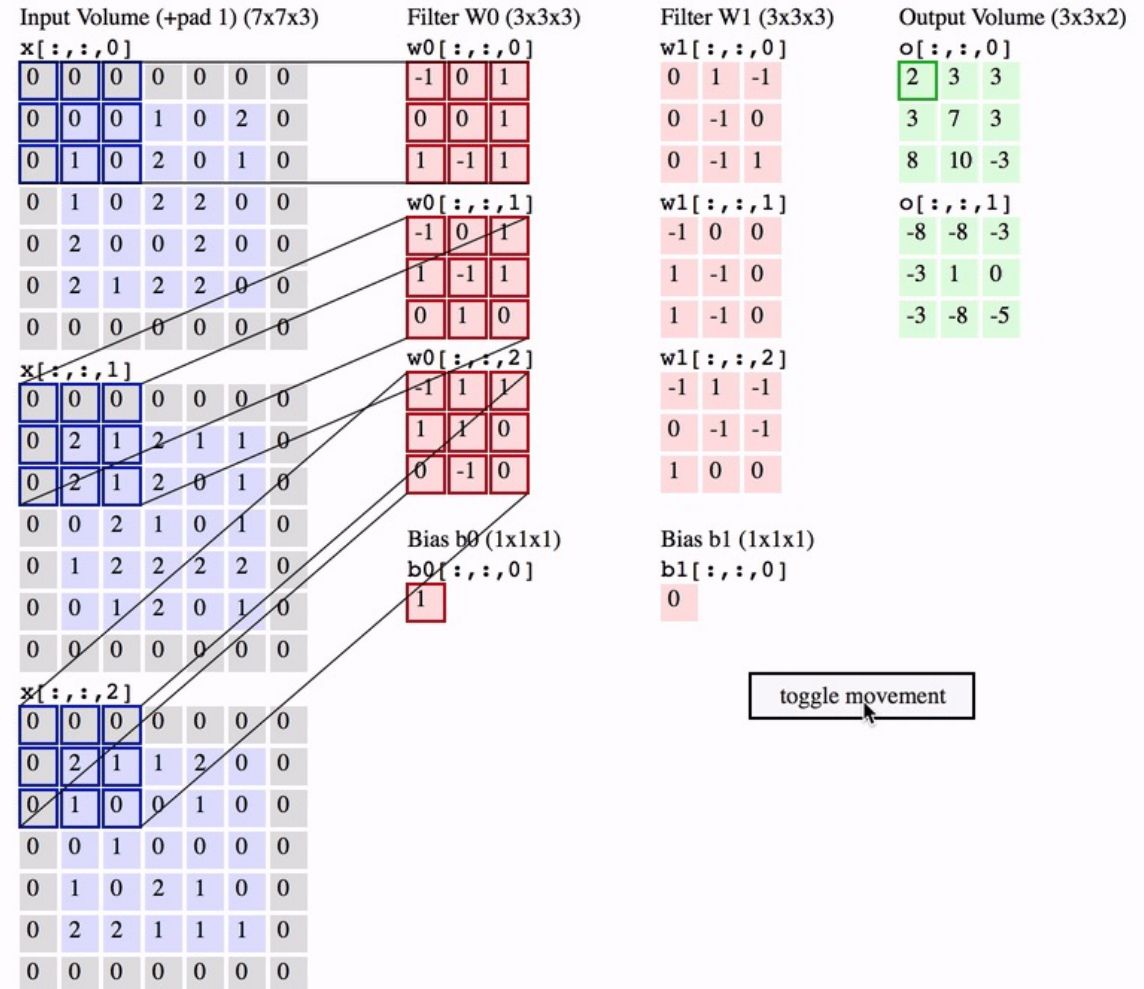
# Padding and Strides

## Padding

- To prevent the transformation from down-sampling (reducing the size of the matrix during convolution to output), we can pad the edges of the image with 0's.

## Strides

- Often, we will pass the filter over every pixel cell, but we don't have to; we might pass over every other cell.
- This is what strides refers to (skipping).



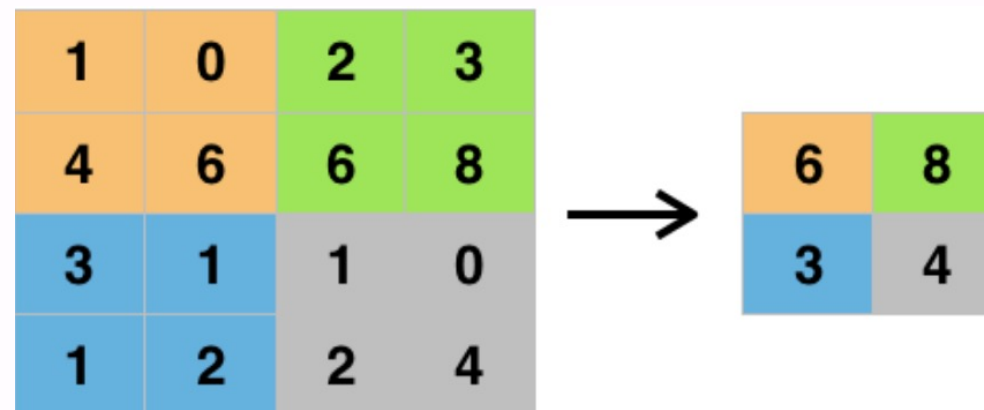
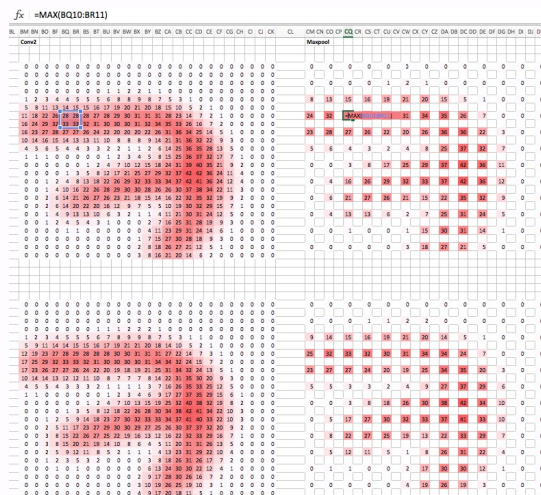
# What is Pooling?

## Down-sampling Detected Features

- The idea is to compress the resulting data down into a coarser representation, to reduce model complexity, and to also force attention toward a broader section of the original image (helps reduce overfitting).

## Forcing Attention to Larger Blocks of the Original Image

- Because we typically use stride = pool width, the pooling output is aggregating over segments of the input.



# What is GlobalAveragePooling?

**Think About This as Something Like a Global 'Feature Confidence' Map**

- Essentially, this create a map of the whole

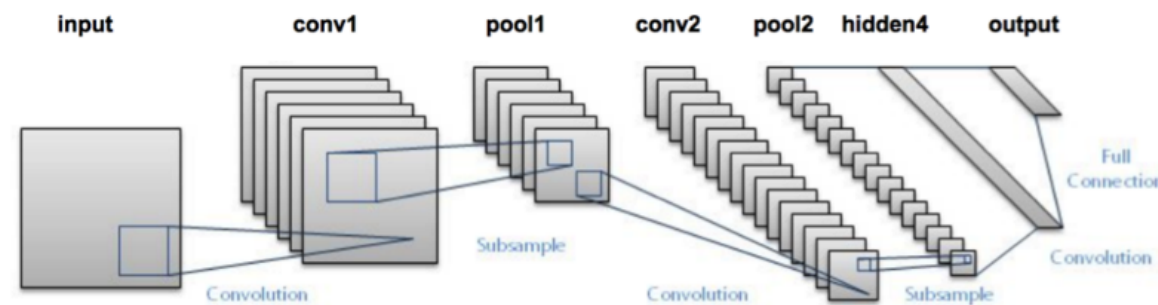
# Basic Image Labeling Topology

## Progressively More Filters

- As you move through the network, the number of features rises exponentially.
- More filters as you move along means it allows more permutations / combinations

## Progressively Smaller Filter Maps

- Smaller filter map arises from the pooling steps, which means that each element of the final map distills features (high level features, derived from low level features, derived from raw pixels) derived from a larger segment of the original picture.

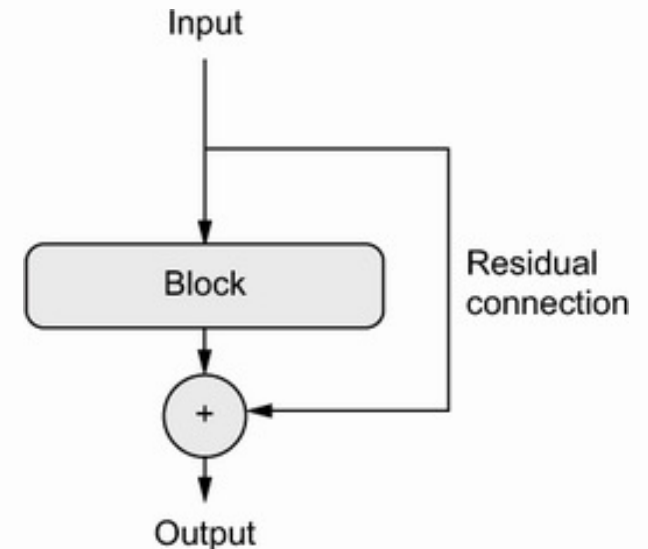


# Residual Connections

## Construct Repeating 'Blocks' of Layers

- Convolution, Pooling, Residual Connection, Convolution, Pooling, Residual Connection...
- In deep networks, vanishing gradient problem arises because there is noise in every step, and if we go far enough back the noise overwhelms the signal. The residual connections help ensure we have more signal than noise.

```
1 inputs = keras.Input(shape=(32, 32, 3))
2 x = layers.Rescaling(1./255)(inputs)
3
4 def residual_block(x, filters, pooling=False):
5     residual = x
6     x = layers.Conv2D(filters, 3, activation="relu", padding="same")(x)
7     x = layers.Conv2D(filters, 3, activation="relu", padding="same")(x)
8     if pooling:
9         x = layers.MaxPooling2D(2, padding="same")(x)
10        residual = layers.Conv2D(filters, 1, strides=2)(residual)
11    elif filters != residual.shape[-1]:
12        residual = layers.Conv2D(filters, 1)(residual)
13    x = layers.add([x, residual])
14    return x
15
16 x = residual_block(x, filters=32, pooling=True)
17 x = residual_block(x, filters=64, pooling=True)
18 x = residual_block(x, filters=128, pooling=False)
19
20 x = layers.GlobalAveragePooling2D()(x)
21 outputs = layers.Dense(1, activation="sigmoid")(x)
22 model = keras.Model(inputs=inputs, outputs=outputs)
23 model.summary()
```



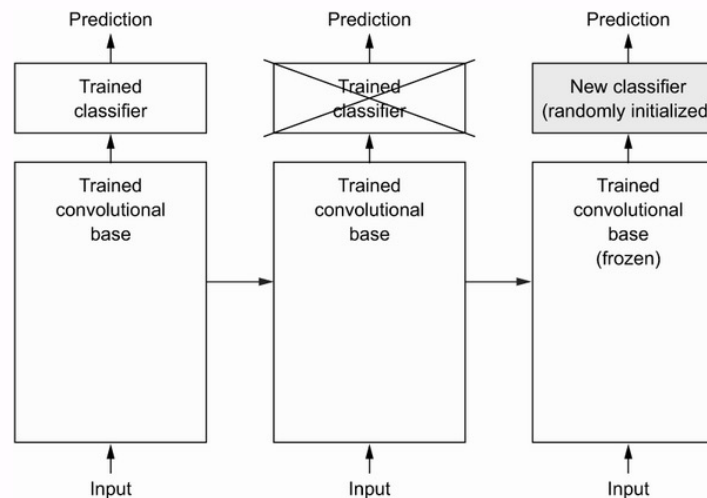


# Pre-Trained Models: Feature Extraction

Take the convolutional base layers from someone else's model, then...

## Two Options

- Feed Data Through Model Base: feed your images through convolutional base, take the outputs, and then use those as your predictors, feeding them into a network of dense layers.
- Freeze Model Base and Include in Network: Take the convolutional base layers from someone else's model and freeze them (make parameters non-trainable), then stack your (trainable) Dense layers onto the end. This lets you add data-augmentation to the front of the model.



# Pre-Trained Models: Fine Tuning

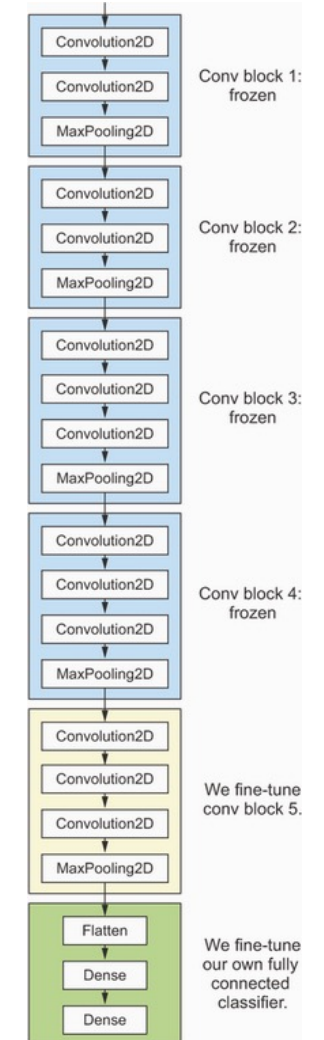
Take the convolutional base layers from someone else's model, then...

## Freeze Only the First Several Layers

- Allow your network to modify / update the last few convolutional base layers as part of training, along with your own Dense layers...
- Iterate over the layers in the network and set the last few to be trainable.

**Listing 8.27** Freezing all layers until the fourth from the last

```
1 conv_base.trainable = True
2 for layer in conv_base.layers[:-4]:
3     layer.trainable = False
```



# Homework Assignment #2

## Multi-modal Prediction Exercise

- I am providing you with a sample of real estate data that includes photos of houses (four of each), along with home characteristics as numeric features.
- Your goal is to use this multi-modal data to predict home price.

## Deliverable

- Produce a Jupyter notebook documenting your work (include the names of the contributors at the top of your notebook).
- Feel free to incorporate heavy comments so we (the TA and I) are clear what you were trying to do! Feel free to re-use code from class examples, but you should not collaborate with anyone else (outside your pair, if you decide to submit as a pair).
- Pre-process the data as / if necessary; define your models, explain your thought process, e.g., around the choice of loss function, activation functions, network topology, etc.
- Tune the model; describe what 'experiments' you ran exploring different configurations and report ultimate performance. You should work some sort of holdout sample here (e.g., cross-validation).

# Questions?