# Computer Vision

# What's In the News?



**Gizmodo**  LATEST  NEWS  REVIEWS  IO9  SCIENCE  DEALS  DOWNLOADS  NEWSLETTERS

ARTIFICIAL INTELLIGENCE

## Pentagon Considers Designating Anthropic AI as a 'Supply Chain Risk': Report

The move would require anyone doing business with the U.S. military to cut ties with the AI company.

THE DAILY  The New York Times  GIVE THE TIMES  Account

The Daily  Subscribe: Apple Podcasts · Google Podcasts

Feb. 18, 2026

## Can A.I. Already Do Your Job?

Tools like Claude Code generate computer code when people type prompts, so those with no coding experience can create their own programs and apps.

# This Week

## Convolutional Neural Networks (CNNs)

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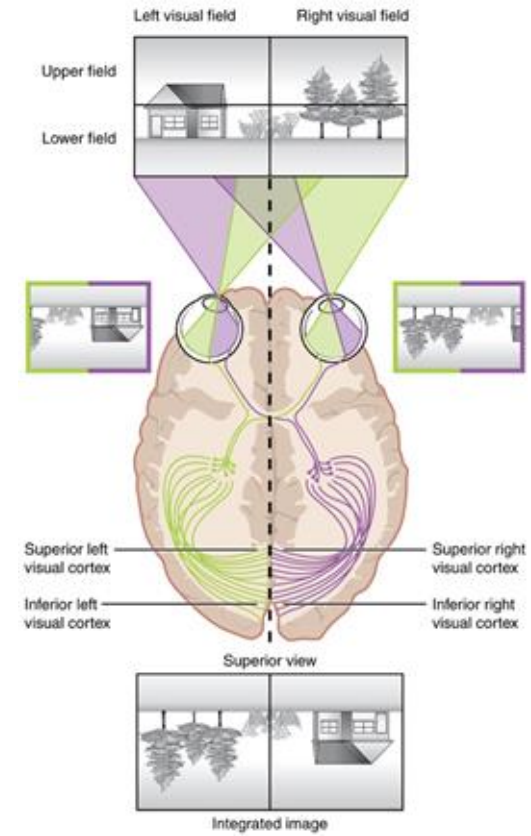- What is pooling?
  - Max, min, avg pooling.

## Other Stuff

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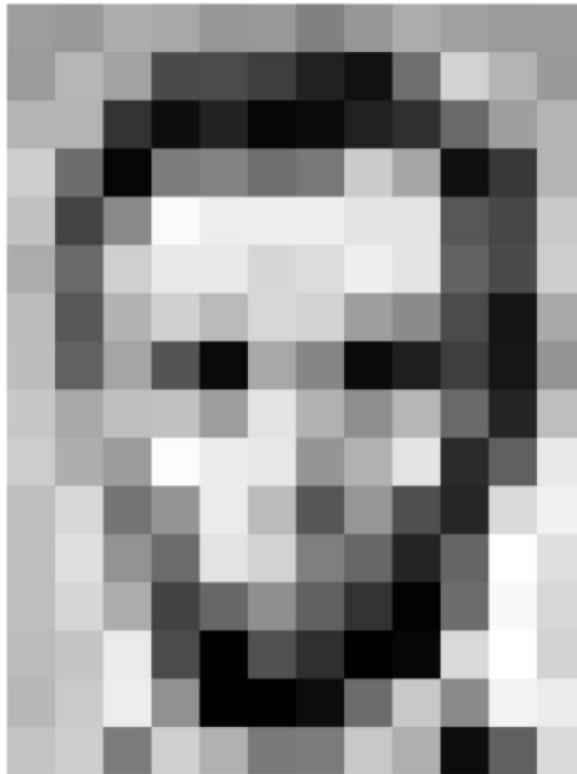- Adapting pre-trained models (transfer learning).

# 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.

# Images are Numbers

What the computer sees



An image is just a matrix of numbers [0,255]!  i.e., 1080x1080x3 for an RGB image
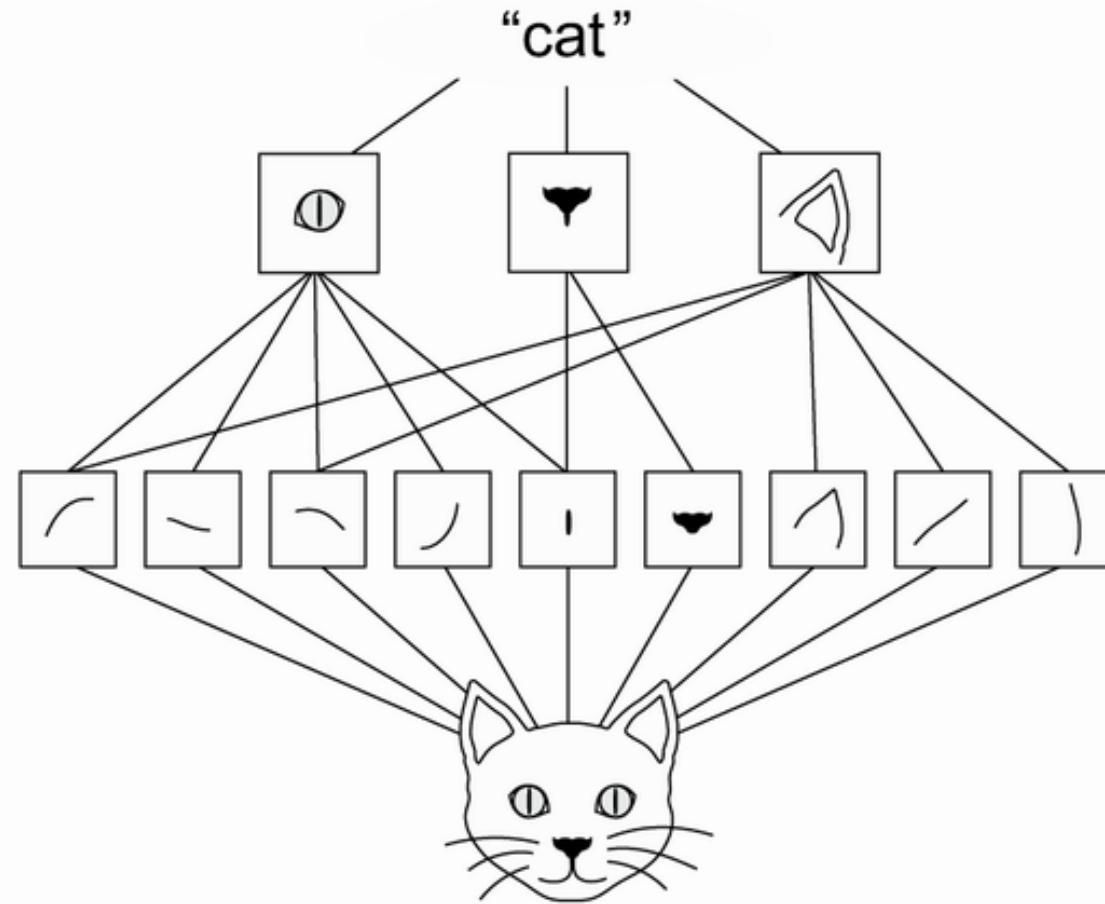
# We Don't Detect Features 'Manually'

## 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!
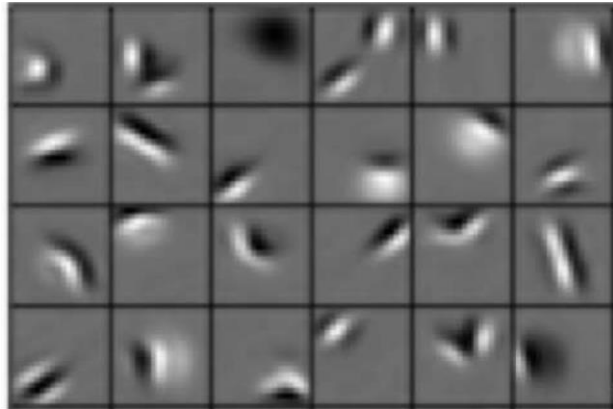
# Feature Detection / Aggregation
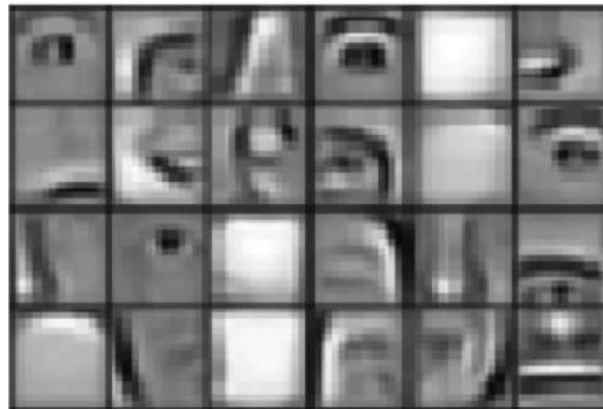
# Learning Feature Representations

Can we learn a **hierarchy of features** directly from the data instead of hand engineering?

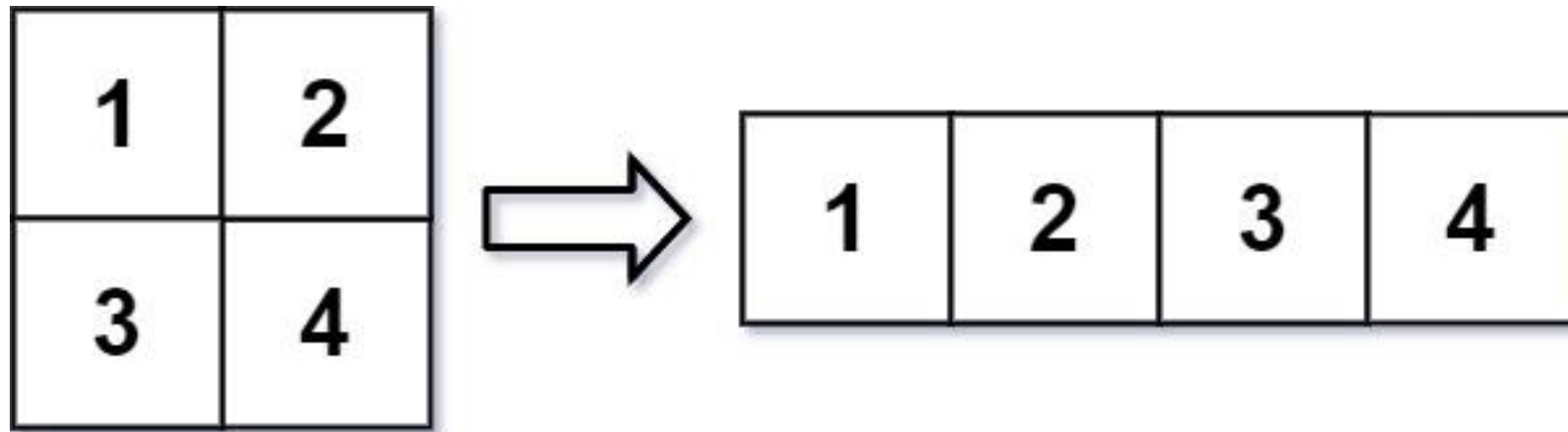| Low level features | Mid level features | High level features |
|---|---|---|



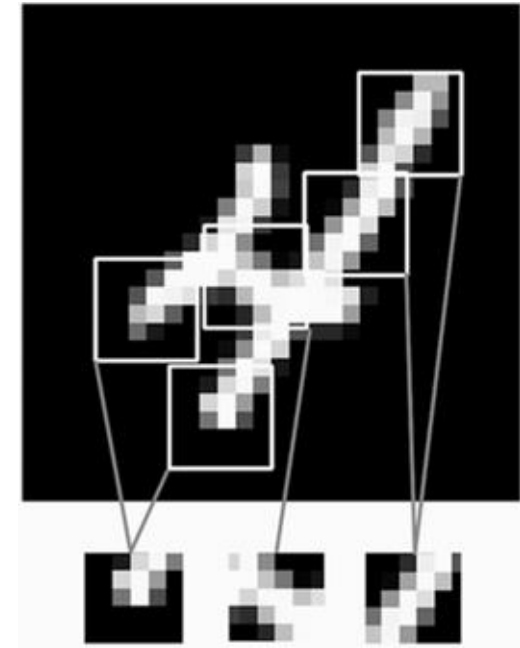| Edges, dark spots | Eyes, ears, nose | Facial structure |
|---|---|---|

# Review: Why are Fully Connected Dense Layers (MLPs) Not Good for This?

# We Need to Preserve Spatial Arrangements

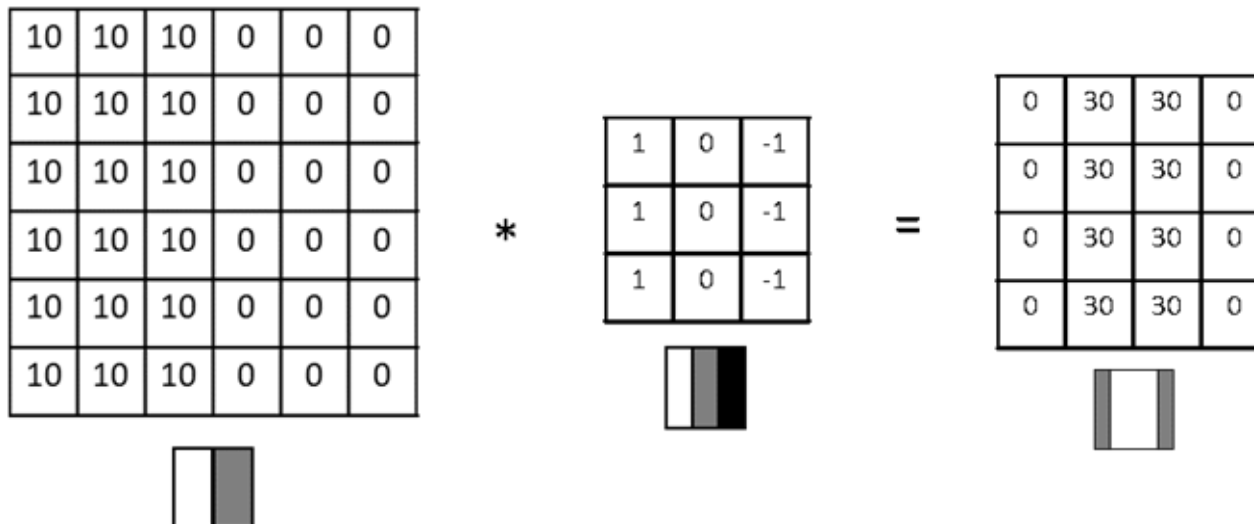## Hone-in On Sub-sections of the Image

- 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
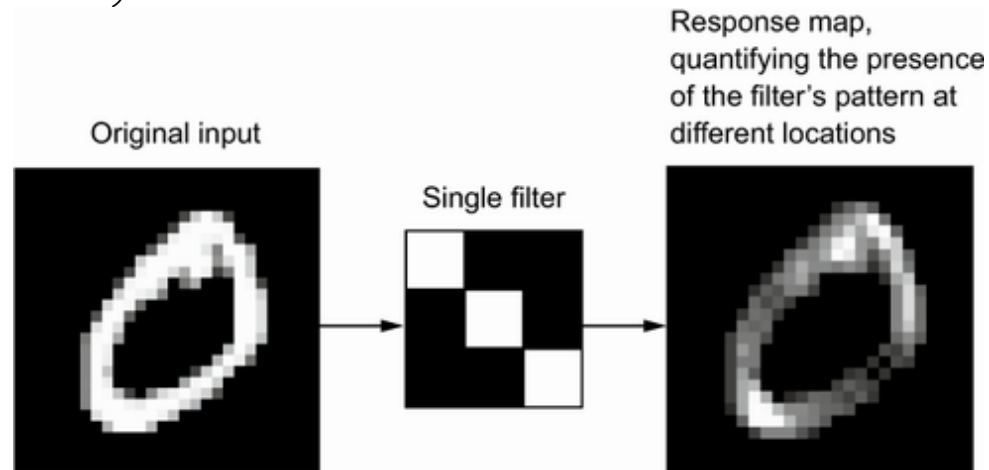
## Consider in Matrix Representation

- 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
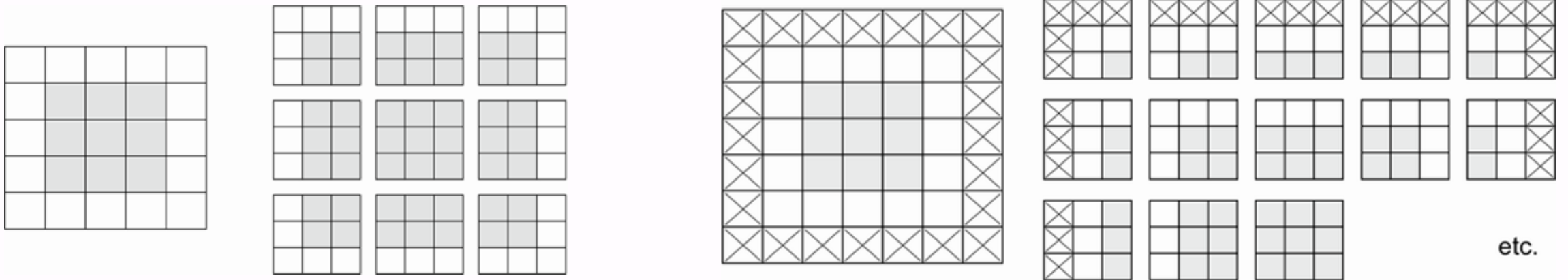
## Consider in Matrix Representation

- 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).



Original input — Single filter — Response map, quantifying the presence of the filter's pattern at different locations
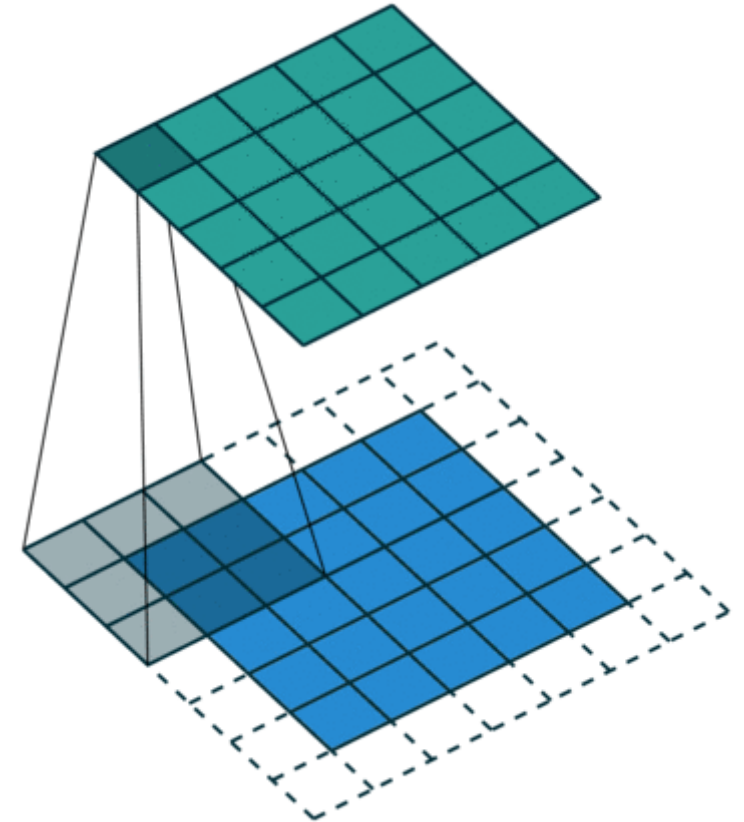
# Padding

## 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.

# Padding Technique



Zero-padding added to image



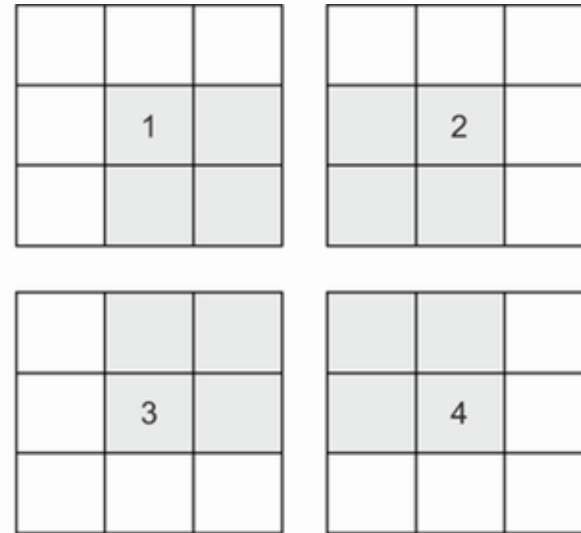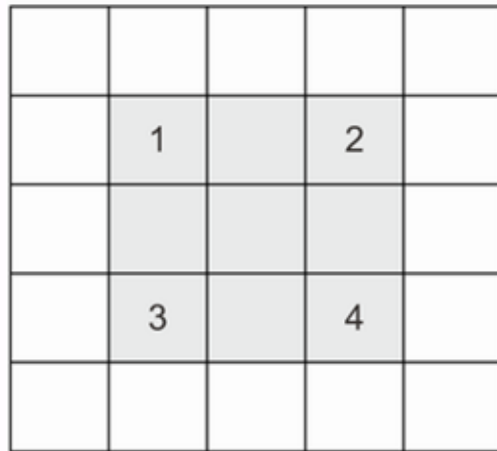Add appropriate padding to keep the dimension after filtering

$m \times m$ input image with $ff \times ff$ filters
add $p = (ff - 1)/2$ padding keep the same dimension of input image

# 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).
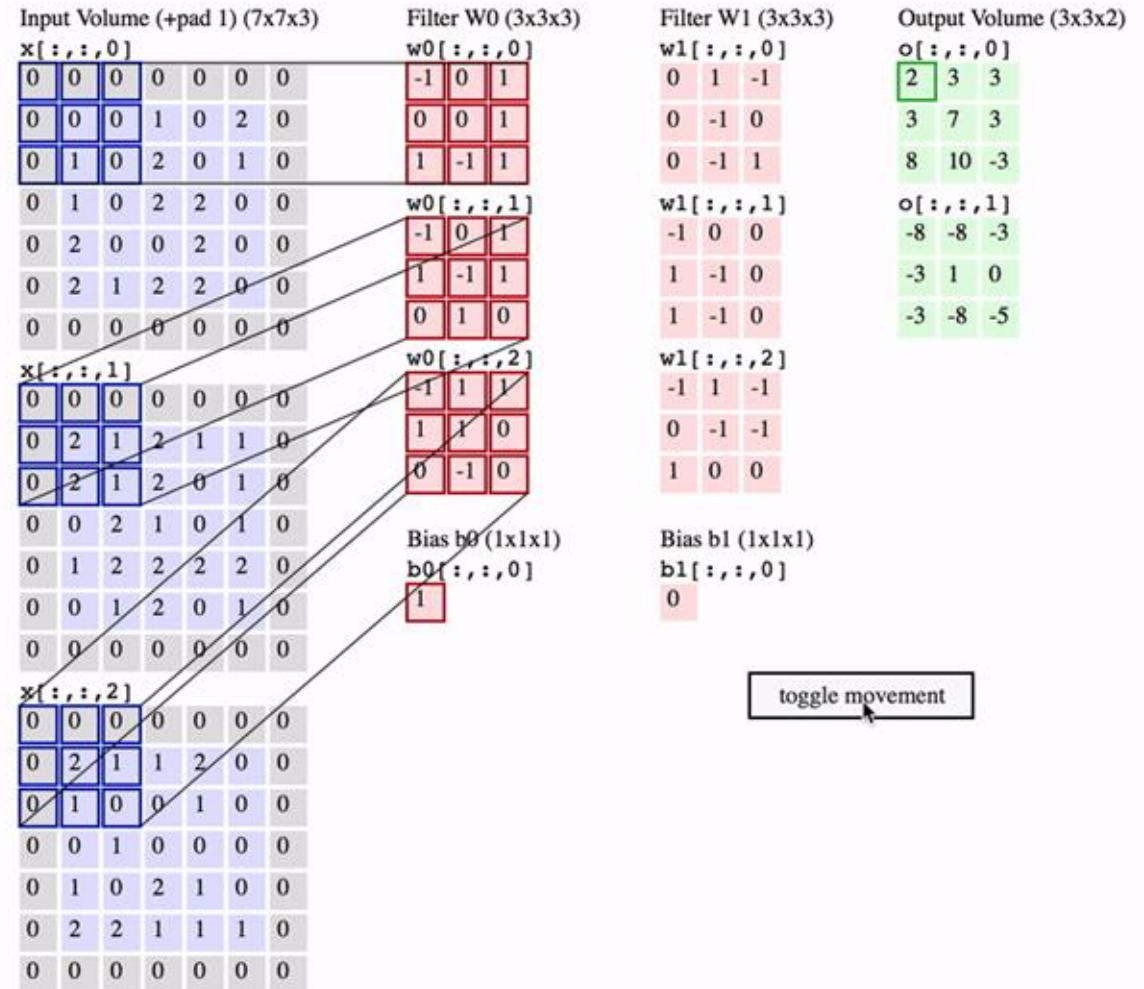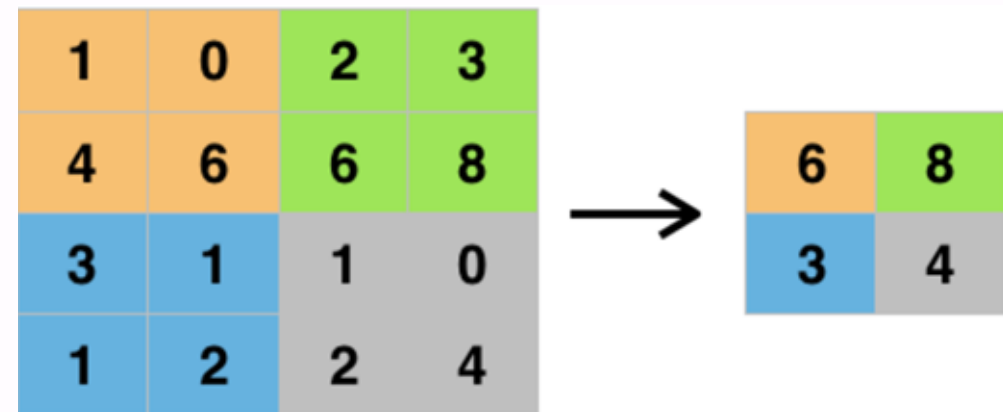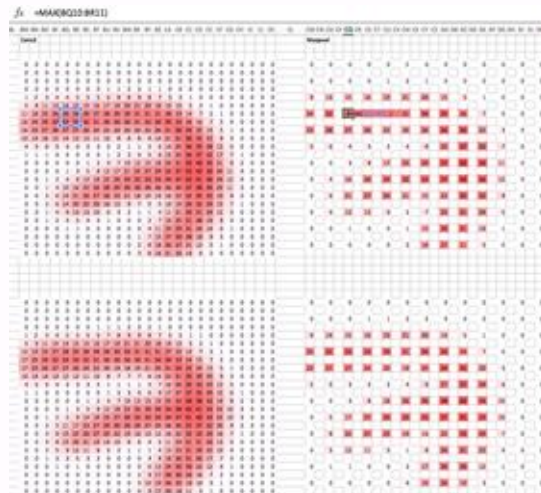
# Padding and Strides

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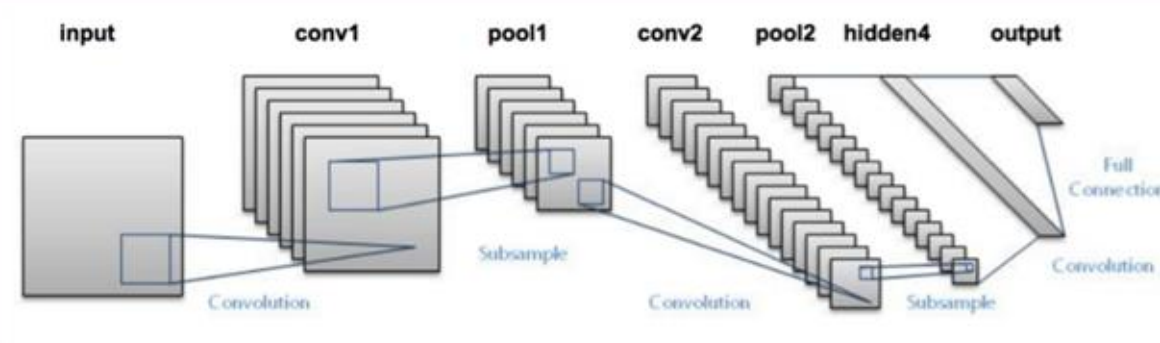- Because we typically use stride = pool width, the pooling output is aggregating over segments of the input.

# A Basic Image Labeling Architecture

## Convolution Layers Apply (Multiple) Filters
- Each filter is a matrix of trainable weights.
- As you move through the network, the number of features typically rises exponentially.
- More filters as you move along means it allows more permutations / combinations
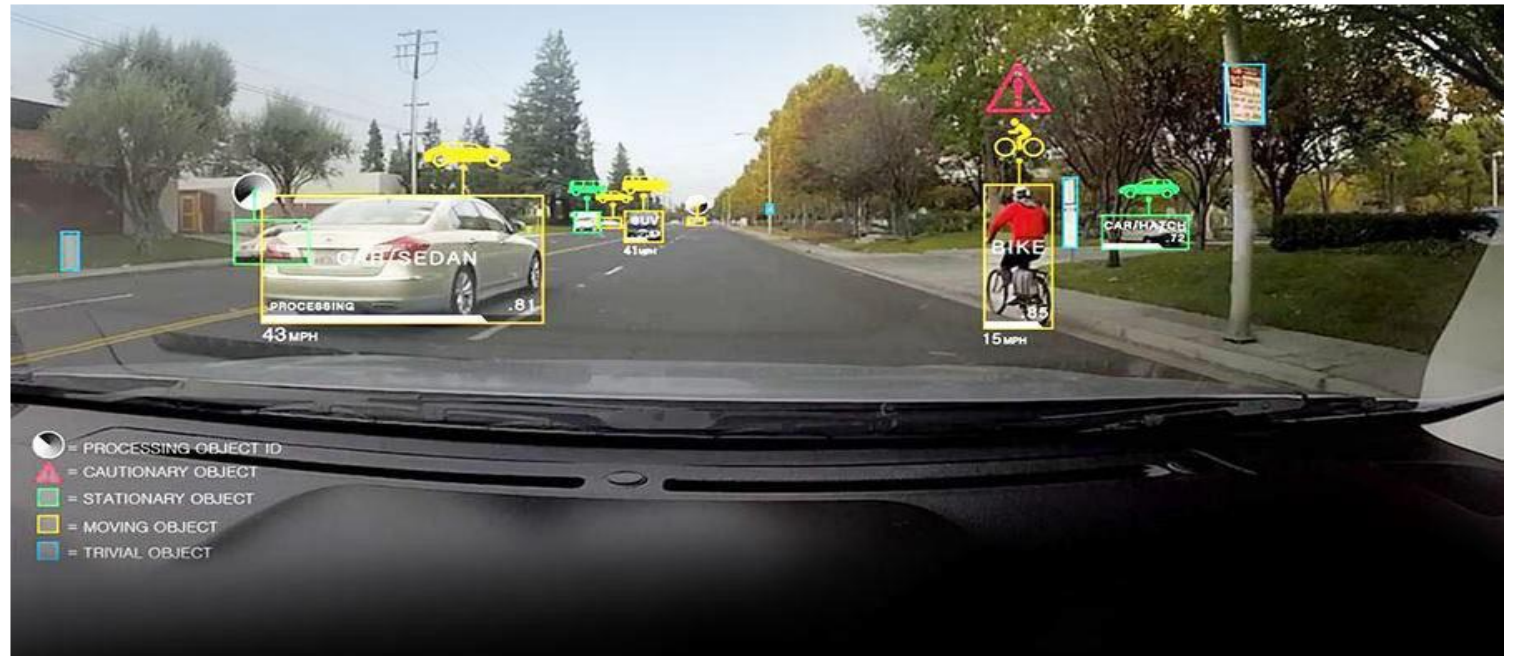
## Progressively Smaller Filter Maps
- Smaller filter map arises from the pooling steps, which means that each element of the final map distills features (high level features, derived from low level features, derived from raw pixels) derived from a larger segment of the original picture.

# Object Recognition Tasks
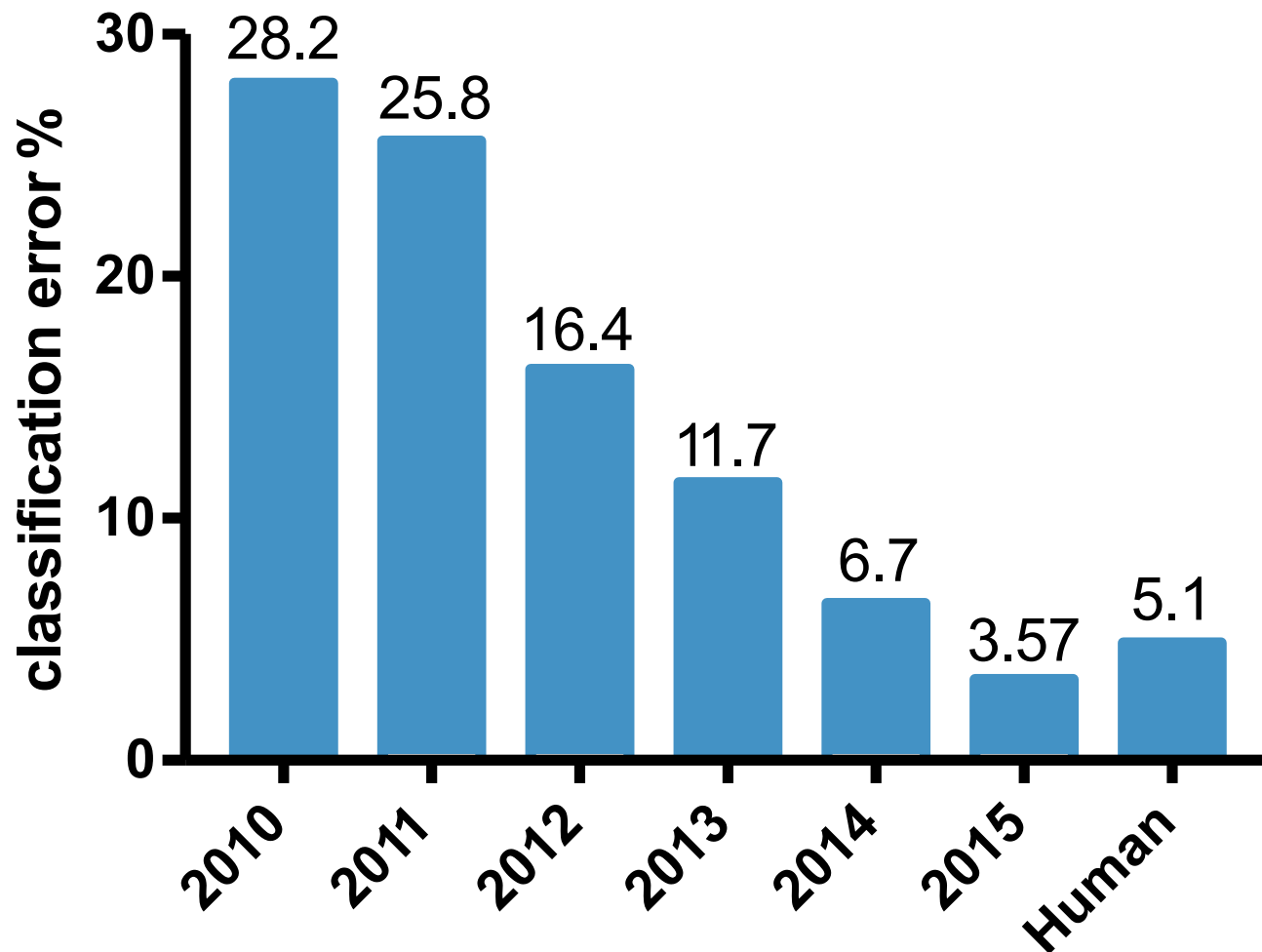


**CAT**, **DOG**, **DUCK**

# ImageNet Challenge



Classification task: "Top 5 error among 1000 categories": rate at which the model does not output correct label in top 5 predictions

# ImageNet Challenge: Classification Task



2012: AlexNet. First CNN to win.
- 8 layers, 61 million parameters

2013: ZFNet
- 8 layers, more filters
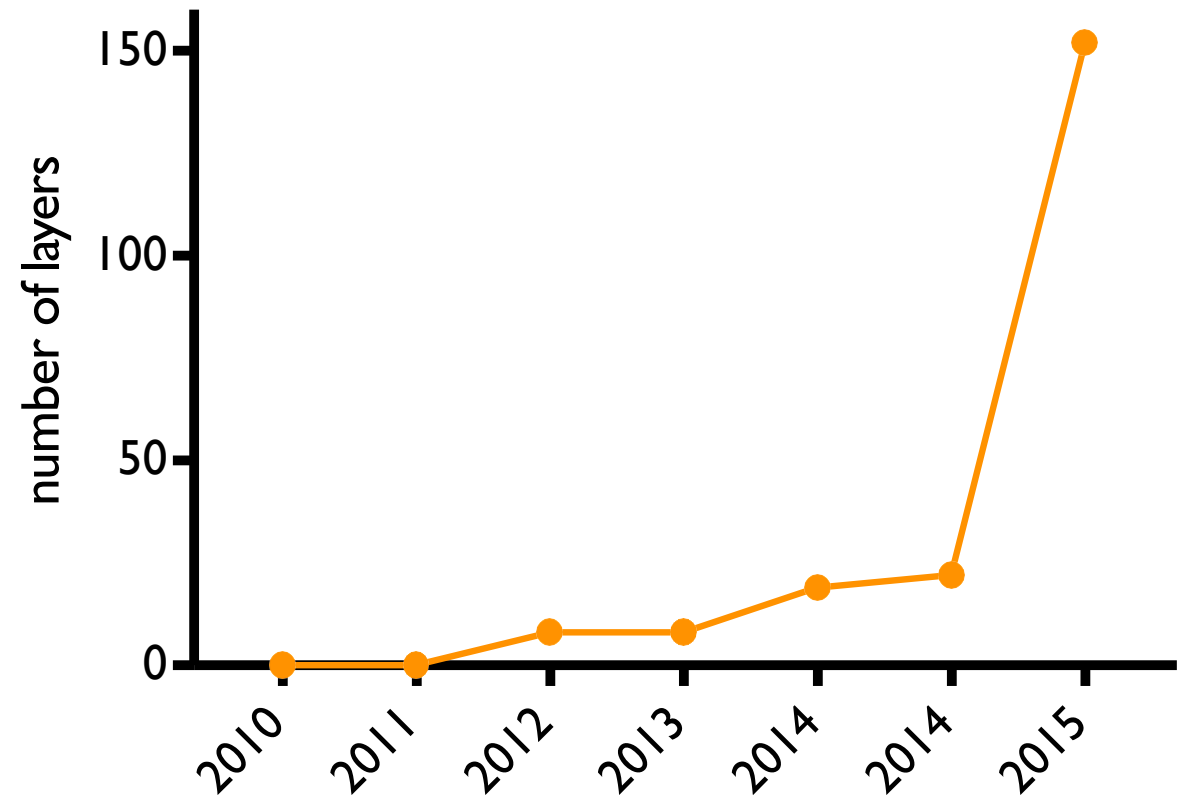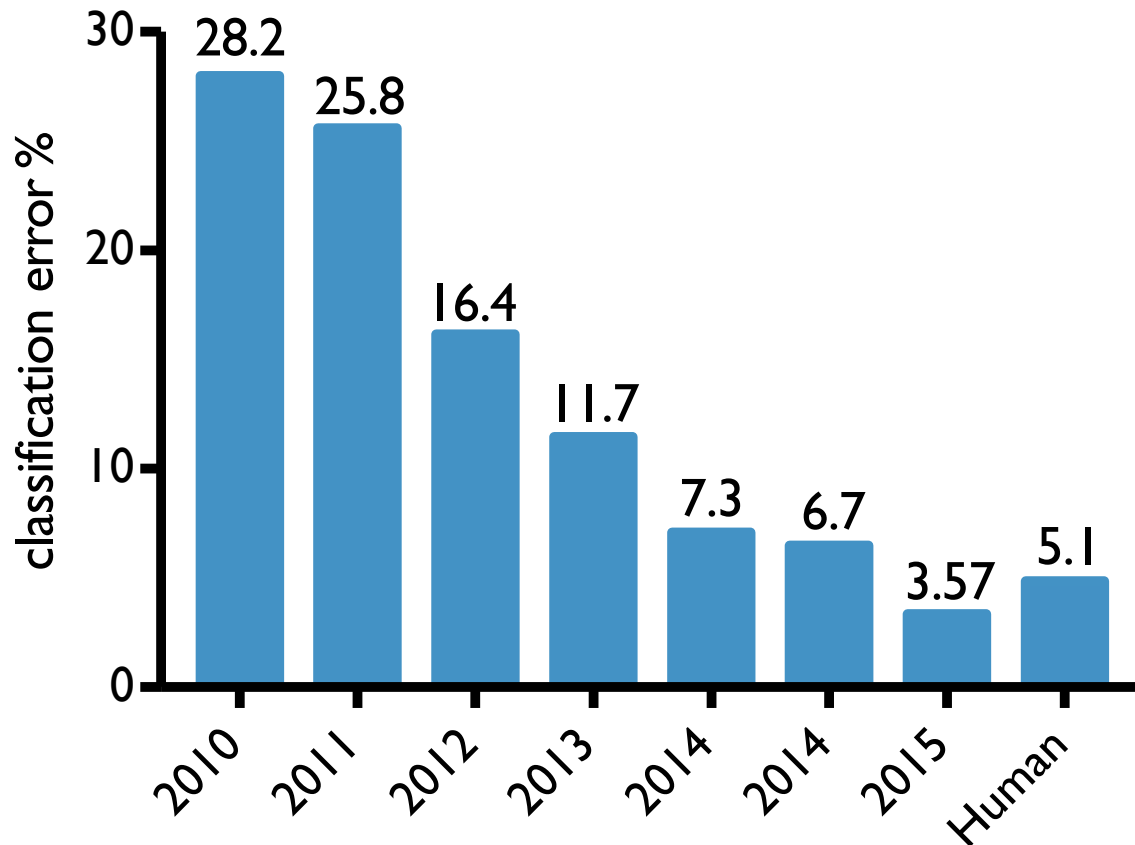
2014: VGG
- 19 layers

2014: GoogLeNet
- "Inception" modules
- 22 layers, 5 million parameters

2015: ResNet
- 152 layers

# ImageNet Challenge: Classification Task

# ImageNet Challenge Was Ended After 2017

**NewScientist**

Enter search keywords

Computer vision is ready for its next big test: seeing in 3D. The ImageNet Challenge, which has boosted the development of image-recognition algorithms, will be replaced by a new competition next year that aims to help robots see the world in all its depth.

Since 2010, researchers have trained image recognition algorithms on the ImageNet database, a go-to set of more than 14 million images hand-labelled with information about the objects they depict. The algorithms learn to classify the objects in the photos into different categories, such as house, steak or Alsatian. Almost all computer vision systems are trained like this before being fine-tuned on a more specific set of images for different tasks.

Every year, participants in the ImageNet Large Scale Visual Recognition Challenge try to code algorithms that can categorise these images with as few errors as possible. Seven years ago, this was a difficult task, but now computer vision is great at categorising images.

# Shifted Focus to Harder Problems

# Shifted Focus to Harder Problems
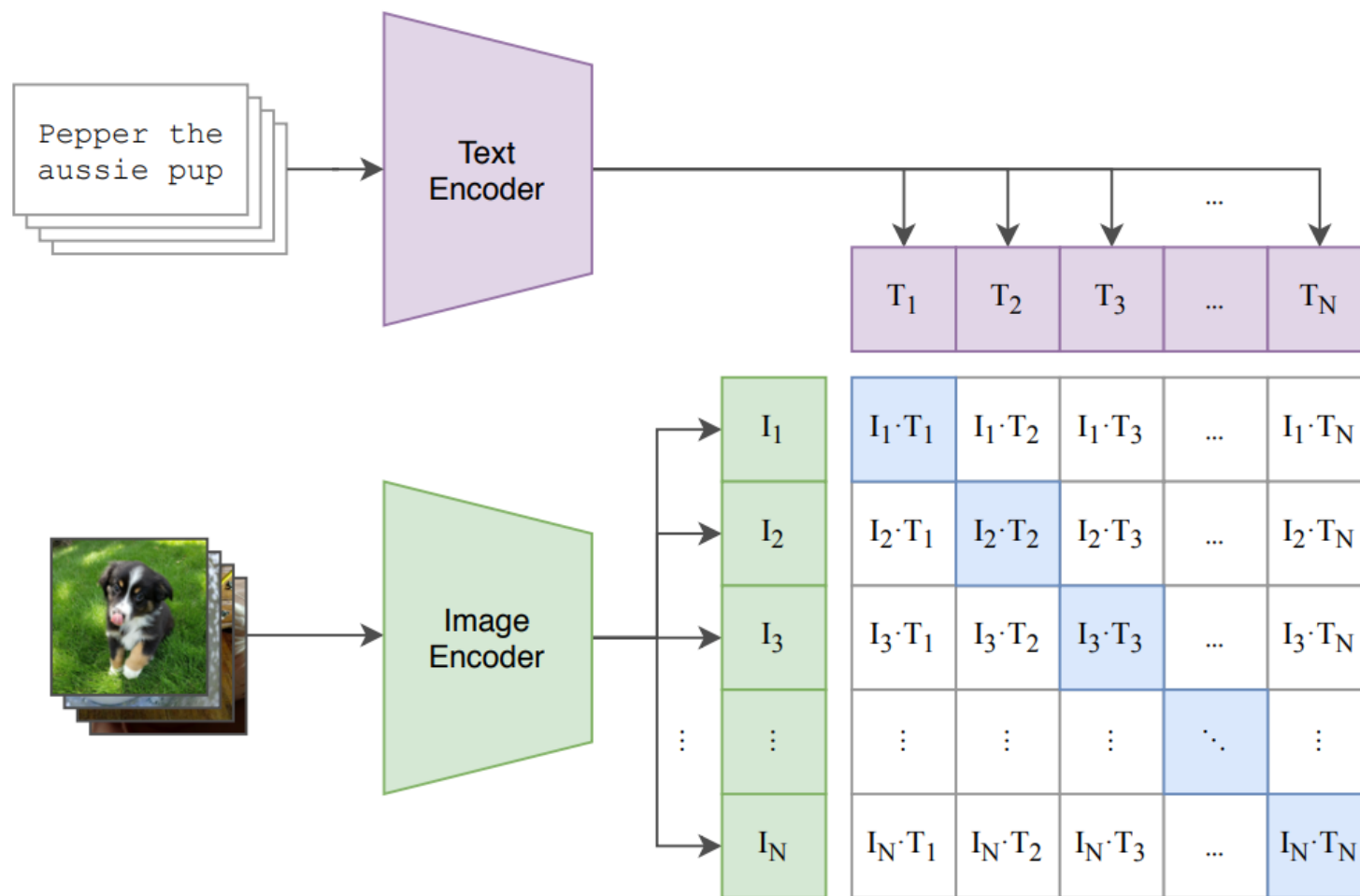


© Gordon Burtch, 2026

# Visual Question and Answer (VQA) challenge



VISUAL QUESTION ANSWERING (VQA) CHALLENGE: ACCURACY
Source: VQA Challenge, 2020 | Chart: 2021 AI Index Report

# Contrastive Language-Image Pre-training (CLIP)

# Contrastive Loss

The goal here is *not* prediction. The goal is to train a model that can encode inputs into numeric vector representations in a shared embedding space, such that 'similar' vectors are conceptually related, and 'dissimilar' vectors are conceptually unrelated.

We have a special loss function that we minimize for this purpose, where it is once again a 2-part loss, depending on whether the 'pair' of inputs is positive (related) or negative (unrelated).
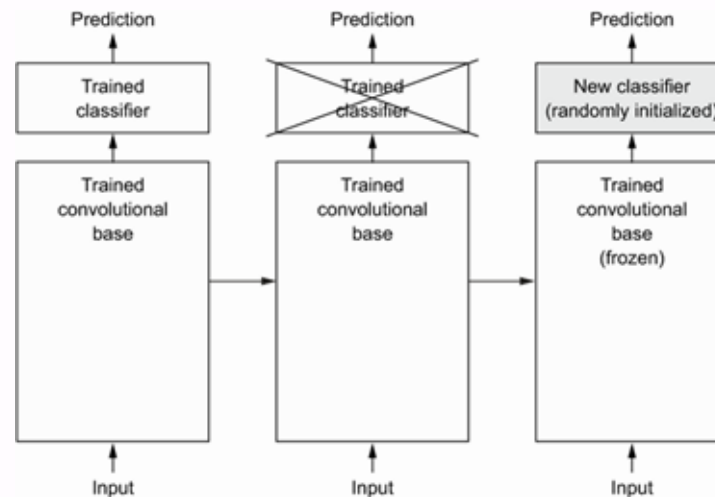
$$L = y \cdot D^2 + (1 - y) \cdot \max(0, m - D)^2$$

# 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.
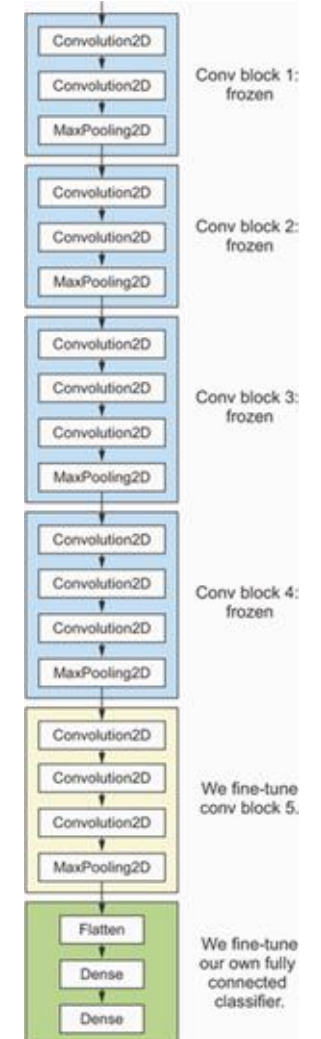
# Pre-Trained Models: Fine Tuning

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
```

# What Are People Tackling Now?
# 3D Scene Reconstruction

# Self-Evaluation & Peer Call-outs

# Questions?