

Lecture 6: Perception II

Professor Katie Driggs-Campbell

February 1, 2024

ECE484: Principles of Safe Autonomy

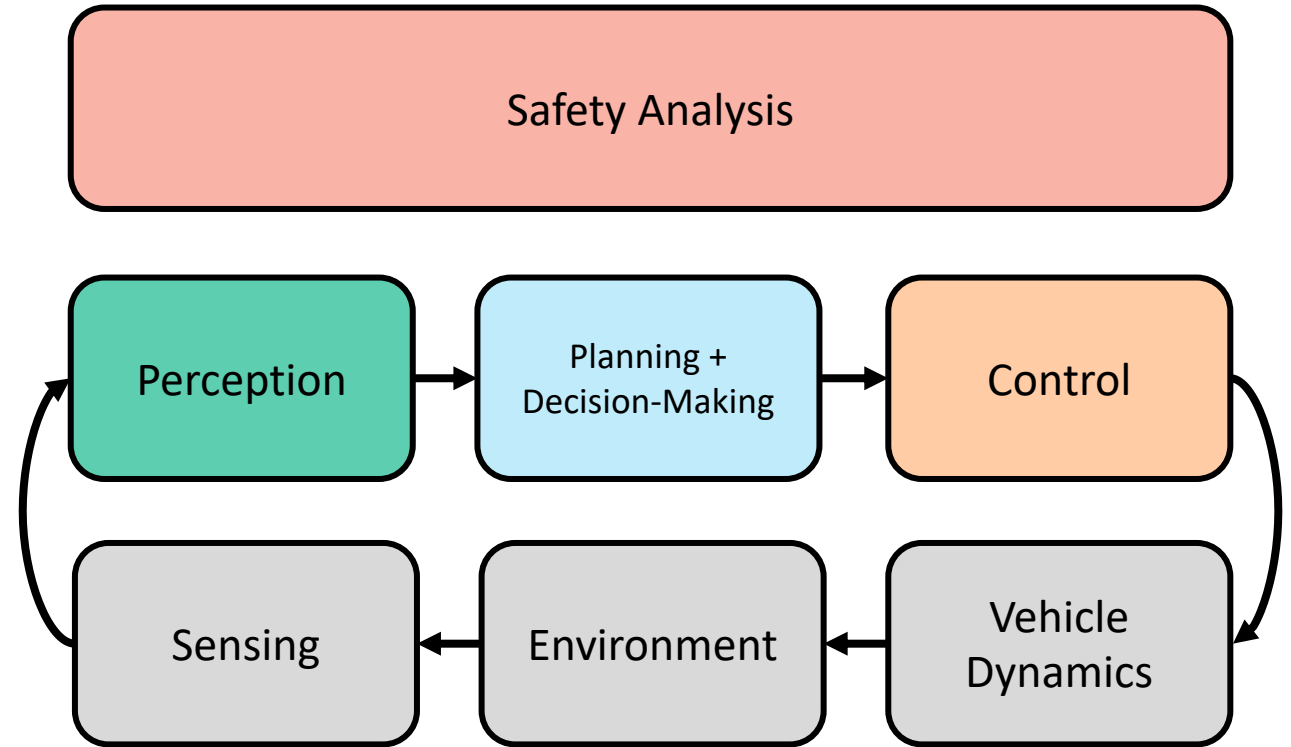


Administrivia

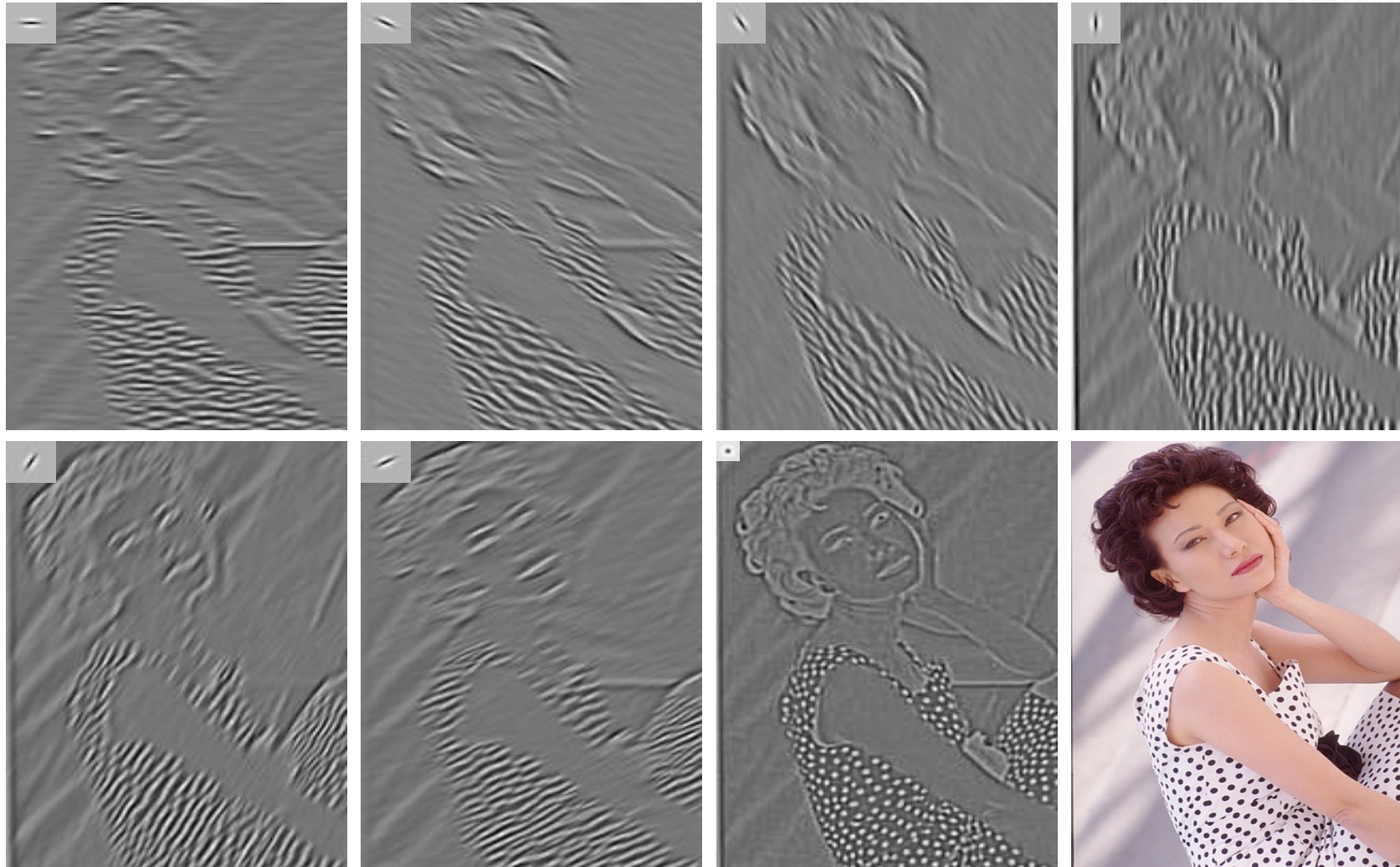
- Repeat pop quiz comments
- MP1 released **next** week (2/9)
- Upcoming due dates:
 - HW0 and MP0 due Friday 2/9
 - HW1 and MP1 due Friday 2/23



Autonomous GEM Vehicle

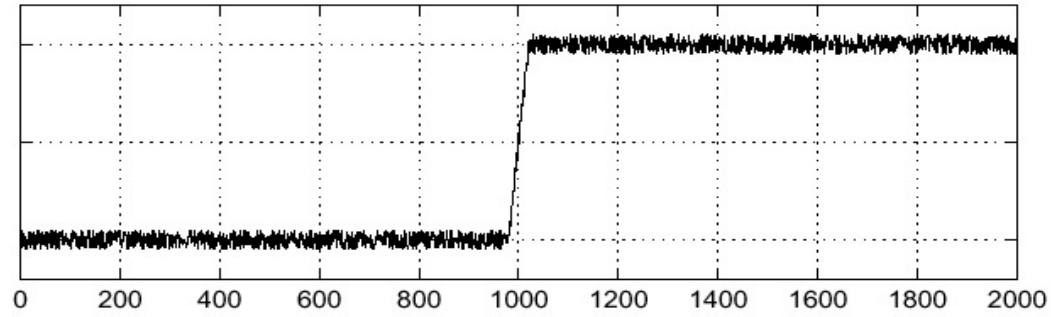


Filter Outputs

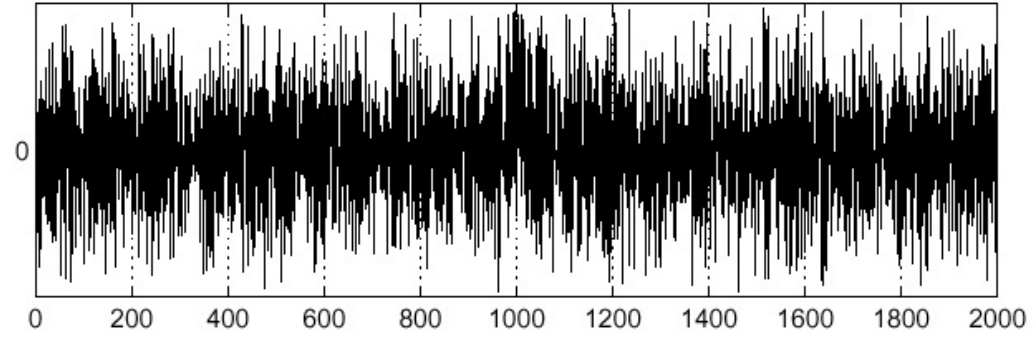


Impact of Noise

f

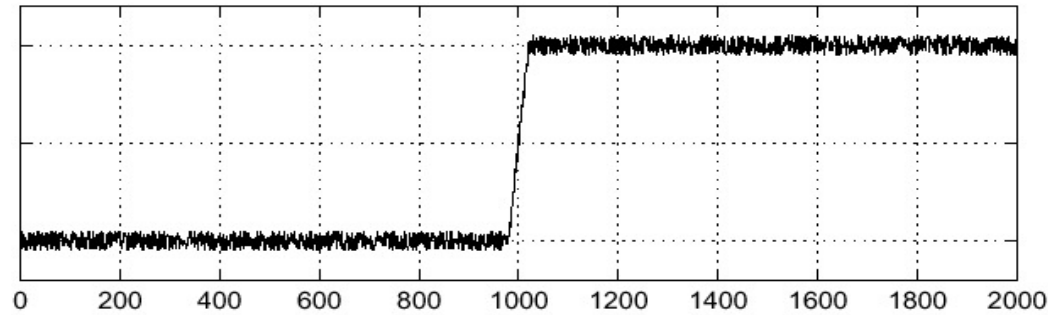


$\frac{d}{dx} f(x)$

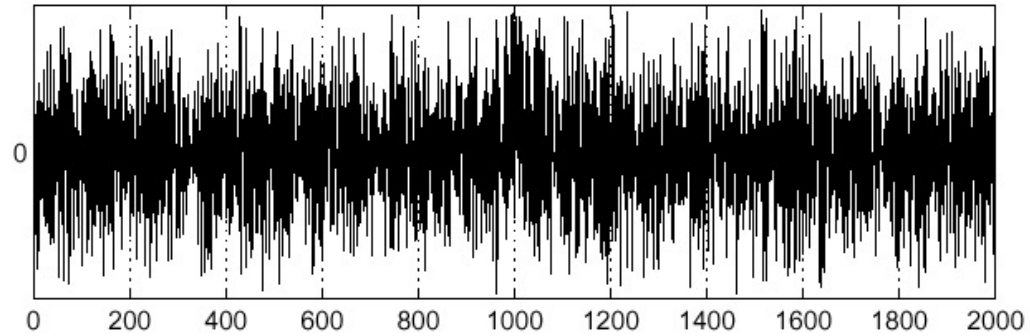


Impact of Noise

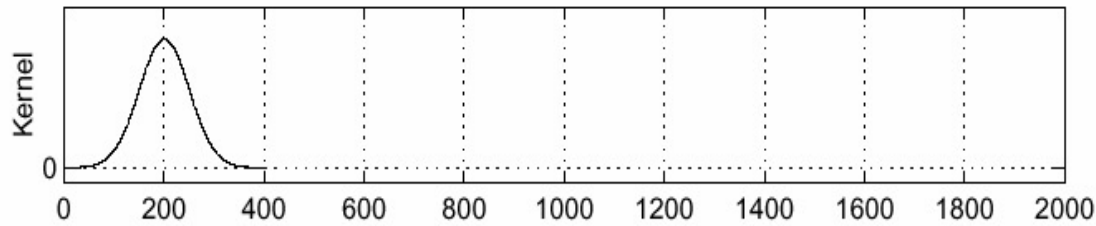
f



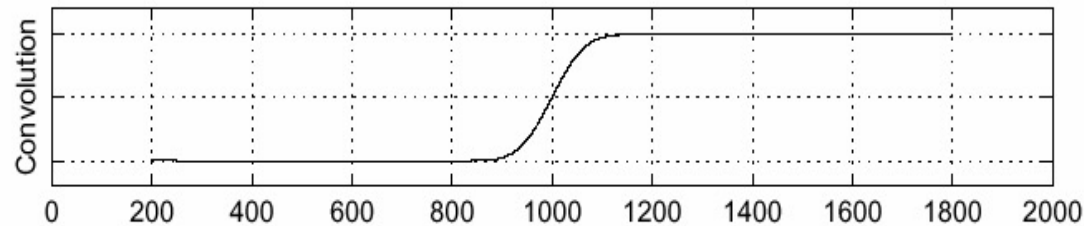
$\frac{d}{dx} f(x)$



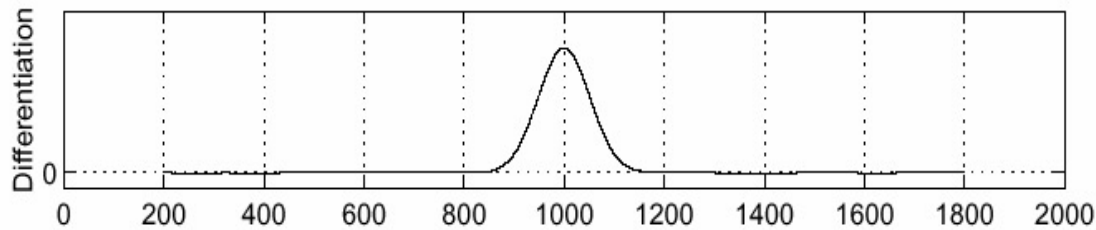
g



$f * g$

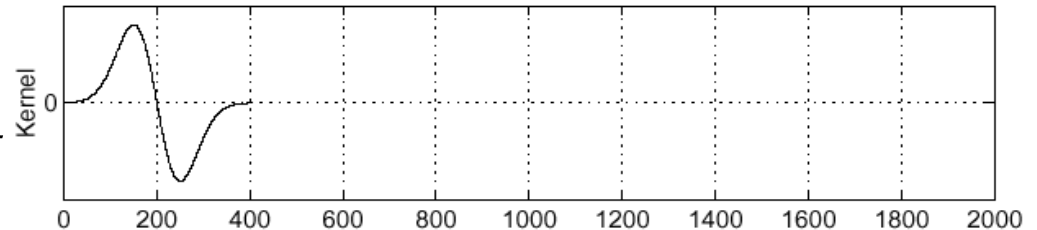


$\frac{d}{dx} (f * g)$

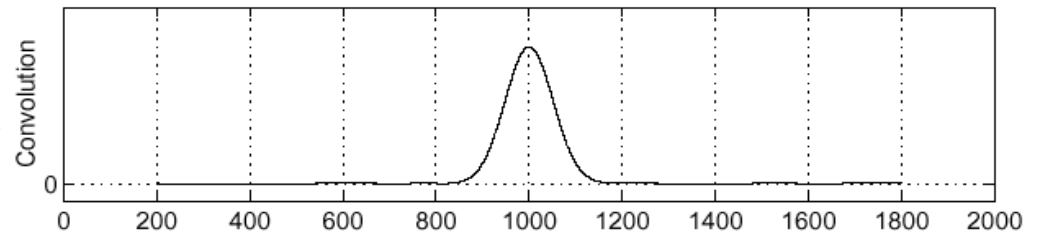


$$\frac{d}{dx} (f * g) = f * \frac{d}{dx} g$$

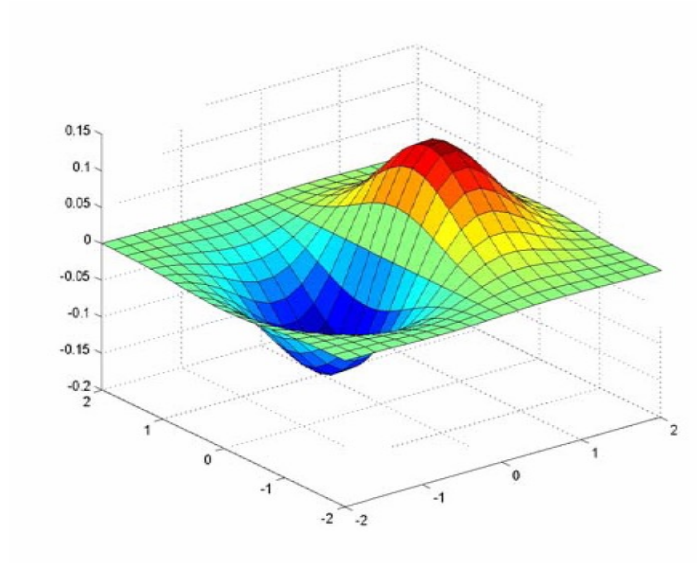
$\frac{d}{dx} g$



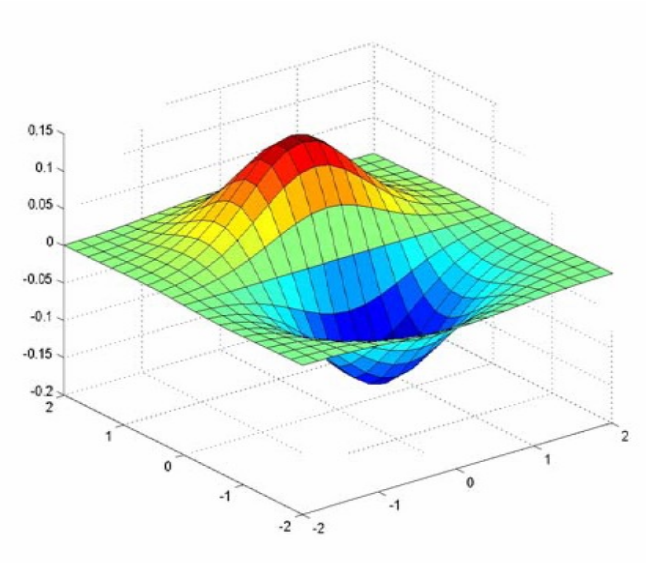
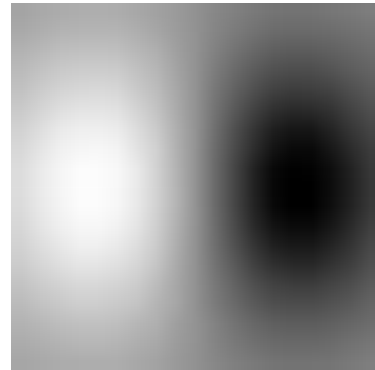
$f * \frac{d}{dx} g$



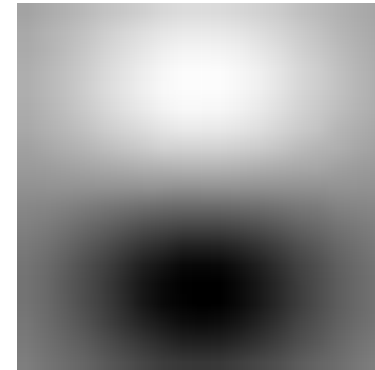
Derivative of Gaussian filters



x-direction



y-direction



Building an edge detector

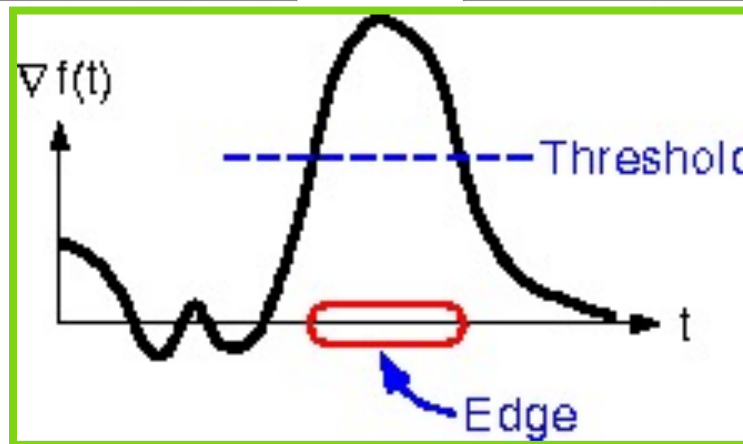
Original Image



Edge Image



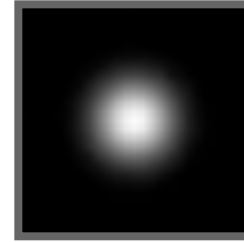
$$\|\nabla f\| = \sqrt{\left(\frac{\partial f}{\partial x}\right)^2 + \left(\frac{\partial f}{\partial y}\right)^2}$$



Review: smoothing vs. derivative filters

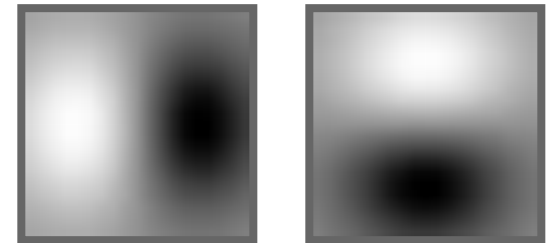
- Smoothing filters

- Gaussian: remove “high-frequency” components; “low-pass” filter
- What should the values sum to?
 - One: constant regions are not affected by the filter



- Derivative filters

- Derivatives of Gaussian
- What should the values sum to?
 - Zero: no response in constant regions



Object Recognition

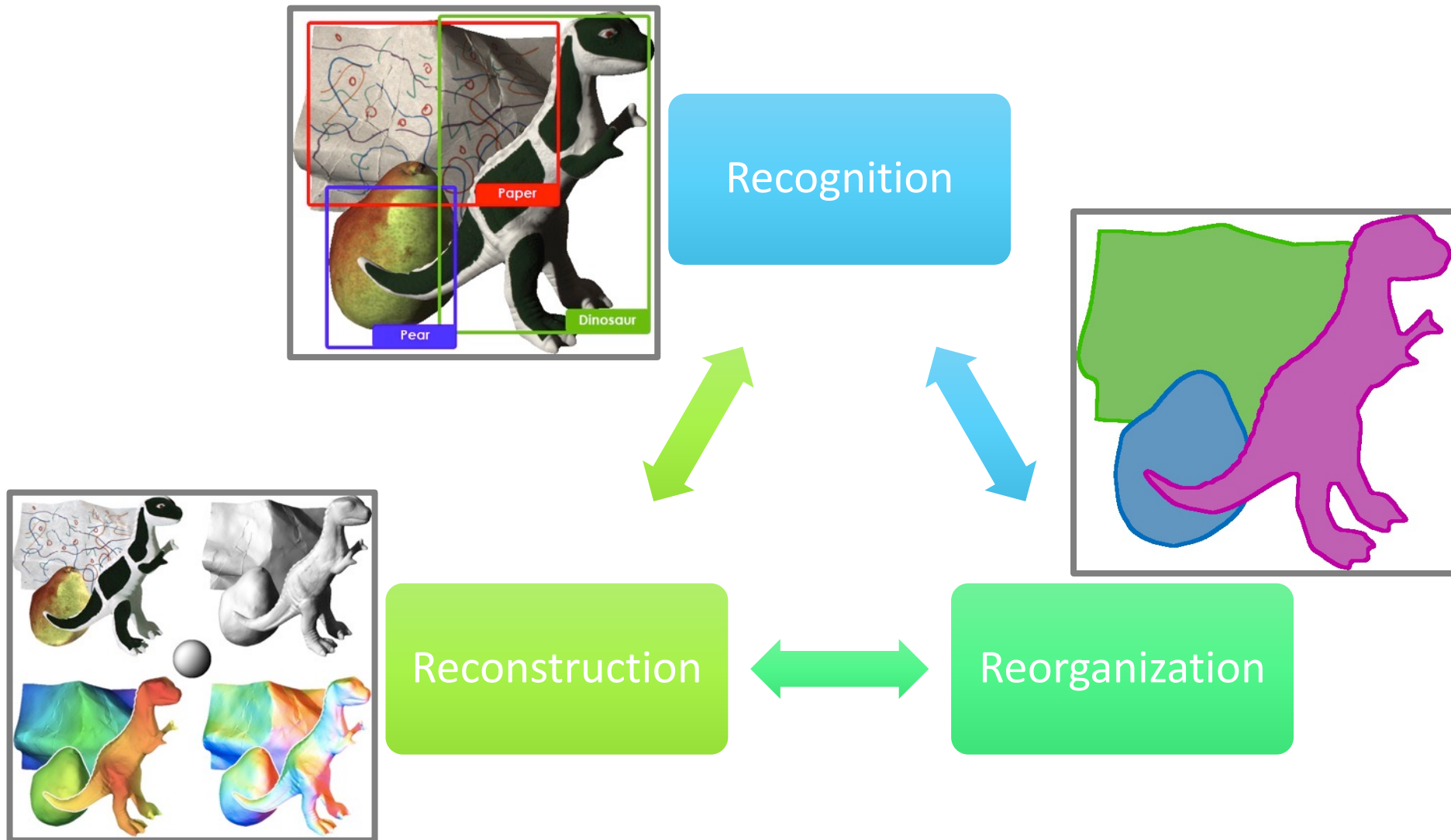


Today's Plan

- Computer vision overview
- Object recognition
 - Feature representations
 - Classification
- (Convolutional) Neural Networks



The Three R's of Computer Vision





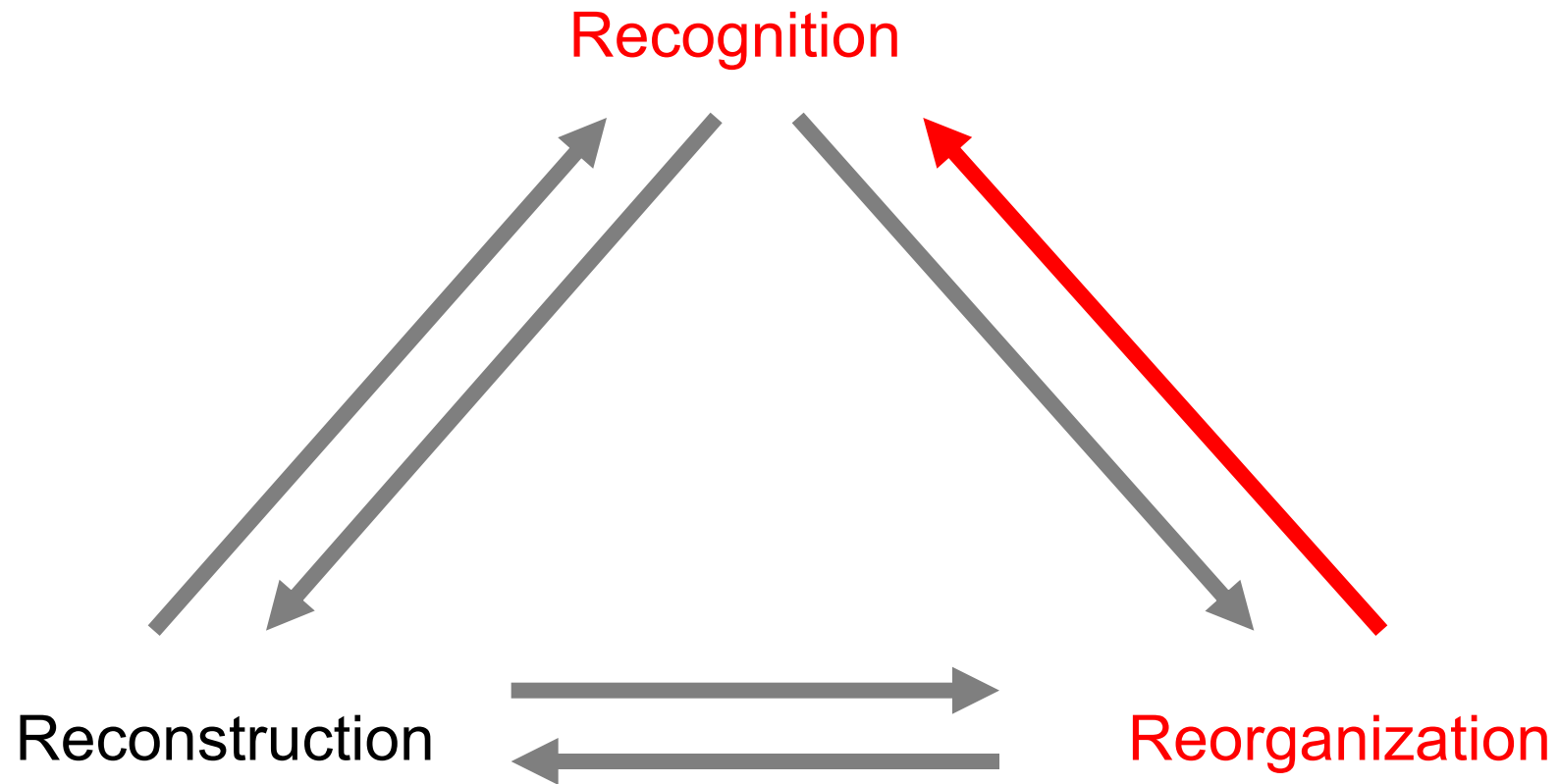
What we would like to infer...



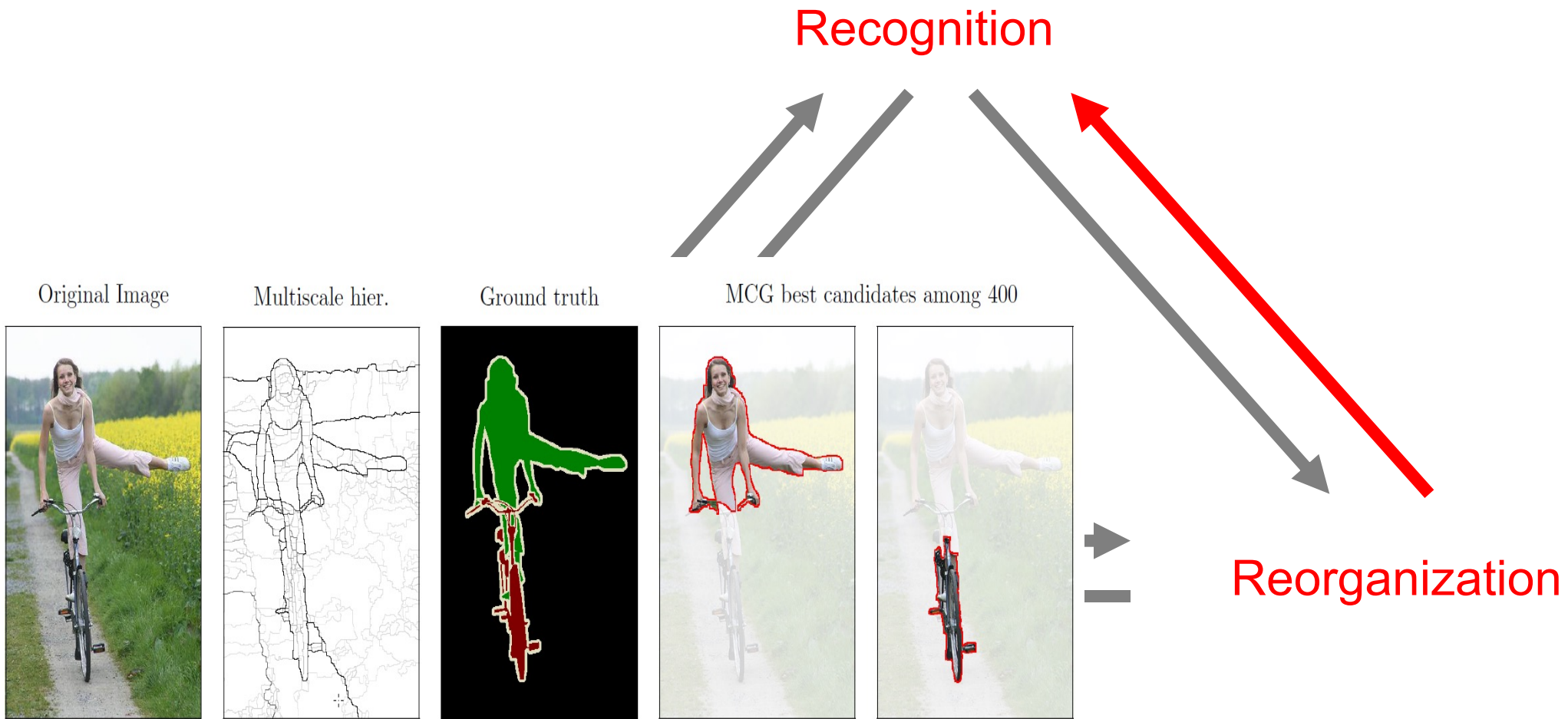
Will person B put some money into Person C's tip bag?



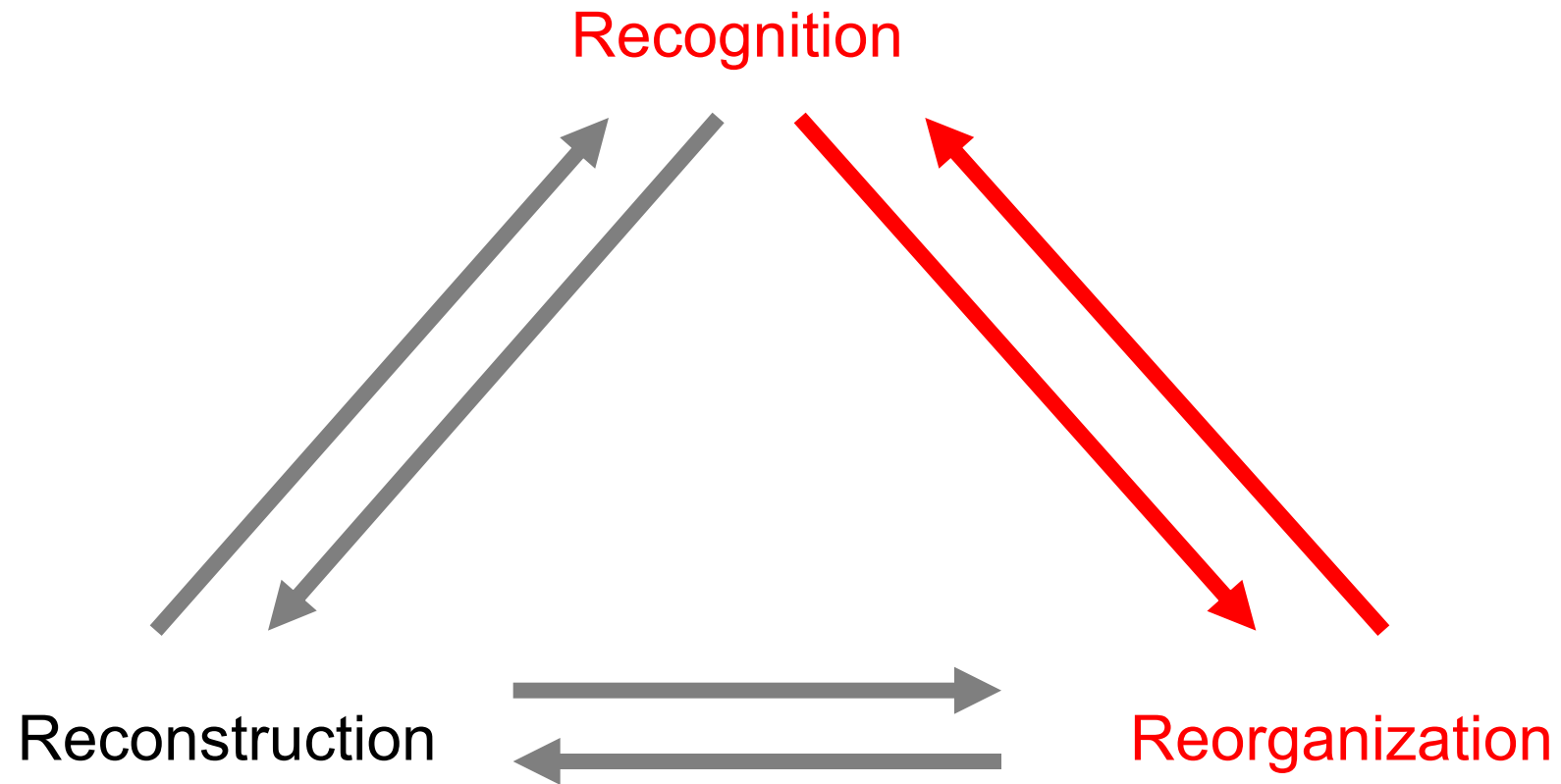
The Three R's of Vision



The Three R's of Vision



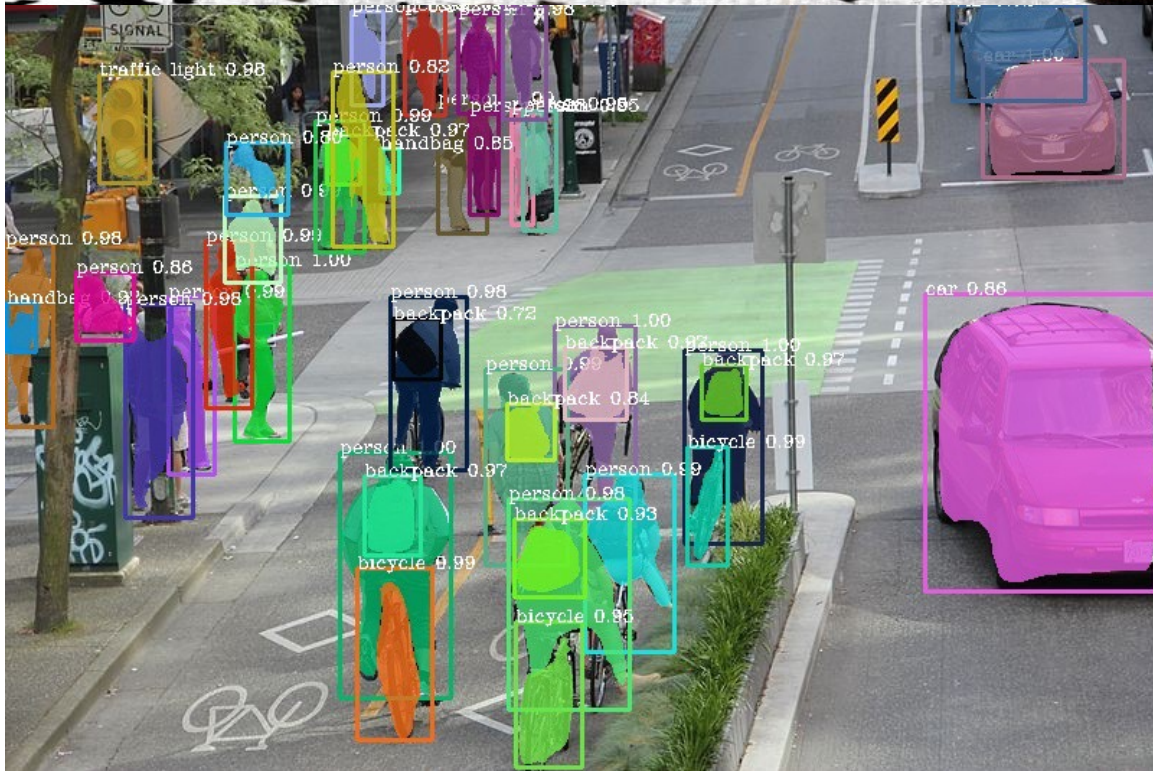
How about the other direction...





other direction...

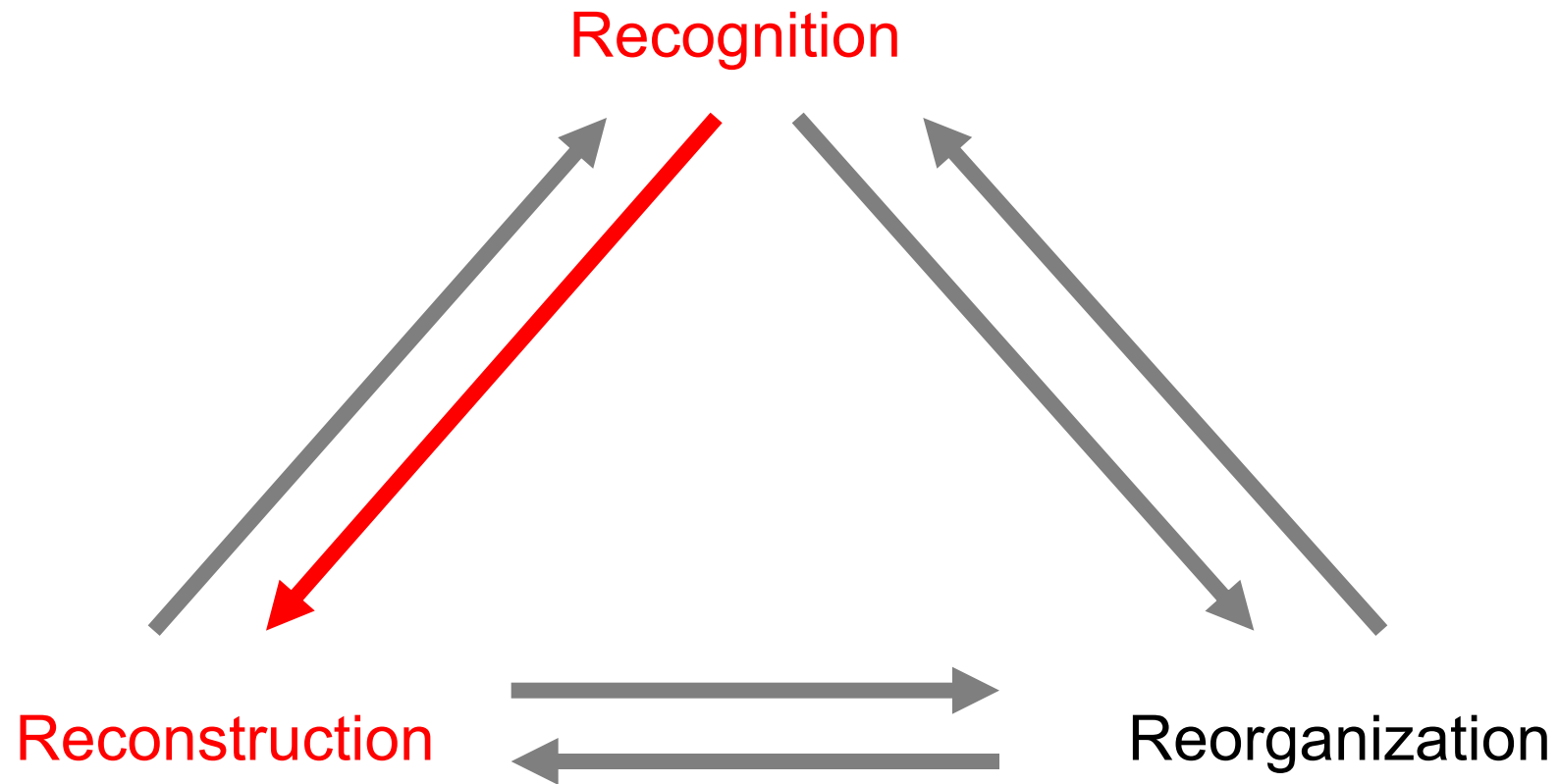
cognition



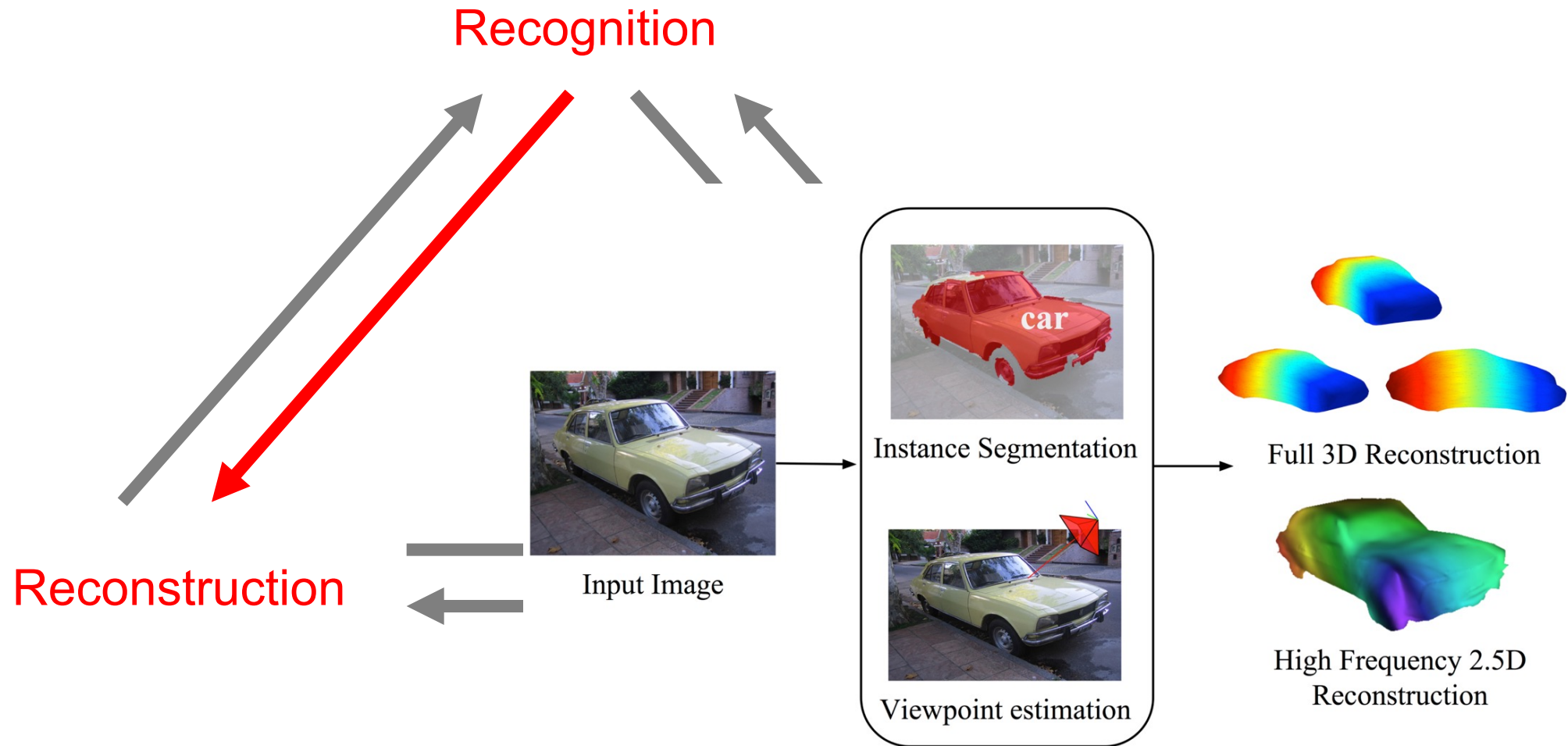
Reorganization



The Three R's of Vision



The Three R's of Vision



What is image recognition?

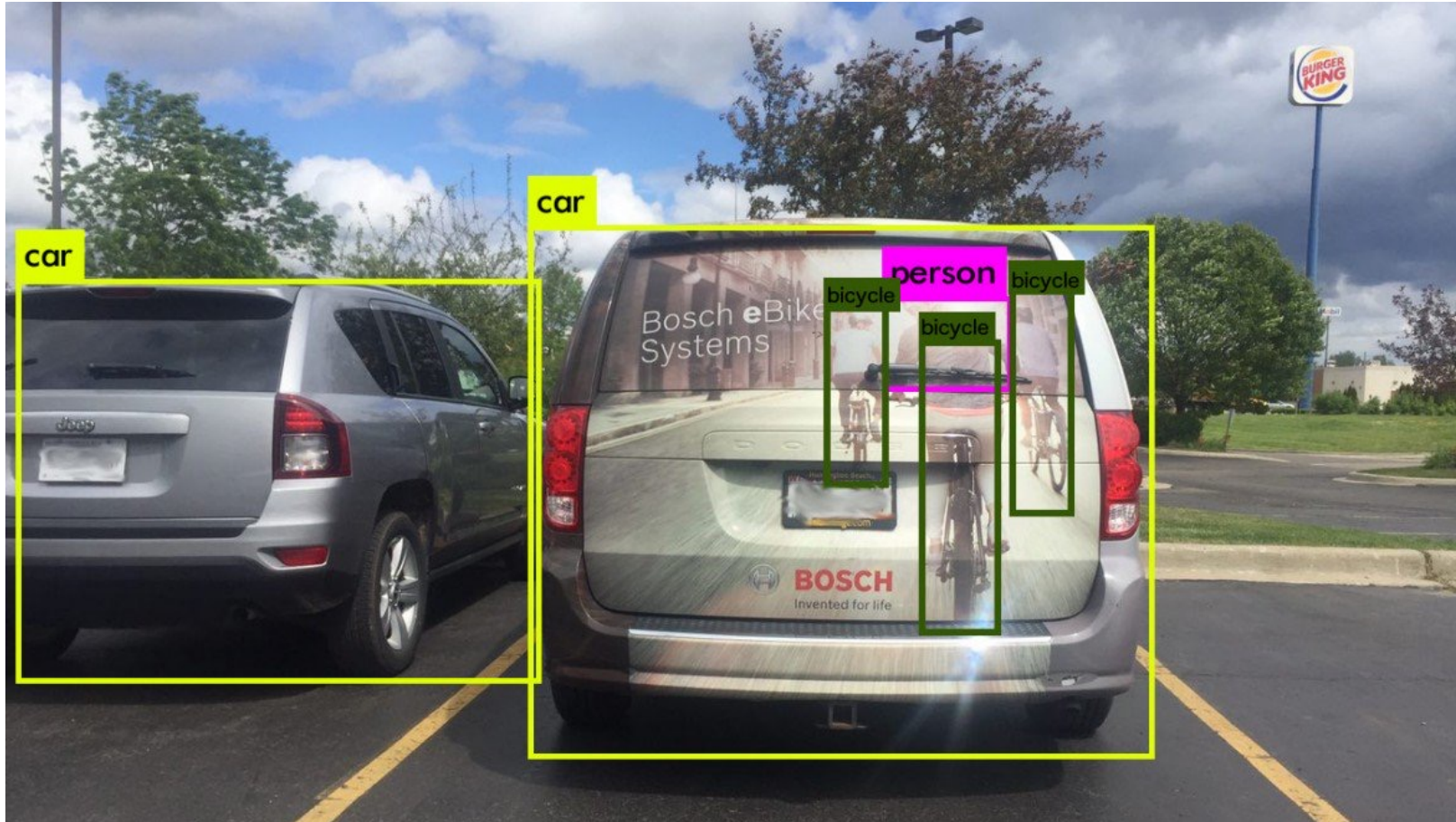


Image Recognition



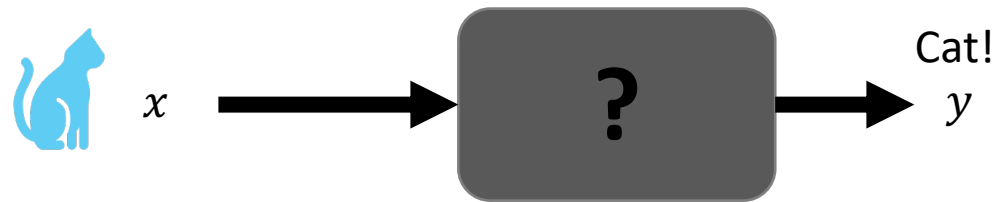
$f(\text{apple}) = \text{"apple"}$

$f(\text{tomato}) = \text{"tomato"}$

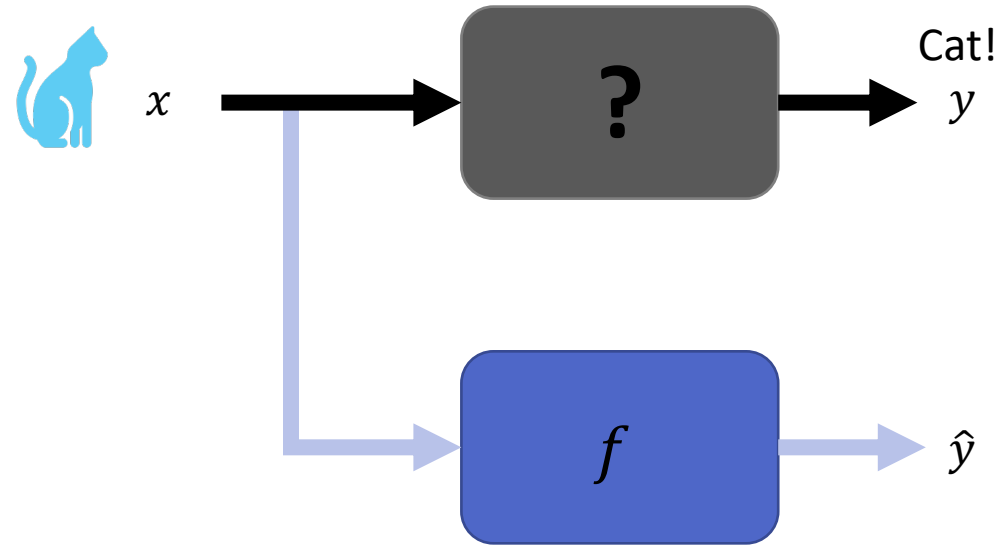
$f(\text{cow}) = \text{"cow"}$



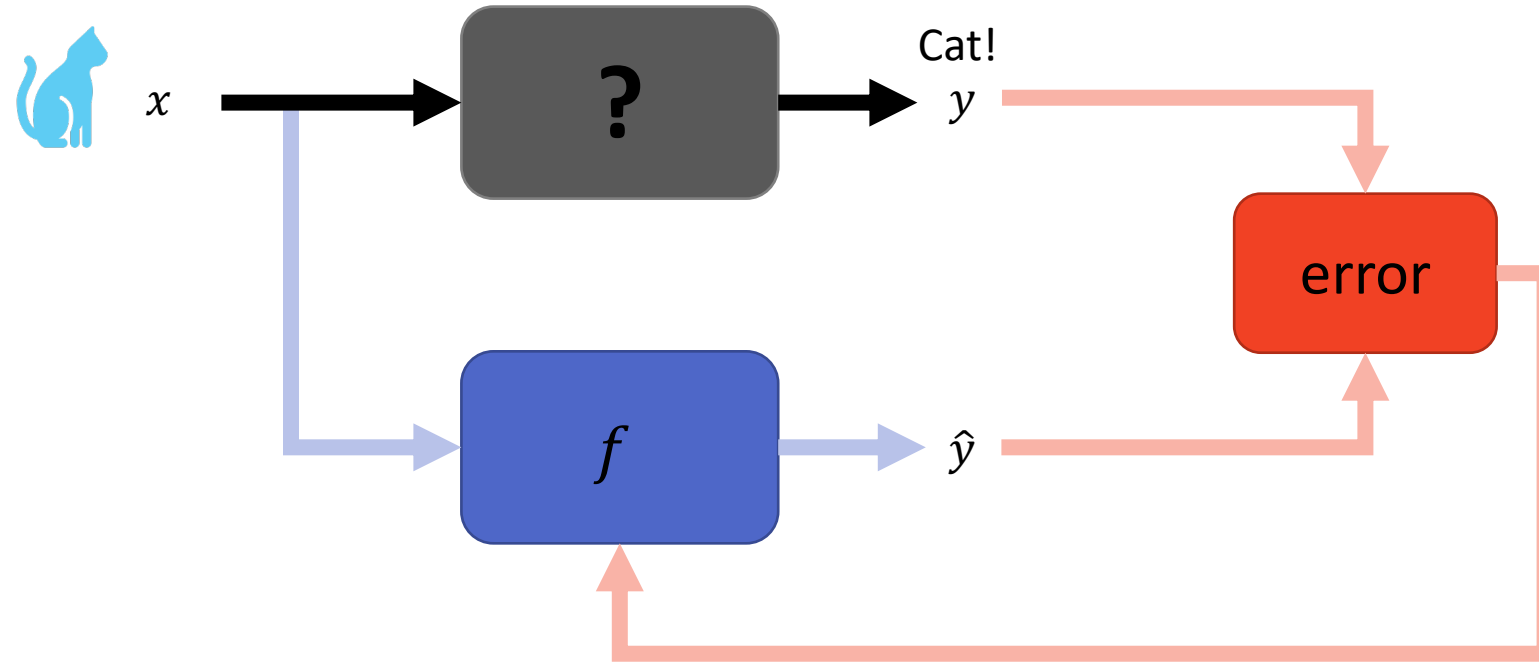
Learning Models



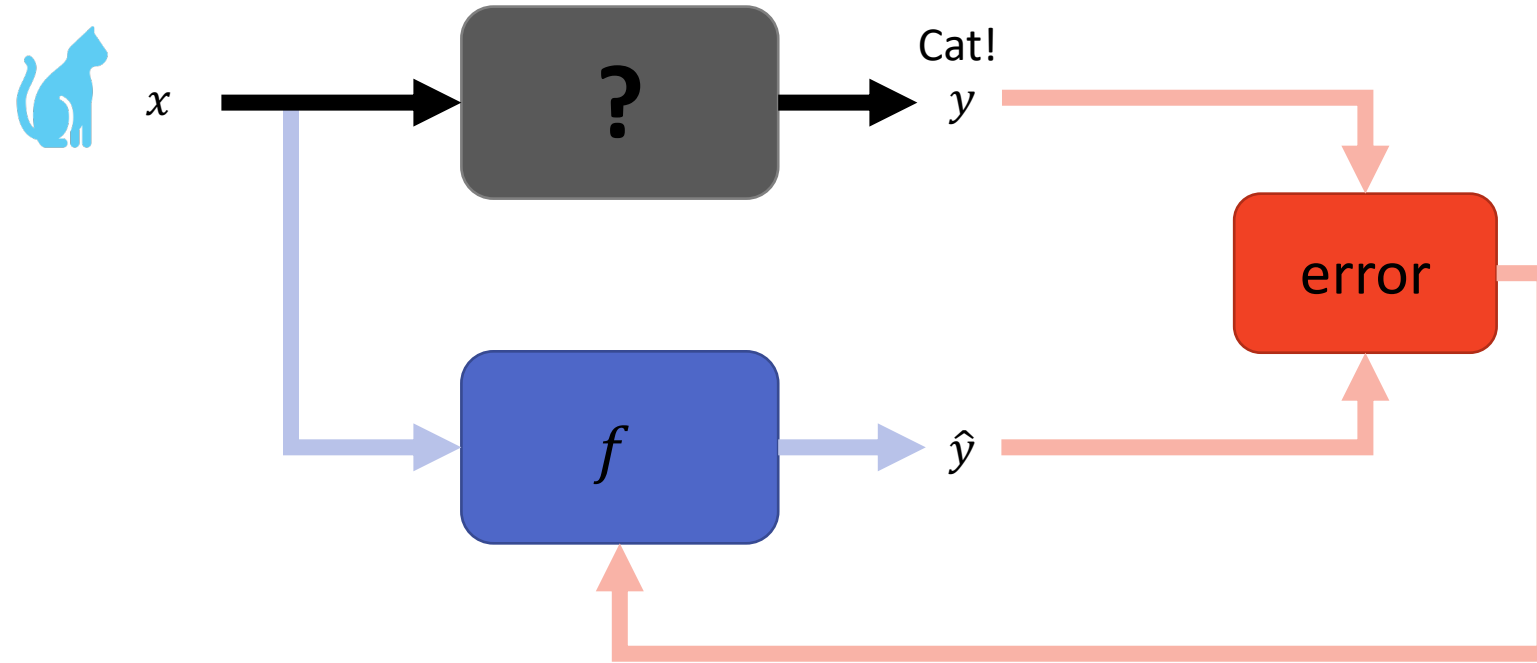
Learning Models



Learning Models



Learning Models



Learned models (like neural networks) are good when:

- Your system needs to learn and adapt
- Original is highly nonlinear / multi-variable
- Physics / model-based approaches are not available or are too computationally expensive



Training



Image Features



Training



Learned model

Training Labels



Training



Image Features



Training



Learned model

Training Labels



Testing



Image Features



Prediction

Learned model



Naïve classification

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

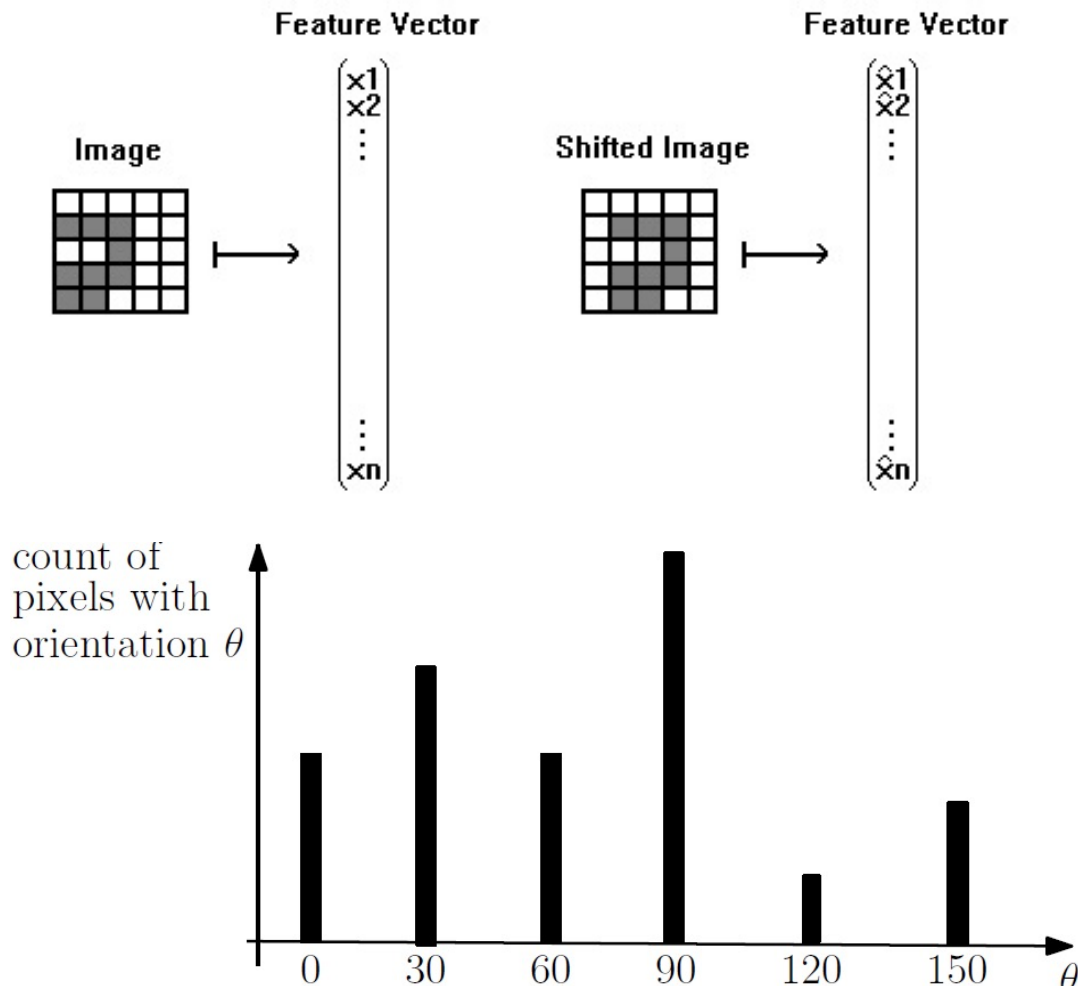
Fig. 4. Size-normalized examples from the MNIST database.



Classification in Feature Space



Naïve Features \rightarrow Orientation Histograms

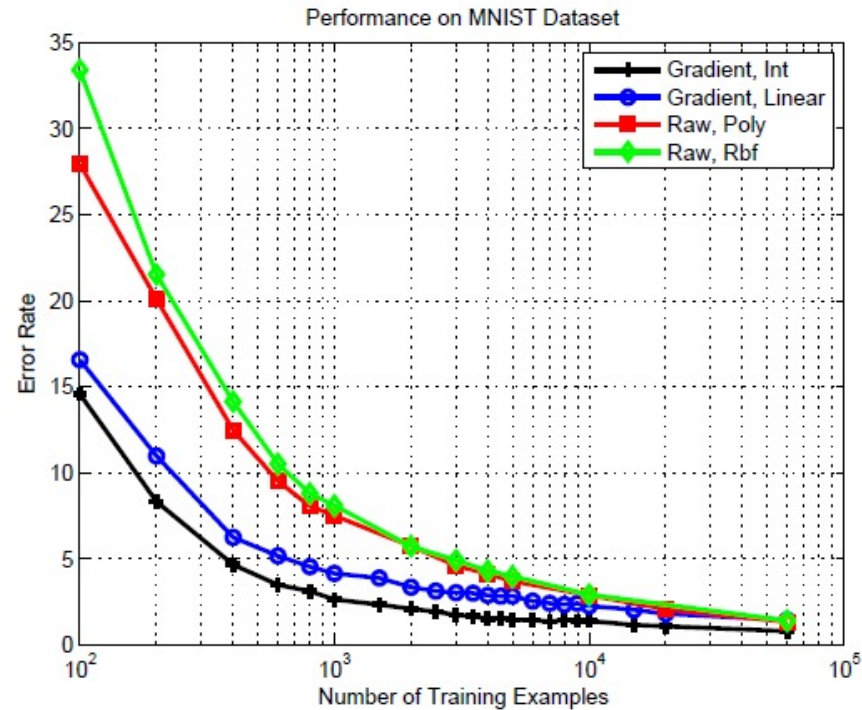


- Orientation histograms can be computed on blocks of pixels, so we can obtain tolerance to small shifts of a part of the object
- For gray-scale images of 3d objects, the process of computing orientations, gives partial invariance to illumination changes
- Small deformations when the orientation of a part changes only by a little causes no change in the histogram



Histogram of Oriented Gradients

Error rates vs. training examples



Misclassifications

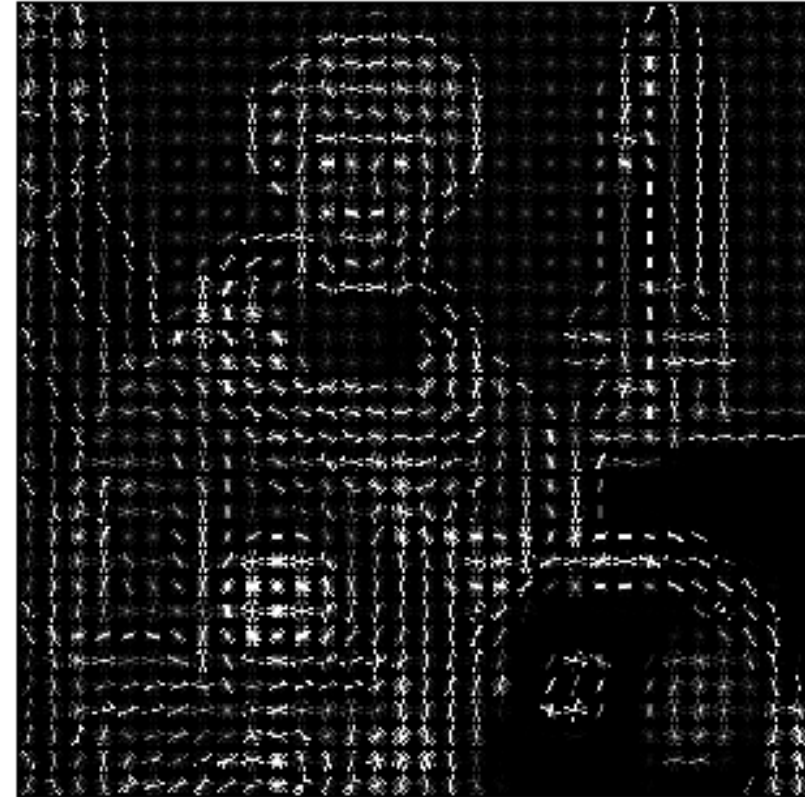


Histogram of Oriented Gradients

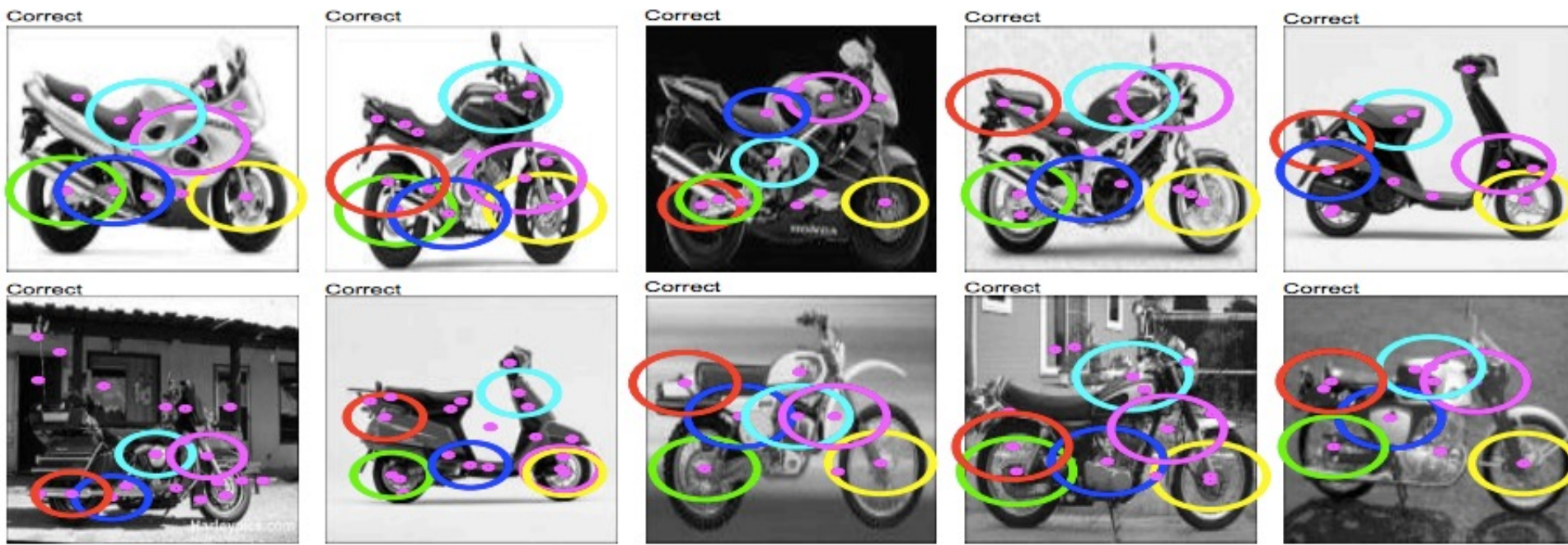
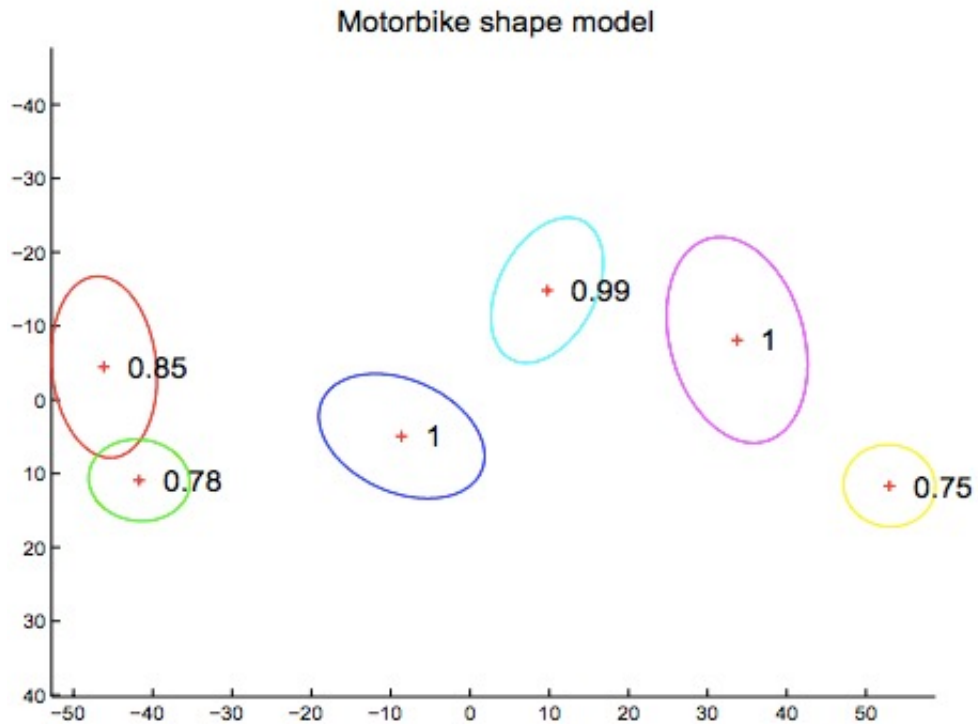
Input image



Histogram of Oriented Gradients



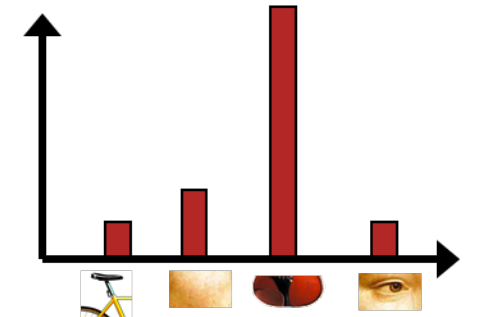
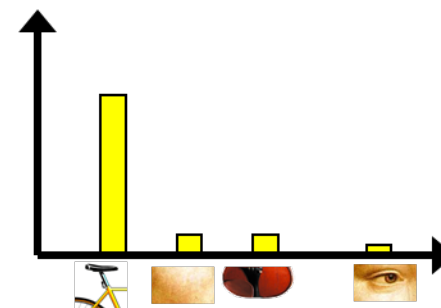
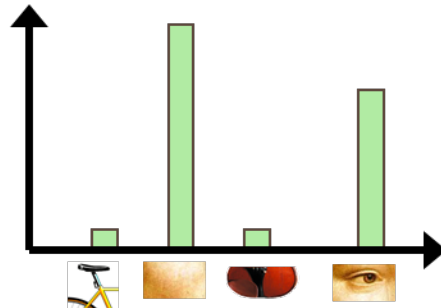
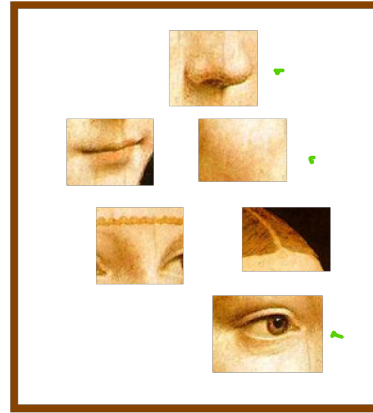
Part-based models



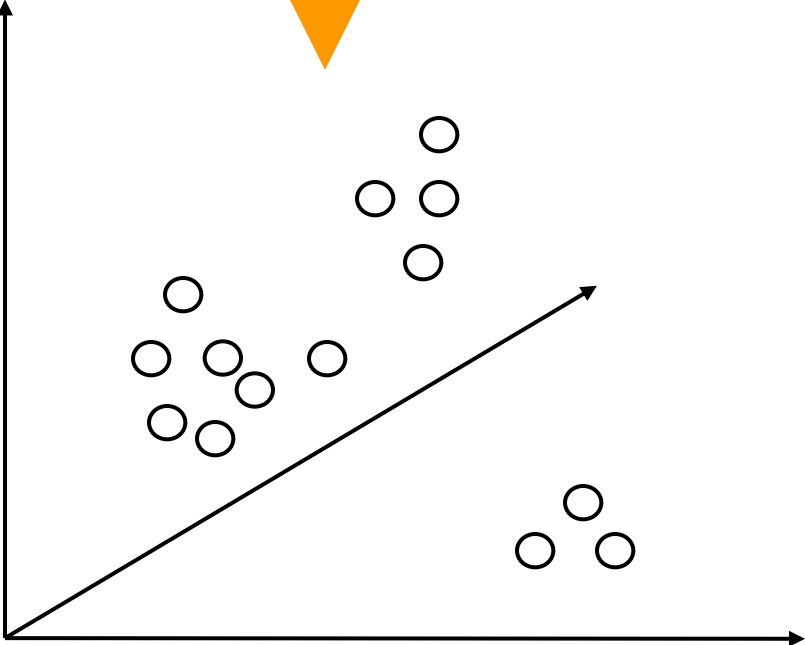
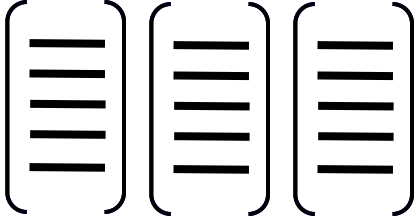
Weber, Welling &
Perona (2000),
Fergus, Perona &
Zisserman (2003)

Bag of features

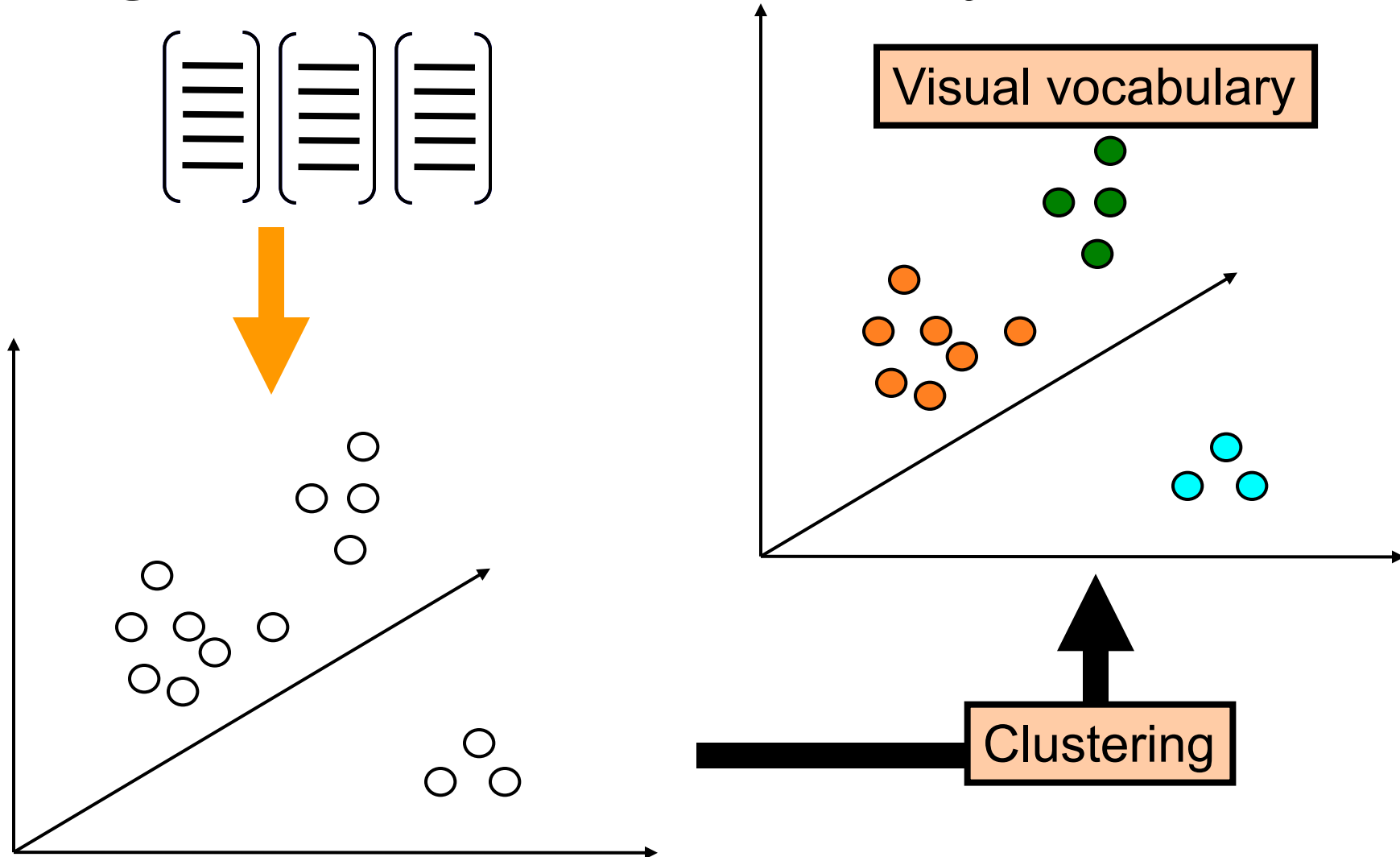
1. Extract local features
2. Learn “visual vocabulary”
3. Quantize local features using visual vocabulary
4. Represent images by frequencies of “visual words”



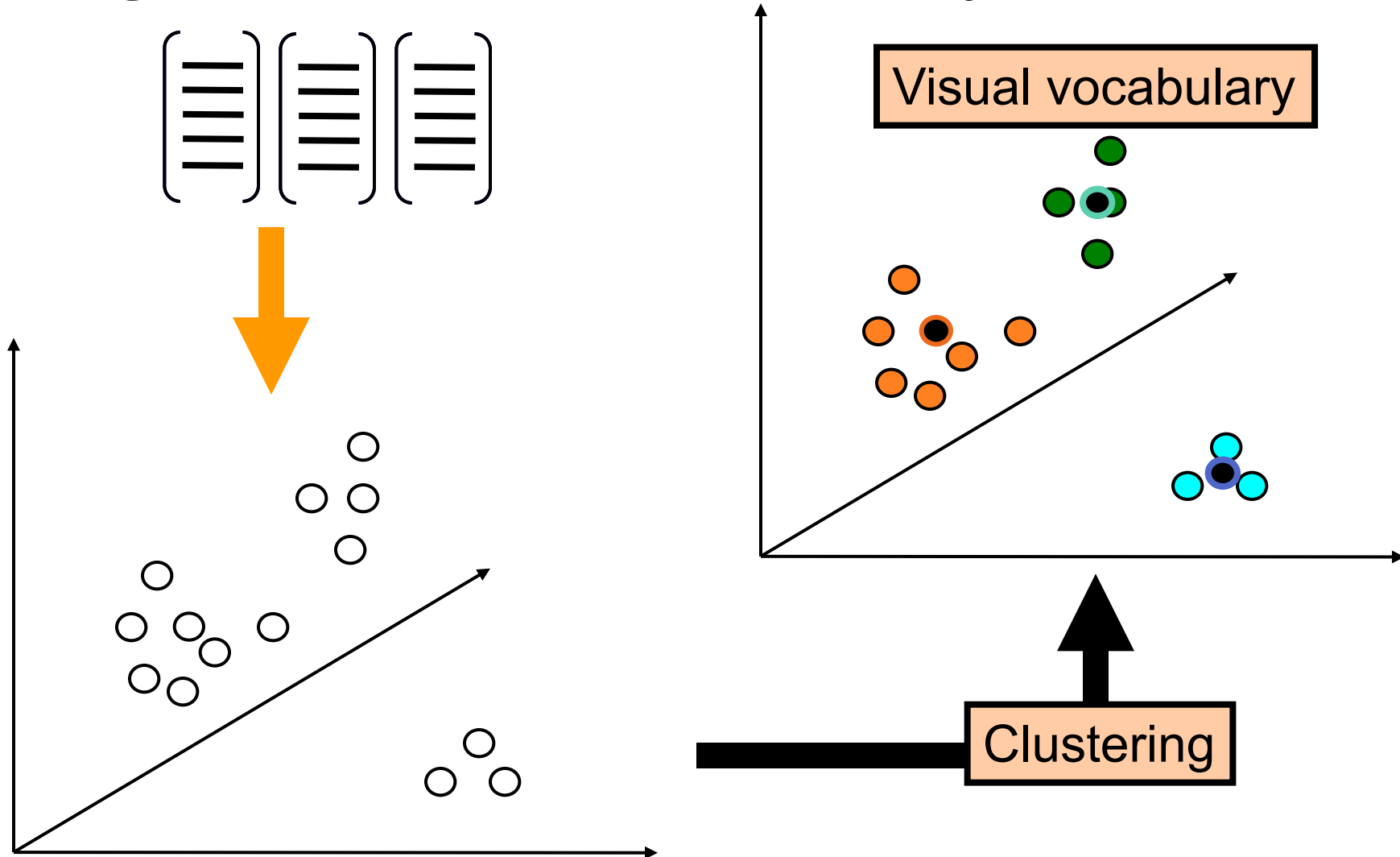
Learning a Visual Vocabulary



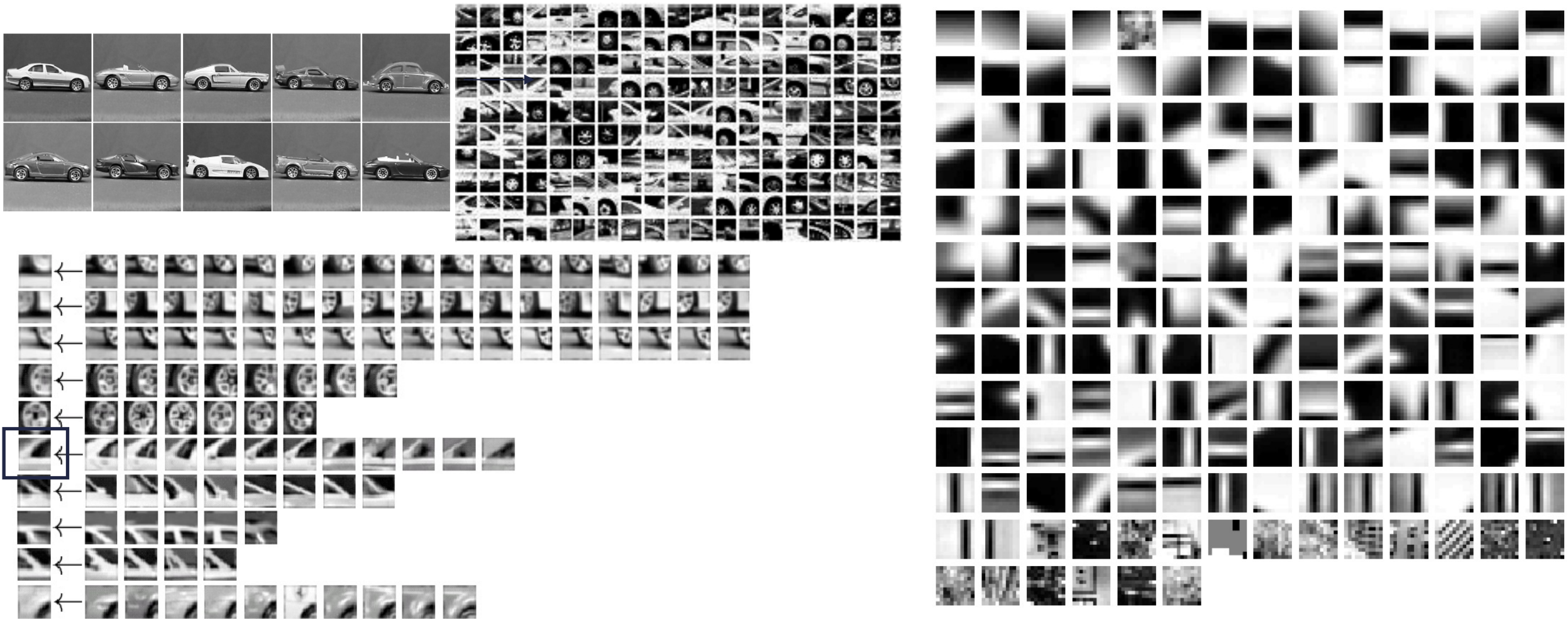
Learning a Visual Vocabulary



Learning a Visual Vocabulary

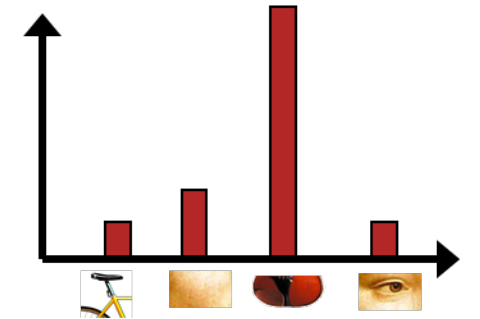
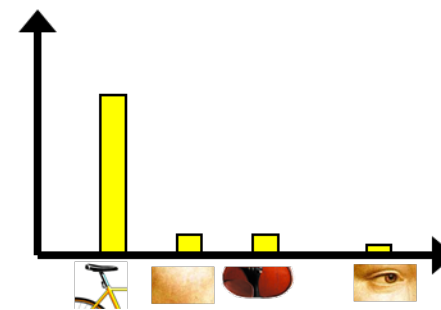
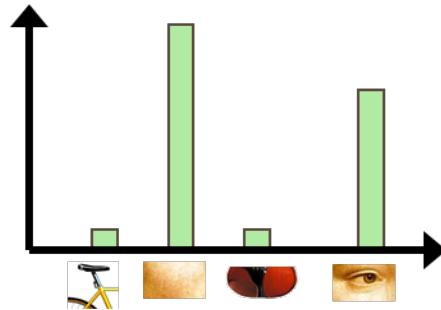
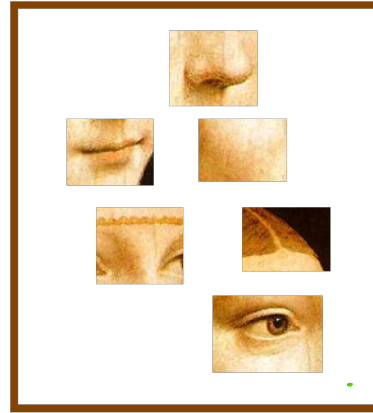


Example Visual Codebooks



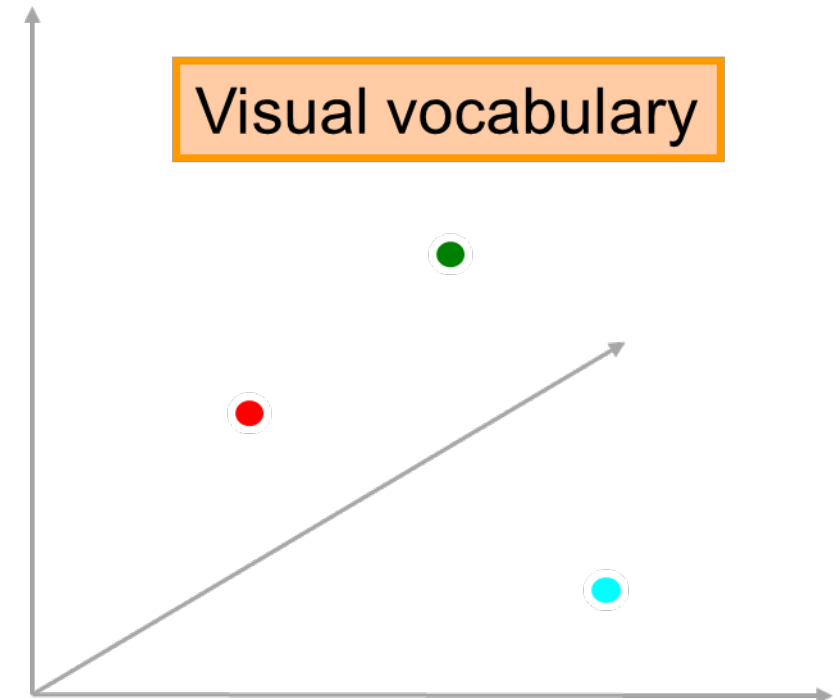
Bag of features

1. Extract local features
2. Learn “visual vocabulary”
3. Quantize local features using visual vocabulary
4. Represent images by frequencies of “visual words”



Bag of features

1. Extract local features
2. Learn “visual vocabulary”
3. Quantize local features using visual vocabulary
4. Represent images by frequencies of “visual words”



Images as Histogram of Patches



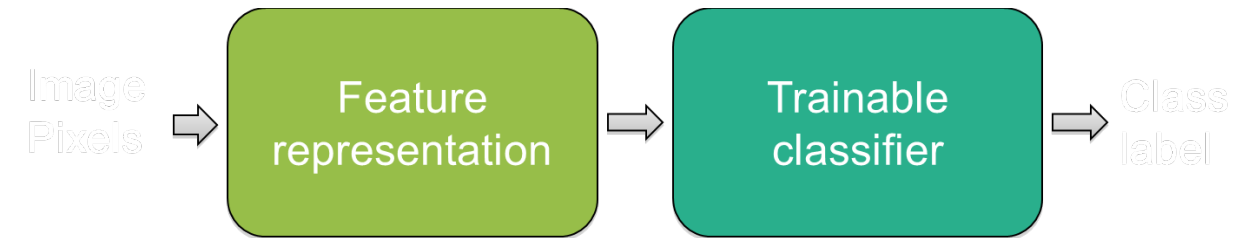
Today's Plan

- Computer vision overview
- Object recognition
 - Feature representations
 - Classification
- (Convolutional) Neural Networks

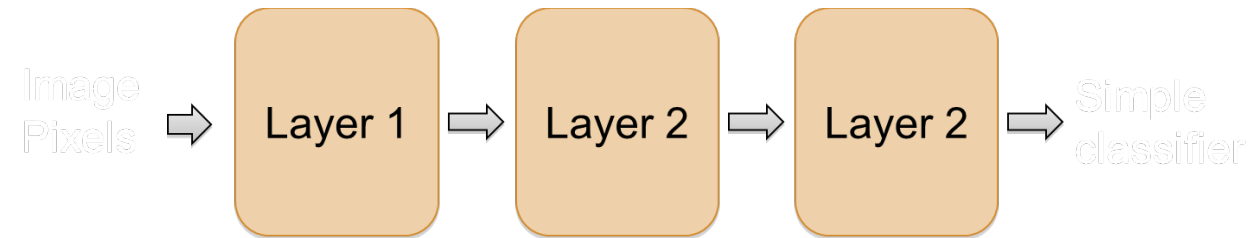


From Shallow to Deep Learning

Traditional “Shallow” Pipeline



“Deep” Recognition Pipeline



Multi-Layer Perceptron (MLP)



Activation Functions

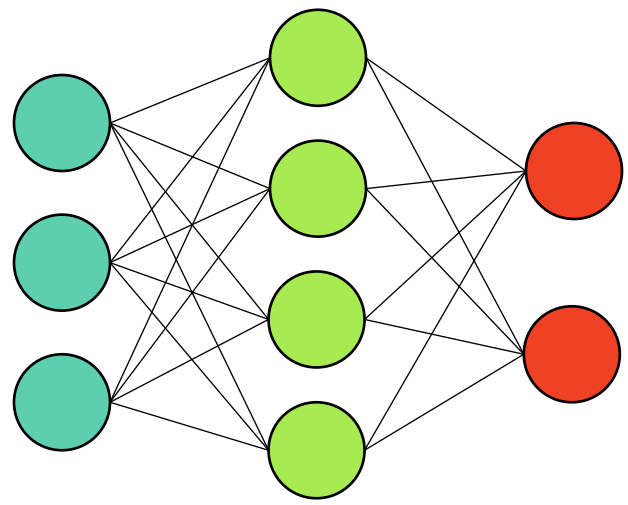
- Sigmoid
 - Homage to the original formulation
 - Not very popular nowadays as they tend to saturate and kills gradients
- Tanh
 - This is a scaled sigmoid and is almost always preferred
- ReLU – Rectified Linear Unit
 - Fast computation, doesn't saturate, might lead to better convergence rates
 - Tends to be fragile in training
- Maxout: $\max(w_1^T x + b_1, w_2^T x + b_2)$
 - Extension of ReLU that does not die
 - Doubles the number of parameters for every unit



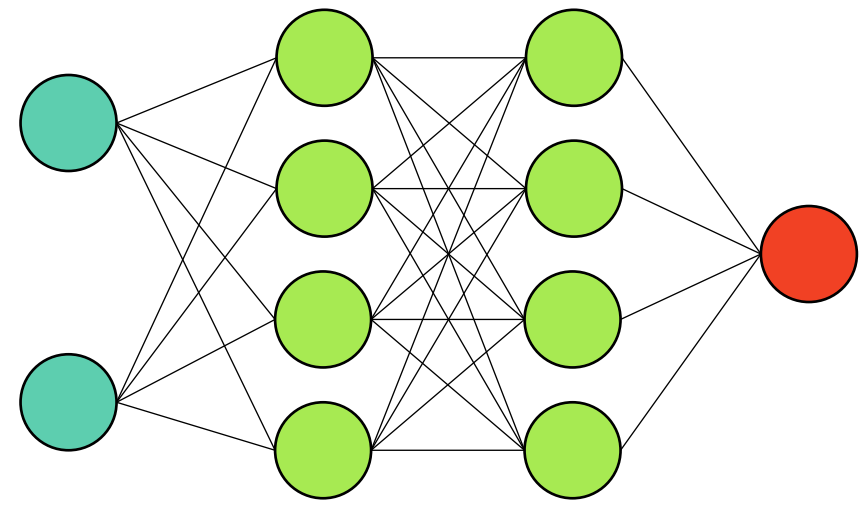
Network Architectures and Sizes

For regular networks, most commonly use fully connected network

Sizing of networks determined by layers, number of units, and/or number of parameters



6 units, 6 biases
 $[3*4]+[4*2]=20$ weights
26 learnable parameters



9 units, 9 biases
 $[2*4]+[4*4]+[4*1]=28$ weights
37 learnable parameters

For context, convolutional networks typically have on the order of 100 million parameters



Universal Function Approximators

Let $\varphi : \mathbb{R} \rightarrow \mathbb{R}$ be a nonconstant, **bounded**, and **continuous** function. Let I_m denote the m -dimensional **unit hypercube** $[0, 1]^m$. The space of real-valued continuous functions on I_m is denoted by $C(I_m)$. Then, given any $\varepsilon > 0$ and any function $f \in C(I_m)$, there exist an integer N , real constants $v_i, b_i \in \mathbb{R}$ and real vectors $w_i \in \mathbb{R}^m$ for $i = 1, \dots, N$, such that we may define:

$$F(x) = \sum_{i=1}^N v_i \varphi(w_i^T x + b_i)$$

as an approximate realization of the function f ; that is,

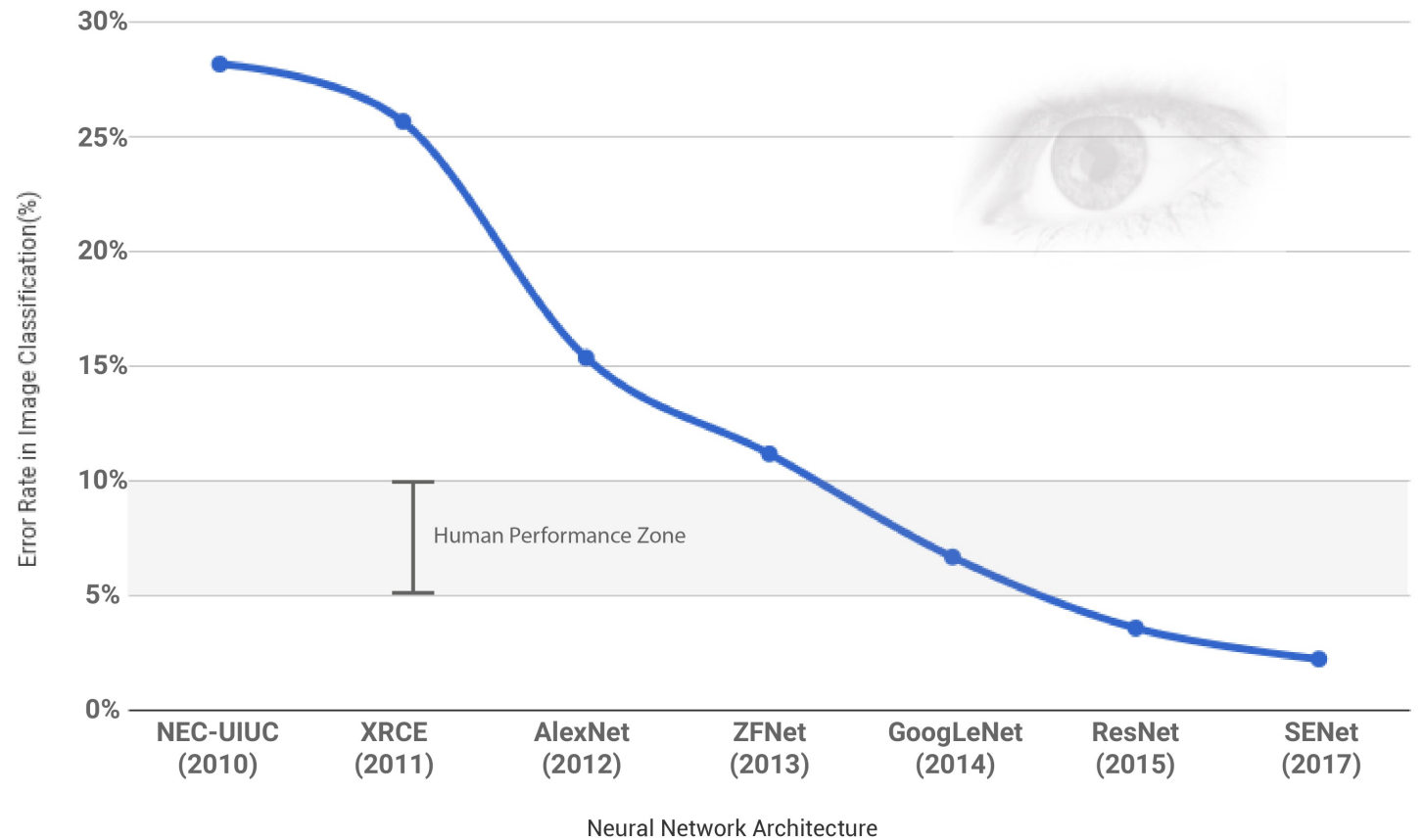
$$|F(x) - f(x)| < \varepsilon$$

for all $x \in I_m$. In other words, functions of the form $F(x)$ are **dense** in $C(I_m)$.

A feedforward network with a single hidden layer containing a finite number of units can approximate continuous functions on compact subsets of \mathbb{R}^n , under mild assumptions on the activation function.



Classification Improvements



Neural Networks

- Pros:
 - + Flexible and general function approximation framework
 - + Generally successful in high dimensional and model free problems



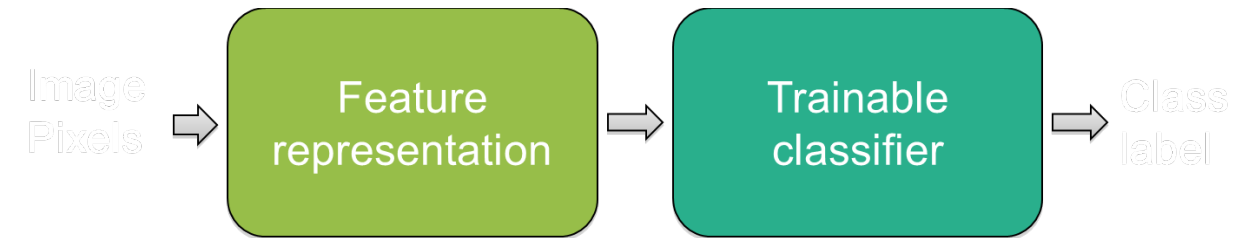
Neural Networks

- Pros:
 - + Flexible and general function approximation framework
 - + Generally successful in high dimensional and model free problems
- Cons
 - Very few theoretical guarantees
 - Training is prone to local optima and unstable
 - Large amount of training data and computing power are required
 - Huge variety of implementation choices need to be hand tuned (network architectures, parameters, etc.)

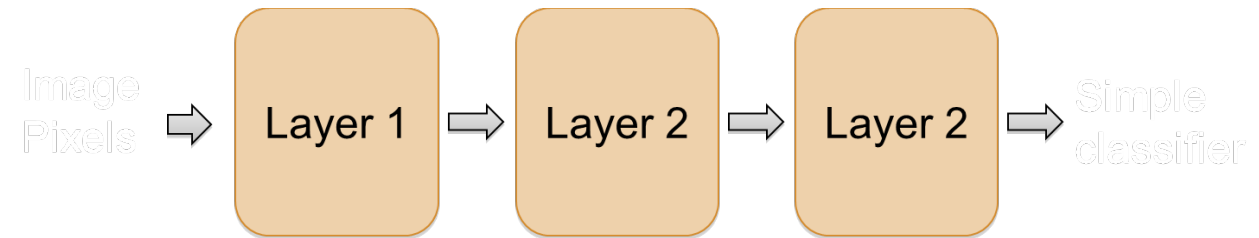


From Shallow to Deep Learning

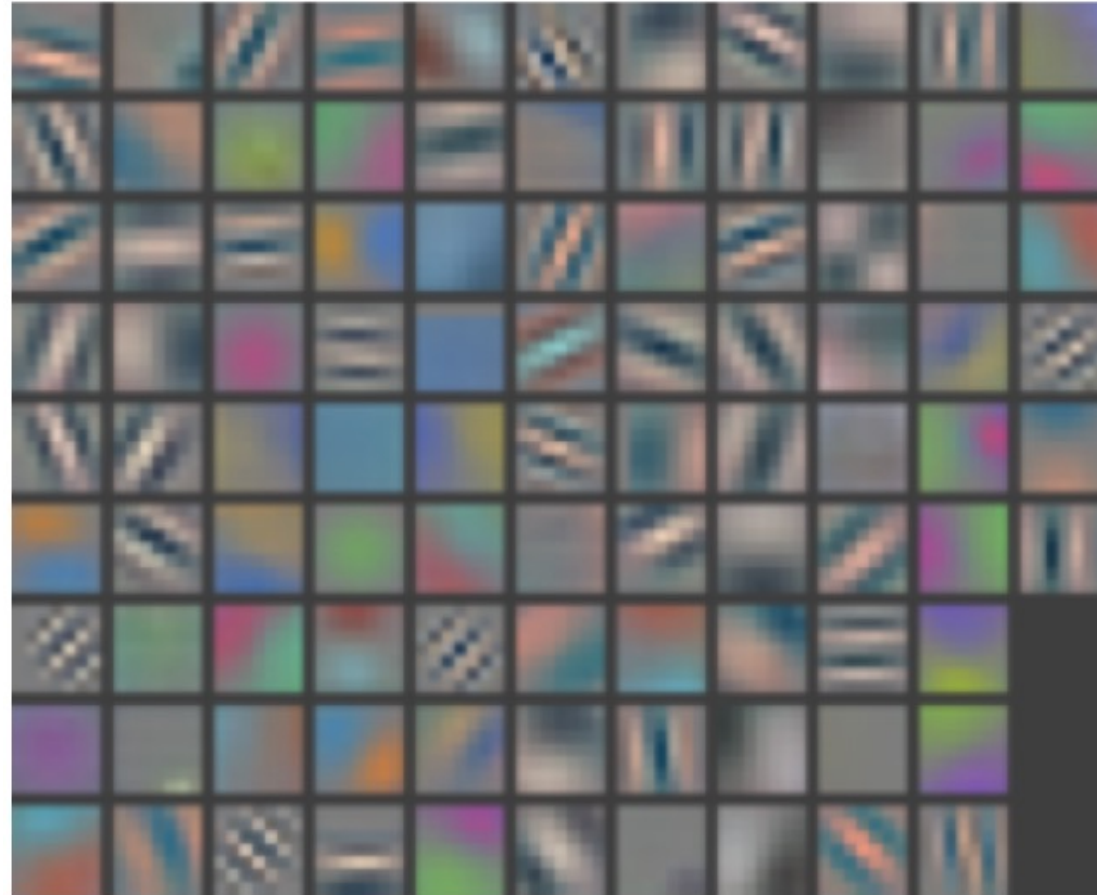
Traditional “Shallow” Pipeline



“Deep” Recognition Pipeline



Layers as filters

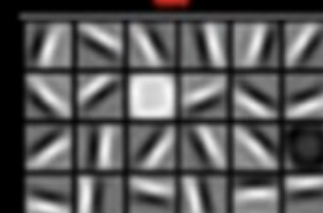
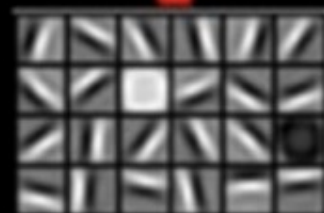
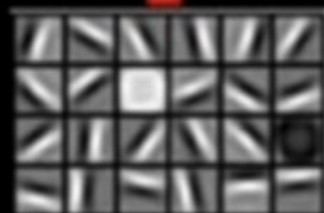
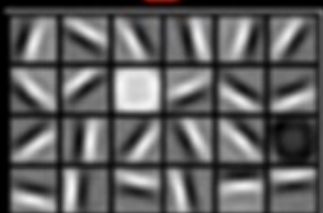
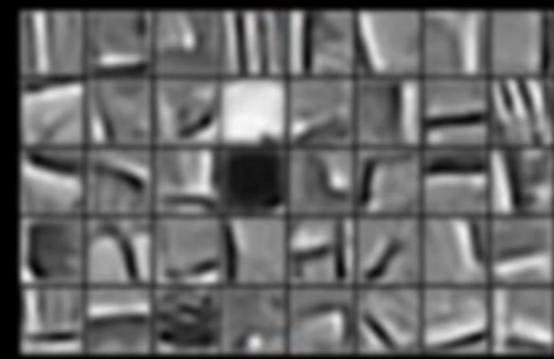
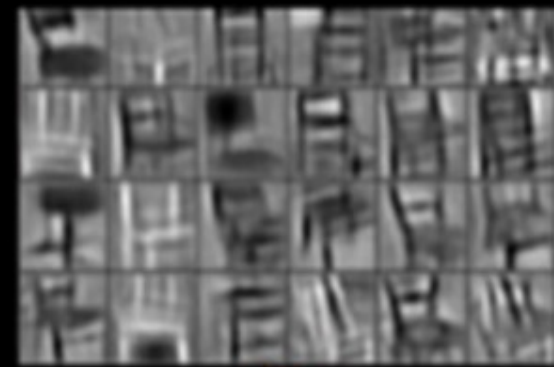


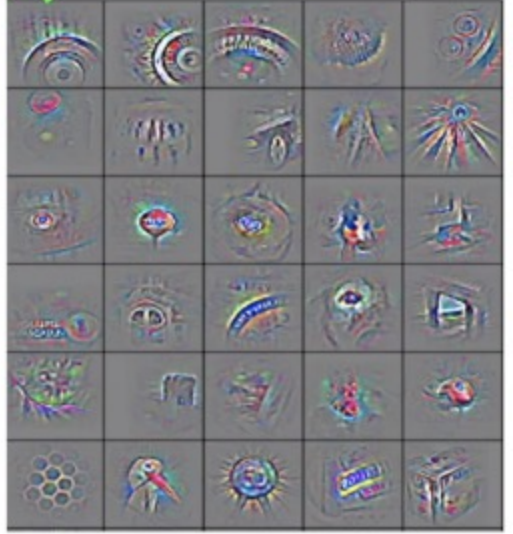
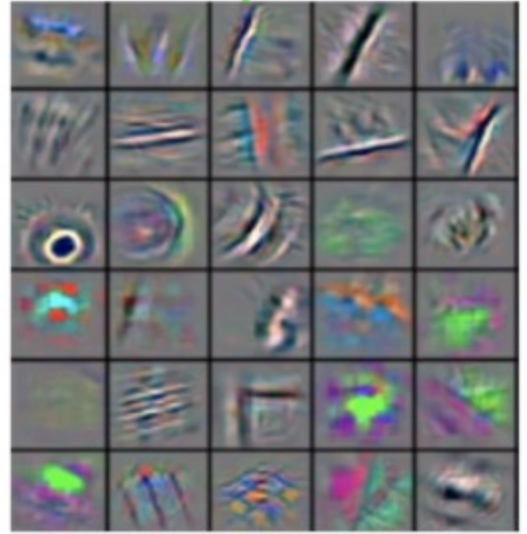
Faces

Cars

Elephants

Chairs





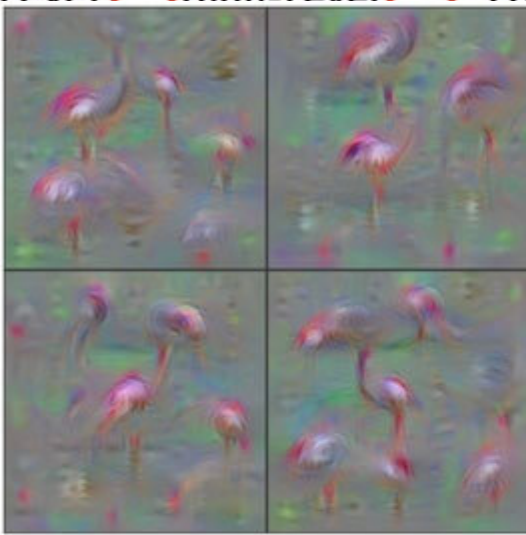
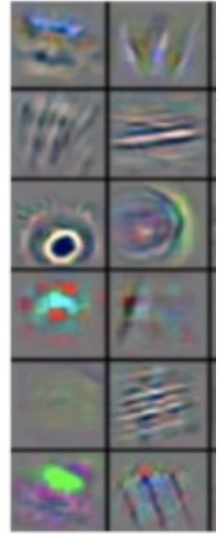


Low-Level Feature

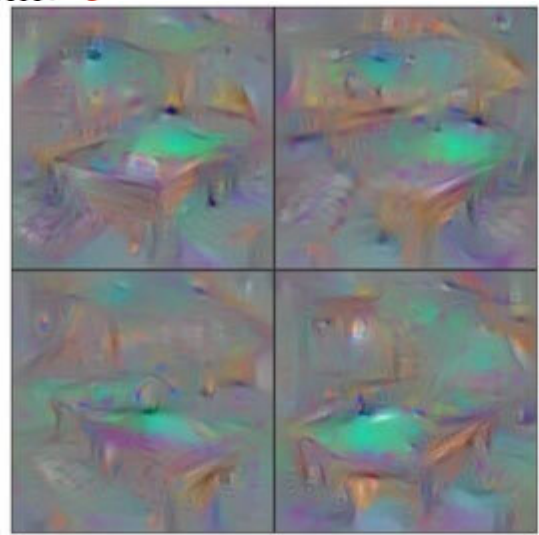
Mid-Level Fea

High-Level

Trainable



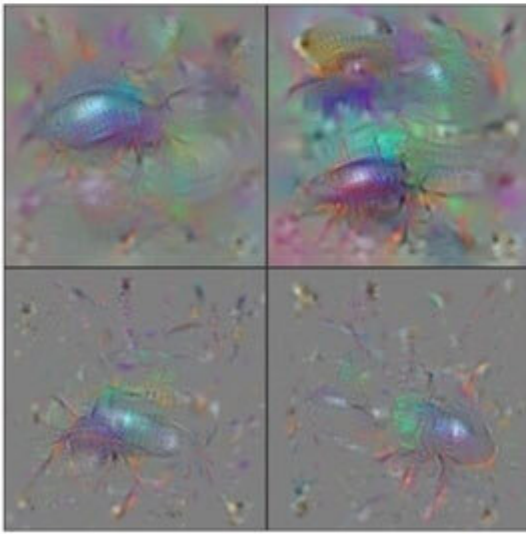
Flamingo



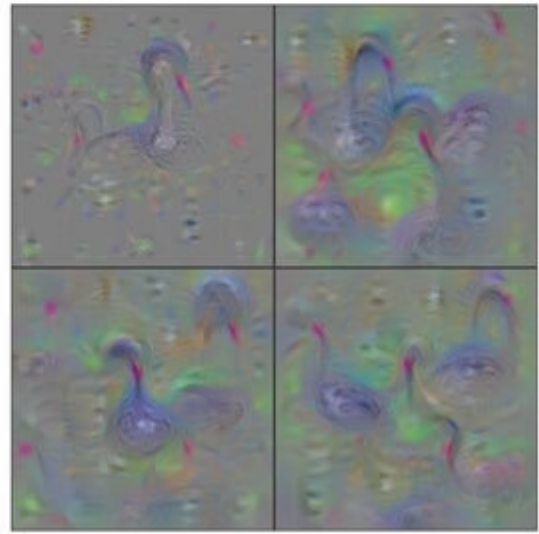
Billiard Table



School Bus



Ground Beetle



Black Swan

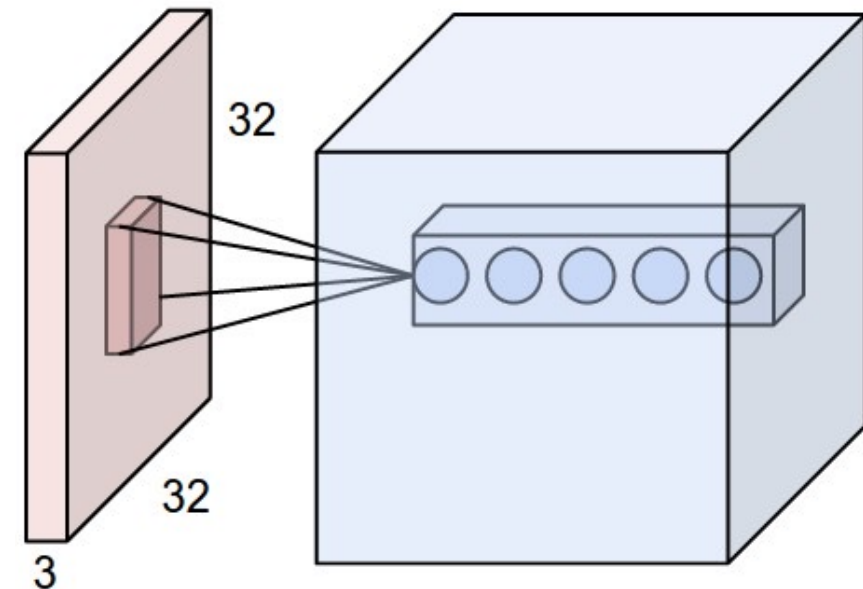
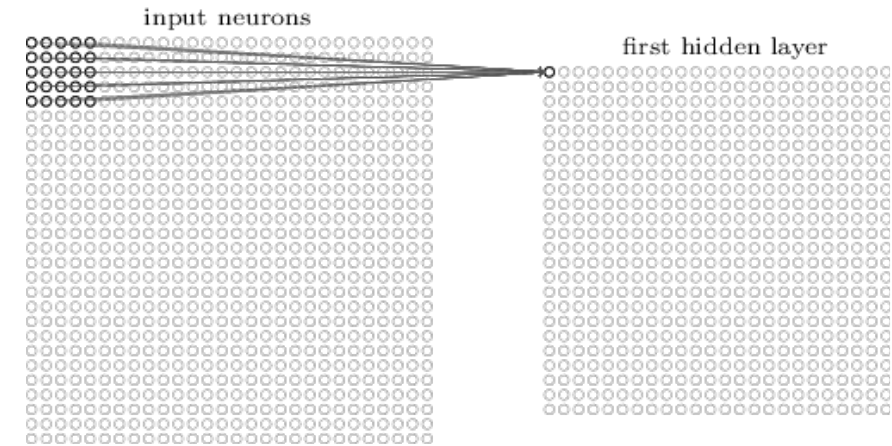


Tricycle



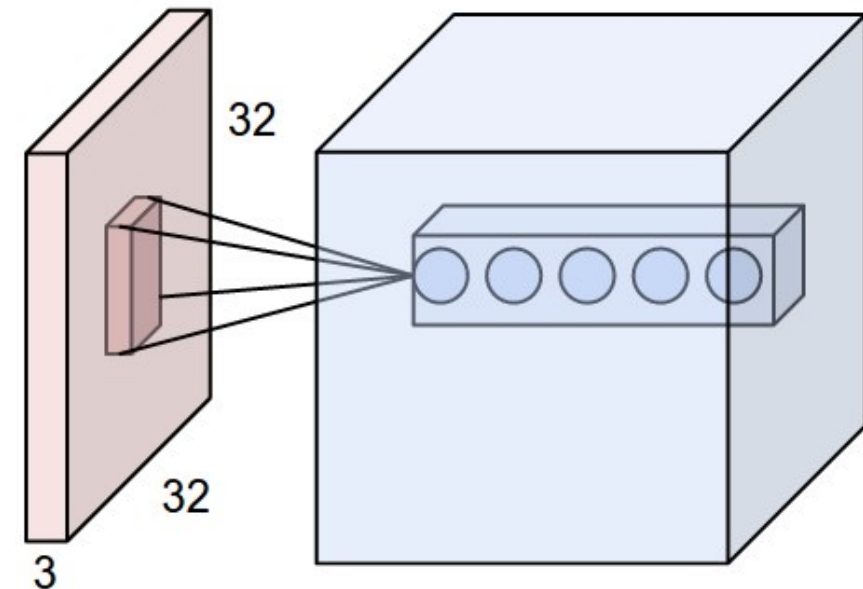
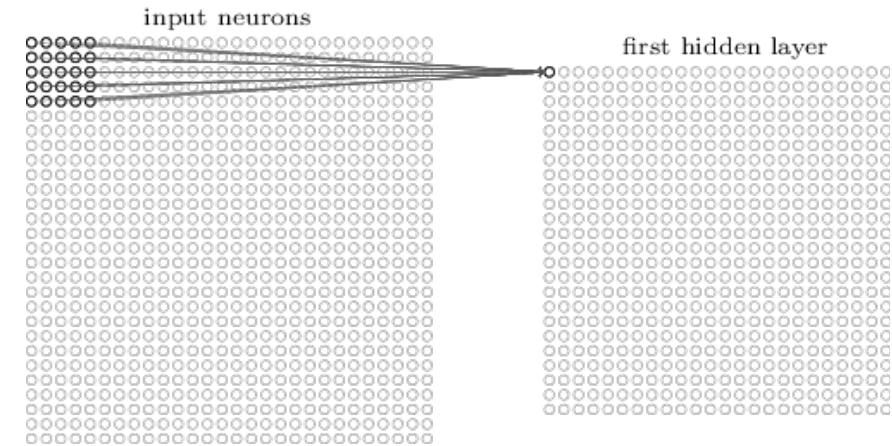
Convolutional Layers

- Each unit has a receptive field that connects it to a small local region of the input
 - If all units in a depth slice use identical weights, then the forward pass of this layer can be computed as a *convolution* of the weights with the input volume



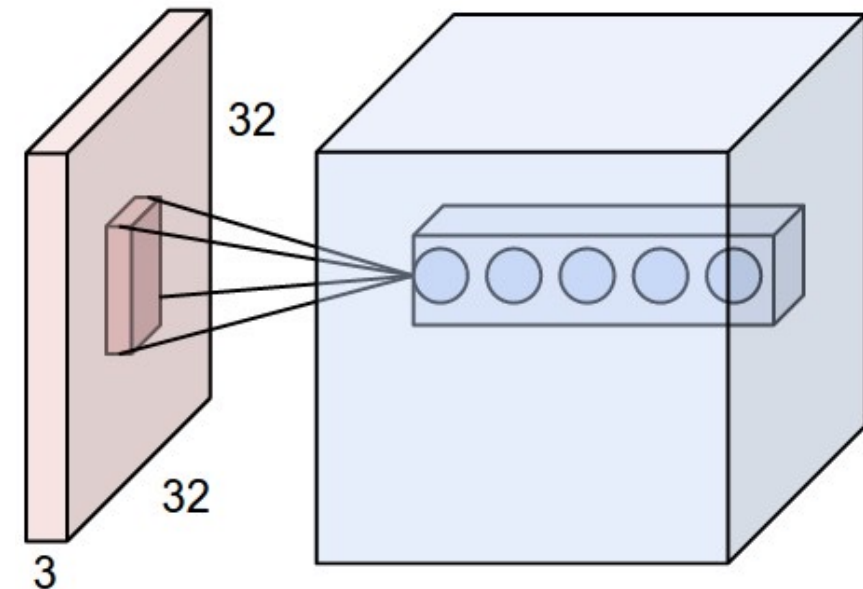
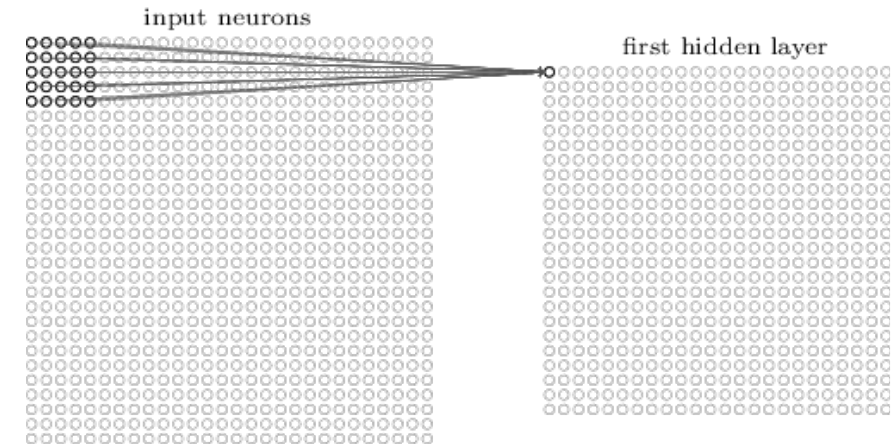
Convolutional Layers

- Each unit has a receptive field that connects it to a small local region of the input
 - If all units in a depth slice use identical weights, then the forward pass of this layer can be computed as a *convolution* of the weights with the input volume
- Each conv layer acts like a learnable filter that activates for some type of visual feature (e.g., edge, corner, eye, cat)



Convolutional Layers

- Each unit has a receptive field that connects it to a small local region of the input
 - If all units in a depth slice use identical weights, then the forward pass of this layer can be computed as a *convolution* of the weights with the input volume
- Each conv layer acts like a learnable filter that activates for some type of visual feature (e.g., edge, corner, eye, cat)
- Recall: large ConvNets have *a ton of* parameters
 - Parameter sharing restricts the weights along one *slice* of the depth, reducing the parameters down to ~35,000 (see first point)



Convolutional Network Components

ConvNets transform the original image layer by layer from the original pixel values to the final class scores

This is done via convolutional layers, pooling, ReLUs, and fully connected (FC) layers



Convolutional Network Components

ConvNets transform the original image layer by layer from the original pixel values to the final class scores

This is done via convolutional layers, pooling, ReLUs, and fully connected (FC) layers

- Conv/FC layers perform transformations that are a function of trainable parameters
 - Ex: CIFAR-10 images are size $32 \times 32 \times 3$, so one fully-connected unit in a first hidden layer of a regular NN would have $32 \times 32 \times 3 = 3072$ weights
- ReLU/Pool layers are fixed and not trained



Pooling Layers

The goal of pooling is to progressively reduce the spatial size of the representation and the amount of parameters and computation

- operates independently on depth slice and resizes it spatially, often using the max operation

Generally speaking:

- Accepts a volume of size $W1 \times H1 \times D1$
- Requires two hyperparameters: their spatial extent F , the stride S ,
- Produces a volume of size $W2 \times H2 \times D2$ where:

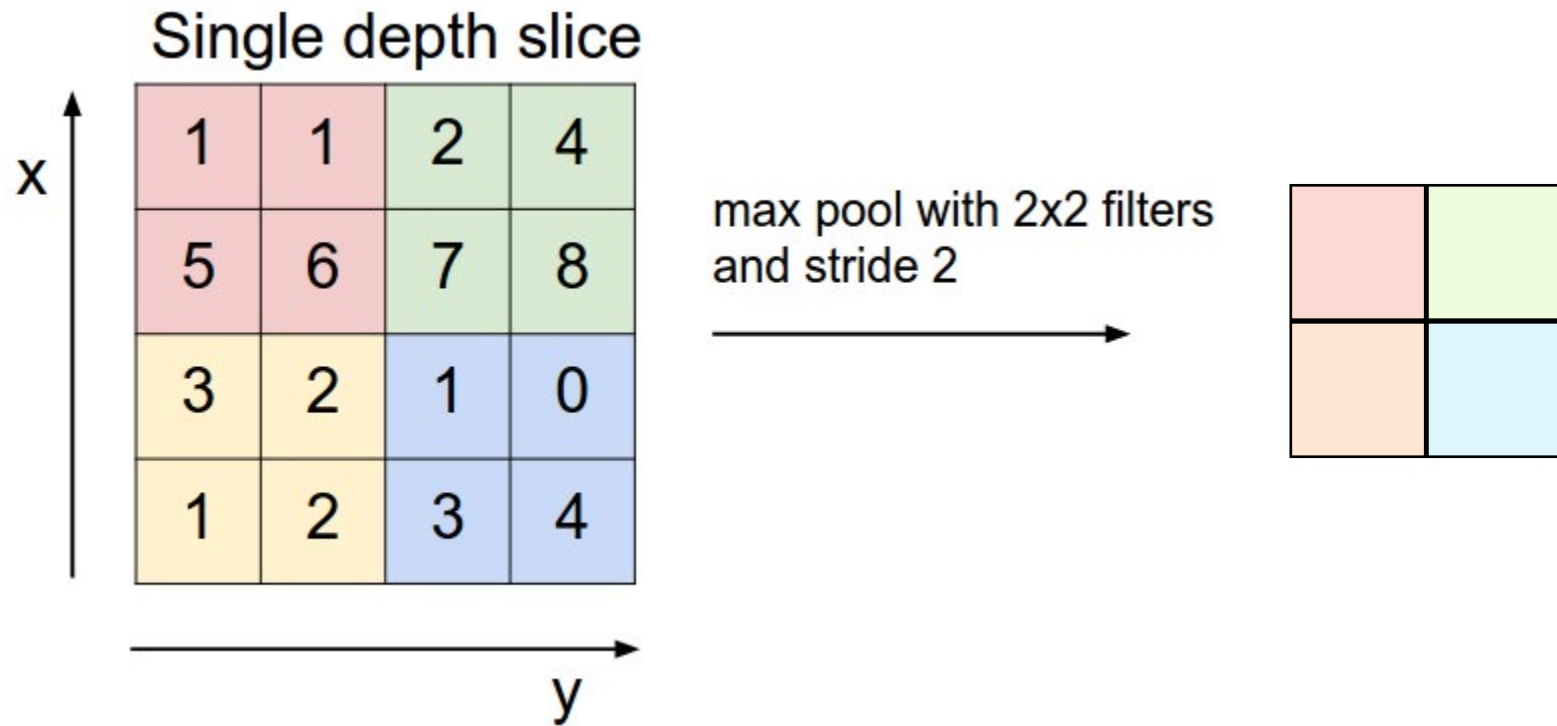
$$W2 = (W1 - F) / S + 1$$

$$H2 = (H1 - F) / S + 1$$

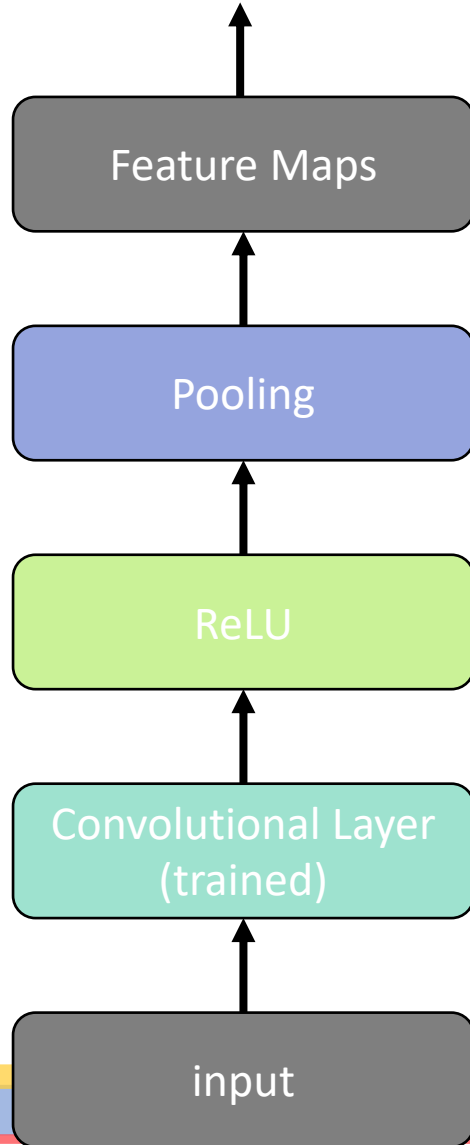
$$D2 = D1$$



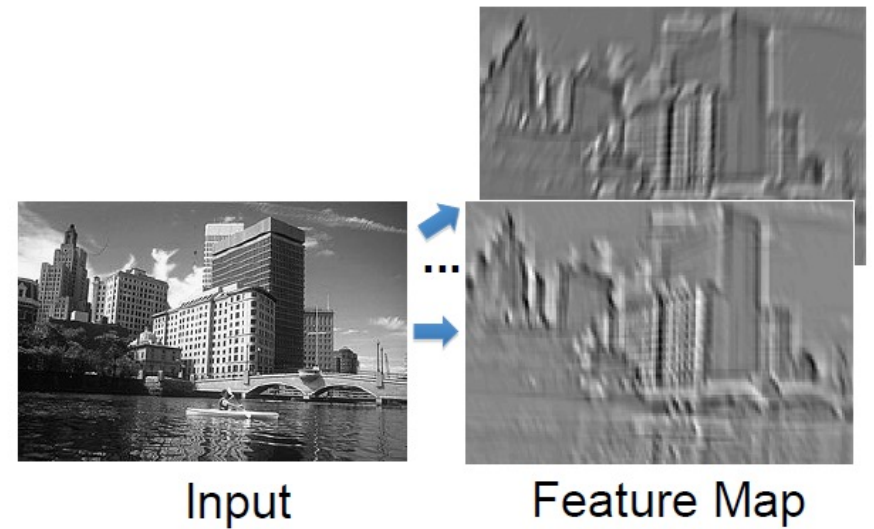
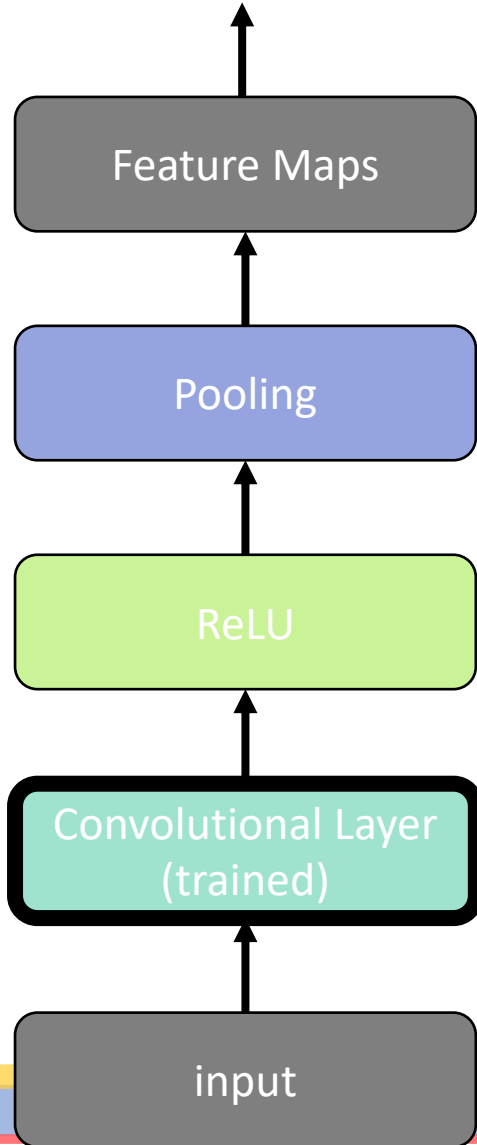
Example: Max Pooling



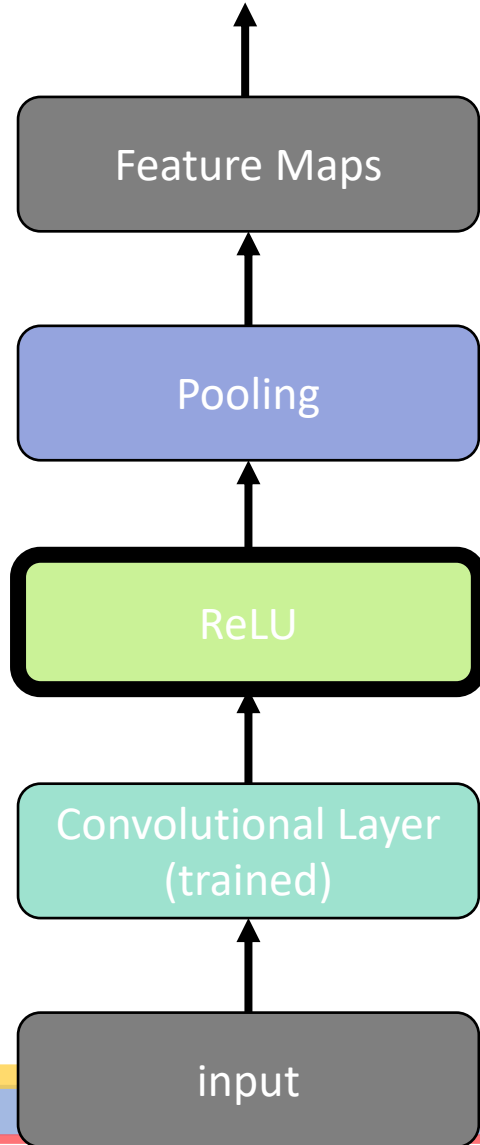
ConvNet Recap



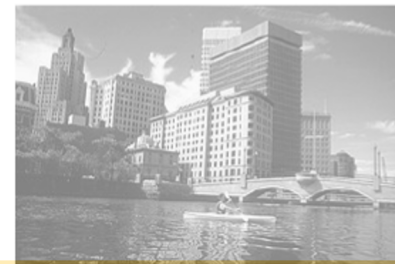
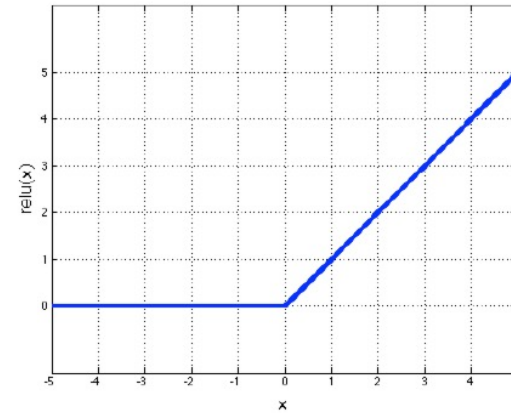
ConvNet Recap



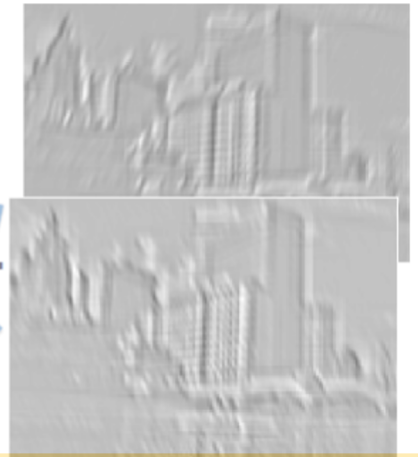
ConvNet Recap



Rectified Linear Unit (ReLU)



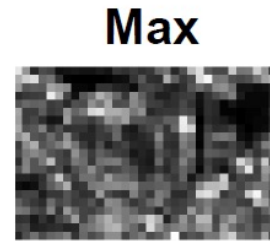
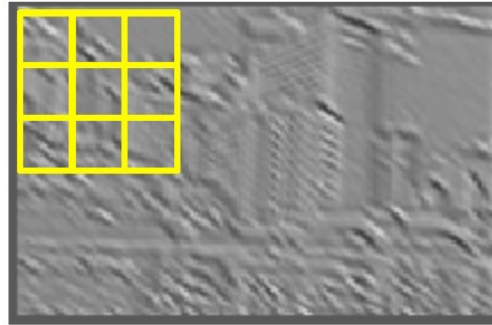
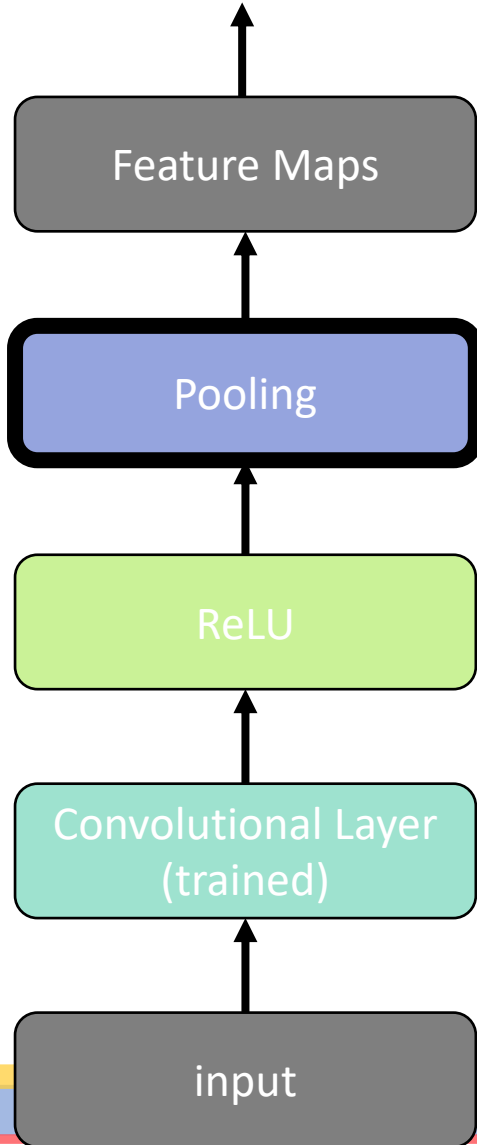
Input



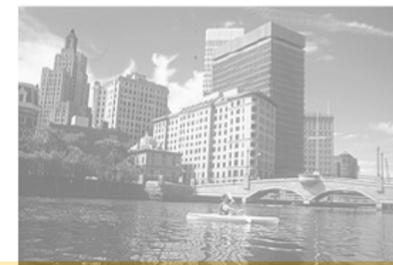
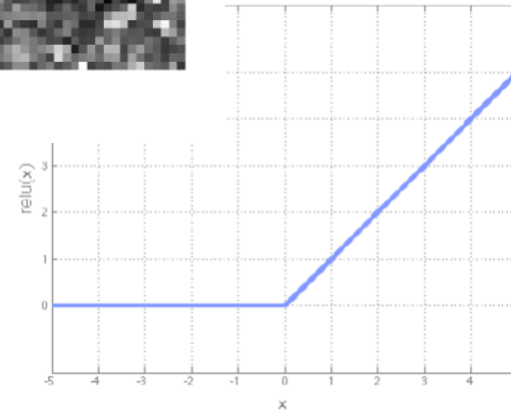
Feature Map



ConvNet Recap



Linear Unit (ReLU)

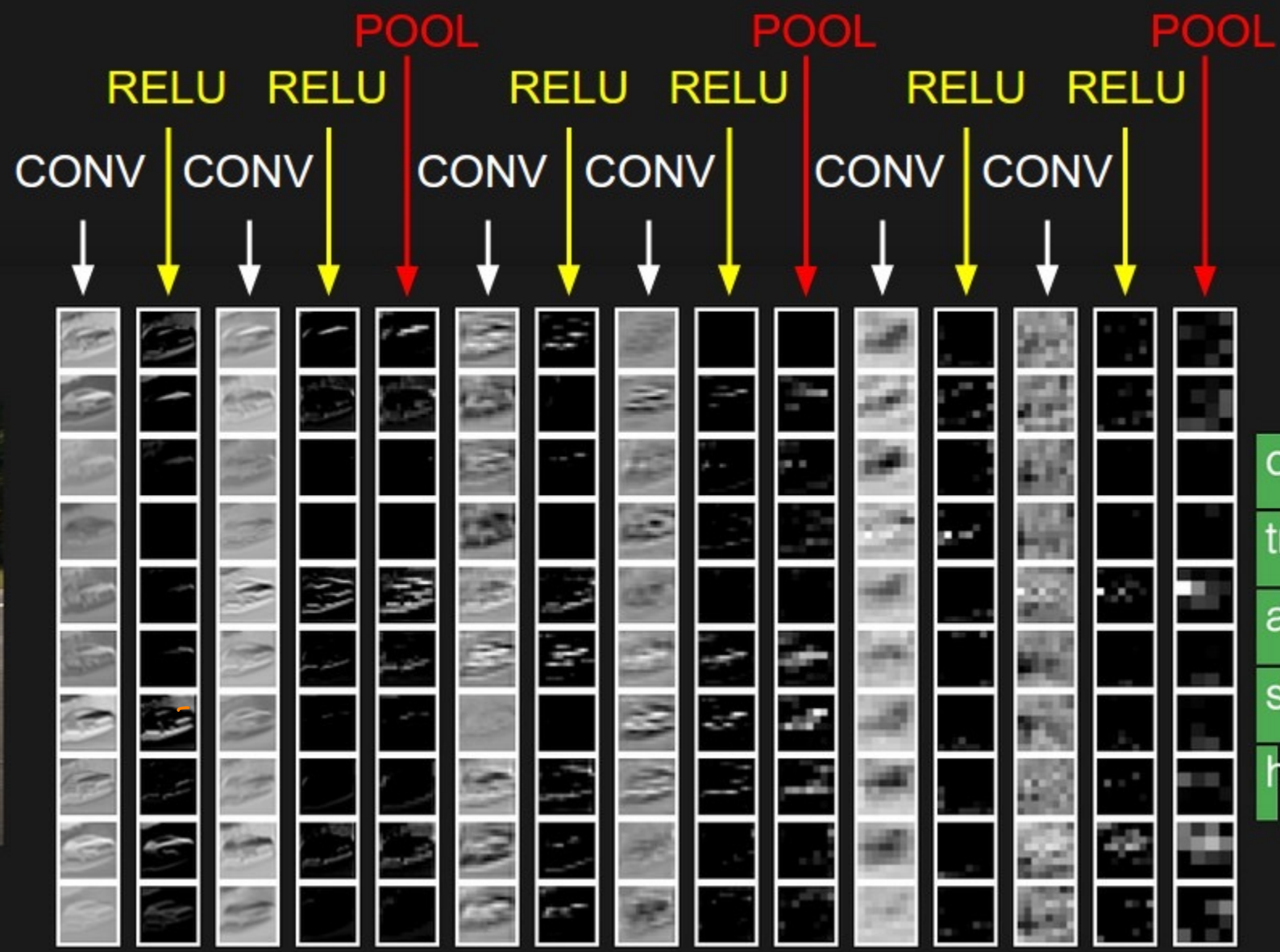


Input



Feature Map





- FC
- car
 - truck
 - airplane
 - ship
 - horse

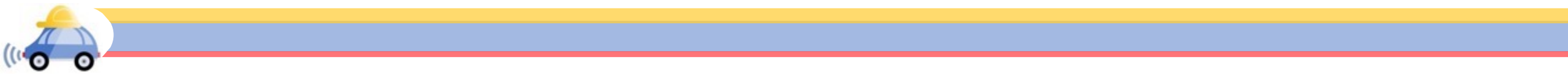


Summary

- Crash course in computer vision
 - Recognition, reconfiguration, and reconstruction
 - Traditional features vs. learned features
- Introduced the basics of neural networks
 - Did not discuss: backpropagation or training methods
 - Did not discuss: state-of-the-art object detection architectures
- Next time: we'll look at modeling and control of vehicles!



Extra Slides



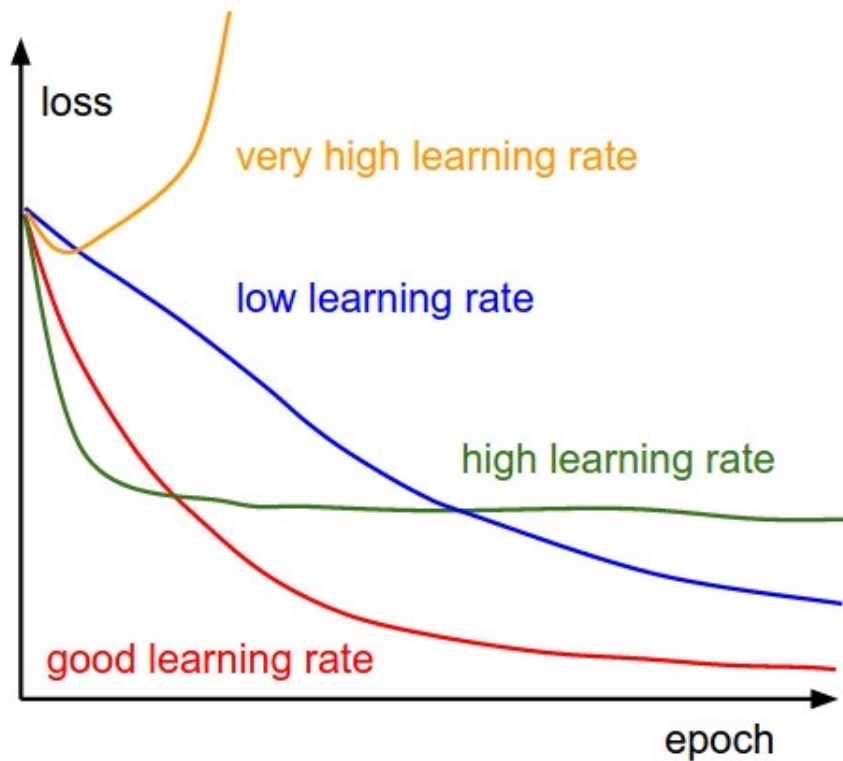
So you want to train a neural net.

1. Pre-process data
 - Zero-center and scale by standard deviation
2. Initialize network
 - Initializing weights can be difficult due to instabilities
 - Small random numbers from normal distribution
3. Set up your regularization (penalty term, dropout)
4. Pick a loss function
 - Depends on problem, but try to shoot for softmax whenever possible
5. You are ready to train your network!
 - Initially try to overfit on a tiny subset of your data ~20 samples. Make sure you get zero loss.
6. Sweep over hyperparameters
 - Initial learning rate and decay schedule, regularization strength
 - Use cross-validation techniques and be prepared to wait. This can take weeks for large networks.

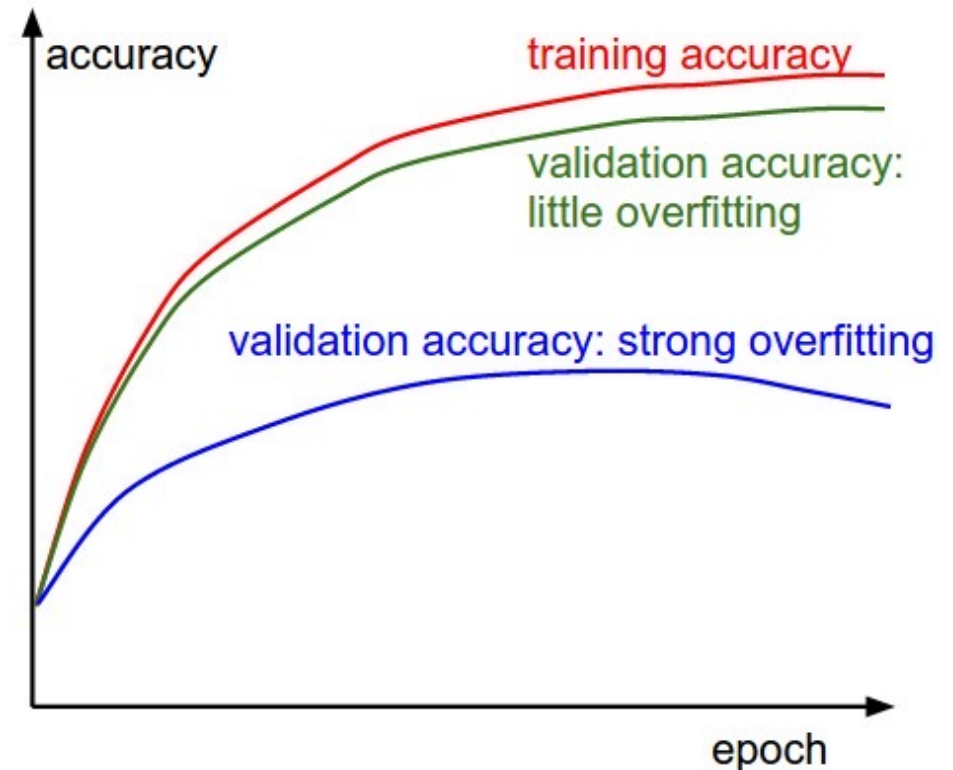


What to watch during training

Loss Rates



Training vs. Validation

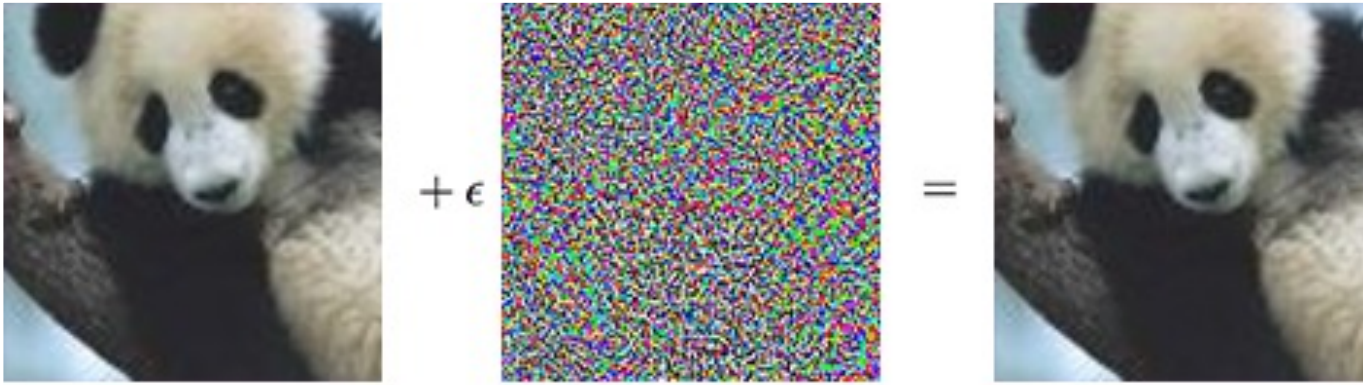


Risks for Autonomy

- What is it about neural networks that are particularly difficult?
 - Training stability is a problem
 - Large amounts of data are required
 - Huge amounts of computation are required



Adversarial Examples



"panda"
57.7% confidence

"gibbon"
99.3% confidence



(a) Image from dataset



(b) Clean image



(c) Adv. image, $\epsilon = 4$



(d) Adv. image, $\epsilon = 8$

| Classification Results |
|--|
| washer: 0.5398173 |
| safe: 0.34602574 washer: 0.22088042 |
| safe: 0.3719305 loudspeaker: 0.24184975 |



A few things to keep in mind

1. Machine Learning is not always the answer – try simple methods first.
 - However, use ML over a complex heuristic. A simple heuristic can only get you so far, while a complex heuristic is unmaintainable.
2. When picking features, make sure they are generalizable!
3. Watching out for data imbalances or other quirks with your data.
4. You may skew your data by causing a discrepancy between how you handle data in the training and testing.
5. Cross validation is key – never peek at your testing data. You may create a feedback loop between your model and your algorithm.

