# Introduction to Recognition

# Saurabh Gupta

# Computer Vision
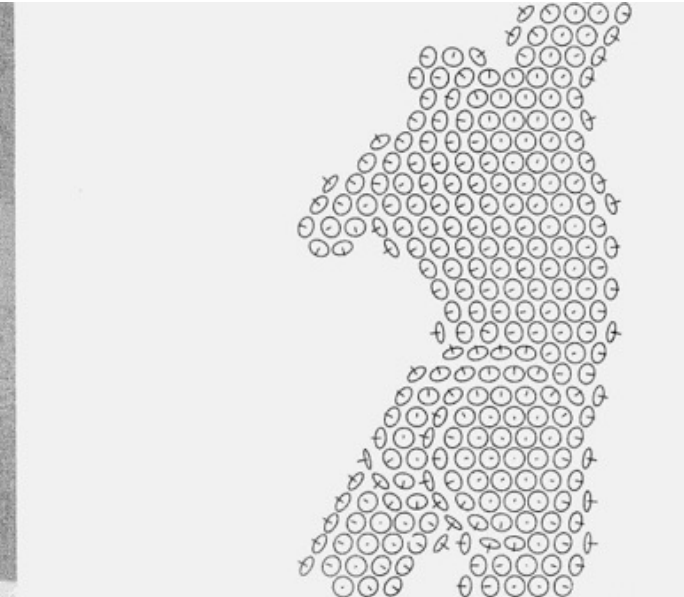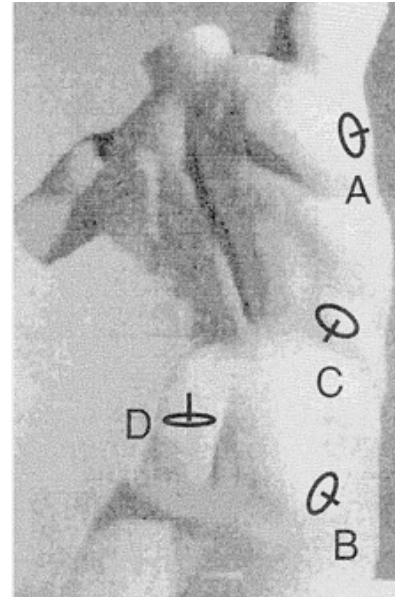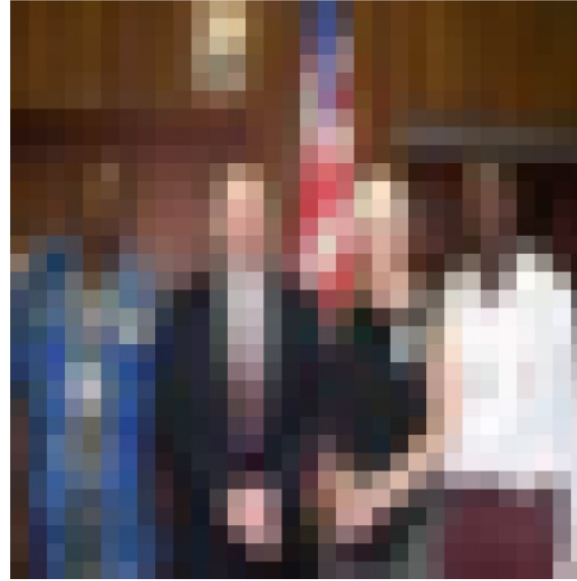
To extract "meaning" from pixels



person, motorcycle, car, chair

*Meaning* can take different forms:

- Geometric Inferences

- Semantic Inferences

- Inferences about actions

- …

# Computer vision is easy for humans

- Effortlessly analyze images for a variety of tasks

- Infer semantics even from severely ablated

- Can also make precise inference about certain geometric properties

# Yet has proven very hard for computers

- Computer vision research easily goes back 60 years …

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

PROJECT MAC

Artificial Intelligence Group          July 7, 1966
Vision Memo. No. 100.

THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".
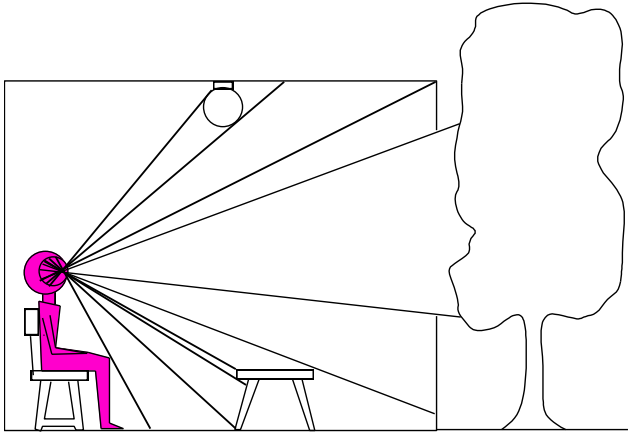


- Entirely true as of 2014 (or so) when this xkcd was published
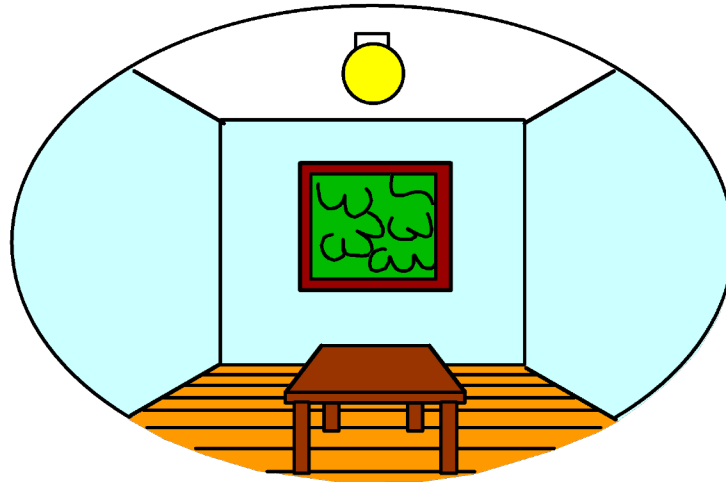
# Why is computer vision hard?

- Images are a lossy projection of the world

**3D world**

**2D image**



Point of observation

Geometry information is lost

# Why is computer vision hard?

- Images are a lossy projection of the world

What color is the dress?

A) Black and blue

B) White and gold?

Appearance information is also lost

https://www.wired.com/2015/02/science-one-agrees-color-dress/

# Why is computer vision hard?

- Images are a lossy projection of the world



Might cause objects to blend

# Why is computer vision hard?

- Images are a lossy projection of the world (geometry, appearance, ... are lost)
- Visual world is diverse



Viewpoint variation

Shape variation

# Why is computer vision hard?

- Images are a lossy projection of the world (geometry, appearance, ... are lost)
- Visual world is diverse



Background clutter



Occlusion

# Why is computer vision hard?

- Images are a lossy projection of the world (geometry, appearance, … are lost)
  - need some priors to interpret what you are seeing

- Visual world is diverse
  - can't write down these priors by hand



*John's Diner with John's Chevelle*, 2007

Enter machine learning

# Why machine learning?

- Good old-fashioned AI (GOFAI) answer:
  Program expertise into the agent



Figure 3. Structure of description built by the system.

Figure 2. Semantic network for knowledge representation.

Figure 5-a. Digitized color scene.

5-b. Result of preliminary segmentation.

5-c. Plan image.

5-d. Result of semantic segmentation.

Y. Ohta, T. Kanade and T. Sakai. An Analysis System for Scenes Containing objects with Substructures. Proc. of the Fourth International Joint Conference on Pattern Recognition, pp. 752-754, 1978

# Why machine learning?

- Good old-fashioned AI (GOFAI) answer:
  Program expertise into the agent



Figure 4-a. "Building" region and "windows".



Figure 4-b. The production for analyzing "windows".

# Why machine learning?

- Good old-fashioned AI (GOFAI) answer:
  Program expertise into the agent

# Why machine learning?

- Good old-fashioned AI (GOFAI) answer:
  Program expertise into the agent
  - Never worked (in general)

# Why machine learning?

- Good old-fashioned AI (GOFAI) answer:
  Program expertise into the agent

- Modern answer: Program into the agent the *ability to improve performance based on experience*
  - Experience should come from *training data* or *demonstrations*
  - We want to optimize the performance of the agent on the training data, with the hope that it will *generalize* to unseen inputs
  - This is the *statistical learning* viewpoint

# The basic ML framework (for supervised learning)

**Training time**

Labeled training data

"apple"
"pear"
"tomato"
"cow"
"dog"
"horse"

Training → Learned model

**Test time**

Unlabeled test sample

Learned model

↓

Inference → Label prediction

"apple"

# The basic ML framework (for supervised learning)

$$y = f(x)$$

output   prediction   input
         function



- **Training** (or **learning**): given a *training set* of labeled examples $\{(x_1, y_1), \ldots, (x_N, y_N)\}$, instantiate a predictor $f$
- **Testing** (or **inference**): apply $f$ to a new *test example* $x$ and output the predicted value $y = f(x)$



- Rather than hand-defining how 2D projections of apples are different from pears, $f$ will learn this from the data.

# Deep Learning

- A general way to model function $f$ as composition (layers) of simple functions, very loosely inspired by the brain.

# Lecture overview

- Different recognition problems in computer vision

- Supervised classification

- Taxonomy of learning problems

# Different Recognition Problems



**Classification**: Assign image to one of a fixed set of categories

**Object Detection**: Put a bounding box around each instance of a class

**Instance Segmentation**: Mark pixels for each instance of a class

**Semantic Segmentation**: Label each pixels with its category

# Different Recognition Problems



**Image Captioning**: Man riding a horse on a beach



**Depth Prediction:** how far is each pixel in the image



**Keypoint prediction**



**Pose Prediction:** Rotation that aligns object to a canonical pose

# The basic ML framework (for supervised learning)

**Training time**

Labeled training data

"apple"
"pear"
"tomato"
"cow"
"dog"
"horse"



→ Training → Learned model

**Test time**

Unlabeled test sample



Learned model

↓

→ Inference →

Label prediction

"apple"

# The basic ML framework (for supervised learning)

$$y = f(x)$$



output      prediction      input
function

- **Training** (or **learning**): given a *training set* of labeled examples $\{(x_1, y_1), \ldots, (x_N, y_N)\}$, instantiate a predictor $f$
- **Testing** (or **inference**): apply $f$ to a new *test example* $x$ and output the predicted value $y = f(x)$



- Rather than hand-defining how 2D projections of apples are different from pears, $f$ will learn this from the data.

# Is an image classifier all you need?

- Image Classification
- Object Detection
- Instance Segmentation
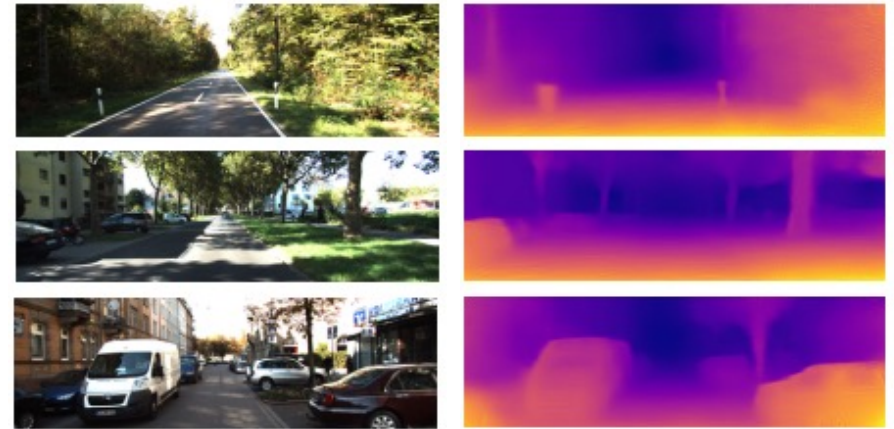- Semantic Segmentation
- Image Captioning
- Depth Prediction
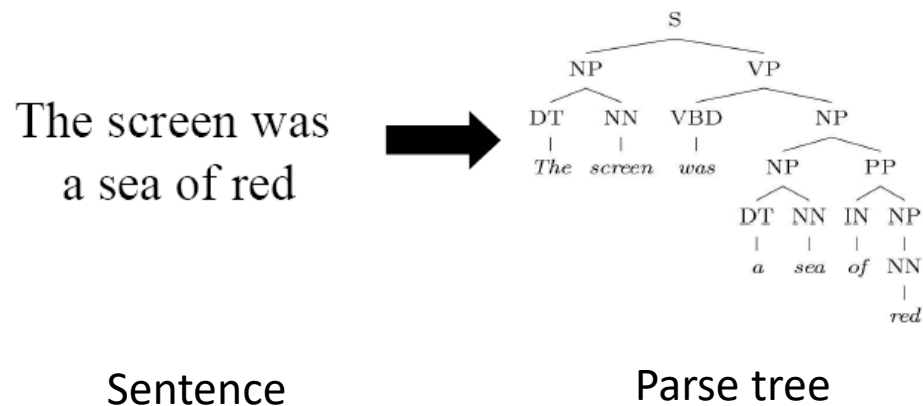- Keypoint Prediction
- Pose Prediction
- …

# Taxonomy of learning problems

- **Type of output**
  - Classification
  - Regression
    - $y = f(x)$. $y$ is an arbitrary scalar and not a class label.
  - Structured prediction
    - $y = f(x)$. $y$ is a structured object.



**Depth Prediction:** how far is each pixel in the image

Several computer vision problems have structure in the output space, but often solving a classification problem with some simple post-processing (or even without) ends up being sufficient.



The screen was a sea of red $\longrightarrow$

Sentence      Parse tree