

# Advanced Computer Vision

Elisa Ricci & Nicu Sebe

# Top 10 AI technology trends



## Deep learning theory

The information bottleneck principle explains how a deep neural network learns.



## Capsule networks

New type of deep neural network that learns with fewer errors and less data, by preserving key hierarchical relationships.



## Deep reinforcement learning

This technique combines reinforcement learning with deep neural networks to learn by interacting with the environment.



## Generative adversarial networks

A type of unsupervised deep learning system, implemented as two competing neural networks, enabling machine learning with less human intervention.



## Lean and augmented data learning

Different techniques that enable a model to learn from less data or synthetic data.



## Probabilistic programming

A high-level language that makes it easy for developers to define probability models.



## Hybrid learning models

Approach that combines different types of deep neural networks with probabilistic approaches to model uncertainty.



## Automated machine learning

Technique for automating the standard workflow of machine learning.



## Digital twin

A virtual model used to facilitate detailed analysis and monitoring of physical or psychological systems.



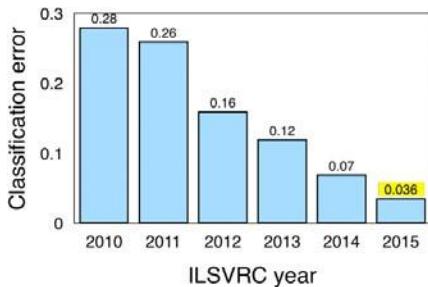
## Explainable artificial intelligence

Machine learning techniques that produce more explainable models while maintaining high performance.

# Successes of Deep Learning in AI

## The New York Times

A Learning Advance in Artificial Intelligence Rivals Human Abilities



IMAGENET

Deep Learning for self-driving cars



Google's DeepMind Masters Atari Games

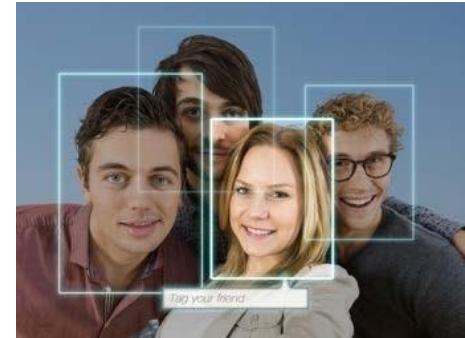


Google Translate

English Russian Chinese (Simplified) ▾

Time flies like an arrow

时间过得很快像箭



Face Recognition

# So is AI solved?

pedestrian detection FAIL (Volvo S60 Pedestrian Detection System Test)



[https://www.youtube.com/watch?  
v=w2pwxv8rFkU](https://www.youtube.com/watch?v=w2pwxv8rFkU)

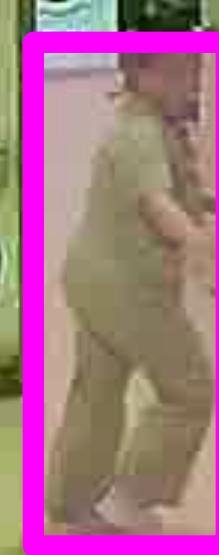
Pedestrian detection is only active at <22mph and >2mph. Once the brakes have been engaged, they do release afterwards assuming the driver who wasn't initially paying attention is now fully aware of the audible alarm and flashing red lights on the windshield. This is why this guy almost got run over: the system was working as designed but he was not behind the wheel to take back!!!

# Challenges



What does computer vision involve?

Detection: are there people?





indoor scene

long term care  
facility

walker

chair

Objects and scenes: where are they?

floor



stand

run

Action recognition: what are they doing?

fall

squat

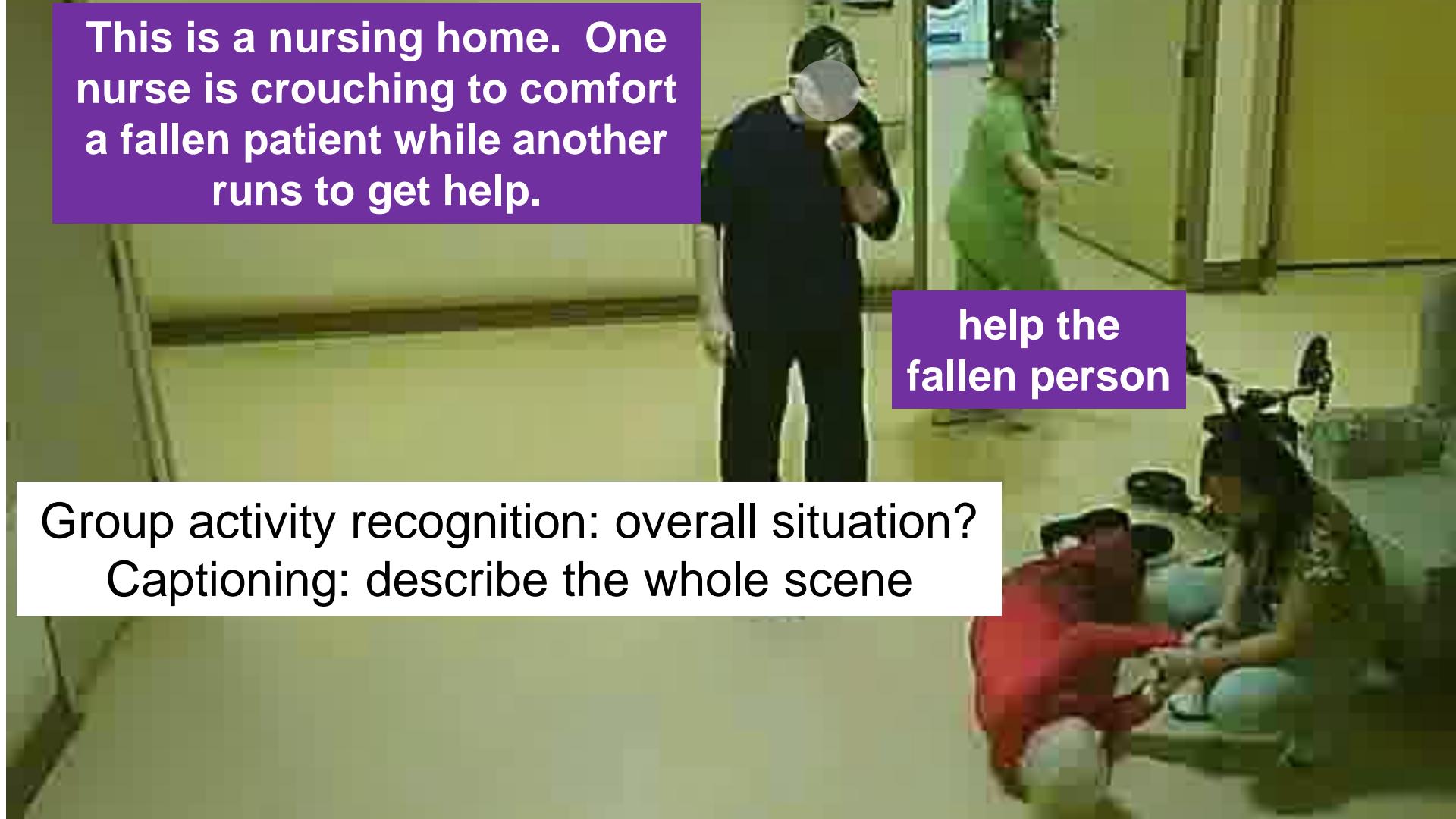




Intention/social role: why are they doing this?

comfort

This is a nursing home. One nurse is crouching to comfort a fallen patient while another runs to get help.



help the  
fallen person

Group activity recognition: overall situation?  
Captioning: describe the whole scene



indoor scene

This is a nursing home. One nurse is crouching to comfort a fallen patient while another runs to get help.

long term care facility

get help

run

help the fallen person

walker

chair

“AI-complete”: full semantic understanding necessary for success

floor

fall

comfort

squat

# Why is computer vision important?

## User videos



~300 hours of videos per minute

- Video indexing and retrieval

## Monitoring cameras



Streaming videos 24/7

- Surveillance
- Patient/elderly monitoring

## Media



## Wearables/robots



Content analysis, experience enrichment

- Recommendation systems
- Advertising
- Sports analytics

Streaming videos to be analyzed in real-time

- Lifelogging
- Robot operations and actions

# Why is it difficult?

- **Large variations in appearance:** occlusions, non-rigid motion, view-point changes, clothing, etc.

Action *Hugging*:



...

- **Manual collection of training samples is prohibitive:** many classes, rare occurrence



...

- **Vocabulary is not well-defined**

Action *Open*:



...

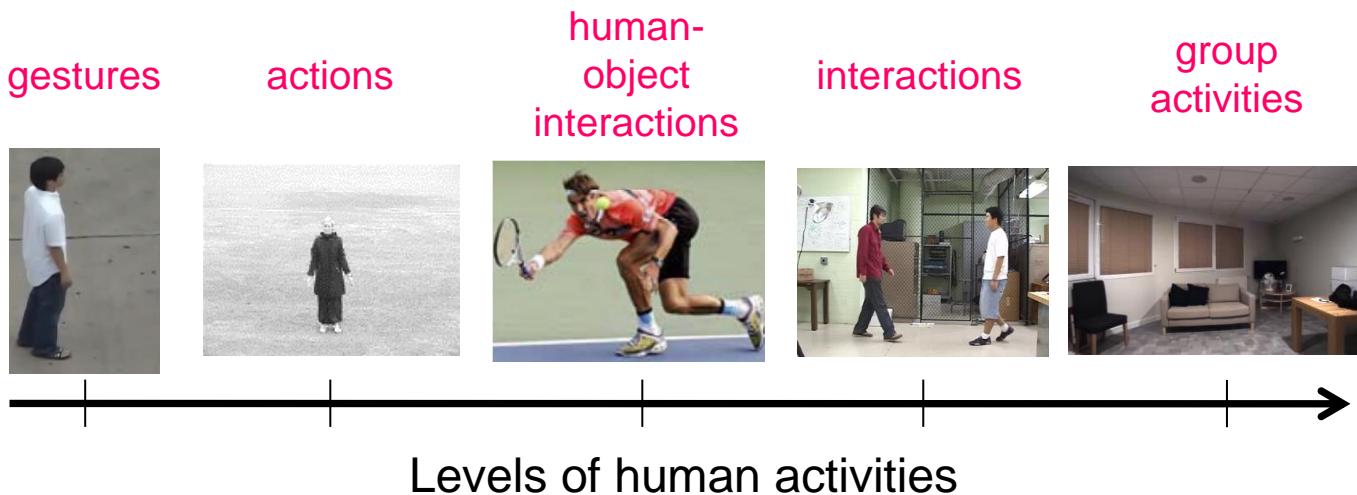
# Why is it difficult?

## Humans Outperform Computers

- Humans **generalize** better than deep learning systems
- Humans can handle **out of distribution** data
- Humans are more robust to **domain shifts**
- Humans can handle **new** situations, objects, environments
- Humans are better in building **abstractions**

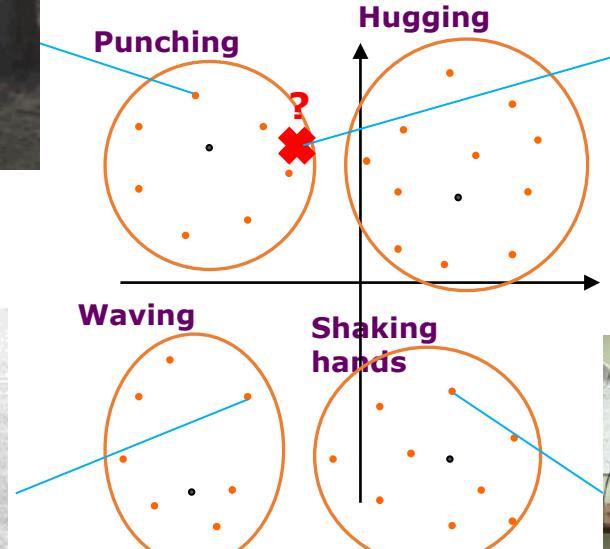
# Example: Human activity recognition

- There are various types/levels of activities
  - The ultimate goal is to make computers recognize all of them reliably

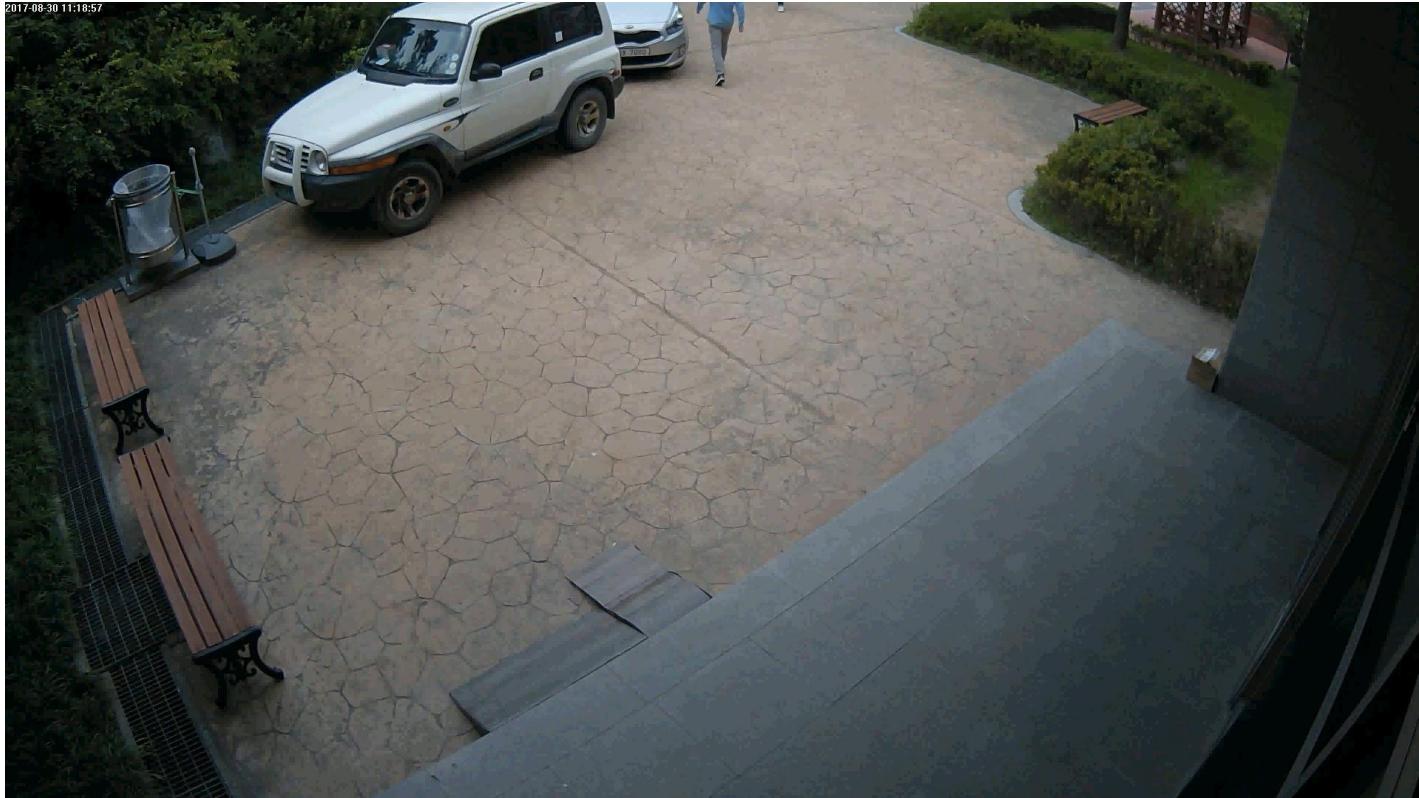


# Example: Activity classification

- Categorization of segmented videos
  - Input = a video segment containing only one activity

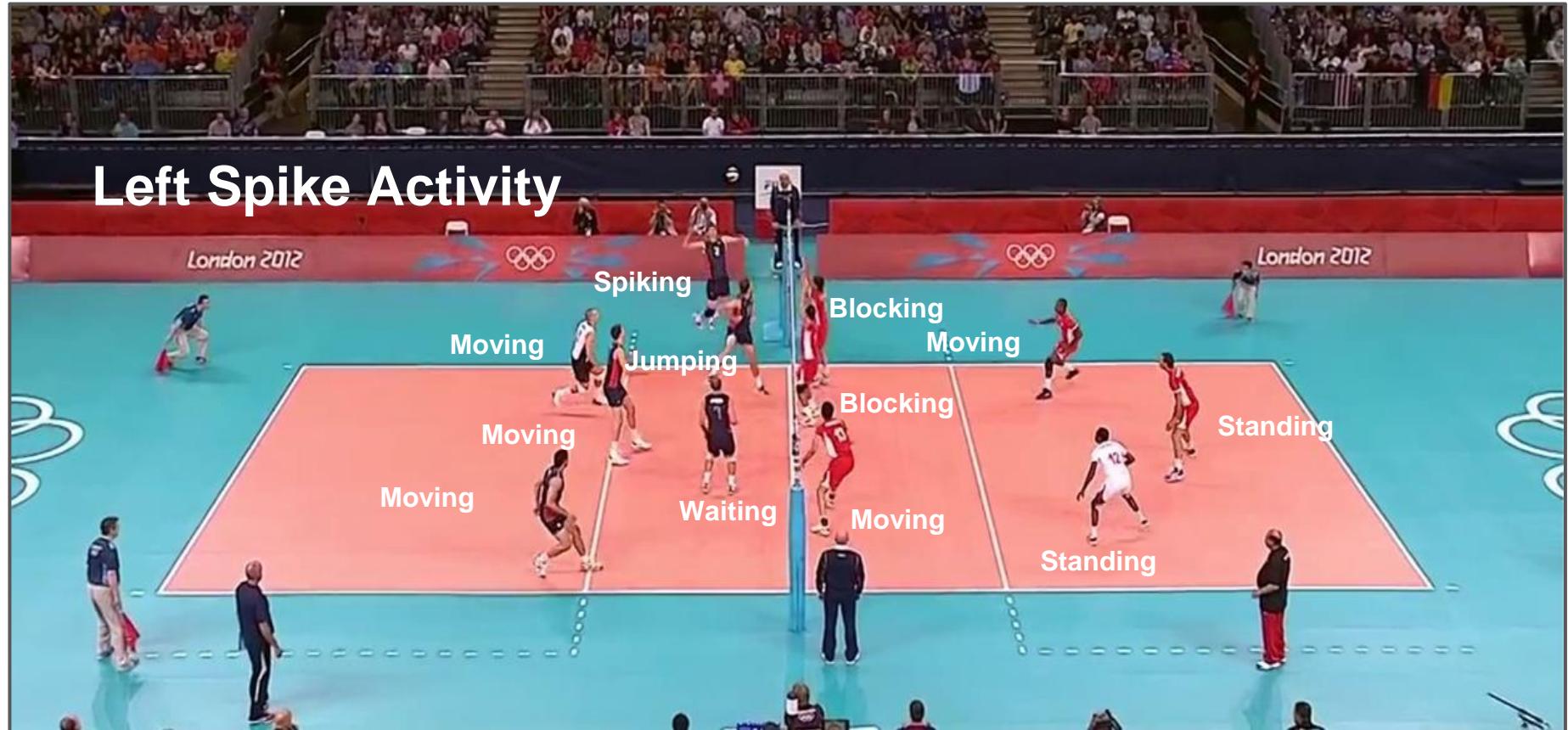


# Example: Action Recognition



# Group activity recognition

## Left Spike Activity



# Relationships?

**Left team:**

- **Spike** by a player
- **Fake jump** by another player



**Right team:**

- 2 players jump to **block**

# Learning Vision: Human vs. Computer

## **Computer Vision: Static Images (or Videos)**

- Directly learn to detect objects
- Learn to segment

## **Human children: Sequence of related images**

- Learn unity of objects
- Learn object permanence
- Learn to solve the binding problem

# Major limitation of deep learning

**Not data efficient:** Learning requires **millions** of labeled examples, models do not generalize well to new domains; not like humans!



# “What you saw is not what you get”



What your net is trained on



What it is asked to label

**“Dataset Bias”**  
**“Domain Shift”**  
**“Domain Adaptation”**  
**“Domain Transfer”**

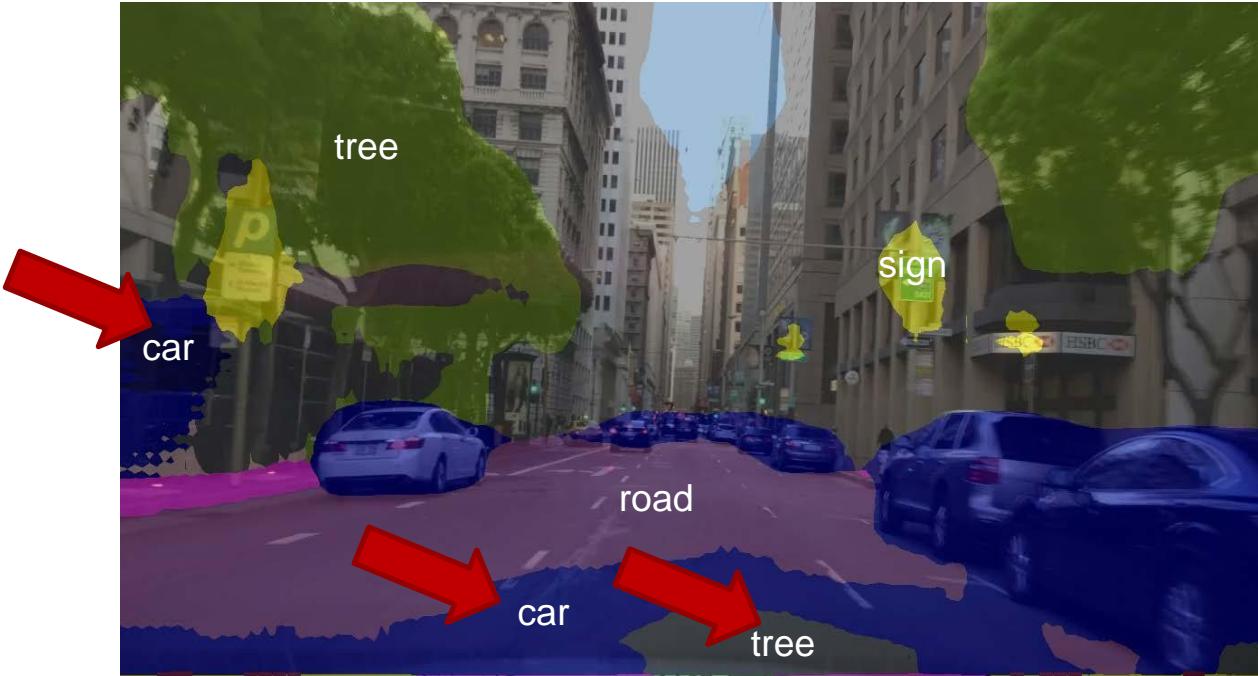
# Example: Scene segmentation



Train on Cityscapes, Test on [Cityscapes](#)

J. Hoffman, D. Wang, F. Yu, T. Darrell, FCNs in the Wild: Pixel-level Adversarial and Constraint-based Adaptation, Arxiv 2016

# Domain shift: Cityscapes to SF

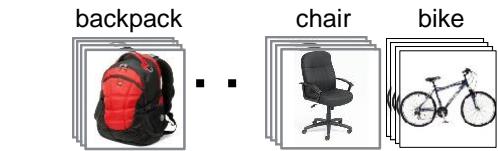
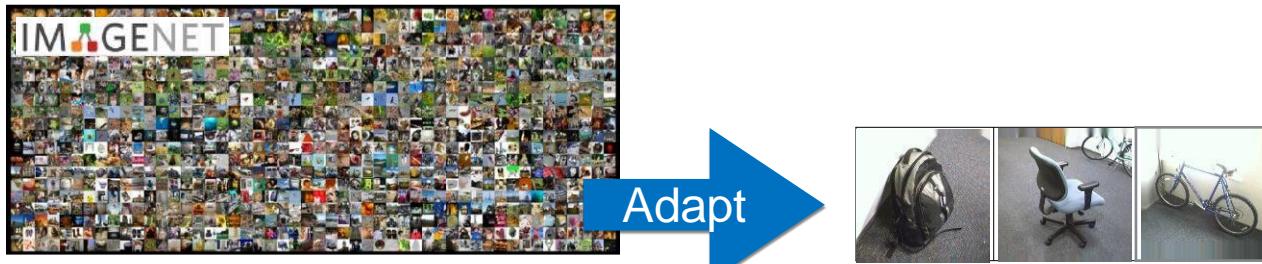


Train on Cityscapes, Test on San Francisco Dashcam

# No tunnels in CityScapes?...



# Domain Adaptation from source to target distribution



Source Domain  $\sim P_S(X, Y)$   
lots of **labeled** data

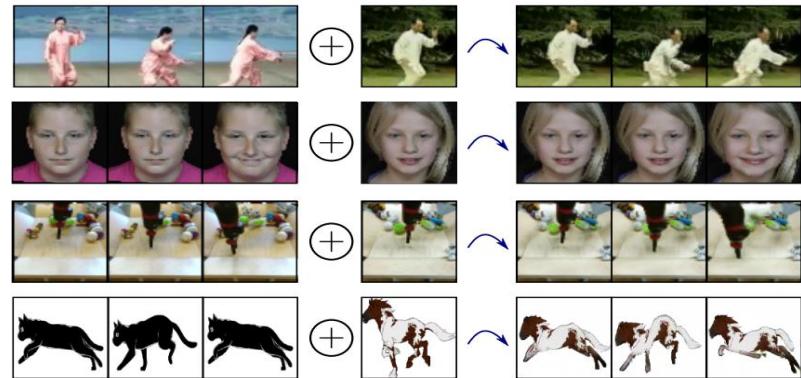
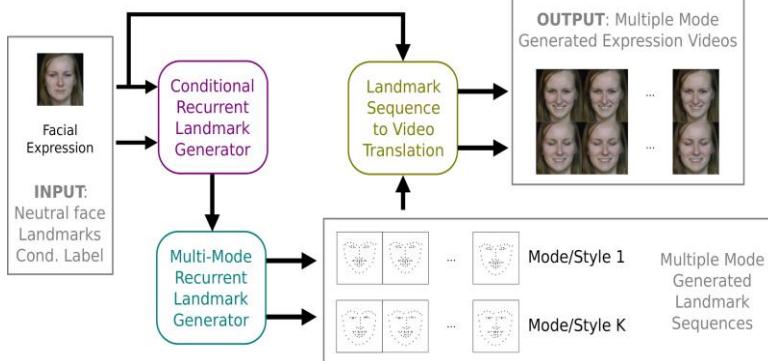
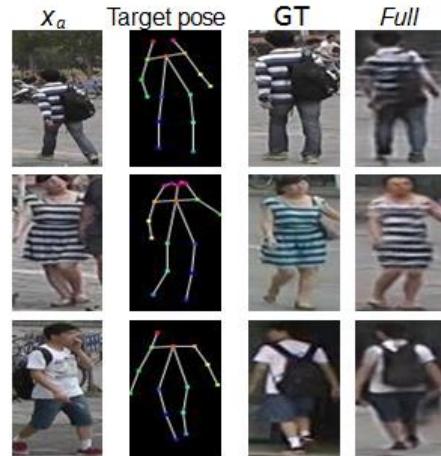
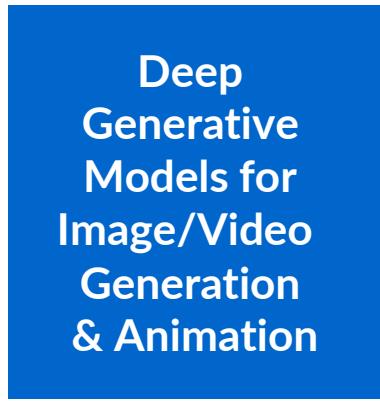


Target Domain  $\sim P_T(Z, H)$   
unlabeled or limited labels

$$D_S = \{(\mathbf{x}_i, y_i), \forall i \in \{1, \dots, N\}\}$$

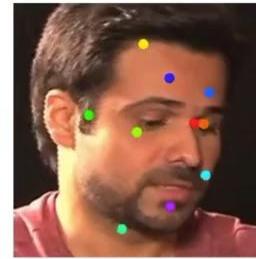
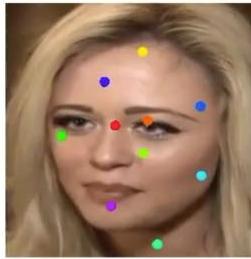
$$D_T = \{(\mathbf{z}_j, ?), \forall j \in \{1, \dots, M\}\}$$

# Domain Shift: Image and Video Generation



Arbitrary Object Animation without 3D modeling

# Learned Keypoints





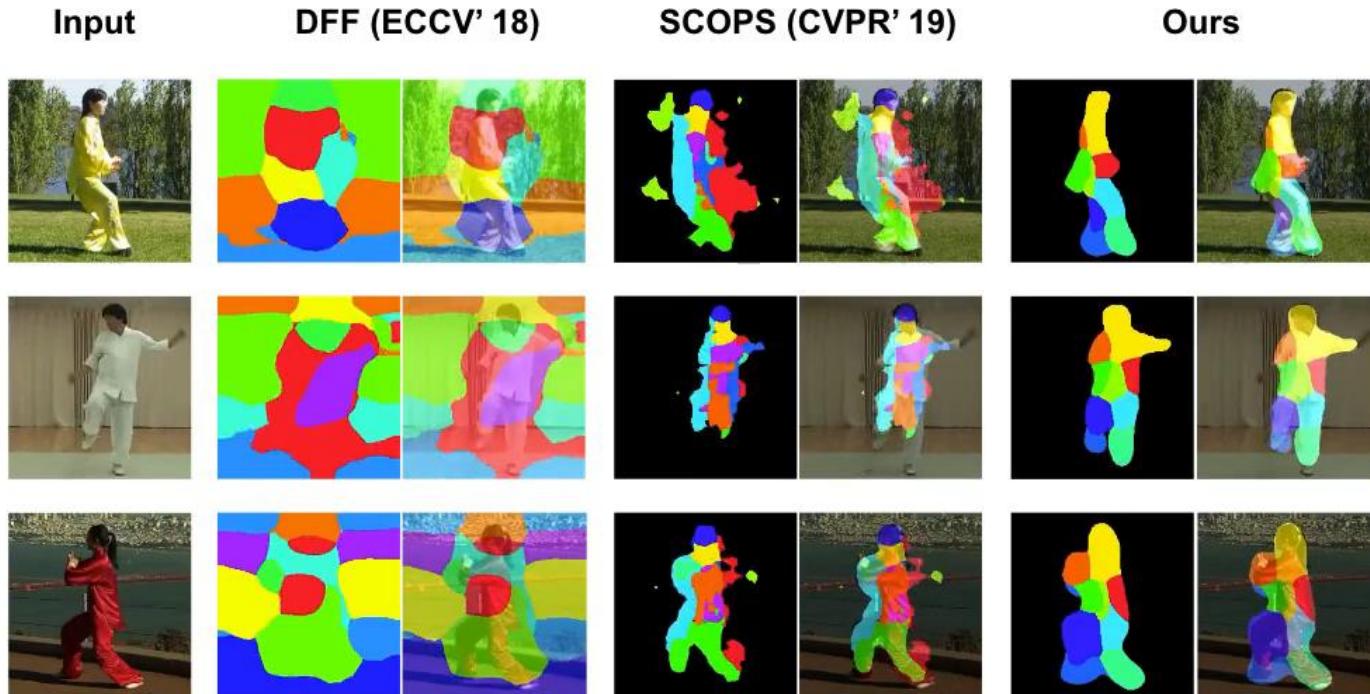
Kuhal

@KuhalCodesAI

# Deep Fakes ...



# Tai-Chi-HD



humans | SCALE



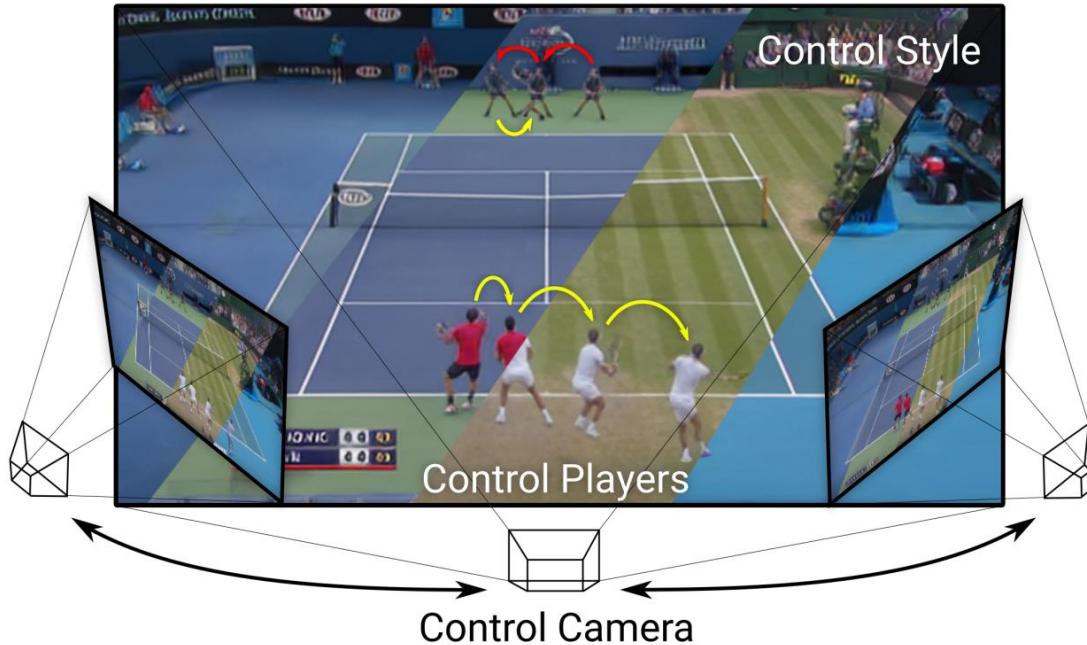
# Playable Video Generation



Menapace, et al., “Playable Video Generation”, in CVPR 2021

<https://github.com/willi-menapace/PlayableVideoGeneration>

# Playable Environments



- Learn a model that represents the observed environment
- Allow the user to input actions to the model through a controller at test time

# Playable Environments



Menapace et al, "Playable Environments: Video Manipulation in Space and Time", in CVPR 2022

<https://github.com/willi-menapace/PlayableEnvironments>

# Playable Environments



Menapace et al, "Playable Environments: Video Manipulation in Space and Time", in CVPR 2022

<https://github.com/willi-menapace/PlayableEnvironments>

# Playable Environments



Menapace et al, "Playable Environments: Video Manipulation in Space and Time", in CVPR 2022

<https://github.com/willi-menapace/PlayableEnvironments>



# Topics

- Scene recognition and understanding
- Object detection/localization/recognition
- Pixel-level prediction
  - Semantic segmentation
  - Depth estimation
- Motion analysis and activity recognition
- Biometrics: face, gesture, body pose, emotions
- Social signals processing

# Topics

- Image and video generation
- Tiny ML for computer vision
- Transfer learning and deep domain adaptation
- Explainable AI
- Privacy, transparency and ethics in vision

# Course organization

- Alternate oral presentations with practical sessions:
  - We will present in detail some of the basic code implementation to provide details on how things are done in practice
- Final evaluation
  - Students (groups) will be asked to choose from a list of recent important articles
  - Briefly present the main ideas
  - Propose extensions that could improve/generalize the approach

# Course timeline

- Classes on:
  - September 12, 13, 19, 20, 26, 27
  - October 3, 4, 17, 18 (pause for ACM MM and ECCV)
  - November 7, 8, 14, 15, 22, 28, 29
- Student presentations on:
  - December 6, 12, 13

# Scene Understanding and Object Recognition

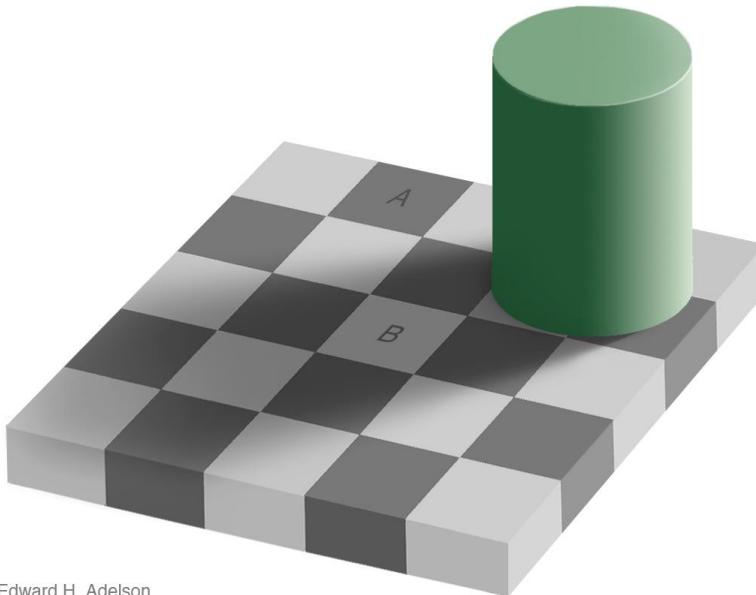
# Why is Vision Interesting?

- Psychology
  - ~ 35% of cerebral cortex is for vision
  - Vision is how we experience the world
- Engineering
  - Want machines to interact with the world
  - Digital images are everywhere

# Measurement vs. Perception

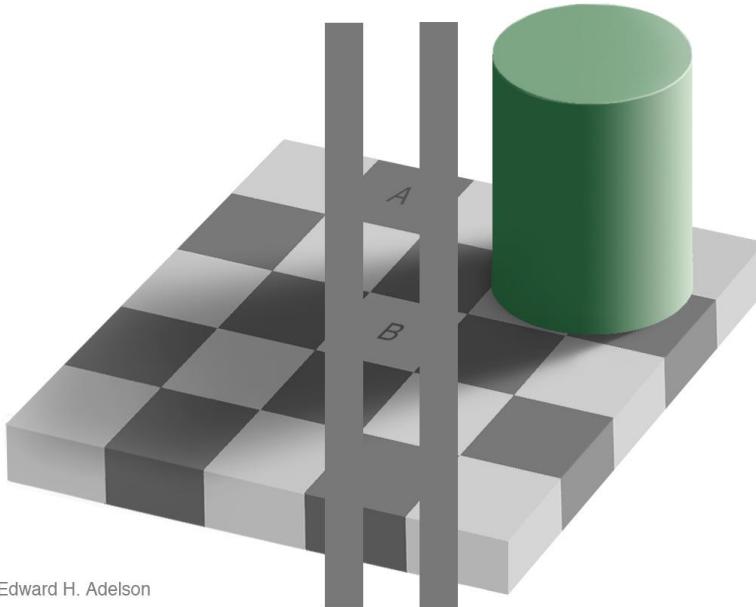


# Brightness: Measurement vs. Perception



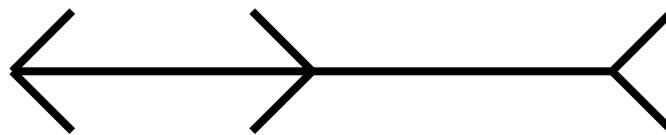
Edward H. Adelson

# Brightness: Measurement vs. Perception

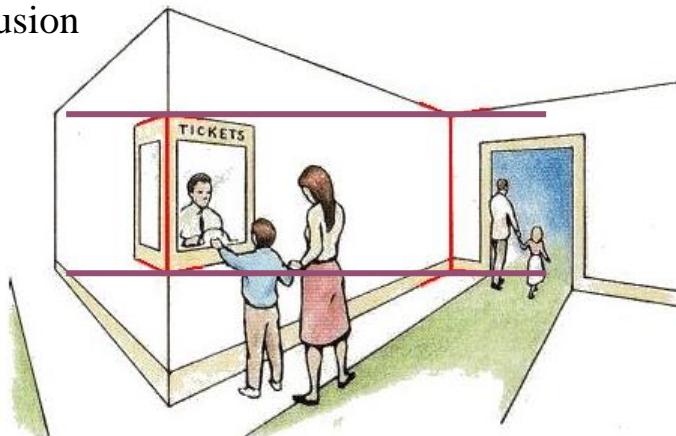


Proof!

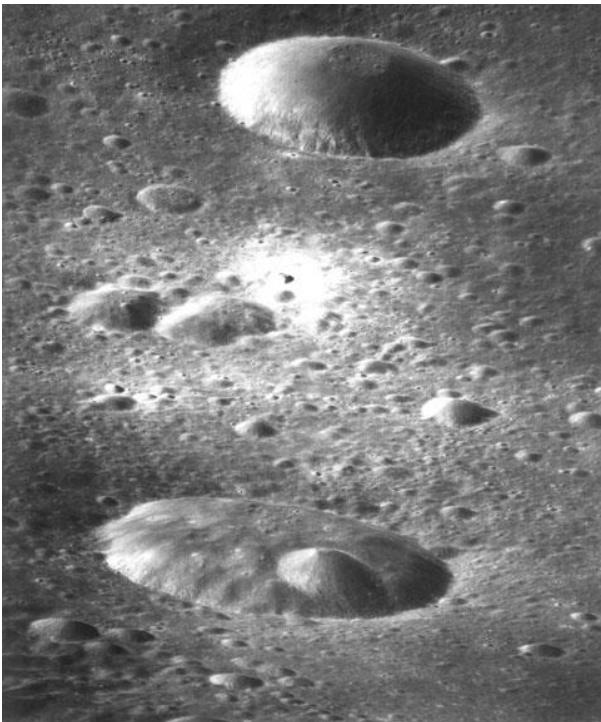
# Lengths: Measurement vs. Perception



Müller-Lyer Illusion



# Vision is Inferential: Prior Knowledge



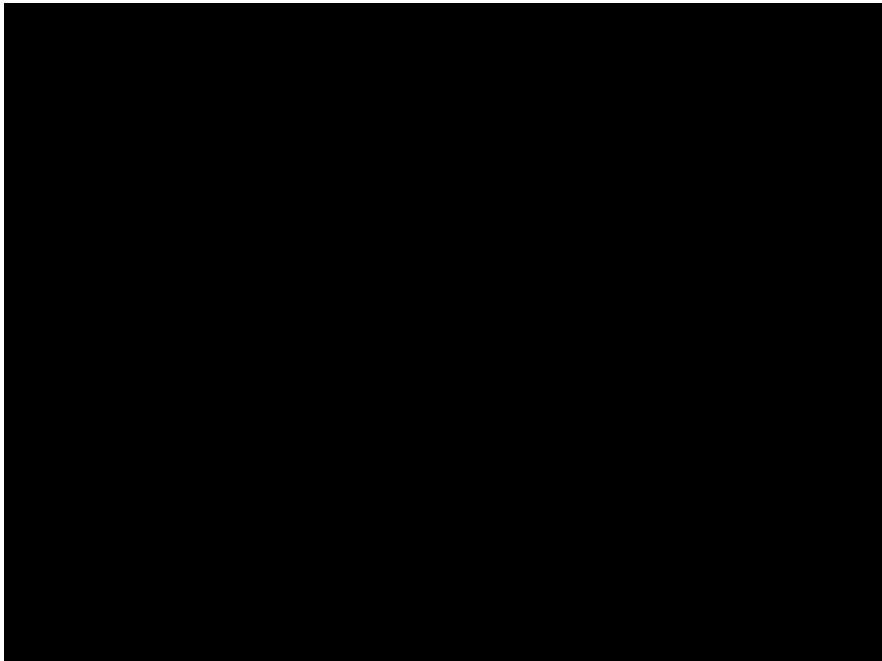
# Vision is Inferential: Prior Knowledge

- Our visual system is very sophisticated
- Humans can interpret images successfully under a wide range of conditions – even in the presence of very limited cues



# Playing with Perspective

- Perspective gives us very strong depth cues  
⇒ hence we can perceive a 3D scene by viewing its 2D representation (i.e. image)
- An example where perception of 3D scenes is misleading:



“Ames room”

A clip from "The computer  
that ate Hollywood"  
documentary. Dr.  
Vilayanur S.  
Ramachandran.

- “*What if I don’t care about this wishy-washy human perception stuff? I just want to make my robot go!*”
- Small Reason:
  - For measurement, other sensors are often better (in DARPA Grand Challenge, vision was barely used!)
  - For navigation, you still need to learn!
- Big Reason:
- The goals of computer vision (**what + where**) are in terms of what humans care about

# So what do humans care about?



# Verification: is that a bus?



# Detection: are there cars?



# Identification: is that a picture of Mao?

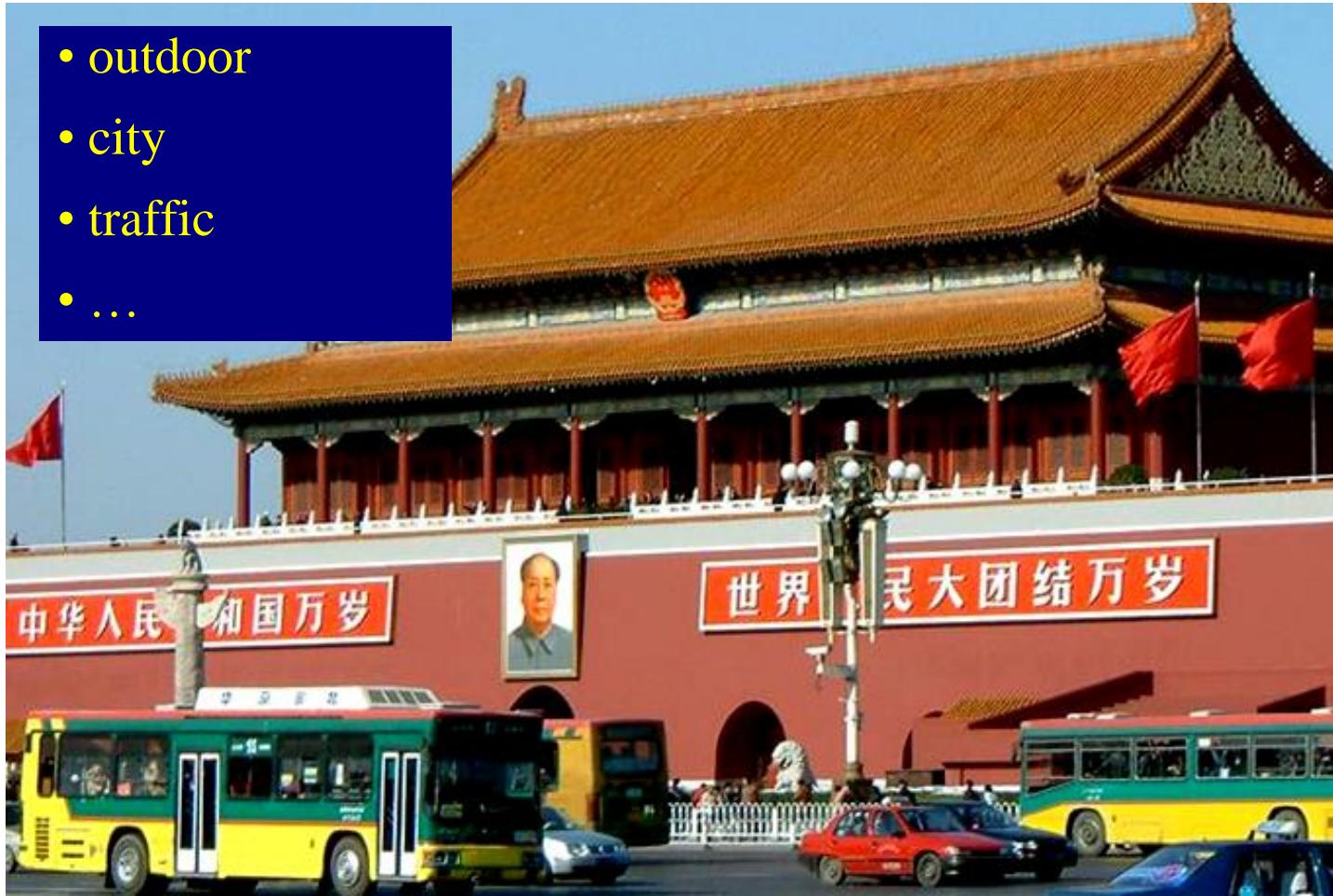


# Object categorization



# Scene and context categorization

- outdoor
- city
- traffic
- ...



# Rough 3D layout, depth ordering



# Challenges 1: view point variation

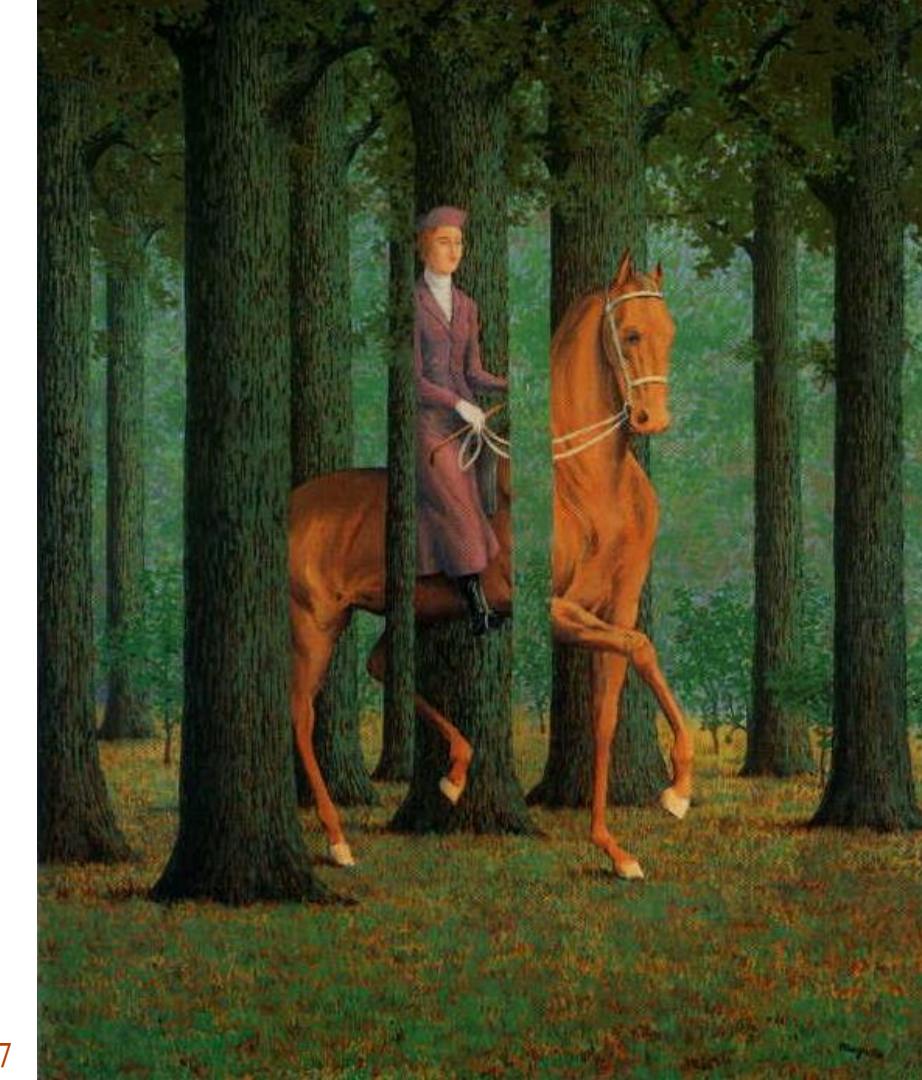


Michelangelo 1475-1564

## Challenges 2: illumination



## Challenges 3: occlusion

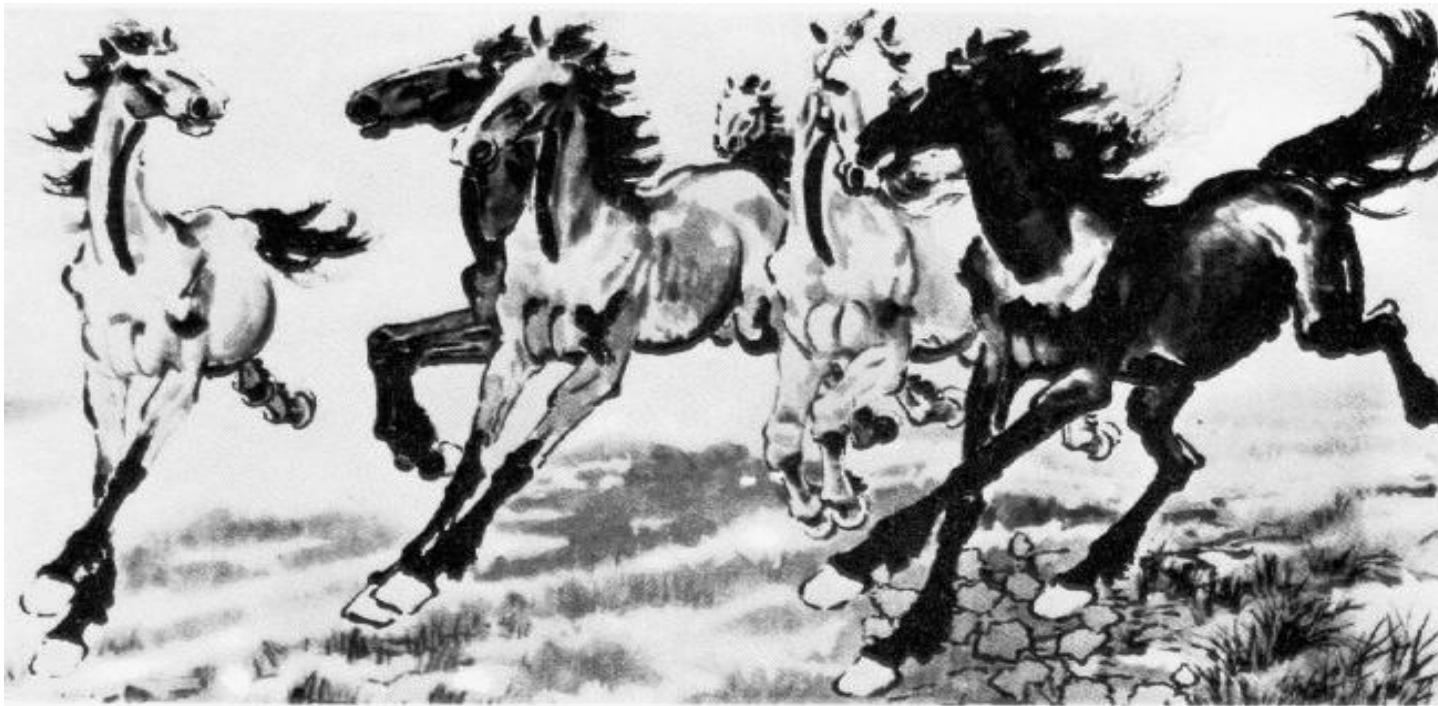


Magritte, 1957

## Challenges 4: scale



## Challenges 5: deformation



Xu, Beihong 1943

## Challenges 6: background clutter

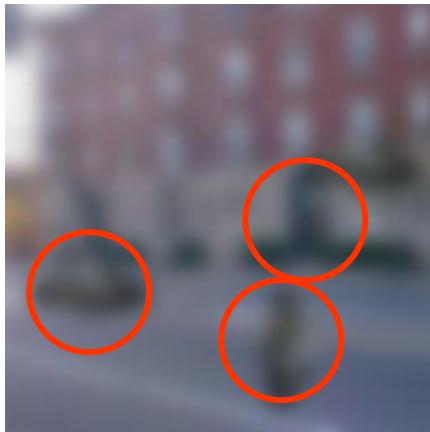
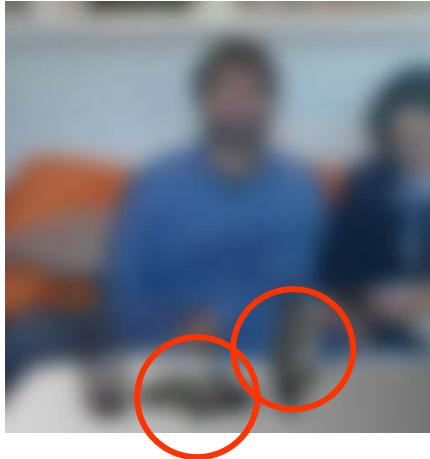


Klimt, 1913

# Challenges 7: object intra-class variation



## Challenges 8: local ambiguity



# Object detection and localization

# Computer Vision Tasks

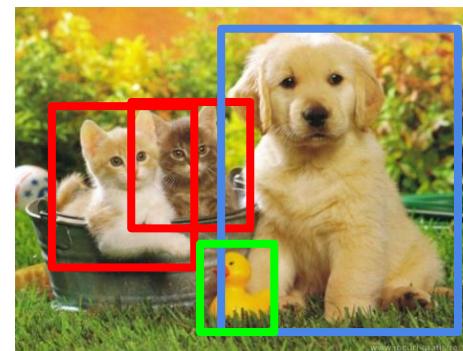
**Classification**



**Classification + Localization**



**Object Detection**



**Instance Segmentation**



CAT

CAT

CAT, DOG, DUCK

CAT, DOG, DUCK

Single object

Multiple objects

# Faces

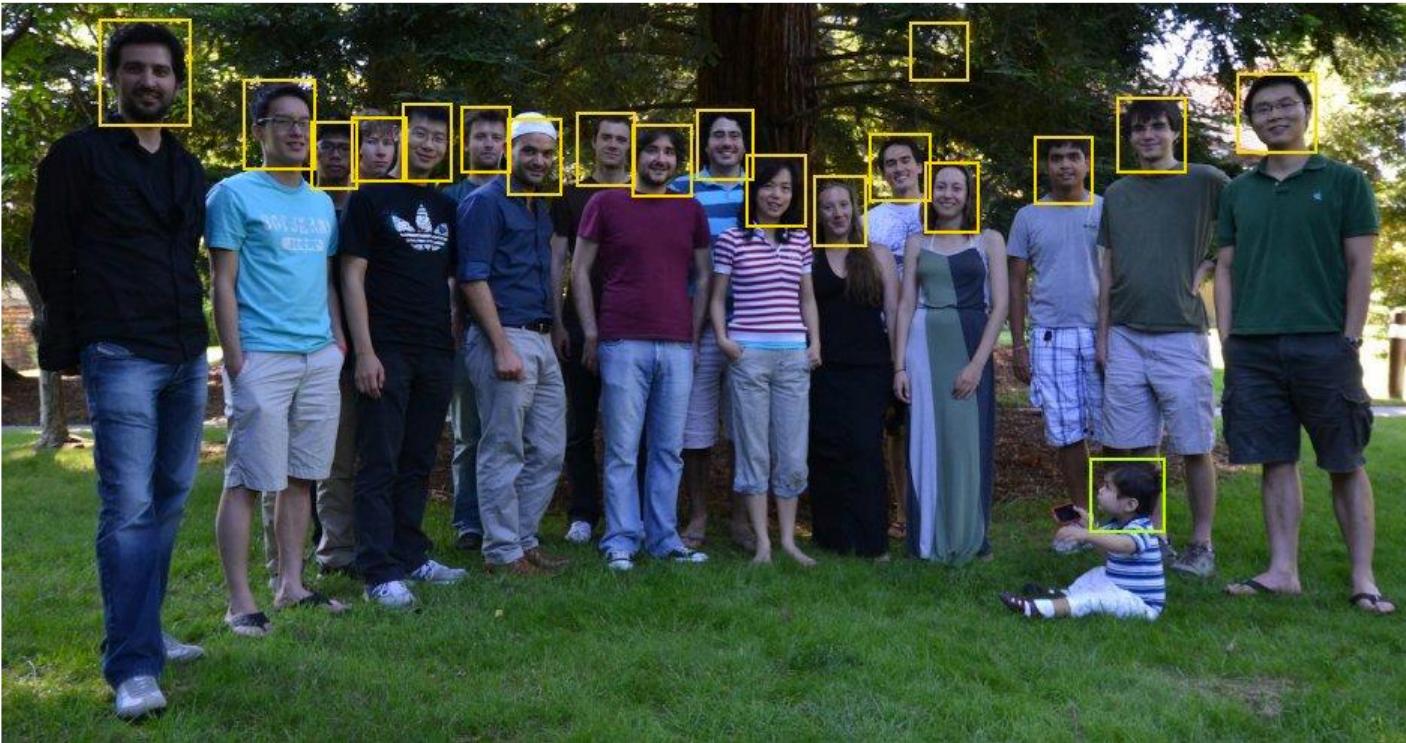
- Face detection
- One category: face
- Frontal faces
- Fairly rigid, unoccluded



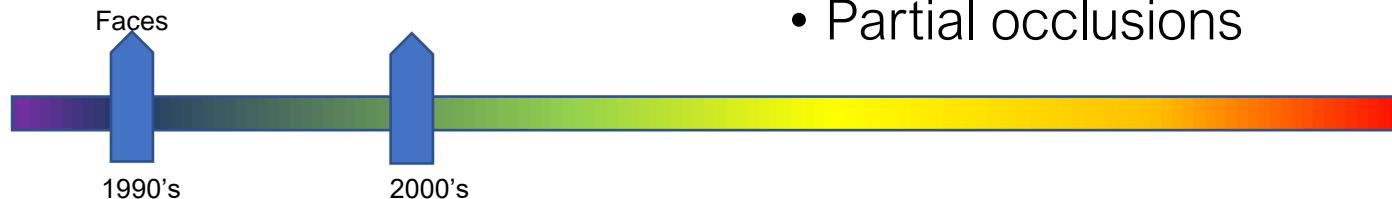
1990's

Human Face Detection in Visual Scenes. H. Rowley, S. Baluja, T. Kanade. 1995.

# Faces

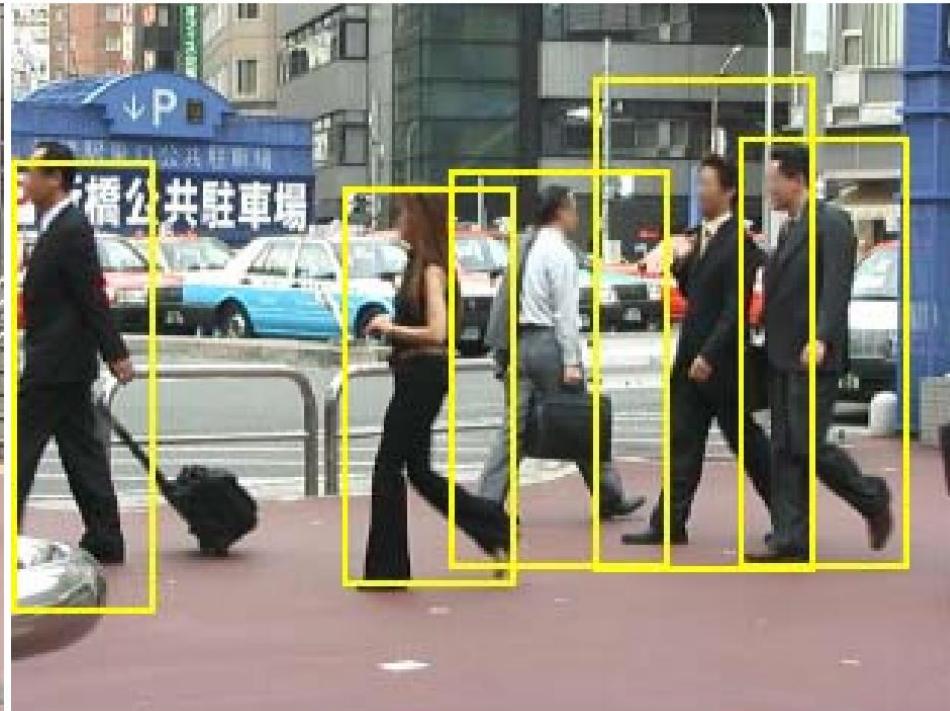
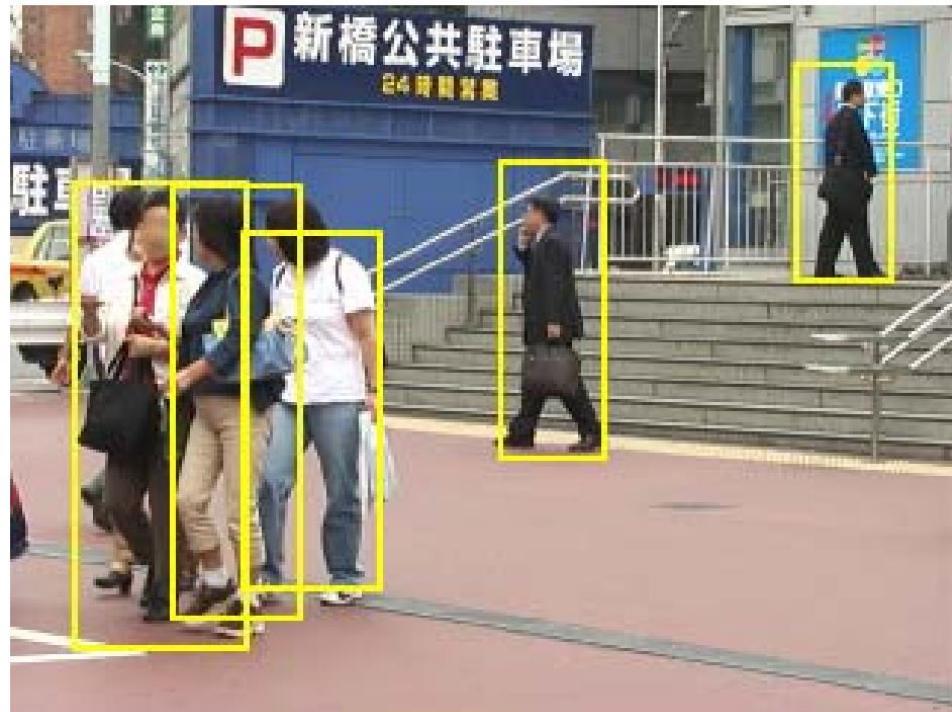


# Pedestrians



- One category: pedestrians
- Slight pose variations and small distortions
- Partial occlusions

# Pedestrians



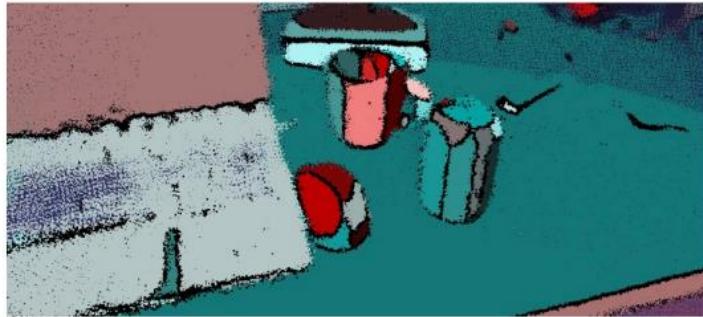
# Applications: Tagging People



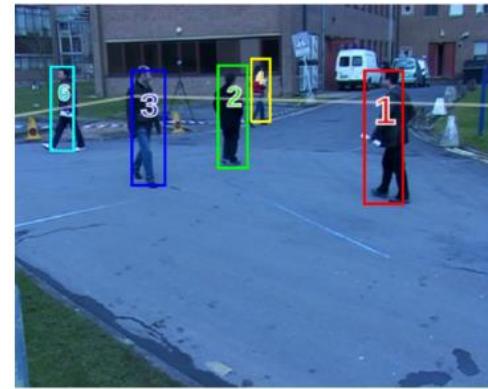
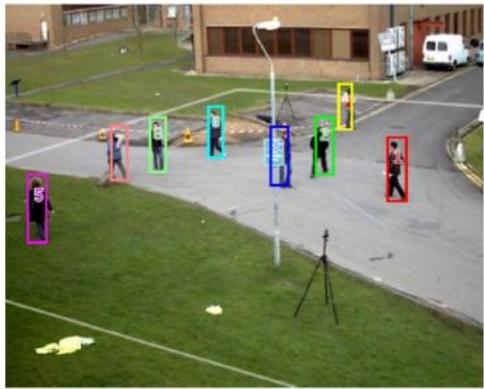
# Applications: Autonomous Driving



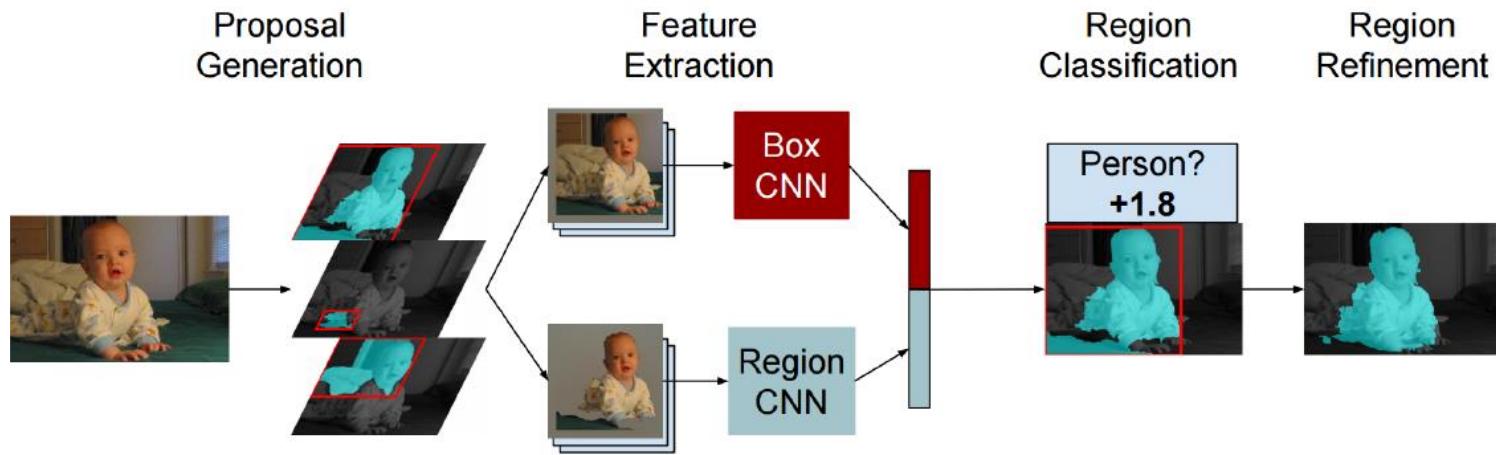
# Applications: Robotics



# Applications: Tracking

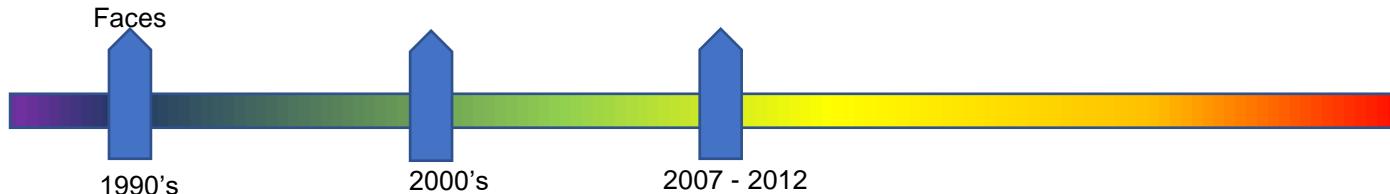
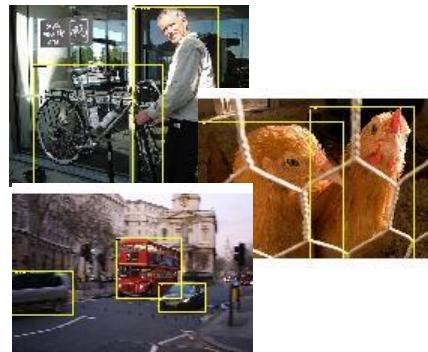


# Applications: Semantic Segmentation



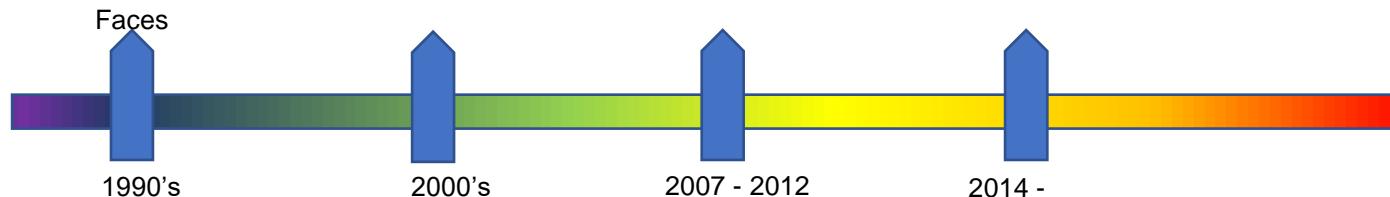
# PASCAL VOC

- 20 categories
- 10K images
- Large pose variations, heavy occlusions
- Generic scenes
- Cleaned up performance metric

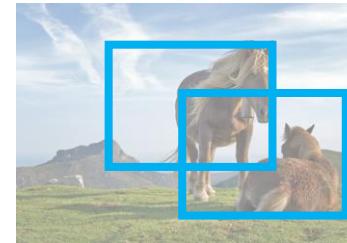
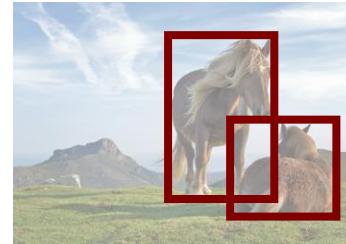


# Coco

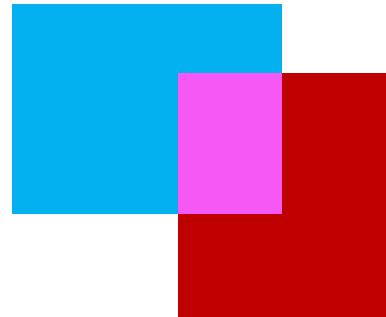
- 80 diverse categories
- 100K images
- Heavy occlusions, many objects per image, large scale variations



# Evaluation metric



# Matching detections to ground truth



$$IoU(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

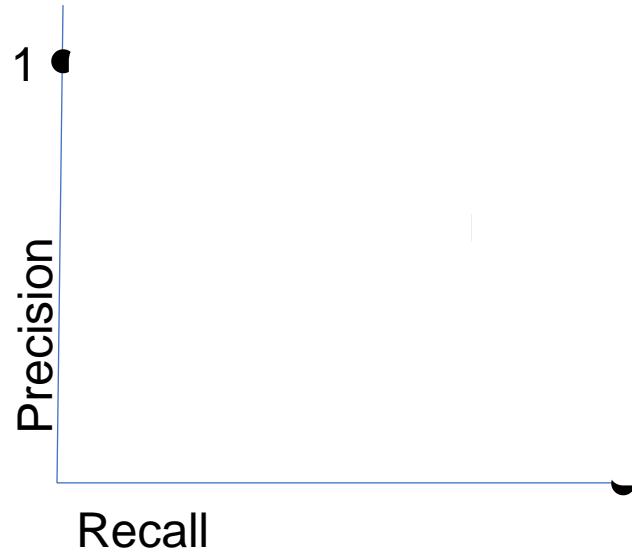
# Matching detections to ground truth

- Match detection to most similar ground truth
  - highest IoU
- If  $\text{IoU} > 50\%$ , mark as correct
- If multiple detections map to same ground truth, mark only one as correct
- **Precision** = #correct detections / total detections
- **Recall** = #ground truth with matched detections / total ground truth

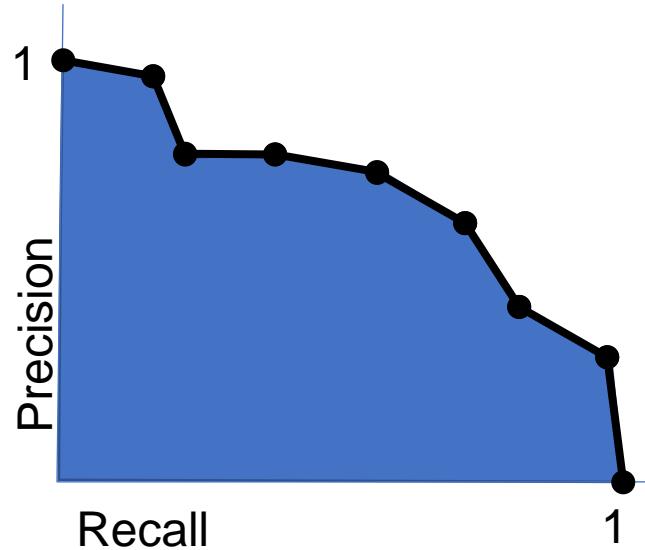
# Tradeoff between precision and recall

- ML usually gives scores or probabilities, so we need to threshold
- Too low threshold → too many detections → low precision, high recall
- Too high threshold → too few detections → high precision, low recall
- Right tradeoff depends on application
  - Detecting cancer cells in tissue: need high recall
  - Detecting edible mushrooms in forest: need high precision

# Average precision



# Average precision



# *Average average precision*

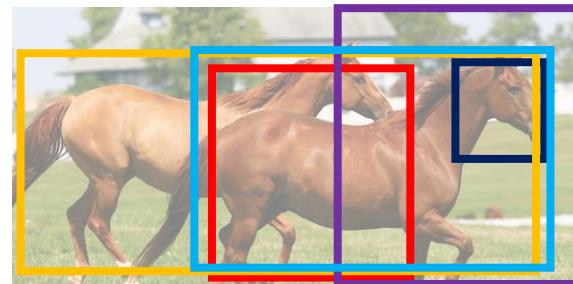
- AP marks detections with overlap  $> 50\%$  as correct
- But may need better localization
- Average AP across multiple overlap thresholds
- Confusingly, still called average precision
- Introduced in COCO

# Mean and category-wise AP

- Every category evaluated independently
- Typically report mean AP averaged over all categories
- Confusingly called “mean Average Precision”, or “mAP”

# Why is detection hard(er)?

- Precise localization



# Why is detection hard(er)?

- Much larger impact of pose



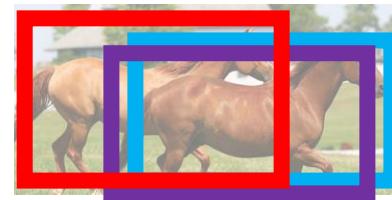
# Why is detection hard(er)?

- Occlusion makes localization difficult



# Why is detection hard(er)?

- Counting



# Why is detection hard(er)?

- Small objects

