SIC Saarland Informatics Campus

# High Level Computer Vision

# Introduction
# @ April 10, 2019

**Bernt Schiele & Mario Fritz**

**www.mpi-inf.mpg.de/hlcv/**

**Max Planck Institute for Informatics & Saarland University,**
**Saarland Informatics Campus Saarbrücken**

# Computer Vision and Multimodal Computing Group @ Max-Planck-Institute for Informatics

**Gerard Pons-Moll**
Real Virtual Humans

**Paul Swoboda**
Combinatorial Vision Group

**Bernt Schiele**
Computer Vision

**Zeynep Akata**
Multimodal Deep Learning
U Amsterdam

**Mario Fritz**
Scalable Learning & Perception
CISPA Helmholtz Center i.G.

**Max Planck Institute for Informatics & Saarland University,
Saarland Informatics Campus Saarbrücken**

# Computer Vision

- Lecturer:
  - ▸ Bernt Schiele (schiele@mpi-inf.mpg.de)
  - ▸ Mario Fritz (mfritz@mpi-inf.mpg.de)
- Assistants:
  - ▸ Yang He (yang@mpi-inf.mpg.de)
  - ▸ Rakshith Shetty (rshetty@mpi-inf.mpg.de)

- Language:
  - ▸ English
- mailing list for announcements etc.
  - ▸ send email (see instructions on the web)
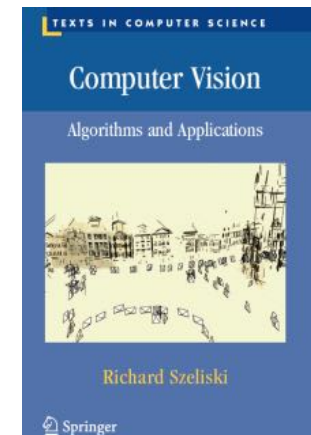    Rakshith Shetty <rshetty@mpi-inf.mpg.de>

# Lecture & Exercise

- Officially: 2V (lecture) + 2Ü (exercise)
  - ▸ Lecture:  Wed: 10:15am - 12pm (room 024)
  - ▸ Exercise: Mon: 10:15am - 12pm (room 024)

- typically 1 exercise sheet every 1-2 weeks
  - ▸ part of the final grade
  - ▸ some pencil and paper, mostly practical including a project
  - ▸ larger project in second half of lecture
    - we/you propose projects, mentoring, final presentation

- 1. exercise is Python tutorial

- Exam
  - ▸ oral exam (grading 50% oral exam and 50% exercises)
  - ▸ after the SS - there will be proposed dates

# Material

- For "non-deep-learning" parts of the lecture:

  ▸ available online
    http://szeliski.org/Book

- Background on deep learning:
  Deep Learning Book

  ▸ available online
    http://deeplearning.org

# Why Study Computer Vision

- Science
  - ▸ Foundations of perception. How do WE as humans see?
  - ▸ computer vision to explore "computational model of human vision"

- Engineering
  - ▸ How do we build systems that perceive the world
  - ▸ computer vision to solve real-world problems
    (e.g. self-driving cars to detect pedestrians)

- Applications
  - ▸ medical imaging (computer vision to support medical diagnosis, visualization)
  - ▸ surveillance (to follow/track people at the airport, train-station, ...)
  - ▸ entertainment (vision-based interfaces for games)
  - ▸ graphics (image-based rendering, vision to support realistic graphics)
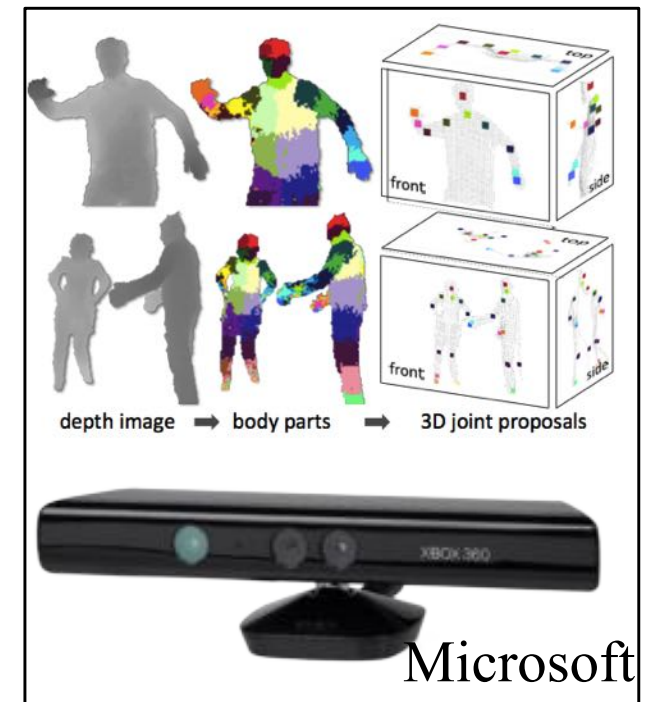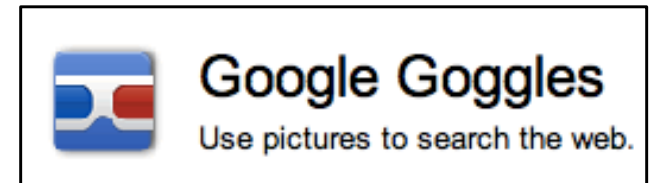  - ▸ car-industry (lane-keeping, pre-crash intervention, …)
  - ▸ …

# Some Applications

- License Plate Recognition

  ‣ London Congestion Charge

  ‣ http://www.cclondon.com/imagingandcameras.html

  ‣ http://en.wikipedia.org/wiki/London_congestion_charge

- Surveillance

  ‣ Face Recognition

  ‣ Airport Security (People Tracking)

- Medical Imaging

  ‣ (Semi-)automatic segmentation and measurements

- Autonomous Driving & Robotics

# More Applications

- Vision on Cellphones:
  - ▸ e.g. Google Goggles

- Vision for Interfaces:
  - ▸ e.g. Microsoft Kinect

- Reconstruction



Google Goggles
Use pictures to search the web.



depth image ➡ body parts ➡ 3D joint proposals

Microsoft



Photo Tourism
Exploring photo collections in 3D

Microsoft

(a)          (b)          (c)

# Goals of today's lecture

- First intuitions about

  ▶ What is computer vision?

  ▶ What does it mean to see and how do we (as humans) do it?

  ▶ How can we make this computational?

- Applications & Appetizers

- Role of Deep Learning

  - with several slides taken from Fei-Fei Li, Justin Johnson, Serena Yeung @ Stanford

- 2 case studies:

  ▶ Recovery of 3D structure

  - slides taken from Michael Black @ Brown University / MPI Intelligent Systems

  ▶ Object Recognition

  - intuition from human vision...

# Applications & Appetizers

… work from our group

# Detection & Recognition of Visual Categories



Challenges:
- multi-scale
- multi-view
- multi-class
- varying illumination
- occlusion
- cluttered background
- articulation
- high intraclass variance
- low interclass variance

# Challenges of Visual Categorization

- high intra-class variation
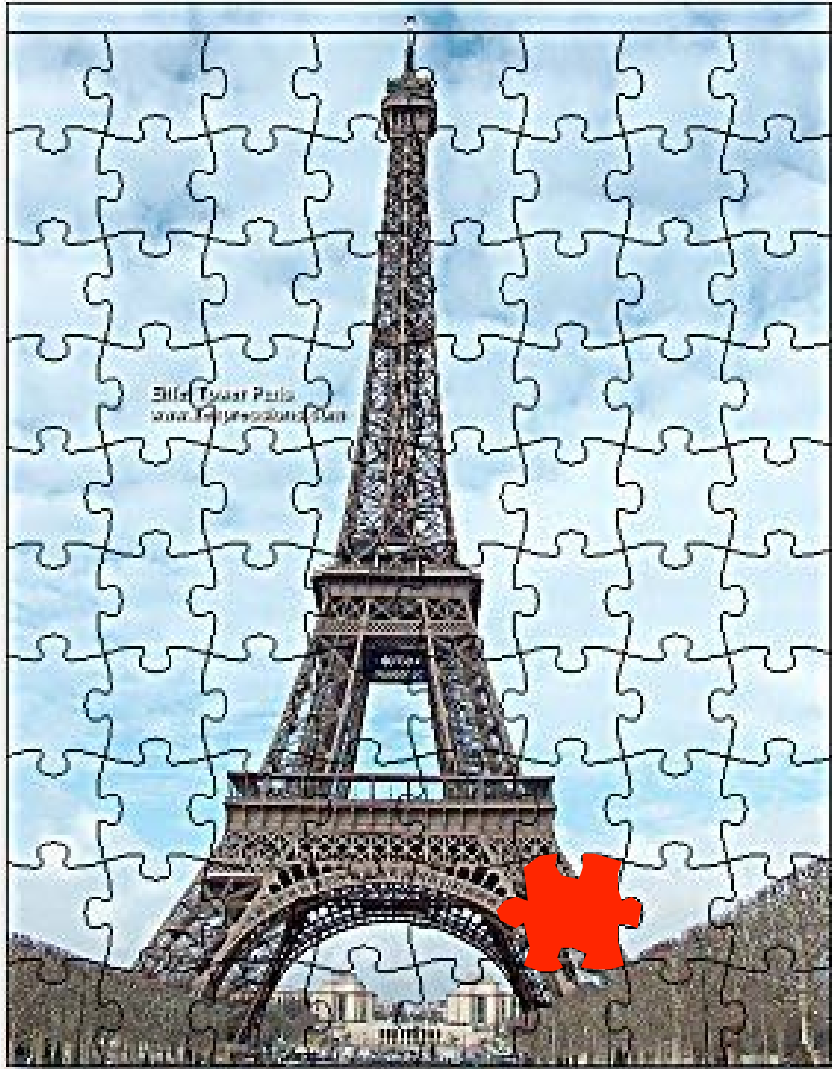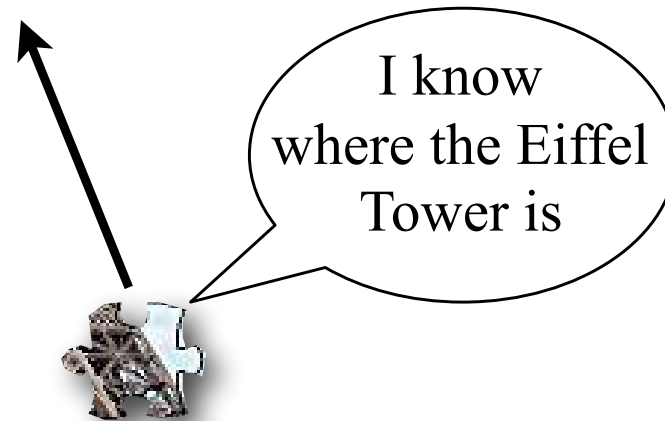


- low inter-class variation
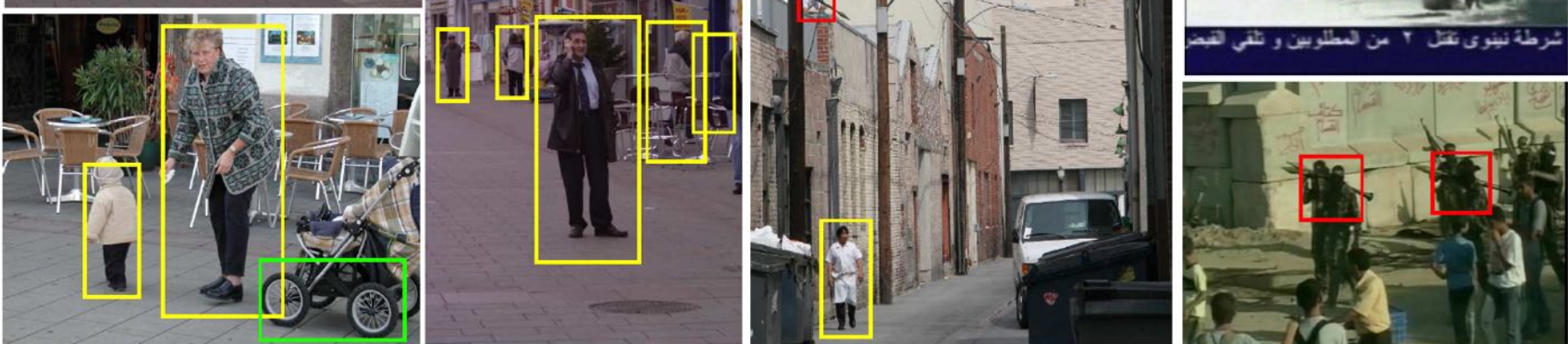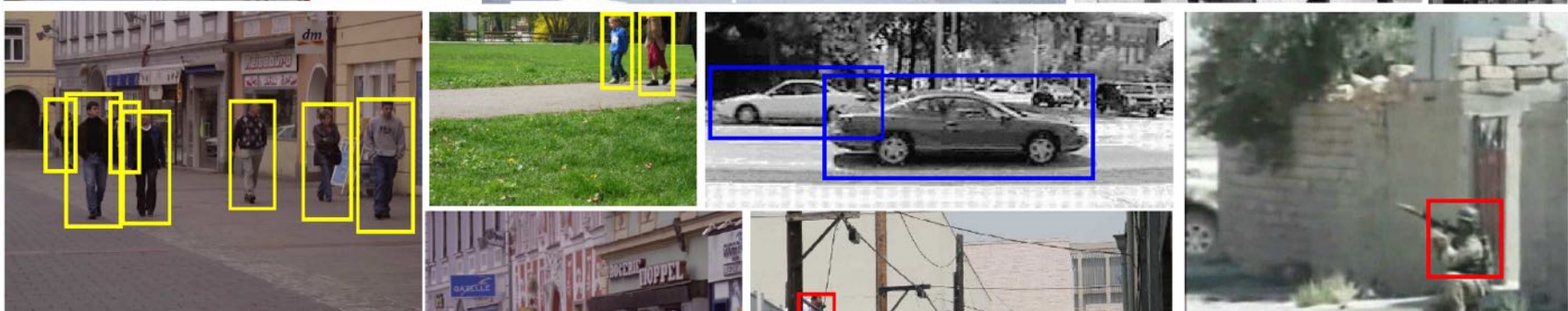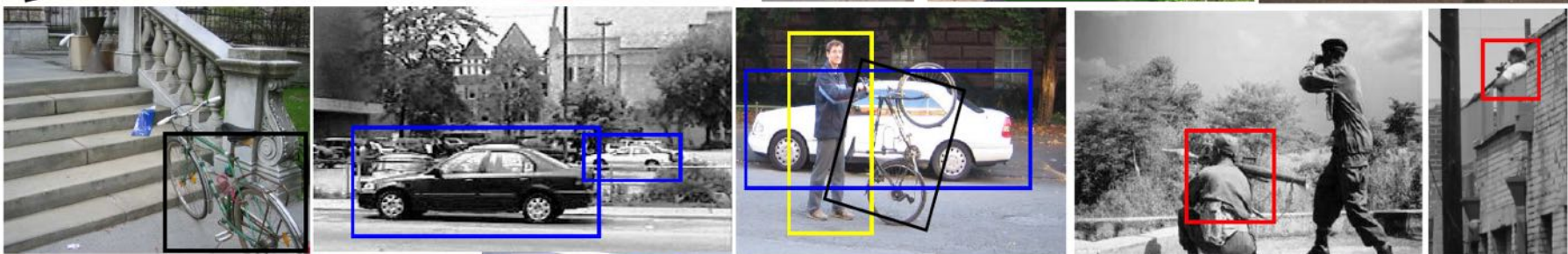
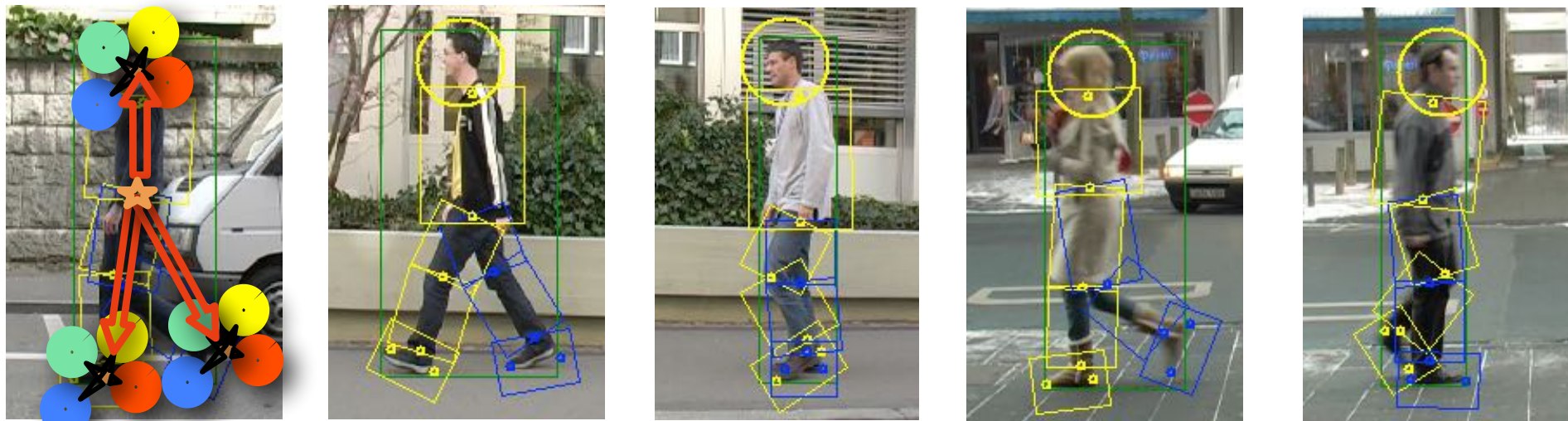- high intra-class variation

# Sample Category: Motorbikes

# Basic Idea



global

local

I know where the Eiffel Tower is

# Video...

# Articulation Model

- Assume uniform position prior for the whole body

- Learn the conditional relation between part position and body center from data:
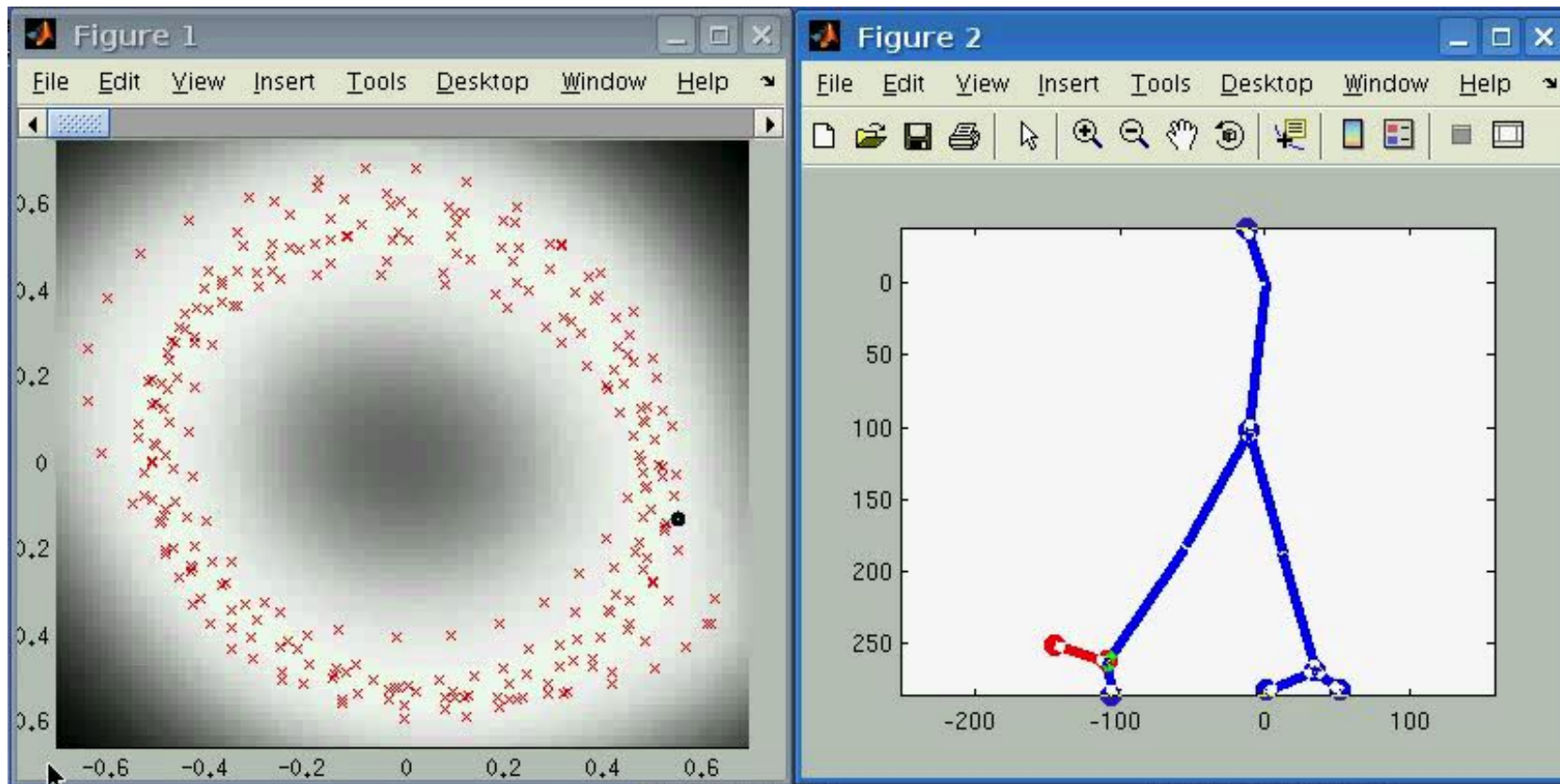
$$p(L|a) = p(\mathbf{x}^o) \prod_{i=1}^{N} p(\mathbf{x}^i | \mathbf{x}^o, a)$$
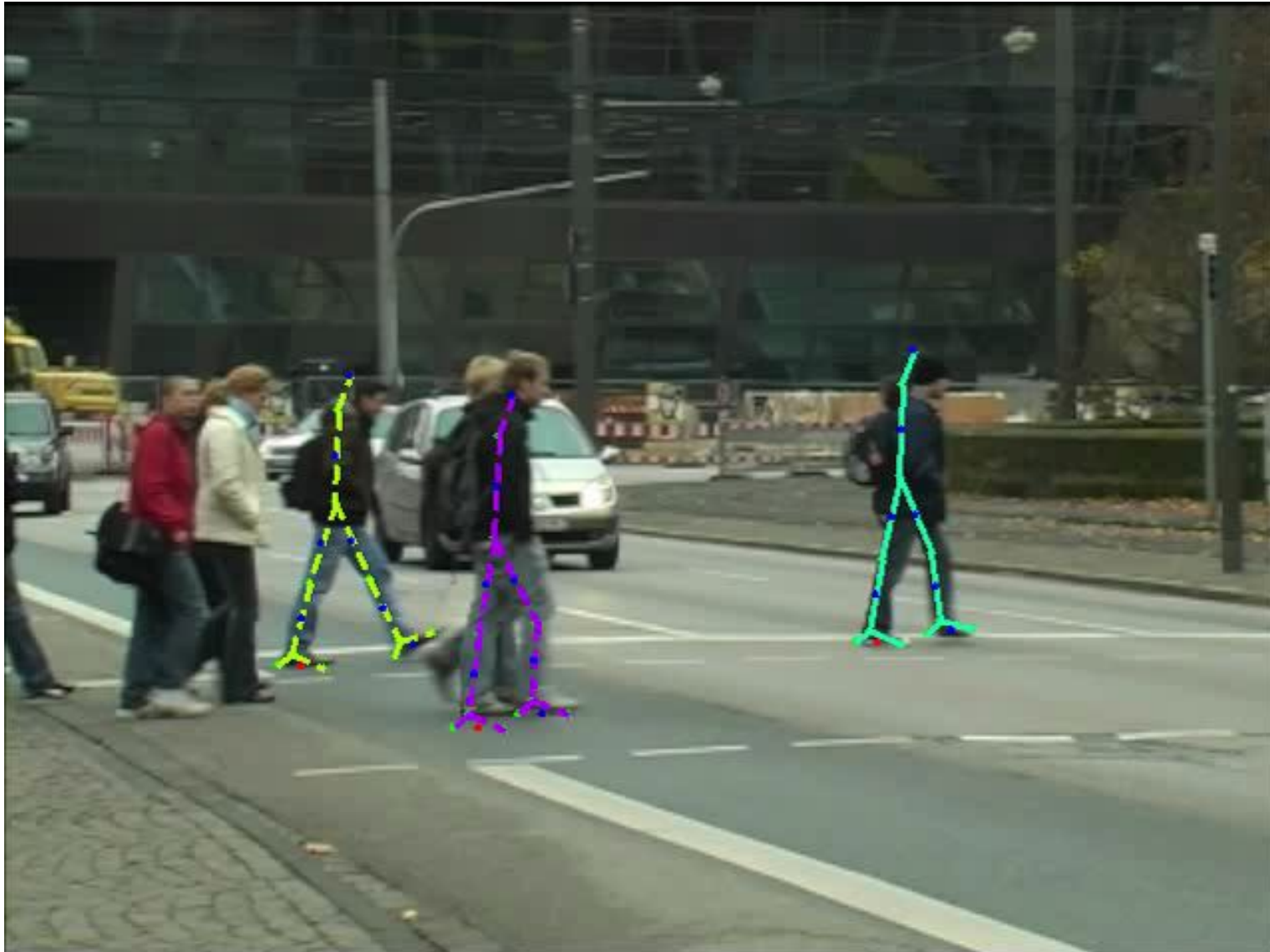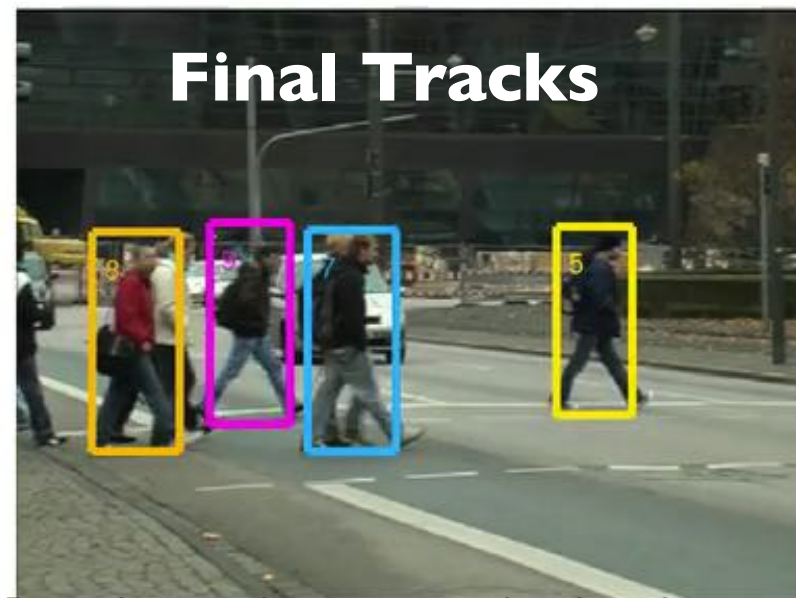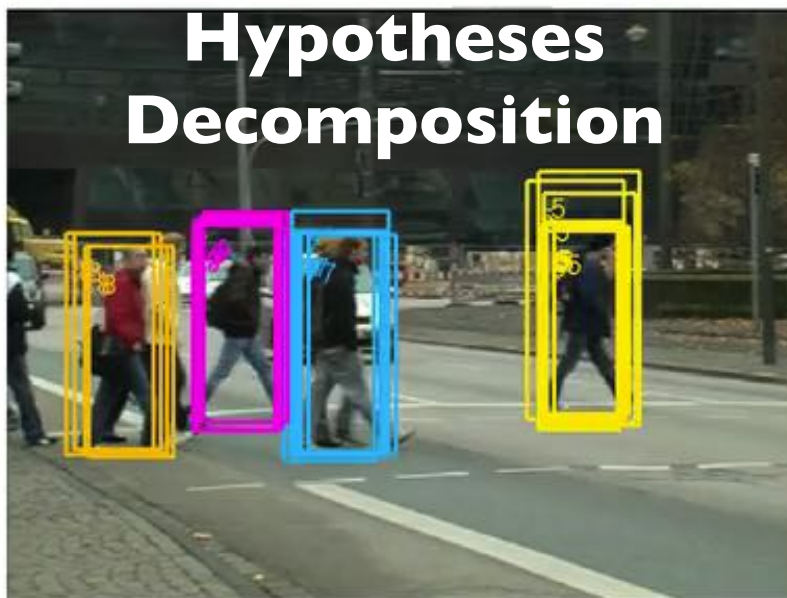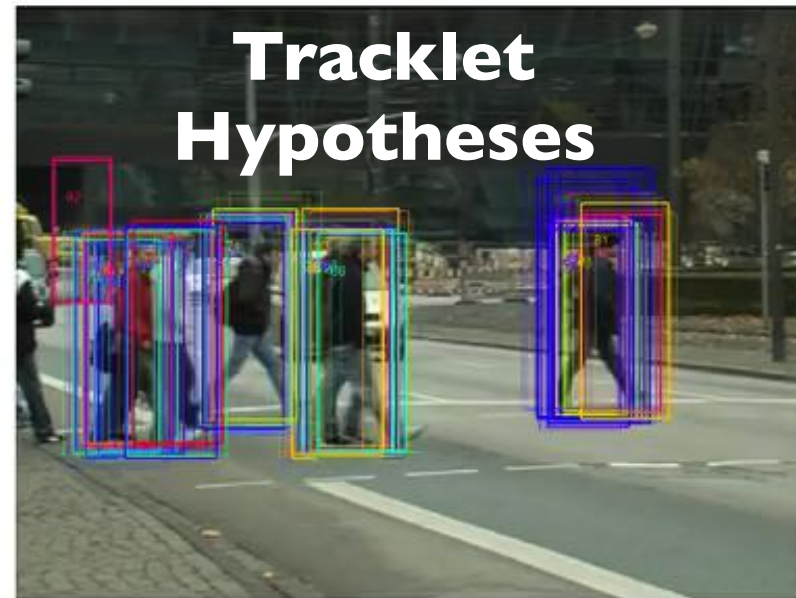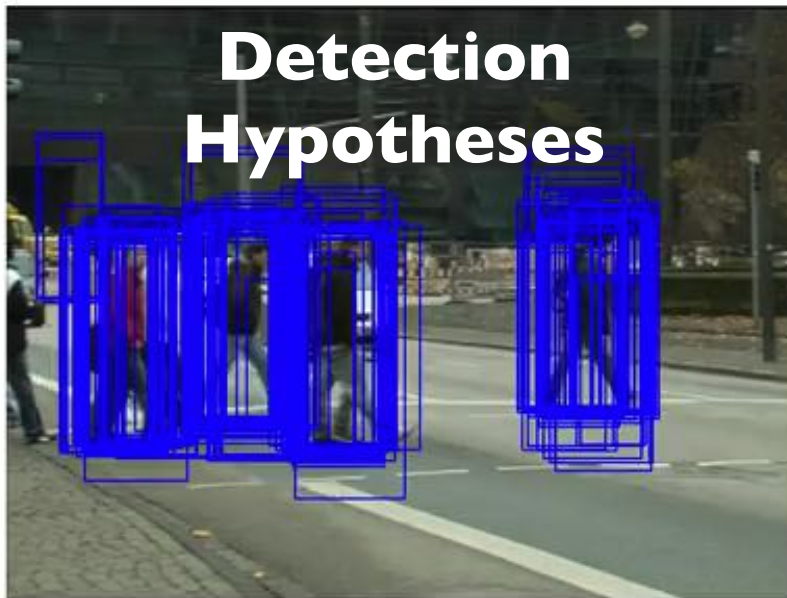


400 annotated training images

# Modeling Body Dynamics

- Visualization of the hierarchical Gaussian process latent variable model (hGPLVM)

# Our Subgraph Multicut Tracking Results



Detection Hypotheses

Tracklet Hypotheses

Hypotheses Decomposition

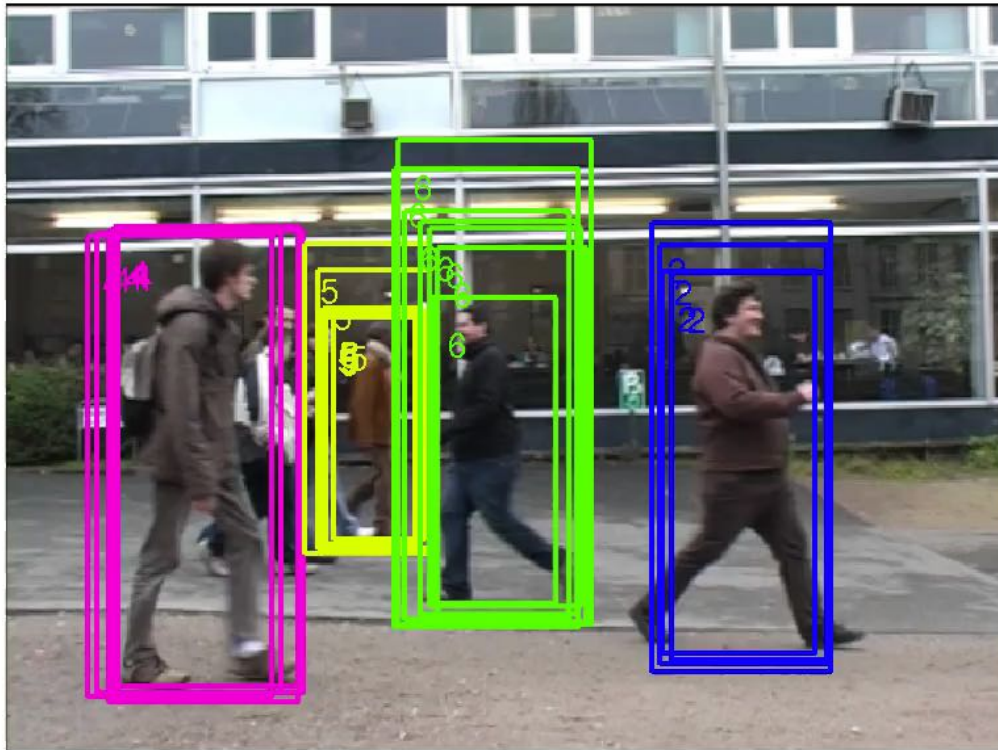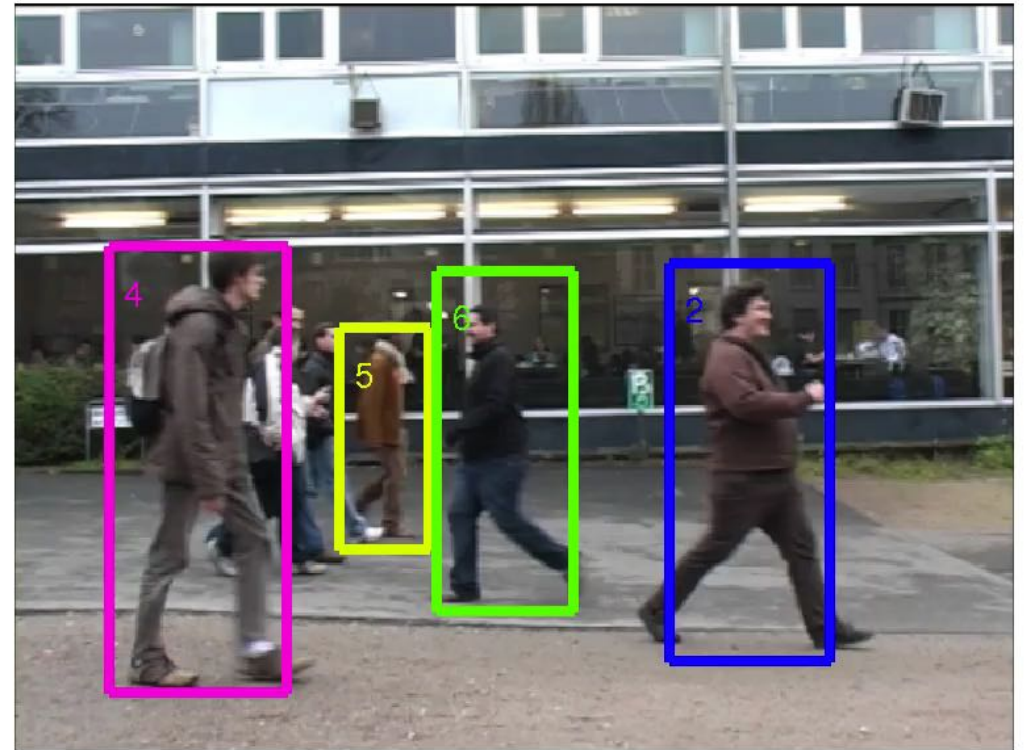Final Tracks

Dotted rectangles are interpolated tracks.

# More Results



**Decompositions (clusters)**

Dotted rectangles are interpolated tracks.



**Tracks**

# More Results



Dotted rectangles are interpolated tracks.

**Decompositions (clusters)**

**Tracks**

**Deep Learning
have become an important tool
for object recognition**

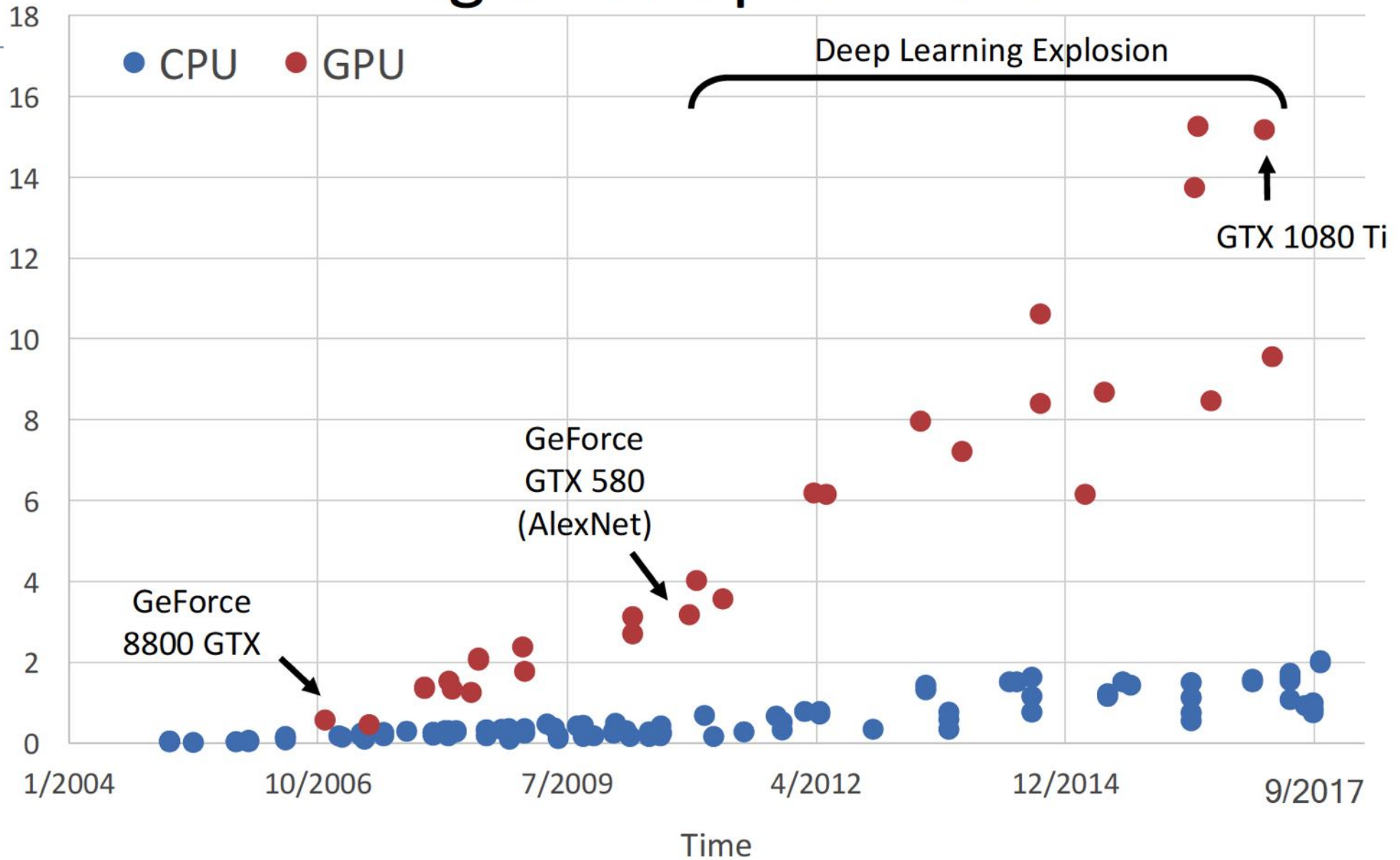(and other computer vision tasks)

Let's briefly discuss CNNs
(Convolutional Neural Networks)
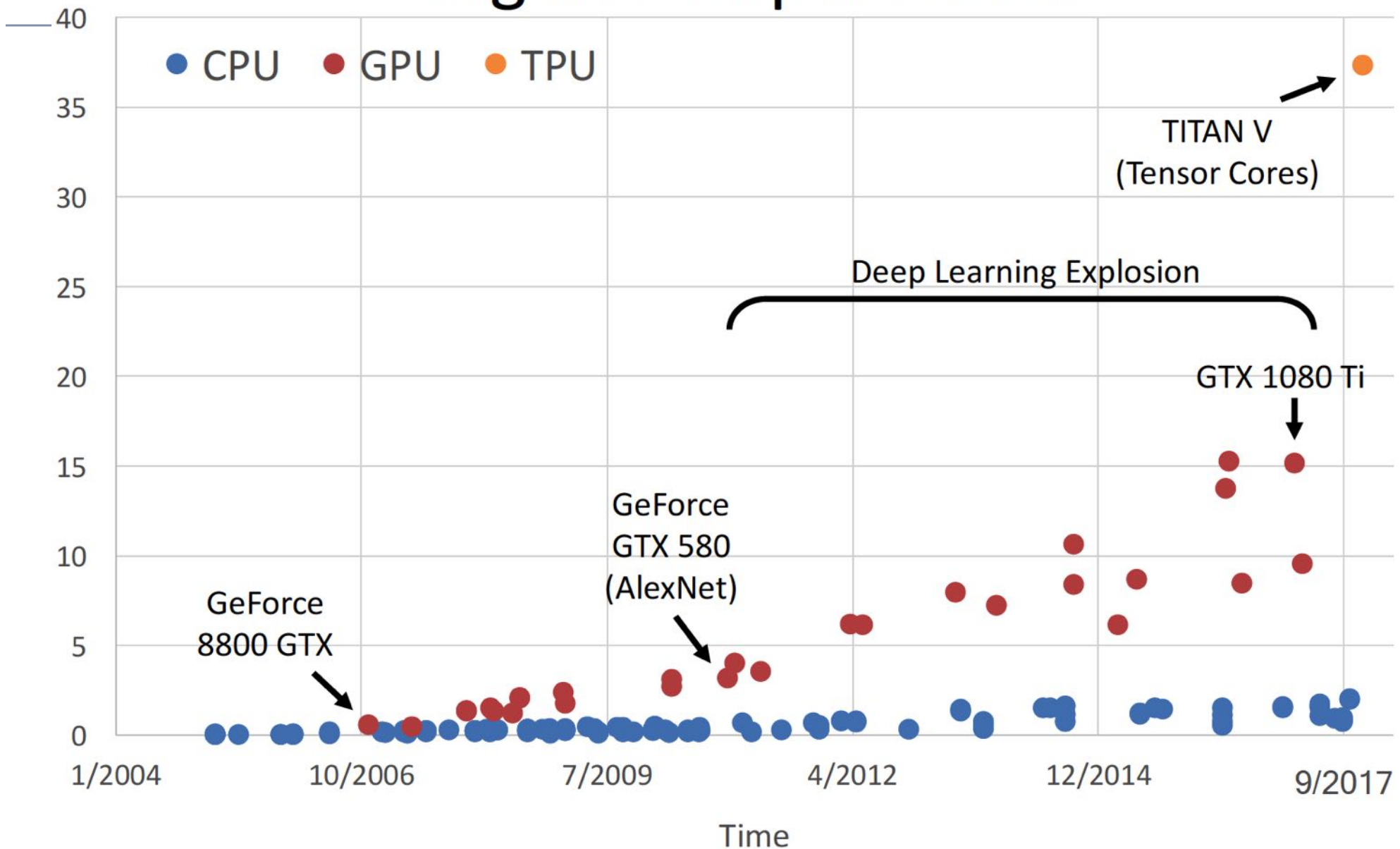
# Ingredients for Deep Learning



slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# GigaFLOPs per Dollar



slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# GigaFLOPs per Dollar



slide credit: Fei-Fei, Justin Johnson, Serena Yeung

**IMAGENET**

www.image-net.org

**22K** categories and **14M** images

- Animals
  - Bird
  - Fish
  - Mammal
  - Invertebrate
- Plants
  - Tree
  - Flower
  - Food
  - Materials
- Structures
- Artifact
  - Tools
  - Appliances
  - Structures
- Person
- Scenes
  - Indoor
  - Geological Formations
- Sport Activities

Deng, Dong, Socher, Li, Li, & Fei-Fei, 2009

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

**IMAGENET Large Scale Visual Recognition Challenge**

The Image Classification Challenge:
1,000 object classes
1,431,167 images

Output:
Scale
T-shirt
**Steel drum**
Drumstick
Mud turtle
✔

Output:
Scale
T-shirt
Giant panda
Drumstick
Mud turtle
✘

Russakovsky et al. IJCV 2015

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Validation classification

# Validation classification

**IMAGENET Large Scale Visual Recognition Challenge**

The Image Classification Challenge:
1,000 object classes
1,431,167 images

Russakovsky et al. IJCV 2015

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# IMAGENET Large Scale Visual Recognition Challenge

## Year 2010
### NEC-UIUC



Dense descriptor grid:
HOG, LBP

↓

Coding: local coordinate,
super-vector

↓

Pooling, SPM

↓

Linear SVM

[Lin CVPR 2011]

Lion image by Swissfrog is
licensed under CC BY 3.0

## Year 2012
### SuperVision



[Krizhevsky NIPS 2012]

Figure copyright Alex Krizhevsky, Ilya
Sutskever, and Geoffrey Hinton, 2012.
Reproduced with permission.

## Year 2014
### GoogLeNet

- Pooling
- Convolutio
- n
- Softmax
- Other



[Szegedy arxiv 2014]

### VGG

| Image |
| conv-64 |
| conv-64 |
| maxpool |
| conv-128 |
| conv-128 |
| maxpool |
| conv-256 |
| conv-256 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| conv-512 |
| conv-512 |
| maxpool |
| fc-4096 |
| fc-4096 |
| fc-1000 |
| softmax |

[Simonyan arxiv 2014]

## Year 2015
### MSRA



[He ICCV 2015]

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# How deep is enough?

AlexNet (2012)

5 convolutional layers

3 fully-connected layers

# How deep is enough?

AlexNet (2012)    VGG-M (2013)    VGG-VD-16 (2014)    GoogLeNet (2014)

# How deep is enough?

GoogLeNet (2014)

VGG-VD-16 (2014)

VGG-M (2013)

AlexNet (2012)

ResNet 50 (2015)

ResNet 152 (2015)

16 convolutional layers

50 convolutional layers

152 convolutional layers

Krizhevsky, I. Sutskever, and G. E. Hinton. *ImageNet classification with deep convolutional neural networks*. In Proc. NIPS, 2012.

C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. *Going deeper with convolutions*. In Proc. CVPR, 2015.

K. Simonyan and A. Zisserman. *Very deep convolutional networks for large-scale image recognition*. In Proc. ICLR, 2015.

K. He, X. Zhang, S. Ren, and J. Sun. *Deep residual learning for image recognition*. In Proc. CVPR, 2016.

# Convolutional Neural Networks (CNNs)
# were not invented overnight...

1998
LeCun et al.

Input

Image Maps

Convolutions

Subsampling

Fully Connected

Output

# of transistors
$10^6$
pentium II

# of pixels used in training
$10^7$ NIST

2012
Krizhevsky et al.

Figure copyright Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton, 2012. Reproduced with permission.

# of transistors
$10^9$
intel Xeon processor

GPUs
nvidia

# of pixels used in training
$10^{14}$ IMAGENET

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Try it out yourself

- Caffe ist an open implementation from the Berkeley Vision Group

  ▸ http://caffe.berkeleyvision.org

  ▸ http://demo.caffe.berkeleyvision.org

**Deep Learning
have become an important tool
for object recognition / image classification**

but there exist many other computer vision tasks
where Deep Learning is also an essential ingredient

a few examples...

# Human Pose Estimation

- **Single Person Pose Estimation** - two "phases"

  ▶ **Phase 1**: **pictorial structures models** e.g.
    [Felzenszwalb&Huttenlocher@ijcv05],
    [Andriluka&al@ijcv11], [Yang&Ramanan@pami13],
    [Pishchulin&al@iccv13], …

  ▶ **Phase 2**: using **deep learning** e.g.
    [Thoshev,Szegedy@cvpr14], [Thompson&al@nips14],
    [Chen&Yuille@nips14], [Carreira&al@cvpr16],
    [Hu&Ramanan@cvpr16], [Wei&al@cvpr16],
    [Newell&al@cvpr16], …

# MPII Human Pose Dataset: Dataset demo

- 410 human activities (after merging similar activities)

- over 40,000 annotated poses

- over 1.5M video frames



**http://human-pose.mpi-inf.mpg.de/**

# Analysis - overall performance

**Best Methods today: deep learning "takes" over**

**PCKh total, MPII Single Person**



Legend:
- Newell et al., arXiv'16
- Wei et al., CVPR'16
- Insafutdinov et al., arXiv'16
- Gkioxary et al., arXiv'16
- Lifshitz et al., arXiv'16
- Pishchulin et al., CVPR'16
- Hu&Ramanan, CVPR'16
- Tompson et al., CVPR'15
- Carreira et al., CVPR'16
- Tompson et al., NIPS'14
- Pishchulin et al., ICCV'13

x-axis: Normalized distance
y-axis: Detection rate, %

**Best Method as of ICCV'13**

✓ since CVPR'14, dataset has become **de-facto standard benchmark**

✓ **large training set** facilitated development of **deep learning methods**

# Cityscapes: Large-Scale Datasets for Semantic Labeling of Street Scenes



- Joint effort of:

DAIMLER    max planck institut informatik    TECHNISCHE UNIVERSITÄT DARMSTADT

Towards 3D Visual Scene

| Class | Group | Class | Group | Class | Group |
|---|---|---|---|---|---|
| road | ground | building | | person[1] | human |
| sidewalk | | wall | | rider[1] | |
| car[1] | | fence | | tree | nature |
| truck[1] | | traffic sign | infra- | terrain | |
| bus[1] | | traffic light | structure | ground | |
| on rails[1] | vehicle | pole | | dynamic | void |
| motorcycle[1] | | bridge[2] | | static | |
| bicycle[1] | | tunnel[2] | | | |
| license plate[2] | | sky | sky | | |

[1]Single instance annotation available
[2]Not included in fine label set challenges

# Image Description



A female tennis player in action on the court.

A group of young men playing a game of soccer.

A man riding a wave on top of a surfboard.

# Image Description



**Ours**: a person on skis jumping over a ramp

**Ours**: a skier is making a turn on a course

**Ours**: a cross country skier makes his way through the snow

**Ours**: a skier is headed down a steep slope

[Rakshith'17]

**Baseline**: a man riding skis down a snow covered slope

# Towards a Visual Turing Challenge

Q: What is the object on the counter in the corner?
A: micro wave

...bject in the scene?

A: brown

QA: (what is beneath the candle holder, decorative plate)
Some annotators use variations on spatial relations that are similar, e.g. 'beneath' is closely related to 'below'.

QA: (what is in front of the wall divider?, cabinet)
Annotators use additional properties to clarify object references (i.e. wall divider). Moreover, the perspective plays an important role in these spatial relations interpretations

Q How many lights are on?
A: 6

1449 RG... ...(NYU depth dataset)

available

QA1: (what is in front of the curtain behind the armchair?, guitar)
QA2: (what is in front of the curtain?, guitar)
Spatial relations matter more in complex environments where reference resolution becomes more relevant. In cluttered scenes, pragmatism starts playing a more important...

QA: (How many drawers are there?, 8)
The annotators use their common-sense knowledge for amodal completion. Here the annotator infers the 8th drawer from the context

QA: (What is the shape of the green chair?, horse shaped)
In this example, an annotator refers to a "horse shaped chair" which requires a quite abstract reasoning about the shapes.

The annotators are using different names to call the same things. The names of the brown object near the bed include 'night stand', 'stool', and 'cabinet'.

Some objects, like the table on the left of image, are severely occluded or truncated. Yet, the annotators refer to them in the questions.

QA: (What is in front of toilet?, door)
Here the 'open door' is clearly visible, yet...

QA1:(How many doors are in the image?, 1)
QA2:(How many doors are in the image?, 5)
Different interpretation of 'door' results in different counts: 1 door at the end of the hall vs. 5 doors including lockers.

QA: (What is the object on the counter in the corner?, micro wave)
References like 'corner' are difficult to...

QA: (How many doors are open?, 1)
Notion of 'open' or 'close' is poorly not well captured by current vision techniques.

QA: (Where is oven?, on the right side of refrigerator)
On some occasions, the annotators prefer to...

A    B    C

0.7   0.8   0.9   1

# Question Answering Results



**What is on the right side of the cabinet?**

Vision + Language:          **bed**

Language Only:               **bed**



**What objects are found on the bed?**

Vision + Language:          **bed sheets, pillow**

Language Only:               **doll, pillow**



**How many burner knobs are there?**

Vision + Language: **4**

Language Only:      **6**

# Video Object Segmentation

Goal: Separating a specific **foreground object** from **background** in a video given its **1st frame mask annotation**.



1st frame ———————————————→ *t*

DAVIS 2016
[Perazzi et al.'16]

# MaskTrack - Proposed Approach

➔ we process video per-frame, using guidance from previous frame



**MaskTrack**

DeepLab
[Chen et al., ICLR'15]

Frame *t*
output mask

Frame *t-1*
output mask

Frame *t*
input

➔ **we want to train from static images only**

# Qualitative Results



**https://www.mpi-inf.mpg.de/masktrack**

# Basic Concepts and Terminology

Computer Vision vs. Computer Graphics

# Pinhole Camera (Model)

- (simple) standard and abstract model today

  ▸ box with a small hole in it

# Camera Obscura

- around 1519, Leonardo da Vinci (1452 - 1519)

  ▸ http://www.acmi.net.au/AIC/CAMERA_OBSCURA.html

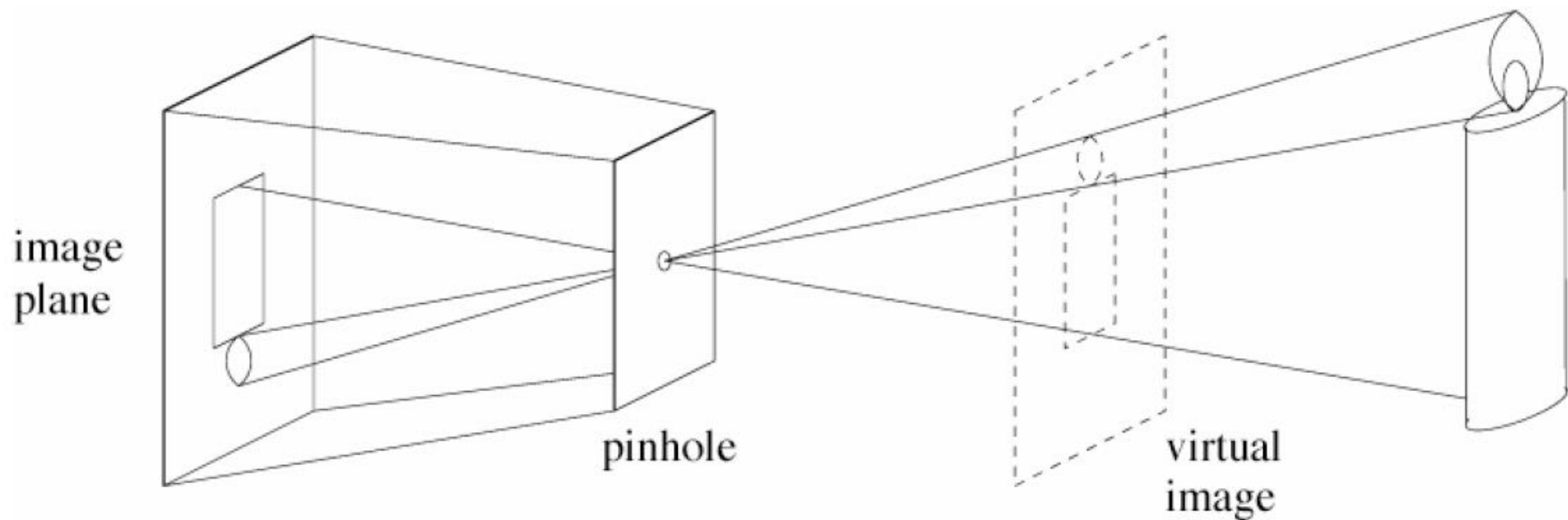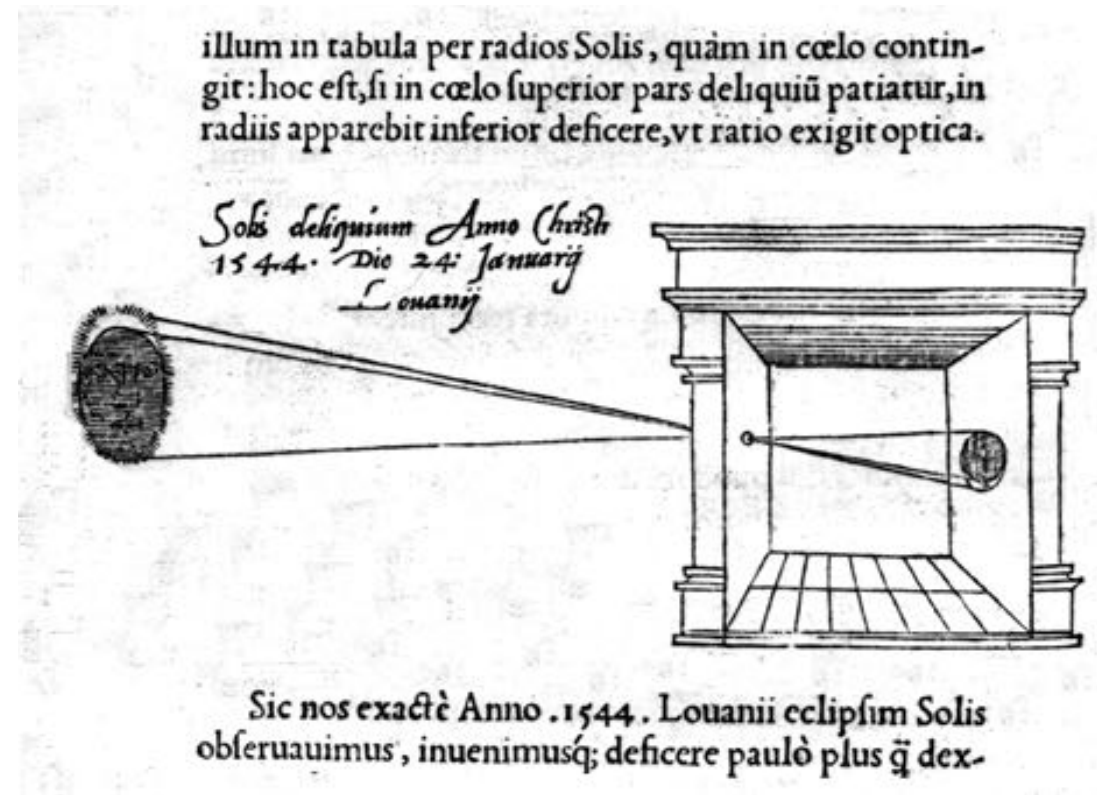  ▸ "when images of illuminated objects ... penetrate through a small hole into a very dark room ... you will see [on the opposite wall] these objects in their proper form and color, reduced in size ... in a reversed position owing to the intersection of the rays"



illum in tabula per radios Solis, quàm in cœlo contingit: hoc eft, fi in cœlo fuperior pars deliquiũ patiatur, in radiis apparebit inferior deficere, vt ratio exigit optica.

Solis deliquium Anno Chrifti 1544. Die 24. Januarij Louanij

Sic nos exactè Anno .1544. Louanii eclipfim Solis obferuauimus, inuenimusq; deficere paulò plus q̃ dex-

# Principle of pinhole....
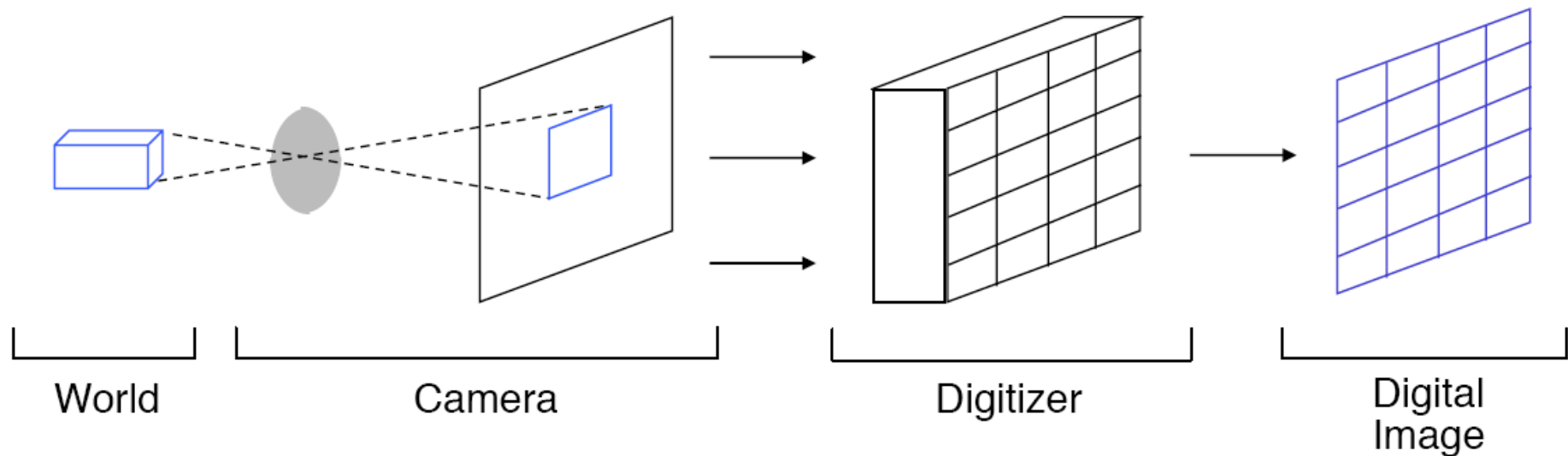
- ...used by artists
  - ▶ (e.g. Vermeer 17th century, dutch)
- and scientists

# Digital Images

- Imaging Process:
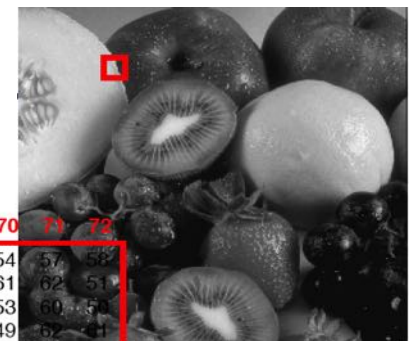  - ▶ (pinhole) camera model
  - ▶ digitizer to obtain digital image



World     Camera     Digitizer     Digital Image

# (Grayscale) Image

- 'Goals' of Computer Vision

  ‣ how can we recognize fruits from an array of (gray-scale) numbers?

  ‣ how can we perceive depth from an array of (gray-scale) numbers?

  ‣ …

- 'Goals' of Graphics

  ‣ how can we generate an array of (gray-scale) numbers that looks like fruits?

  ‣ how can we generate an array of (gray-scale) numbers so that the human observer perceives depth?

  ‣ …

- computer vision = the problem of 'inverse graphics' …?

# Visual Cues for Image Analysis

… in art and visual illusions

# 1. Case Study:
# Human & Art - Recovery of 3D Structure

# 1. Case Study:
# Human & Art - Recovery of 3D Structure

# 1. Case Study:
# Human & Art - Recovery of 3D Structure



Vincent van Gogh *Interior of a Restaurant at Arles* 1888

# 1. Case Study:
# Human & Art - Recovery of 3D Structure



Vincent van Gogh *Snowy Landscape with Arles in the Background* 1888

# 1. Case Study
# Computer Vision - Recovery of 3D Structure

- take all the cues of artists and
  'turn them around'

  ‣ exploit these cues to **infer**
    the structure of the world

  ‣ need **mathematical** and
    **computational models** of these cues

- sometimes called
  'inverse graphics'



http://www.vrvis.at/ar2/adm/shading/

# A 'trompe l'oeil'

- depth-perception
  - ▸ movement of ball stays the same
  - ▸ location/trace of shadow changes

# Another 'trompe l'oeil'

- illusory motion
  - ▶ only shadows changes
  - ▶ square is stationary

# Color & Shading

# Color & Shading

# 2. Case Study: Computer Vision & Object Recognition

- is it more than inverse graphics?

- how do you recognize
  - ▸ the banana?
  - ▸ the glas?
  - ▸ the towel?

- how can we make computers to do this?

- ill posed problem:
  - ▸ missing data
  - ▸ ambiguities
  - ▸ multiple possible explanations

# Image Edges:
# What are edges? Where do they come from?



- Edges are changes in pixel brightness

# Image Edges:
# What are edges? Where do they come from?



- Edges are changes in pixel brightness
    - **Foreground/Background Boundaries**
    - **Object-Object-Boundaries**
    - **Shadow Edges**
    - **Changes in Albedo or Texture**
    - **Changes in Surface Normals**

# Line Drawings:
# Good Starting Point for Recognition?

MASSACHUSETTS INSTITUTE OF TECHNOLOGY

PROJECT MAC

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

## THE SUMMER VISION PROJECT

Seymour Papert

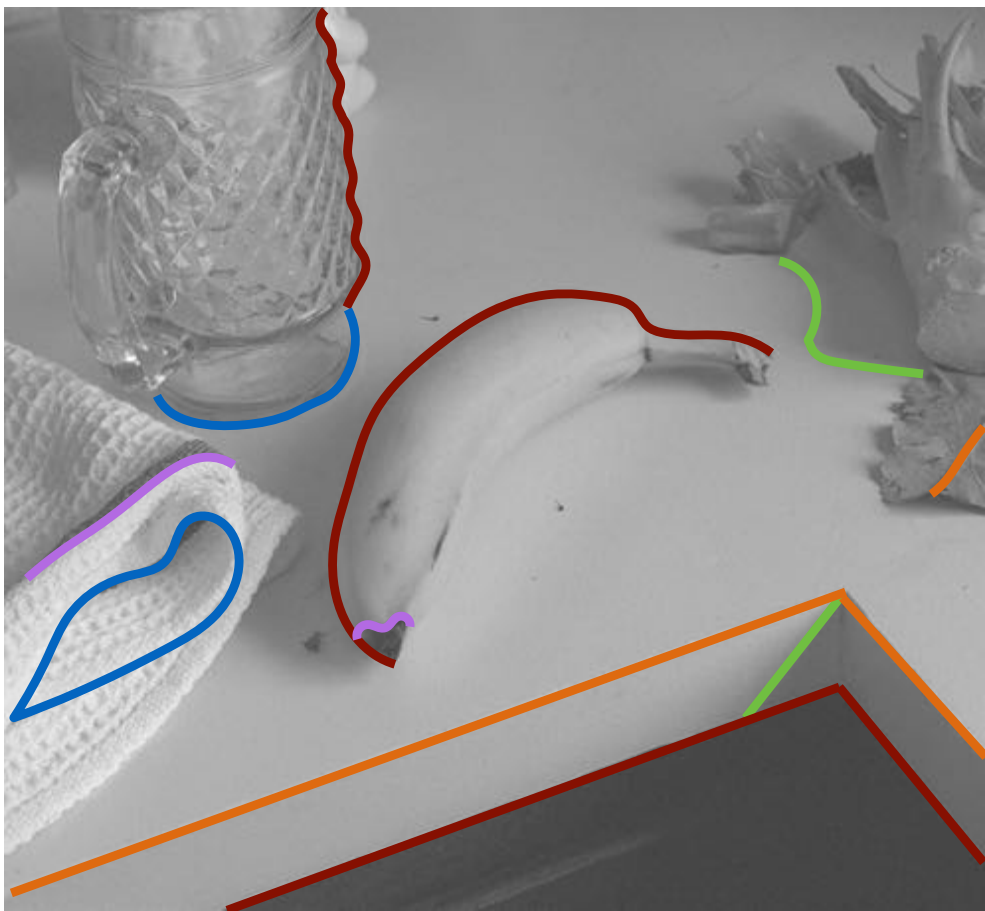The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the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

VISION

David Marr

FOREWORD BY
Shimon Ullman

AFTERWORD BY
Tomaso Poggio

David Marr, 1970s

Input image | Edge image | 2 ½-D sketch | 3-D model

| Input Image | Primal sketch | 2 ½-D | 3-D Model Representation |
|---|---|---|---|
| Perceived intensities | Zero crossings, blobs, edges, bands, lines, curves boundaries | Local surface orientation and discontinuities in surface orientation | 3-D models hierarchically organized in terms of surface and volumetric primitives |

Stages of Visual Representation, David Marr, 1970s

slide credit: Fei-Fei, Justin Johnson, Serena Yeung

# Complexity of Recognition

# Complexity of Recognition

# Complexity of Recognition

# Recognition: the Role of Context

- Antonio Torralba

# Recognition: the role of Prior Expectation

- Guiseppe Arcimboldo

# Complexity of Recognition

# One or Two Faces ?

# Class of Models: Pictorial Structure

- Fischler & Elschlager 1973

- Model has two components
  - ▸ parts
    (2D image fragments)
  - ▸ structure
    (configuration of parts)

# Deformations



A

B

C

D

# Clutter

# Example

# Recognition, Localization, and Segmentation

a few terms

... let's briefly define what we mean by that

# Object Recognition:
# First part of this Computer Vision class

- Different Types of Recognition Problems:

  ▶ Object **Identification**

    - recognize your pencil, your dog, your car

  ▶ Object **Classification**

    - recognize any pencil, any dog, any car

    - also called: generic object recognition, object categorization, …
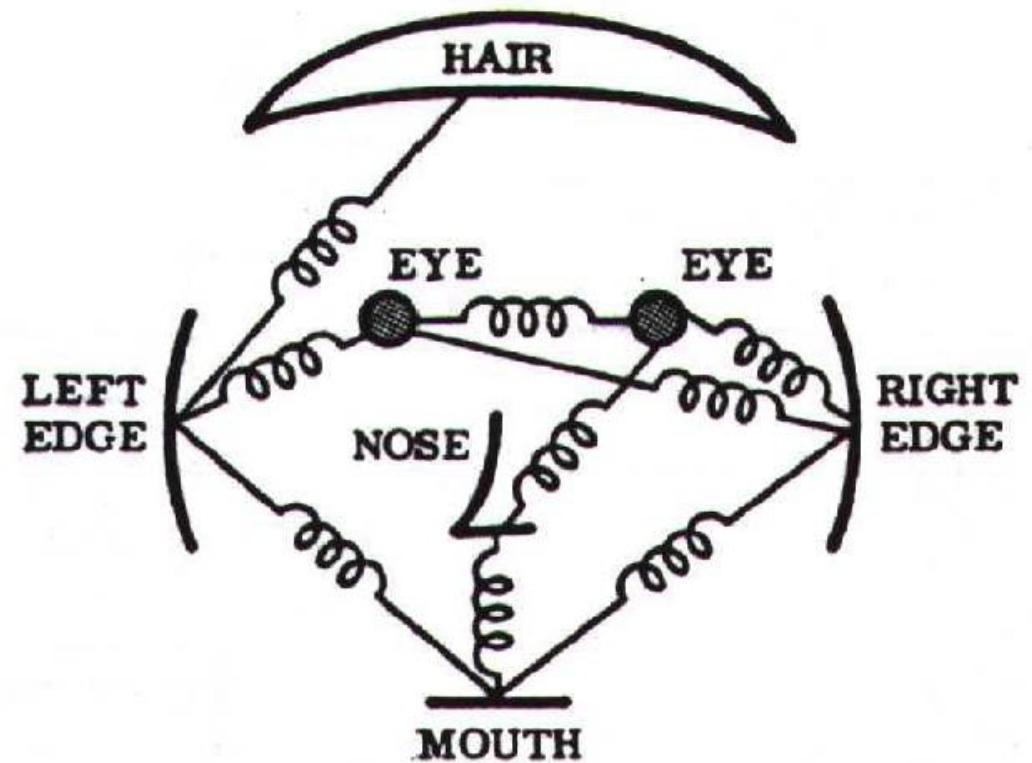
- Recognition and

  ▶ **Segmentation**: separate pixels belonging to the foreground (object) and the background

  ▶ **Localization/Detection**: position of the object in the scene, pose estimate (orientation, size/scale, 3D position)

# Object Recognition:
# First part of this Computer Vision class

- Different Types of Recognition Problems:

  - ▸ Object **Identification**

    - recognize your apple, your cup, your dog

  - ▸ Object **Classification**

    - recognize any apple, any cup, any dog

    - also called:
      generic object recognition, object categorization, …

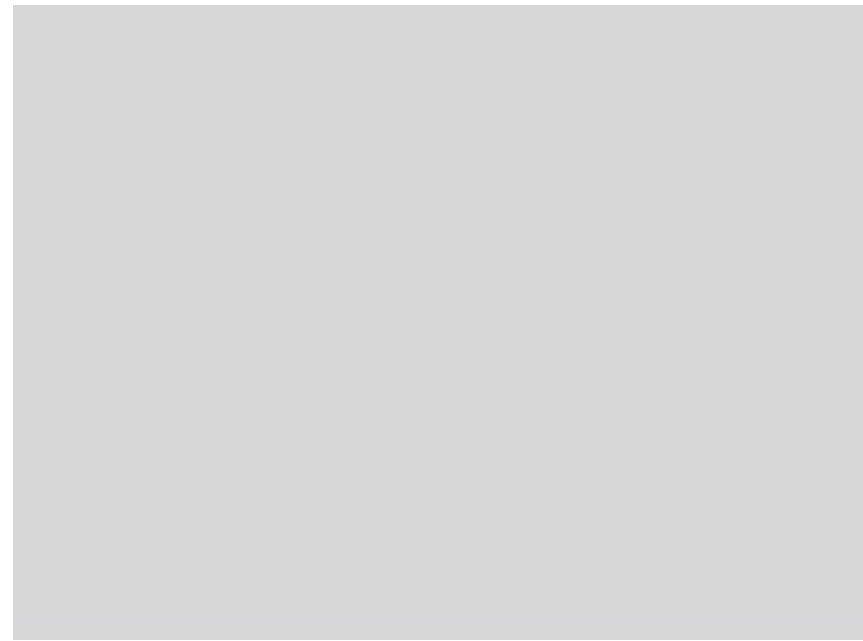    - typical definition:
      'basic level category'

# Which Level is right for Object Classes?

- Basic-Level Categories
  - the highest level at which category members have **similar perceived shape**
  - the highest level at which a **single mental image** can reflect the entire category
  - the highest level at which a person uses similar **motor actions** to interact with category members
  - the level at which human subjects are usually **fastest** at identifying category members
  - the first level named and understood by **children**

  - (while the definition of basic-level categories depends on culture there exist a remarkable consistency across cultures...)

- Most recent work in object recognition has focused on this problem
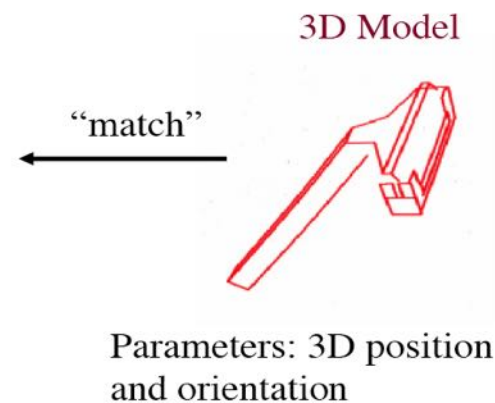  - we will discuss several of the most successful methods in the lecture :-)
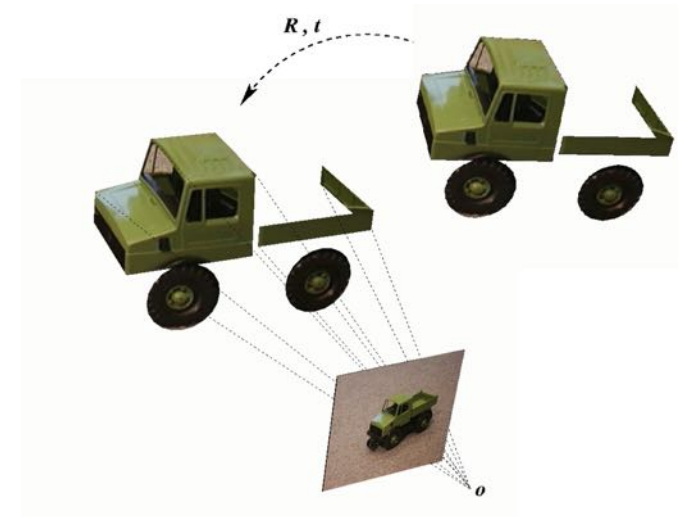
# Object Recognition & Segmentation

- Recognition and

  ▸ **Segmentation**: separate pixels belonging to the foreground (object) and the background

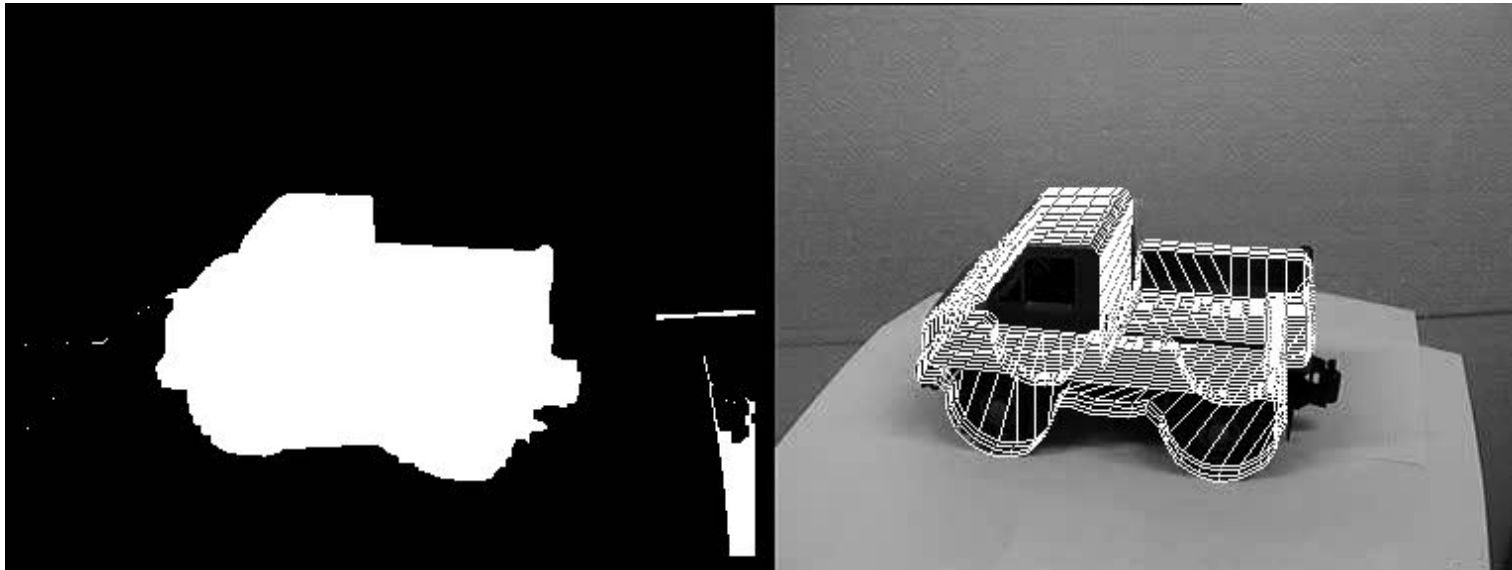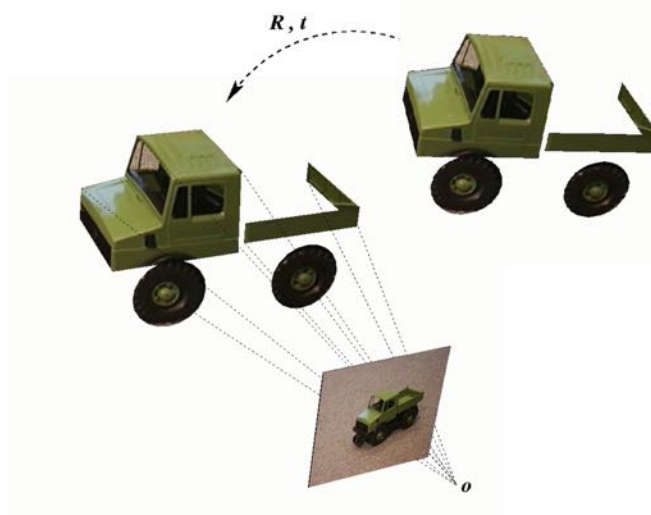# Object Recognition & Localization

- Recognition and

  ▸ **Localization**: to position the object in the scene, estimate the object's pose (orientation, size/scale, 3D position)

  

  ▸ Example from David Lowe:

# Localization: Example Video 1

# Localization: Example Video 2

# Object Recognition

- Different Types of Recognition Problems:

  ▶ Object **Identification**

    - recognize your pencil, your dog, your car

  ▶ Object **Classification**

    - recognize any pencil, any dog, any car

    - also called: generic object recognition, object categorization, …

- Recognition and

  ▶ **Segmentation**: separate pixels belonging to the foreground (object) and the background

  ▶ **Localization**: position the object in the scene, estimate pose of the object (orientation, size/scale, 3D position)

# Goals of today's lecture

- First intuitions about
  - ▸ What is computer vision?
  - ▸ What does it mean to see and how do we (as humans) do it?
  - ▸ How can we make this computational?

- Applications & Appetizers

- Role of Deep Learning
  - – with several slides taken from Fei-Fei Li, Justin Johnson, Serena Yeung @ Stanford

- 2 case studies:
  - ▸ Recovery of 3D structure
    - – slides taken from Michael Black @ Brown University / MPI Intelligent Systems
  - ▸ Object Recognition
    - – intuition from human vision...