



计算机视觉表征与识别

Chapter 1: Introduction to computer vision

王利民

媒体计算课题组

<http://mcg.nju.edu.cn/>



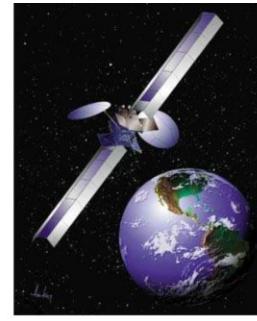
Overview



- What is computer vision about?
- Computer vision is useful.
- Computer vision is difficult.
- History and progress of computer vision
- Course overview



Introduction



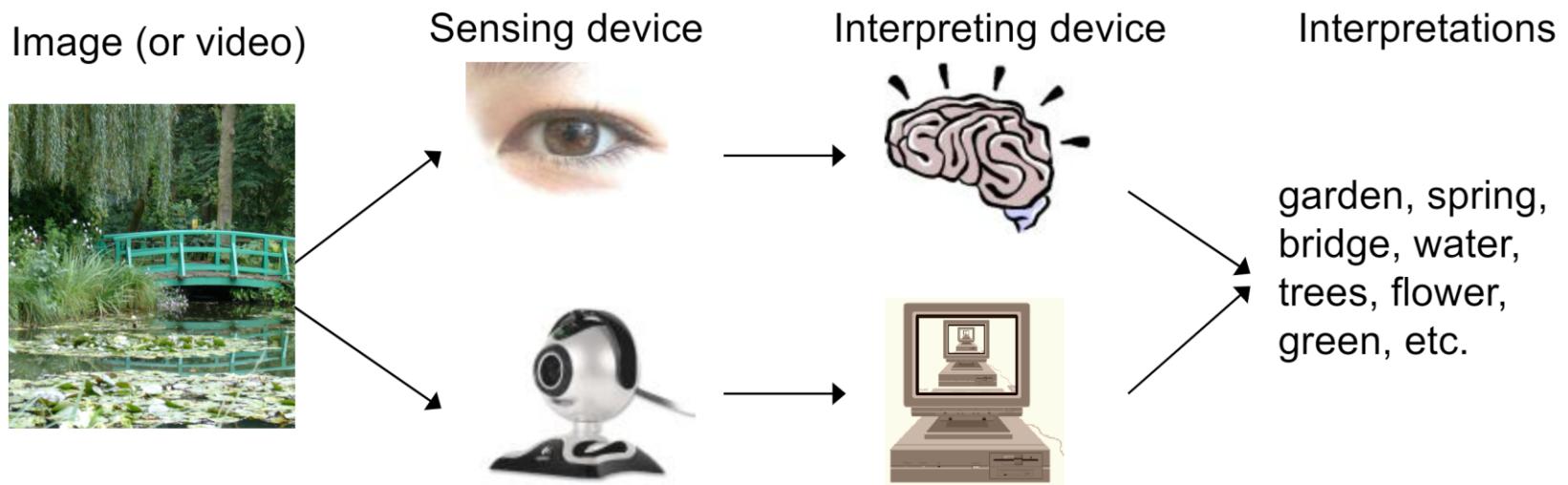


What is computer vision





What is computer vision





The goal of computer vision



- To extract “meaning” from pixels



What we see

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

What a computer sees

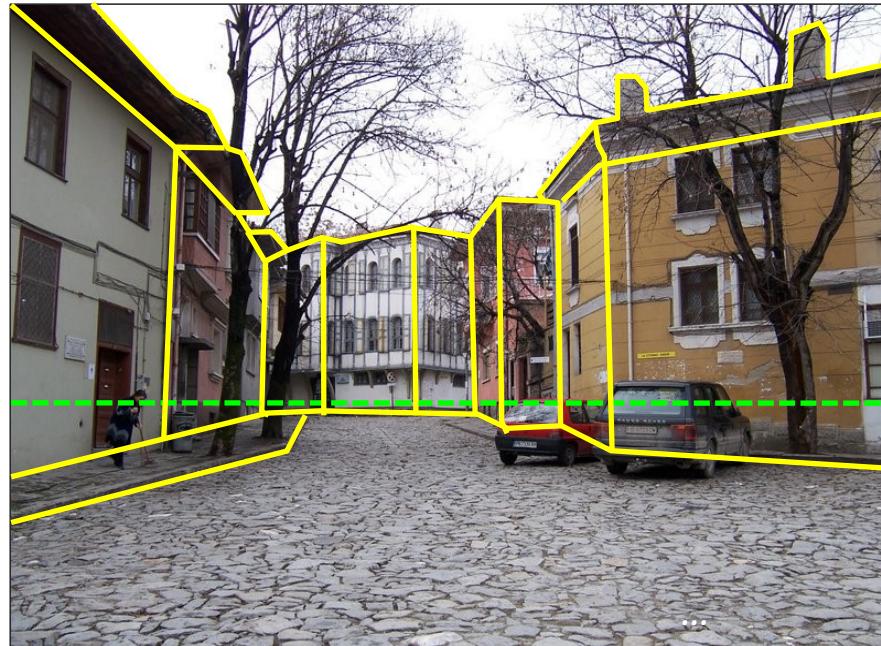


What kind of information





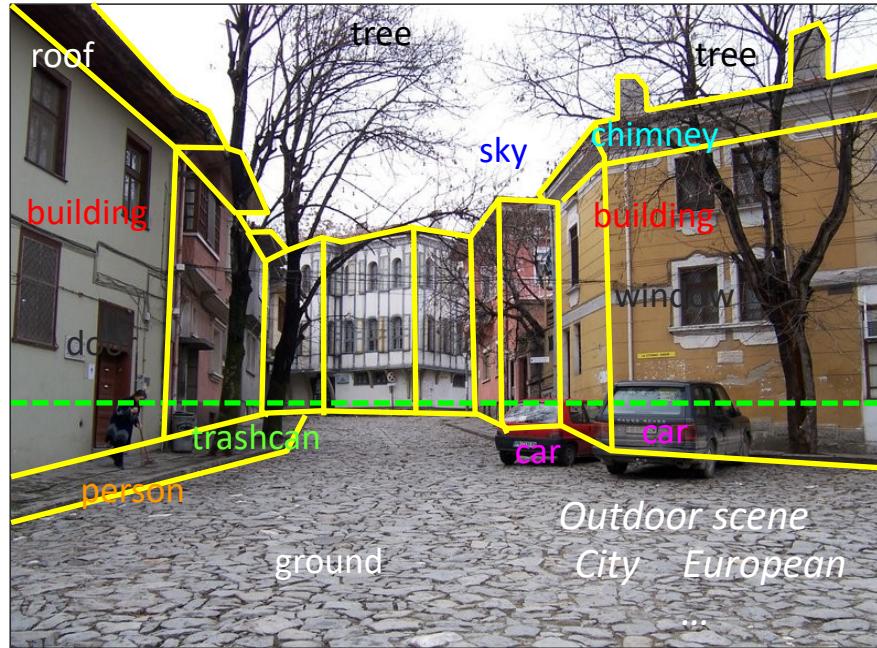
What kind of information



Geometric information



What kind of information



Geometric information
Semantic information



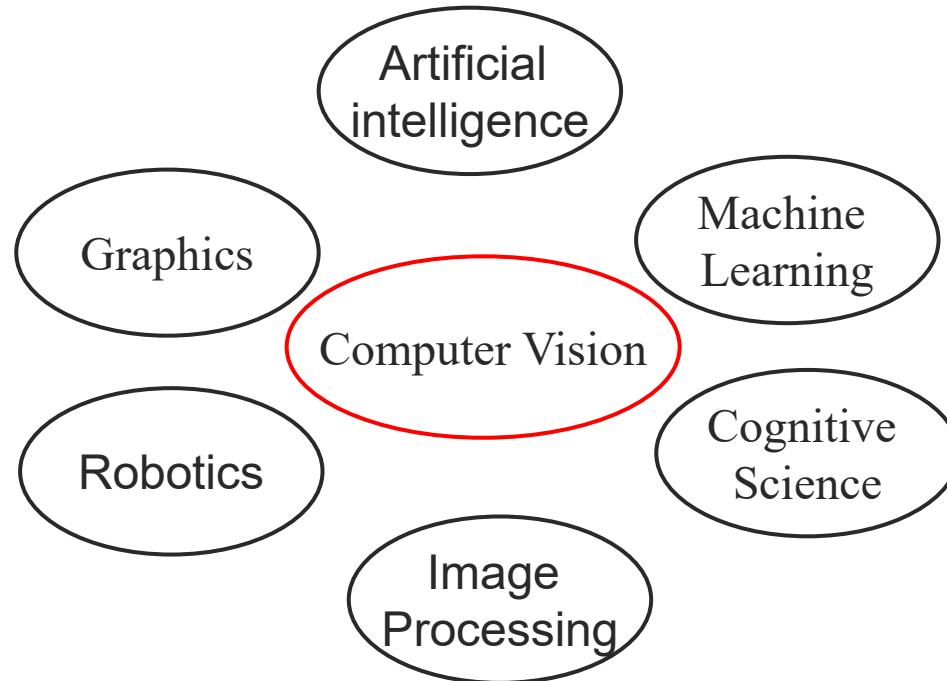
Computer vision



- Automatic understanding of images and video
 - 1. Computing properties of the 3D world from visual data (**measurement**)
 - 2. Algorithms and representations to allow a machine to recognize objects, scene, and people (**perception and interpretation**)

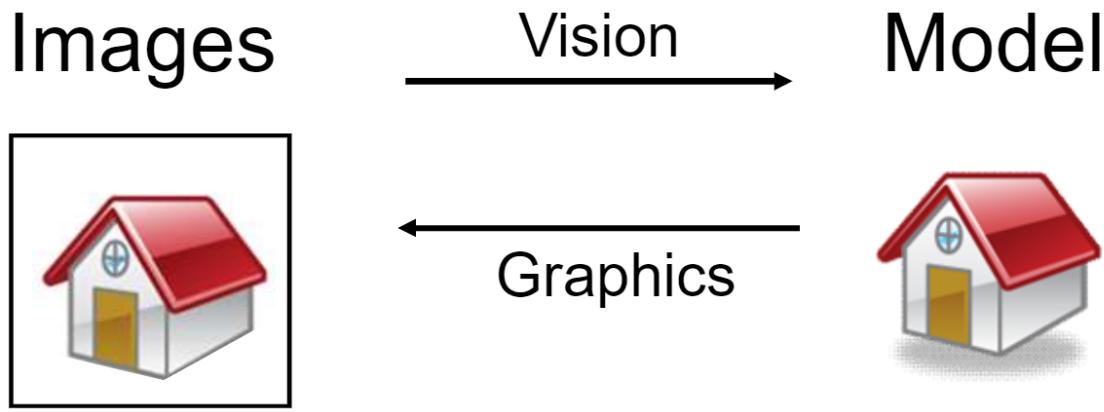


Related disciplines





Vision and graphics



Inverse problems: analysis and synthesis.



Why study computer vision



- Computer vision is useful
- Computer vision is difficult
- Computer vision is fast developing



Vision is useful



- For people vision is their most crucial sense, for good reason:
 - half our brain is devoted to it
 - developed many times during evolution
 - it can be implemented with high resolution
 - yields color, texture, depth, motion, shape



Faces and digital cameras



Camera waits for everyone to smile to take a photo [Canon]



Setting camera focus via face detection

Revisions





Pig face recognition



A system made by Yingzi Techology, a small Chinese company, scanning a barn to recognize pig faces.
Yingzi Technology



Video-based interfaces



Human joystick, NewsBreaker Live



Assistive technology systems
Camera Mouse, Boston College



Microsoft Kinect



THE
SOCIAL
MEDIA SHOW

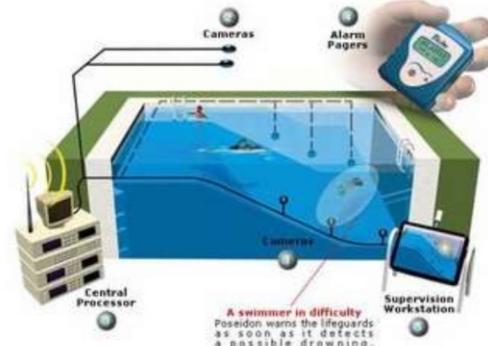




Safety & security



Navigation,
driver safety



Monitoring pool
(Poseidon)



Pedestrian detection
MERL, Viola et al.



Surveillance



The system converts image data taken by 4 super-wide angle cameras, to display a virtual image of the vehicle from above.

AI LEARNING

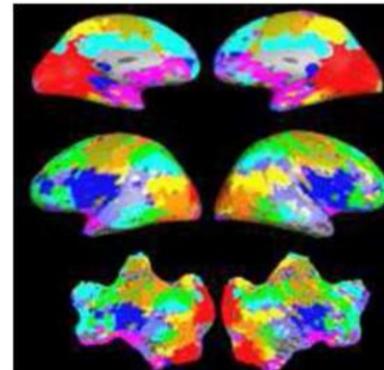




Vision for medical & neuroimages



Image guided surgery
MIT AI Vision Group

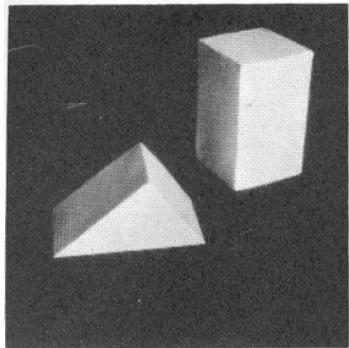


fMRI data
Golland et al.

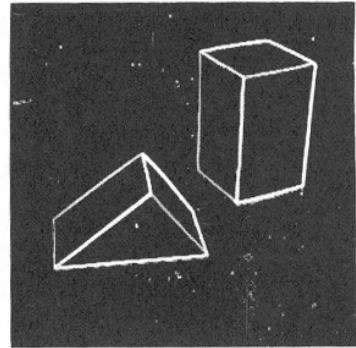




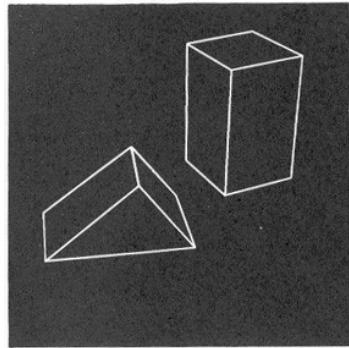
Origins of computer vision



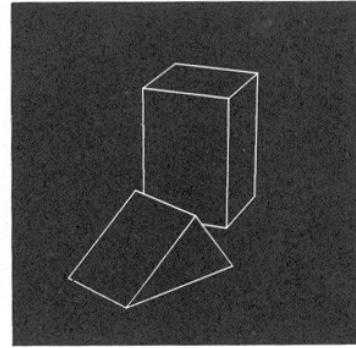
(a) Original picture.



(b) Differentiated picture.



(c) Line drawing.



(d) Rotated view.

L. G. Roberts *Machine Perception of Three Dimensional Solids*
Ph.D. thesis, MIT Department of Electrical Engineering, 1963.

- 25 - 4445(a-d)



Origins of computer vision



MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group
Vision Memo. No. 100.

July 7, 1966

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".



Computer vision is difficult



0	120	120	120	120	120	120	120	121	121	121	121	122	122	121	121	121	121	122	123	1
1	82	87	88	88	88	88	89	89	88	89	89	89	89	89	89	90	90	90	90	90
2	95	94	94	95	95	95	95	94	95	95	95	95	95	95	96	96	96	96	96	96
3	89	92	92	92	92	92	93	93	93	93	93	93	93	93	93	93	93	93	93	93
4	98	98	98	98	99	99	99	99	99	100	100	99	100	101	101	101	101	100	101	101
5	98	98	98	98	99	99	99	99	99	100	100	99	100	101	101	101	101	100	101	101
6	90	93	94	94	94	94	93	93	93	94	94	94	94	94	94	94	94	95	95	96
7	99	99	100	100	100	100	99	100	101	101	101	102	102	102	102	102	102	102	102	102
8	91	95	94	96	96	94	95	95	95	95	95	94	95	95	95	95	95	96	96	97
9	101	101	102	102	102	102	102	102	102	102	103	103	103	103	103	103	103	103	103	103
10	92	97	96	97	97	97	97	97	97	97	97	97	97	97	97	98	98	98	98	98
11	103	102	103	103	103	103	103	103	104	104	104	104	104	104	104	104	104	104	105	1
12	97	98	98	98	98	98	98	98	98	98	99	99	98	98	98	99	99	99	99	99
13	103	103	104	104	104	104	104	104	105	106	106	106	106	106	106	106	105	105	105	1
14	95	99	98	99	99	99	99	99	100	100	100	100	100	101	101	101	100	100	100	1
15	104	104	105	105	105	105	106	107	106	107	107	107	107	107	107	107	108	108	108	1
16	96	100	100	100	102	101	102	101	102	102	102	102	102	102	102	102	103	102	103	1
17	107	107	107	107	107	107	107	108	108	108	108	109	108	109	108	109	109	109	109	1
18	98	102	102	102	103	103	103	103	103	103	103	103	103	103	103	103	103	104	104	1
19	108	108	108	109	108	109	109	109	109	109	109	110	109	111	111	111	111	111	111	1
20	100	103	103	103	103	104	104	104	104	104	104	104	104	104	104	105	105	105	105	1
21	109	109	110	109	110	110	110	110	111	111	111	112	112	112	112	112	113	113	113	1
22	101	104	104	105	106	106	105	105	106	105	106	107	107	107	107	107	107	107	107	1
23	111	111	112	112	112	112	112	112	113	112	112	112	113	113	113	114	114	114	114	1
24	102	106	107	107	106	107	107	106	107	108	107	108	108	108	108	108	108	108	108	1
25	113	113	113	113	113	113	113	113	114	114	114	113	114	114	114	114	115	115	115	1
26	105	108	108	108	108	108	108	109	109	109	109	109	109	110	110	109	109	109	109	1
27	114	114	114	114	114	114	114	114	115	116	115	115	115	116	117	116	117	117	117	1
28	106	109	109	109	110	110	110	110	111	111	110	111	111	112	112	111	111	112	112	1
29	116	116	116	116	116	116	116	117	117	117	117	117	118	118	118	118	118	118	118	1
30	108	111	111	110	110	111	111	112	112	113	113	113	113	113	113	113	113	113	113	1
31	117	117	118	118	118	118	118	118	119	118	120	120	119	119	119	119	119	119	119	1
32	110	113	113	113	113	113	113	114	114	114	114	114	114	114	114	114	115	115	115	1
33	119	119	119	119	119	119	120	120	120	120	120	120	120	121	121	121	121	121	122	1
34	114	114	114	114	115	115	115	115	115	116	116	117	116	117	117	117	117	117	117	1
35	121	120	120	121	121	121	121	121	121	121	121	121	122	122	123	123	123	123	123	1
36	116	116	116	117	117	117	117	117	118	118	118	118	118	118	118	118	118	118	118	1
37	122	122	123	123	123	123	123	123	124	124	124	124	124	124	125	125	125	124	124	1
38	114	117	118	118	119	119	119	120	120	120	120	120	119	119	120	120	120	120	120	1
39	124	124	124	124	124	125	125	125	126	125	125	126	127	125	126	127	127	127	127	1
40	116	119	119	120	120	120	121	119	137	123	122	121	122	121	121	122	122	122	122	1
41	125	125	126	126	126	127	127	127	127	127	127	127	126	126	128	128	129	129	129	1
42	117	120	121	122	122	123	123	121	145	128	122	123	124	124	124	124	124	124	124	1

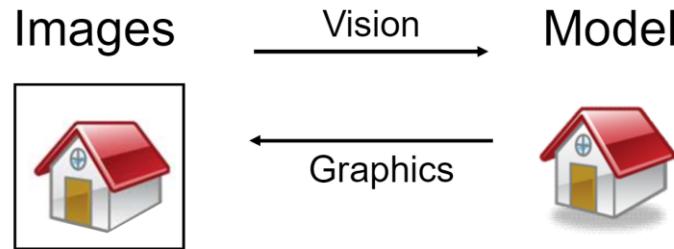
Gap between low level signal and high level meanings



Challenges: ill posed problem



- Ill-posed problem: real world much more complex than what we can measure in images.
- Impossible to literally “invert” image formation process.



Inverse problems: analysis and synthesis.



Challenges: large variations



Illumination



Object pose



Clutter



Occlusions



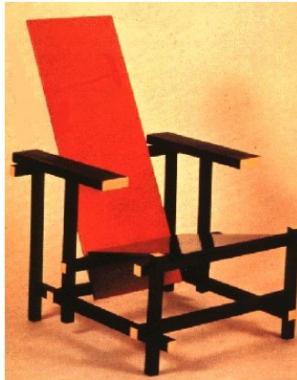
**Intra-class
appearance**



Viewpoint



Challenge: intra-class variation



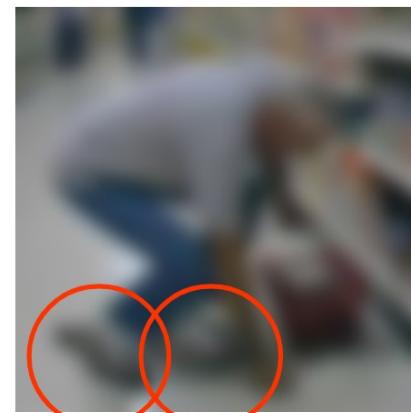
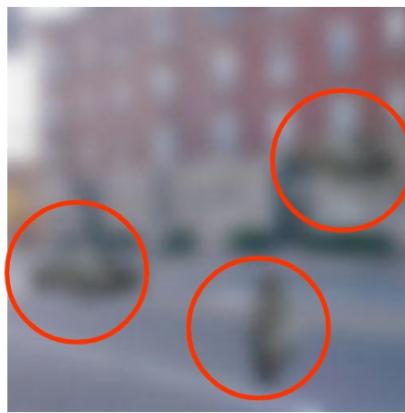
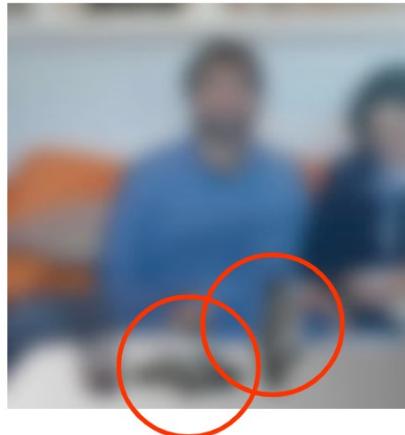
slide credit: Fei-Fei, Fergus & Torralba



Challenge: context



All encircled
patterns
are identical:





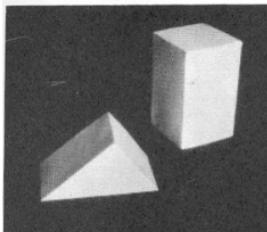
Challenges: complexity



- High dimension of input data
 - Millions of pixels in an image
- Large number of categories
 - 30,000 recognizable object categories
- Large number of visual data
 - Billions of images online
 - 144K hours of new video on YouTube daily
- About half of the cerebral cortex in primates is devoted to processing visual information



Progress charted by dataset



Roberts 1963



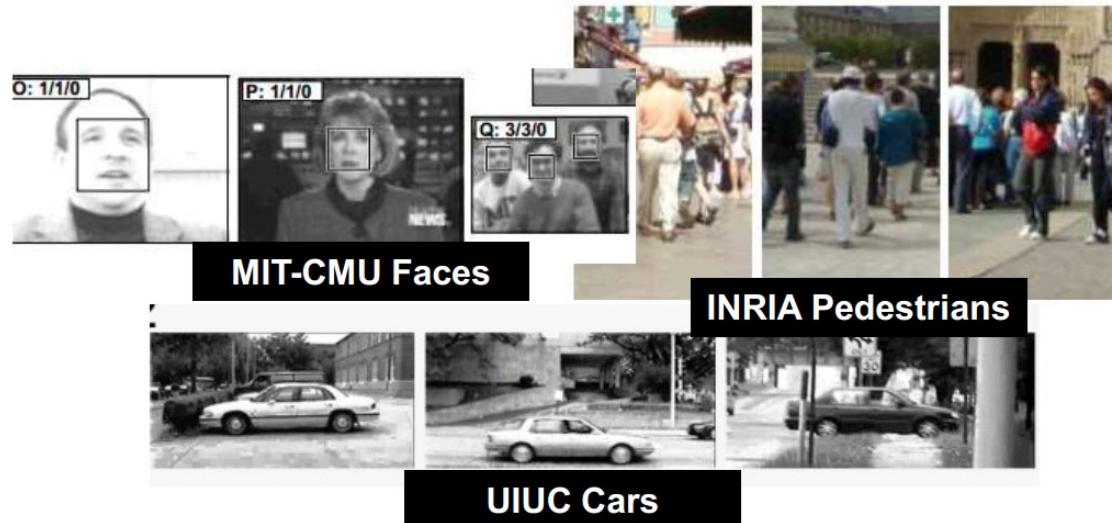
COIL



1963 ... 1996



Progress charted by dataset





Progress charted by dataset



MSRC 21 Objects



Caltech-101



Caltech-256



1963 ... 1996

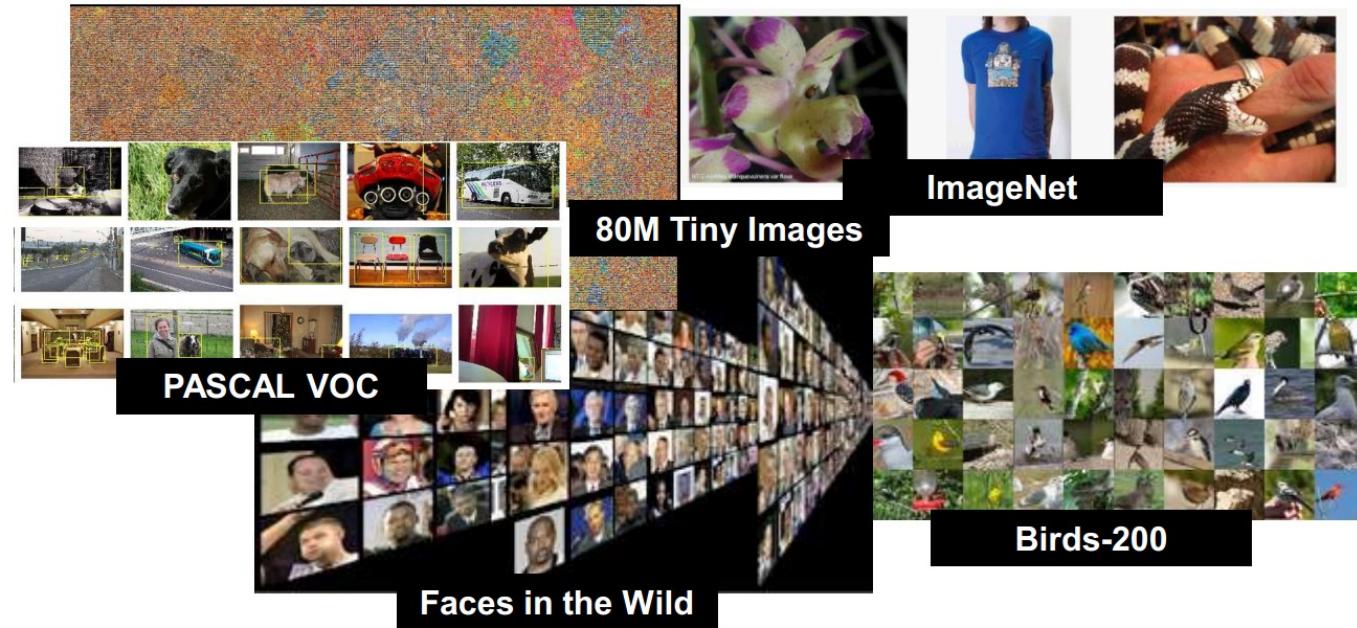
2000

2005





Progress charted by dataset

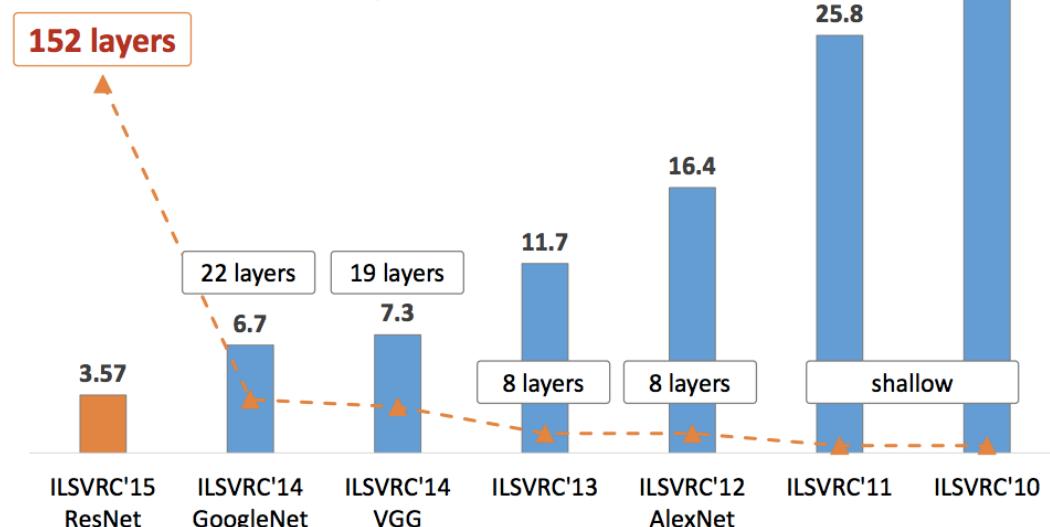
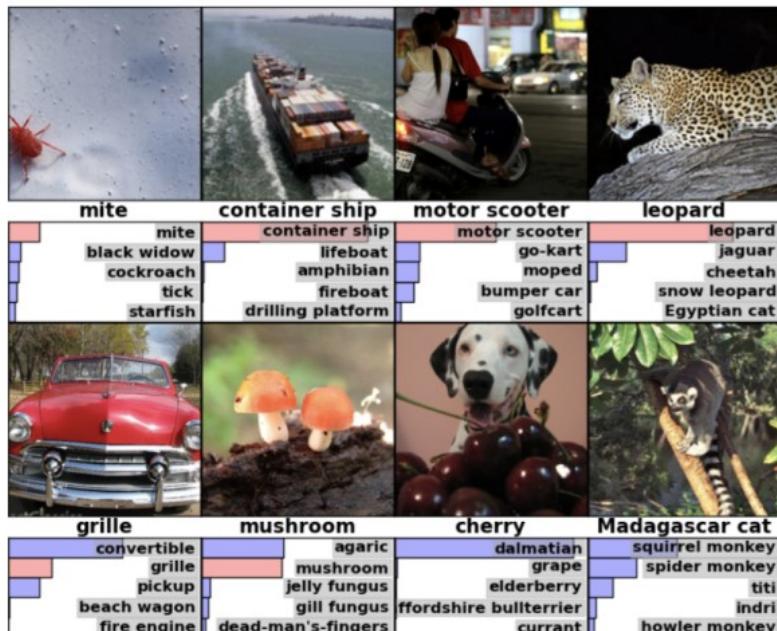




Recognition: General categories

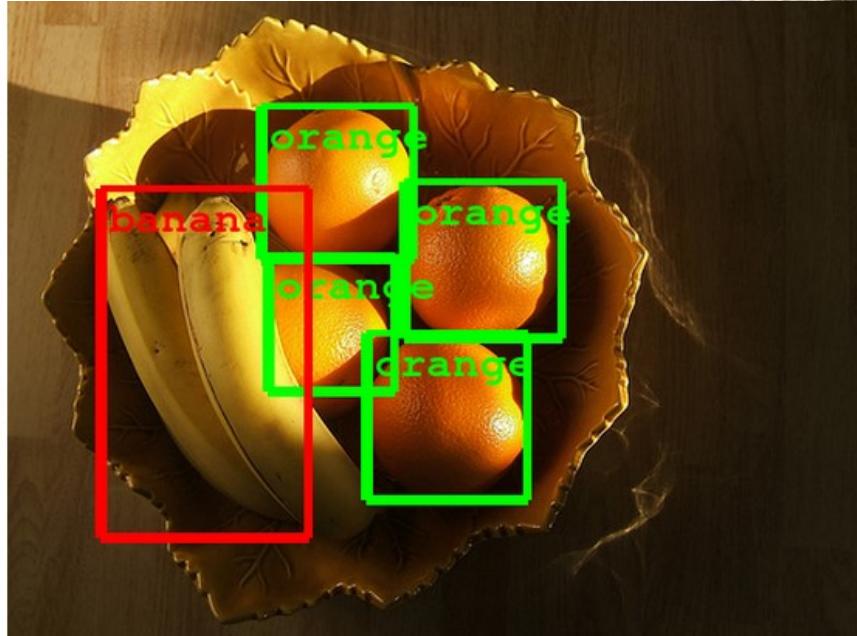
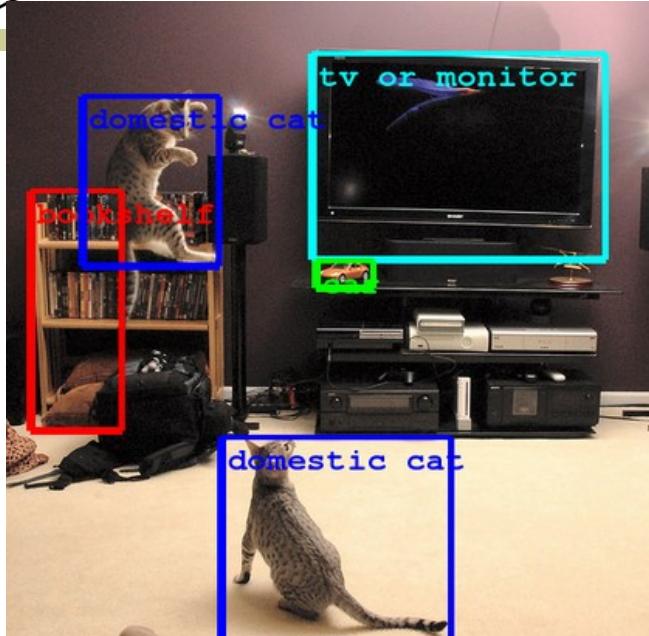


- ImageNet challenge





Recognition: General categories



- [Computer Eyesight Gets a Lot More Accurate](#),
NY Times Bits blog, August 18, 2014
- [Building A Deeper Understanding of Images](#),
Google Research Blog, September 5, 2014





Large scale recognition



clarifai

ABOUT

TECHNOLOGY

API

NEWS

BLOG

CAREERS

CONTACT

Paste a url here...

USE THE URL

CHOOSE A FILE INSTEAD

*By using the demo you agree to our terms of service



Predicted Tags

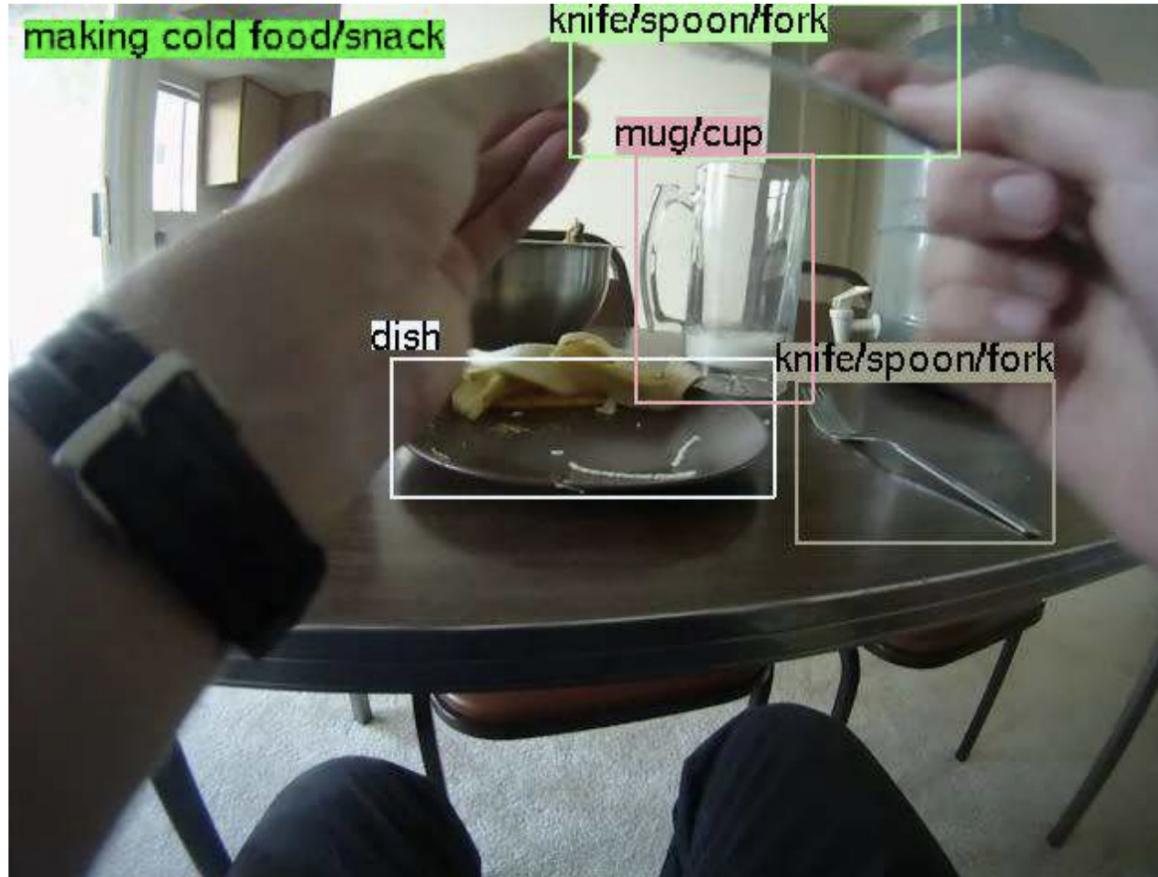
mammal livestock cattle
pasture agriculture bovine
farm nobody meadow grass

Similar Images



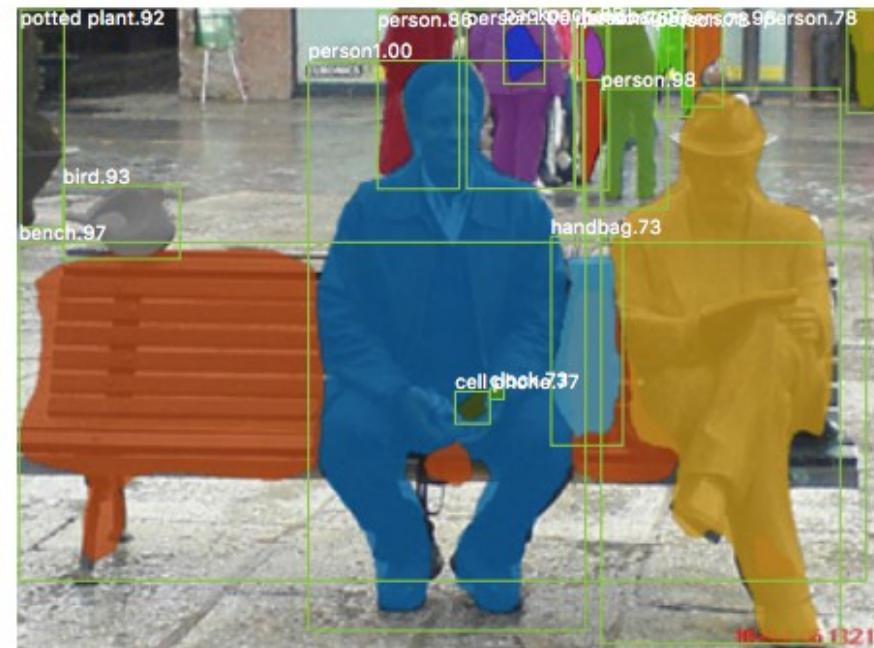


Recognition in first-person view





Object detection, instance segmentation



K. He, G. Gkioxari, P. Dollar, and R. Girshick, [Mask R-CNN](#),
ICCV 2017 (Best Paper Award)



Image captioning



Describes without errors	Describes with minor errors	Somewhat related to the image	Unrelated to the image
			
<p>A person riding a motorcycle on a dirt road.</p>	<p>Two dogs play in the grass.</p>	<p>A skateboarder does a trick on a ramp.</p>	<p>A dog is jumping to catch a frisbee.</p>
			
<p>A group of young people playing a game of frisbee.</p>	<p>Two hockey players are fighting over the puck.</p>	<p>A little girl in a pink hat is blowing bubbles.</p>	<p>A refrigerator filled with lots of food and drinks.</p>
			
<p>A herd of elephants walking across a dry grass field.</p>	<p>A close up of a cat laying on a couch.</p>	<p>A red motorcycle parked on the side of the road.</p>	<p>A yellow school bus parked in a parking lot.</p>



Image generation



- Faces: 1024x1024 resolution, CelebA-HQ



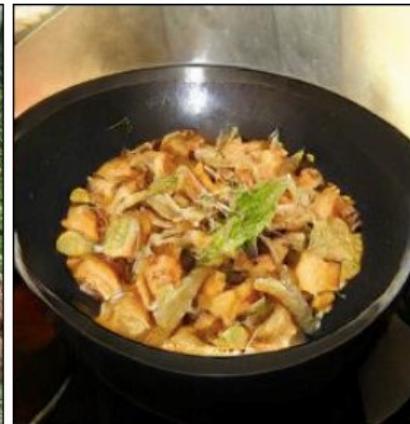
T. Karras, T. Aila, S. Laine, and J. Lehtinen, [Progressive Growing of GANs for Improved Quality, Stability, and Variation](#), ICLR 2018



Image generation



- BigGAN: 512 x 512 resolution, ImageNet



A. Brock, J. Donahue, K. Simonyan, [Large scale GAN training for high fidelity natural image synthesis](#), arXiv 2018



Image generation



- BigGAN: 512 x 512 resolution, ImageNet



A. Brock, J. Donahue, K. Simonyan, [Large scale GAN training for high fidelity natural image synthesis](#), arXiv 2018

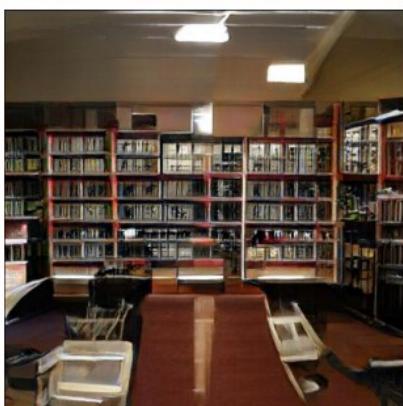


Image generation



- BigGAN: 512 x 512 resolution, ImageNet

Easy classes



Difficult classes



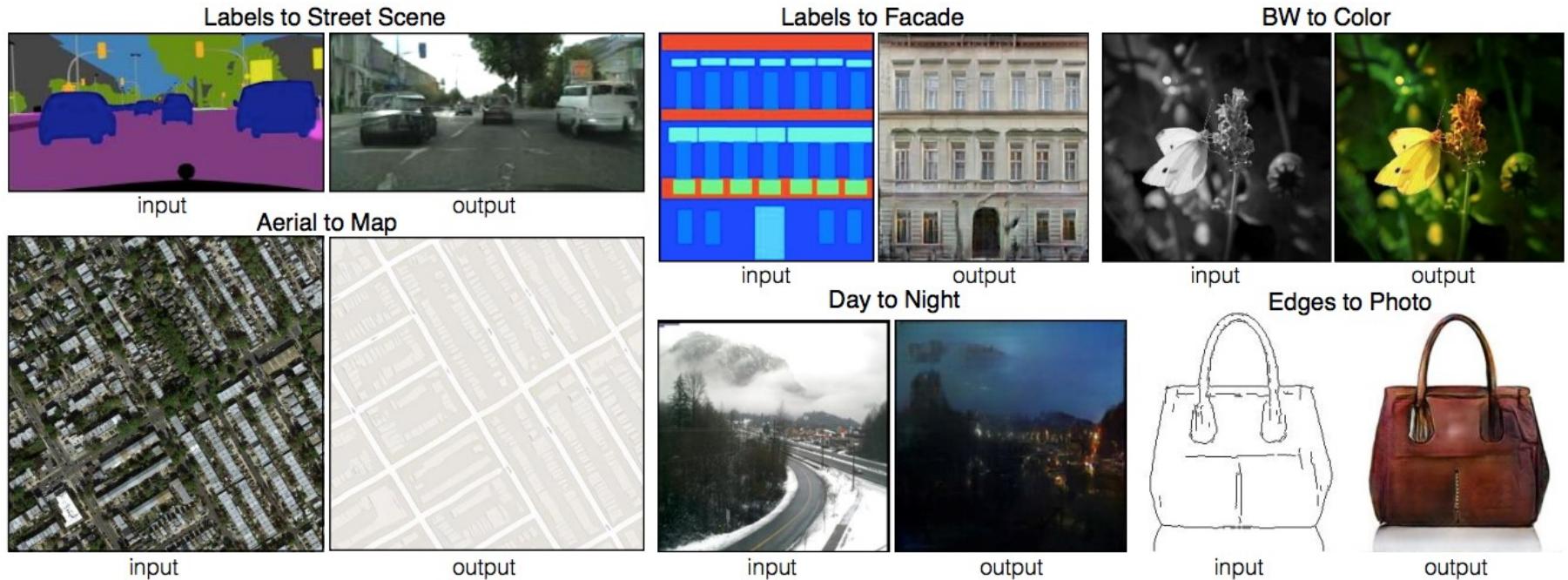
A. Brock, J. Donahue, K. Simonyan, [Large scale GAN training for high fidelity natural image synthesis](#), arXiv 2018



Image generation



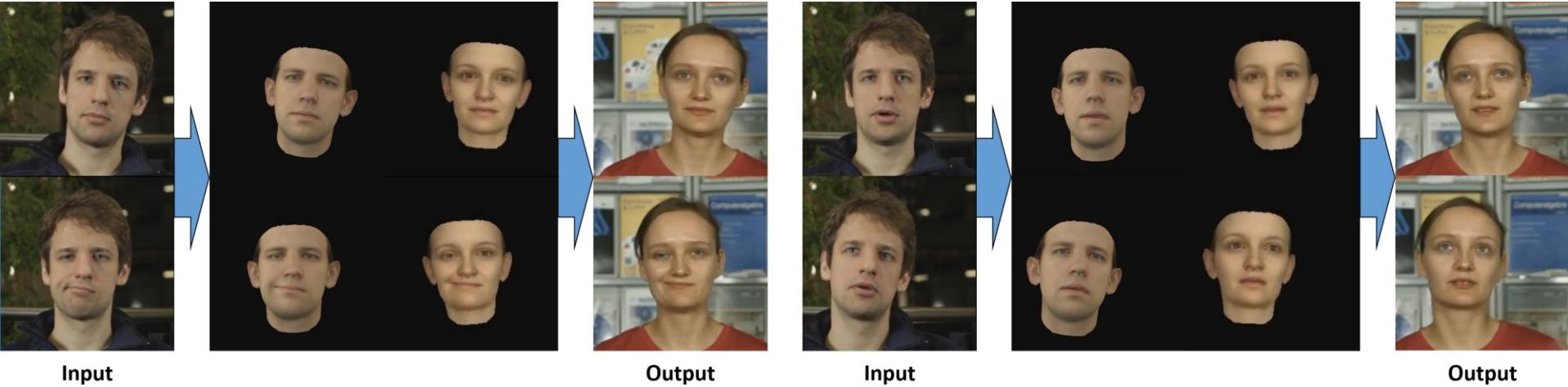
- Image-to-image translation



P. Isola, J.-Y. Zhu, T. Zhou, A. Efros, [Image-to-Image Translation with Conditional Adversarial Networks](#),
CVPR 2017



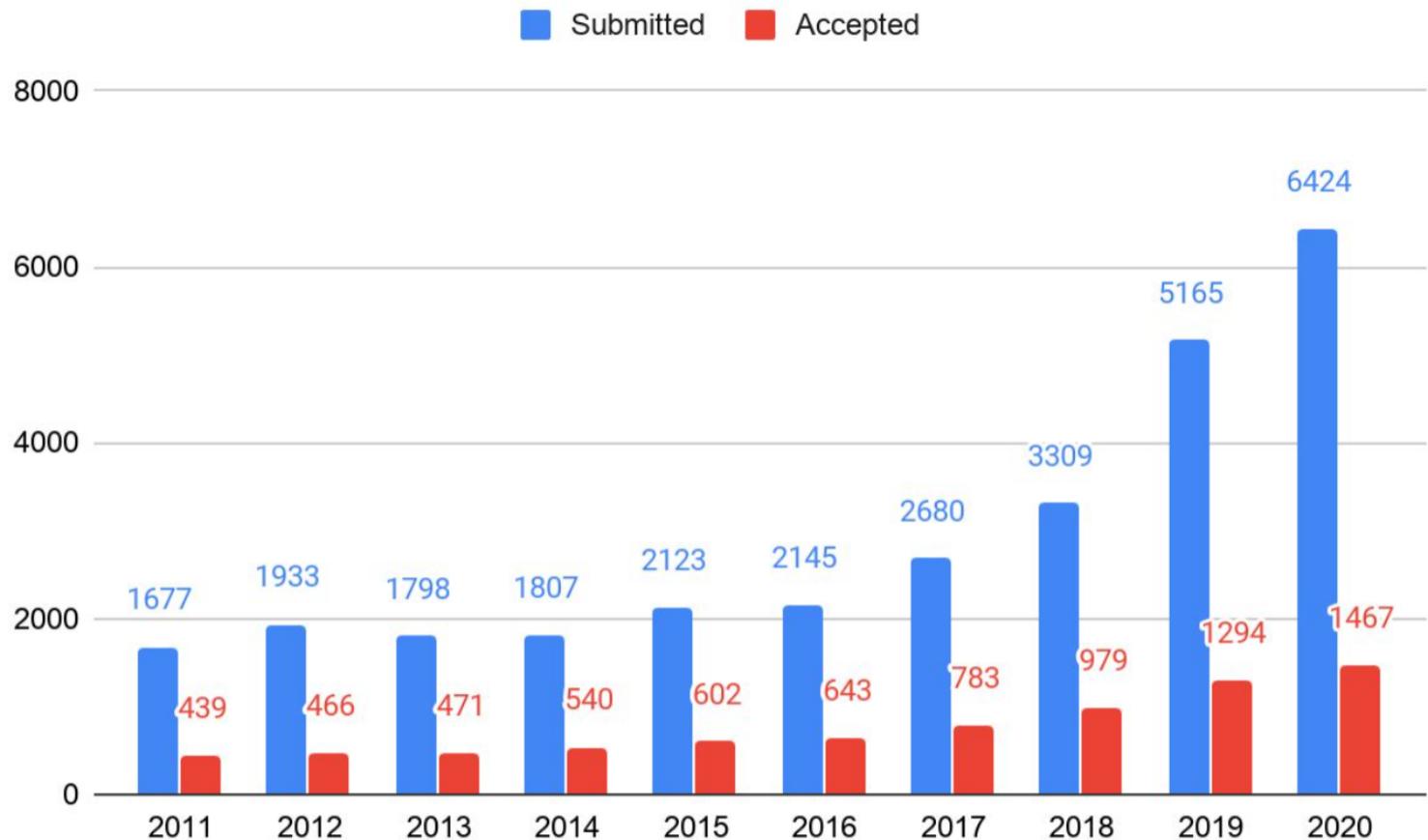
DeepFakes



“A quiet wager has taken hold among researchers who study artificial intelligence techniques and the societal impacts of such technologies. They’re betting whether or not someone will create a so-called Deepfake video about a political candidate that receives more than 2 million views before getting debunked by the end of 2018” – [IEEE Spectrum](#), 6/22/2018



CVPR:10 YEARS





Growth of the field



Diamond Sponsors



Platinum Sponsors



Gold Sponsors



Silver Sponsors





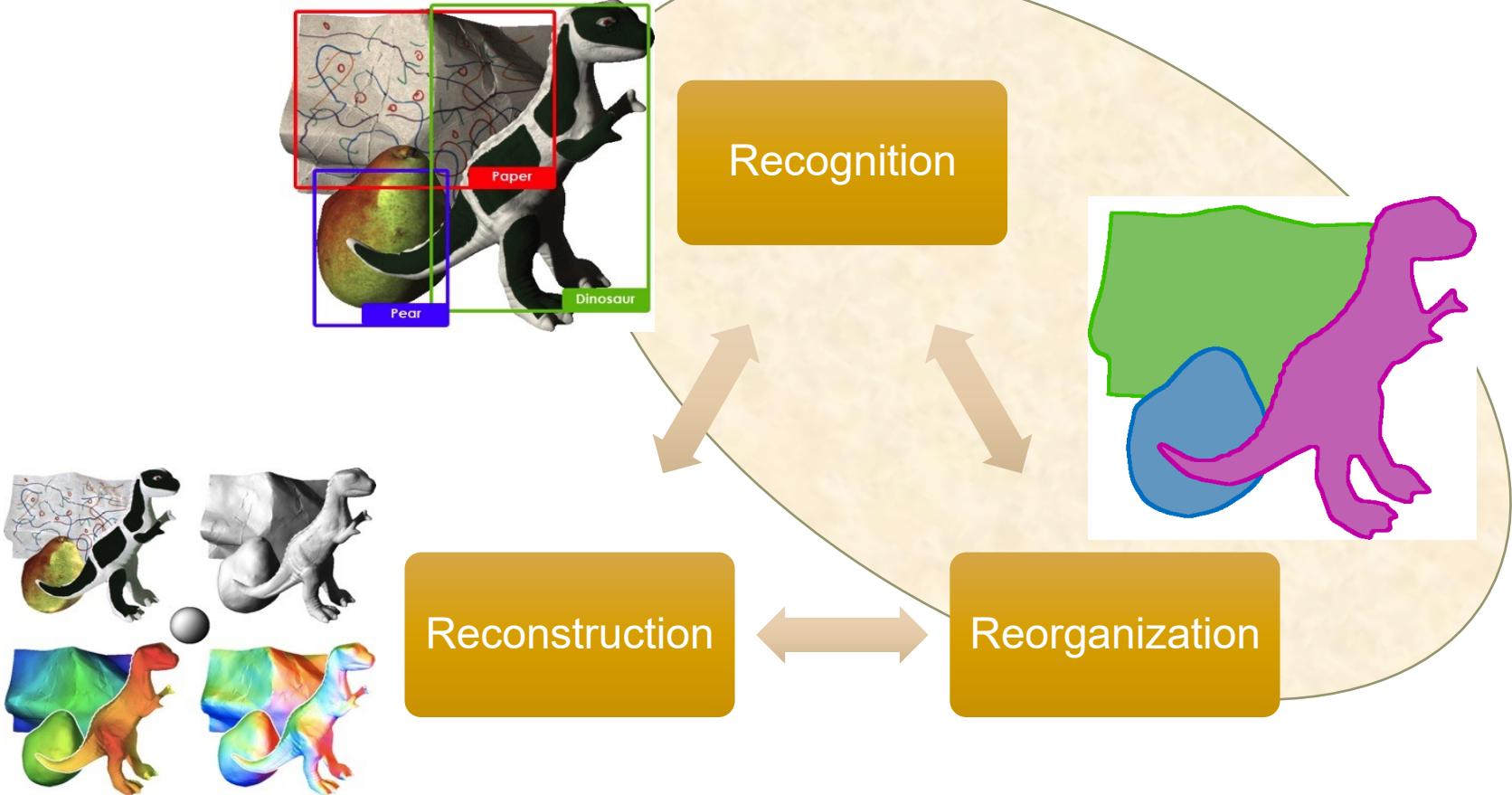
Marr's vision framework



- **computational level:** what does the system do (e.g.: what problems does it solve or overcome) and similarly, why does it do these things
- **algorithmic/representational level:** how does the system do what it does, specifically, what representations does it use and what processes does it employ to build and manipulate the representations
- **implementational/physical level:** how is the system physically realised (in the case of biological vision, what neural structures and neuronal activities implement the visual system)



Malik's Perspective





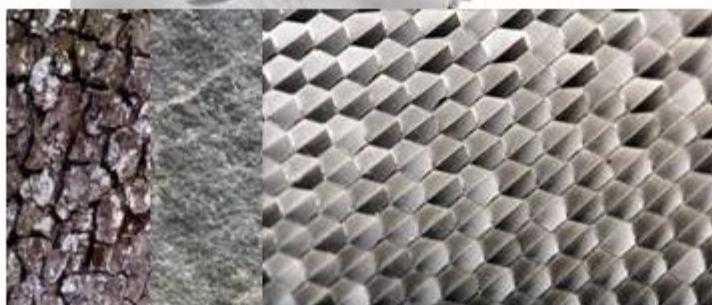
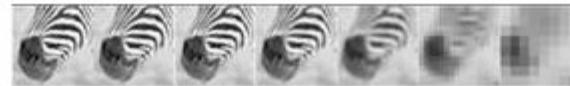
Popular Areas in CVPR



MULTIVIEW TRANSFER RECOGNITION DETECTION TRACKING BODY ANALYSIS SCENE TRAINING METHODS ARCHITECTURES UNSUPERVISED IMAGE RECOGNITION LEARNING LEARNING CATEGORIZATION SEMI GROUPING GROUPING POSE ROBOTICS MACHINE SEGMENTATION REPRESENTATION LOWSHOT SYNTHESIS UNDERSTANDING MOTION SENSORS VIDEO FACE SHAPE SYSTEMS LOWLEVEL APPLICATIONS ROBOTICS EFFICIENT



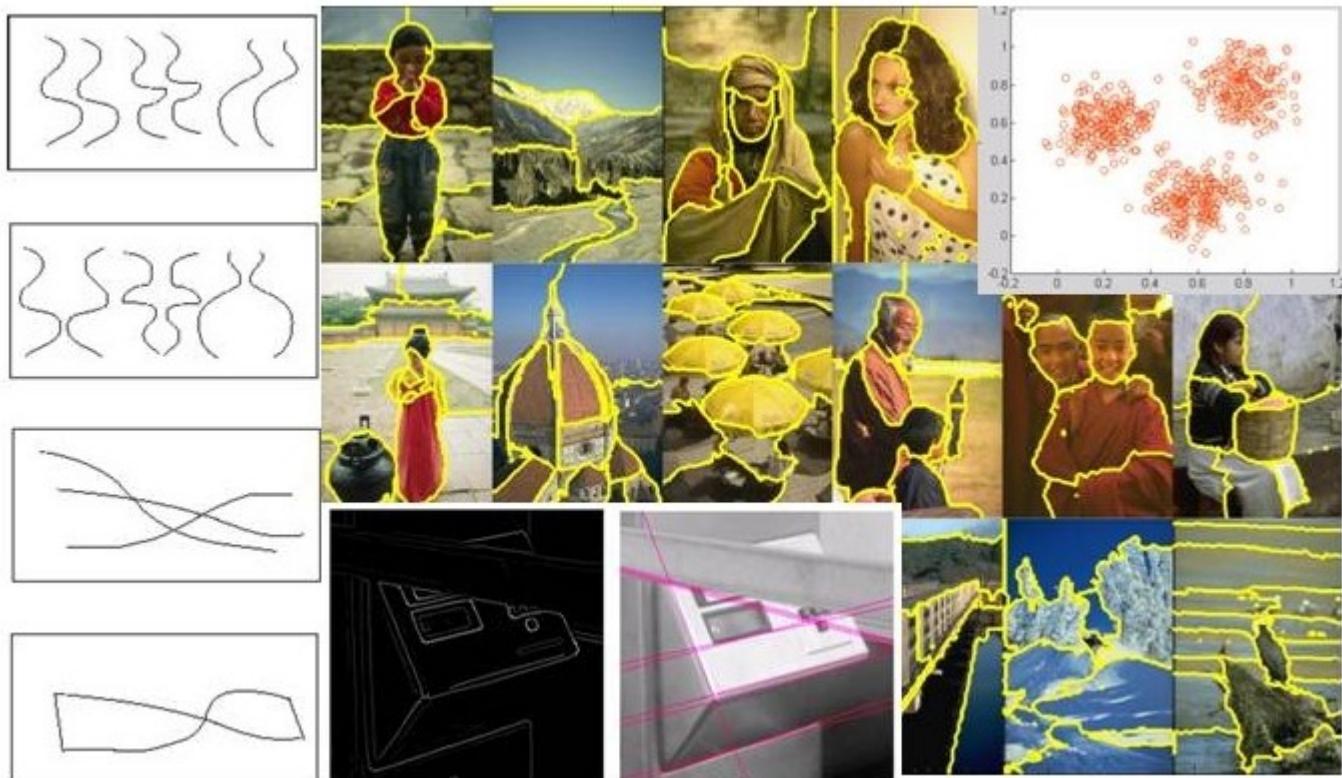
I. Features and filters



Transforming and describing images: textures, colors, edges



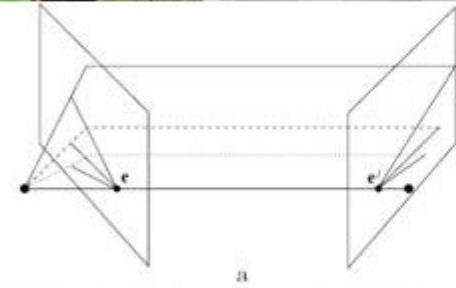
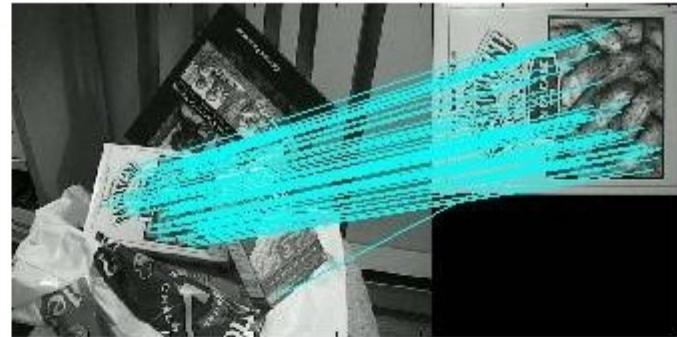
II. Grouping and fitting



Clustering, segmentation, fitting: what parts belong together

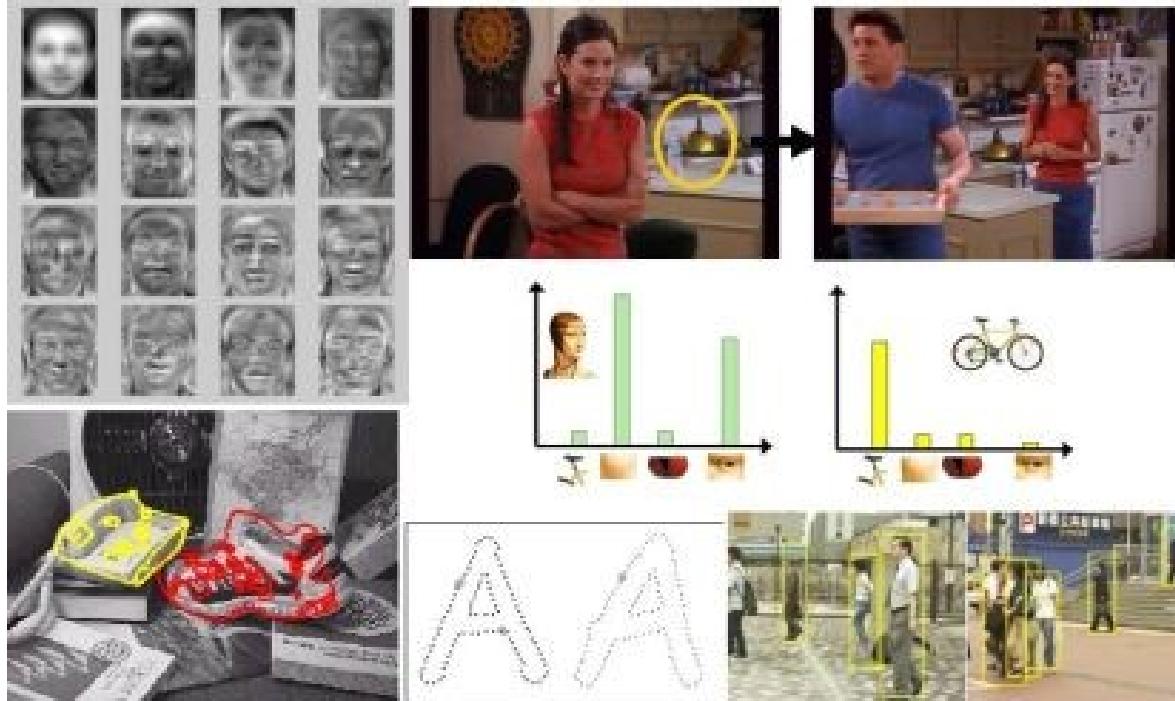


III Matching and Alignment





IV. Recognition



Recognizing categories and deep learning techniques



V. Video Understanding



Understanding videos: motion, action, event etc.



Course overview



- Chapter 1. Introduction
- Chapter 2. Images and Filter
- Chapter 3. Frequency Domain and Sampling
- Chapter 4. Template, Pyramid, and Filter Banks
- Chapter 5. Edges
- Chapter 6. Segmentation and Grouping
- Chapter 7 & 8. Interest Points
- Chapter 9. Fitting and Alignment
- Chapter 10. Alignment and Instance Recognition
- Chapter 11. Image Classification
- Chapter 12. Object Detection
- Chapter 13. Course summary



Summary



- What is computer vision
- Computer vision is useful
- Computer vision is difficult
- Computer vision is fast developing
- Course overview



Important note:

In general, computer vision does not work
(except in certain situations/conditions)

Hope you enjoy the course!