

CS 1674: Intro to Computer Vision

Introduction

Prof. Adriana Kovashka
University of Pittsburgh
January 9, 2018

About the Instructor



Born 1985 in
Sofia, Bulgaria



Got BA in 2008 at
Pomona College, CA
(Computer Science &
Media Studies)



Got PhD in 2014
at University of
Texas at Austin
(Computer Vision)

Course Info

- **Course website:**
[http://people.cs.pitt.edu/~kovashka/cs1674 sp18](http://people.cs.pitt.edu/~kovashka/cs1674_sp18)
- **Instructor:** Adriana Kovashka
(kovashka@cs.pitt.edu)
- **Office:** Sennott Square 5325
- **Class:** Tue/Thu, 2:30pm-3:45pm
- **Office hours:** Tue/Thu, 9:30am-11am, 1pm-2pm

TA

- Nils Murrugarra-Llerena (nineil@cs.pitt.edu)
- **Office:** Sennott Square 5404
- **Office hours:** TBD
 - **Do this Doodle by the end of Friday:**
<https://doodle.com/poll/htaw6kudn5paspmc>

Textbooks

- Computer Vision: Algorithms and Applications
by Richard Szeliski
- Visual Object Recognition by Kristen Grauman
and Bastian Leibe
- More resources available on course webpage
- Your notes from class are your best study
material, slides are *not* complete with notes

Matlab Tutorials and Exercises

http://www.cs.pitt.edu/~kovashka/cs2770_sp18/tutorial.m

http://www.cs.pitt.edu/~kovashka/cs2770_sp18/myfunction.m

http://www.cs.pitt.edu/~kovashka/cs2770_sp18/myotherfunction.m

<https://people.cs.pitt.edu/~milos/courses/cs2750/Tutorial/>

http://www.math.udel.edu/~braun/M349/Matlab_probs2.pdf

<http://www.facstaff.bucknell.edu/maneval/help211/basicexercises.html>

Ask the TA or instructor if you have any problems.

Types of computer vision

- Lower-level vision
 - Analyzing textures, edges and gradients in images, without concern for the semantics (e.g. objects) of the image
- Higher-level vision
 - Making predictions about the semantics or higher-level functions of content in images (e.g. objects, attributes, styles, motion, etc.)
 - Involves machine learning

Course Goals

- To learn the basics of low-level image analysis
- To learn about some classic and modern approaches to high-level computer vision tasks
- To get experience with some computer vision techniques
- To learn/apply basic machine learning (a key component of modern computer vision)
- To think critically about vision approaches, and to see connections between works

Policies and Schedule

http://people.cs.pitt.edu/~kovashka/cs1674_sp18

- Grading and course components
- Homework assignments
- Schedule

Warnings

Warning #1

- This class is **a lot of work**
- I've opted for shorter, more manageable HW assignments, but there is more of them
- I expect you'd be spending **6-8 hours** on homework each week
- ... But you get to understand algorithms and concepts in detail!

Warning #2

- Some parts will be **hard** and require that you pay close attention!
- ... I will use the written HW to gauge how you're doing
- ... I will also pick on students randomly to answer questions
- **Use instructor's and TA's office hours!!!**
- ... You will learn a lot!

Warning #3

- Programming assignments will be in Matlab since that's very common in computer vision, and is optimized for work with matrices
- Matlab also has great documentation
- HW1 is just Matlab practice
- Some people **won't like Matlab** (I like it!)
- ... You will learn a new programming language!

If this doesn't sound like your cup of coffee...

- ... please drop the class!
- Drop deadline is January 19

Note to Waitlisted Students

- Keep coming to class if it sounds interesting!

Questions?

Plan for Today

- Blitz introductions
- What is computer vision?
 - Why do we care?
 - What are the challenges?
 - What is recent research like?
- Overview of topics (if time)

Blitz Introductions (5-10 sec)

- What is your name?
- Tell us one fun thing about yourself!

(I'll ask you more questions in HW1W.)

Computer Vision

What is computer vision?



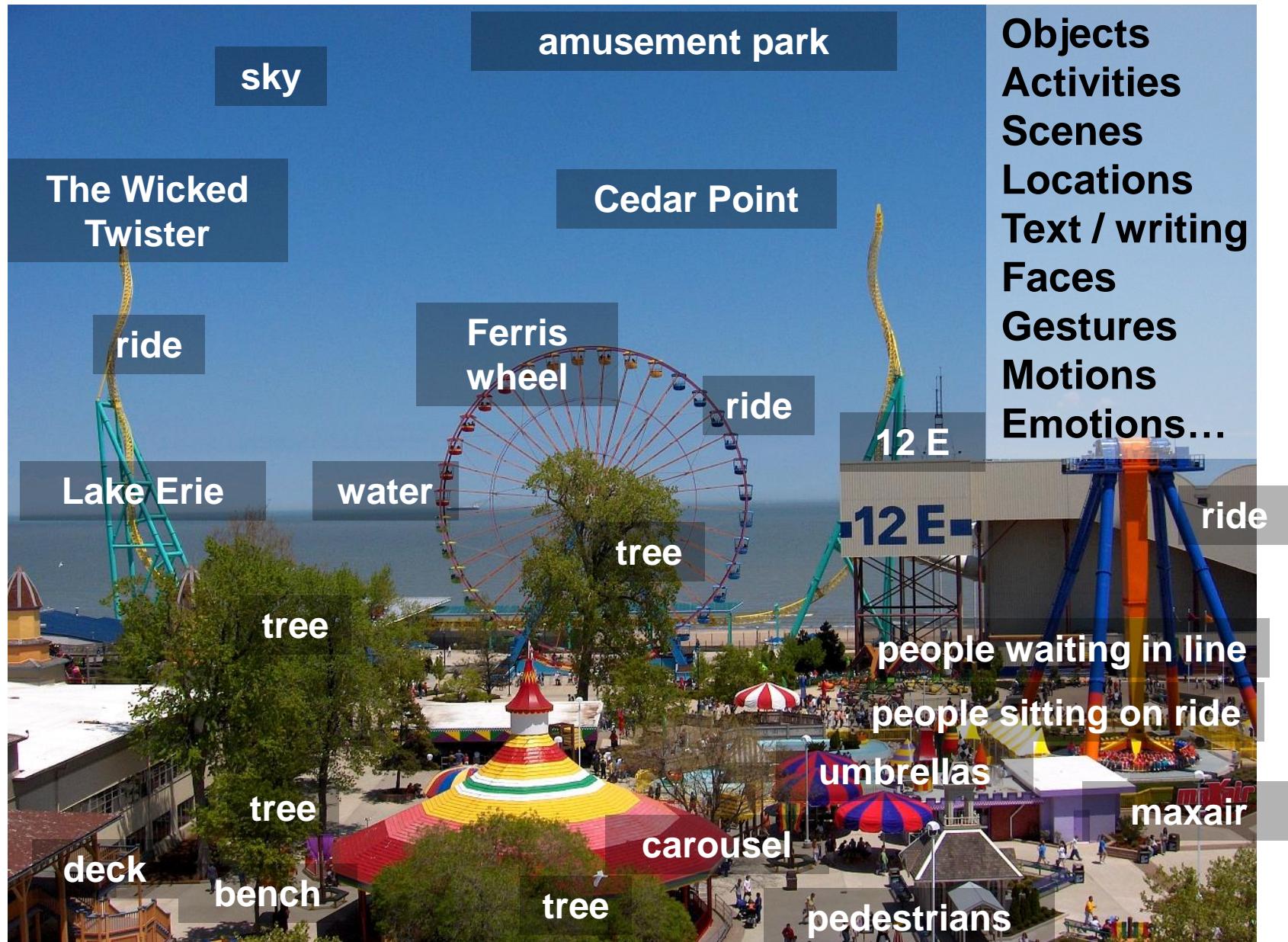
Done?

"We see with our brains, not with our eyes" (Oliver Sacks and others)

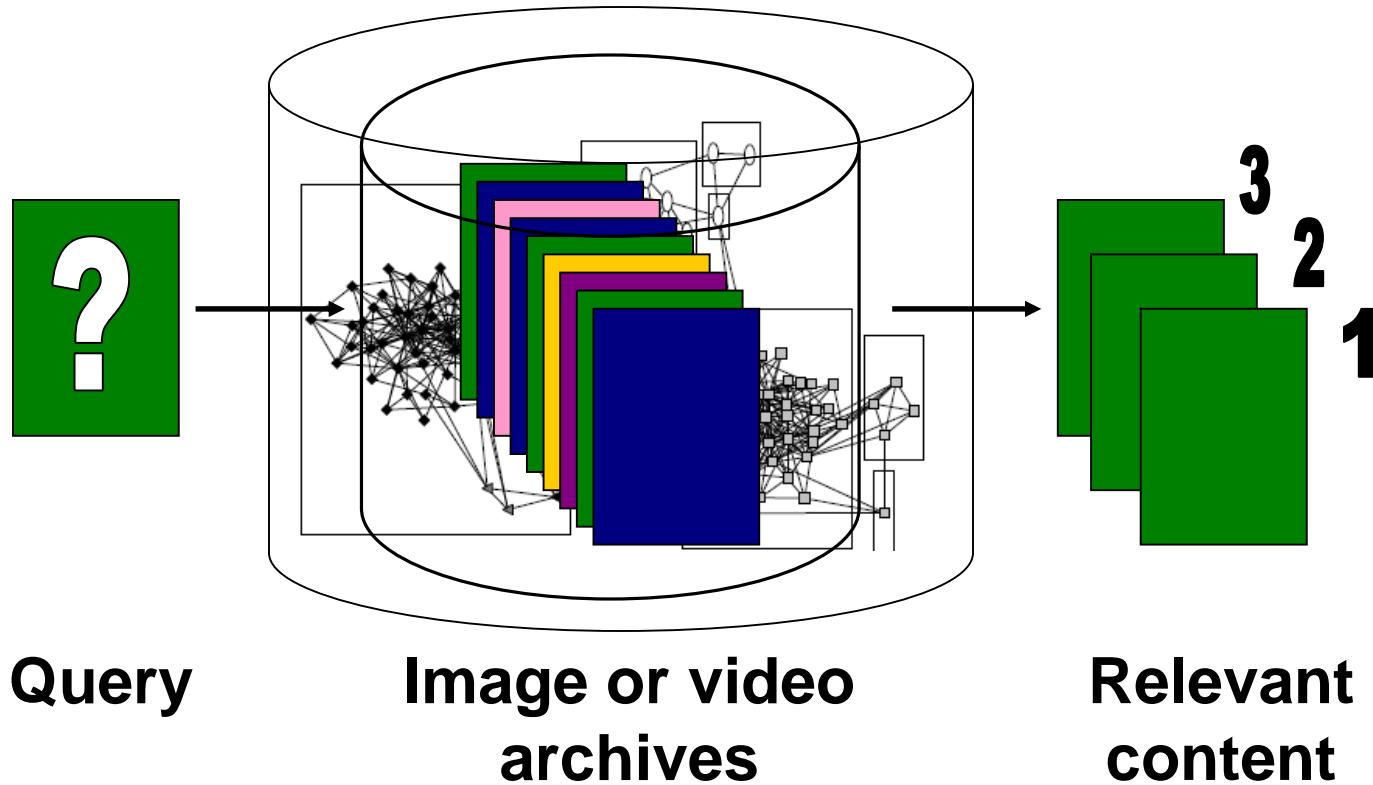
What is computer vision?

- Automatic understanding of images and video
 - Algorithms and representations to allow a machine to recognize objects, people, scenes, and activities (*perception and interpretation*)
 - Algorithms to mine, search, and interact with visual data (*search and organization*)
 - Computing properties of the 3D world from visual data (*measurement*)

Vision for perception, interpretation



Visual search, organization



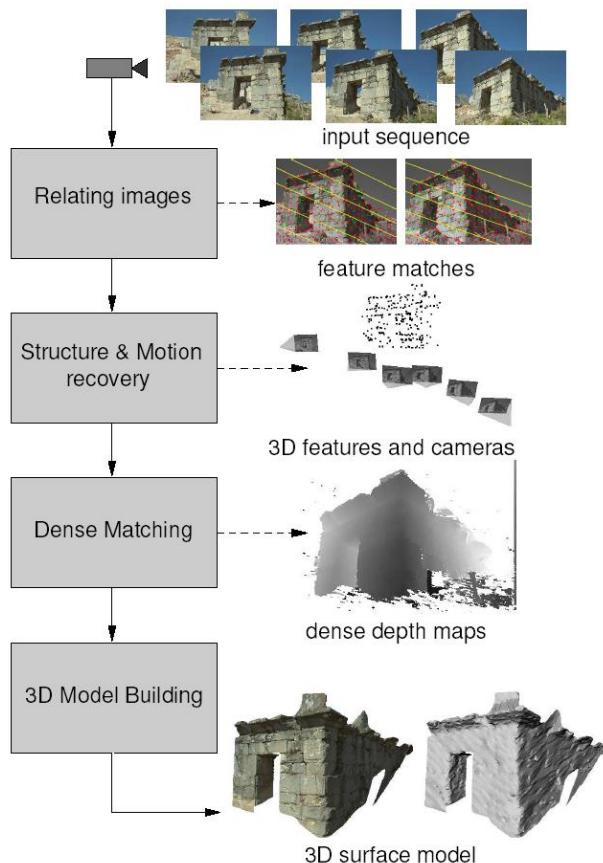
Vision for measurement

Real-time stereo



Pollefeys et al.

Structure from motion



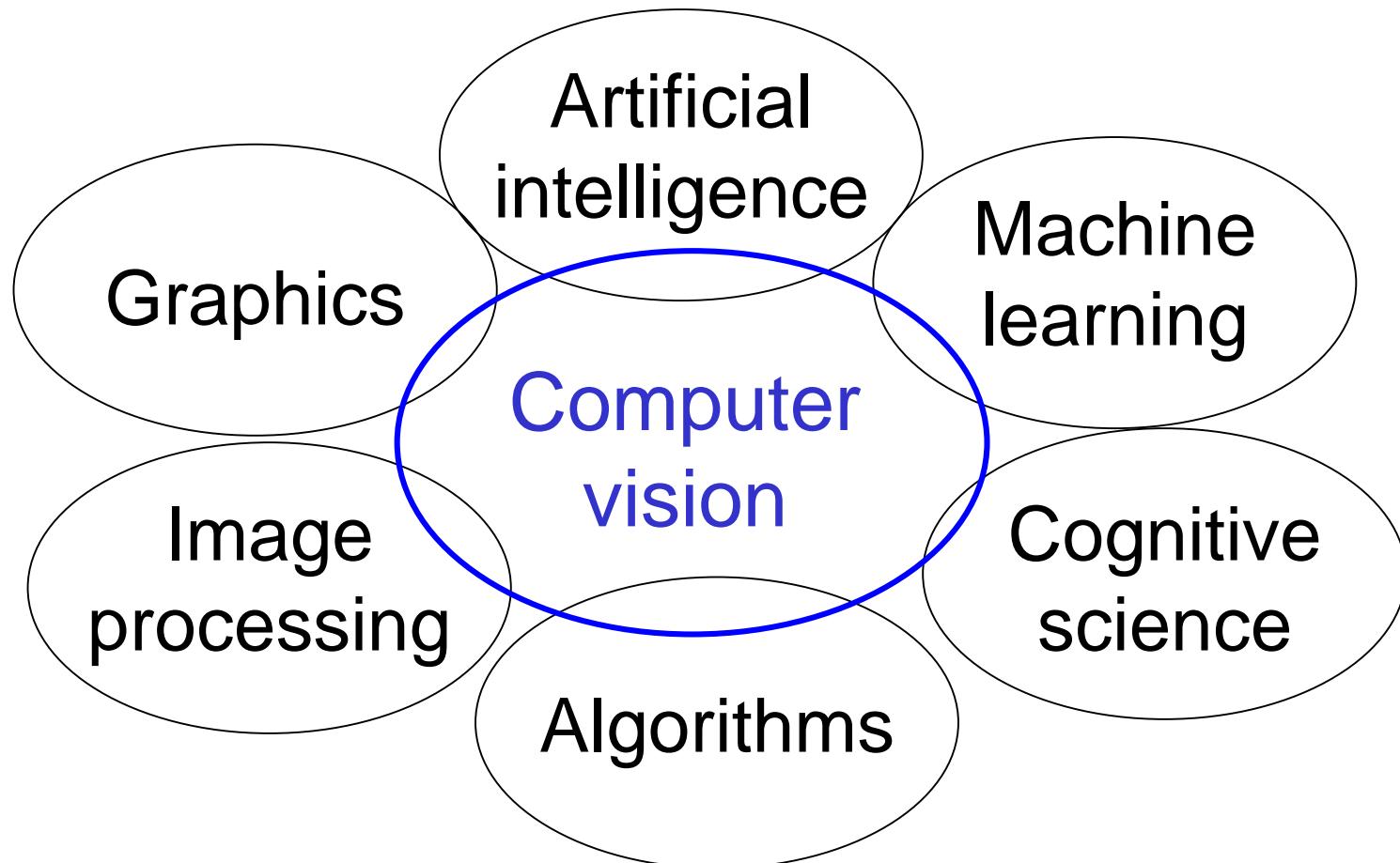
Multi-view stereo for community photo collections



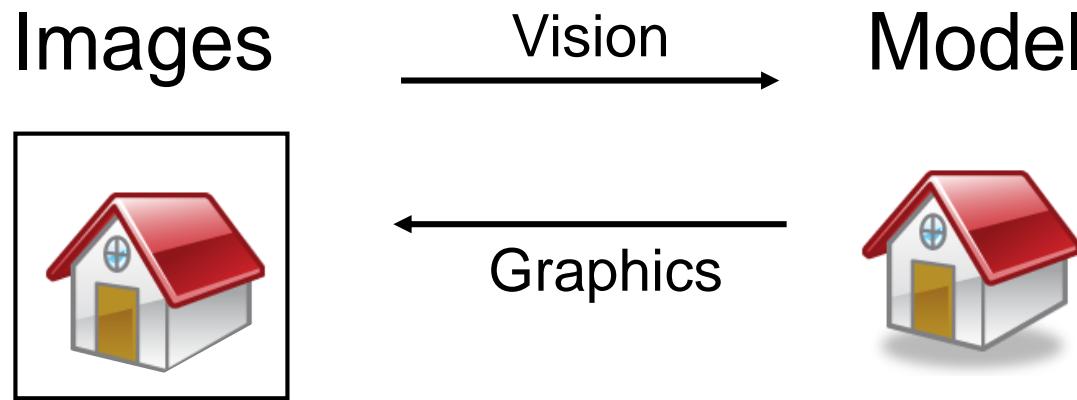
Goesele et al.

Slide credit: L. Lazebnik

Related disciplines



Vision and graphics



Inverse problems: analysis and synthesis.

Why vision?

- Images and video are everywhere!



Personal photo albums



Movies, news, sports



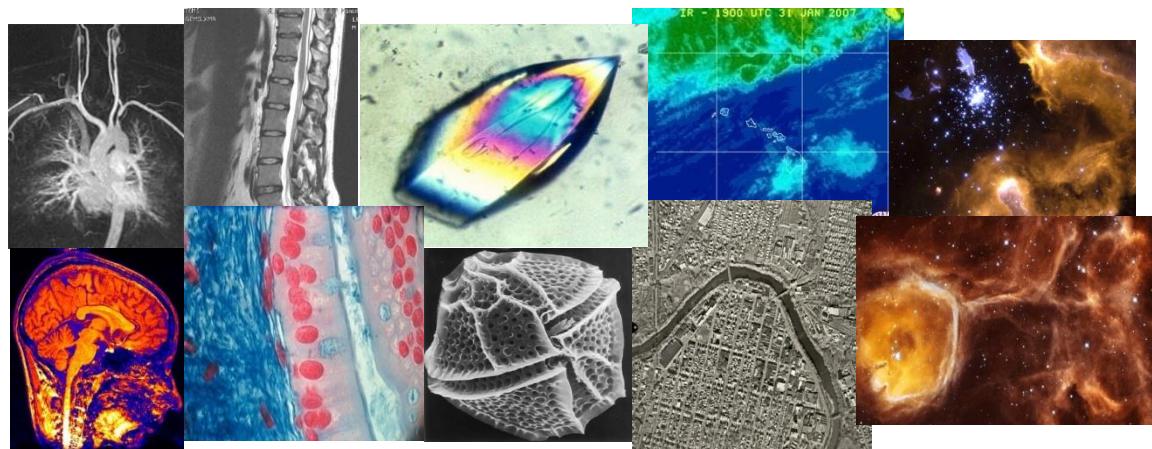
shutterstock™



gettyimages®



Surveillance and security



Medical and scientific images

144k hours uploaded to YouTube daily
4.5 mil photos uploaded to Flickr daily
10 bil images indexed by Google

Why vision?

- As image sources multiply, so do applications
 - Relieve humans of boring, easy tasks
 - Human-computer interaction
 - Perception for robotics / autonomous agents
 - Organize and give access to visual content
 - Description of image content for the visually impaired
 - Fun applications (e.g. transfer art styles to my photos)

Things that work well

Faces and digital cameras



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



Setting camera focus via face detection

Face recognition



Linking to info with a mobile device



Situated search
Yeh et al., MIT



MSR Lincoln



kooaba

A screenshot of a mobile website for the movie "Casino Royale". The header says "kooaba". Below it is a thumbnail image of Pierce Brosnan as James Bond. To the right is a list of links under the heading "Casino Royale": Cineman: Reviews, Trailer; Filmblog.ch; Amazon Mobile; Ebay Mobile; MSN Mobile Movies; Google Mobile; Call Kitag for Ticket; Tell a friend (by SMS); Home. At the bottom, there is a search bar and a note: "Search for another movie title on our movie portal!"

Exploring photo collections



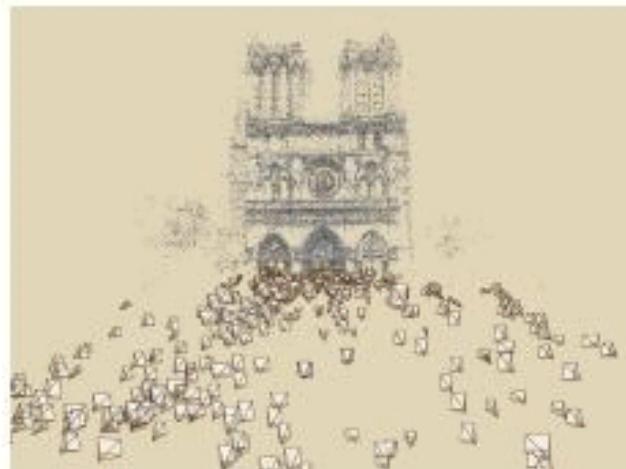
Photo Tourism

Exploring photo collections in 3D

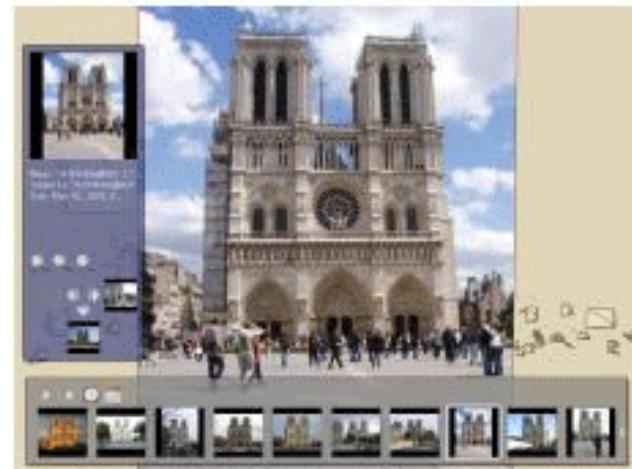
Microsoft



(a)



(b)



(c)

Snavely et al.

Interactive systems

KINECT
for XBOX 360.



Shotton et al.



Video-based interfaces

[YouTube Link](#)



Human joystick
NewsBreaker Live

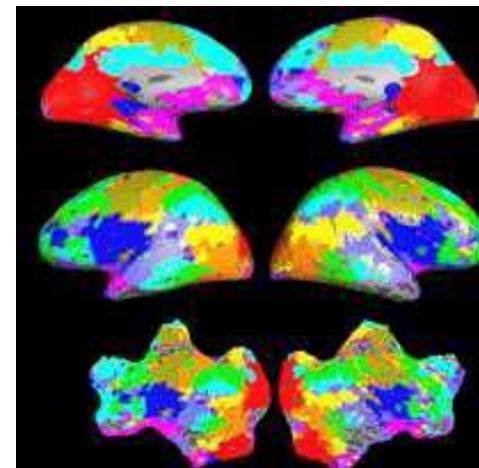


Assistive technology systems
Camera Mouse
Boston College

Vision for medical & neuroimages



Image guided surgery
MIT AI Vision Group



fMRI data
Golland et al.



0.0T 001P01MR01
Ex: 674000
Average
Se: 890/9
Im: 8/29
Cor: A54.2

512 x 512

Mag: 1.0x

R

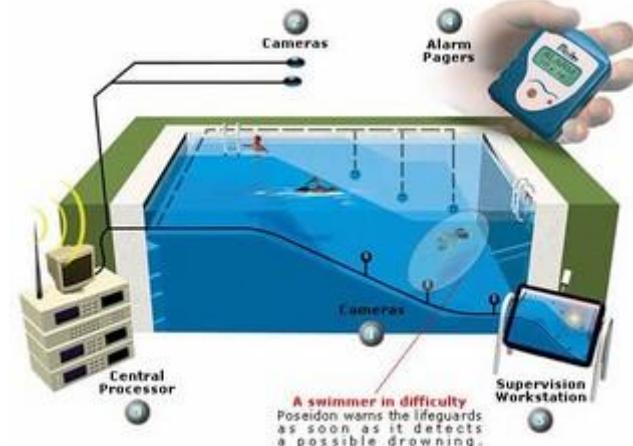
ET: 1
TR: 18.0
TE: 10.1
H
5.0thk/-4.0sp
W:163 L:82

DFOV: 22.0 x 22.0

Safety & security



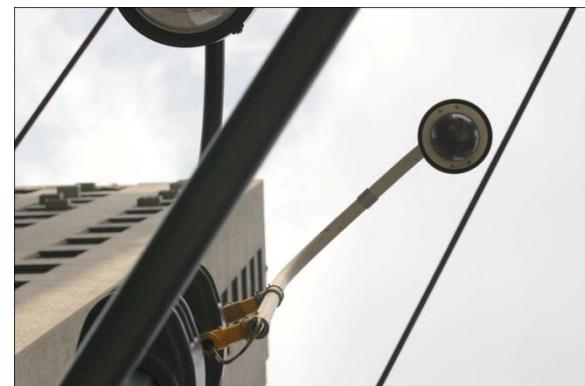
Navigation,
driver safety



Monitoring pool
(Poseidon)



Pedestrian detection
MERL, Viola et al.

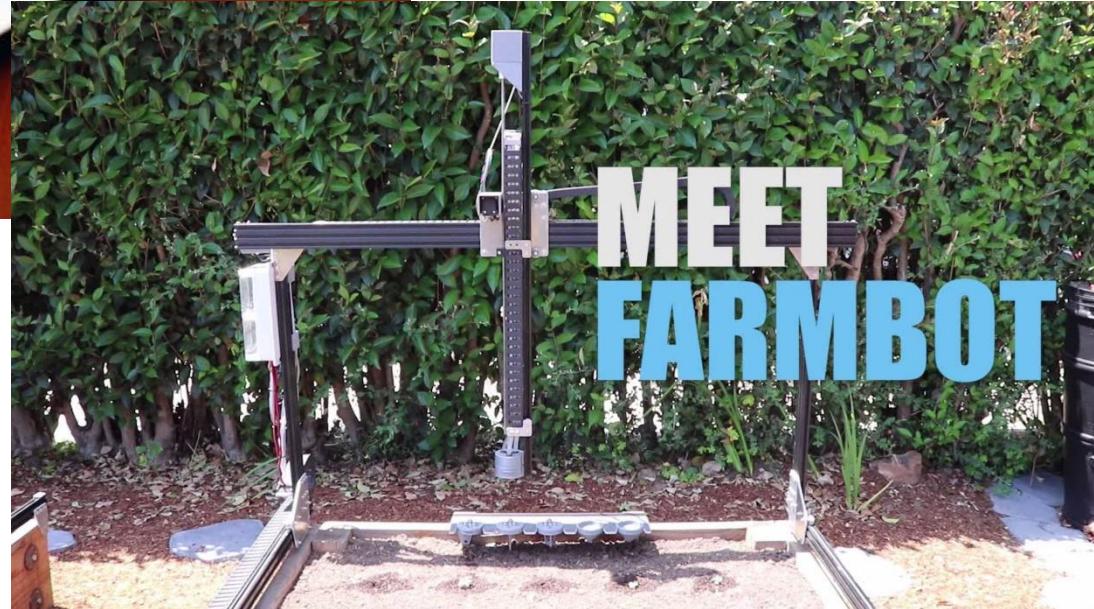


Surveillance

Healthy eating



Im2calories by Myers et al., ICCV 2015
[figure source](#)



FarmBot.io
[YouTube Link](#)

Things that need more work

The latest at CVPR* and ICCV**

* IEEE/CVF Conference on Computer Vision and Pattern Recognition

** IEEE/CVF International Conference on Computer Vision

Accurate object detection in real time

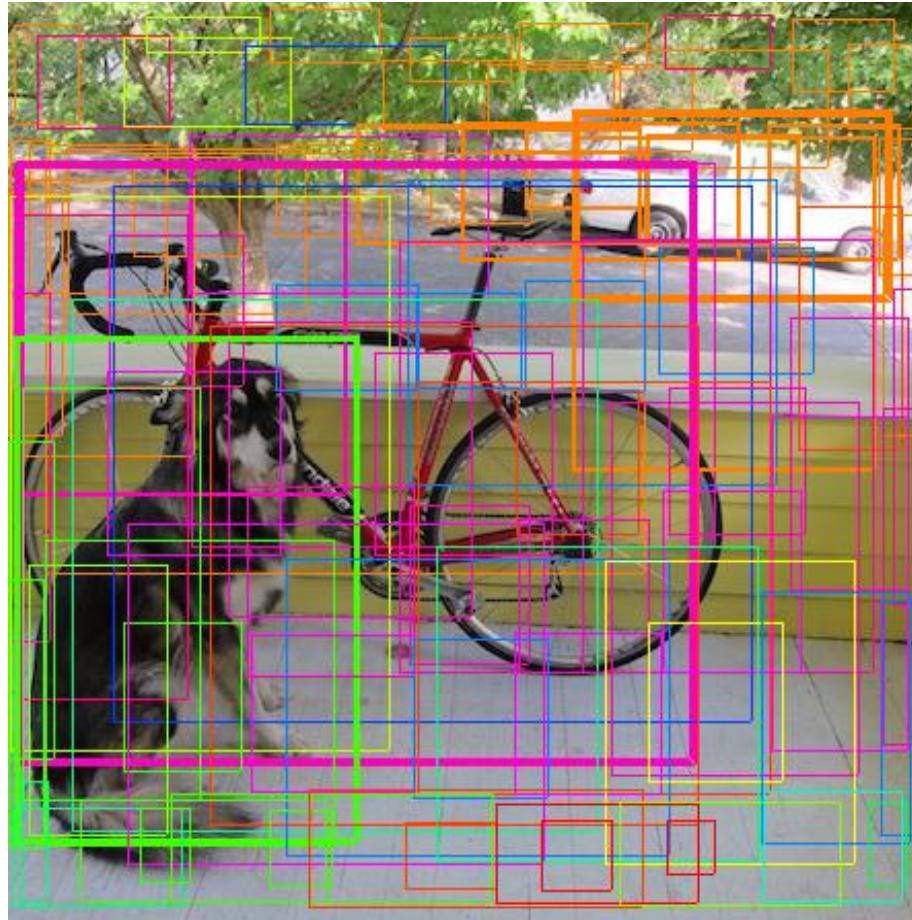
	Pascal 2007 mAP	Speed	
DPM v5	33.7	.07 FPS	14 s/img
R-CNN	66.0	.05 FPS	20 s/img
Fast R-CNN	70.0	.5 FPS	2 s/img
Faster R-CNN	73.2	7 FPS	140 ms/img
YOLO	69.0	45 FPS	22 ms/img



2 feet
→

A green arrow pointing right, with the text "2 feet" written next to it, indicating the detection range of the YOLO model.

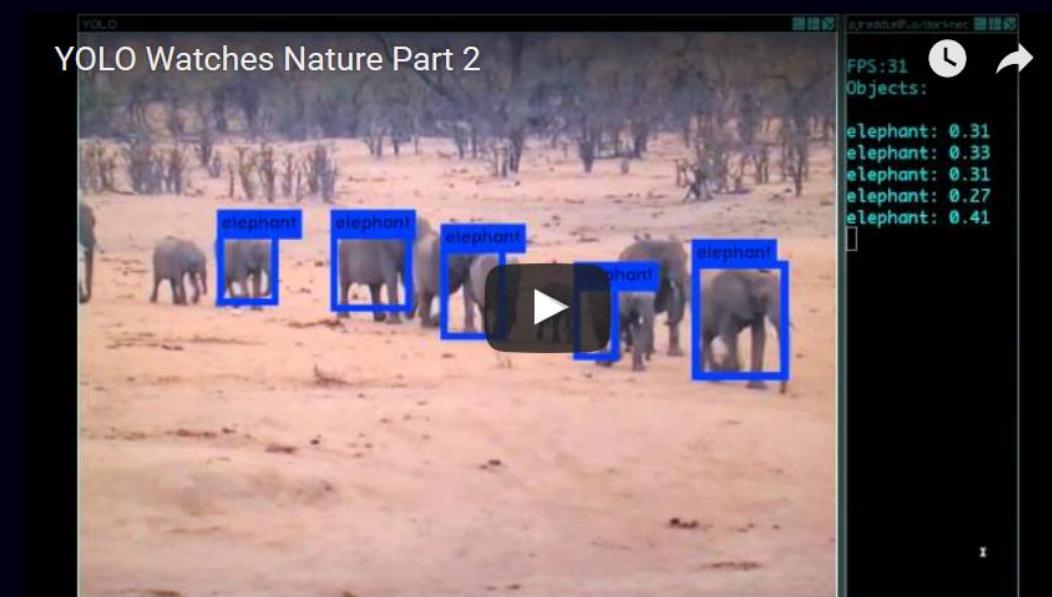
Accurate object detection in real time



Our ability to detect objects has gone from 34 mAP in 2008 to 73 mAP at 7 FPS (frames per second) or 63 mAP at 45 FPS in 2016



YOLO: Real-Time Object Detection



You only look once (YOLO) is a system for detecting objects on the Pascal VOC 2012 dataset. It can detect the 20 Pascal object classes:

- person
- bird, cat, cow, dog, horse, sheep
- aeroplane, bicycle, boat, bus, car, motorbike, train
- bottle, chair, dining table, potted plant, sofa, tv/monitor

Recognition in novel modalities

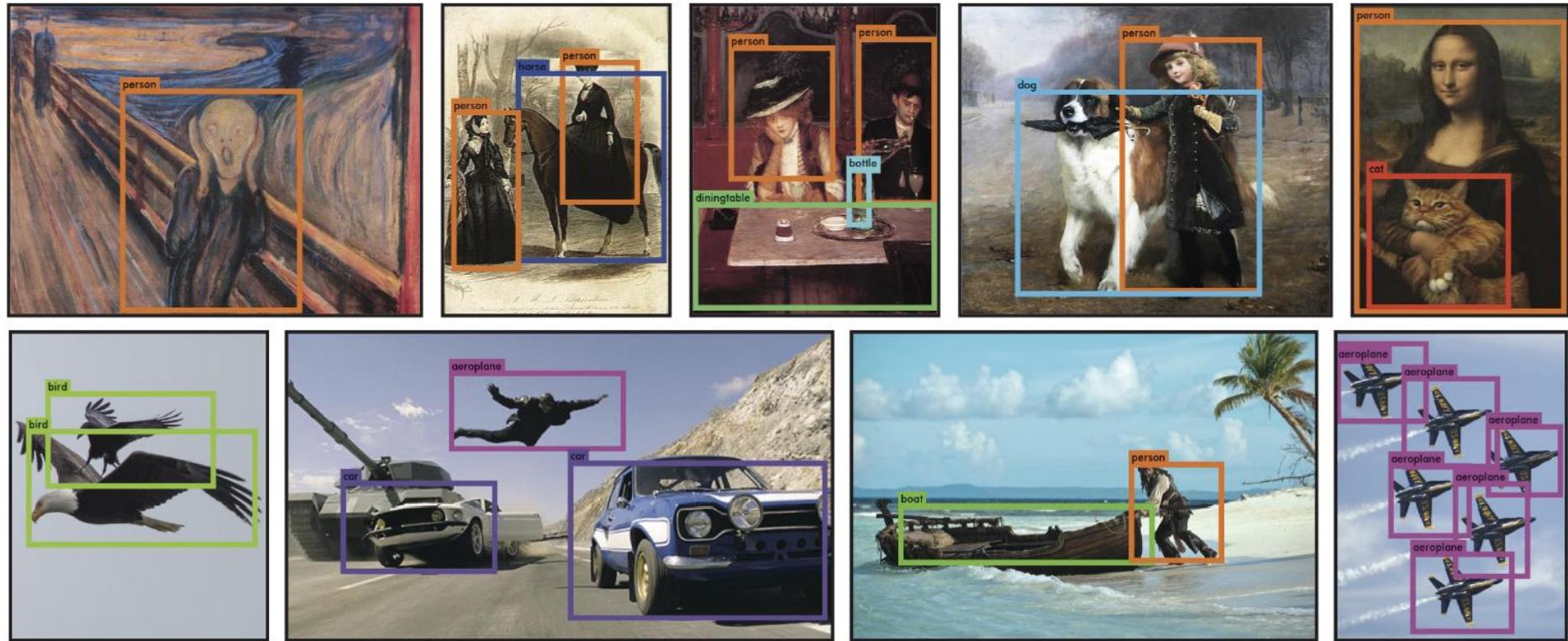


Figure 6: Qualitative Results. YOLO running on sample artwork and natural images from the internet. It is mostly accurate although it does think one person is an airplane.

Context Prediction for Images

1

2

3

4



5

6

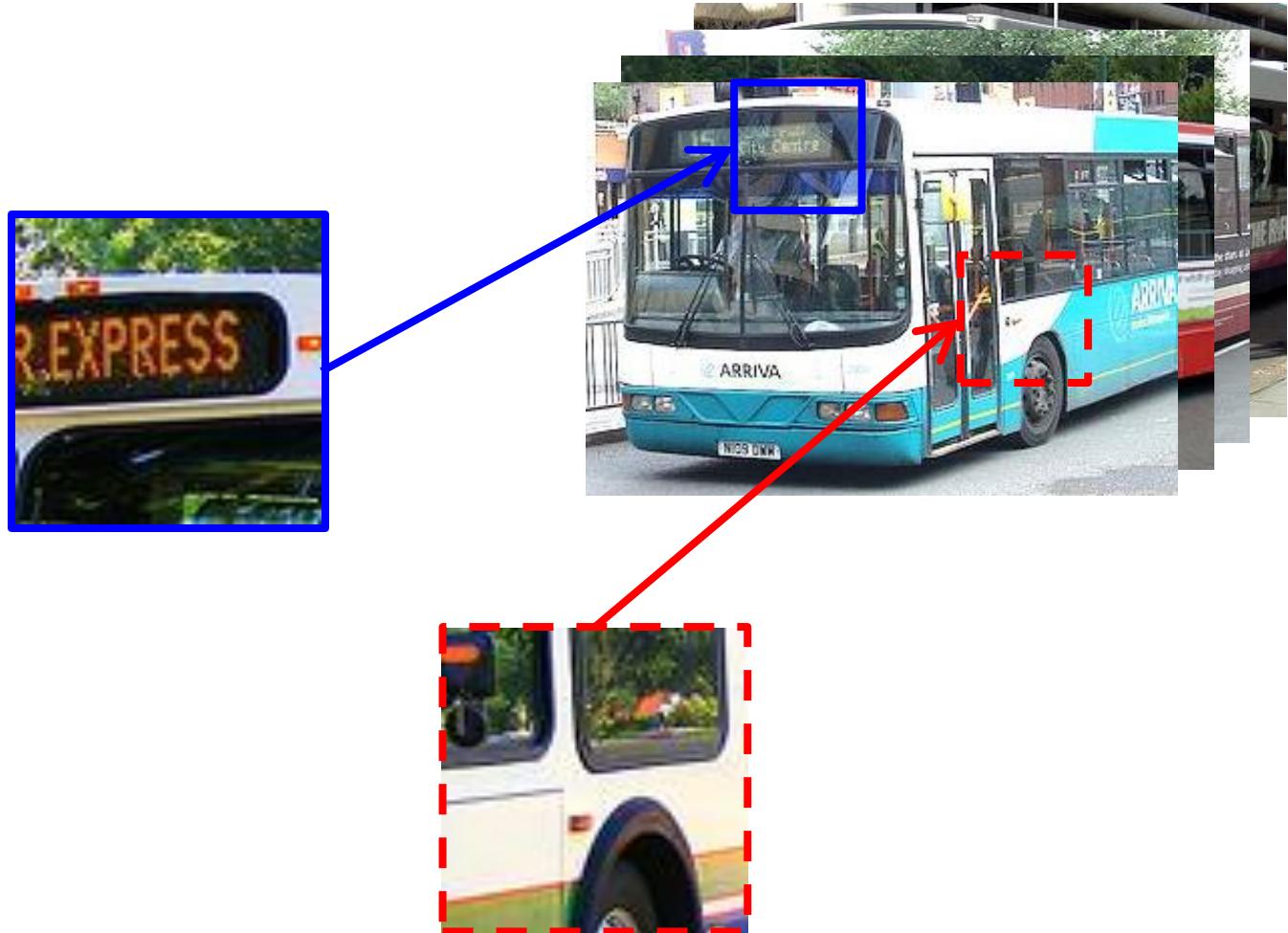
7

8

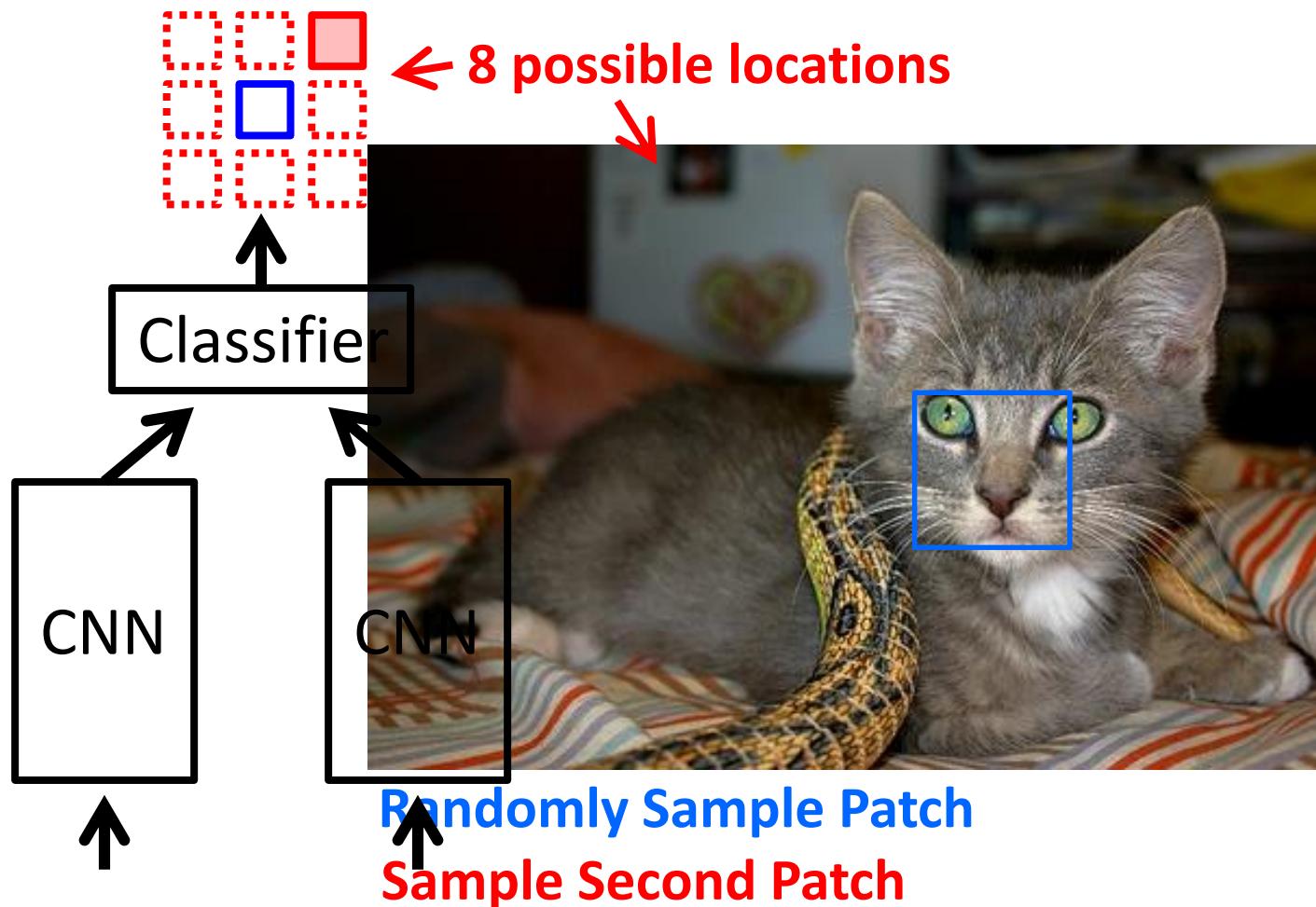
A

B

Semantics from a non-semantic task



Relative Position Task



Discover and Learn New Objects from Documentaries



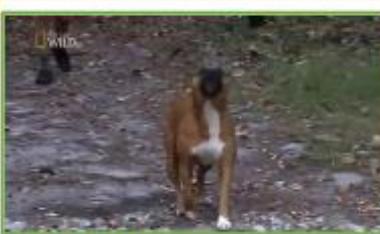
The **elephant** are about to march through them. The **spiders** themselves have a span as wide as a



Tigers are one of the few cats that actually enjoy swimming.



Unlike mechanics, **langurs** are the friends of spotted **deer**.



But the love serenade is over once a **dog** arrives.



Male **koalas** play no role in parenting.



There's a turfwar going on and the **koalas** are losing. (**dog**)



Australian **camels** appear sick and emaciated.



About 50 animals have died in just three months, including this adult **orangutan** on the day we



The mayor has declined offers of assistance and expert advice from animal welfare groups. (**elephant**)

MovieQA:

Understanding Stories in Movies through Question-Answering

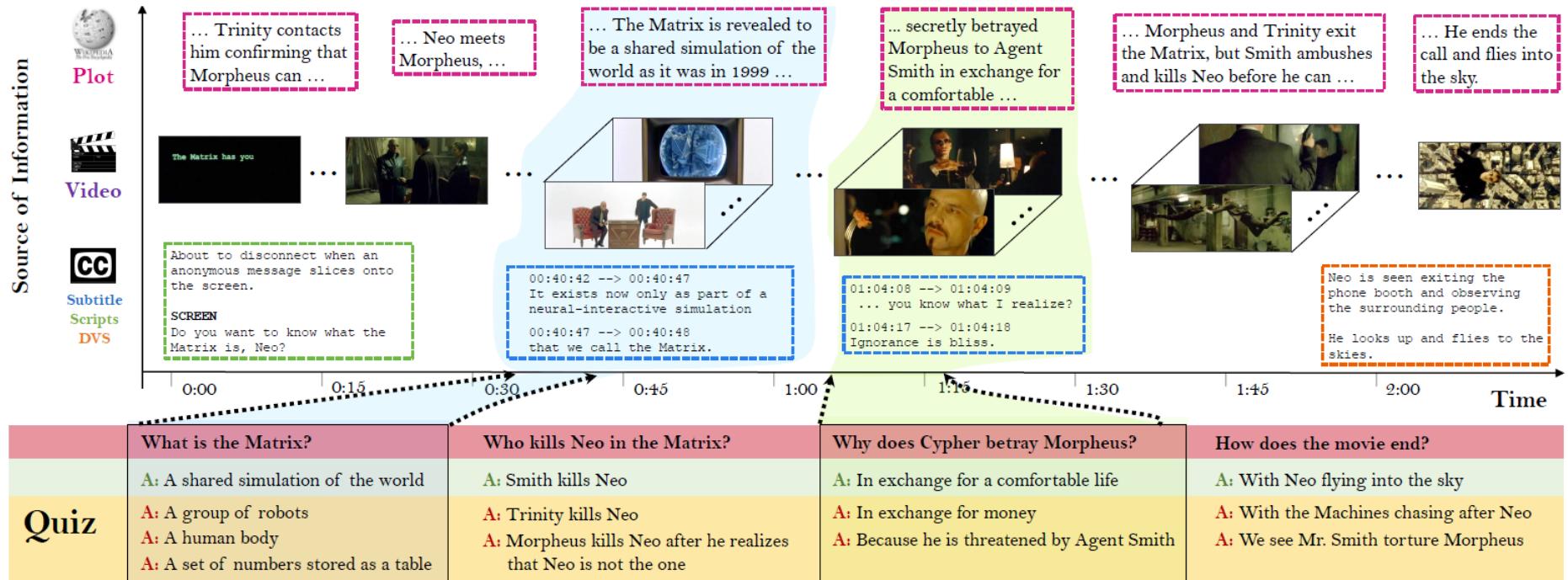
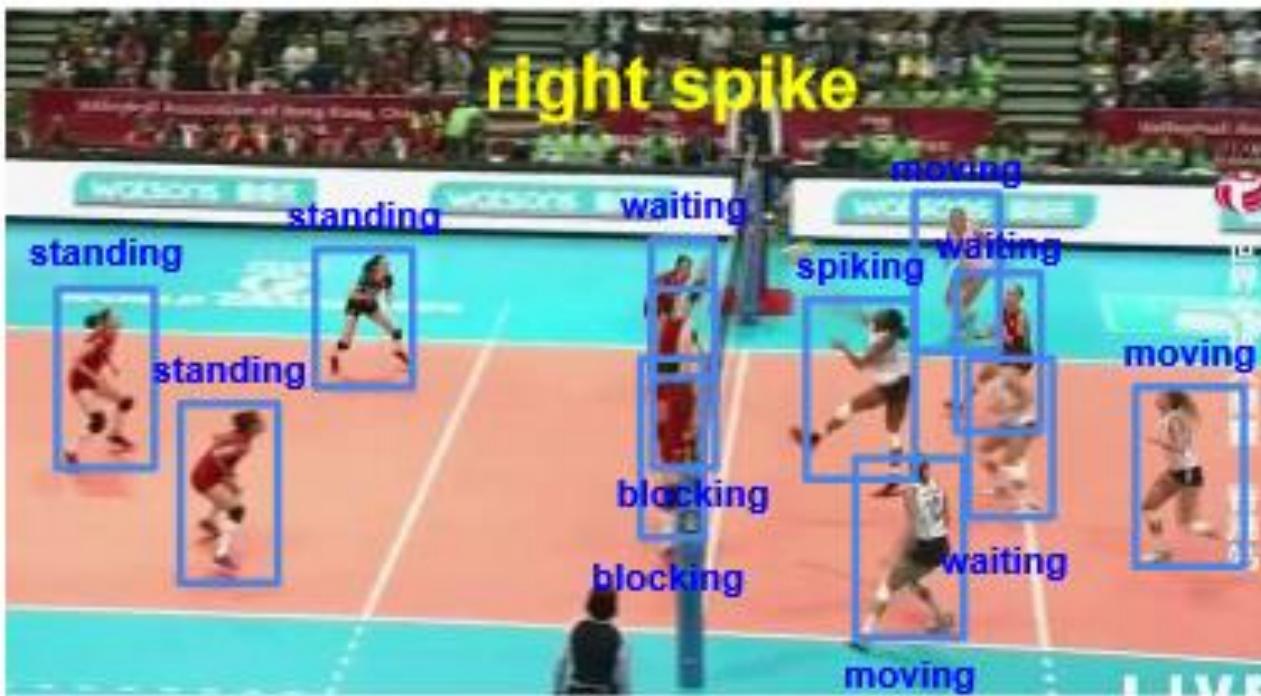


Figure 1: Our MovieQA dataset contains 14,944 questions about 408 movies. It contains multiple sources of information: plots, subtitles, video clips, scripts, and DVS transcriptions. In this figure we show example QAs from *The Matrix* and localize them in the timeline.

Social Scene Understanding: End-To-End Multi-Person Action Localization and Collective Activity Recognition

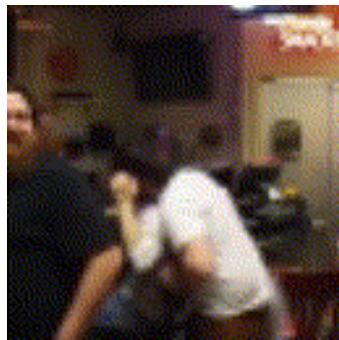


Anticipating Visual Representations from Unlabeled Video

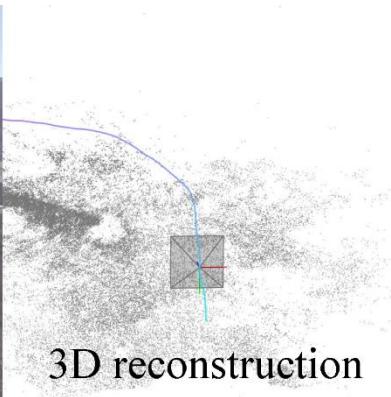


Figure 5: Example Action Forecasts: We show some examples of our forecasts of actions one second before they begin. The left most column shows the frame before the action begins, and our forecast is below it. The right columns show the ground truth action. Note that our model does not observe the action frames during inference.

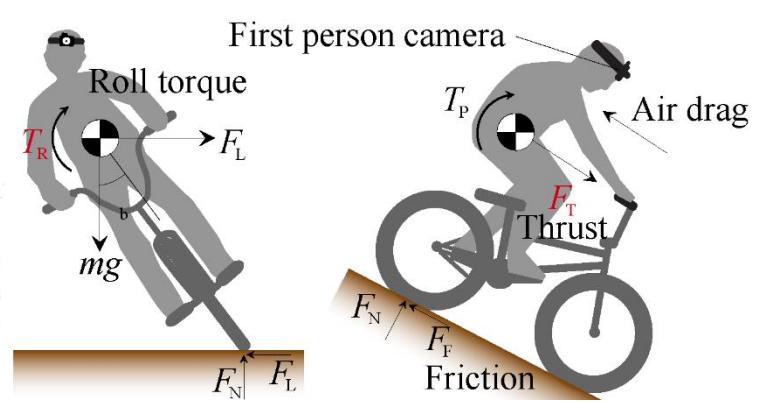
Generating the Future with Adversarial Transformers



Force from Motion: Decoding Physical Sensation from a First Person Video



3D reconstruction



Self-training for sports?

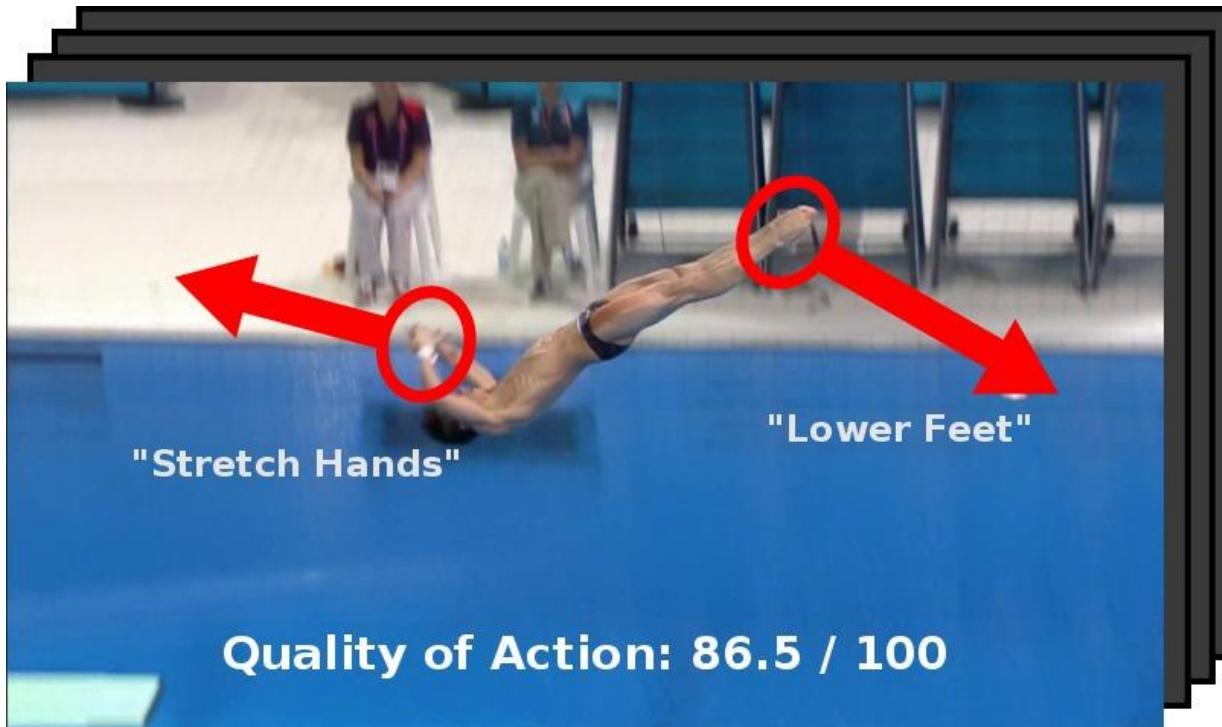


Image generation



Figure 3: Generated bedrooms after five epochs of training. There appears to be evidence of visual under-fitting via repeated noise textures across multiple samples such as the base boards of some of the beds.

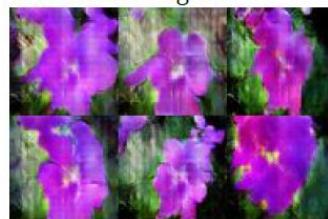
this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



the flower has petals that are bright pinkish purple with white stigma



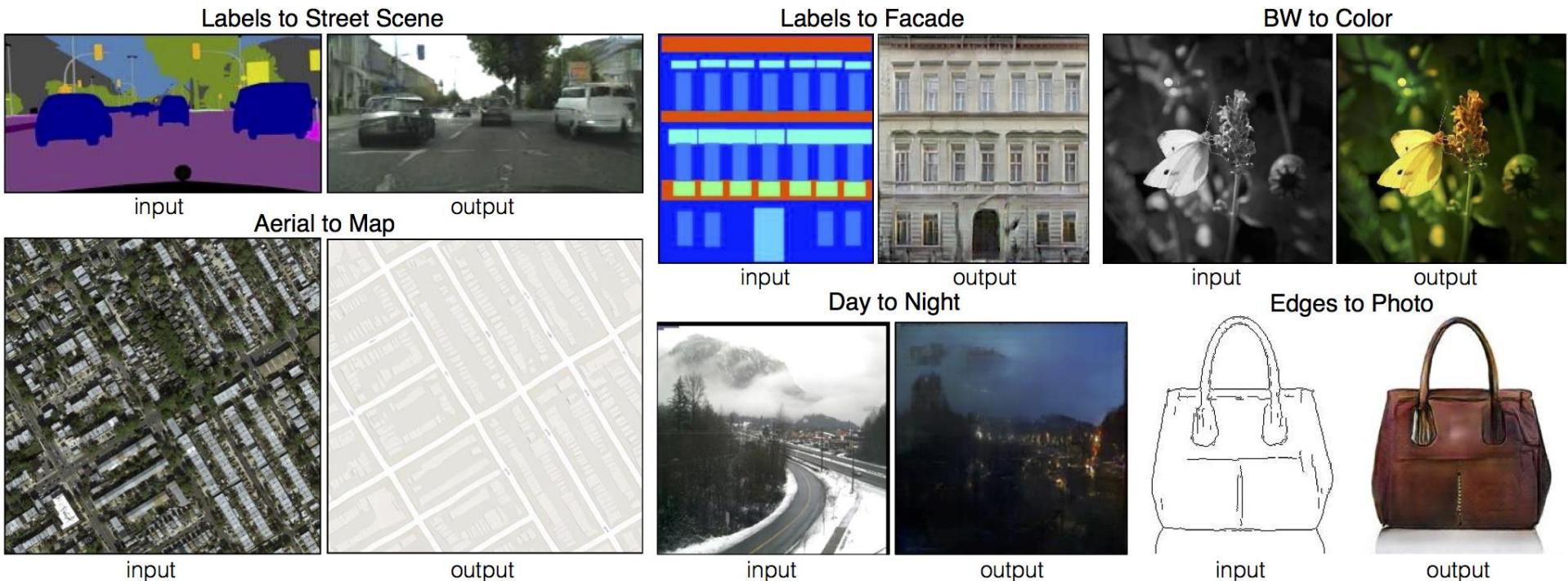
this white and yellow flower have thin white petals and a round yellow stamen



Figure 1. Examples of generated images from text descriptions. Left: captions are from zero-shot (held out) categories. Right: captions are from training set categories.

Reed et al., ICML 2016

Image-to-Image Translation with Conditional Adversarial Nets

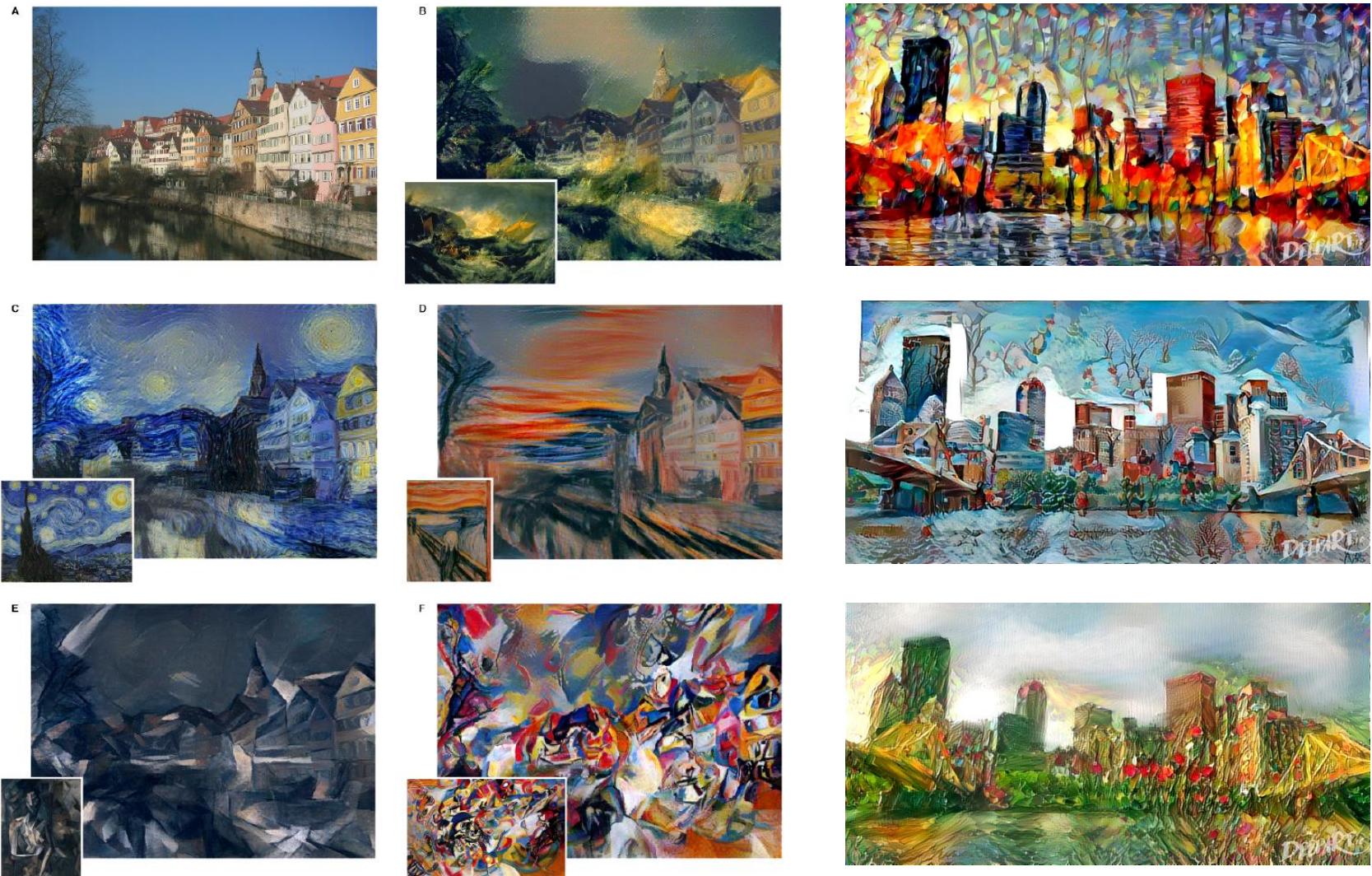


Scribbler: Controlling Deep Image Synthesis with Sketch and Color



Figure 1. A user can sketch and scribble colors to control deep image synthesis. On the left is an image generated from a hand drawn sketch. On the right several objects have been deleted from the sketch, a vase has been added, and the color of various scene elements has been constrained by sparse color strokes. For best resolution and additional results, see scribbler.eye.gatech.edu

Image Style Transfer Using Convolutional Neural Networks



DeepArt.io – try it for yourself!

Automatic Understanding of Image and Video Advertisements



Zaeem Hussain, Mingda Zhang, Xiaozhong Zhang, Keren Ye, Christopher Thomas,
Zuha Agha, Nathan Ong, Adriana Kovashka

University of Pittsburgh



Understanding advertisements is more challenging than simply recognizing physical content from images, as ads employ a variety of strategies to persuade viewers.



Symbolism



Atypical Objects



Culture/Memes

We collect an advertisement dataset containing 64,832 images and 3,477 videos, each annotated by 3-5 human workers from Amazon Mechanical Turk.

Image	Topic	204,340	Strategy	20,000
	Sentiment	102,340	Symbol	64,131
	Q+A Pair	202,090	Slogan	11,130
Video	Topic	17,345	Fun/Exciting	15,380
	Sentiment	17,345	English?	17,374
	Q+A Pair	17,345	Effective	16,721

Here are some sample annotations in our dataset.



New Caddy Maxi Life, infinitely bigger.



For the love of automobiles

What's being advertised in this image?

Cars, automobiles

What sentiments are provoked in the viewer?

Amused, Creative, Impressed, Youthful, Conscious

What strategies are used to persuade viewer?

Symbolism, Contrast, Straightforward, Transferred qualities

What should the viewer do, and why should they do this?

- I should buy Volkswagen because it can hold a big bear.
- I should buy VW SUV because it can fit anything and everything in it.
- I should buy this car because it can hold everything I need.

More information available at <http://cs.pitt.edu/~kovashka/ads>

Is computer vision solved?

- Given an image, we can guess with 96% accuracy what object categories are shown (ResNet)
- ... but we only answer “why” questions about images with 14% accuracy!

Why does it seem like it's solved?

- Deep learning makes excellent use of massive data (labeled for the task of interest?)
 - But it's hard to understand *how* it does so
 - It doesn't work well when massive data is not available and your task is different than tasks for which data is available
- Sometimes the manner in which deep methods work is not intellectually appealing, but our “smarter” / more complex methods perform worse

Seeing AI

[YouTube link](#)



Microsoft Cognitive Services: Introducing the Seeing AI project

Obstacles?

MASSACHUSETTS INSTITUTE OF TECHNOLOGY
PROJECT MAC

Artificial Intelligence Group
Vision Memo. No. 100.

July 7, 1966

THE SUMMER VISION PROJECT

Seymour Papert

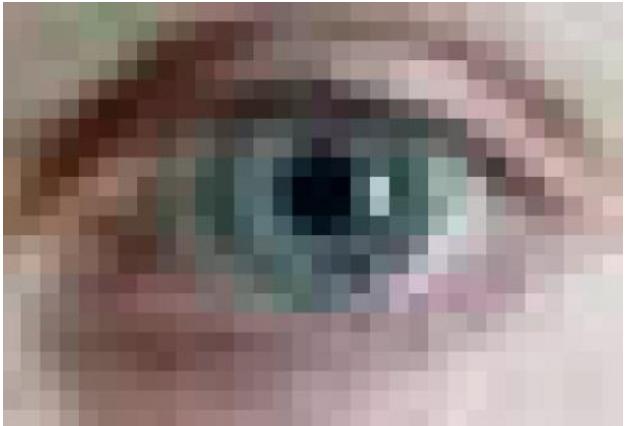
The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the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

Read more about the history: Szeliski Sec. 1.2

Why is vision difficult?

- Ill-posed problem: real world much more complex than what we can measure in images
 - $3D \rightarrow 2D$
- Impossible to literally “invert” image formation process with limited information
 - Need information outside of this particular image to generalize what image portrays (e.g. to resolve occlusion)

What the computer gets



153	156	148	152	149	147	139	146	142	150	146	144	137	125	120	119	136	146	151	164	172	175	183	188	196	200	205	208	214	214	219	217	
155	151	150	148	140	138	139	129	119	104	86	82	89	97	107	115	118	130	128	132	128	144	160	168	179	188	200	208	213	220	212	214	
149	146	153	147	147	146	132	99	73	78	96	95	105	126	138	151	145	157	163	171	165	161	146	126	157	184	190	201	215	212	214	214	
145	150	154	148	148	126	93	67	72	78	96	107	117	127	131	134	127	154	166	167	183	194	200	195	143	140	175	190	197	203	206	207	
151	153	151	147	120	85	67	75	84	83	94	92	81	78	98	91	83	117	126	144	178	200	201	203	206	175	127	159	185	196	195	206	
146	144	139	123	79	66	74	83	79	69	64	62	58	50	46	54	54	66	60	86	108	141	191	184	200	187	123	144	175	198	199		
135	130	115	87	64	77	90	79	78	85	81	63	55	57	56	53	70	62	61	68	59	58	84	105	168	194	196	183	131	151	185	197	
116	112	92	71	82	94	103	101	83	101	88	66	70	90	80	42	39	53	88	73	76	82	116	87	97	144	188	195	190	166	171	203	
135	120	84	83	108	127	135	115	100	92	79	49	85	74	59	0	0	0	0	50	69	52	79	157	141	100	84	136	187	206	204	189	200
144	103	91	115	139	147	127	91	87	80	72	44	61	84	25	0	0	0	50	181	45	69	142	164	167	113	93	130	193	199	208	203	
139	102	123	143	137	131	109	85	93	84	68	47	77	86	31	0	3	0	51	156	53	75	141	169	199	151	171	108	143	181	199	208	
141	135	153	142	114	104	97	93	98	77	42	77	96	79	21	0	23	58	46	56	77	155	199	212	161	194	193	164	187	202	205		
160	172	164	141	128	112	98	95	100	96	91	73	68	86	75	73	64	65	54	69	77	115	190	212	193	181	174	188	210	194	202	207	
179	189	160	140	139	116	97	97	100	103	110	99	75	80	72	83	50	55	54	95	98	174	205	185	179	188	185	190	193	217	224		
189	183	152	130	121	105	105	117	114	108	107	115	110	81	85	85	87	81	81	124	183	202	175	180	179	171	173	204	225	215	229		
178	161	149	135	120	115	122	129	137	145	131	121	125	115	109	91	92	111	132	159	173	170	184	176	184	190	191	217	210	226	228	233	
187	159	139	127	125	115	118	121	121	131	133	134	140	137	134	139	140	152	141	154	170	163	195	194	176	198	216	209	219	224	223	226	
185	164	140	122	116	110	109	108	113	118	115	116	123	127	135	148	154	162	165	170	171	160	183	198	201	210	223	216	221	222	221	226	
188	175	150	130	118	117	113	110	108	115	117	123	130	132	138	150	157	158	174	182	189	186	198	221	224	221	227	221	223	218	228		
187	179	158	141	124	127	125	127	126	129	130	135	139	141	150	165	175	172	185	195	207	210	212	226	229	222	224	224	223	218	219	221	
184	182	172	158	138	135	133	143	143	143	144	146	145	147	160	174	181	191	199	207	211	213	217	244	227	223	221	221	218	224	233		
183	181	187	174	153	148	136	140	147	145	148	157	162	160	158	165	174	181	188	201	210	212	216	228	234	226	226	215	217	215	204	204	

Why is this problematic?

Challenges: many nuisance parameters



Illumination



Object pose



Clutter



Occlusions



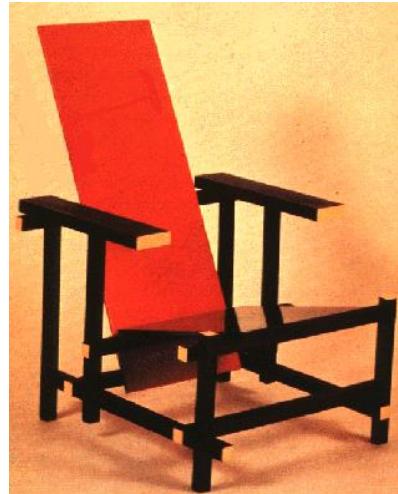
**Intra-class
appearance**



Viewpoint

Think again about the pixels...

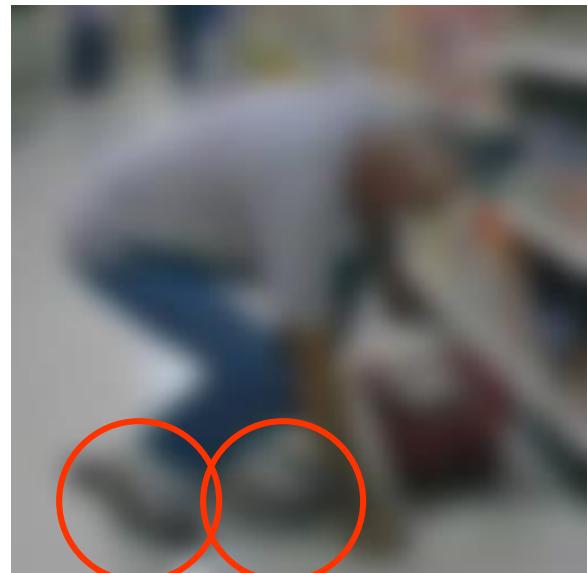
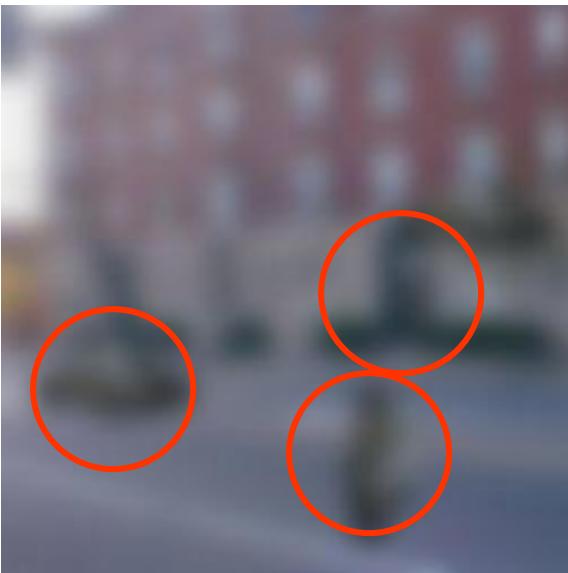
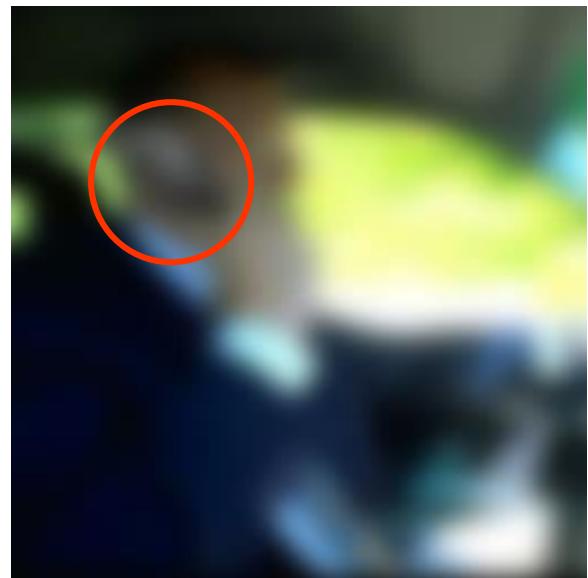
Challenges: intra-class variation



CMOA Pittsburgh



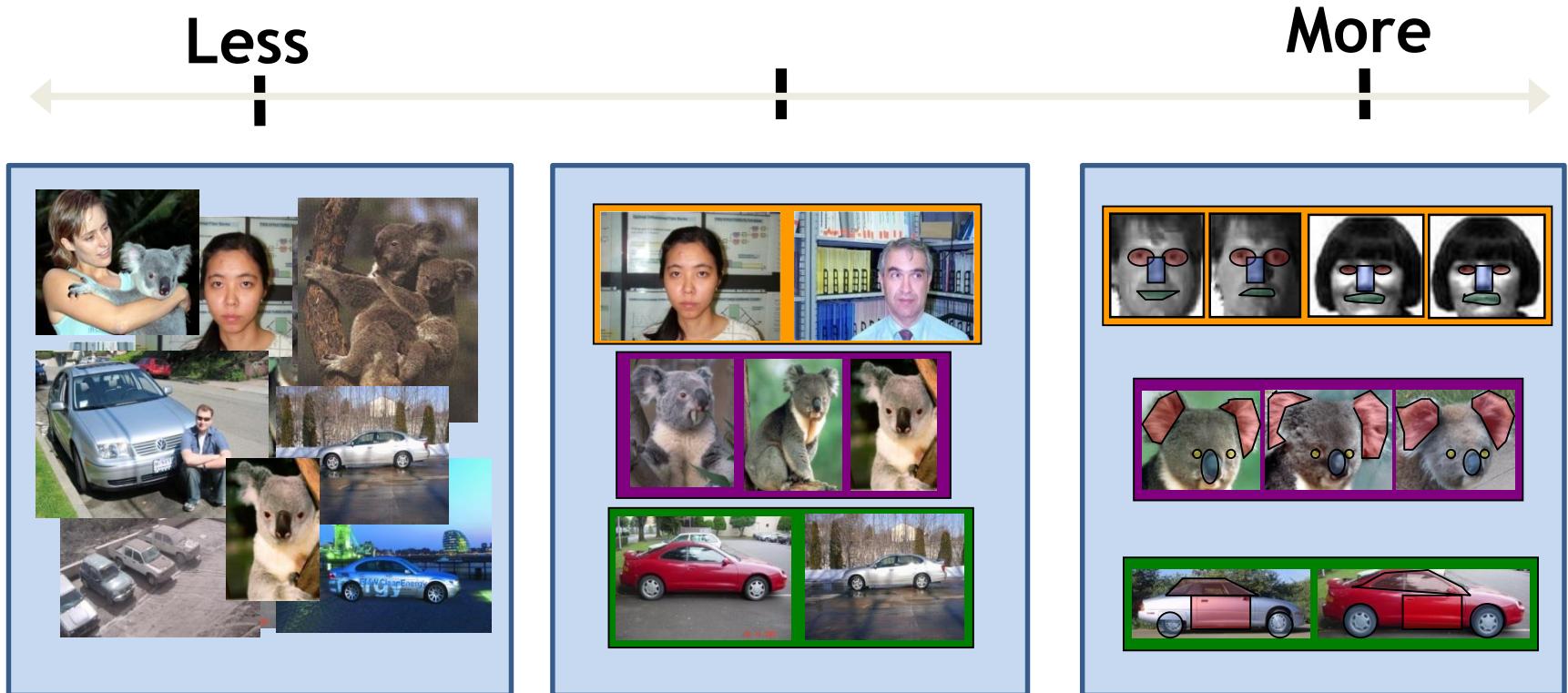
Challenges: importance of context



Challenges: Complexity

- Thousands to millions of pixels in an image
- 3,000-30,000 human recognizable object categories
- 30+ degrees of freedom in the pose of articulated objects (humans)
- Billions of images indexed by Google Image Search
- 1.424 billion smart camera phones sold in 2015
- About half of the cerebral cortex in primates is devoted to processing visual information [Felleman and van Essen 1991]

Challenges: Limited supervision



Unlabeled,
multiple objects

Classes labeled,
some clutter

Cropped to object,
parts and classes
labeled

Challenges: Vision requires reasoning



What color are her eyes?
What is the mustache made of?



How many slices of pizza are there?
Is this a vegetarian pizza?



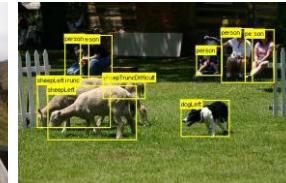
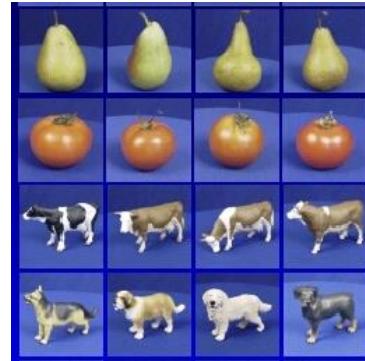
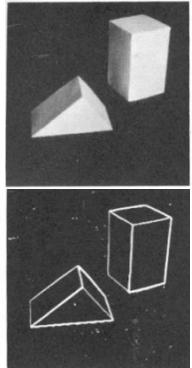
Is this person expecting company?
What is just under the tree?



Does it appear to be rainy?
Does this person have 20/20 vision?

Evolution of datasets

- Challenging problem → active research area



PASCAL:
20 categories, 12k images



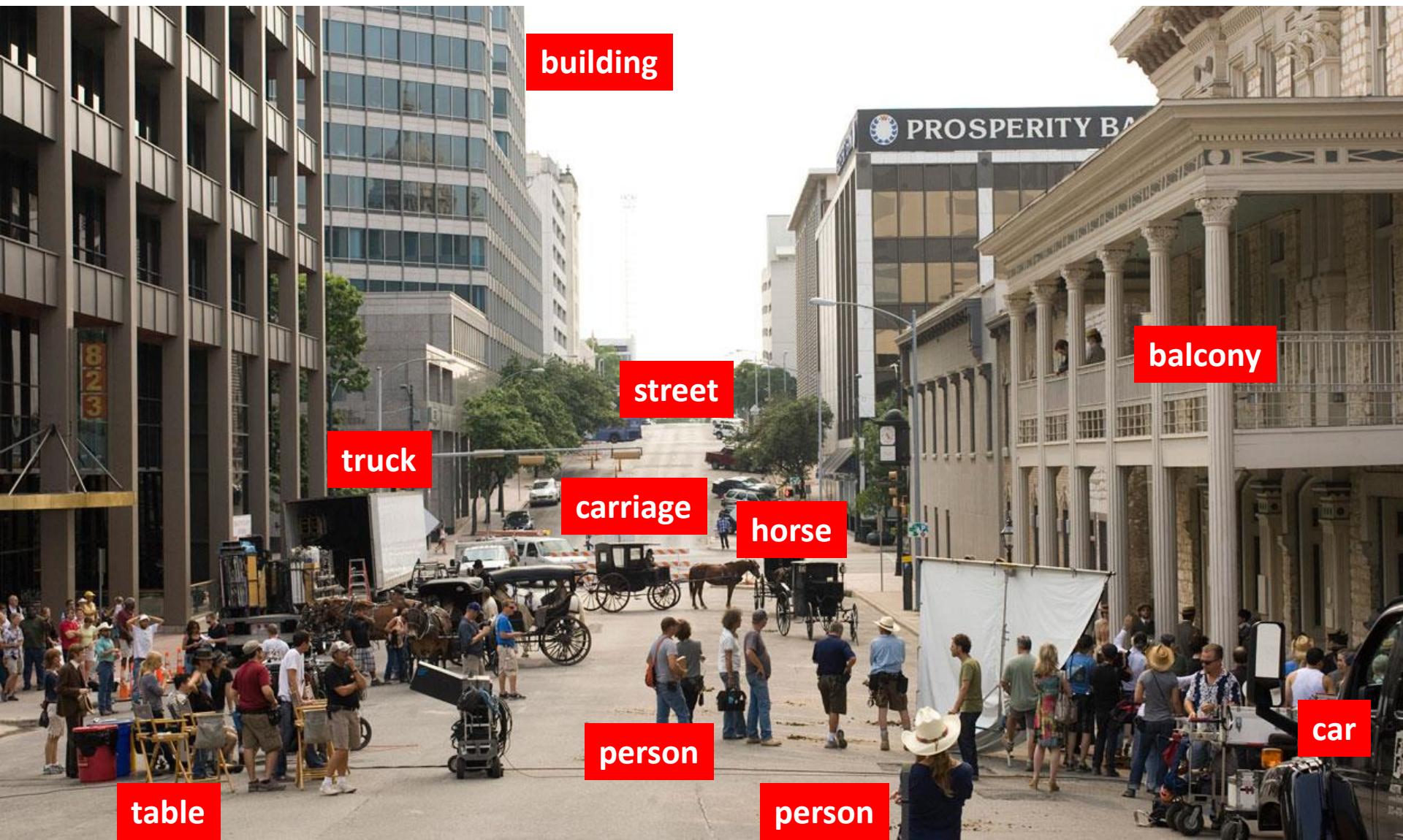
ImageNet:
22k categories, 14mil images

Microsoft COCO:
80 categories, 300k images

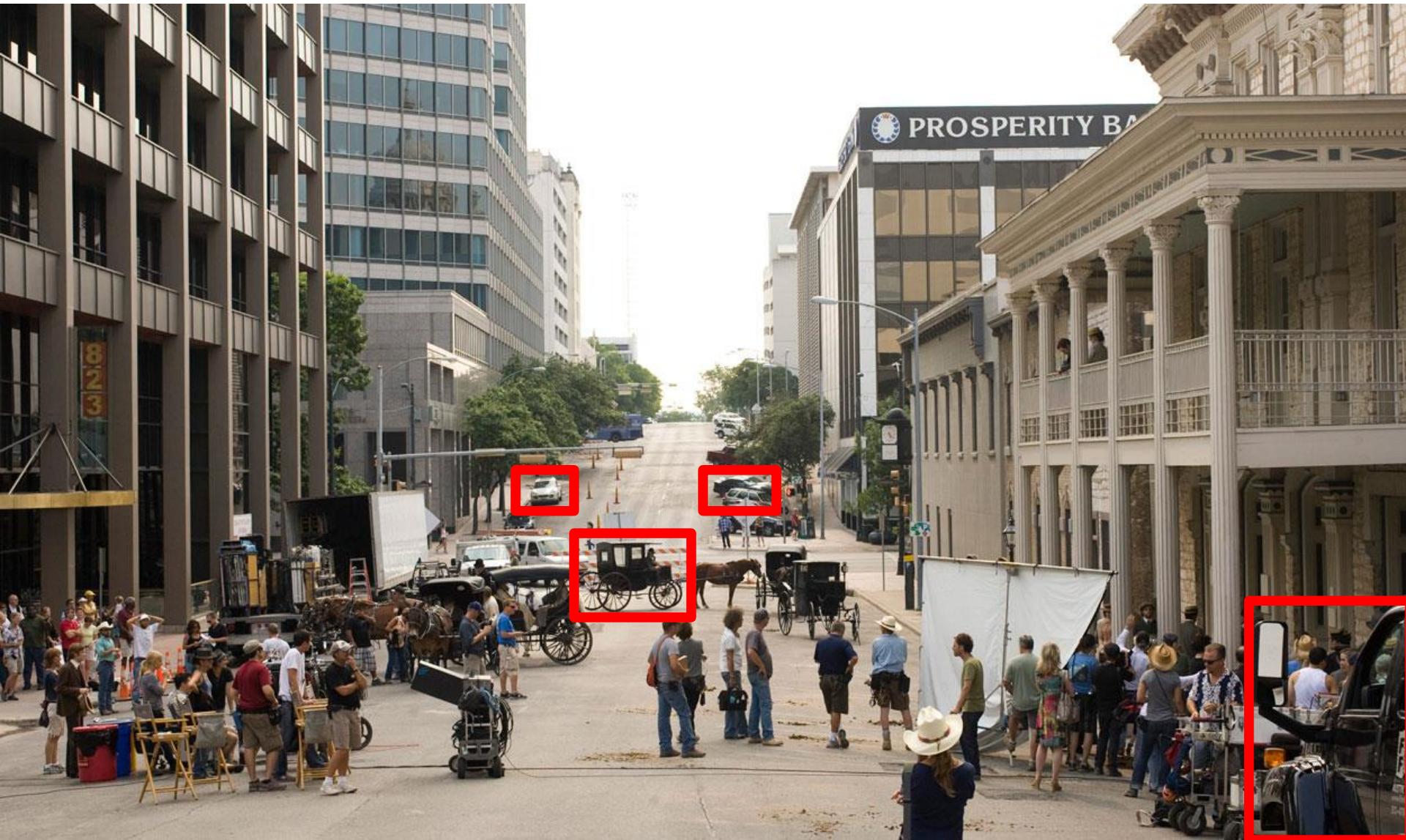
Some Visual Recognition Problems: Why are they challenging?



Recognition: What objects do you see?



Detection: Where are the cars?



Activity: What is this person doing?



Scene: Is this an indoor scene?



Instance: Which city? Which building?

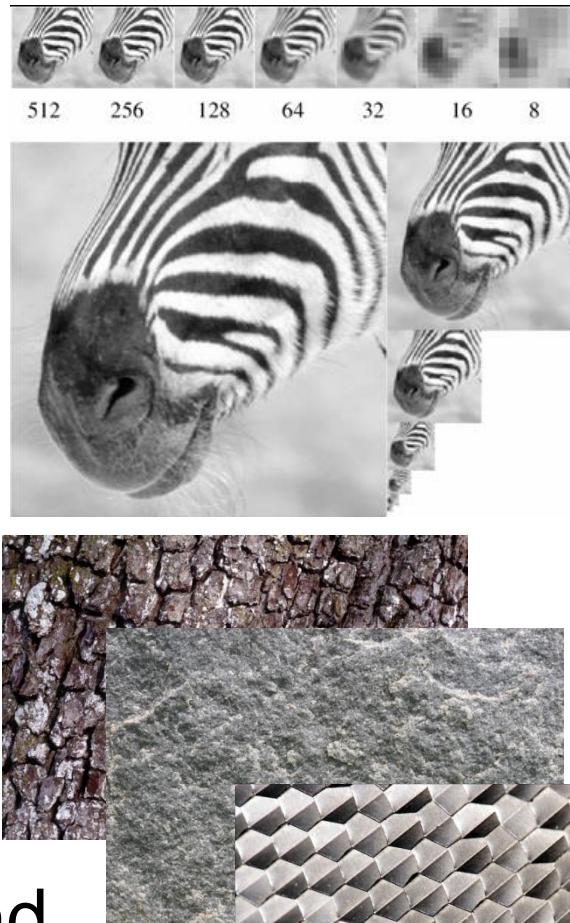
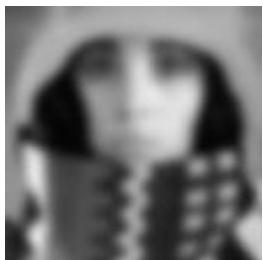


Visual question answering: Why is there a carriage in the street?



Overview of topics

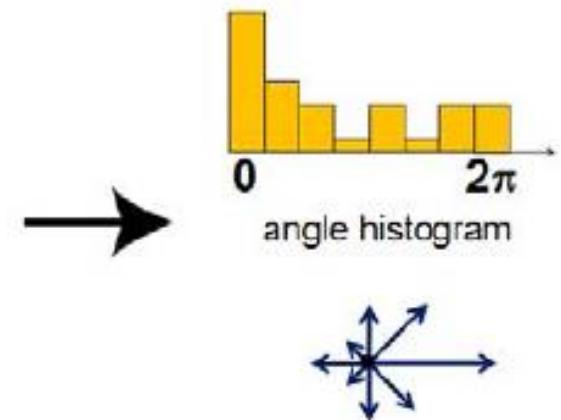
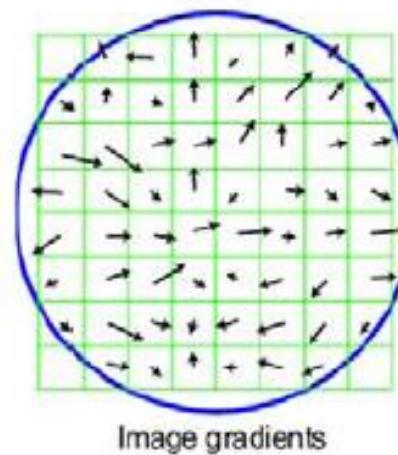
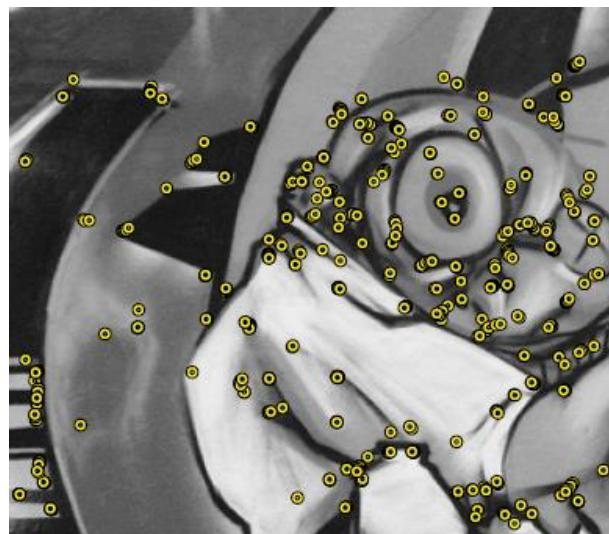
Features and filters



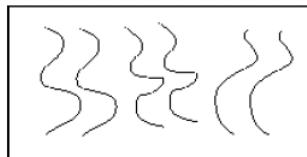
- Transforming and describing images; textures, colors, edges

Features and filters

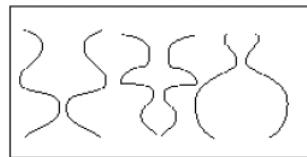
- Detecting distinctive + repeatable features
- Describing images with local statistics



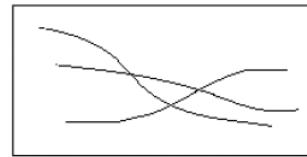
Grouping and fitting



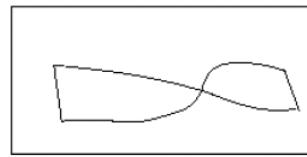
Parallelism



Symmetry



Continuity

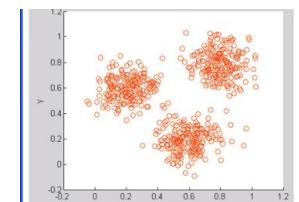
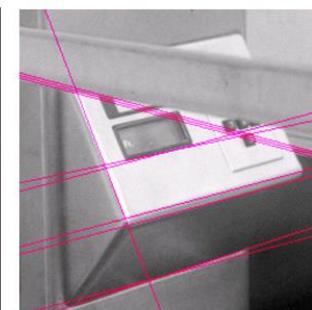
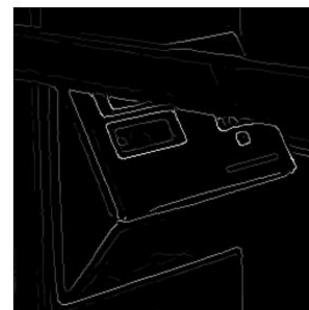


Closure



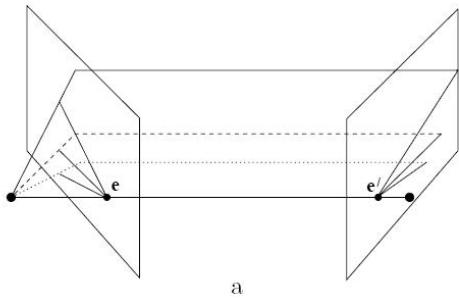
[fig from Shi et al]

- Clustering, segmentation, fitting; what parts belong together?



Multiple views

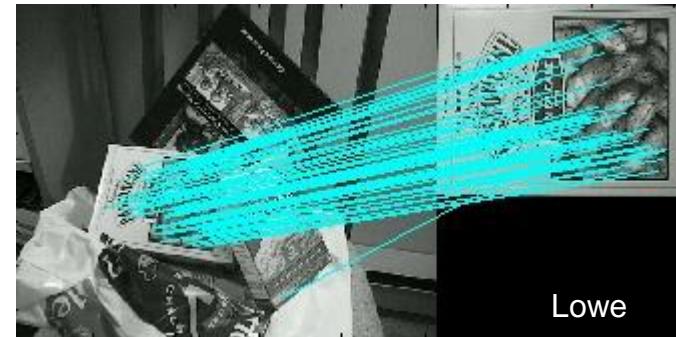
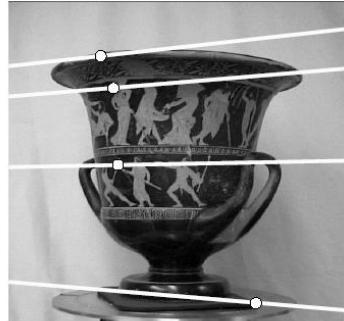
- Multi-view geometry, matching, invariant features, stereo vision



a



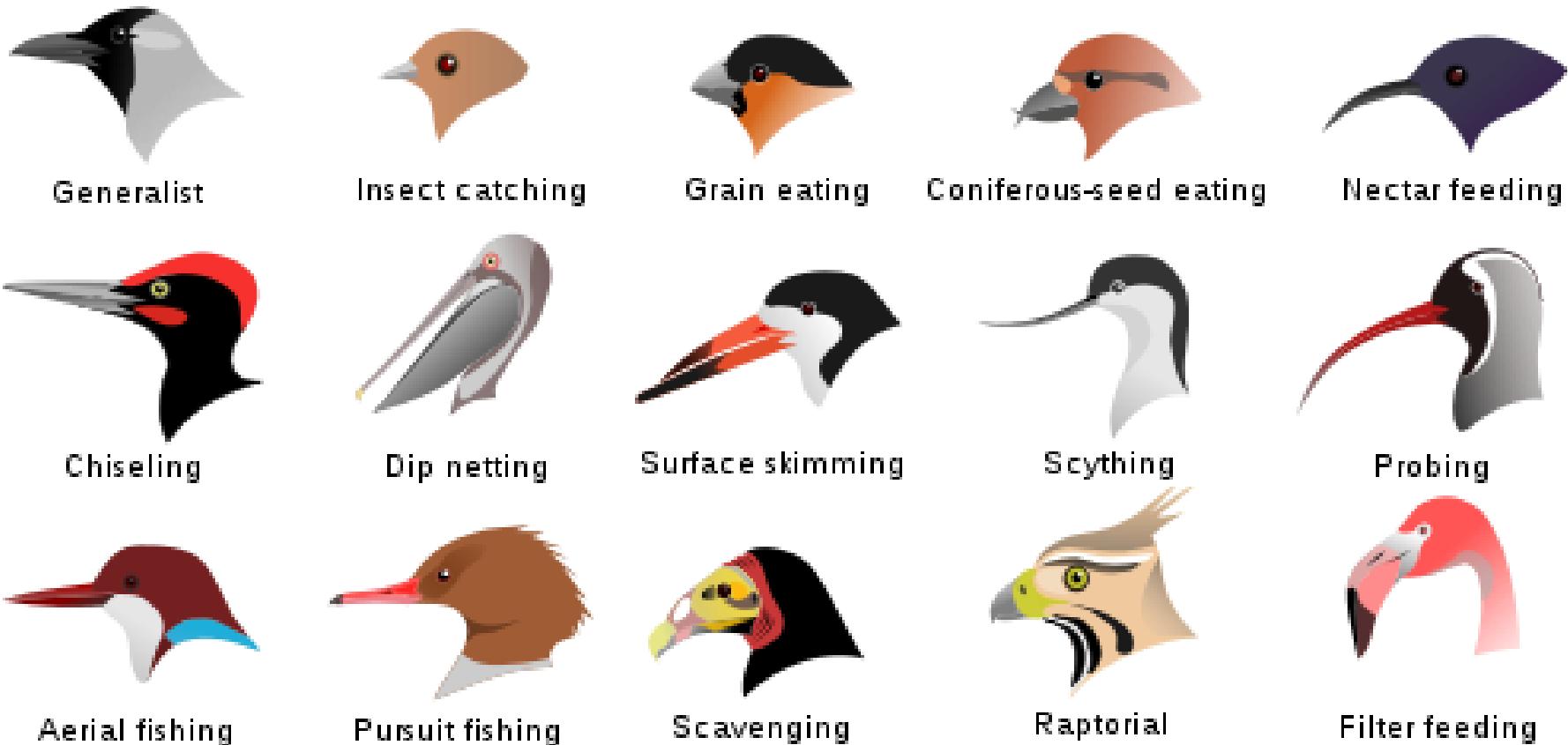
Hartley and Zisserman



Fei-Fei Li

Image categorization

- Fine-grained recognition



Region categorization

- Material recognition



[Bell et al. CVPR 2015]

Slide credit: D. Hoiem

Image categorization

- Image style recognition



HDR



Macro



Baroque



Rococo



Vintage



Noir



Northern Renaissance



Cubism



Minimal



Hazy



Impressionism



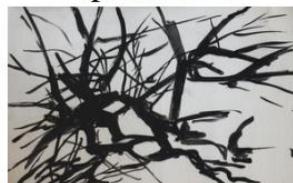
Post-Impressionism



Long Exposure



Romantic



Abs. Expressionism

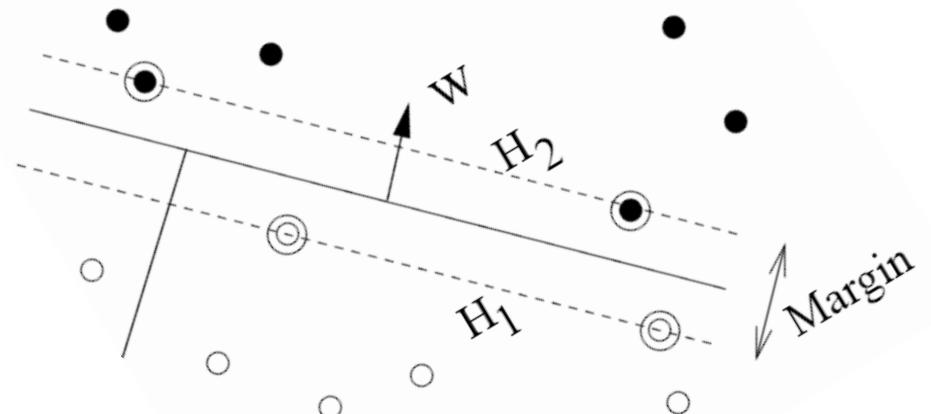
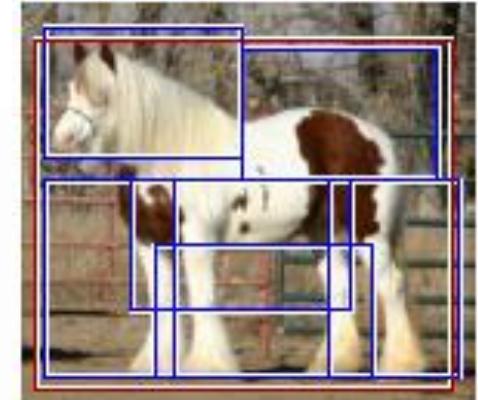
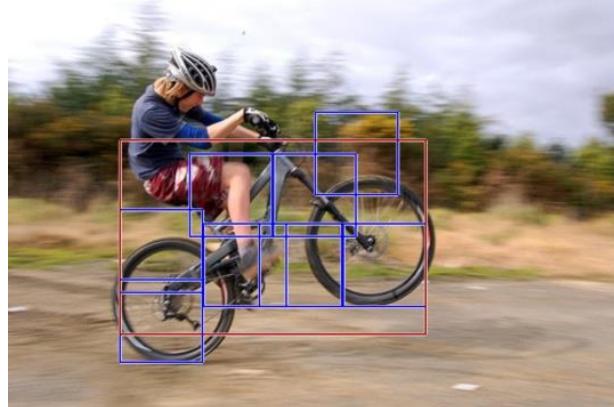
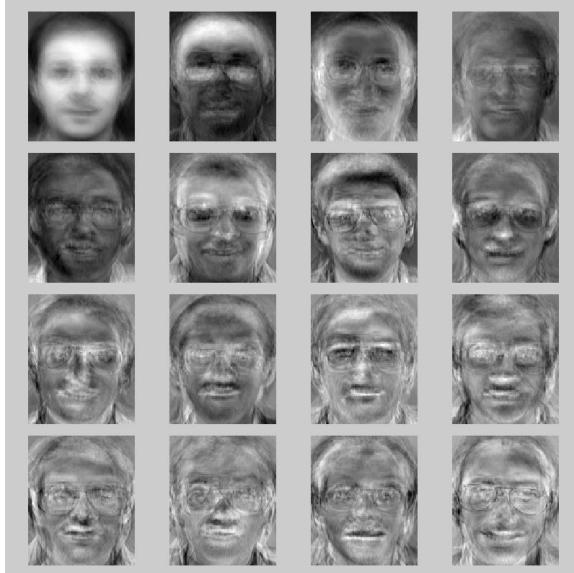


Color Field Painting

Flickr Style: 80K images covering 20 styles.

Wikipaintings: 85K images for 25 art genres.

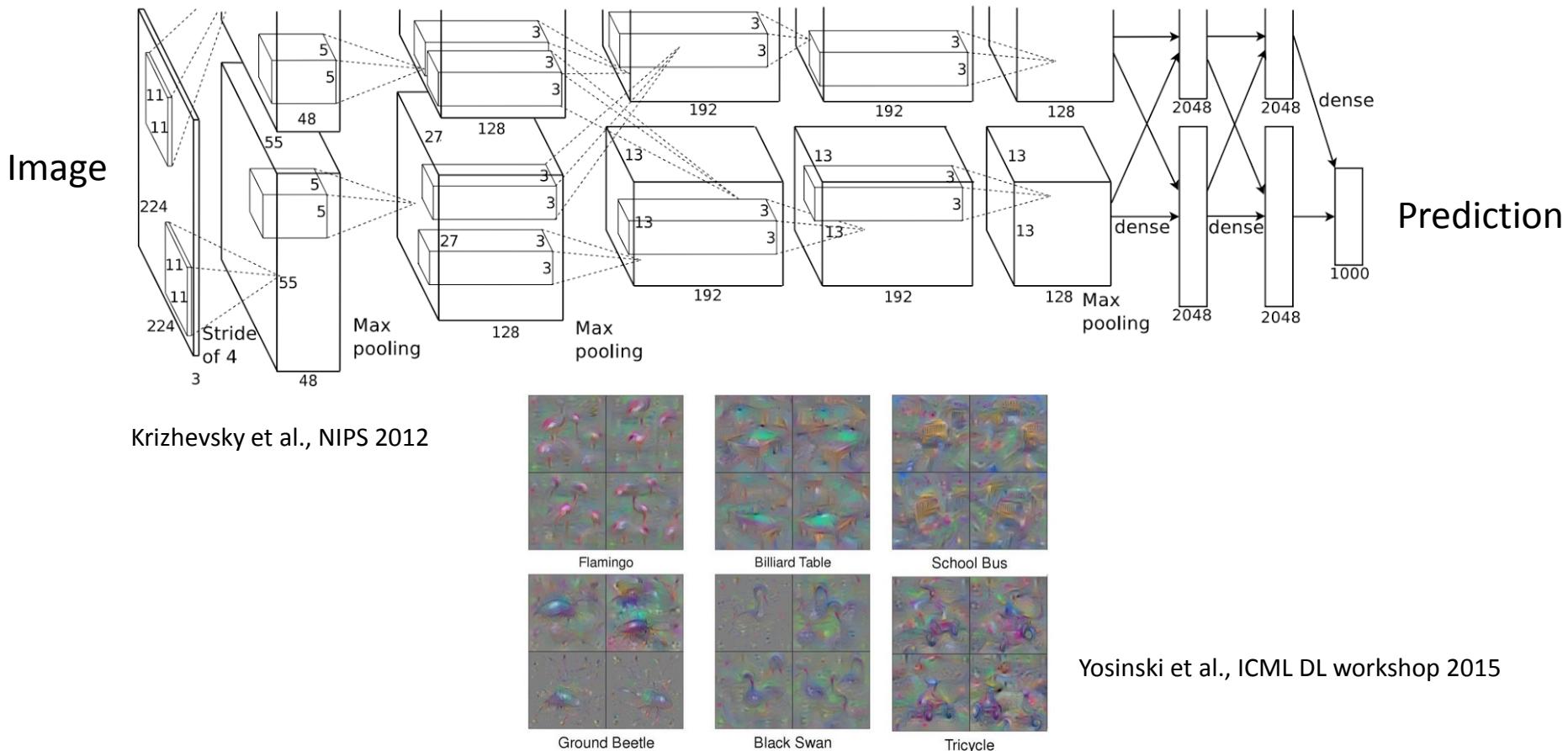
Visual recognition and SVMs



- Recognizing objects and categories, learning techniques

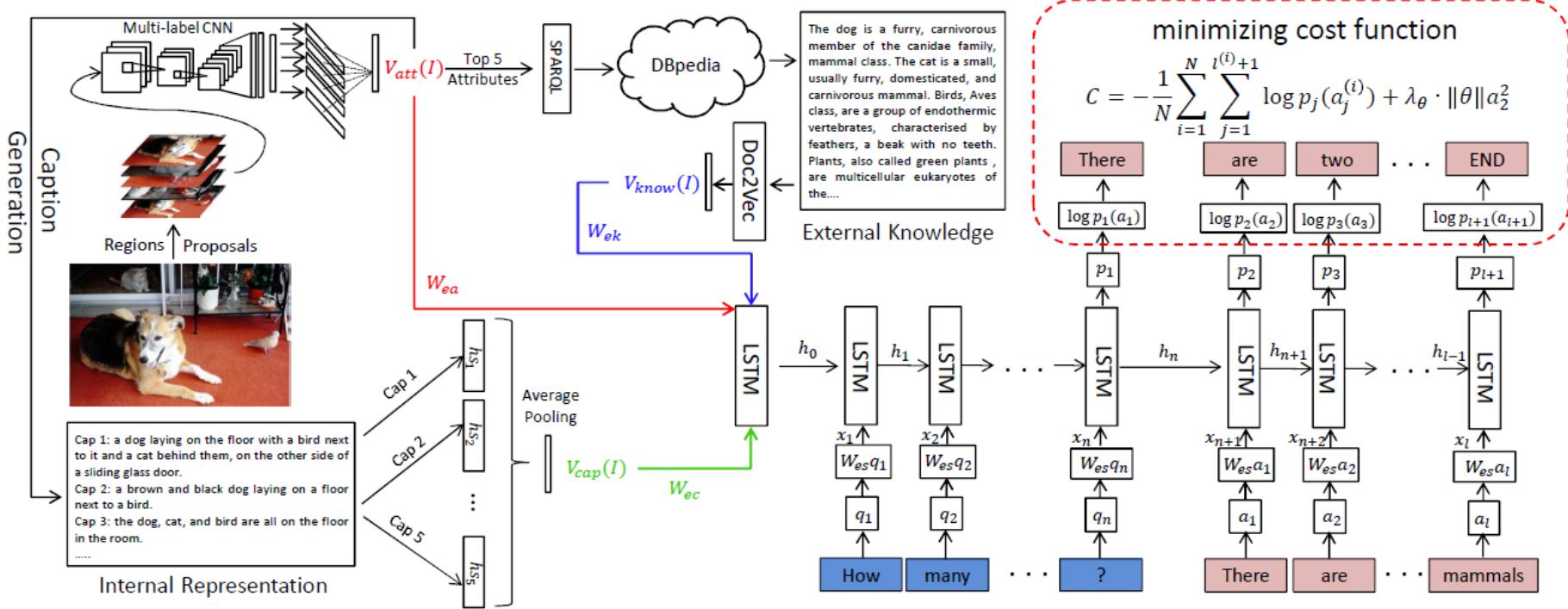
Convolutional neural networks (CNNs)

- State-of-the-art on many recognition tasks



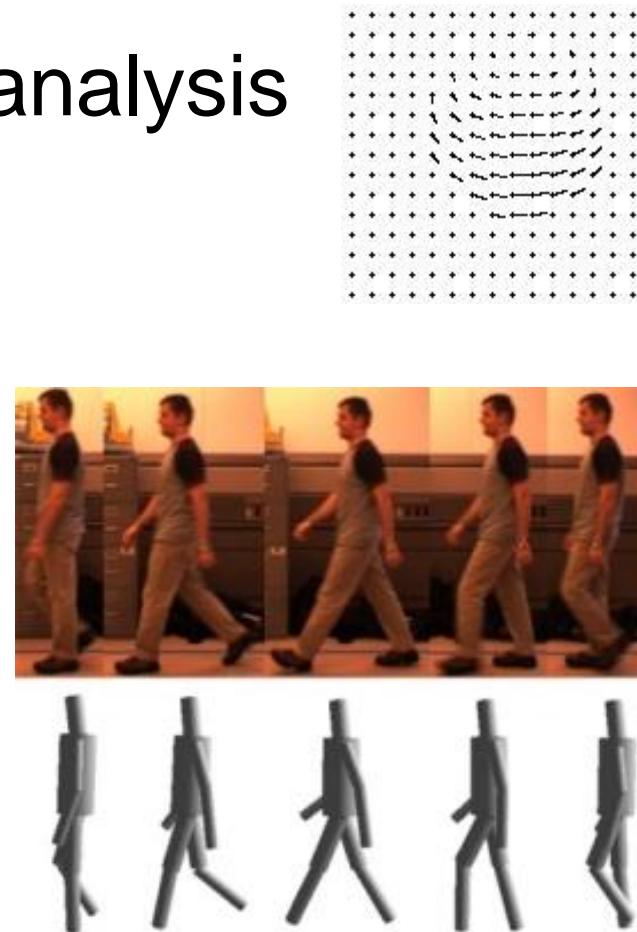
Recurrent neural networks

- Sequence processing, e.g. question answering



Motion and tracking

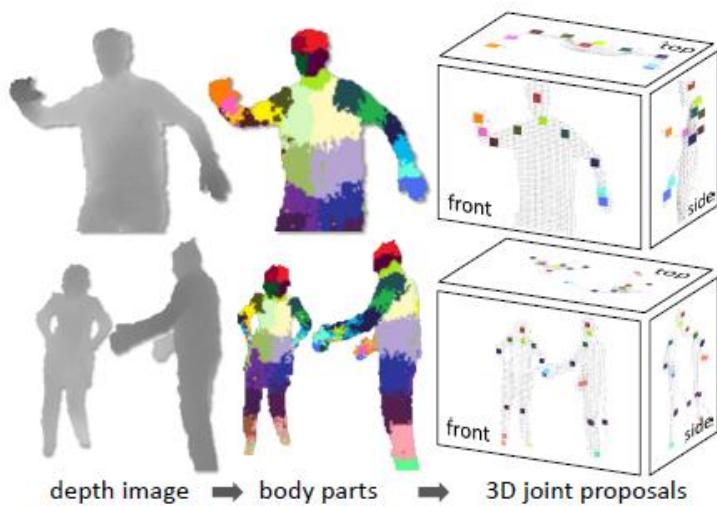
- Tracking objects, video analysis



Tomas Izo

Pose and actions

- Automatically annotating human pose (joints)
- Recognizing actions in first-person video



Your Homework

- Fill out Doodle
- Read entire course website
- Do first reading