# Computer Vision

Course webpage URL: http://luthuli.cs.uiuc.edu/~daf

And follow links

# Vision group at Illinois

**David Forsyth**
- Fulton Watson Copp chair; Marr prize, 1993; 2 ex students with Marr prizes; IEEE Tech. Achievement, Fellow; ACM Fellow; EIC IEEE TPAMI; 4 patents

**Jim Rehg**
- HCESC director; multiple famous ex-students, best paper awards; 26 patents

**Derek Hoiem**
- best paper, CVPR 2006; ACM Doctoral Dissertation honorable mention; Sloan Fellow; PAMI-TC Young Researcher; major startup

**Lana Lazebnik**
- Microsoft Faculty Fellow; Sloan Fellow; Koenderink Prize (2016); EIC IJCV

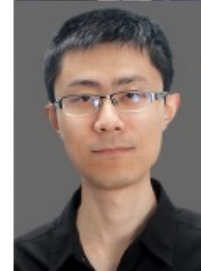**Alex Schwing**
- Visual learning, segmentation and GAN models

**Saurabh Gupta**
- Linking visual sensing to motion
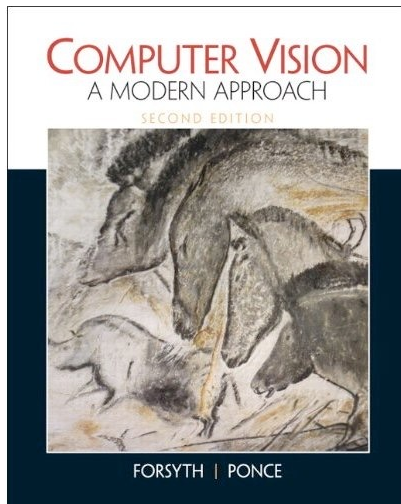
**Liangyan Gui**
- Understanding human movement

**Shenlong Wang**
- Simulation and sensing for autonomous vehicles

**Yuxiong Wang**
- Learning to detect and classify with very little data

# Vision group

**COMPUTER VISION**
A MODERN APPROACH
SECOND EDITION
FORSYTH | PONCE

The New Computer Vision

D.A. Forsyth

Likely about ~~202x~~ 2026

Cover design opportunity!

**Well-known ex-students:**

**Lana Lazebnik (UIUC)**

**Tamara Berg (UNC)**

**Pinar Duygulu (Hacettepe U.)**

**Ian Endres**

**Ali Farhadi (UW)**

**Varsha Hedau**

**Nazli Ikizler (Hacettepe U.)**

**Brett Jones**

**Kevin Karsch**

**Zicheng Liao**

**Deva Ramanan (CMU)**

**Raj Sodhi**

**Gang Wang (now Alibaba)**

**Amin Sadeghi**

**Zicheng Liao (Zhejiang U.)**

**Startups:**

**Lightform**

**Revery.ai**

**Reconstruct**

**Depix**

David Forsyth

Probability and Statistics for Computer Science
Springer

David Forsyth

Applied Machine Learning
Springer

## Top publications

| Categories | | English | |
| --- | --- | --- | --- |
| | Publication | h5-index | h5-median |
| 1. | Nature | 488 | 745 |
| 2. | IEEE/CVF Conference on Computer Vision and Pattern Recognition | 440 | 689 |
| 3. | The New England Journal of Medicine | 434 | 897 |
| 4. | Science | 409 | 633 |
| 5. | Nature Communications | 375 | 492 |
| 6. | The Lancet | 368 | 678 |
| 7. | Neural Information Processing Systems | 337 | 614 |
| 8. | Advanced Materials | 327 | 420 |
| 9. | Cell | 320 | 482 |
| 10. | International Conference on Learning Representations | 304 | 584 |
| 11. | JAMA | 298 | 498 |
| 12. | Science of The Total Environment | 297 | 436 |
| 13. | IEEE/CVF International Conference on Computer Vision | 291 | 484 |
| 14. | Angewandte Chemie International Edition | 281 | 361 |
| 15. | Nature Medicine | 274 | 474 |
| 16. | Journal of Cleaner Production | 272 | 359 |
| 17. | International Conference on Machine Learning | 268 | 424 |
| 18. | Chemical Reviews | 267 | 461 |
| 19. | Proceedings of the National Academy of Sciences | 267 | 405 |
| 20. | IEEE Access | 266 | 364 |

.com

Home / Best Conferences - Computer Science

# Best Computer Science Conferences

The ranking of leading conferences for Computer Science was prepared by Research.com, one of the primary websites for Computer Science research offering accurate data on scientific output since 2014.

The spot on the list is based on Impact Score data collected on 21-11-2023. It was based on a careful analysis of as much as 3,040 conference profiles and websites. Show more

| All resear | All publisl | All countr | Paper submission open |
|---|---|---|---|

| Rank | Conference Details | Impact Score |
|---|---|---|
| 1 Rank ◆IEEE | **Computer Vision and Pattern Recognition** 18-06-2023 - 22-06-2023 - **Vancouver** | 62.60 |
| 2 Rank ◆IEEE | **International Conference on Computer Vision** 11-10-2021 - 11-10-2021 - **Montreal** | 43.80 |
| 3 Rank | **Neural Information Processing Systems** 12-12-2023 - 14-12-2023 - **New Orleans** | 40.60 |
| 4 Rank Open Review .net | **International Conference on Learning Representations** 01-05-2023 - 05-05-2023 - **Kigali** | 38.10 |
| 5 Rank | **AAAI Conference on Artificial Intelligence** 07-02-2023 - 14-02-2023 - **Washington DC** | 35.00 |
| 6 Rank Karolinska Institutet | **International Conference on Machine Learning** 17-07-2022 - 23-07-2022 - **Baltimore** | 30.10 |
| 7 Rank Springer | **European Conference on Computer Vision** 29-09-2024 - 04-10-2024 - **Milan** | 29.60 |
| 8 Rank | **Meeting of the Association for Computational Linguistics** 22-05-2022 - 27-05-2022 - **Dublin** | 28.90 |
| 9 Rank | **Empirical Methods in Natural Language Processing** 06-12-2023 - 10-12-2023 - **Singapore City** | 21.90 |
| 10 Rank IJCAI | **International Joint Conference on Artificial Intelligence** 23-07-2022 - 29-07-2022 - **Vienna** | 21.10 |

# Outline

- Key tasks
- Why it is hard
- History of computer vision
- Current state of the art
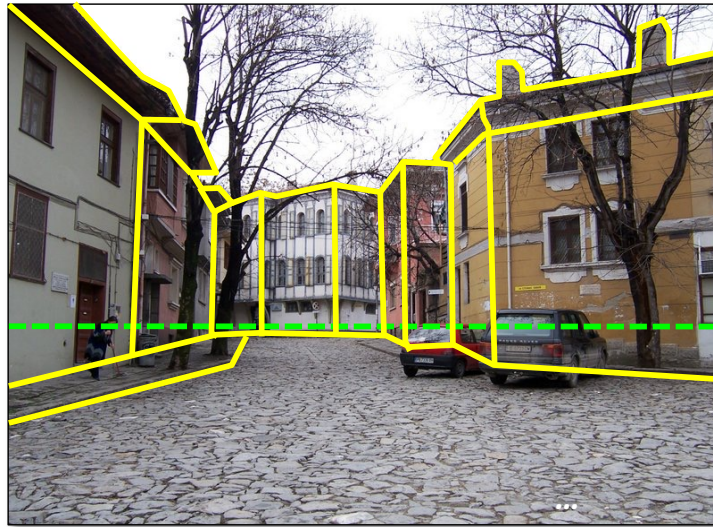- Topics covered in class

# What does vision do?

- *Classification: What is it?*

- *Localization: Where is it?*

- *Detection: Where and what?*

- *Tracking: Where is it going?*

- *Odometry: How have I moved?*

- *Navigation: Where am I?*

- *Modelling: What is the world like?*

- *Control: What should I do?*

- *Speculation: What will it be like if?*

# What kind of information can be extracted from an image?

# What kind of information can be extracted from an image?



**Geometric** information

# What kind of information can be extracted from an image?



**Geometric** information
**Semantic** information

# What kind of information can be extracted from an image?



**Geometric** information

**Semantic (?)** information – *affordances*

# What kind of information can be extracted from an image?



**Geometric** information
**Semantic** information
*Vision for action*

# Classification



Image

Some neural stuff;
differentiable wrt
parameters, input

Cat
Dog
·
·
·
Car

# Detection



FIGURE 18.8: *Faster RCNN uses two networks. One uses the image to compute "objectness" scores for a sampling of possible image boxes. The samples (called "anchor boxes") are each centered at a grid point. At each grid point, there are nine boxes (three scales, three aspect ratios). The second is a feature stack that computes a representation of the image suitable for classification. The boxes with highest objectness score are then cut from the feature map, standardized with ROI pooling, then passed to a classifier. Bounding box regression means that the relatively coarse sampling of locations, scales and aspect ratios does not weaken accuracy.*

# Detection and localization in 2D



YOLOv11 documentation by Ultralytics

# Localization in 3D from detection



Wu et al, *6D-VNet: End-to-end 6DoF Vehicle Pose Estimation from Monocular RGB Images*

# Lane detection

US 9081385

Waymo and Google  2012

Strategy:  detect markers (reflective paint), join up

exercise in robust fitting of curves



FIGURE 4A

# Tracking Anything Results



**Cheng et al.; DEVA: Tracking anything with decoupled video segmentation; 2023**

# Tracking Anything Results



**Cheng et al.; DEVA: Tracking anything with decoupled video segmentation; 2023**

# Tracking Anything Results



**Cheng et al.; DEVA: Tracking anything with decoupled video segmentation; 2023**

# Visual odometry



https://github.com/MAC-VO/MAC-VO/blob/main/asset/ICRAvideo.gif

# Extreme odometry

https://www.youtube.com/watch?v=fBiataDpGIo

# Goal: To extract useful information from pixels



| What we see | What a computer sees |
|---|---|

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| 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 |

# Trad. View: Images are fundamentally ambiguous!

# Humans are remarkably good at vision…

# …still, vision is hard even for humans

# …still, vision is hard even for humans



Figure from Marr (1982), attributed to R. C. James

# Outline

- Logistics, requirements
- Goal of computer vision and why it is hard
- History of computer vision

How it started



(a) Original picture.    (b) Differentiated picture.
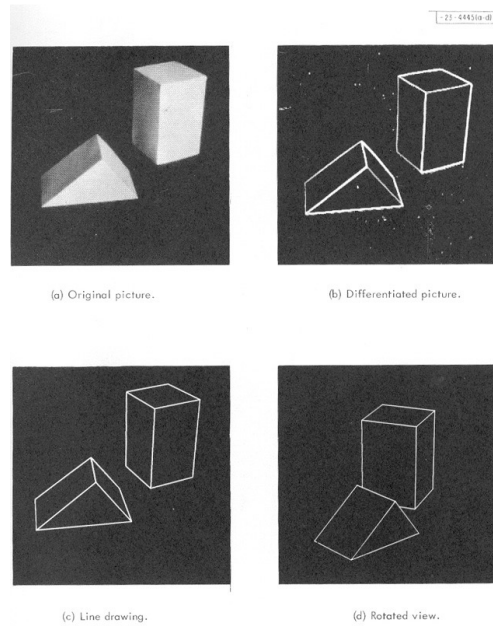
(c) Line drawing.    (d) Rotated view.

L. G. Roberts, 1963

How it's going



OpenAI DALL-E, 2020

# Origins



Hough, 1959



(a) Original picture.    (b) Differentiated picture.

(c) Line drawing.    (d) Rotated view.

Roberts, 1963



PICTURE PROCESSING BY COMPUTER

AZRIEL ROSENFELD

COMPUTER SCIENCE AND APPLIED MATHEMATICS

Rosenfeld, 1969



Pattern Classification and Scene Analysis

Richard O. Duda and Peter E. Hart

Duda & Hart, 1972

# Decade by decade

- **1960s**: Blocks world, image processing and pattern recognition
- **1970s**: Key recovery problems defined: structure from motion, stereo, shape from shading, color constancy. Attempts at knowledge-based recognition
- **1980s**: Fundamental and essential matrix, multi-scale analysis, corner and edge detection, optical flow, geometric recognition as alignment
- **1990s**: Multi-view geometry, statistical and appearance-based models for recognition, first approaches for (class-specific) object detection
- **2000s**: Local features, generic object recognition and detection
- **2010s**: Deep learning, big data

- For much more detail: see Prof Lazebnik's   historical overview

Adapted from J. Malik

# Growth of the field: CVPR papers



More *accepted* papers in 2022
than *submitted* papers in 2012!

Source: CVPR 2022 opening sides

# Growth of the field: CVPR attendance



Source: CVPR 2022 opening sides

# Introduction: Outline

- Logistics, requirements
- Goal of computer vision and why it is hard
- History of computer vision
- **Current(ish) state of the art**
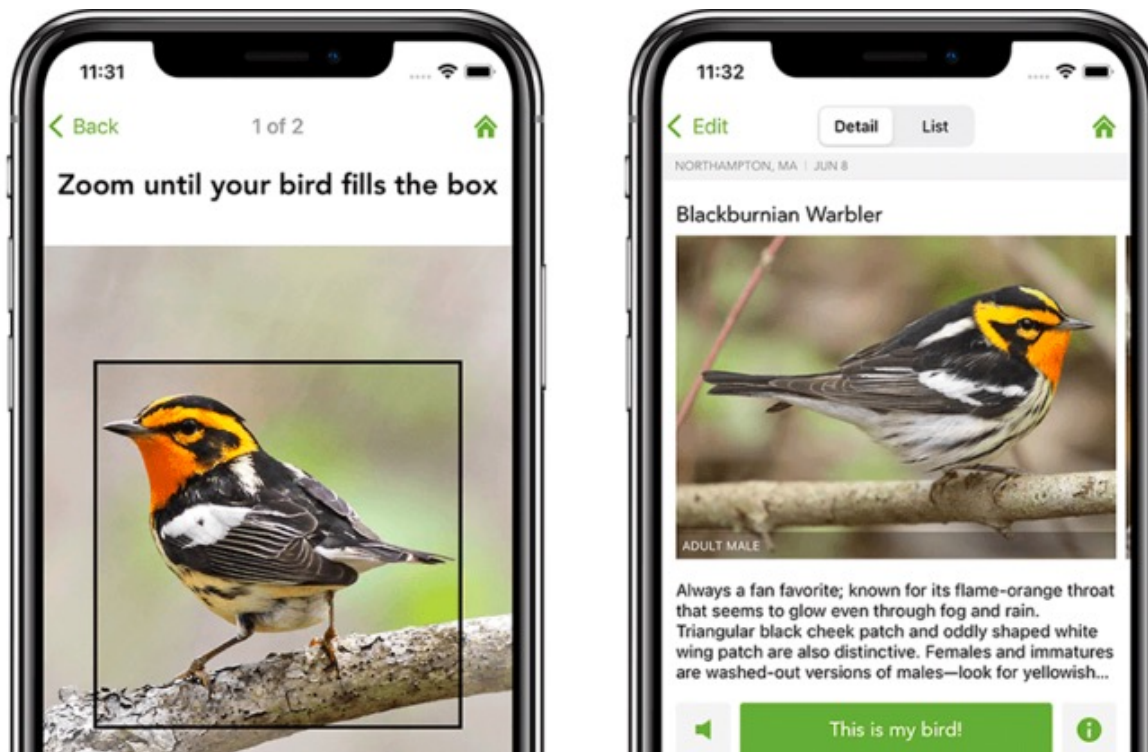
# What can computer vision do today?



In the 60s, Marvin Minsky assigned a couple of undergrads to spend the summer programming a computer to use a camera to identify objects in a scene. He figured they'd have the problem solved by the end of the summer. Half a century later, we're still working on it.

https://xkcd.com/1425/

(September 24, 2014)

# What can computer vision do today?

- It's 2025 now…



https://merlin.allaboutbirds.org/

# What can computer vision do today?

- It's 2025 now…

# What can computer vision do today?

- Reconstruction
- Recognition
- *Reconstruction meets recognition, or 3D scene understanding*
- *Image generation*
- *Vision for action*

# Regression

- We must make image-like things from images
- Examples:
  - depth map from image
  - normal map from image
  - derained image from rainy image
  - defogged image from foggy image
- Train with pairs (image, depth)
  - or (image, normal), etc
  - Loss
    - Squared error +abs value of error+other terms as required

**Depth (omnimap, current best depth est)**  **Normal (omnimap, current best normal est)**

# Correspondence yields 3D configuration



Camera 1

Camera 2

# Reconstruction: 3D from photo collections



Colosseum, Rome, Italy          San Marco Square, Venice, Italy

Q. Shan, R. Adams, B. Curless, Y. Furukawa, and S. Seitz, The Visual Turing Test for Scene Reconstruction, 3DV 2013

YouTube Video

# Reconstruction: 4D from photo collections



**Figure 1:** *We mine Internet photo collections to generate time-lapse videos of locations all over the world. Our time-lapses visualize a multitude of changes, like the retreat of the Briksdalsbreen Glacier in Norway shown above. The continuous time-lapse (bottom) is computed from hundreds of Internet photos (samples on top).* Photo credits: `Aliento Más Allá, jirihnidek, mcxurxo, elka-cz,` Juan Jesús Orío, Klaus Wißkirchen, `Daikrieg, Free the image, dration` and Nadav Tobias.

R. Martin-Brualla, D. Gallup, and S. Seitz, Time-Lapse Mining from Internet Photos, SIGGRAPH 2015

YouTube Video

# Reconstruction: 4D from depth cameras



Figure 1: Real-time reconstructions of a moving scene with DynamicFusion; both the person and the camera are moving. The initially noisy and incomplete model is progressively denoised and completed over time (left to right).

R. Newcombe, D. Fox, and S. Seitz, DynamicFusion:
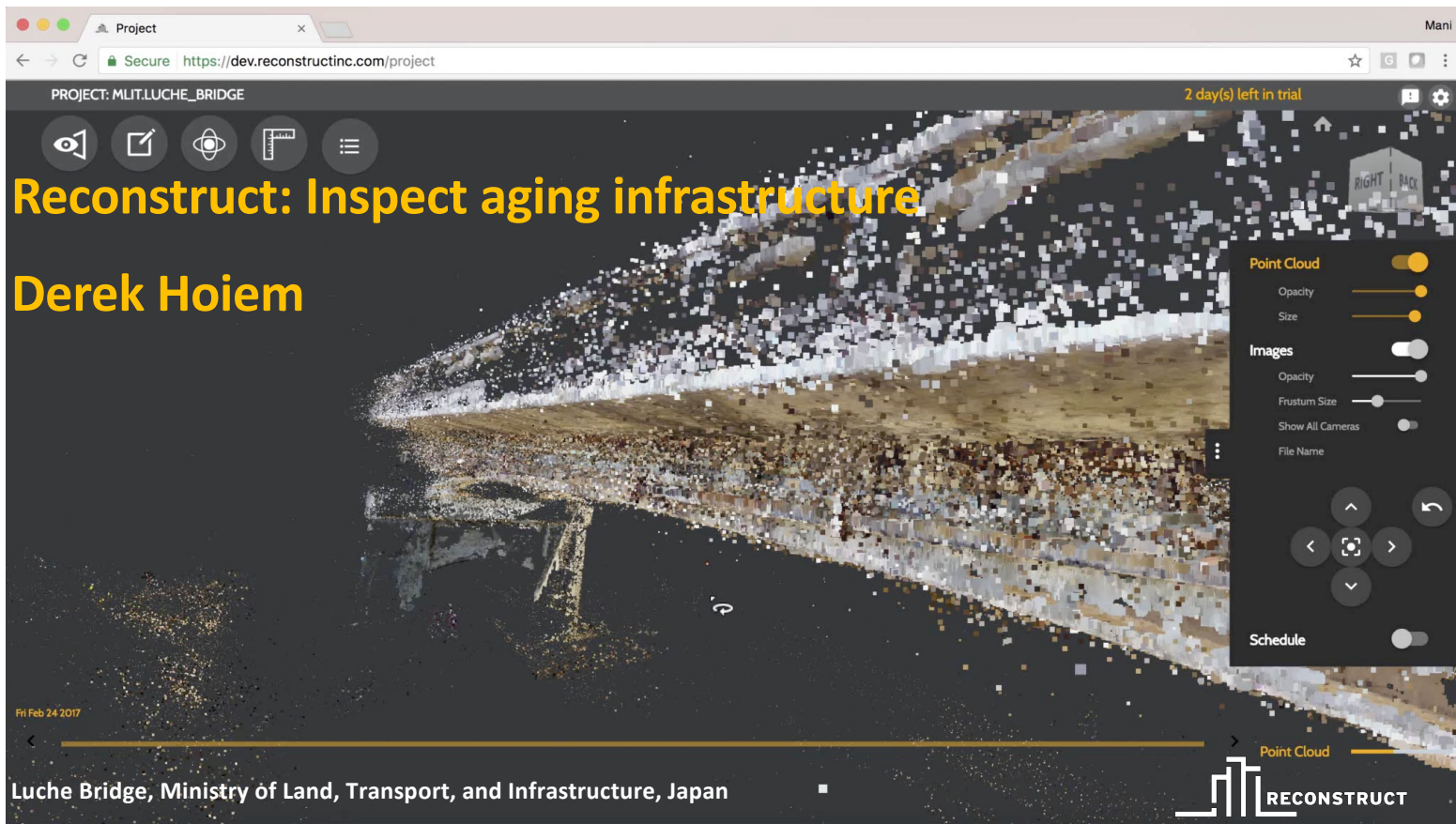Reconstruction and Tracking of Non-rigid Scenes in Real-Time,
CVPR 2015

YouTube Video

# Reconstruction: Commercial applications



https://www.zillow.com/z/3d-home/

**Reconstruct: Inspect aging infrastructure**

**Derek Hoiem**

Luche Bridge, Ministry of Land, Transport, and Infrastructure, Japan

Reconstruct: Align reality to plans for construction management

Derek Hoiem

# Reconstruction: Commercial applications



RECONSTRUCT INTEGRATES REALITY AND PLAN

**Visual Asset Management**

Reconstruct 4D point clouds and organize images and videos from smartphones, time-lapse cameras, and drones around the project schedule. View, annotate, and share anywhere with a web interface.

**4D Visual Production Models**

Integrate 4D point clouds with 4D BIM, review "who does what work at what location" on a daily basis and improve coordination and communication among project teams.

**Predictive Visual Data Analytics**

Analyze actual progress deviations by comparing Reality and Plan and predict risk with respect to the execution of the look-ahead schedule for each project location, to offer your project team with an opportunity to tap off potential delays before they surface on your jobsite.

**reconstructinc.com**

Source: D. Hoiem

# Recognition: "Simple" patterns

# Recognition: Faces

# Recognition: Faces





[How China Uses High-Tech Surveillance to Subdue Minorities](#) – New York Times, 5/22/2019

[The Secretive Company That Might End Privacy As We Know It](#) – New York Times, 1/18/2020
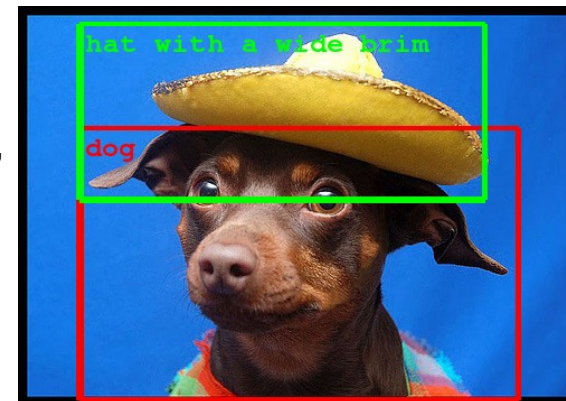
[Wrongfully Accused by an Algorithm](#) – New York Times, 6/24/2020

[Facial Recognition Goes to War](#) – New York Times, 4/7/2022

# 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

# Recognition: General categories

# Object detection, instance segmentation



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

# 3D scene understanding: NERFs



$$(x, y, z, \theta, \phi) \rightarrow \boxed{|||} \rightarrow (RGB\sigma)$$
$$F_\Theta$$

B. Mildenhall et al., Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV 2020

# 3D scene understanding: NERFs



B. Mildenhall et al., Representing Scenes as Neural Radiance Fields for View Synthesis, ECCV 2020

# 3D scene understanding: Single-view reconstruction



Figure 2: Our model (MarrNet) has three major components: (a) 2.5D sketch estimation, (b) 3D shape estimation, and (c) a loss function for reprojection consistency. MarrNet first recovers object normal, depth, and silhouette images from an RGB image. It then regresses the 3D shape from the 2.5D sketches. In both steps, it uses an encoding-decoding network. It finally employs a reprojection consistency loss to ensure the estimated 3D shape aligns with the 2.5D sketches. The entire framework can be trained end-to-end.

J. Wu, Y. Wang, T. Xue, X. Sun, W. Freeman, J. Tenenbaum, MarrNet: 3D Shape Reconstruction via 2.5D Sketches, NeurIPS 2017

# Image generation: Faces

- 1024x1024 resolution, CelebA-HQ dataset



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

Follow-up work

# Image generation: DeepFakes



**Harrison Ford Is Young Han In Solo Deepfake Video**

Thanks to deepfake technology, the maligned Solo: A Star Wars Story now stars Harrison Ford instead of Alden Ehrenreich as the young Han.

BY DAN ZINSKI
2 DAYS AGO

Just a random recent example…

https://screenrant.com/star-wars-han-solo-movie-harrison-ford-video-deepfake/
https://www.youtube.com/watch?v=bC3uH4Xw4Xo

**https://en.wikipedia.org/wiki/Deepfake**

# Image generation: OpenAI DALL-E, DALL-E 2



vibrant portrait painting of Salvador Dalí with a robotic half face

a shiba inu wearing a beret and black turtleneck

a close up of a handpalm with leaves growing from it

an espresso machine that makes coffee from human souls, artstation

panda mad scientist mixing sparkling chemicals, artstation

a corgi's head depicted as an explosion of a nebula

A. Ramesh et al., Zero-Shot Text-to-Image Generation, ICML 2021. https://openai.com/blog/dall-e/

A. Ramesh et al., Hierarchical Text-Conditional Image Generation with CLIP Latents, arXiv 2022. https://openai.com/dall-e-2/

# Vision for action: Visuomotor learning



**Overview video**,

**training video**



S. Levine, C. Finn, T. Darrell, P. Abbeel, End-to-end training of deep visuomotor policies, JMLR 2016

# Does computer vision matter for action?



"Our main finding is that computer vision does matter. Models equipped with intermediate representations train faster, achieve higher task performance, and generalize better to previously unseen environments."

B. Zhou, P. Krähenbühl, and V. Koltun, Does Computer Vision Matter for Action? Science Robotics, 4(30), 2019 (video)

# Vision for action: Learning skills from video



Fig. 1. Simulated characters performing highly dynamic skills learned by imitating video clips of human demonstrations. **Left:** Humanoid performing cartwheel B on irregular terrain. **Right:** Backflip A retargeted to a simulated Atlas robot.

**Video**

X. B. Peng, A. Kanazawa, J. Malik, P. Abbeel, S. Levine, SFV: Reinforcement Learning of Physical Skills from Videos, SIGGRAPH Asia 2018

# Outline

- Logistics, requirements
- Goal of computer vision and why it is hard
- History of computer vision
- Current state of the art
- Topics covered in class

# Topics covered in class

I. Elementary Image Representations:

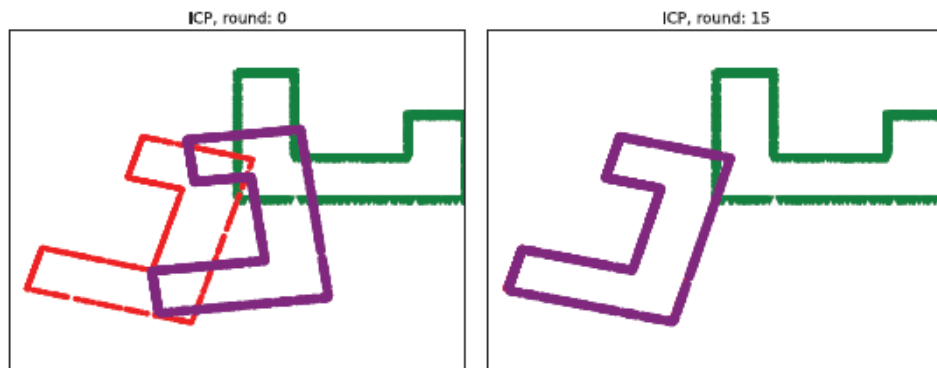Point transformations, geometric transformations, filters, denoising, edges, interest points

II. Mid-level vision:

Voting, Fitting, Registration

III. Learned Image Representations:

Learned denoising, Mapping images to images, Classification, Detection

IV. Image formation and geometric vision

Cameras, Light, Color, Calibration

V. Pairs of Cameras and more:

Geometry, Odometry, Optic Flow, Stereopsis, Structure from Motion, Tracking

# Elementary image representations



Basic image processing



Upsampling and downsampling

downsampled by 8



Linear filtering
Edge detection



Feature extraction

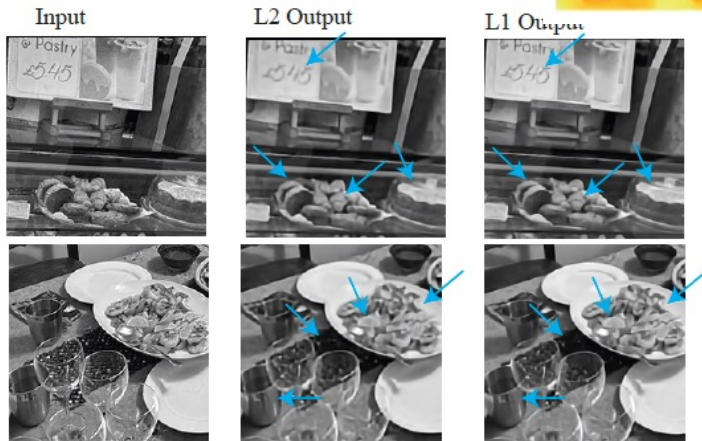# Mid-level Vision


Voting


Fitting
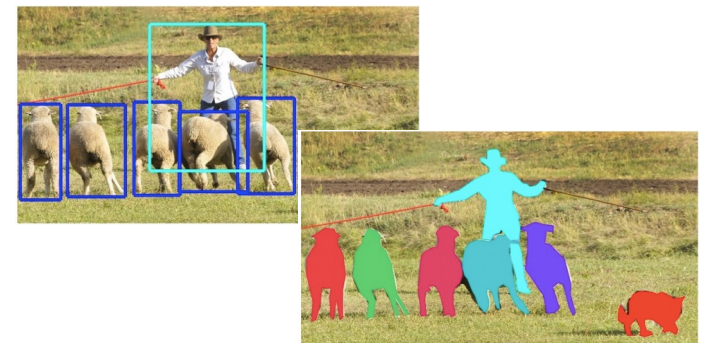

Registration


Mosaics

# Learned Image Representations
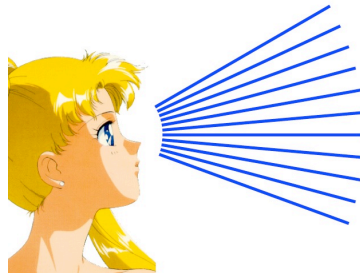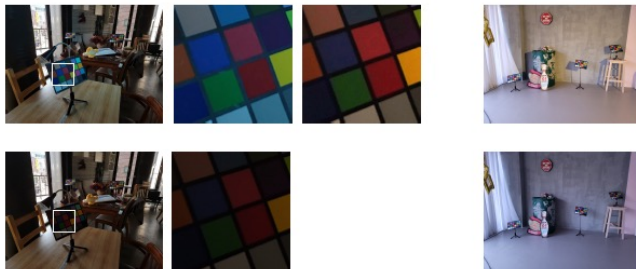


Depth, normal, etc

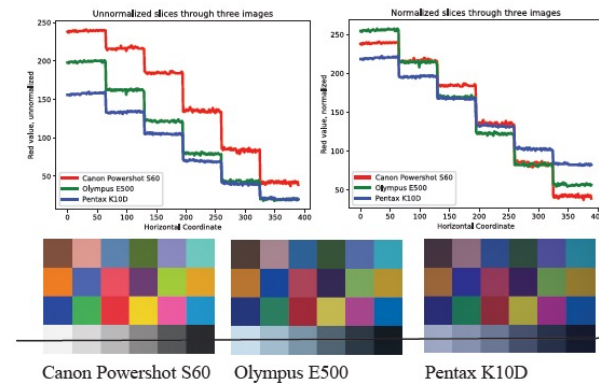From images



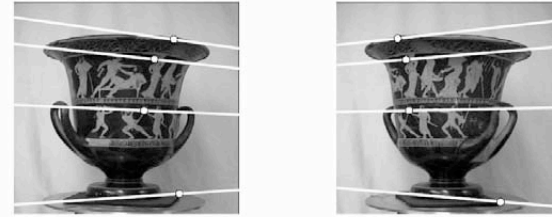Learned denoising



Object detection and segmentation

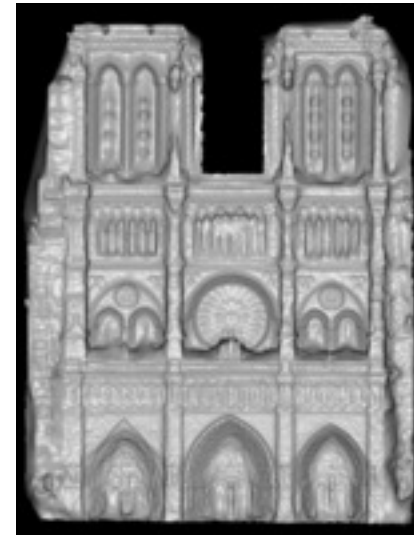# IV. Image formation and geometric vision


Cameras and sensors

# V. Pairs of cameras and more


Two-view geometry, stereo


Structure from motion


Multi-view stereo