

CS 6476-A: Computer Vision

Instructor: James Hays

TAs: Yiming Chen, Jim James, (head TAs),
Sirish Gambhira, Kartik Garg, Keerthi Kaashyap, Amogh Palasamudram,
Esther Shen, Mohit Talreja, Sanchit Tanwar, Haotian Xue

Today's Class

- Who am I?
- What is Computer Vision?
- Specifics of this course
- Geometry of Image Formation
- Questions

A bit about me



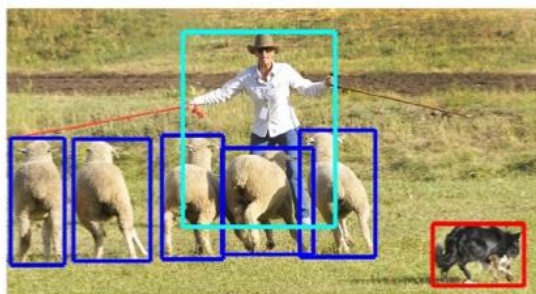
What type of stuff do I work on?

Microsoft COCO: Common Objects in Context

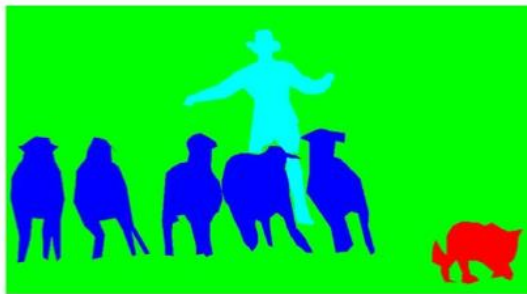
Tsung-Yi Lin¹, Michael Maire², Serge Belongie¹, James Hays³, Pietro Perona²,
Deva Ramanan⁴, Piotr Dollár⁵, and C. Lawrence Zitnick⁵



(a) Image classification



(b) Object localization

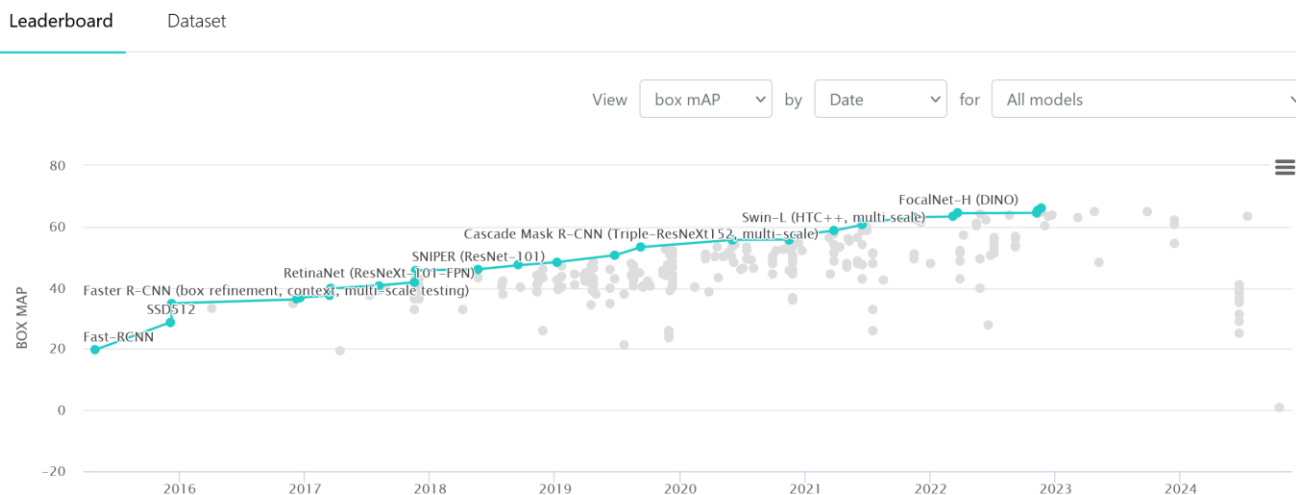


(c) Semantic segmentation



(d) This work

Object Detection on COCO test-dev



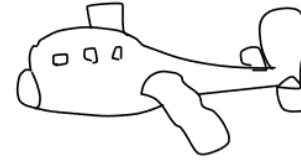
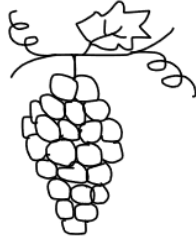
Winner of ECCV 2024
Koenderink Prize and
PAMI Mark
Everingham Prize

How Do Humans Sketch Objects?

Mathias Eitz*
TU Berlin

James Hays†
Brown University

Marc Alexa‡
TU Berlin



| | | | | | |
|---------------------|---------------------|------------------------|----------------------|-----------------------|---------------------|
| t-shirt 100% | snake 99% | comb 99% | flower 99% | eyeglasses 98% | elephant 98% |
| leaf 98% | sun 98% | wrist-watch 96% | pineapple 96% | trousers 96% | ladder 96% |
| apple 96% | airplane 96% | butterfly 96% | umbrella 96% | chair 95% | key 95% |

| | | | | | |
|--------------------------|-----------------------|---------------------|------------------|----------------------|-----------------------|
| seagull 2.5% | panda 11% | armchair 13% | tire 21% | ashtray 24% | snowboard 25% |
| flying bird 47% | bear 44% | chair 89% | wheel 44% | cigarette 30% | skateboard 32% |
| standing bird 24% | teddy bear 30% | couch 3% | donut 16% | bowl 15% | knife 7% |
| pigeon 14% | dog 8% | bench 1% | fan 6% | bathtub 11% | canoe 3% |

Siggraph 2012. **Won Siggraph Test of Time Award 2024.**

Personalized Residuals



Concept



"A rusty V toy gnome in a post-apocalyptic landscape"*



Concept



"V plushie oil painting Ghibli inspired"*



Concept



"V cat wearing sunglasses"*

Personalized Residuals
+ LAG Sampling



Concept



"V action figure riding a motorcycle"*



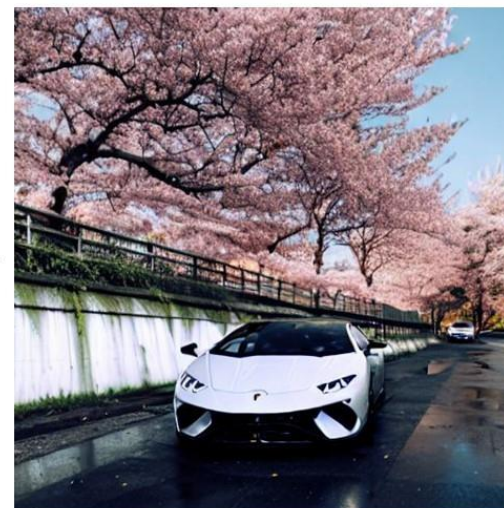
Concept



"The V lighthouse surrounded by a tranquil lake"*



Concept

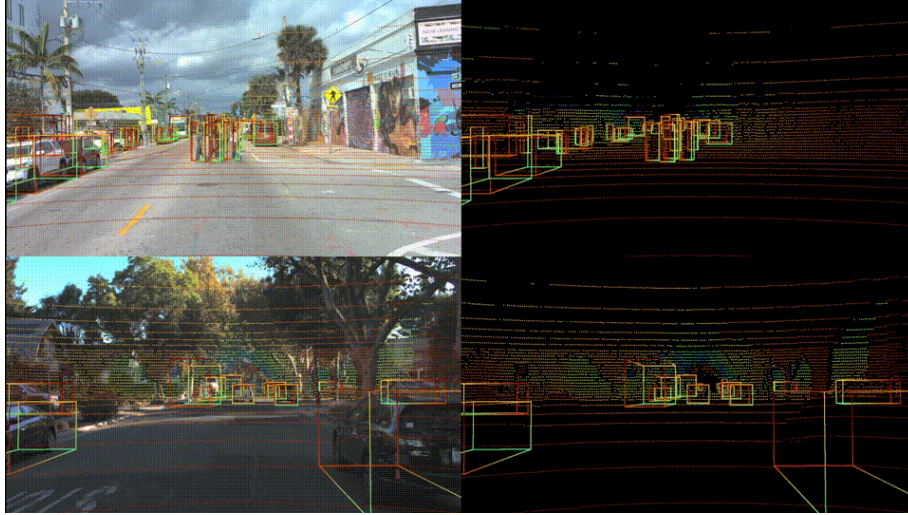


"A V car resting beneath the cherry blossoms in full bloom"*

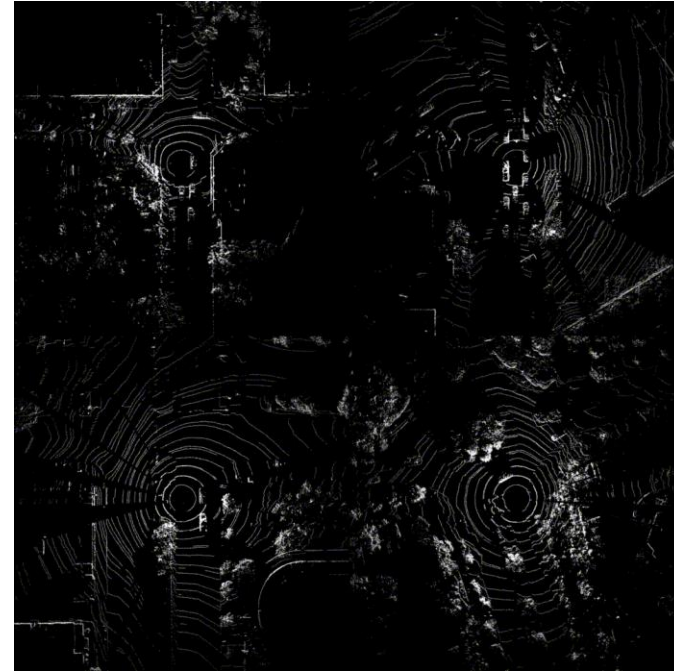
Personalized Residuals for Concept-Driven Text-to-Image Generation. Cusuh Ham, Matthew Fisher, James Hays, Nicholas Kolkin, Yuchen Liu, Richard Zhang, Tobias Hinz. CVPR 2024

Argoverse 2 (AV2) : Four Datasets

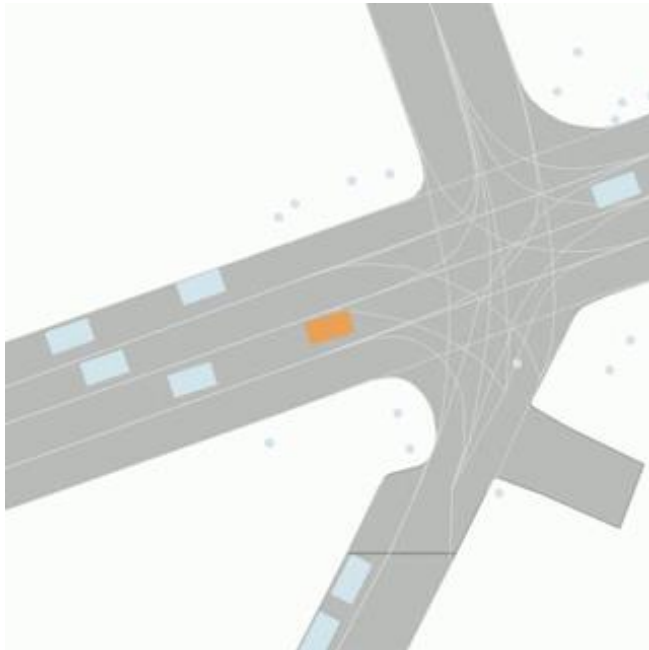
Sensor



Lidar



Motion
Forecasting



Map
Change

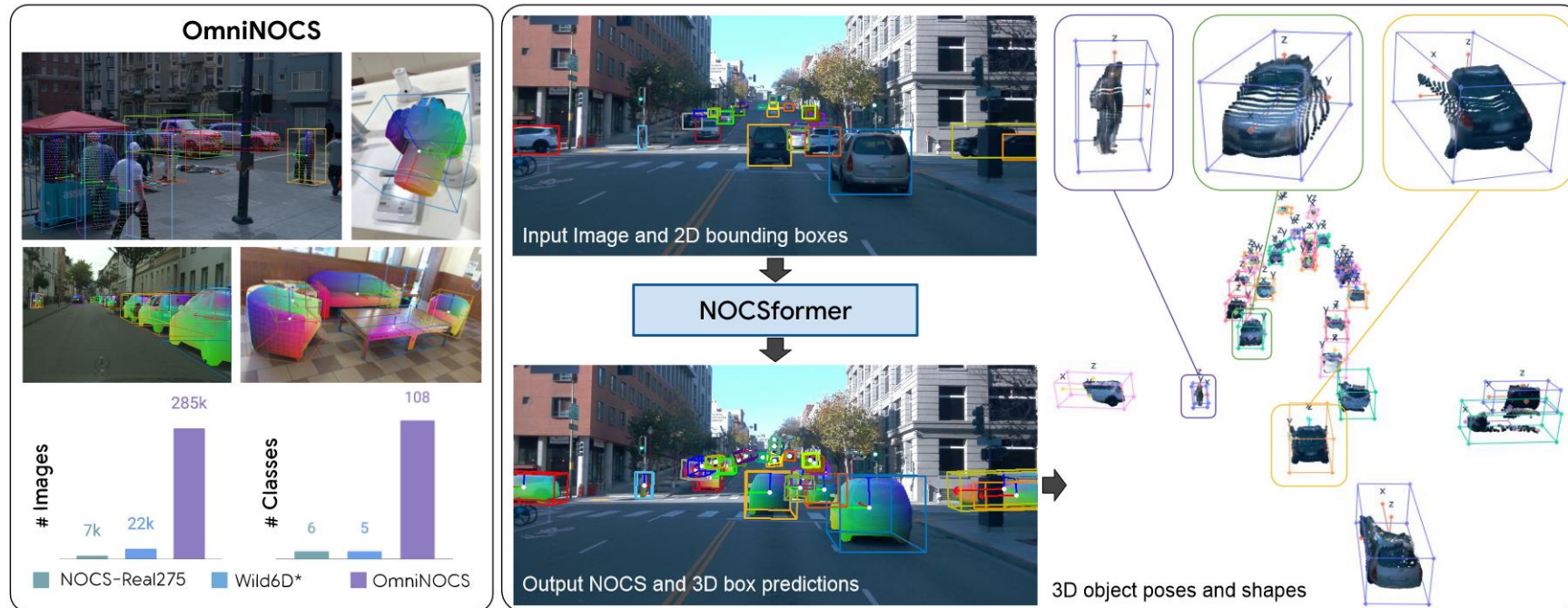


OmniNOCS: A unified NOCS dataset and model for 3D lifting of 2D objects

Akshay Krishnan^{1,2} , Abhijit Kundu¹ , Kavis-Kokitsi Maninis¹ , James Hays² , and Matthew Brown¹ 

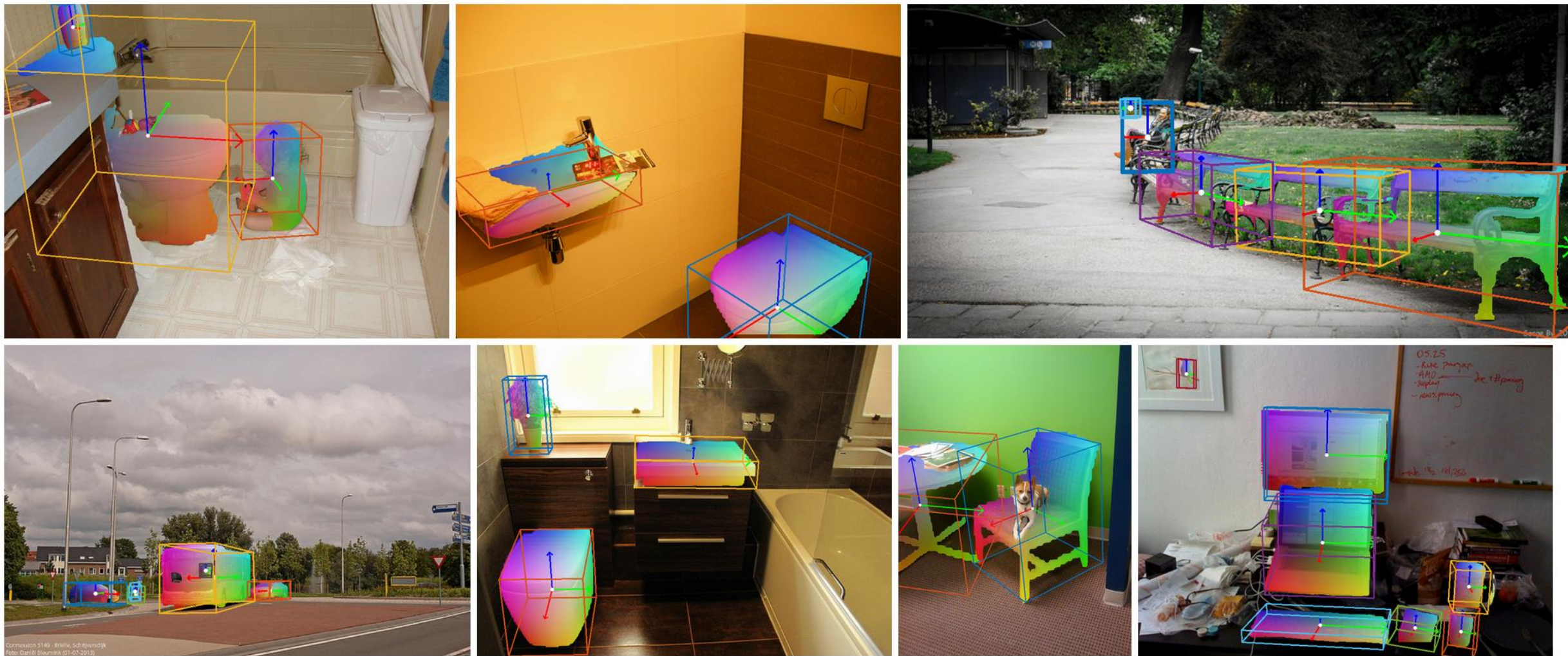
¹ Google Research[†]

² Georgia Institute of Technology



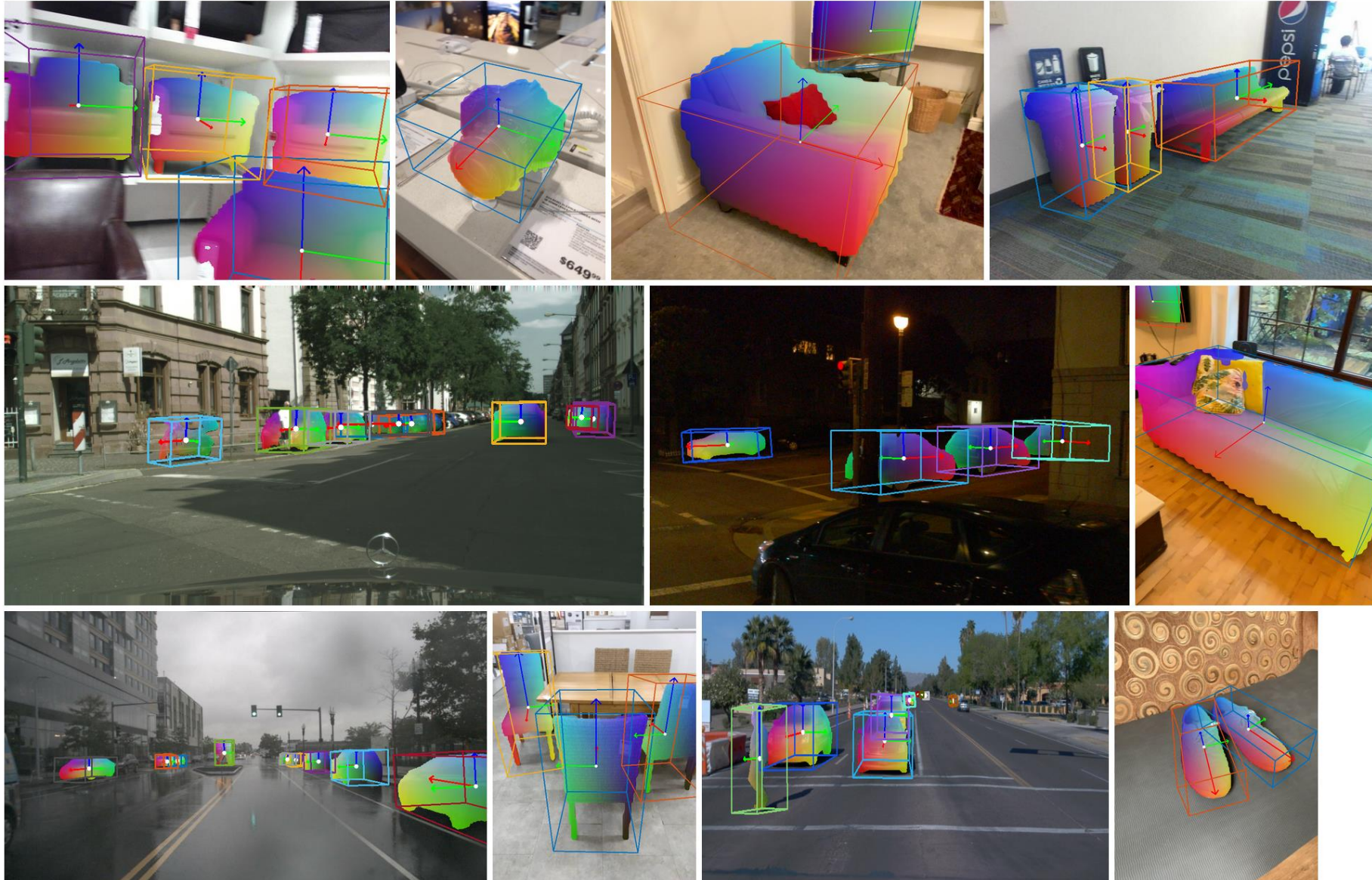
NOCS predictions on COCO objects

NOCSformer can generalize to in-the-wild objects in COCO images when trained on OmniNOCS.

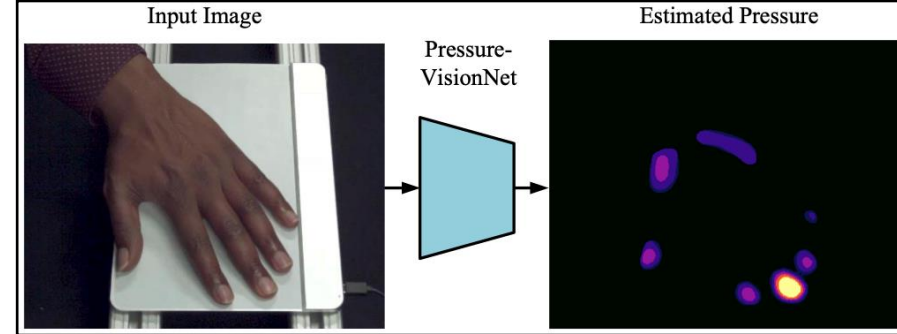
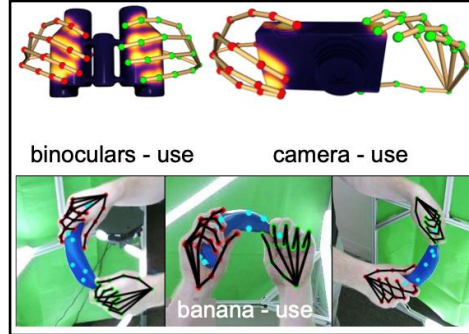
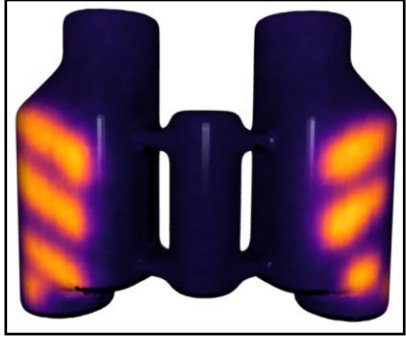


NOCS predictions across OmniNOCS

NOCSformer generalizes to the wide range of object classes and domains in OmniNOCS, including indoor and outdoor scenes, as well as object-centric images.



Creative Sensing for People and Robots



Presented by



James
Hays



Samarth
Brahmbhatt



Patrick
Grady



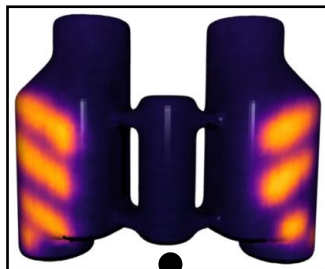
Mengyu
Yang



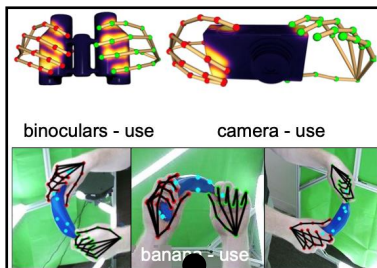
Charles C.
Kemp

And collaborators Cusuh Ham, Chengcheng Tang, Christopher D. Twigg, Minh Vo, Chengde Wan, Ankur Handa, Dieter Fox, Jeremy Collins

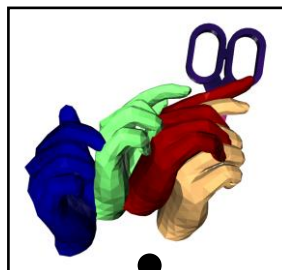
Brahmbhatt et al
CVPR '19 (oral)
Best Paper finalist



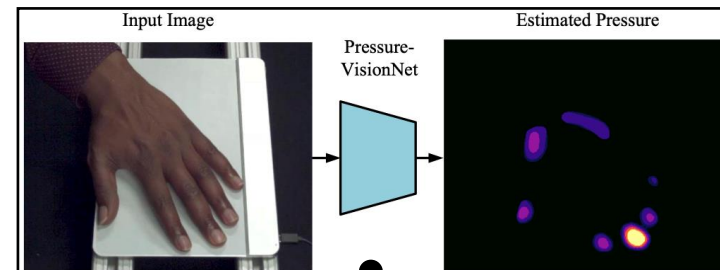
Brahmbhatt et al
ECCV '20



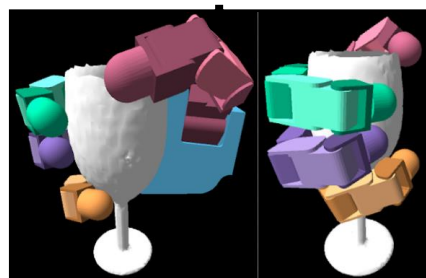
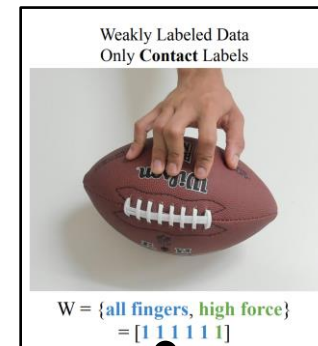
Grady et al
CVPR '21 (oral)



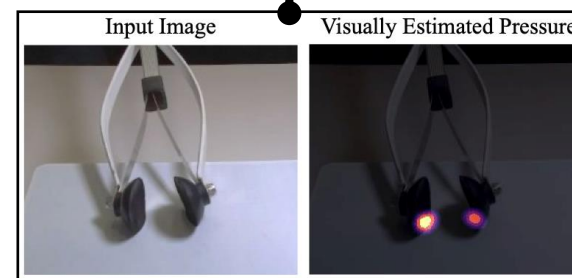
Grady et al
ECCV '22 (oral)



Grady et al.
WACV 2024



Brahmbhatt et al
IROS '19



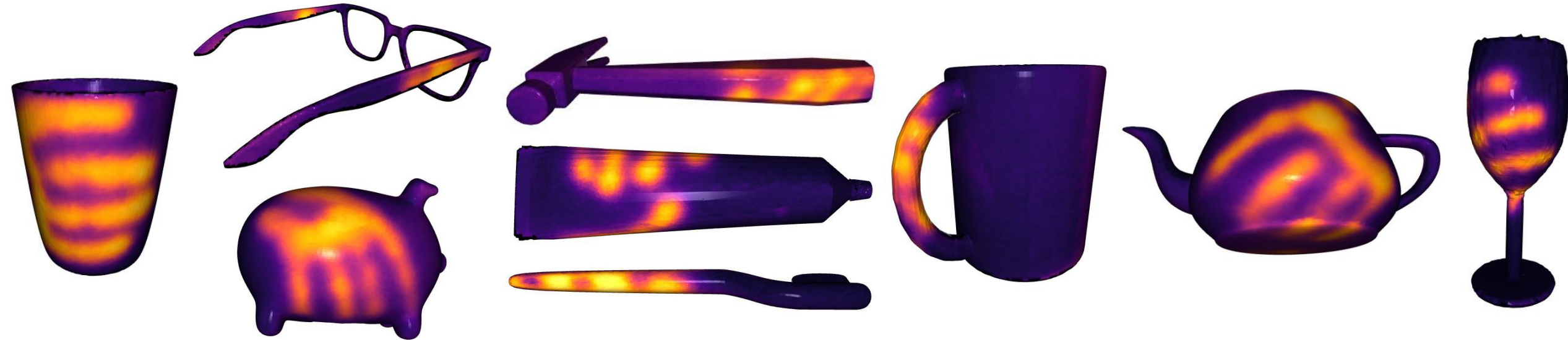
Grady et al
IROS '22

Why is observing contact difficult?

Occlusion



ContactDB: Analyzing and Predicting Grasp Contact via Thermal Imaging



Samarth
Brahmabhatt



Cusuh Ham



Charles C.
Kemp



James
Hays



2 seconds



5 seconds



10 seconds



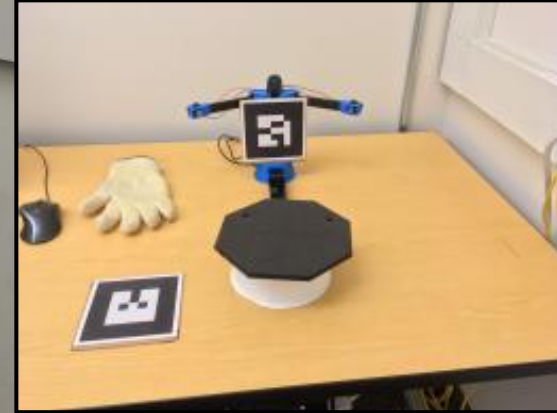
Computer

Turntable

Camera

A dark wooden table with a yellow ruler and a black pen. The text "Table with 3D printed objects" is overlaid in white.

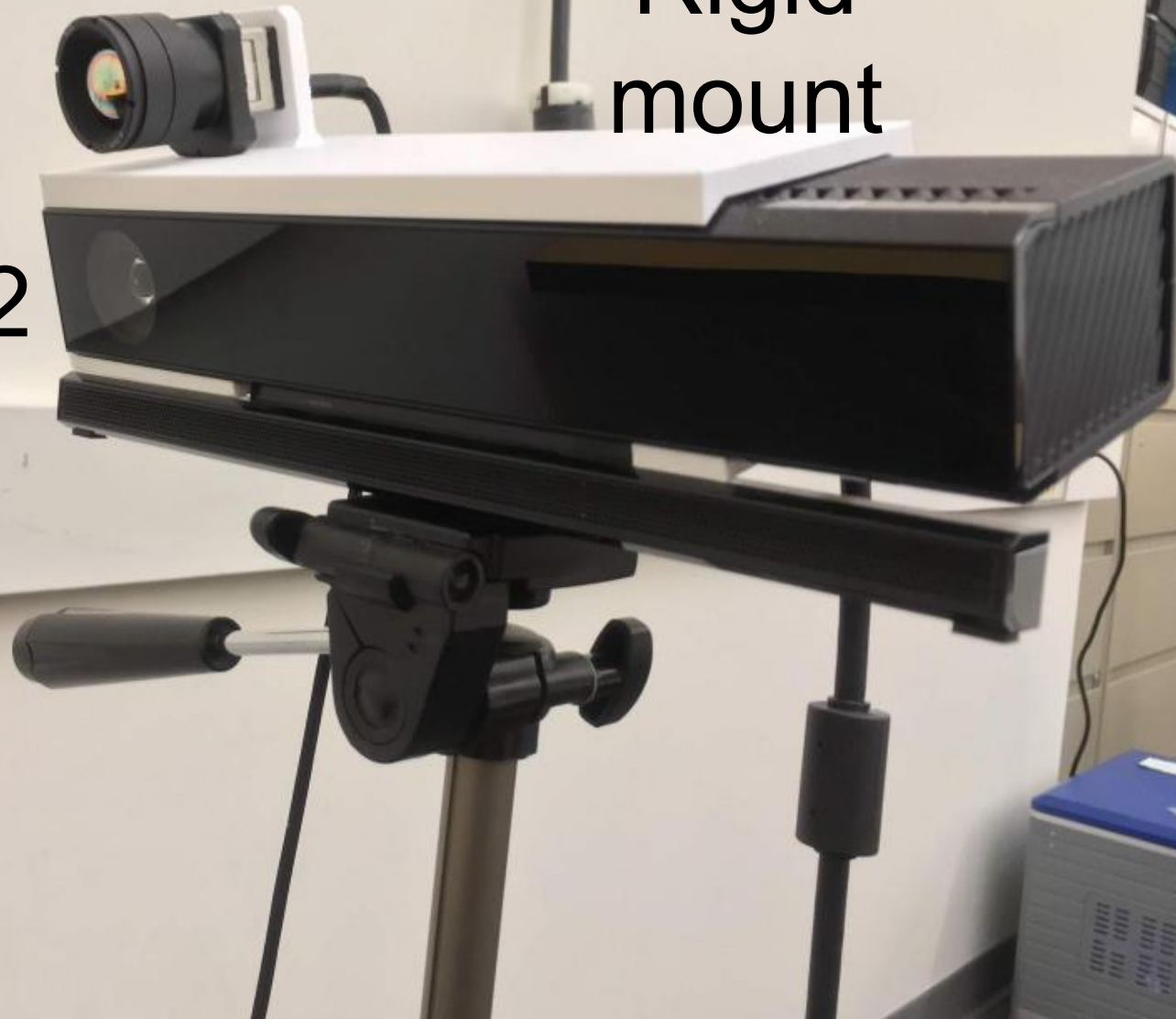
Table with 3D printed objects

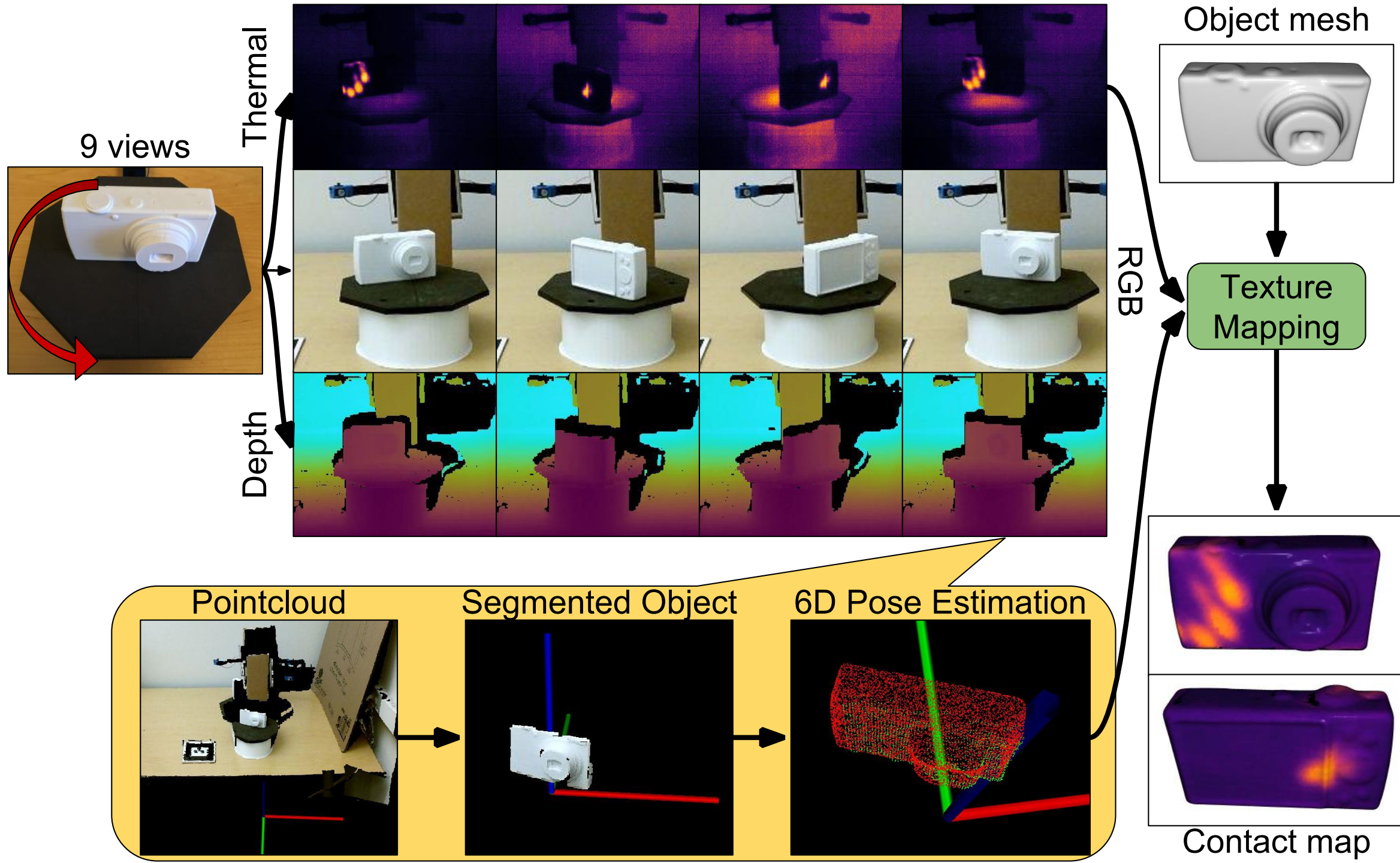


FLIR Boson 640
Thermal camera

Rigid
mount

Kinect v2
RGB-D
camera

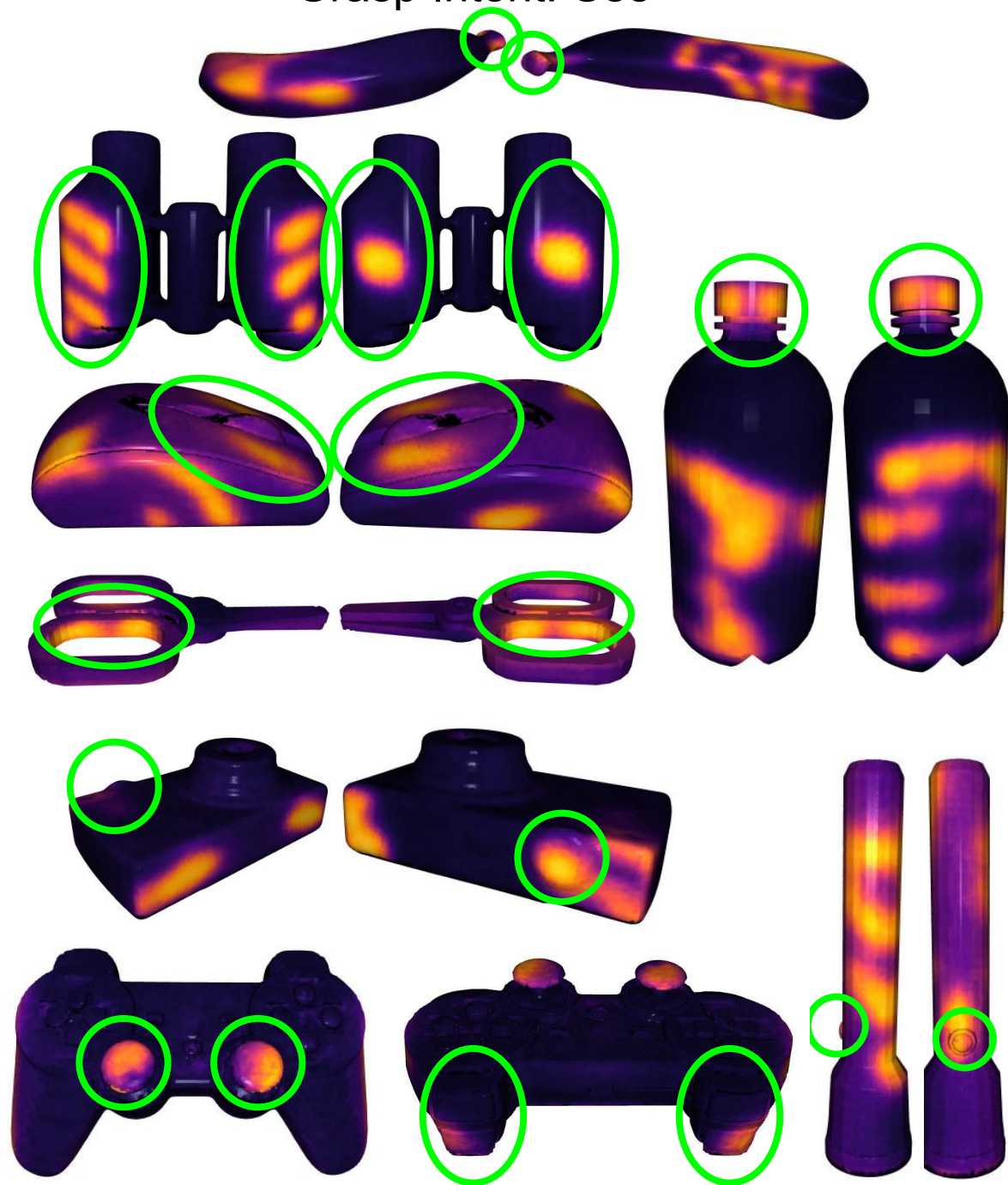




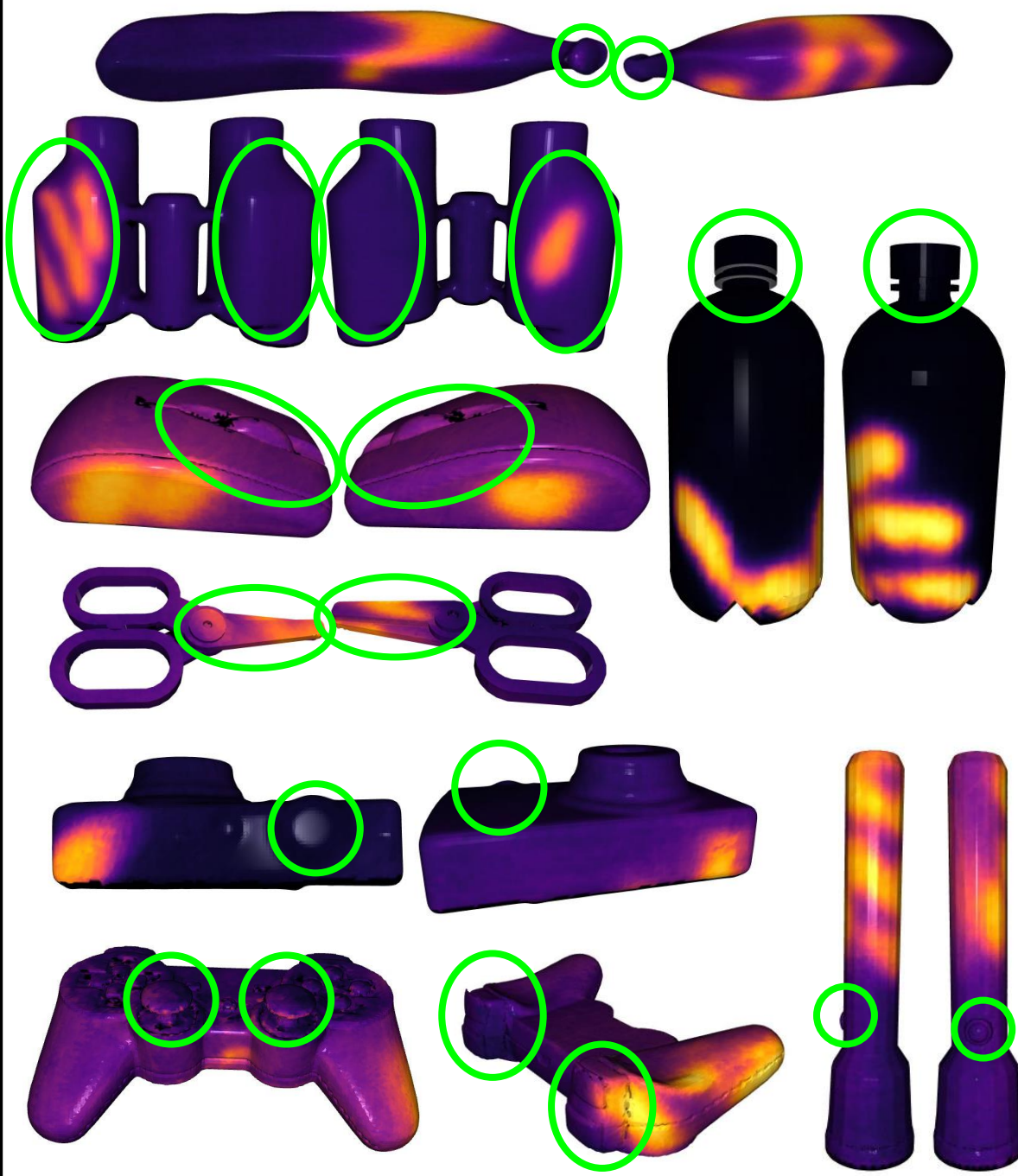


Contact map

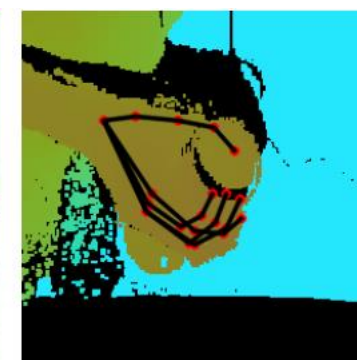
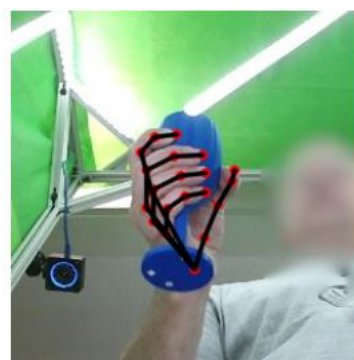
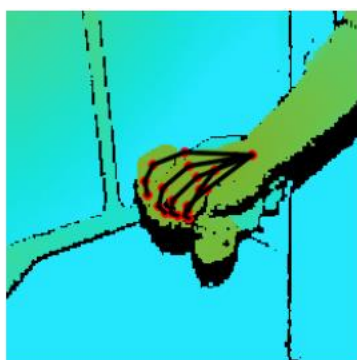
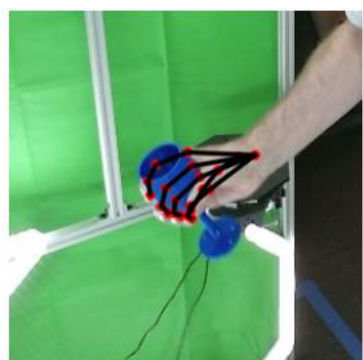
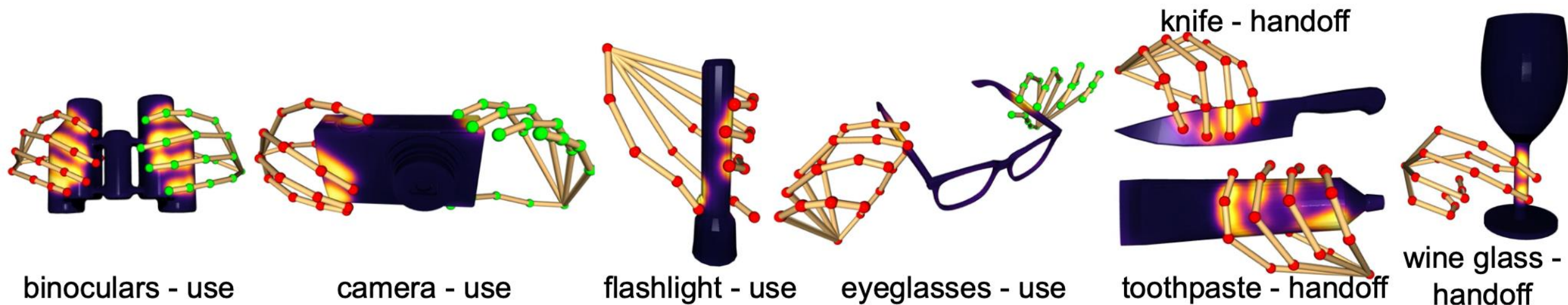
Grasp Intent: Use



Grasp Intent: Handoff



ContactPose: Capturing Contact + Hand Pose



Kinect1-color

Kinect1-depth

Kinect2-color

Kinect2-depth

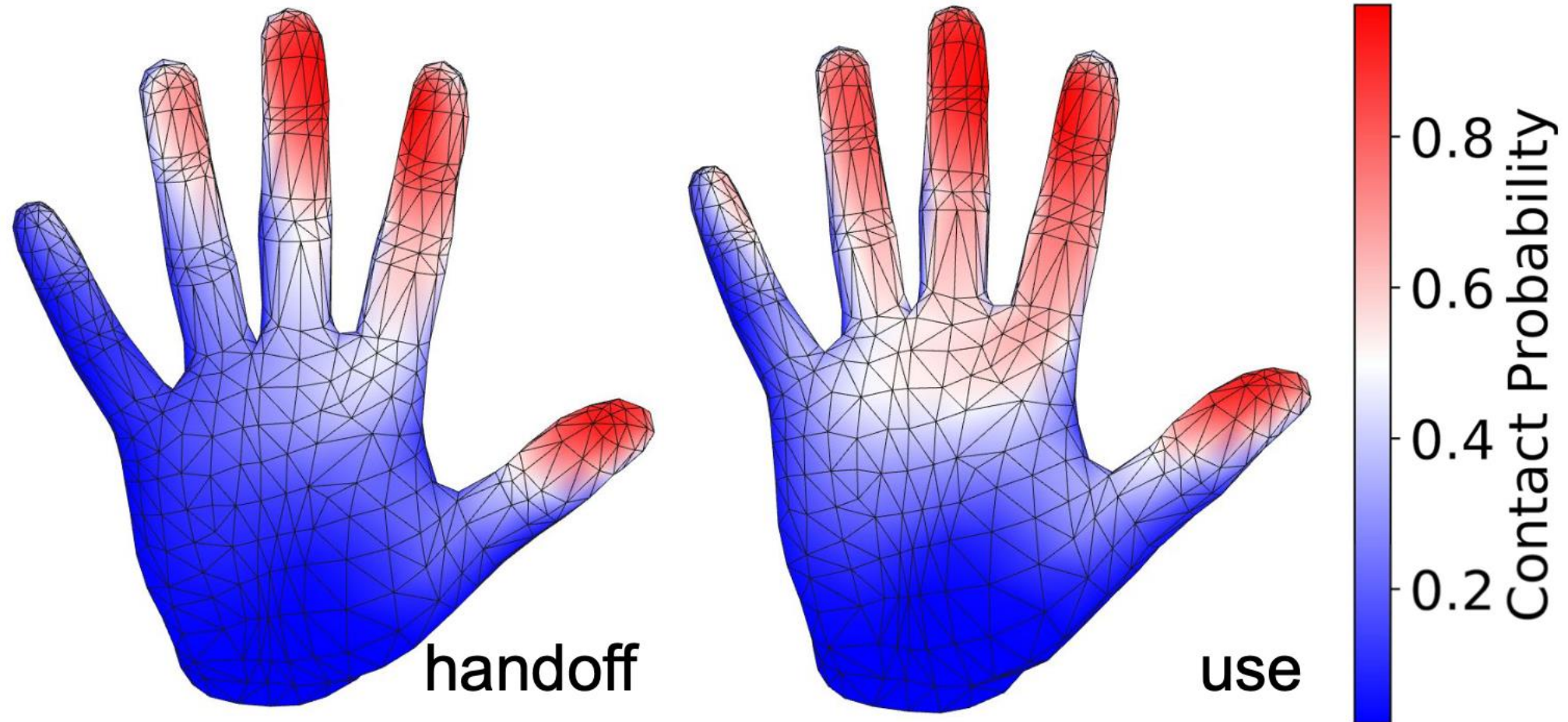
Kinect3-color

Kinect3-depth

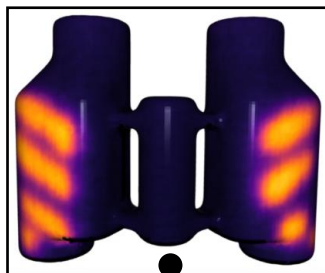
Samarth Brahmbhatt, Chengcheng Tang, Christopher D. Twigg, Charles C. Kemp, and James Hays

ECCV 2020

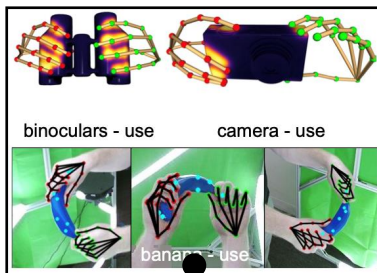
Hand Contact Probability



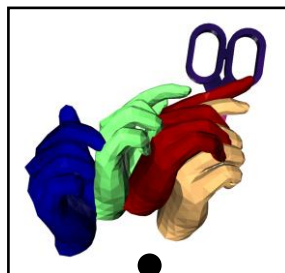
Brahmbhatt et al
CVPR '19 (oral)
Best Paper finalist



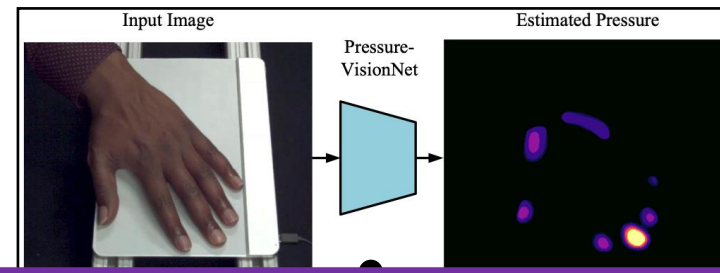
Brahmbhatt et al
ECCV '20



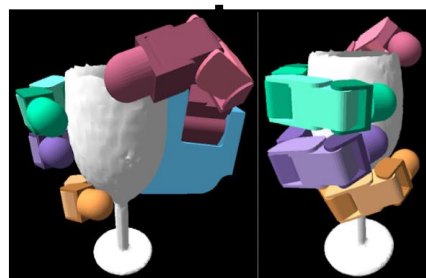
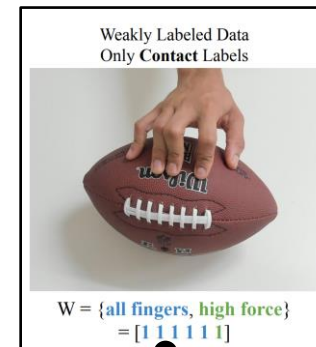
Grady et al
CVPR '21 (oral)



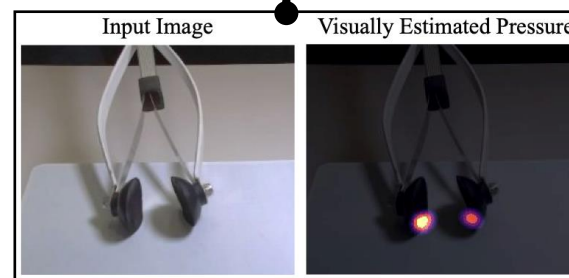
Grady et al
ECCV '22 (oral)



Grady et al.
WACV 2024



Brahmbhatt et al
IROS '19



Grady et al
IROS '22



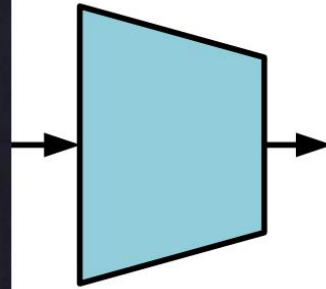


PressureVision: Estimating Hand Pressure from a Single RGB Image

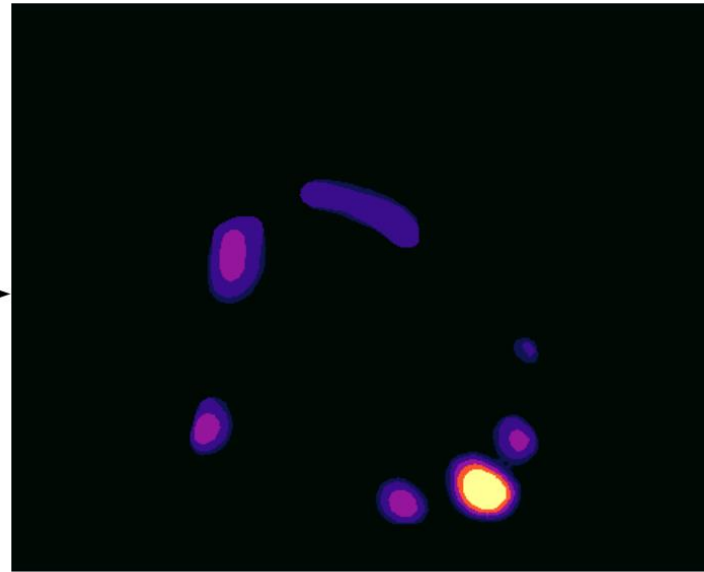
Input Image



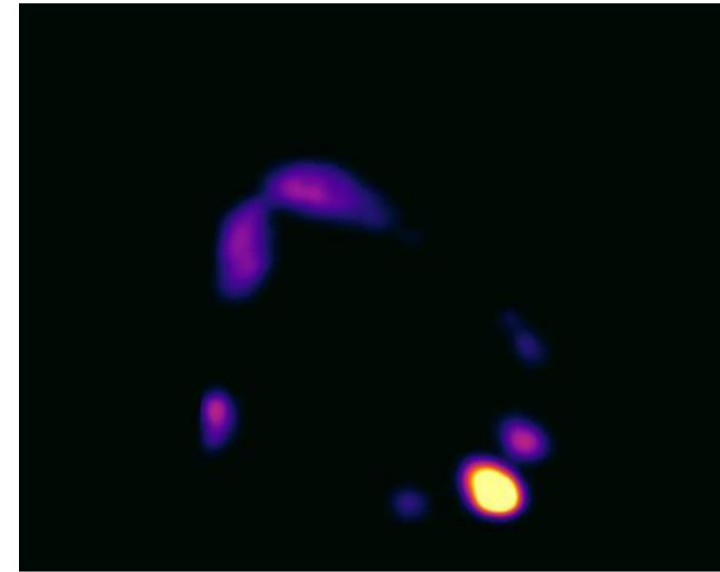
Pressure-
VisionNet



Estimated Pressure

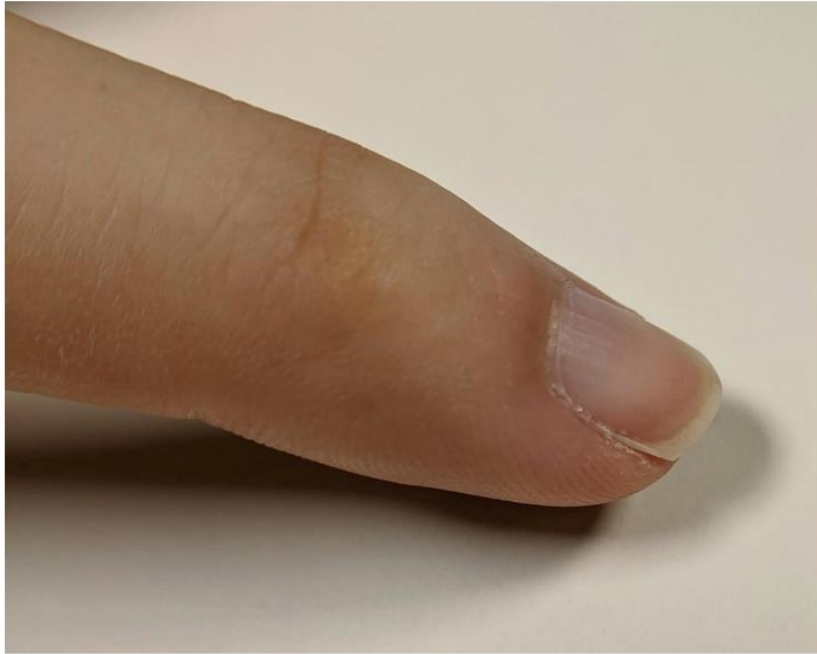


Ground Truth Pressure

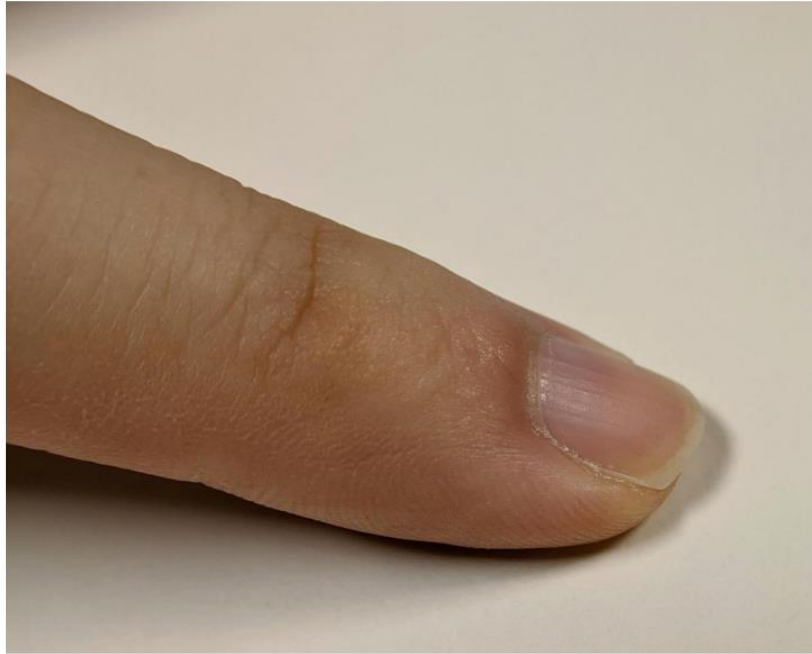


Patrick Grady, Chengcheng Tang, Samarth Brahmbhatt, Christopher D. Twigg,
Chengde Wan, James Hays, and Charles C. Kemp

No Contact



Low Force



High Force



We train a deep network, PressureVisionNet,
to estimate pressure from a single RGB image.

The pressure for each frame is calculated independently.

PressureVision++: Estimating Fingertip Pressure From Diverse RGB Images

Input Image



Estimated Pressure



Prompt: Press thumb and index, low force

Input Image



Estimated Pressure



Prompt: Press all fingers, high force

Patrick Grady, Jeremy Collins, Chengcheng Tang,
Christopher D. Twigg, James Hays, and Charles C. Kemp
WACV 2024

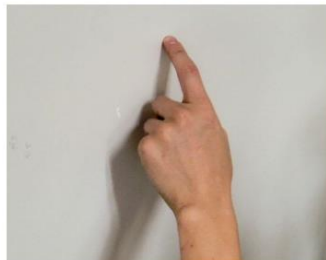
Surface/Prompt

Image

ContactLabelNet

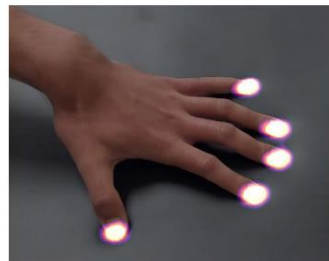
Wall

*Press index
Low force*



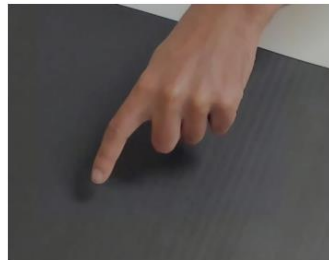
Foam mat

*Press all fingers
High force*



Foam mat

No contact



Mirror

*Press middle
Low force*



Mirror

*Press all fingers
High force*



Surface/Prompt

Image

ContactLabelNet

Football

*Press pinky
High force*



Notebook

*Press all fingers
Low force*



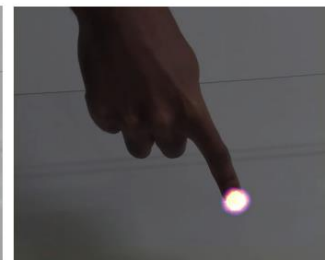
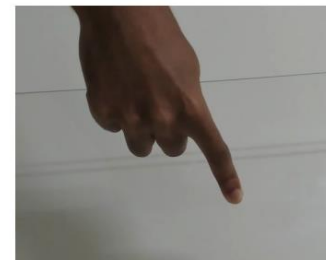
Notebook

*Press index, thumb
Low force*



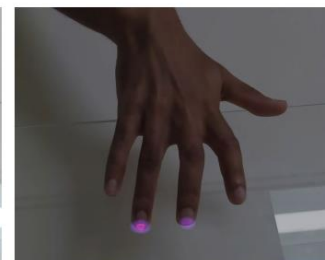
Glass

*Press index
High force*



Glass

*Press ring
Low force*





ICRA2024
YOKOHAMA | JAPAN

The Un-Kidnappable Robot: Acoustic Localization of Sneaking People

Mengyu Yang, Patrick Grady, Samarth Brahmabhatt,
Arun Balajee Vasudevan, Charles C. Kemp, James Hays





ICRA2024
YOKOHAMA | JAPAN

How easy is it to **sneak**
up on a robot?



ICRA2024
YOKOHAMA | JAPAN

We train robots to detect people using *only* the **subtle and incidental sounds** they produce as they move around, even when they try to be **quiet**.

Real Time

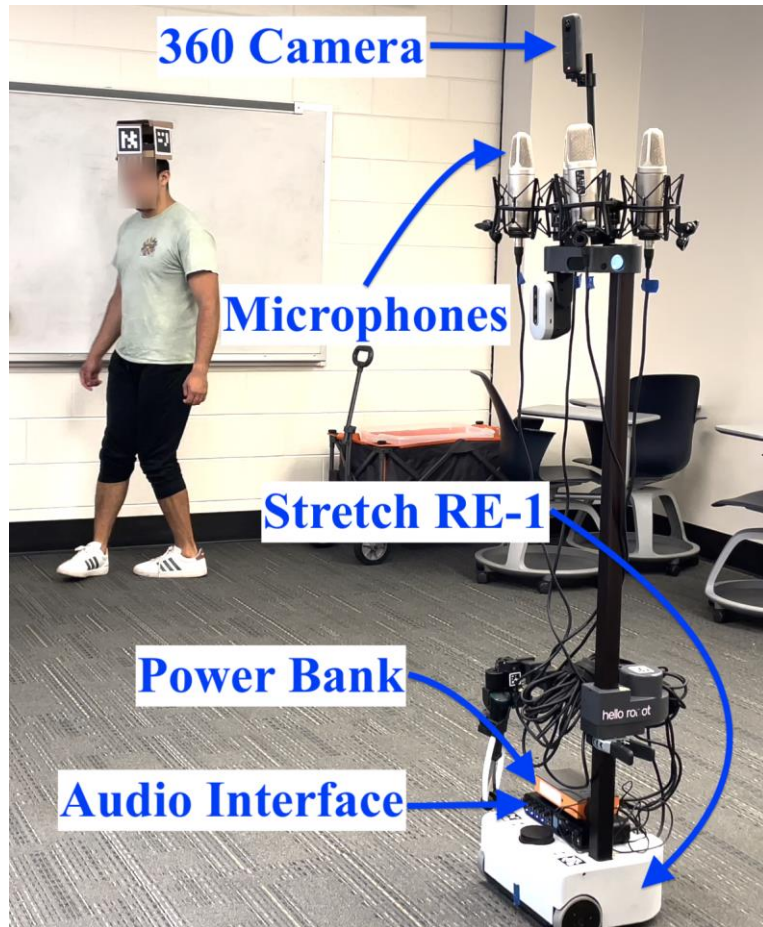




The Robot Kidnapper Dataset



ICRA2024
YOKOHAMA | JAPAN



- 4-channel audio
- 360 degree egocentric video
- 12 participants in 8 indoor settings:
 - Standing still
 - Walking quietly
 - Walking normally
 - Walking loudly

Data Annotation



ICRA2024
YOKOHAMA | JAPAN



Azimuthal angle

Data Annotation



ICRA2024
YOKOHAMA | JAPAN

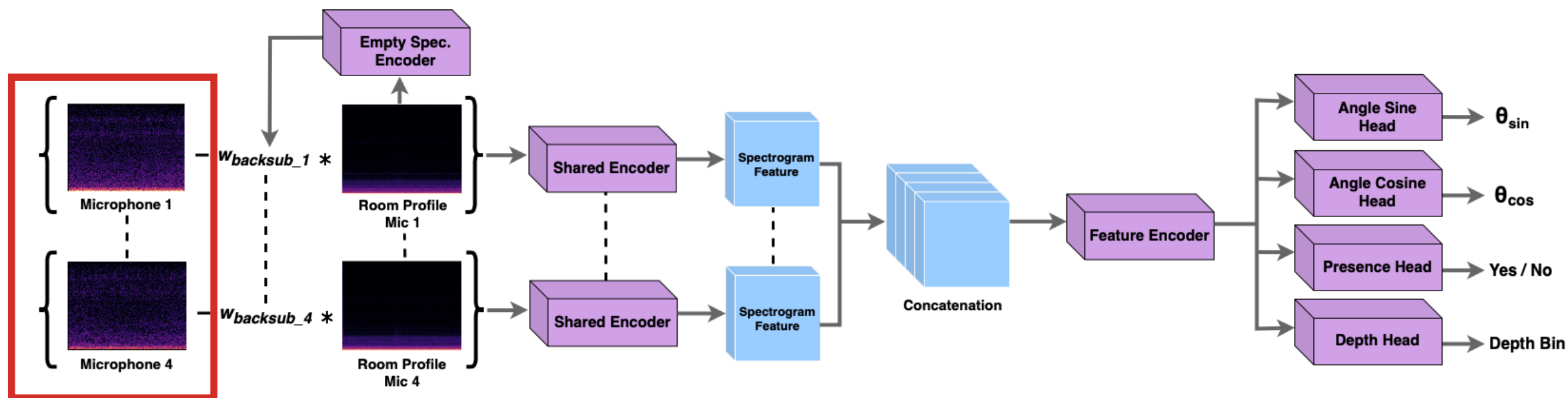


Radial distance from detecting ArUco markers

Architecture



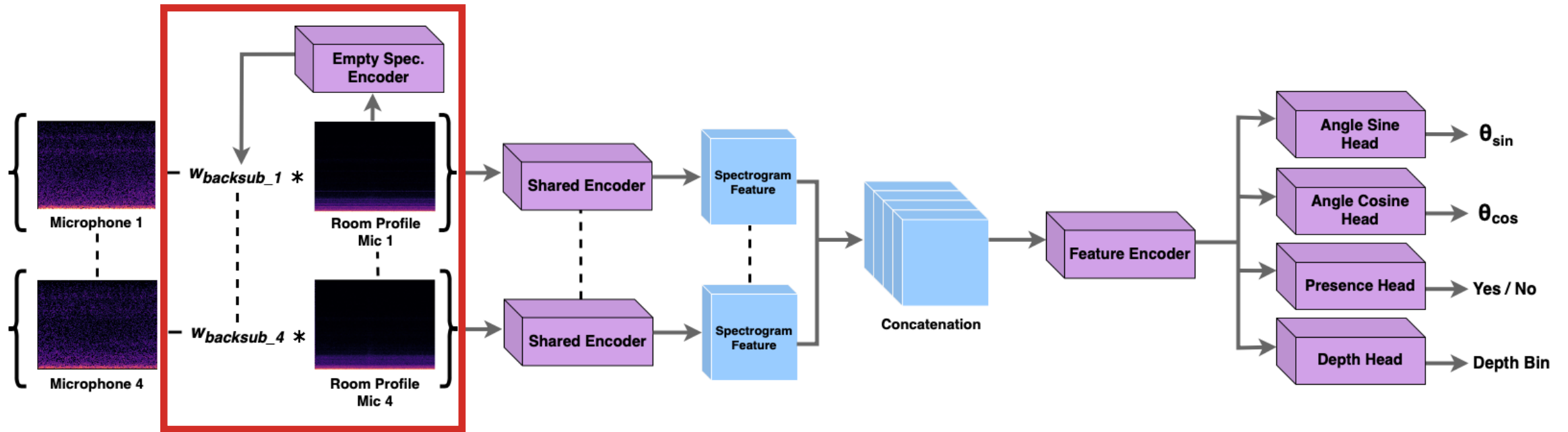
ICRA2024
YOKOHAMA | JAPAN



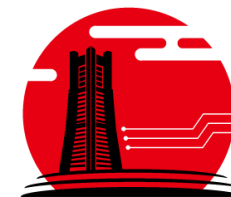
Architecture



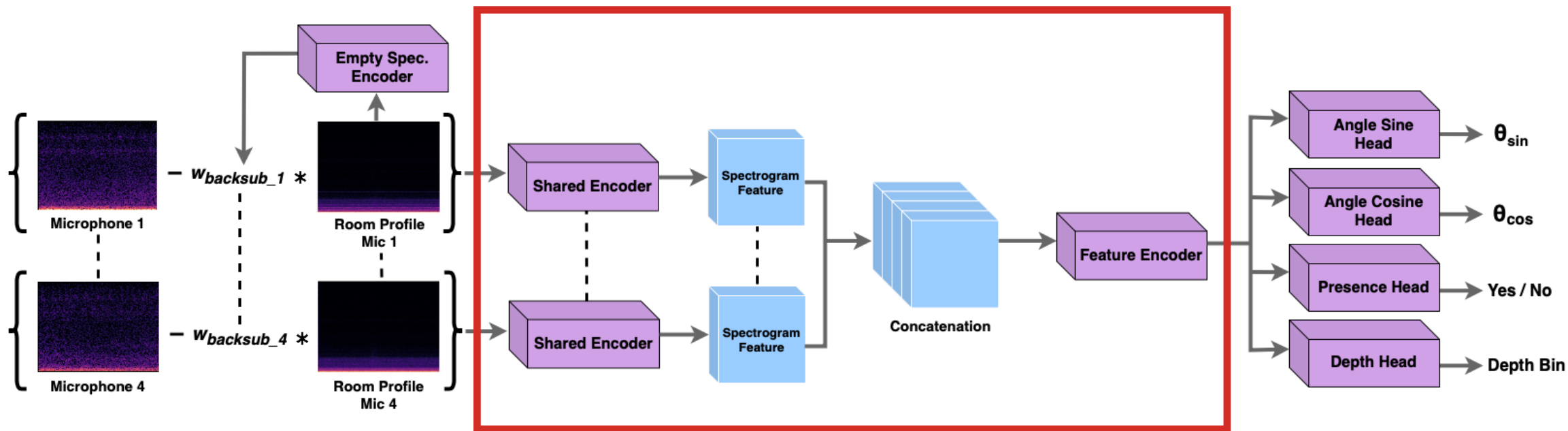
ICRA2024
YOKOHAMA | JAPAN



Architecture



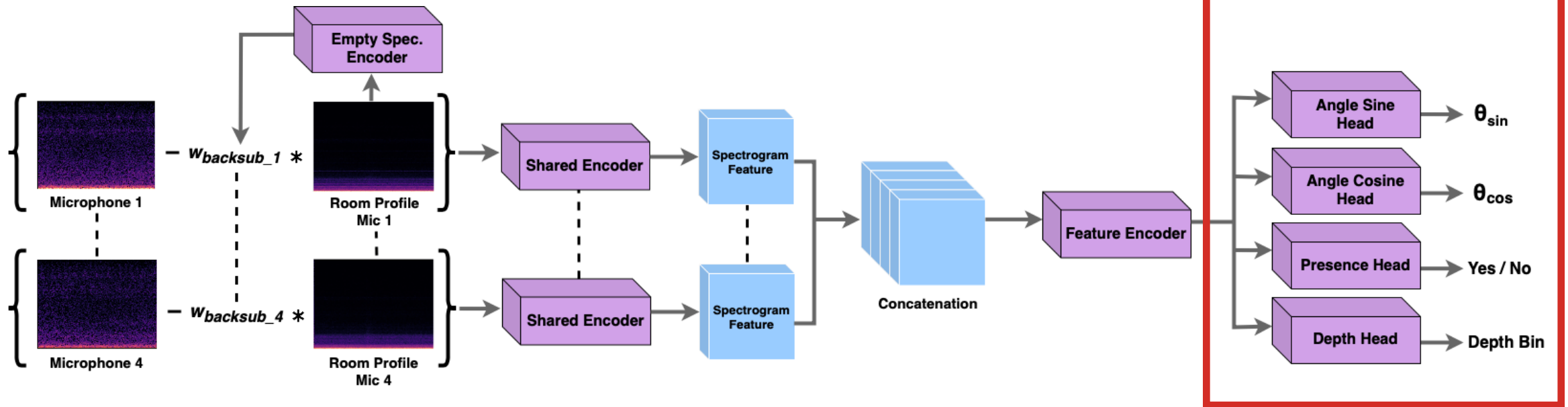
ICRA2024
YOKOHAMA | JAPAN



Architecture



ICRA2024
YOKOHAMA | JAPAN



Azimuthal Angle Prediction



ICRA2024
YOKOHAMA | JAPAN

| CATEGORY | MODEL | <i>Quiet</i> | | <i>Normal</i> | | <i>Loud</i> | |
|-----------------|----------------|--------------|-----------|---------------|-----------|-------------|-----------|
| | | Sta. | Dyn. | Sta. | Dyn. | Sta. | Dyn. |
| Random | Uniform 360° | 90 | 90 | 90 | 90 | 90 | 90 |
| Oracle Mic Pair | Constant Front | 50 | 43 | 50 | 46 | 50 | 43 |
| | GCC-PHAT [25] | 44 | 47 | 45 | 43 | 46 | 47 |
| | StereoCRW [26] | 52 | 46 | 51 | 48 | 37 | 34 |
| Ours | 1 Mic | 67 | 75 | 64 | 71 | 64 | 74 |
| | 2 Mics | 37 | 54 | 37 | 48 | 36 | 47 |
| | Base 4 Mics | 47 | 55 | 50 | 48 | 49 | 47 |
| | 4 Mics | 21 | 26 | 22 | 24 | 19 | 22 |

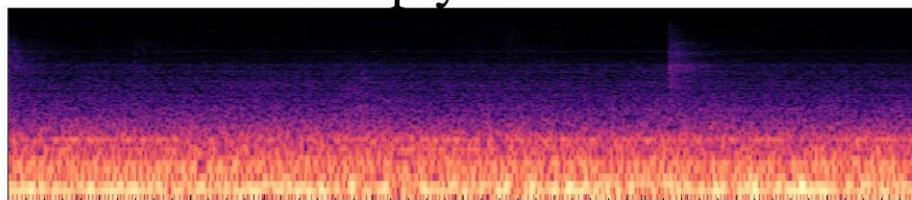
Mean absolute error (MAE) in degrees

Qualitative Comparisons

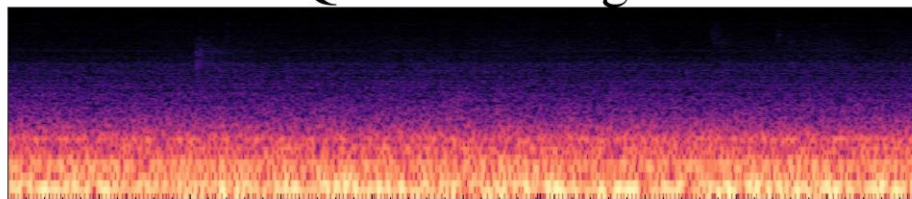


ICRA2024
YOKOHAMA | JAPAN

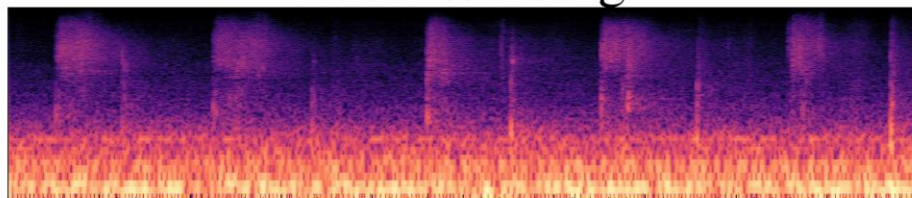
Empty Room



Quiet Walking

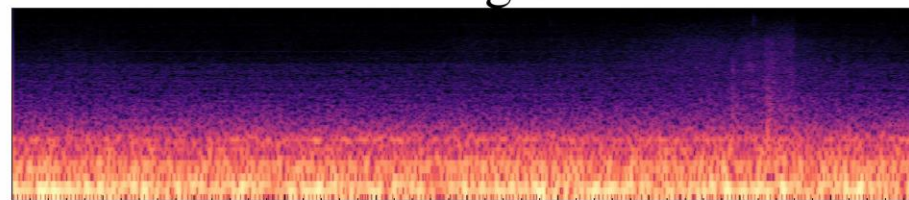


Loud Walking

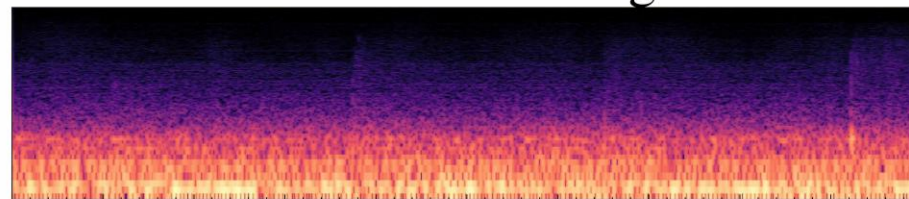


0 0.5 1 1.5 2 2.5 3 3.5 4 4.5 5

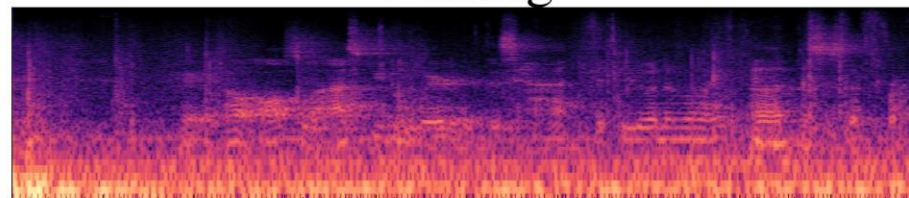
Standing Still



Normal Walking



Talking



0 0.5 1 1.5 2 2.5 3 3.5 4 4.5 5

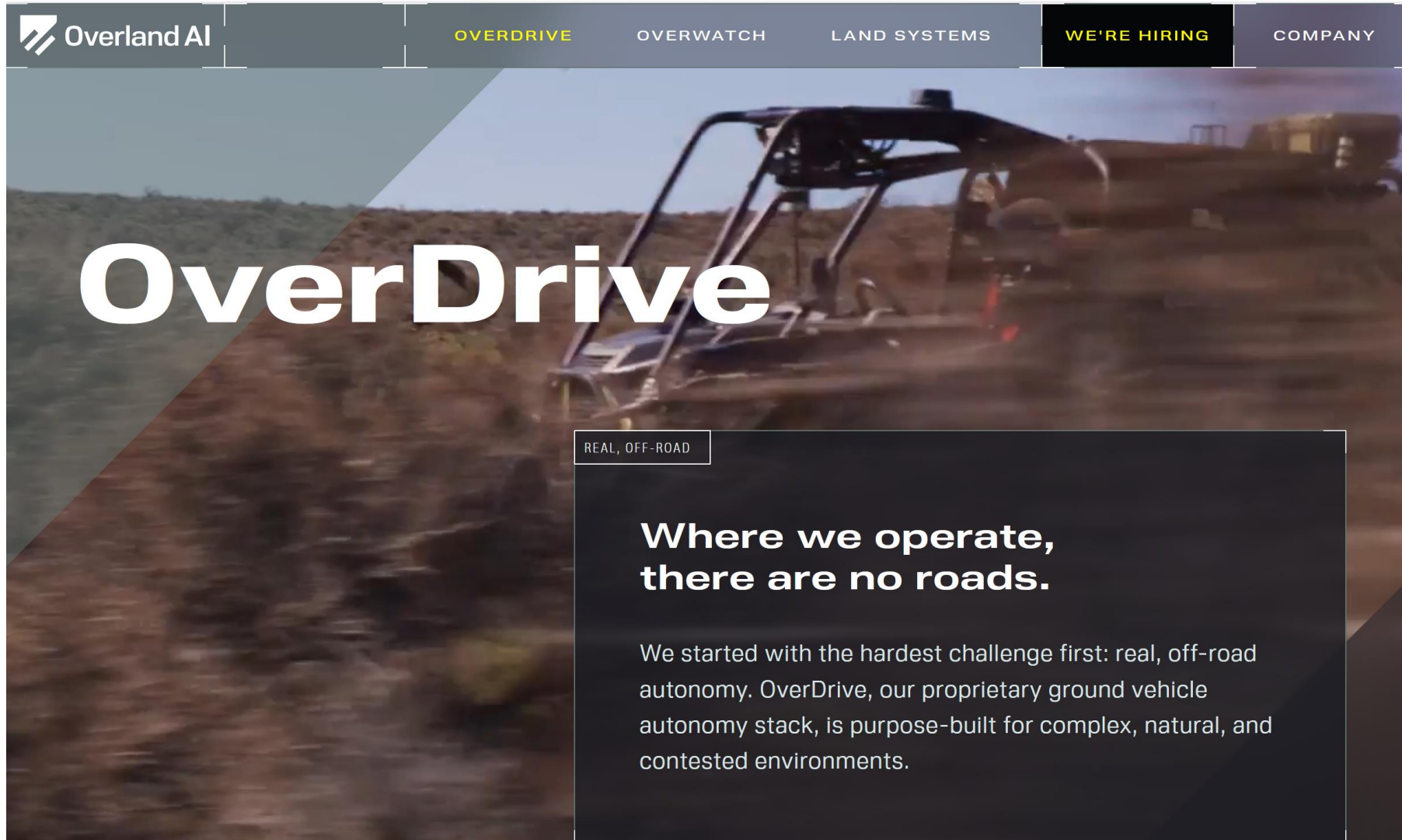
Conclusion



ICRA2024
YOKOHAMA | JAPAN

- Human detection with only subtle incidental sounds of them moving
- Robot Kidnapper dataset collected on robot in real-world indoor environments
- Our model outperforms previous sound localization methods
- Real-time detection on robot

I still work on autonomous vehicles



OVERDRIVE

OVERWATCH

LAND SYSTEMS

WE'RE HIRING

COMPANY

OverDrive

REAL, OFF-ROAD

**Where we operate,
there are no roads.**

We started with the hardest challenge first: real, off-road autonomy. OverDrive, our proprietary ground vehicle autonomy stack, is purpose-built for complex, natural, and contested environments.

Today's Class

- ~~Who am I?~~
- What is Computer Vision?
- Specifics of this course
- Geometry of Image Formation
- Questions

What is Computer Vision?

Derogatory summary of computer vision:
Machine learning applied to visual data



<http://cocodataset.org/#explore?id=550719>

http://farm5.staticflickr.com/4106/4990212001_a18aeffa3d_z.jpg





Detect all of the objects in this photo

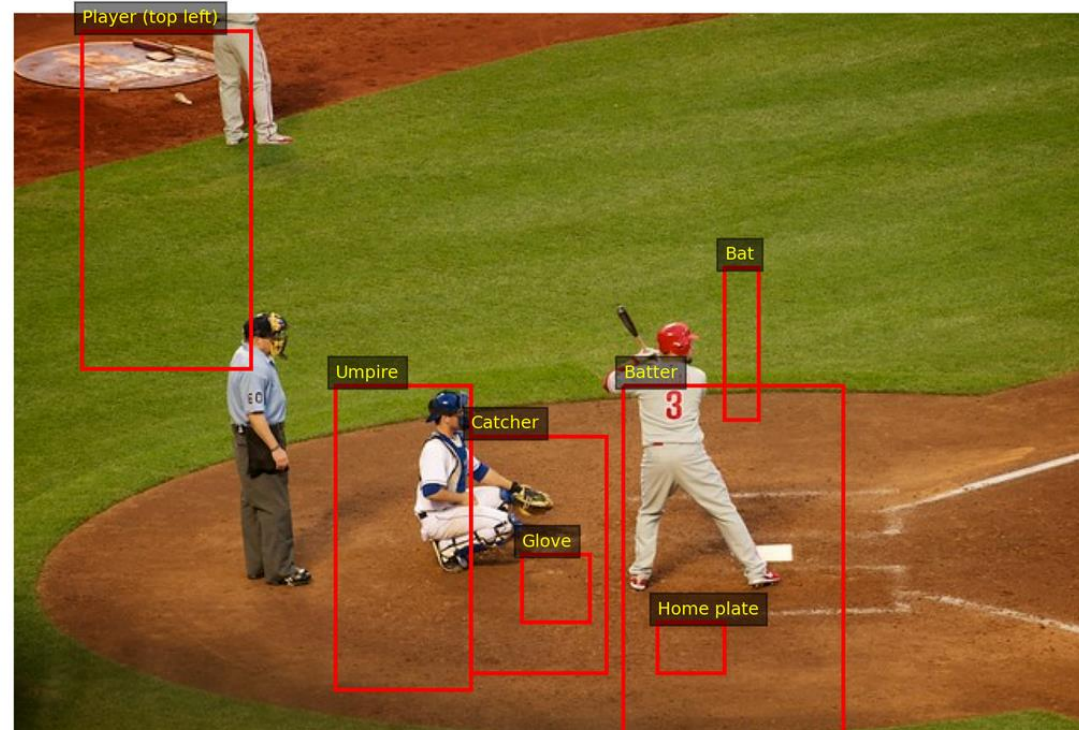


Here are the main objects detected in the photo:

- **Baseball batter** (wearing a red helmet and white uniform with number 3)
- **Catcher** (wearing a blue helmet and protective gear, crouched behind home plate)
- **Umpire** (in a blue shirt, black pants, and protective gear, standing to the left of the catcher)
- **Another player** (partially visible at the top left, standing on the dirt near the baseline)
- **Baseball bat** (held by the batter)
- **Baseball glove** (worn by the catcher)
- **Home plate** (white plate on the ground where the batter stands)
- **Batting box chalk lines** (white chalk lines marking the batter's box)
- **Grass field** (green outfield/turf area)
- **Dirt infield** (brown/red dirt around home plate area)



Matplotlib Chart



Here's the image with bounding boxes drawn around the detected objects:

- **Batter** (red helmet, white uniform, number 3)
- **Catcher** (in blue gear, crouched)
- **Umpire** (in blue shirt, standing)
- **Player (top left)** (partially visible)
- **Bat** (held by batter)
- **Glove** (worn by catcher)
- **Home plate**



Computer Vision

- Automatic understanding of images and video
 1. Computing properties of the 3D world from visual data
(*measurement*)

1. Vision for measurement

Real-time stereo



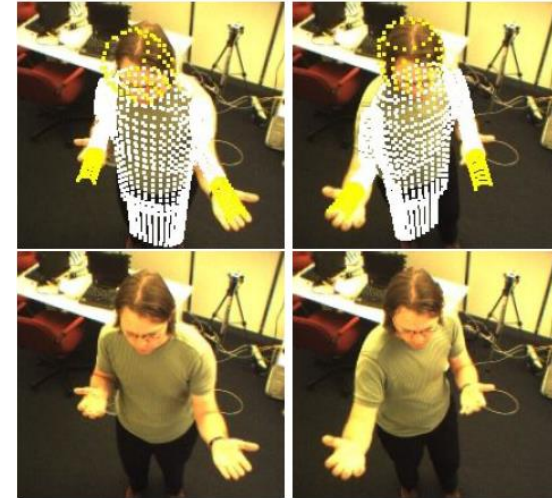
Wang et al.

Structure from motion



Snavely et al.

Tracking

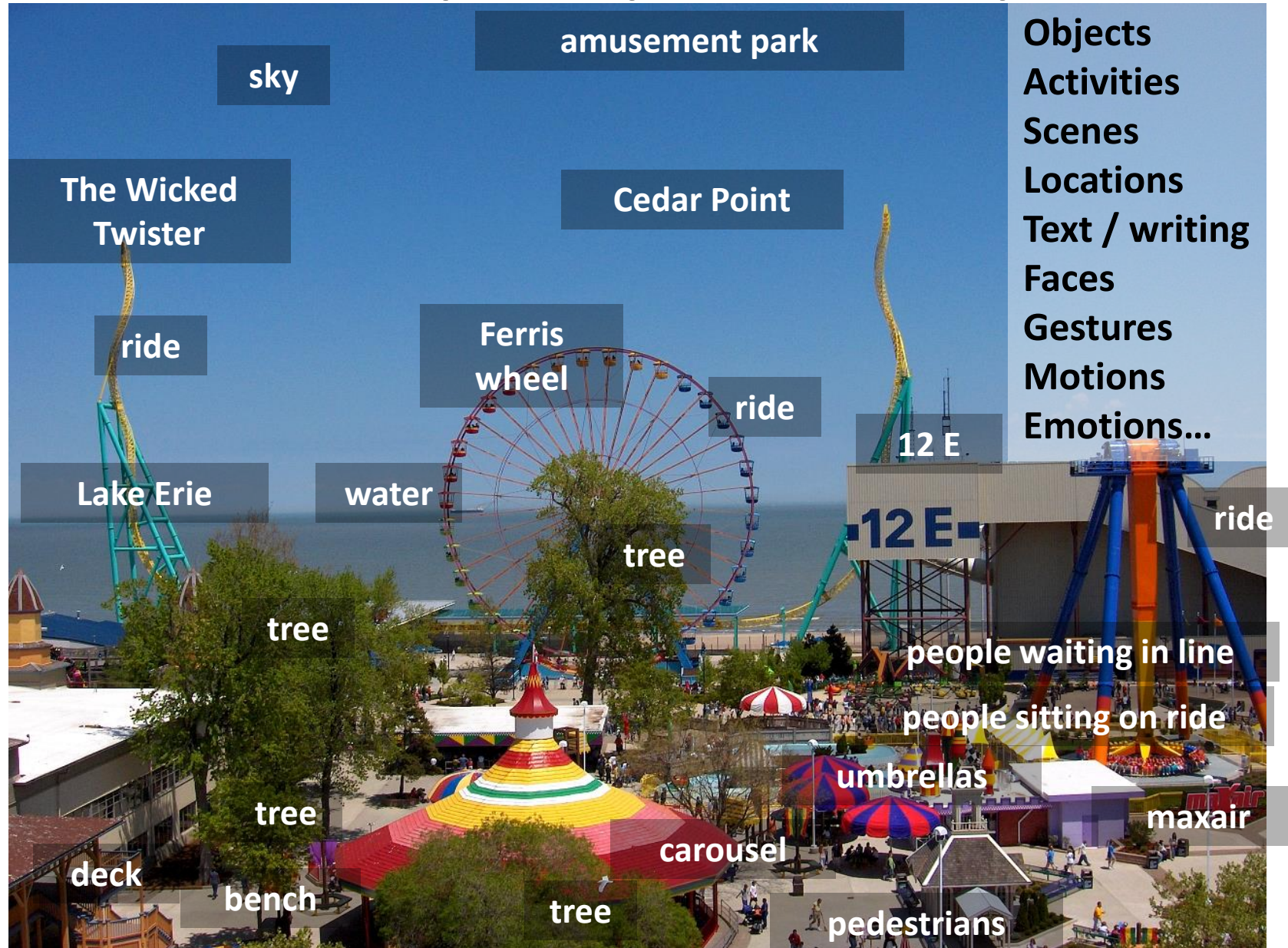


Demirdjian et al.

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, people, scenes, and activities.
(perception and interpretation)

2. Vision for perception, interpretation



Slide credit: Kristen Grauman

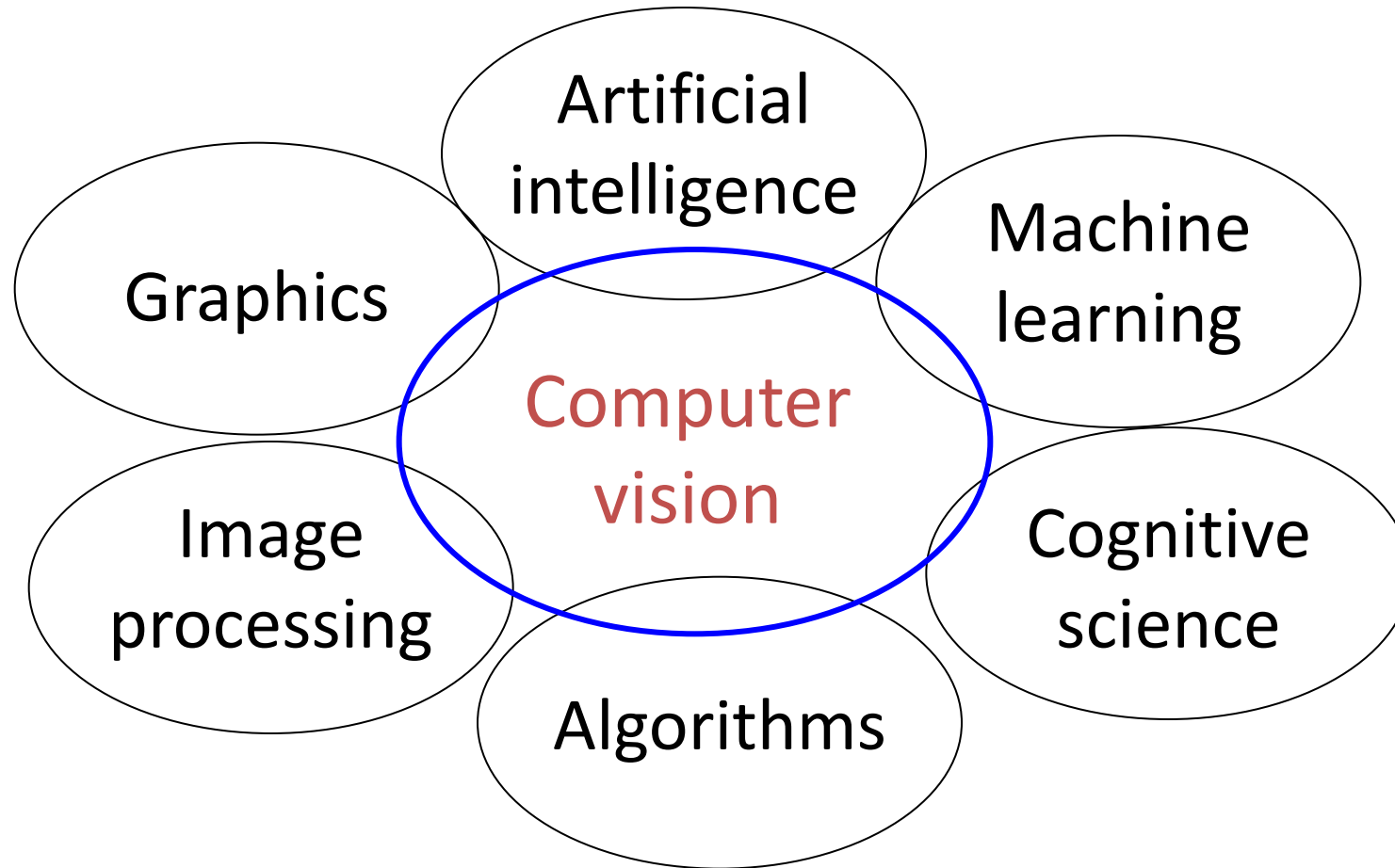
Computer Vision

- 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, people, scenes, and activities. (*perception and interpretation*)
 3. Algorithms to mine, search, and interact with visual data (*interaction*)

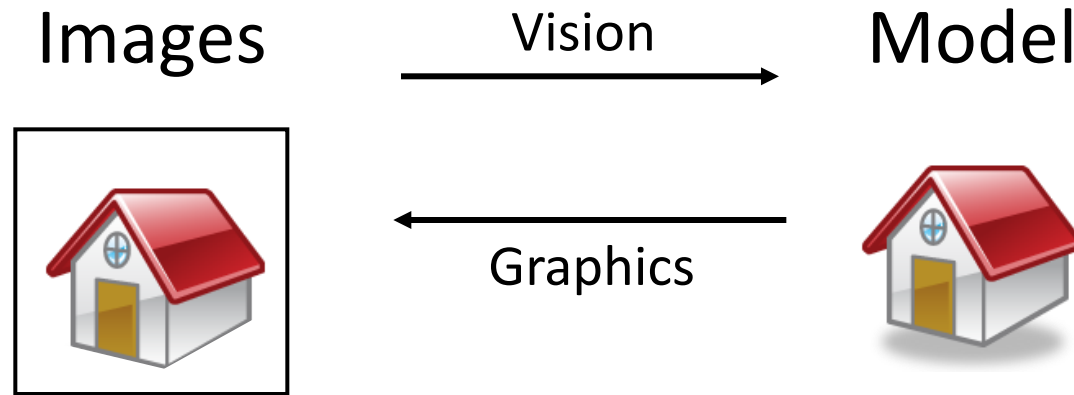
3. Interaction



Related disciplines

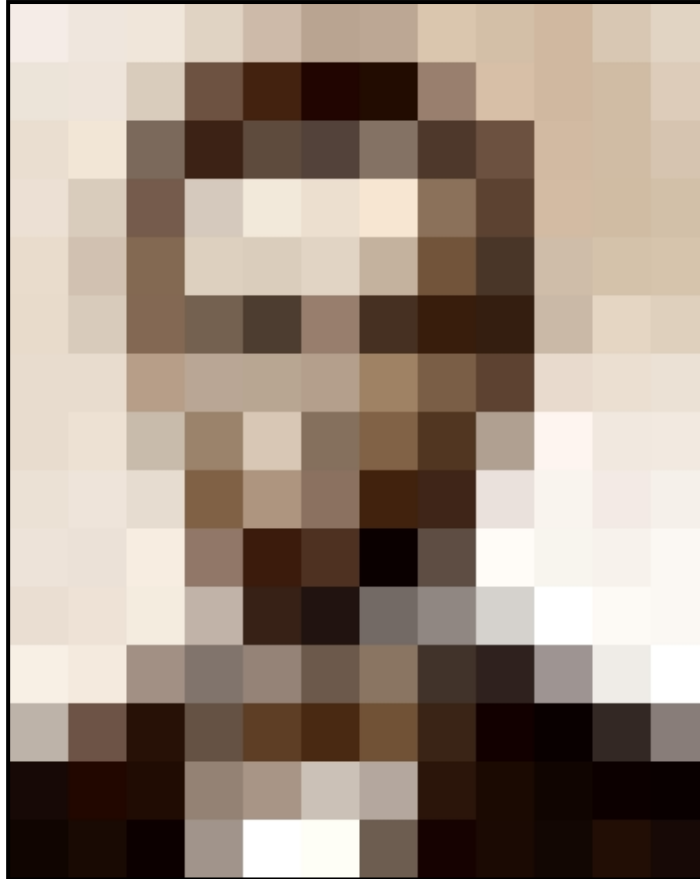


Vision and graphics

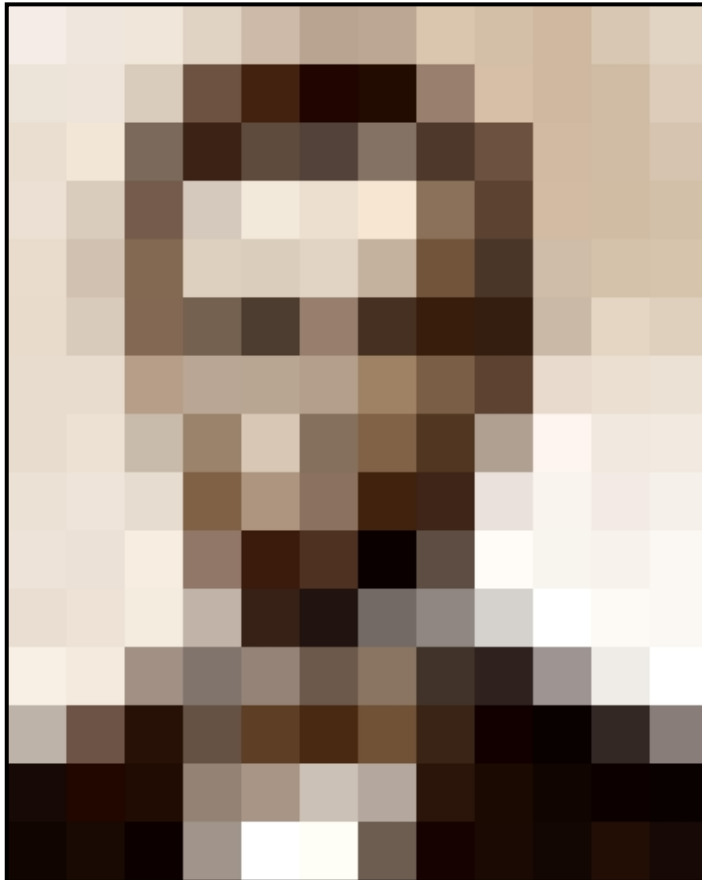


Inverse problems: analysis and synthesis.

What humans see



What computers see

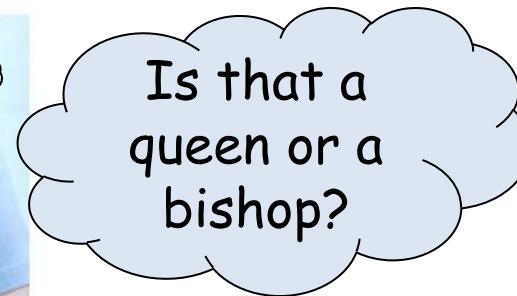


What do humans see?



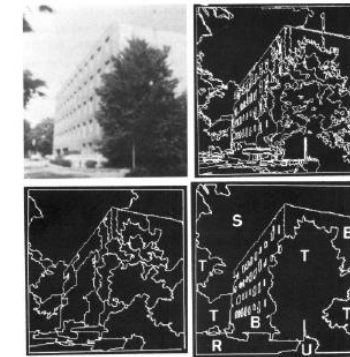
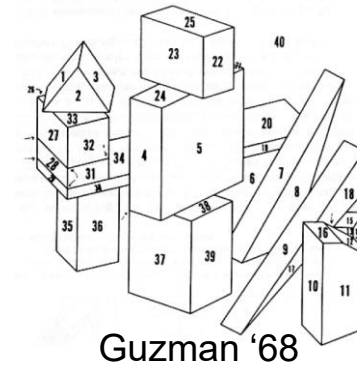
Vision is really hard

- Vision is an amazing feat of natural intelligence
 - Visual cortex occupies about 50% of Macaque brain
 - One third of human brain devoted to vision (more than anything else)



Ridiculously brief history of computer vision

- 1966: Minsky assigns computer vision as an undergrad summer project
- 1960's: interpretation of synthetic worlds
- 1970's: some progress on interpreting selected images
- 1980's: ANNs come and go; shift toward geometry and increased mathematical rigor
- 1990's: face recognition; statistical analysis in vogue
- 2000's: broader recognition; large annotated datasets available; video processing starts
- 2010's: Deep learning with ConvNets
- 2020's: Widespread autonomous vehicles?
- 2030's: robot uprising?



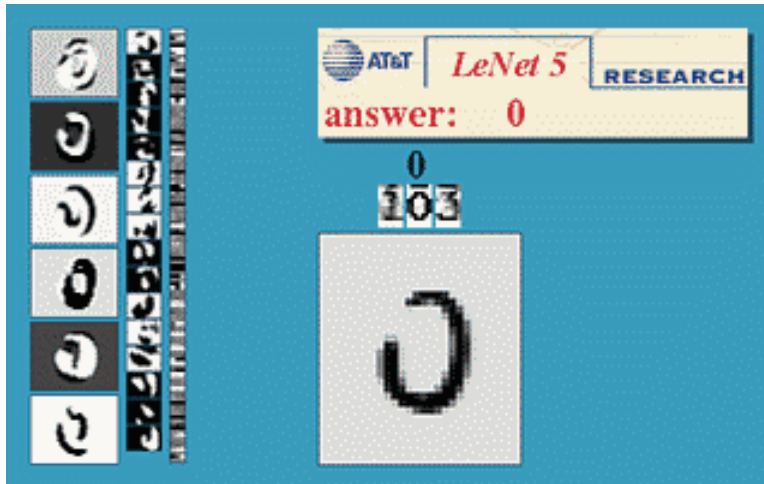
How vision is used now

- Examples of real-world applications

Optical character recognition (OCR)

Technology to convert scanned docs to text

- If you have a scanner, it probably came with OCR software



Digit recognition, AT&T labs

<http://www.research.att.com/~yann/>



License plate readers

http://en.wikipedia.org/wiki/Automatic_number_plate_recognition

Optical character recognition (OCR)

- Most US postal service mail is automatically read.
- In 1997, there were 55 offices reviewing images of 19 billion pieces of mail that OCR failed on.
- Today, there is 1 office, and they only looked at 1.2 billion pieces of mail this year.



<https://www.youtube.com/watch?v=XxCha4Kez9c>

Face detection



- Digital cameras detect faces

Vision in space

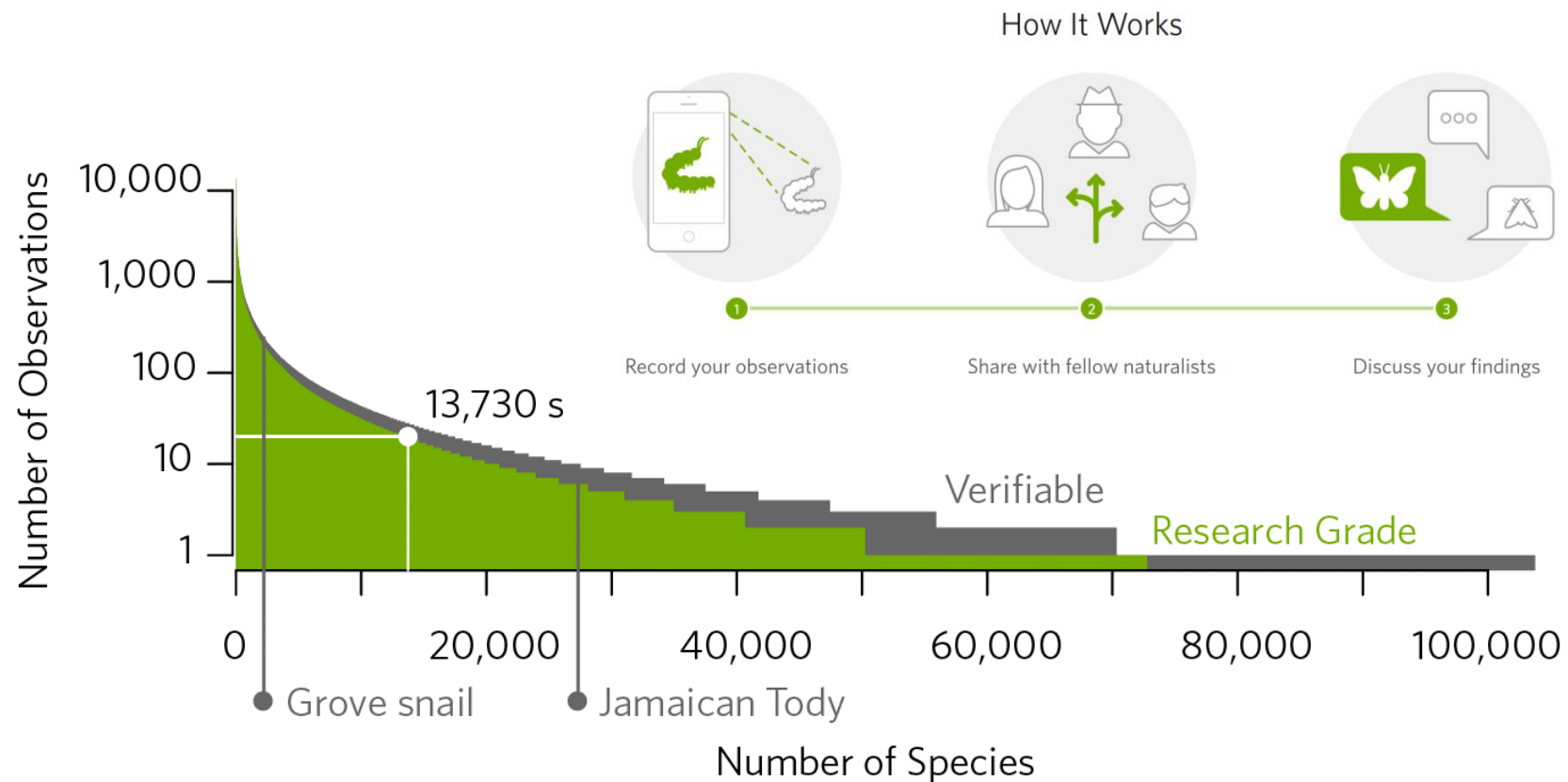


[NASA'S Mars Exploration Rover Spirit](#) captured this westward view from atop a low plateau where Spirit spent the closing months of 2007.

Vision systems (JPL) used for several tasks

- Panorama stitching
- 3D terrain modeling
- Obstacle detection, position tracking
- For more, read “[Computer Vision on Mars](#)” by Matthies et al.

iNaturalist



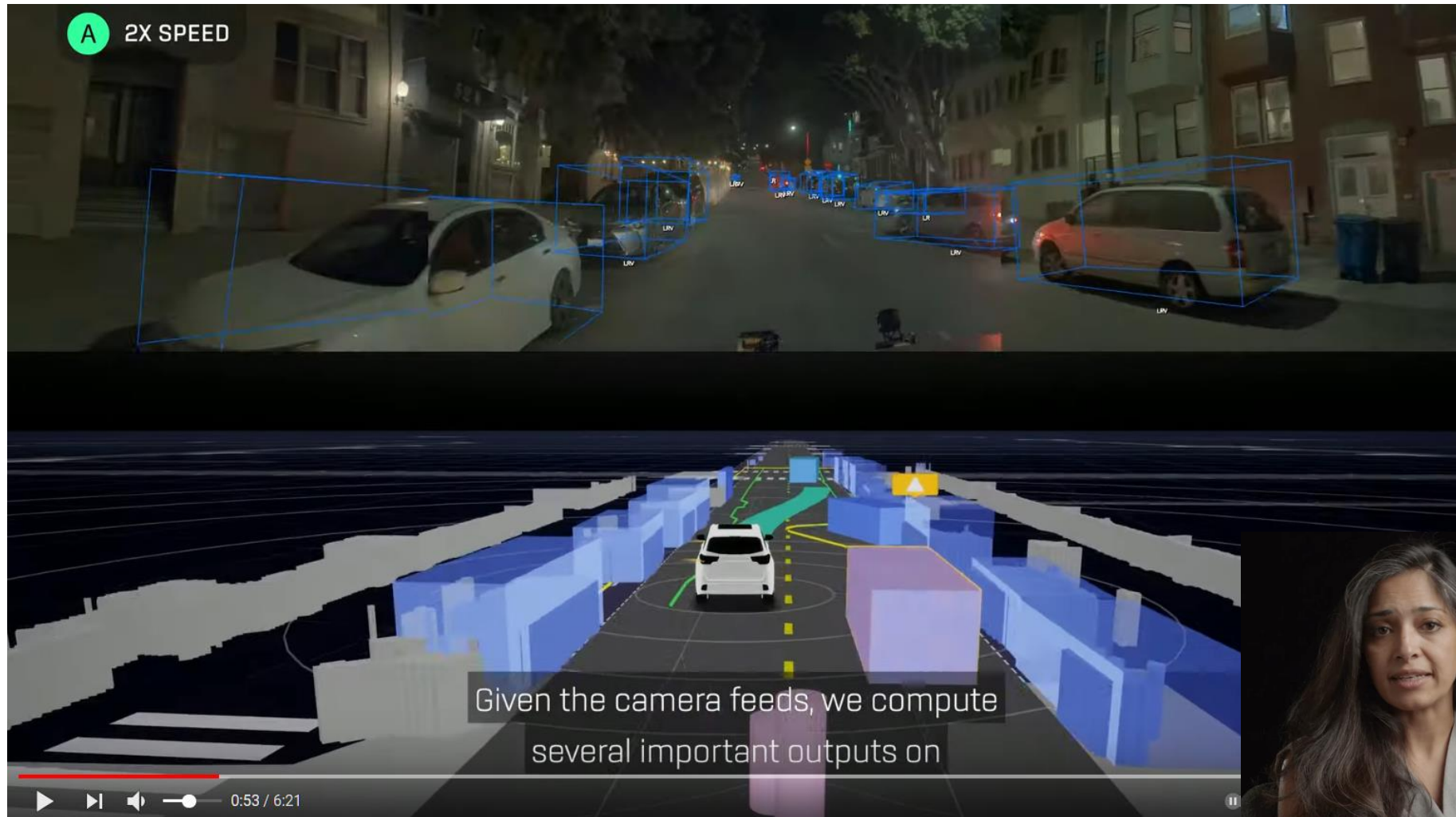
https://www.inaturalist.org/pages/computer_vision_demo

Skydio



<https://www.skydio.com/>

Zoox Computer Vision Demo



<https://www.youtube.com/watch?v=BVRMh9NO9Cs>



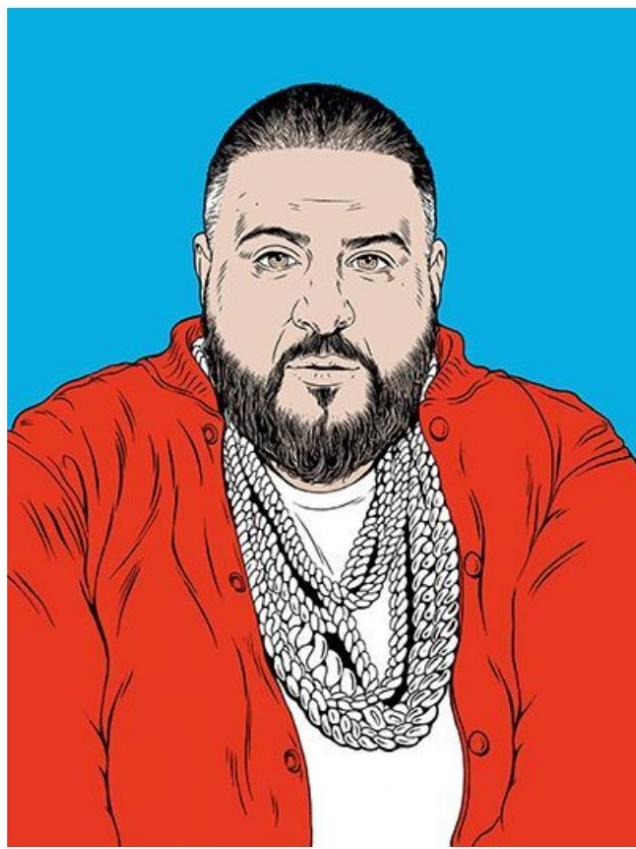
State of the art today?

With enough training data, computer vision ~~now~~ nearly matches human vision at most recognition tasks

Deep learning has been an enormous disruption to the field. More and more techniques are being “deepified”.

WIRED 100

WHO'S SHAPING THE DIGITAL WORLD?



DJ Khaled

Credit **Louise Zergaeng Pomeroy**

73. DJ Khaled

Snapchat icon; DJ and producer

Louisiana-born Khaled Mohamed Khaled, aka DJ Khaled, cut his musical chops in the early 00s as a host for Miami urban music radio WEDR. He proceeded to build a solid if not dazzling career as a mixtape DJ and music producer (he founded his label We The Best Music Group in 2008, and was appointed president of Def Jam South in 2009).

69. Geoffrey Hinton

Psychologist, computer scientist; researcher, Google Toronto

British-born Hinton has been dubbed the "godfather of deep learning". The Cambridge-educated cognitive psychologist and computer scientist started being an ardent believer in the potential of neural networks and deep learning in the 80s, when those technologies enjoyed little support in the wider AI community.

But he soldiered on: in 2004, with support from the Canadian Institute for Advanced Research, he launched a University of Toronto programme in neural computation and adaptive perception, where, with a group of researchers, he carried on investigating how to create computers that could behave like brains.

Hinton's work – in particular his algorithms that train multilayered neural networks – caught the attention of tech giants in Silicon Valley, which realised how deep learning could be applied to voice recognition, predictive search and machine vision.

The spike in interest prompted him to launch a free course on neural networks on e-learning platform Coursera in 2012. Today, 68-year-old Hinton is chair of machine learning at the University of Toronto and moonlights at Google, where he has been using deep learning to help build internet tools since 2013.

63. Yann Lecun

Director of AI research, Facebook, Menlo Park

LeCun is a leading expert in deep learning and heads up what, for Facebook, could be a hugely significant source of revenue: understanding its user's intentions.

62. Richard Branson

Founder, Virgin Group, London

Branson saw his personal fortune grow £550 million when Alaska Air bought Virgin America for \$2.6 billion in April. He is pressing on with civilian space travel with [Virgin Galactic](#).

61. Taylor Swift

Entertainer, Los Angeles

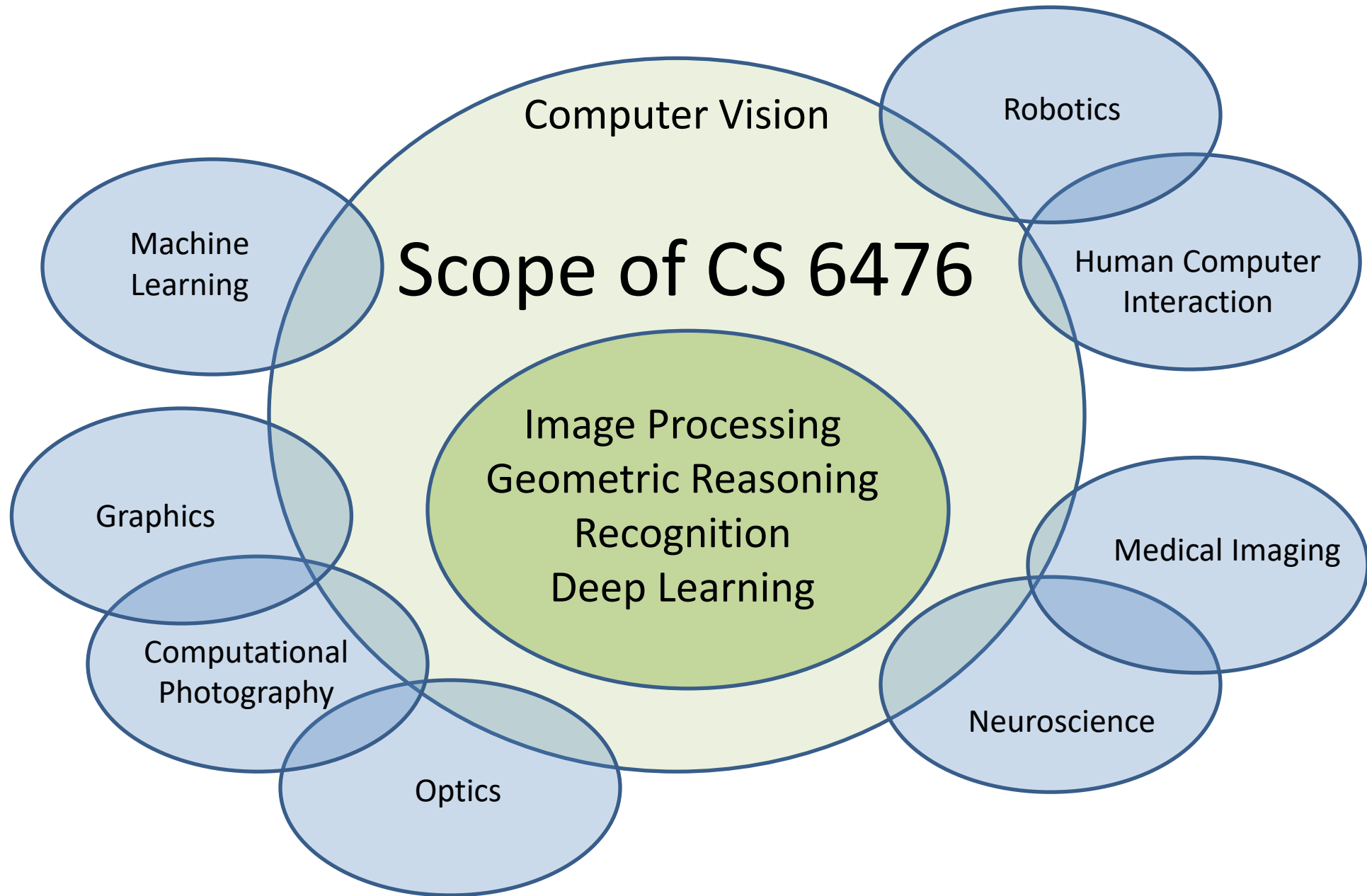


Today's Class

- ~~Who am I?~~
- ~~What is Computer Vision?~~
- Specifics of this course
- Geometry of Image Formation
- Questions

Grading

- 70% programming projects (6 total)
- 30% 2 or 3 quizzes in class
- We will have no final exam. The last project might extend into the final exam period.



Textbook

Computer Vision: Algorithms and Applications, 2nd ed.

© 2020 [Richard Szeliski](#), Facebook



<http://szeliski.org/Book/>

Prerequisites

- **Linear algebra**, basic calculus, and probability
- Experience with image processing will help but is not necessary
- Experience with Python or Python-like languages will help

Projects

- (project 0 to test environment setup and handin)
- Image Filtering and Hybrid Images
- Local Feature Matching and Ransac
- Image Classification with Deep Learning
- Semantic Segmentation with Deep Learning
- Point cloud classification with PointNet
- Neural Radiance Fields (NeRF)

You may want to buy a month or two of Google Colab Pro near the end of the semester

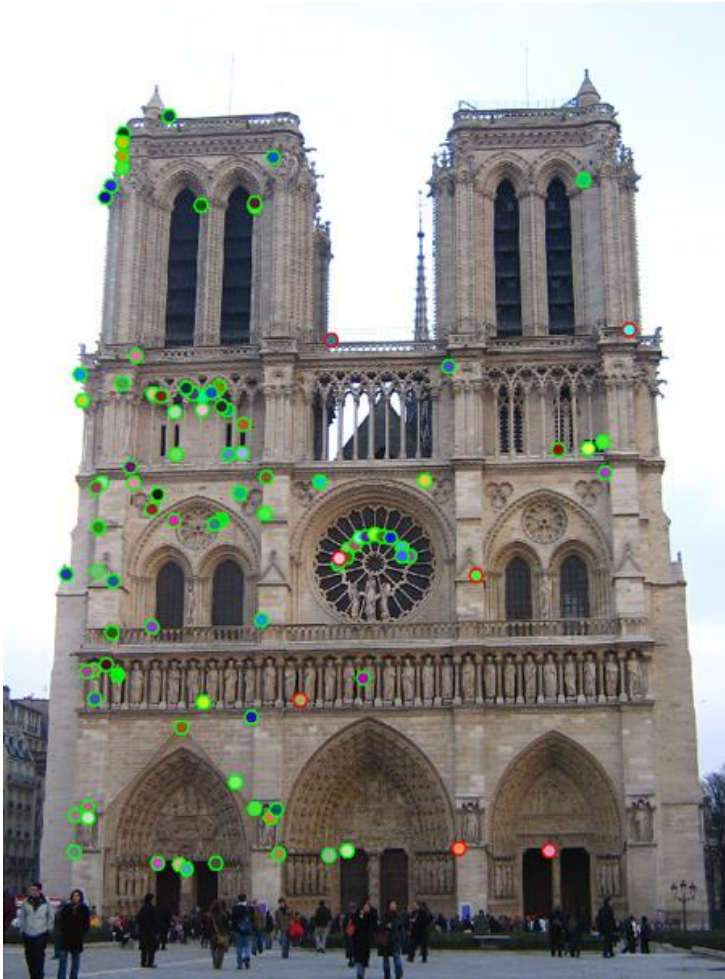
Proj1: Image Filtering and Hybrid Images

- Implement image filtering to separate high and low frequencies
- Combine high frequencies and low frequencies from different images to create an image with scale-dependent interpretation



Proj2: Local Feature Matching

- Implement interest point detector, SIFT-like local feature descriptor, and simple matching algorithm.



Course Syllabus (tentative)

<https://faculty.cc.gatech.edu/~hays/compvision/>

Code of Conduct

Your work must be your own. We'll look for cheating. Don't talk at the level of code with other students.

Canvas Survey

Today's Class

- ~~Who am I?~~
- ~~What is Computer Vision?~~
- ~~Specifics of this course~~
- Geometry of Image Formation
- Questions