

Anticipating the Unseen and Unheard for Embodied Perception

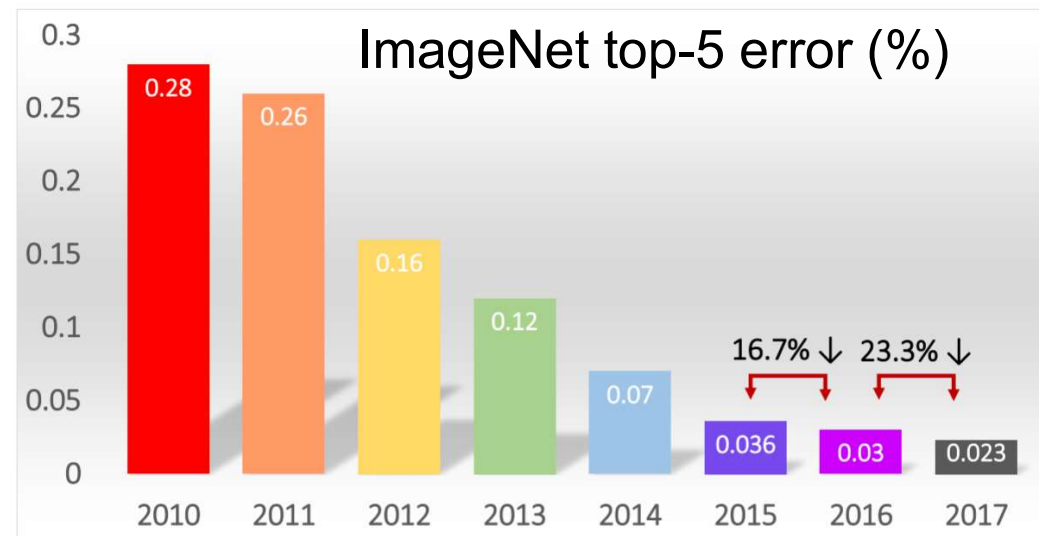
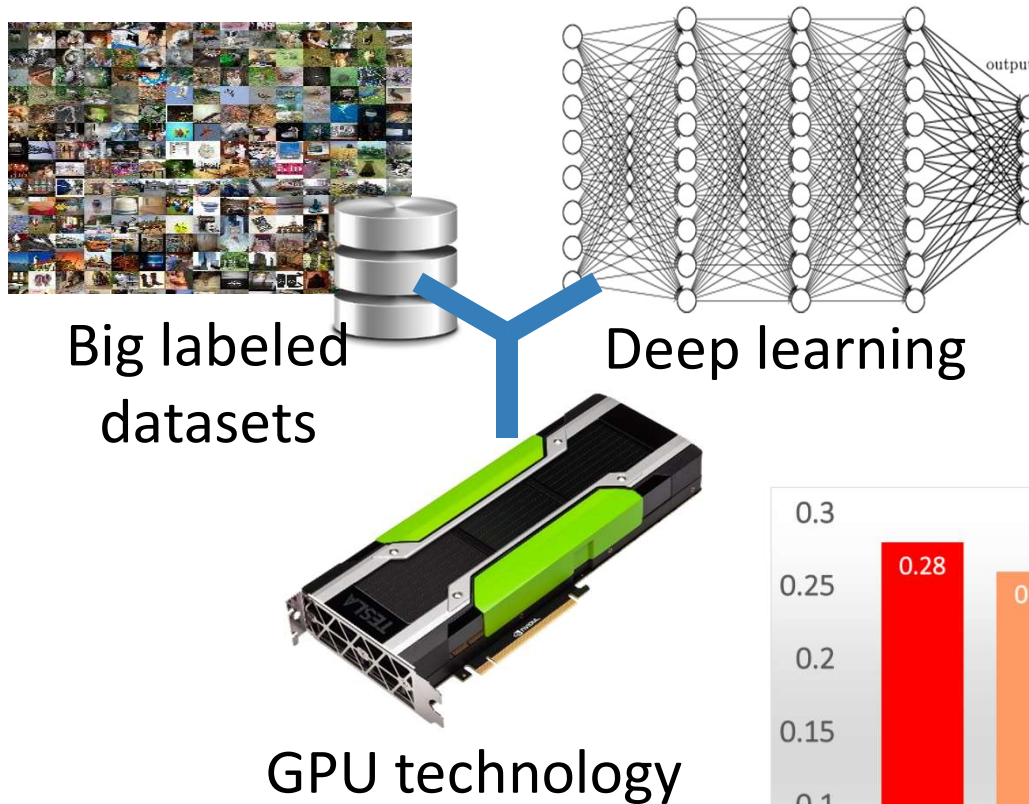
Kristen Grauman

University of Texas at Austin

Facebook AI Research

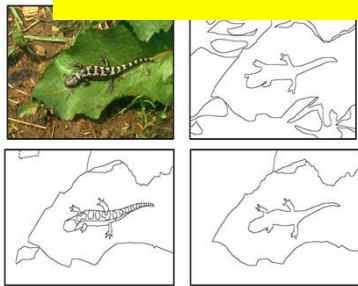


Visual recognition: significant recent progress



The Web photo perceptual experience

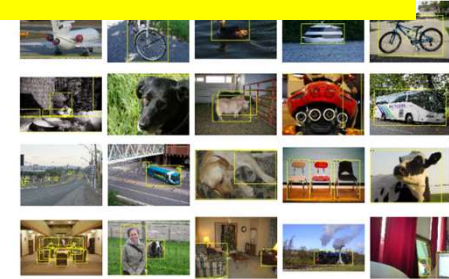
A “disembodied” well-curated moment in time



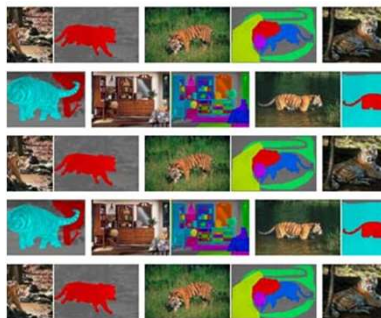
BSD (2001)



Caltech 101 (2004), Caltech 256 (2006)



PASCAL (2007-12)



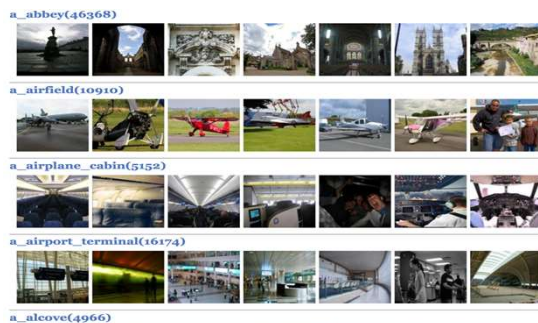
LabelMe (2007)



ImageNet (2009)



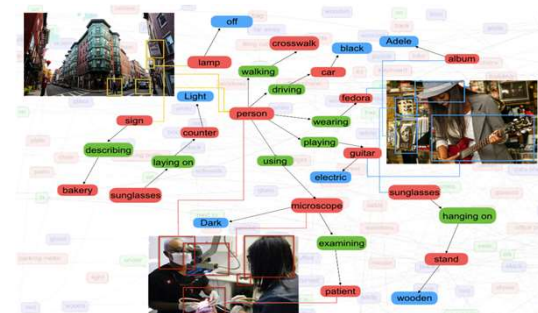
SUN (2010)



Places (2014)



MS COCO (2014)



Visual Genome (2016)

Egocentric perceptual experience

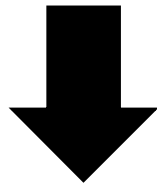
A tangle of relevant and irrelevant multi-sensory information



Big picture goal: Embodied perception

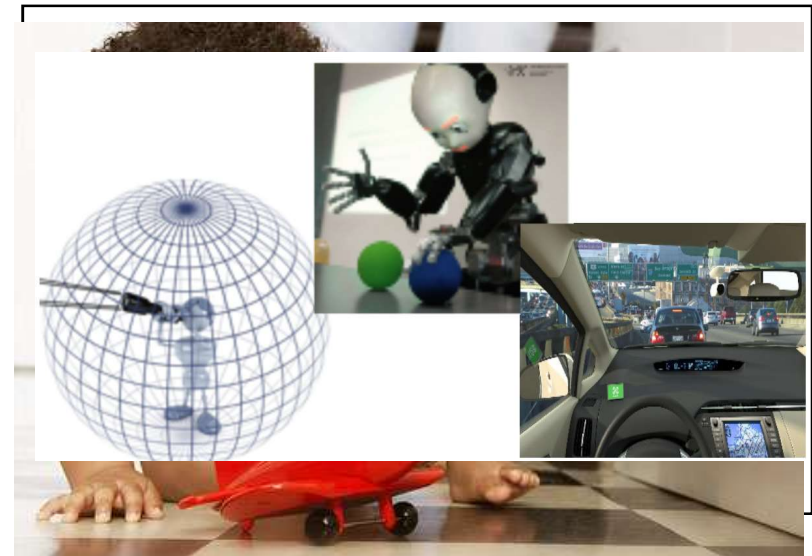
Status quo:

Learning and inference with
“disembodied” snapshots.



On the horizon:

Visual learning in the
context of **action, motion,**
and **multi-sensory**
observations.



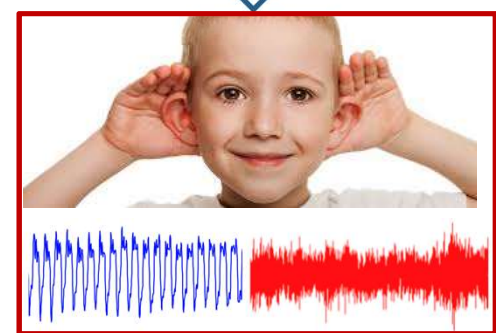
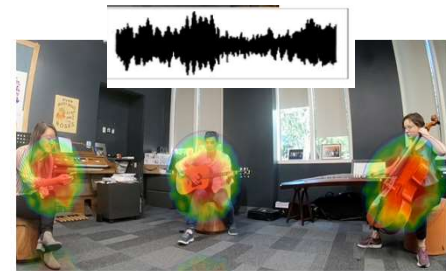
Anticipating the unseen and unheard



**Look-around
policies**



**Affordance
learning**

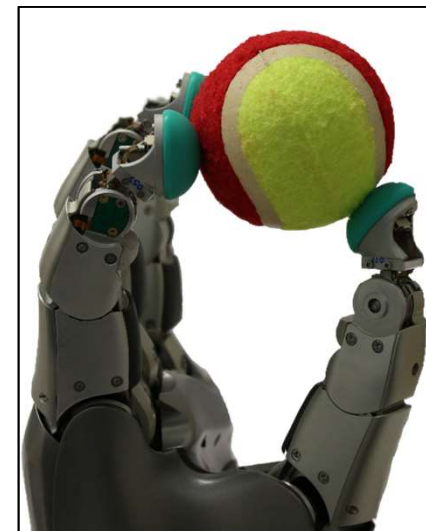


**Audio-visual
learning**

Towards embodied perception

Active perception

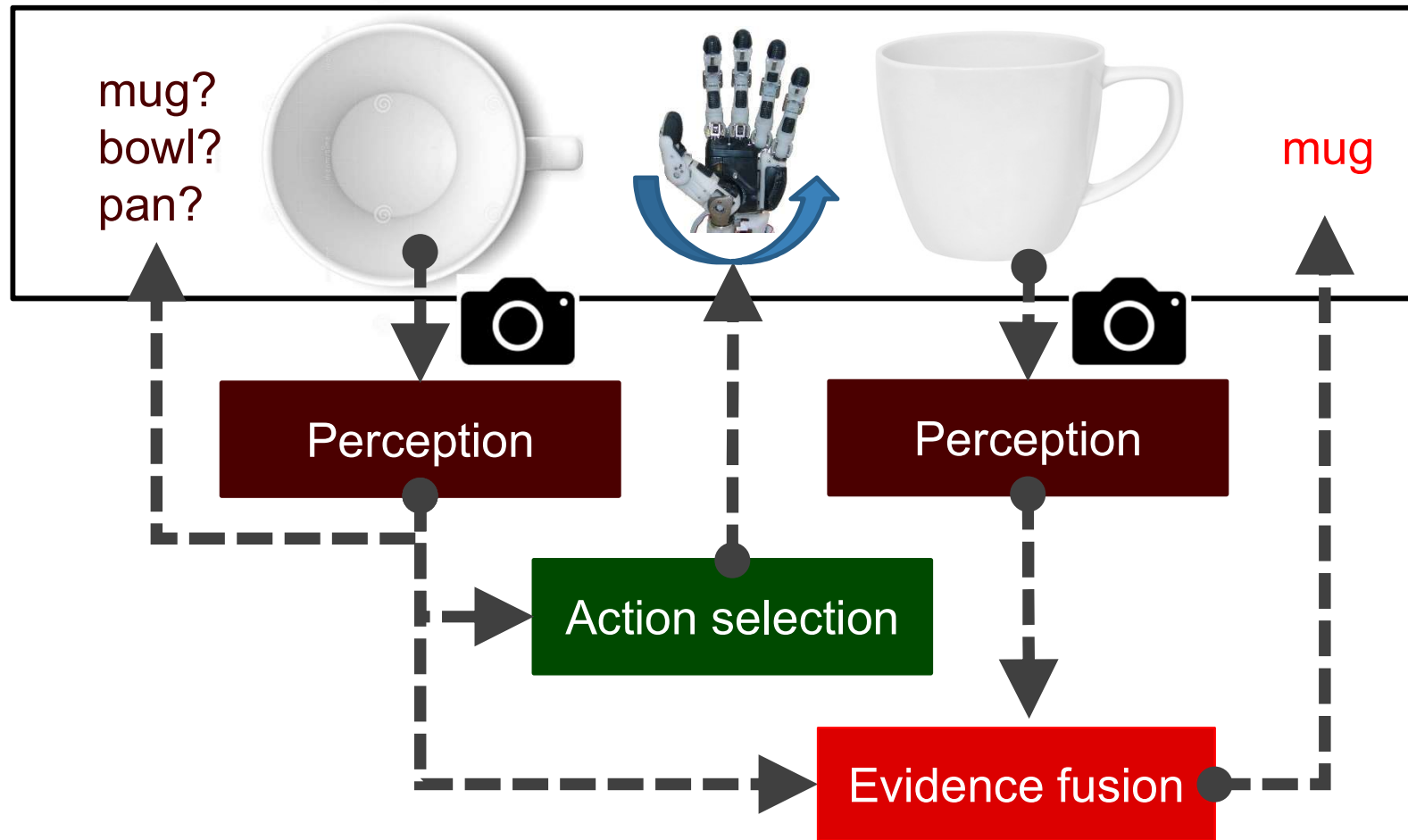
From learning *representations* to learning *policies*



*Bajcsy 1985, Aloimonos 1988, Ballard 1991, Wilkes 1992, Dickinson 1997,
Schiele & Crowley 1998, Tsotsos 2001, Denzler 2002, Soatto 2009,
Ramanathan 2011, Borotschnig 2011, ...*

End-to-end active recognition

Main idea: Deep reinforcement learning approach that anticipates visual changes as a function of egomotion



End-to-end active recognition

Predicted
label:



T=1



T=2



T=3

[Jayaraman and Grauman, ECCV 2016, PAMI 2018]

Goal: Learn to “look around”



recognition

vs.



reconnaissance



search and rescue

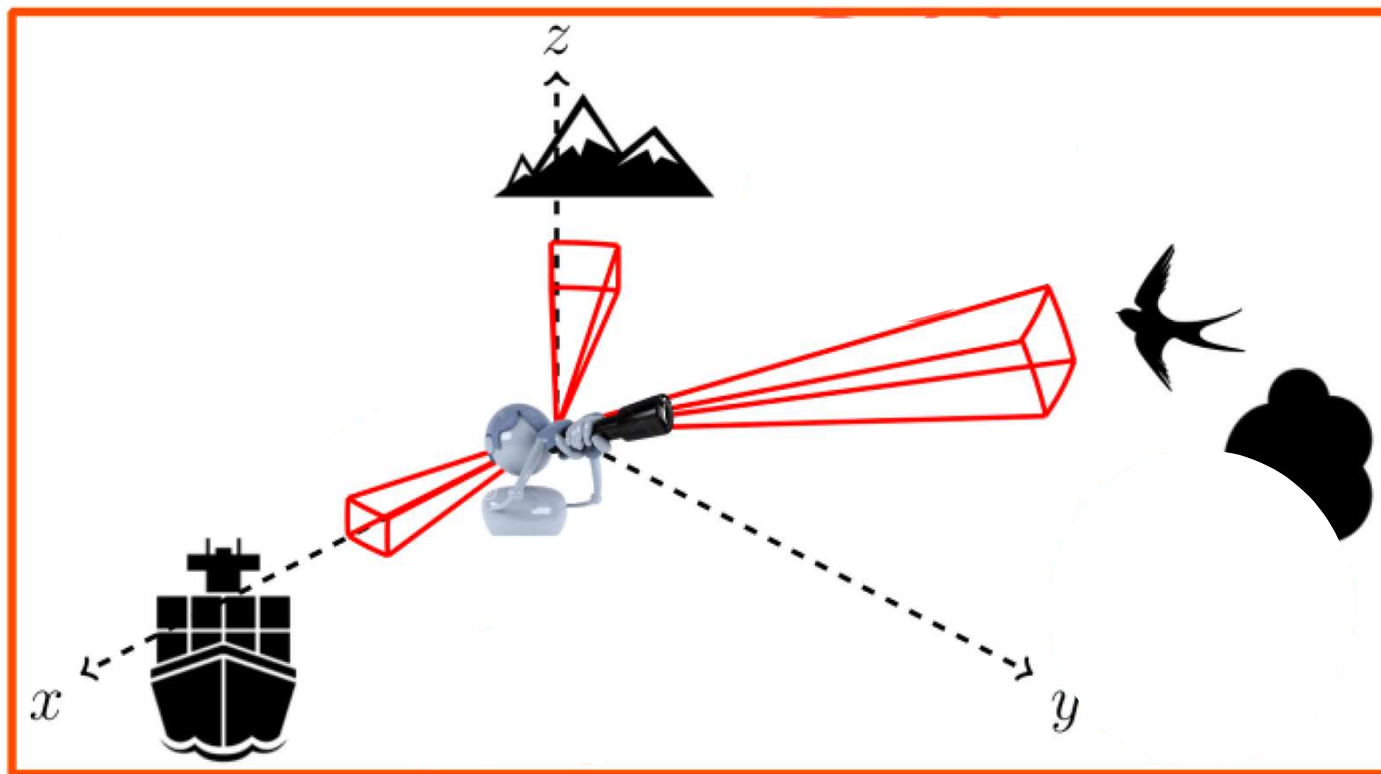
task predefined

task unfolds dynamically

Can we learn **look-around policies** for visual agents that are curiosity-driven, exploratory, and generic?

Key idea: Active observation completion

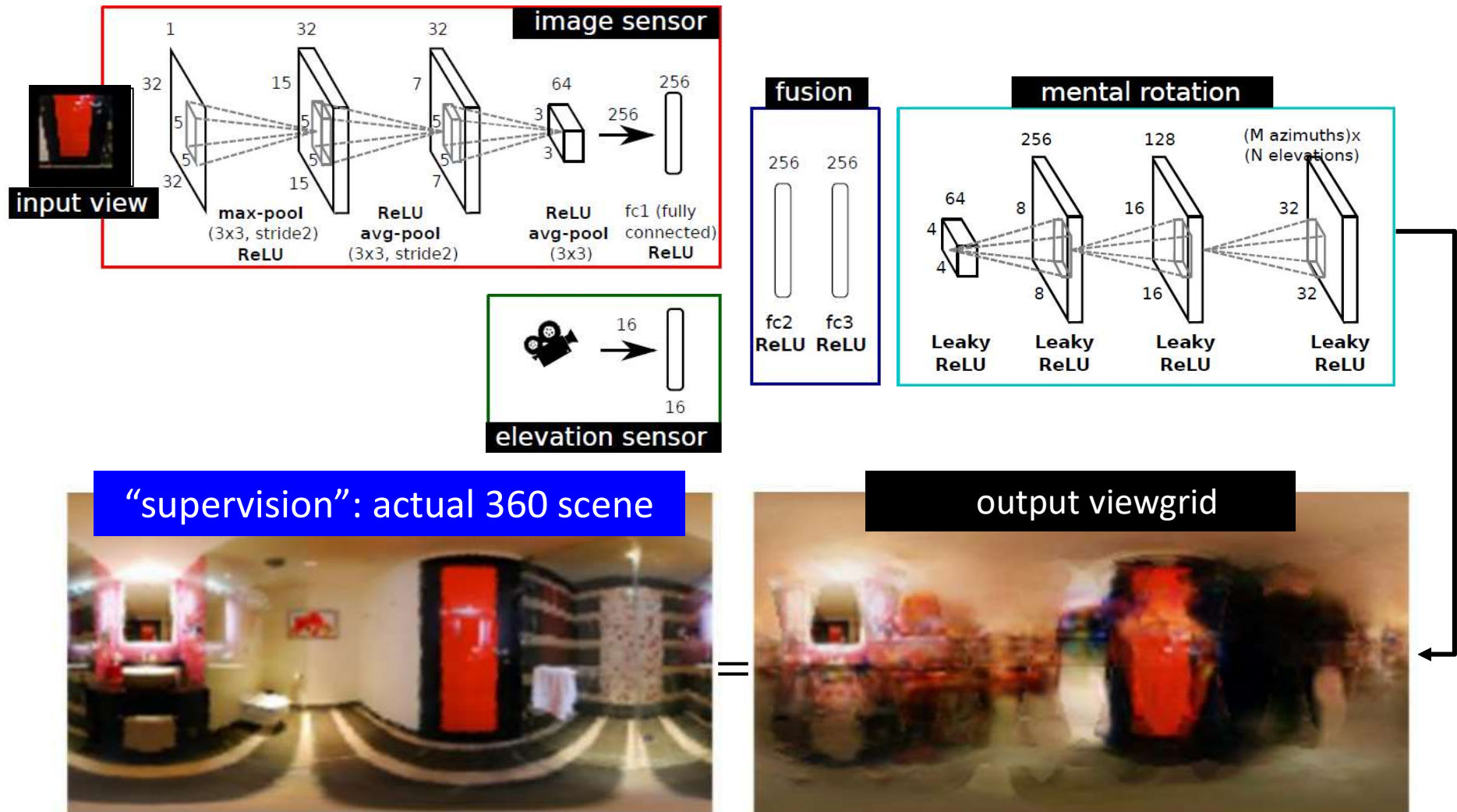
Completion objective: Learn policy for efficiently inferring (pixels of) all yet-unseen portions of environment



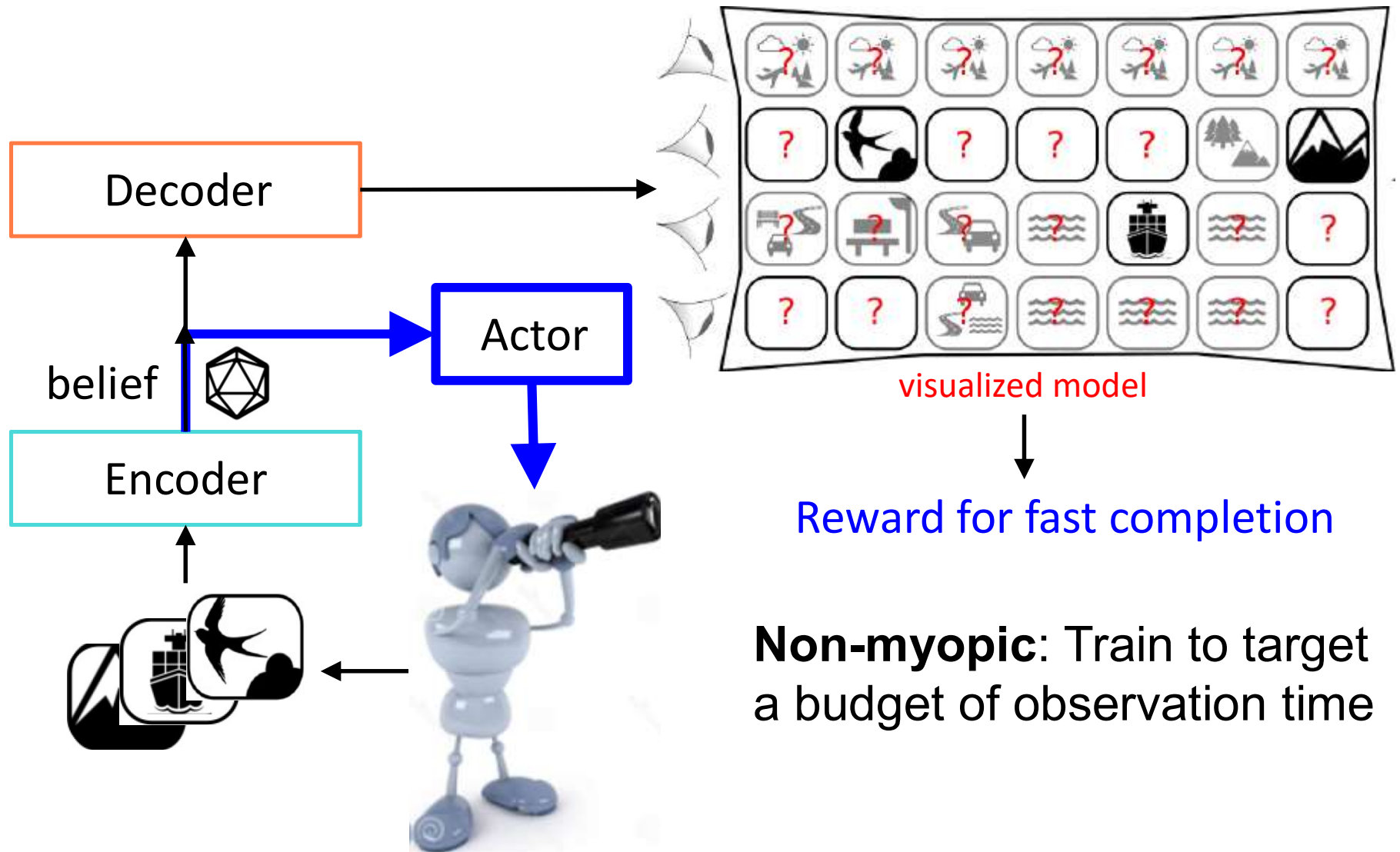
Agent must choose where to look *before* looking there.

Completing unseen views

Encoder-decoder model to infer unseen viewpoints



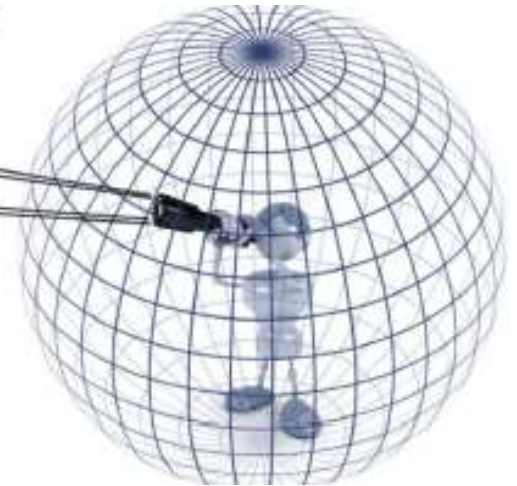
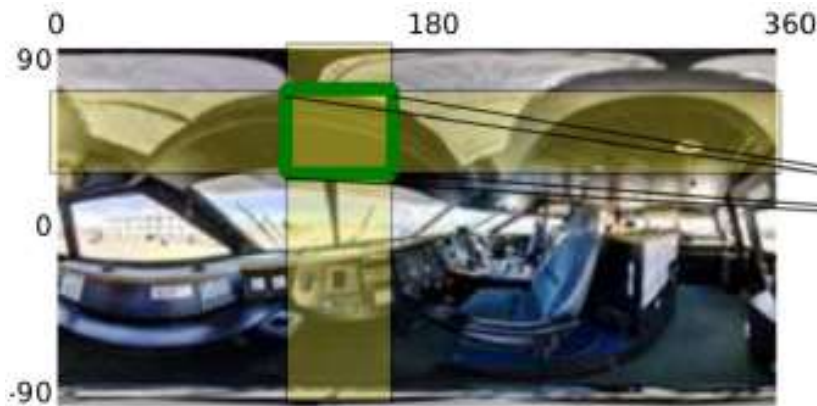
Actively selecting observations



Two scenarios

Where to look next?

agent

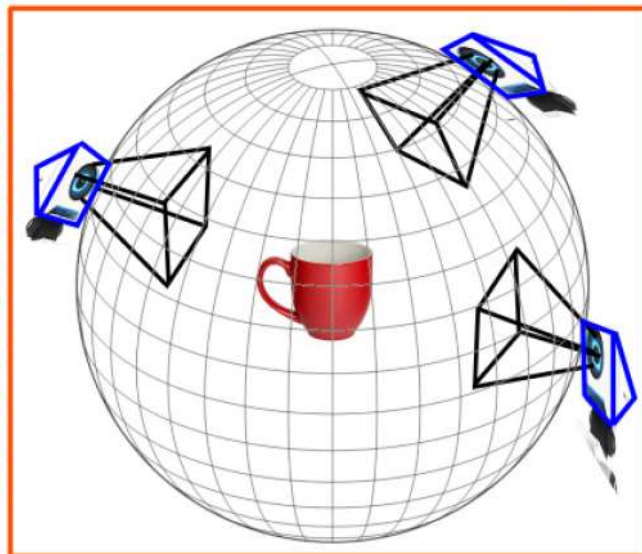


How to manipulate?

agent



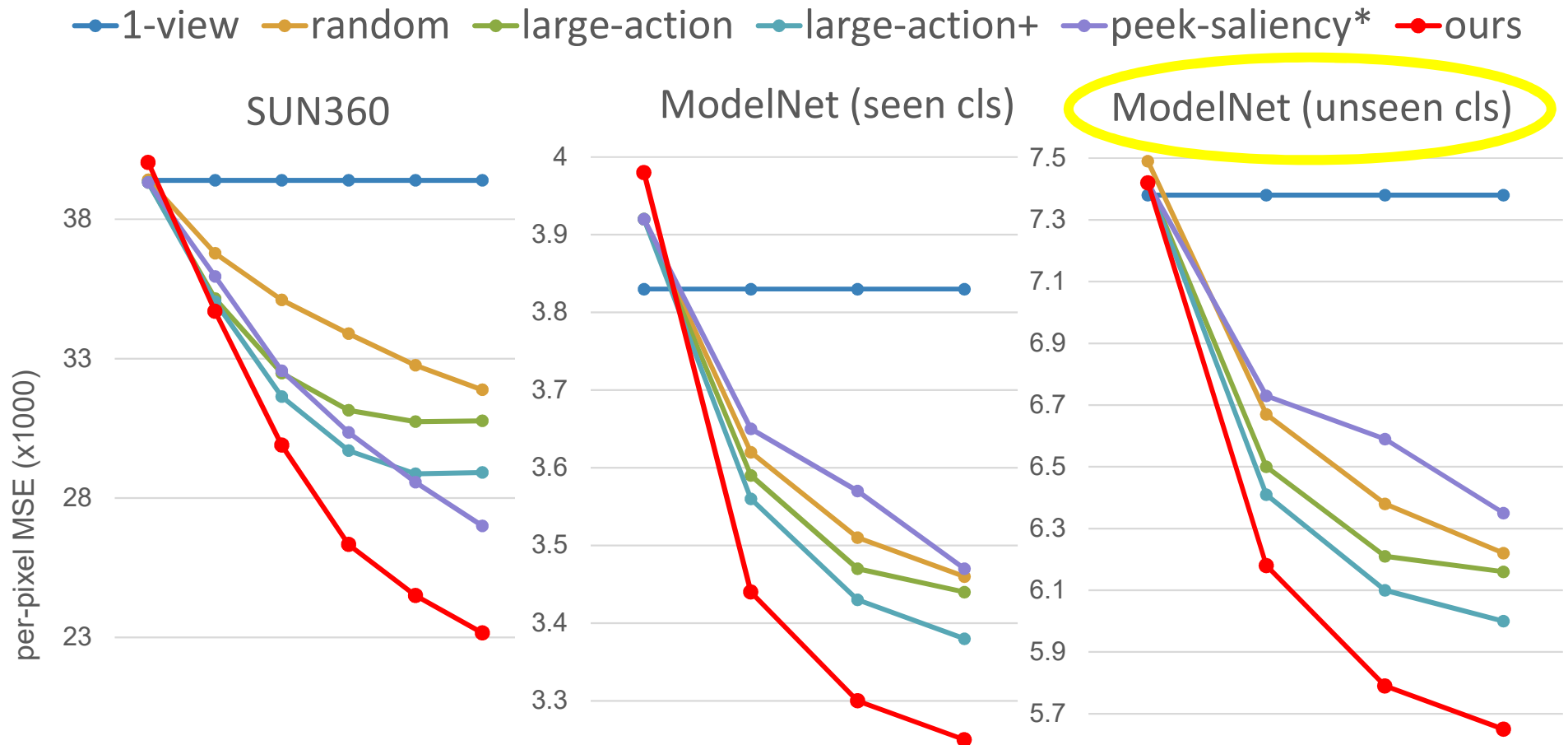
environment



observations



Active “look around” results



Learned active look-around policy: quickly grasp environment independent of a specific task

Active “look around” results

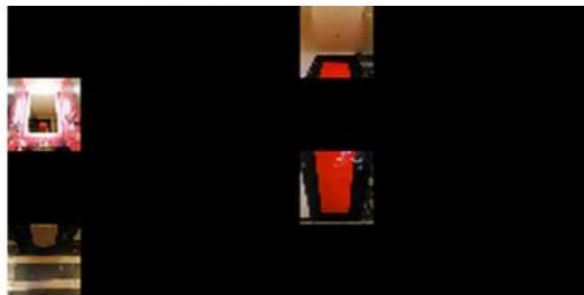
Ground Truth



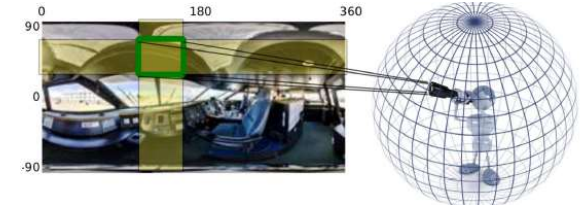
Agent inputs



Reconstruction



Active “look around”

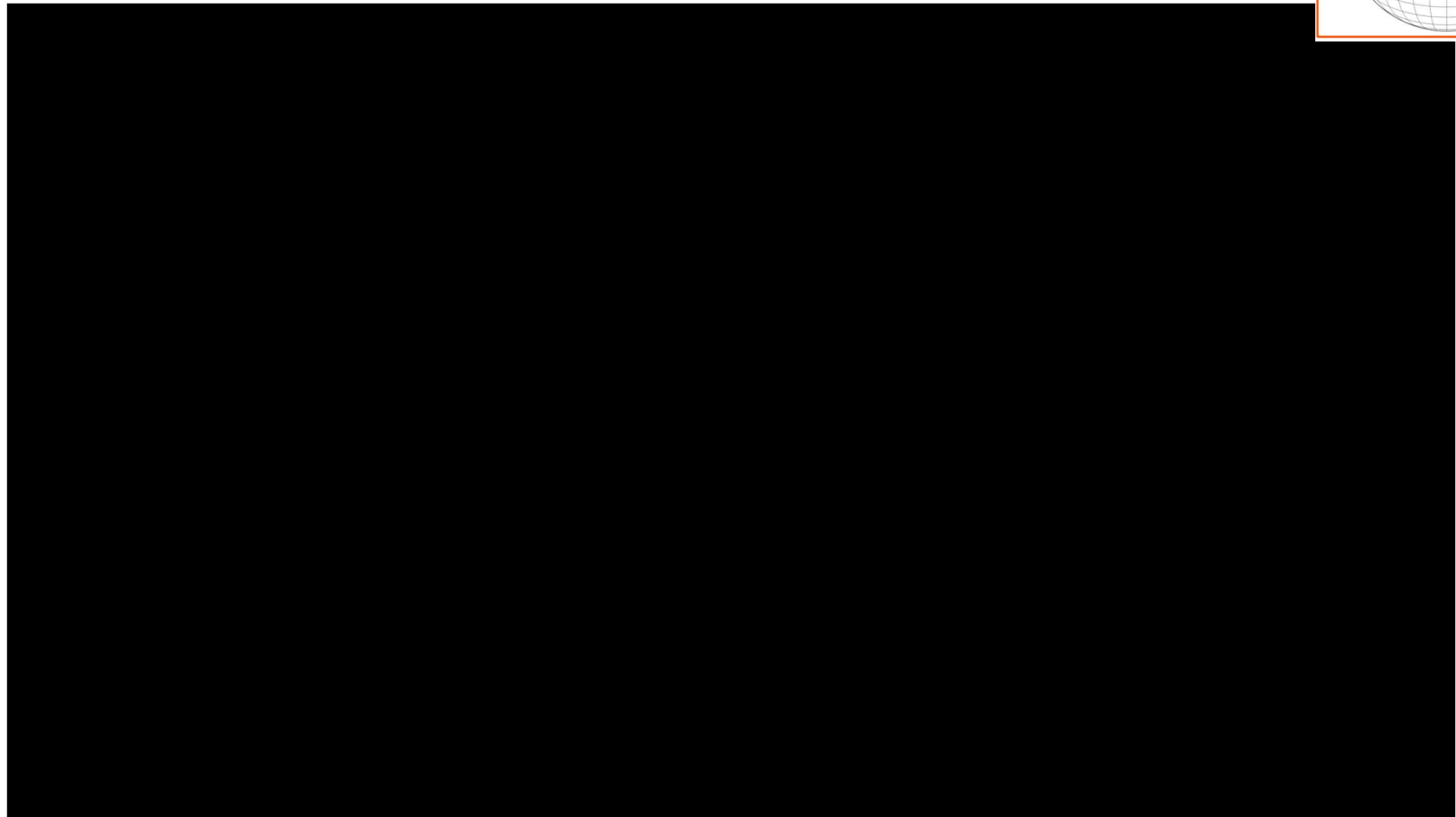
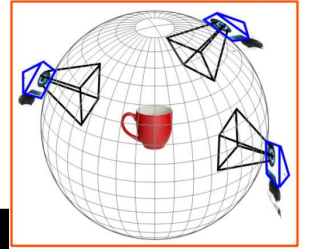


Example #1

Agent's mental model for 360 scene evolves with
actively accumulated glimpses

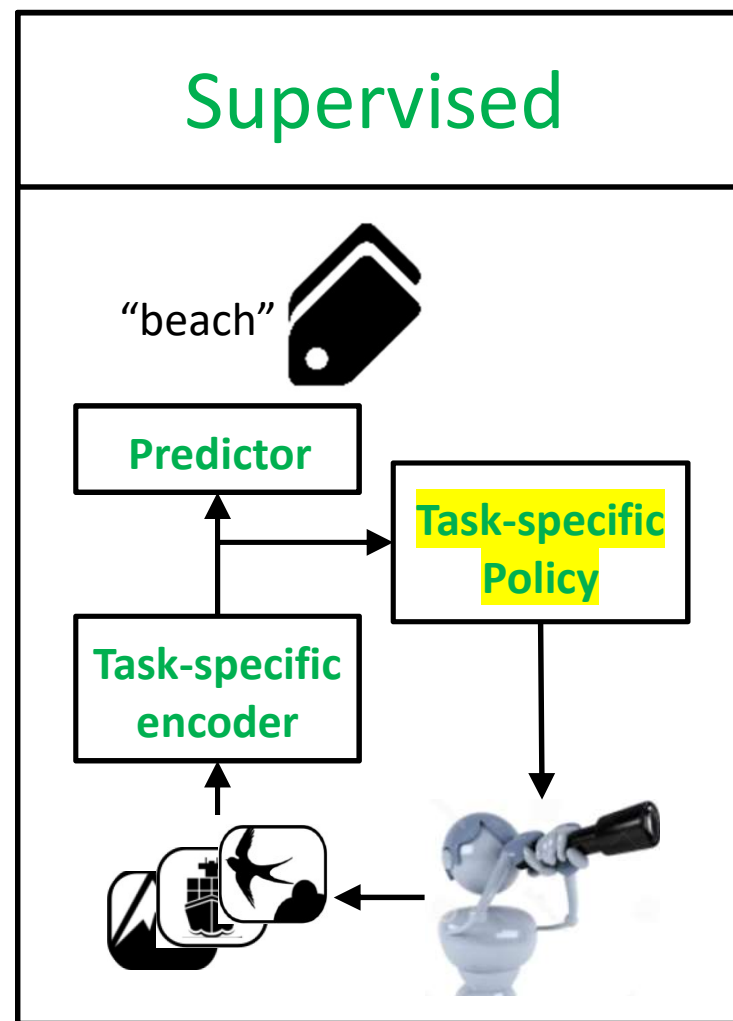
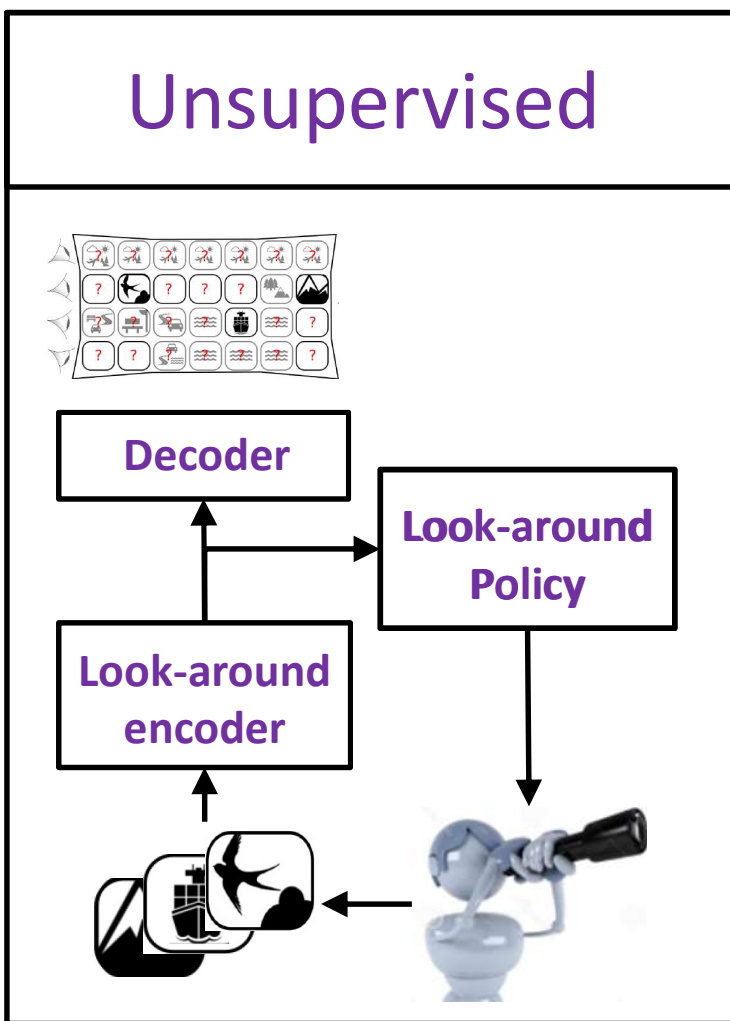
Jayaraman and Grauman, CVPR 2018; Ramakrishnan & Grauman, ECCV 2018

Active “look around”



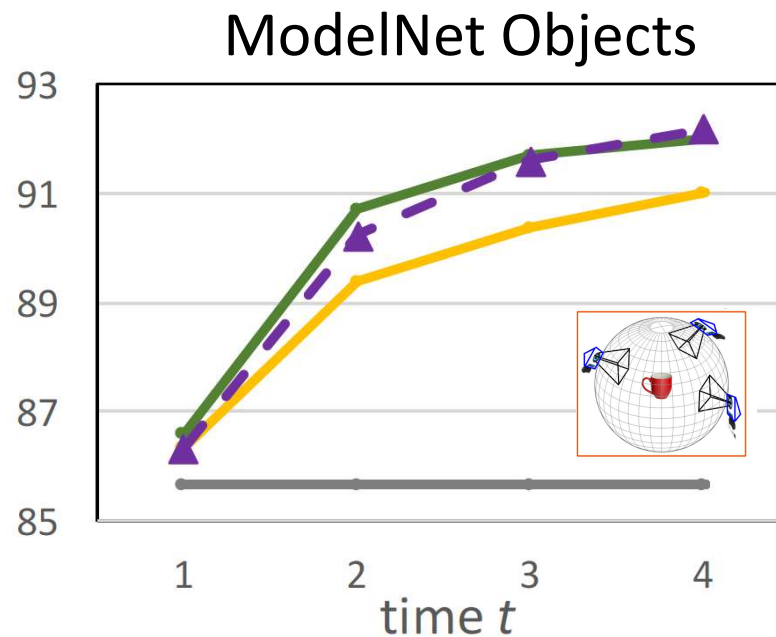
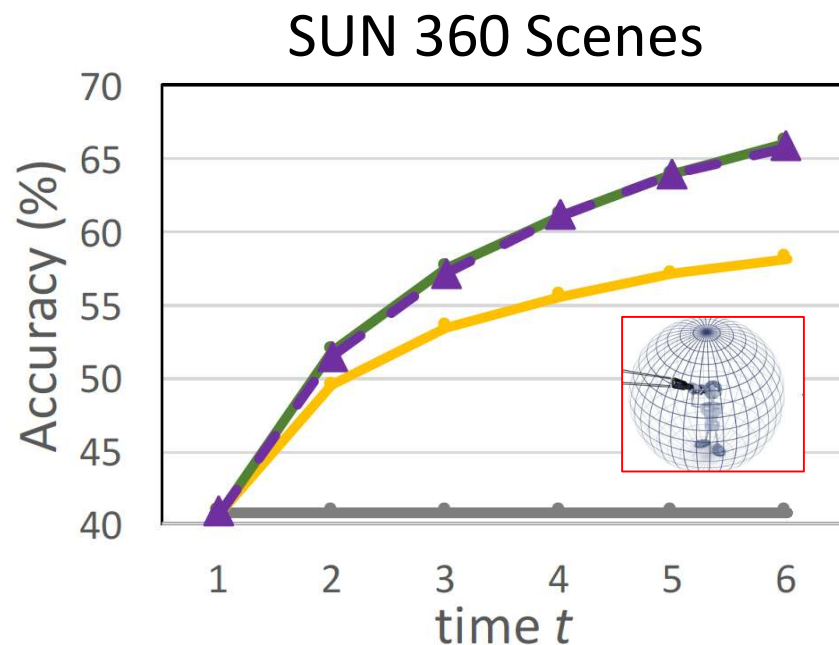
Agent's mental model for 3D object evolves with
actively accumulated glimpses

Look-around policy transfer



Plug observation completion policy in for **new** task

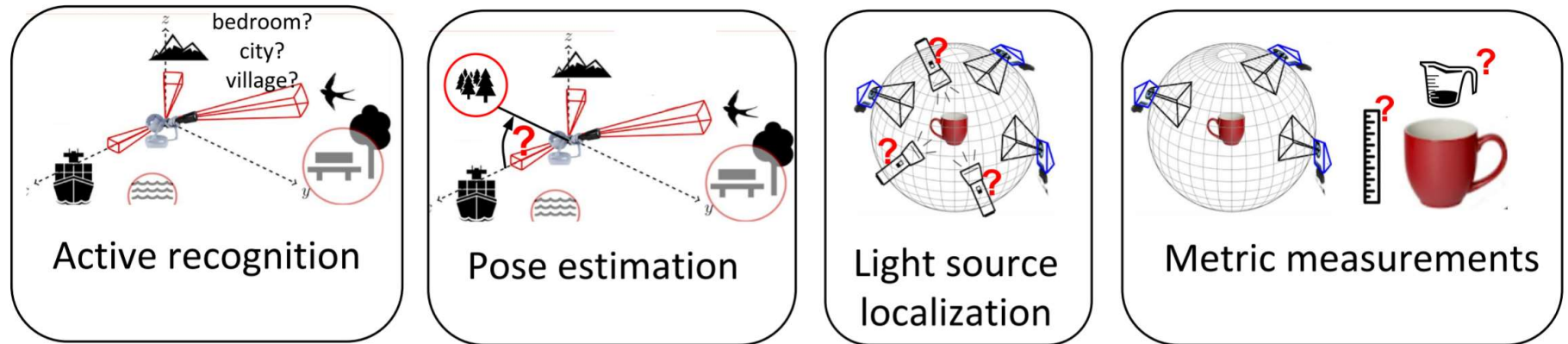
Look-around policy transfer



— 1-view — random-policy — sup-policy ▲ ours (policy transfer)

Plug **Unsupervised exploratory policy** approaches task
supervised task-specific policy accuracy!

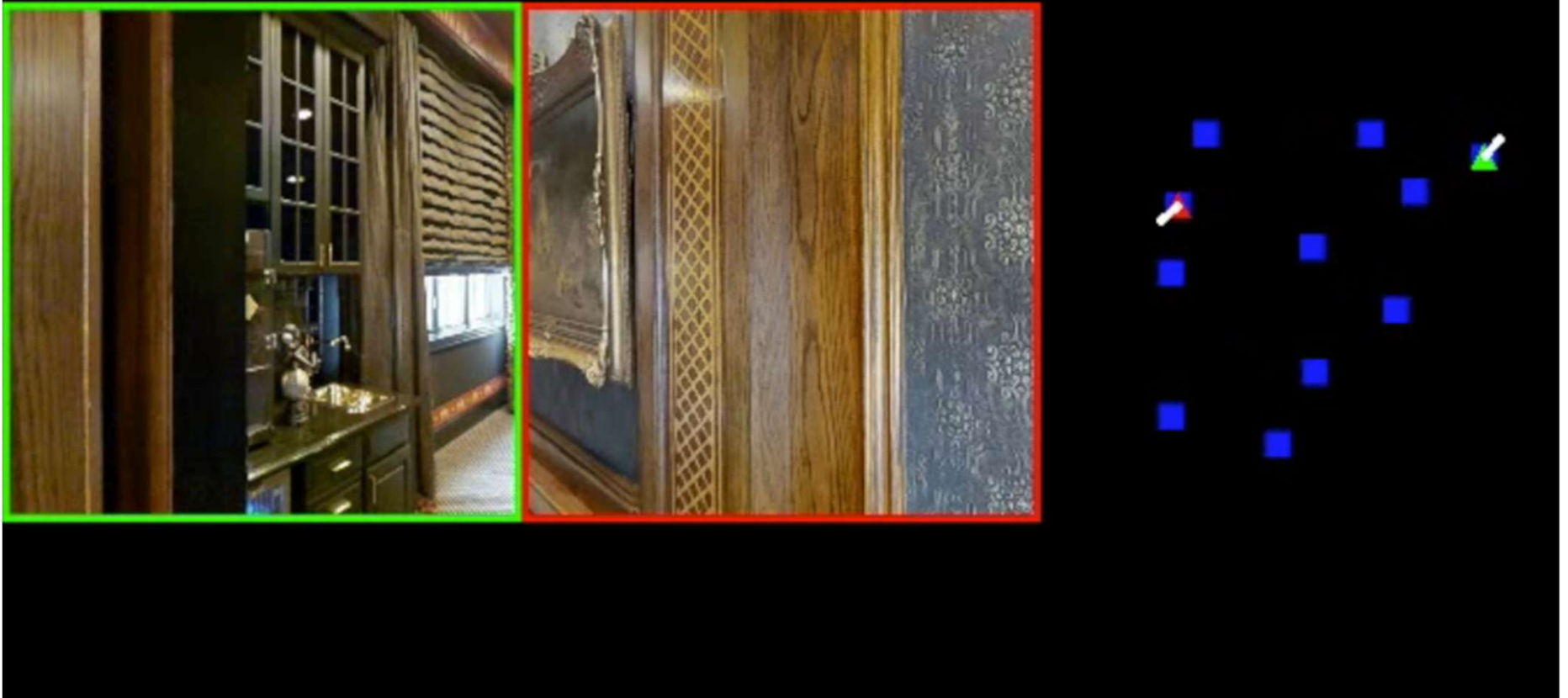
Look-around policy transfer



Multiple perception tasks

| Method | Task | SUN360 | | | ModelNet | | |
|----------------|------|---------------------------------------|--|---|---------------------------------------|--|---|
| | | Active recogn. Accuracy \uparrow | Pose estimation. AE azimuth. \downarrow | Pose estimation. AE elev. \downarrow | Active recogn. Accuracy \uparrow | Light source loc. Accuracy \uparrow | Surface area RMSE $\times 100\downarrow$ |
| one-view | | 51.94 | 75.74 | 30.32 | 83.60 | 58.74 | 21.22 |
| rnd-actions | | 62.90 | 66.18 | 19.53 | 88.46 | 72.97 | 19.04 |
| large-action | | 63.73 | 67.57 | 19.94 | 89.05 | 75.14 | 18.38 |
| peek-saliency | | 64.20 | 65.46 | 19.76 | 88.74 | 71.19 | 18.85 |
| supervised | | 68.21 | 51.36 | 9.81 | 88.58 | 86.30 | 18.43 |
| lookaround | | 68.89 | 50.00 | 9.94 | 89.00 | 83.29 | 18.82 |
| lookaround+spl | | 69.32 | 47.13 | 9.36 | 89.38 | 83.08 | 18.14 |

Look-around policy transfer



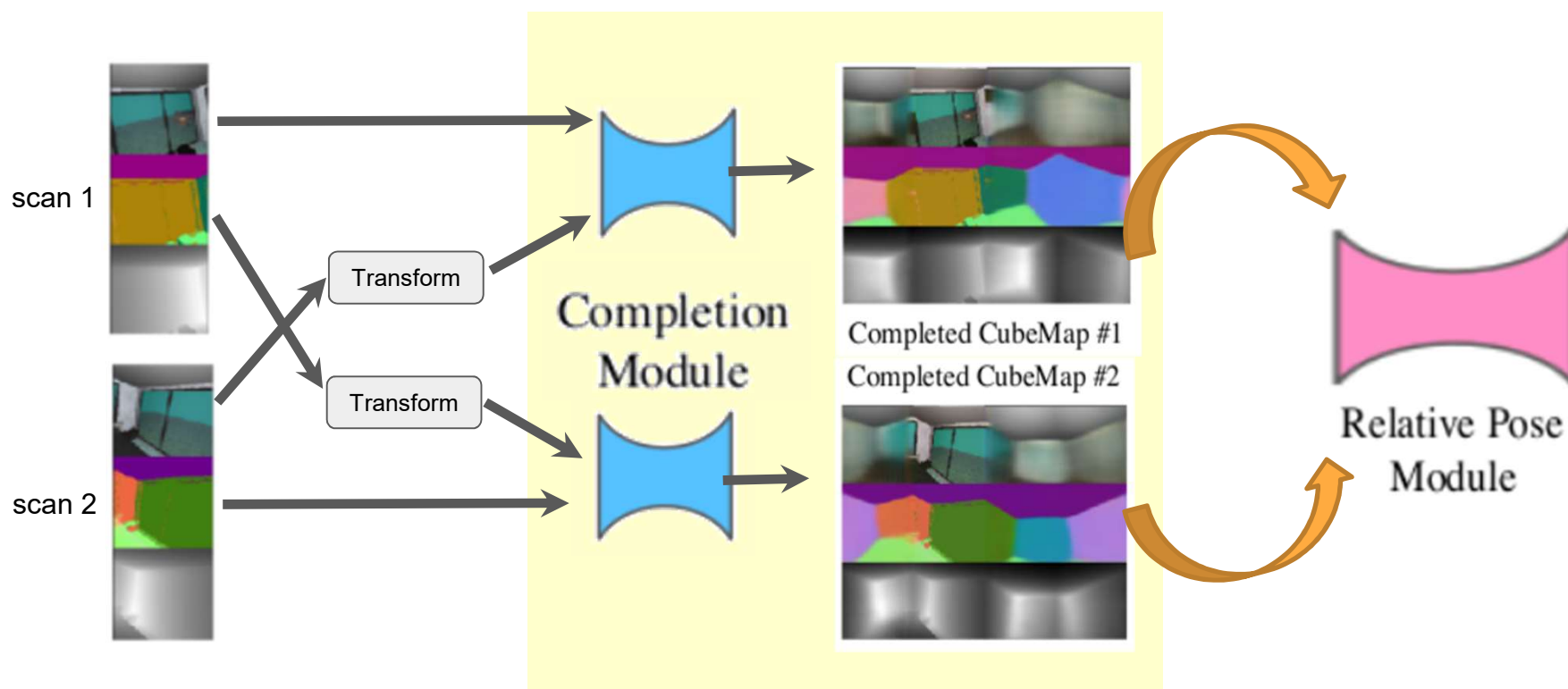
Agent navigates 3d environment leveraging active exploration

Kristen Grauman

Extreme relative pose from RGB-D scans

Input: Pair of RGB-D scans with little or *no* overlap

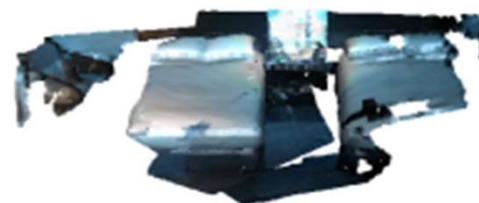
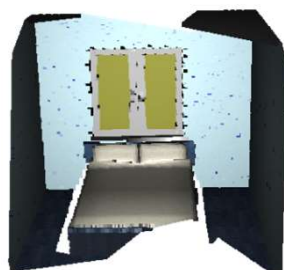
Output: Rigid transformation (R, t) that separates them



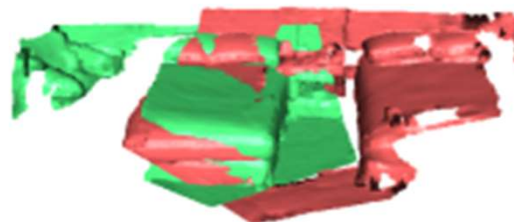
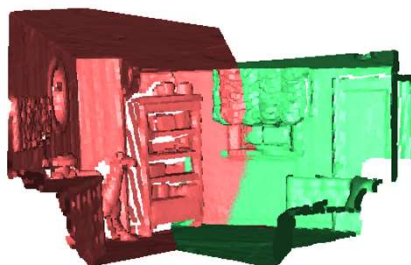
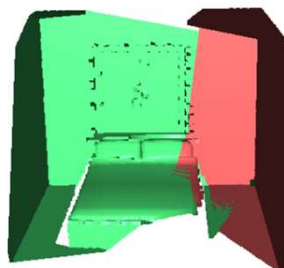
Approach: Alternate between completion and matching

Extreme relative pose from RGB-D scans

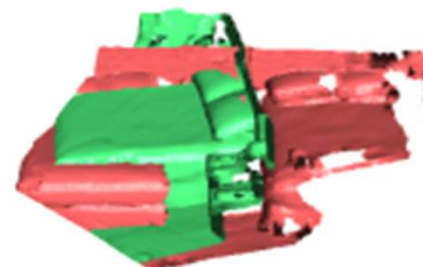
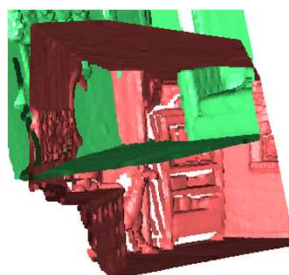
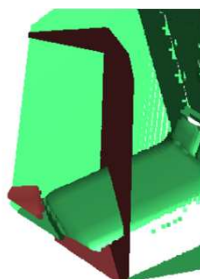
GT



Ours



4PCS



Outperform existing methods on SUNCG / Matterport / ScanNet, particularly for small overlap case (10% to 50%)

360° video: a “look around” problem for people



Control by mouse



Where to look when?

AutoCam

Input 360° Video



Output NFOV Video



Automatically select FOV and viewing direction

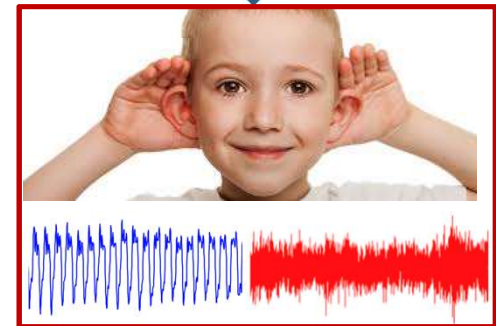
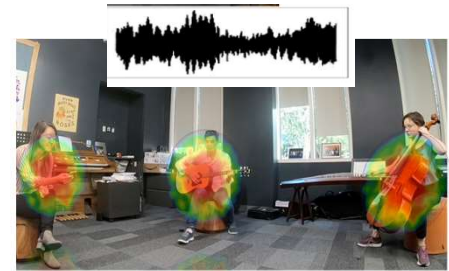
Anticipating the unseen and unheard



**Look-around
policies**



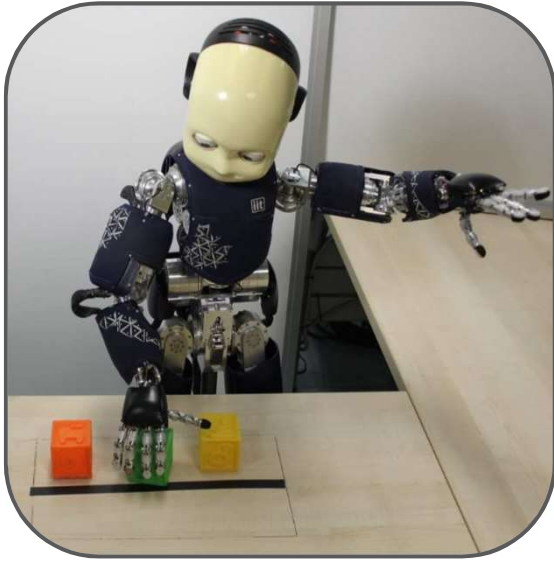
**Affordance
learning**



**Audio-visual
learning**

Towards embodied perception

Object interaction



Embodied
perception system



Object
manipulation

Turn on

Increase
height

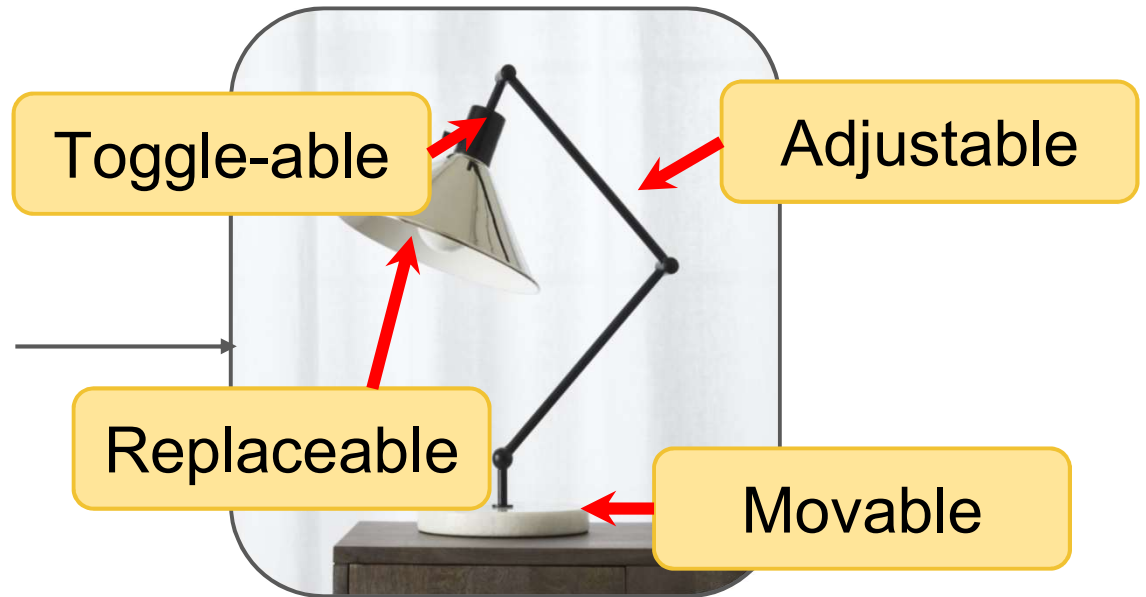
Move
lamp

Replace
lightbulb

What actions does an object *afford*?

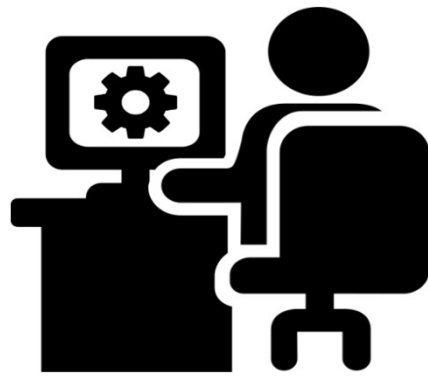


Embodied
perception system

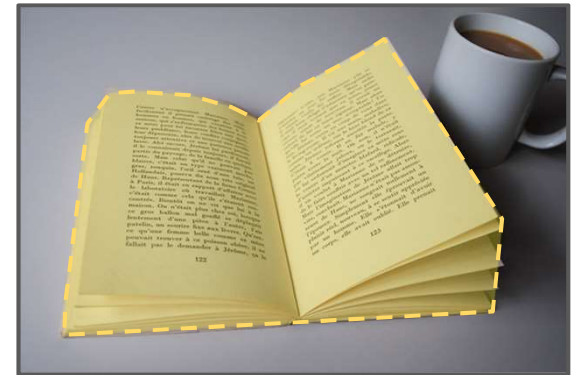


Object
manipulation

Current approaches: affordance as semantic segmentation



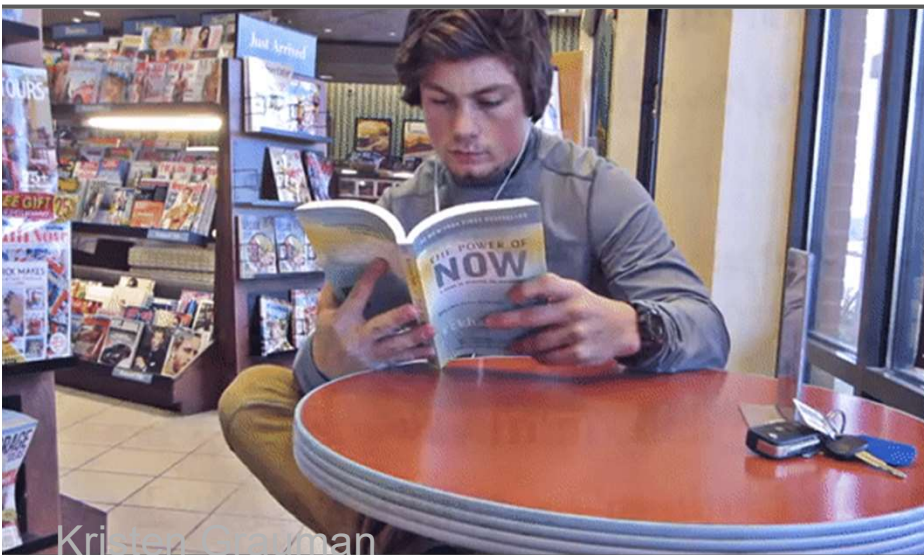
Label
“holdable”
regions



Captures annotators' expectations of what is important

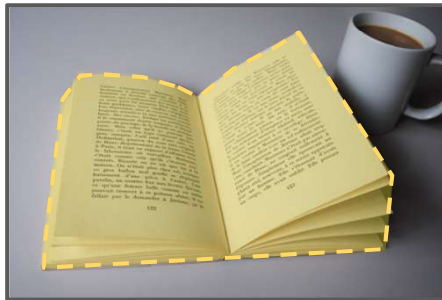
*Sawatzky et al. (CVPR 17), Nguyen et al. (IROS 17),
Kristen Grauman Roy et al. (ECCV 16), Myers et al. (ICRA 15), ...*

...but real human behavior is complex



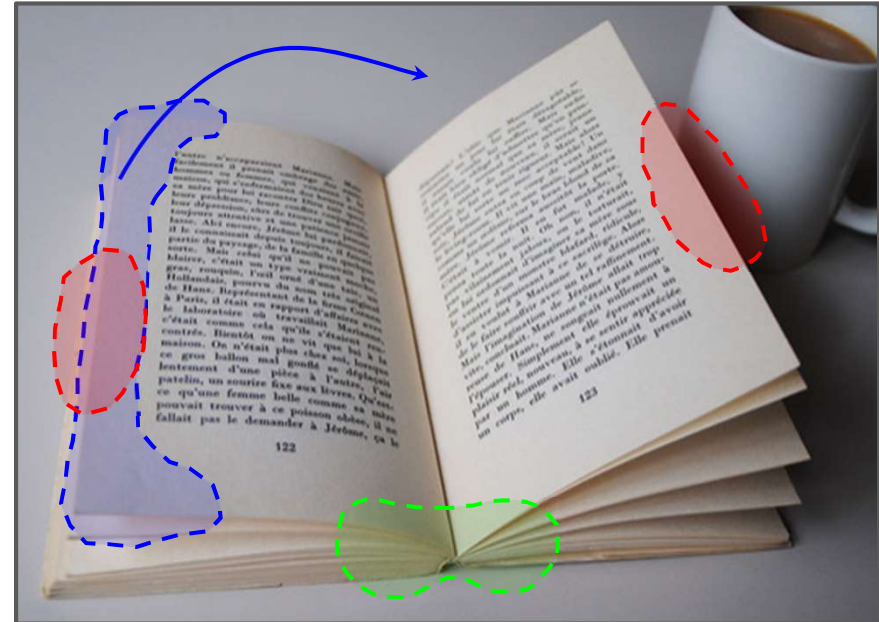
Kristen Grayman

How to learn object affordances?



Manually
curated
affordances

V
S.



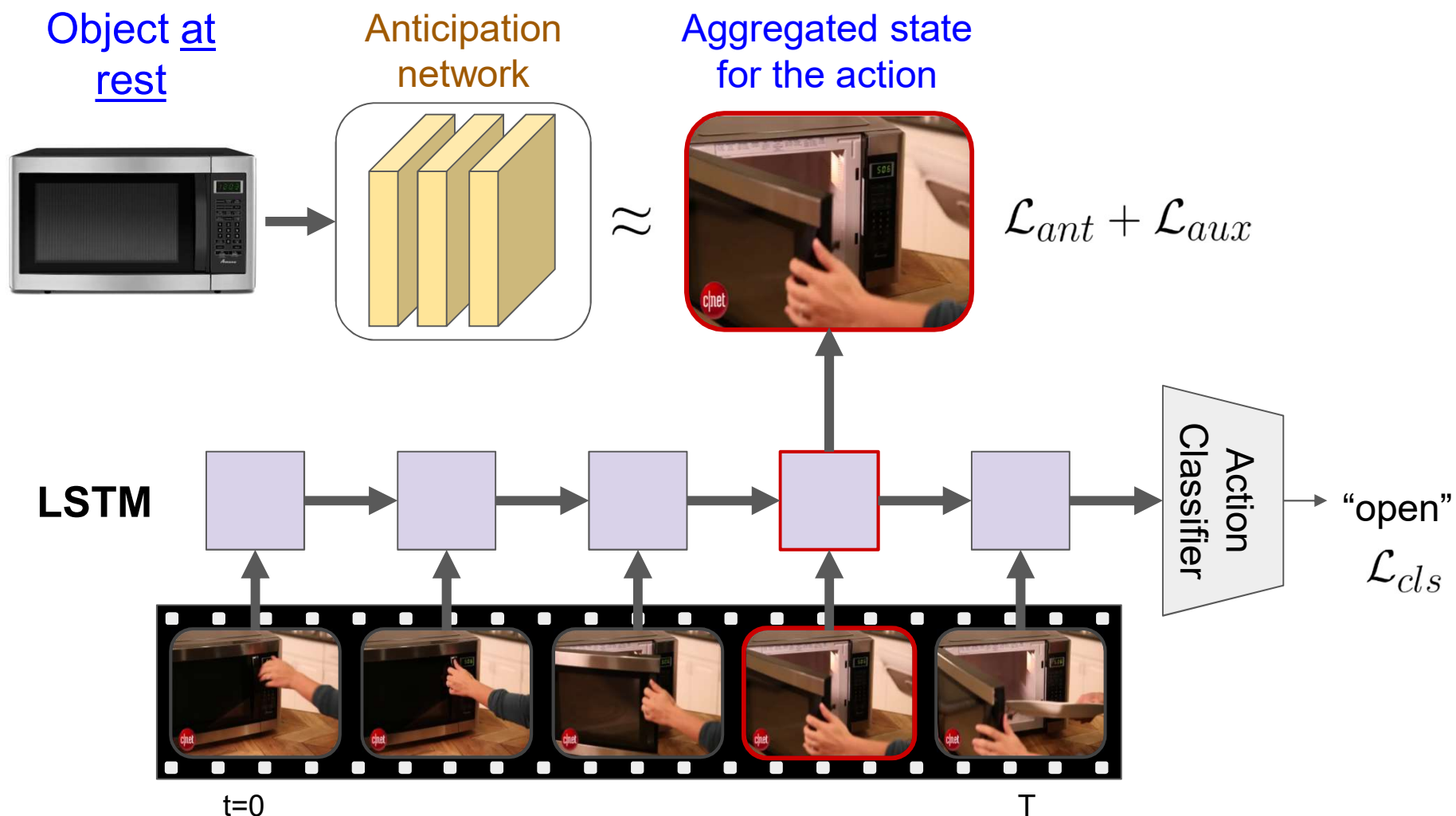
Real human
interactions?

Sawatzky et al. (CVPR 17), Nguyen et al. (IROS 17), Roy et al. (ECCV 16), Myers et al. (ICRA 15), ...

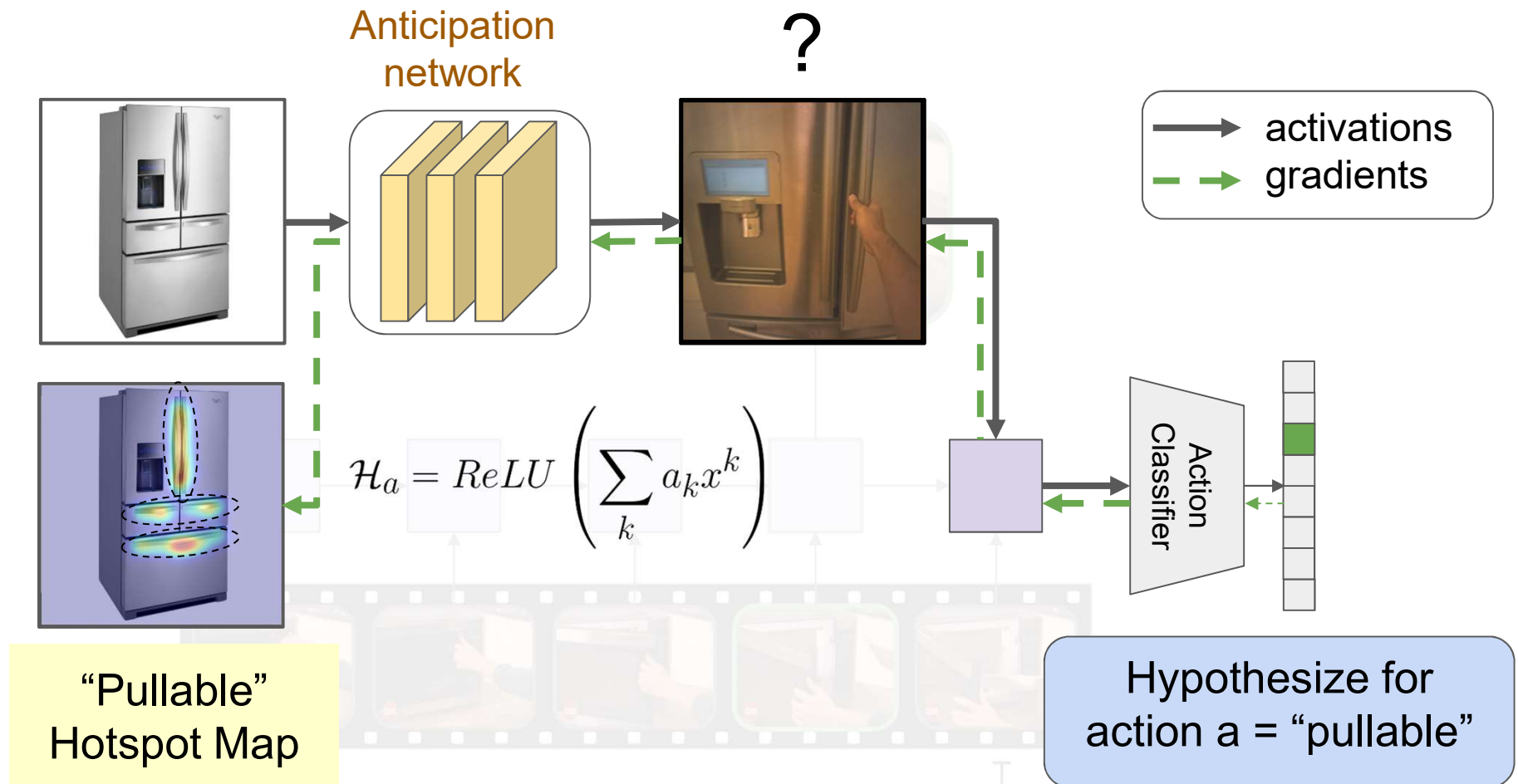
Our idea: Learn directly by watching people (video)



Learning affordances from video



Extracting interaction hotspot maps



Activation mapping to identify responsible spatial regions

Wait, is this just action recognition?

Action
recognition +
Grad-CAM



Ours



No: Hotspot anticipation model maps
object at rest to potential for interaction

Evaluating interaction hotspots

OPRA

(Fang et al., CVPR 18)



EPIC Kitchens

(Damen et al., ECCV 18)



MS COCO

(Lin et al., ECCV 14)

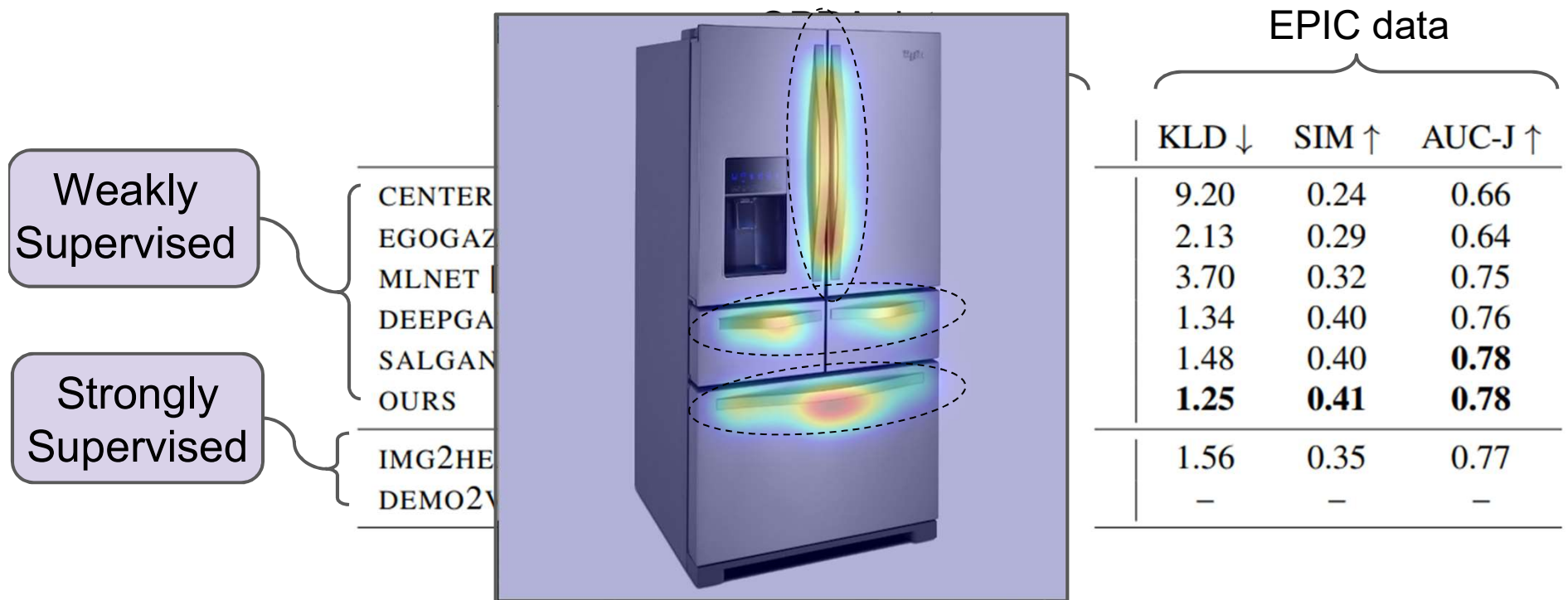


Train on video datasets, generate
heatmaps on novel images---
even from unseen categories



Results: interaction hotspots

Given static image of object at rest, infer affordance regions



Up to 24% increase vs. weakly supervised methods

Results: interaction hotspots



Results: hotspots for recognition

| N → | COCO | | | |
|-------------|-------------------|-------------------|-------------------|-------------------|
| | 5 | 25 | 100 | 3300 (all) |
| VANILLA | 44.3 ± 0.3 | 56.6 ± 0.2 | 65.6 ± 0.4 | 75.2 ± 0.1 |
| AUTOENCODER | 39.4 ± 0.4 | 51.2 ± 0.2 | 59.1 ± 0.2 | 72.8 ± 0.3 |
| OURS | 46.8 ± 0.3 | 57.9 ± 0.1 | 63.2 ± 0.2 | 73.9 ± 0.3 |



Better low-shot object recognition
by **anticipating object function**

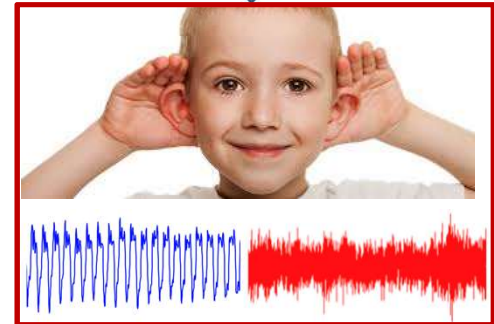
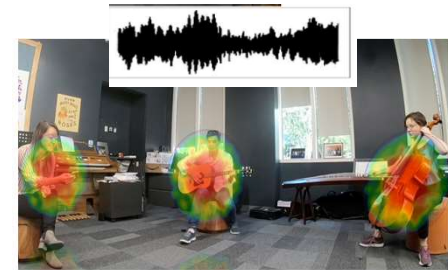
Anticipating the unseen and unheard



**Look-around
policies**



**Affordance
learning**



**Audio-visual
learning**

Towards embodied perception

Listening to learn



woof



meow



ring



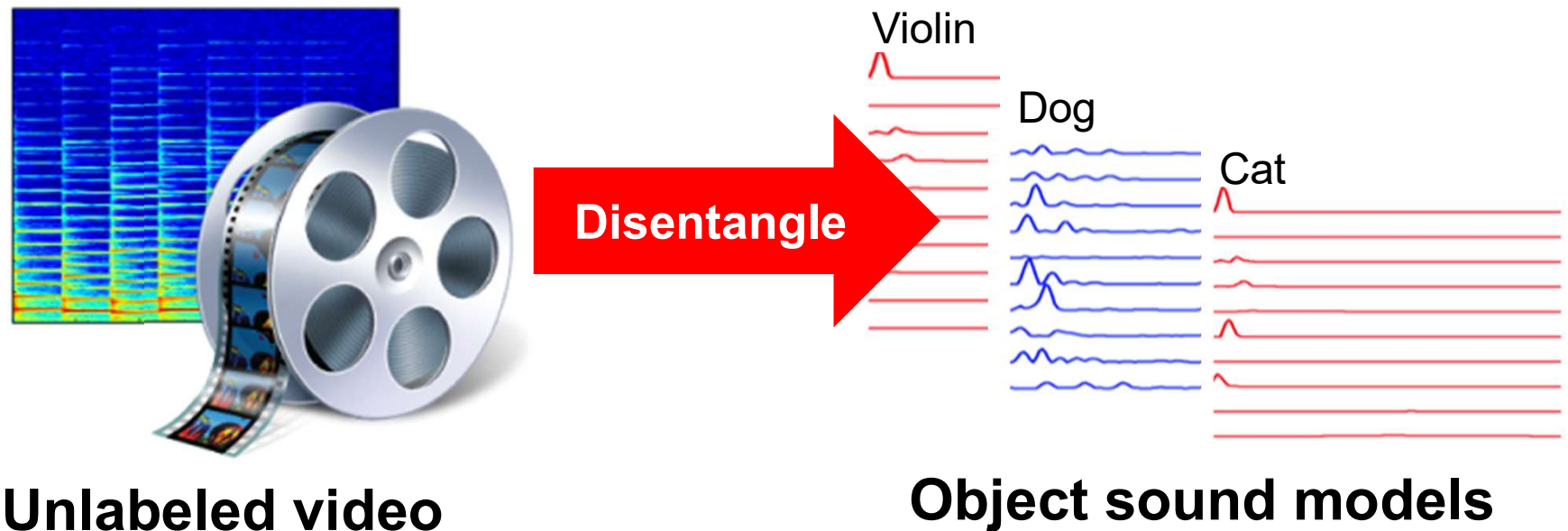
clatter

Goal: a repertoire of objects and their sounds

Challenge: a single audio channel mixes
sounds of multiple objects

Learning to separate object sounds

Our idea: Leverage visual objects to learn from *unlabeled* video with *multiple* audio sources



Apply to **separate** simultaneous sounds in novel videos

Results: audio-visual source separation

Train on 100,000 unlabeled multi-source video clips, then separate audio for novel video



original video
(before separation)

visual predictions:
violin & acoustic guitar

Dataset: AudioSet [Gemmeke et al. 2017]

Results: audio-visual source separation

Train on 100,000 unlabeled multi-source video clips, then separate audio for novel video



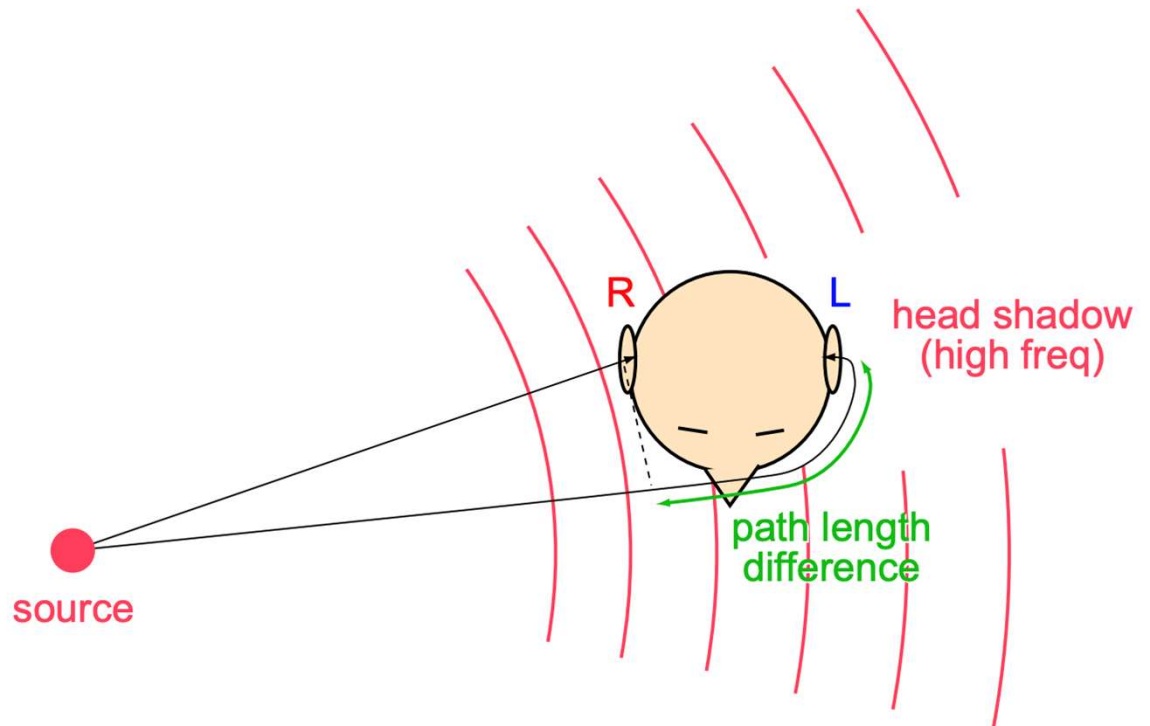
original video
(before separation)

visual predictions:
dog & violin

Spatial effects in audio



**Spatial effects
absent in
monaural audio**

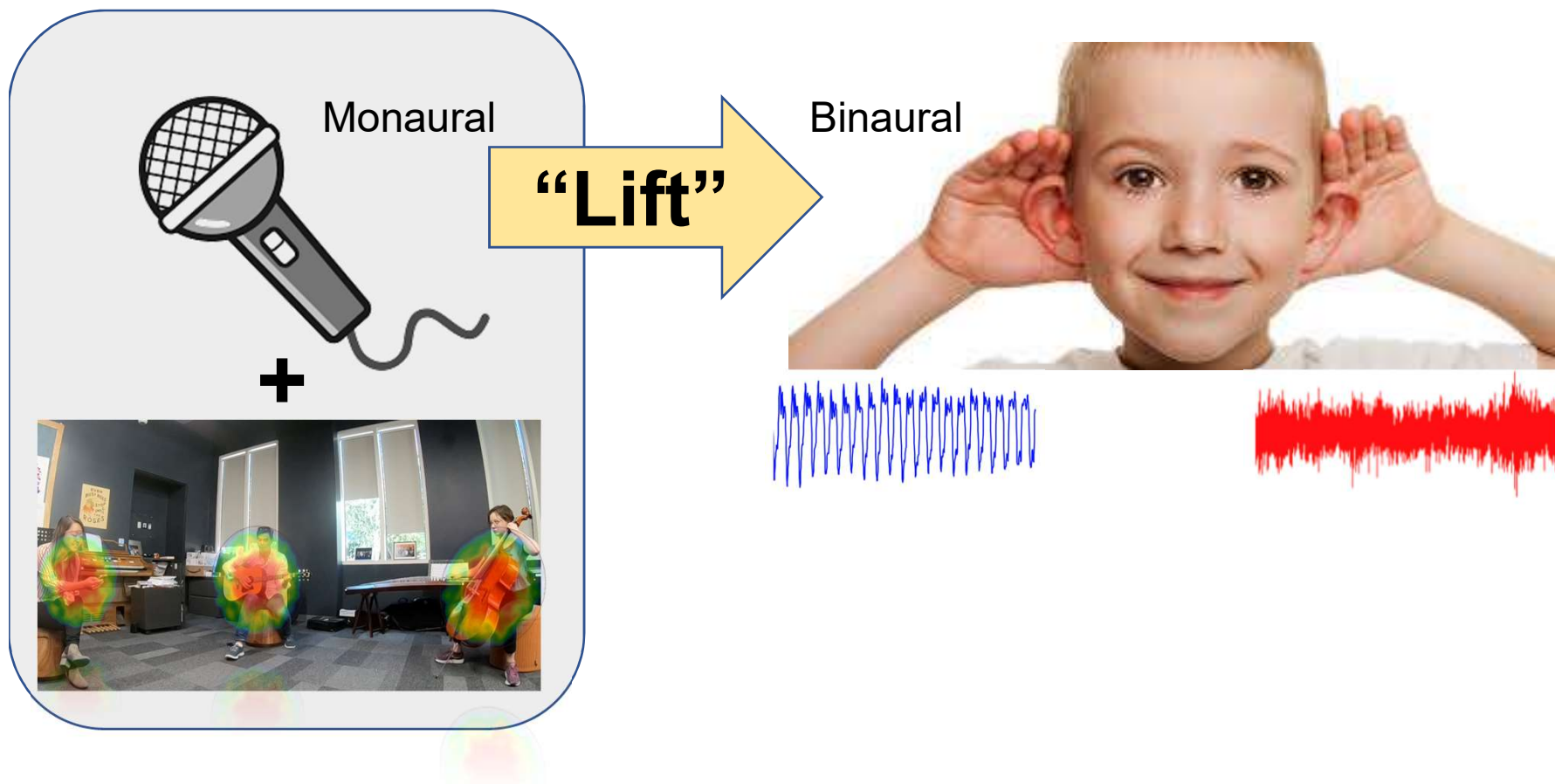


Cues for spatial hearing:

- Interaural time difference (ITD)
- Interaural level difference (ILD)
- Spectral detail (from pinna reflections)

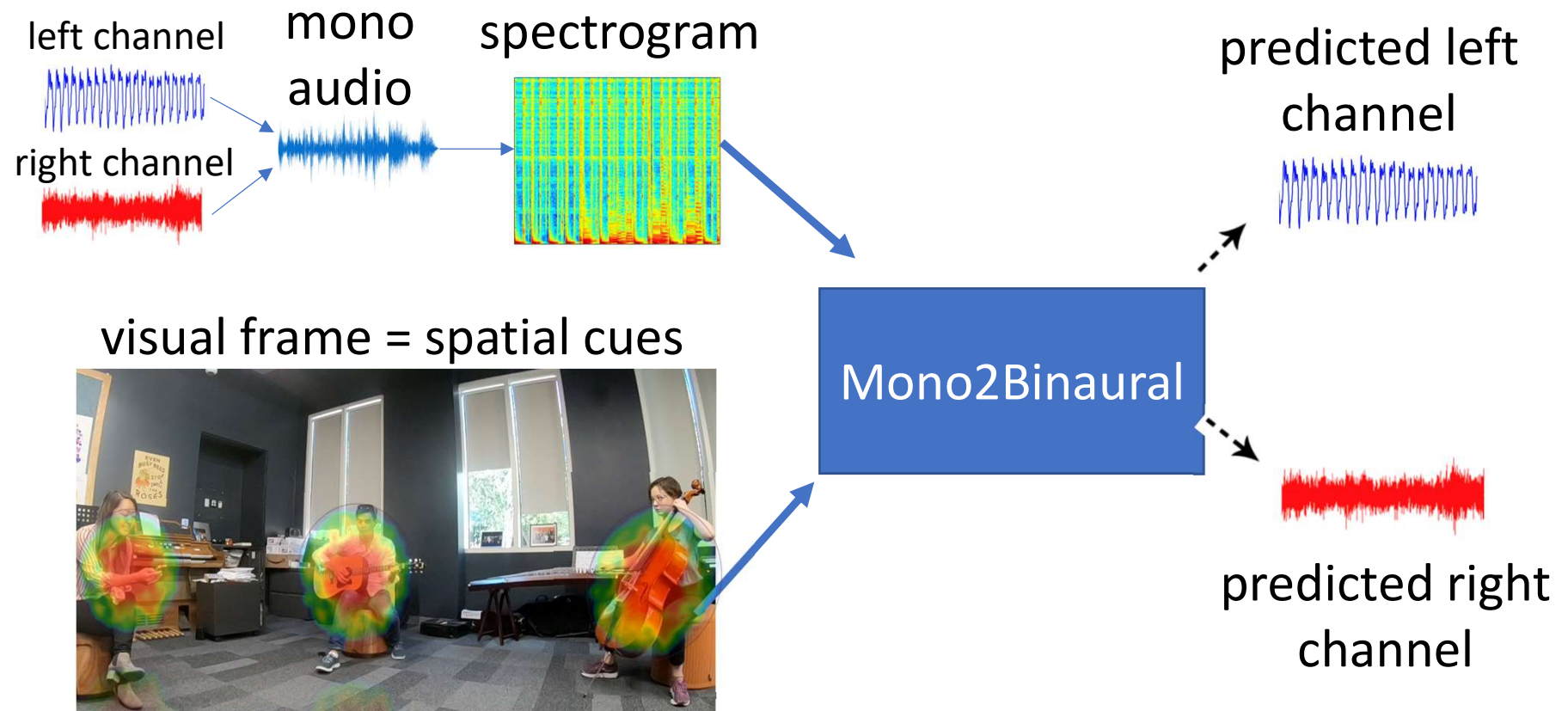
Our idea: 2.5D visual sound

“Lift” mono audio to spatial audio via visual cues



Our idea: 2.5D visual sound

“Lift” mono audio to spatial audio via visual cues



New: FAIR-Play dataset

<https://github.com/facebookresearch/FAIR-Play>

Binaural
microphone rig
linked to camera
and monoaural mic



Capture ~5 hours
video and binaural
sound in music room

Kristen Grauman

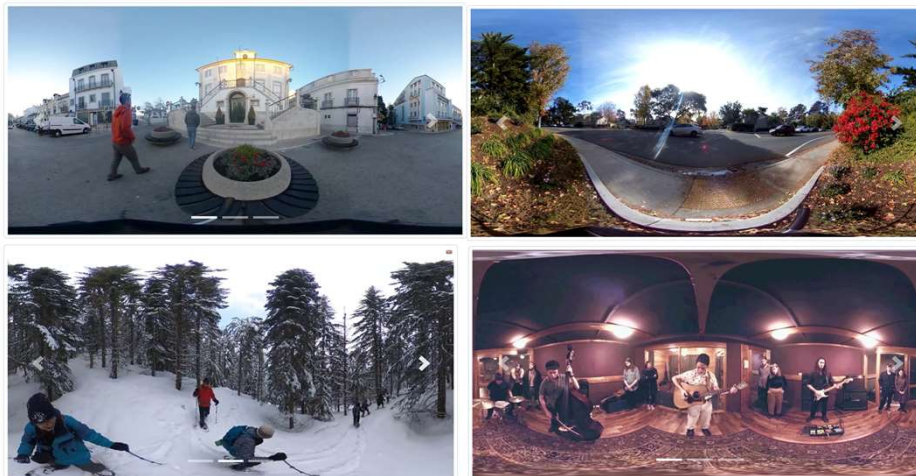
[Gao & Grauman, CVPR 2019]

Datasets



FAIR-Play Binaural

- 10 musical instruments, e.g., cello, guitar, harp, ukulele, trumpet, etc.
- ~5 hours of performances

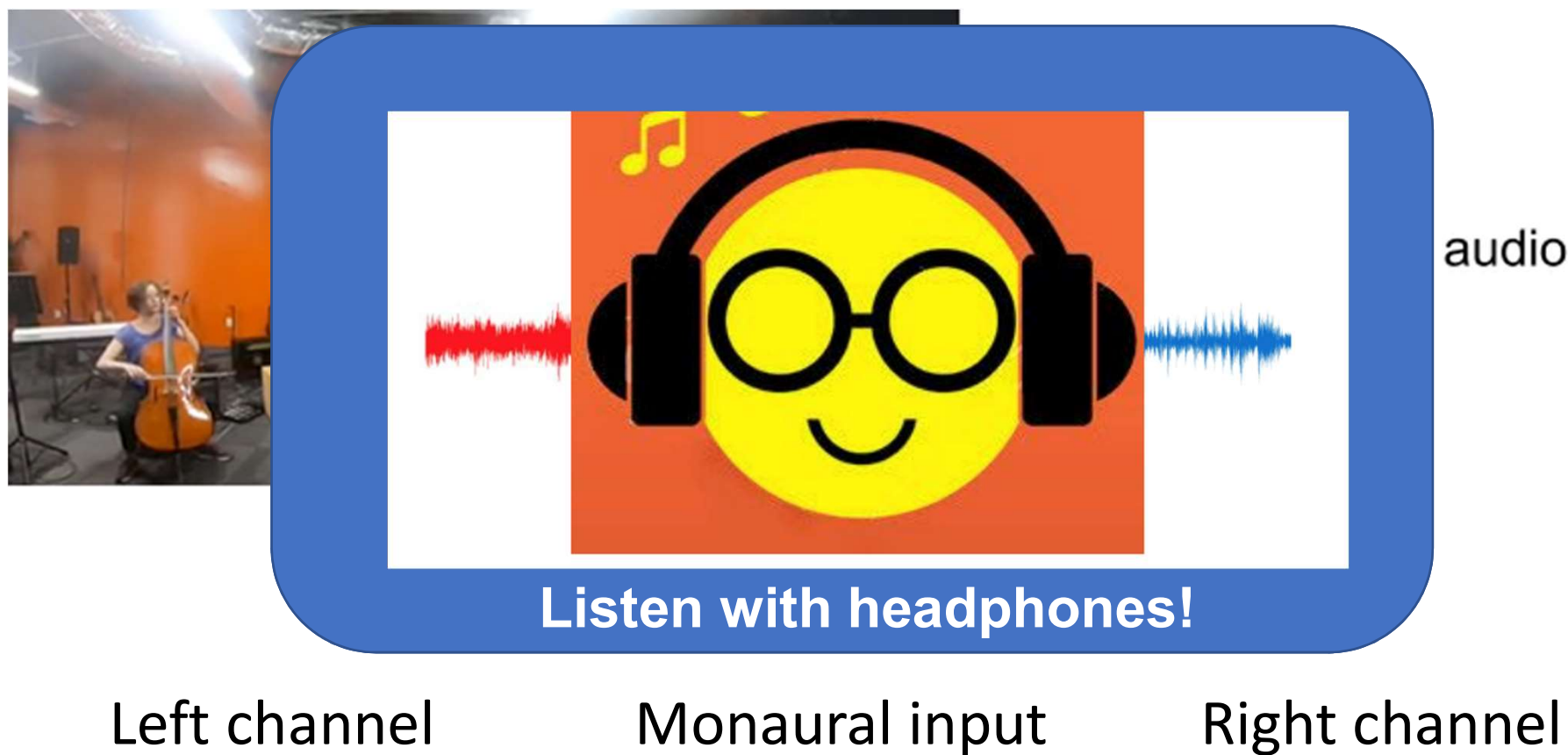


Ambisonics Datasets

[Morgado *et al.* NIPS 2018]

- Streets, random YouTube
- ~1000 360° video clips
- Converted to binaural audio using decoder

Results: 2.5D visual sound



vision.cs.utexas.edu/projects/2.5D_visual_sound/

Results: 2.5D visual sound



input monaural audio



Left channel

Monaural input

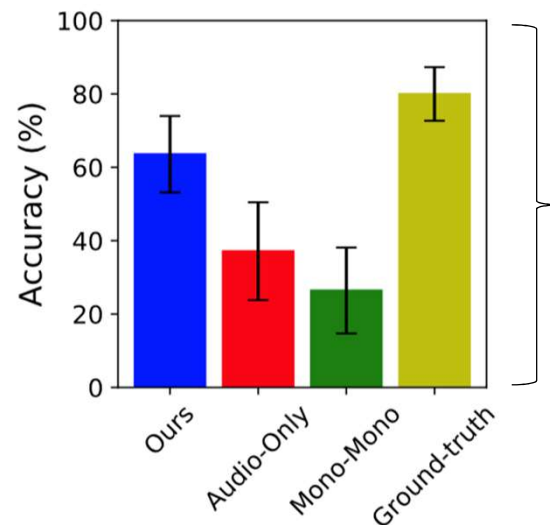
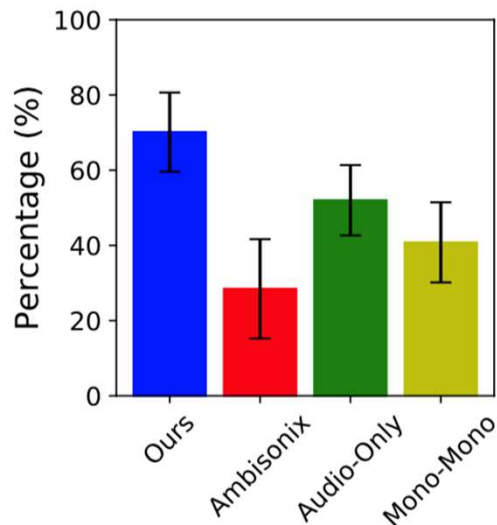
Right channel

vision.cs.utexas.edu/projects/2.5D_visual_sound/

Results: 2.5D visual sound

| | BINAURAL-MUSIC-ROOM | | REC-STREET | | YT-CLEAN | | YT-MUSIC | |
|-----------------------------|---------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|
| | STFT | ENV | STFT | ENV | STFT | ENV | STFT | ENV |
| Ambisonics [29] | - | - | 0.744 | 0.126 | 1.435 | 0.155 | 1.885 | 0.183 |
| Audio-Only | 1.022 | 0.143 | 0.590 | 0.114 | 1.065 | 0.131 | 1.553 | 0.167 |
| Flipped-Visual | 1.136 | 0.148 | 0.658 | 0.123 | 1.095 | 0.132 | 1.590 | 0.165 |
| Mono-Mono | 1.141 | 0.152 | 0.774 | 0.136 | 1.369 | 0.153 | 1.853 | 0.184 |
| MONO2BINAURAL (Ours) | 0.875 | 0.133 | 0.565 | 0.109 | 1.027 | 0.130 | 1.451 | 0.156 |

Binaural audio generation error, all four datasets



User studies:
perceived realism

Binaural audio offers “embodied” 3D sensation.
...and improves sound source separation!

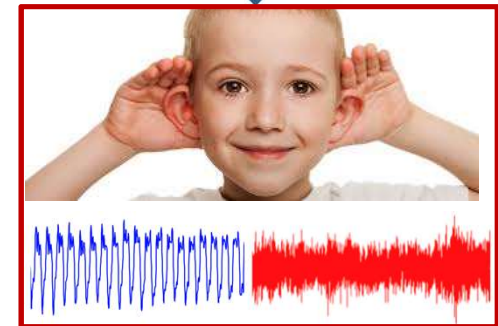
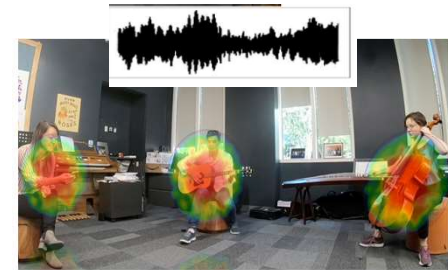
Anticipating the unseen and unheard



**Look-around
policies**



**Affordance
learning**



**Audio-visual
learning**

Towards embodied perception

Summary

Towards embodied perception

- self-supervised learning via anticipation
- learning to autonomously direct the camera
- multi-sensory observations (audio, motion, visual)
- object interaction from video



Ruohan
Gao



Tushar
Nagarajan



Dinesh
Jayaraman



Santhosh
Ramakrishnan



Christoph
Feichtenhofer