# Visual Learning with Unlabeled Video and Look-Around Policies
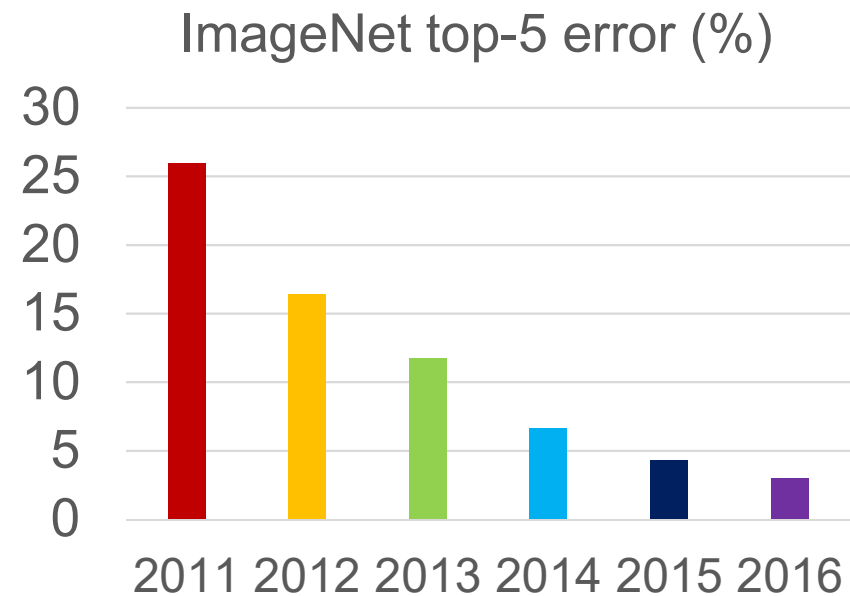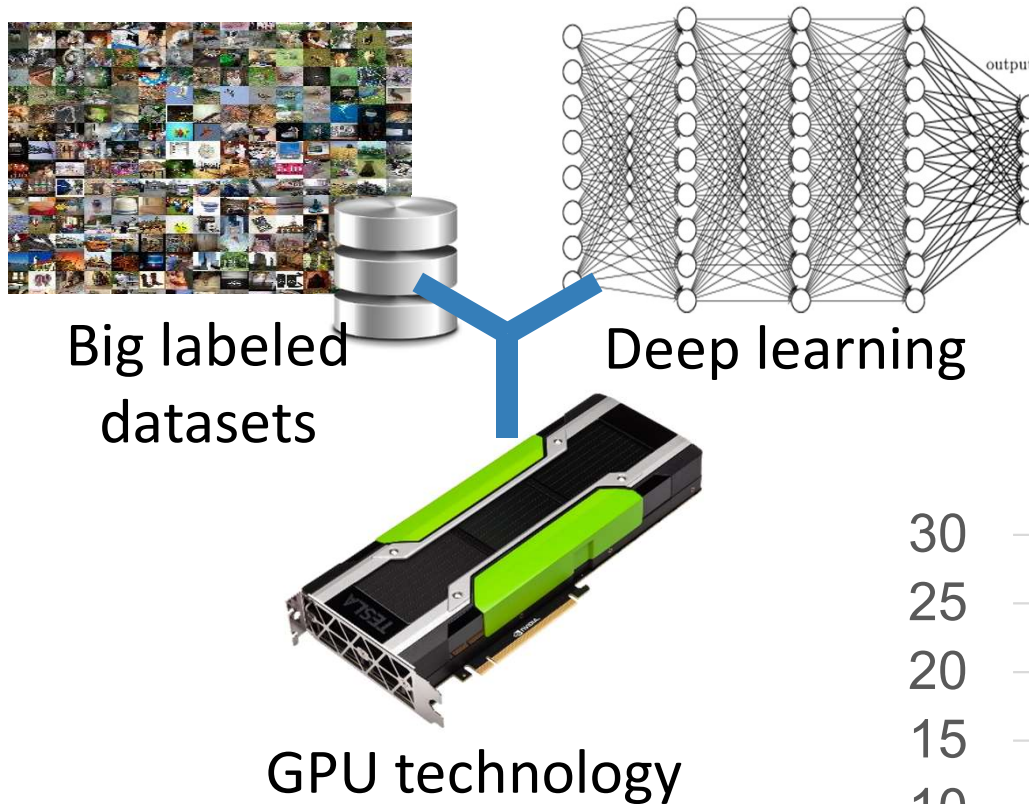
Kristen Grauman

Department of Computer Science

University of Texas at Austin

THE UNIVERSITY OF

TEXAS

AT AUSTIN

# Visual recognition: significant recent progress



Big labeled datasets

Deep learning
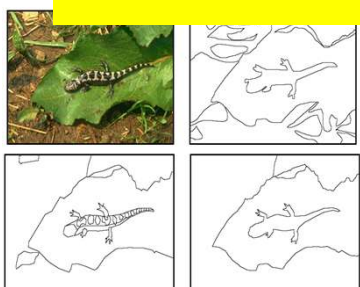
GPU technology

ImageNet top-5 error (%)

# How do systems typically learn about objects today?
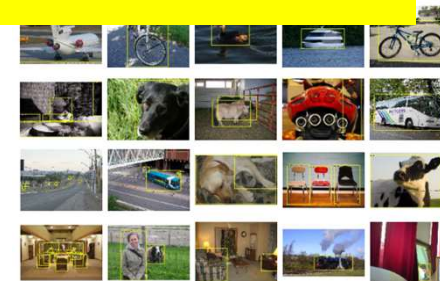


dog

boat

# Recognition benchmarks

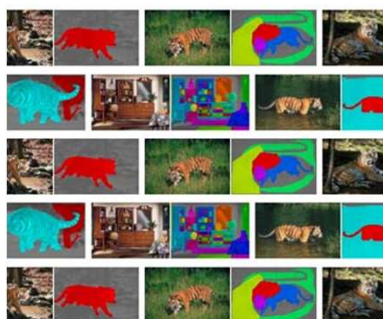A "disembodied" well-curated moment in time


BSD (2001)


Caltech 101 (2004), Caltech 256 (2006)


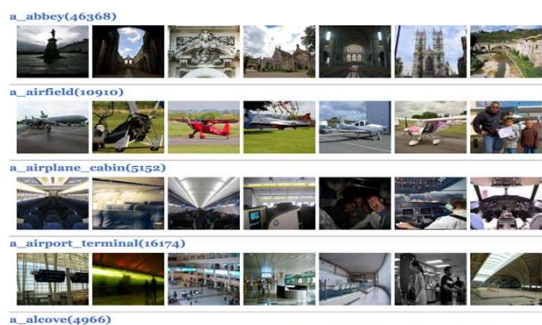PASCAL (2007-12)


LabelMe (2007)
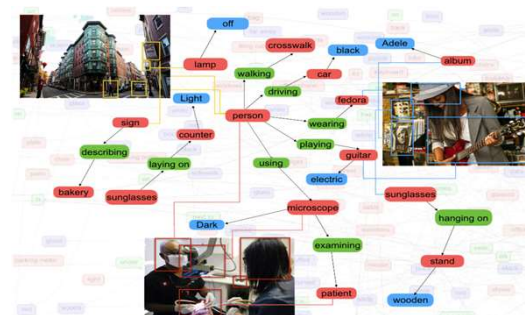

ImageNet (2009)


SUN (2010)


Places (2014)


MS COCO (2014)


Visual Genome (2016)

# Egocentric perceptual experience

A tangle of relevant and irrelevant multi-sensory information

Kristen Grauman, UT Austin

# Big picture goal: Embodied visual learning

**Status quo**:

Learn from "disembodied" bag of labeled snapshots.



**On the horizon:**

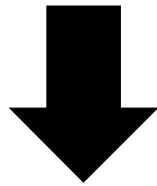Visual learning in the context of acting and moving in the world.



Kristen Grauman, UT Austin

# Big picture goal: Embodied visual learning

**Status quo**:

Learn from "disembodied" bag of labeled snapshots.



**On the horizon:**

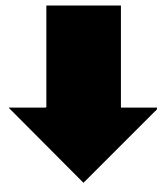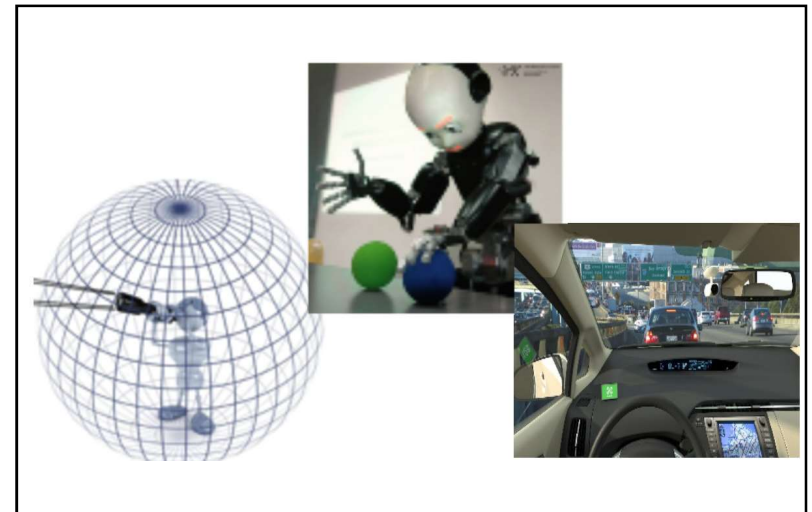Visual learning in the context of acting and moving in the world.
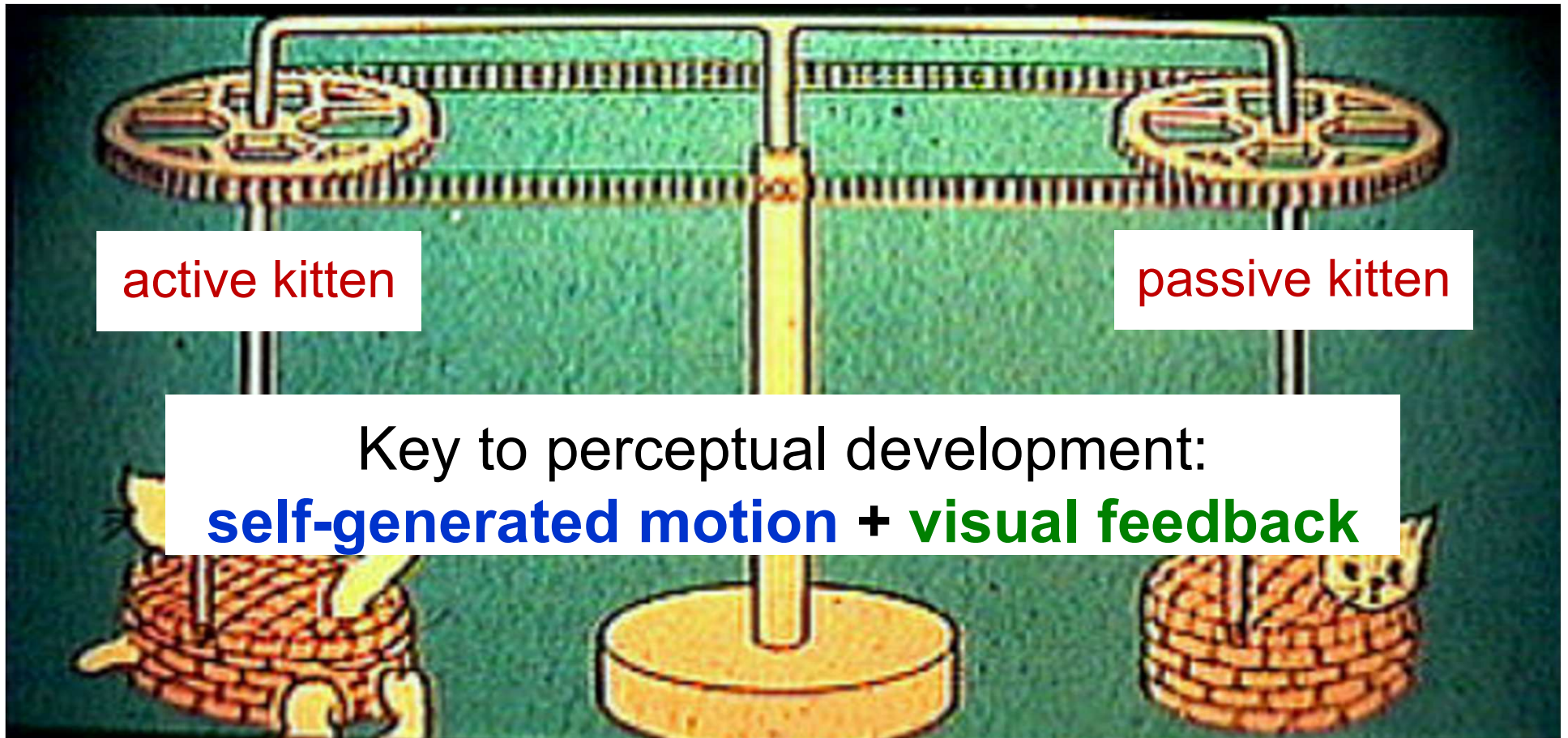


Kristen Grauman, UT Austin

# This talk

Towards embodied visual learning

1. Learning from unlabeled video and multiple sensory modalities

2. Learning policies for how to move for recognition and exploration

# The kitten carousel experiment
## [Held & Hein, 1963]



active kitten

passive kitten

Key to perceptual development:
**self-generated motion** + **visual feedback**

# Idea: Ego-motion ↔ vision

**Goal:** Teach computer vision system the connection:
"how I move" ↔ "how my visual surroundings change"



**Ego-motion motor signals**          **Unlabeled video**

*[Jayaraman & Grauman, ICCV 2015, IJCV 2017]*
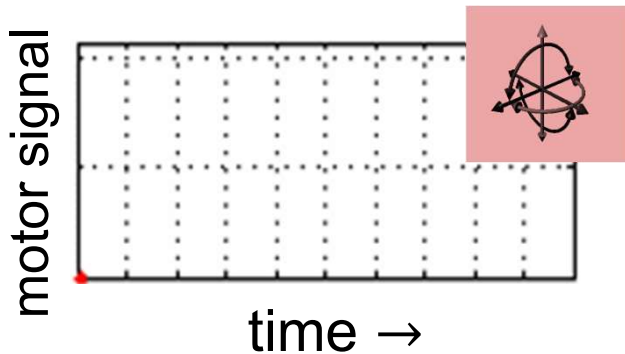
# Ego-motion ↔ vision: view prediction



**After moving:**

# Approach: Ego-motion equivariance

**Training data**
Unlabeled video +
motor signals

**Equivariant embedding**
organized by ego-motions



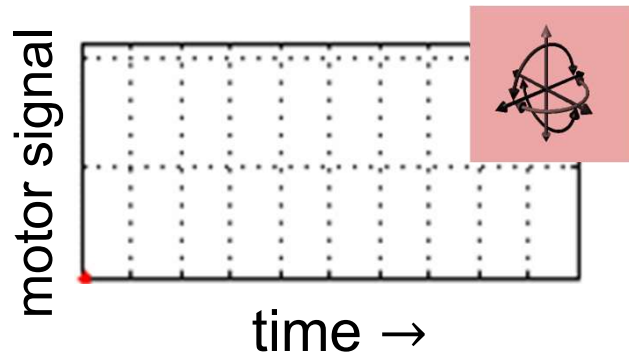$$\mathbf{z}(g\mathbf{x}) \approx M_g \mathbf{z}(\mathbf{x})$$

Learn

Pairs of frames related by
similar ego-motion should
be related by same
feature transformation

*[Jayaraman & Grauman, ICCV 2015, IJCV 2017]*

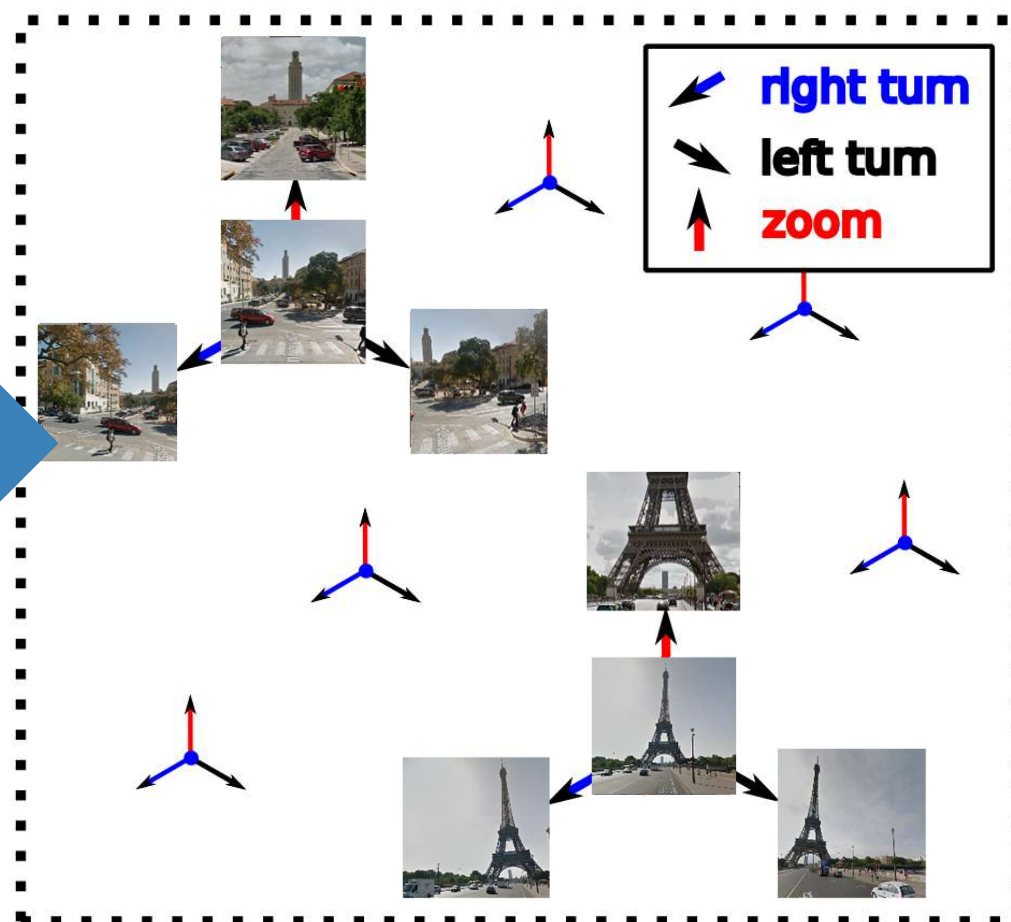# Approach: Egomotion equivariance

**Training data**
Unlabeled video +
motor signals

**Equivariant embedding**
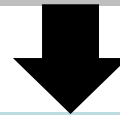organized by egomotions



Learn

right turn
left turn
zoom

*[Jayaraman & Grauman, ICCV 2015, IJCV 2017]*

# Example result: Recognition

Learn from *unlabeled* **car video** (KITTI)



Geiger et al, IJRR '13

Exploit features for **static scene classification**
(SUN, 397 classes)



Apse    Window se...    ...ardhouse

**30% accuracy increase**
when labeled data scarce

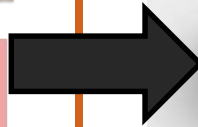, CVPR '10

# Passive → complete ego-motions

Pre-recorded video



Moving around to inspect

# One-shot reconstruction

Infer representation

View from other views



**Key idea:** One-shot reconstruction as a proxy task to learn semantic shape features.

# One-shot reconstruction



[Snavely et al, CVPR '06]



[Sinha et al, ICCV'93]

Shape from many views
geometric problem

Shape from one view
semantic problem

# Approach: ShapeCodes



- Implicit 3D shape representation
- No "canonical" azimuth to exploit
- Category agnostic

*[Jayaraman & Grauman, arXiv 2017, ECCV 2018]*

# ShapeCodes for recognition



ModelNet
[Wu et al 2015]

ShapeNet
[Chang et al 2015]

Accuracy (%)

■ Pixels  ■ Random wts  ■ DrLIM*  ■ Autoencoder**  ■ LSM^  ■ Ours

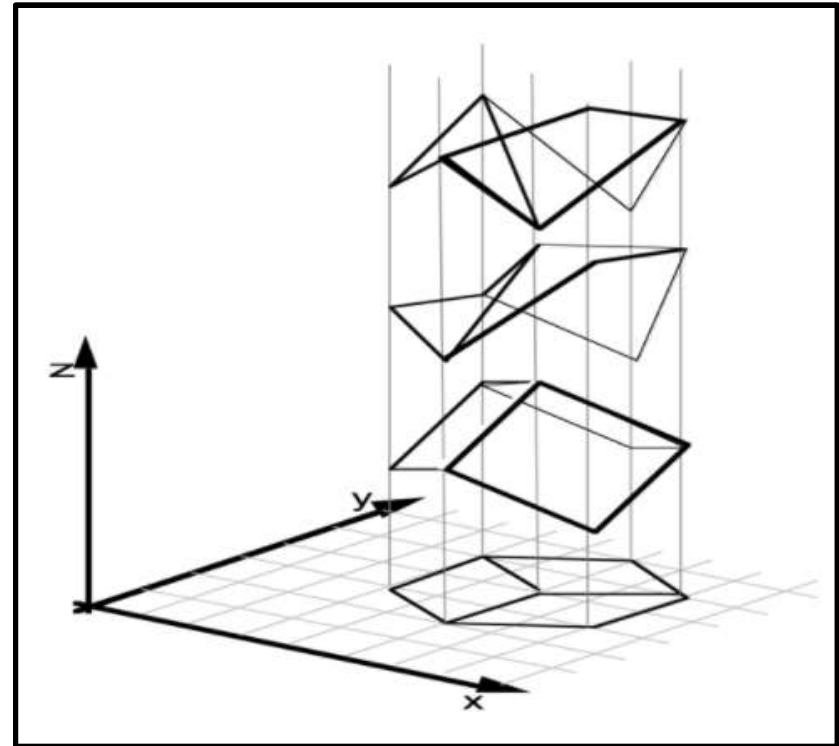*Hadsell et al, Dimensionality reduction by learning an invariant mapping, CVPR 2005
** Masci et al, Stacked Convolutional Autoencoders for Hierarchical Feature Extraction, ICANN 2011
^Agrawal, Carreira, Malik, Learning to See by Moving, ICCV 2015

# Ego-motion and implied body pose

Learn relationship between egocentric scene motion and 3D human body pose



**Input:**
egocentric video

**Output:**
sequence of 3d joint positions

*[Jiang & Grauman, CVPR 2017]*

# Ego-motion and implied body pose

Learn relationship between egocentric scene motion and 3D human body pose



**Wearable camera video**  **Inferred pose of camera wearer**

*[Jiang & Grauman, CVPR 2017]*

# This talk

Towards embodied visual learning

1. Learning from unlabeled video and multiple sensory modalities
   a) Egomotion / motor signals
   b) Audio signals

2. Learning policies for how to move for recognition and exploration

# Recall: Disembodied visual learning



dog

boat

# Listening to learn

# Listening to learn

# Listening to learn



**woof**   **meow**   **ring**   **clatter**

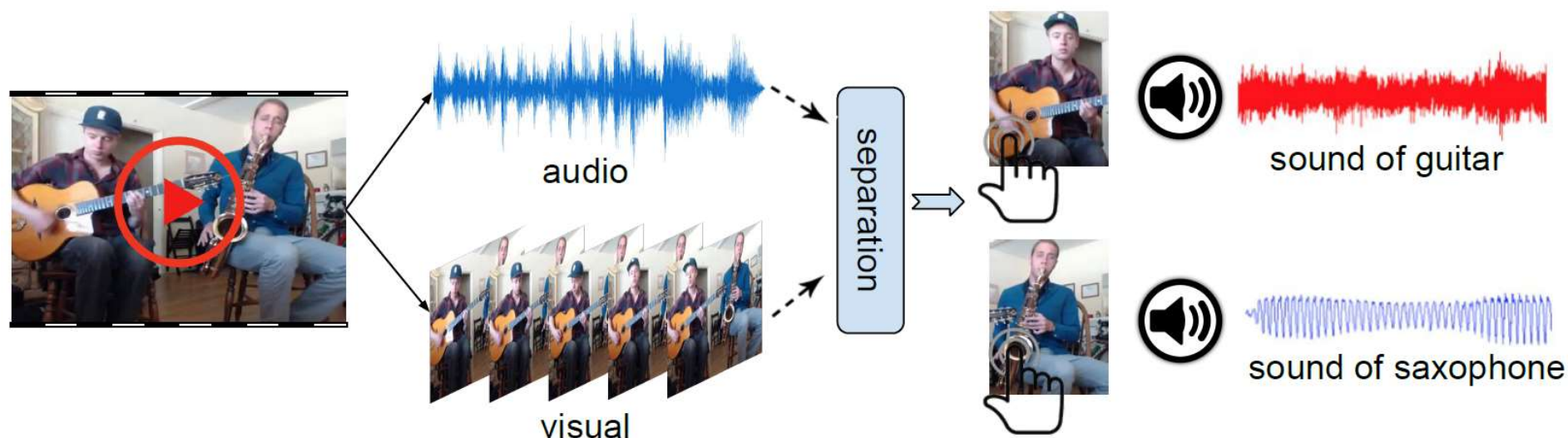**Goal**: A repertoire of objects and their sounds

# Visually-guided audio source separation



**Traditional approach:**
- Detect low-level correlations within a single video
- Learn from clean *single audio source* examples

*[Darrell et al. 2000; Fisher et al. 2001; Rivet et al. 2007; Barzelay & Schechner 2007; Casanovas et al. 2010; Parekh et al. 2017; Pu et al. 2017; Li et al. 2017]*

# Learning to separate object sounds

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



**Unlabeled video**

**Disentangle**

Violin

Dog

Cat

**Object sound models**

*[Gao, Feris, & Grauman, arXiv 2018]*

# Our approach: learning

Deep multi-instance multi-label learning (MIML) to disentangle which visual objects make which sounds



**Output: Group of audio basis vectors per object class**

# Our approach: inference

Given a novel video, use **discovered object sound models** to guide audio source separation.

# Results: Separating object sounds

## Train on 100,000 unlabeled video clips, then separate audio for novel video
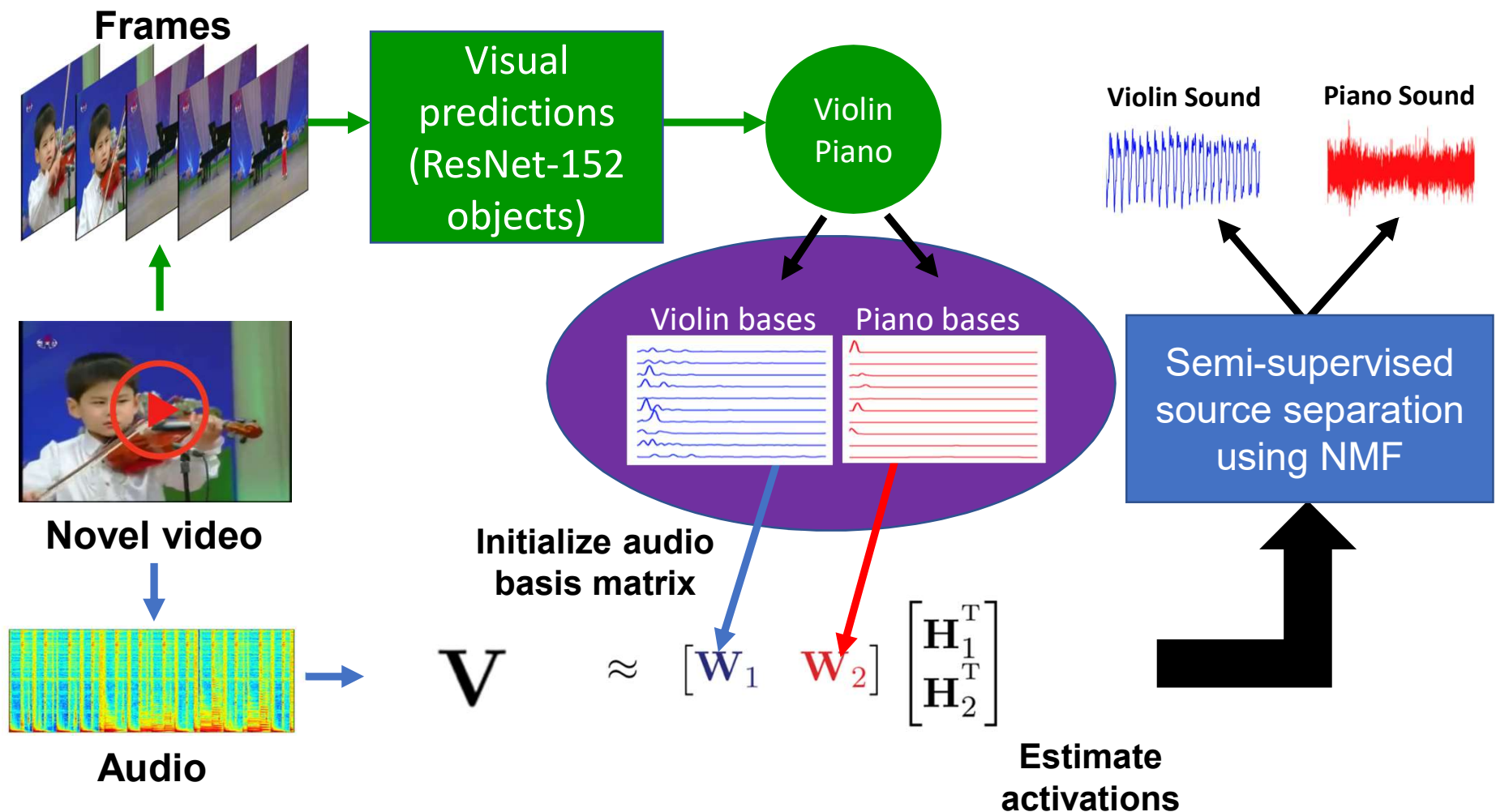


original video
(before separation)

visual predictions:
acoustic guitar & harmonica

Baseline: M. Spiertz, Source-filter based clustering for monaural blind source separation. International Conference on Digital Audio Effects, 2009

*[Gao, Feris, & Grauman, arXiv 2018]*

# Results: Separating object sounds

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



Failure case

original video
(before separation)

visual predictions:
acoustic guitar & electric guitar

Failure cases

[Gao, Feris, & Grauman, arXiv 2018]

# Results: Separating object sounds

| | Instrument Pair | Animal Pair | Vehicle Pair | Cross-Domain Pair |
|---|---|---|---|---|
| Upper-Bound | 2.05 | 0.35 | 0.60 | 2.79 |
| K-means Clustering | -2.85 | -3.76 | -2.71 | -3.32 |
| MFCC Unsupervised [65] | 0.47 | -0.21 | -0.05 | 1.49 |
| Visual Exemplar | -2.41 | -4.75 | -2.21 | -2.28 |
| Unmatched Bases | -2.12 | -2.46 | -1.99 | -1.93 |
| Gaussian Bases | -8.74 | -9.12 | -7.39 | -8.21 |
| Ours | **1.83** | **0.23** | **0.49** | **2.53** |

**Visually-aided audio source separation (SDR)**

| | Wooden Horse | Violin Yanni | Guitar Solo | Average |
|---|---|---|---|---|
| Sparse CCA (Kidron et al. [43]) | 4.36 | 5.30 | 5.71 | 5.12 |
| JIVE (Lock et al. [50]) | 4.54 | 4.43 | 2.64 | 3.87 |
| Audio-Visual (Pu et al. [56]) | 8.82 | 5.90 | **14.1** | 9.61 |
| Ours | **12.3** | **7.88** | 11.4 | **10.5** |

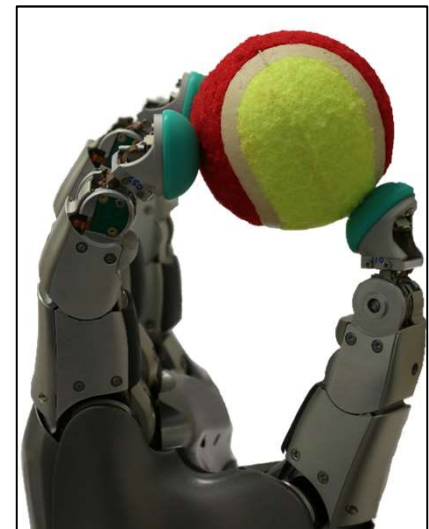**Visually-aided audio denoising (NSDR)**

Train on 100K unlabeled video clips from AudioSet [Gemmeke et al. 2017]

# This talk

Towards embodied visual learning

1. Learning from unlabeled video and multiple sensory modalities

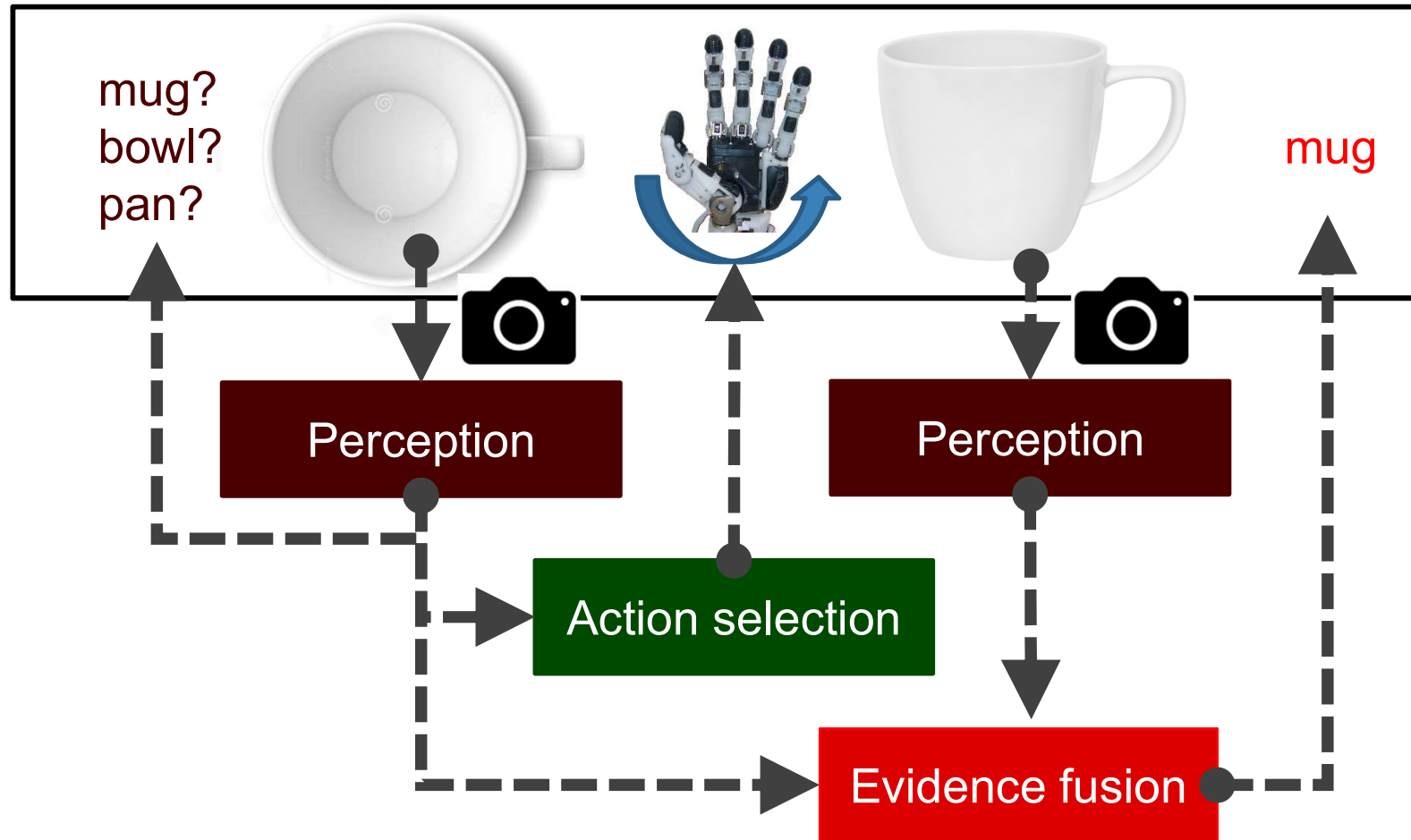2. Learning policies for how to move for recognition and exploration

# Moving to recognize



Time to revisit active recognition in challenging settings!

*Bajcsy 1985, Aloimonos 1988, Ballard 1991, Wilkes 1992, Dickinson 1997, Schiele & Crowley 1998, Tsotsos 2001, Denzler 2002, Soatto 2009, Ramanathan 2011, Borotschnig 2011, …*
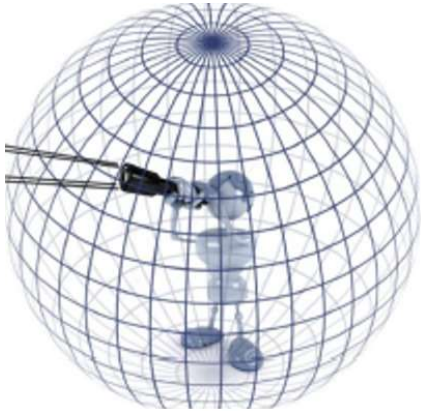
# End-to-end active recognition



mug?
bowl?
pan?

mug

Perception

Perception

Action selection

Evidence fusion

*[Jayaraman and Grauman, ECCV 2016, PAMI 2018]*

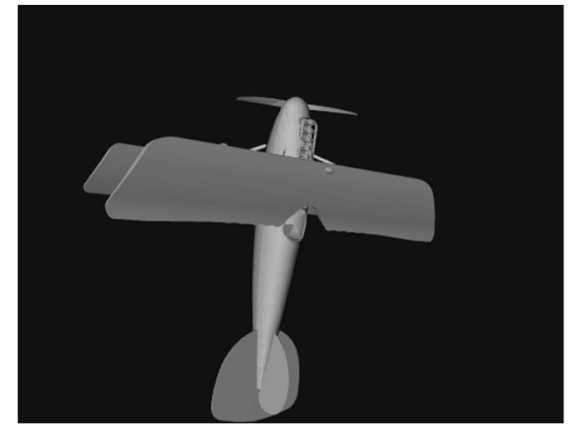# End-to-end active recognition

Look around scene

Manipulate object

Move around an object
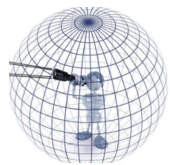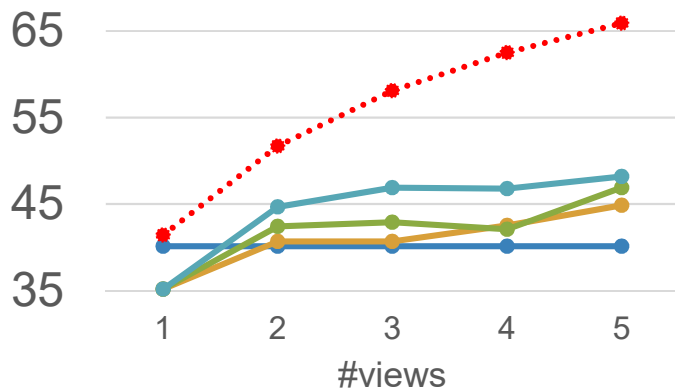


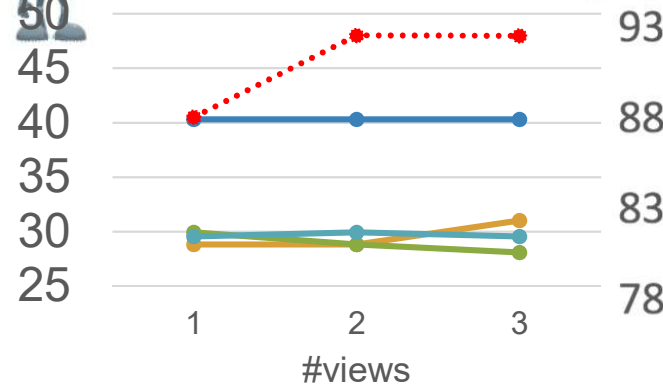*[Jayaraman and Grauman, ECCV 2016, PAMI 2018]*

# End-to-end active recognition



**SUN 360**

**GERMS**

**ModelNet-10**

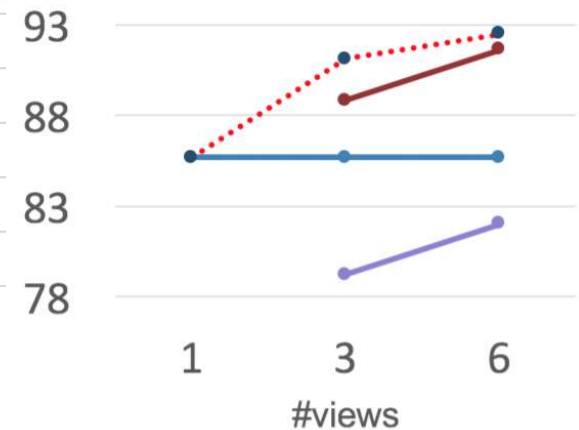Passive neural net
Transinformation [Schiele98]
SeqDP [Denzler03]
Transinformation+SeqDP
Ours

Passive neural net
Transinformation [Schiele98]
SeqDP[Denzler03]
Transinformation+SeqDP
Ours

Passive neural net
ShapeNets [Wu15]
Pairwise [Johns 16]
Ours

Agents that learn to look around intelligently can recognize things faster.

*[Jayaraman and Grauman, ECCV 2016, PAMI 2018]*

# End-to-end active recognition: example

P("Plaza courtyard"):  (6.28)          (11.95)          (68.38)

Top 3 guesses:
- Restaurant          Theater          Plaza courtyard
- Train interior          Restaurant          Street
- Shop          Plaza courtyard          Theater



*[Jayaraman and Grauman, ECCV 2016, PAMI 2018]*

# End-to-end active recognition: example

Predicted
label:



T=1                T=2                T=3

GERMS dataset: Malmir et al. BMVC 2015

*[Jayaraman and Grauman, ECCV 2016, PAMI 2018]*

# Goal: Learn to "look around"



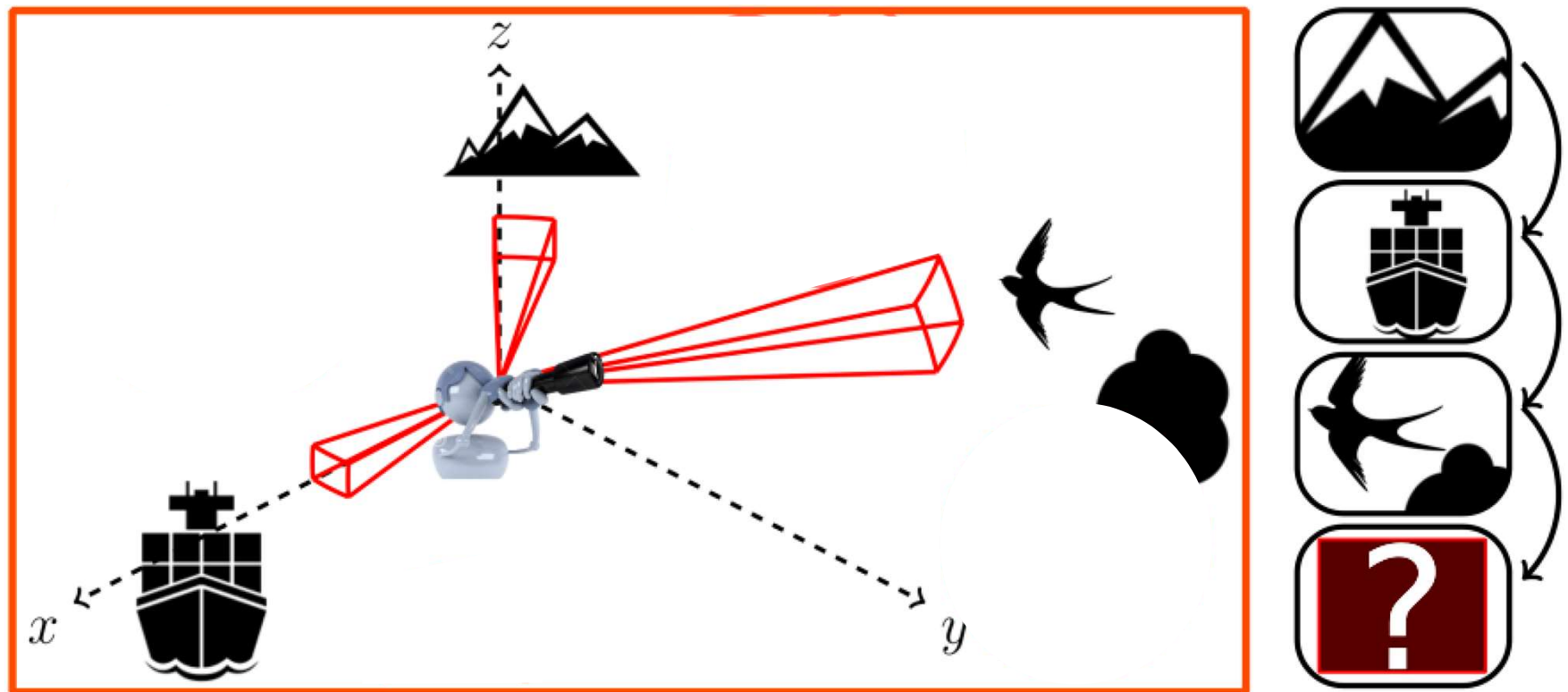recognition       **vs.**       reconnaissance       search and rescue

task predefined       task unfolds dynamically

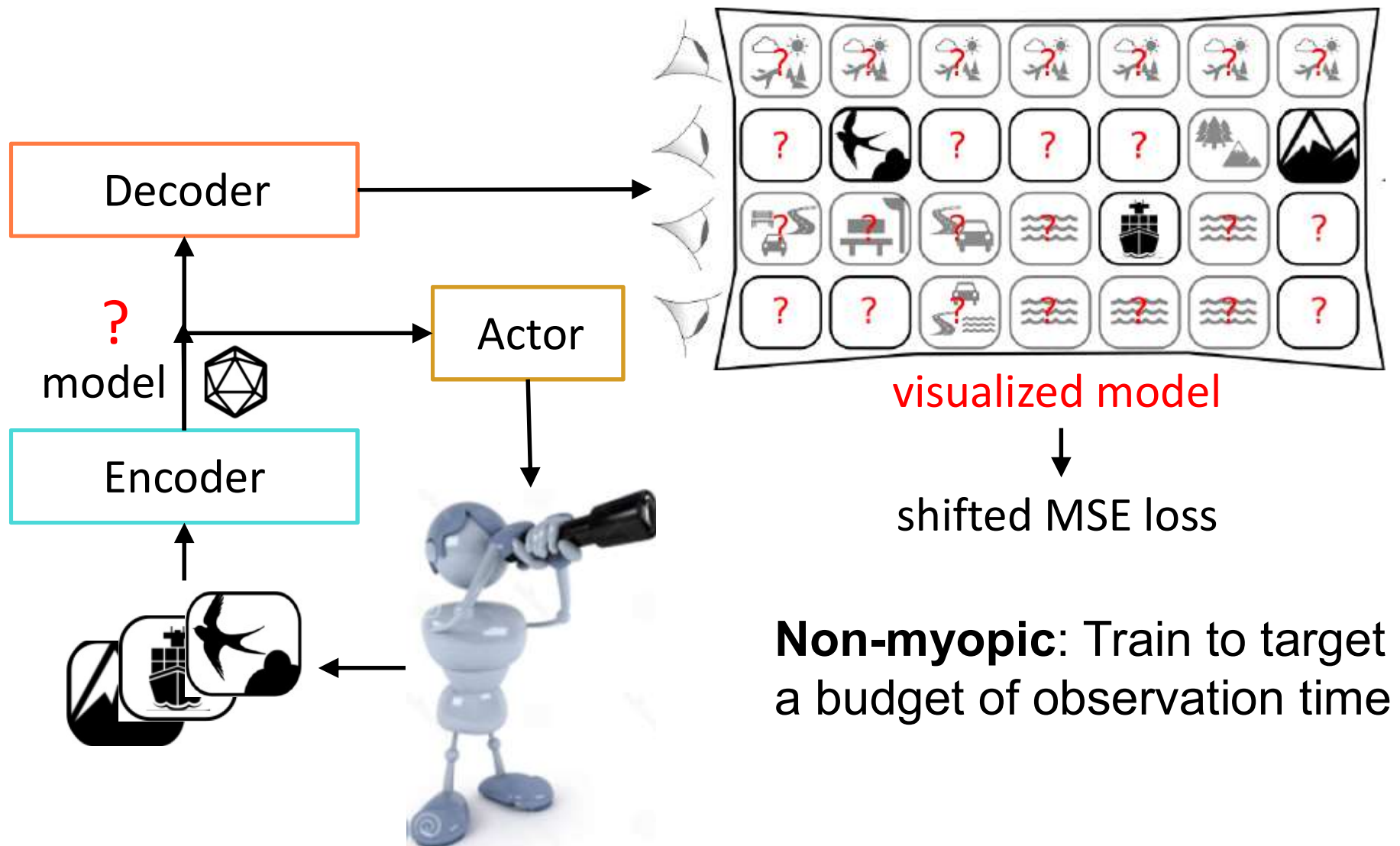Can we learn look-around policies for visual agents that are curiosity-driven, exploratory, and generic?

# Key idea: Active observation completion

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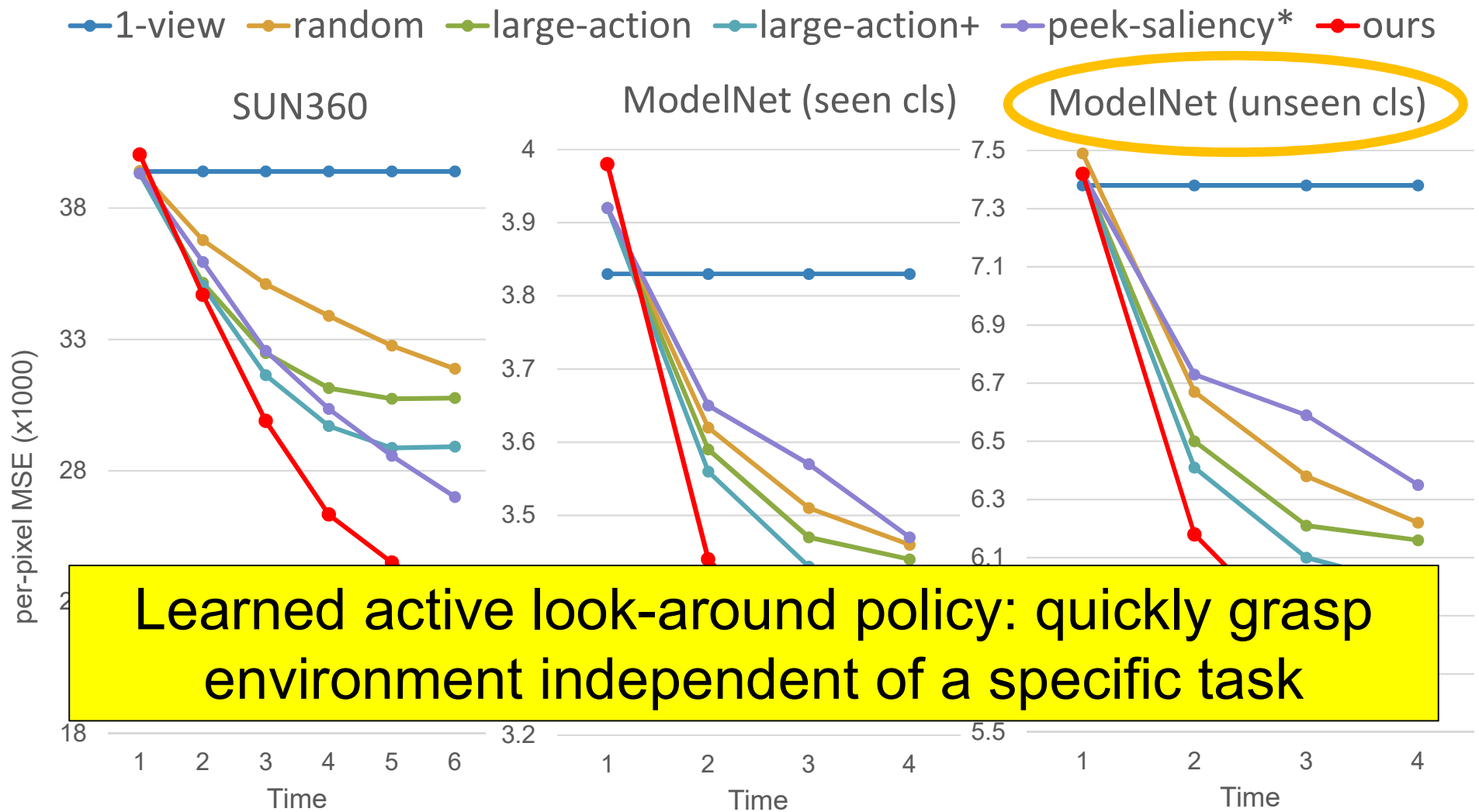


**Agent must choose where to look *before* looking there.**

# Approach: Active observation completion



visualized model

shifted MSE loss

**Non-myopic**: Train to target a budget of observation time

*Jayaraman and Grauman, CVPR 2018*
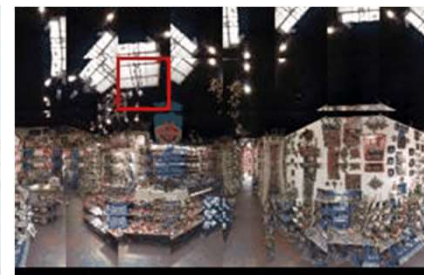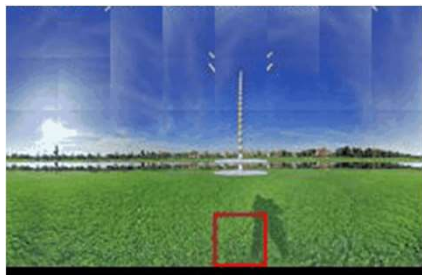
# Active "look around" results



*Saliency -- Harel et al, Graph based Visual Saliency, NIPS'07     *Jayaraman and Grauman, CVPR 2018*
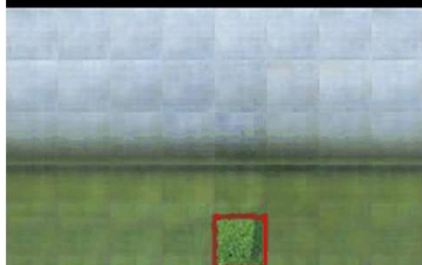
# Active "look around" visualization
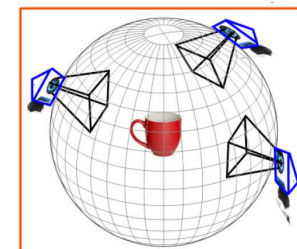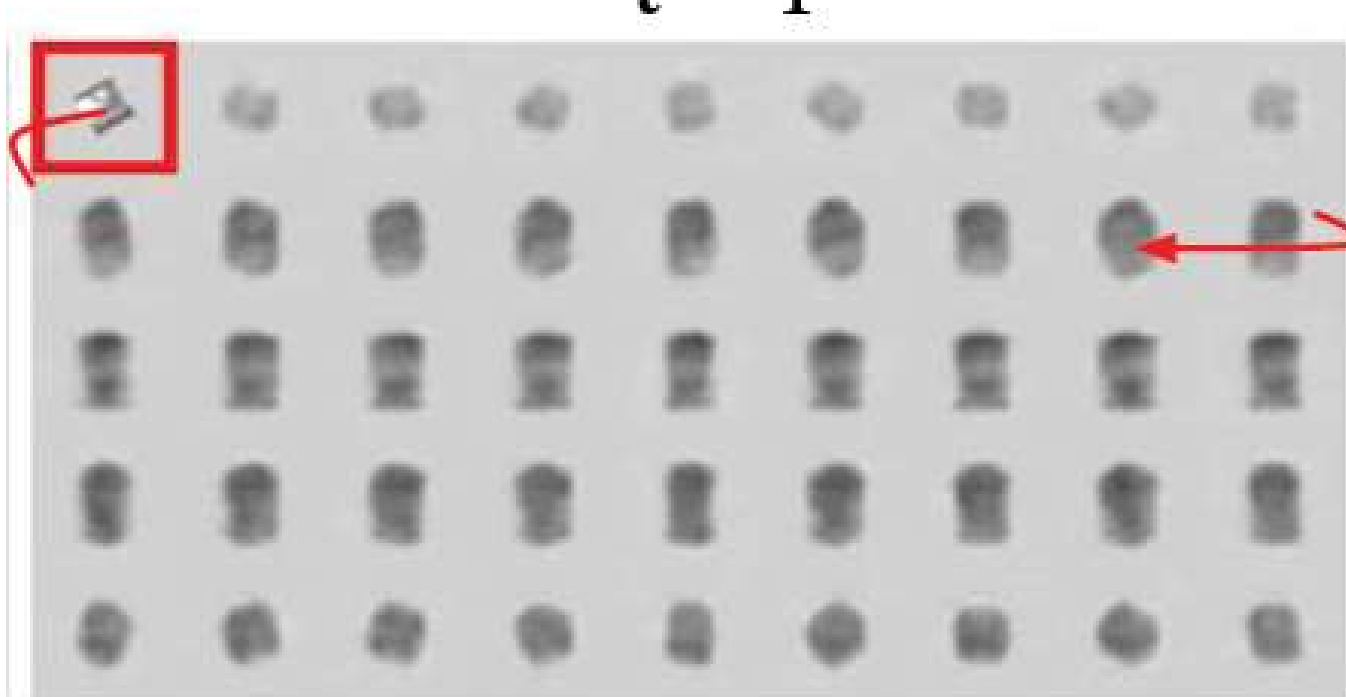


Complete 360 scene (ground truth)

Inferred scene

☐ = observed views

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

*Jayaraman and Grauman, CVPR 2018*

# Active "look around" visualization

$$t = 1$$



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

*Jayaraman and Grauman, CVPR 2018*

# Active "look around" visualization



$t = 2$

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

*Jayaraman and Grauman, CVPR 2018*

# Active "look around" visualization



$$t = 3$$

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

*Jayaraman and Grauman, CVPR 2018*
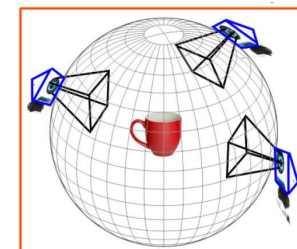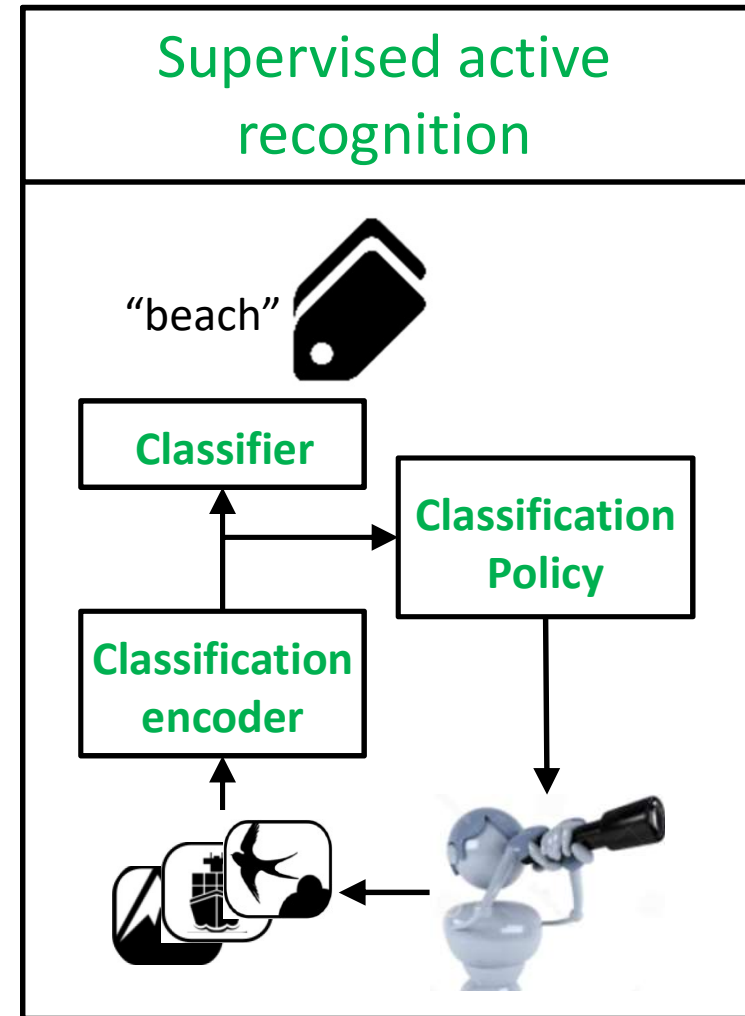
# Motion policy transfer



Unsupervised observation completion

Decoder

Look-around Policy

Look-around encoder

Supervised active recognition

"beach"

Classifier

Classification Policy

Classification encoder

Plug observation completion policy in for new task

# Motion policy transfer



SUN 360 Scenes

ModelNet Objects

Unsupervised exploratory policy approaches supervised task-specific policy accuracy!

*Jayaraman and Grauman, CVPR 2018*

# Summary



THE UNIVERSITY OF
# TEXAS
AT AUSTIN

- Visual learning benefits from

  – context of action and motion in the world

  – continuous unsupervised observations

- New ideas:

  – Embodied feature learning via visual and motor signals

  – Learning to separate object sound models from unlabeled video

  – Active policies for view selection and camera control



Dinesh
Jayaraman



Ruohan
Gao

Kristen Grauman, UT Austin

# Papers/code/videos

- **Learning to Separate Object Sounds by Watching Unlabeled Video**. R. Gao, R. Feris, and K. Grauman. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, Sept 2018. (Oral) [pdf] [videos]

- **ShapeCodes: Self-Supervised Feature Learning by Lifting Views to Viewgrids**. D. Jayaraman, R. Gao, and K. Grauman. In Proceedings of the European Conference on Computer Vision (ECCV), Munich, Germany, Sept 2018. [pdf]

- **End-to-end Policy Learning for Active Visual Categorization**. D. Jayaraman and K. Grauman. To appear, Transactions on Pattern Analysis and Machine Intelligence (PAMI), 2018. [pdf]

- **Learning to Look Around: Intelligently Exploring Unseen Environments for Unknown Tasks**. D. Jayaraman and K. Grauman. In Proceedings of IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Salt Lake City, June 2018. [pdf] [animations]

- **Learning Image Representations Tied to Egomotion from Unlabeled Video**. D. Jayaraman and K. Grauman. International Journal of Computer Vision (IJCV), Special Issue for Best Papers of ICCV 2015, Mar 2017. [pdf] [preprint] [project page, pretrained models]