

# Learning image representations from unlabeled video

Kristen Grauman

Department of Computer Science

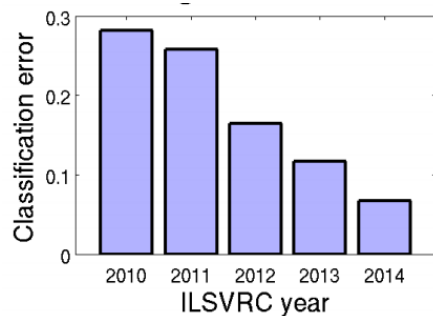
The University of Texas at Austin

Work with Dinesh Jayaraman



# Learning visual categories

- Recent major strides in category recognition



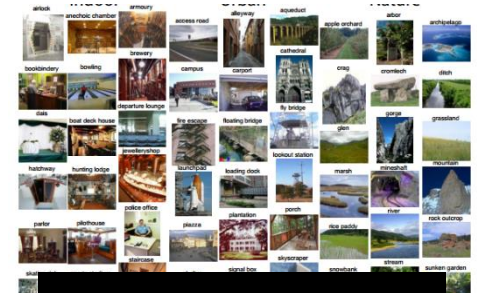
- Facilitated by large labeled datasets



**ImageNet**  
[Deng et al.]



**80M Tiny Images**  
[Torralba et al.]



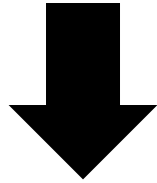
**SUN Database**  
[Xiao et al.]

[Papageorgiou & Poggio 1998, Viola & Jones 2001, Dalal & Triggs 2005, Grauman & Darrell 2005, Lazebnik et al. 2006, Felzenszwalb et al. 2008, Krizhevsky et al. 2012, Russakovsky IJCV 2015...]

# Big picture goal: Embodied vision

## Status quo:

Learn from “disembodied”  
bag of labeled snapshots.



## Our goal:

Learn in the context of **acting**  
and **moving** in the world.



# Beyond “bags of labeled images”?



**Visual development in nature**  
is based on:

- *continuous observation*
- *multi-sensory feedback*
- *motion and action*

... in an environment.

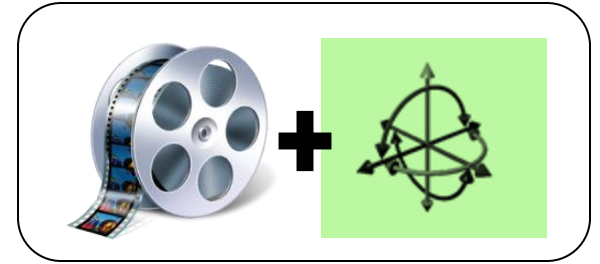
**Inexpensive, and unrestricted in scope**

Evidence from: psychology, evolutionary biology,  
cognitive science.

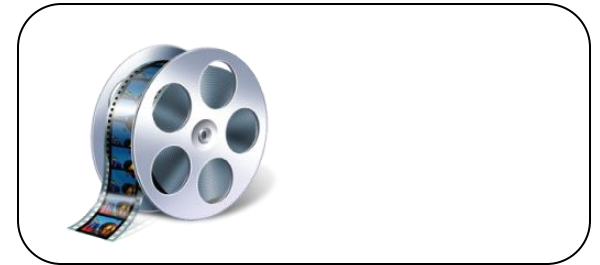


# Talk overview

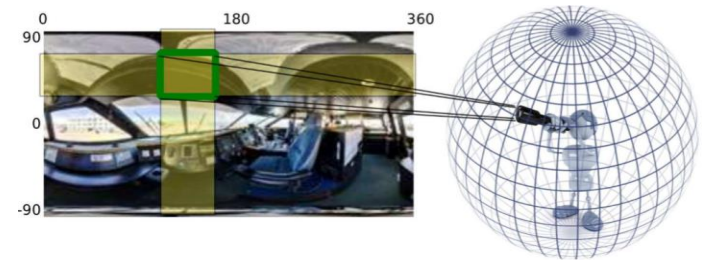
1. Learning representations tied to ego-motion



2. Learning representations from unlabeled video

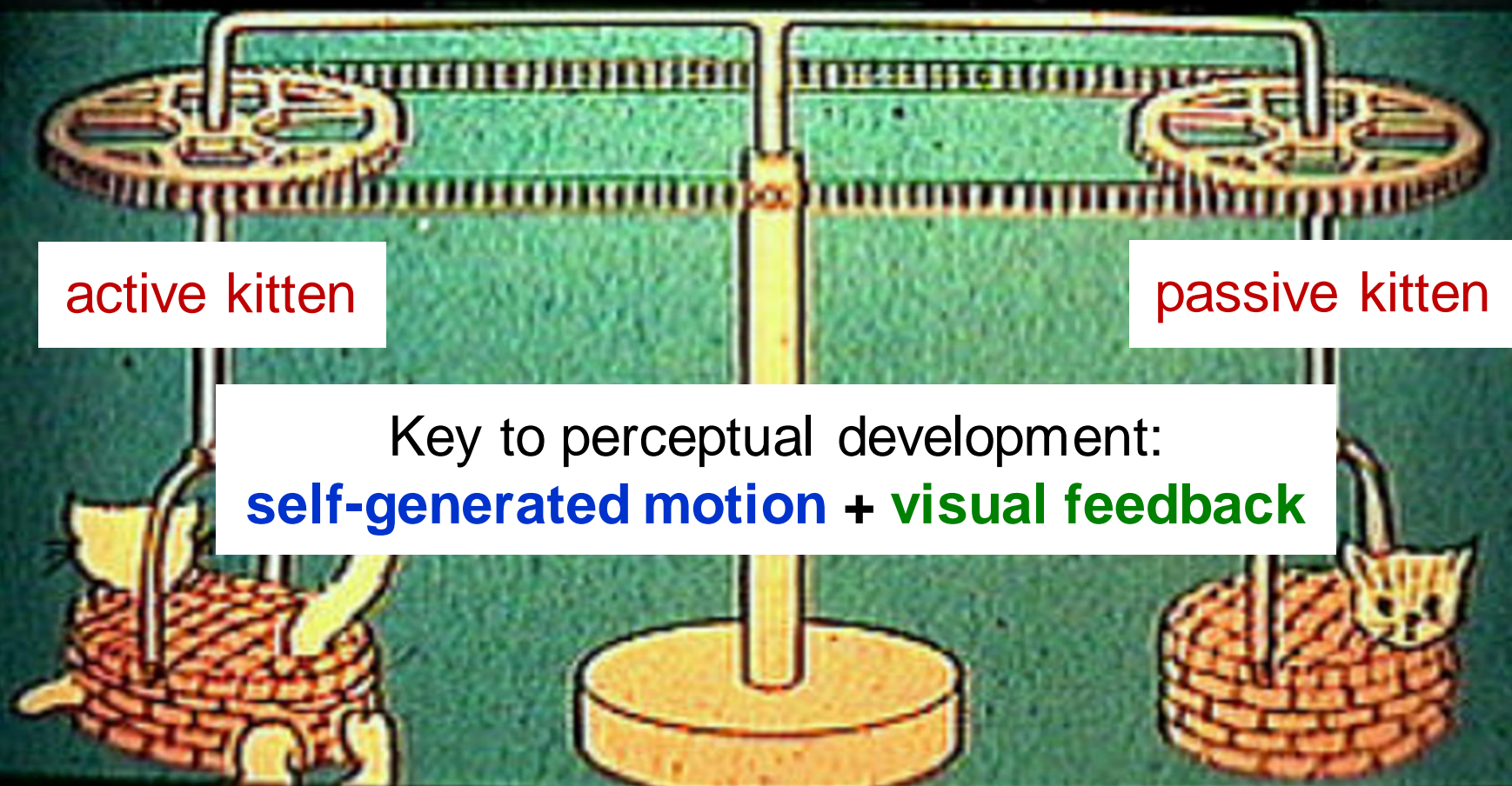


3. Learning how to move and where to look



# The kitten carousel experiment

[Held & Hein, 1963]



# Our idea: **Ego-motion** $\leftrightarrow$ **vision**

**Goal:** Teach computer vision system the connection:  
“**how I move**”  $\leftrightarrow$  “**how my visual surroundings change**”



**Ego-motion motor signals**

+



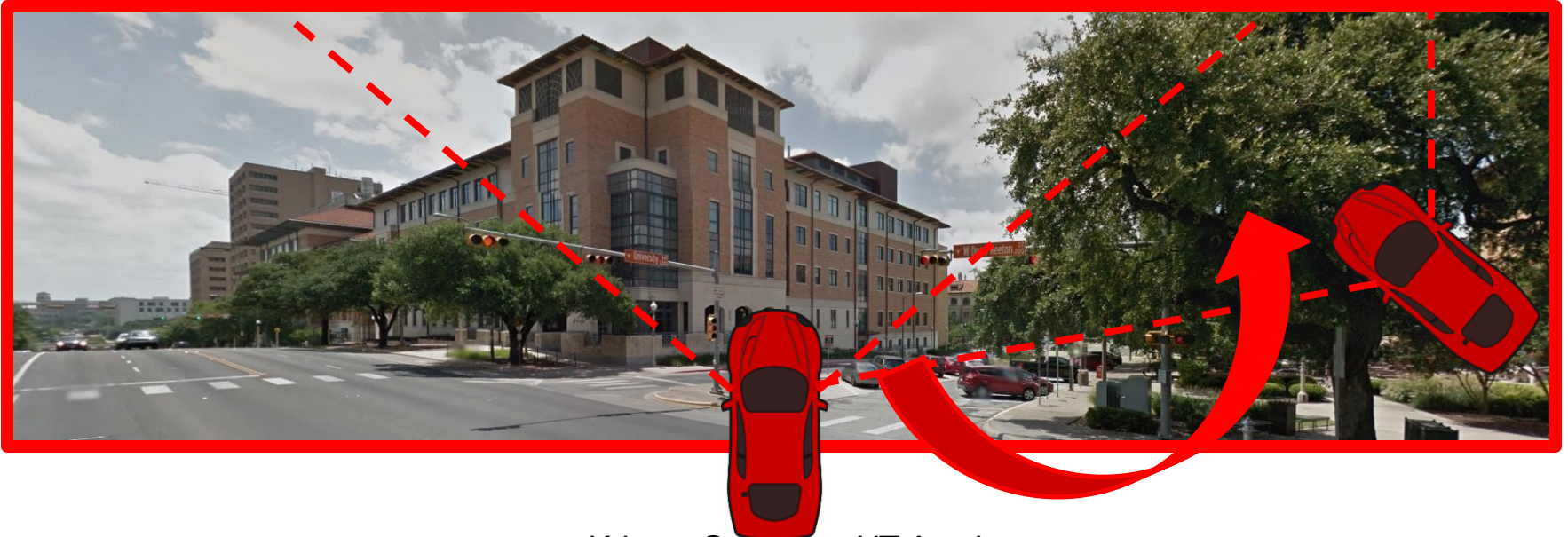
**Unlabeled video**



# Ego-motion $\leftrightarrow$ vision: view prediction



After moving:





# Ego-motion $\leftrightarrow$ vision for recognition

Learning this connection requires:

- Depth, 3D geometry
- Semantics
- Context



Also key to  
recognition!

Can be learned without manual labels!

**Our approach:** unsupervised feature learning  
using egocentric video + motor signals

# Approach idea: Ego-motion equivariance

**Invariant features:** unresponsive to some classes of transformations

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

Simard et al, Tech Report, '98

Wiskott et al, Neural Comp '02

Hadsell et al, CVPR '06

Mobahi et al, ICML '09

Zou et al, NIPS '12

Sohn et al, ICML '12

Cadieu et al, Neural Comp '12

Goroshin et al, ICCV '15

Lies et al, PLoS computation biology '14

...

# Approach idea: Ego-motion equivariance

**Invariant features:** unresponsive to some classes of transformations

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

**Equivariant features:** *predictably* responsive to some classes of transformations, through simple mappings (e.g., linear)

$$\mathbf{z}(g\mathbf{x}) \approx \overset{\text{“equivariance map”}}{\mathbf{M}_g} \mathbf{z}(\mathbf{x})$$

Invariance discards information;  
equivariance organizes it.



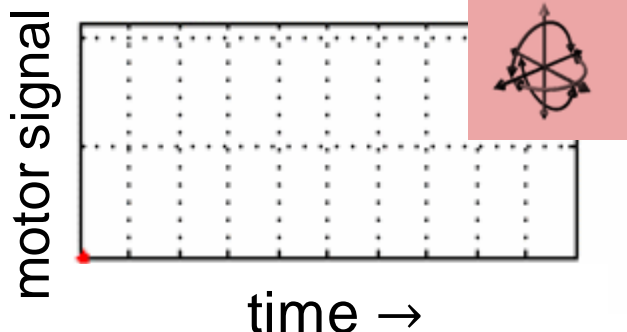
# Approach idea: Ego-motion equivariance

## Training data

Unlabeled video +  
motor signals

## Equivariant embedding

organized by ego-motions



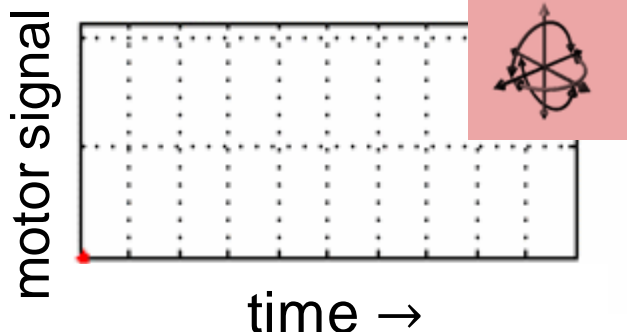
Learn

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

# Approach idea: Ego-motion equivariance

## Training data

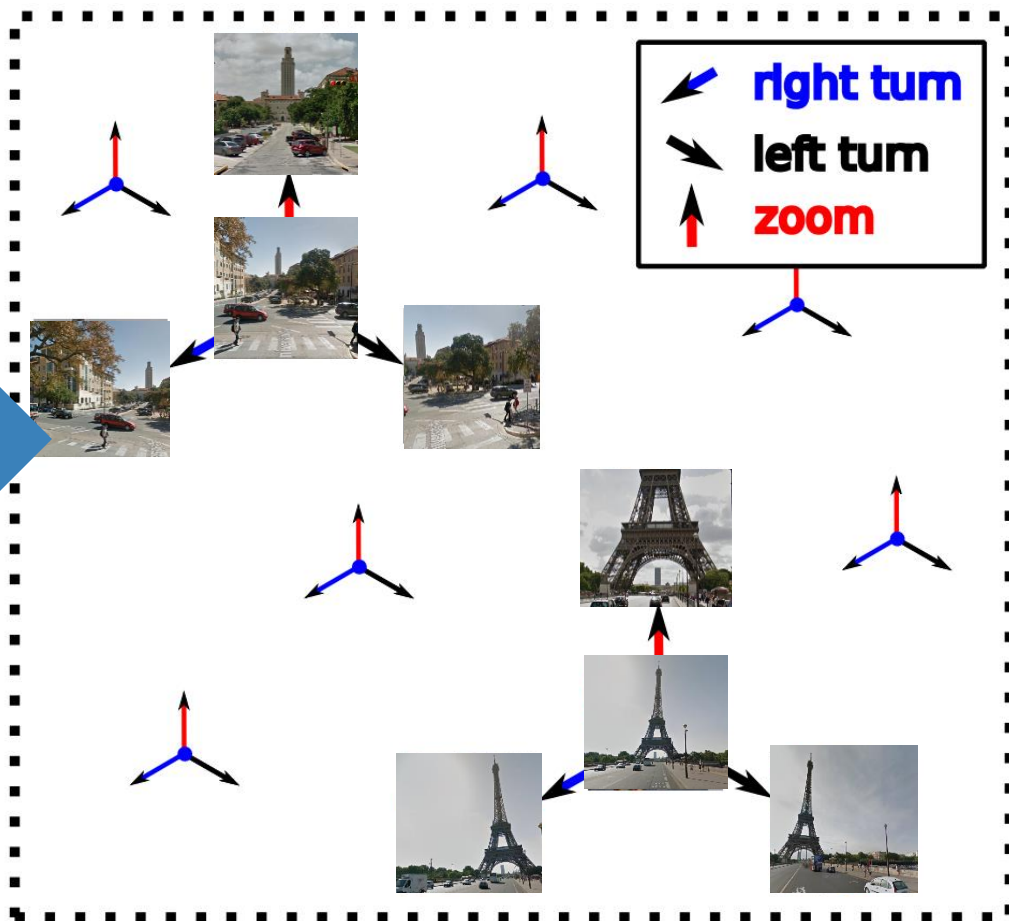
Unlabeled video +  
motor signals



Learn

## Equivariant embedding

organized by ego-motions



# Approach overview

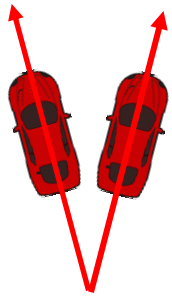
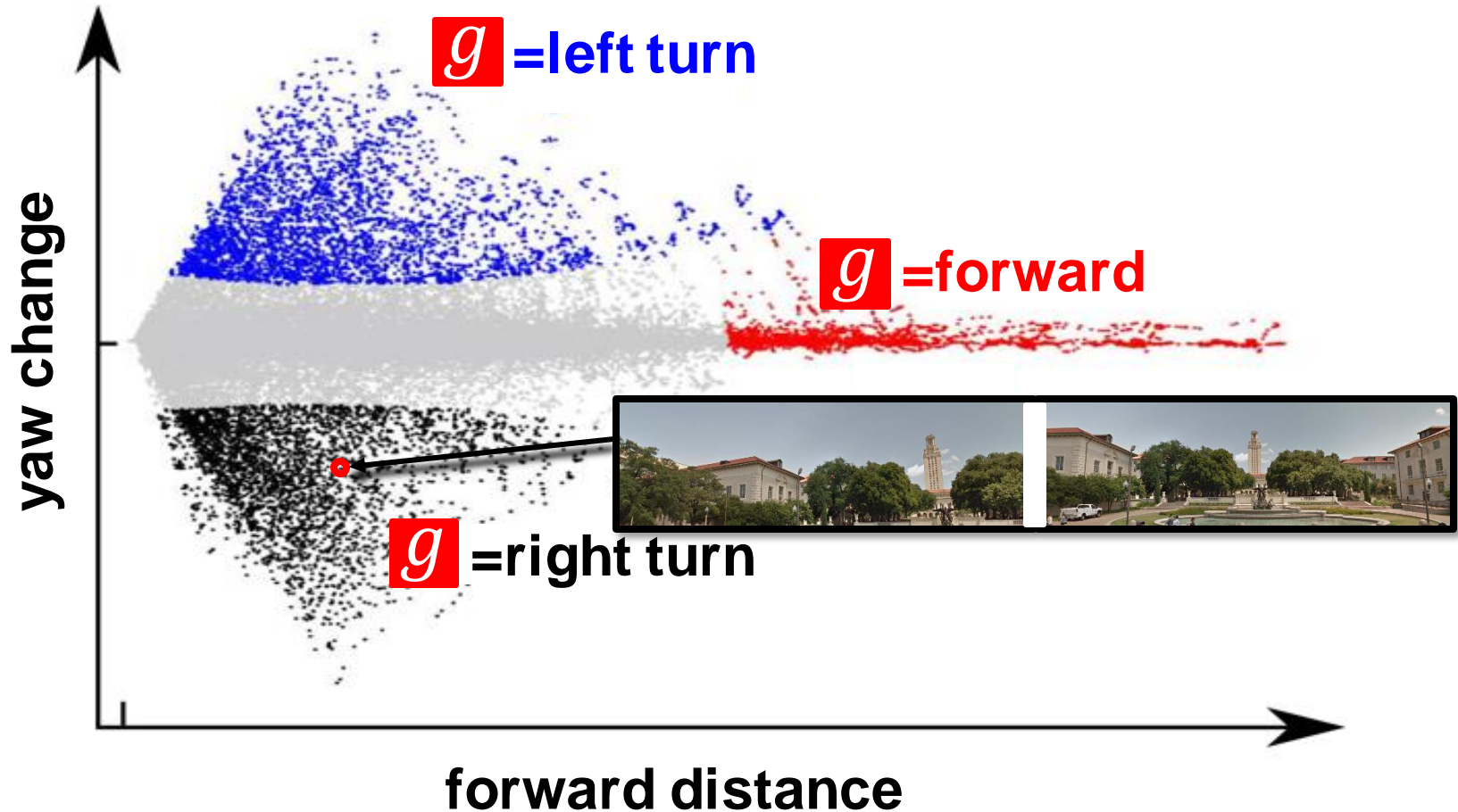
**Our approach:** unsupervised feature learning using egocentric video + motor signals

1. Extract training frame pairs from video
2. Learn ego-motion-equivariant image features
3. Train on target recognition task in parallel



# Training frame pair mining

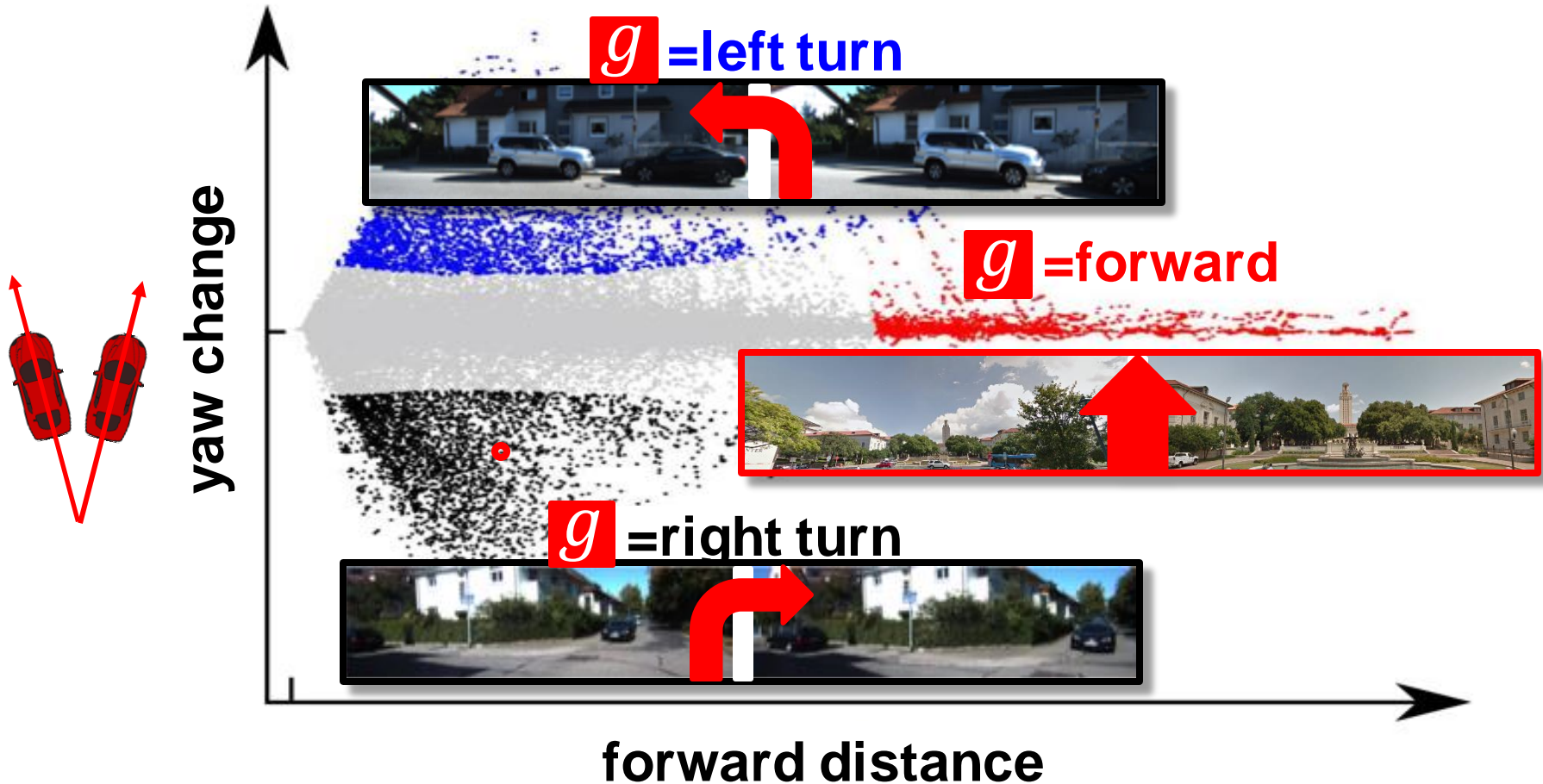
## Discovery of ego-motion clusters



Kristen Grauman, UT Austin

# Training frame pair mining

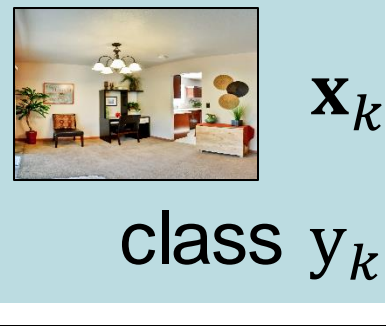
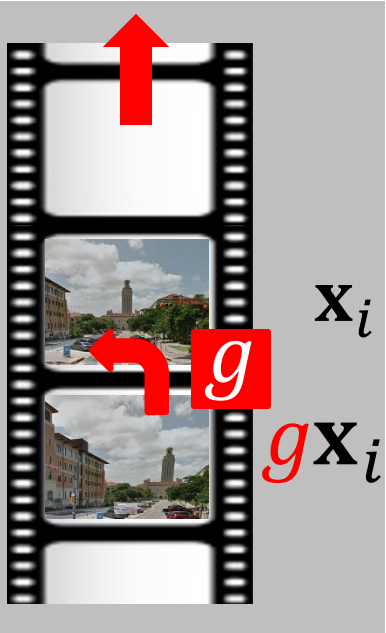
## Discovery of ego-motion clusters



Kristen Grauman, UT Austin

# Ego-motion equivariant feature learning

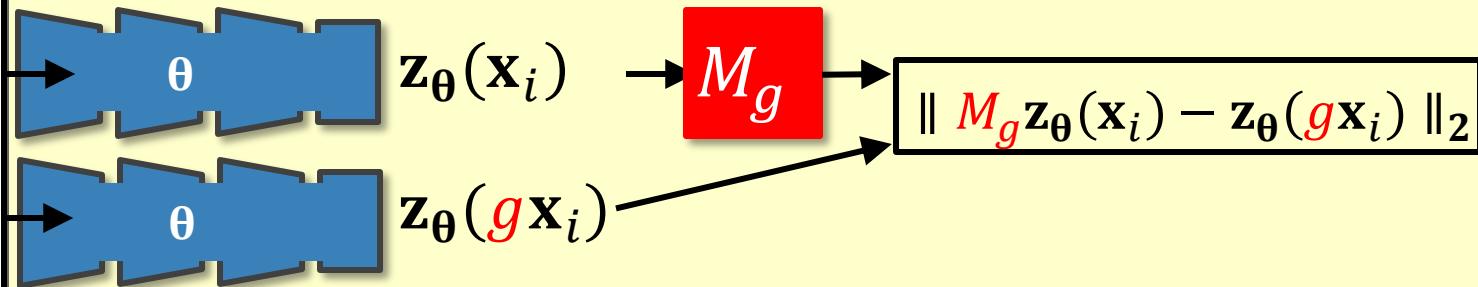
**Given:**



**Desired:** for all motions  $g$  and all images  $\mathbf{x}$ ,

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

**Unsupervised training**



**Supervised training**



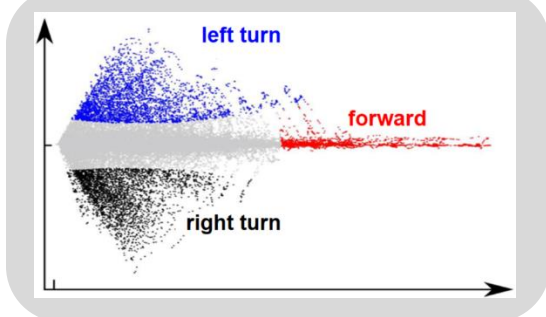
$\theta$ ,  $M_g$  and  $W$  jointly trained



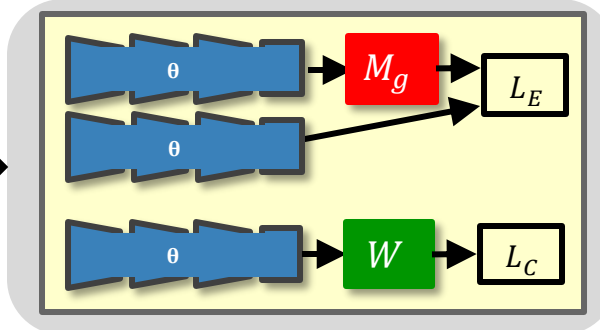
# Method recap

APPROACH

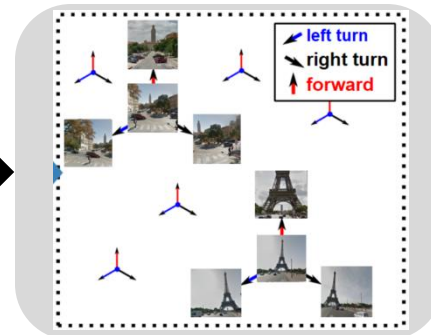
Ego-motion training pairs



Neural network training



Equivariant embedding



Scene and object recognition



Football field?  
Pagoda?  
Airport?  
Cathedral?  
Army base?

Next-best view selection



cup

frying pan

RESULTS

# Datasets

## KITTI video

[Geiger et al. 2012]

Car platform

Egomotions: yaw and forward distance



## SUN images

[Xiao et al. 2010]

Large-scale scene classification task with 397 categories (static images)

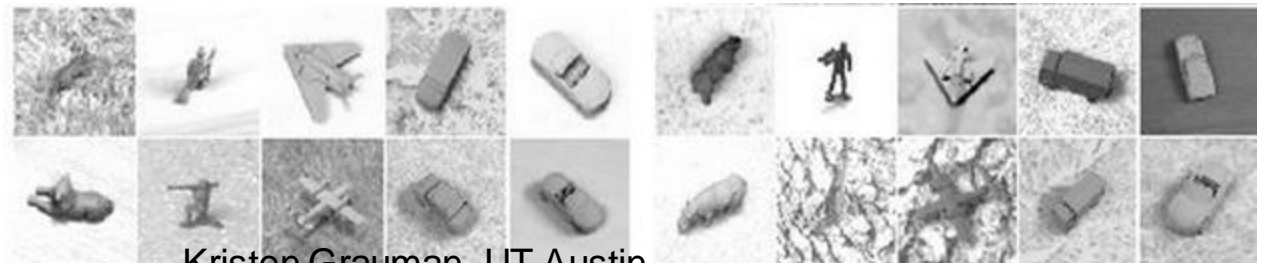


## NORB images

[LeCun et al. 2004]

Toy recognition

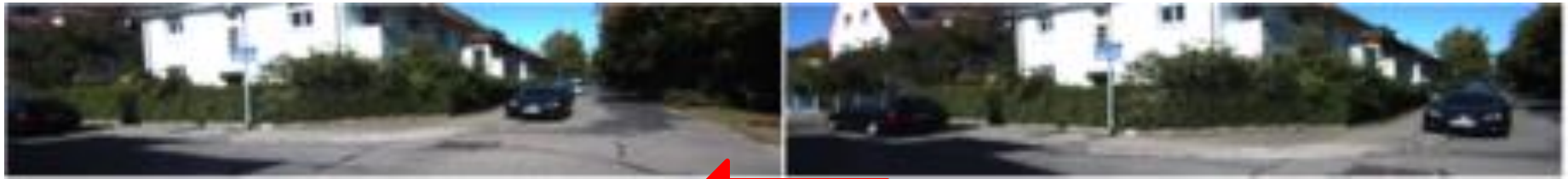
Egomotions: elevation and azimuth



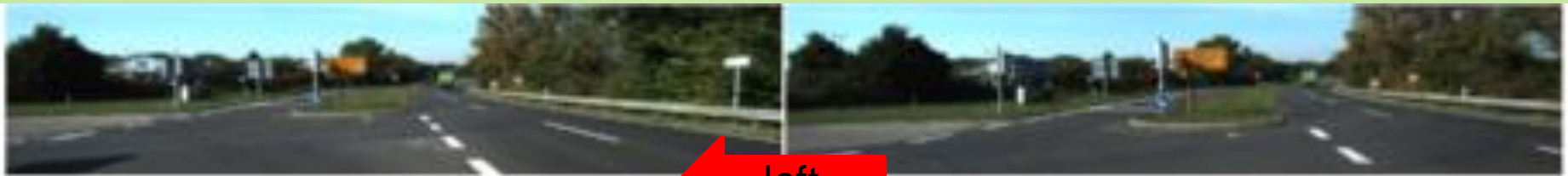
Kristen Grauman, UT Austin

# Results: Equivariance check

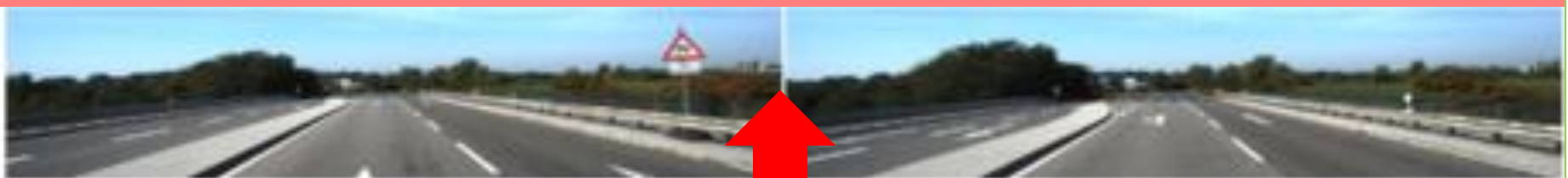
Visualizing how well equivariance is preserved



**Query pair**



**Neighbor pair (our features)**



**Neighbor pair (pixel space)**





# Results: Equivariance check

How well is equivariance preserved?

Methods↓	atomic	composite
random	1.0000	1.0000
CLSNET	0.9239	0.9145
TEMPORAL [19]	0.7587	0.8119
DRLIM [7]	0.6404	0.7263
EQUIV	<b>0.6082</b>	<b>0.6982</b>
EQUIV+DRLIM	<b>0.5814</b>	<b>0.6492</b>

Recognition loss only →

Temporal coherence ↗ ↘

Ours ↗ ↘

Normalized error:

$$\rho_g = E \left[ \|\mathbf{z}_\theta(\mathbf{x}) - M'_g \mathbf{z}_\theta(g\mathbf{x})\|_2 / \|\mathbf{z}_\theta(\mathbf{x}) - \mathbf{z}_\theta(g\mathbf{x})\|_2 \right]$$

Temporal coherence: Hadsell et al. CVPR 2006, Mohabi et al. ICML 2009

# Results: Recognition

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



Geiger et al, IJRR '13

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



Apse

Window seat

Art school

Library

Auditorium

Bus interior

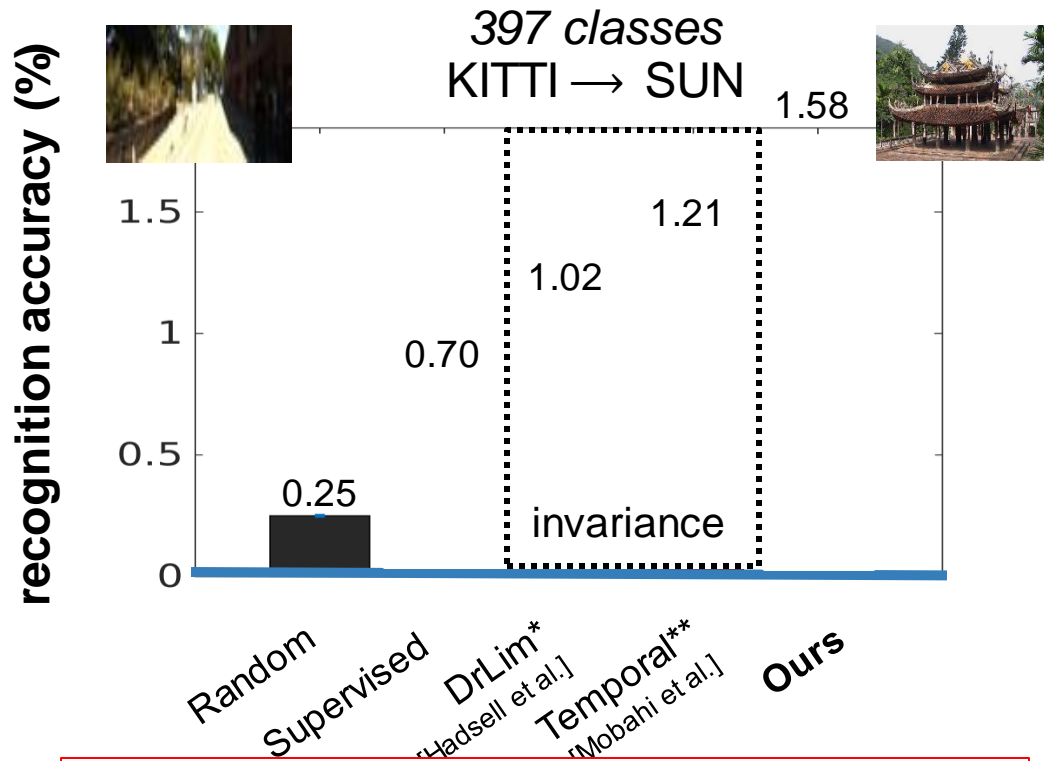
Cathedral

Freeway

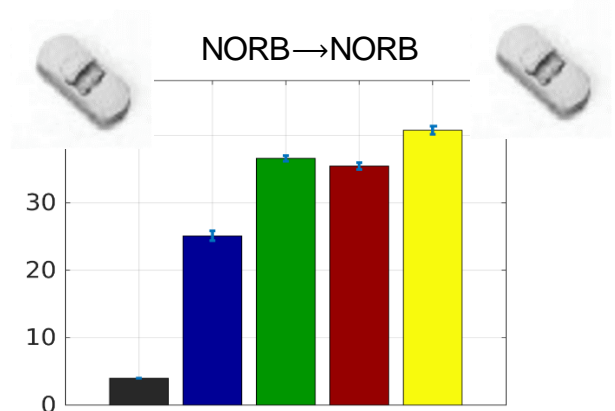
Guardhouse

# Results: Recognition

Do ego-motion equivariant features improve recognition?



6 labeled training examples per class



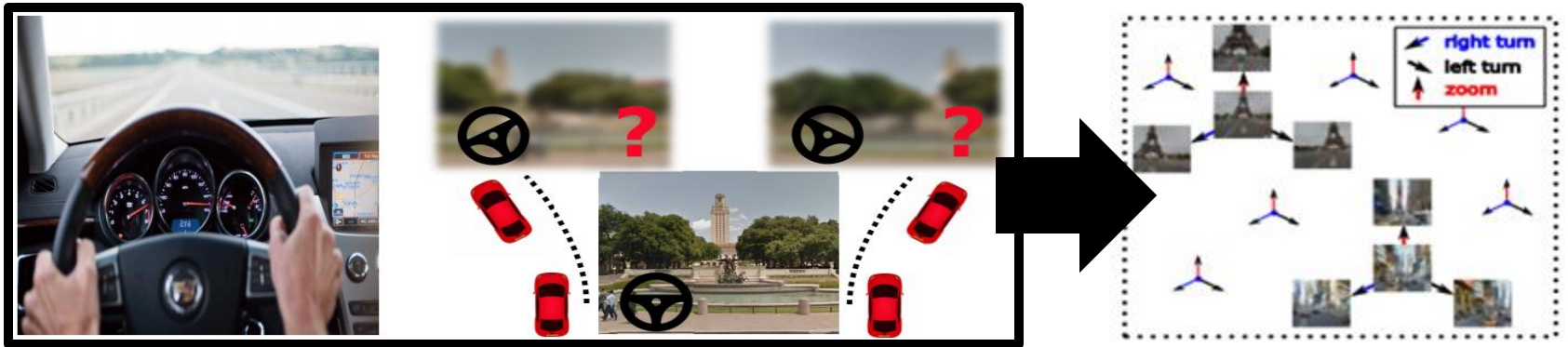
**Up to 30% accuracy increase  
over state of the art!**

\*Hadsell et al., Dimensionality Reduction by Learning an Invariance

\*\*Mobahi et al., Deep Learning from Temporal Coherence in Video, ICML'09

Kristen Grauman, UT Austin

# Recap so far

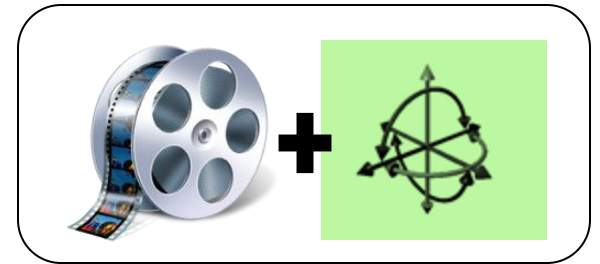


<http://vision.cs.utexas.edu/projects/egoequiv/>

- New *embodied visual feature learning* paradigm
- *Ego-motion equivariance* boosts performance across multiple challenging recognition tasks
- Future work: volition at training time too

# Talk overview

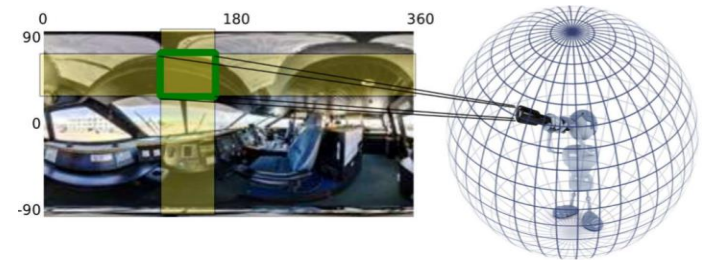
1. Learning representations tied to ego-motion



2. Learning representations from unlabeled video

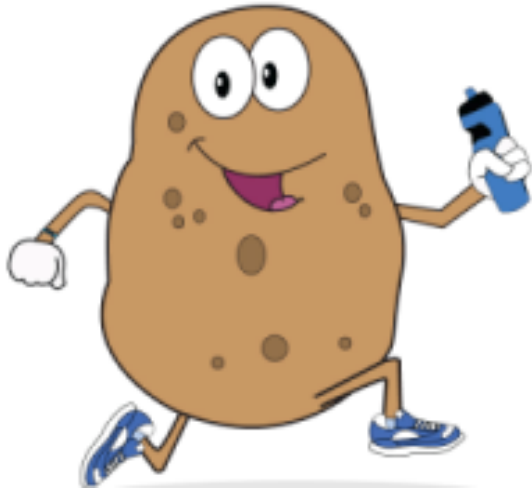


3. Learning how to move and where to look





# Learning from arbitrary unlabeled video?

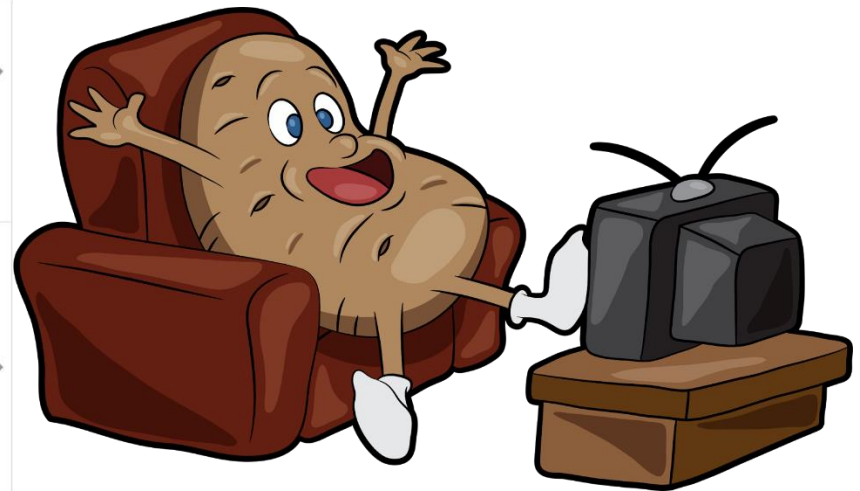
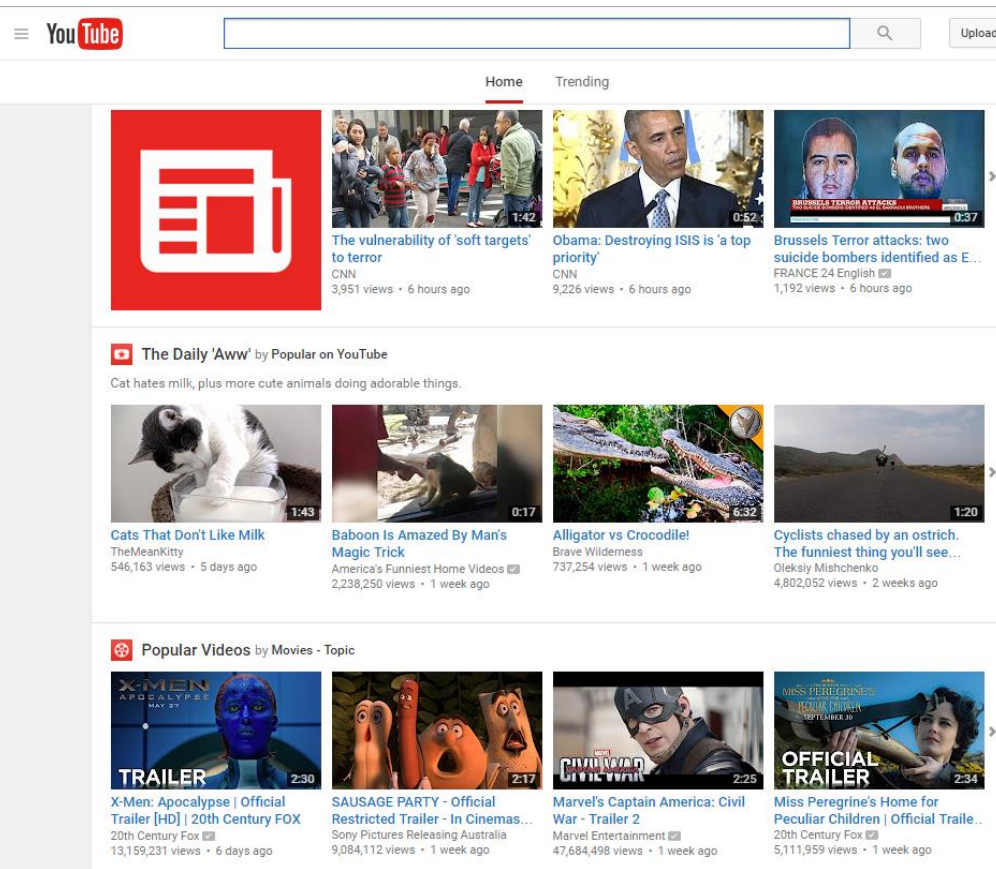


**Unlabeled video  
+ ego-motion**



**Unlabeled video**

# Learning from arbitrary unlabeled video?

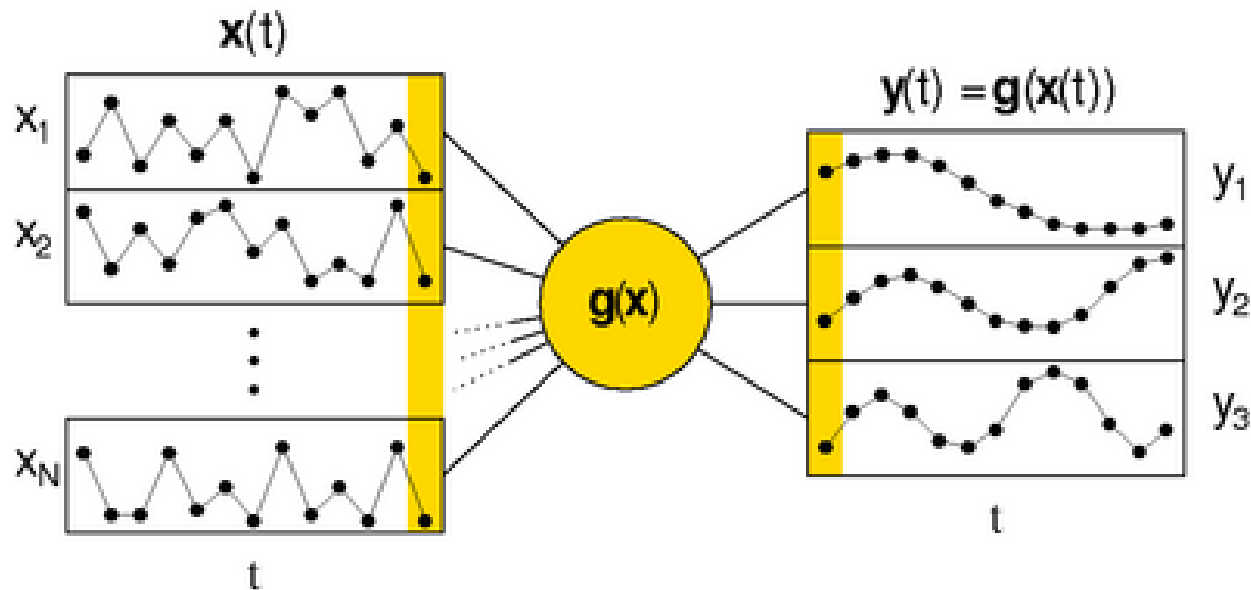


Unlabeled video

# Background: Slow feature analysis

[Wiskott & Sejnowski, 2002]

Find functions  $\mathbf{g}(\mathbf{x})$  that map



quickly varying input  
signal  $\mathbf{x}(t)$

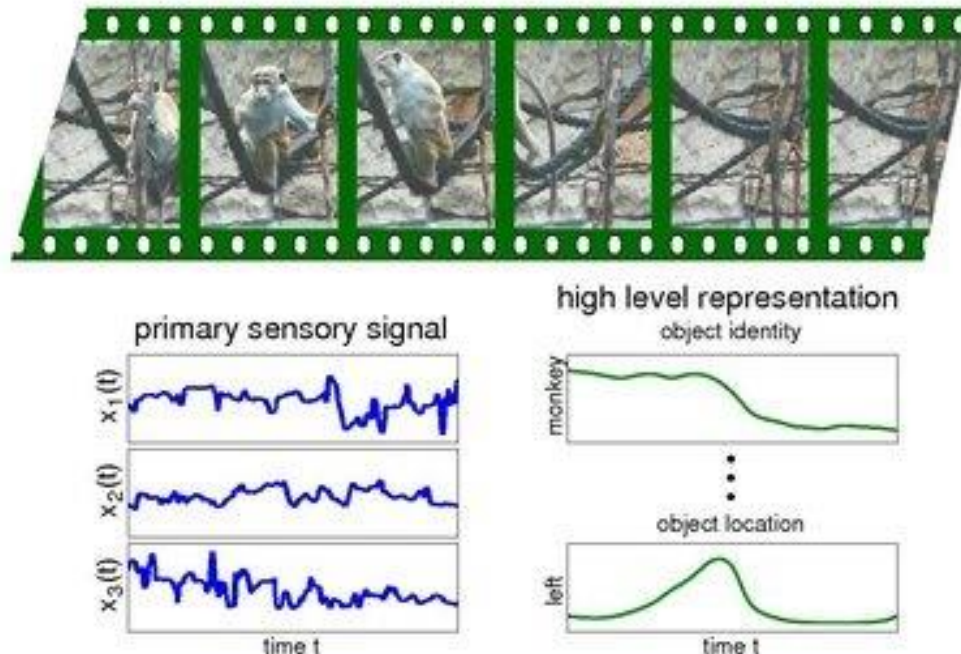


slowly varying  
features  $\mathbf{y}(t)$

# Background: Slow feature analysis

*[Wiskott & Sejnowski, 2002]*

Find functions  $\mathbf{g}(\mathbf{x})$  that map



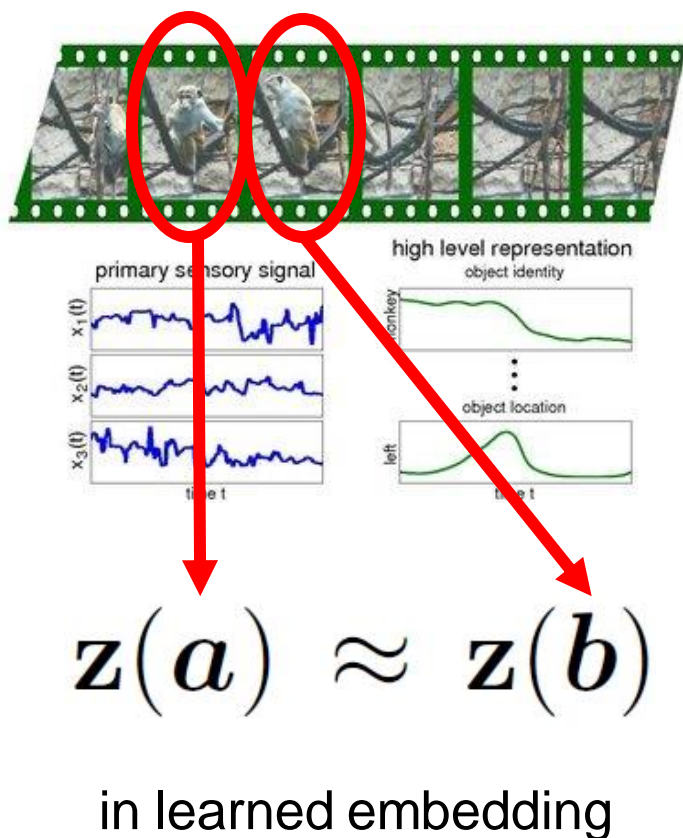
quickly varying input  
signal  $\mathbf{x}(t)$



slowly varying  
features  $\mathbf{y}(t)$

# Background: Slow feature analysis

[Wiskott & Sejnowski, 2002]



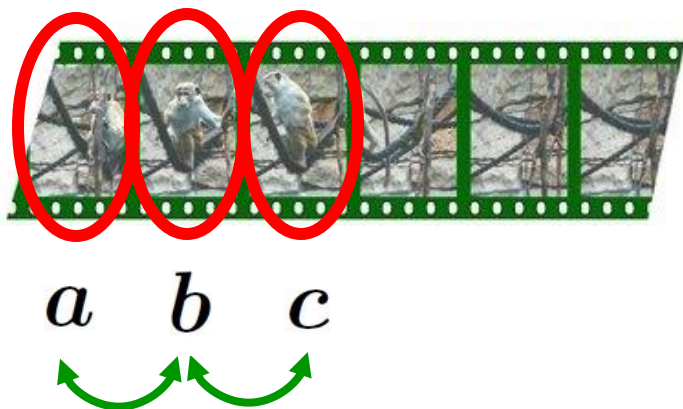
- Existing work exploits “slowness” as **temporal coherence** in video → learn invariant representation

[Hadsell et al. 2006; Mobahi et al. 2009; Bergstra & Bengio 2009; Goroshin et al. 2013; Wang & Gupta 2015, ...]

- Fails to capture *how* visual content changes over time



# Our idea: **Steady** feature analysis



- Higher order temporal coherence in video → learn equivariant representation

Second order slowness operates on frame triplets:

$$\mathbf{z}(b) - \mathbf{z}(a) \approx \mathbf{z}(c) - \mathbf{z}(b)$$

in learned embedding

[Jayaraman & Grauman, CVPR 2016]

# Approach: Steady feature analysis

Learn classifier  $W$  and representation  $\theta$  jointly,

$$(\theta^*, W^*) = \arg \min_{\theta, W} L_s(\theta, W, \mathcal{S}) + \lambda L_u(\theta, \mathcal{U})$$

with unsupervised regularization loss:

$$L_u(\theta, \mathcal{U}) = \underbrace{R_2(\theta, \mathcal{U})}_{\text{slow}} + \lambda' \underbrace{R_3(\theta, \mathcal{U})}_{\text{steady}}$$

slow

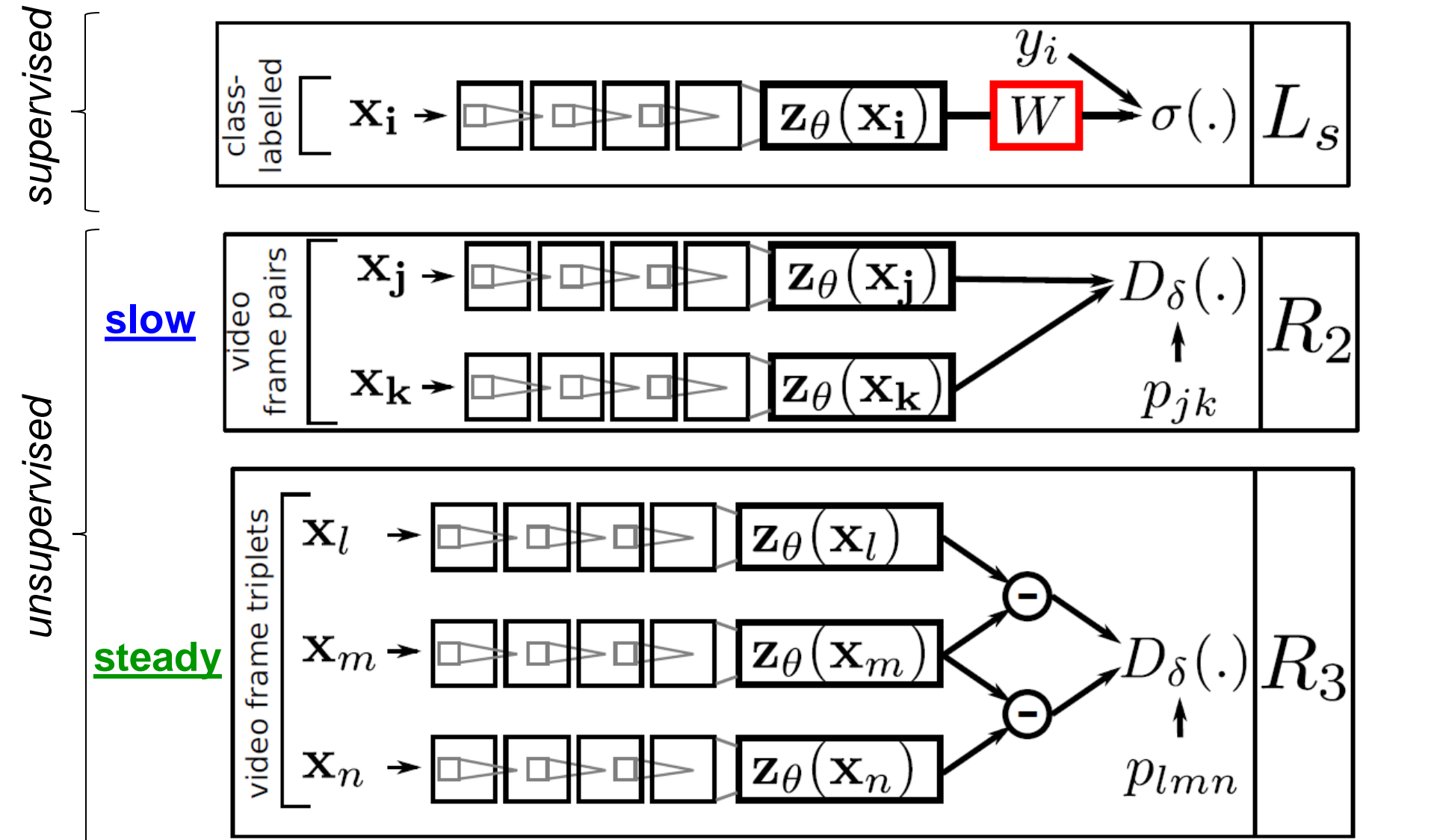
$$\sum_{(j,k) \in \mathcal{U}_2} D_\delta(\mathbf{z}_\theta(\mathbf{x}_j), \mathbf{z}_\theta(\mathbf{x}_k), p_{jk})$$

steady

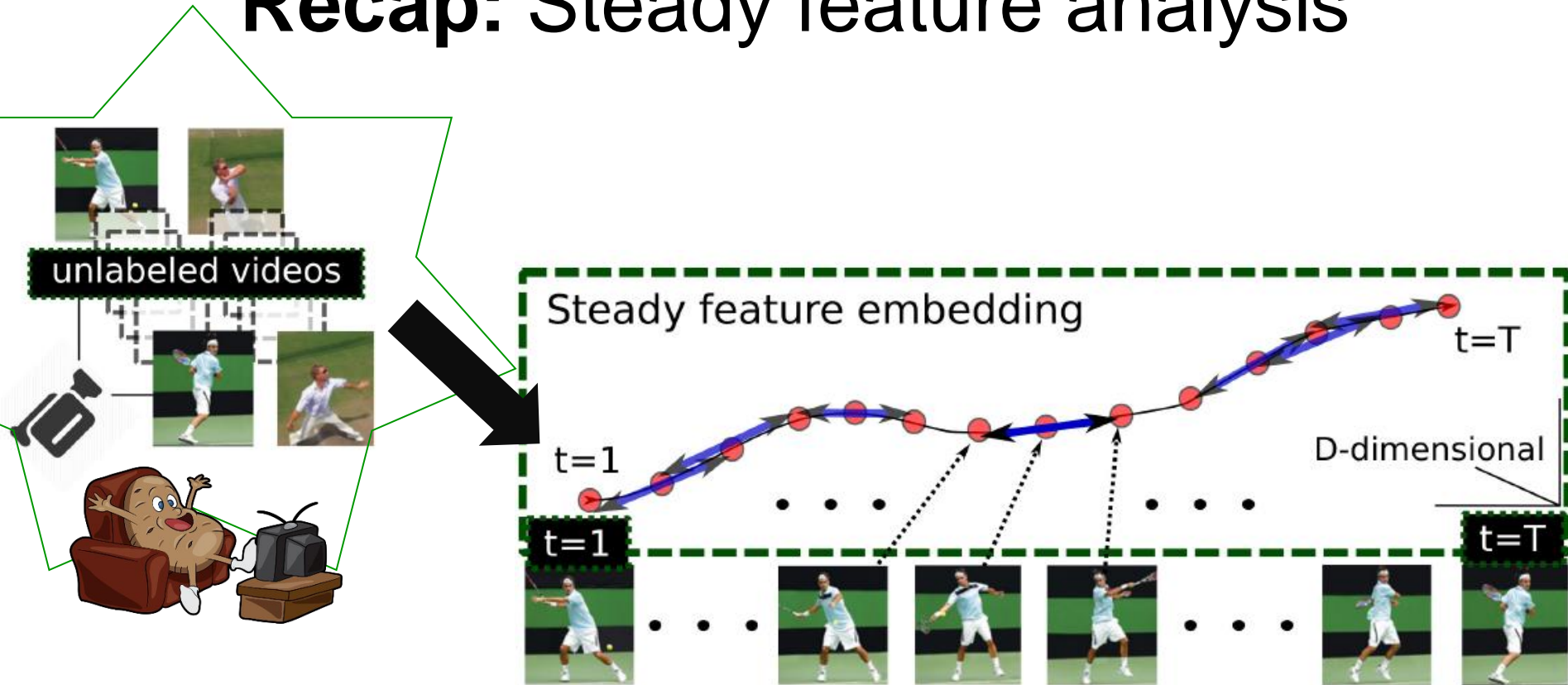
$$\sum_{(l,m,n) \in \mathcal{U}_3} D_\delta(\mathbf{z}_{\theta_l} - \mathbf{z}_{\theta_m}, \mathbf{z}_{\theta_m} - \mathbf{z}_{\theta_n}, p_{lmn})$$

**Contrastive loss**  
that also exploits  
“negative” tuples

# Approach: Steady feature analysis



# Recap: Steady feature analysis



Equivariance  $\approx$  “steadily” varying frame features!

$$d^2 \mathbf{z}_\theta(\mathbf{x}_t) / dt^2 \approx \mathbf{0}$$

[Jayaraman & Grauman, CVPR 2016]

# Datasets

## Unlabeled video



**Human Motion  
Database (HMDB)**



**KITTI Video**



**NORB**

## Target task (few labels)



**PASCAL 10 Actions**



**SUN 397 Scenes**



**NORB 25 Objects**

32 x 32 images or 96 x 96 images



# Results: Sequence completion

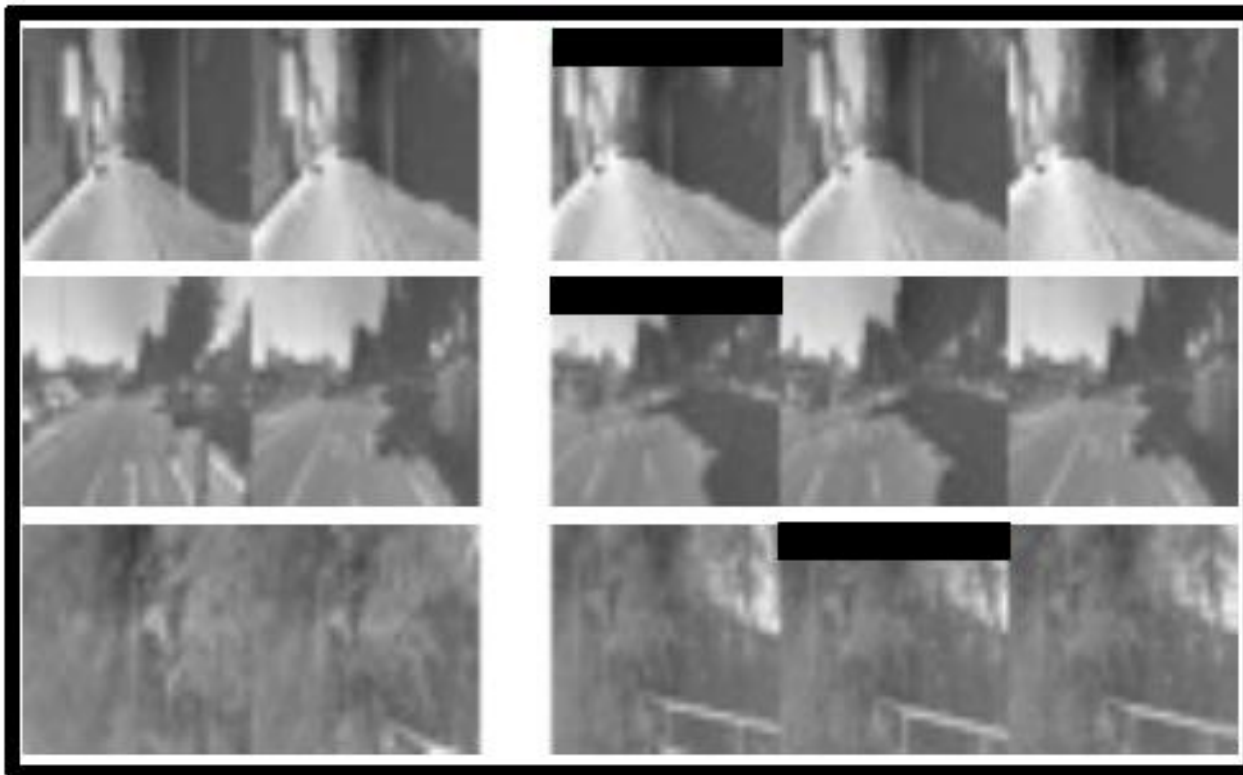
Given sequential pair, infer next frame (embedding)

$$\tilde{\mathbf{z}}_{\theta}(\mathbf{x}_3) = 2\mathbf{z}_{\theta}(\mathbf{x}_2) - \mathbf{z}_{\theta}(\mathbf{x}_1)$$

$\mathbf{x}_1$

$\mathbf{x}_2$

Our top 3 estimates for  $\mathbf{x}_3$



# Results: Sequence completion

Given sequential pair, infer next frame (embedding)

	Datasets→	NORB	KITTI	HMDB
<u>slow</u>	SFA-1 [30] *	0.95	31.04	2.70
<u>slow</u>	SFA-2 [14] **	0.91	8.39	2.27
<u>slow &amp; steady</u>	SSFA (ours)	<b>0.53</b>	<b>7.79</b>	<b>1.78</b>

Percentile rank of correct completion (lower is better)

\*Hadsell et al., Dimensionality Reduction by Learning an Invariant Mapping, CVPR'06

\*\*Mobahi et al., Deep Learning from Temporal Coherence in Video, ICML'09

Kristen Grauman, UT Austin

# Results: Recognition



Task type→	Objects	Scenes		Actions
Datasets→	NORB→NORB	KITTI→SUN		HMDB→PASCAL-10
Methods↓	[25 cls]	[397 cls]	[397 cls, top-10]	[10 cls]
random	4.00	0.25	2.52	10.00
UNREG	24.64±0.85	0.70±0.12	6.10±0.67	15.34±0.28
SFA-1 [30]*	37.57±0.85	1.21±0.14	8.24±0.25	19.26±0.45
SFA-2 [14]**	39.23±0.94	1.02±0.12	6.78±0.32	19.04±0.24
SSFA (ours)	<b>42.83±0.33</b>	<b>1.65±0.04</b>	<b>9.19±0.10</b>	<b>20.95±0.13</b>

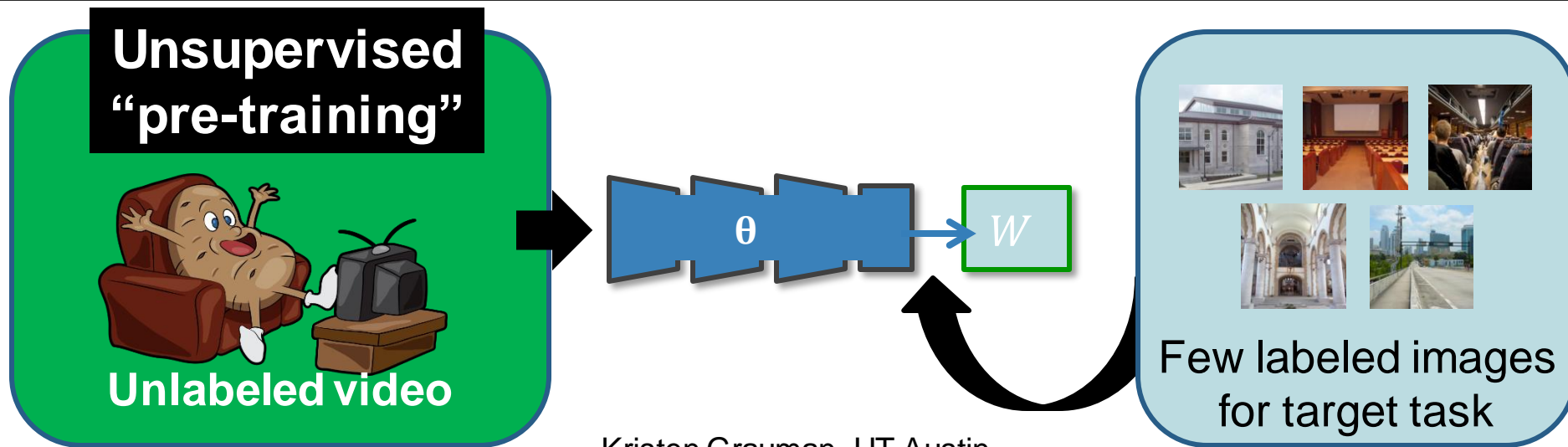
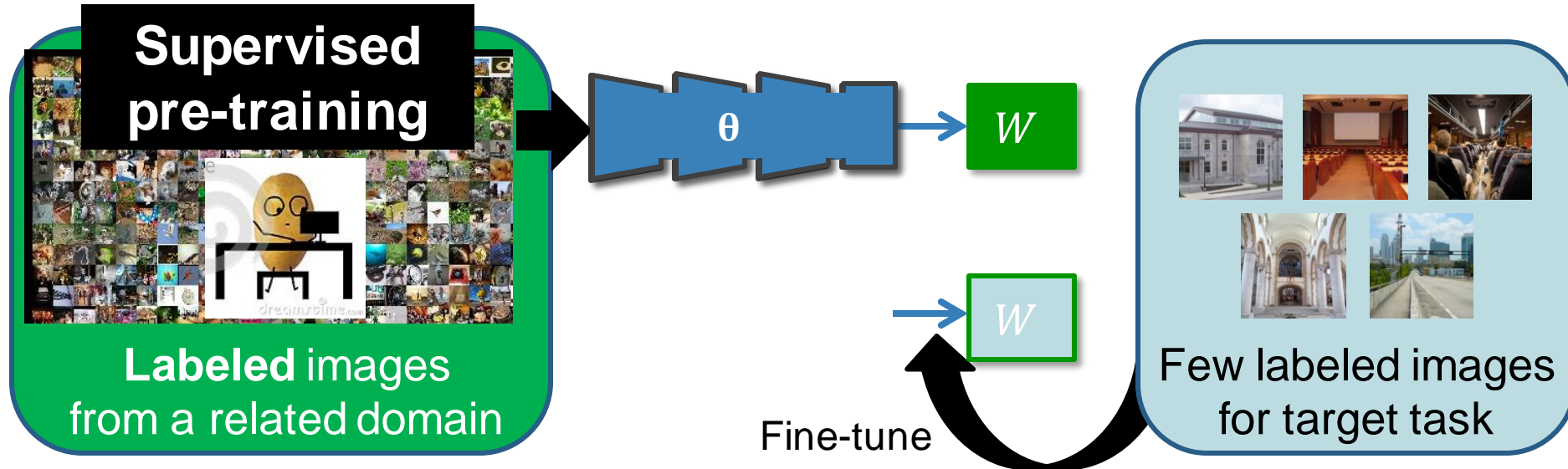
Multi-class recognition accuracy

\*Hadsell et al., Dimensionality Reduction by Learning an Invariant Mapping, CVPR'06

\*\*Mobahi et al., Deep Learning from Temporal Coherence in Video, ICML'09

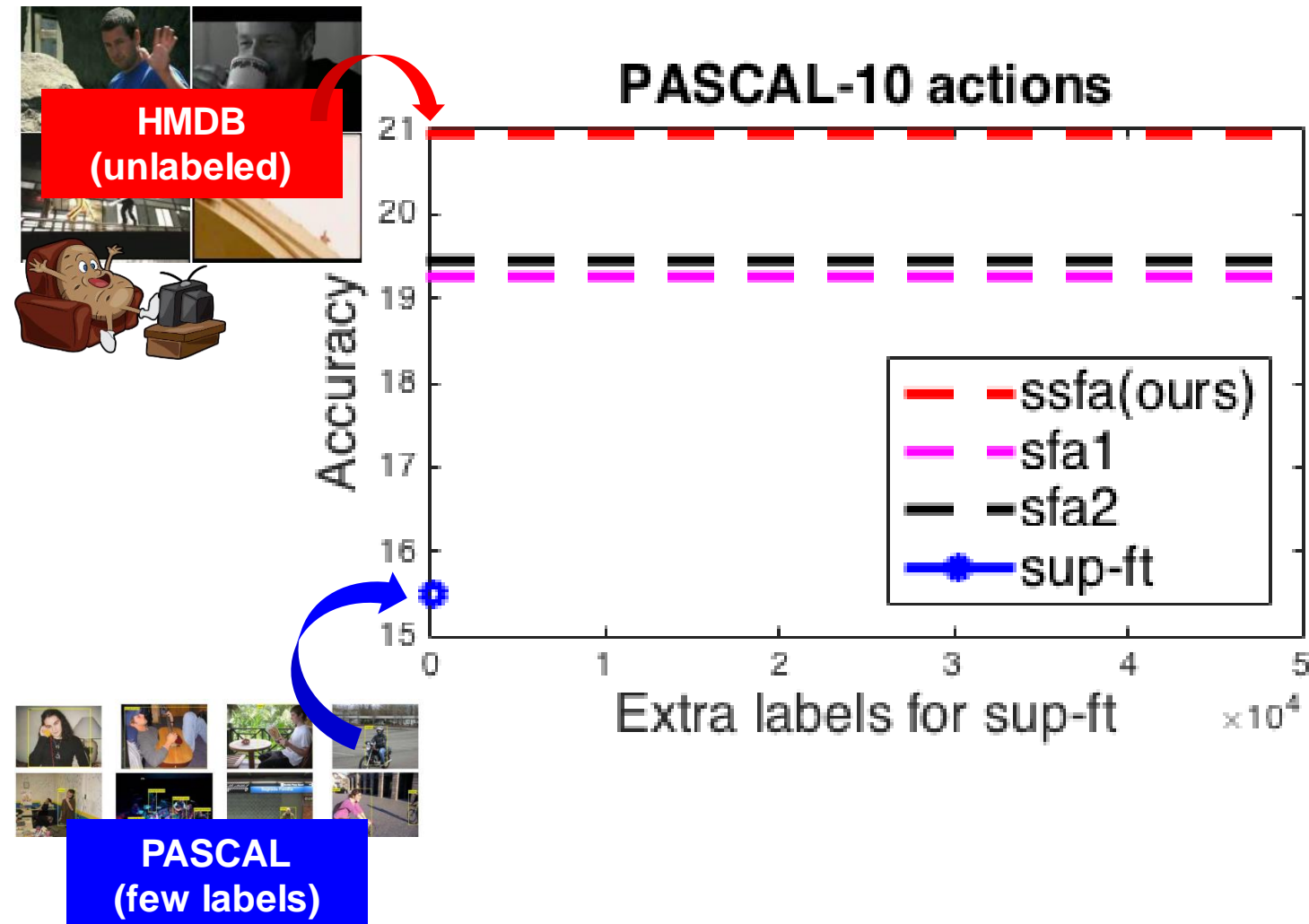
Kristen Grauman, UT Austin

# Pre-training a representation

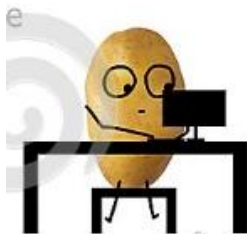
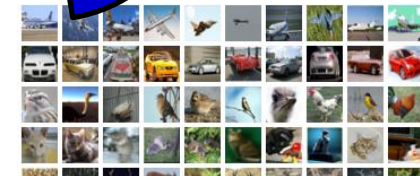
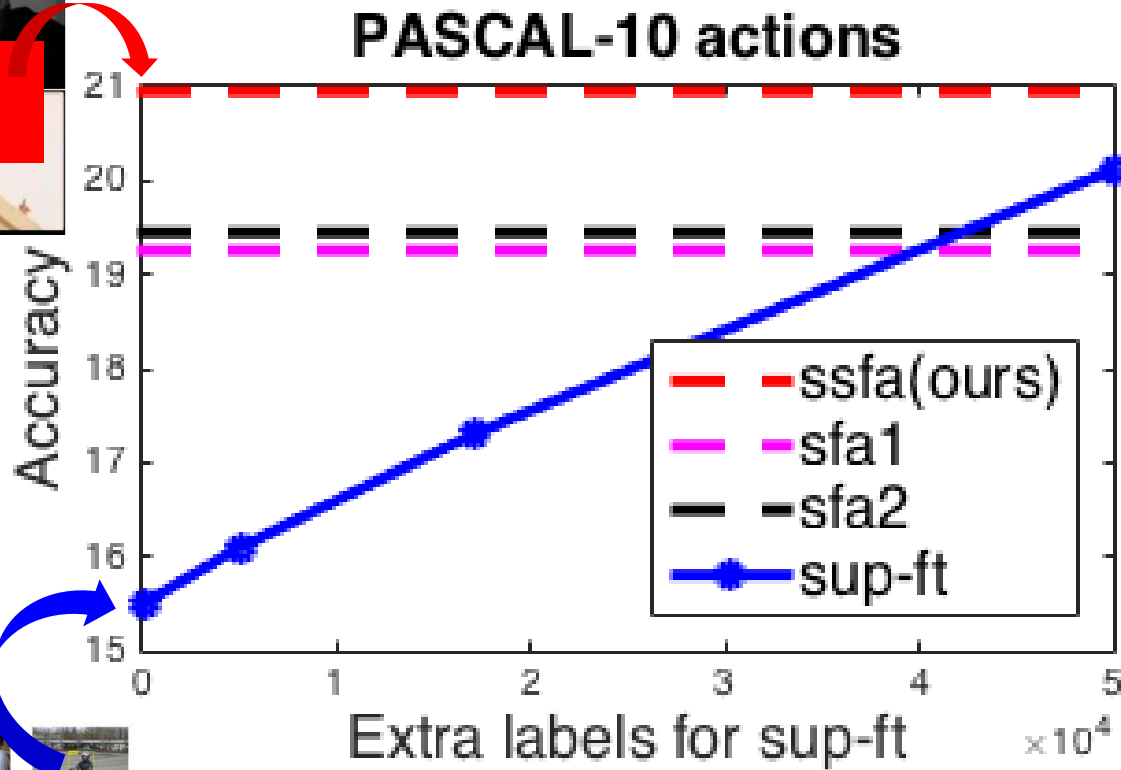




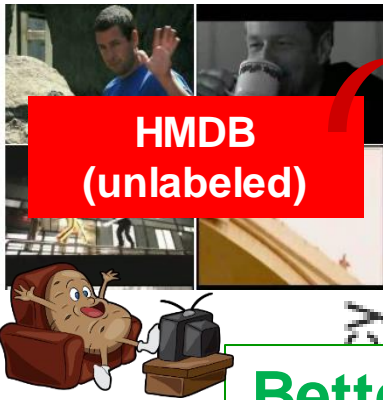
# Results: Can we learn *more* from unlabeled video than “related” labeled images?



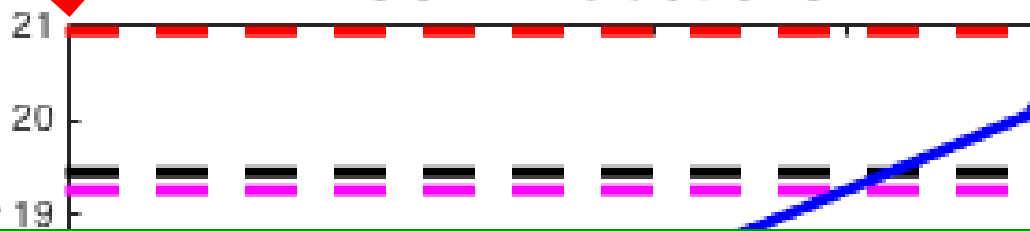
# Results: Can we learn *more* from unlabeled video than “related” labeled images?



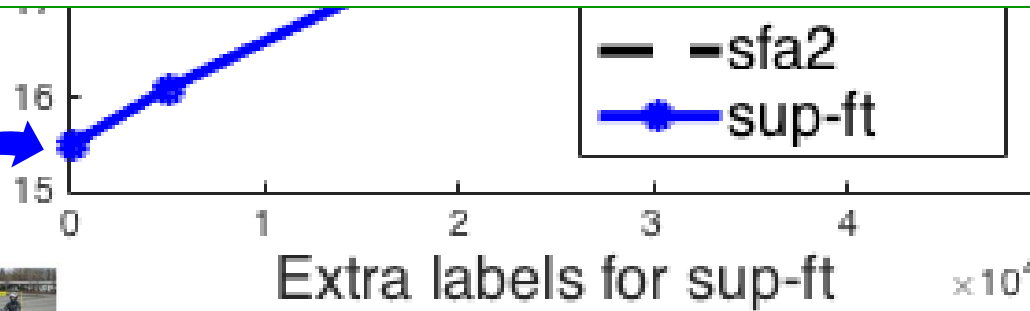
# Results: Can we learn *more* from unlabeled video than “related” labeled images?



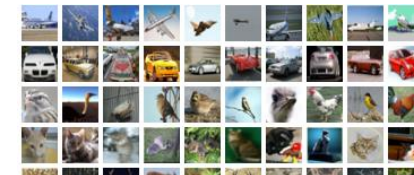
PASCAL-10 actions



Better even than providing 50,000 extra manual labels for auxiliary classification task!



Extra labels for sup-ft  $\times 10^4$

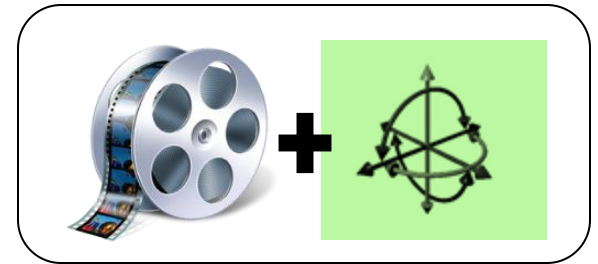


CIFAR-100 (labeled for other categories)



# Talk overview

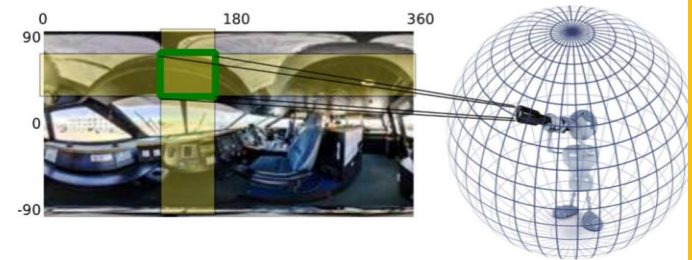
1. Learning representations tied to ego-motion



2. Learning representations from unlabeled video



3. Learning how to move and where to look



# Learning how to move for recognition



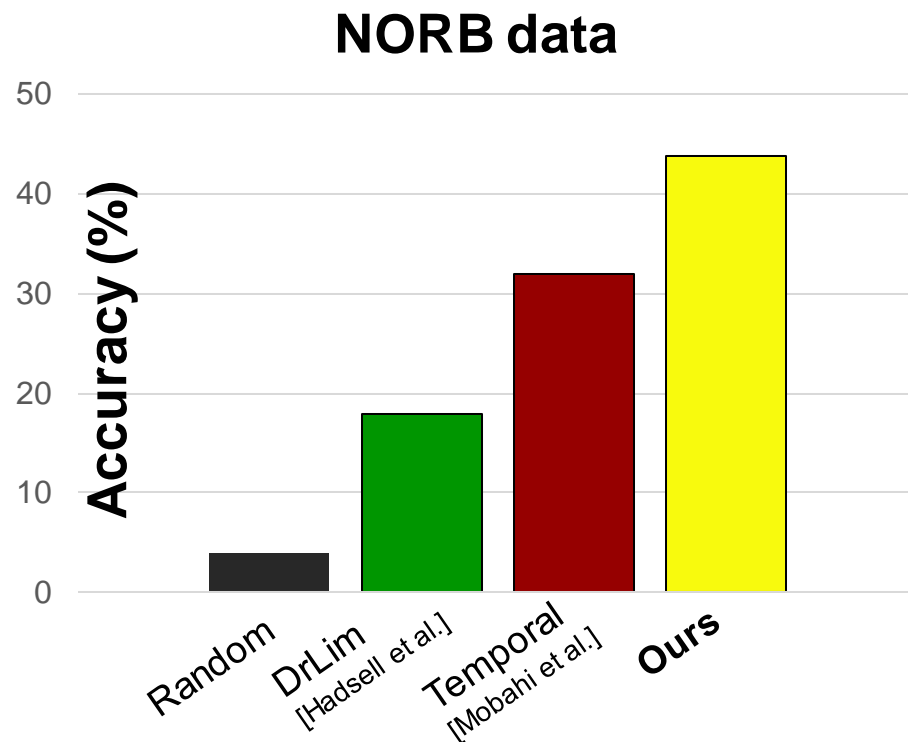
Time to revisit **active recognition** in  
challenging settings!

[Bajcsy 1985, Schiele & Crowley 1998, Dickinson et al. 1997, Tsotsos et al. 2001, Soatto 2009,...]



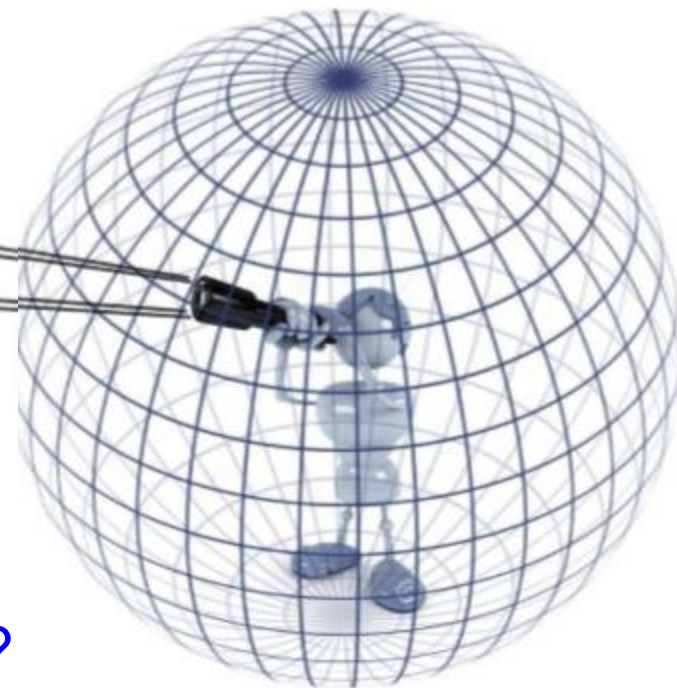
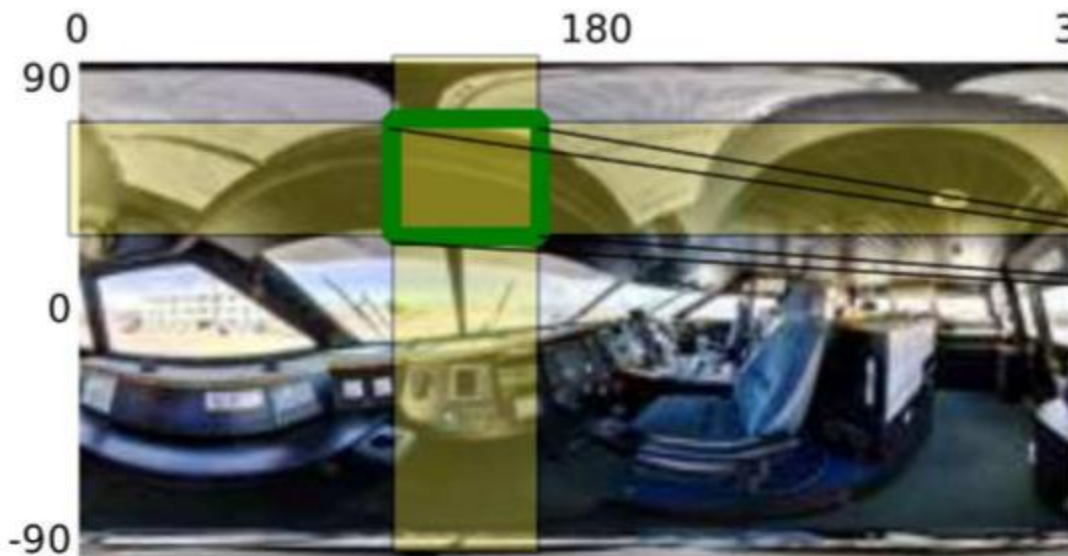
# Learning how to move for recognition

Leverage proposed ego-motion equivariant  
embedding to **select next best view**



[Jayarman & Grauman, ICCV 2015]

# Learning how to move for recognition



Best sequence of glimpses in 3D scene?

## Requires:

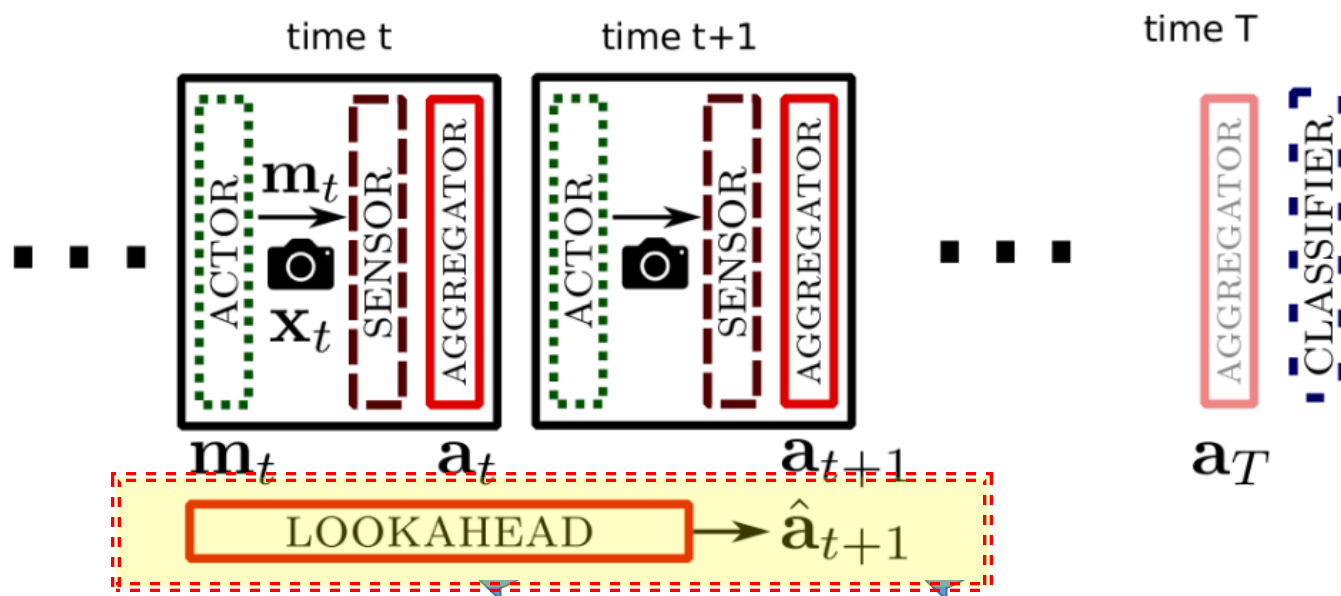
- Action selection
- Per-view processing
- Evidence aggregation
- Look-ahead prediction
- Final class belief prediction

**Learn all end-to-end**

Jayaraman and Grauman, UT TR AI15-06

Kristen Grauman, UT Austin

# Active visual recognition



Requires several separate functionalities:

- Action selection
- Per-view processing
- Across-view evidence aggregation
- Next-view prediction
- Final class belief prediction

Learn all end-to-end

# Active recognition: example results

P("Plaza courtyard"):

Top 3 guesses:

(6.28)

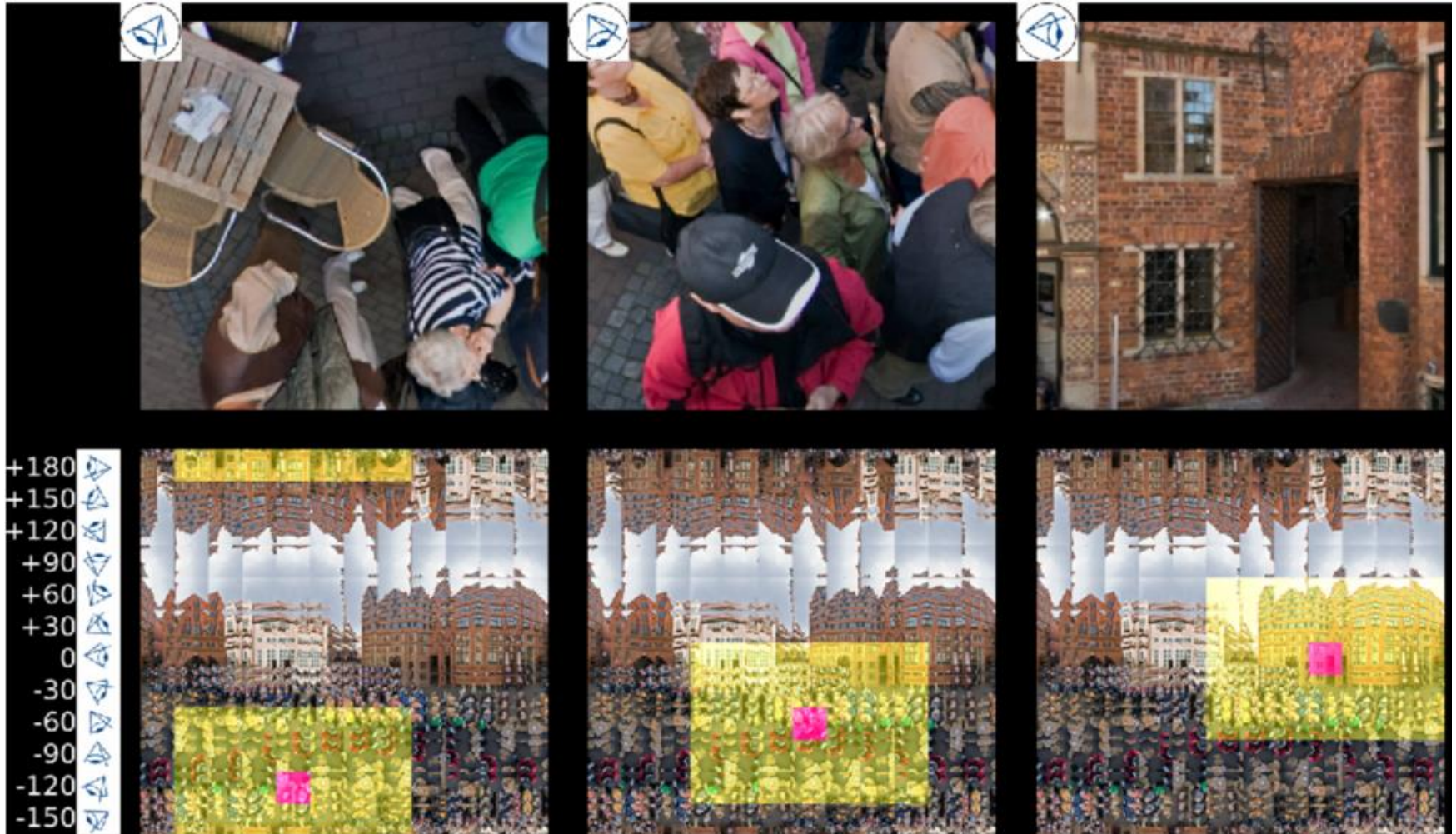
Restaurant  
Train interior  
Shop

(11.95)

Theater  
Restaurant  
Plaza courtyard

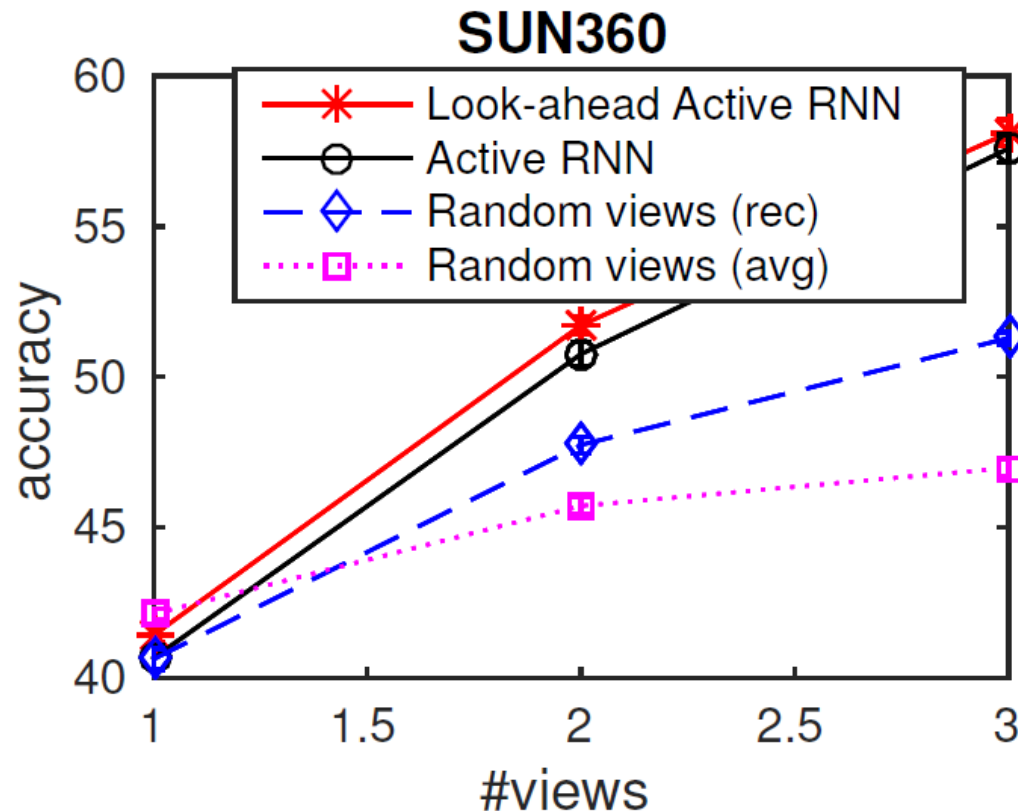
(68.38)

Plaza courtyard  
Street  
Theater





# Active recognition: Results



Active selection + look-ahead → better scene categorization from sequence of glimpses in 360 panorama



# Summary

- Visual learning requires
  - context of action and motion in the world
  - with continuous self-acquired feedback
- New ideas:
  - “Embodied” feature learning using both visual and motor signals
  - Feature learning from unlabeled video via higher order temporal coherence
  - Steps towards active view selection in 360 scenes

# References

- Learning Image Representations Tied to Ego-Motion. D. Jayaraman and K. Grauman. In Proceedings of the IEEE International Conference on Computer Vision (ICCV), Santiago, Chile, Dec 2015. (Oral)
- Slow and Steady Feature Analysis: Higher Order Temporal Coherence in Video. D. Jayaraman and K. Grauman. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), Las Vegas, June 2016. (Spotlight)
- Look Ahead Before You Leap: End-to-End Active Recognition by Forecasting the Effect of Motion. D. Jayaraman and K. Grauman. UT Tech Report A115-06, Dec 2015.