# Slow and steady feature analysis: higher order temporal coherence in video

## Dinesh Jayaraman and Kristen Grauman
### UT Austin

THE UNIVERSITY OF TEXAS AT AUSTIN

## Problem

Learning *unsupervised* generic visual features from unlabeled video



**Status quo** — Learning from "bags of labeled images"
- Expensive
- Limited data
- Task-specific

**Desired** — Learning from continuous observation
- Free
- Unlimited data
- Generic

## Learning from natural world temporal dynamics

### Prior work

**Slowness:** "concepts only change slowly over time", i.e., given $k = 1$ video frame, next frame is close in a semantic space $z$:

$k = 1$
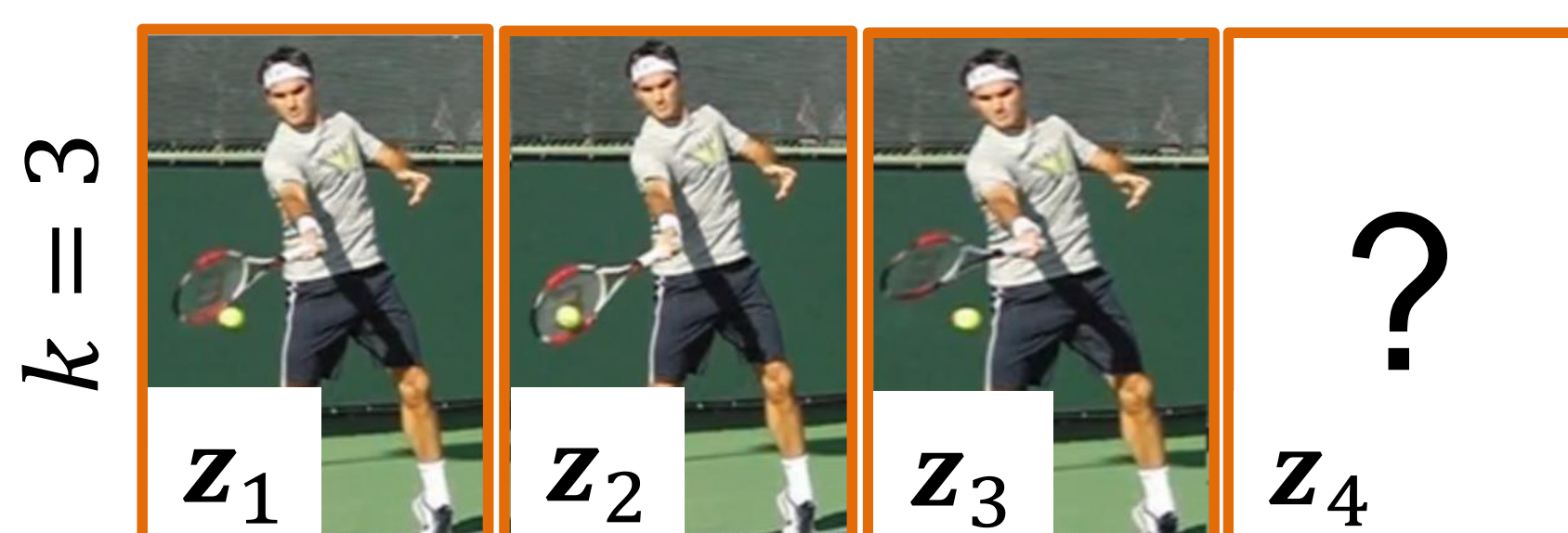


$z_1$ $z_2$

Predict: $\widehat{z_2} \approx z_1$

[Wiskott 2002], [Hadsell 2006], [Mobahi 2009], [Bergstra 2009], [Goroshin 2013], [Wang 2015] …

### Our idea

**Steadiness:** "concepts evolve '*steadily*' over time", i.e., given $k$ frames, perform $k^{th}$ order extrapolation to guess next frame:

$k = 2$



$z_1$ $z_2$ $z_3$

Predict: $\widehat{z_3} \approx z_2 + (z_2 - z_1)$

$k = 3$



$z_1$ $z_2$ $z_3$ $z_4$

Predict: $\widehat{z_4} \approx z_3 + (z_3 - z_2) + ((z_3 - z_2) - (z_2 - z_1))$

**Idea:** train a feature mapping $z(.)$ to be a "steady" semantic space i.e. $k^{th}$ order extrapolation must predict future frames well:
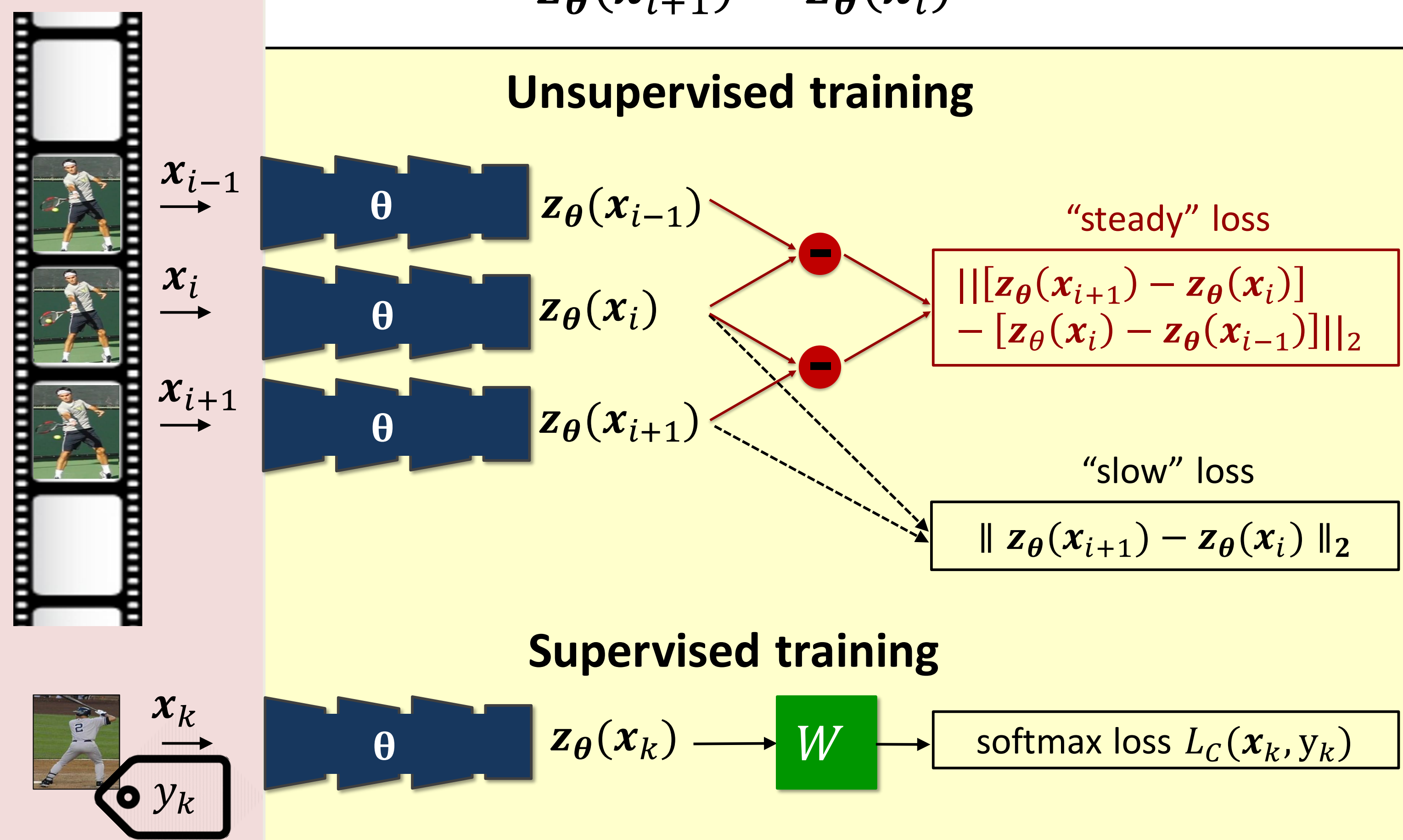- Induces desirable feature properties such as invariance ($k = 1$), equivariance ($k = 2$) and so on.

## Target representation

**Input**

unlabeled videos

**Output**

Steady feature embedding

$t=1$ … $t=T$

D-dimensional

$t=1$ … $t=T$

"$K^{th}$ order steady" features $z \Rightarrow \dfrac{dz}{dt} \approx \mathbf{0}, \dfrac{d^2 z}{dt^2} \approx \mathbf{0}, \dots, \dfrac{d^K z}{dt^K} \approx \mathbf{0}$

## Steady representation learning

**Desired:** Feature space with second-order temporal coherence:

$$z_\theta(x_{i+1}) - z_\theta(x_i) \approx z_\theta(x_i) - z_\theta(x_{i-1}), \text{ and}$$
$$z_\theta(x_{i+1}) \approx z_\theta(x_i)$$

**Given:**

**Unsupervised training**

$x_{i-1} \rightarrow \theta \rightarrow z_\theta(x_{i-1})$

$x_i \rightarrow \theta \rightarrow z_\theta(x_i)$

$x_{i+1} \rightarrow \theta \rightarrow z_\theta(x_{i+1})$

"steady" loss
$$\|[z_\theta(x_{i+1}) - z_\theta(x_i)] - [z_\theta(x_i) - z_\theta(x_{i-1})]\|_2$$

"slow" loss
$$\| z_\theta(x_{i+1}) - z_\theta(x_i) \|_2$$

**Supervised training**

$x_k \rightarrow \theta \rightarrow z_\theta(x_k) \rightarrow W \rightarrow$ softmax loss $L_C(x_k, y_k)$

$y_k$

**Training within a Siamese triplet neural network architecture**

## Datasets

unlabeled video → target task



NORB → NORB 25 Objects
KITTI Video → SUN 397 Scenes
Human Motion (HMDB) → PASCAL 10 Actions

| Task | Img/frame dims | #Classes | Recog. Task | #Train | #Test | Unsup. Input Type | #Pairs (1:3) | #Triplets (1:1) |
|---|---|---|---|---|---|---|---|---|
| NORB→NORB | 96×96×1 | 25 | object | 150 | 8100 | pose-reg. images | 50,000 | 75,000 |
| KITTI→SUN | 32×32×1 | 397 | scene | 2382 | 7940 | car-mounted video | 100,000 | 100,000 |
| HMDB→PASCAL-10 | 32×32×3 | 10 | action | 50 | 2000 | web video | 100,000 | 100,000 |

## Sequence completion



NORB — Query — Top-3 Results
KITTI — Query — Top-3 Results
HMDB — Query — Top-3 Results

**Problem:** Given two frames, find third in sequence, as: $z_3 = z_2 + (z_2 - z_1)$

**Score:** $\eta = E\left[\dfrac{rank}{\#candidates}\right] \times 100$

| Datasets | NORB | KITTI | HMDB |
|---|---|---|---|
| SFA-1 [1] | 0.95 | 31.04 | 2.70 |
| SFA-2 [2] | 0.91 | 8.39 | 2.27 |
| SSFA (Ours) | **0.53** | **7.79** | **1.78** |

**Qualitative and quantitative feature "steadiness" verification**

## Regularized category recognition from few samples



**25 objects** NORB → NORB
**397 scenes** KITTI → SUN
**10 actions** HMDB → PASCAL-10

- Clsnet (black)
- SFA-1 [1] (blue)
- SFA-2 [2] (green)
- Ours (SSFA) (red)

Slowness baselines:
[1] Mobahi, ICML'09
[2] Hadsell, CVPR'06

**Strong and consistent accuracy gains for *higher-order* temporal coherence vs. *slow* feature learning methods**

## Unsupervised vs. supervised pretraining and finetuning

Pretraining on unlabeled video vs. labeled CIFAR-100 images:



**PASCAL-10 Actions** — Accuracy vs. Extra supervision for SUP-FT ×10⁴
**SUN Scenes** — Accuracy vs. Extra supervision for SUP-FT ×10⁴

- SUP-FT
- SSFA (ours)
- SFA-2
- SFA-1

**Our unsupervised features can even surpass supervised pretraining with up to 50,000 additional class labels for an auxiliary task!**