# LEARNING CORRESPONDENCE FROM THE CYCLE-CONSISTENCY OF TIME

Xiaolong Wang    Carnegie Mellon University
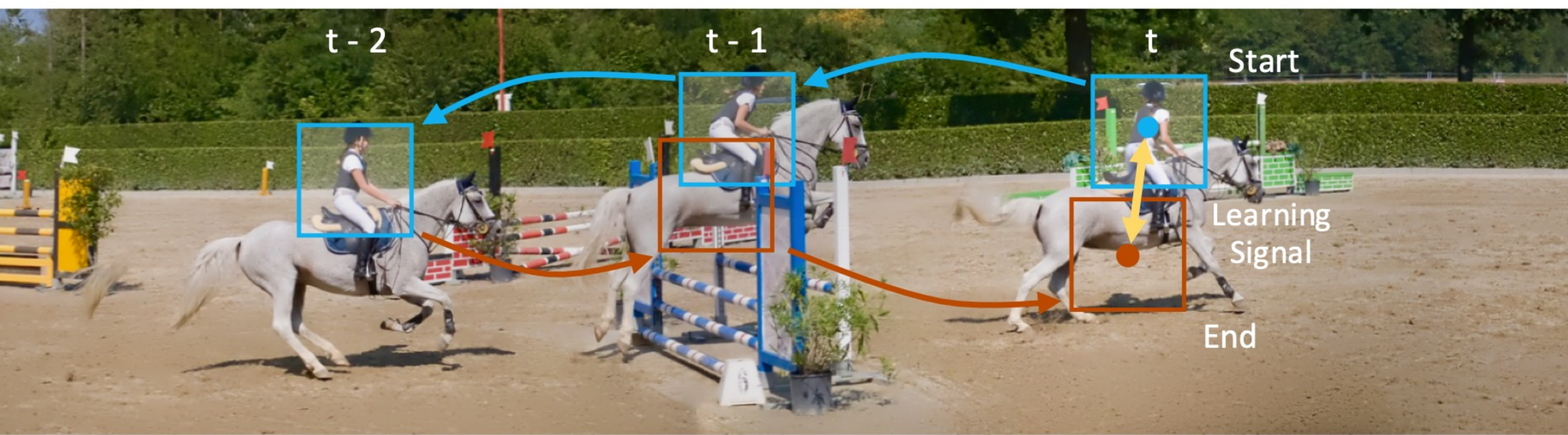
Allan Jabri*        UC Berkeley

Alexei A. Efros    UC Berkeley

# CORRESPONDENCE

- Pixel level: Optical Flow

  - Training data limitation: Synthetic datasets may not match real images

- Object level: Tracking

  - Training data limitation: human annotation of objects

# CYCLE CONSISTENCY

- Concept: define a powerful feature descriptor network $\phi$ and a weak tracking operator $\mathcal{T}$ that together track a patch through an image

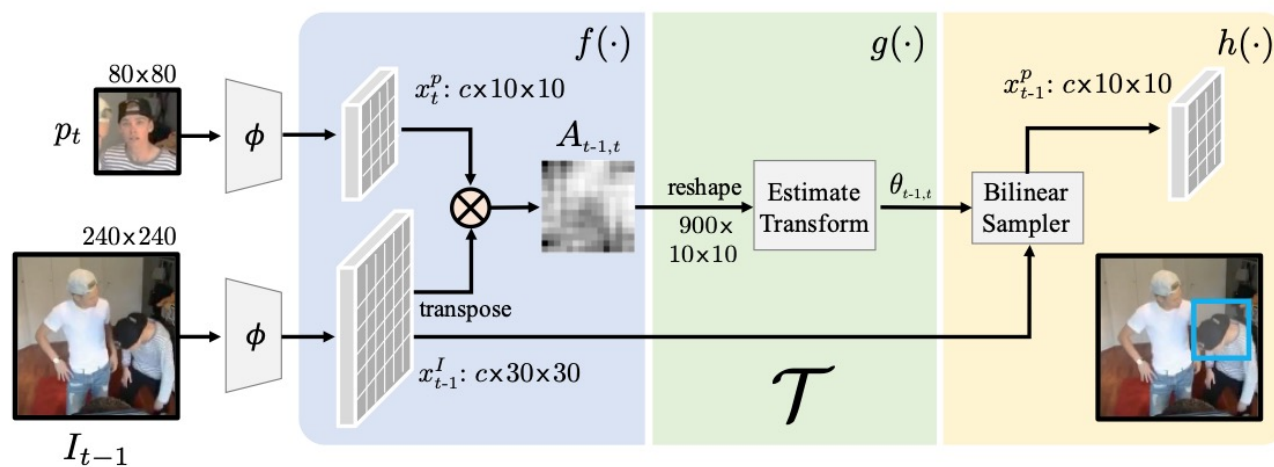- Loss: how "cycle consistent" is $\phi$ when combined with $\mathcal{T}$

# NOTATION:

- $I_{t-k:t}$ — A sequence of k+1 images from a video

- $p_t$ — A patch from image t

- $\phi$ — An encoder that produces a grid of feature vectors

- $x^I_{t-k:t}$ — $\phi(I_{t-k:t})$ — k+1 x c x 30 x 30

- $x^p_t$ — $\phi(p_t)$ — c x 10 x 10

- $\mathcal{T}$ — $x^I_s \times x^p_t \rightarrow x^p_s$

- $\mathcal{T}$ finds the patch in $x^I_s$ that is most similar to $x^p_t$

# TRAINING PROCESS

$$\mathcal{T}$$

- $A: 900 \times 100$

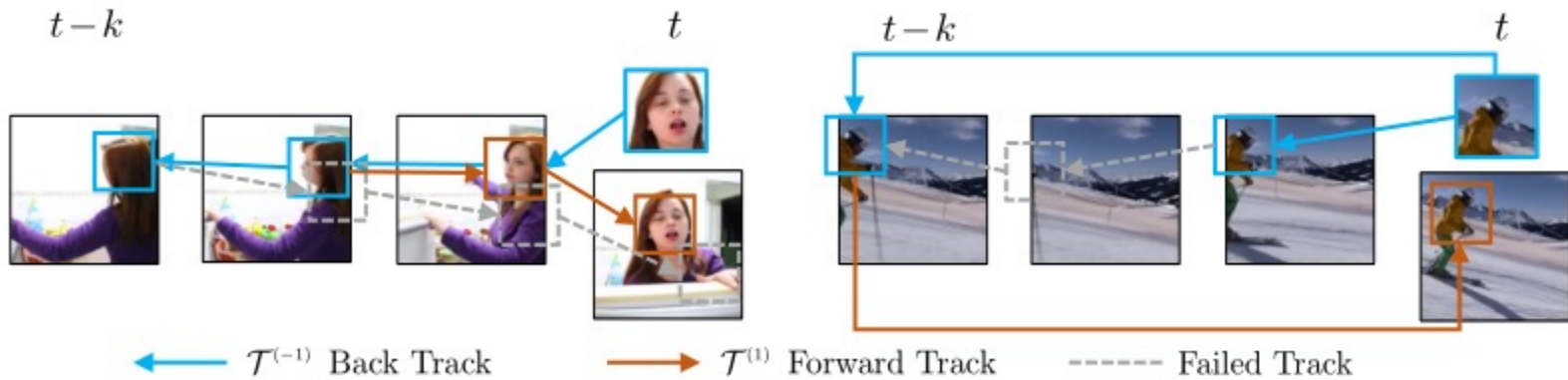- $A(i,j) = \dfrac{e^{x^I(j) \cdot x^p(i)}}{\sum_j e^{x^I(j) \cdot x^p(i)}}$



(b) Differentiable Tracking Operation $\mathcal{T}$

# ITERATE OPERATOR $\mathcal{T}$

$$\mathcal{T}^{(i)}(x_{t-i}^I, x^p) = \mathcal{T}(x_{t-1}^I, \mathcal{T}(x_{t-2}^I, ... \mathcal{T}(x_{t-i}^I, x^p)))$$

$$\mathcal{T}^{(-i)}(x_{t-1}^I, x^p) = \mathcal{T}(x_{t-i}^I, \mathcal{T}(x_{t-i+1}^I, ... \mathcal{T}(x_{t-1}^I, x^p)))$$
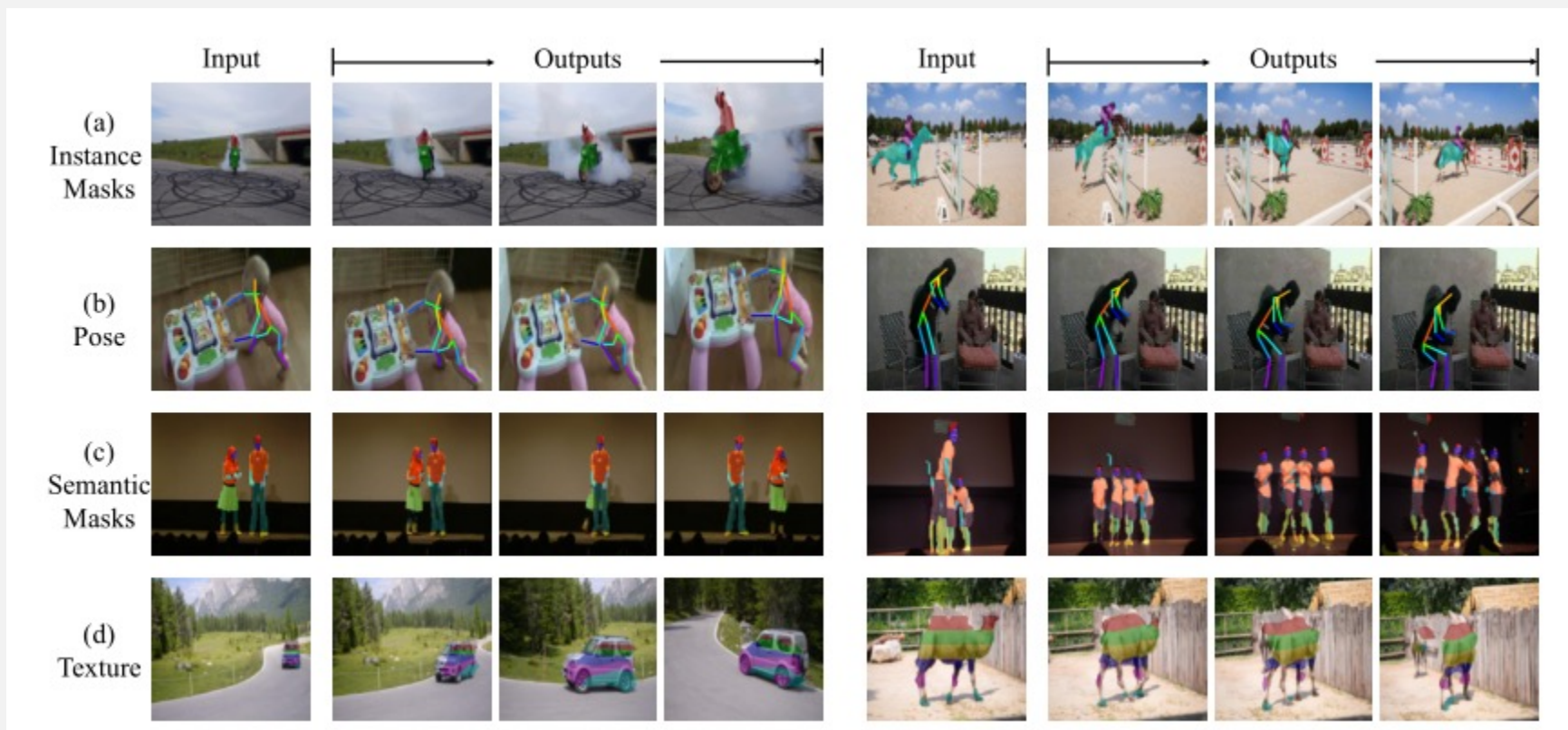
# FULL LOSS



$\mathcal{L}^3_{long}$        $\mathcal{L}^3_{skip}$

- $\mathcal{L} = \sum_i \mathcal{L}^i_{sim} + \lambda\mathcal{L}^i_{skip} + \lambda\mathcal{L}^i_{long}$
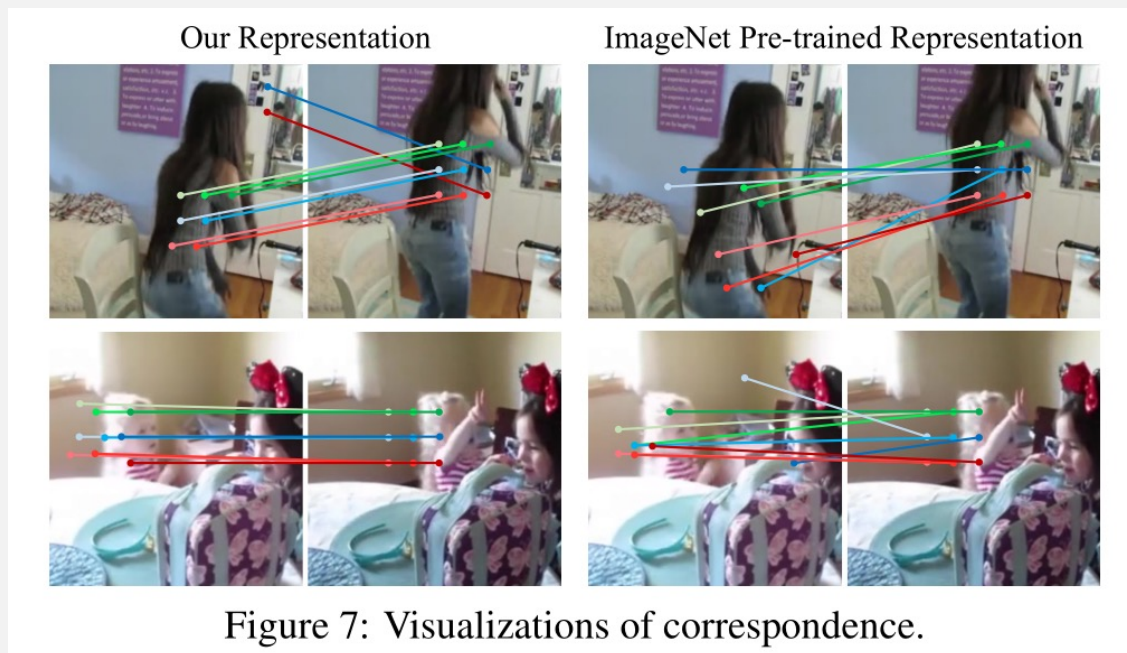
# INFERENCE PROCESS

- Drop $\mathcal{T}$ entirely! Just use features from $\phi$

- Propagate labels:

  - Now A is over whole image instead of patch

  - $A(j,i) = \dfrac{e^{x^I_{t-1}(j)\cdot x^I_t(i)}}{\sum_j e^{x^I_{t-1}(j)\cdot x^I_t(i)}}$

  - $y_i = \sum_j A_{t-1,t}(j,i)\, y_j$

  - Finally, label map upsampled

# VISUAL RESULTS

# VISUAL RESULTS



Figure 7: Visualizations of correspondence.

# NUMERICAL RESULTS

- DAVIS mask propagation

JHMDB pose propagation

| model | Supervised | $\mathcal{J}$(Mean) | $\mathcal{F}$(Mean) |
|---|---|---|---|
| Identity | | 22.1 | 23.6 |
| Random Weights (ResNet-50) | | 12.4 | 12.5 |
| Optical Flow (FlowNet2) [22] | | 26.7 | 25.2 |
| SIFT Flow [39] | | 33.0 | 35.0 |
| Transitive Inv. [74] | | 32.0 | 26.8 |
| DeepCluster [8] | | 37.5 | 33.2 |
| Video Colorization [69] | | 34.6 | 32.7 |
| Ours (ResNet-18) | | 40.1 | 38.3 |
| Ours (ResNet-50) | | **41.9** | **39.4** |
| ImageNet (ResNet-50) [18] | ✓ | 50.3 | 49.0 |
| Fully Supervised [81, 7] | ✓ | 55.1 | 62.1 |

| model | Supervised | PCK@.1 | PCK@.2 |
|---|---|---|---|
| Identity | | 43.1 | 64.5 |
| Optical Flow (FlowNet2) [22] | | 45.2 | 62.9 |
| SIFT Flow [39] | | 49.0 | 68.6 |
| Transitive Inv. [74] | | 43.9 | 67.0 |
| DeepCluster [8] | | 43.2 | 66.9 |
| Video Colorization [69] | | 45.2 | 69.6 |
| Ours (ResNet-18) | | 57.3 | 78.1 |
| Ours (ResNet-50) | | **57.7** | **78.5** |
| ImageNet (ResNet-50) [18] | ✓ | 58.4 | 78.4 |
| Fully Supervised [59] | ✓ | 68.7 | 92.1 |

Table 2: Evaluation on pose propagation on JHMDB [26]. We report the PCK in different thresholds.

# DISCUSSION QUESTIONS:

- Why do we use a neural network to interpret A(j, i) during training, but use it directly during test time?

  - Why do we learn rotation during training but not use it during test time?

- Do the authors make a convincing argument that their process is better than pretraining on ImageNet?