# Using Video to Learn about Visual Correspondence







Alexei (Alyosha) Efros UC Berkeley



target





Why video?

# 1. Richer signal







Jackson Pollock Number 21 (detail)

"It irritated him that the "dog" of 3:14 in the afternoon, seen in profile, should be indicated by the same noun as the dog of 3:15, seen frontally..."

-- from *Funes the Memorious* 



2. Correspondences



**Jorge Luis Borges** 

### Continuity crucial for visual development





Wood 2013, 2016, 2018

# 3. Ordering?

### "Time is what keeps everything from happening at once."



#### -- Ray Cummings (1919)

larverst of or an an inclusion of a state at



Canadity connections an arreful to ABOX CYNDIAL XUN DEMANDED INF. рюзния

второй нелекинай

Financers of parameters equidant test, which are built and ROMANN STREET CODESS. DESCRIPTION.

# Video in the "old days"



# Space-time XYT volumes



#### Clean and beautiful story

# Recognition with XYT volumes



"braided patterns" of cyclic motion [Niyogi & Adelson, 1994]



histograms of spatio-temporal gradients [Zelnik-Manor & Irani, 2001]



3D Convolutions [Tran et al, 2015]

# XYT segmentation [Shi & Malik, 98]







time->



# Video Textures

Arno Schödl Richard Szeliski David Salesin Irfan Essa

#### SIGGRAPH 2000

## Problem statement



#### video clip



#### video texture

# Text Synthesis

- [Shannon,'48] proposed a way to generate English-looking text using N-grams:
  - Assume a generalized Markov model
  - Use a large text to compute prob. distributions of each letter given N-1 previous letters
  - Starting from a seed repeatedly sample this
     Markov chain to generate new letters
  - Also works for whole words

#### WE NEED TO EAT CAKE

# Texture Synthesis: Markov chain on pixels



#### Synthesizing a pixel

Efros & Leung, ICCV 1999

non-parametric sampling



# Video Textures: Markov chain on frames

#### • How do we find good transitions?



# Finding good transitions

#### Compute L<sub>2</sub> distance $D_{i, j}$ between all frames



#### Similar frames make good transitions



# • Transition from i to j if successor of i is similar to j





# Transition costs

• Cost function:  $C_{i \rightarrow j} = D_{i+1, j}$ 



# Preserving dynamics



# Preserving dynamics



# Preserving dynamics

# • Cost for transition $i \rightarrow j$ • $C_{i \rightarrow j} = \sum_{k = -N}^{N-1} \sum_{k = -N}^$



•  $C_{i \rightarrow j} = \sum_{k=1}^{N-1} w_k D_{i+k+1, j+k}$ 

# Preserving dynamics – effect

Cost for transition  $i \rightarrow j$ 



•  $C_{i \rightarrow j} = \sum_{k=1}^{N-1} w_k D_{i+k+1, j+k}$ k = -N

### User-controlled video textures



Slow

#### User selects target frame range





fast

# Video sprite extraction



#### blue screen matting and velocity estimation





#### • Augmented transition cost:

# Animation $C_{i \rightarrow j} = \alpha C_{i \rightarrow j} + \beta \text{ angle } \langle c_{i \rightarrow j} \rangle$ Similarity term Control term

## Video sprite control

# vector to mouse pointer velocity vector

# Interactive fish





# Trouble with XYT volumes

- Time is not just "another dimension" – Very sparse sampling in t
- Implicit correspondence assumed  $-e.g.(x,y,t) \rightarrow (x,y,t+1)$  (fixed camera) or other simple models



# Explicit Correspondences

### **Optical Flow**



- + dense correspondences
  short-range (2 frames)
- local

### **Object Tacking**



- + longer-range & mid-level
- one object at a time
- not stable
  - "time to failure" metric



- Enter: learning
- detect in all other frames
- Now fancy off-the-shelf detectors
- But no explicit correspondence



"Strike a Pose" [Ramanan, Forsyth, Zisserman, 2005]

• Train simple person detector on one frame, then

### Data association

# • Temporal correspondence across detections - Step 1: detect objects (requires supervision) (but divorced from pixels)



- Step 2: find correspondences across detections



# Using Time as Supervision

### Temporal Self-supervised Feature Learning tasks

Inputs



#### Outputs



aligned vs. not-aligned	
3D Convolution	
1D Convolution	
1D Convolution	
	-

#### Predict Audio-visual Shifts [Owens & Efros, 2018]



#### Predict Arrow of Time [Misra et al, 2016; Wei et al, 2018]



#### Predict Color in Time [Vondrick et al, 2018]

# Using Tracking to Learn Features







Tracking → Similarity [Wang et al, 2015; Pathak et al, 2017]

Tracking → Similarity [Wang et al, 2015; Pathak et al, 2017]

# Limited by Off-the-shelf Trackers





CNN

# Using Tracking to Learn Features









#### Similarity requires tracking





#### Tracking requires similarity



# Let's jointly learn both!



# Learning Correspondence from the Cycle-consistency of Time



#### Xiaolong Wang





Allan Jabri CVPR 2019

Alexei Efros

# Learning to Track

#### F: a deep tracker


# Supervision: Cycle-Consistency in Time

### Track backwards



### Track forwards, back to the future

# Supervision: Cycle-Consistency in Time



### Backpropagation through time along the cycle

# Visualization of Training



### Limitations



- One patch at a time
- Winner-takes-all tracking
- Complex tracker (Spatial Transformer)
- Does not improve with longer cycles (6+ frames)
- Does not improve with more training data
- Fresh new work addresses these...



# Space-Time Correspondence as a Contrastive Random Walk



(on ArXiv next week!)

Allan Jabri, Andrew Owens, Alexei A. Efros

# Aim:

# Learn a similarity representation for correspondence from unlabeled video

**Focus:** Simplicity and Scalability

# Supervised correspondences



# Supervised correspondences























## Latent correspondences













# Palindrome sequence

































# <u>Self-supervised</u> correspondences









t



# Contrastive Random Walk







#### Video as a Graph



 $I_t$ 

Pixels

Nodes

 $\mathbf{q}_t$ 

#### Video as a Graph



**q**<sub>t</sub>

 $\mathbf{q}_{t+1}$ 

$$A_{ij} = \frac{e^{d_{\phi}(q_t^i, q_{t+1}^j)/\tau}}{\sum_l e^{d_{\phi}(q_t^i, q_{t+1}^l)/\tau}}$$
$$= P(X_{t+1} = j | X_t = i)$$

where  $d_{\phi}(x, y) = \phi(x)^{\mathsf{T}} \phi(y)$ 

 $X_t$  is the position of walker at time t







#### Correspondence as a Random Walk



### Learn embedding $\phi$ = Fitting transition probabilities



### One Step Contrastive Learning



#### Maximize

### $P(X_{t+1} = target | X_t = query)$

### **One Step Contrastive Learning**



#### Maximize

 $P(X_{t+1} = pos | X_t = query) = \frac{e^{\phi(query)^{\mathsf{T}}\phi(pos)}}{\sum_l e^{\phi(query)^{\mathsf{T}}\phi(neg_l)}}$ 



 $l_2$ -normalized Embedding

**Contrastive Learning** 

#### with cross entropy loss

Dosovitsky et al (2014)

Isola et al (2015)

Tian et al (2019)

Wu et al (2018)

van den Oord et al (2019)









target

#### *k*-step Transition Matrix

$$\overline{A}_{t}^{t+k} = \prod_{i=0}^{k-1} A_{t+i}^{t+i+1} \qquad \text{Sums ove}$$
intermediation
$$= P(X_{t+k}|X_{t})$$





target

#### Easy case:

#### Correspondence is obvious



target

#### Harder case:

Ambiguity

#### Considers multiple paths



#### Supervised $\rightarrow$ Self-Supervised



### Train on Palindromes

$$\mathcal{L}_{cyc}^k = \mathcal{L}_{CE}(\bar{A}_t^{t+k}\bar{A}_{t+k}^t,$$

Contrastive learning with latent correspondences





### Simple



#### **Algorithm 1** Pseudocode in a PyTorch-like style.

```
# load a minibatch x with B sequences
 # B x C x T x H x W -> B x C x T x P x h x w
 v = 12\_norm(resnet(x)) \# Embed patches (B x C x T x P)
 # Transitions from t to t+1 (B \times T \times P \times P)
 A = einsum("bcti,bctj->bcij", v[:,:,:-1], v[:,:,1:])
 # Transition similarities for palindrome graph
AA = cat((A, A[:,::-1].transpose(-1,-2), 1)
  # Perform random walk
for t in range(2*T):
  At = bmm(softmax(dropedge(AA[:,t]),dim=-1), At)
  # Target is the original node
  loss = cross_ent_loss(At, labels=[range(P)]*B)
```

## Evaluation

#### Semantic Part Propagation



#### Pose Propagation 15 Keypoints

Use  $\phi$  for

Label Propagation



#### **Object Propagation**



#### 20 Parts



#### **VIP Benchmark**



#### JHMDB Benchmark



1-4 Objects

#### **DAVIS Benchmark**



#### Label Propagation



#### Using more context, i.e. last *m* frames

Wang et al. (2019) Li et al. (2019) Lai et al. (2019



## Semantic Part Propagation

on Video Instance Parsing (VIP) dataset











# Pose Tracking

### on JHMDB
















#### "Common Fate" by Edge Dropout



t

t + k

# Higher-level correspondence

#### by modeling correlated paths

Laws of organization in perceptual forms. Wehrtheimer (1938)



# Edge Dropout

# Video Object Propagation on DAVIS dataset

#### Ours



#### Ours



# Ours



### Still Frame Comparison



UVC [Li et al., 2019]



#### Ours



#### Quantitative Comparison: Self-Supervised State-of-the-Art



#### Quantitative Comparison: Semantic Part Propagation



#### Quantitative Comparison: Pose Propagation



Sun et al, ICML 2020

# (Self-supervised Test-time Training)

### Pre-Train $\rightarrow$ Self-sup Train on X $\rightarrow$ Test on X

So far...

## Pre-Train → Test on Example X

(No Adaptation)

# Test-time Training (DAVIS)



# Path Length at Training



# Learning Curve



#### DAVIS J&F Mean

# Conclusions

- We proposed a simple and effective formulation for learning correspondence in a scalable way from unlabeled video.
- The method outperforms self-supervised state of the art methods despite being more simple.