# Space-Time Correspondence
# as a Contrastive Random Walk

**Allan A. Jabri**
UC Berkeley

**Andrew Owens**
University of Michigan
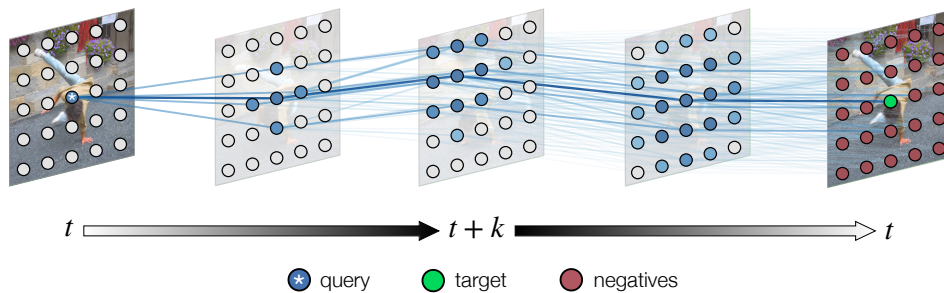
**Alexei A. Efros**
UC Berkeley

Figure 1: We represent video as a graph, where nodes are image patches, and edges are affinities (in some feature space) between nodes of neighboring frames. Our aim is to learn features such that temporal correspondences are represented by strong edges. We learn to find paths through the graph by performing a random walk between query and target nodes. A contrastive loss encourages paths that reach the target, implicitly supervising latent correspondence along the path. Learning proceeds *without labels* by training on a *palindrome* sequence, walking from frame $t$ to $t + k$, then back to $t$, using the initial node itself as the target. Please see our webpage for videos.

## Abstract

This paper proposes a simple self-supervised approach for learning a representation for visual correspondence from raw video. We cast correspondence as prediction of links in a space-time graph constructed from video. In this graph, the nodes are patches sampled from each frame, and nodes adjacent in time can share a directed edge. We learn a representation in which pairwise similarity defines transition probability of a random walk, so that long-range correspondence is computed as a walk along the graph. We optimize the representation to place high probability along paths of similarity. Targets for learning are formed without supervision, by cycle-consistency: the objective is to maximize the likelihood of returning to the initial node when walking along a graph constructed from a palindrome of frames. Thus, a single path-level constraint implicitly supervises chains of intermediate comparisons. When used as a similarity metric without adaptation, the learned representation outperforms the self-supervised state-of-the-art on label propagation tasks involving objects, semantic parts, and pose. Moreover, we demonstrate that a technique we call edge dropout, as well as self-supervised adaptation at test-time, further improve transfer for object-centric correspondence.

## 1   Introduction

There has been a flurry of advances in self-supervised representation learning from still images, yet this has not translated into commensurate advances in learning from video. Video is often treated as a simple extension of an image into time, modeled as a spatio-temporal $XYT$ volume [72, 118, 14]. Yet, treating time as yet another dimension is limiting [26]. One practical issue is the sampling rate

mismatch between $X$ and $Y$ vs. $T$. But a more fundamental problem is that a physical point depicted at position $(x, y)$ in frame $t$ might not have any relation to what we find at that same $(x, y)$ in frame $t + k$, as the object or the camera will have moved in arbitrary (albeit smooth) ways. This is why the notion of *temporal correspondence* — "what went where" [113] — is so fundamental for learning about objects in dynamic scenes, and how they inevitably change.

Recent approaches for self-supervised representation learning, such as those based on pairwise similarity learning [17, 22, 100, 88, 114, 45, 96, 41, 16], are highly effective when pairs of matching views are assumed to be known, e.g. constructed via data augmentation. Temporal correspondences, however, are *latent*, leading to a chicken-and-egg problem: we need correspondences to train our model, yet we rely on our model to find these correspondences. An emerging line of work aims to address this problem by bootstrapping an initially random representation to infer which correspondences should be learned in a self-supervised manner e.g. via cycle-consistency of time [110, 106, 58]. While this is a promising direction, current methods rely on complex and greedy tracking that may lead to local optima, especially when applied recurrently in time.

In this paper, we learn to associate features across space and time by formulating correspondence as pathfinding on a space-time graph. The graph is constructed from a video, where nodes are image patches and only nodes in neighboring frames share an edge. The strength of the edge is determined by similarity under a learned representation, whose aim is to place weight along paths linking visually corresponding patches (see Figure 1). Learning the representation amounts to fitting the transition probabilities of a walker stepping through time along the graph, reminiscent of the classic work of Meila and Shi [68] on learning graph affinities with a local random walk. This learning problem requires supervision — namely, the target that the walker should reach. In lieu of ground truth labels, we use the idea of cycle-consistency [122, 110, 106], by turning training videos into *palindromes*, e.g. sequences where the first half is repeated backwards. This provides every walker with a target — returning to its starting point. Under this formulation, we can view each step of the walk as a contrastive learning problem [17], where the walker's target provides supervision for entire chains of intermediate comparisons.

The central benefit of the proposed model is efficient consideration and supervision of many paths through the graph by computing the expected outcome of a random walk. This lets us obtain a learning signal from all views (patches) in the video simultaneously, and handling ambiguity in order to learn from harder examples encountered during training. Despite its simplicity, the method learns a representation that is effective for a variety of correspondence tasks. When used as a similarity metric without any adaptation, the representation outperforms state-of-the-art self-supervised methods on video object segmentation, pose keypoint propagation, and semantic part propagation. The model scales and improves in performance as the length of walks used for training increases. We also show several extensions of the model that further improve the quality of object segmentation, including an edge dropout [92] technique that encourages the model to group "common-fate" [112] nodes together, as well as test-time adaptation.

## 2 Contrastive Random Walks on Video

We represent each video as a directed graph where nodes are patches, and weighted edges connect nodes in neighboring frames. Let $\mathbf{I}$ be a set of frames of a video and $\mathbf{q}_t$ be the set of $N$ nodes extracted from frame $\mathbf{I}_t$, e.g. by sampling overlapping patches in a grid. An encoder $\phi$ maps nodes to $l_2$-normalized $d$-dimensional vectors, which we use to compute a pairwise similarity function $d_\phi(q_1, q_2) = \langle \phi(q_1), \phi(q_2) \rangle$ and an embedding matrix for $\mathbf{q}_t$ denoted $Q_t \in \mathbb{R}^{N \times d}$. We convert pairwise similarities into non-negative affinities by applying a softmax (with temperature $\tau$) over edges departing from each node. For timesteps $t$ and $t + 1$, the stochastic matrix of affinities is

$$A_t^{t+1}(i, j) = \texttt{softmax}(Q_t Q_{t+1}^\top)_{ij} = \frac{\exp(d_\phi(\mathbf{q}_t^i, \mathbf{q}_{t+1}^j)/\tau)}{\sum_{l=1}^N \exp(d_\phi(\mathbf{q}_t^i, \mathbf{q}_{t+1}^l)/\tau)}, \qquad (1)$$

where the $\texttt{softmax}$ is row-wise. Note that this describes only the *local* affinity between the patches of two video frames, $\mathbf{q}_t$ and $\mathbf{q}_{t+1}$. The affinity matrix for the entire graph, which relates all nodes in the video as a Markov chain, is block-sparse and composed of local affinity matrices.
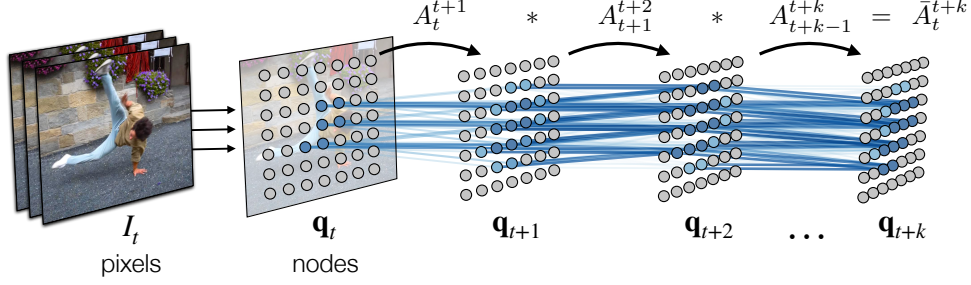
2

Figure 2: **Correspondence as a Random Walk**. We build a space-time graph by extracting nodes from each frame and allowing directed edges between nodes in neighboring frames. The transition probabilities of a random walk along this graph are determined by pairwise similarity in a learned representation.
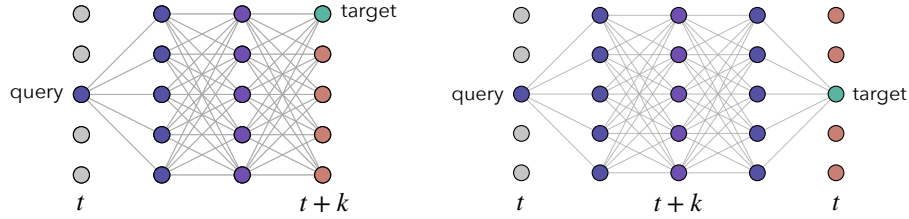


Figure 3: **Learning to Walk on Video.** (a) Specifying a target multiple steps in the future provides implicit supervision for *latent* correspondences along each path *(left)*. (b) We can construct targets for free by choosing palindromes as sequences for learning *(right)*.

Given the spatio-temporal connectivity of the graph, a step of a random walker on this graph can be viewed as performing tracking by *contrasting* similarity of neighboring nodes (using encoder $\phi$). Let $X_t$ be the state of the walker at time $t$, with transition probabilities $A_t^{t+1}(i,j) = P(X_{t+1} = j|X_t = i)$, where $P(X_t = i)$ is the probability of being at node $i$ at time $t$. With this view, we can formulate long-range correspondence as walking multiple steps along the graph (Figure 2):

$$\bar{A}_t^{t+k} = \prod_{i=0}^{k-1} A_{t+i}^{t+i+1} = P(X_{t+k}|X_t). \tag{2}$$

**Guiding the walk.**   Our aim is to train the embedding to encourage the random walker to follow paths of corresponding patches as it steps through time. While ultimately we will train without labels, for motivation suppose that we did have ground-truth correspondence between nodes in two frames of a video, $t$ and $t + k$ (Figure 3a). We can use these labels to fit the embedding by maximizing the likelihood that a walker beginning at a *query* node at $t$ ends at the *target* node at time $t + k$:

$$\mathcal{L}_{sup} = \mathcal{L}_{CE}(\bar{A}_t^{t+k}, Y_t^{t+k}) = -\sum_{i=1}^{N} \log P(X_{t+k} = Y_t^{t+k}(i)|X_t = i), \tag{3}$$

where $\mathcal{L}_{CE}$ is cross entropy loss and $Y_t^{t+k}$ are correspondence labels for matching time $t$ to $t + k$. Given the way transition probabilities are computed, the walk can be viewed as a chain of contrastive learning problems. Providing supervision at *every* step amounts to maximizing similarity between query and target nodes adjacent in time, while minimizing similarity to all other neighbors.

The more interesting case is supervision of longer-range correspondence, i.e. $k > 1$. In this case, the labels of $t$ and $t + k$ provide *implicit* supervision for intermediate frames $t + 1, ..., t + k - 1$, assuming that latent correspondences exist to link $t$ and $t + k$. Recall that in computing $P(X_{t+k}|X_t)$, we marginalize over all intermediate paths that link nodes in $t$ and $t + k$. By minimizing $\mathcal{L}_{sup}$, we shift affinity to paths that link the query and target. In easier cases (e.g. smooth videos), the paths that the walker takes from each node will not overlap, and these paths will simply be reinforced. In more ambiguous cases – e.g. deformation, multi-modality, or one-to-many matches – transition probability may be split across latent correspondences, such that we consider distribution over paths with higher entropy. The embedding should capture similarity between nodes in a manner that hedges probability over paths to overcome ambiguity, while avoiding transitions to nodes that lead the walker astray.

3

## 2.1 Self-Supervision

How to obtain query-target pairs that are known to correspond, without human supervision? We can consider training on graphs in which correspondence between the first and last frames are known, by construction. One such class of sequences are *palindromes*, i.e. sequences that are identical when reversed, for which targets are known since the first and last frames are identical. Given a sequence of frames $(I_t, ..., I_{t+k})$, we form training examples by simply concatenating the sequence with a temporally reversed version of itself: $(I_t, ...I_{t+k}, ...I_t)$. Treating each query node's position as its own target (Figure 3b), we obtain the following cycle-consistency objective:

$$\mathcal{L}_{cyc}^k = \mathcal{L}_{CE}(\bar{A}_t^{t+k}\bar{A}_{t+k}^t, I) = -\sum_{i=1}^{N} \log P(X_{t+2k} = i | X_t = i) \tag{4}$$

By leveraging structure in the graph, we can generate supervision for chains of contrastive learning problems that can be made arbitrarily long. As the model computes a soft attention distribution at every time step, we can backpropagate error across – and thus learn from – the many alternate paths of similarity that link query and target nodes.

**Contrastive learning with latent views.** To better understand the model, we can interpret it as contrastive learning with latent views. The popular InfoNCE formulation [75] draws the representation of two views of the same example closer by minimizing the loss $\mathcal{L}_{CE}(U_1^2, I)$, where $U_1^2 \in \mathbb{R}^{n \times n}$ is the normalized affinity matrix between the vectors of the first and second views of $n$ examples, as in Equation 1. Suppose, however, that we do not know *which* views should be matched with one another, merely that there should be a soft one-to-one alignment between them. We can formulate this as contrastive learning guided by a 'one-hop' cycle-consistency constraint, composing $U_1^2$ with the "transposed" stochastic similarity matrix $U_2^1$, to produce the loss $\mathcal{L}_{CE}(U_1^2 U_2^1, I)$, akin to Equation 4.

This task becomes more challenging with multiple hops, as avoiding spurious features that lead to undesirable diffusion of similarity across the graph becomes more important. While there are other ways of learning to align sets of features – e.g. by assuming soft bijection [18, 31, 115, 86] – it is unclear how they should extend to the multi-hop setting, where such heuristics may not always be desirable at each intermediate step. The proposed objective avoids the need to explicitly infer intermediate latent views, instead imposing a sequence-level constraint based on long-range correspondence known by construction.

## 2.2 Edge Dropout

One might further consider correspondence on the level of broader *segments*, where points within a segment have strong affinity to all other points in the segment. This inspires a trivial extension of the method – randomly dropping edges from the graph, thereby forcing the walker to consider alternative paths. We apply dropout [92] (with rate $\delta$) to the transition matrix $A$ to obtain $\tilde{A} = \texttt{dropout}(A, \delta)$, and then re-normalize. The resulting transition matrix $B$ and noisy cycle loss are:

$$B_{ij} = \frac{\tilde{A}_{ij}}{\sum_l \tilde{A}_{il}} \qquad \mathcal{L}_{c\tilde{y}c}^k = \mathcal{L}_{CE}(B_t^{t+k}B_{t+k}^t, I).$$

Edge dropout affects the task by randomly obstructing paths, thus encouraging hedging of mass to paths correlated with the ideal path – i.e. paths of *common fate* [112] – similar to the effect in spectral-based segmentation [90, 68]. In practice, we apply edge dropout before normalizing affinities, by setting values to a negative constant. We will see in Section 3.2 that edge dropout improves object-centric correspondence.

## 2.3 Implementation

We now describe how we construct the graph and parameterize the node embedding $\phi$. Algorithm 1 provides complete pseudocode for the method.

**Pixels to Nodes.** At training time, we follow [44], where patches of size $64 \times 64$ are sampled on a $7 \times 7$ grid from a $256 \times 256$ image (i.e. 49 nodes per frame). Patches are spatially jittered to prevent matching based on borders. At test time, we found that we could reuse the convolutional feature map between patches instead of processing the patches independently [60], making the features computable with only a single feed-forward pass of our network.[1]

---

[1]Using a single convolutional feature map for training was susceptible to shortcut solutions; see Appendix C.

**Object Propagation** 1-4 Objects   **Pose Propagation** 15 Keypoints

**Semantic Part Propagation** 20 Parts

Figure 4: Qualitative results for label propagation under our model for object, pose, and semantic part propagation tasks. The first frame is indicate with a blue outline. **Please see our webpage for video results**, as well as a qualitative comparison with other methods.

**Encoder** $\phi$. We create an embedding for each image patch using a convolutional network, namely ResNet-18 [43]. We apply a linear projection and $l_2$ normalization after average pooling, obtaining a 128-dimensional vector. We reduce the stride of last two residual blocks (`res3` and `res4`) to be 1. Please see Appendix G for details.

**Shorter paths.** During training, we consider paths of multiple lengths. For a sequence of length $T$, we optimize all *sub*-cycles: $\mathcal{L}_{train} = \sum_{i=1}^{T} \mathcal{L}_{cyc}^i$. This loss encourages the sequence of nodes visited in the walk to be a palindrome, i.e. on a walk of length $N$, the node visited at step $t$ should be the same node as $N - t$. It induces a curriculum, as short walks are easier to learn than long ones. This can be computed efficiently, since the losses share affinity matrices.

**Training.** We train $\phi$ using the (unlabeled) videos from Kinetics400 [14], with Algorithm 1.

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

```
for x in loader: # x: batch with B sequences
  # Split image into patches
  # B x C x T x H x W -> B x C x T x N x h x w
  x = unfold(x, (patch_size, patch_size))
  x = spatial_jitter(x)
  # Embed patches (B x C x T x N)
  v = l2_norm(resnet(x))

  # Transitions from t to t+1 (B x T-1 x N x N)
  A = einsum("bcti,bctj->btij",
             v[:,:,:-1], v[:,:,1:]) / temperature

  # Transition energies for palindrome graph
  AA = cat((A, A[:,::-1].transpose(-1,-2), 1)
  AA[rand(AA) < drop_rate] = -1e10 # Edge dropout
  At = eye(P)                      # Init. position

  # Compute walks
  for t in range(2*T-2):
      At = bmm(softmax(AA[:,t]), dim=-1), At)

  # Target is the original node
  loss = At[[range(P)]*B]].log()
```

bmm: batch matrix multiplication; `eye`: identity matrix; `cat`: concatenation.; `rand`: random tensor drawn from $(0, 1)$.

We used the Adam optimizer [50] for two million updates with a learning rate of $1 \times 10^{-4}$. We use a temperature of $\tau = 0.07$ in Equation 1, following [114] and resize frames to $256 \times 256$ (before extracting nodes, as above). Except when indicated otherwise, we report results with edge dropout rate 0.1 and a videos of length 10. Please find more details in Appendix E.

## 3 Experiments

We evaluate the learned representation on video label propagation tasks involving objects, keypoints, and semantic parts, by using it as a similarity metric. We also study the effects of edge dropout, training sequence length, and self-supervised adaptation at test-time. In addition to comparison with the state-of-the-art, we consider a baseline of label propagation with strong pre-trained features. Please find additional details, comparisons, ablations, and qualitative results in the Appendices.

### 3.1 Transferring the Learned Representation

We transfer the trained representation to label propagation tasks involving objects, semantic parts, and human pose. To isolate the effect of the representation, we use a simple inference algorithm based on $k$-nearest neighbors. Qualitative results are shown in Figure 4.

5

| Method | Resolution | Train Data | $\mathcal{J}\&\mathcal{F}_{\mathrm{m}}$ | $\mathcal{J}_{\mathrm{m}}$ | $\mathcal{J}_{\mathrm{r}}$ | $\mathcal{F}_{\mathrm{m}}$ | $\mathcal{F}_{\mathrm{r}}$ |
|---|---|---|---|---|---|---|---|
| VINCE [32] | 1× | Kinetics | 60.4 | 57.9 | 66.2 | 62.8 | 71.5 |
| CorrFlow* [57] | 2× | OxUvA | 50.3 | 48.4 | 53.2 | 52.2 | 56.0 |
| MAST* [56] | 2× | OxUvA | 63.7 | 61.2 | 73.2 | 66.3 | 78.3 |
| MAST* [56] | 2× | YT-VOS | 65.5 | 63.3 | 73.2 | 67.6 | 77.7 |
| TimeCycle [110] | 1× | VLOG | 48.7 | 46.4 | 50.0 | 50.0 | 48.0 |
| UVC+track* [58] | 1× | Kinetics | 59.5 | 57.7 | 68.3 | 61.3 | 69.8 |
| UVC [58] | 1× | Kinetics | 60.9 | 59.3 | 68.8 | 62.7 | 70.9 |
| **Ours**  w/ dropout | 1× | Kinetics | **67.6** | **64.8** | **76.1** | **70.2** | **82.1** |
| w/ dropout & adaptation | 1× | Kinetics | 68.3 | 65.5 | 78.6 | 71.0 | 82.9 |

Table 1: **Video object segmentation results on DAVIS 2017 val set** Comparison of our method (2 variants), with previous self-supervised approaches and strong pretrained feature baselines. *Resolution* indicates if the approach uses a high-resolution (2x) feature map. *Train Data* indicates which dataset was used for pre-training. $\mathcal{F}$ is a boundary alignment metric, while $\mathcal{J}$ measures region similarity as IOU between masks. * indicates that our label propagation algorithm is not used.

**Label propagation.**    All evaluation tasks considered are cast as video label propagation, where the task is to predict labels for each pixel in *target* frames of a video given only ground-truth for the first frame (i.e. the *source*). We use the representation as a similarity function for prediction by $k$-nearest neighbors, which is natural under our model and follows prior work for fair comparison [110, 58].

Say we are given source nodes $\mathbf{q}_s$ with labels $L_s \in \mathbb{R}^{N \times C}$, and target nodes $\mathbf{q}_t$. Let $K_t^s$ be the matrix of transitions between $\mathbf{q}_t$ and $\mathbf{q}_s$ (Equation 1), with the special property that only the top$-k$ transitions are considered per target node. Labels $L_t$ are propagated as $L_t = K_t^s L_s$, where each row corresponds to the soft distribution over labels for a node, predicted by $k$-nearest neighbor in $d_\phi$.

To provide temporal context, as done in prior work [110, 57, 58], we use a queue of the last $m$ frames. We also restrict the set of source nodes considered to a spatial neighborhood of the query node for efficiency (i.e. *local* attention). The source set includes nodes of the first labeled frame, as well as the nodes in previous $m$ frames, whose predicted labels are used for auto-regressive propagation. The softmax computed for $K_t^s$ is applied over all source nodes. See Appendix F for further discussion and hyper-parameters.

**Baselines.**    All baselines use ResNet-18 [43] as the backbone, modified to increase spatial resolution of the feature map by reducing the stride of the last two residual blocks to be 1. For consistency across methods, we use the output of the penultimate residual block as node embeddings at test-time.

Pre-trained visual features: We evaluate pretrained features from strong image- and video-based representation learning methods. For a strongly supervised approach, we consider a model trained for classification on **ImageNet** [20]. We also consider a strong self-supervised method, **MoCo** [41]. Finally, we compare with a video-based contrastive learning method, **VINCE** [32], which extends MoCo to videos (Kinetics) with views from data augmentation *and* neighbors in time.

Task-specific approaches: **Wang et al.** [110] uses cycle-consistency to train a spatial transformer network as a deterministic patch tracker. We also consider methods based on the **Colorization** approach of Vondrick et al. [105], including high-resolution methods: **CorrFlow** [57] and **MAST** [56]. CorrFlow combines cycle consistency with colorization. MAST uses a deterministic region localizer and memory bank for high-resolution colorization, and performs multi-stage training on [99]. Notably, both [57, 56] use feature maps that are significantly higher resolution than other approaches (2×) by removing max pooling from the network. Finally, **UVC** [58] jointly optimizes losses for colorization, grouping, pixel-wise cycle-consistency, and patch tracking with a deterministic patch localizer.

### 3.1.1   Video Object Segmentation

We evaluate our model on DAVIS 2017 [84], a popular benchmark for video object segmentation, for the task of semi-supervised multi-object (i.e. 2-4) segmentation. Following common practice, we evaluate on 480p resolution images. We apply our label propagation algorithm for all comparisons, except CorrFlow and MAST [57, 56], which require 4× more GPU memory. We report mean (m) and recall (r) of standard boundary alignment ($\mathcal{F}$) and region similarity ($\mathcal{J}$) metrics, detailed in  [79].

As shown in Table 1, our approach outperforms other self-supervised methods, without relying on machinery such as localization modules or multi-stage training. We also outperform [56] despite being

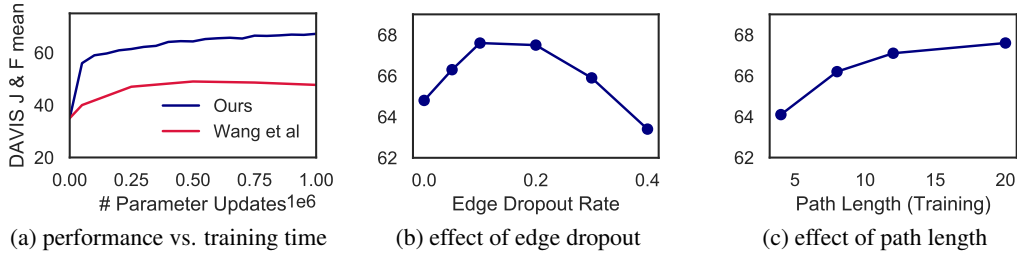| (a) performance vs. training time | (b) effect of edge dropout | (c) effect of path length |

Figure 5: **Variations of the Model.** (a) Downstream task performance as a function of training time. (b) Moderate edge dropout improves object-level correspondences. (c) Training on longer paths is beneficial. All evaluations are on the DAVIS segmentation task.

more simple at train and test time, and using a lower-resolution feature map. We found that when combined with a properly tuned label propagation algorithm, the more generic pretrained feature baselines fare better than more specialized temporal correspondence approaches. Our approach outperformed approaches such as MoCo [41] and VINCE [32], suggesting that it may not always be optimal to choose views for contrastive learning by random crop data augmentation of frames. Finally, our model compares favorably to many supervised approaches with architectures designed for dense tracking [79, 12, 107] (see Appendix B).

### 3.1.2 Pose Tracking

We consider pose tracking on the JHMDB benchmark, which involves tracking 15 keypoints. We follow the evaluation protocol of [58], using $320 \times 320$px images. As seen in Table 2, our model outperforms existing self-supervised approaches, including video colorization models that directly optimize for fine-grained matching with pixel-level objectives [58]. We attribute this success to the fact that our model sees sufficiently hard negative samples drawn from the same image at training time to learn features that discriminate beyond color. Note that our inference procedure is naive in that we propagate keypoints independently, without leveraging relational structure between them.

### 3.1.3 Video Part Segmentation

We consider the semantic part segmentation task of the Video Instance Parsing (VIP) benchmark [121], which involves propagating labels of 20 parts — such as arm, leg, hair, shirt, hand — requiring more precise correspondence than DAVIS. The sequences are longer and sampled at a lower frame rate. We follow the evaluation protocol of [58], using $560 \times 560$px images and $m = 1$. The model outperforms existing self-supervised methods, and when using more temporal context (i.e. $m = 4$), outperforms the baseline supervised approach of [121].

|  | Parts | Pose | |
| Method | mIoU | PCK@0.1 | PCK@0.2 |
| TimeCycle [110] | 28.9 | 57.3 | 78.1 |
| UVC [58] | 34.1 | 58.6 | 79.6 |
| Ours | 36.0 | 59.0 | 83.2 |
| Ours + context | **38.6** | **59.3** | **84.9** |
| ImageNet [43] | 31.9 | 53.8 | 74.6 |
| ATEN [121] | **37.9** | – | – |
| Yang et al. [116] | – | **68.7** | **92.1** |

Table 2: **Part and Pose Propagation tasks**, with the VIP and JHMDB benchmarks, respectively. For comparison, we show supervised methods below.

### 3.2 Variations of the Model

**Edge dropout.** We test the hypothesis (Figure 5b) that edge dropout should improve performance on the object segmentation task, by training our model with different edge dropout rates: {0, 0.05, 0.1, 0.2, 0.3, 0.4}. Moderate edge dropout yields a significant improvement on the DAVIS benchmark. Edge dropout simulates partial occlusion, forcing the network to consider reliable context.

**Path length.** We also asked how important it is for the model to see longer sequences during training, by using clips of length 2, 4, 6, or 10 (resulting in paths of length 4, 8, 12, or 20). Longer sequences yield harder tasks due to compounding error. We find that longer training sequences accelerated convergence as well as improved performance on the DAVIS task (Figure 5c). This is in contrast to prior work [110]; we attribute this success to considering multiple paths at training time via soft-attention, which allows for learning from longer sequences, despite ambiguity.

**Improvement with training** We found that the model's downstream performance on DAVIS improves as more data is seen during self-supervised training (Figure 5a). Compared to Wang et al [110], there is less indication of saturation of performance on the downstream task.

7

### 3.3 Self-supervised Adaptation at Test-time

A key benefit of not relying on labeled data is that training need not be limited to the training phase, but can continue during deployment [3, 71, 82, 94]. Our approach is especially suited for such adaptation, given the non-parametric inference procedure. We ask whether the model can be improved for object correspondence by fine-tuning the representation *at test time* on a novel video. Given an input video, we can perform a small number of iterations of gradient descent on the self-supervised loss (Algorithm 1) *prior* to label propagation. We argue it is most natural to consider an online setting, where the video is ingested as a stream and fine-tuning is performed continuously on the sliding window of $k$ frames around the current frame. Note that only the raw, unlabeled video is used for this adaptation; we do not use the provided label mask. As seen in Table 1, test-time training improves object propagation. Interestingly, we see most improvement in the recall of the region similarity metric $\mathcal{J}_{recall}$ (which measures how often more than 50% of the object is segmented). More experiment details can be found in Appendix E.

## 4 Related Work

**Temporal Correspondence.** Many early methods represented video as a spatio-temporal $XYT$ volume, where patterns, such as lines or statistics of spatio-temporal gradients, were computed for tasks like gait tracking [72] and action recognition [118]. Because the camera was usually static, this provided an implicit temporal correspondence via $(x, y)$ coordinates. For more complex videos, optical flow [63] was used to obtain short-range explicit correspondences between patches of neighboring frames. However, optical flow proved too noisy to provide long-range composite correspondences across many frames. Object tracking was meant to offer robust long-range correspondences for a given tracked object. But after many years of effort (see [27] for overview), that goal was largely abandoned as too difficult, giving rise to "tracking as repeated detection" paradigm [85], where trained object detectors are applied to each frame independently. In the case of multiple objects, the process of "data association" resolves detections into coherent object tracks. Data association is often cast as an optimization problem for finding paths through video that fulfill certain constraints, e.g. appearance, position overlap, etc. Approaches include dynamic programming, particle filtering, various graph-based combinatorial optimization, and more recently, graph neural networks [119, 87, 8, 83, 15, 117, 49, 48, 54, 11, 13, 52]. Our work can be seen as contrastive data association via soft-attention, as a means for learning representations directly from pixels.

**Graph Neural Networks and Attention.** Representing inputs as graphs has led to unified deep learning architectures. Graph neural networks – versatile and effective across domains [101, 38, 53, 102, 108, 6, 13] – can be seen as learned message passing algorithms that iteratively update node representations, where propagation of information is dynamic, contingent on local and global relations, and often implemented as soft-attention. Iterative routing of information encodes structure of the graph for downstream tasks. Our work uses cross-attention between nodes of adjacent frames to learn to propagate node identity through a graph, where the task – in essence, instance discrimination across space and time – is designed to induce representation learning.

**Graph Partitioning.** Graphs have been widely used in image and video segmentation as a data structure. Given a video, a graph is formed by connecting pixels in spatio-temporal neighborhoods, followed by spectral clustering [89, 90, 28] or MRF/GraphCuts [10]. Most relevant is the work of Meila and Shi [68], which poses Normalized Cuts as a Markov random walk, describing an algorithm for learning an affinity function for segmentation by fitting the transition probabilities to be uniform within segments and zero otherwise. More recently, there has been renewed interest in the problem of unsupervised grouping [35, 51, 34, 25, 52]. Many of these approaches can be viewed as end-to-end neural architectures for graph partitioning, where entities are partitions of images or video inferred by learned clustering algorithms or latent variable models implemented with neural networks. While these approaches explicitly group without supervision, they have mainly considered simpler data. Our work similarly aims to model groups in dynamic scenes, but does so implicitly so as to scale to real, large-scale video data. Incorporating more explicit entity estimation is an exciting direction.

**Graph Representation Learning.** Graph representation learning approaches solve for distributed representations of nodes and vertices given connectivity in the graph [39]. Most relevant are similarity learning approaches, which define neighborhoods of positives with fixed (i.e. $k$-hop neighborhood) or stochastic (i.e. random walk) heuristics [80, 36, 95, 38], while sampling negatives at random. Many of these approaches can thus be viewed as fitting shallow graph neural networks with tasks reminiscent of Mikolov et al. [69]. Backstrom et al. [5] learns to predict links by supervising a

random walk on social network data. While the above consider learning representations given a single graph, others have explored learning node embeddings given multiple graphs. A key challenge is inferring correspondence between graphs, which has been approached in prior work [115, 86] with efficient optimal transport algorithms [91, 18, 81]. We use graph matching as a means for representation learning, using cycle-consistency to supervise a chain of matches, without inferring correspondence between intermediate pairs of graphs. In a similar vein, cycle-consistency has also been shown to be a useful constraint for solving large-scale optimal transport problems [62].

**Self-supervised Visual Representation Learning.** Most work in self-supervised representation learning can be interpreted as data imputation: given an example, the task is to predict a part — or *view* — of its data given another view [7, 19, 17]. Earlier work leveraged unlabeled visual datasets by constructing *pretext* prediction tasks [21, 73, 120]. For video, temporal information makes for natural pretext tasks, including future prediction [33, 93, 67, 61, 65], arrow of time [70, 111], motion estimation [1, 47, 98, 59] or audio [77, 2, 76, 55]. The use of off-the-shelf tools to provide supervisory signal for learning visual similarity has also been explored [109, 29, 78]. Recent progress in self-supervised learning has focused on improving techniques for large-scale deep similarity learning, e.g. by combining the cross-entropy objective with negative sampling [37, 69]. Sets of corresponding views are constructed by composing combinations of augmentations of the same instance [22, 9, 114], with domain knowledge being crucial for picking the right data augmentations. Strong image-level visual representations can be learned by heuristically choosing views that are close in space [100, 45, 4, 41, 16], in time [88, 82, 40, 97, 32] or both [46, 96], even when relying on noisy negative samples. However, forcing random crops to be similar is not always desirable because they may not be in correspondence. In contrast, we implicitly determine which views to bring closer – a sort of automatic view selection.

**Self-supervised Correspondence and Cycle-consistency.** Our approach builds on recent work that uses cycle-consistency [122, 23] in time as supervisory signal for learning visual representations from video [110, 106]. The key idea in [110, 106] is to use self-supervised tracking as a pretext task: given a patch, first track forward in time, then backward, with the aim of ending up where it started, forming a cycle. These methods rely on trackers with hard attention, which limits them to sampling, and learning from, one path at a time. In contrast, our approach computes soft-attention at every time step, considering many paths to obtain a dense learning signal and overcome ambiguity. Li et al. [58] combines patch tracking with other losses including color label propagation [105], grouping, and cycle-consistency via an orthogonality constraint [30], considering pairs of frames at a time. Lai et al. [57, 56] refine architectural and training design decisions that yield impressive results on video object segmentation and tracking tasks. While colorization is a useful cue, the underlying assumption that corresponding pixels have the same color is often violated, e.g. due to lighting or deformation. In contrast, our loss is discriminative and permits association between regions that may have significant differences in their appearance.

## 5   Discussion

While data augmentation can be tuned to induce representation learning tasks involving invariance to color and local context, changes in other important factors of variation – such as physical transformations – are much harder to simulate. We presented a self-supervised approach for learning representations for space-time correspondence from unlabeled video data, based on learning to walk on a space-time graph. Under our formulation, a simple path-level constraint provides implicit supervision for a chain of contrastive learning problems. Our learning objective aims to leverage the natural data augmentation of dynamic scenes, *i.e.* how objects change and interact over time, and can be combined with other learning objectives. Moreover, it builds a connection between self-supervised representation learning and unsupervised grouping [68]. As such, we hope this work is a step toward learning to discover and describe the structure and dynamics of natural scenes from large-scale unlabeled video.

## 6   Broader Impact

Research presented in the paper has a potential to positively contribute to a number of practical applications where establishing temporal correspondence in video is critical, among them pedestrian safely in automotive settings, patient monitoring in hospitals and elderly care homes, video-based animal monitoring and 3D reconstruction, etc. However, there is also a potential for the technology to be used for nefarious purposes, mainly in the area of unauthorized surveillance, especially by

autocratic regimes. As partial mitigation, we commit to not entering into any contracts involving this technology with any government or quasi-governmental agencies of countries with an *EIU Democracy Index* [24] score of $4.0$ or below ("authoritarian regimes"), or authorizing them to use our software.

# References

[1] Pulkit Agrawal, Joao Carreira, and Jitendra Malik. Learning to see by moving. In *ICCV*, 2015.

[2] Relja Arandjelovic and Andrew Zisserman. Look, listen and learn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 609–617, 2017.

[3] Michal Irani Assaf Shocher, Nadav Cohen. "zero-shot" super-resolution using deep internal learning. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018.

[4] Philip Bachman, R Devon Hjelm, and William Buchwalter. Learning representations by maximizing mutual information across views. In *Advances in Neural Information Processing Systems*, pages 15535–15545, 2019.

[5] Lars Backstrom and Jure Leskovec. Supervised random walks: predicting and recommending links in social networks. In *Proceedings of the fourth ACM international conference on Web search and data mining*, pages 635–644, 2011.

[6] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018.

[7] Suzanna Becker and Geoffrey E Hinton. Self-organizing neural network that discovers surfaces in random-dot stereograms. *Nature*, 355(6356):161–163, 1992.

[8] Jerome Berclaz, Francois Fleuret, Engin Turetken, and Pascal Fua. Multiple object tracking using k-shortest paths optimization. *IEEE transactions on pattern analysis and machine intelligence*, 33(9):1806–1819, 2011.

[9] Piotr Bojanowski and Armand Joulin. Unsupervised learning by predicting noise. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 517–526. JMLR. org, 2017.

[10] Yuri Boykov and Gareth Funka-Lea. Graph cuts and efficient nd image segmentation. *International journal of computer vision*, 70(2):109–131, 2006.

[11] Guillem Brasó and Laura Leal-Taixé. Learning a neural solver for multiple object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6247–6257, 2020.

[12] Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and Luc Van Gool. One-shot video object segmentation. In *CVPR*, 2017.

[13] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. *arXiv preprint arXiv:2005.12872*, 2020.

[14] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *Computer Vision and Pattern Recognition (CVPR)*, 2017.

[15] Albert YC Chen and Jason J Corso. Temporally consistent multi-class video-object segmentation with the video graph-shifts algorithm. In *2011 IEEE Workshop on Applications of Computer Vision (WACV)*, pages 614–621. IEEE, 2011.

[16] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.

[17] Sumit Chopra, Raia Hadsell, and Yann LeCun. Learning a similarity metric discriminatively, with application to face verification. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)*, volume 1, pages 539–546. IEEE, 2005.

[18] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in neural information processing systems*, pages 2292–2300, 2013.

[19] Virginia R de Sa. Learning classification with unlabeled data. In *Advances in neural information processing systems*, pages 112–119, 1994.

[20] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *Computer Vision and Pattern Recognition (CVPR)*, 2009.

[21] Carl Doersch, Abhinav Gupta, and Alexei A. Efros. Unsupervised visual representation learning by context prediction. In *ICCV*, 2015.

[22] Alexey Dosovitskiy, Philipp Fischer, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE transactions on pattern analysis and machine intelligence*, 38(9):1734–1747, 2015.

[23] Debidatta Dwibedi, Yusuf Aytar, Jonathan Tompson, Pierre Sermanet, and Andrew Zisserman. Temporal cycle-consistency learning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1801–1810, 2019.

[24] EIU.com. Democracy index 2019 a year of democratic setbacks and popular protest. https://www.eiu.com/public/topical_report.aspx?campaignid=democracyindex2019, 2019.

[25] Martin Engelcke, Adam R Kosiorek, Oiwi Parker Jones, and Ingmar Posner. Genesis: Generative scene inference and sampling with object-centric latent representations. *arXiv preprint arXiv:1907.13052*, 2019.

[26] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 6202–6211, 2019.

[27] David A. Forsyth and Jean Ponce. *Computer Vision - A Modern Approach, Second Edition.* Pitman, 2012.

[28] Charless Fowlkes, Serge Belongie, Fan Chung, and Jitendra Malik. Spectral grouping using the nystrom method. *IEEE transactions on pattern analysis and machine intelligence*, 26(2):214–225, 2004.

[29] Ruohan Gao, Dinesh Jayaraman, and Kristen Grauman. Object-centric representation learning from unlabeled videos. In *Asian Conference on Computer Vision*, pages 248–263. Springer, 2016.

[30] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. A neural algorithm of artistic style. *arXiv preprint arXiv:1508.06576*, 2015.

[31] Aude Genevay, Gabriel Peyre, and Marco Cuturi. Learning generative models with sinkhorn divergences. In *International Conference on Artificial Intelligence and Statistics*, pages 1608–1617, 2018.

[32] Daniel Gordon, Kiana Ehsani, Dieter Fox, and Ali Farhadi. Watching the world go by: Representation learning from unlabeled videos, 2020.

[33] Ross Goroshin, Joan Bruna, Jonathan Tompson, David Eigen, and Yann LeCun. Unsupervised learning of spatiotemporally coherent metrics. *ICCV*, 2015.

[34] Klaus Greff, Raphaël Lopez Kaufman, Rishabh Kabra, Nick Watters, Chris Burgess, Daniel Zoran, Loic Matthey, Matthew Botvinick, and Alexander Lerchner. Multi-object representation learning with iterative variational inference. *arXiv preprint arXiv:1903.00450*, 2019.

[35] Klaus Greff, Sjoerd Van Steenkiste, and Jürgen Schmidhuber. Neural expectation maximization. In *Advances in Neural Information Processing Systems*, pages 6691–6701, 2017.

[36] Aditya Grover and Jure Leskovec. node2vec: Scalable feature learning for networks. In *Proceedings of the 22nd ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 855–864, 2016.

[37] Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation: A new estimation principle for unnormalized statistical models. In *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, pages 297–304, 2010.

[38] Will Hamilton, Zhitao Ying, and Jure Leskovec. Inductive representation learning on large graphs. In *Advances in neural information processing systems*, pages 1024–1034, 2017.

[39] William L Hamilton, Rex Ying, and Jure Leskovec. Representation learning on graphs: Methods and applications. *arXiv preprint arXiv:1709.05584*, 2017.

[40] Tengda Han, Weidi Xie, and Andrew Zisserman. Video representation learning by dense predictive coding. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019.

[41] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *CVPR*, 2020.

[42] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *International Conference on Computer Vision (ICCV)*, 2017.

[43] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR)*, 2016.

[44] Olivier J Hénaff, Aravind Srinivas, Jeffrey De Fauw, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019.

[45] R Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *International Conference on Learning Representations*, 2018.

[46] Phillip Isola, Daniel Zoran, Dilip Krishnan, and Edward H Adelson. Learning visual groups from co-occurrences in space and time. *arXiv preprint arXiv:1511.06811*, 2015.

[47] Dinesh Jayaraman and Kristen Grauman. Learning image representations tied to egomotion. In *ICCV*, 2015.

[48] Neel Joshi, Wolf Kienzle, Mike Toelle, Matt Uyttendaele, and Michael F Cohen. Real-time hyperlapse creation via optimal frame selection. *ACM Transactions on Graphics (TOG)*, 34(4):1–9, 2015.

[49] Armand Joulin, Kevin Tang, and Li Fei-Fei. Efficient image and video co-localization with frank-wolfe algorithm. In *European Conference on Computer Vision*, pages 253–268. Springer, 2014.

[50] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv*, 2014.

[51] Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural relational inference for interacting systems. *arXiv preprint arXiv:1802.04687*, 2018.

[52] Thomas Kipf, Elise van der Pol, and Max Welling. Contrastive learning of structured world models. In *International Conference on Learning Representations*, 2020.

[53] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.

[54] Shu Kong and Charless Fowlkes. Multigrid predictive filter flow for unsupervised learning on videos. *arXiv preprint arXiv:1904.01693*, 2019.

[55] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *Advances in Neural Information Processing Systems*, 2018.

[56] Zihang Lai, Erika Lu, and Weidi Xie. Mast: A memory-augmented self-supervised tracker. *arXiv preprint arXiv:2002.07793*, 2020.

[57] Zihang Lai and Weidi Xie. Self-supervised learning for video correspondence flow. *arXiv preprint arXiv:1905.00875*, 2019.

[58] Xueting Li, Sifei Liu, Shalini De Mello, Xiaolong Wang, Jan Kautz, and Ming-Hsuan Yang. Joint-task self-supervised learning for temporal correspondence. In *Advances in Neural Information Processing Systems*, pages 317–327, 2019.

[59] Yin Li, Manohar Paluri, James M. Rehg, and Piotr Dollár. Unsupervised learning of edges. In *CVPR*, 2016.

[60] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.

[61] William Lotter, Gabriel Kreiman, and David Cox. Deep predictive coding networks for video prediction and unsupervised learning. *arXiv preprint arXiv:1605.08104*, 2016.

[62] Guansong Lu, Zhiming Zhou, Jian Shen, Cheng Chen, Weinan Zhang, and Yong Yu. Large-scale optimal transport via adversarial training with cycle-consistency. *arXiv preprint arXiv:2003.06635*, 2020.

[63] Bruce D Lucas and Takeo Kanade. An iterative image registration technique with an application to stereo vision. *IJCAI*, 1981.

[64] Jonathon Luiten, Paul Voigtlaender, and Bastian Leibe. Premvos: Proposal-generation, refinement and merging for video object segmentation. In *Asian Conference on Computer Vision*, pages 565–580. Springer, 2018.

[65] Zelun Luo, Boya Peng, De-An Huang, Alexandre Alahi, and Li Fei-Fei. Unsupervised learning of long-term motion dynamics for videos. 2017.

[66] K. K. Maninis, S. Caelles, Y. Chen, J. Pont-Tuset, L. Leal-Taixé, D. Cremers, and L. Van Gool. Video object segmentation without temporal information. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 41(6):1515–1530, 2019.

[67] Michaël Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv*, 2015.

[68] Meila, Marina and Shi, Jianbo. Learning segmentation by random walks. In *Advances in neural information processing systems*, pages 873–879, 2001.

[69] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

[70] Ishan Misra, C. Lawrence Zitnick, and Martial Hebert. Shuffle and Learn: Unsupervised Learning using Temporal Order Verification. In *ECCV*, 2016.

[71] Ravi Teja Mullapudi, Steven Chen, Keyi Zhang, Deva Ramanan, and Kayvon Fatahalian. Online model distillation for efficient video inference. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3573–3582, 2019.

[72] Sourabh A Niyogi and Edward H Adelson. Analyzing gait with spatiotemporal surfaces. In *Proceedings of 1994 IEEE Workshop on Motion of Non-rigid and Articulated Objects*, pages 64–69. IEEE, 1994.

[73] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *European Conference on Computer Vision*, pages 69–84. Springer, 2016.

[74] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9226–9235, 2019.

[75] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.

[76] Andrew Owens and Alexei A Efros. Audio-visual scene analysis with self-supervised multi-sensory features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 631–648, 2018.

[77] Andrew Owens, Phillip Isola, Josh McDermott, Antonio Torralba, Edward H Adelson, and William T Freeman. Visually indicated sounds. In *Computer Vision and Pattern Recognition (CVPR)*, 2016.

[78] Deepak Pathak, Ross Girshick, Piotr Dollár, Trevor Darrell, and Bharath Hariharan. Learning features by watching objects move. In *CVPR*, 2017.

[79] F. Perazzi, J. Pont-Tuset, B. McWilliams, L. Van Gool, M. Gross, and A. Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *Computer Vision and Pattern Recognition*, 2016.

[80] Bryan Perozzi, Rami Al-Rfou, and Steven Skiena. Deepwalk. *Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining - KDD '14*, 2014.

[81] Gabriel Peyre, Marco Cuturi, and Justin Solomon. Gromov-wasserstein averaging of kernel and distance matrices. In *International Conference on Machine Learning*, pages 2664–2672, 2016.

[82] Soren Pirk, Mohi Khansari, Yunfei Bai, Corey Lynch, and Pierre Sermanet. Online object representations with contrastive learning. *arXiv preprint arXiv:1906.04312*, 2019.

[83] Hamed Pirsiavash, Deva Ramanan, and Charless C. Fowlkes. Globally-optimal greedy algorithms for tracking a variable number of objects. In *CVPR 2011*, pages 1201–1208, 2011.

[84] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv:1704.00675*, 2017.

[85] Deva Ramanan, David A Forsyth, and Andrew Zisserman. Strike a pose: Tracking people by finding stylized poses. In *CVPR*, 2005.

[86] Paul-Edouard Sarlin, Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich. Super-glue: Learning feature matching with graph neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4938–4947, 2020.

[87] Steven M Seitz and Simon Baker. Filter flow. In *2009 IEEE 12th International Conference on Computer Vision*, pages 143–150. IEEE, 2009.

[88] Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, Sergey Levine, and Google Brain. Time-contrastive networks: Self-supervised learning from video. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 1134–1141. IEEE, 2018.

[89] Jianbo Shi and Jitendra Malik. Motion segmentation and tracking using normalized cuts. In *Sixth International Conference on Computer Vision (IEEE Cat. No. 98CH36271)*, pages 1154–1160. IEEE, 1998.

[90] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000.

[91] Richard Sinkhorn and Paul Knopp. Concerning nonnegative matrices and doubly stochastic matrices. *Pacific Journal of Mathematics*, 21(2):343–348, 1967.

[92] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, pages 1929–1958, 2014.

[93] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhutdinov. Unsupervised learning of video representations using LSTMs. *arXiv*, 2015.

[94] Yu Sun, Xiaolong Wang, Zhuang Liu, John Miller, Alexei A. Efros, and Moritz Hardt. Test-time training with self-supervision for generalization under distribution shifts. In *International Conference on Machine Learning (ICML)*, 2020.

[95] Jian Tang, Meng Qu, Mingzhe Wang, Ming Zhang, Jun Yan, and Qiaozhu Mei. Line: Large-scale information network embedding. In *Proceedings of the 24th international conference on world wide web*, pages 1067–1077, 2015.

[96] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *CoRR*, abs/1906.05849, 2019.

[97] Michael Tschannen, Josip Djolonga, Marvin Ritter, Aravindh Mahendran, Neil Houlsby, Sylvain Gelly, and Mario Lucic. Self-supervised learning of video-induced visual invariances. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.

[98] Hsiao-Yu Tung, Hsiao-Wei Tung, Ersin Yumer, and Katerina Fragkiadaki. Self-supervised learning of motion capture. In *NIPS*, 2017.

[99] Jack Valmadre, Luca Bertinetto, Joao F Henriques, Ran Tao, Andrea Vedaldi, Arnold WM Smeulders, Philip HS Torr, and Efstratios Gavves. Long-term tracking in the wild: A benchmark. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 670–685, 2018.

[100] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018.

[101] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Neural Information Processing Systems (NIPS)*, 2017.

[102] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. In *International Conference on Learning Representations*, 2018.

[103] Paul Voigtlaender, Michael Krause, Aljosa Osep, Jonathon Luiten, Berin Balachandar Gnana Sekar, Andreas Geiger, and Bastian Leibe. MOTS: multi-object tracking and segmentation. *CoRR*, abs/1902.03604, 2019.

[104] Paul Voigtlaender and Bastian Leibe. Online adaptation of convolutional neural networks for video object segmentation. *arXiv*, 2017.

[105] Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking emerges by colorizing videos. In *ECCV*, 2017.

[106] Ning Wang, Yibing Song, Chao Ma, Wengang Zhou, Wei Liu, and Houqiang Li. Unsupervised deep tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1308–1317, 2019.

[107] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip HS Torr. Fast online object tracking and segmentation: A unifying approach. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1328–1338, 2019.

[108] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *CVPR*, 2018.

[109] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *ICCV*, 2015.

[110] Xiaolong Wang, Allan Jabri, and Alexei A Efros. Learning correspondence from the cycle-consistency of time. In *CVPR*, 2019.

[111] Donglai Wei, Joseph Lim, Andrew Zisserman, and William T. Freeman. Learning and using the arrow of time. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.

[112] Max Wertheimer. Laws of organization in perceptual forms. In *A source book of Gestalt psychology*, pages 71–88. Routledge & Kegan Paul, London, 1938.

[113] Josh Wills, Sameer Agarwal, and Serge Belongie. What went where. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 1, pages 37–44, Madison, WI, 2003.

[114] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3733–3742, 2018.

[115] Hongteng Xu, Dixin Luo, Hongyuan Zha, and Lawrence Carin. Gromov-wasserstein learning for graph matching and node embedding. *arXiv preprint arXiv:1901.06003*, 2019.

[116] Linjie Yang, Yanran Wang, Xuehan Xiong, Jianchao Yang, and Aggelos K Katsaggelos. Efficient video object segmentation via network modulation. 2018.

[117] Amir Roshan Zamir, Afshin Dehghan, and Mubarak Shah. Gmcp-tracker: Global multi-object tracking using generalized minimum clique graphs. In *European Conference on Computer Vision*, pages 343–356. Springer, 2012.

[118] Lihi Zelnik-Manor and Michal Irani. Event-based analysis of video. In *Proceedings of the 2001 IEEE Computer Society Conference on Computer Vision and Pattern Recognition. CVPR 2001*, volume 2, pages II–II. IEEE, 2001.

[119] Li Zhang, Yuan Li, and Ramakant Nevatia. Global data association for multi-object tracking using network flows. In *2008 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008.

[120] Richard Zhang, Phillip Isola, and Alexei A Efros. Split-brain autoencoders: Unsupervised learning by cross-channel prediction. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1058–1067, 2017.

[121] Qixian Zhou, Xiaodan Liang, Ke Gong, and Liang Lin. Adaptive temporal encoding network for video instance-level human parsing. In *ACM MM*, 2018.

[122] Tinghui Zhou, Philipp Krahenbuhl, Mathieu Aubry, Qixing Huang, and Alexei A Efros. Learning dense correspondence via 3d-guided cycle consistency. In *CVPR*, 2016.

# A  Label Noise: Effect of Identical Patches

Here, we show that false negatives that are identical to the positive – for example, patches of the sky – do not change the sign of gradient associated with the positive. Let $q$ be the query, $u$ be the positive, $V$ be the set of negatives. W.l.o.g, let the softmax temperature $\tau = 1$. The loss and corresponding gradient can be expressed as follows, where $Z$ is the partition function:

$$L(q, u, V) = u^\top q - \log[\exp u^\top q + \sum_{v \in V} \exp v^\top q] = u^\top q - \log Z$$

$$\nabla_q L(q, u, V) = u - \frac{\exp u^\top q}{Z} u - \sum_{v \in V} \frac{\exp v^\top q}{Z} v = (1 - \frac{\exp u^\top q}{Z}) u - \sum_{v \in V} \frac{\exp v^\top q}{Z} v$$

Let $V^-$ be the set of false negatives, such that $V^- \subseteq V$ and $V^+ = V \setminus V^-$. Consider the worst case, whereby $v_- = u, \forall v_- \in V^-$, so that false negatives are exactly identical to the positive:

$$\nabla_q L(q, u, V) = (1 - \frac{\exp u^\top q}{Z}) u - \sum_{v_- \in V^-} \frac{\exp v_-^\top q}{Z} v_- - \sum_{v_+ \in V^+} \frac{\exp v_+^\top q}{Z} v_+$$

$$= \underbrace{\left(1 - \frac{(1 + |V^-|) \exp u^\top q}{Z}\right)}_{\lambda_u} u - \sum_{v_+ \in V^+} \frac{\exp v_+^\top q}{Z} v_+$$

We see that the contribution of the negatives that are identical to the positive do not flip the sign of the positive gradient, i.e. $\lambda_u \geq 0$, so that in the worse case the gradient vanishes:

$$\lambda_u = 1 - \frac{(1 + |V^-|) \exp u^\top q}{Z}$$

$$= 1 - \frac{(1 + |V^-|) \exp u^\top q}{(1 + |V^-|) \exp u^\top q + \sum_{v_+ \in V^+} \exp v_+^\top q}$$

$$\geq 0$$

# B  Comparison to Supervised Methods on DAVIS-VOS

The proposed method outperforms many supervised methods for video object segmentation, despite relying on a simple label propagation algorithm, not being trained for object segmentation, and not training on the DAVIS dataset. We also show comparisons to pretrained feature baselines with larger networks.

| Method | Backbone | Train Data (#frames) | $\mathcal{J}\&\mathcal{F}_m$ | $\mathcal{J}_m$ | $\mathcal{J}_r$ | $\mathcal{F}_m$ | $\mathcal{F}_r$ |
|---|---|---|---|---|---|---|---|
| OSMN [116] | VGG-16 | I/C/D (1.2M + 227k) | 54.8 | 52.5 | 60.9 | 57.1 | 66.1 |
| SiamMask [107] | ResNet-50 | I/V/C/Y (1.2M + 2.7M) | 56.4 | 54.3 | 62.8 | 58.5 | 67.5 |
| OSVOS [12] | VGG-16 | I/D (1.2M + 10k) | 60.3 | 56.6 | 63.8 | 63.9 | 73.8 |
| OnAVOS [104] | ResNet-38 | I/C/P/D (1.2M + 517k) | 65.4 | 61.6 | 67.4 | 69.1 | 75.4 |
| OSVOS-S [66] | VGG-16 | I/P/D (1.2M + 17k) | 68.0 | 64.7 | 74.2 | 71.3 | 80.7 |
| FEELVOS [103] | Xception-65 | I/C/D/Y (1.2M + 663k) | 71.5 | 69.1 | 79.1 | 74.0 | 83.8 |
| PReMVOS [64] | ResNet-101 | I/C/D/P/M (1.2M + 527k) | 77.8 | 73.9 | 83.1 | 81.8 | 88.9 |
| STM [74] | ResNet-50 | I/D/Y (1.2M + 164k) | 81.8 | 79.2 | - | 84.3 | - |
| ImageNet [42] | ResNet-50 | I (1.2M) | 66.0 | 63.7 | 74.0 | 68.4 | 79.2 |
| MoCo [41] | ResNet-50 | I (1.2M) | 65.4 | 63.2 | 73.0 | 67.6 | 78.7 |
| **Ours** | ResNet-18 | K (20M unlabeled) | 67.6 | 64.8 | 76.1 | 70.2 | 82.1 |

Table 3: **Video object segmentation results on DAVIS 2017 val set**. We show results of state-of-the-art **supervised** approaches in comparison to our unsupervised one (see main paper for comparison with unsupervised methods). Key for *Train Data* column: I=ImageNet, K=Kinetics, V = ImageNet-VID, C=COCO, D=DAVIS, M=Mapillary, P=PASCAL-VOC Y=YouTube-VOS. $\mathcal{F}$ is a boundary alignment metric, while $\mathcal{J}$ measures region similarity as IOU between masks.

## C    Using a Single Feature Map for Training

We follow the simplest approach for extracting nodes from an image without supervision, which is to simply sample patches in a convolutional manner. The most efficient way of doing this would be to only encode the image once, and pool the features to obtain region-level features [60].

We began with that idea and found that the network could cheat to solve this dense correspondence task even across long sequences, by learning a shortcut. It is well-known that convolutional networks can learn to rely on boundary artifacts [60] to encode position information, which is useful for the dense correspondence task. To control for this, we considered: 1) removing padding altogether; 2) reducing the receptive field of the network to the extent that entries in the center crop of the spatial feature map do not see the boundary; we then cropped the feature map to only see this region; 3) randomly blurring frames in each video to combat space-time compression artifacts; and 4) using *random* videos made of noise. Surprisingly, the network was able to learn a shortcut in each case. In the case of random videos, the shortcut solution was not nearly as successful, but we still found it surprising that the self-supervised loss could be optimized at all.

## D    Frame-rate Ablation

**Effect of frame-rate at training time**    We ablate the effect of frame-rate (i.e. frames per second) used to generate sequences for training, on downstream object segmentation performance. The case of infinite frame-rate corresponds to the setting where the *same* image is used in each time step; this experiment is meant to disentangle the effect of data augmentation (spatial jittering of patches) from the natural "data augmentation" observed in video. We observe that spatio-temporal transformations is beneficial for learning of representations that transfer better for object segmentation.

| Frame rate | $\mathcal{J}\&\mathcal{F}_\mathrm{m}$ |
|---|---|
| 2 | 65.9 |
| 4 | 67.5 |
| 8 | 67.6 |
| 30 | 62.3 |
| $\infty$ | 57.5 |

## E    Hyper-parameters

We list the key hyper-parameters and ranges considered at training time. Due to computational constraints, we did not tune the patch extraction strategy, nor several other hyper-parameters. The hyper-parameters varied, namely edge dropout and video length, were ablated in Section 3 (shown in bold). Note that the effective training path length is twice that of the video sequence length.

| *Train* Hyper-parameters | Values |
|---|---|
| Learning rate | 0.0001 |
| Temperature $\tau$ | 0.07 |
| Dimensionality $d$ of embedding | 128 |
| Frame size | 256 |
| Video length | **2, 4, 6, 10** |
| Edge dropout | **0, 0.05, 0.1, 0.2, 0.3** |
| Frame rate | **2, 4, 8, 30** |
| Patch Size | 64 |
| Patch Stride | 32 |
| Spatial Jittering (crop range) | (0.7, 0.9) |

We tuned test hyper-parameters with the ImageNet baseline. In general, we found performance to increase given more context. Here, we show hyper-parameters used in reported experiments; we largely follow prior work, but for the case of DAVIS, we used 20 frames of context.

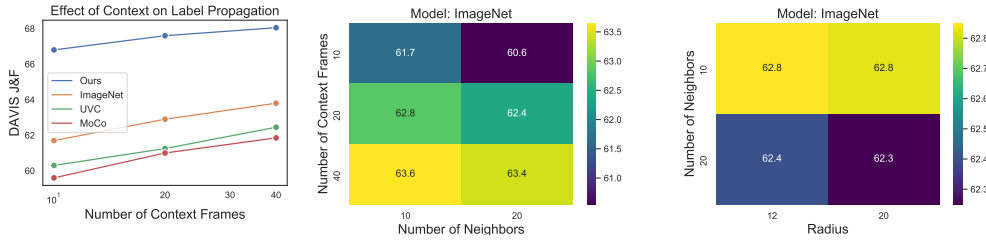| *Test* Hyper-parameters | Values |
|---|---|
| Temperature $\tau$ | 0.07 |
| Number of neighbors $k$ | **10**, 20 |
| Number of context frames $m$ | Objects: 20 |
| | Pose: 7 |
| | Parts: 4 |
| Spatial radius of source nodes | **12**, 20 |

# F   Label Propagation

We found that the performance of baselines can be improved by carefully implementing label propagation by $k$-nearest neighbors. When compared to baseline results reported in [58] and [56], the differences are:

1. Restricting the set of source nodes (context) considered for each target node, on the basis of spatial locality, i.e. *local* attention. This leads to a gain of $+4\%$ J&F for the ImageNet baseline.

   Many of the task-specific approaches for temporal correspondence incorporate restricted attention, and we found this rudimentary form to be effective and reasonable.

2. Computing attention over all source nodes at once and selecting the top-$k$, instead of independently selecting the top-$k$ from each frame. This leads to a gain of $+3\%$ J&F for the ImageNet baseline.

   This is more natural than computing nearest neighbors in each frame individually, and can be done efficiently if combined with local attention. Note that the softmax over context can be performed after nearest neighbors retrieval, for further efficiency.

## F.1   Effect of Label Propagation Hyper-parameters



We study the effect of hyper-parameters of the label propagation algorithm, when applied with strong baselines and our method. The key hyper-parameters are the length of context $m$, the number of neighbors $k$, and the search radius $r$. In the figures above, we see the benefit of adding context (see left, with $k = 10, r = 12$), effect of considering more neighbors (middle, with $r = 12$), and effect of radius (right, with $m = 20$).

# G   Encoder Architecture

We use the ResNet-18 network architecture, modified to increase the resolution of the output convolutional feature map. Specifically, we modify the stride of convolutions in the last two residual blocks from 2 to 1. This increases the resolution by a factor of four, so that the downsampling factor is $1/8$. Please refer to Table 4 for a detailed description.

For evaluation, when applying our label propagation algorithm, we report results using the output of `res3` as node embeddings, for fair comparison to pretrained feature baselines ImageNet, MoCo, and VINCE, which were trained with stride 2 in `res3` and `res4`. We also found that `res3` features compared favorably to `res4` features.

| Layer | Output | Details |
|---|---|---|
| input | $H \times W$ | |
| conv1 | $H/2 \times W/2$ | $7 \times 7$, 64, stride 2 |
| maxpool | $H/4 \times W/4$ | stride 2 |
| res1 | $H/4 \times W/4$ | $\begin{bmatrix} 3 \times 3, 64 \\ 3 \times 3, 64 \end{bmatrix} \times 2$, stride 1 |
| res2 | $H/8 \times W/8$ | $\begin{bmatrix} 3 \times 3, 128 \\ 3 \times 3, 128 \end{bmatrix} \times 2$, stride 2 |
| res3 | $H/8 \times W/8$ | $\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$, stride 2̶ 1 |
| res4 | $H/8 \times W/8$ | $\begin{bmatrix} 3 \times 3, 512 \\ 3 \times 3, 512 \end{bmatrix} \times 2$, stride 2̶ 1 |

Table 4: Modified ResNet-18 Architecture. Our modifications are shown in blue.

## H   Test-time Training Details

We adopt the same hyper-parameters for optimization as in training: we use the Adam optimizer with learning rate 0.0001. Given an input video $I$, we fine-tune the model parameters by applying Algorithm 1 with input frames $\{I_{t-m}, ..., I_t, ..., I_{t+m}\}$, *prior* to propagating labels to $I_t$. For efficiency, we only finetune the model every 5 timesteps, applying Adam for 100 updates. In practice, we use $m = 10$, which we did not tune.

## I   Utility Functions used in Algorithm 1

**Algorithm 2** Utility functions.

```
// psize : size of patches to be extracted

import torch
import kornia.augmentation as K

# Turning images into list of patches
unfold = torch.nn.Unfold((psize, psize), stride=(psize//2, psize//2))

# l2 normalization
l2_norm = lambda x: torch.nn.functional.normalize(x, p=2, dim=1)

# Slightly cropping patches once extracted
spatial_jitter = K.RandomResizedCrop(size=(psize, psize), scale=(0.7, 0.9), ratio=(0.7, 1.3))
```