# Time Does Tell: Self-Supervised *Time-Tuning* of Dense Image Representations

Mohammadreza Salehi, Efstratios Gavves, Cees G. M. Snoek, Yuki M. Asano

QUVA Lab, University of Amsterdam

`(s.salehidehnavi, e.gavves, c.g.m.snoek, y.m.asano)@uva.nl`

## Abstract

*Spatially dense self-supervised learning is a rapidly growing problem domain with promising applications for unsupervised segmentation and pretraining for dense downstream tasks. Despite the abundance of temporal data in the form of videos, this information-rich source has been largely overlooked. Our paper aims to address this gap by proposing a novel approach that incorporates temporal consistency in dense self-supervised learning. While methods designed solely for images face difficulties in achieving even the same performance on videos, our method improves not only the representation quality for videos – but also images. Our approach, which we call time-tuning, starts from image-pretrained models and fine-tunes them with a novel self-supervised temporal-alignment clustering loss on unlabeled videos. This effectively facilitates the transfer of high-level information from videos to image representations. Time-tuning improves the state-of-the-art by 8-10% for unsupervised semantic segmentation on videos and matches it for images. We believe this method paves the way for further self-supervised scaling by leveraging the abundant availability of videos. The implementation can be found here : https://github.com/SMSD75/Timetuning*

## 1. Introduction

Dense self-supervised learning, whereby meaningful deep features for each pixel or patch of input are learned in an unsupervised manner, has recently received increasing attention [74, 24, 57, 51]. By learning spatially consistent features for different views of an input, strong gains in unsupervised semantic segmentation have been achieved using unlabeled images. However, so far, an even more information-rich source for unsupervised training has been largely overlooked: videos. With their additional time dimension and being the most rapidly growing form of digital content, they are well-suited to scaling dense self-supervised learning even further.

Some efforts have already been made to learn from the video domain by using different frames from a video as aug-
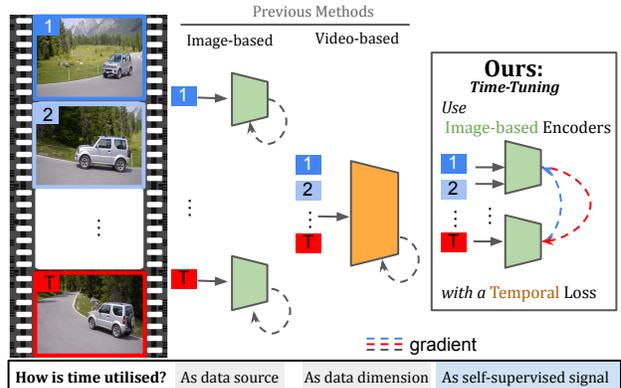


Figure 1: **Time-tuning compared to previous methods.** Unlike existing methods that ignore or utilize expensive 3D models to implicitly model temporal information, the proposed method explicitly incorporates temporal consistency in dense feature representations using a temporal self-supervised loss. The method starts with a 2D encoder pretrained on images and fine-tunes it using unlabeled videos. This approach leads to improved performance not only for videos but also for images.

mentations [18, 63, 49] or by involving temporal correspondence [71, 58]; however, they mostly did it in a supervised way [29, 73, 19, 40, 65, 38, 31, 28, 41, 37, 54, 44], which is not scalable specifically for dense tasks where the number of targets can increase significantly as the number of pixels grows. To this end, self-supervised learning approaches offer a solution by reducing the need for supervision. However, these methods typically rely on the notion of "views", which involves learning similar features for corresponding locations over time. This usually leads to a chicken-and-egg problem, where the correspondences are required for learning dense features – which in turn enable good correspondences [30].

In images, the challenge is trivially solved by considering the correspondence function based on the augmentation function. For instance, In the case of color augmentations, this correspondence is simply given by the identity. However, shifts in time cannot be viewed as mere augmenta-

tions. As we demonstrate through a new evaluation protocol, using image models on frames alone is not nearly as effective. A similar finding has also been reported, albeit for non-dense works [20, 34], which assumed the passage of time as an image augmentation. These works have generally reported reduced performances, even when compared to simple image-based pretraining methods. Similarly, video-level tasks [17, 18, 55] assume sufficiently similar semantics between different frames. This is also not true for dense tasks, as static features can only be assumed where nothing is moving – which is rare due to possible object, camera, and background motion between the frames.

To address this challenge, we propose to model the additional time-dimension explicitly to identify which pixels should retain similar embeddings and which should not. We propose two separate modules to tackle the correspondence and the dense learning, respectively. For the former, we introduce the Feature-Forwarder (FF) module, which breaks the mentioned chicken-and-egg loop by leveraging the good tracking performance of pretrained image models, and allows an approximate second "view" that can then be treated as a target for the further dense self-supervised loss. On top of this, we introduce a spatio-temporally dense clustering module, which learns unsupervised clusters across samples, locations and time. Using these two components and starting from image-pretrained features, our proposed method allows *time-tuning* (TIMET) the dense representation in a self-supervised manner, see Figure 1.

Finally, we demonstrate that TIMET paves the way for further scaling of self-supervised learning by leveraging the abundant availability of video datasets and transferring their knowledge to the image domain. This results in consistently achieving state-of-the-art performances not only for the task of unsupervised semantic segmentation of videos, but also for unsupervised *image* semantic segmentation, a feat previously out of reach for methods trained on videos.

Overall, this paper makes the following contributions:

- We show that image-based unsupervised dense segmentation models applied to videos exhibit degraded performance and lack temporal consistency in their segmentation maps.

- Building on this observation, we propose a novel dense self-supervised learning method that utilizes temporal consistency as a learning signal.

- We demonstrate that our method enables the scaling of self-supervised learning by leveraging abundant video datasets and effectively transferring knowledge to the image domain. Our approach consistently achieves state-of-the-art performance for both images and videos, opening up new opportunities in the field.

## 2. Related Works

**Dense self-supervised learning.** These methods build upon image-level self-supervised representation learning by incorporating existing losses to enhance the spatial features, demonstrating a commendable advancement. DenseCL [61] works on spatial features by constructing dense correspondences across views using the contrastive objective given in MoCo [23], while PixPro [62] utilizes the augmentation wrapper to get the spatial correspondence of the pixel intersection between two views. Similarly, MaskContrast [57], Leopart [74], DetCon [24] and Odin [25] also ensure spatial feature similarities via contrastive learning and spatial correspondences. Self-Patch [67] treats the spatial neighbors of the patch as positive examples for learning more semantically meaningful relations among patches. Inspired by SelfPatch, AD-CLR [69] proposes patch-level contrasting via query crop and cross-attention mechanism. Pursuing the same objective through none end-to-end approaches. Both [68] and [59] propose an unsupervised salient object segmentation pipeline that extracts noisy object masks from the inputs and fine-tunes a specifically designed object segmentation head with several self-training steps to make an unsupervised semantic segmentation model. Unlike the existing image-based approaches, we propose a method that improves the dense prediction performance of a pretrained encoder by explicitly modeling the temporal dimension and learning from diverse natural dynamics and variations found in *videos*.

**Video to image knowledge transfer.** Learning from videos instead of images has recently received attention since they contain far more information than still images and hold the potential for learning rich representations of the visual world. To this end, several *non-dense* works have been released [46, 4, 27]. VITO [46] shows that naively applying image domain self-supervised learning methods on videos can lead to a performance drop coming from the distribution shift, which by applying data processing techniques is relieved but not fully solved. Following the same way, [27] trains masked autoencoders [22] with contrastive learning [10] on video datasets and shows a decent performance on both image and video tasks. While video to image knowledge transfer is the end goal of such methods and our paper; our dense task is different and is as of yet unaddressed.

**Unsupervised video object segmentation.** Unsupervised video object segmentation methods do not require any manual annotations. However, they are only designed to tackle a *foreground/background* segmentation task, which refers to the segmentation of the most prominent, general objects in video sequences [66, 43, 50, 1, 39]. For instance, [66]
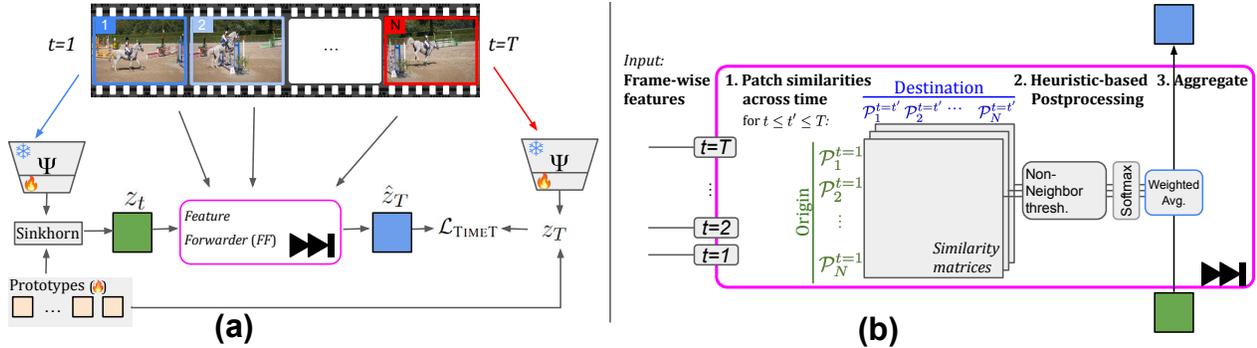
Figure 2: **Time-Tuning overview.** We conduct self-supervised video semantic segmentation by tuning image-pretrained models with temporal information from videos. **(a)**: Our general pipeline of adapting an image-pretrained ViT transformer on video data using our dense clustering loss. The encoder is specified by $\Psi$ and kept frozen except for the last two layers. Also, $z_t$ shows the output of Sinkhorn-Knopp algorithm, which is forwarded to time-step $T$ to be compared with the features that are obtained using the last time-step. **(b)**: Detailed view into our Feature-Forwarder module that is used for aligning cluster-maps from past frames to logit-features of future ones.

uses an AutoEncoder-based architecture based on slot attention so that pixels that show the same motion flows based on optical-flow are grouped together. In this way, dominant moving objects can be detected and separated from the background. In our experiments, we benchmark two representative instances of this line of work, *i.e.*, [66, 1], to provide better insights into the differences of unsupervised object detection methods and our proposed benchmark for video semantic segmentation.

To conclude, we propose the first self-supervised video semantic segmentation method, which tunes image-trained models such as [9, 72] based on the temporal information that exists in unlabeled videos to enhance their effectiveness in dense prediction tasks.

## 3. Method

Our method works by densely clustering features in a manner that is consistent with time. At a high level, it works by transforming features from a past time $t'=t$ to the current frame $t'=T$ and forcing these to be consistent with the currently observed ones at $t'=T$. The overall learning signal comes from the dense temporal clustering of forwarded features and allows "timetuning" image-pretrained models to learn from videos. Figure 2 shows the proposed method's architecture. In the following, we describe each component in detail.

### 3.1. Feature Forwarding

The goal of this component, the *feature-forwarder* (FF), is to propagate the dense features in time as accurately as possible, given some RGB input frames. More formally, let $f$ be the feature forwarding function, $I_t$ be a frame at time $t$ for a given video, and $z_t$ the features corresponding to $I_t$,

then the task is to propagate $z_t$ to a future time $T$,

$$f([I_t, \ldots, I_T], z_t) \rightarrow \hat{z}_T. \qquad (1)$$

Note that $\hat{z}_T$ and $z_T$ are not necessarily the same for $T > t$; since $\hat{z}_T$ represents an approximation of $z_T$ using the previous features. Obviously, the further $T$ is from $t$, the more variance is seen and therefore a better learning signal might be provided; however, for such larger intervals, the idea of feature forwarding encounters two main challenges: error accumulation and object non-permanence.

**Challenge 1: Error accumulation.** While computing how features have changed simply based on the start and end frames is the most straight-forward (*i.e.*, $I_t$ and $I_T$), it does not fully leverage the knowledge contained in the intermediate frames. Instead, it is common [30, 60, 36] to forward features for every time-step $\delta t$, *i.e.*,

$$f([I_t, .., I_T], z_t) = \bigcirc_{t'=t+\delta t}^{T} f([I_{t'-\delta t}, I_{t'}], z_t), \qquad (2)$$

where $\bigcirc$ is a composition operator such as element-wise multiplication. However, this results in the accumulation of errors over time.

**Challenge 2: Object non-permanence.** One difficulty that arises when using videos as training data is the fact that objects sometimes simply disappear from the screen. Compared to images, which carry a heavy photographer's bias, videos are more prone to natural variation due to camera and object motion resulting in temporary object *non-permanence* in a particular video clip. A method that simply assumes an object (or a feature) has to be present in any given frame simply because it was present in the previous

one might therefore easily fail when going towards videos that exist in-the-wild. A case in point are simple occlusions that arise from an object being fully- or partially-covered by another one further in the background. Taking longer training clips has been attempted to resolve the issue. However, this increases the training complexity, especially for dense-tasks that necessitate a higher number of predictions for each input.

*We address these challenges using a novel, yet simple component, which we call the* **Stabilizing Feature-Forwarder.** The first choice in designing this module is the function and data used for composing information across time. For instance, a simple approach would be to use pixel-wise optical flow and aggregate the flow at the feature level to predict how the features change over time. However, as we also show in our experiments, optical flow is prone to accumulating error with time and also cannot easily recover from occlusions [6]. Instead, we utilize the fact that we wish to pretrain the visual encoder and *recycle* part of it for the purpose of feature forwarding. This not only speeds up the process, as no additional encoder/modality is required, but also produces a positive feedback loop of better forwarding and better feature learning.

Concretely, we first compute the L2-normalised dense features of a pretrained visual encoder $\Phi$ to compute spatial similarities across time (the $'$ in $t'$ left out for clarity):

$$F_{ij}^t = \langle \Phi(I_t)_i, \Phi(I_T)_j \rangle / \tau, \tag{3}$$

which yields matrices $F^{\{t,...T-\delta t\}} \in R^{N \times N}$ with values between [0,1] indicating semantic similarities between the $N$ spatial features that are sharpened by a temperature $\tau$. Note that because we compare each frame with the final frame, the effect of object non-permanence can be minimized: even if the object is only present in every other frame, there is enough signal to forward its features.

Next to propagate from time $t$ to $T$, these similarities are stacked and normalised along time, as follows:

$$\tilde{F}_{ij}^t = \exp(\mathcal{N}(F_{ij}^t)) / \sum_{i,t'} \exp(\mathcal{N}(F_{ij}^{t'})), \tag{4}$$

where $t \leq t' \leq T$ and $\mathcal{N}$ is a neighborhood thresholding method which forces short-term spatial-smoothness, *i.e.*, $\mathcal{N}(\tilde{F}_{ij}^t)=0$ if $i$ is not within a local window around $j$ with the size $k$. Finally, the propagated feature are computed $\hat{z}_T$ as:

$$\hat{z}_T(j) = \sum_{t',i} \tilde{F}_{ij}^{t'} \hat{z}_{t'}(i), \quad j \in \{1,..,N\}. \tag{5}$$

This means that for arriving at the target feature $\hat{z}_T$, we not only use the previous frame as the source, but instead use the past $(T-t)/\delta t$ frames and aggregate these.

**Relation to mask-propagation methods.** While previous methods such as DINO [9] and STC [30] have utilized a similar technique for propagating ground-truth masks for evaluating, for example on the DAVIS dataset, there are three key differences to our forwarding method. First, instead of propagating binary maps of foreground-vs-background, our feature-forwarder has as inputs soft, noisy and multi-label self-supervised segmentation maps, which require the forwarder to tolerate overlapping of different object probabilities throughout training. Second, previous methods have used this approach mainly for inference; we, however, use this module as a trainable component and show how our loss improves this forwarder with training time. Finally, we do not follow a typical re-normalizing step across feature dimensions (typically done after Eq. (5)) as this harms the scale of logits that are being propagated and leads to heavily diluted target distributions.

## 3.2. Self-supervised dense clustering

While the Feature-Forwarder produces target features that include information about the dynamics with time, its computation utilizes $\Phi(I_T)$ and could lead to trivial solutions, *i.e.*, $f([I_t,..,I_T], z_t) = \Phi(I_T)$. To counteract this, we propose a self-supervised clustering task across views in time. For this, we utilize the basic online clustering algorithm based optimal-transport [12], utilized in works such as SeLa [2], SwAV [8] and DINOv2 [45]. In particular, let $\Psi$ be a visual encoder with a clustering head $g$ that yields a $K$ dimensional output, then the clustering loss is given by:

$$\mathcal{L}(x_i) = -\tilde{y}_i \log(g(\Psi(x)_i)), \tag{6}$$

which is a standard cross-entropy loss with regards to self-supervised pseudo-label $\tilde{y}_i$. These labels, in turn, are generated by solving an entropy-regularised optimal transport problem on the batch $\mathcal{B}$ [2]:

$$\min_{\tilde{y}} \langle \tilde{y}, -\log g(\Psi(x)) \rangle + \frac{1}{\lambda} \text{KL}(\tilde{y} \| rc^\top), \tag{7}$$

$$\text{with } r = \frac{1}{K} \cdot \mathbb{1}, \quad c = \frac{1}{|\mathcal{B}|} \cdot \mathbb{1}. \tag{8}$$

Here $\lambda$ controls the entropy regularisation and $r, c$ are the marginals for the prototypes and the batch, respectively. Note that solving this problem can be done extremely quickly on the GPU and yields soft pseudo-labels $\tilde{Y}$, such that $\text{argmax}(\tilde{Y}) = \tilde{y}$.

n, these might yield sufficiently similar semantics [17] this is not true for dense tasks, as static features can only be assumed for videos where nothing is moving, which is rare.

## 3.3. Overall TIMET **loss**

We combine the previous two modules to arrive at our full method (as shown in Fig. 2): First, self-supervised

Sinkhorn-Knopp clustering is conducted on early features $g(\Psi(I_t))$, yielding soft pseudo-labels $\mathrm{SK}(g(\Psi(I_t)))=\tilde{Y}_t$. These are then forwarded in time using our Feature-Forwarder to arrive at dense targets $\mathrm{FF}(\tilde{Y}_t)$, which are used in the final loss. Compactly:

$$\mathcal{L}_{\mathrm{TIMET}}(I_T) = -\sum_{i,j} \mathrm{FF}(\tilde{Y}_t) \log(g(\Psi(I_T))). \quad (9)$$

## 4. Experiments

### 4.1. Setup

**Datasets.** We train our method and baselines on *YTVOS* [64], one of the largest video segmentation datasets available, and evaluate on *DAVIS17* [48] and YTVOS. For YTVOS the ground truth masks are only available for the first frames of the test and validation sets, and therefore, a fixed random 20% of the training set is used for testing. For transfer learning experiments, we use the validation set of *Pascal VOC 2012* [16], As the dataset has been commonly used as a main reference for recent works in dense self-supervised image segmentation [74, 57, 61]. For completeness, we also report the performance on egocentric datasets that have less object-centric bias and are prevalent in real-world scenarios. For egocentric experiments, we train on *EPIC-KITCHENS-100* [13] and evaluate on *VISOR* [14]. Further details are provided in Appendix A.

**Models and baselines.** Currently, there is a lack of available unsupervised semantic video semantic segmentation methods. Nevertheless, we have included a comprehensive comparison of our method with state-of-the-art techniques in both image and video domains for unsupervised image semantic segmentation and unsupervised video object segmentation. To evaluate the image-based models, we utilized either official reported numbers or provided pretrained models of STEGO [21] and Leopart [74]. We have taken every measure to ensure fairness in our comparison. To do so, all the used pretrained models have the same pretraining dataset (ImageNet-1k), and the number of their backbone parameters is roughly similar. For those models that use extra datasets, for instance, Leopart, we select the pretrained backbones that closely match the specification of YTVOS training dataset. Additionally, we trained image-based models on the same video datasets, where we converted video data into their corresponding image data and reported their performance wherever possible. To ensure comprehensive analysis, we have included the recent concurrent work Flowdino [71], which utilizes optical flow to refine DINO features on unlabelled videos, in our comparison benchmarks as well. This is done only in cases where there is a shared experimental setup.

**Evaluation procedure.** Despite unsupervised object segmentation being a well-established evaluation in the image domain [57, 74], evaluating unsupervised video multi-label object segmentation is challenging due to the absence of an established evaluation protocol for video object semantic segmentation *without* supervision. In this regard, we propose a set of evaluation protocols for unsupervised video multi-label object segmentation, which exploits existing video object segmentation datasets for evaluation purposes (see details in Appendix B). In our experiments, we discard any projection head used during training and evaluate ViT's spatial tokens directly, similar to [57, 32, 61], using four methods: classification with a linear head, classification with an FCN head, overclustering, and clustering with as many clusters as the number of ground truth objects. To fine-tune different heads on top of the frozen spatial tokens, we follow [57]. For unsupervised semantic segmentation, we apply $K$-Means on all spatial tokens, where $K$ is chosen to be higher or equal to the ground-truth number of classes, following the common protocol in image clustering [32, 56]. For grouping clusters to ground-truth classes, we match them either by pixel-wise precision or Hungarian matching on merged cluster maps [35]. For video datasets, we allow the matching to be per-frame, per-video, and per-dataset. The evaluation metric is permutation-invariant, following [32], and the results are averaged over five different seeds. Clustering evaluations are preferred as they require less supervision and fewer hyperparameters than training a linear classifier and operate directly on the learned embedding spaces. We report our results in mean Intersection over Union (mIoU) unless otherwise specified.

**Model training.** We train a ViT-Small with patch size 16 and initialized from ImageNet-pretrained DINO weights [9] using the proposed self-supervised loss function. For the ablations and the main experiments, our models are trained for 12 and 30 epochs respectively. Further training details are provided in Appendix A. Code will be made available.

### 4.2. Ablations

We first examine the essential parameters of our method by training TIMET on YTVOS and assessing its ability to perform semantic segmentation on Pascal VOC. In addition to presenting the results through clustering and overclustering, we also demonstrate linear classification outcomes.

**Number of prototypes.** We first ablate the influence of the number of prototypes on downstream segmentation performance. Results in Table 1a indicate a sharp increase in performance with a rise in the number of prototypes, but once a sufficiently large number is reached, performance stabilizes. We observe peak performance with a moderate number of 200 prototypes. The stability of our method over

| (a) Ablating # prototypes $K$. | | | |
|---|---|---|---|
| $P$ | LC | K=21 | K=500 |
| 10 | 37.5 | 6.1 | 24.1 |
| 50 | 58.2 | 8.2 | 40.5 |
| **200** | **59.7** | **9.2** | **42.8** |
| 300 | 59.4 | 9.0 | 42.3 |

| (b) Ablating time interval $\delta T$. | | | | |
|---|---|---|---|---|
| # frames | $\delta T$ | LC | K=21 | K=500 |
| 1 | 0s | 50.7 | 5.6 | 26.2 |
| 4 | 0.2s | 56.5 | 7.4 | 36.5 |
| 4 | 0.5s | **59.7** | **9.2** | **42.8** |
| 4 | 1.0s | 57.1 | 8.3 | 38.1 |

| (c) Ablating number of frames used. | | | | |
|---|---|---|---|---|
| # frames | $T$ | LC | K=21 | K=500 |
| 1 | 0s | 50.7 | 5.6 | 26.2 |
| 2 | 2.0s | 52.6 | 6.2 | 37.1 |
| 4 | 2.0s | **59.7** | **9.2** | **42.8** |
| 8 | 2.0s | 59.7 | 9.0 | 42.3 |

Table 1: **Ablations of the key parameters of our method.** The model is trained for 12 epochs on Pascal VOC, and results for unsupervised segmentation with clustering ($K$=21), overclustering ($K$=500), and linear pixel-wise classification (LC) are shown. The stability of our method over a range of prototypes (50-300), inter-frame time intervals ($\delta T \in$[0.5s-1.0s]), and the number of training frames (4-8) at a fixed clip duration ($T$) shows the robustness of the method.

a range of prototypes (50-300) suggests that our approach is robust to the change of this parameter.

**Understanding FF.** Next, we ablate the temporal dimensions that influence the working of the Feature-Forwarder. In Table 1b, we vary the time interval $\delta T$ between frames whilst keeping a fixed number of four frames. We generally find an increase in performance when increasing the time-interval, with performance peaking at 0.5s. One of the most critical observations is the clear difference in clustering performance from 26.2% to 42.8% between training on single frames (first row) *vs.* training with multiple frames. This clearly indicates the significant benefit of utilizing temporal information from the Feature-Forwarder during training.

Next, in Table 1c we vary the number of frames given a fixed clip duration of 2s. Again, we find that a moderate number of frames between the source and target time is most favorable. While additional frames do not degrade the performance much, they do add computational complexity, so we prefer using 4 frames in our main method.

**Choice of propagation feature.** Finally, in Table 2, we vary the *type of feature* that we forward to a future time. In the first row, we report the performance when not using any feature forwarding – thus solely training on single-frame inputs (*i.e.*, the same as row 1 in Table 1c). In the second row, "Identity" shows the performance that is obtained when the feature map from the source frame is simply forwarded to the future without any changes, which shows an increase in the performance compared to the first row. This shows that often, training videos are not very dynamic, and a static assumption can already lead to some gains. However, compared to forwarding the Sinkhorn-Knopp (SK) clustered features, these gains are small (+4.6% *vs.* + 17%). Importantly, our Feature-Forwarder relies on forwarding SK-sharpened feature maps and does not work when simply forwarding the network's output logits $\Phi(x)$, followed by a clustering step. The reason for obtaining considerably lower numbers here can be traced back to using an un-entropy-

| FF type | LC | K=21 | K=500 |
|---|---|---|---|
| None | 50.7 | 5.6 | 26.2 |
| Identity | 55.3 | 7.4 | 36.1 |
| $\Psi(x)$ | 53.1 | 6.6 | 35.5 |
| **SK$(\Psi(x))$** | **59.7** | **9.2** | **42.8** |

Table 2: **Propagating different features in FF** on Pascal VOC. We find that our Sinkhorn-Knopp (SK) based module outperforms static training ('None' or 'Identity'). Moreover, propagating logits $\Psi(x)$ does not work as well as SK-regularised features.

| | | At Init | | +TIMET | |
|---|---|---|---|---|---|
| Pretrain | Backbone | K=500 | LC | K=500 | LC |
| MSN [3] | ViT-S/16 | 26.0 | 55.4 | 48.3$_{\uparrow 22.3}$ | 67.2$_{\uparrow 11.8}$ |
| iBOT [72] | ViT-S/16 | 32.1 | 62.1 | 47.1$_{\uparrow 15.0}$ | 67.1$_{\uparrow 5.0}$ |
| DINO [9] | ViT-S/8 | 22.5 | 55.8 | 53.9$_{\uparrow 31.4}$ | 69.5$_{\uparrow 13.7}$ |
| DINO [9] | ViT-B/16 | 28.9 | 59.1 | 52.7$_{\uparrow 23.8}$ | 70.2$_{\uparrow 11.1}$ |

Table 3: **Applying TIMET to different pretrainings** on Pascal VOC. TIMET can boost ($\uparrow$) the performance of different backbones with different initialization by a considerable margin, showing the generality of our approach.

regularized clustering algorithm. In this case, the logits tend to cause a few prototypes to dominate the cluster centers, exacerbated after the propagation, resulting in highly noisy and uninformative propagated logits. This shows that careful design of the Feature-Forwarder is indeed required.

**Applying TIMET to different models.** As shown in Table 3, our method is generalisable to different backbones and self-supervised learning initializations, enabling it to be an effective method to transfer the knowledge of the videos to images in an unsupervised way, which reduces the cost of labeling for different downstream tasks.
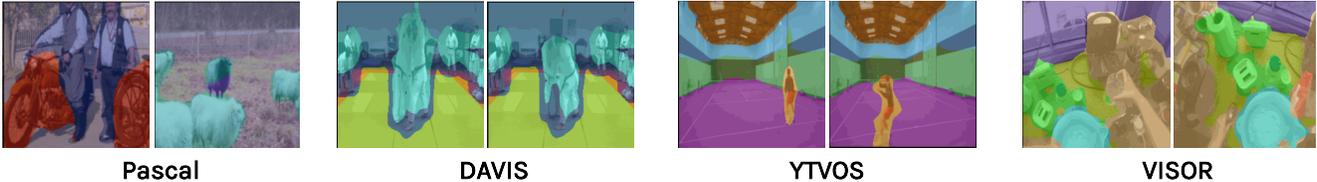
| Pascal | DAVIS | YTVOS | VISOR |

Figure 3: **TIMET unsupervised segmentations.** As we regularize DINO's backbone to be consistent across time on YTVOS, it obtains strong performance on both image and video segmentation datasets, yielding high class consistencies (indicated by the segmentation colors) and tight borders. We provide more qualitative results in Appendix D.

| | Clustering | | | | | | Overclustering | | | | | |
| | YTVOS | | | DAVIS | | | YTVOS | | | DAVIS | | |
| | *F* | *C* | *D* | *F* | *C* | *D* | *F* | *C* | *D* | *F* | *C* | *D* |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| *Trained on Images* | | | | | | | | | | | | |
| Resnet50 | 44.0 | 43.4 | 1.7 | 39.3 | 37.4 | 4.2 | 55.6 | 52.8 | 3.1 | 46.6 | 44.2 | 8.4 |
| SwAV [8] | 39.5 | 38.2 | 3.2 | 32.0 | 29.6 | 7.3 | 59.8 | 58.1 | 5.8 | 50.5 | 50.1 | 25.7 |
| DINO [9] | 39.1 | 37.9 | 1.9 | 30.2 | 31.0 | 1.6 | 66.2 | 65.4 | 4.0 | 56.9 | 54.9 | 17.9 |
| Leopart [74] | 39.2 | 37.9 | 11.7 | 30.3 | 30.2 | 16.5 | 64.5 | 62.8 | 15.5 | 54.9 | 54.4 | 26.7 |
| *Trained on Videos* | | | | | | | | | | | | |
| STEGO* | 41.5 | 40.3 | 2.0 | 31.9 | 31.0 | 3.2 | 58.1 | 54.3 | 5.1 | 47.6 | 46.3 | 10.4 |
| DINO* | 37.2 | 36.1 | 1.2 | 29.3 | 29.2 | 2.4 | 53.1 | 50.9 | 1.3 | 45.4 | 44.0 | 8.6 |
| Leopart* | 41.5 | 40.5 | 7.7 | 37.5 | 36.5 | 12.6 | 60.8 | 59.8 | 6.8 | 53.7 | 53.1 | 16.8 |
| TIMET(ours) | **52.5** | **51.3** | **13.3** | **53.7** | **53.0** | **20.5** | **68.6** | **66.8** | **15.8** | **59.8** | **61.5** | **31.8** |

Table 4: **Unsupervised video semantic segmentation results** in mIoU for clustering (K=GT) and overclustering. TIMET not only gets better numbers on YTVOS, but also achieves considerably better results on DAVIS. This shows the better quality of the learned features at transferability to other video datasets. The clip-length is set to 16 and 4 for DAVIS and YTVOS, respectively. For clustering, the Hungarian algorithm [35] matches the unsupervised segmentation clusters (K) with the ground truth (GT) per frame (*F*), clip (*C*) or across the whole dataset (*D*). For overclustering, we use K=10 for the frame-wise (*F*) and clip-wise (*C*) evaluations, and for dataset-wise (*D*) evaluation, K=200 for DAVIS and K=500 for YTVOS. *denotes finetuning pretrained models on the same video frames as input as our method. The matching protocol is greedy many-to-one, see Appendix B for details.

### 4.3. Large-scale experiments

**Unsupervised video semantic segmentation.** In this section, we train our TIMET method on the YTVOS dataset and evaluate it for unsupervised video object segmentation on both DAVIS and YTVOS. The results are shown in Table 4. It should be noted that the number of prototypes we analyzed in Table 1a is not the number of clusters used to report accuracies for Table 4, as the actual number is usually unknown in the case of unsupervised learning. For the clustering experiment, we set K to the number of ground truth objects, which makes evaluation metrics easier to interpret and is commonly done in image clustering and segmentation. This evaluation is repeated but for larger numbers of clusters, in the "overclustering" scenario. Note that in dense self-supervised learning, overclustering can be particularly important because the learned representations are

used as a feature extractor for downstream tasks such as semantic segmentation or object detection. By using more fine-grained representations, the network may be able to extract more discriminative features such as object parts [74], which can lead to better performances.

From Table 4, we observe a clear trend: our method, trained on YTVOS not only achieves superior performances on YTVOS, but also beats existing image-trained models on DAVIS by a large margin. In particular, the state-of-the-art self-supervised clustering method, Leopart, has 4% lower performance on per dataset DAVIS evaluation for K=GT and 1.6% lower performance on YTVOS. The gap even becomes larger when Leopart is trained on the same video dataset with 8% to 16% lower numbers in per dataset and per clip numbers across different datasets. This means that when grouping objects of the same class over time and

the whole dataset, the image-based self-supervised methods have trouble generalising to videos, where objects are not centered in the frame, can appear in varying poses, and are more difficult in general. A strong contender is Leopart, which matches the performance of DINO for per frame and per clip clustering, however, improves considerably with the per dataset clustering. We attribute this to the fact that their dense learning objective improves mainly the generalisation capabilities of the learned representation, however falling short of generalising to temporal variations. We make similar observations for overclustering, improving the DINO baseline which we utilize as initialisation by 12% to 14% across different datasets in the per dataset metric.

We conclude that the proposed method outperforms all other methods in learning robust and discriminative representations for dense video understanding, and image-based self-supervised learning may not have sufficient generalisation capacity. This shows that if used right, time is a particularly strong inductive bias for self-supervised learning. Figure 3 shows segmentations returned by the proposed model trained on YTVOS and tested on YTVOS, DAVIS, and VISOR respectively. The proposed method groups objects accurately, and importantly, the frame segmentations are considerably consistent over time (see more examples in Appendix D). This highlights the importance and relevance of temporal fine-tuning, not just for higher accuracy, but crucially for consistency and robustness.

**Transfer from video training to images.** Despite the common belief that features learned from videos perform worse when transferred to images [34, 20], the results in Table 5 demonstrate our method achieves high performance that match state-of-the-art methods directly trained on images, specifically for the K=500, FCN, and LC metrics. We also compare our performance gains in videos against models designed for unsupervised image segmentation. The results show that models highly biased towards image datasets cannot provide the same performance when trained on videos, lagging behind our approach by 7% to 30%. Our results demonstrate that achieving high transfer performance on challenging tasks through video-based self-supervised learning is not only feasible but can also maintain high performance across modalities. These findings suggest that our method can drive further advances in self-supervised learning and inspire new directions for research in this field. Figure 3 shows the qualitative results on Pascal VOC. For more visualisations we refer to Appendix D.

**Salient object segmentation.** In Table 6, we compare foreground masks obtained with various DINO ViT based methods. We use the cluster-based foreground extraction protocol from [74] (details provided in Appendix C). First, we find that our method outperforms the original DINO at-

| | Pascal VOC | | | |
|---|---|---|---|---|
| | K=21 | K=500 | LC | FCN |
| *Trained on Images* | | | | |
| ResNet-50 | 4.5 | 36.5 | 53.8 | - |
| DINO [9] | 5.5 | 17.4 | 50.6 | 60.6 |
| SwAV [8] | 11.6 | 35.7 | 50.7 | - |
| MaskContrast [57] | 35.0 | 45.4 | 49.2 | - |
| DenseCL [61] | - | 43.6 | 49.0 | 69.4 |
| STEGO [21] | 7.0 | 19.5 | 59.1 | 63.5 |
| CrOC [52] | 20.6 | - | 61.6 | - |
| Leopart [74] | **36.6** | 50.5 | **68.0** | 70.1 |
| *Trained on Videos* | | | | |
| STEGO* | 4.0 | 15.5 | 51.1 | 55.5 |
| Leopart* | 14.9 | 21.2 | 53.2 | 63.2 |
| Flowdino† [70] | - | - | 59.4 | - |
| TIMET (ours) | 34.5 | **53.2** | **68.0** | **70.6** |

Table 5: **Transfer from video training to images.** Numbers taken from [74, 52, 70]. *: finetuning pretrained models on the same video frames as input as our method. †: uses a 40% larger superset of our train data. Details on the different datasets and methods are provided in the Appendix.

| | Pascal VOC | DAVIS | YTVOS |
|---|---|---|---|
| DINO [9] | 52.1 | 34.5 | 32.1 |
| Leopart [74] | 59.6 | 37.3 | 38.6 |
| STEGO [21] | 49.1 | 30.4 | 32.1 |
| TIMET (ours) | **63.9** | **44.5** | **43.5** |

Table 6: **Salient object segmentation.** We report performance using the Jaccard score [9] and use the official pretrained models for evaluation.

tention maps consistently by 10-11%. We also surpass the results of Leopart [74] and STEGO [21], works that rely on the same pretrained backbone. Even for the evaluation on Pascal VOC, our method achieves higher performances despite the domain shift of having trained on videos.

**Generalisation to egocentric datasets.** We train our method on EPIC-Kitchens-100 [13] and evaluate on VISOR [14]. We use the official code to convert VISOR to a DAVIS-like structure, in which we can report our numbers for per frame and per clip evaluation. Note that, after conversion, the object IDs do not maintain global consistency in the whole dataset; therefore, we cannot report the per dataset number. As Table 7 shows, TIMET outperforms image-based competitors on egocentric datasets as well.

|          | VISOR | |
|----------|------|------|
|          | *F* | *C* |
| DINO [9] | 24.8 | 18.7 |
| Leopart [74] | 24.1 | 18.5 |
| TIMET (ours) | **26.5** | **21.5** |

Table 7: **Generalisation to egocentric datasets.** Unsupervised semantic segmentation results in mIoU for clustering (K=GT) on VISOR [14] after training onn EPIC-Kitchens-100 [13]. The clip-length is set to 4. Our method gets better results on egocentric datasets as well.

**Visual In-Context Learning evaluation.** Here, we contrast our approach with a recently introduced benchmark, which assesses the in-context reasoning capabilities of models within the vision domain. Unlike linear or FCN classification methods, visual in-context learning evaluation obviates the need for fine-tuning or end-to-end training. Instead, it constructs validation segmentation maps using patch-level, feature-wise nearest neighbor similarities between validation images (referred to as "queries") and training samples (termed "keys"). This approach mirrors strategies in the NLP domain, aiming to evaluate the proficiency of models in learning tasks from limited examples presented as demonstrations. The results are shown by Table 8. Given that the vast majority of extant models in the domain utilize ViT-S16 as their backbone, we conducted a re-evaluation of their checkpoints to furnish a consistent and directly comparable evaluation table. Subsequently, we re-trained our model employing ViT-B16 and benchmarked it against other baselines as presented by [5]. The results, as depicted in the table, indicate that even though our model was exclusively trained on videos, it registers performance metrics in line with Leopart [74]. Furthermore, it surpasses the results of the leading method, CrOC [52], which was trained on images. Contrary to Hummingbird [5], TIMET is not custom-fitted to the specific evaluation setup and boasts superior computational efficiency. Notably, it requires only a single GPU for training, in contrast to the 16 TPUs demanded by Hummingbird.

## 5. Conclusion

This paper has aimed to learn dense representations that are learned from videos; yet can be generalised to the image domain as well. As video content is growing rapidly and contains more information compared to images, learning generalisable knowledge from them facilitates further scaling of self-supervised learning methods. To this effect, we have proposed a self-supervised clustering loss to encourage temporally consistent features between different frames of a video. To efficiently find corresponding views between

|          | Encoder | Params | mIoU |
|----------|---------|--------|------|
| *Trained on Images* | | | |
| Supervised | ViT-S16 | 21M | 35.1 |
| MoCo-v3* [11] | ViT-S16 | 21M | 19.5 |
| DINO* [9] | ViT-S16 | 21M | 47.9 |
| CrOC* [52] | ViT-S16 | 21M | 50.0 |
| Leopart* [74] | ViT-S16 | 21M | 63.6 |
| DINO [9] | ViT-B16 | 86M | 55.9 |
| MoCo-v3 [11] | ViT-B16 | 86M | 37.2 |
| MAE [22] | ViT-B16 | 86M | 6.6 |
| LOCA [7] | ViT-B16 | 86M | 57.5 |
| Hummingbird [5] | ViT-B16 | 86M | **70.5** |
| *Trained on Videos* | | | |
| TIMET (ours) | ViT-S16 | 21M | 61.6 |
| TIMET (ours) | ViT-B16 | 86M | 65.5 |

Table 8: **Visual In-Context Learning evaluation.** Numbers are taken from [5]. *: numbers are produced by this paper. This appeoach is equivalent to a non-parametric nearest-neighbor based evaluation on Pascal VOC. The details of the evaluation benchmark can be found in the provided implementation codes.

clip frames, we have proposed to recycle pretrained transformer features to leverage their natural tracking ability at each clip. Our empirical results indicate that this method achieves significant gains on the challenging task of video object segmentation across three different evaluation protocols and three datasets. Moreover, by transferring the learned model to the image domain, we have demonstrated the generalisability of the learned features by surpassing or matching the state-of-the-art for unsupervised image segmentation.

## 6. Acknowledgement

## References

[1] Nikita Araslanov, Simone Schaub-Meyer, and Stefan Roth. Dense unsupervised learning for video segmentation. *Advances in Neural Information Processing Systems*, 34:25308–25319, 2021. 2, 3, 15

[2] Yuki M. Asano, Christian Rupprecht, and Andrea Vedaldi. Self-labelling via simultaneous clustering and representation learning. In *International Conference on Learning Representations (ICLR)*, 2020. 4

[3] Mahmoud Assran, Mathilde Caron, Ishan Misra, Piotr Bojanowski, Florian Bordes, Pascal Vincent, Armand Joulin, Mike Rabbat, and Nicolas Ballas. Masked siamese networks for label-efficient learning. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXI*, pages 456–473. Springer, 2022. 6

[4] Arthur Aubret, Markus R. Ernst, Céline Teulière, and Jochen Triesch. Time to augment self-supervised visual representation learning. In *The Eleventh International Conference on Learning Representations*, 2023. 2

[5] Ivana Balažević, David Steiner, Nikhil Parthasarathy, Relja Arandjelović, and Olivier J Hénaff. Towards in-context scene understanding. *arXiv preprint arXiv:2306.01667*, 2023. 9

[6] Thomas Brox and Jitendra Malik. Large displacement optical flow: Descriptor matching in variational motion estimation. *IEEE TPAMI*, 33(3):500–513, 2011. 4

[7] Mathilde Caron, Neil Houlsby, and Cordelia Schmid. Location-aware self-supervised transformers. *arXiv preprint arXiv:2212.02400*, 2022. 9

[8] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020. 4, 7, 8

[9] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9650–9660, 2021. 3, 4, 5, 6, 7, 8, 9, 15

[10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020. 2

[11] Xinlei Chen*, Saining Xie*, and Kaiming He. An empirical study of training self-supervised vision transformers. *ICCV*, 2021. 9

[12] Marco Cuturi. Sinkhorn distances: Lightspeed computation of optimal transport. In *Advances in Neural Information Processing Systems*, volume 26, 2013. 4

[13] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, , Antonino Furnari, Jian Ma, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Rescaling egocentric vision: Collection, pipeline and challenges for epic-kitchens-100. *International Journal of Computer Vision (IJCV)*, 130:33–55, 2022. 5, 8, 9

[14] Ahmad Darkhalil, Dandan Shan, Bin Zhu, Jian Ma, Amlan Kar, Richard Higgins, Sanja Fidler, David Fouhey, and Dima Damen. Epic-kitchens visor benchmark: Video segmentations and object relations. In *Proceedings of the Neural Information Processing Systems (NeurIPS) Track on Datasets and Benchmarks*, 2022. 5, 8, 9

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 13

[16] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 5, 13, 14

[17] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3299–3309, 2021. 2, 4

[18] Basura Fernando, Hakan Bilen, Efstratios Gavves, and Stephen Gould. Self-supervised video representation learning with odd-one-out networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3636–3645, 2017. 1, 2

[19] Raghudeep Gadde, Varun Jampani, and Peter V Gehler. Semantic video cnns through representation warping. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4453–4462, 2017. 1

[20] Daniel Gordon, Kiana Ehsani, Dieter Fox, and Ali Farhadi. Watching the world go by: Representation learning from unlabeled videos. *arXiv preprint arXiv:2003.07990*, 2020. 2, 8

[21] Mark Hamilton, Zhoutong Zhang, Bharath Hariharan, Noah Snavely, and William T Freeman. Unsupervised semantic segmentation by distilling feature correspondences. *arXiv preprint arXiv:2203.08414*, 2022. 5, 8

[22] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022. 2, 9

[23] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020. 2

[24] Olivier J Hénaff, Skanda Koppula, Jean-Baptiste Alayrac, Aaron Van den Oord, Oriol Vinyals, and João Carreira. Efficient visual pretraining with contrastive detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10086–10096, 2021. 1, 2

[25] Olivier J Hénaff, Skanda Koppula, Evan Shelhamer, Daniel Zoran, Andrew Jaegle, Andrew Zisserman, João Carreira, and Relja Arandjelović. Object discovery and representation networks. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXVII*, pages 123–143. Springer, 2022. 2

[26] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint arXiv:1606.08415*, 2016. 13

[27] Jefferson Hernandez, Ruben Villegas, and Vicente Ordonez. Visual representation learning from unlabeled video using contrastive masked autoencoders. *arXiv preprint arXiv:2303.12001*, 2023. 2

[28] Ping Hu, Fabian Caba, Oliver Wang, Zhe Lin, Stan Sclaroff, and Federico Perazzi. Temporally distributed networks

for fast video semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8818–8827, 2020. 1

[29] Junhwa Hur and Stefan Roth. Joint optical flow and temporally consistent semantic segmentation. In *Computer Vision–ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part I 14*, pages 163–177. Springer, 2016. 1

[30] Allan Jabri, Andrew Owens, and Alexei Efros. Space-time correspondence as a contrastive random walk. *Advances in neural information processing systems*, 33:19545–19560, 2020. 1, 3, 4

[31] Samvit Jain, Xin Wang, and Joseph E Gonzalez. Accel: A corrective fusion network for efficient semantic segmentation on video. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8866–8875, 2019. 1

[32] Xu Ji, Joao F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9865–9874, 2019. 5

[33] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547, 2019. 13

[34] Thomas Kipf, Gamaleldin F Elsayed, Aravindh Mahendran, Austin Stone, Sara Sabour, Georg Heigold, Rico Jonschkowski, Alexey Dosovitskiy, and Klaus Greff. Conditional object-centric learning from video. *arXiv preprint arXiv:2111.12594*, 2021. 2, 8

[35] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 5, 7

[36] Xueting Li, Sifei Liu, Shalini De Mello, Xiaolong Wang, Jan Kautz, and Ming-Hsuan Yang. Joint-task self-supervised learning for temporal correspondence. *Advances in Neural Information Processing Systems*, 32, 2019. 3

[37] Xiangtai Li, Wenwei Zhang, Jiangmiao Pang, Kai Chen, Guangliang Cheng, Yunhai Tong, and Chen Change Loy. Video k-net: A simple, strong, and unified baseline for video segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18847–18857, 2022. 1

[38] Yule Li, Jianping Shi, and Dahua Lin. Low-latency video semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5997–6005, 2018. 1

[39] Runtao Liu, Zhirong Wu, Stella Yu, and Stephen Lin. The emergence of objectness: Learning zero-shot segmentation from videos. *Advances in Neural Information Processing Systems*, 34:13137–13152, 2021. 2

[40] Si Liu, Changhu Wang, Ruihe Qian, Han Yu, Renda Bao, and Yao Sun. Surveillance video parsing with single frame supervision. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 413–421, 2017. 1

[41] Yifan Liu, Chunhua Shen, Changqian Yu, and Jingdong Wang. Efficient semantic video segmentation with per-frame inference. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X 16*, pages 352–368. Springer, 2020. 1

[42] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 13

[43] Xiankai Lu, Wenguan Wang, Jianbing Shen, Yu-Wing Tai, David J Crandall, and Steven CH Hoi. Learning video object segmentation from unlabeled videos. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 8960–8970, 2020. 2

[44] Yuval Nirkin, Lior Wolf, and Tal Hassner. Hyperseg: Patch-wise hypernetwork for real-time semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4061–4070, 2021. 1

[45] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023. 4

[46] Nikhil Parthasarathy, SM Eslami, João Carreira, and Olivier J Hénaff. Self-supervised video pretraining yields strong image representations. *arXiv preprint arXiv:2210.06433*, 2022. 2

[47] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Advances in neural information processing systems*, 32, 2019. 13

[48] Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alex Sorkine-Hornung, and Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv preprint arXiv:1704.00675*, 2017. 5, 13

[49] Kanchana Ranasinghe, Muzammal Naseer, Salman Khan, Fahad Shahbaz Khan, and Michael S Ryoo. Self-supervised video transformer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2874–2884, 2022. 1

[50] Sucheng Ren, Wenxi Liu, Yongtuo Liu, Haoxin Chen, Guoqiang Han, and Shengfeng He. Reciprocal transformations for unsupervised video object segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 15455–15464, 2021. 2

[51] Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Dominik Zietlow, Tianjun Xiao, Carl-Johann Simon-Gabriel, Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox, et al. Bridging the gap to real-world object-centric learning. *arXiv preprint arXiv:2209.14860*, 2022. 1

[52] Thomas Stegmüller, Tim Lebailly, Behzad Bozorgtabar, Tinne Tuytelaars, and Jean-Philippe Thiran. Croc: Cross-view online clustering for dense visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7000–7009, 2023. 8, 9

[53] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *Proceedings of the IEEE/CVF conference on*

*computer vision and pattern recognition*, pages 8934–8943, 2018. 15

[54] Guolei Sun, Yun Liu, Henghui Ding, Thomas Probst, and Luc Van Gool. Coarse-to-fine feature mining for video semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3126–3137, 2022. 1

[55] Fida Mohammad Thoker, Hazel Doughty, Piyush Bagad, and Cees GM Snoek. How severe is benchmark-sensitivity in video self-supervised learning? In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXIV*, pages 632–652. Springer, 2022. 2

[56] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, Marc Proesmans, and Luc Van Gool. Scan: Learning to classify images without labels. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part X*, pages 268–285. Springer, 2020. 5

[57] Wouter Van Gansbeke, Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Unsupervised semantic segmentation by contrasting object mask proposals. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10052–10062, 2021. 1, 2, 5, 8, 13

[58] Junke Wang, Dongdong Chen, Zuxuan Wu, Chong Luo, Chuanxin Tang, Xiyang Dai, Yucheng Zhao, Yujia Xie, Lu Yuan, and Yu-Gang Jiang. Look before you match: Instance understanding matters in video object segmentation. *arXiv preprint arXiv:2212.06826*, 2022. 1

[59] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra. Cut and learn for unsupervised object detection and instance segmentation. *arXiv preprint arXiv:2301.11320*, 2023. 2

[60] Xiaolong Wang, Allan Jabri, and Alexei A Efros. Learning correspondence from the cycle-consistency of time. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2566–2576, 2019. 3

[61] Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3024–3033, 2021. 2, 5, 8, 13

[62] Zhenda Xie, Yutong Lin, Zheng Zhang, Yue Cao, Stephen Lin, and Han Hu. Propagate yourself: Exploring pixel-level consistency for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16684–16693, 2021. 2

[63] Dejing Xu, Jun Xiao, Zhou Zhao, Jian Shao, Di Xie, and Yueting Zhuang. Self-supervised spatiotemporal learning via video clip order prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10334–10343, 2019. 1

[64] Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint arXiv:1809.03327*, 2018. 5, 13

[65] Yu-Syuan Xu, Tsu-Jui Fu, Hsuan-Kung Yang, and Chun-Yi Lee. Dynamic video segmentation network. In *Proceed-*

*ings of the IEEE conference on computer vision and pattern recognition*, pages 6556–6565, 2018. 1

[66] Charig Yang, Hala Lamdouar, Erika Lu, Andrew Zisserman, and Weidi Xie. Self-supervised video object segmentation by motion grouping. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7177–7188, 2021. 2, 3, 15

[67] Sukmin Yun, Hankook Lee, Jaehyung Kim, and Jinwoo Shin. Patch-level representation learning for self-supervised vision transformers. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8354–8363, 2022. 2

[68] Andrii Zadaianchuk, Matthaeus Kleindessner, Yi Zhu, Francesco Locatello, and Thomas Brox. Unsupervised semantic segmentation with self-supervised object-centric representations. *arXiv preprint arXiv:2207.05027*, 2022. 2

[69] Shaofeng Zhang, Feng Zhu, Rui Zhao, and Junchi Yan. Patch-level contrasting without patch correspondence for accurate and dense contrastive representation learning. In *The Eleventh International Conference on Learning Representations*, 2023. 2

[70] Xinyu Zhang and Abdeslam Boularias. Optical flow boosts unsupervised localization and segmentation. *arXiv preprint arXiv:2307.13640*, 2023. 8

[71] Yurong Zhang, Liulei Li, Wenguan Wang, Rong Xie, Li Song, and Wenjun Zhang. Boosting video object segmentation via space-time correspondence learning. *arXiv preprint arXiv:2304.06211*, 2023. 1, 5

[72] Jinghao Zhou, Chen Wei, Huiyu Wang, Wei Shen, Cihang Xie, Alan Yuille, and Tao Kong. ibot: Image bert pre-training with online tokenizer. *ICLR*, 2022. 3, 6

[73] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. Deep feature flow for video recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2349–2358, 2017. 1

[74] Adrian Ziegler and Yuki M Asano. Self-supervised learning of object parts for semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14502–14511, 2022. 1, 2, 5, 7, 8, 9, 13, 14, 15, 16, 18

# Appendix

## A. Implementation Details

Code is provided in the supplementary materials and will be open-sourced.

**Training and evaluation datasets.**

**Video datasets.** *DAVIS17* [48] is designed for video object segmentation, comprising 150 videos, with 60 allocated for training, 30 for validation, and 60 for testing. Only the first frames of the test set have ground truth foreground masks, so the validation set is used for evaluation. *YTVOS* [64] is another dataset for video object segmentation and is significantly larger than DAVIS17. It consists of 4,453 videos that are annotated with 65 object categories. As with DAVIS17, ground truth masks are only available for the first frames of the test and validation sets, and therefore, a fixed 20% of the training set is randomly sampled for the evaluation phase, details are provided in the supplementary material. Additionally, meta information is utilized to ensure objects in the same category have the same class id throughout the dataset for semantic, category-level assessments. Figure 4 shows the distribution of objects in YTVOS.

**Image datasets.** *Pascal VOC 2012* [16] is an object recognition dataset with 20 object categories and one background class. It includes pixel-level segmentation, bounding box, and object class annotations for each image, and has been extensively used as a benchmark for object detection, semantic segmentation, and classification tasks. The dataset is split into three subsets, with 1,464 images allocated for training, 1,449 for validation, and a private testing set. As the dataset has been commonly used as a main reference for recent works in dense self-supervised image segmentation [74, 57, 61], we also use its validation set as one of the evaluation datasets.

**Model training.** We use batches of size 128 on 1 NVIDIA GeForce RTX 3090, and the optimizer is AdamW [42] with learning rate equal to 1e-4 for the projection head and the backbone's learning rate is 1e-5. We freeze the backbone model except for the last two blocks for fine-tuning. Our model is implemented in torch [47]. We use Faiss [33] for K-Means clustering. We chose to train a ViT-Small [15] image because it has roughly the same number of parameters as a ResNet-50 (21M vs. 23M). The projection head learning rate is 1e-4 and the backbone's learning rate is 1e-5. The projection head consists of three linear layers with hidden dimensionality of 2048 and Gaussian error linear units
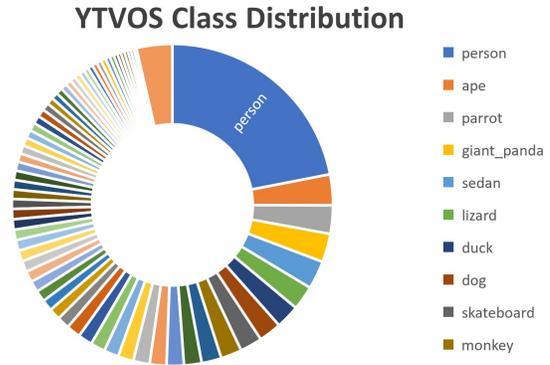


Figure 4: The distribution of classes in YTVOS. Some of the more dominant classes are labeled.

```
# mm : torch.mm
# exp : torch.exp
# bmm : torch.bmm
# Reshape : In-place operation to change the
    input shape
# Normalize : torch.Normalize
# F[i] : Shows the ith feature map
# C-Map[i] : Shows the ith cluster map

prev_feat = []    # (nmb-context, dim, h*w)
prev_maps = []
For i in range(N-1):
    prev_feat.append(F[i])
    prev_maps.append(C-Map[i])
src_feat = Stack(prev_feat)
trgt_feat = F[N] # (1, dim, h*w)
trgt_feat= Normalize(trgt_feat, dim=1, p=2)
src_feat = Normalize(src_feat, dim=1, p=2)
aff = exp(bmm(trgt_feat, src_feat) / 0.1)
Reshape(aff, (nmb-context * h*w, h*w))
aff = aff / sum(aff, keepdim=True, axis=0)
aff = mask-neighborhood(aff)
prev_maps = Stack(prev_maps)    # (nmb-context, C,
    h, w)
Reshape(prev_maps, (C, nmb-context*h*w))
trgt_cmap = mm(prev_maps, aff)
```

Listing 1: **FF component**. The pytorch implementation of **FF** is shown.

as activation function [26]. We set the temperature to 0.1 and use Adam as an optimizer with a cosine weight decay schedule. The augmentations used are random color-jitter, Gaussian blur, grayscale, and random cropping.

**Evaluation details.** Since we evaluate the pre-GAP *layer4* features or the spatial tokens, their output resolution does not match the mask resolution. To fix that, we bilinearly interpolate before applying the linear head; or directly interpolate the clustering results by nearest neighbor upsampling. For a fair comparison between ResNets and ViTs,

**Algorithm 1** Evaluation Pipeline Pseudocode. $M$ is the model, $C$ is the clustering algorithm, MA is the matching algorithm by which the clusters are scored, and GT is the given ground-truth.

---

1: input = input.reshape(bs * $n\_f$, c, h, w)
2: $F_b$ = M(input)
3: $F_b$ = F.reshape(bs, $n\_f$, num-patch, dim)
4: score_list = []
5: **if** Per frame **then**
6:     **for** $F_c$ In $F_b$ **do**
7:         **for** $F_f$ In $F_c$ **do**
8:             C_Map = $C(F_f)$
9:             score = $\mathrm{MA}(\mathrm{C\_Map}, \mathrm{GT}_f)$
10:             score_list.append(score)
11:         **end for**
12:     **end for**
13: **else if** Per clip **then**
14:     **for** $F_c$ In $F_{all}$ **do**
15:         C_Maps = $C(F_c)$
16:         score = $\mathrm{MA}(\mathrm{C\_Maps}, \mathrm{GT}_c)$
17:         score_list.append(score)
18:     **end for**
19: **else if** Per dataset **then**
20:     C_Maps = $C(F_b)$
21:     score = $\mathrm{MA}(\mathrm{C\_Maps}, \mathrm{GT}_b)$
22:     score_list.append(score)
23: **end if**
24: **return**(score_list.mean())

---

we use dilated convolution in the last bottleneck layer of the ResNet such that the spatial resolution of both network architectures match (28x28 for 448x448 input images). All overclustering results were computed using downsampled 100x100 masks to speed up the Hungarian matching as we found that the results do not differ from using full-resolution masks.

## B. Evaluation Protocol for Unsupervised Video Object Semantic Segmentation

Here, we provide details for the evaluation protocols for unsupervised video multi-label object segmentation. To be consistent with the image domain [74], a clustering algorithm is applied to the features extracted from frozen encoders to craft dense assignment maps of pseudo-labels. To produce scores, based on each evaluation protocol, the crafted maps are matched with the ground truth, and their MIOU is reported. Suppose the matching algorithm is specified by $M$(labels, ground-truth), clustering algorithm by $K$, dataset features by $F \in R^{N \times n\_f \times d \times h \times w}$, where $N$, $n\_f$, $d$, $h$, and $w$ stand for dataset size, number of frames per clip, feature dimension, and feature spatial resolutions, re-

spectively. we introduce three evaluation protocols that are specific to the video domain.

**Per frame evaluation (*F*).**

$$F_{\mathrm{frame}}[i,j] = F[i,j] \tag{10}$$

$$\text{score} = \frac{1}{N \times n\_f} \sum_{i=1}^{N} \sum_{j=1}^{n\_f} \mathrm{MIOU}(M(K(F_{\mathrm{frame}}[i,j]), \mathrm{GT}[i,j])) \tag{11}$$

this measures a basic alignment of a given feature map with the ground-truth.

**Per clip evaluation (*C*).**

$$F_{\mathrm{clip}}[i] = (F[i,1], \cdots, F[i, nf]) \tag{12}$$

$$\text{score} = \frac{1}{N} \sum_{i=1}^{N} \mathrm{MIOU}(M(K(F_{\mathrm{clip}}[i]), \mathrm{GT}[i])) \tag{13}$$

This evaluation tests whether the assigned pseudo-labels remain consistent over time for each clip.

**Per dataset evaluation (*D*).**

$$F_{\mathrm{dataset}} = (F[1,1], \cdots, F[N, n\_f]) \tag{14}$$

$$\text{score} = \mathrm{MIOU}(M(K(F_{\mathrm{dataset}}), \mathrm{GT})) \tag{15}$$

This evaluation measures the most difficult ability of generating not only temporally stable features of objects across time but across videos.

## C. Additional Experiments

To provide a complete evaluation of our method compared to the baseline on Pascal VOC [16], we show the per-class performance in Figure 6. As the figure shows, we improve the class "person" by more than 40%, which could be beacuse of the high number of such objects in YTVOS as Figure 4 shows. The classes "cat" and "dog" also show a significant improvement since they are of the further dominant classes after "person".

**Unsupervised video object segmentation and tracking.**
Although such methods are designed for salient object detection or unsupervised mask propagation, and not specifically for unsupervised semantic segmentation, we evaluate their performance on our proposed evaluation protocols to highlight their strengths and limitations. The results are shown in Table 9. As it is shown, TIMET outperforms such methods on all the evaluation protocols by a margin between 12% to 24%. Not being specifically designed for semantic segmentation tasks may explain their inferior performance.

| Method | $F$ | $C$ | $D$ |
|---|---|---|---|
| DUL [1] | 28.2 | 27.4 | 2.4 |
| Motion Grp [66] | 32.0 | 30.7 | 1.5 |
| TIMET (ours) | **56.5** | **55.5** | **14.1** |

Table 9: **Comparison to video unsupervised object segmentation methods.** Evaluation on DAVIS with K=GT.

**Comparing to unsupervised video object segmentation and tracking.** Although such methods are designed for salient object detection or unsupervised mask propagation, and not specifically for unsupervised semantic segmentation, we evaluate their performance on our proposed benchmark to highlight their strengths and limitations. We conducted a comprehensive comparison of our method with state-of-the-art techniques, such as Motion-Grp [66] and DUL [1], which aim to learn unsupervised features for propagating a given first frame's mask in test time or separating foreground from background using motion flows, as previously mentioned. To ensure consistency, we trained and tested all models on DAVIS, a widely used benchmark. Results in Table 9 demonstrate that our method outperforms these techniques across all reported metrics. It is worth noting that these methods were not specifically designed for semantic segmentation tasks, which may explain their inferior performance.

**Analyzing per class results.** In Figure 6 and 5 the per class improvements on Pascal VOC are reported. As the figures show, for the classes that make over 95% of the number of objects existing in YTVOS, the performance almost always improves considerably. The reason that the class "bird" does not behave the same as the others might be due to the small size of this object in YTVOS.

**Component Contributions.** In Table 10, we show the effect of using different components on the Pascal clustering results. As it is shown, TIMETimproves DINO [9] by 12%. The results are improved by another 18% after applying CBFE [74] to the features, showing high overclustering performance on this dataset. The number of clusters used for CBFE is 300 for this experiment.

**Comparing different propagators.** Complementary to the previous ablations, in Table 11, we conduct a detailed study on the performance of forwarding foreground masks using various temporal intervals and comparing the use of optical flow [53] with our Feature-Forwarder (FF). We find that across all values of $\delta t$, FF outperforms flow by a large margin of +10% or more. This superiority has the added benefit of our FF module not having to expensively com-

| | MIOU |
|---|---|
| K=150 | 48.2 |
| DINO | 4.6 |
| +TimeT | 16.5 |
| +CBFE | 34.5 |

Table 10: **Component contributions.** We show the gains that each individual component brings for Pascal VOC segmentation and K=21.

pute flow but instead reusing the activations that are processed for the clustering step.

| $\delta t$ | $GT_0$ | Flow | FF |
|---|---|---|---|
| 0.1s | 26.4 | 29.8 | 39.8 |
| 0.2s | 25.7 | 28.5 | **40.3** |
| 0.4s | 24.5 | 26.1 | 40.1 |
| 0.8s | 21.8 | 22.7 | 37.5 |

Table 11: Comparing FF with other forwarding methods. The numbers of reported on DAVIS validation set.

## D. Additional Visualisations

**Clustering with K=GT.** In Figure 8, we show some qualitative results on Pascal VOC when K is set to the number of ground-truth objects, similar to Figure 3 in the paper. While the method is trained on YTVOS with a temporal loss, it obtains strong performances on an image segmentation dataset yielding high class consistencies (indicated by the segmentation colors) and tight borders.

We also have provided further qualitative results with the same setting on DAVIS and YTVOS in Figure 7 and the attached HTML file. As the videos show, we get considerably more consistent and structured visualizations compared to DINO [9] and Leopart [74]. This further supports the effectiveness of temporal fine-tuning compared to the models solely trained on images.

**Overclustering results by merging 500 clusters using ground-truth labels.** In the attached HTML file, we show the visualizations of our method on Pascal in the overslutering setting as well. As depicted, objects from different classes can be segmented precisely with different colors, showing that the learned patch features are semantic.
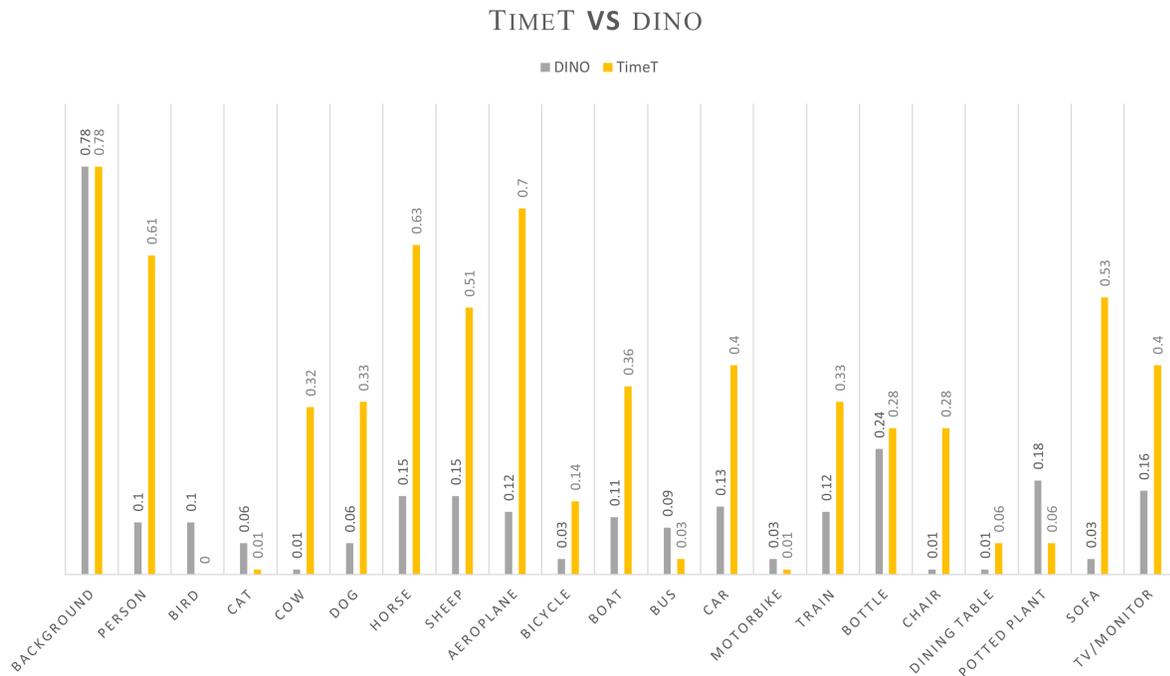
Figure 5: The per class performance of DINO and TIMET is shown for the clustering experiment with K=GT. Cluster-based foreground extraction [74] has been applied to both methods. As it is seen, this paper almost always improves the baseline performance for this evaluation as well. Pascal VOC is used in this experiment.
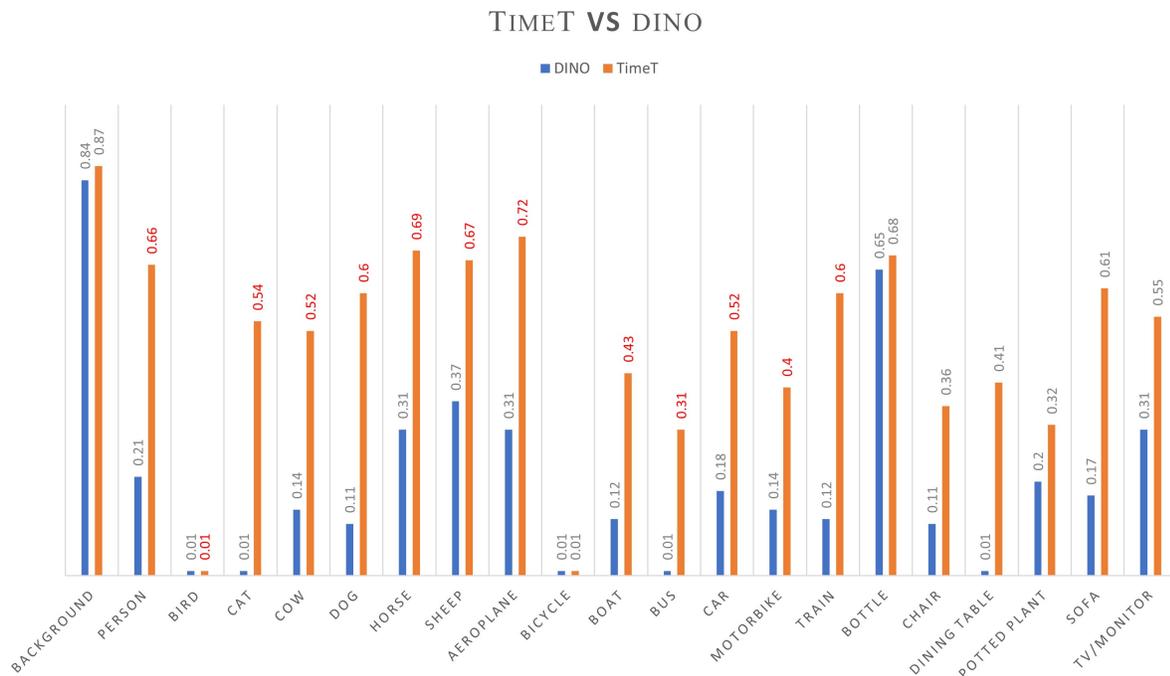


Figure 6: The per class performance of DINO and TIMET is shown for the overclustering experiment with K=500. As it is seen, this paper consistently improves the baseline performance. The numbers for the dominant shared classes between YTVOS and Pascal VOC are shown by the color red.
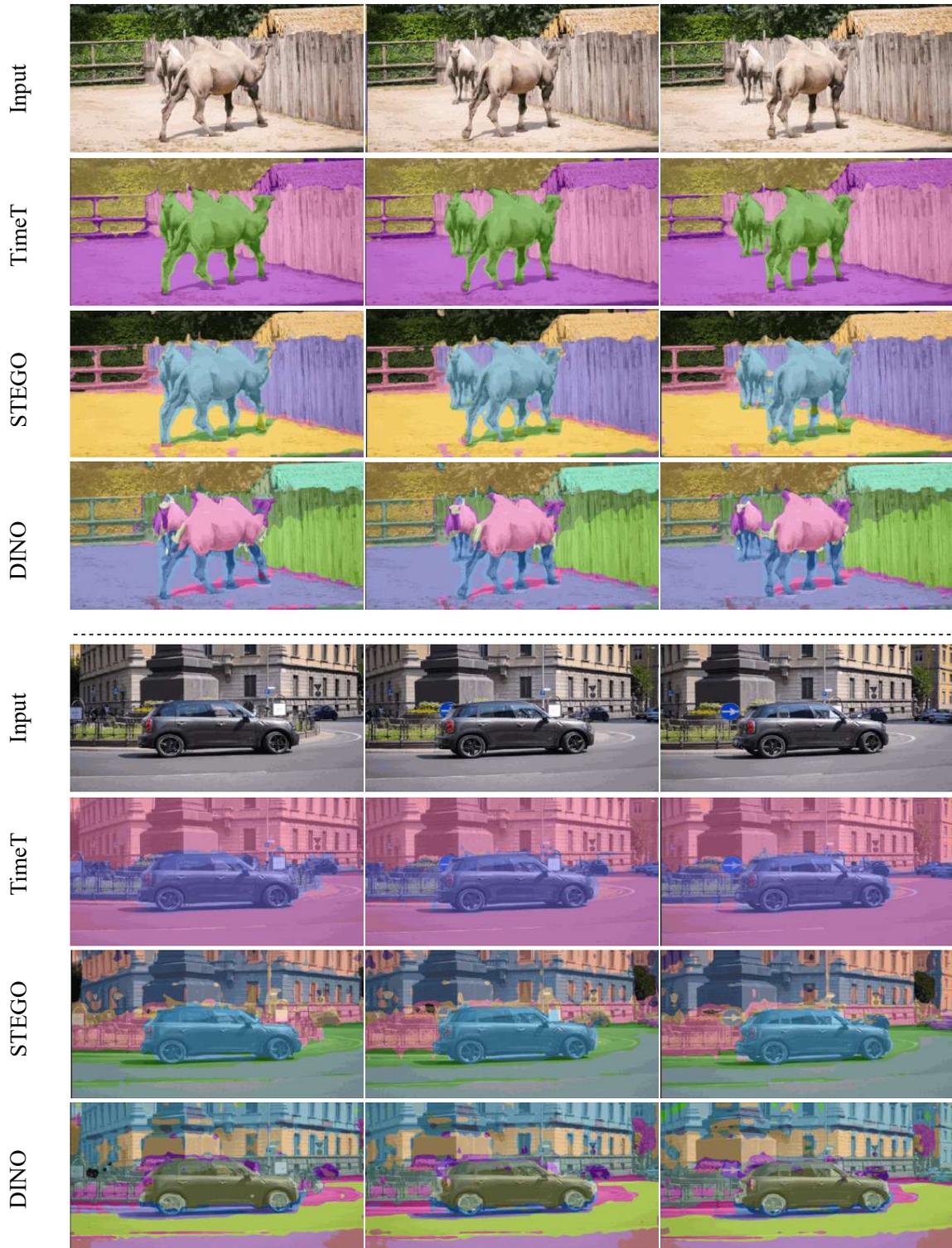
Figure 7: **TIMET segmentations on DAVIS with K=GT.** Here, we compare the performance of DINO, STEGO, and TIMET on the task of unsupervised video semantic segmentations. TIMET has a clear advantage over both DINO and STEGO in terms of providing tight segmentation boundaries and specifying different objects with different category IDs. Different colors in the figure specify different IDs.

Figure 8: **TIMET segmentations on Pascal VOC with K=21.** We use CBFE [74] to focus on the foreground objects. While the method is trained on YTVOS with a temporal loss, it obtains strong performances on an image segmentation dataset yielding high class consistencies (indicated by the segmentation colors) and tight borders.
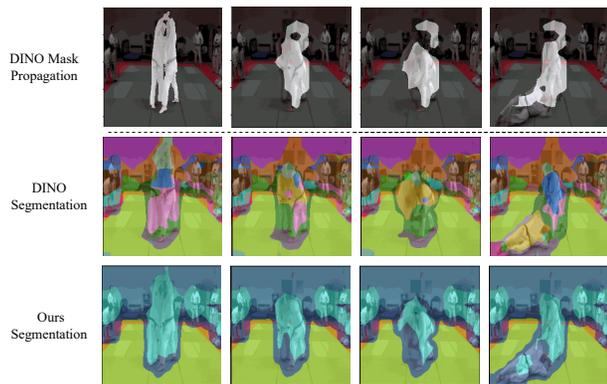
Figure 9: DINO's features jump around across time, leading to inconsistent cluster maps. Our proposed TIMET-trained model observes more temporal consistency.