

# EMERGING TRACKING FROM VIDEO DIFFUSION

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

We find video diffusion models, renowned for their generative capabilities, surprisingly excel at pixel-level object tracking without any explicit training for this task. We introduce a simple and effective method to extract motion representations from video diffusion models, achieving state-of-the-art tracking results. Our approach enables the tracking of identical objects, overcoming limitations of previous methods reliant on intra-frame appearance correspondence. Visualizations and empirical results show that our approach outperforms recent self-supervised tracking methods, including the state-of-the-art, by up to 6 points. Our work demonstrates video generative models can learn intrinsic temporal dynamics of video, and excel in tracking tasks beyond original video synthesis.

## 1 INTRODUCTION

*“What I cannot create, I do not understand.”*

– Richard Feynman

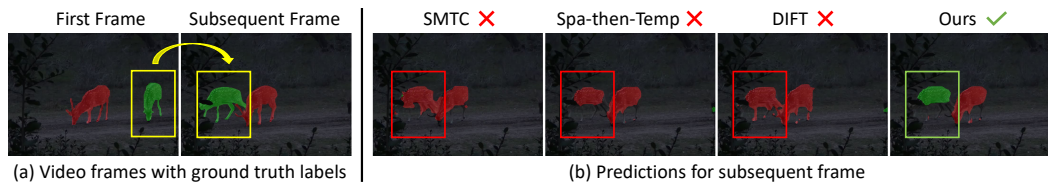


Figure 1: **Predictions from video label propagation task.** State-of-the-art models fail to find the correct temporal correspondence when multiple objects look similar in a video, such as SMTC (Qian et al., 2023), Spa-then-Temp (Li & Liu, 2023), and DIFT (Tang et al., 2023). For instance, the deer with green segmentation map labels in (a) are mislabeled as red by existing models, as highlighted by the red boxes in (b). By introducing latent representations from pretrained video diffusion models, our method captures temporal motions and correctly identifies the deer, highlighted by the green box in (b). Our work significantly improves tracking performance across various scenarios.

The ability of temporal relational reasoning over time (Yi et al., 2019) is crucial for visual intelligence. Rather than performing simple appearance correspondence, people often rely on temporal relational reasoning to track moving objects in complex situations (Yi et al., 2019; Gerstenberg et al., 2015; Ullman, 2015). For example, given the two moving deer in Figure 1(a), we can easily reason and track different deer even after they change their relative positions.

Learning video representations for temporal correspondence is essential for tasks like video object segmentation (Caron et al., 2021). Appearance-based correspondence methods have been used for tracking (Wang et al., 2021; Hu et al., 2022), including the recent state-of-the-art DIFT (Tang et al., 2023) that uses latent representations from image diffusion models (Rombach et al., 2021; Dhariwal & Nichol, 2021). Some research also integrates temporal information in model training (Wang et al., 2019; Jabri et al., 2020). However, existing methods often have low accuracy because they fail to capture temporal context in complex scenarios, see Figure 1(b), where state-of-the-art models (Qian et al., 2023; Li & Liu, 2023; Tang et al., 2023) fail to differentiate between two deer.

In this paper, we demonstrate that representations from video diffusion models can improve tracking across various scenarios, including those with multiple objects of similar appearance. Video diffusion

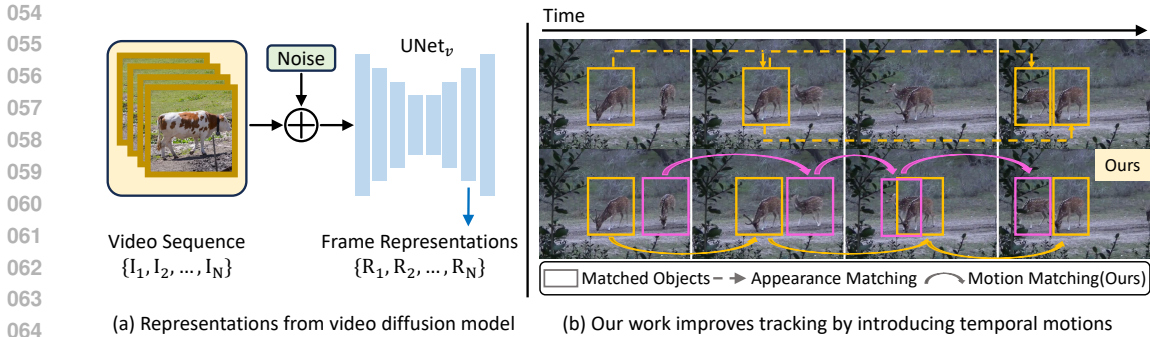


Figure 2: **Framework.** Our work focuses on the video label propagation task, which uses frame representations to transfer the first frame’s label to subsequent frames. We term the representation for the  $j$ th frame  $I_j$  as  $R_j$ . Unlike existing methods that often extract  $R_j$  by a 2D image model, we introduce the 3D UNet backbone from video diffusion models, which includes a temporal axis and processes the entire video sequence as input (see (a)). Our approach improves tracking by integrating temporal motions as shown in (b), where different colors indicate different matching pairs. In (b), the first row shows traditional appearance matching, which relies on visual similarity across frames and may misidentify objects, such as incorrectly matching two deer in the last frame to the same deer in the first frame. In contrast, our work (second row) captures motion patterns among frames, resulting in more accurate tracking. We term our **Temporal Enhanced Diffusion** tracking method as **TED**. Experiments demonstrate that our TED improves tracking performance across diverse video scenarios, including those with similar-looking objects. (**Best viewed when zoomed in.**)

models, trained to generate consistent videos across frames, capture both the object appearance and the temporal relationships between objects. We show that without additional training, the internal layer outputs of UNet from video diffusion models introduce the temporal reasoning capability that aids tracking in complex situations. For example, as shown in the last column of Figure 1(b), our video diffusion representations successfully track two deer even as they change their positions relative to each other in the video. We term our **Temporal Enhanced Diffusion** tracking method as **TED**.

Experimental results show that our TED method outperforms 23 popular baseline models, achieving state-of-the-art performance in self-supervised pixel-level object tracking. On the DAVIS dataset for semi-supervised video object segmentation, our TED significantly outperforms SFC (Hu et al., 2022) by 6.4%, SMTC (Qian et al., 2023) by 4.6%, Spa-then-Temp (Li & Liu, 2023) by 3.5%, and DIFT (Tang et al., 2023) by 1.9%. Furthermore, we introduce a challenging task of tracking similar-looking objects, and a new real-world dataset for evaluation, termed YouTube-Similar. Benefiting from the temporal reasoning ability, our TED improves upon DIFT (Tang et al., 2023) by 5.3%. Moreover, our approach achieves state-of-the-art results in human pose tracking. Our work is the first to show that temporal motions learned from video diffusion models can solve perception challenges and significantly improve perception performance. We will release our code and data.

## 2 RELATED WORK

Learning temporal correspondence is crucial for visual tracking (Tao et al., 2016; Xu & Wang, 2021; Li & Liu, 2023). Due to limited annotations, prior studies have developed methods to learn correspondence in a self-supervised manner (Caron et al., 2021; Qian et al., 2023). Our work contributes to this field of self-supervised correspondence, and we discuss related work as follows.

**Temporal correspondence learned from images.** Self-supervised learning that trains on image datasets has achieved great success in downstream tasks, including temporal correspondence. Pioneering work, such as MoCo (He et al., 2020) and DINO (Caron et al., 2021), adopt instance discrimination as pretext task which learns similar representations for different augmentations of the same image. DenseCL (Wang et al., 2021) and PixPro (Xie et al., 2021b) further apply contrastive learning to pixel-level, which improve dense prediction tasks. SFC (Hu et al., 2022) boosts performance on temporal correspondence further by fusing image-level and pixel-level representations. Recently, DIFT (Tang et al., 2023) achieves state-of-the-art results in temporal correspondence task

by leveraging internal representations from image diffusion models (Rombach et al., 2021). These methods learn intra-frame information and rely on appearance for pixel-level tracking. Our work highlights the limitations of using appearance alone for temporal correspondence and significantly improves tracking by introducing temporal reasoning capabilities.

**Temporal correspondence learned from videos.** Temporal information in videos provides supervision signals to learn video representations during training. Two widely used pretext tasks for model training are frame reconstruction and cycle-consistency over time. Frame reconstruction tasks involve reconstructing a frame from adjacent frames (Vondrick et al., 2018; Lai & Xie, 2019; Li et al., 2019; Lai et al., 2020), while cycle-consistency tasks track a patch backwards and forward in time to align start and end points (Wang et al., 2019; Jabri et al., 2020). However, these methods often overlook spatial features crucial for creating discriminative and robust representations (Li & Liu, 2023). Recent research integrates spatial with temporal information in model training, such as Spa-then-Temp (Li & Liu, 2023) and SMTC (Qian et al., 2023). Despite incorporating temporal information during model training, our work reveals that existing methods still face challenges in complex scenarios, such as tracking multiple similar-looking objects, as shown in Figure 1. By introducing temporal reasoning ability from video diffusion models to tracking, our approach significantly improves performance across various video scenarios, including those involving similar-looking objects.

**Video diffusion models.** Diffusion models have significantly advanced image and video generation (Ho et al., 2020; Saharia et al., 2022; Ho et al., 2022; Ruiz et al., 2023). Text-to-image diffusion models (Nichol et al., 2021; Ramesh et al., 2022) allow precise control over generated image content via text prompts, with Stable Diffusion (Rombach et al., 2021) improving generation efficiency and quality by performing diffusion process in latent space. To generate videos with consistent frames, video diffusion models are created by inserting temporal blocks into image diffusion models, which are then trained on video datasets (Blattmann et al., 2023b; Zhang et al., 2023). Representative video diffusion models include Sora (Brooks et al., 2024), ModelScope (Wang et al., 2023), I2VGen-XL (Zhang et al., 2023), and Stable Video Diffusion (Blattmann et al., 2023a). Our work is the first to demonstrate that temporal dynamics learned by video diffusion models can significantly improve tracking performance. Our work highlights the potential of video generative models in tracking tasks beyond their original use in video synthesis.

### 3 METHOD

We focus on the video label propagation task and first introduce the background in Section 3.1. We then discuss the challenges faced by previous methods in tracking identical objects in Section 3.2. In Section 3.3, we show how our approach addresses these challenges and improves tracking performance by leveraging temporal context. Our implementation details are provided in Section 3.4.

#### 3.1 BACKGROUND

**Video label propagation task** aims to transfer ground truth labels, such as segmentation maps, from the first frame to subsequent frames (Vondrick et al., 2018), as shown in Figure 3. The key is training models to represent frames and establish pixel-level mapping among them (Hu et al., 2022). Due to limited annotations, prior work trains the models in a self-supervised manner with various pretext tasks (Jabri et al., 2020; Li & Liu, 2023). DIFT (Tang et al., 2023) significantly improves tracking performance using latent representations from image diffusion models. We first introduce diffusion models and then discuss how DIFT uses them for tracking.

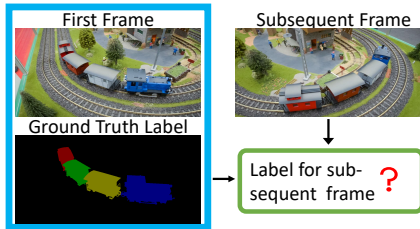


Figure 3: **Video label propagation task** transfers the ground truth label of the first frame to subsequent frames.

**Diffusion models** have achieved unprecedented success in generating images and videos with rich content (Rombach et al., 2022; Brooks et al., 2024). They are probabilistic models that learn the data distribution  $p(\mathbf{x})$  and generate  $\mathbf{x}$  from a random Gaussian variable (Nichol et al., 2021), where  $\mathbf{x}$  is the image for image diffusion models.

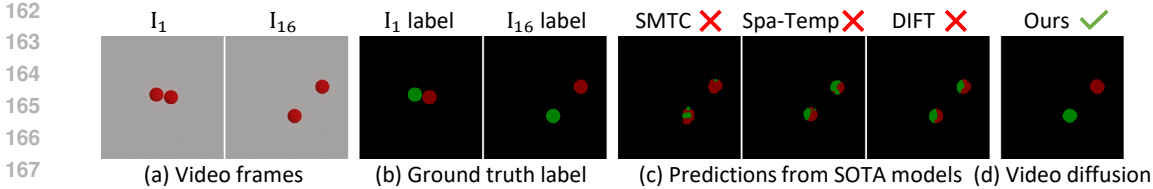


Figure 4: **Video diffusion representations enable tracking objects with identical appearances.** We conduct a controlled study, that we perform object label propagation on videos featuring two independently moving and identical-looking balls, with frames and their ground truth labels shown in (a) and (b). State-of-the-art methods (Qian et al., 2023; Li & Liu, 2023; Tang et al., 2023) fail to distinguish the two balls, leading to incorrect predictions (c). In contrast, our video diffusion representations accurately track both balls despite their identical appearance, as shown in (d).

Diffusion models learn rich visual concepts by recovering signals from corrupted data  $\mathbf{x}_\tau$  at varying noise levels (Choi et al., 2022), with loss defined in Eqn. 1:

$$L = \mathbb{E}_{\mathbf{x}, \epsilon \sim \mathcal{N}(0,1), \tau} \left[ \|\epsilon - \epsilon_\theta(\mathbf{x}_\tau, t)\|_2^2 \right] \quad (1)$$

where  $\epsilon$  is the actual noise corrupting the clean data and  $\epsilon_\theta(\mathbf{x}_\tau, t)$  is the noise predicted by the denoising model  $\epsilon_\theta$ . UNet (Ronneberger et al., 2015) is commonly used as the denoising model  $\epsilon_\theta$ .

Noisy  $\mathbf{x}_\tau$  is generated by adding noise from a Gaussian distribution  $\mathcal{N}(0, 1)$  to the clean data  $\mathbf{x}_0$  according to the noise scheduler  $\alpha_t$  (Ho et al., 2020), defined as:

$$\mathbf{x}_\tau = \sqrt{\alpha_\tau} \mathbf{x}_0 + \sqrt{1 - \alpha_\tau} \epsilon, \quad \epsilon \sim \mathcal{N}(0, 1) \quad (2)$$

Here,  $\tau$  represents the timestep in diffusion process, with larger  $\tau$  indicating higher noise levels.

**Tracking by image diffusion representations.** DIFT (Tang et al., 2023) improves video propagation performance using latent representations from image diffusion models (Rombach et al., 2021; Dhariwal & Nichol, 2021). It leverages outputs from internal layers of UNet backbone, defined as:

$$\mathbf{R} = \text{UNet}(\mathbf{x}_\tau, n) \quad (3)$$

where  $n$  is the layer index. Specifically,  $\mathbf{R}_i = \text{UNet}_i(\mathbf{x}_\tau, n)$ , where subscript  $i$  indicates image diffusion models. The input  $\mathbf{x}_\tau$  is generated using Eqn. 2 at a chosen timestep  $\tau$ . Since  $\text{UNet}_i$  processes a single image at a time, DIFT treats video frames as independent images and extracts  $\mathbf{R}_i$  for each frame with a single forward pass through the UNet model.

### 3.2 CHALLENGES FOR TRACKING IDENTICAL OBJECTS

Prior studies have achieved impressive results in video label propagation by establishing pixel-level mappings among frames based on frame representations (Jabri et al., 2020; Li & Liu, 2023). For videos with a single object, pixel-level mapping often relies on object appearances, such as the semantic information used in SFC (Hu et al., 2022). However, in videos with multiple similar-looking objects, like tanks with similar fish, establishing accurate correspondence remains challenging and is underexplored in the video label propagation task.

**Controlled toy example.** We begin with a controlled toy example that tracks two independently moving, identical-looking balls in a video, as shown in Figure 4(a). We use the Kubric simulator (Greff et al., 2022) to create a video dataset with random ball sizes and motions, termed Kubric-Similar. In this dataset, we propagate the segmentation map of each ball from the first frame (see Figure 4(b)) to the subsequent frames. We follow the label propagation procedures in prior studies (Jabri et al., 2020; Li & Liu, 2023), with implementation details in Section 3.4.

We evaluate state-of-the-art models on Kubric-Similar, with results reported in Figure 4(c) and Table 1. Figure 4(c) shows that existing methods struggle with object identity, leading to poor tracking. This aligns with Table 1, where many methods, including DIFT (Tang et al., 2023) that uses image diffusion representations, achieve a  $J_m$  around 50%. Note that a  $J_m$  around 50% indicates performance no better than random guessing due to the identical size of the two balls. These findings highlight the difficulty of tracking multiple similar-looking objects in temporal correspondence tasks.



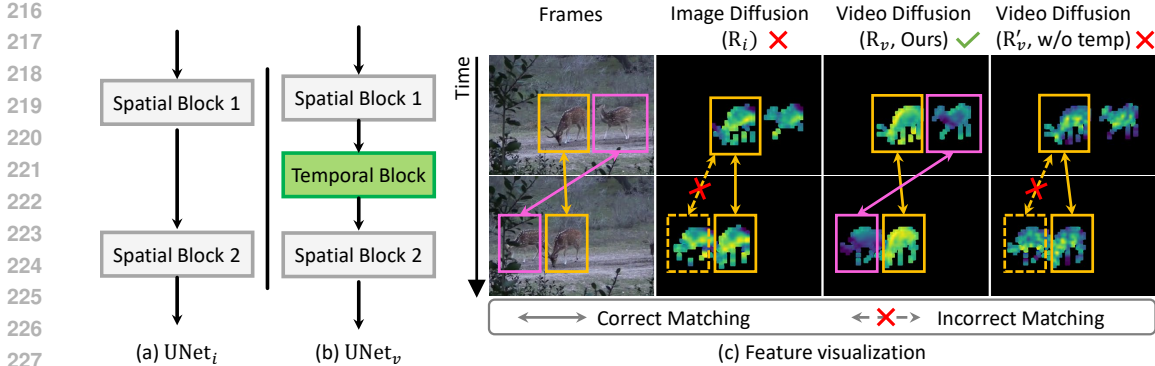


Figure 5: **Track by temporal context.** To understand why video diffusion representation ( $\mathbf{R}_v$ ) excels in tracking similar-looking objects, we compare the UNet backbones of video and image diffusion models. (a) UNet<sub>i</sub> from image diffusion models consists of spatial layer blocks that process each image independently. (b) UNet<sub>v</sub> is constructed by inserting temporal layer blocks to UNet<sub>i</sub> to ensure frame consistency. In (c), we perform principal component analysis (PCA) on the representations from different frames of each model, such as  $\mathbf{R}_v^s, \mathbf{R}_v^t = \text{PCA}(\mathbf{R}_v^s \parallel \mathbf{R}_v^t)$  for  $\mathbf{R}_v$ , where  $s$  and  $t$  represent different frames. The results reveal that image diffusion representation ( $\mathbf{R}_i$ ) from DIFT (Tang et al., 2023) learns similar features for both deer, leading to incorrect matching. In contrast, our  $\mathbf{R}_v$  learns distinguishing features that achieve correct matching. Removing temporal layers from UNet<sub>v</sub> results in losing its distinguishing capability, shown in  $\mathbf{R}'_v$ (w/o temp). By integrating information across frames,  $\mathbf{R}_v$  enhances tracking by incorporating temporal context, outperforming  $\mathbf{R}_i$  from DIFT (Tang et al., 2023), which is limited to intra-frame information and appearance-based tracking.

**Video diffusion representations achieve significant improvement in tracking identical objects.** To improve label propagation for identical objects, we replace the image diffusion representations ( $\mathbf{R}_i$ ) in DIFT (Tang et al., 2023) with video diffusion representations ( $\mathbf{R}_v$ ). Specifically,  $\mathbf{R}_v$  is obtained by applying UNet<sub>v</sub> from video diffusion models following Eqn. 3. Thus,  $\mathbf{R}_v = \text{UNet}_v(\mathbf{x}_\tau, n)$ , where  $\mathbf{x}_\tau$  represents the video sequence of multiple frames. The process of obtaining  $\mathbf{R}_v$  is illustrated in Figure 2(a), with additional implementation details in Section 3.4. Figure 4(d) shows that our  $\mathbf{R}_v$  accurately tracks both balls, despite their identical appearance, outperforming existing methods. Table 1 further confirms our approach’s effectiveness, outperforming DIFT (Tang et al., 2023) by 38.3% in  $\mathcal{J}\&\mathcal{F}_m$ .

Table 1: **Results in tracking identical objects.** We perform object label propagation on videos featuring two independently moving, identical-looking balls, as shown in Figure 4(a). Our video diffusion representations achieve state-of-the-art results in tracking identical objects. Colors of the numbers highlight the **best** results.

Model	$\mathcal{J}\&\mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$
MOCO (He et al., 2020)	51.9	61.8	56.8
SimSiam (Chen & He, 2021)	52.8	63.9	58.3
TimeCycle (Wang et al., 2019)	42.6	55.6	49.1
UVC (Li et al., 2019)	58.0	52.6	43.9
CRW (Jabri et al., 2020)	49.5	59.7	54.6
SFC (Hu et al., 2022)	41.8	51.1	46.5
SMTc (Qian et al., 2023)	72.6	68.5	76.7
Spa-then-Temp (Li & Liu, 2023)	44.5	39.8	49.2
DIFT <sub>sd</sub> (Tang et al., 2023)	46.3	43.4	49.3
DIFT <sub>adm</sub> (Tang et al., 2023)	52.6	50.3	54.8
<b>Video diffusion (<math>\mathbf{R}_v</math>, ours)</b>	<b>90.9</b>	<b>87.2</b>	<b>94.5</b>

### 3.3 TRACK OBJECTS BY TEMPORAL CONTEXT

Motivated by the success of video diffusion representations when tracking identical objects in the toy example above, we will first investigate where the tracking capability comes from. We will then capitalize on the findings and propose a simple and effective method for better tracking.

**Where does the ability to track similar-looking objects come from?** We hypothesize that this tracking capability stems from temporal context learned during video synthesis. Video diffusion models (Blattmann et al., 2023a) insert temporal layers into image diffusion models to learn temporal dynamics such as motion, ensuring frame consistency in generated videos (see Figure 5 (a)(b)). We denote the UNet representations from video and image diffusion models as  $\mathbf{R}_v$  and  $\mathbf{R}_i$ , respectively.

Table 2: **Results for pixel-level object tracking.** We evaluate our TED method on semi-supervised video object segmentation, and compare it with 24 baseline models, including self-supervised and supervised approaches. Our TED achieves state-of-the-art tracking performance on both the DAVIS and Youtube-Similar datasets, outperforming recent methods by up to 6%. We visualize the tracking results in Figure 6. These results demonstrate the effectiveness of our method in object tracking, even when multiple objects have similar appearances. Colored numbers indicate the **best** results. TED refers to our default setting using  $\mathbf{R}_f$ , while TED<sup>†</sup> denotes the setting using  $\mathbf{R}_v$ .

Super-vised	Method	Dataset	DAVIS			Youtube-Similar		
			$\mathcal{J} \& \mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$	$\mathcal{J} \& \mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$
	InstDis (Wu et al., 2018)		66.4	63.9	68.9	-	-	-
	MoCo (He et al., 2020)		65.9	63.4	68.4	48.0	48.5	47.4
	SimCLR (Chen et al., 2020)		66.9	64.4	69.4	37.5	36.9	38.1
	BYOL (Grill et al., 2020)	ImageNet w/o labels	66.5	64.0	69.0	47.1	47.7	46.5
	SimSiam (Chen & He, 2021)		67.2	64.8	68.8	47.4	47.9	47.0
	DINO (Caron et al., 2021)		71.4	67.9	74.9	63.9	63.0	64.7
	DetCo (Xie et al., 2021a)		65.7	63.3	68.1	41.4	42.0	40.9
	DenseCL (Wang et al., 2021)		61.4	60.0	62.9	46.4	46.5	46.3
	PixPro (Xie et al., 2021b)		57.5	56.6	58.3	45.6	45.9	45.3
	DIFT <sup>adm</sup> (Tang et al., 2023)		75.7	72.7	78.6	60.7	59.8	61.7
	DIFT <sup>sd</sup> (Tang et al., 2023)	LAION	70.0	67.4	72.5	56.3	55.8	56.7
	Colorization (Vondrick et al., 2018)		34.0	34.6	32.7	-	-	-
	VINCE (Gordon et al., 2020)		65.2	62.5	67.8	44.9	45.4	44.3
	VFS (Xu & Wang, 2021)	Kinetic	68.9	66.5	71.3	57.3	57.1	57.5
	UVC (Li et al., 2019)		60.9	59.3	62.7	49.7	49.8	49.7
	CRW (Jabri et al., 2020)		67.6	64.8	70.2	52.0	52.3	51.6
	CorrFlow (Lai & Xie, 2019)	OxUvA	50.3	48.4	52.2	39.6	40.0	39.3
	TimeCycle (Wang et al., 2019)	VLOG	48.7	46.4	50.0	39.8	41.3	38.2
	MAST (Lai et al., 2020)	YT-VOS	65.5	63.3	67.6	-	-	-
	SMTC (Qian et al., 2023)		73.0	69.4	76.6	57.5	57.2	57.9
	SFC (Hu et al., 2022)	ImageNet, YT-VOS	71.2	68.3	74.0	55.5	55.3	55.7
	Spa-then-Temp (Li & Liu, 2023)		74.1	71.1	77.1	59.6	59.2	60.1
	TED <sup>†</sup> (Ours, $\mathbf{R}_v$ )	Web-Vid	66.3	63.4	69.1	62.0	61.5	62.5
	TED (Ours, $\mathbf{R}_f$ )	ImageNet, Web-Vid	<b>77.6</b>	<b>74.4</b>	<b>80.8</b>	<b>66.0</b>	<b>65.1</b>	<b>67.0</b>
✓	OSVOS (Caelles et al., 2017)	ImageNet, DAVIS	60.3	56.6	63.9	-	-	-
	OnAVOS (Voigtlaender & Leibe, 2017)	ImageNet, DAVIS	65.4	61.6	69.1	-	-	-
	CFBI+ (Yang et al., 2020)	YT-VOS, DAVIS	82.8	80.1	85.5	-	-	-

We examine the properties of  $\mathbf{R}_v$  and  $\mathbf{R}_i$  using principal component analysis (PCA), as shown in Figure 5(c), where two moving deer change their relative positions over time. We will show that  $\mathbf{R}_v$  learns distinguishing features for two deer even if they have similar appearances, while  $\mathbf{R}_i$  learns similar features for both deer.

We perform PCA on pairs of frames for each model, such as  $\tilde{\mathbf{R}}_v^s, \tilde{\mathbf{R}}_v^t = \text{PCA}(\mathbf{R}_v^s \parallel \mathbf{R}_v^t)$  for  $\mathbf{R}_v$  where  $s$  and  $t$  represent different frames. Figure 5(c) shows that  $\mathbf{R}_i$  of DIFT (Tang et al., 2023) learn similar features for both deer, leading to incorrect matching. In contrast,  $\mathbf{R}_v$  learns distinguishing features for the two deer that enable correct matching. We then remove the temporal blocks from UNet<sub>v</sub> and recomputed  $\mathbf{R}_v$ , termed  $\mathbf{R}'_v$ . Interestingly, Figure 5(c) shows that  $\mathbf{R}'_v$  loses the distinguishing features between the two deer. Unlike  $\mathbf{R}_i$  which only uses intra-frame information, the temporal layers in UNet<sub>v</sub> (like temporal attention layers) enable  $\mathbf{R}_v$  to integrate information across multiple video frames, introducing temporal motion to tracking. We compare our temporal motion matching using  $\mathbf{R}_v$  to the appearance matching of  $\mathbf{R}_i$  in Figure 2(b). Our results and discussions demonstrate the superiority of  $\mathbf{R}_v$  in using temporal context for tracking.

**Using  $\mathbf{R}_v$  for better tracking.** Our investigations show that video diffusion representations ( $\mathbf{R}_v$ ) capture temporal context, crucial for tracking identical objects. Since temporal context is orthogonal to appearance information, it complements prior tracking methods like image diffusion ( $\mathbf{R}_i$ ). As shown in Eqn. 4, we employ a simple concatenation of the representations from video and image diffusion models in later experiments:

$$\mathbf{R}_f = \text{Concat}(\alpha \|\mathbf{R}_v\|_2, (1 - \alpha) \|\mathbf{R}_i\|_2) \quad (4)$$

where  $\|\cdot\|$  denotes L2 normalization and  $\alpha$  is a hyperparameter between 0 and 1. We term our Temporal Enhanced Diffusion tracking method as TED. We use  $\mathbf{R}_f$  by default and denote the setting that uses  $\mathbf{R}_v$  as TED<sup>†</sup> for distinction. We will show our TED achieves state-of-the-art tracking results.

### 3.4 IMPLEMENTATION DETAILS

**Video label propagation.** Our work follows prior studies (Jabri et al., 2020; Caron et al., 2021; Tang et al., 2023) for the evaluation protocol of label propagation, which includes representation extraction and label prediction stage, as shown in Algorithm 1. We first obtain frame representations  $R_f$  using video and image diffusion models. To predict the label of current frame, similar pixel pairs between current and previous frames are identified by computing the similarities of their representations. Each pixel in the current frame is then labeled by aggregating the labels of similar pixels from previous frames, weighted by their pixel similarity. More experimental setups are detailed in Appendix B.1.

**Appearance representations.** Following (Tang et al., 2023), we use the output from the internal layers of UNet<sub>i</sub> as the appearance representation  $\mathbf{R}_i$ , following  $\mathbf{R}_i = \text{UNet}_i(\mathbf{x}_\tau, n)$ .  $\mathbf{x}_\tau$  represents each video frame, generated according to Eqn. 2 with an empirically determined  $\tau$ . We process each frame through a single forward pass of UNet<sub>i</sub>. Our framework accommodates any pre-trained image diffusion model for  $\mathbf{R}_i$ , using ADM (Dhariwal & Nichol, 2021) by default. We also investigate other models such as Stable Diffusion (Rombach et al., 2022).

**Temporal representations.** We obtain  $\mathbf{R}_v$  following  $\mathbf{R}_v = \text{UNet}_v(\mathbf{x}_\tau, n)$  as shown in Fig 2(a). The key difference in obtaining  $\mathbf{R}_v$  compared to  $\mathbf{R}_i$  is using UNet<sub>v</sub> from video diffusion models, which process video sequence of multiple frames as  $\mathbf{x}_\tau$ . Since current video diffusion models accept limited frames as input, long videos are split into subsequences.  $\mathbf{R}_v$  is then obtained for each subsequence through a one-pass forward process in UNet<sub>v</sub>. Our framework supports any off-the-shelf pre-trained video diffusion model for  $\mathbf{R}_v$ , using I2VGen-XL (Zhang et al., 2023) by default. We also explore additional models like Stable Video Diffusion (Blattmann et al., 2023a).

## 4 EXPERIMENTS

### 4.1 EXPERIMENTAL SETUP

We evaluate our TED method on the video label propagation task, and compare it with 24 baseline models. Our work uses video representations from pre-trained diffusion models, and does not require additional training.

- **Pretrained self-supervised learning models.** We evaluate 9 self-supervised models pre-trained on ImageNet known for strong temporal correspondence performance: 6 instance discrimination models like MoCo (He et al., 2020) and 3 dense contrastive learning models, such as DenseCL (Wang et al., 2021).
- **Image diffusion model representations.** We compare with DIFT (Tang et al., 2023), which leverages representations from image diffusion models for temporal correspondence.
- **Task-specific models (self-supervised).** We include 11 self-supervised models tailored for temporal correspondence tasks, trained by pretext tasks like frame reconstruction (e.g., UVC (Li et al., 2019)), cycle consistency (e.g., CRW (Jabri et al., 2020)), and video contrastive learning (e.g., VFS (Xu & Wang, 2021)). We also include recent methods such as SMTC (Qian et al., 2023) and Spa-then-Temp (Li & Liu, 2023).
- **Task-specific models (supervised).** We compare our method with 3 supervised approaches that utilize labeled data during training, such as CFBI+ (Yang et al., 2020).

**Evaluation datasets.** We evaluate TED on the semi-supervised video object segmentation task, which propagates the object segmentation from the first frame to subsequent frames. We evaluate widely-used DAVIS-2017 (Pont-Tuset et al., 2017) which includes 30 videos from various scenarios. To test the tracking ability for similar-looking objects, we introduce the Youtube-Similar dataset, composed of 28 videos from Youtube-VOS (Xu et al., 2018) that feature multiple similar-looking objects. Following (Tang et al., 2023), we report region-based similarity ( $J_m$ ) and contour-based accuracy ( $F_m$ ). More dataset details are provided in the Appendix B.3.

### 4.2 EXPERIMENTAL RESULTS

**Quantitative results.** We compared our TED to 24 baseline models on the DAVIS and Youtube-Similar dataset, with results detailed in Table 2. Our TED achieves the **state-of-the-art** tracking

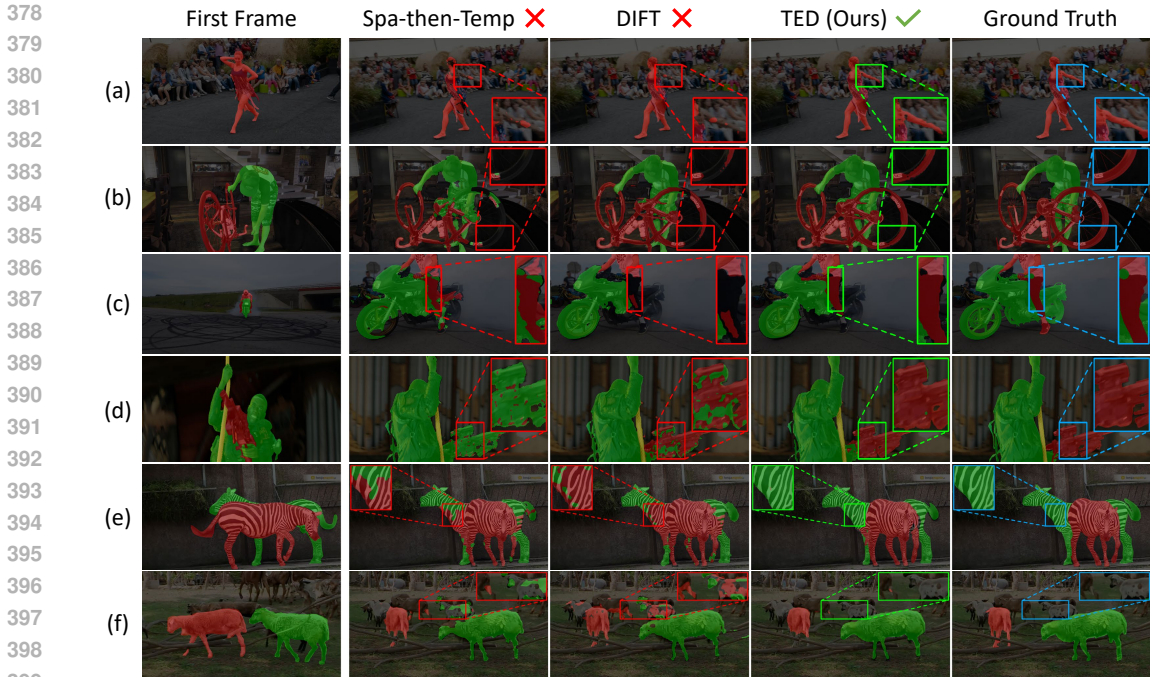


Figure 6: **Predictions for pixel-level object tracking.** We evaluate our TED method on semi-supervised video object segmentation, which propagates object segmentation maps from the first frame to subsequent frames. Our TED consistently outperforms state-of-the-art methods (Li & Liu, 2023; Tang et al., 2023) on the DAVIS (Figure a-d) and YouTube-Similar (Figure e-f) datasets, aligning with the results in Table 2. Notably, our TED delivers more accurate predictions in scenarios with complex deformations (a) and viewpoint changes (b), while Spa-then-Temp (Li & Liu, 2023) and DIFT (Tang et al., 2023) struggle with tracking completeness, such as the missing arm in (a). Our TED also excels in multi-object scenarios, delivering superior tracking for interacting objects (c-d) and similar-looking objects (e-f). In contrast, Spa-then-Temp (Li & Liu, 2023) and DIFT (Tang et al., 2023) have mislabeling issues, such as incorrect labels for the gun in (d) and misaligned labels for sheep in the background (f). These results show that our TED significantly improves tracking performance, highlighting the benefits of incorporating temporal reasoning into tracking. (**Best viewed when zoomed in.**)

performance on both datasets, surpassing recent methods by up to **6%**. Specifically, on the DAVIS dataset, our method outperforms SFC (Caron et al., 2021) by 6.4%, SMTc (Qian et al., 2023) by 4.6%, Spa-then-Temp (Li & Liu, 2023) by 3.5%, and DIFT (Tang et al., 2023) by 1.9%. On the Youtube-Similar dataset, our TED shows an even greater improvement, exceeding Spa-then-Temp (Li & Liu, 2023) by 6.4% and DIFT (Tang et al., 2023) by 5.3%. These improvements highlight the effectiveness of our method in object tracking, even for challenging settings with multiple similar-looking objects.

**Visualizations.** We present our tracking results alongside those from state-of-the-art methods in Figure 6, with results for DAVIS shown in Figure 6(a-d) and for YouTube-Similar in Figure 6(e-f). Our TED outperforms existing studies on both datasets, aligning with Table 2. Our TED effectively handles complex deformations (a) and viewpoint changes (b), outperforming Spa-then-Temp (Li & Liu, 2023) and DIFT (Tang et al., 2023), which struggle with tracking elements like the human arm in Figure 6(a). Additionally, our TED excels in multiple-object scenarios, such as interacting objects (c-d) and similar-looking objects (e-f), whereas Spa-then-Temp and DIFT often confuse different objects, leading to incorrect label assignments. For instance, in Figure 6(d), Spa-then-Temp (Li & Liu, 2023) incorrectly labels the gun as a human, and DIFT (Tang et al., 2023) shows significant errors in the predicted contour. In Figure 6 (f), featuring multiple sheep, both Spa-then-Temp (Li & Liu, 2023) and DIFT (Tang et al., 2023) mistakenly align the object label to the sheep in the background. Our TED consistently achieves more accurate tracking results across these scenarios, demonstrating significant performance improvements through enhanced temporal reasoning.



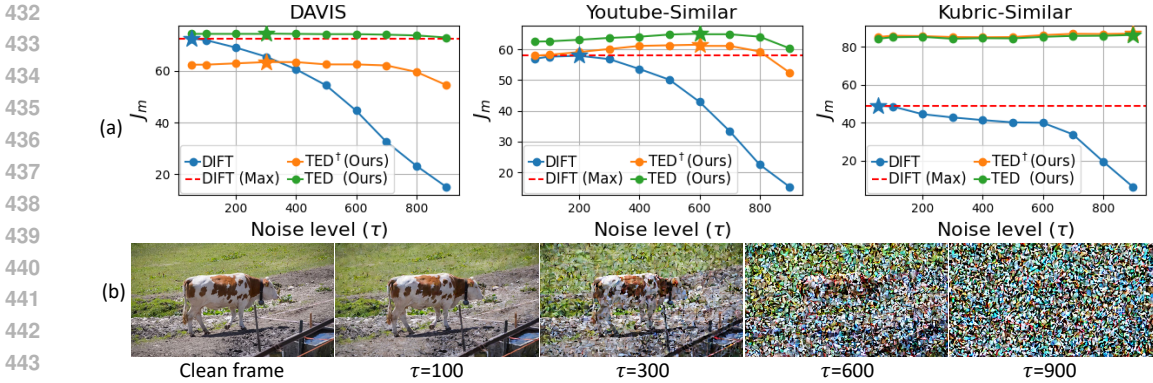


Figure 7: **How and why does  $t$  influence tracking.** We present tracking results in (a) using diffusion representations obtained at varying noise levels  $\tau$  (see Eqn. 2 and Eqn. 3), with higher  $\tau$  indicating more noise (b). TED uses combined  $\mathbf{R}_f$  defined in Eqn. 4, and TED<sup>†</sup> uses video diffusion representation  $\mathbf{R}_v$ . The best result for each method is marked with a star and the best result for DIFT (Tang et al., 2023) across all  $\tau$  is indicated as a red dashed line. Using image diffusion representations ( $\mathbf{R}_i$ ), DIFT peaks at low noise ( $\tau \leq 200$ ) and deteriorates as noise increases. This is due to its reliance on appearance for tracking, which becomes almost unavailable at high noises. In contrast, TED<sup>†</sup> (using  $\mathbf{R}_v$ ) excels at higher  $\tau$  values, peaking at  $\tau=600$  on Youtube-Similar and  $\tau=900$  on Kubric-Similar where the input video is heavily corrupted (b). **The high accuracy at high noise levels is because  $\mathbf{R}_v$  learns coarse-grained motions that enable tracking similar-looking objects, such as object positions. When the video input is less noisy, the diffusion model is trained to denoise appearance details, where motion feature may not be so prioritized, leading to performance decrease at low noise levels.** Our TED consistently outperforms DIFT (Tang et al., 2023) across various  $\tau$  values on all datasets, demonstrating the superiority of incorporating temporal information into tracking.

### 4.3 ABLATION STUDIES AND DISCUSSIONS

**How and why does  $\tau$  influence tracking.** We obtain frame representations from diffusion models as defined in Eqn. 3, with the UNet input  $x_\tau$  generated according to Eqn. 2. Following DIFT (Tang et al., 2023), we empirically determine the noise level  $t$  to produce  $x_\tau$ . We investigate the impact of noise level  $\tau$  on tracking performance in Figure 7(a), where a higher  $\tau$  indicates more noise (Figure 7(b)). In Figure 7(a), Kubric-Similar is a dataset featuring independently moving and identical-looking balls, defined in Section 3.2. We mark the best result for each method with a star. TED uses combined  $\mathbf{R}_f$  defined in Eqn. 4, and TED<sup>†</sup> uses video diffusion representation  $\mathbf{R}_v$ . Using image diffusion representations ( $\mathbf{R}_i$ ), DIFT (Tang et al., 2023) achieves the best result at low noise ( $\tau \leq 200$ ) and decreases rapidly as noise increases due to diminishing availability of appearance information. In contrast, our TED<sup>†</sup> with  $\mathbf{R}_v$  peaks at a higher  $t$  and maintains robust tracking over a much broader range of  $\tau$ . Notably, TED<sup>†</sup> reaches its best performance at  $\tau=600$  on Youtube-Similar and  $\tau=900$  on Kubric-Similar, where the input video is heavily corrupted and appearance information is almost unavailable as shown in Figure 7(b). These results suggest that  $\mathbf{R}_v$  encodes temporal motion that can be used for tracking at higher noise levels. Moreover, our TED with  $\mathbf{R}_f$  consistently outperforms DIFT (Tang et al., 2023) across a wide  $\tau$  range, demonstrating the effectiveness of our TED by integrating temporal dynamics into tracking.

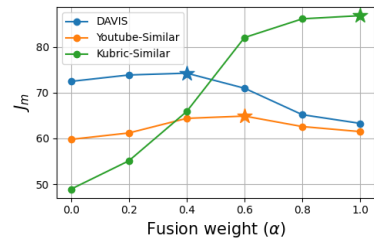


Figure 8: **Fusion weight  $\alpha$ .** Our TED outperforms DIFT (Tang et al., 2023) ( $\alpha=0.0$ ) on all datasets.

Diffusion models solve different tasks at different noise levels during training (Choi et al., 2022). When the video input is corrupted at high noise levels, video diffusion models are trained to solve the hard task that learning coarse-grained signals in the video, such as motion (like the change of object positions among frames). Therefore, its representation encodes rich motion information that enables

tracking similar-looking objects. When the video input is less noisy, the diffusion model is trained to denoise appearance details, where motion features may not be so prioritized, leading to performance decrease at low noise levels.

**Fusion weight  $\alpha$ .** To utilize both temporal motion and appearance for better tracking, our TED combines  $\mathbf{R}_v$  and  $\mathbf{R}_i$  into  $\mathbf{R}_f$  as defined in Eqn. 4. Figure 8 shows the tracking results with varying fusion weight  $\alpha$ , where a higher  $\alpha$  increases the contribution of  $\mathbf{R}_v$ .  $\mathbf{R}_f$  reduces to  $\mathbf{R}_i$  when  $\alpha=0.0$  and to  $\mathbf{R}_v$  when  $\alpha=1.0$ . We mark the best result on each dataset with a star in Figure 8. Our TED achieves the best results with medium  $\alpha$  values of 0.4 for DAVIS and 0.6 for Youtube-Similar, demonstrating that the integration of appearance and temporal information improves tracking performance. For Kubric-Similar, TED performs best with  $\alpha=1.0$ , reflecting the dataset’s unique characteristics of containing identical objects where appearance information from  $\mathbf{R}_i$  does not provide additional value for tracking. On all datasets, our TED consistently outperforms DIFT (Tang et al., 2023) ( $\alpha=0.0$ ), highlighting the advantage of our work by introducing temporal motions to tracking.

**Feature layers for video diffusion representations.** We use  $\mathbf{R}_v$  from internal layers of the UNet in video diffusion models for video label propagation, as illustrated in Figure 2 and Eqn. 3. Following (Tang et al., 2023), we use the decoder representations from UNet and report the tracking results of TED<sup>†</sup> on DAVIS using  $\mathbf{R}_v$  from different decoder blocks in Table 3. Table 3 shows that the medium block (block 2) yields the best performance among all blocks.

**Different diffusion models.** We evaluate the tracking results of TED using  $\mathbf{R}_f$  obtained from different video and image diffusion models on the DAVIS dataset, as shown in Table 4. We investigate video diffusion models like Stable Video Diffusion (SVD) (Blattmann et al., 2023a) and I2VGen-XL (I2V) (Zhang et al., 2023), image diffusion models like Stable Diffusion (SD) (Rombach et al., 2022) and ADM (Dhariwal & Nichol, 2021). Our TED achieves the best tracking performance when using video diffusion representations from I2VGen-XL (Zhang et al., 2023) and image diffusion representations from ADM (Dhariwal & Nichol, 2021), which is used as the default setting in the paper.

**Results on human pose tracking.** In addition to video object segmentation, we test our method on the JHMDB benchmark (Jhuang et al., 2013), which tracks 15 human pose keypoints in 268 videos. We follow the evaluation protocol of prior studies (Li et al., 2019; Jabri et al., 2020; Li & Liu, 2023), and report the percentage of correctly tracked keypoints (PCK) for JHMDB. We compare our method with baseline models in Table 5. Table 5 shows that our approach achieves state-of-the-art performance in the human pose tracking task.

## 5 CONCLUSION

In this work, we leverage latent representations from video diffusion models for pixel-level tracking. Benefiting from video diffusion models’ ability to incorporate information across multiple frames, our work introduces temporal reasoning to the tracking tasks. Without additional training, our method improves tracking performance in various video scenarios, even enabling tracking of similar-looking objects where previous methods struggle. Experimental results show that our approach achieves state-of-the-art tracking performance, outperforming recent studies by up to 6 points. Our work highlights the potential of video generative models in tracking applications beyond their original use in video synthesis task.

Table 3: **Ablation study on UNet blocks.** TED<sup>†</sup> achieves the best tracking results using  $\mathbf{R}_v$  from block 2.

Block	$J_m \& F_m$	$J_m$	$F_m$
0	24.8	28.2	21.4
1	47.6	52.7	42.5
2	66.3	63.4	69.1
3	31.5	27.2	35.8

Table 4: **Pretrained diffusion models for TED.** Our TED achieves the best tracking results using representations from I2VGen-XL and ADM.

Video	Image	$\mathcal{J} \& \mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$
SVD	SD	71.5	68.9	74.1
SVD	ADM	76.6	73.6	79.7
I2V	SD	71.7	69.0	74.5
I2V	ADM	<b>77.6</b>	<b>74.4</b>	<b>80.8</b>

Table 5: **Results on JHMDB dataset.** Our method achieves state-of-the-art performance in human pose tracking.

Method	PCK@0.1	PCK@0.2
SFC (Hu et al., 2022)	61.9	83.0
SMTc (Qian et al., 2023)	62.5	84.1
DIFT (Tang et al., 2023)	63.4	84.3
Spa-then-Temp (Li & Liu, 2023)	66.4	84.4
<b>TED (Ours)</b>	<b>68.3</b>	<b>85.8</b>

## REFERENCES

- 540  
541  
542 Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and  
543 image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference*  
544 *on Computer Vision*, pp. 1728–1738, 2021.
- 545 Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik  
546 Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling  
547 latent video diffusion models to large datasets. *arXiv preprint arXiv:2311.15127*, 2023a.
- 548  
549 Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and  
550 Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models.  
551 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
552 22563–22575, 2023b.
- 553 Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe  
554 Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video  
555 generation models as world simulators. 2024. URL [https://openai.com/research/  
556 video-generation-models-as-world-simulators](https://openai.com/research/video-generation-models-as-world-simulators).
- 557  
558 Sergi Caelles, Kevis-Kokitsi Maninis, Jordi Pont-Tuset, Laura Leal-Taixé, Daniel Cremers, and  
559 Luc Van Gool. One-shot video object segmentation. In *Proceedings of the IEEE conference on*  
560 *computer vision and pattern recognition*, pp. 221–230, 2017.
- 561 Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and  
562 Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the*  
563 *IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- 564  
565 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for  
566 contrastive learning of visual representations. In *International conference on machine learning*, pp.  
567 1597–1607. PMLR, 2020.
- 568 Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of*  
569 *the IEEE/CVF conference on computer vision and pattern recognition*, pp. 15750–15758, 2021.
- 570  
571 Jooyoung Choi, Jungbeom Lee, Chaehun Shin, Sungwon Kim, Hyunwoo Kim, and Sungroh Yoon.  
572 Perception prioritized training of diffusion models. In *Proceedings of the IEEE/CVF Conference*  
573 *on Computer Vision and Pattern Recognition*, pp. 11472–11481, 2022.
- 574 Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances*  
575 *in neural information processing systems*, 34:8780–8794, 2021.
- 576  
577 Tobias Gerstenberg, Noah D Goodman, David A Lagnado, and Joshua B Tenenbaum. How, whether,  
578 why: Causal judgments as counterfactual contrasts. In *CogSci*, 2015.
- 579 Daniel Gordon, Kiana Ehsani, Dieter Fox, and Ali Farhadi. Watching the world go by: Representation  
580 learning from unlabeled videos. *arXiv preprint arXiv:2003.07990*, 2020.
- 581  
582 Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J  
583 Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, et al. Kubric: A scalable dataset  
584 generator. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
585 pp. 3749–3761, 2022.
- 586  
587 Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena  
588 Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar,  
589 et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in neural*  
590 *information processing systems*, 33:21271–21284, 2020.
- 591 Yuwei Guo, Ceyuan Yang, Anyi Rao, Zhengyang Liang, Yaohui Wang, Yu Qiao, Maneesh Agrawala,  
592 Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models  
593 without specific tuning. In *The Twelfth International Conference on Learning Representations*,  
2024.

- 594 Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for  
595 unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on*  
596 *computer vision and pattern recognition*, pp. 9729–9738, 2020.
- 597 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in*  
598 *neural information processing systems*, 33:6840–6851, 2020.
- 600 Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J  
601 Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646,  
602 2022.
- 603 Yingdong Hu, Renhao Wang, Kaifeng Zhang, and Yang Gao. Semantic-aware fine-grained corre-  
604 spondence. In *European Conference on Computer Vision*, pp. 97–115. Springer, 2022.
- 606 HuggingFace. Diffusers documentation. [https://huggingface.co/docs/diffusers/](https://huggingface.co/docs/diffusers/index)  
607 [index](https://huggingface.co/docs/diffusers/index), 2024.
- 608 Allan Jabri, Andrew Owens, and Alexei Efros. Space-time correspondence as a contrastive random  
609 walk. *Advances in neural information processing systems*, 33:19545–19560, 2020.
- 611 Hueihan Jhuang, Juergen Gall, Silvia Zuffi, Cordelia Schmid, and Michael J Black. Towards  
612 understanding action recognition. In *Proceedings of the IEEE international conference on computer*  
613 *vision*, pp. 3192–3199, 2013.
- 614 Zihang Lai and Weidi Xie. Self-supervised learning for video correspondence flow. In *BMVC*, 2019.
- 616 Zihang Lai, Erika Lu, and Weidi Xie. Mast: A memory-augmented self-supervised tracker. In  
617 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
618 6479–6488, 2020.
- 619 Rui Li and Dong Liu. Spatial-then-temporal self-supervised learning for video correspondence.  
620 In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp.  
621 2279–2288, 2023.
- 622 Xueting Li, Sifei Liu, Shalini De Mello, Xiaolong Wang, Jan Kautz, and Ming-Hsuan Yang. Joint-task  
623 self-supervised learning for temporal correspondence. *Advances in Neural Information Processing*  
624 *Systems*, 32, 2019.
- 626 Andrzej Maćkiewicz and Waldemar Ratajczak. Principal components analysis (pca). *Computers &*  
627 *Geosciences*, 19(3):303–342, 1993.
- 628 Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew,  
629 Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with  
630 text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021.
- 632 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,  
633 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning  
634 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
- 635 Jordi Pont-Tuset, Federico Perazzi, Sergi Caelles, Pablo Arbeláez, Alexander Sorkine-Hornung, and  
636 Luc Van Gool. The 2017 davis challenge on video object segmentation. *arXiv:1704.00675*, 2017.
- 637 Rui Qian, Shuangrui Ding, Xian Liu, and Dahua Lin. Semantics meets temporal correspondence:  
638 Self-supervised object-centric learning in videos. In *Proceedings of the IEEE/CVF International*  
639 *Conference on Computer Vision*, pp. 16675–16687, 2023.
- 641 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-  
642 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2):3, 2022.
- 643 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
644 resolution image synthesis with latent diffusion models, 2021.
- 646 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-  
647 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*  
*ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.



- 648 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical  
649 image segmentation. In *Medical image computing and computer-assisted intervention–MICCAI*  
650 *2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceedings, part III*  
651 *18*, pp. 234–241. Springer, 2015.
- 652 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.  
653 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceed-*  
654 *ings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–22510,  
655 2023.
- 656 Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar  
657 Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic  
658 text-to-image diffusion models with deep language understanding. *Advances in Neural Information*  
659 *Processing Systems*, 35:36479–36494, 2022.
- 660 Mohammadreza Salehi, Efstratios Gavves, Cees GM Snoek, and Yuki M Asano. Time does tell:  
661 Self-supervised time-tuning of dense image representations. In *Proceedings of the IEEE/CVF*  
662 *International Conference on Computer Vision*, pp. 16536–16547, 2023.
- 663 Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent  
664 correspondence from image diffusion. *Advances in Neural Information Processing Systems*, 36:  
665 1363–1389, 2023.
- 666 Ran Tao, Efstratios Gavves, and Arnold WM Smeulders. Siamese instance search for tracking. In  
667 *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1420–1429,  
668 2016.
- 669 Tomer David Ullman. *On the nature and origin of intuitive theories: learning, physics and psychology*.  
670 PhD thesis, Massachusetts Institute of Technology, 2015.
- 671 Paul Voigtlaender and Bastian Leibe. Online adaptation of convolutional neural networks for video  
672 object segmentation. *arXiv preprint arXiv:1706.09364*, 2017.
- 673 Carl Vondrick, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. Tracking  
674 emerges by coloring videos. In *Proceedings of the European conference on computer vision*  
675 *(ECCV)*, pp. 391–408, 2018.
- 676 Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Mod-  
677 elscope text-to-video technical report. *arXiv preprint arXiv:2308.06571*, 2023.
- 678 Xiaolong Wang, Allan Jabri, and Alexei A Efros. Learning correspondence from the cycle-consistency  
679 of time. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*,  
680 pp. 2566–2576, 2019.
- 681 Xinlong Wang, Rufeng Zhang, Chunhua Shen, Tao Kong, and Lei Li. Dense contrastive learning  
682 for self-supervised visual pre-training. In *Proceedings of the IEEE/CVF conference on computer*  
683 *vision and pattern recognition*, pp. 3024–3033, 2021.
- 684 Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-  
685 parametric instance discrimination. In *Proceedings of the IEEE conference on computer vision*  
686 *and pattern recognition*, pp. 3733–3742, 2018.
- 687 Enze Xie, Jian Ding, Wenhai Wang, Xiaohang Zhan, Hang Xu, Peize Sun, Zhenguo Li, and Ping Luo.  
688 Detco: Unsupervised contrastive learning for object detection. In *Proceedings of the IEEE/CVF*  
689 *international conference on computer vision*, pp. 8392–8401, 2021a.
- 690 Zhenda Xie, Yutong Lin, Zheng Zhang, Yue Cao, Stephen Lin, and Han Hu. Propagate yourself:  
691 Exploring pixel-level consistency for unsupervised visual representation learning. In *Proceedings*  
692 *of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16684–16693, 2021b.
- 693 Jiarui Xu and Xiaolong Wang. Rethinking self-supervised correspondence learning: A video  
694 frame-level similarity perspective. In *Proceedings of the IEEE/CVF International Conference on*  
695 *Computer Vision*, pp. 10075–10085, 2021.

702 Ning Xu, Linjie Yang, Yuchen Fan, Dingcheng Yue, Yuchen Liang, Jianchao Yang, and Thomas  
703 Huang. Youtube-vos: A large-scale video object segmentation benchmark. *arXiv preprint*  
704 *arXiv:1809.03327*, 2018.

705 Zongxin Yang, Yuhang Ding, Yunchao Wei, and Yi Yang. Cfbi+: Collaborative video object  
706 segmentation by multi-scale foreground-background integration. *The 2020 DAVIS Challenge on*  
707 *Video Object Segmentation - CVPR Workshops*, 2020.

708 Kexin Yi, Chuang Gan, Yunzhu Li, Pushmeet Kohli, Jiajun Wu, Antonio Torralba, and Joshua B  
709 Tenenbaum. Clevrer: Collision events for video representation and reasoning. *arXiv preprint*  
710 *arXiv:1910.01442*, 2019.

711 Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun,  
712 and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot  
713 semantic correspondence. *Advances in Neural Information Processing Systems*, 36, 2024.

714 Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qing, Xiang Wang,  
715 Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded  
716 diffusion models. 2023.

717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

## 756 A DISCUSSIONS

### 757 A.1 ADVANTAGES OF OUR WORK OVER DIFT

758 We clarify and highlight the advantages of our work over the state-of-the-art DIFT (Tang et al., 2023)  
759 as follows.

- 760 • **We solve a task that tracks similar-looking objects which DIFT cannot solve.** Tracking  
761 similar-looking objects in label propagation is a very fundamental task in the field. Since  
762 DIFT learns only appearance features, it fails to track similar-looking objects.
- 763 • **Our work uses temporal motions learned from video diffusion models in tracking,  
764 providing new insights into how motion-based tracking emerges.** Our experiments and  
765 analysis show that temporal layers in video diffusion models enable motion-aware features  
766 necessary for tracking similar-looking objects, which are absent in DIFT.
- 767 • **Improved tracking accuracy across various scenarios.** Our work significantly outper-  
768 forms DIFT in tracking performance in various videos, such as those with severe object  
769 deformation, achieving 1.9% higher accuracy on DAVIS and 5.3% on YouTube-Similar.

### 770 A.2 ADVANTAGE OF OUR WORK IN LEARNING TEMPORAL FEATURES

771 We clarify and highlight the advantages of our work over prior studies in learning temporal features  
772 as follows.

- 773 • **Better representations obtained by solving a harder generative task.** Previous methods  
774 are trained on easier tasks that always have shortcuts. For example, mismatched patches  
775 in Wang et al. (2019); Jabri et al. (2020) or objects in Gordon et al. (2020); Xu & Wang  
776 (2021) with similar appearances can also yield low training loss. In contrast, our video  
777 diffusion models are trained to fully reconstruct every pixel from noisy inputs, enabling  
778 better representation learning.
- 779 • **Advanced temporal attention vs. simple pairwise correlation.** During training, prior  
780 methods learn temporal features by simple correlations between spatial features across  
781 frames (Vondrick et al., 2018; Wang et al., 2019; Li et al., 2019; Lai & Xie, 2019; Lai et al.,  
782 2020; Qian et al., 2023; Li & Liu, 2023), which fail to distinguish similar-looking objects.  
783 In contrast, our video diffusion model uses temporal attention layers to integrate multiple  
784 frames, enabling advanced reasoning in complex scenarios like the deer with changing  
785 positions in Figure 1.
- 786 • **Significantly improved tracking accuracy.** Quantitative results and visualization show  
787 that our method significantly improves the tracking performance compared to prior studies,  
788 by more than 3.5% on DAVIS and 6.4% on YouTube-Similar.

## 798 B EXPERIMENTAL SETUPS

### 799 B.1 VIDEO LABEL PROPAGATION

800 In this work, we evaluate the video label propagation task, which predicts pixel-level labels for  
801 subsequent video frames given the ground-truth label of the first frame. We follow the evaluation  
802 protocol of prior studies (Jabri et al., 2020; Caron et al., 2021; Tang et al., 2023), as detailed in  
803 Algorithm 1. Pixel-level labels for each frame are predicted based on the frame representations and  
804 labels of previous frames. For the current frame, we first identify similar pixel pairs between this  
805 frame and the previous frames by computing the similarities of their pixel representations. Labels for  
806 the current frame are then predicted by aggregating the labels of similar pixels from previous frames,  
807 weighted by their pixel similarity. A key advantage of TED over prior studies is that it generates  
808 frame representations by inputting the video sequence into the 3D UNet<sub>v</sub>, which encodes the temporal  
809

810 motions learned in video generation and significantly improves tracking accuracy. To handle videos  
 811 longer than the sequence length limit of UNet<sub>v</sub>, we split each video into multiple sequences and  
 812 process each sequence separately. An optional technique to improve accuracy is allowing overlapping  
 813 frames among sequences. Another optional technique is using a batch of random noise to obtain an  
 814 averaged representation for each video following DIFT (Tang et al., 2023).

---

816 **Algorithm 1: Temporal Enhanced Diffusion tracking (TED)**

---

817 **Input:** Video frames  $I_1$  to  $I_N$ ; Ground-truth label for the first frame  $L_1$ ; Video diffusion model  
 818 with UNet<sub>v</sub>; Image diffusion model with UNet<sub>i</sub>.

820 **Output:** Label predictions  $L_2$  to  $L_N$  for frames  $I_2$  to  $I_N$ .

```

821 1 Let  $d$  be the sequence length defined by UNetv, split all video frames to  $n = \lfloor \frac{N}{d} \rfloor$  sequences;
822 2 Initialize queue  $Q = \emptyset$  to store the representations and labels of the previous  $p$  frames;
823 3 for each sequence  $j = 0$  to  $n - 1$  do
824 4   Select the frames  $I_{1+j \cdot d}$  to  $I_{(j+1) \cdot d}$  as the current sequence;
825
826   Step 1: Computation of Frame Representations
827
828   6 Compute the video diffusion representation  $R_v$  using a single forward pass of UNetv:
829   7    $R_{v,1+j \cdot d}, \dots, R_{v,(j+1) \cdot d} = \text{UNet}_v(I_{1+j \cdot d}, \dots, I_{(j+1) \cdot d})$ ;
830   8 Compute the image diffusion representation  $R_i$  using  $d$  forward passes of UNeti:
831   9   for each frame  $I_k$  in  $I_{1+j \cdot d}$  to  $I_{(j+1) \cdot d}$  do
832   10     $R_{i,k} = \text{UNet}_i(I_k)$ ;
833   11 Compute the fused representation  $R_f$  following Eqn. 4:
834   12    $R_f = \text{Concat}(\alpha \|R_v\|_2, (1 - \alpha) \|R_i\|_2)$ ;
835
836   Step 2: Label Prediction
837
837   14 if  $j = 0$  then
838   15    $\lfloor$  Add  $(R_{f,1}, L_1)$  of the first frame to the queue  $Q$ ;
839   16 for each frame  $I_k$  in the sequence from  $I_{1+j \cdot d}$  to  $I_{(j+1) \cdot d}$  do
840   17   if  $k = 1$  then
841   18      $\lfloor$  Skip the first frame since the ground truth label is already provided ;
842   19   Compute the pixel similarity matrix  $A$  between pixel representations of current frame
843   20    $R_{f,k}$  and previous frames  $R \in Q$ ;
844   21   for each pixel in the frame  $I_k$  do
845   22     Retain only the similarities for spatially neighboring pixels in  $A$ ;
846   23     Apply top- $\kappa$  filtering to retain the strongest similarities and set the remaining values
847     in  $A$  to zero;
848   24   Predict the labels for the current frame  $k$  by propagating the labels from the most similar
849   25   pixels in previous frames, weighted by their pixel similarity:
850   26    $L_k = A \cdot (\text{labels } L \in Q)$ ;
851   27   Add  $(R_{f,k}, L_k)$  to the queue  $Q$ ;
852   28   if the size of  $Q$  exceeds the maximum allowed size  $p$  then
853   29      $\lfloor$  Remove the oldest entry from the queue  $Q$ ;
854
855 28 return  $L_2$  to  $L_N$ .

```

---

## 858 B.2 PRETRAINED DIFFUSION MODELS

859  
 860 In our work, we utilize pretrained diffusion models without additional training. Our framework  
 861 supports any pretrained diffusion model and we use open-sourced checkpoints for our experiments.  
 862 For video diffusion models, we use the official weights of Stable Video Diffusion (Blattmann et al.,  
 863 2023a) and I2VGen-XL (Zhang et al., 2023) available on Hugging Face (HuggingFace, 2024). For  
 image diffusion models, we use pretrained weights from Hugging Face for Stable Diffusion (Rombach



et al., 2022) (version 2-1) and from the official GitHub repository for ADM (Dhariwal & Nichol, 2021). We follow the configurations of DIFT (Tang et al., 2023) and summarize as in Table 6.

Table 6: Experimental setups of TED in video label propagation task.

Dataset	Video diffusion			Image diffusion			Fusion weight	Softmax temp	Propagation radius	k for top-k
	Model	Timestep	Block	Model	Timestep	Block				
DAVIS	I2VGen-XL	300	2	ADM	51	7	0.4	0.2	15	10
Youtube-Similar	I2VGen-XL	600	2	ADM	51	7	0.6	0.1	15	10

### B.3 EVALUATION DATASETS

We evaluate TED on the semi-supervised video object segmentation task using three datasets: DAVIS-2017 (Pont-Tuset et al., 2017), Youtube-Similar, and Kubric-Similar. Figure 9 shows video examples from each dataset.

- **DAVIS-2017** (Pont-Tuset et al., 2017): A widely used benchmark for semi-supervised object segmentation. Following prior work (Caron et al., 2021; Tang et al., 2023), we use the val subset, which includes 30 videos with 2023 frames and 59 annotated objects.
- **Youtube-Similar**: We propose this benchmark to evaluate tracking on multiple similar-looking objects. It includes 28 videos from Youtube-VOS (Xu et al., 2018) with 839 frames and 69 annotated objects.
- **Kubric-Similar**: We use Kubric simulator (Greff et al., 2022) to generate this dataset for tracking identical objects. Each of the 14 videos contains two identical balls with random sizes and movements, totaling 224 frames and 28 objects.



Figure 9: **Dataset examples.** We present video examples from various evaluation datasets. Following prior work (Caron et al., 2021; Tang et al., 2023), we evaluate our method on the widely-used DAVIS-2017 dataset (Pont-Tuset et al., 2017), shown in the first two columns of the figure. For the first time, we propose the challenging task of tracking multiple similar-looking objects in video label propagation. To assess model performance in this setting, we introduce two new datasets: Youtube-Similar (the third and fourth columns) and Kubric-Similar (the fifth column).

### B.4 FEATURE VISUALIZATION

In Section 3.3, we use PCA (Maćkiewicz & Ratajczak, 1993) to reduce the dimension of pixel representations for visualization. Figure 5(c) visualizes the representations after PCA, where similar colors indicate similar pixel representations. If different objects have distinct pixel colors, it indicates they are successfully distinguished from each other. Figure 5(c) shows that our work succeeds in distinguishing and tracking similar-looking objects (third column), unlike DIFT which learns similar pixel representations for different objects and fails in tracking (second column). These results highlight the effectiveness of temporal motions in our work for tracking, which DIFT lacks.

## C ADDITIONAL RESULTS

### C.1 COMPUTATION COST ANALYSIS

We compare computation cost of our method with DIFT (Tang et al., 2023) in Table 7. We track a 100-frame video, reporting average time per frame and maximum GPU memory. Our TED (efficient) outperforms DIFT (best) by 1.5% in accuracy with similar speed and slightly higher memory use, while TED (best) achieves higher accuracy at greater computation cost. In real applications, users can choose the version based on their requirements on accuracy and efficiency.

We introduce the setups for computation cost analysis as follows. We test the model on a single NVIDIA TITAN RTX GPU using a 100-frame DAVIS video. Following DIFT, we introduce two TED versions, efficient and best, based on whether to use the optional techniques. For DIFT, the optional technique involves averaging representations using a batch of noise. For TED, it includes both averaged representations and overlapping frames among sequences, as discussed in Appendix B.1.

Our work demonstrates, for the first time, that motions learned from video diffusion models can solve perception challenges and achieve state-of-the-art results. Our work offers new insights for diffusion and tracking, benefiting both communities. We believe our method can be further accelerated with future research on diffusion model acceleration as well as advances in computing and resources.

Table 7: **Computation cost analysis.** Our TED (efficient) outperforms DIFT (best) by 1.5% in accuracy with similar speed and slightly higher memory use, while TED (best) achieves higher accuracy at greater computation cost. Here, the time refers to the duration required to track a single image.

Model	Version	Optional Techniques	Accuracy	Time (s)	FPS	Memory (GB)
DIFT (Tang et al., 2023)	Efficient	No	74.7	0.73	1.37	5.53
	Best	Yes	75.7	1.36	0.74	9.25
TED(ours)	Efficient	No	77.2	1.21	0.82	11.65
	Best	Yes	77.6	2.24	0.47	15.20

### C.2 DISCUSSIONS ON THE TRAINING DATASET

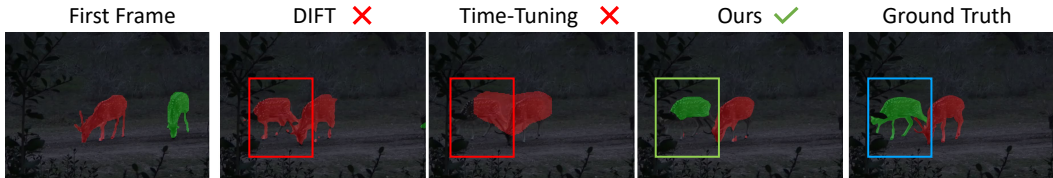
To investigate the influence of training dataset on the tracking results, we train image diffusion model from DIFT (Tang et al., 2023) on the same training dataset as our video diffusion model for comparison. Table 8 shows that without temporal modeling, training on additional video data fails to track similar-looking objects, indicated by a low  $\mathcal{J}\&\mathcal{F}_m$  of 43.8% on Kubric-Similar. Web-Vid (Bain et al., 2021) has lower individual image quality (Guo et al., 2024), such as motion blur and watermarks. Fine-tuning DIFT on Web-Vid even reduces performance. In contrast, our TED achieves significant improvements using video diffusion models and effectively distinguishes similar-looking objects, demonstrating the importance of learning temporal motions from video diffusion models for tracking.

Table 8: **Fine-tune DIFT’s image diffusion models on video datasets.** DIFT fails to track similar-looking objects even when trained on the same datasets as our video diffusion models. This is because image diffusion models learn only appearance features from video datasets, lacking the temporal motion information critical for tracking.

Version	Model	Dataset	Kubric-Similar			Youtube-Similar			DAVIS		
			$\mathcal{J}\&\mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$	$\mathcal{J}\&\mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$	$\mathcal{J}\&\mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$
Original	DIFT	ImageNet	52.6	50.3	54.8	60.7	59.8	61.7	75.7	72.7	78.6
Finetune on Web-Vid	DIFT	ImageNet, Web-Vid	43.8	40.1	47.4	58.9	58.3	59.4	72.9	70.1	75.6
Ours	TED	ImageNet, Web-Vid	<b>90.4</b>	<b>86.9</b>	<b>94.0</b>	<b>66.0</b>	<b>65.1</b>	<b>67.0</b>	<b>77.6</b>	<b>74.4</b>	<b>80.8</b>

972 C.3 RESULTS OF TIME-TUNING METHOD  
973

974 We use Time-Tuning features (Salehi et al., 2023) for video label propagation task and find that  
975 it fails to distinguish similar-looking objects in our work, as shown in Figure 10. This failure is  
976 because Time-Tuning is trained to learn semantic features for semantic segmentation task, as shown  
977 in Figure 3 of the original paper (Salehi et al., 2023), which lacks object motions needed in tracking  
978 similar-looking objects.



980  
981  
982  
983  
984  
985  
986 Figure 10: Time-Tuning fails to distinguish multiple similar-looking objects.  
987

988 C.4 RESULTS WITH ADDITIONAL DINO FEATURES  
989

990 Prior work (Zhang et al., 2024) shows that the combination of Stable Diffusion and DINOv2 (Oquab  
991 et al., 2023) features significantly improves performance in semantic correspondence task. Follow-  
992 ing (Zhang et al., 2024), we add DINOv2 features to our TED and report the tracking results in  
993 Table 9. Table 9 shows that incorporating additional DINOv2 features in our TED does not further  
994 improve tracking performance.

995  
996 Table 9: **TED with additional DINOv2 features.** We introduce additional DINOv2 features as a  
997 complementary to our TED method following Zhang et al. (2024). We find that additional DINOv2  
998 features do not further improve the performance of our TED in the tracking task.

999

Model	Features	Kubric-Similar			Youtube-Similar			DAVIS		
		$\mathcal{J} \& \mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$	$\mathcal{J} \& \mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$	$\mathcal{J} \& \mathcal{F}_m(\uparrow)$	$\mathcal{J}_m(\uparrow)$	$\mathcal{F}_m(\uparrow)$
TED (With DINOv2 features)	ADM, I2VGen-XL, DINOv2	90.0	86.6	93.5	65.9	65.0	66.7	77.3	74.2	80.5
TED (Ours)	ADM, I2VGen-XL	<b>90.4</b>	<b>86.9</b>	<b>94.0</b>	<b>66.0</b>	<b>65.1</b>	<b>67.0</b>	<b>77.6</b>	<b>74.4</b>	<b>80.8</b>

1000  
1001  
1002  
1003  
1004  
1005  
1006  
1007  
1008  
1009  
1010  
1011  
1012  
1013  
1014  
1015  
1016  
1017  
1018  
1019  
1020  
1021  
1022  
1023  
1024  
1025