# VidSeg: Training-free Video Semantic Segmentation based on Diffusion Models

Qian Wang[1]  Abdelrahman Eldesokey[1]  Mohit Mendiratta[2]  Fangneng Zhan[2]
Adam Kortylewski[2]  Christian Theobalt[2]  Peter Wonka[1]
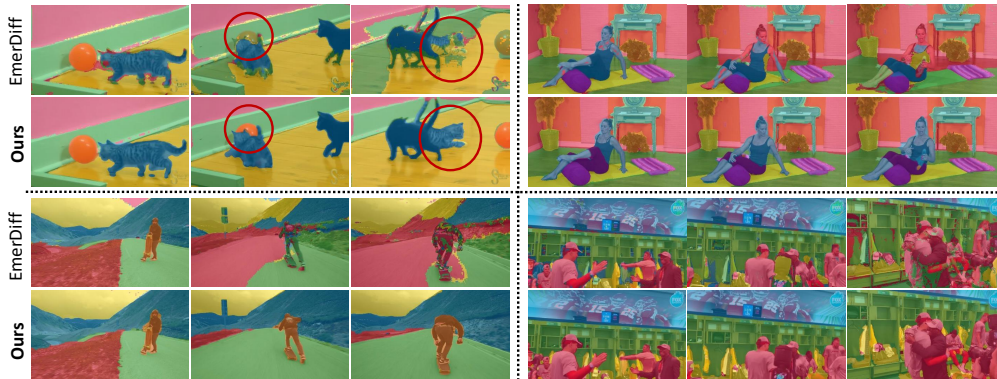[1]KAUST
[2]Max Planck Institute for Informatics

Figure 1. We propose the first *training-free* diffusion-based approach for Video Semantic Segmentation (VSS). Our approach produces temporally consistent predictions compared to the diffusion-based image segmentation method EmerDiff [15].

## Abstract

*We introduce the first* training-free *approach for Video Semantic Segmentation (VSS) based on pre-trained diffusion models, termed* VidSeg. *A growing research direction attempts to employ diffusion models to perform downstream vision tasks by exploiting their deep understanding of image semantics. Yet, the majority of these approaches have focused on image-related tasks like semantic segmentation, with less emphasis on video tasks such as VSS. Ideally, diffusion-based image semantic segmentation approaches can be applied to videos in a frame-by-frame manner. However, we find their performance on videos to be subpar due to the absence of any modeling of temporal information inherent in the video data. To this end, we tackle this problem and introduce a framework tailored for VSS based on pre-trained image and video diffusion models. We propose building a scene context model based on the diffusion features, where the model is autoregressively updated to adapt to scene changes. This context model predicts per-frame coarse segmentation maps that are temporally consistent. To refine these maps further, we propose a correspondence-based refinement strategy that aggregates predictions temporally, resulting in more confi-dent predictions. Finally, we introduce a masked modulation approach to upsample the coarse maps to a high-quality full resolution. Experiments show that our proposed approach significantly outperforms existing training-free image semantic segmentation approaches on various VSS benchmarks without any training or fine-tuning. Moreover, it rivals supervised VSS approaches on the VSPW dataset despite not being explicitly trained for VSS.*

## 1. Introduction

Diffusion models [2, 9, 17, 20, 21] have showcased remarkable capabilities in learning complex data distributions effectively. This was achieved by exploiting their scalability to train on large-scale datasets [1, 22], allowing them to generate high-quality images and videos with soaring diversity and fidelity. Interestingly, those models could learn a profound understanding of images and their semantics as an indirect consequence of their large-scale training. This entitled them to be considered as *Foundation Models* with high degrees of generalizability and comprehension of images. As a result, a growing research direction attempts to use the internal representations of diffusion models to perform
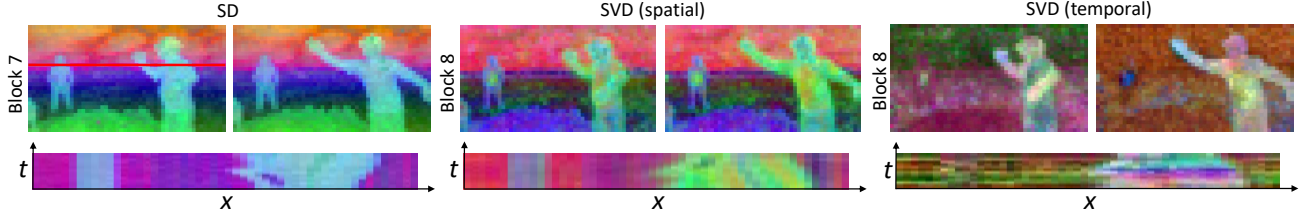
Figure 2. A visualization of the first three PCA components for the features of two video frames extracted from the most semantically rich blocks in SD (Block 7) and SVD (Bock 8). In the second row, we show the $x$-$t$ slice of an image row (highlighted in the red line in the leftmost image) horizontally across the PCA visualization ($x$-axis) and stack it chronologically across the full batch of video frames ($t$-axis). The plot shows that the spatial features of both SD and SVD are temporally more consistent between video frames compared to the features of temporal layers in SVD.

various downstream image vision tasks. For instance, pre-trained diffusion models were used to perform *image vision tasks* such as semantic correspondence [12, 23], keypoints detection [8], and semantic segmentation [13, 15].

Video vision tasks, on the other hand, have not received the same attention compared to their image counterparts. A simple approach is adopting image-based approaches to solve video tasks. To investigate this, we test the diffusion-based image semantic segmentation approach EmerDiff [15] on the task Video Semantic Segmentation (VSS) in a frame-by-frame manner. VSS aims to predict a semantic class for every pixel in each frame according to the pre-defined categories in a video. Our initial experiments show that EmerDiff performs poorly in segmenting videos in terms of temporal consistency, as shown in Figure 1. This can be attributed to the lack of modeling of video temporal information, causing inconsistent predictions across frames.

An intuitive solution for enhancing the temporal consistency of these approaches is employing a video diffusion model, *e.g.*, Stable Video Diffusion (SVD) [2] as a backbone. The SVD model is initially trained on images, then expanded with temporal layers and fine-tuned on videos. Ideally, the temporal features from SVD should exhibit better temporal consistency. To investigate this hypothesis, we visualize the temporal features of a pre-trained SVD in Figure 2 to examine their temporal consistency. The figure shows that the temporal features are surprisingly unstable and tend to change significantly between video frames. On the other hand, the spatial features encode the structure of the scene, similar to Image Stable Diffusion [19] (SD). Based on these observations, we capitalize on the spatial features of either SD or SVD and attempt to enhance them temporally.

In this paper, we introduce a diffusion-based training-free approach for VSS, called *VidSeg*. First, we build a *scene context model* based on the features of a pre-trained image (SD) or video (SVD) diffusion model. This context model predicts per-frame coarse segmentation maps and is autoregressively updated to accommodate scene changes throughout the video. To further enhance the temporal and spatial con-

sistency of these coarse maps, we propose a correspondence-based refinement (CBR) strategy encompassing a pixel-wise voting scheme between the video frames. Finally, we propose a masked modulation process to reconstruct the full-resolution segmentation maps that are more stable and less noisy than Namekata et al. [15]. Experiments show that our proposed approach significantly outperforms training-free image semantic segmentation methods on various VSS benchmarks. More specifically, we improve mIOU over image semantic approaches by at least 9.3% on VSPW, CityScapes, and Camvid datasets. Remarkably, our approach performs comparably well as supervised VSS approaches on the diverse VSPW dataset. We also show that for currently released models, SD features lead to a higher quality result than SVD features, but this trend may reverse when SVD training considers larger datasets in the future.

## 2. Related Work

### 2.1. Diffusion Models' Features

The large-scale training of diffusion models on the LAION-5B dataset [22] with 5 billion images allowed them to learn semantically rich image features. As a result, a growing direction of research attempts to employ these features to perform downstream vision tasks. Several approaches [12, 23, 26] investigated using these features to perform semantic correspondence. They observed that the features are semantically meaningful and generalize well across different objects and styles of images. For example, a human head in any arbitrary image will have features similar to those of any human head in other images, regardless of the scene variations. Those features even generalize across similar classes of objects like animal heads. At the same time, the features will differ from those of unrelated object classes like buildings, landscapes, vehicles, *etc*. EmerDiff [15] capitalized on this observation to perform image semantic segmentation. Since the diffusion features are distinct for different objects, they can easily be clustered to separate those objects and produce a coarse segmentation map. Then, they proposed a

modulation strategy to produce fine segmentation maps at a remarkable quality. However, we observed that the produced segmentation maps by EmerDiff are not temporally consistent, as illustrated in Figure 1, making it unsuitable for Video Semantic Segmentation (VSS). Therefore, we propose a diffusion-based pipeline tailored for VSS with a focus on improving temporal consistency.

## 2.2. Video Semantic Segmentation

Video semantic segmentation (VSS) [5, 10, 24, 25, 27–29] is a spatiotemporal variation of image segmentation on videos that aims to predict a pixel-wise label across the video frames. Those predictions should be temporally consistent under object deformations and camera motion, making VSS more challenging than its image counterpart. Recent approaches attempted to exploit the temporal correlation between video features to produce temporally consistent predictions. Several approaches [7, 30] incorporated optical flow prediction to model motion between frames. Other approaches [11] proposed a temporal consistency loss on the per-frame segmentation predictions as an efficient replacement for optical flow. TMANet [24] utilized a temporal attention module to capture the relations between the current frame and a memory bank. DVIS [27] further improved the efficiency by treating VSS as a first-frame-segmentation followed by a tracking problem. Recent work UniVS [10] proposed a single unified model for multiple video segmentation tasks by considering the features from previous frames as visual prompts for the consecutive frames. GvSeg [4] proposed a general model and adapted it for segmentation tasks with different properties. Despite their remarkable performance, these supervised approaches do not generalize well on unseen datasets [28]. Therefore, it is desirable to have an approach that generalizes well across datasets. Inspired by the success of EmerDiff [15] on image semantic segmentation, we attempt to exploit the diffusion features to propose the first temporally consistent training-free VSS approach.

## 3. Preliminaries

### 3.1. Stable Diffusion architecture

Stable Diffusion (SD) [16, 19] is one of the prominent latent diffusion models that achieves a good tradeoff between efficiency and quality. It is trained to approximate the image data distribution by adding noise to the latents of data samples until they converge to pure Gaussian isotropic noise. During sampling, it performs a series of Markovian denoising steps starting from pure noise to recover a noise-free latent that is decoded to produce a synthetic image. SD utilizes a UNet architecture to predict either the noise or some other signal at each time step. This UNet encompasses multiple blocks for the encoder and the decoder at different resolutions ranging from $8 \times 8$ to $64 \times 64$, where every

block has residual blocks, self-attention, and cross-attention modules. The attention is computed as:

$$f \left( \sigma \left( \frac{QK^T}{\sqrt{d}} \right) \cdot V \right) \tag{1}$$

where $Q, K, V$ are the query, key, and value vectors in the attention layers, $\sigma$ is the Softmax activation, and $f$ denotes a fully connected layer. The query is always computed from the image features, while the key and the value are computed from the image in self-attention and a conditional signal (*e.g.* textual prompt) in cross-attention. Since the semantically rich features are located in the decoder [15, 23], we only consider the decoder blocks. The decoder has 12 blocks over 4 resolutions, where we refer to the first block as 0, with a resolution of $8 \times 8$, and the last block as 11, with a resolution of $64 \times 64$.

### 3.2. Emerging Image Semantic Segmentation from Diffusion Models

EmerDiff [15] observed that it is possible to extract semantically rich features from some of the UNet blocks and use them to produce coarse semantic segmentation maps. Given an RGB image $X$ with a resolution of $H \times W$, the spatial resolution of the low-dimensional UNet features becomes $H/S_i \times W/S_i$, where $S_i$ is the scale factor of block $i$ determined by both the size of the latent representation and the downsampling factor of that block. By applying K-Means clustering on the low-dimensional feature maps from Block $b_k$ at timestep $t_k$, we obtain a set of binary masks $\mathcal{M} = \{M_1, M_2, ..., M_L\}$, where $M \in \mathbb{R}^{H/S_i \times W/S_i}$, and $L$ is the number of distinct clusters.

A *modulation* strategy is used to obtain fine-grained image-resolution segmentation maps. This is achieved by modulating the attention module at block $b_m$ and denoising timestep $t_m$ for each binary mask $M_l$ based on the following formula:

$$f \left( \sigma \left( \frac{QK^T}{\sqrt{d}} \right) \cdot V \pm \lambda M_l \right), \tag{2}$$

where $\lambda$ controls the degree of modulation. The intuition behind this process is to add or subtract a certain amount of perturbation $\lambda$ on the region specified by mask $M_l$ and then continue the denoising process to reconstruct a modulated image. By applying $+\lambda$ and $-\lambda$, we get two different modulated images denoted as $I_l^+$ and $I_l^-$ respectively. Then, a difference map is computed as $D_l = \|I_l^+ - I_l^-\|^2$, where $D_l \in \mathbb{R}^{H \times W}$. This is repeated for all masks in $\mathcal{M}$ to get a set of difference maps $\mathcal{D} = \{D_1, D_2, ..., D_L\}$. Finally, the full-resolution segmentation map is computed as $Y = \arg\max_l \mathcal{D}$.

## 4. Method

Our method encompasses three main components, as illustrated in Figure 3. First, we construct a *scene context model*
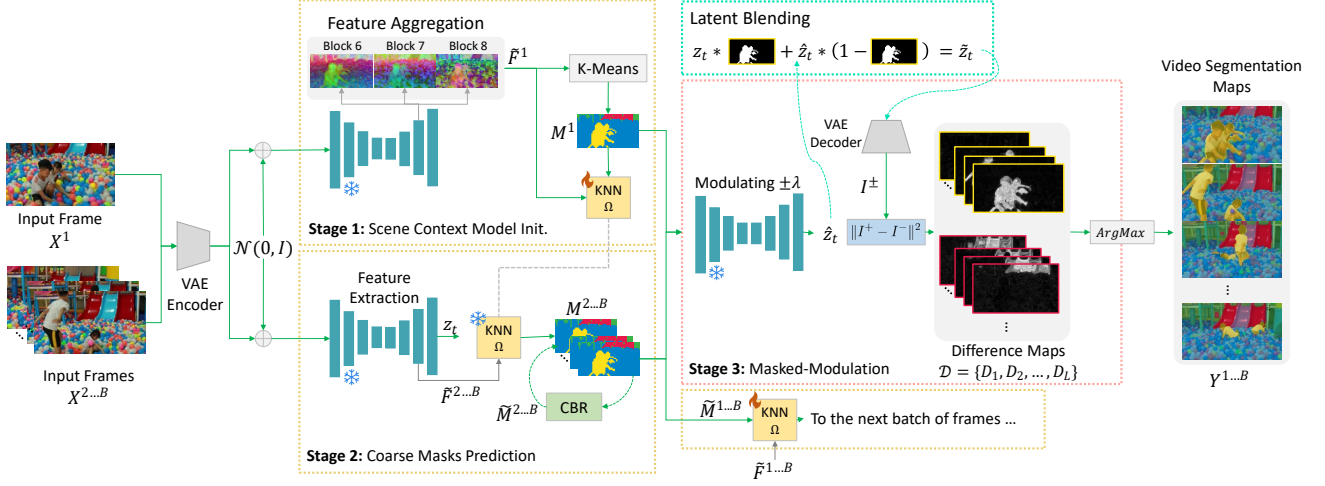
Figure 3. Our Video Semantic Segmentation (VSS) approach encompasses three stages. In Stage 1, we initialize a Scene Context Model $\Omega$ as a KNN classifier with the aggregated diffusion features $\widetilde{F}^1$ of the first frame and a coarse mask $M^1$ produced by K-Means clustering. In Stage 2, we use the context model $\Omega$ to predict coarse masks for the remaining frames in the batch $M^{2...B}$. We refine the coarse maps $M^{1...B}$ using our correspondence-based refinement (CBR). In Stage 3, we use the refined coarse masks to modulate the attention layers of the diffusion process with factor $\pm\lambda$ to obtain a modulated latent $\hat{z}_t$. Then, we blend $\hat{z}_t$ with the original unmodulated latent $z_t$ using the coarse masks to obtain a less noisy latent $\tilde{z}_t$. Finally, the latent $\tilde{z}_t$ is decoded to obtain images $I^+, I^-$ that are used to compute a set of difference maps per segment $l \in L$. The final predictions are made by applying an $\arg\max$ operation over the difference maps similar to [15]. The process is repeated for the following batch of frames where the context model is updated in an autoregressive manner using the coarse masks $M^{1...B}$ and their corresponding features $\widetilde{F}^{1...B}$.

to produce coarse segmentation maps based on the diffusion features (Section 4.1). Then, we introduce a *correspondence-based refinement* strategy to curate the coarse segmentation maps (Section 4.2). Finally, we propose a *masked modulation* approach that produces less noisy and more stable full-resolution segmentation maps (4.3). We also provide some details on adapting our approach to employ features from Stable Video Diffusion (SVD) in Section 4.4.

## 4.1. Scene Context Model

Image segmentation algorithms are designed for segmenting individual images and can only process videos in a frame-by-frame manner. This is not ideal for videos, as the per-frame predictions will be completely independent and consequently temporally inconsistent. To address this limitation, we propose to create a scene context model that is initialized at the first frame and then updated throughout the video in an autoregressive manner.

Given a video sequence $\mathcal{X} = \{X^1, X^2, \ldots, X^N\}$ with $N$ frames, we extract diffusion features $F_i^n$ for all frames in $[1, N]$, where $i$ is the decoder block. Since different blocks have different information, we aggregate features from multiple blocks by averaging to produce an aggregated feature $\widetilde{F}^n$. We found that aggregating blocks 6, 7, and 8 attains the best results (see Section 5.5). Note that these blocks share the same resolution and number of channels.

Then, we process the video in batches of length $B$.

For the first batch, we use K-Means to extract the initial coarse segmentation map $M^1$ for the first frame based on the aggregated features $\widetilde{F}^1$. Given the diffusion features $\widetilde{F}^1$ and the coarse map $M^1$ as labels, we train a KNN classifier $\Omega(M^1, \widetilde{F}^1)$ as a context model to discriminate between different clusters. Then, we use the context model $\Omega$ to predict the coarse segmentation maps for the remaining frames in the first batch $\{M^2, M^3, \ldots, M^B\}$. We refine the coarse maps further using the correspondence-based refinement (Section 4.2) and use them alongside their aggregated diffusion features to update the context model as $\Omega([\widetilde{M}^1, \widetilde{M}^2, \ldots, \widetilde{M}^B], [\widetilde{F}^1, \widetilde{F}^2, \ldots, \widetilde{F}^B])$. The context model is then used for the next batch. This strategy ensures that the context model $\Omega$ adapts to changes in the video in an auto-regressive manner.

## 4.2. Correspondence-Based Refinement

Since the context model operates purely in the feature space, it is unaware of the spatial arrangement of clusters or how they develop temporally. This might cause inconsistencies between clusters, especially across borders between objects. To alleviate this, we propose a refinement strategy based on the semantic correspondence between consecutive frames. We compute per-pixel correspondences (in the coarse map resolution) similar to Tang et al. [23] based on the diffusion
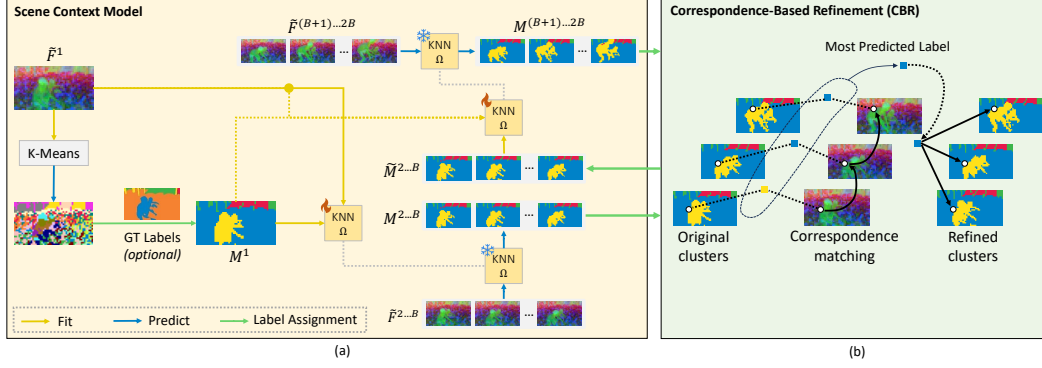
Figure 4. A detailed illustration of (a) the scene context model and (b) correspondence-based refinement.

features of block $c$ between images $j$ and $j+1$ to produce a correspondence-based coarse segmentation $\hat{M}^j$. First, we compute a trajectory $T$ for each pixel $p$ in frame $j$ that maps to the most similar pixel $q$ in the following frame $j+1$ as follows:

$$T^j[p] = \arg\max_q \Gamma(F_c^j[p], F_c^{j+1}[q]), \quad \text{with } \| p-q \|^2 \le \mathcal{T}$$
(3)

where $\Gamma$ is a distance metric that we choose to be the cosine similarity. The threshold $\mathcal{T}$ discards faulty matches that are spatially implausible. These correspondences are computed for all pixels $p$ and over all frames within the batch. Then, we define a recursive tracking function $\texttt{TRACK}$ that follows the trajectory from frame to frame to fetch the corresponding class label:

$$\texttt{TRACK}(p, j, J) = \begin{cases} \texttt{TRACK}(T^j[p], j+1, J), & \text{if } j \le J \\ p, & \text{if } j > J \end{cases},$$
(4)

where $J$ is the number of frames. Afterward, we employ this function to query class labels for each pixel on the trajectory of pixel $p$. We perform a majority voting for all pixels over the temporal axis to produce the final coarse segmentation map $\widetilde{M}^j$:

$$\widetilde{M}^j[p] = \arg\max_{l \in L} \sum_{k=j}^{B-1} \mathbf{1}\left(M^j[\texttt{TRACK}(p,j,k)] = l\right),$$
(5)

where $L$ is the number of clusters. We compute the counting using the indicator function $\mathbf{1}$, which is equal to one if the condition $M^j[\texttt{TRACK}(p,j,k)] = l$ is True, and zero otherwise. This interplay between the context model and the correspondence-based refinement leads to more accurate predictions. The context model encodes the feature space

of the video batch, while the refinement strategy spatially and temporally regularizes the predictions. This process is illustrated in Figure 4.

### 4.3. Masked Modulation

The modulation process aims to upsample the coarse segmentation masks to the full resolution of the video frames. When applying the modulation process, it is only the modulated regions that are expected to change, as explained in Section 3.2. However, in practice, the modulation process produces noise outside that region, causing discrepancies when computing the final segmentation labels. Therefore, we propose a masked modulation process that employs the coarse segmentation map to mask out both the latents and the difference maps outside the modulated region. For a denoising timestep $t$, we blend the latents as follows:

$$\widetilde{z}_t = z_t * (1 - M_l) + \hat{z}_t * M_l,$$
(6)

where $z_t$ is the latent from the unmodified sampling step, $\hat{z}_t$ is the latent from the modulated sampling process, and $M_l$ is the low-resolution mask we are modulating. We bilinearly resize $M_l$ to match the resolution of $z_t$. Even though the modulation is only performed at timestep $t_m$, we apply latent blending after $t_m$ until timesteps $t_f$, as once the attention map of a certain timestep is modified, it will influence all the following denoising timesteps.

To further suppress the noise in the difference map, we can apply the same blending strategy to the difference maps. We compute the filtered difference map $\widetilde{D}_l$ as:

$$\widetilde{D}_l = D_l * M_l + s \cdot D_l * (1 - M_l),$$
(7)

where $s$ is a scaling hyperparameter that controls the filtering strength. We bilinearly resize $M_l$ to match the size of $D_l$.

### 4.4. Adapting SVD for VSS

As it is natural to explore applying video models for video tasks, we investigate the possibility of using Stable Video

Table 1. Quantitative performance comparison.

| Method | Backbone | Training | VSPW | | | Cityscapes | Camvid | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | mIoU | $mVC_8$ | $mVC_{16}$ | mIoU | mIoU | $mVC_8$ | $mVC_{16}$ |
| TMANet[24] | ResNet-50 | Supervised | – | – | – | 80.3 | 76.5 | – | – |
| UniVS[10] | Swin-T | Supervised | 59.8 | 92.3 | – | – | – | – | – |
| DVIS++[27] | VIT-L | Supervised | 63.8 | 95.7 | 95.1 | – | – | – | – |
| CLIPpy | T5 + DINO | Unsupervised | 17.7 | 72.4 | 68.4 | 4.7 | 2.6 | 44.4 | 35.6 |
| EmerDiff (SVD) | SVD | Training-free | 39.7 | 82.1 | 78.5 | 11.0 | 7.3 | 71.7 | 64.4 |
| EmerDiff (SD) | SD 2.1 | Training-free | 43.4 | 68.9 | 64.3 | 21.5 | 6.9 | 39.8 | 32.9 |
| Ours(SVD) | SVD | Training-free | 53.2 | 89.3 | 88.0 | 36.2 | 16.6 | 87.4 | 85.8 |
| Ours(SD) | SD 2.1 | Training-free | **60.6** | **90.7** | **89.6** | **37.3** | **20.6** | **92.3** | **91.9** |

Diffusion (SVD)[2] features to perform VSS. Unlike SD, which is text-conditioned, SVD is image-conditioned. SVD adapts the SD 2.1 architecture and finetunes on a highly-curated large video dataset. Moreover, it employs a video VAE encoder to encode and decode the videos. SVD extends the SD architecture design mainly in two parts: 1) A video residual network that consists of a set of convolutional blocks that handle every frame of the video latent independently. 2) A temporal attention layer is applied on top of the output of the spatial attention layer, which computes the full attention of the features along the temporal axis on each spatial location. The output of the temporal attention and the output of the spatial attention are then mixed with a learnable weight. We use SVD as a backbone model to extract the features as well as do modulation following similar steps described in previous sections.

## 5. Experiments

To the best of our knowledge, no training-free diffusion-based video semantic segmentation (VSS) approach exists. Therefore, we compare against existing training-free image segmentation methods by adapting them to the VSS setting.

### 5.1. Evaluation Protocol

The clusters generated by K-Means in the first frame are class-agnostic. To be able to compare our predicted segmentation maps against the GT, we use the labels from the GT of the first frame to match existing segmentation clusters. Unlike common practice in Video Object Segmentation (VOS), which often relies on direct propagation of the first-frame GT to guide segmentation across subsequent frames, we simply associate exiting clusters with labels based on an initial alignment with the GT. This step is a post-processing procedure used solely to evaluate the segments against GT. In scenarios where labeling is not required, this step can be entirely omitted without impacting the segmentation itself. Therefore, our method can still be considered as **zero-shot** video segmentation method.

### 5.2. Experimental Setup

**Implementation details.** We evaluate our method with both SD 2.1 and SVD backbones, and we denote them as *Ours (SD)* and *Ours (SVD)*. Given the video frames, we encode them using VAE, and we add a certain level of noise

that corresponds to timestep $t_{inv}$ and start the denoising process at this noise level. As we need to perform modulation at timestep $t_m$, the value of $t_{inv}$ must be larger than $t_m$. Initially, we set the modulating coefficient $\lambda$ to 10 as in [15]. However, we observed that latent blending suppresses the modulating strength; thus, we increase $\lambda$ to 50 when latent blending is enabled. We set the batch size $B = 14$, which is also the original training batch size of SVD. We show more implementation details in the Supplementary Materials.

**Baselines.** We compare against unsupervised image segmentation approach CLIPpy [18] and training-free image segmentation approach EmerDiff [15]. For EmerDiff, we adapt it to our VSS setup. We use the *exact same* label assignment strategy for the first frame as explained in Sec. 5.1, but then we train a KNN classifier to predict the next frame in an autoregressive manner, similar to our approach. We also provide results of the supervised approaches DVIS++ [27], UniVS [10], and TMANet [24] to showcase where our training-free approach stands compared to them.

**Dataset.** We evaluate on the validation set of three commonly used VSS datasets: VSPW [14] with diverse videos, Cityscapes [6] and CamVid [3] with driving videos. We provide details of the settings for individual datasets in the Appendix.

**Metrics.** We report mean Intersection-over-Union (mIoU) and mean Video Consistency (mVC) [14] as quantitative metrics similar to existing VSS approaches [24, 27, 28]. mIoU describes the mean intersection-over-union between the predicted and ground-truth pixels, while mVC computes the mean categories' consistency over the long-range adjacent frames. In mVC, the common area is denoted as the area over frames where the ground-truth label does not change, and mVC is computing the ratio between the intersection of predictions in the common area over the common area. We denote mVC evaluated under 8 and 16 video frames as $mVC_8$ and $mVC_{16}$, respectively. We use both metrics together to showcase the segmentation quality on individual images as well as the overall temporal consistency.

### 5.3. Quantitative results

We provide the quantitative results in Table 1. Our approach with both SD and SVD backbones performs the best

Table 2. Ablation study. The videos we use here are the first 30 videos from VSPW validation set. CBR refers to Correspondence-Based refinement.

| Batch size $B$ | Masked Modulation | Feature Aggregation | CBR | SD2.1 | | | SVD | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | mIoU | mVC$_8$ | mVC$_{16}$ | mIoU | mVC$_8$ | mVC$_{16}$ |
| 1 | ✗ | ✗ | ✗ | 33.4 | 70.6 | 60.2 | 26.6 | 82.2 | 78.3 |
| 14 | ✗ | ✗ | ✗ | 43.4 | 76.9 | 73.8 | 38.9 | 90.2 | 88.7 |
| 14 | ✓ | ✗ | ✗ | 45.6 | 87.5 | 85.5 | 38.6 | 90.5 | 89.0 |
| 14 | ✗ | ✓ | ✗ | 45.2 | 78.4 | 75.6 | 39.0 | 89.8 | 88.6 |
| 14 | ✗ | ✗ | ✓ | 44.6 | 79.3 | 76.8 | 38.4 | 90.4 | 89.1 |
| 14 | ✓ | ✓ | ✗ | 46.6 | 87.6 | 85.7 | **39.4** | 90.2 | 89.2 |
| 14 | ✓ | ✗ | ✓ | **47.4** | 89.2 | 87.5 | 37.1 | 91.5 | 90.3 |
| 14 | ✓ | ✓ | ✓ | 46.5 | **89.8** | **88.4** | 38.2 | **92.3** | **91.3** |



Figure 5. Qualitative comparison of different unsupervised / training-free methods. Note that the color of a segmentation cluster only represents the relative index of the cluster when the video is processed. The color itself does not map to an absolute label.

in terms of all evaluation metrics on all datasets amongst the unsupervised / training-free approaches. More specifically, it improves over EmerDiff for both SD and SVD backbones in terms of mIoU with 17.2%, 13.5% on the VSPW dataset, 15.8%, 25.2% on the CityScapes dataset, and 13.7% and 9.3% on the Camvid dataset. Furthermore, our approach performs similarly to the supervised method UniVS and DVIS++ in terms of mIoU despite not being explicitly trained for this task. On the CityScapes and Camvid datasets, our approach outperforms the unsupervised / training-free methods by a huge margin.

However, there is still a performance gap of our approach compared to the supervised approaches. We attribute this to two main reasons. First, CityScapes and Camvid datasets are for driving scenarios and have challenging lighting conditions, which poses a challenge when inverting the video frames (see appendix). Secondly, we observed that for CityScapes and Camvid, our method may not be able to segment very small objects. As CityScapes at the original resolution of $2048 \times 1024$ is computationally expensive, we downsample the video frames of CityScapes by a factor of 16. Furthermore, Stable Diffusion works in latent space, necessitating a compression factor of 64 from the pixel space to

the latent space in our setup. This level of compression poses challenges for segmenting small objects, such as pedestrians and traffic poles, as many may not be adequately represented in the latent space of diffusion models. Here we show the per-class mIoUs to validate our hypothesis. For backone SD 2.1 on CityScapes dataset, the mIoU for big objects such as building, sky and vegetation are 70.5, 63.7, 68.5, respectively, while the mIoU for small objects such as traffic light, person and bicycle are 3.69, 24.5, 25.7, respectively. In the VSS evaluation setting, missing small objects can lead to significant penalties, even though these objects occupy relatively small areas in pixel space. We leave it for future work to adopt a tiled approach for high-resolution video segmentation.

### 5.4. Qualitative results

We show some qualitative results of different unsupervised / training-free methods in Figure 5. CLIPpy struggles to locate the boundaries accurately and produces coarse segments. EmerDiff can segment the first frame accurately but struggles to preserve the masks temporally. Our method with SD backbone produces the best segmentation maps in terms of segmentation quality and temporal consistency. The maps have sharper boundaries and clean clusters com-
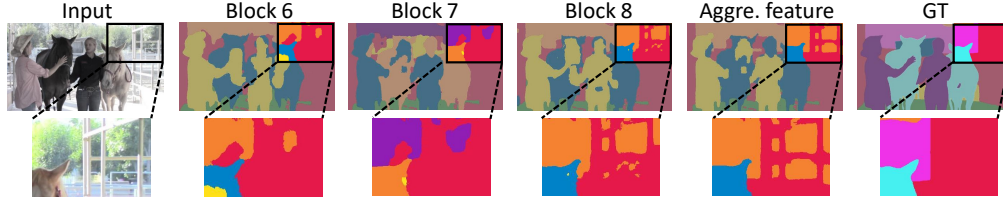
Figure 6. High-resolution segmentation map generated by aggregated feature has more details than features from a single block. We omit the low-resolution segmentation map for a better visual comparison.
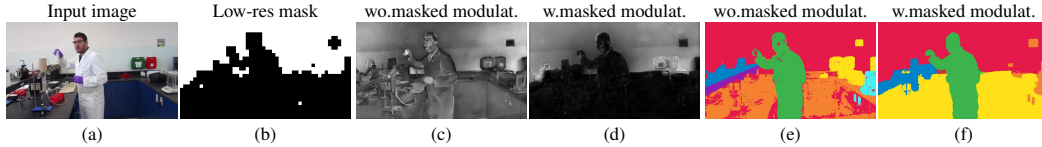


Figure 7. We show that given a low-resolution mask of a sub-region (b), the corresponding difference map without masked modulation (c) can have high activation outside the masked sub-region, which will result in spatial inconsistency on the final segmentation map (e). By applying masked modulation on the intermediate latents, the high activation on the irrelevant regions of the difference map will be removed (d), therefore producing a cleaner segmentation map (f).

pared to other methods. Ours (SVD) performs better than its EmerDiff counterpart. However, the overall segmentation quality obtained by the SVD backbone is worse than that of SD. This can be attributed to the degraded feature representation of SVD compared to SD as a result of training on a relatively small video dataset compared to SD.

### 5.5. Ablation Analysis

We show an ablation of our newly proposed components in Table 2. Incorporating more frames in training our context model greatly improves the mIoU as well as mVC. Enabling the masked modulation and the feature aggregation improves all metrics further. The best results in terms of mIoU are attained by enabling the correspondence-based refinement and disabling the feature aggregation. Despite the fact that correspondence-based refinement itself does not bring substantial improvement in terms of mIoU, it improves overall video consistency. The feature aggregation could negatively impact the mIoU due to the coarse ground truth of the VSPW dataset. Figure 6 shows an example where our approach remarkably predicts the fine details of the window as one class and the background as another class, while the ground truth annotates the whole region as a window. Finally, the best trade-off between the segmentation quality (mIoU) and temporal consistency (mVC) is achieved with all components. We provide more ablations in the Appendix.

**Feature Aggregation.** To validate the efficacy of feature aggregation, we show a visual comparison of the segmentation maps produced by features in different blocks in Figure 6. The aggregated features can encode more spatial details, which could enhance the coarse masks and, consequently, the high-resolution segmentation maps.

**Masked Modulation.** We show the qualitative comparison with or without the latent blending in Figure 7. The figure

shows that without the latent blending, the difference map in (c) contains high activations outside of the modulated sub-region indicated by the coarse binary mask. The existence of activation in these regions can lead to a false assignment of the segmentation labels, as shown in (e). For example, a part of the lab table is classified as a wall. After applying latent blending, we remove activations outside of the mask region and obtain a cleaner segmentation mask, as shown in (f).

## 6. Limitations and Future Work

One of the limitations of our approach is its dependency on the quality of the image inversion method. Moreover, fine image details are likely to be discarded due to the compression from the VAE encoder. Therefore, our approach can benefit from future research improving both VAE encoding and image inversion. It is also beneficial to investigate if other video diffusion models have semantically higher quality features than SVD. Another limitation of our approach is that it is instance-agnostic, *i.e.* it groups all objects of the same class into the same cluster. This is due to the inherent nature of diffusion features that group similar semantics. For future work, our approach can be extended to perform Video Instance or Panoptic segmentation.

## 7. Conclusions

In this work, we introduced the first *training-free* method for Video Semantic Segmentation (VSS) using pre-trained diffusion models. We proposed a pipeline tailored for VSS that leverages image and video diffusion features, and attempts to enhance their temporal consistency. Experiments showed that our proposed approach significantly outperforms training-free image semantic segmentation methods on several VSS benchmarks, and performs comparably well to supervised methods on VSPW dataset.

# References

[1] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*, 2021. 1

[2] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large datasets, 2023. 1, 2, 6

[3] Gabriel J. Brostow, Julien Fauqueur, and Roberto Cipolla. Semantic object classes in video: A high-definition ground truth database. *Pattern Recognition Letters*, 30(2):88–97, 2009. Video-based Object and Event Analysis. 6

[4] Mu Chen, Liulei Li, Wenguan Wang, Ruijie Quan, and Yi Yang. General and task-oriented video segmentation. In *ECCV*, 2024. 3

[5] Bowen Cheng, Anwesa Choudhuri, Ishan Misra, Alexander Kirillov, Rohit Girdhar, and Alexander G. Schwing. Mask2former for video instance segmentation, 2021. 3

[6] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proc. of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. 6

[7] Raghudeep Gadde, Varun Jampani, and Peter V. Gehler. Semantic video cnns through representation warping, 2017. 3

[8] Eric Hedlin, Gopal Sharma, Shweta Mahajan, Xingzhe He, Hossam Isack, Abhishek Kar Helge Rhodin, Andrea Tagliasacchi, and Kwang Moo Yi. Unsupervised keypoints from pretrained diffusion models. *arXiv preprint arXiv:2312.00065*, 2023. 2

[9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 1

[10] Minghan Li, Shuai Li, Xindong Zhang, and Lei Zhang. Univs: Unified and universal video segmentation with prompts as queries, 2024. 3, 6

[11] Yifan Liu, Chunhua Shen, Changqian Yu, and Jingdong Wang. Efficient semantic video segmentation with per-frame inference, 2020. 3

[12] Grace Luo, Lisa Dunlap, Dong Huk Park, Aleksander Holynski, and Trevor Darrell. Diffusion hyperfeatures: Searching through time and space for semantic correspondence, 2023. 2

[13] Pablo Marcos-Manchón, Roberto Alcover-Couso, Juan C. SanMiguel, and Jose M. Martínez. Open-vocabulary attention maps with token optimization for semantic segmentation in diffusion models, 2024. 2

[14] Jiaxu Miao, Yunchao Wei, Yu Wu, Chen Liang, Guangrui Li, and Yi Yang. Vspw: A large-scale dataset for video scene parsing in the wild. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4133–4143, 2021. 6

[15] Koichi Namekata, Amirmojtaba Sabour, Sanja Fidler, and Seung Wook Kim. Emerdiff: Emerging pixel-level semantic knowledge in diffusion models. In *The Twelfth International Conference on Learning Representations*, 2024. 1, 2, 3, 4, 6

[16] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis, 2023. 3

[17] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 1(2): 3, 2022. 1

[18] K. Ranasinghe, B. McKinzie, S. Ravi, Y. Yang, A. Toshev, and J. Shlens. Perceptual grouping in contrastive vision-language models. In *2023 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 5548–5561, Los Alamitos, CA, USA, 2023. IEEE Computer Society. 6

[19] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2021. 2, 3

[20] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695, 2022. 1

[21] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. *Advances in Neural Information Processing Systems*, 35:36479–36494, 2022. 1

[22] Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in Neural Information Processing Systems*, 35:25278–25294, 2022. 1, 2

[23] Luming Tang, Menglin Jia, Qianqian Wang, Cheng Perng Phoo, and Bharath Hariharan. Emergent correspondence from image diffusion, 2023. 2, 3, 4

[24] Hao Wang, Weining Wang, and Jing Liu. Temporal memory attention for video semantic segmentation. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 2254–2258. IEEE, 2021. 3, 6

[25] Shusheng Yang, Xinggang Wang, Yu Li, Yuxin Fang, Jiemin Fang, Liu, Xun Zhao, and Ying Shan. Temporally efficient vision transformer for video instance segmentation. In *Proc. IEEE Conf. Computer Vision and Pattern Recognition (CVPR)*, 2022. 3

[26] Junyi Zhang, Charles Herrmann, Junhwa Hur, Luisa Polania Cabrera, Varun Jampani, Deqing Sun, and Ming-Hsuan Yang. A tale of two features: Stable diffusion complements dino for zero-shot semantic correspondence. *Advances in Neural Information Processing Systems*, 36, 2024. 2

[27] Tao Zhang, Xingye Tian, Yu Wu, Shunping Ji, Xuebo Wang, Yuan Zhang, and Pengfei Wan. Dvis: Decoupled video instance segmentation framework, 2023. 3, 6

[28] Tao Zhang, Xingye Tian, Yikang Zhou, Shunping Ji, Xuebo Wang, Xin Tao, Yuan Zhang, Pengfei Wan, Zhongyuan Wang, and Yu Wu. Dvis++: Improved decoupled framework for universal video segmentation, 2023. 3, 6

[29] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6230–6239, 2017. 3

[30] Xizhou Zhu, Yuwen Xiong, Jifeng Dai, Lu Yuan, and Yichen Wei. Deep feature flow for video recognition, 2017. 3