

# Zero-Shot Video Editing Using Off-The-Shelf Image Diffusion Models

Wen Wang<sup>1\*</sup> Kangyang Xie<sup>1\*</sup> Zide Liu<sup>1\*</sup> Hao Chen<sup>1</sup> Yue Cao<sup>2</sup> Xinlong Wang<sup>2</sup> Chunhua Shen<sup>1†</sup>

<sup>1</sup>Zhejiang University <sup>2</sup>Beijing Academy of Artificial Intelligence

[Frame-wise Image Editing Result] A dog is walking on the ground, Van Gogh style



[Input Video] A dog is walking on the ground.



[Our Editing Result] A dog is walking on the ground, Van Gogh style



Figure 1: Video editing using an **off-the-shelf** image diffusion model. The straightforward frame-wise image editing results in severe flickering effects (first row). In contrast, our method achieves temporally-consistent editing results (third row).

## Abstract

*Large-scale text-to-image diffusion models achieve unprecedented success in image generation and editing. However, how to extend such success to video editing is unclear. Recent initial attempts at video editing require significant text-to-video data and computation resources for training, which is often not accessible. In this work, we propose vid2vid-zero, a simple yet effective method for zero-shot video editing. Our vid2vid-zero leverages off-the-shelf image diffusion models, and doesn't require training on any video. At the core of our method is a null-text inversion module for text-to-video alignment, a cross-frame modeling module for temporal consistency, and a spatial regularization module for fidelity to the original video. Without any training, we leverage the dynamic nature of the attention mechanism to enable bi-directional temporal modeling at test time. Experiments and analyses show promising*

*results in editing attributes, subjects, places, etc., in real-world videos. Code will be made available at <https://github.com/baaivision/vid2vid-zero>.*

## 1. Introduction

Text-to-image diffusion models such as DALL·E 2 [33], Imagen [38], and Stable Diffusion [36] are capable of generating unprecedented diverse and realistic images with complex objects and scenes, opening a new era of image generation. As an important application, image editing built on top of pre-trained diffusion models has also made significant progress [7, 10, 22, 29, 45]. These methods allow users to edit the input images through simple text prompts, and can achieve satisfying alignment with the target prompt and fidelity to the original image.

However, it is still unclear how to extend such success to the video editing realm. Given the input video and text prompts, the text-driven video editing algorithm is required to output an edited video that satisfies (1) text-to-video alignment: the generated edited video should be aligned to

\*Equal contribution. †Corresponding author. Part of this work was done when Wen Wang was an intern at Beijing Academy of Artificial Intelligence.

the description of the text prompt; (2) fidelity to the original video: each frame of the edited video should be consistent in content with the corresponding frame of the original video; and (3) Quality: the generated video should be temporal consistent and in high quality.

Drawing on the success of the text-to-image diffusion model, one way is to build video editing algorithms on top of video diffusion models pre-trained on large-scale video datasets [9, 28]. However, the training requires both significant amounts of paired text-video data and computational resources, which is often inaccessible. Instead, a recent work Tune-A-Video [50] fine-tunes the pre-trained text-to-image diffusion model on a single video for video synthesis. While efficient in training, it still requires more than 70M parameters per video for fine-tuning. In addition, few-shot training updates the pre-trained weights and could lead to degradation in generation quality.

Different from these methods, we aim at performing zero-shot video editing using off-the-shelf image diffusion models, without training on any video. While it may appear easy to achieve video editing via frame-wise image editing techniques, severe flickering occurs even with the content-preserving DDIM inversion [41] and cross-attention guidance [10], as shown in the first row in Fig. 1. The frame-wise image editing produces temporal inconsistent results when alternating the video style, due to the lack of temporal modeling.

To tackle this problem, we propose a simple yet effective pipeline, termed vid2vid-zero, for zero-shot video editing. Our method directly uses the pre-trained image diffusion models and is completely training-free. Specifically, it contains three major components, including a null-text inversion module for text-to-video alignment, a spatial regularization module for video-to-video fidelity, and a cross-frame modeling module for temporal consistency. To achieve the balance between effective training-free temporal modeling and reducing discrepancy to the sampling process of the pre-trained text-to-image diffusion model, a spatial-temporal attention module is proposed to model bi-directional temporal information at test time for video editing. Without bells and whistles, our method shows promising video editing results on real-world videos, as shown in the third row in Fig. 1. To summarize, the contribution of this paper is listed as follows:

- To our knowledge, we present the first zero-shot video editing method that directly uses the pre-trained text-to-image diffusion model, without any further training.
- We leverage the dynamic nature of the attention mechanism and propose a simple yet effective dense spatial-temporal attention module that achieves bi-directional temporal modeling for video editing.
- Experiments and analyses show promising results in

editing attributes, subjects, places, *etc.*, in real-world videos.

## 2. Related Work

### 2.1. Diffusion Models for Generation

With the powerful capability of approximating data distribution, diffusion models have achieved unprecedented success in the image generation domain [12, 41–43]. The magical generation fidelity and diversity outperform previous state-of-the-art generative models like Generative Adversarial Networks (GANs) [8, 21, 23, 37]. Driven by pre-trained large language models [31, 32], text-to-image diffusion models can generate high fidelity images that are consistent with the text description [33, 34, 36, 38].

Compared to image generation, video generation is a more challenging problem due to its higher-dimensional complexity and the lack of high-quality datasets [16, 46, 48, 49]. Encouraged by diffusion models’ superiority in modeling complex higher-dimensional data and their success in image generation, recent literature attempts to extend diffusion models to the video generation task [9, 15, 38, 50]. As one of the pioneers, Video Diffusion Models [15] design a new architecture using 3D U-Net [6] with factorized space-time attention [2, 3, 13] for generating temporally-coherent results. Imagen Video [38] further scales up the Video Diffusion Model to achieve photo-realistic video generation. However, the difficulty of collecting large-scale text-video datasets and the expensive computational overheads are not feasible for most researchers.

To this end, MagicVideo [51] optimizes the efficiency by performing the diffusion process in the latent space. Leveraging a pre-trained text-to-image diffusion model, Make-A-Video [40] performs video generation through a spatial-temporal decoder trained only on unlabeled video data. Recently, Tune-A-Video [50] proposes a one-shot fine-tuning strategy on text-to-image pre-trained diffusion models for text-to-video generation.

### 2.2. Diffusion Models for Editing

Besides the great success made in image generation, text-to-image diffusion models have also broken the dominance of the previous state-of-the-art models [18–20, 24] in text-driven image editing [1, 35, 44]. Building on top of a pre-trained text-to-image diffusion model, recent works [7, 10, 22, 29, 45] achieve impressive performance in text-driven image editing. Based on the key observation that the spatial layout and geometry information of generated images are retained in the text-image cross-attention map, Prompt-to-Prompt [45] achieves a fine-grained control of the spatial layout in the edited image by directly manipulating the cross-attention maps in the generating process. However, directly replacing the cross-attention maps may limit the

editing scope of various motion changes. To tackle this problem, pix2pix-zero [29] applies cross-attention guidance in generating process to ensure the structural information stays close to the input image but still has the flexibility to be changed under the guidance of the input prompt. Furthermore, they propose to compute the edit direction based on multiple sentences in the text embedding space which allows for more robust and better performance in text-driven image editing. With similar motivation, Plug-and-Play [45] achieves the state-of-the-art text-driven image editing performance with satisfying semantic and spatial correspondence by replacing spatial features in generating process with extracted corresponding counterparts of the input image.

While text-to-image diffusion models have made great success in image editing, few explorations have been made for extending them to the video editing realm. Dreamix [28] pioneers the application of a diffusion-based model for text-driven video editing. They fine-tune the video diffusion model [15] on both original videos and unordered frames set with a mixed objective to improve the performance of motion edits and temporal modeling. GEN-1 [9] achieves impressive video-editing performance under the guidance of the input prompt without fine-tuning on individual input videos. However, it requires training a latent video diffusion model conditioned on depth maps and input prompts to maintain the structure consistency. Different from these methods that are built on top of large-scale trained video diffusion models, we completely bypass the requirements of expensive computational resources and the significant amount of text-video training data, and instead achieve zero-shot video editing using the off-the-shelf image diffusion models.

Several concurrent works [5, 26, 30, 39] also make attempts to tackle the video editing problem. However, they either require fine-tuning the image diffusion model on the input video [26, 39] or rely on a video diffusion model to perform zero-shot video editing [30]. Probably the most similar work to our vid2vid-zero is Pix2Video [5], which relies on a pre-trained structure-guided (*e.g.*, depth) diffusion model and utilizes features from the previous and the first edited frames to edit the current frame. Differently, we do not assume the image diffusion model is structure-guided and can take advantage of the bi-directional temporal modeling for temporal consistency.

### 3. Method

In text-driven video editing, the user provides an input video  $X_0 = \{x_0^i \mid i \in [1, F]\}$  with  $F$  frames, with the corresponding text description (source prompt)  $c$ , and then queries the video editing algorithm with a target prompt  $\hat{c}$  to produce the edited video. To tackle this problem, we propose a simple yet effective method vid2vid-zero,

---

#### Algorithm 1 vid2vid-zero algorithm

---

```

Input:  $X_0$ : input video
 $c$ : source prompt
 $\hat{c}$ : target prompt
 $\tau_{null}, \tau_M$ : injection thresholds

▷ Real Video Inversion
 $\{X_t^{\text{inv}}\}_{t=1}^T = \text{DDIM-INV}(X_0)$ 
 $\{\mathcal{O}_t\}_{t=1}^T = \text{NULL-OPT}(\{X_t^{\text{inv}}\}_{t=1}^T, c)$ 

▷ Computing reference cross-attention maps
 $X_T \leftarrow X_T^{\text{inv}}$ 
for  $t = T \dots 1$  do
    ▷ null-text embedding is injected
     $z_t, M_t \leftarrow \hat{\epsilon}_\theta(X_t, t, c; \mathcal{O}_t)$ 
     $X_{t-1} = \text{DDIM-SAMP}(X_t, z_t, t)$ 
end for

▷ Editing with guidance
 $\hat{X}_T \leftarrow X_T^{\text{inv}}$ 
for  $t = T \dots 1$  do
    if  $t < \tau_M$  then  $M_t \leftarrow \emptyset$  ▷ without injection
    if  $t < \tau_{null}$  then  $\mathcal{O}_t \leftarrow \emptyset$ 
     $\hat{z}_t, \dots \leftarrow \hat{\epsilon}_\theta(\hat{X}_t, t, \hat{c}; \mathcal{O}_t, M_t)$ 
     $\hat{X}_{t-1} = \text{DDIM-SAMP}(\hat{X}_t, \hat{z}_t, t)$ 
end for
Output:  $\hat{X}_0$ : edited video

```

---

that edits video in a zero-shot manner using the publicly available text-to-image diffusion model, as shown in Figure 2. Specifically, our method contains three major components, a video inversion module for text-to-video alignment (Sec. 3.1), a spatial regularization module for video-to-video fidelity (Sec. 3.3), and a cross-frame modeling module for temporal consistency (Sec. 3.2).

#### 3.1. Real Video Inversion

**DDIM Inversion.** A common practice in image editing is to find the variables in the latent space that corresponds to the image. Afterward, image editing can be achieved by finding editing directions in the latent space. Motivated by their success in image editing [27, 29, 29], we first inverse each frame in the input video to the noise space, through the commonly used deterministic inversion method, DDIM inversion [41]. The inversion trajectory over  $T$  timesteps can be denoted as  $\{X_t^{\text{inv}}\}_{t=1}^T$ .

**Null-text Optimization.** Although using DDIM inversion preserves the information of each frame in the video, the obtained latent noise may not be aligned to the user-provided text description  $c$ . In other words, when sampling with the latent  $X_T^{\text{inv}}$  and the source prompt, the reconstructed video may be significantly different from the

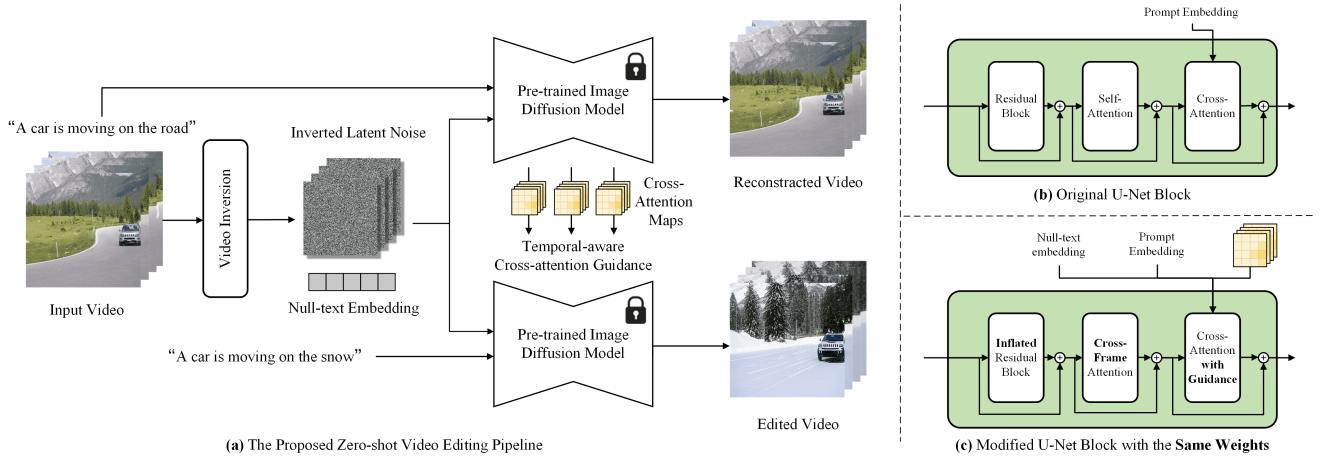


Figure 2: Illustration of the proposed vid2vid-zero for zero-shot video editing. (a) Our framework first inverts the input video to obtain the latent noise and null-text embedding, then use the inversion results to generate the edited video, under cross-attention guidance. Temporal modeling is achieved by replacing self-attention in the original U-Net block (b) with the cross-frame attention in (c) that shares the same model weights.

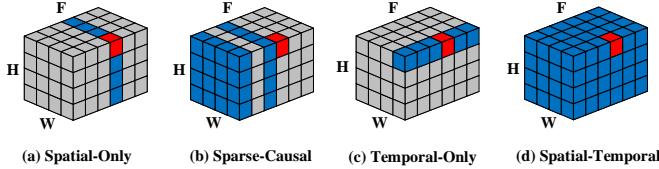


Figure 3: Illustration of different attention mechanisms. The query and keys are highlighted in red and blue, respectively.  $H$ ,  $W$ , and  $F$  represent the height, width, and temporal dimension of the input video.

original video. To solve the above problem, we resort to prompt-tuning [25] to learn a soft text embedding that aligns with the video content. Although more sophisticated methods in prompt-tuning can be used, we find that optimizing the null-text embedding [27] can achieve promising text-to-video alignment. Taking the DDIM inversion trajectory  $\{X_t^{\text{inv}}\}_{t=1}^T$  and the source prompt  $c$  as input, the optimized null-text embedding trajectory can be written as  $\{\emptyset_t\}_{t=1}^T$ . Specifically, the null-text embedding  $\emptyset_t$  is updated according to the following equation:

$$\min_{\emptyset_t} \|X_{t-1}^{\text{inv}} - f_\theta(\bar{X}_t, t, c; \emptyset_t)\|_2^2, \quad (1)$$

where  $f_\theta$  denotes applying DDIM sampling on top of the noise predicted by pre-trained image diffusion model  $\epsilon_\theta$ ,  $\bar{X}_t$  is the sampled latent code during null-text inversion. In practice, we share the same null-text embedding for different video frames, in order to preserve consistent information across video frames.

### 3.2. Temporal Modeling

Temporal modeling is essential for video editing. The dynamic computation nature makes the attention mechanism a suitable choice for temporal modeling at test time. To this end, Sparse-Causal Attention [50] (SC-Attn) is recently proposed to replace self-attention in pre-trained diffusion models, to capture the causal temporal information required for video generation. As shown in Fig. 3(b), it queries tokens in both the previous and the first frames for temporal modeling. While effective for generation, it is not optimal for the video editing task, as bi-directional temporal modeling is necessary to ensure temporally-coherent video editing results. An alternative choice for bi-directional temporal modeling is the temporal-only attention [17, 40], as shown in Fig. 3(c). However, it completely abandons spatial modeling, leading to a significant discrepancy with the self-attention in pre-trained image diffusion models (Fig. 3(a)). Our experiments in Sec. 4.3 indicate that this discrepancy leads to severe degradation in editing performance.

To reach a trade-off between bi-directional temporal modeling and alignment to the sampling process of the pre-trained image diffusion model, we propose spatial-temporal attention (ST-Attn) for video editing. As shown in Fig. 3(d), each frame  $x_i$  attends to all frames in the video. Specifically, query features are computed from all spatial features in the query frame  $x_i$  while key and value features are computed from all spatial features across all frames  $x_{1:T}$ .

Mathematically, the query, key, and value in the proposed ST-Attn can be written as:  $Q = W^Q x_i$ ,  $K = W^K x_{1:T}$ , and  $V = W^V x_{1:T}$ , respectively. Here,  $W^Q$ ,  $W^K$ , and  $W^V$  are the pre-trained projection weights in the self-attention layers, and are shared by all tokens in differ-

ent spatial and temporal locations.

Attending both previous frames and future frames with the proposed spatial-temporal attention provides bi-directional temporal modeling capability, while querying with spatial features across locations alleviates the discrepancy to the sampling process of pre-trained image diffusion models.

As shown in Fig. 2(c), we replace the pre-trained self-attention modules with cross-frame attention that shares exactly the same weights. In practice, We find that using ST-Attn at the beginning of the down-sampling, middle, and up-sampling block, and replacing the other self-attention layers with SC-Attn already works well for bi-directional temporal modeling. Besides, the original 2D residual block is inflated to 3D by copying the weights to each frame, to enable the inference on video inputs. We denote the model with updated U-Net blocks as  $\hat{\epsilon}_\theta$ .

### 3.3. Spatial Regularization

Maintaining fidelity to the original input video is an important aspect of the video editing task. To this end, we introduce spatial regularization to the editing process. The video inversion in Sec. 3.1 can largely ensure the reconstruction of the original video when sampling with the inversion results and the source prompt. And the cross-attn maps generated during reconstruction contain the spatial information of the original video [10]. For this reason, we can use the cross-attention maps as a spatial regularization and force the model to focus on the prompt-related areas via attention map injection. Formally, the noise prediction process can be written as:

$$\hat{z}_t = \hat{\epsilon}_\theta \left( \hat{X}_t, t, \hat{c}; \emptyset_t, M_t \right), \quad (2)$$

where  $M_t$  is the attention mask produced during the reconstruction of the input video. Here, we use  $\epsilon_\theta(\cdot; \emptyset_t, M_t)$  to denote the denoising steps with injected cross-attention map  $M_t$  and null-text embedding  $\emptyset_t$ . Thanks to our zero-shot temporal modeling methods in Sec. 3.2, the guidance cross-attention mask is temporal-aware, thus preventing sudden changes in cross-attention maps over different frames.

## 4. Experiments

### 4.1. Implementation Details

In experiments, we implement our method on top of the latent diffusion models [36] and use the publicly available Stable Diffusion model weights<sup>1</sup> by default. More results with different text-to-image diffusion model weights are presented in the Appendix. Following Tune-A-Video [50], each video contains 8 frames in  $512 \times 512$  resolution by default. During inference, we use DDIM sampler [41] with

<sup>1</sup><https://huggingface.co/CompVis/stable-diffusion-v1-4>

50 steps and classifier-free guidance [14] with a guidance scale of 7.5. For null-text inversion, we follow the hyper-parameters in [27], except that we set the inner step as 1, which we found already works well. The cross-attention injection threshold and the null-text embedding injection threshold are set as 0.8 and 0.5, respectively.

## 4.2. Main Results

We showcase the effectiveness of our zero-shot video editing method in Fig. 4. Specifically, two sets of video examples are included, whose text descriptions are “a man is running” and “A car is moving on the road”, respectively. More demonstrations are provided in the Appendix.

**Editing Style.** The style of a video is reflected by the global spatial and temporal characteristics of the video. In the 2nd row in Fig. 4, we add “anime style” to the input video. As can be seen, our method is capable of transforming all frames in the video to the target style, without alternating the semantics of the original video.

**Editing Attributes.** In the 3rd and 7th rows in Fig. 4, we demonstrate the ability of vid2vid-zero to edit attributes. For example, turning a young man into an old man and changing the model of the car. As can be seen, the edits are highly consistent with the target prompt. What’s more, they also maintain the other attributes of the original video, like the color and the pose of the car in the 7th row.

**Editing Background.** Rows 3, 4, 6, and 7 of Fig. 4 show the results of vid2vid-zero on editing background. As can be seen, it successfully turns the backgrounds to beach, complex urban street scenes, and the desert. These results demonstrate the ability of vid2vid-zero to retain the generative power of the pre-trained Stable-Diffusion, enabling creative video editing creations in a zero-shot manner. we found that when editing the background, the model will also change accordingly to make the overall effect more realistic. As row 6 in Fig. 4, when changing the prompt to “A car is moving on the snow”, the model also adds snow on the top of the car according to the scene.

**Replacing Subjects.** To demonstrate the effectiveness of our method on editing subjects, we turn the man into Stephen Curry in the 4th row of Fig. 4, and replace the horse with the dog in the 4th row of Fig. 6. As can be seen, our editing results not only align well with the text description but also maintain fidelity to the original videos. What’s more, our method can edit multiple properties at the same time, for example, replacing subjects (replacing the man with Curry) and editing background (changing the background to Time Square) at the same time.

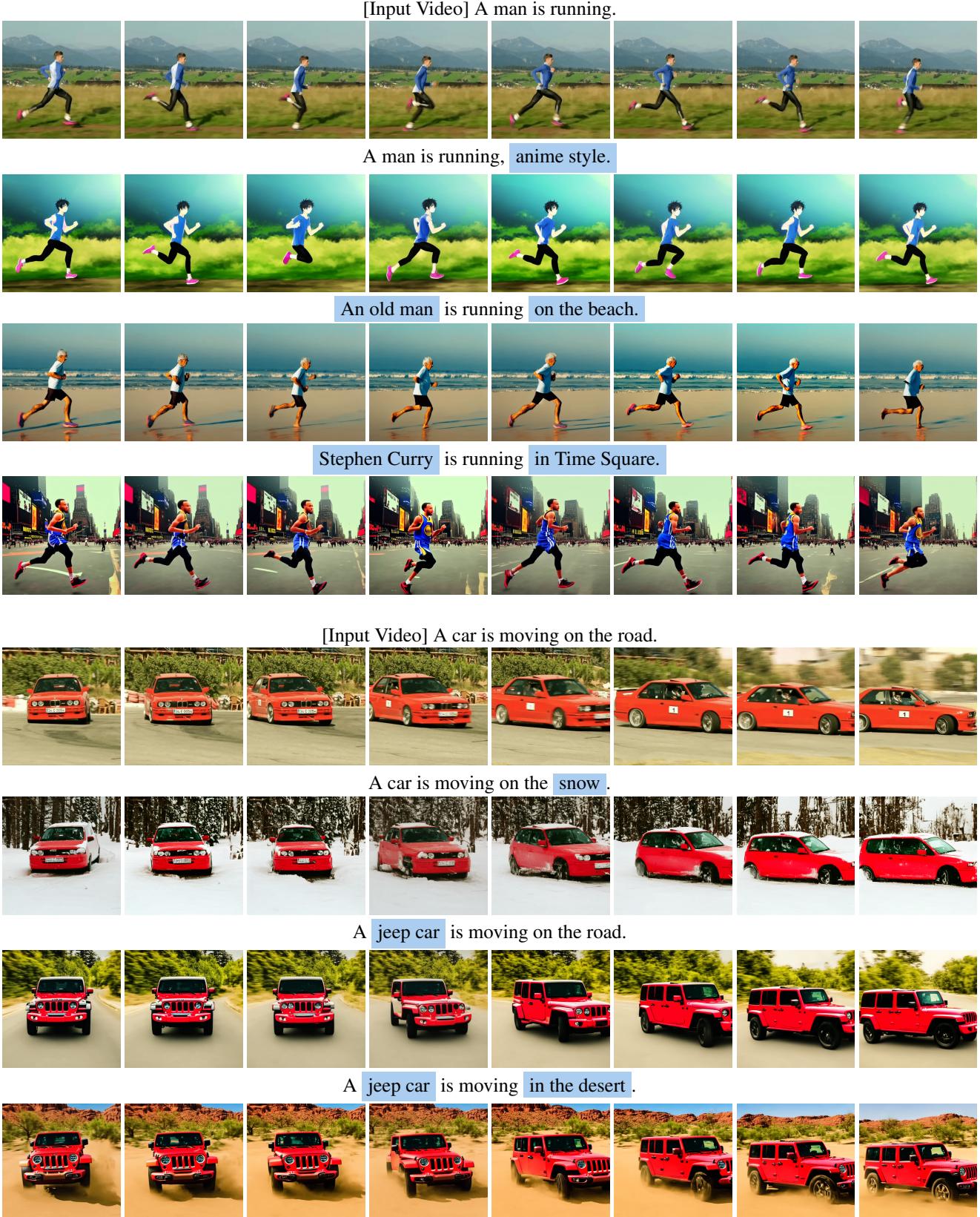


Figure 4: Videos editing results from various input videos and prompts. Our vid2vid-zero generates temporal consistent videos that align with the semantics of text prompts and faithfully to the original video.

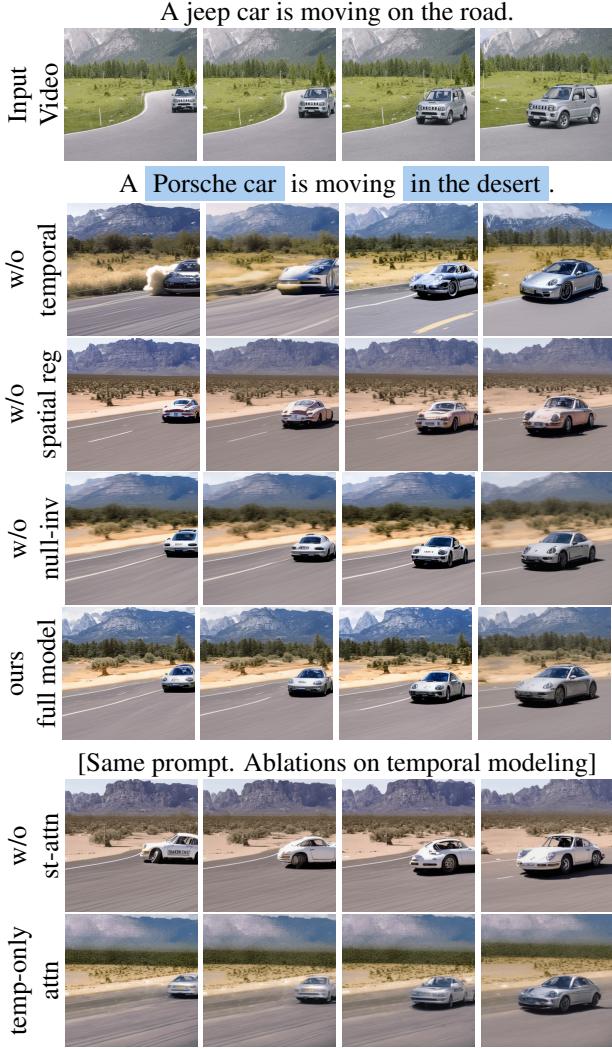


Figure 5: Ablations on the effectiveness of each component in vid2vid-zero (2nd ~ 5th row), and the temporal modeling design (6th ~ 7th row).

### 4.3. Analysis

**Ablations on Each Module.** We conduct ablation studies by isolating each component in vid2vid-zero, as shown in the 2nd to 4th row in Fig. 5. We have the following observation. Firstly, without the proposed temporal modeling, both the foreground car and the background mountain become inconsistent over time. Secondly, without spatial attention guidance, the edited video became less faithful to the input video. For example, the color of the car changes, and the background trees are missing. Thirdly, without null-text inversion, the background mountain and trees become blurry, since the optimized null-text embedding contains fine-grained details that align with the input video. Lastly, the above three components are complemen-



Figure 6: Comparison to other methods. Our vid2vid-zero achieves both temporal consistency and fidelity to the input video.

tary to each other, and combining them together gives rise to the best video editing result.

**Ablations on Temporal Attention Modeling.** Effective temporal modeling is the key to achieving zero-shot video editing. To this end, we perform further ablation experiments on the temporal modeling design. First, we consider removing the dense spatial-temporal attention that enables bidirectional temporal modeling and replacing it with Sparse-Causal Attention. As shown in the 6th row of Fig. 5, the car in the first few frames of the video is deformed. This is mainly because SC-Attention over-emphasizes the previous frames of the video, and the editing errors in the first few video frames propagate to the subsequent frames, leading to serious artifacts. Besides, we also present the results of replacing Dense-attn with Temporal-only Attention, which only focuses on temporal modeling while ignoring other spatial locations. As can be seen in the last row of Fig. 5, Temporal-only Attention can also alleviate the error propagation caused by the first few frames. However, it has a large gap compared with the pre-training self-attn, which on contrary only focuses on spatial location. As a result, there is a large distribution gap between the edited video frames and the output space of the pre-training text-to-image diffusion model, as indicated by the blurry video frames.



Figure 7: Visualization of the spatial-temporal attention. See Fig. 4 for the input and edited videos. The query is located near the right car light area in the third frame, highlighted by a white box.

**Spatial-Temporal Attention Visualization.** To gain a better understanding of the behavior of the proposed dense spatial-temporal attention, we visualize the attention maps in Fig. 7. The target prompt is “A jeep car is moving in the desert”, and the corresponding input video and edited video are shown in Fig. 4. The query is located at the right car light region in the third frame, highlighted by the white box. As can be seen, it successfully attends the car light areas in both previous and future frames, which demonstrates the effectiveness of our method in capturing bidirectional temporal and spatial information.

#### 4.4. Comparison

**Compared Methods.** *Tune-A-Video (TAV)* [50] realizes video generation via fine-tuning a pre-trained text-to-image diffusion model on a single video, and has shown potential applications in subject replacement, background change, attribute modification, and style transfer. *Plug-and-Play (PnP)* [45] is a state-of-the-art image editing method, which achieves high-quality image editing via feature replacement. We use PnP to edit the input video frame-by-frame and concatenate them to obtain the edited video.

**Qualitative Comparison.** The results of different methods on editing subjects are shown in Fig. 6. As can be seen, each frame produced by PnP shares similar poses to the corresponding input frame. However, it suffers from flickering effects like the sudden appearance of white bands in the first column and the deformation and repeated tails in the 2nd and 3rd columns. On the contrary, TAV is able to achieve relatively better temporal consistency, thanks to its temporal modeling and one-shot tuning. However, the pose and the motion in the edited video are not faithful to the original video. Our method can take the best of both worlds, and successfully turn the horse into a dog while maintaining both temporal consistency and faithfulness to the input video.

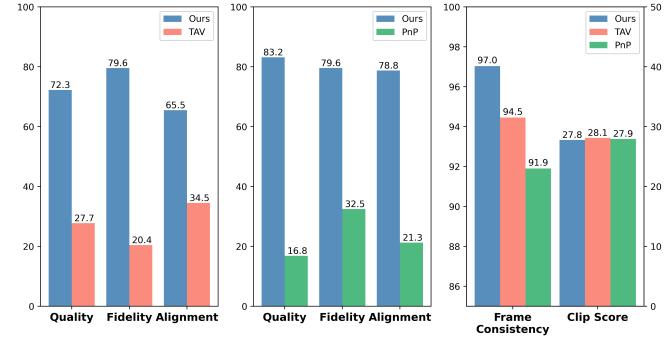


Figure 8: Quantitative results. User preference studies compared to TAV and PnP are shown in the first and second columns, respectively. CLIP score and frame consistency are presented in the third column.

**Quantitative Comparison.** (a) *User Preference.* Following [40, 50], we conduct user preference studies on 32 videos to compare our method with TAV and PnP, respectively. The users are asked to select editing better results based on (1) quality, (2) text-to-video alignment, and (3) fidelity to the original video. The results are shown in the first two columns in Fig. 8. As can be seen, we achieve higher scores in all three indicators when compared to PnP and TAV. (b) *CLIP score* [9, 11] measures text-to-video alignment by computing the similarity between the target prompt and the edited videos. As can be seen in the third column in Fig. 8, our method obtains an on-par CLIP score to that of TAV finetuned on the input video. We should note that the automatic evaluation metrics like CLIP score are not always consistent with human perception [28], and should only be used as an imperfect score for reference. (c) *Frame consistency* [9] measures the consistency of the generated video. We follow [9] to calculate the cosine feature similarity of consecutive frame pairs, based on their CLIP embedding, and take the average result as the score. As can be seen in the third column of Fig. 8, our method obtains the best score in these three methods, which demonstrates the superiority of our method in maintaining temporal consistency.

## 5. Conclusion and Discussion

In this paper, we propose vid2vid-zero, a method that achieves zero-shot video editing using off-the-shelf image diffusion models, without training on any video data. We leverage the dynamic nature of the attention mechanism to enable effective test-time temporal modeling. Experiments show that vid2vid-zero preserves the creativity of the image diffusion model and can obtain high-quality, text-aligned, and faithful video editing results. With the rapid development of generative models, we hope to bring novel insights

into the applications of these models.

**Limitations.** Since our method directly uses the pre-trained weights of the image diffusion model, it may inherit the limitations of the off-the-shelf image generation model. For example, our method lacks temporal and motion priors, because the off-the-shelf image diffusion model is not trained on any video data. As a result, it cannot be directly used to edit actions in videos, as reflected by our inability to effectively edit verbs in the source prompts.

## 6. Acknowledgments

This work was supported by National Key R&D Program of China (No. 2022ZD0118700).

## References

- [1] Yuval Alaluf, Or Patashnik, and Daniel Cohen-Or. Restyle: A residual-based stylegan encoder via iterative refinement. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6711–6720, 2021. [2](#)
- [2] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021. [2](#)
- [3] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *ICML*, volume 2, page 4, 2021. [2](#)
- [4] Daniel Bolya, Cheng-Yang Fu, Xiaoliang Dai, Peizhao Zhang, Christoph Feichtenhofer, and Judy Hoffman. Token merging: Your vit but faster. *arXiv preprint arXiv:2210.09461*, 2022. [11](#)
- [5] Duygu Ceylan, Chun-Hao Paul Huang, and Niloy J Mitra. Pix2video: Video editing using image diffusion. *arXiv preprint arXiv:2303.12688*, 2023. [3](#)
- [6] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2016: 19th International Conference, Athens, Greece, October 17–21, 2016, Proceedings, Part II* 19, pages 424–432. Springer, 2016. [2](#)
- [7] Guillaume Couairon, Jakob Verbeek, Holger Schwenk, and Matthieu Cord. Diffedit: Diffusion-based semantic image editing with mask guidance. *arXiv preprint arXiv:2210.11427*, 2022. [1, 2](#)
- [8] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021. [2](#)
- [9] Patrick Esser, Johnathan Chiu, Parmida Atighchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. *arXiv preprint arXiv:2302.03011*, 2023. [2, 3, 8](#)
- [10] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. [1, 2, 5](#)
- [11] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning. *arXiv preprint arXiv:2104.08718*, 2021. [8](#)
- [12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020. [2](#)
- [13] Jonathan Ho, Nal Kalchbrenner, Dirk Weissenborn, and Tim Salimans. Axial attention in multidimensional transformers. *arXiv preprint arXiv:1912.12180*, 2019. [2](#)
- [14] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022. [5](#)
- [15] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv:2204.03458*, 2022. [2, 3](#)
- [16] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022. [2](#)
- [17] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers, 2022. [4](#)
- [18] Yongcheng Jing, Xiao Liu, Yukang Ding, Xinchao Wang, Errui Ding, Mingli Song, and Shilei Wen. Dynamic instance normalization for arbitrary style transfer. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 4369–4376, 2020. [2](#)
- [19] Yongcheng Jing, Yang Liu, Yezhou Yang, Zunlei Feng, Yizhou Yu, Dacheng Tao, and Mingli Song. Stroke controllable fast style transfer with adaptive receptive fields. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 238–254, 2018. [2](#)
- [20] Yongcheng Jing, Yining Mao, Yiding Yang, Yibing Zhan, Mingli Song, Xinchao Wang, and Dacheng Tao. Learning graph neural networks for image style transfer. In *Computer Vision—ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part VII*, pages 111–128. Springer, 2022. [2](#)
- [21] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakkio Lehtinen, and Timo Aila. Analyzing and improving the image quality of StyleGAN. In *Proc. CVPR*, 2020. [2](#)
- [22] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022. [1, 2](#)
- [23] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in neural information processing systems*, 34:21696–21707, 2021. [2](#)
- [24] Huan Ling, Karsten Kreis, Daiqing Li, Seung Wook Kim, Antonio Torralba, and Sanja Fidler. Editgan: High-precision semantic image editing. *Advances in Neural Information Processing Systems*, 34:16331–16345, 2021. [2](#)
- [25] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroyuki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys*, 55(9):1–35, 2023. [4](#)

- [26] Shaoteng Liu, Yuechen Zhang, Wenbo Li, Zhe Lin, and Jiaya Jia. Video-p2p: Video editing with cross-attention control. *arXiv preprint arXiv:2303.04761*, 2023. 3
- [27] Ron Mokady, Amir Hertz, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Null-text inversion for editing real images using guided diffusion models. *arXiv preprint arXiv:2211.09794*, 2022. 3, 4, 5
- [28] Eyal Molad, Eliahu Horwitz, Dani Valevski, Alex Rav Acha, Yossi Matias, Yael Pritch, Yaniv Leviathan, and Yedid Hoshen. Dreamix: Video diffusion models are general video editors. *arXiv preprint arXiv:2302.01329*, 2023. 2, 3, 8
- [29] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. *arXiv preprint arXiv:2302.03027*, 2023. 1, 2, 3
- [30] Chenyang Qi, Xiaodong Cun, Yong Zhang, Chenyang Lei, Xintao Wang, Ying Shan, and Qifeng Chen. Fatezero: Fusing attentions for zero-shot text-based video editing. *arXiv preprint arXiv:2303.09535*, 2023. 3
- [31] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 2
- [32] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020. 2
- [33] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. 1, 2
- [34] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *ICML*, pages 8821–8831. PMLR, 2021. 2
- [35] Elad Richardson, Yuval Alaluf, Or Patashnik, Yotam Nitzan, Yaniv Azar, Stav Shapiro, and Daniel Cohen-Or. Encoding in style: a stylegan encoder for image-to-image translation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2287–2296, 2021. 2
- [36] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, pages 10684–10695, 2022. 1, 2, 5
- [37] Chitwan Saharia, William Chan, Huiwen Chang, Chris Lee, Jonathan Ho, Tim Salimans, David Fleet, and Mohammad Norouzi. Palette: Image-to-image diffusion models. In *ACM SIGGRAPH 2022 Conference Proceedings*, pages 1–10, 2022. 2
- [38] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022. 1, 2
- [39] Chaehun Shin, Heeseung Kim, Che Hyun Lee, Sang-gil Lee, and Sungroh Yoon. Edit-a-video: Single video editing with object-aware consistency. *arXiv preprint arXiv:2303.07945*, 2023. 3
- [40] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792*, 2022. 2, 4, 8
- [41] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020. 2, 3, 5, 11
- [42] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. In Hugo Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin, editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, 2020. 2
- [43] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. In *International Conference on Learning Representations*, 2021. 2
- [44] Omer Tov, Yuval Alaluf, Yotam Nitzan, Or Patashnik, and Daniel Cohen-Or. Designing an encoder for stylegan image manipulation. *ACM Transactions on Graphics (TOG)*, 40(4):1–14, 2021. 2
- [45] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. *arXiv preprint arXiv:2211.12572*, 2022. 1, 2, 3, 8, 11
- [46] Ruben Villegas, Mohammad Babaeizadeh, Pieter-Jan Kindermans, Hernan Moraldo, Han Zhang, Mohammad Taghi Saffar, Santiago Castro, Julius Kunze, and Dumitru Erhan. Phenaki: Variable length video generation from open domain textual description. *arXiv preprint arXiv:2210.02399*, 2022. 2
- [47] Wenhui Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyramid vision transformer: A versatile backbone for dense prediction without convolutions. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 568–578, 2021. 11
- [48] Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021. 2
- [49] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. NUwa: Visual synthesis pre-training for neural visual world creation, 2021. 2
- [50] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Weixian Lei, Yuchao Gu, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. *arXiv preprint arXiv:2212.11565*, 2022. 2, 4, 5, 8
- [51] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022. 2

## A. More Implementation Details

We provide more implementation details on the proposed vid2vid-zero. Specifically, for DDIM inversion [41], we perform inversion frame by frame and set the number of steps as 50. All experiments are conducted on a single Tesla V100 GPU with 32G memory. For the visualization of ST-Attn in Fig. 7 in the submission, we use features from the up-sampling block in the 40th DDIM sampling step, which we found is more semantically meaningful.

## B. More Analyses

### B.1. More Ablation Studies

In this section, we provide more detailed ablation studies on our method, as shown in Fig. 9. We use two samples for illustration, including “A Porsche car is moving in the desert.” and “A jeep car is moving in the snow” for the ablated experiments. The input video is shown in Fig. 5.

**Ablations on ST-Attn location.** As detailed in Sec. 3.2, we only use three ST-Attn modules at the first stage of the down-sampling, middle, and up-sampling blocks, respectively. In this section, we perform more detailed studies on this design choice. Specifically, we consider four other experiment setups, including (1) using ST-Attn modules at the middle stage of down-sampling, middle, and up-sampling blocks; (2) using ST-Attn modules at the last stage of down-sampling, middle, and up-sampling blocks; (3) using ST-Attn only in the first three stages in the down-sampling block; and (4) using ST-Attn only in the last three stages in the up-sampling block. The results are shown in Fig. 9. As shown, there are few differences among those settings, while our default setting shows slightly better temporal consistency and fewer artifacts.

**Ablations on null-text injection ratio.** As previously detailed, we set the null-text injection threshold as 0.5 by default. In this section, we show the results obtained by using different null-text injection thresholds. As can be seen, when the injection threshold is too high, the editing results suffer from severe artifacts, for example, the shape of the Porsche car is distorted when the injection threshold equals 1.0. We suspect that the null-text embedding contains information for reconstructing the input video, and injecting too much null-text embedding feature limited the editability of the proposed method. By contrast, when the injection threshold goes too low, the outputs get blurry, since the high-frequency details preserved in the null-text embedding are not fully utilized. To achieve a trade-off between editability and fidelity to the original video, we set the injection threshold as 0.5.

**Ablations on DDIM inversion with temporal modeling.** By default, we do not include temporal modeling during DDIM inversion. In this section, we compare the default setting with the results obtained from DDIM inversion with

temporal modeling. As shown in Fig. 9, excluding temporal modeling in DDIM inversion better preserves the spatial structure in each frame.

### B.2. More Visualisations on ST-Attn

In this section, we provide more visualization for the proposed ST-Attn, as shown in Fig. 10. The 5th to 7th rows show the visualization of features from down-sampling, middle, and up-sampling blocks, respectively. The query is located in the middle area of the tiger’s body, in the 3rd frame. The last row shows the visualization of the features from the up-sampling layer and the query is located in the middle area of the watermelon in the 3rd frame. We have the following observation. Firstly, the feature in the up-sampling block is more semantically meaningful, and can accurately locate the subject of interest. By contrast, the features in shallow layers are slightly worse at grouping features from semantically similar areas. This is consistent with the observations made in PnP [45]. Secondly, as can be seen in the last two rows, the ST-Attn map can accurately locate either the tiger or the watermelon when the corresponding query is related to each of these subjects. These results further prove that our ST-Attn can successfully capture the bi-directional relationship across different temporal frames and spatial locations.

### B.3. Discussions on ST-Attn Complexity

As detailed in Sec. 3.2, we use 3 ST-Attn modules in our vid2vid-zero model, and it achieves an ideal trade-off between effective temporal modeling and reducing gaps to the training time. While the dense spatial-temporal attention is not computationally efficient due to the quadratic complexity of the attention mechanism, we find that with the efficient xformer<sup>2</sup> attention implementation, the increase in memory is negligible. Without bells and whistles, our vid2vid-zero can be extended to input videos with 24 frames on a single Tesla V100 GPU with 32G memory.

For very long videos, several improvements can be included to reduce computational complexity while maintaining effective bi-directional temporal modeling. For example, (1) we can down-sample the key and value in spatial size [47] to significantly reduce the number of tokens for attention computation; (2) we can apply test-time token merge strategies like ToMe [4] to reduce the computational cost; (3) we can select neighborhood frames with fixed window size or select context video frames with strides.

## C. More Sample Results

In this section, we provide more visualization of the video editing results obtained by our vid2vid-zero. The results with are shown in Fig. 11, Fig. 12, and Fig. 13.

<sup>2</sup><https://github.com/facebookresearch/xformers>

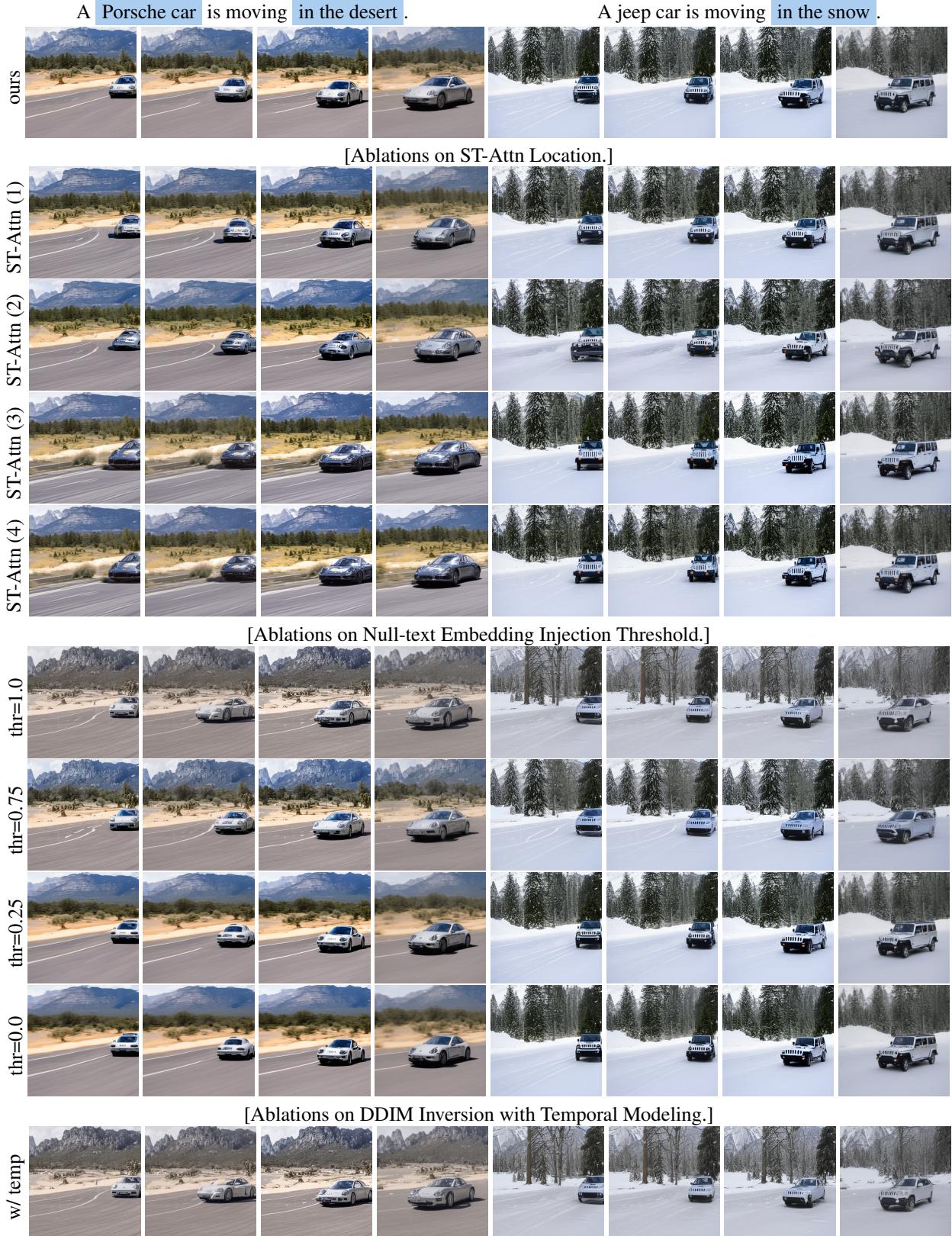
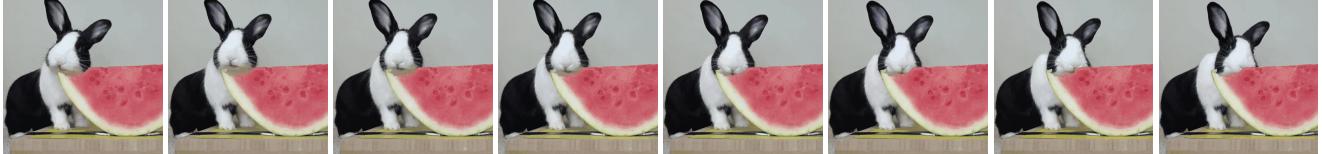
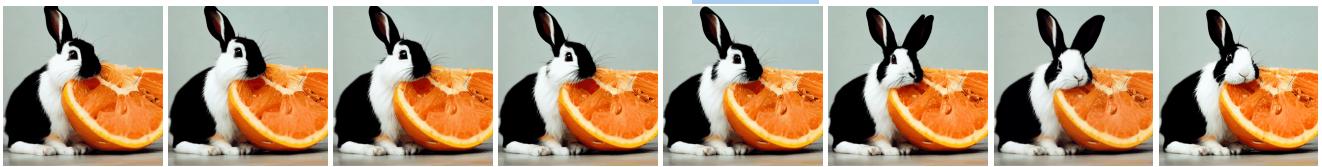


Figure 9: Ablations on ST-Attn location (2nd ~ 5th row), null-text injection threshold (6th ~ 9th row), and DDIM inversion with temporal modeling (the last row). The input video is shown in Fig. 5 in the main text. See the text for detailed settings.

[Input Video] A rabbit is eating a watermelon.



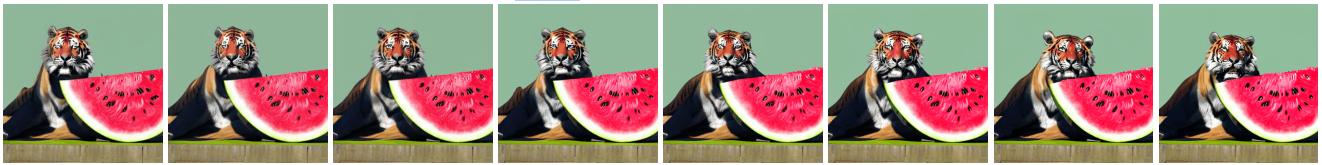
A rabbit is eating [an orange].



A [puppy] is eating [an orange].



A [tiger] is eating a watermelon.



[ST-Attn Visualizations on “A tiger is eating a watermelon.”]

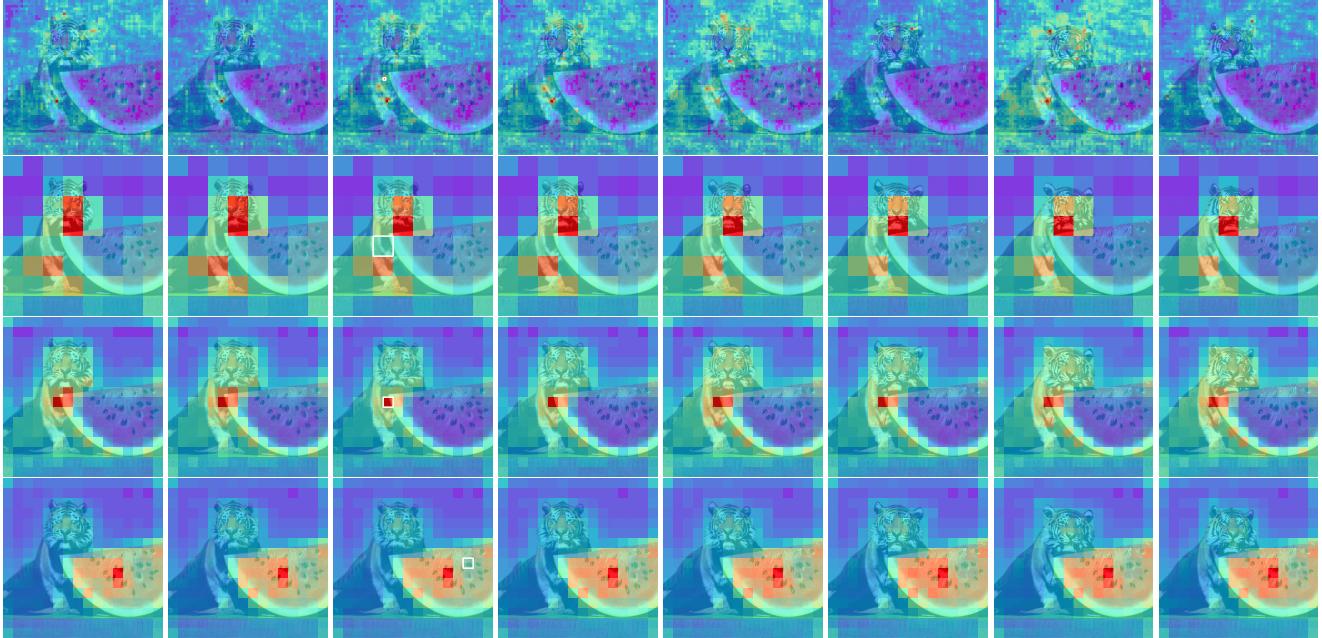
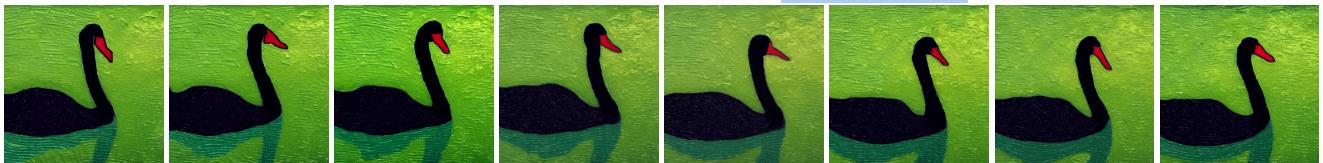


Figure 10: Videos editing results from various prompts (2nd ~ 4th row) and ST-Attn visualization (5th ~ 8th row). The 7th and 8th rows are visualization of the ST-Attn in the up-sampling block, while the 5th and 6th rows are visualization of the ST-Attn in the down-sampling and middle blocks, respectively. The query location is highlighted by the small white boxes.

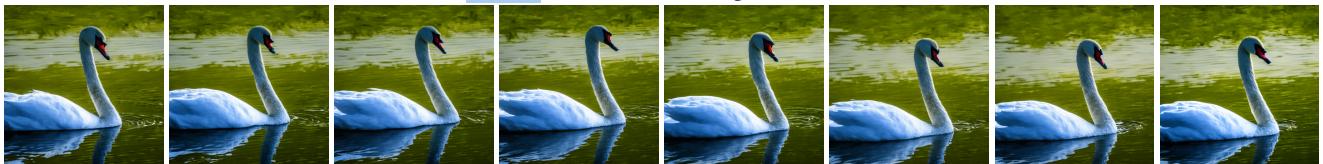
[Input video] A blackswan is swimming in the water.



A blackswan is swimming in the water, Van Gogh style.



A white swan is swimming in the water.



[Input video] A child is riding a bike on the road.



A lego child is riding a bike on the road.



A child is riding a bike on the flooded road.



[Input video] A man with a dog skateboarding on the road.



A man with a dog skateboarding on the desert.



Figure 11: Videos editing results from various input videos and prompts.



Figure 12: Videos editing results from various input videos and prompts.

[Input video] A man is surfing.



A boy is surfing.



Iron Man is surfing.



[Input video] A brown bear is playing on the ground.



A polar bear is playing on the grass.



A black bear is playing on the grass.



[Input video] A man is playing skateboard on the road.



A boy is playing skateboard on the road.



Figure 13: Videos editing results from various input videos and prompts.