# CCEdit: Creative and Controllable Video Editing via Diffusion Models

Ruoyu Feng[1,2] [*], Wenming Weng[1,2], Yanhui Wang[1,2],
Yuhui Yuan[2], Jianmin Bao[2], Chong Luo[2] [†], Zhibo Chen[1] [†], Baining Guo[2]
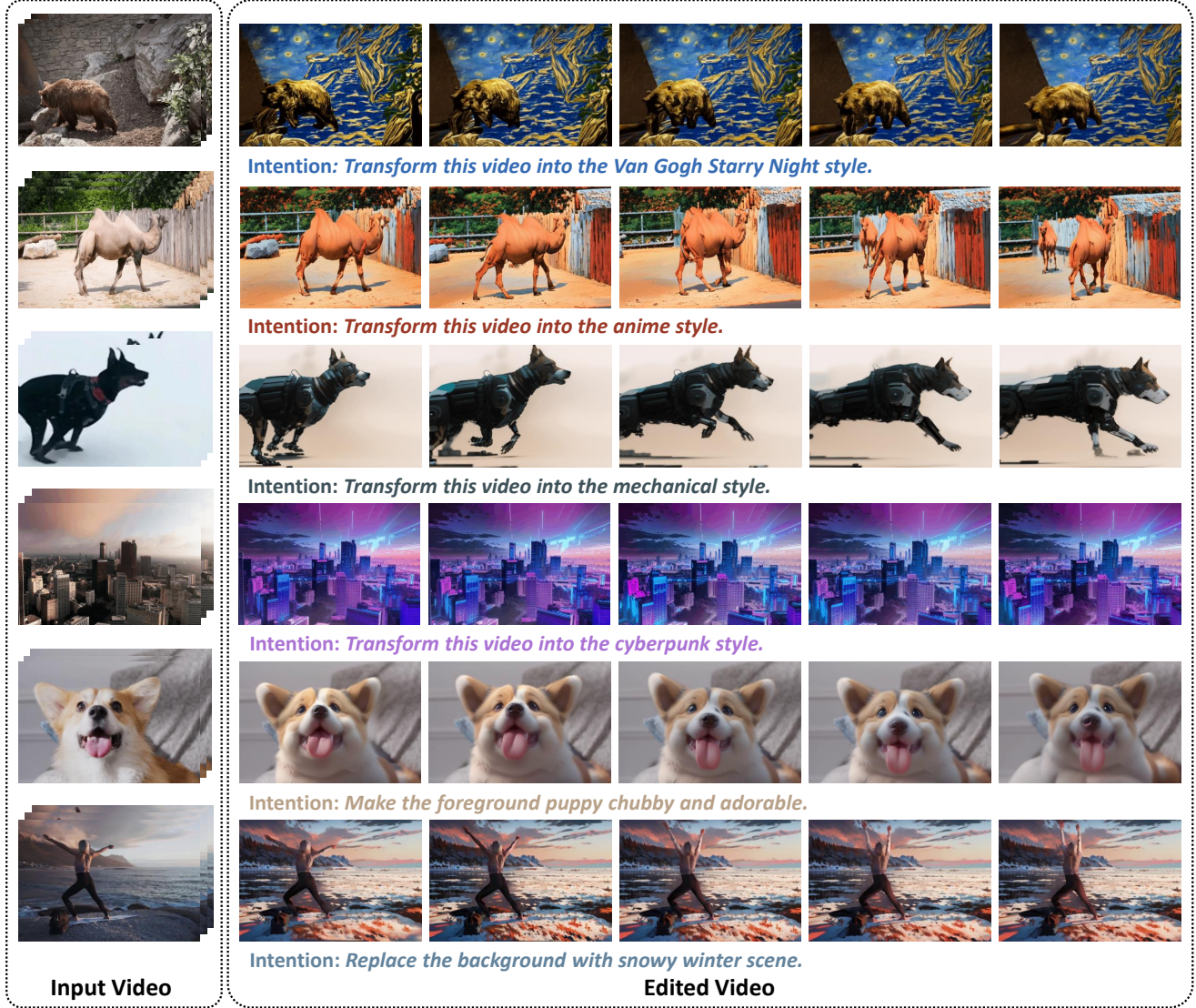[1]University of Science and Technology of China [2]Microsoft Research Asia

Figure 1. **The CCEdit framework accommodates a wide spectrum of user editing demands**: the first four rows showcase global style translation driven by various intentions, the fifth row illustrates foreground editing, and the sixth row demonstrates background modification.

---

[*] This work is done when Ruoyu Feng is an intern with MSRA.
[†] Corresponding author.

## Abstract

*In this work, we present CCEdit, a versatile framework designed to address the challenges of creative and controllable video editing. CCEdit accommodates a wide spectrum of user editing requirements and enables enhanced creative control through an innovative approach that decouples video structure and appearance. We leverage the foundational ControlNet architecture to preserve structural integrity, while seamlessly integrating adaptable temporal modules compatible with state-of-the-art personalization techniques for text-to-image generation, such as DreamBooth and LoRA. Furthermore, we introduce reference-conditioned video editing, empowering users to exercise precise creative control over video editing through the more manageable process of editing key frames. Our extensive experimental evaluations confirm the exceptional functionality and editing capabilities of the proposed CCEdit framework. Demo video is available at https://www.youtube.com/watch?v=UQw4jq-igN4.*

## 1. Introduction

In recent years, the domain of visual content creation and editing has undergone a profound transformation, driven by the emergence of diffusion-based generative models [8, 13, 34]. A large body of prior research has demonstrated the exceptional capabilities of diffusion models in generating diverse and high-quality images [23, 27, 30] and videos [3, 12, 33], conditioned by text prompts. These advancements have naturally paved the way for innovations in generative video editing [4, 20, 22, 38, 40, 41, 44].

Generative video editing, although in its nascent stages, faces a series of significant challenges. These challenges include accommodating diverse editing requests, achieving fine-grained control over the editing process, and harnessing the creative potential of generative models. Diverse editing requirements encompass tasks such as stylistic alterations, foreground replacements, and background modifications. Generative algorithms, while powerful, may not always align perfectly with the editor's intentions or artistic vision, resulting in a lack of precise control. Balancing creativity and controllability is crucial for a successful generative video editing system.

In response to these challenges, this paper introduces CCEdit[1], an innovative and comprehensive framework meticulously designed to strike a harmonious balance between controllability and creativity while accommodating a wide range of editing requirements. CCEdit unleashes the potential of video editing by decoupling the control of video structure and appearance during the editing process.

Specifically, we leverage both text prompt and personalized text-to-image (T2I) diffusion models [27] for appearance control. Building upon these pre-trained weights, we employ ControlNet [42] to inherit structural information from the source video. Users can choose from a variety of structural information types, such as line drawings [5], PIDI boundaries [35], and depth maps. For each type of structure information, we train distinct temporal layers within the main diffusion network while keeping the pretrained T2I model unaltered. These temporal layers seamlessly integrate structural control into the generation network while maintaining temporal consistency across edited video frames. This approach allows the temporal modules to function as adaptable components, accommodating various personalized model versions [14, 29] derived from the same text-to-image model, thereby enabling users to control appearance by selecting their desired image model.

However, there are occasions when users desire finergrained control over the appearance of the target video, considering more nuanced aspects. Furthermore, we empirically observed compatibility issues between the personalized T2I models and the newly trained temporal modules, especially with certain stylish LoRA models. In such cases, the video editing output may not perfectly maintain the style of the adopted image model. To address this, our framework introduces reference-aware video editing, an elegant mechanism that enables users to manually or AI-assistedly edit a keyframe and propagate those changes to other frames in the same video. Technically, we introduce an additional model to extract features from the edited reference image which are then incorporated into the features of the center frame extracted by the main network. The temporal modules propagate this newly introduced information throughout the entire video sequence, ensuring the generation of a coherent video.

To complete the system, we incorporate a frame interpolation method for editing high-frame-rate videos. Qualitative results demonstrate that CCEdit strikes a delicate balance between artistic creativity and operational controllability, empowering editors to tailor the editing process precisely to their unique requirements. The framework's compatibility with a diverse spectrum of generative models and the reference-aware approach underscore its adaptability and efficiency, heralding a future of enhanced generative video editing experiences.

## 2. Related Work

### 2.1. Diffusion-based Image and Video Generation

Diffusion models (DM) [8, 13, 34] have demonstrated exceptional capabilities in the field of image synthesis. These

---

[1] CCEdit is currently a research project, and there are no immediate intentions to integrate it into a product or extend public accessibility. Any future research endeavor will adhere to Microsoft's AI principles.

models indeed help by learning to approximate a data distribution through the iterative denoising of a diffused input. What makes DMs truly practical is the incorporation of text prompt as condition to control the output image during the generative process [19, 24, 27, 30]. Apart from the proliferation of advanced techniques in the field of image synthesis, DMs have also excelled in video generation [3, 12, 33]. This is accomplished by integrating modulated spatial-temporal modules, enabling the synthesis of high-quality videos while maintaining temporal consistency.

## 2.2. Video Editing with Diffusion Models

Recent studies leverage the inherent generative priors in DMs for image editing [2, 7, 11, 17, 21, 36]. The same idea is also applied in the field of video editing. Unlike image editing, video editing involves not only the manipulation of appearance-based attributes but also requires the meticulous preservation of temporal coherence throughout the video sequence. A lapse in maintaining this temporal coherence can result in visual artifacts, such as flickering and degradation. Some methods [22, 37, 41, 43] attempt to achieve temporal consistency without training by transitioning from spatial self-attention mechanisms within text-to-image (T2I) diffusion models to temporal-aware cross-frame attention techniques. Furthermore, certain methods [16, 32, 39, 44] optimize the parameters of the pre-trained T2I model based on the input video to attain temporal coherence of the target video. Nonetheless, the process of optimization for each input video can be time-consuming. Inadequate tuning of the temporal modules can result in suboptimal temporal coherence. Recent studies [10, 15, 40] have introduced supplementary trainable temporal layers to construct text-to-video (T2V) generative models. These models are trained on extensive text-video paired datasets, and their versatility is demonstrated in both video generation and video editing tasks. Different from previous methods, this paper disentangle the video editing process into distinct components, namely structure control and appearance control. Structure control aims to maintain the structure information of the input video, encompassing a wide range types from coarse to fine-grained, such as line drawings, PIDI boundaries [35], and depth maps. In terms of appearance control, with the goal of addressing creators' preferences comprehensively, we have introduced a spectrum of tools. These include text prompts, personalized T2I models, and edited reference images, allowing for meticulous control over the appearance of the output video. Additionally, we develop dedicated temporal layers to achieve precise alignment and coherence in the temporal dimension.

## 2.3. Personalized Text-to-Image models

The Stable Diffusion [27] model is trained on a huge dataset that encompasses a broad spectrum of domains [31].

Although the Stable Diffusion model is highly versatile and capable of generating a wide array of images, it occasionally falls short in specific details, particularly when it comes to generating human faces and hands, where subtle variations can markedly influence the overall perception. Additionally, it often struggles to precisely meet users' expectations for specific content, styles, and attributes. Therefore, personalized T2I models are designed to address these challenges. Two respective methods are DreamBooth [29] and LoRA [14]. The former uses a unique string as an indicator to represent the corresponding domain or concept during training. Once trained, this indicator can be employed to transfer the expectations to the fine-tuned T2I model. DreamBooth faces challenges due to the extensive weight parameters, making communication less convenient. To use less parameters and inherent the generalization of the base model, LoRA is proposed to fine-tune the model by freezing all original parameters and introducing the weight residuals $\Delta W$ to update the weights $W$. This process is formulated as $W' = W + \alpha \Delta W$, where $\alpha$ is the hyperparameter that controls the significance of the added $\Delta W$. Typically, the parameters of $\Delta W$ are significantly fewer than those of $W$. Finally, two additional methods for creating robust personalized T2I base models are fine-tuning the entire model directly on the self-collected datasets and blending parameters from various models.

Personalized T2I models play a pivotal role in establishing the ecosystem of contemporary AI content generation. They empower both novices and experienced artists, as well as enthusiasts, to rapidly and autonomously create stunning images and develop new models. One of the most significant objectives of our framework is to ensure compatibility with personalized T2I models, allowing creators to freely combine and perform highly creative edits on videos using models from the community. This is primarily attributed to our inheritance of the original spatial weights from Stable Diffusion model and the additional introduction of reference-aware editing operations.

## 3. Approach

### 3.1. Preliminary

**Diffusion Models** [13] are probabilistic generative models that approximate a data distribution $p(\mathbf{x})$ by gradually denoising a normally distributed variable. Specifically, DMs aim to learn the reverse dynamics of a predetermined Markov chain with a fixed length of $T$. The forward Markov chain can be conceptualized as a procedure of injecting noise into a pristine image, thereby transforming it into a stochastic representation, typically conforming to a standard Gaussian distribution. Empirically, DMs can be interpreted as an equally weighted sequence of denoising autoencoders $\epsilon_\theta(\mathbf{x}_t, t)$ where $t = 1, ..., T$. These autoen-
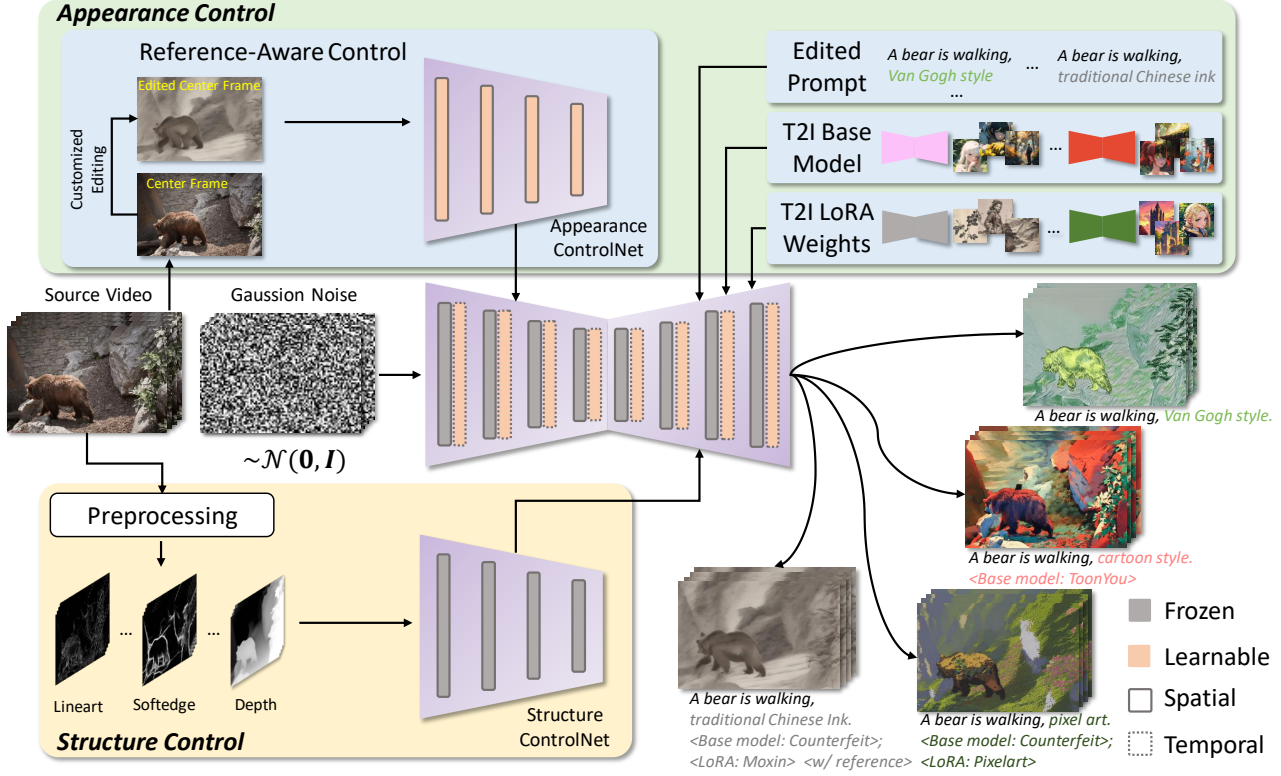
Figure 2. **Illustration of our overall framework.** Structure and appearance information in the target video are modulated independently. Structure control is conducted via the pre-trained ControlNet [42], encompassing a range from coarse to fine-grained. Appearance control is achieved by the flexible combination of text prompts, personalized T2I models, and edited center frame as reference. We omit the autoencoder and iterative denoising process for simplicity. Best viewed in color.

coders are trained to predict a denoised variant of the input $\mathbf{x}_t$. The corresponding objective can be simplified to:

$$\mathbb{E}_{\mathbf{x}, \epsilon \sim \mathcal{N}(0,1), t}[\|\epsilon - \epsilon_\theta(\mathbf{x}_t, t)\|_2^2]. \qquad (1)$$

**Latent Diffusion Models** (LDMs) are trained in the learned latent representation space $\mathbf{z}_t$, instead of the pixel space $\mathbf{x}_t$. The bridge between this latent space and the original pixel-level domain is established via a perceptual compression model. The perceptual compression model is composed of an encoder $\mathcal{E}$ and a decoder $\mathcal{D}$, where $\mathbf{z} = \mathcal{E}(\mathbf{x})$ and $\mathbf{x} \approx \mathcal{D}(\mathcal{E}(\mathbf{x}))$. Then the optimization objective in Equation (1) is modified as:

$$\mathbb{E}_{\mathcal{E}(\mathbf{x}), \epsilon \sim \mathcal{N}(0,1), t}[\|\epsilon - \epsilon_\theta(\mathbf{z}_t, t)\|_2^2]. \qquad (2)$$

### 3.2. Overall Framework

The architecture of our proposed framework is illustrated in Figure 2. Our objective is to create a comprehensive and user-friendly framework for video editing. To this end, we decouple the editing process into two distinct modules: structure control and appearance control.

Specifically, our framework is built on the Stable Diffusion [27] model. During the training phase, we freeze all weights of the original models to preserve both its robust visual synthesis and language comprehension capabilities. Structure control is achieved by incorporating visual priors such as line drawings [5], PIDI boundaries [35], and depth maps extracted from the source video. We use ControlNet to extract features and insert them into the base T2I model. For appearance control, we employ three different approaches: text prompts, personalized model weights, and customized center frame as reference. Most importantly, users can flexibly combine these elements according to their specific requirements to achieve the purpose of customizing the target video. Finally, the temporal coherence of the generated video is preserved through our introduced temporal consistency modules. This ensures a seamless and high-quality visual experience across the video sequence. Additionally, we also implement the temporal interpolation method in our framework for high frame rate.

Section 3.3 provides the implementation of structure control. Section 3.4 introduces the appearance control used in our framework. The detailed implementation of the newly introduced temporal consistency modules is presented in Section 3.5. Section 3.6 describes the method of

4

temporal interpolation.

## 3.3. Structure Control

We utilize the pre-trained ControlNet [42] to achieve the inherence of structure information of the input video. The ControlNet is initialized by copying the encoder of the T2I model before its training. The features extracted by ControlNet are integrated into the decoder of the frozen T2I base model, thereby controlling the structure of the generated image. Besides, the projection-in and projection-out layers are initialized with zeros before training to ensure stable and progressive updating. In our framework, we employ various types of structure information and corresponding ControlNets, ranging from coarse to fine, to ensure the flexibility of controlling the structure from varying degrees. To maintain the strong capability of the pre-trained ControlNet, we freeze the weights of it during training. Specifically, as illustrated in Figure 2, each frame of the input video undergoes preprocessing to obtain a structure representation. The representation is then input into ControlNet, which extracts features for each frame. These features are subsequently incorporated into the decoder of the main network.

## 3.4. Appearance Control

We propose different approaches to control the appearance of the target video, from the perspective of coarse-grained to fine-grained. Section 3.4.1 discusses the coarse-grained methods, which involve text prompts and personalized models. In Section 3.4.2, we delve into the fine-grained technique, *i.e.*, controlling the appearance of the video at the pixel level using the edited center frame as a reference.

### 3.4.1 Coarse-Grained Appearance Control

**Text Prompts.** Stable Diffusion [27] is a text-to-image (T2I) model trained on the large and diverse image-text paired dataset LAION [31], possessing strong language understanding capabilities. During the training process, we retain all weights of the original layers to preserve its robust abilities. As a result, the simplest method of video editing involves directly editing the text prompts. For instance, in Figure 2, the original text prompt is "a bear is walking." The edited text could be "a bear is walking, Van Gogh style." The appearance of the target video is controlled by the text prompt, the structure consistency is maintained by the pre-trained structure ControlNet, and temporal coherence is preserved by temporal consistency modules.

**Personalized Models.** Stable Diffusion excels as a general-purpose generative model for synthesizing realistic images, but it underperforms when tasked with generating contents within a specific domain or depicting a particular identity.
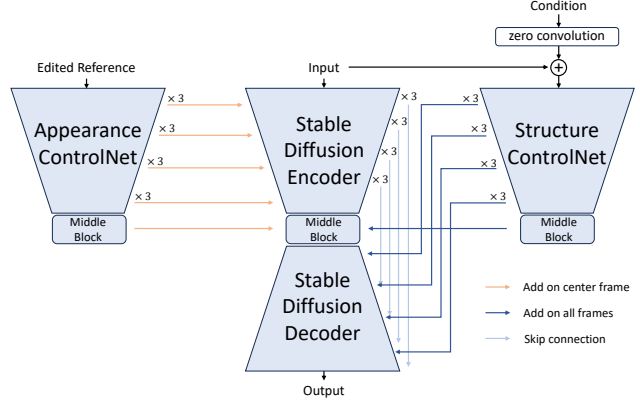


Figure 3. **Illustration of the reference-aware pipeline in our framework.** Text prompts and time embedding are incorporated into all modules, which are omitted for simplicity.

DreamBooth [29] and LoRA [14] are two representative techniques for generating custom styles, concepts, identities, *etc*. Recently, the Stable Diffusion community has seen rapid advancements, with contributors fine-tuning the T2I model on specialized datasets and sharing the pre-trained weights online [6, 9]. Users can readily substitute the original T2I model's weights with those personalized ones to generate images tailored to specific domains. As long as the personalized models and LoRAs are trained from the same base model (*e.g.*, Stable Diffusion-v1.5), they can exhibit good compatibility with each other. Different LoRAs can function like plug-ins on various base models. In our framework, as shown in Figure 2, the flexibility to change the weights of the T2I base model and LoRAs serves as another avenue for appearance control. It's also achieved by freezing the weights of the original Stable Diffusion model during the training phase. This enables the generation of video content tailored to specific styles and domains.

### 3.4.2 Fine-Grained Appearance Control

**Customized Center Frame.** At times, professional video editors desire to edit and produce videos strictly according to their own requirements. Additionally, we empirically find that the trained temporal modules can sometimes suppress the performance of the personalized T2I base model and LoRAs. Therefore, we propose the reference-aware mechanism in our framework for controlling the appearance of the target video from the fine-grained perspective. In the workflow, users can initiate editing on the center frame initially [11, 18, 28, 42], allowing for meticulous and comprehensive control over the detailed appearance. More specifically, as shown in Figure 3, the center frame of the original video is edited and fed into the Appearance ControlNet. Subsequently, the features extracted from each layer are added to the corresponding features of the center frame in the base model's encoder side. Consequently, the appear-

ance information within the center frame is propagated to all frames through the temporal modules, achieving the stable and controllable editing. During training, parameters of newly introduced Appearance ControlNet and temporal consistency modules are jointly optimized. It's worth noting that different appearance control methods are not necessarily used in isolation. They can be flexibly selected and combined according to the requirements to achieve different purposes.

### 3.5. Temporal Consistency Modules

**Pseudo-3D Layers Design.** Inspired by previous works that turn the pre-trained Stable Diffusion [27] into video generative models and meanwhile protecting the strong content synthesising capability and language understanding capability [3, 10], we freeze weights of all layers of original image generation model and only update newly introduced temporal consistency modules to achieve the motion smoothness and content coherence. Specifically, 2D Resblocks and attention layers are transformed into pseudo-3D ones. As illustrated in Figure 4 (a) and (b), a 1D version of corresponding 2D transform layer is appended, which is conducted on temporal dimension. Formally, suppose the spatial layer and appended temporal layer are represented by $l_\theta^i$ and $l_\phi^i$ parameterized by parameters $\theta$ and $\phi$, where $i$ indicates the index of corresponding layer. Given the feature $\mathbf{z}^i \in \mathbb{R}^{B \times C \times T \times H \times W}$ of the video clip, where $B$ indicates batch size, $C$ denotes feature channels, $T$ denotes sequence length, $H$ and $W$ presents the height and width of the latent feature. The operation can be presented as follows (using einops [26] notation):

$$\mathbf{z}^i \leftarrow \texttt{rearrange}(\mathbf{z}^i, \texttt{b c t h w} \rightarrow \texttt{(b t) c h w})$$
$$\mathbf{z}^i \leftarrow l_\theta^i(\mathbf{z}^i, \mathbf{c}^i)$$
$$\mathbf{z}^i \leftarrow \texttt{rearrange}(\mathbf{z}^i, \texttt{(b t) c h w} \rightarrow \texttt{(b h w) c t})$$
$$\mathbf{z}^i \leftarrow l_\phi^i(\mathbf{z}^i, \mathbf{c}^i)$$
$$\mathbf{z}^i \leftarrow \texttt{rearrange}(\mathbf{z}^i, \texttt{(b h w) c t} \rightarrow \texttt{b c t h w}),$$

where $\mathbf{c}^i$ indicates the condition including time step and text prompt.

Besides, we also apply skip connection and zero-initialized projection out layer for each newly introduced module for progressive and stable updating.

**Anchor-Aware Attention Module.** For the technique of taking customized center frame as reference for appearance control in Section 3.4.2, we additionally introduce the anchor-aware attention module to enhance the control by the reference center frame, which is shown in Figure 4 (c). Specifically, each token calculates the corresponding key and value features based on the concatenation of the current frame and the center frame. Formally, suppose that the
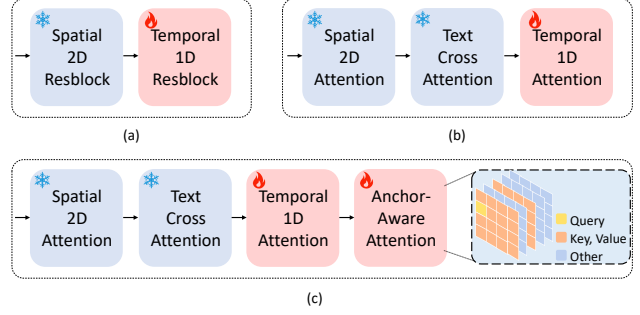


Figure 4. **Details of proposed temporal modules.** (a) The modified temporal-aware Resblock. (b) The modified temporal-aware attention block. (c) The modified anchor-aware attention block. We apply skip connection and zero-initialized projection-out layer for each newly introduced module.

latent feature $\mathbf{z}$ (we omit the layer index $i$ for simplicity), $\mathbf{z}_j$ and $\mathbf{z}_k$ represent features of the $j$-th frame and the center frame. The anchor-aware attention is formulated as follows:

$$Q = W^Q \mathbf{z}_j, K = W^K[\mathbf{z}_j; \mathbf{z}_k], V = W^V[\mathbf{z}_j; \mathbf{z}_k], \quad (3)$$

where $W^Q$, $W^K$, $W^V$ are projection layers, and $[\cdot]$ indicates concatenation operation. Besides, we copy the original spatial self-attention weights as initialization, with the exception of zero-initializing the projection-out layer.
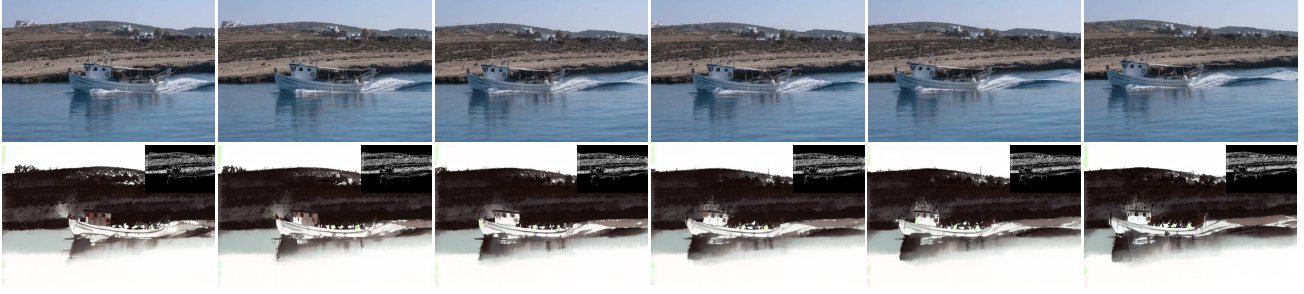
### 3.6. Temporal Interpolation

For video editing, a high frame rate of the edited video plays a crucial role in enhancing visual continuity. The techniques described in Section 3.2, Section 3.4, and Section 3.5 are capable of conducting edits on *key frames* with customized changes. However, this editing process is constrained by relatively low frame rates due to memory limitations. To achieve this, we further introduce the temporal interpolation model to interpolate intermediate frames between every two key frame. Similarly to the information incorporation mechanism for the edited center frame as a reference introduced in Section 3.4.2, we use the Appearance Controlnet to extract the features of the two given key frames and then add them to the features of the first and last frame, respectively. The structure information of both key frames and intermediate frames is input into the Structure ControlNet as structural guidance. Then the temporal modules propagate appearance features from both ends towards the center.

## 4. Experiments

### 4.1. Implementation Details

**Training.** Stable Diffusion-v1.5 is used as the base model to train temporal consitency modules and Appearance

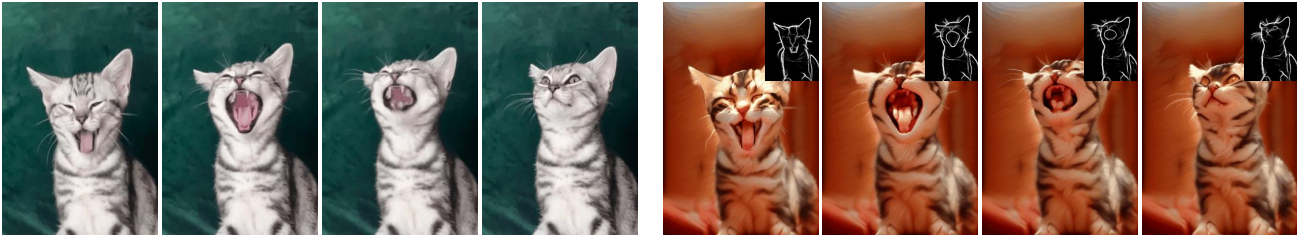*<ReV Animated, Moxin>* A boat is sailing, *Chinese traditional ink style.*

*< MeinaMix, Pixel Art Style>* A bear is walking, *pixel art style.*

*<ReV Animated, CAGAlienSpaceships> An alien spacecraft* is traveling through space.

*< A-Zovya Photoreal, mechanical dog> A mechanical Doberman* sits solemnly on the ground.

*< hellonijicute25d>* A little yawning cat, *anime style.*

Figure 5. **Results of video style translation.** Our method can conduct style transferring by flexibly combining different kinds of structure information inherence and appearance control. $\langle \cdot \rangle$ indicate the personalized T2I model we used.

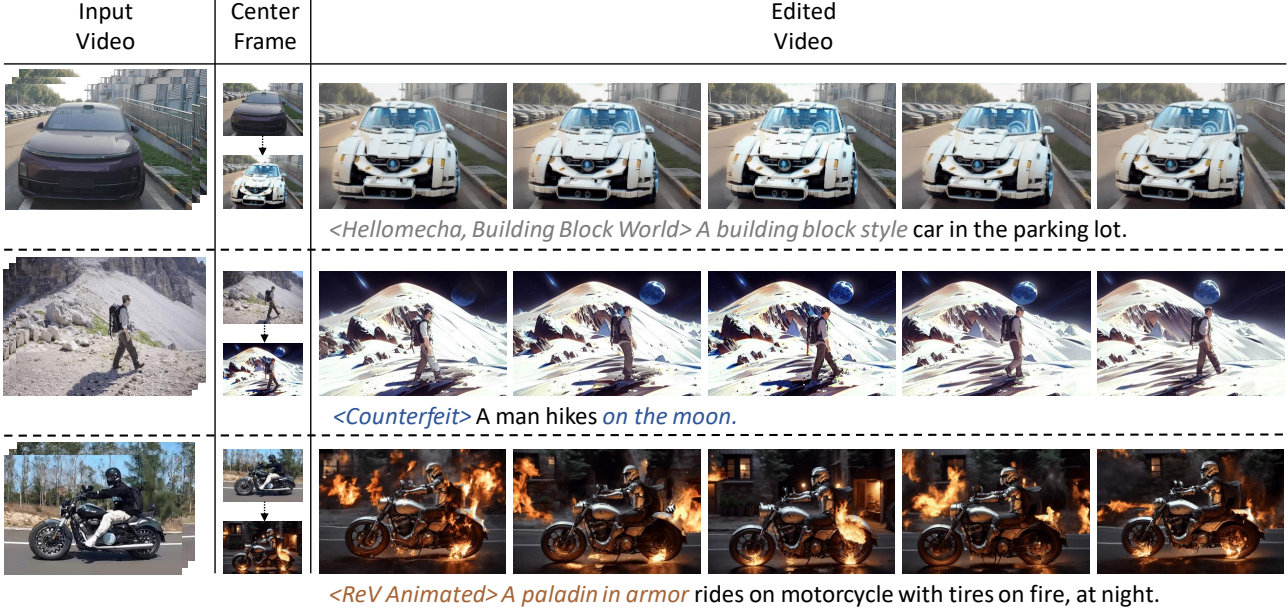|  | Input Video | Center Frame | Edited Video |
|--|--|--|--|

Figure 6. **Video editing results with customized center frame as reference.** The first row corresponds to customizing foreground, the second row corresponds to customizing background, and the third row is taking given reference image to affect the entire picture. ⟨·⟩ indicate the personalized T2I model we used.



Figure 7. **Video editing results under different structural guidance.** The types of structure information illustrated in this figure include "Lineart (line drawings [5])", "Softedge (PIDI boundaries [35])", "Scribble (human sketching)", "Depth". We use the personalized T2I model "*ReV Animated*" and LoRA "*kMechAnimal*". The prompt for all samples is "*Mech4nim4lAI, tiger, robotic.*"

ControlNet. We use the pre-trained ControlNet [42] for the structure information guidance. The training dataset is a self-collected private dataset. We trained the temporal consistency modules and Appearance ControlNet for various types of structural information, including line drawings [5], PIDI boundaries [35], depth maps detected by Midas [25], and human scribbles. During the training process, we first resize the shorter side to 384 pixels, followed by a ran-
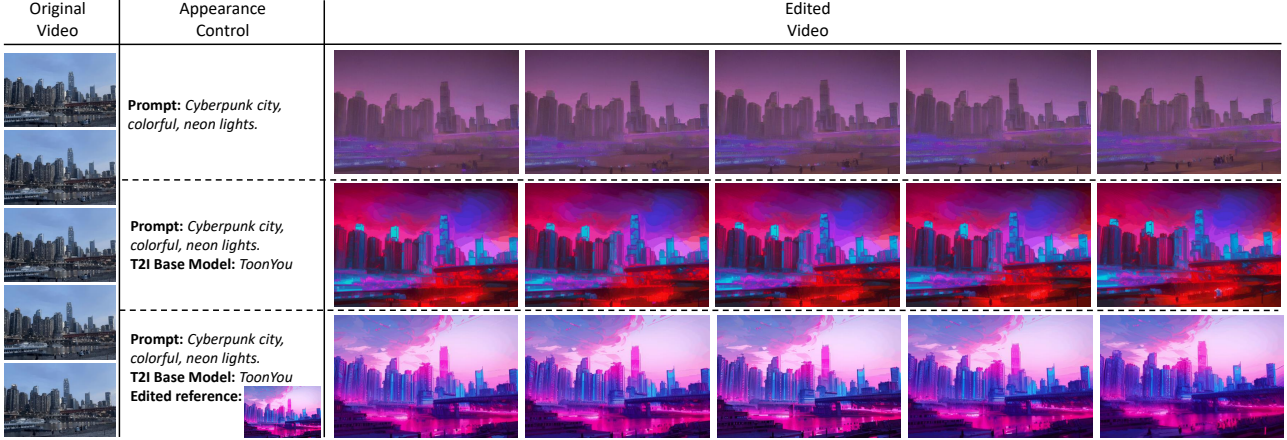
8

Figure 8. **Video editing results under different appearance control.** Dependence on the original Stable Diffusion pre-trained model alone failed to produce satisfactory results. Even after transitioning to a personalized T2I base model, the generated outcomes remained suboptimal. Only when we initially edited the center frame and take it as the appearance reference that we achieved the expected results.

| Model Name | Type |
|---|---|
| Counterfeit | T2I Base Model |
| ToonYou | T2I Base Model |
| ReV Animated | T2I Base Model |
| HelloMecha | T2I Base Model |
| hellonijicute25d | T2I Base Model |
| A-Zovya Photoreal | LoRA |
| kMechAnimal | LoRA |
| Pixel Art Style | LoRA |
| fat animal | LoRA |
| Building Block World | LoRA |
| MoXin | LoRA |
| mechanical dog | LoRA |

Table 1. Personalized models utilized in the evaluation, all sourced from CivitAI [6].

dom crop to obtain video clips with a size of 384×576. 17 frames at 4 fps are sampled from each video. The batch size is 16 and the learning rate is $3e - 5$.

**Evaluation.** We validate the effectiveness and functionality of our approach through the use of diverse combinations of structure control and appearance control. Structure control encompasses line drawings, PIDI boundaries [35], depth maps detected by Midas, and human scribbles, in a fine-to-coarse manner. For appearance control, we collect several personalized T2I base models and LoRA weights from CivitAI [6] and explored different combinations, which are illustrated in Table 1. Similar to previous work [10], we employ the "trigger words" to activate these personalized models.

## 4.2. Applications

**Video Style Translation.** The basic function of our framework lies in video style translation. By flexibly combining various types of structural and appearance control, users are empowered to translate videos into styles that align with their preferences. We show some samples in Figure 5, including anime style, pixel art style, and mechanical style, *etc*. These results demonstrate the superior coordination of our method in integrating various types of structure control and appearance control, affirming the compatibility and controllability of the framework.

**Customized Video Editing.** Sometimes, users require stronger control over the content they want to generate. For example, they may want to change only the foreground, alter just the background, or edit the texture content of a video in a specific way. In our framework, this can be achieved by first using image editing techniques to modify the content of the intermediate frame and then using it as a reference for editing the entire video. As depicted in Figure 6, we first edit the center frames of the videos by Stable Diffusion Web UI [1], followed by utilizing these edited central frames as guides for the video editing process. Thanks to the joint training of Appearance ControlNet and temporal consistency modules, our method can effectively propagate the edited information of the center frame to the whole video in a coherent manner.

## 4.3. Ablation Study

**Structure Control.** As shown in Figure 7, we also experimented with how different structure control methods affect the results. Line drawings (2nd row) and PIDI boundaries (3rd row) provide detailed control over the structure, making them more suitable for use cases that require a signifi-

9

cant preservation of the structure. On the other hand, scribbles (4th row) and depth maps (5th row) offer a coarser level of control over the structure, making them suitable for more imaginative use cases.

**Appearance Control.** Figure 8 illustrate the importance of taking the customized center frame as a reference in certain scenarios. Initially, translating the video scenes into "`cyberpunk`" style (1st row) solely through prompt adjustments appears challenging, as this word may be unfamiliar to the pre-trained Stable Diffusion and the temporal consistency modules. Even after loading a T2I base model that is more familiar with such terminology, it still doesn't seem to achieve the desired outcome effectively. This could potentially be due to the fact that the temporal consistency modules may not be highly compatible with just any T2I base model. The solution is illustrated in the third row in Figure 8, that is, first edit the center frame and then combine it with other appearance control factors to manage the appearance of the video.

## 5. Limitation and Future Works

**Out-of-Distribution Personalized Models.** As shown in Figure 9, our approach encounters challenges in generating satisfactory results when using personalized models that are too out of distribution. This issue can probably be tackled by collecting relevant text-video paired data and fine-tuning the temporal modules.
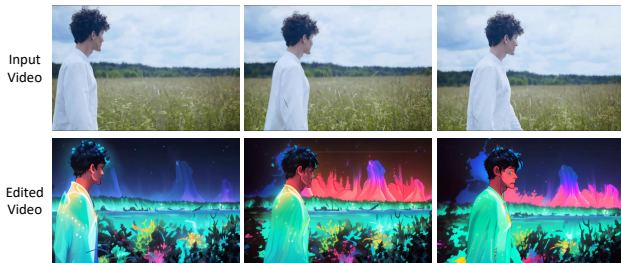


Figure 9. Poor temporal coherence can occur when the target video deviates significantly from the training data distribution. In this experiment, we employ the T2I base model of *"ToonYou"*, and the text prompt provided is: *"A man walking in the wonderland, lush, bioluminescent, vibrant, dense, unique, towering, floating, colorful, enormous, intricate, expansive, rich biodiversity, exotic, mystical."*

**Heavy ControlNet.** In this paper, we directly utilize the pre-trained ControlNet to perform structure control. However, due to ControlNet's evolution from the original Stable Diffusion encoder, it possesses a substantial parameter count (more than 400M). Furthermore, for reference-aware scenarios, two ControlNets are employed, with each being responsible for extracting structure and appearance information, further increasing the parameters of the overall

model. This may be unnecessary and could lead to issues such as increased GPU memory consumption and longer inference time. In the future, we plan to explore the adoption of more lightweight structure control networks [18, 28] to address these concerns.

**Flickering Problem.** We observed flickering in videos with higher frame rates or after frame interpolation, especially noticeable in high-frequency fine texture details. This is primarily attributed to our video editing operations being performed in the latent domain encoded by the 2D autoencoder. Introducing additional temporal layers in the autoencoder [3] is an promising way to solve this problem.

## 6. Conclusion

In this paper, we propose a unified and practical framework for diffusion-based video editing. Our key design choice is to decouple the video editing process into structure control and appearance control. We utilize the pre-trained ControlNet for structure control, and we achieve appearance control through three different ways, *i.e.*, text prompt, personalized model weights, and customized center frame as reference. For the appearance control of using a customized center frame, we innovatively introduced the Appearance ControlNet. This network is responsible for extracting features from the edited center frame and seamlessly integrating them into the base model. Additionally, we introduce temporal consistency modules to maintain temporal coherence. Last but not least, we also propose the temporal interpolation model to generate high frame rate videos. Our work offers a comprehensive and flexible set of tools for creative and controllable video editing. We hope that our work can benefits practitioners in the context of creative and controllable video editing.

## References

[1] AUTOMATIC1111. Stable Diffusion Web UI, Aug. 2022. 9

[2] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18208–18218, 2022. 3

[3] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 2, 3, 6, 10

[4] Wenhao Chai, Xun Guo, Gaoang Wang, and Yan Lu. Stablevideo: Text-driven consistency-aware diffusion video editing. *arXiv preprint arXiv:2308.09592*, 2023. 2

[5] Caroline Chan, Frédo Durand, and Phillip Isola. Learning to generate line drawings that convey geometry and semantics. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7915–7925, 2022. 2, 4, 8

[6] Civitai. Civitai. https://civitai.com/, 2022. 5, 9

[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. 3

[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] Hugging Face. Hugging face. https://huggingface.co/, 2022. 5

[10] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning. *arXiv preprint arXiv:2307.04725*, 2023. 3, 6, 9

[11] 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. 3, 5

[12] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2, 3

[13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020. 2, 3

[14] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021. 2, 3, 5

[15] Jun Hao Liew, Hanshu Yan, Jianfeng Zhang, Zhongcong Xu, and Jiashi Feng. Magicedit: High-fidelity and temporally coherent video editing. *arXiv preprint arXiv:2308.14749*, 2023. 3

[16] 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

[17] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. *arXiv preprint arXiv:2108.01073*, 2021. 3

[18] Chong Mou, Xintao Wang, Liangbin Xie, Jian Zhang, Zhongang Qi, Ying Shan, and Xiaohu Qie. T2i-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion models. *arXiv preprint arXiv:2302.08453*, 2023. 5, 10

[19] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. *arXiv preprint arXiv:2112.10741*, 2021. 3

[20] Hao Ouyang, Qiuyu Wang, Yuxi Xiao, Qingyan Bai, Juntao Zhang, Kecheng Zheng, Xiaowei Zhou, Qifeng Chen, and Yujun Shen. Codef: Content deformation fields for temporally consistent video processing. *arXiv preprint arXiv:2308.07926*, 2023. 2

[21] Gaurav Parmar, Krishna Kumar Singh, Richard Zhang, Yijun Li, Jingwan Lu, and Jun-Yan Zhu. Zero-shot image-to-image translation. In *ACM SIGGRAPH 2023 Conference Proceedings*, pages 1–11, 2023. 3

[22] 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. 2, 3

[23] 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. 2

[24] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *International Conference on Machine Learning*, pages 8821–8831. PMLR, 2021. 3

[25] René Ranftl, Katrin Lasinger, David Hafner, Konrad Schindler, and Vladlen Koltun. Towards robust monocular depth estimation: Mixing datasets for zero-shot cross-dataset transfer. *IEEE transactions on pattern analysis and machine intelligence*, 44(3):1623–1637, 2020. 8

[26] Alex Rogozhnikov. Einops: Clear and reliable tensor manipulations with einstein-like notation. In *International Conference on Learning Representations*, 2021. 6

[27] 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. 2, 3, 4, 5, 6

[28] Denis Zavadski Carsten Rother. Controlnet-xs. 2023. 5, 10

[29] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22500–22510, 2023. 2, 3, 5

[30] 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. 2, 3

[31] 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. 3, 5

[32] 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

[33] 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, 3

[34] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 2

[35] Zhuo Su, Wenzhe Liu, Zitong Yu, Dewen Hu, Qing Liao, Qi Tian, Matti Pietikäinen, and Li Liu. Pixel difference networks for efficient edge detection. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 5117–5127, 2021. 2, 3, 4, 8, 9

[36] Narek Tumanyan, Michal Geyer, Shai Bagon, and Tali Dekel. Plug-and-play diffusion features for text-driven image-to-image translation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1921–1930, 2023. 3

[37] Wen Wang, Kangyang Xie, Zide Liu, Hao Chen, Yue Cao, Xinlong Wang, and Chunhua Shen. Zero-shot video editing using off-the-shelf image diffusion models. *arXiv preprint arXiv:2303.17599*, 2023. 3

[38] Yuanzhi Wang, Yong Li, Xin Liu, Anbo Dai, Antoni Chan, and Zhen Cui. Edit temporal-consistent videos with image diffusion model. *arXiv preprint arXiv:2308.09091*, 2023. 2

[39] 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. 3

[40] Zhen Xing, Qi Dai, Han Hu, Zuxuan Wu, and Yu-Gang Jiang. Simda: Simple diffusion adapter for efficient video generation. *arXiv preprint arXiv:2308.09710*, 2023. 2, 3

[41] Shuai Yang, Yifan Zhou, Ziwei Liu, and Chen Change Loy. Rerender a video: Zero-shot text-guided video-to-video translation. *arXiv preprint arXiv:2306.07954*, 2023. 2, 3

[42] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *arXiv preprint arXiv:2302.05543*, 2023. 2, 4, 5, 8

[43] Yabo Zhang, Yuxiang Wei, Dongsheng Jiang, Xiaopeng Zhang, Wangmeng Zuo, and Qi Tian. Controlvideo: Training-free controllable text-to-video generation. *arXiv preprint arXiv:2305.13077*, 2023. 3

[44] Min Zhao, Rongzhen Wang, Fan Bao, Chongxuan Li, and Jun Zhu. Controlvideo: Adding conditional control for one shot text-to-video editing. *arXiv preprint arXiv:2305.17098*, 2023. 2, 3