

# CustomVideo: Customizing Text-to-Video Generation with Multiple Subjects

Zhao Wang<sup>1</sup>, Aoxue Li<sup>2\*</sup>, Lingting Zhu<sup>3</sup>, Yong Guo<sup>2</sup>, Qi Dou<sup>1</sup>, Zhenguo Li<sup>2</sup>

<sup>1</sup>The Chinese University of Hong Kong, <sup>2</sup>Huawei Noah Ark's Lab, <sup>3</sup>The University of Hong Kong  
 {zwang21@cse., qidou}@cuhk.edu.hk, lax@pku.edu.cn,  
 ltzhu99@connect.hku.hk, guoyongcs@gmail.com, li.zhenguo@huawei.com



Figure 1: **Customized text-to-video generation results** of our proposed CustomVideo given multiple subjects (left) and text prompts (below). Our approach can disentangle highly similar subjects, e.g., *cat* v.s. *dog*, simultaneously preserving the fidelity of subjects and smooth motions.

## Abstract

Customized text-to-video generation aims to generate high-quality videos guided by text prompts and subject references. Current approaches for personalizing text-to-video generation suffer from tackling multiple subjects, which is a more challenging and practical scenario. In this work, our aim is to promote multi-subject guided text-to-video customization. We propose CustomVideo, a novel framework that can generate identity-preserving videos with the guidance of multiple subjects. To be specific, firstly, we encourage the co-occurrence of multiple subjects via composing them in a single image. Further, upon a basic text-to-video diffusion model, we design a simple yet effective attention control strategy to disentangle different subjects in the latent space of diffusion model. Moreover, to help the model focus on the specific area of the object, we segment the object from given reference images and provide a corresponding object mask for attention learning. Also, we collect a multi-subject text-to-video generation dataset as a comprehensive benchmark, with 63 individual subjects from 13 different categories and 68 meaningful pairs. Extensive qualitative, quantitative, and user study results demonstrate

\*Corresponding author.

the superiority of our method compared to previous state-of-the-art approaches. The project page is <https://kyfayd.wang/projects/customvideo>.

## 1 Introduction

Text-to-video (T2V) generation [1, 2, 3, 4, 5] has achieved fantastic progress taking advantages of diffusion models [6, 7, 8]. Recently, artists have dreamed of generating videos with their own belongings, *e.g.*, pets, which directs a new research direction named customized T2V generation. Although existing methods [9, 10, 11] have been proposed to generate videos from a single object, tackling multiple objects still remains a difficult scenario. The key challenge is to ensure the co-occurrence of multiple objects in the generated video and retain their corresponding identities.

A recent work, VideoDreamer [12], proposes disen-mix finetuning and human-in-the-loop re-finetuning strategies based on Stable Diffusion [13], aiming to generate videos from multiple subjects. However, VideoDreamer falls short in guaranteeing the co-occurrence of multiple subjects and disentangling similar subjects due to its inconsistent object mixing strategy. As shown in Figure 3, both the car and cat can not be consistently generated across all frames. Another straight forward solution for multi-subject driven T2V generation is by combining multi-subject driven Text-to-Image (T2I) and Image-to-Video (I2V) models. However, this naive approach fails from two aspects, including inaccurate T2I customization and lack of motion (see Figure 3 ‘DisenDiff [14] + SVD [15]’).

In contrast, in our approach, we ensure the co-occurrence of multiple objects during model training, which encourages the model to capture the presence of different subjects simultaneously, thereby facilitating co-occurrence during inference. Additionally, we propose an attention control mechanism to disentangle multiple subjects during training, effectively guiding the model to focus on the corresponding subject area while disregarding irrelevant parts of the image. To facilitate this process, we incorporate a ground truth object mask, obtained through segmentation either from a model like SAM [16] or provided by human annotators, as supervision during optimization. Our attention mechanism consists of two key designs that contribute to the disentanglement of subjects. Firstly, we highlight the corresponding subject area using the ground truth object mask on the cross-attention map, promoting better preservation of subject identity. Secondly, we optimize the cross-attention map towards a slight negative value, excluding the desired subject and mitigating the influence of irrelevant area in the input image. To comprehensively evaluate our proposed approach, we have curated a diverse dataset covering a wide range of 13 categories. Beyond the MultiStudioBench [12] which includes 12 subject pairs, our dataset comprises 63 individual subjects and 68 subject pairs, notably featuring some challenging scenarios involving visually similar objects. Through extensive experiments on this benchmark dataset, we provide qualitative, quantitative, and user study results that demonstrate the superiority of our method in generating high-quality videos with customized subjects. In summary, our contributions are as follows:

- We propose CustomVideo, a novel multi-subject driven T2V generation framework, powered by a simple yet effective co-occurrence and attention control mechanism.
- We collect a multi-subject T2V dataset and build it as a comprehensive benchmark, which covers a wide range of subject categories and diverse subject pairs over them.
- Our method consistently outperforms previous state-of-the-art (SOTA) approaches in qualitative, quantitative and human preference evaluation, *e.g.*, 11.99% *CLIP Image Aligement* and 23.39% *DINO Image Aligement* better than the previous best method.

## 2 Related Work

**Text-to-Video Generation.** Text-to-video generation has made significant advancements in recent years [17, 18, 1]. Early approaches [19, 20, 21, 22] employed GANs [23] and VQ-VAE [24], while more recent works have explored diffusion models to generate high-quality videos [25, 26, 27, 28]. Make-A-Video [29] utilizes a pre-trained image diffusion model with additional temporal attention finetuning. VideoLDM [30] introduces a multi-stage alignment approach in the latent space to generate high-resolution and temporally consistent videos. Other methods [31, 32] generate videos with an image as the first frame and randomly initialized subsequent frames. To enhance controllability, VideoComposer [33] incorporates additional guidance signals, such as depth maps,

to produce desired videos alongside text inputs. VideoDirectorGPT [34] aims to control the video generation with the guidance from different scenes and specific layouts generated by GPT4 in the temporal axis. GEST [35] models the representation between the text and video via a graph of events in the spatial and temporal space, in which the video is generated following a timeline. Tune-A-Video [36] proposes a temporal self-attention module that fine-tunes a pre-trained image diffusion model, achieving successful generation of videos with specific text guidance. Furthermore, diffusion-based video-to-video editing approaches [37, 38] have also been proposed for practical usages.

**Subject-driven Customization.** There has been a growing interest in customizing pre-trained image and video diffusion models for personalized generation [39, 40, 41, 42, 43, 44, 11]. Customization involves generating images and videos with specific subjects, typically based on a few reference images. For image diffusion customization, Textual Inversion [39] represents a specific object as a learnable text token using only a few reference images. This learned text token can then be integrated into a sentence to generate personalized images during inference. Additionally, DreamBooth [45] fine-tunes the weights of the diffusion model to improve fidelity in image generation. Several works [46, 14, 47, 48] have explored personalized image diffusion with multiple subjects, focusing on parameter-efficient finetuning and text embedding learning. While there have been initial attempts to customize video diffusion, such as VideoAssembler [9], VideoBooth [10], and ID-Animator [49], which use reference images to personalize the video diffusion model while preserving subject identity, and DreamVideo [11], which decouples the learning process for subject and motion customization, these methods are designed for single objects and cannot handle multiple subjects when given. A recent work, VideoDreamer [12], proposes multi-subject driven video customization through disen-mix finetuning strategy with LoRA [50]. However, the generated videos do not guarantee the co-occurrence of multiple subjects, or disentanglement of different subjects. In this work, we propose a simple yet effective co-occurrence and attention mechanism that disentangles multiple subjects using masks as guidance while preserving the co-occurrence of subjects in the generated videos.

### 3 Method

#### 3.1 Preliminary: Text-to-Video

Video diffusion models (VDMs) [1, 18, 17] generate videos by gradually denoising a randomly sampled Gaussian noise  $\epsilon$ , following an iterative denoising process that resembles a reverse procedure of a fixed-length Markov Chain. This iterative denoising allows VDMs to capture the temporal dependencies presented in video data. To be specific, a video diffusion model  $\theta$  predicts the added noise at each timestep  $t$  given a text condition  $c$ , where  $t \in \{1, 2, \dots, T\}$ . The training objective for this process can be expressed as a reconstruction loss:

$$\mathcal{L}_{recon} = \mathbb{E}_{\epsilon, z, c, t} \left[ \|\epsilon - \epsilon_{\theta}(z_t, \mathcal{T}(c), t)\|_2^2 \right], \quad (1)$$

where  $z \in \mathbb{R}^{B \times L \times H \times W \times D}$  is the latent code of the input video with batch size  $B$ , video length  $L$ , height  $H$ , width  $W$ , and latent dimension  $D$ .  $\epsilon_{\theta}$  is the noise prediction from the model.  $\mathcal{T}$  is a pre-trained text encoder.  $z_t$  is obtained by adding noise to the ground truth  $z_0$  with  $z_t = \alpha_t z_0 + \sqrt{1 - \alpha_t^2} \epsilon$ , where  $\alpha_t$  is a diffusion hyperparameter. In this work, we utilize zeroscope [51] T2V model as our base model, which is built upon a 3D U-Net, with spatial and temporal modeling for generating high-quality videos.

#### 3.2 CustomVideo with Multiple Subjects

An overview of our proposed CustomVideo framework is illustrated in Figure 2. For customization perspective, we introduce a new learnable word token for every subject, *e.g.*, ‘<new1>’ for cat and ‘<new2>’ for dog, representing the corresponding identity. During the training stage, we employ a concatenation technique to ensure the co-occurrence of multiple subjects. Specifically, we combine these subjects into a single image, facilitating model’s learning of multi-subject patterns. To address the challenge of entanglement among highly similar subjects, we propose an attention control mechanism. This mechanism ensures that the learnable word tokens align with their corresponding regions on the cross-attention map (see Figure 6). By achieving this alignment, we enable the model to disentangle different subjects and enhance the quality of the generated videos. The training process focuses on training the subject-related learnable word tokens, as well as the key and value

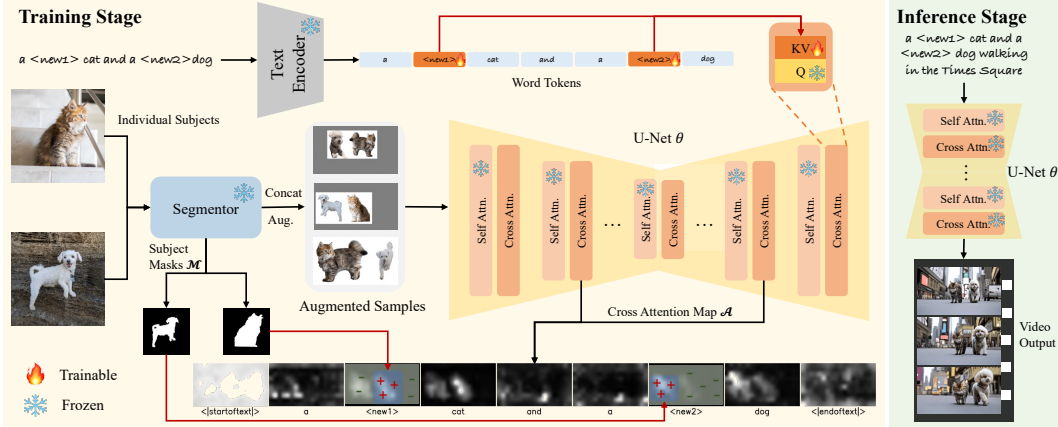


Figure 2: **The overview of our proposed CustomVideo.** We propose a simple yet effective co-occurrence and attention control mechanism with mask guidance to preserve the fidelity of subjects for multi-subject driven text-to-video generation. During the training stage, only the key and value weights in the cross-attention layers are fine-tuned, together with a learnable word token for every subject. In the inference stage, given a text prompt integrated with learned word tokens, we can easily obtain high-quality videos with specific subjects.

weights within the cross-attention layers of the U-Net architecture, adopting a parameter-efficient finetuning approach. During the inference stage, users only need to provide a text prompt description integrated with the corresponding learned word tokens to generate high-quality videos aligned with their preferences. In the following, we will delve into the detailed workings of CustomVideo.

**Co-occurrence Control.** In the context of multi-subject driven T2V generation, ensuring that the model consistently learns to generate videos with multiple subjects is crucial. In a previous work, VideoDreamer [12], the authors proposed a disen-mix strategy for multi-subject driven generation by fine-tuning the model using both single subject images and concatenated images of multiple subjects. However, we have identified that this mixing strategy can confuse the model due to the inconsistent number of subjects present in a single image. Consequently, it leads to unstable generation of videos with multiple subjects (see line ‘both single and concat’ in Figure 5). In our approach, we have found that fine-tuning the model using only concatenated images of multiple subjects is sufficient to ensure the co-occurrence of multiple subjects in the generated videos. Additionally, we have observed that providing clear subjects without background aids the model in learning the specific characteristics of the subjects more effectively. To achieve this, we perform background removal on the subject images. This can be accomplished manually or by employing an automatic tool such as the SAM model [16].

**Attention Control.** Ensuring the co-occurrence of multiple subjects in generated videos is achieved through the above. However, a more challenging task is preserving the distinct characteristics of each subject without interference. Simply fine-tuning the model with concatenated images can lead to confusion among subject characteristics. As shown in Figure 5 (line ‘w/o pos. attn.’), the generated cat predominantly resembles the shape of the provided dog sample, despite having similar color and texture to the provided cat. Thus, disentangling the multiple subjects becomes crucial to generate high-quality videos that faithfully represent each subject’s characteristics.

In our approach, the learnable word tokens operate in the cross-attention layers of the diffusion model. To effectively regulate the subject learning process, we can directly leverage the cross-attention map. As illustrated in Figure 2, we employ an automatic segmentor to obtain the ground truth mask  $\mathcal{M}^p$  for each subject, which indicates the spatial position of the subjects. During the training stage, we extract the cross-attention map  $\mathcal{A}$  from the cross-attention layer using the activations of each word in the given text prompt. For each learnable word token, such as ‘<new1>’, we enhance the corresponding area on the cross-attention map  $\mathcal{A}$  by encouraging the alignment between the subject area and the ground truth subject mask with a loss function as the following:

$$\mathcal{L}_{attn} = \frac{1}{B} \sum_{i=1}^B \frac{1}{N} \sum_{j=1}^N \|\mathcal{A}_{i,j} - \mathcal{M}_{i,j}^p\|_2^2 \quad (2)$$

where  $N$  is the number of subjects, the corresponding area of subject in the mask  $\mathcal{M}^p$  is filled by value 1 while the remaining area is filled by 0. By employing this positive attention mechanism, the model is compelled to allocate more attention to the correct subject area, leading to an effective learning of the corresponding subject characteristics.

The positive-style attention mechanism mentioned above effectively enhances the learning of specific characteristics for each subject. However, there exist issues with the generated subjects, particularly when irrelevant area is present. For example, in Figure 5 (line ‘w/o neg. attn.’), we can observe that the generated dog’s legs are affected by color information from the given cat, which is not desirable. To address this problem, we introduce negative guidance by considering the area outside the subject. In addition to the positive guidance within the ground truth subject mask, we incorporate a small negative value, denoted as  $\eta$ , into the region outside the subject within the mask  $\mathcal{M}^p$ . This modified mask, denoted as  $\mathcal{M}^{[p,\eta]}$ , is then used in Eq. (2) to regulate the subject learning process. By integrating negative guidance, we can alleviate the issue of irrelevant area and improve the fidelity of the generated subjects.

### 3.3 Model Training and Inference

**Training Strategy.** During training, we only fine-tune the weights of key and value in all of the cross-attention layers. We use a more flexible approach by fine-tuning the model with subject images extended to single-frame videos as both input and supervision, rather than using subject videos. This approach guarantees the co-occurrence of multiple subjects because subject images are much easier to obtain and smoother to concatenate compared to subject videos. Moreover, fine-tuning with subject images saves lots of computation cost. Inspired by previous T2I personalization [45, 40], we conduct class-specific prior preservation to improve the diversity of generated videos and alleviate the issue of language drift. The loss for prior preservation is

$$\mathcal{L}_{recon}^{pr} = \mathbb{E}_{\epsilon', z', c', t'} \left[ \|\epsilon' - \epsilon'_\theta(z', \mathcal{T}(c'), t')\|_2^2 \right], \quad (3)$$

where  $z'$  is the latent code of the input class image,  $\epsilon'_\theta$  is the noise prediction from the model  $\theta$ ,  $\epsilon'$  is the randomly sampled Gaussian noise,  $c'$  is the text condition for class image, and  $t'$  is the sampling timestep. To this end, we train our CustomVideo via an end-to-end manner, with the following overall training objective:

$$\mathcal{L} = \mathcal{L}_{recon} + \alpha \cdot \mathcal{L}_{attn} + \beta \cdot \mathcal{L}_{recon}^{pr}, \quad (4)$$

where  $\alpha$  and  $\beta$  are two hyper-parameters to control the weight of attention control and prior preservation, respectively.

**Inference.** During inference, CustomVideo only requires a specific text prompt with corresponding learned word token integrated to generate a required video. To be note that ground truth masks of subjects are not required during inference. Although the model is fine-tuned with single-frame videos extended from subject images, our proposed CustomVideo can generate videos with high diversity while not losing the the fidelity of subjects and motion smoothness (see Figure 1 and Figure 3).

## 4 Experiment

### 4.1 Experimental Setup

**Dataset.** We collect a dataset CustomStudio with 63 individual subjects and construct 68 meaningful pairs for multi-subject driven T2V generation. These subjects are adapted from DreamBooth [45], CustomDiffusion [40], and Mix-of-Show [52], in which they cover a wide range of 13 categories. Each pair of subjects has 10 different text prompts for evaluation, which are designed with different contexts, actions, and so on. More details of dataset can be found in Sec. A.

**Implementation Details.** We train CustomVideo for 500 steps with AdamW [53] optimizer (batch size 2, learning rate 4e-5, and weight decay 1e-2). For class-specific prior preservation, we generate 200 class images with Stable Diffusion v2.1 [13] for each subject, in which the generating prompts are obtained from Claude-3-Opus [54]. Note that the corresponding class images are also concatenated during training phase. The negative value  $\eta$  in the mask  $\mathcal{M}$  is set as -1e-8. The weight parameters  $\alpha$  and  $\beta$  in Eq. (4) are set as 0.2 and 1.0, respectively. During inference, we perform 50 steps denoising with DPM-Solver [8] sampler and classifier-free guidance [55]. The resolution of generated 24-frame

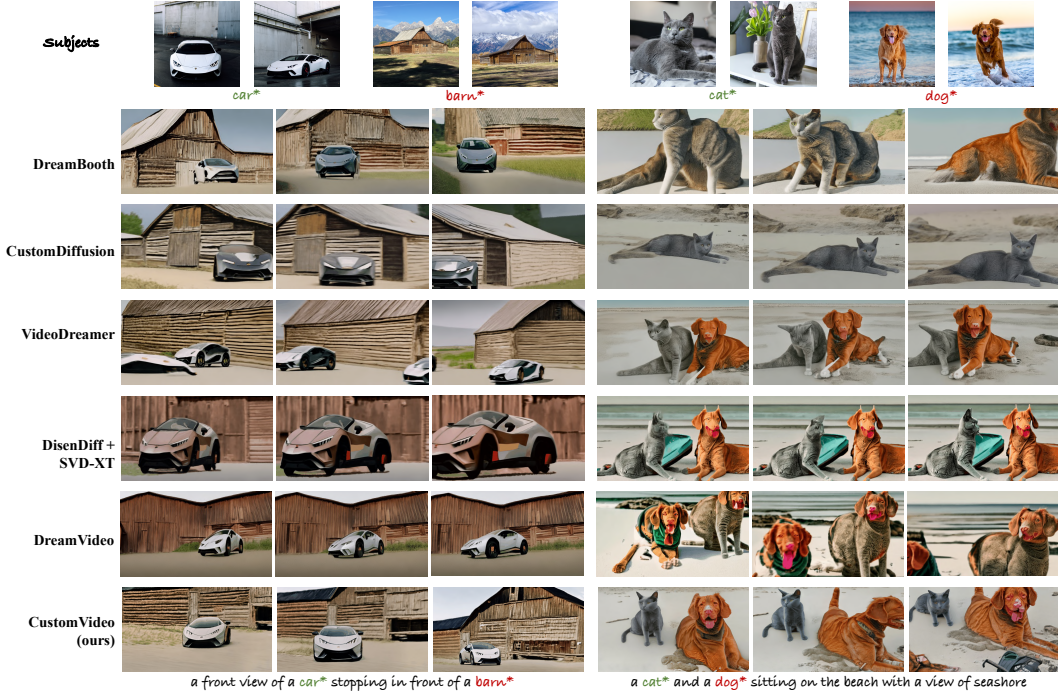


Figure 3: **Qualitative results of our CustomVideo with comparison to SOTA methods**, including DreamBooth, CustomDiffusion, VideoDreamer, DisenDiff + SVD-XT, and DreamVideo. The first line indicates the given subjects, while each line indicates the frames generated by the corresponding method. The bottom line is the text prompt for inference. We can observe that our CustomVideo can generate videos with better fidelity of subjects compared with previous SOTA methods.

videos are  $576 \times 320$  with 8 fps. Interestingly, we find that the trained weights and word tokens learned from low resolution T2V model can be directly loaded for high resolution zeroscope T2V model [56] to generate personalized videos with  $1024 \times 576$  resolution, in which no additional training computation cost is required. Our CustomVideo is implemented based on Diffusers [57]. The training phase takes about 8 minutes for a subject pair on 1 RTX 3090 GPU. Meanwhile, it takes about 1 minute and 2 minutes to generate a low and high resolution video on 1 RTX 3090 GPU, respectively.

**Comparison Methods.** Except for most relevant VideoDreamer [12], we also consider adapting previous SOTA image-based multi-subject driven methods to a video scenario for comparison, including DreamBooth [45] and CustomDiffusion [40]. Also, we compare our method with single-subject driven video personalization method DreamVideo [11]. Moreover, we consider an alternative way for multi-subject T2V customization, that is, the combination of multi-concept T2I and I2V diffusion models. For this case, we utilize the latest multi-concept driven T2I method DisenDiff [14] for customization and I2V method SVD-XT [15] for animating the generated image.

**Evaluation Metrics.** Following previous works [11, 12], we quantitatively evaluate our CustomVideo with the following 4 metrics: 1) *CLIP Textual Alignment* computes the average cosine similarity between the generated frames and text prompt with CLIP [58] ViT-B/32 [59] image and text models; 2) *CLIP Image Alignment* calculates the average cosine similarity between the generated frames and subject images with CLIP ViT-B/32 image model; 3) *DINO Image Alignment* measures the average visual similarity between generated frames and reference images with DINO [60] ViT-S/16 model; 4) *Temporal Consistency* [61] evaluates the average cosine similarity of all consecutive frame pairs in the generated videos.

## 4.2 Main Results

**Qualitative Results.** We present the qualitative comparison results in Figure 3. From these results, we observe that our CustomVideo approach effectively ensures the co-occurrence of multiple subjects

Method	Reference	CLIP-T $\uparrow$	CLIP-I $\uparrow$	DINO-I $\uparrow$	T. Cons. $\uparrow$
DreamBooth [45]	CVPR'23	0.6451	0.6079	0.3109	0.7084
CustomDiffusion [40]	CVPR'23	0.6524	0.6206	0.3164	0.7342
VideoDreamer [12]	arXiv'23	0.6638	0.6297	0.3479	0.7267
DisenDiff [14] + SVD [15]	CVPR'24	0.6592	0.6179	0.3318	0.7095
DreamVideo [11]	CVPR'24	0.6646	0.6128	0.3228	0.7498
<b>CustomVideo (ours)</b>	/	<b>0.7075</b>	<b>0.6863</b>	<b>0.3983</b>	<b>0.7960</b>

Table 1: **Quantitative results of our CustomVideo with comparison to SOTA methods**, including DreamBooth, CustomDiffusion, VideoDreamer, DisenDiff + SVD-XT, and DreamVideo. Our proposed CustomVideo consistently outperforms previous SOTA methods for all 4 evaluation metrics.

and successfully disentangles different subjects. However, both DreamBooth and CustomDiffusion fail to capture the structural color information of the car provided in our experiments. The dominant black color from the car window obscures the entire car, resulting in low fidelity. Moreover, the generated frames from VideoDreamer lack consistency, as some frames depict two cars while others only show one car. Additionally, VideoDreamer fails to capture accurate color information, such as the color of the car door. The alternative approach, DisenDiff + SVD-XT can not neither retain the identities of the given subjects nor generate videos with rich motions. Without specific design for tackling multiple subjects, DreamVideo can not distinguish the similar pair ‘cat’ and ‘dog’. In contrast, our CustomVideo method excels in handling such challenging scenes and foreground subject scenarios, effectively capturing the intricate structural details of the provided car for video generation. Similarly, when considering the case of ‘cat’ and ‘dog’, our approach also demonstrates superior capability in generating high-quality videos. More results can be found in Figure 9 and Figure 10.

**Quantitative Results.** We conduct quantitative experiments on our collected dataset. To ensure a thorough analysis, we generate videos using 10 individual prompts and 4 random seeds for each pair of subjects, resulting in a total of 2,720 generated videos for each method. We then evaluate the quality of the generated videos using four metrics, and the results are presented in Table 1. The table clearly demonstrates that our proposed method is capable of generating videos that are better aligned with the given subjects, outperforming the most recent video diffusion model, DreamVideo, by 11.99% and 23.39% in terms of *CLIP Image Alignment* and *DINO Image Alignment*, respectively. Meanwhile, the videos generated from our approach achieve better alignment with the desired text prompts, exceeding DreamVideo by 6.46% *CLIP Textual Alignment*. These improvements can be attributed to our specially designed co-occurrence and attention control mechanisms, which effectively disentangle and preserve the fidelity of the subjects. Furthermore, our CustomVideo generates videos with significantly higher temporal consistency compared to SOTA methods, as indicated in Table 1. For instance, CustomVideo surpasses DreamVideo by 6.16% in terms of *Temporal Consistency*.

**Human Preference Study.** To further validate our method, we conduct human evaluations on our CustomVideo with comparison to 5 SOTA methods. In this study, we collect 1500 answers from 25 independent human raters with the following questions: 1) which one is aligned to the text prompt best? 2) which one is aligned to the subject images best? 3) which one has the best overall quality? The results are shown in Figure 4. Our CustomVideo is found to be the most preferred option based on human evaluations in all three dimensions.

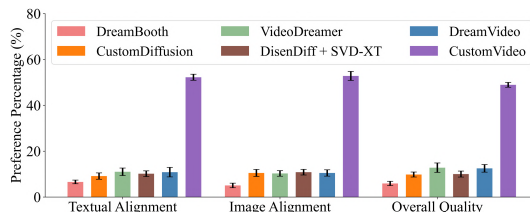


Figure 4: **User study.** Our CustomVideo achieves the best human preference compared with 5 SOTA comparison methods in terms of *Textual Alignment*, *Image Alignment*, and *Overall Quality*.

### 4.3 Ablation Studies

We conduct ablation analysis of our method from 3 aspects: 1) the effect of each component in our method; 2) the behavior of attention control; 3) the effect of hyper-parameters in attention control.

**Component Analysis.** We conduct a thorough analysis of each component in our CustomVideo, presenting both qualitative and quantitative results in Figure 5 and Table 2. One crucial observation

is the significance of ensuring co-occurrence during the fine-tuning process. We notice a significant drop in performance when multiple subjects are not concatenated into a single image (line ‘w/o concat’ in Table 2). The absence of subject concatenation leads to the domination of one single subject, which is deemed unacceptable (line ‘w/o concat’ in Figure 5). Moreover, we conduct a study on fine-tuning T2V model with both single and concatenated subjects to examine if the single subject could aid in learning corresponding characteristics. Surprisingly, the results reveal that adding the single subjects to the training process proved detrimental, resulting in inconsistent generation during inference (line ‘both single and concat’ in Figure 5).

Our attention control mechanism plays a crucial role in preserving the identities of subjects. The positive attention guidance in particular improves the *CLIP Image Alignment* metric by 11.45%, as demonstrated in Table 2. The negative attention guidance also proved beneficial in promoting better image alignment. Moreover, the attention control mechanism significantly enhances the temporal consistency of the generated videos, a crucial factor in video generation. The videos generated without positive and negative attention guidance clearly showcase the positive impact of both components (lines ‘w/o pos. attn.’ and ‘w/o neg. attn.’ in Figure 5). The positive guidance specifically preserves the unique characteristics of a subject, such as distinguishing between a cat and a dog, while the negative guidance weakens the influence of other subjects on a specific subject. By incorporating both positive and negative guidance, our CustomVideo excels in generating videos with high subject fidelity and remarkable temporal consistency. Furthermore, we observe that removing the background from the given subject images significantly improves the generation of high-quality videos (line ‘w/o remove bg’ in Figure 5). By eliminating the background, the T2V model can focus solely on learning the characteristics of the given subjects.

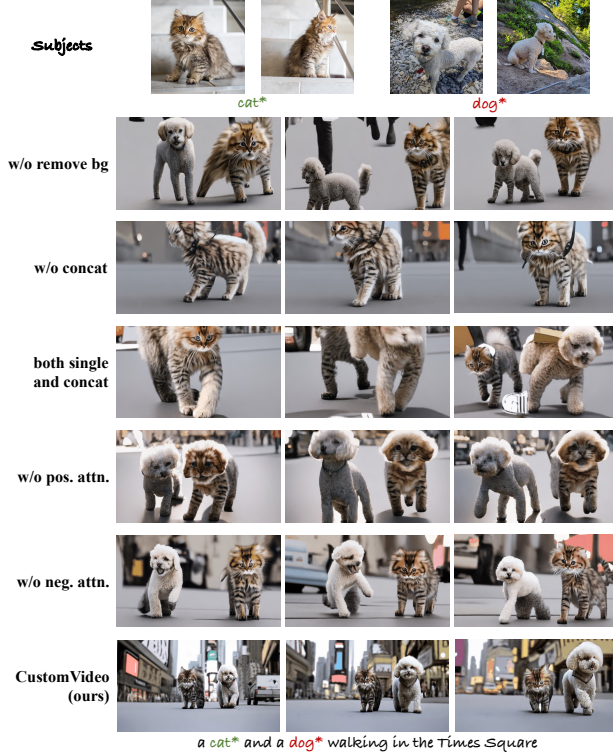


Figure 5: **Qualitative results for component analysis** of our proposed CustomVideo. We find that ensuring concatenating subjects during training is effective for guaranteeing the co-occurrence in the generated video. Moreover, our attention mechanism could disentangle different subjects.

**Comparison of Cross-attention Maps.** We investigate the cross-attention maps during training, and study the behavior of our attention control mechanism. The visualization results are shown in Figure 6. Without attention control, it can be observed that the learnable word token ‘<new1>’ for cat is aligned to the region of dog, which is unacceptable. In contrast, equipped with attention control mechanism, the learnable word token ‘<new1>’ and ‘<new2>’ are better aligned to the correct areas of cat and dog, respectively. Such results strongly demonstrate the effectiveness of our approach.

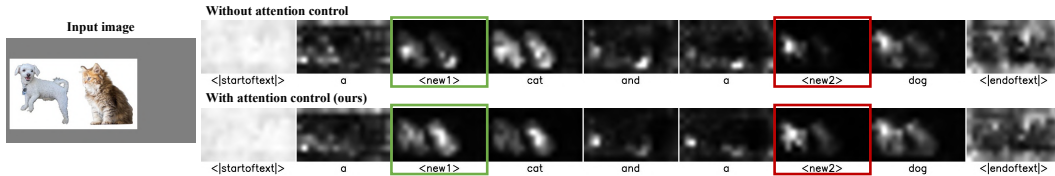


Figure 6: **Comparison results of cross-attention maps.** We observe that the learnable word tokens better match the corresponding correct areas of the subjects with our proposed attention control.



Method	CLIP-T $\uparrow$	CLIP-I $\uparrow$	DINO-I $\uparrow$	T. Cons. $\uparrow$
w/o remove bg	0.6682	0.6752	0.3636	0.7609
w/o concat	0.6115	0.6611	0.3386	0.7034
both single and concat	0.6248	0.6096	0.3054	0.7604
w/o pos. attn.	0.6851	0.6158	0.3215	0.7213
w/o neg. attn.	0.6961	0.6716	0.3613	0.7886
<b>CustomVideo (ours)</b>	<b>0.7075</b>	<b>0.6863</b>	<b>0.3983</b>	<b>0.7960</b>

Table 2: **Quantitative results for component analysis** of our proposed CustomVideo. We observe a significant performance drop when removing subjects concatenation or background removal. Moreover, training the model with both single and concatenated subjects will produce unstable results, leading to performance degradation. In practice, we suggest simultaneously using positive and negative attention mechanisms to obtain the best results.

$\alpha$	CLIP-T $\uparrow$	CLIP-I $\uparrow$	DINO-I $\uparrow$	T. Cons. $\uparrow$
1.0	0.6909	0.6496	0.3773	0.7812
<b>0.2</b>	<b>0.7075</b>	<b>0.6863</b>	<b>0.3983</b>	<b>0.7960</b>
0.01	0.6924	0.6710	0.3715	0.7886

(a) Weight  $\alpha$  for attention loss

$\eta$	CLIP-T $\uparrow$	CLIP-I $\uparrow$	DINO-I $\uparrow$	T. Cons. $\uparrow$
-1e-5	0.6786	0.6392	0.3719	0.7604
<b>-1e-8</b>	<b>0.7075</b>	<b>0.6863</b>	<b>0.3983</b>	<b>0.7960</b>
-1e-11	0.6926	0.6737	0.3613	0.7812

(b) Negative value  $\eta$  for GT mask

Table 3: **Quantitative results of the weight of attention loss and the negative value.** Interestingly,  $\alpha=0.2$  and  $\eta=-1e-8$  yield the best results. Practically, we adopt this setting in all the experiments.

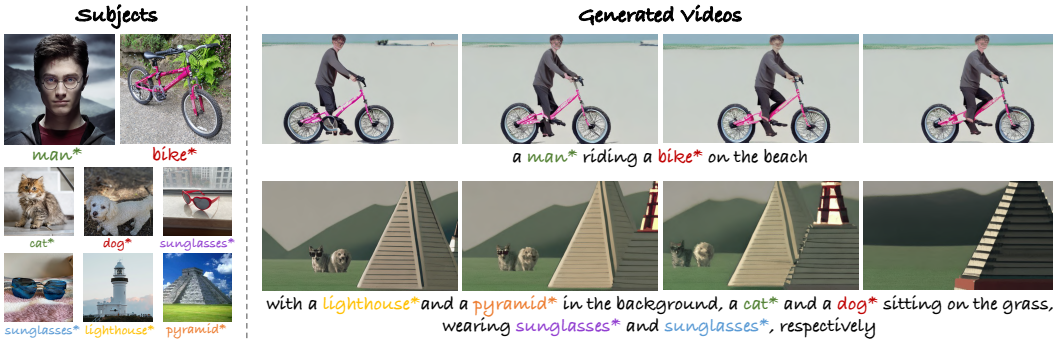


Figure 7: **Failure cases of our CustomVideo.** Our CustomVideo fails to generate vivid facial contents with a global view, such as ‘man’ and ‘bike’ (first line). Also, our approach can not tackle too many subjects, e.g., customizing T2V generation with 6 different subjects (second line).

**Effect of Hyper-parameters in Attention Control.** We also conduct investigations into the effects of two important parameters: the weight of attention loss ( $\alpha$ ) and the negative value ( $\eta$ ) used in the guidance mask ( $\mathcal{M}$ ). The quantitative results of these investigations are presented in Table 3a and Table 3b. We observe that the best performance is achieved when the weight  $\alpha$  is set to 0.2. This value results in optimal alignment between the generated videos and the given subjects, indicating the importance of appropriately balancing the attention guidance during training. For the negative value in the guidance mask, our findings reveal that even a slight negative value is sufficient to enhance the quality of T2V generation. This implies that incorporating negative attention guidance can effectively suppress the influence of other subjects on a specific subject, improving the generation quality.

#### 4.4 Limitations

In Figure 7, we present some failure cases encountered during our experiments. Since our method relies on the ability of the base model, it would fail if the base model could not generate, such as small faces in a global view. Currently, Our approach can not tackle too many subjects, in which the spatial positions of different objects in the video conflict with each other. To address such case, we can control the positions of different subjects via conditioning on the pre-defined spatial bounding boxes during inference sampling.

## 5 Conclusions

This paper provides a novel framework CustomVideo for multi-subject driven T2V generation, powered by a simple yet effective co-occurrence and attention control mechanism. It successfully disentangles different similar subjects and preserves the corresponding identities. We design a flexible and efficient training process, only requiring subject images rather than subject videos. And during

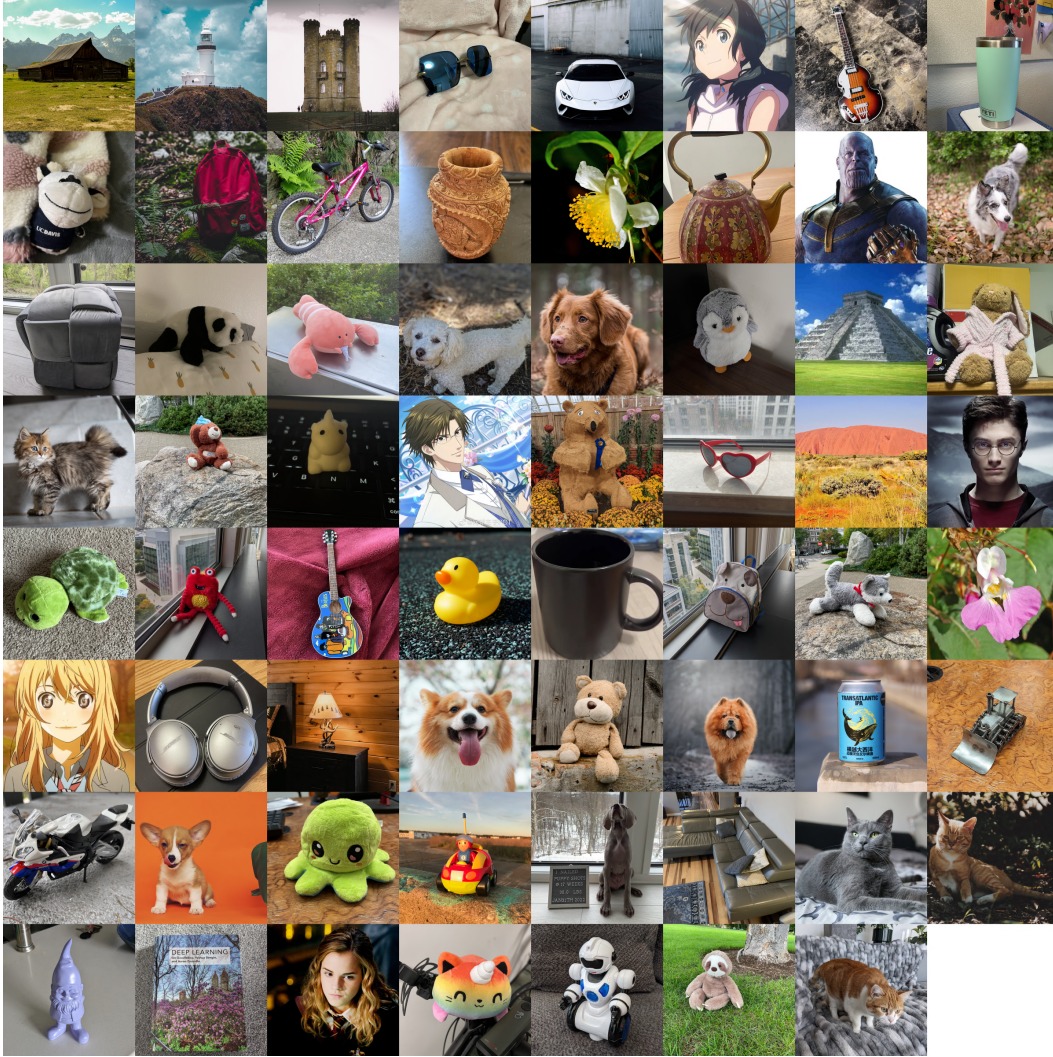


Figure 8: **The overview of our CustomStudio dataset.** Samples in our CustomStudio dataset cover a wide range of 13 object categories, including anime character, decorate item, furniture, instrument, person, pet, plant, plush, scene, thing, toy, transport, and wearable item.

inference, a desired video can be easily generated just by providing a text prompt from the user. We collect a comprehensive dataset CustomStudio for evaluation. Extensive quantitative and qualitative experiments, together with a user study, demonstrates the effectiveness of our CustomVideo, outperforming SOTA methods significantly. We hope our method could serve as a strong baseline and motivate further research on subject-driven applications, especially for multi-subject scenarios.

## A CustomStudio Dataset

**Dataset Info.** We collect a dataset CustomStudio for multi-subject driven video generation. This dataset is composed of 63 individual objects and 68 meaningful pairs. These objects cover a wide range of 13 diverse categories, including anime character, decorate item, furniture, instrument, person, pet, plant, plush, scene, thing, toy, transport, and wearable item. For each subject, 4-21 images are provided as the reference images. The images are adopted from DreamBooth [45], CustomDiffusion [40] and Mix-of-Show [52]. An overview of the samples in CustomStudio dataset is shown in Figure 8. A meaningful pair is composed of different subjects may appear in a video together in practice. For example, ‘cat’ and ‘dog’ can be a meaningful pair, while ‘book’ and ‘barn’ can not arise in the

Two-subject Pair	[decorate item, plant], [person, transport], [pet, pet], [furniture, thing], [toy, toy] [decorate item, furniture], [instrument, person], [person, person], [plush, plush] [anime character, anime character], [toy, transport], [person, wearable item] [person, scene], [scene, transport], [instrument, anime character]
Three-subject Pair	[anime character, anime character, scene], [anime character, instrument, scene] [toy, toy, toy], [plush, plush, plush], [person, instrument, scene] [person, person, scene], [pet, wearable item, scene], [toy, transport, scene] [pet, pet, scene], [decorate item, furniture, furniture]

Table 4: **Subject pairs** constructed upon categories of the subjects in our CustomStudio dataset.

<b>Prompt Template for 2-subject Pair</b>	
a <c1> and a <c2> sitting on an antique table	
a <c1> and a <c2> sitting on beach with a view of seashore	
a <c2> and a <c1> side by side on a mountaintop, overlooking a sunrise	
a <c1> and a <c2> on a surfboard together in the middle of a clear blue ocean	
a <c1> playing with <c2>	
a <c1> playing with a robot toy <c2>	
a <c2> playing with a robot toy <c1>	
a plush toy replica of a <c1> and a <c2> sitting beside it	
a <c1> and a <c2> walking in the Times Square	
a <c1> and a <c2> walking on the Great Wall	
<b>Prompt Template for 3-subject Pair</b>	
a <c1> wearing a <c2>, with a <c3> in the background	
a <c1> wearing a <c2> and giving a speech, with a <c3> in the background	
a wide shot of a <c1> wearing a <c2> with boston city in background, with a <c3> in the background	
a long shot of a <c1> walking their golden retriever, wearing a <c2>, with a <c3> in the background	
a <c1> sitting on the sidewalk wearing a <c2>, with a <c3> in the background	
a <c1> standing at a graffiti wall, showcasing a new <c2>, with a <c3> in the background	
close shot of a <c1> walking on a ramp wearing a <c2>, with a <c3> in the background	
a <c1> stepping out of a taxi in a <c2>, with a <c3> in the background	
close shot of a <c1> walking in rain wearing a <c2>, with a <c3> in the background	
long shot of a <c1> sitting on the edge of a roof wearing <c2>, with a <c3> in the background	

Table 5: **Prompt templates** used for generating videos during inference.

most cases. We set up 25 types of subject pairs, which are constructed based on the corresponding categories, as shown in Table 4.

**Prompt Templates.** The prompt templates used for generating videos during inference are shown in Table 5. For example, we want to generate a customized video with a ‘cat’ and a ‘dog’, the prompt for inference can be ‘a <new1> cat and a <new2> dog walking in the Times Square’.

## B Additional Experimental Results

**Layers of Cross-attention Maps.** We investigate the effect of how we extract the cross-attention maps from the cross-attention layers. There are four levels of cross-attention layers in the U-Net, in which the sizes of cross-attention maps are different. Taking the resolution  $576 \times 320$  of our generated low resolution video as an example, the sizes of cross-attention maps are  $72 \times 40$  ( $\ell_1$ ),  $36 \times 20$  ( $\ell_2$ ),  $18 \times 10$  ( $\ell_3$ ), and  $9 \times 5$  ( $\ell_4$ ). Here, we study the effect of different levels of cross-attention maps, the results are shown in Table 6. We observe a significant performance drop with too large ( $\ell_1$ ) or small ( $\ell_4$ ) cross-attention maps. Large cross-attention maps may, on one hand, lead to the loss of subtle subject features, while on the other hand, small cross-attention maps can result in inaccurate optimization between the maps and the downsampled ground truth masks. Regarding the cross-attention maps with middle size from  $\ell_2$  and  $\ell_3$ ,  $\ell_3$  works better. Thus, in our experiments, we extract cross-attention maps from  $\ell_3$  for attention control.

$\ell_1$	$\ell_2$	$\ell_3$	$\ell_4$	CLIP-T $\uparrow$	CLIP-I $\uparrow$	DINO-I $\uparrow$	T. Cons. $\uparrow$
✓	✗	✗	✗	0.6111	0.6335	0.3964	0.7315
✗	✓	✗	✗	0.6080	0.6641	0.3055	0.7242
✗	✗	✓	✗	<b>0.7075</b>	<b>0.6863</b>	<b>0.3983</b>	<b>0.7960</b>
✗	✗	✗	✓	0.6134	0.6420	0.3215	0.7793
✓	✓	✗	✗	0.6234	0.6375	0.3497	0.7443
✗	✗	✓	✓	0.7047	0.6007	0.3642	0.7380
✗	✓	✓	✗	0.6730	0.6364	0.3977	0.7433
✓	✗	✗	✓	0.6508	0.6360	0.3055	0.7469
✓	✓	✓	✓	0.6432	0.6738	0.3594	0.6975

Table 6: **Quantitative results of cross-attention levels.** We find that only using cross-attention map from  $\ell_3$  works best and adopt this setting by default.

---

A fluffy white Persian cat with blue eyes, sitting on a plush red cushion
A playful orange tabby kitten chasing a ball of yarn across a hardwood floor
A sleek black cat with green eyes, perched on a windowsill, looking out at a rainy day
A cute Siamese cat with a pink collar, napping in a wicker basket
A majestic Maine Coon cat with a bushy tail, walking through a field of wildflowers
A curious Sphinx cat with wrinkled skin, exploring a cardboard box
A lazy gray British Shorthair cat, curled up on a soft blanket, basking in the sunlight
A mischievous Calico cat, peeking out from behind a potted plant
A regal Russian Blue cat with silver fur, sitting on a velvet armchair
A friendly Ragdoll cat with bright blue eyes, being petted by a gentle hand
A playful Bengal cat with spots, pouncing on a feather wand toy
A curious Abyssinian cat with large ears, investigating a paper bag
A fluffy Himalayan cat with a flat face, lounging on a fuzzy rug
A sleepy Birman cat with white paws, dozing off on a cozy bed
A graceful Siberian cat with long fur, walking along a wooden fence
A cute Munchkin cat with short legs, playing with a crinkly ball
A curious Norwegian Forest cat with a thick coat, climbing a cat tree
A playful Scottish Fold cat with folded ears, batting at a dangling string
A majestic Savannah cat with a tall, slender build, surveying its surroundings from atop a bookshelf
A friendly Tonkinese cat with a unique coat pattern, rubbing against its owner’s leg

---

Table 7: **Class prompts** for generating class images used in class-specific prior preservation. In our experiments, we generate one class image with one independent text prompt. Here we show 20 examples for 20 class images.

**Additional Qualitative Results.** We show additional qualitative results here, including comparison of our CustomVideo with previous SOTA methods in Figure 9 and generated video of CustomVideo in Figure 10. It is clear that our method can generate videos with higher fidelity and better preserve the identities of the reference subjects.

## C Additional Experimental Details

All of the training process is under fp16 mixed precision with accelerate package [62]. We use data augmentations during training, including randomly horizontal flip, randomly crop and resize, together with corresponding prompt change (‘very small’ or ‘close up’ appends before the prompt).

**Class Prompts.** For class-specific prior preservation, we generate 200 class images with Stable Diffusion v2.1 [13] for each subject, in which the generating prompts are obtained from Claude-3-Opus [54]. We use Stable Diffusion v2.1 [13] to generate 200 class images for every subject in class-specific prior preservation. The text prompts for generating are obtained from Claude-3-Opus [54]. We show 20 examples in Table 7.

**Implementation Details.** We adopt the follow setting for in our experiments, including the previous SOTA methods for comparison and our CustomVideo:



Figure 9: **Additional comparison results** of our CustomVideo with previous SOTA methods.

- **DreamBooth** [45]: For efficient fine-tuning, we utilize LoRA [50] to adapt the U-Net under DreamBooth [45]. The rank of LoRA module is set as 4. We use AdamW [53] optimizer, with learning rate  $4e-4$ , weight decay  $1e-2$ , batch size 2, and training steps 1000.
- **CustomDiffusion** [40]: The weights of query and value in all of the cross-attention layers of U-Net are fine-tuned in CustomDiffusion [40]. We use AdamW [53] optimizer, with learning rate  $4e-5$ , weight decay  $1e-2$ , batch size 2, and training steps 500.

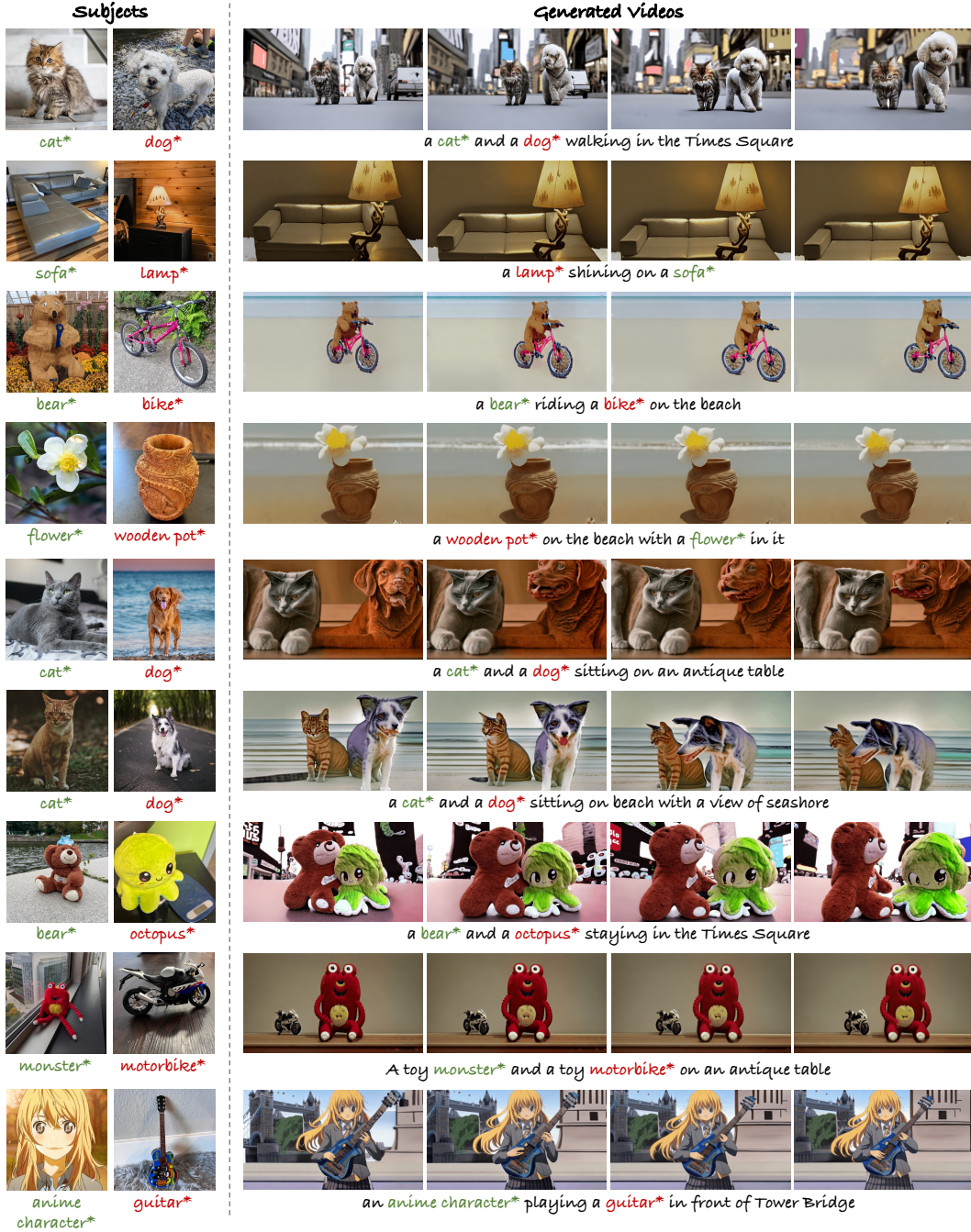


Figure 10: Additional qualitative results of generated videos from our CustomVideo.

- **VideoDreamer** [12]: LoRA [50] is utilized for fine-tuning the U-Net and text encoder in VideoDreamer [12]. The AdamW [53] optimizer is used with learning rate 5e-5, weight decay 1e-2, batch size 2, and training steps 500.
- **DisenDiff + SVD-XT** [14, 15]: DisenDiff is trained for 250 steps with AdamW [53] optimizer (learning rate 4e-5, weight decay 1e-2, batch size 2). The sampling step for SVD-XT is 30.
- **DreamVideo** [11]: We train the learnable textual token and identity adapter for 3000 and 800 steps, respectively, with learning rate 1e-4 and 1e-5. We use AdamW [53] optimizer with a batch size of 2.

Given two subjects for customization, which one is aligned to the text prompt best?

Text Prompt: a **cat\*** and a **dog\*** sitting on the beach with a view of seashore

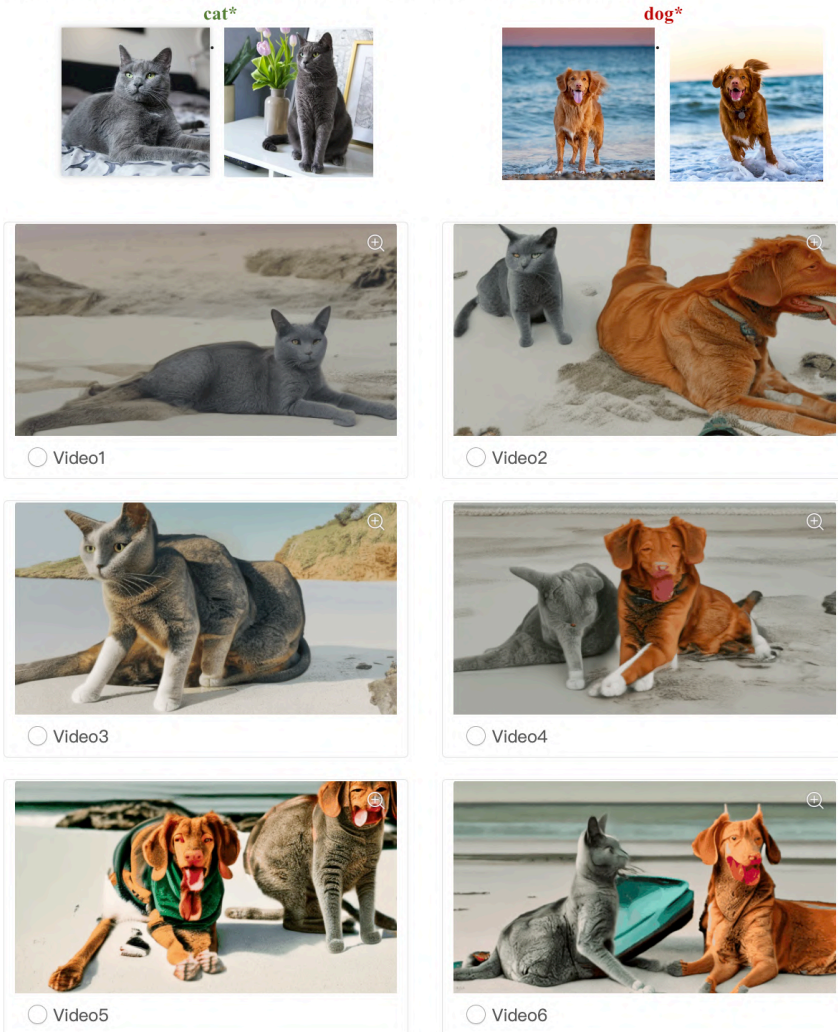


Figure 11: **Human interface** used in the user study. The generated videos of 5 comparison approaches and our CustomVideo are randomly ordered and anonymous for fair comparison.

- **CustomVideo:** The weights of query and value in all of the cross-attention layers (including the temporal and spatial layers) of U-Net are fine-tuned. The sampling step for high resolution video is 30.

**Human Interface for User Study.** The human annotators are all researchers with experience in artificial intelligence. The interface for human evaluation is shown in Figure 11.

## D Broader Impact

Generating realistic videos with human imagination and personal belongings is crucial in practice. Our work makes this poorly explored problem come true by multi-subject driven T2V generation. We believe our work will benefit a wide range of AIGC applications, such as film production and virtual reality. Although our model significantly outperforms other competitive methods, it still struggle to generate promising results for some cases, such as small faces and lots of subjects shown in Figure 7. Please be careful to use our model in the situations where failures will cause severe consequences.

## References

- [1] Jiuniu Wang, Hangjie Yuan, Dayou Chen, Yingya Zhang, Xiang Wang, and Shiwei Zhang. Modelscope text-to-video technical report. *arXiv*, 2023. [2](#), [3](#)
- [2] Yan Zeng, Guoqiang Wei, Jiani Zheng, Jiaxin Zou, Yang Wei, Yuchen Zhang, and Hang Li. Make pixels dance: High-dynamic video generation. *arXiv*, 2023. [2](#)
- [3] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large datasets. *arXiv*, 2023. [2](#)
- [4] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, et al. Videopoet: A large language model for zero-shot video generation. *arXiv*, 2023. [2](#)
- [5] Haoxin Chen, Yong Zhang, Xiaodong Cun, Menghan Xia, Xintao Wang, Chao Weng, and Ying Shan. Videocrafter2: Overcoming data limitations for high-quality video diffusion models, 2024. [2](#)
- [6] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *NeurIPS*, 33:6840–6851, 2020. [2](#)
- [7] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv*, 2020. [2](#)
- [8] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *NeurIPS*, 35:5775–5787, 2022. [2](#), [5](#)
- [9] Haoyu Zhao, Tianyi Lu, Jiayi Gu, Xing Zhang, Zuxuan Wu, Hang Xu, and Yu-Gang Jiang. Videoassembler: Identity-consistent video generation with reference entities using diffusion model. *arXiv*, 2023. [2](#), [3](#)
- [10] Yuming Jiang, Tianxing Wu, Shuai Yang, Chenyang Si, Dahua Lin, Yu Qiao, Chen Change Loy, and Ziwei Liu. Videobooth: Diffusion-based video generation with image prompts. *CVPR*, 2024. [2](#), [3](#)
- [11] Yujie Wei, Shiwei Zhang, Zhiwu Qing, Hangjie Yuan, Zhiheng Liu, Yu Liu, Yingya Zhang, Jingren Zhou, and Hongming Shan. Dreamvideo: Composing your dream videos with customized subject and motion. In *CVPR*, 2024. [2](#), [3](#), [6](#), [7](#), [14](#)
- [12] Hong Chen, Xin Wang, Guanning Zeng, Yipeng Zhang, Yuwei Zhou, Feilin Han, and Wenwu Zhu. Videodreamer: Customized multi-subject text-to-video generation with disen-mix finetuning. *arXiv*, 2023. [2](#), [3](#), [4](#), [6](#), [7](#), [14](#)
- [13] 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. [2](#), [5](#), [12](#)
- [14] Yanbing Zhang, Mengping Yang, Qin Zhou, and Zhe Wang. Attention calibration for disentangled text-to-image personalization. *CVPR*, 2024. [2](#), [3](#), [6](#), [7](#), [14](#)
- [15] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, Varun Jampani, and Robin Rombach. Stable video diffusion: Scaling latent video diffusion models to large datasets, 2024. [2](#), [6](#), [7](#), [14](#)
- [16] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C. Berg, Wan-Yen Lo, Piotr Dollár, and Ross Girshick. Segment anything. *ICCV*, 2023. [2](#), [4](#)
- [17] Haoxin Chen, Menghan Xia, Yingqing He, Yong Zhang, Xiaodong Cun, Shaoshu Yang, Jinbo Xing, Yaofang Liu, Qifeng Chen, Xintao Wang, et al. Videocrafter1: Open diffusion models for high-quality video generation. *arXiv*, 2023. [2](#), [3](#)
- [18] Zhongjie Duan, Lizhou You, Chengyu Wang, Cen Chen, Ziheng Wu, Weining Qian, Jun Huang, Fei Chao, and Rongrong Ji. Diffsynth: Latent in-iteration deflickering for realistic video synthesis. *arXiv*, 2023. [2](#), [3](#)
- [19] Yogesh Balaji, Martin Renqiang Min, Bing Bai, Rama Chellappa, and Hans Peter Graf. Conditional gan with discriminative filter generation for text-to-video synthesis. In *IJCAI*, volume 1, page 2, 2019. [2](#)
- [20] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video generator with the price, image quality and perks of stylegan2. In *CVPR*, pages 3626–3636, 2022. [2](#)
- [21] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv*, 2022. [2](#)
- [22] 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 descriptions. In *ICLR*, 2023. [2](#)
- [23] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NeurIPS*, pages 2672–2680, 2014. [2](#)



- [24] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. In *NeurIPS*, pages 6309–6318, 2017. [2](#)
- [25] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity video generation with arbitrary lengths. *arXiv*, 2022. [2](#)
- [26] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *arXiv*, 2023. [2](#)
- [27] David Junhao Zhang, Jay Zhangjie Wu, Jia-Wei Liu, Rui Zhao, Lingmin Ran, Yuchao Gu, Difei Gao, and Mike Zheng Shou. Show-1: Marrying pixel and latent diffusion models for text-to-video generation. *arXiv*, 2023. [2](#)
- [28] Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv*, 2023. [2](#)
- [29] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *ICLR*, 2023. [2](#)
- [30] 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 *CVPR*, pages 22563–22575, 2023. [2](#)
- [31] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *ICCV*, 2023. [2](#)
- [32] 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*, 2023. [2](#)
- [33] Xiang Wang, Hangjie Yuan, Shiwei Zhang, Dayou Chen, Jiuniu Wang, Yingya Zhang, Yujun Shen, Deli Zhao, and Jingren Zhou. Videocomposer: Compositional video synthesis with motion controllability. *arXiv*, 2023. [2](#)
- [34] Han Lin, Abhay Zala, Jaemin Cho, and Mohit Bansal. Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning. *arXiv*, 2023. [3](#)
- [35] Mihai Masala, Nicolae Cudlenco, Traian Rebedea, and Marius Leordeanu. Explaining vision and language through graphs of events in space and time. In *ICCV*, pages 2826–2831, 2023. [3](#)
- [36] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, 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. In *ICCV*, 2023. [3](#)
- [37] 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. *ICCV*, 2023. [3](#)
- [38] Michal Geyer, Omer Bar-Tal, Shai Bagon, and Tali Dekel. Tokenflow: Consistent diffusion features for consistent video editing. *arXiv*, 2023. [3](#)
- [39] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using textual inversion. In *ICLR*, 2023. [3](#)
- [40] Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept customization of text-to-image diffusion. In *CVPR*, pages 1931–1941, 2023. [3](#), [5](#), [6](#), [7](#), [10](#), [13](#)
- [41] Omri Avrahami, Kfir Aberman, Ohad Fried, Daniel Cohen-Or, and Dani Lischinski. Break-a-scene: Extracting multiple concepts from a single image. In *SIGGRAPH Asia*, pages 1–12, 2023. [3](#)
- [42] Jian Ma, Junhao Liang, Chen Chen, and Haonan Lu. Subject-diffusion: Open domain personalized text-to-image generation without test-time fine-tuning. *arXiv*, 2023. [3](#)
- [43] Zhiheng Liu, Ruili Feng, Kai Zhu, Yifei Zhang, Kecheng Zheng, Yu Liu, Deli Zhao, Jingren Zhou, and Yang Cao. Cones: Concept neurons in diffusion models for customized generation. *ICML*, 2023. [3](#)
- [44] Zhiheng Liu, Yifei Zhang, Yujun Shen, Kecheng Zheng, Kai Zhu, Ruili Feng, Yu Liu, Deli Zhao, Jingren Zhou, and Yang Cao. Cones 2: Customizable image synthesis with multiple subjects. *NeurIPS*, 2023. [3](#)
- [45] Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. Dream-booth: Fine tuning text-to-image diffusion models for subject-driven generation. In *CVPR*, pages 22500–22510, 2023. [3](#), [5](#), [6](#), [7](#), [10](#), [13](#)
- [46] Guangxuan Xiao, Tianwei Yin, William T Freeman, Frédo Durand, and Song Han. Fastcomposer: Tuning-free multi-subject image generation with localized attention. *arXiv*, 2023. [3](#)

- [47] Zhichao Wei, Qingkun Su, Long Qin, and Weizhi Wang. Mm-diff: High-fidelity image personalization via multi-modal condition integration. *arXiv*, 2024. 3
- [48] Jiaxiu Jiang, Yabo Zhang, Kailai Feng, Xiaohe Wu, and Wangmeng Zuo. Mc<sup>2</sup>: Multi-concept guidance for customized multi-concept generation. *arXiv*, 2024. 3
- [49] Xuanhua He, Quande Liu, Shengju Qian, Xin Wang, Tao Hu, Ke Cao, Keyu Yan, Man Zhou, and Jie Zhang. Id-animator: Zero-shot identity-preserving human video generation. *arXiv*, 2024. 3
- [50] 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. In *ICLR*, 2022. 3, 13, 14
- [51] Spencer Sterling. Zeroscope. [https://huggingface.co/cerspense/zeroscope\\_v2\\_576w](https://huggingface.co/cerspense/zeroscope_v2_576w), 2023. 3
- [52] Yuchao Gu, Xintao Wang, Jay Zhangjie Wu, Yujun Shi, Yunpeng Chen, Zihan Fan, Wuyou Xiao, Rui Zhao, Shuning Chang, Weijia Wu, et al. Mix-of-show: Decentralized low-rank adaptation for multi-concept customization of diffusion models. *Advances in Neural Information Processing Systems*, 36, 2023. 5, 10
- [53] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv*, 2017. 5, 13, 14
- [54] Anthropic. Introducing the next generation of claude. <https://www.anthropic.com/news/claude-3-family>, 2024. 5, 12
- [55] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv*, 2022. 5
- [56] Spencer Sterling. Zeroscope xl. [https://huggingface.co/cerspense/zeroscope\\_v2\\_XL](https://huggingface.co/cerspense/zeroscope_v2_XL), 2023. 6
- [57] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. Diffusers: State-of-the-art diffusion models. <https://github.com/huggingface/diffusers>, 2022. 6
- [58] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, pages 8748–8763, 2021. 6
- [59] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 6
- [60] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *ICCV*, pages 9650–9660, 2021. 6
- [61] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. In *ICCV*, pages 7346–7356, 2023. 6
- [62] Sylvain Gugger, Lysandre Debut, Thomas Wolf, Philipp Schmid, Zachary Mueller, Sourab Mangrulkar, Marc Sun, and Benjamin Bossan. Accelerate: Training and inference at scale made simple, efficient and adaptable. <https://github.com/huggingface/accelerate>, 2022. 12