

---

# ModelScope Text-to-Video Technical Report

---

Jiuniu Wang\* Hangjie Yuan\* Dayou Chen\* Yingya Zhang\*†  
Xiang Wang Shiwei Zhang

Alibaba Group

## Abstract

This paper introduces ModelScopeT2V, a text-to-video synthesis model that evolves from a text-to-image synthesis model (*i.e.*, Stable Diffusion). ModelScopeT2V incorporates spatio-temporal blocks to ensure consistent frame generation and smooth movement transitions. The model could adapt to varying frame numbers during training and inference, rendering it suitable for both image-text and video-text datasets. ModelScopeT2V brings together three components (*i.e.*, VQGAN, a text encoder, and a denoising UNet), totally comprising 1.7 billion parameters, in which 0.5 billion parameters are dedicated to temporal capabilities. The model demonstrates superior performance over state-of-the-art methods across three evaluation metrics. The code and an online demo are available at <https://modelscope.cn/models/damo/text-to-video-synthesis/summary>.

## 1 Introduction

Artificial intelligence has expanded the boundaries of content generation in diverse modalities following simple and intuitive instructions. This encompasses textual content [40, 5, 56], visual content [39, 24, 48, 44, 27, 71, 16, 38, 51, 59] and auditory content [4, 32, 31]. In the realm of visual content generation, research efforts have been put into image generation [39, 24, 48, 44] and editing [27], leveraging diffusion models [53].

While video generation [52, 20, 67] continues to pose challenges. A primary hurdle lies in the training difficulty, which often leads to generated videos exhibiting sub-optimal fidelity and motion discontinuity. This presents ample opportunities for further advancements. The open-source image generation methods (*e.g.*, Stable Diffusion [45]) have significantly advanced research in the text-to-image synthesis. Nonetheless, the field of video generation has yet to benefit from a publicly available codebase, which could potentially catalyze further research efforts and progress.

To this end, we propose a simple yet easily trainable baseline for video generation, termed ModelScope Text-to-Video (ModelScopeT2V). This model, which has been publicly available, presents two technical contributions to the field. Firstly, regarding the architecture of ModelScopeT2V, we explore LDM [45] in the field of text-to-video generation by introducing the spatio-temporal block that models temporal dependencies. Secondly, regarding the pre-training technique, we propose a multi-frame training strategy that utilizes both the image-text and video-text paired datasets, thereby enhancing the model’s semantic richness. Experiments have shown that videos generated by ModelScopeT2V perform quantitatively and qualitatively similar or superior to other state-of-the-art methods. We anticipate that ModelScopeT2V can serve as a powerful and effective baseline for future research related to the video synthesis, and propel innovative advancements and exploration.

---

\*Equal contribution.

†Corresponding author.

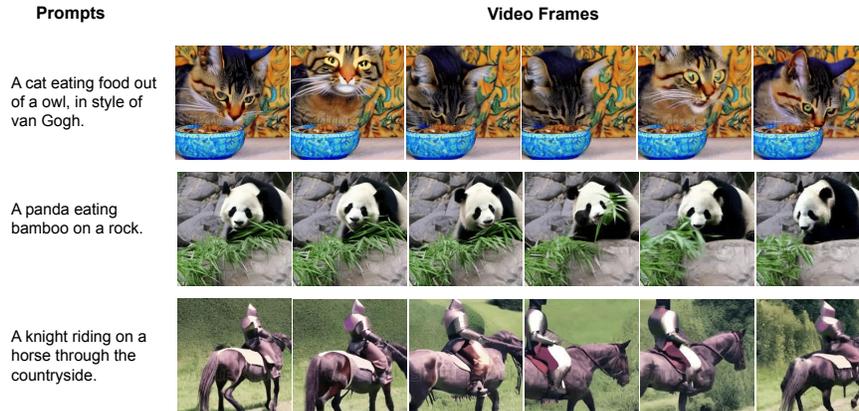


Figure 1: **The qualitative results of our ModelScopeT2V model.** The text-to-video-synthesis model could generate video frames adhering to the given prompts.

## 2 Related work

**Diffusion probabilistic models.** Diffusion Probabilistic Models (DPM) was originally proposed in [53]. The early efforts of utilizing DPM into the image synthesis task at scale has been proven effective [28, 7], surpassing dominant generative models, *e.g.*, generative adversarial networks [13] and variational autoencoders [29], in terms of diversity and fidelity. The original DPM suffers the problem of low-efficiency when adopted for image/video generation due to the iterative denoising process and the high-resolution pixel space. To solve the first obstacle, research efforts focused on improving the sampling efficiency by learning-free sampling [54, 55, 35, 69, 37, 69] and learning based sampling [60, 70, 49]. To address the second obstacle, methods like LDM [46], LSGM [55] and RDM [23] resorted to manifolds with lower intrinsic dimensionality [10, 65]. Our modelScopeT2V follows LDM [46] but modifies it to the video generation task.

**Text-to-image synthesis via diffusion models.** By receiving knowledge from natural language instructions (*e.g.*, CLIP [42] and T5 [43]), diffusion models can be utilized for text-to-image synthesis. LDM [46] designed language-conditioned image generator by augmenting the UNet backbone [47] with cross-attention layers [58]. DALL-E 2 [44] generated image embeddings for a diffusion decoder with CLIP text encoder. The concurrent work, Imagen [48], found the scalability of T5, which means increasing the size of T5 could boost image fidelity and language-image alignment. Building on existing image generation framework [46, 48], Imagic [27] achieves text-based semantic image edits by leveraging intermediate text embeddings that align with the input image and the target text. Composer [24] reformulates images as various compositions, thus enabling image generation from not just texts, but also sketches, masks, depthmaps and more. The ModelScopeT2V initialize the spatial part from Stable Diffusion model [46], and proposes the spatio-temporal block that empowers the capacity of temporal dependencies.

**Text-to-video synthesis via diffusion models.** Generating realistic videos remains challenging due to the difficulty in generating videos with high fidelity and motion continuity [52, 20, 67]. Recent works have utilized diffusion models [53, 37] to generate authentic videos [66, 14, 21, 59]. Text, as a highly intuitive and informative instruction, has been employed to guide video generation. Various approaches have been proposed, such as Video Diffusion [19], which introduces a spatio-temporal factorized 3D Unet [6] with a novel conditional sampling mechanism. Imagen Video [16] synthesizes high definition videos given a text prompt by designing a video generator and a video super-resolution model. Make-A-Video [51] employs off-the-shelf text-to-image generation model combined with spatio-temporal factorized diffusion models to generate high-quality videos without relying on paired video-text data. Instead of modeling the video distribution in the visual space (*e.g.*, RGB format), MagicVideo [71] designed a generator in the latent space with a distribution adapter, which did not require temporal convolution. In order to improve the content and motion performance, VideoFusion [38] decouples per-frame noise into base noise and residual noise, which benefits from a well-pretrained DALL-E 2 and provides better control over content and motion. Targetting

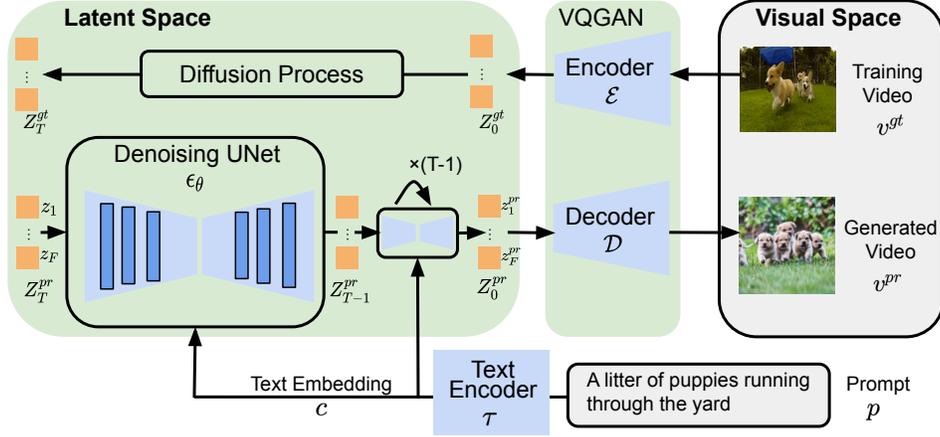


Figure 2: **The overall architecture of ModelScopeT2V.** Here, a text encoder  $\tau$  encodes the prompt  $p$  into text embedding  $c$ . Then, input the embedding  $c$  into the UNet  $\epsilon_\theta$ , directing the denoising process. During training, a diffusion process is performed, transitioning from  $Z_0^{gt}$  to  $Z_T^{gt}$ ; so the denoising UNet could be trained on these latent variables. Conversely, during inference, random noise  $Z_0^{pr}$  is sampled and utilized for the denoising procedure.

decoupling, Gen-1 [8] defines the content latents as CLIP embeddings and the structure latents as the monocular depth estimates, attaining superior decoupled controllability. ModelScopeT2V proposes a simple yet effective training pipeline that benefits from semantic diversity inherent in image-text and video-text paired datasets, further enhancing the learning process and performance in video generation.

### 3 Methodology

In this section, we introduce the overall architecture of ModelScopeT2V (Sec. 3.1), the key spatio-temporal block used in UNet (Sec. 3.2) and a multi-frame training mechanism that stabilizes training (Sec. 3.3).

#### 3.1 ModelScopeT2V

**Structure overview.** Given a text prompt  $p$ , ModelScopeT2V outputs a video  $v^{pr}$  through a latent video diffusion model that conforms to the semantic meaning of the prompt. The architecture of the latent video diffusion model is shown in Fig. 2. As illustrated in the figure, the training video  $v^{gt}$  and generated video  $v^{pr}$  are in the visual space. The diffusion process and denoising UNet  $\epsilon_\theta$  are in the latent space. VQGAN [9] converts the data between visual space and latent space through its encoder  $\mathcal{E}$  and decoder  $\mathcal{D}$ . The latent space in this paper is proposed by VQGAN [9]. Given a training video  $v^{gt} = [f_1, \dots, f_F]$  with  $F$  frames, we could encode the video with VQGAN encoder  $\mathcal{E}$  as,

$$Z_0^{gt} = [\mathcal{E}(f_1), \dots, \mathcal{E}(f_F)], \quad (1)$$

where  $v^{gt} \in \mathbb{R}^{F \times H \times W \times 3}$  is a RGB video,  $Z_0^{gt} \in \mathbb{R}^{F \times \frac{H}{8} \times \frac{W}{8} \times 4}$  is the ground-truth latent variable,  $H$  and  $W$  are the height and width of frames. The denoising UNet works on the latent space. During training, the diffusion process turns  $Z_0^{gt}$  into  $Z_T^{gt}$  by adding gaussian noise  $[\epsilon_1^{gt}, \dots, \epsilon_T^{gt}]$  for  $T$  steps. Therefore, we have  $[Z_0^{gt}, \dots, Z_T^{gt}]$ , which contains less information as the diffusion process proceeds. During inference, the UNet predicts the added noise for each step, so that we finally generate  $Z_0^{pr} = [z_1^{pr}, \dots, z_F^{pr}]$  from a random noise  $Z_T^{pr}$ . Then we could generate a video  $v^{pr}$  by VQGAN decoder  $\mathcal{D}$  as:

$$v^{pr} = [\mathcal{D}(z_1^{pr}), \dots, \mathcal{D}(z_F^{pr})], \quad (2)$$

**Text conditioning mechanism.** ModelScopeT2V aims to generate videos conforming to the given text prompts. Therefore, it is desired to ensure the textual controllability by effectively injecting

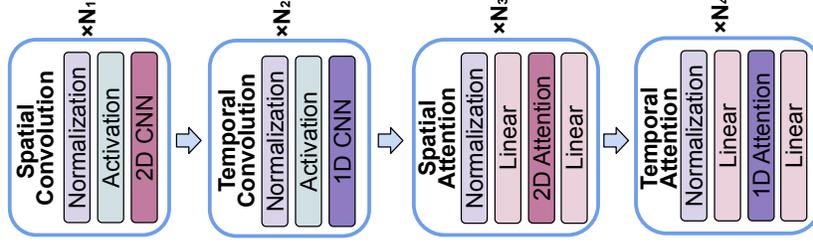


Figure 3: **The structure of the spatio-temporal block.** It includes four modules, *i.e.*, spatial convolution, temporal convolution, spatial attention, and temporal attention. The main layers for these modules are marked in different colors.

textual information into the generative process. Inspired from Stable Diffusion [46], we augment the UNet structure with a cross-attention mechanism, which is an effective approach to condition the visual content on texts [68, 26, 24, 33]. Specifically, we use the text embedding  $c$  of the prompt  $p$  in the spatial attention block as the key and value in the multi-head attention layer. This enables the intermediate UNet features to aggregate text features seamlessly, thereby facilitating an alignment of language and vision embeddings. To ensure a great alignment between language and vision, we utilize the text encoder from pre-trained CLIP ViT-H/14 [42] to convert the prompt  $p$  into the text embedding  $c \in \mathbb{R}^{N_p \times N_c}$  where  $N_p$  represents the maximum token length of the prompt, and  $N_c$  represents the dimension of the token embedding.

**Denoising UNet.** The UNet includes different types of blocks such as the initial block, downsampling block, spatio-temporal block and upsampling block, represented by the dark blue squares in Figure 2. Most parameters of the model are concentrated in the denoising UNet. So the denoising UNet  $\epsilon_\theta$  is considered as the core of the latent video diffusion model, which performs the diffusion process in the latent space. The UNet aims to denoise from  $Z_T$  to  $Z_0$  by predicting the noise of each step. Given a specific step index  $\hat{t} \in [1, 2, \dots, T]$ , the predicted noise  $\epsilon_{\hat{t}}^{pr}$  can be formulated as:

$$c = \tau(p) \quad (3)$$

$$\epsilon_{\hat{t}}^{pr} = \epsilon_\theta(Z_{\hat{t}}, c, \hat{t}) \quad (4)$$

where  $p$  denotes the prompt,  $c$  represent the text embedding, and  $Z_{\hat{t}}$  is the latent variable in  $\hat{t}$ -th step. In this case,  $Z_{\hat{t}}$  denotes  $Z_{\hat{t}}^{gt}$  during training, and denotes the denoised latent video representation  $Z_{\hat{t}}^{pr}$  during inference. The model’s objective is to minimize the discrepancy between the predicted noise  $\epsilon_{\hat{t}}^{pr}$  and ground-truth noise  $\epsilon_{\hat{t}}^{gt}$ . Consequently, the training loss  $\mathcal{L}$  of the UNet can be formulated as:

$$\mathcal{L} = \mathbb{E}_{Z_{\hat{t}}, \epsilon_{\hat{t}}^{gt} \sim \mathcal{N}(0,1), \hat{t}} \left[ \|\epsilon_{\hat{t}}^{gt} - \epsilon_{\hat{t}}^{pr}\|_2^2 \right] \quad (5)$$

which means we hope to decrease the mathematical expectation  $\mathbb{E}$  of the difference between  $\epsilon_{\hat{t}}^{gt}$  and  $\epsilon_{\hat{t}}^{pr}$ .

### 3.2 Spatio-temporal block

Our video diffusion model is built upon a UNet architecture, which consists of four key building blocks: the initial block, the downsampling block, the spatio-temporal block and the upsampling block. The initial block projects the input into the embedding space, while the downsampling and upsampling blocks spatially downsample and upsample the feature maps, respectively. The spatio-temporal block plays a crucial role in capturing complex spatial and temporal dependencies in the latent space, thereby enhancing the quality of video synthesis. To this end, we leverage the power of spatio-temporal convolutions and attentions to comprehensively obtain such complex dependencies.

**Structure overview.** Figure 3 illustrates the architecture of the spatio-temporal block in our video diffusion model. To ensure effectively synthesising videos, we factorise the convolution and the attention mechanism over space and time [41, 1]. Therefore, the spatio-temporal block is composed of four sub-components, namely spatial convolution, temporal convolution, spatial attention and temporal

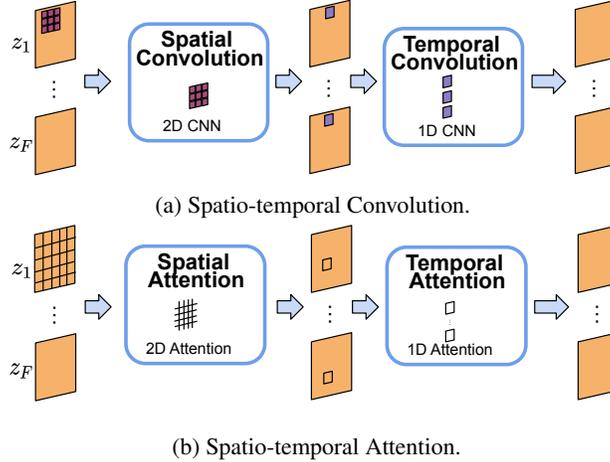


Figure 4: **Diagram of the processing flow for the spatio-temporal block.** Here  $z_i$  denotes the latent variable of the  $i$ -th frame. (a) displays the core structure of spatio-temporal convolution, including a 2D CNN with a kernel size of  $3 \times 3$  and a 1D CNN with a kernel size of 3. (b) shows the variable processing in multi-head attention of spatio-temporal attention, the frame variable is flattened at the spatial scale in spatial attention; while variables at the same positions between frames are grouped together in temporal attention.

attention. The spatio-temporal convolutions capture correlations across frames by convolving over both the spatial and temporal dimensions of the video, while the spatio-temporal attentions capture correlations across frames by selectively attending to different regions and time steps within the video. By utilizing these spatio-temporal building blocks, our model can effectively learn spatio-temporal representations and generate high-quality videos. Specifically, one spatio-temporal block consists of  $N_1$  spatial convolutions,  $N_2$  temporal convolutions,  $N_3$  spatial attention, and  $N_4$  temporal attention operations. In our experiments, we set  $(N_1, N_2, N_3, N_4) = (2, 4, 2, 2)$  by default to achieve a balance between performance and computational efficiency. Regarding the architecture of the two spatial attentions, we instantiate it in two different ways. The first attention is a cross-attention module that conditions the visual features on the textual features, allowing cross-modal interactions. The other attention is a self-attention module that operates solely on visual features, responsible for spatial modeling. While the two temporal attentions are both self-attention module.

**Spatio-temporal convolutions.** Figure 4a illustrates the spatio-temporal convolution, which is composed of both spatial and temporal convolutions. The spatial convolution employs a convolution kernel of size  $3 \times 3$  to extract features from the  $\frac{H}{8} \times \frac{W}{8}$  latent features within each frame, where  $H$  and  $W$  denote the height and width of video frames in visual space. Meanwhile, the temporal convolution adopts a convolution kernel of size 3 to extract features from  $F$  frames, where  $F$  represents the number of frames for each video.

**Spatio-temporal attentions.** Figure 4b displays the spatio-temporal attention, which consists of spatial and temporal attention modules. In detail, the spatial attention operates on the latent features in the spatial dimension of size  $\frac{HW}{64}$ , while the temporal attention operates on the temporal dimension of size  $F$ . We adopt the popular Transformer architecture [58] to instantiate both attention mechanisms.

### 3.3 Multi-frame training

ModelScopeT2V is designed to be trained on large-scale video-text paired datasets, such as WebVid [2], which is domain-aligned with video generation. Nonetheless, the scale of such datasets is orders of magnitude smaller compared to image-text paired datasets, such as LAION [50]. Despite initializing the spatial part of ModelScopeT2V with Stable Diffusion [46], training solely on video-text paired datasets can hinder semantic diversity and lead to catastrophic forgetting of image-domain expertise during training [34, 11, 30, 12]. To overcome this limitation and leverage the strengths of both datasets, we propose a multi-frame training approach. Specifically, one eighth of GPUs for

training are applied to image-text paired datasets, while the remaining GPUs handle video-text paired datasets. Since the model structure could adapt to any frame length, one image could be considered as a video with frame length 1 for those GPUs training on image-text paired datasets.

## 4 Experiments

### 4.1 Implementation details

#### 4.1.1 Datasets

**LAION [50].** We utilize the LAION-5B dataset as image-text pairs, specifically the LAION2B-en subset, as the model focuses on English input. The LAION dataset encompasses objects, people, scenes, and other real-world elements.

**WebVid [2].** The WebVid dataset comprises almost 10 million video-text pairs, with a majority of the videos having a resolution of  $336 \times 596$ . Each video clip has a duration of approximately 30 seconds. During the model training, we selected the middle square portion and randomly picked 16 frames with 3 frames per second as training data.

**MSR-VTT [63].** The MSR-VTT dataset is used to validate our model performance and is not utilized for training. This dataset includes 10k video clips, each of which is annotated with 20 sentences. To obtain FID-vid and FVD metric results, 2,048 video clips were randomly selected from the test set, and one sentence was randomly chosen from each clip to generate videos. When evaluating on CLIPSIM metric, we followed previous works [51, 62] and use nearly 60k sentences from the whole test split as prompts to generate videos.

#### 4.1.2 Model instantiation and hyper-parameters

We employ DDPM [17] with  $T = 1,000$  steps for training and use DDIM sampler [54] in classifier-free guidance [18] with 50 steps for inference by default. ModelScopeT2V primarily consists of three modules: the Text encoder  $\tau$ , VQGAN, and Denoising UNet. The pretrained checkpoint for initializing VQGAN and Denoising UNet are obtained from Stable Diffusion [46] version 2.1<sup>3</sup>. The parameters in VQGAN remains frozen during training and inference. The outputs of temporal convolution and temporal attention are initialized as zeros, enabling ModelScopeT2V to generate meaningful yet temporally discontinuous frames at the beginning of training. As the training progresses, the temporal structures will be optimized to learn the temporal correspondence between frames, thereby synthesising continuous videos.

#### 4.1.3 Training details

We train ModelScopeT2V using the AdamW optimizer [36] with a learning rate of  $5 \times 10^{-5}$ . Our model is trained on 80G NVIDIA A100 GPUs. We perform multi-frame training as detailed in Section 3.3, specifically using a batch size of 1,400 for images and a batch size of 3,200 for videos, and training 267 thousand iterations. The compression factor of VQGAN is 8, meaning that it converts RGB images of size  $256 \times 256$  into latent representations of size  $32 \times 32$ . For the text encoder, we set the maximum text length to  $N_p = 77$ , and embedding dim  $N_c = 768$ , which are consistent with the pre-trained OpenCLIP<sup>4</sup>.

We empirically observe that employing either temporal convolution or temporal attention can augment ModelScopeT2V’s ability to capture temporal dependency. This observation is partly supported by VideoCraft<sup>5</sup> which only contains temporal attention for temporal modeling. We take a step further by employing both the temporal convolution and the temporal attention, which facilitates the ModelScopeT2V to achieve superior temporal modeling. In detail, we use  $N_2 = 4$  temporal convolution blocks and  $N_4 = 2$  temporal attention block for each spatio-temporal block. These temporal blocks account for 552 million parameters out of the total 1,345 million parameters in our UNet, indicating that 39% parameters of the UNet parameters are dedicated to capturing temporal

<sup>3</sup><https://github.com/Stability-AI/stablediffusion>

<sup>4</sup>[https://github.com/mlfoundations/open\\_clip](https://github.com/mlfoundations/open_clip)

<sup>5</sup><https://github.com/VideoCrafter/VideoCrafter>

Models	FID-vid ( $\downarrow$ )	FVD ( $\downarrow$ )	CLIPSIM ( $\uparrow$ )
NÜWA [62]	47.68	-	0.2439
CogVideo (Chinese) [20]	24.78	-	0.2614
CogVideo (English) [20]	23.59	1294	0.2631
MagicVideo [71]	-	1290	-
Video LDM [3]	-	-	0.2929
Make-A-Video [51]	13.17	-	<b>0.3049</b>
ModelScopeT2V (ours)	<b>11.09</b>	<b>550</b>	0.2930

Table 1: **Quantitative comparison with state-of-the-art models on MSR-VTT.** We evaluate the models with three metrics (*i.e.*, FID-vid [15], FVD [57], and CLIPSIM [61]).

information. As a result, the entire ModelScopeT2V model (including VQGAN and the text encoder) comprises approximately 1.7 billion parameters.

We observe that use more layers of temporal convolution would lead to better temporal ability. Since the kernel size of 1D CNN in temporal convolution is 3,  $N_2 = 4$  temporal convolution layers could lead the local receptive field as 81 in each spatio-temporal block, which is enough for 16 output frames per video. For multi-frame training, the temporal convolution and temporal attention mechanisms are still active. Our experiments show it is unnecessary to change the range of parameters for different frame settings.

## 4.2 Main results

### 4.2.1 Quantitative results

ModelScopeT2V is evaluated on MSR-VTT [63] dataset. We conduct the evaluation under a zero-shot setting since ModelScopeT2V is not trained on MSR-VTT. We compare ModelScopeT2V with several state-of-the-art models using FID-vid [15], FVD [57], and CLIPSIM [61] metrics. The FID-vid and FVD are assessed based on 2,048 randomly selected videos from MSR-VTT test split, where we compute the metrics using the middle 16 frames of each video with an FPS of 3. CLIPSIM are evaluated based on all captions from MSR-VTT test split following [51]. The resolution of the generated videos is consistently  $256 \times 256$ .

As shown in Table 1, ModelScopeT2V achieves the best performance on both FID-vid (*i.e.*, 11.09) and FVD (*i.e.*, 550), indicating that our generated videos are visually similar to the ground truth videos. Our model also obtains a competitive score of 0.2930 on CLIPSIM, suggesting that our generated videos are semantically similar to the text prompts. The CLIPSIM score of our model is only marginally lower than that of Make-A-Video [51], while they utilize additional data from HD-VILA-100M [64] for training.

### 4.3 Qualitative results

In this subsection, we compare the qualitative results of ModelScopeT2V with other state-of-the-art methods. To facilitate comparison with Make-A-Video and Imagen Video, generated video frames with the same frame index are presented in the same column with ModelScopeT2V. The videos generated by Make-A-Video [51]<sup>6</sup> and Imagen Video [16]<sup>7</sup> were downloaded from their official webpages. Six frames are uniformly sampled from each video for comparison. This comparison is fair in terms of video duration, as all three methods (*i.e.*, Make-A-Video, Imagen Video, and ModelScopeT2V) generate 16-frame videos aligned with given texts. One difference is that Imagen Video generates videos with a aspect ratio of 2 : 5, while the other two generate videos with a aspect ratio of 1 : 1.

The quantitative comparison between ModelScopeT2V and Make-A-Video is displayed in Figure 5. We can observe that both methods generate videos of high quality, which is consistent with the quantitative results. However, the “robot” in the first example and the “dog” in the second example generated by ModelScopeT2V exhibit a superior degree of realism. We attribute this advantage to our

<sup>6</sup><https://makeavideo.studio>

<sup>7</sup><https://imagen.research.google/video>

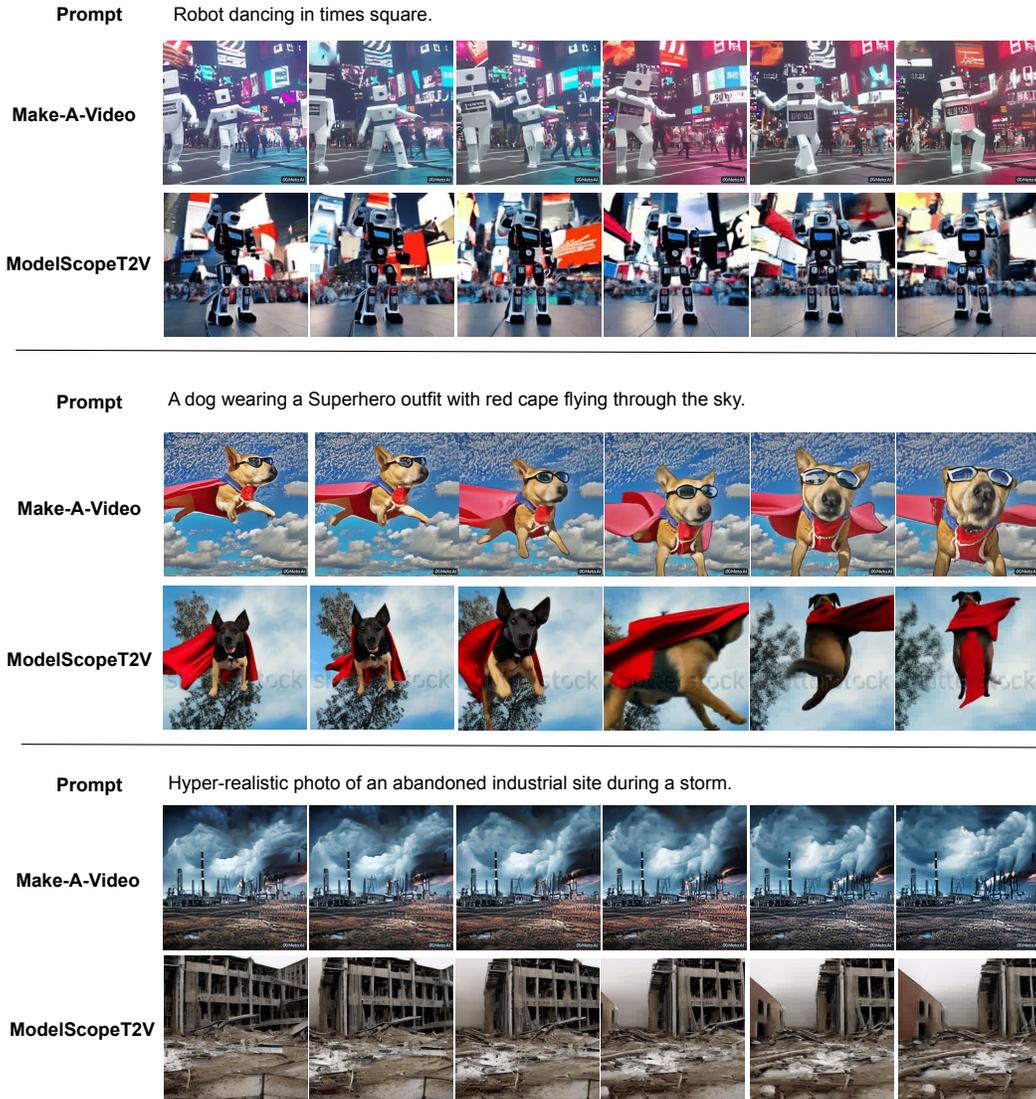


Figure 5: **Qualitative results comparison with Make-A-Video.** We present three examples, each displaying the provided prompt, the video generated by Make-A-Video and the video generated by ModelScopeT2V.

model’s joint training with image-text pairs, which enhances its comprehension of the correspondence between textual and visual data. In the third example, while the “industrial site” generated by Make-A-Video is more closely aligned with the prompt, depicting the overall scene of a “storm”, ModelScopeT2V produces a distinctive interpretation showcasing two “abandoned” factories and a gray sky, captured from various camera angles. This difference stems from Make-A-Video’s use of image CLIP embedding to generate videos, which can result in less dynamic motion information. In general, ModelScopeT2V demonstrates a wider range of motion in its generated videos, distinguishing it from Make-A-Video.

The comparison of our method, ModelScopeT2V, with Imagen Video is illustrated in Figure 6. While Imagen Video generate more vivid and contextually relevant video content, ModelScopeT2V effectively depicts the content of the prompt, albeit with some roughness in the details. For instance, in the first example of Figure 6, it’s noteworthy that Imagen Video generates a video whose second frame illustrates a significantly distorted dog’s tongue, exposing the model’s limitations in accurately rendering the real world. On the other hand, ModelScopeT2V demonstrates its potential in robustly

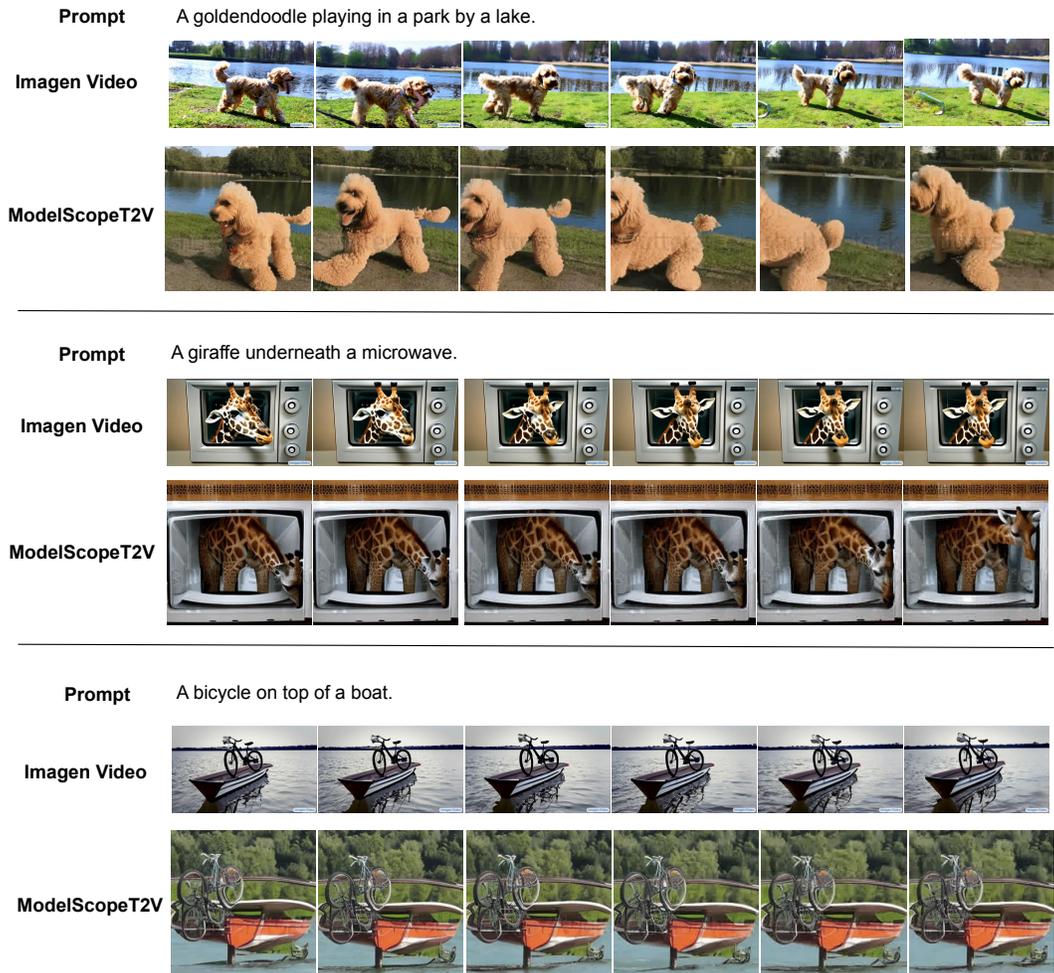


Figure 6: **Qualitative results comparison with Imagen Video.** We present three examples, each displaying the provided prompt, the video generated by Imagen Video and the video generated by ModelScopeT2V.

representing the content described in the prompt. It’s worth noting that the superior performance of Imagen Video can be attributed to the employment of the T5 text encoder, a base model with a larger number of parameters, and a larger-scale training dataset, which is not utilized in ModelScopeT2V. Considering the performance, our ModelScopeT2V lays a strong foundation for future improvements and shows considerable promise in the domain of text-to-video generation.

#### 4.4 Community development

We have made the code for ModelScopeT2V publicly available on the GitHub repositories of ModelScope<sup>8</sup> and Diffuser<sup>9</sup>. Additionally, we have provided online demos of ModelScopeT2V on ModelScope<sup>10</sup> and HuggingFace<sup>11</sup>. The open-source community has actively engaged with our model and uncovered several applications of ModelScopeT2V. Notably, projects such as sd-

<sup>8</sup>[https://github.com/modelscope/modelscope/blob/master/modelscope/models/multi\\_modal/video\\_synthesis](https://github.com/modelscope/modelscope/blob/master/modelscope/models/multi_modal/video_synthesis)

<sup>9</sup><https://huggingface.co/spaces/damo-vilab/modelscope-text-to-video-synthesis/blob/main/app.py>

<sup>10</sup><https://modelscope.cn/studios/damo/text-to-video-synthesis/summary>

<sup>11</sup><https://huggingface.co/spaces/damo-vilab/modelscope-text-to-video-synthesis>

webui-text2video<sup>12</sup> and Text-To-Video-Finetuning<sup>13</sup> have extended the model’s usage and broadened its applicability. Additionally, the video generation feature of ModelScopeT2V has already been successfully utilized for the creation of short videos<sup>14</sup>.

## 5 Conclusion

This paper proposes ModelScopeT2V, the first open-source diffusion-based text-to-video generation model. To enhance the ModelScopeT2V’s ability to modeling temporal dynamics, we design the spatio-temporal block that incorporates spatio-temporal convolution and spatio-temporal attention. Furthermore, to leverage semantics from comprehensive visual content-text pairs, we perform multi-frame training on both text-image pairs and text-video pairs. Comparative analysis of videos generated by ModelScopeT2V and those produced by other state-of-the-art methods demonstrate similar or superior performance quantitatively and qualitatively.

As for future research directions, we expect to adopt additional conditions to enhance video generation quality. Potential strategies include using multi-condition approaches [25] or the LoRA technique [22]. Additionally, an interesting topic to explore could be the generation of longer videos that encapsulate more semantic information.

## References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6836–6846, 2021.
- [2] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *IEEE International Conference on Computer Vision*, 2021.
- [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.
- [4] Zalán Borsos, Raphaël Marinier, Damien Vincent, Eugene Kharitonov, Olivier Pietquin, Matt Sharifi, Olivier Teboul, David Grangier, Marco Tagliasacchi, and Neil Zeghidour. Audioldm: a language modeling approach to audio generation. *arXiv preprint arXiv:2209.03143*, 2022.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [6] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2016: 19th International Conference, Athens, Greece, October 17-21, 2016, Proceedings, Part II 19*, pages 424–432. Springer, 2016.
- [7] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat gans on image synthesis. *Advances in Neural Information Processing Systems*, 34:8780–8794, 2021.
- [8] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. *arXiv preprint arXiv:2302.03011*, 2023.
- [9] Patrick Esser, Robin Rombach, and Björn Ommer. Taming transformers for high-resolution image synthesis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021.

---

<sup>12</sup><https://github.com/deformart/sd-webui-text2video>

<sup>13</sup><https://github.com/ExponentialML/Text-To-Video-Finetuning>

<sup>14</sup><https://youtu.be/Ank49I99EI8>

- [10] Charles Fefferman, Sanjoy Mitter, and Hariharan Narayanan. Testing the manifold hypothesis. *Journal of the American Mathematical Society*, 29(4):983–1049, 2016.
- [11] Tao Feng, Mang Wang, and Hangjie Yuan. Overcoming catastrophic forgetting in incremental object detection via elastic response distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9427–9436, 2022.
- [12] Tao Feng, Hangjie Yuan, Mang Wang, Ziyuan Huang, Ang Bian, and Jianzhou Zhang. Progressive learning without forgetting. *arXiv preprint arXiv:2211.15215*, 2022.
- [13] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 63(11):139–144, 2020.
- [14] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible diffusion modeling of long videos. *arXiv preprint arXiv:2205.11495*, 2022.
- [15] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017.
- [16] 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.
- [17] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in Neural Information Processing Systems*, 33:6840–6851, 2020.
- [18] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *arXiv preprint arXiv:2207.12598*, 2022.
- [19] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022.
- [20] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. *arXiv preprint arXiv:2205.15868*, 2022.
- [21] Tobias Höppe, Arash Mehrjou, Stefan Bauer, Didrik Nielsen, and Andrea Dittadi. Diffusion models for video prediction and infilling. *arXiv preprint arXiv:2206.07696*, 2022.
- [22] 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.
- [23] Chin-Wei Huang, Milad Aghajohari, Joey Bose, Prakash Panangaden, and Aaron C Courville. Riemannian diffusion models. *Advances in Neural Information Processing Systems*, 35:2750–2761, 2022.
- [24] Lianghua Huang, Di Chen, Yu Liu, Yujun Shen, Deli Zhao, and Jingren Zhou. Composer: Creative and controllable image synthesis with composable conditions. *arXiv preprint arXiv:2302.09778*, 2023.
- [25] Lianghua Huang, Di Chen, Yu Liu, Shen Yujun, Deli Zhao, and Zhou Jingren. Composer: Creative and controllable image synthesis with composable conditions. 2023.
- [26] Andrew Jaegle, Felix Gimeno, Andy Brock, Oriol Vinyals, Andrew Zisserman, and Joao Carreira. Perceiver: General perception with iterative attention. In *International conference on machine learning*, pages 4651–4664. PMLR, 2021.
- [27] Bahjat Kawar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. *arXiv preprint arXiv:2210.09276*, 2022.

- [28] Diederik Kingma, Tim Salimans, Ben Poole, and Jonathan Ho. Variational diffusion models. *Advances in Neural Information Processing Systems*, 34:21696–21707, 2021.
- [29] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [30] James Kirkpatrick, Razvan Pascanu, Neil Rabinowitz, Joel Veness, Guillaume Desjardins, Andrei A Rusu, Kieran Milan, John Quan, Tiago Ramalho, Agnieszka Grabska-Barwinska, et al. Overcoming catastrophic forgetting in neural networks. *Proceedings of the national academy of sciences*, 114(13):3521–3526, 2017.
- [31] Zhifeng Kong, Wei Ping, Jiayi Huang, Kexin Zhao, and Bryan Catanzaro. Diffwave: A versatile diffusion model for audio synthesis. In *International Conference on Learning Representations*, 2021.
- [32] Alon Levkovitch, Eliya Nachmani, and Lior Wolf. Zero-shot voice conditioning for denoising diffusion tts models. *arXiv preprint arXiv:2206.02246*, 2022.
- [33] Liunan Harold Li, Pengchuan Zhang, Haotian Zhang, Jianwei Yang, Chunyuan Li, Yiwu Zhong, Lijuan Wang, Lu Yuan, Lei Zhang, Jenq-Neng Hwang, et al. Grounded language-image pre-training. *arXiv preprint arXiv:2112.03857*, 2021.
- [34] Zhizhong Li and Derek Hoiem. Learning without forgetting. *IEEE transactions on pattern analysis and machine intelligence*, 40(12):2935–2947, 2017.
- [35] Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. Pseudo numerical methods for diffusion models on manifolds. In *International Conference on Learning Representations*, 2022.
- [36] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [37] 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. In *Advances in Neural Information Processing Systems*, 2022.
- [38] Zhengxiong Luo, Dayou Chen, Yingya Zhang, Yan Huang, Liang Wang, Yujun Shen, Deli Zhao, Jingren Zhou, and Tieniu Tan. Videofusion: Decomposed diffusion models for high-quality video generation, 2023.
- [39] 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.
- [40] OpenAI. GPT-4 technical report, 2023.
- [41] Zhaofan Qiu, Ting Yao, and Tao Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In *proceedings of the IEEE International Conference on Computer Vision*, pages 5533–5541, 2017.
- [42] 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 *International conference on machine learning*, pages 8748–8763. PMLR, 2021.
- [43] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *The Journal of Machine Learning Research*, 21(1):5485–5551, 2020.
- [44] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- [45] 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.

- [46] 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.
- [47] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention–MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III 18*, pages 234–241. Springer, 2015.
- [48] 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.
- [49] Tim Salimans and Jonathan Ho. Progressive distillation for fast sampling of diffusion models. In *International Conference on Learning Representations*, 2022.
- [50] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- [51] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyan 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.
- [52] Ivan Skorokhodov, Sergey Tulyakov, and Mohamed Elhoseiny. Stylegan-v: A continuous video generator with the price, image quality and perks of stylegan2. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3626–3636, 2022.
- [53] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *International Conference on Computer Vision*, pages 2256–2265. PMLR, 2015.
- [54] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. *arXiv preprint arXiv:2010.02502*, 2020.
- [55] 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.
- [56] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023.
- [57] Thomas Unterthiner, Sjoerd Van Steenkiste, Karol Kurach, Raphael Marinier, Marcin Michalski, and Sylvain Gelly. Towards accurate generative models of video: A new metric & challenges. *arXiv preprint arXiv:1812.01717*, 2018.
- [58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [59] 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 preprint arXiv:2306.02018*, 2023.
- [60] Daniel Watson, William Chan, Jonathan Ho, and Mohammad Norouzi. Learning fast samplers for diffusion models by differentiating through sample quality. In *International Conference on Learning Representations*, 2022.
- [61] Chenfei Wu, Lun Huang, Qianxi Zhang, Binyang Li, Lei Ji, Fan Yang, Guillermo Sapiro, and Nan Duan. Godiva: Generating open-domain videos from natural descriptions. *arXiv preprint arXiv:2104.14806*, 2021.

- [62] Chenfei Wu, Jian Liang, Lei Ji, Fan Yang, Yuejian Fang, Daxin Jiang, and Nan Duan. Nüwa: Visual synthesis pre-training for neural visual world creation. In *European Conference on Computer Vision*, pages 720–736. Springer, 2022.
- [63] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5288–5296, 2016.
- [64] Hongwei Xue, Tiankai Hang, Yanhong Zeng, Yuchong Sun, Bei Liu, Huan Yang, Jianlong Fu, and Baining Guo. Advancing high-resolution video-language representation with large-scale video transcriptions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5036–5045, 2022.
- [65] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Yingxia Shao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *arXiv preprint arXiv:2209.00796*, 2022.
- [66] Ruihan Yang, Prakhar Srivastava, and Stephan Mandt. Diffusion probabilistic modeling for video generation. *arXiv preprint arXiv:2203.09481*, 2022.
- [67] Sihyun Yu, Jihoon Tack, Sangwoo Mo, Hyunsu Kim, Junho Kim, Jung-Woo Ha, and Jinwoo Shin. Generating videos with dynamics-aware implicit generative adversarial networks. *arXiv preprint arXiv:2202.10571*, 2022.
- [68] Hangjie Yuan, Jianwen Jiang, Samuel Albanie, Tao Feng, Ziyuan Huang, Dong Ni, and Mingqian Tang. RLIP: Relational language-image pre-training for human-object interaction detection. In *Advances in Neural Information Processing Systems*, 2022.
- [69] Qinsheng Zhang and Yongxin Chen. Fast sampling of diffusion models with exponential integrator. 2022.
- [70] Huangjie Zheng, Pengcheng He, Weizhu Chen, and Mingyuan Zhou. Truncated diffusion probabilistic models. *stat*, 1050:7, 2022.
- [71] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022.