# MagicVideo: Efficient Video Generation With Latent Diffusion Models

Daquan Zhou*    Weimin Wang*    Hanshu Yan

Weiwei Lv    Yizhe Zhu    Jiashi Feng

ByteDance Inc.

{daquanzhou, weimin.wang, hanshu.yan}@bytedance.com
{vici, yizhe.zhu, jshfeng}@bytedance.com

## Abstract

*We present an efficient text-to-video generation framework based on latent diffusion models, termed MagicVideo. MagicVideo can generate smooth video clips that are concordant with the given text descriptions. Due to a novel and efficient 3D U-Net design and modeling video distributions in a low-dimensional space, MagicVideo can synthesize video clips with 256×256 spatial resolution on a single GPU card, which takes around 64× fewer computations than the Video Diffusion Models (VDM) in terms of FLOPs. In specific, unlike existing works that directly train video models in the RGB space, we use a pre-trained VAE to map video clips into a low-dimensional latent space and learn the distribution of videos' latent codes via a diffusion model. Besides, we introduce two new designs to adapt the U-Net denoiser trained on image tasks to video data: a frame-wise lightweight adaptor for the image-to-video distribution adjustment and a directed temporal attention module to capture temporal dependencies across frames. Thus, we can exploit the informative weights of convolution operators from a text-to-image model for accelerating video training. To ameliorate the pixel dithering in the generated videos, we also propose a novel VideoVAE auto-encoder for better RGB reconstruction. We conduct extensive experiments and demonstrate that MagicVideo can generate high-quality video clips with either realistic or imaginary content. The code will be made public.*

## 1. Introduction

Diffusion-based generative models have shown astonishing achievements over a variety of applications, including text-to-image generation [36, 29, 33] and text-to-3D object generation [26], due to their superior generation quality and scaling capability to large datasets. For example,

DALL-E 2 [29], Imagen [36] and Latent Diffusion [33] models can generate photo-realistic image contents from the given texts after being trained with large-scale image-text datasets (*e.g.*, LAION 400M [38]).

Despite the success in text-to-image generation tasks, using diffusion-based generative models for video generation tasks is still under-explored due to the following difficulties. *1) Data scarcity*. Video data with precise textual descriptions are much harder to collect than image-text data, as videos are more difficult to describe by a single sentence. Besides, different from images carrying compact information, each video may contain some redundant short clips that are less relevant to the textual description. Such information redundancy would limit the effectiveness of video data for model training. *2) Complex temporal dynamics*. Video data contains diverse visual contents across frames and complex temporal dynamics. Therefore, it is much more challenging to model video data distribution compared to static images. *3) High computation cost*. A smooth and informative video may contain more than hundreds of frames. Compared to image generation, directly generating a whole video would consume a huge amount of computational and memory resources.

Existing diffusion-based video generation models propose a cascaded pipeline [13] to deal with the high computational cost. This pipeline generates low-resolution video frames first via an iterative diffusion-based denoising process and then up-samples them by a super-resolution module. Nevertheless, their computational cost is still very high. For example, when generating a coarse video clip of 16 frames and $64 \times 64$ resolution, the recent video diffusion model [13] would take 6-10 seconds with 75G GPU memory[1] for each diffusion iteration. The whole generation process requires tens to hundreds of such iterations for synthesizing a single frame, causing unaffordable time cost.
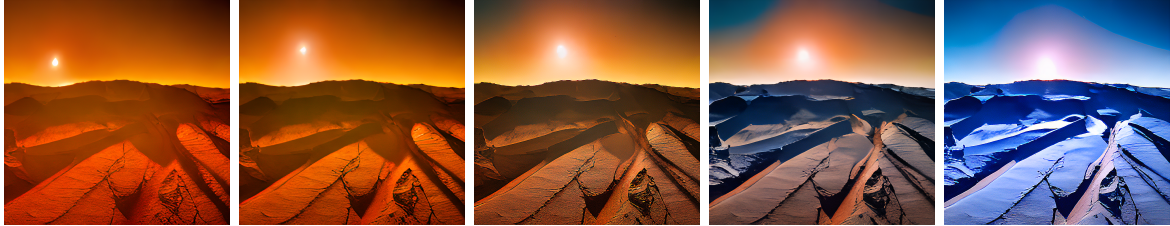
To further *reduce the computational cost of video mod-*

---

*Equal contribution.

[1]We measure the speed on a single Nvidia A100 GPU card

Young attractive woman makeup in the morning.

A beautiful sunrise on mars, Curiosity rover. High definition, timelapse,dramatic colors.

A cute cat.

Melting pistachio ice cream dripping down the cone.

Figure 1: Given various text prompts, MagicVideo produces diverse and temporally-coherent videos that are well-aligned with the prompts. Note that the videos are generated with the key frame model directly without using super-resolution.

*eling*, we alternatively explore using the latent diffusion model (LDM) [33] to learn the distribution of video data, which was developed for image generation and has shown state-of-the-art efficiency in generating images. Specifically, LDM first trains a variational auto-encoder (VAE) [32, 31] to map images into a lower dimensional latent space. Then, it trains a diffusion model to approximate the distributions of the image's latent features instead of the raw RGB images. In this way, the spatial dimension of the latent features undergoing the diffusion-based denoising process is kept low, thus reducing the computational cost significantly. For example, using VAE to reduce the frame resolution by 8 times, the computational cost would be reduced by around $64\times$ for every single frame generation. Motivated by this, we propose adopting the LDM to

build our video generation model, MagicVideo.

To address other aforementioned challenges for video generation, including *data scarcity and complex temporal dynamics modeling*, we introduce the following novel designs to build the first LDM-based video generation model. To *improve data efficiency* and alleviate the demand for paired video-text data, we adopt 2D convolution with temporal computation operators to model the spatial and temporal video features instead of building the model with vanilla 3D [14] or (2+1)D [39] convolutions. This new architecture design allows for initializing 2D convolutions with the parameters of a pre-trained text-to-image model (*e.g.*, LDM [33]) and exploiting its prior knowledge of image generation for facilitating video modeling. Experiments demonstrate that this strategy enables the learning of video

generation even with little video training data. To further reduce the memory cost, we use the same 2D convolutions for generating each frame. However, to avoid deteriorating the temporal consistency in the generated keyframes (*e.g.*, object motion), we introduce a new and lightweight *adaptor* module to adjust the feature distribution per frame. The adaptor only consists of a few scalar parameters yet performs well as it exploits the correlation of the video frames. It effectively alleviates the need for independent 2D convolution blocks for modeling different frames [21].

To model the *temporal dynamics*, our model leverages the *directed self-attention* mechanism. It calculates the features of the future frame based on all the preceding frames while keeping the previous frames unaffected by the future ones. This improves the motion consistency over conventional bi-directional self-attention modules that existing generative models widely use [39, 36, 33, 29]. Furthermore, we propose a novel VideoVAE containing a decoder block that is dedicated to reducing frame generation artifacts (*e.g.*, pixel dithering). Different from the original VAE used in SD that treats each frame independently, VideoVAE considers the temporal relations during the decoding phase, leading to a more consistent high-frequency content.

We conduct extensive experiments to verify the effectiveness of MagicVideo in generating high-resolution videos. MagicVideo can generate photo-realistic video frames with smooth motion and consistent object identity, as shown in Fig. 1, achieving higher quality and efficiency than recent strong methods. We also present its applications to image-to-video and video-to-video generations to show its versatility for various conditional video generation schemes.

## 2. Related works

**Diffusion based generative models.** Denoising diffusion probabilistic models (DDPMs) have achieved great success in image generation [12, 40] and editing [22, 24, 10, 20, 18]. For example, DALL-E-2 [29] uses a generative model (*e.g.* autoregressive or DDPM) to learn the distribution of images' CLIP embeddings and then train a DDPM to synthesize RGB images by conditioning on the sampled embeddings. Alternatively, Imagen [36] directly models the distribution of low-dimensional RGB images via a DDPM and uses cascaded super-resolution models to enhance image qualities. However, due to the intrinsic iterative sampling, the computational overhead of DDPMs gets very high and impedes their applications. To improve the efficiency, Rombach *et al*. 2022 [33] proposed the latent diffusion model (LDM) that models the data distribution in a low-dimensional latent space. Denoising noisy data in a lower dimension may reduce the computational cost in the generation process. Specifically, LDM first trains an au-

toencoder to map images into a low-dimensional space and reconstruct images from latent features. Then, a DDPM with a time-conditional U-Net backbone is used to model the distribution of the latent representations.

**Video generation.** Various video generation methods have been proposed in the past, including GAN-based methods [47, 6, 44] and auto-regressive one [30, 17, 2, 19, 23, 15]. Recently, the success of diffusion-based models for image generation also triggered significant interest in exploring their applications in video modeling. For *unconditional video generation*, Ho *et al*. [14] extended image DDPM models to the video domain by developing a 3D U-Net architecture. Harvey *et al*. [9] proposed to model the distribution of subsequent video frames in an auto-regressive manner. In this work, we are interested in synthesizing videos in a controllable manner, i.e., *text-conditional video generation*. Along this line, Hong *et al*. [16] proposed an auto-regressive framework, CogVideo, that models the video sequence by conditioning itself for the given text and the previous frames. Ho *et al*. [11] proposed a diffusion-based cascaded pipeline, Imagen Video that consists of one base text-to-video module and three spatial and temporal super-resolution modules. Concurrently, Singer *et al*. [39] propose a multi-stage text-to-video generation method, termed Make-A-Video, that first exploits a text-to-image model to generate image embeddings and then trains a low-resolution video generation model with conditioning on the image embeddings, which are then up-sampled via super-resolution models. Both Imagen Video and Make-A-Video model the video distribution in the RGB space. Differently, we explore a more efficient way for video generation by synthesizing videos in a low-dimensional latent space.

## 3. Method

In this section, we introduce the MagicVideo in details (Fig. 2). MagicVideo models video clip distribution in a low-dimension latent space. During the inference stage, MagicVideo first generates key frames in the latent space; then it interpolates key frames to smoothing the frame sequence temporally and maps the latent sequence back to RGB space. Finally, MagicVideo upsamples the obtained video to a high-resolution space for better visual quality.

**Notation.** In this paper, we use $\boldsymbol{x}_t$ to denote a sequence of video frames corrupted with Gaussian noise at intermediate time step $t$. $\boldsymbol{x}_t$ is short for $\boldsymbol{x}_t = [\boldsymbol{x}_t^1, ..., \boldsymbol{x}_t^F]$, where $\boldsymbol{x}_t^i$ represents the $i^{\text{th}}$ frame in the sequence. The encoder and decoder of the proposed video variational auto-encoder (VideoVAE) are denoted by $\mathcal{E}(\cdot)$ and $\mathcal{D}(\cdot)$, respectively. The video frames are mapped into the latent space one by one, *i.e.*, $\boldsymbol{z}_t = [\mathcal{E}(\boldsymbol{x}_t^1), ..., \mathcal{E}(\boldsymbol{x}_t^F)]$. We use CLIP [27] to encode the given text prompt $\boldsymbol{y}$, and the obtained embedding is denoted as $\tau(\boldsymbol{y})$. We use $\epsilon_\theta(\boldsymbol{z}_t, t, \tau(\boldsymbol{y}))$ to represent the
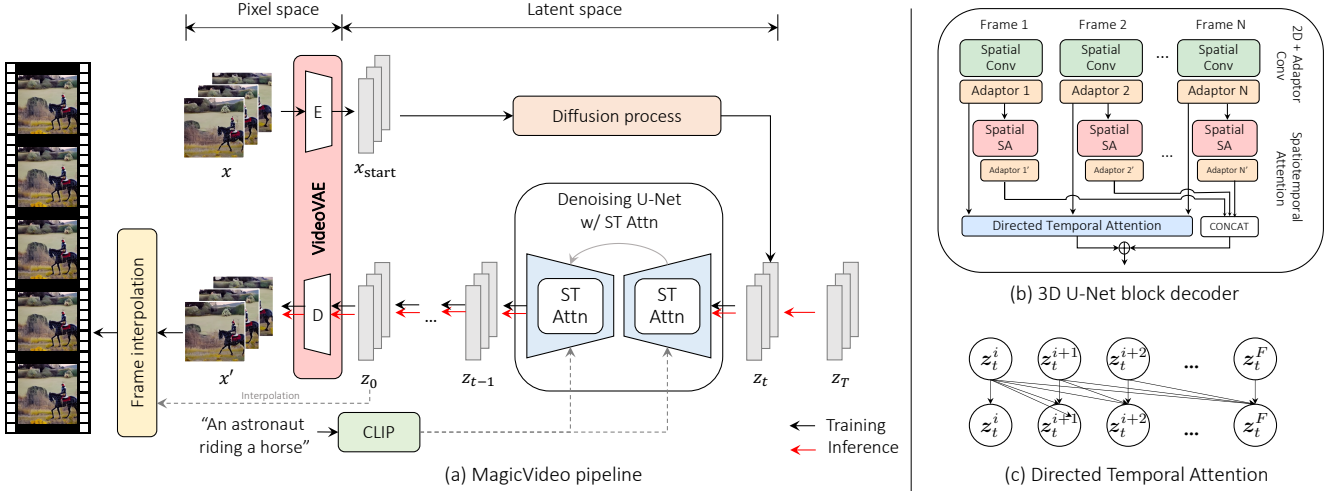
Figure 2: **The overall framework of MagicVideo.** (a) The data flow of both the training and inference phases: during the training phase, a timestep $t$ will be sampled randomly from $[0, T]$ and the input video frames are corrupted via the diffusion process, and a U-Net decoder is used to learn to reconstruct the video frames. Gaussian noise is randomly sampled during inference, and the denoising process is repeated $T$ times. The denoised latent vector $z$ is then fed into a VAE decoder and converted to the RGB space. (b) The structure of the spatiotemporal attention (ST-Attn) module. (c) The directed attention used in the ST-Attn. For the details of the VideoVAE decoder design, please refer to Fig. 3(b).

denoiser of the diffusion model in the latent space.

## 3.1. Key frame generation

The most crucial step of MagicVideo is key frame generation. We use a diffusion model to approximate the distribution of 16 key frames in a low-dimensional latent space. In specific, we design a novel 3D U-Net decoder with an efficient video distribution adaptor and a directed temporal attention module, for video generation. We follow LDM [33] for image generation to add text conditioning via cross-attention, where the text embeddings are used for computing the value and key embeddings, and the intermediate representations of U-Net are used for the query embeddings.

### 3.1.1 Video distribution adaptor

The conventional operator in a neural network model for video data processing is the 3D convolution [5]. However, the computation complexity of 3D convolution is significantly higher than that of 2D convolution. Thus, to reduce the high computational cost, recent video processing models typically replace 3D convolution with a 2D convolution along the spatial dimension followed by a 1D convolution [43] along the temporal dimension (termed "2D+1D").

In this work, we further simplify the operators from "2D+1D" to "2D+adaptor", where the *adaptor* is an even simpler operator compared to the 1D convolution. Specifically, given a sequence of $F$ video frames, we apply a shared 2D convolution for all the frames to extract their spatial features. After that, we adjust the mean and variance for

the intermediate features of every single frame via:

$$z_t^i = S^i \cdot \text{Conv2d}(z_t^i) + B^i, \tag{1}$$

where $z_t^i$ denotes the feature of the $i^{\text{th}}$ frame at denoising time step $t$, and $S, B \in \mathbb{R}^{F \times C}$ are two groups of learnable parameters. This operation is inspired by the observation that frames within each video clip are semantically similar. The small difference among frames may not be necessary for a dedicated 1D convolution layer. Instead, we model those differences via a small group of parameters. The details of the adaptor are shown in Fig. 2(b).

### 3.1.2 Spatial and directed temporal attention

Following previous works [29, 36, 14], within the U-Net, we adopt self-attention modules after the down-sampling blocks that reduce the feature spatial resolution by $4\times$, $8\times$ and $16\times$. The attention operations are conducted along the spatial and temporal dimensions separately. The output of the two parallel attention modules is added and passed to the following modules:

$$z_t = \text{S-Attn}(z_{t-1}) + \text{T-Attn}(z_{t-1}), \tag{2}$$

where S-Attn denotes the attention calculated along the spatial dimension (i.e., to aggregate the frame-wise feature tokens), and T-Attn denotes the self-attention conducted along the temporal dimension. Concretely, the spatial attention is

calculated following previous works [36, 29] via:

$$\text{S-Attn} = \text{Cross-Attn}(\text{LN}(\text{MHSA}(\text{LN}(\boldsymbol{z}_{t-1}))), \tau(\boldsymbol{y})), \tag{3}$$

where MHSA is a standard Multi-head Self-attention module used in vision transformers [7, 51], LN denotes the layer normalization [1], and cross-Attn denotes the cross self-attention module where the attention matrix is calculated between the frame tokens $\boldsymbol{z}_{t-1}$ and the text embedding $\tau(\boldsymbol{y})$. Different from recent VDM [14], we introduce a novel directed self-attention module to better model the video temporal dynamics for the denoising decoder.

**Directed temporal attention.** Recent video generation frameworks [39, 14] mostly use a conventional (i.e., bi-directional) self-attention along the temporal dimension for the motion learning in the video dataset. We notice that the self-attention matrix missed a critical feature of the video data: the motions are directional. In videos, the frames are expected to change in a regular pattern along the temporal dimension. We propose a directed self-attention mechanism to inject the temporal dependency among the frames.

Given a set of given video frames features, $\boldsymbol{z}_t \in \mathbb{R}^{F \times C \times H \times W}$ where $C, F, H, W$ denotes the batch size, number of feature channels, number of frames and the spatial dimension of the features respectively. We first reshape $\boldsymbol{z}_t$ into shape $HW \times \#Heads \times F \times \frac{C}{\#Heads}$ and treat each pixel of each frame as a token, where *#Heads* denotes the number of attention heads. The temporal attention is applied to the tokens of the exact spatial location across different frames to model their dynamics. Specifically, we obtain their query $Q_t$, key $K_t$, and value $V_t$ embeddings for self-attention via three linear transformations, then calculate the temporal attention matrix, $A_t \in \mathbb{R}^{\#Heads \times F \times F}$ via:

$$A_t = \text{Softmax}(Q_t K_t^\top / \sqrt{d}) \odot M, \tag{4}$$

where $d$ is the dimension of embeddings per head and $M$ is an lower triangular matrix with $M_{p,q} = 0$ if $p > q$ else 1. With the mask, the present token is only affected by the previous tokens and independent from the future tokens. Fig. 2(c) illustrates this process.

### 3.1.3 Training strategy

**Frame sampling and training objective.** During training, we first randomly sample a small portion of successive frames (length $L_s$) from each video and read out its frame-per-second (FPS) metadata. Then, we sample 16 frames uniformly from the selected subset as training data. The length of the selected small portion implicitly indicates the speed of motion changing observed within the sampled 16 frames, *i.e.* the longer the subset is, the faster the scene changes. Thus, we compute $\nu = \frac{16}{L_s} \cdot \text{FPS}$ as the new FPS
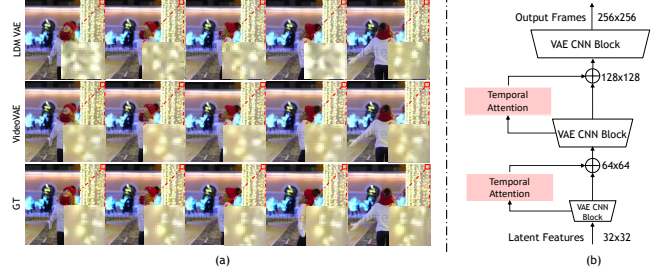


Figure 3: (a) Our proposed VideoVAE for decoding latent features to video frames can effectively reducing pixel dithering artifacts, compared with the conventional VAE model. (b) The architecture of our proposed VideoVAE.

of the 16 frames and use it as an input embedding to MagicVideo. Specifically, we use two linear layers to transform the new FPS $\nu$ into an embedding of dimension $C$:

$$\text{emb}_\nu = \text{Linear}(\text{SiLU}(\text{Linear}(\text{Sin}(\nu)))), \tag{5}$$

where $\text{Sin}(\cdot)$ denotes the sinusoidal position embedding [46]. SiLU denotes the sigmoid ReLU [8]. The embedding $\text{emb}_\nu \in \mathbb{R}^C$ will be added to the video frame features, $\boldsymbol{z} \in \mathbb{R}^{F \times H \times W \times C}$.

We directly use the frame-wise reconstruction loss for the model training. Given a sequence of video frames, the loss of a certain sequence is computed as follows,

$$\mathcal{L}(\boldsymbol{z}_0) = \mathbb{E}_t \sum_{i=0}^{F} \|\boldsymbol{z}_0^i - \epsilon_\theta(\boldsymbol{z}_t^i, t, \boldsymbol{y})\|_2^2. \tag{6}$$

**Unsupervised training scheme.** Text-video pairs are scarce in practice. On the contrary, it is easy to collect abundant high-quality video-only data. Motivated by [39], we adopt an unsupervised training strategy where the embeddings of video frames are used as proxies of text conditions to pretrain the model. The embeddings are extracted using the vision encoder of CLIP [28]. After the unsupervised stage, we finetune the model on a well-annotated video-text paired dataset. The unsupervised and supervised training use the same training objective as defined in Eqn. (6).

### 3.2. Frame interpolation

To increase the temporal resolution and make the generated video smoother, we train a separate frame-interpolation network to synthesize unseen frames between two adjacent key frames. The interpolation model is also trained in the latent space under a similar pipeline as the key frame generation. The difference is that the generation of the interpolated frame features $\boldsymbol{z}$ is conditioned on the adjacent two frames. The conditioning embeddings of the adjacent frames are extracted by CLIP's vision encoder and injected into the cross-attention layers. Besides, we also concatenate

the adjacent two frames' latent embeddings to the randomly sampled noise as input to the interpolation model. We initialize the interpolation U-Net with the key frame generation model for faster convergence. For each pair of two adjacent frames, the interpolation network predicts 3 new intermediate frames between them.

### 3.3. VideoVAE decoder

In LDM [33], RGB images are synthesized by decoding latent features via a pre-trained VAE decoder. In practice, we observe pixel dithering in the generated video frames if we reconstruct videos frame-by-frame via the VAE decoder, leading to visually aesthetic degradation, as shown in Fig. 3.

We empirically find the appearance of dithering relates to the spatial dimension of the latent features: using features with higher dimension suffers less dithering. However, the computational cost will increase if naively increase the feature spatial dimensions. To improve the visual quality without incurring much computational overhead, we keep low dimension of the latent features while adding two temporal directed attention layers in the decoder to build a VideoVAE decoder, as shown in Fig. 3. We find it effectively alleviates the dithering artifacts.

### 3.4. Super-resolution

To generate high-resolution videos, we train a diffusion-based super-resolution (SR) model [37] in RGB space to upsample videos from $256\times 256$ to $1024\times 1024$. The SR model is trained only on image datasets because large-scale high-resolution video datasets are not publicly available and hard to collect. To reduce its computational and memory cost, we train the SR model on $512\times512$ random crops of $1024\times1024$ images with $128\times128$ input frames. During inference, we feed the generated $256\times256$ frames as the input to generate frames with $1024\times1024$ dimension. Chitwan et al. [35] observed noise conditioning augmentation on super-resolution is critical for generating high-fidelity images. Thus, following [35], we degrade low-resolution images with random Gaussian noise and add the noise level as another conditioning signal of the diffusion model.

## 4. Experiments

**Datasets.** We use the weights of LDM [34] pre-trained on Laion 5B [38] to initialize our video 3D U-net denoising decoder. We then conduct unsupervised training on a subset (10M videos) of HD-VILA-100M [50] and Webvid-10M [3]. We fine-tune the video generation model on a subset of self-collected 7M video-text samples. For the ablation study, we randomly sample a subset of 50k videos from Webvid10M to save the computational cost. When comparing with other methods, we evaluate the zero-shot performance with text prompt from the test dataset of UCF-

101 [41], MSR-VTT [49] and calculate the Frechet Inception Distance (FID) [25] and Frechet Video Distance (FVD) [45] with reference to the images in their test dataset. For implementation details, please refer to the sumpplementary material due to space limit.

### 4.1. Analysis

**Ablation studies.** We first investigate the impacts of each proposed component via ablation studies. We randomly sample a subset of 1000 video clips from the test dataset of the Webvid-10M dataset and extract 16 frames from each video to form the reference dataset. The results are shown in Fig. 7. We can observe the directed attention can substantially reduce the FVD. The adaptor not only saves computation cost but also benefits the video generation quality. The unsupervised pre-training can significantly improve the quality—the FVD is further reduced by around 60.

**Spatial and temporal attention.** Different from concurrent works [39], we use a dual self-attention design in parallel: one for spatial self-attention and temporal attention learning in parallel. The detailed architecture for the attention block is shown in Fig. 2(c). To understand the difference between spatial and temporal attention, we visualize the extracted features from each branch separately, and the results are shown in Fig. 4. We find that spatial attention helps generate diverse video frames, and temporal attention generates consistent output among the frames. Combining the spatial and temporal attention guarantees the content diversity and consistency for the generated frames.

### 4.2. Results

**Qualitative evaluation.** We first evaluate our video generation model on qualitative generation performance and compare it with recent state-of-the-art models. The visual results are shown in Fig. 5, where we compare with three strong baselines. We want to highlight that *Make-A-Video is a concurrent work*. Compared with CogVideo [16] and VDM [13], both Make-A-Video and our model can generate videos with richer details. For example, with "Busy freeway at night" as the text input, the videos generated by CogVideo and VDM only show abstract scenes with motion flow without any clear objects (*e.g.*, the cars). Differently, our MagicVideo can generate complex highway objects such as cars with headlights. Moreover, MagicVideo can even generate the perspective phenomenon—the video from our model shows clearer vehicles near the camera. **More samples of the generated videos are provided in the supplementary material**.

**Quantitative evaluation.** We also evaluate MagicVideo quantitatively. Specifically, we pre-train the model on the Webvid-10M dataset. Then we use the text descriptions

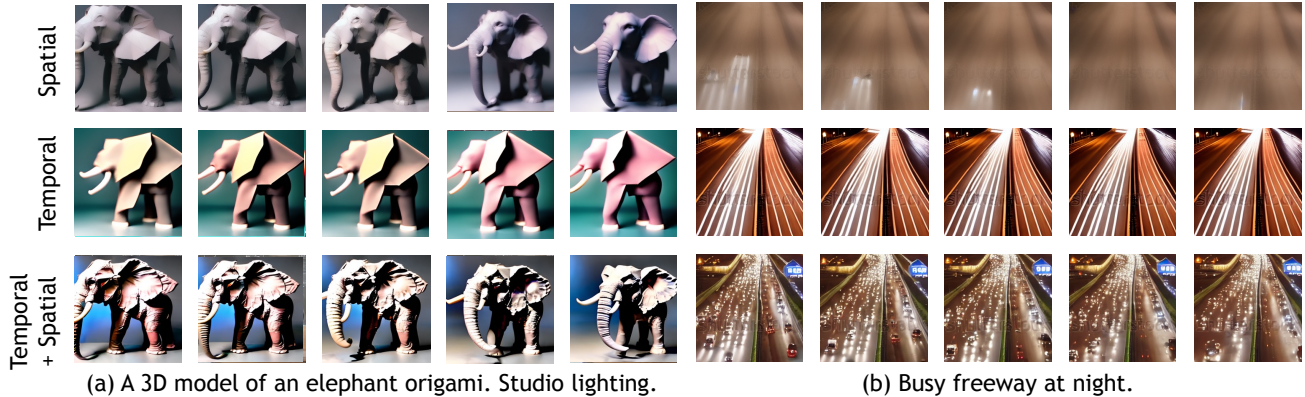(a) A 3D model of an elephant origami. Studio lighting.　(b) Busy freeway at night.

Figure 4: **Illustration of effects of the attention mechanisms.** Spatial attention helps generate diverse frame contents and temporal attention tends to guarantee cross-frame consistency.
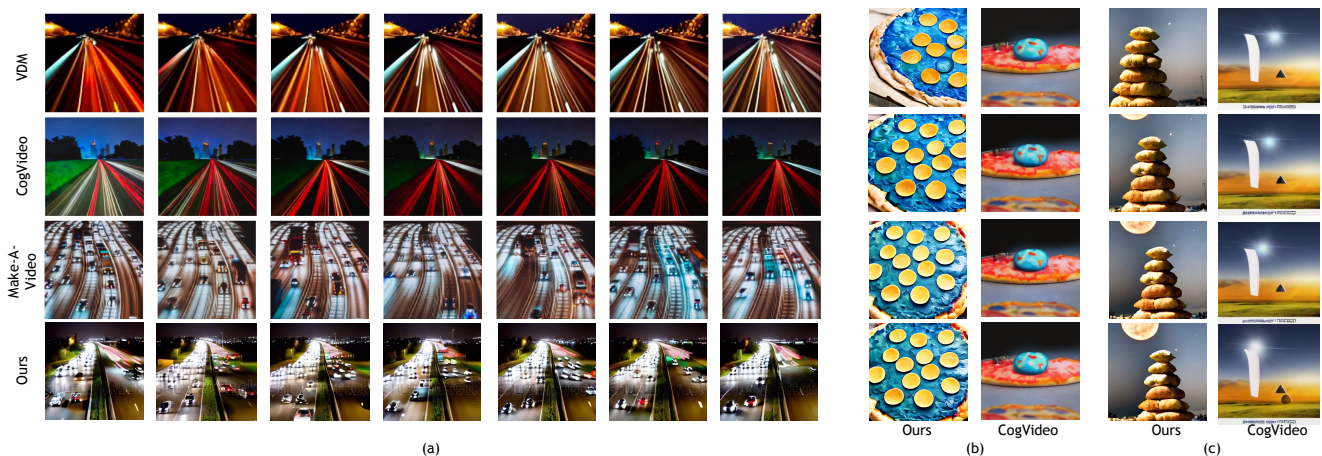


Figure 5: **Qualitative results comparison with recent strong methods**: VDM [14], CogVideo [16], and Make-A-Video [39]. (a) The sample videos are generated with text input "Busy freeway at night", where the samples of the VDM are taken from [39]; (b) & (c) the samples are generated with the text "A blue colored dog" and "A pyramid made of falafel with a partial solar eclipse in the background" respectively. See Tab. 3 for detailed comparison with CogVideo.

Table 1: Video generation evaluation on MSR-VTT.

| Method | Zero-Shot | FID ↓ | FVD ↓ |
|---|---|---|---|
| NÜWA [48] | No | 47.7 | – |
| CogVideo [16] | Yes | 49.0 | 1294 |
| MagicVideo (ours) | Yes | 36.5 | 998 |

Table 2: Video generation evaluation on UCF-101.

| Method | Zero-Shot | FID ↓ | FVD ↓ |
|---|---|---|---|
| MoCoGAN-HD[42] | No | - | 700 |
| CogVideo [16] (Chinese) | Yes | 185 | 751 |
| CogVideo [16] (English) | Yes | 179 | 702 |
| MagicVideo (ours) | Yes | 145 | 655 |

of the test data of MSR-VTT and the class label of UCF-101 validation data as the text prompts to generate 16 key frames for each text prompt without fine-tuning. The comparison between MagicVideo and other recent SOTA methods is shown in Tab. 1 and Tab. 2.

**Human Evaluation.** We also compare to CogVideo (the state-of-the-art open sourced model) on DrawBench [36] by inviting multiple raters. The results are shown in Tab. 3. MagicVideo performs much better than CogVideo [16] with

significantly faster speed.

### 4.3. Applications

We present three applications based on MagicVideo: i) image-to-video generation: given the an input reference image, generating the videos based on the image; ii) video variations: generating a similar video frame sequence based on the input video frames; and iii) video editing: changing the video frame contents based on the input text prompts.

(a) Video-to-video translation

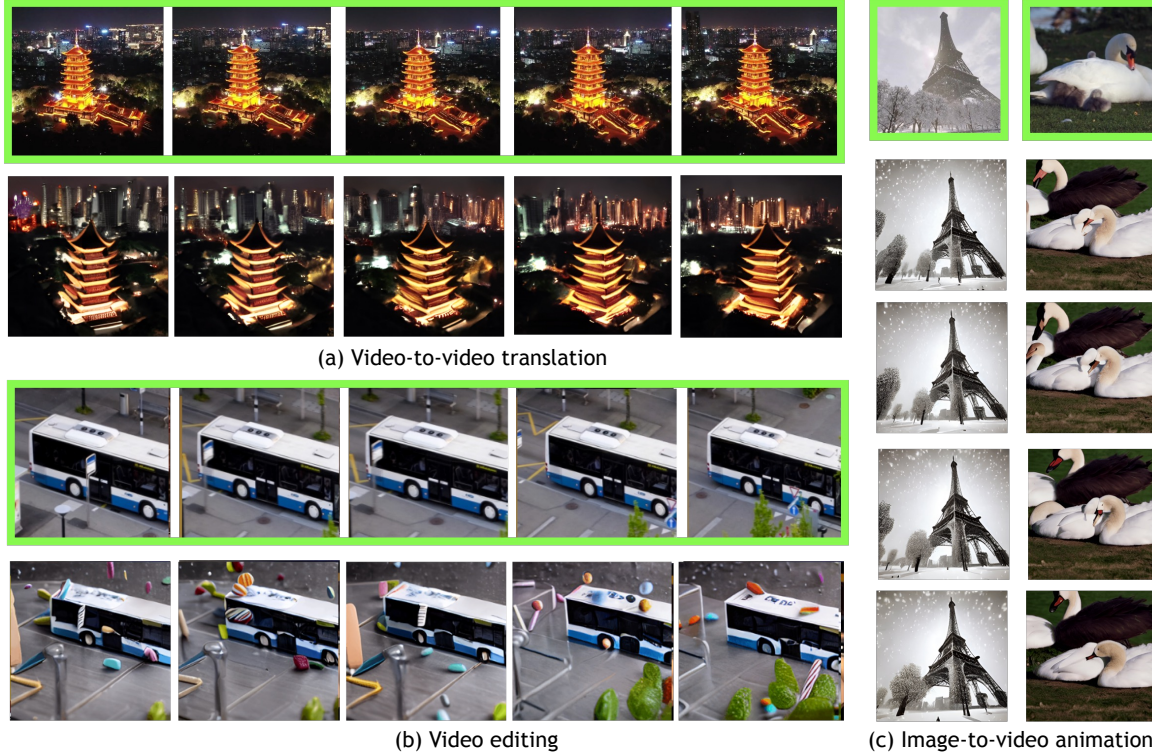(b) Video editing

(c) Image-to-video animation

Figure 6: **Applications based on MagicVideo.** The **green bounding box** denotes the source image/video frames. (a) The video variation results. (b) Given video clip of a driving bus, we use text prompt of "candies falling onto the ground" to MagicVideo for video editing. (c) Given an input image, MagicVideo can generate a short relevant video.
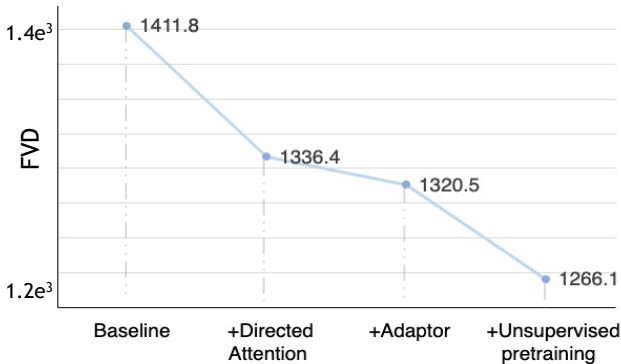


Figure 7: Impacts of different components on the model performance (FVD).

Table 3: Human evaluation comparison between MaigicVideo (ours) and CogVideo [16] on DrawBench [36] on 200 prompts. We evaluate the quality of the generated videos based on four aspects: Realism, Faithfulness, Smoothness, and Efficiency (time cost for generating 16 frames). The numbers in the table show the percentage of samples that the results of our MagicVideo model are better or equal to the samples generated from CogVideo.

| Realism | Faithfulness | Smoothness | Time cost (mins/16 frames) |
|---------|--------------|------------|---------------------------|
| 68 | 86 | 96 | 0.3 vs 21 |

## 5. Conclusions

In this paper, we stepped toward solving the video generation challenge. In particular, we focused on improving the data and computational efficiency of the video generation models. We leveraged the recent latent diffusion model and developed the video generation framework, MagicVideo, in a low-dimensional latent space. Additionally, we introduced several new designs, including the directional attention and the adaptor module, to sufficiently utilize pretrained image generation models. Finally, we demonstrated MagicVideo indeed generates realistic and smooth videos from a text description efficiently.

As shown in Fig. 6(a), with a given image input, MagicVideo is able to generate coherent video frames that are closely related to the main context of the single image input. Fig. 6(b) demonstrates that MagicVideo is able to generate variants of a given video input and Fig. 6(c) shows that by adding some text prompt, MagicVideo can be used to edit a given video. More detailed descriptions on the settings of the three applications are put in the supplementary material.

**Ethical impact.** Video generation may have significant ethical impacts. Besides the applications on generative models for entertainment and art creation, video generation methods are also applicable for malicious purposes by editing videos. However, current deep fake detection technology can detect the fake contents. Another potential issue is using the pre-trained weights from Stable Diffusion [34], which was trained on the LAION dataset [38]. Therefore, it may inherit the LAION dataset contents with ethical issues [4].

# References

[1] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. *arXiv preprint arXiv:1607.06450*, 2016.

[2] Mohammad Babaeizadeh, Mohammad Taghi Saffar, Suraj Nair, Sergey Levine, Chelsea Finn, and Dumitru Erhan. FitVid: Overfitting in Pixel-Level Video Prediction, June 2021. arXiv:2106.13195.

[3] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *ICCV*, pages 1728–1738, 2021.

[4] Abeba Birhane, Vinay Uday Prabhu, and Emmanuel Kahembwe. Multimodal datasets: misogyny, pornography, and malignant stereotypes. *arXiv preprint arXiv:2110.01963*, 2021.

[5] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017.

[6] Aidan Clark, Jeff Donahue, and Karen Simonyan. Adversarial Video Generation on Complex Datasets, Sept. 2019. arXiv:1907.06571.

[7] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

[8] Stefan Elfwing, Eiji Uchibe, and Kenji Doya. Sigmoid-weighted linear units for neural network function approximation in reinforcement learning. *Neural Networks*, 107:3–11, 2018.

[9] William Harvey, Saeid Naderiparizi, Vaden Masrani, Christian Weilbach, and Frank Wood. Flexible Diffusion Modeling of Long Videos. Technical Report arXiv:2205.11495, arXiv, May 2022. arXiv:2205.11495.

[10] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-Prompt Image Editing with Cross Attention Control, Aug. 2022. arXiv:2208.01626.

[11] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen Video: High Definition Video Generation with Diffusion Models, Oct. 2022. arXiv:2210.02303.

[12] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models, Dec. 2020. arXiv:2006.11239.

[13] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *arXiv preprint arXiv:2204.03458*, 2022.

[14] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J. Fleet. Video Diffusion Models. 2022.

[15] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. CogVideo: Large-scale Pretraining for Text-to-Video Generation via Transformers.

[16] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers, 2022.

[17] Nal Kalchbrenner, Aäron Oord, Karen Simonyan, Ivo Danihelka, Oriol Vinyals, Alex Graves, and Koray Kavukcuoglu. Video Pixel Networks. In *Proceedings of the 34th International Conference on Machine Learning*, pages 1771–1779. PMLR, July 2017. ISSN: 2640-3498.

[18] 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, Oct. 2022. arXiv:2210.09276.

[19] Manoj Kumar, Mohammad Babaeizadeh, Dumitru Erhan, Chelsea Finn, Sergey Levine, Laurent Dinh, and Durk Kingma. VideoFlow: A Conditional Flow-Based Model for Stochastic Video Generation. Mar. 2020.

[20] Jun Hao Liew, Hanshu Yan, Daquan Zhou, and Jiashi Feng. MagicMix: Semantic Mixing with Diffusion Models, Oct. 2022. arXiv:2210.16056.

[21] Ziyi Lin, Shijie Geng, Renrui Zhang, Peng Gao, Gerard de Melo, Xiaogang Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Frozen clip models are efficient video learners. *arXiv preprint arXiv:2208.03550*, 2022.

[22] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Re-Paint: Inpainting Using Denoising Diffusion Probabilistic Models. page 11.

[23] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error, Feb. 2016. arXiv:1511.05440.

[24] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations. Technical Report arXiv:2108.01073, arXiv, Jan. 2022. arXiv:2108.01073.

[25] Gaurav Parmar, Richard Zhang, and Jun-Yan Zhu. On aliased resizing and surprising subtleties in gan evaluation. In *CVPR*, 2022.

[26] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988*, 2022.

[27] 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. PMLR, 2021.

[28] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021.

[29] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022.

[30] MarcAurelio Ranzato, Arthur Szlam, Joan Bruna, Michael Mathieu, Ronan Collobert, and Sumit Chopra. Video (language) modeling: a baseline for generative models of natural videos, May 2016. arXiv:1412.6604.

[31] Ali Razavi, Aaron Van den Oord, and Oriol Vinyals. Generating diverse high-fidelity images with vq-vae-2. *Advances in neural information processing systems*, 32, 2019.

[32] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. In *International conference on machine learning*, pages 1278–1286. PMLR, 2014.

[33] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-Resolution Image Synthesis with Latent Diffusion Models, Apr. 2022. arXiv:2112.10752.

[34] 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.

[35] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S Sara Mahdavi, Rapha Gontijo Lopes, et al. Photorealistic text-to-image diffusion models with deep language understanding. *arXiv preprint arXiv:2205.11487*, 2022.

[36] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, S. Sara Mahdavi, Rapha Gontijo Lopes, Tim Salimans, Jonathan Ho, David J Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding, 2022.

[37] Chitwan Saharia, Jonathan Ho, William Chan, Tim Salimans, David J Fleet, and Mohammad Norouzi. Image super-resolution via iterative refinement. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[38] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Theo Coombes, Cade Gordon, Aarush Katta, Robert Kaczmarczyk, and Jenia Jitsev. LAION-5B: laion-5b: A new era of open large-scale multi-modal datasets. https://laion.ai/laion-5b-a-new-era-of-open-large-scale-multi-modal-datasets/, 2022.

[39] 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, Sept. 2022. arXiv:2209.14792.

[40] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. page 36, 2021.

[41] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.

[42] Yu Tian, Jian Ren, Menglei Chai, Kyle Olszewski, Xi Peng, Dimitris N Metaxas, and Sergey Tulyakov. A good image generator is what you need for high-resolution video synthesis. *ICLR*, 2021.

[43] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look

at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018.

[44] Sergey Tulyakov, Ming-Yu Liu, Xiaodong Yang, and Jan Kautz. MoCoGAN: Decomposing Motion and Content for Video Generation, Dec. 2017. arXiv:1707.04993.

[45] Thomas Unterthiner, Sjoerd van Steenkiste, Karol Kurach, Raphaël Marinier, Marcin Michalski, and Sylvain Gelly. Fvd: A new metric for video generation. 2019.

[46] 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.

[47] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Generating Videos with Scene Dynamics, Oct. 2016. arXiv:1609.02612.

[48] 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 *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XVI*, pages 720–736. Springer, 2022.

[49] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. Msr-vtt: A large video description dataset for bridging video and language. In *CVPR*, pages 5288–5296, 2016.

[50] 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 *CVPR*, pages 5036–5045, 2022.

[51] Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animashree Anandkumar, Jiashi Feng, and Jose M Alvarez. Understanding the robustness in vision transformers. In *International Conference on Machine Learning*, pages 27378–27394. PMLR, 2022.