

# MagicVideo-V2: Multi-Stage High-Aesthetic Video Generation

Weimin Wang\*   Jiawei Liu\*   Zhijie Lin   Jiangqiao Yan   Shuo Chen   Chetwin Low  
Tuyen Hoang   Jie Wu   Jun Hao Liew   Hanshu Yan   Daquan Zhou   Jiashi Feng

Bytedance Inc.

<https://magicvideov2.github.io/>

## Abstract

*The growing demand for high-fidelity video generation from textual descriptions has catalyzed significant research in this field. In this work, we introduce MagicVideo-V2 that integrates the text-to-image model, video motion generator, reference image embedding module and frame interpolation module into an end-to-end video generation pipeline. Benefiting from these architecture designs, MagicVideo-V2 can generate an aesthetically pleasing, high-resolution video with remarkable fidelity and smoothness. It demonstrates superior performance over leading Text-to-Video systems such as Runway, Pika 1.0, Morph, Moon Valley and Stable Video Diffusion model via user evaluation at large scale.*

## 1. Introduction

The proliferation of Text-to-Video (T2V) models has marked a significant advancement [6, 9, 11, 15], driven by the recent diffusion-based models. In this work, we propose MagicVideo-V2, a novel multi-stage T2V framework that integrates Text-to-Image (T2I), Image-to-Video (I2V), Video-to-Video (V2V) and Video Frame Interpolation (VFI) modules into an end-to-end video generation pipeline.

The T2I module sets the foundation by producing an initial image from the text prompt, capturing the aesthetic essence of the input. Then the I2V module takes the image as input and outputs low-resolution keyframes of the generated video. The subsequent V2V module increases the resolution of the keyframes and enhances their details. Finally, the frame interpolation module adds smoothness to the motion in the video.

## 2. MagicVideo-V2

The proposed MagicVideo-V2 is a multi-stage end-to-end video generation pipeline capable of generating high-

aesthetic videos from textual description. It consists of the following key modules:

- **Text-to-Image** model that generates an aesthetic image with high fidelity from the given text prompt.
- **Image-to-Video** model that uses the text prompt and generated image as conditions to generate keyframes.
- **Video to video** model that refines and performs super-resolution on the keyframes to yield a high-resolution video.
- **Video Frame Interpolation** model that interpolates frames between keyframes to smoothen the video motion and finally generates a high resolution, smooth, highly aesthetic video.

The following subsections will explain each module in more details.

### 2.1. The Text-to-Image Module

The T2I module takes a text prompt from users as input and generates a  $1024 \times 1024$  image as the reference image for video generation. The reference image helps describe the video contents and the aesthetic style. The proposed MagicVideo-V2 is compatible with different T2I models. Specifically, we use an internally developed diffusion-based T2I model in MagicVideo-V2 that could output high aesthetic images.

### 2.2. The Image-to-Video Module

The I2V module is built on a high-aesthetic SD1.5 [12] model, that leverages human feedback to improve model capabilities in visual quality and content consistency. The I2V module inflates this high-aesthetic SD1.5 with a motion module inspired by [10], both of which were trained on internal datasets.

The I2V module is augmented with a reference image embedding module for utilizing the reference image. More specifically, we adapt an appearance encoder to extract the reference image embeddings and inject them into the I2V module via a cross-attention mechanism. In this way, the image prompt can be effectively decoupled from the text

\*Equal contribution.

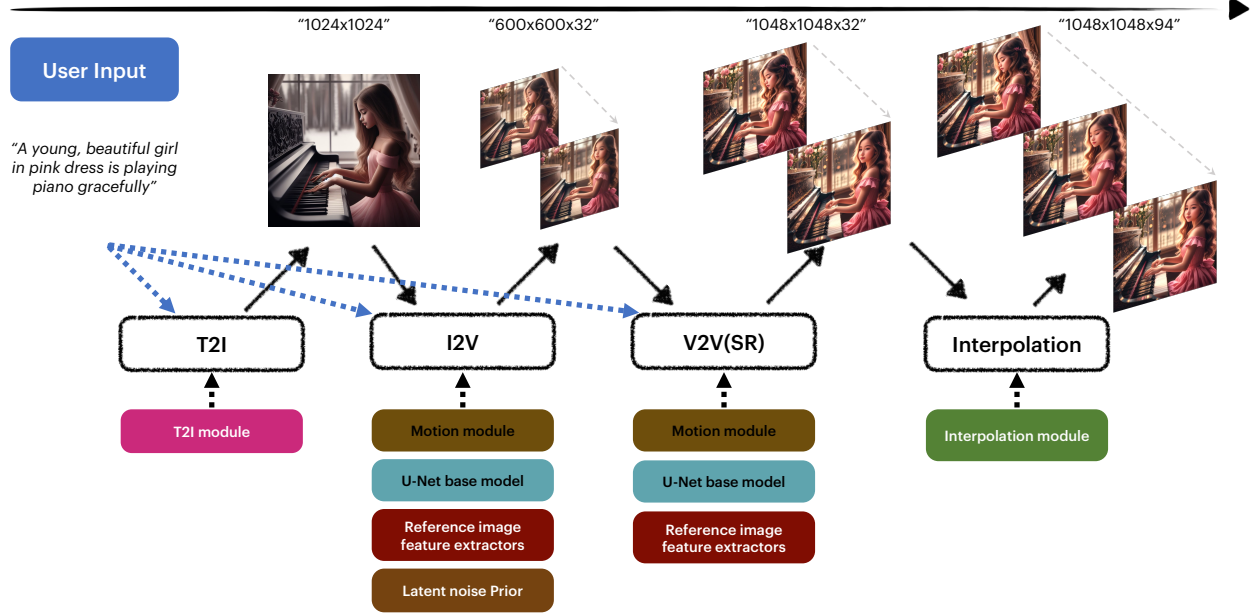


Figure 1. Overview of MagicVideo-V2. The T2I module creates a  $1024 \times 1024$  image that encapsulates the described scene. Subsequently, the I2V module animates this still image, generating a sequence of  $600 \times 600 \times 32$  frames, with the latent noise prior ensuring continuity from the initial frame. The V2V module enhances these frames to a  $1048 \times 1048$  resolution while refining the video content. Finally, the interpolation module extends the sequence to 94 frames, getting a  $1048 \times 1048$  resolution video that exhibits both high aesthetic quality and temporal smoothness.

prompts and provide stronger image conditioning. In addition, we employ a latent noise prior strategy to provide layout condition in the starting noisy latents. The frames are initialized from standard Gaussian noise whose means have shifted from zeros towards the value of reference image latent. With a proper noise prior trick, the image layout could be partially retained and the temporal coherence across frames could also be improved. To further enhance layout and spatial conditioning, we deploy a ControlNet [14] module to directly extract RGB information from the reference image and apply it to all frames. These techniques align the the frames with the reference image well while allowing the model to generate clear motion.

We employ an image-video joint training strategy for training the I2V module, where the images are treated as single-frame videos. The motivation here for joint training is to leverage our internal image datasets of high quality and aesthetics, to improve frame quality of generated videos. The image dataset part also serves as a good compensation for our video datasets that are lacking in diversity and volume.

### 2.3. The Video-to-Video Module

The V2V module has a similar design as the I2V module. It shares the same backbone and spatial layers as in I2V module. Its motion module is separately finetuned using a high-resolution video subset for video super-resolution.

The image appearance encoder and ControlNet module are also used here. This turns out to be crucial, as we are generating video frames at a much higher resolution. Leveraging the information from the reference image helps guide the video diffusion steps by reducing structural errors and failure rates. In addition, it could also enhance the details generated at the higher resolution.

### 2.4. Video Frame Interpolation (VFI)

The VFI module uses an internally trained GAN based VFI model. It employs an Enhanced Deformable Separable Convolution (EDSC) head [7] paired with a VQ-GAN based architecture, similar to the autoencoder model used in the research conducted by [8]. To further enhance its stability and smoothness, we used a pretrained lightweight interpolation model proposed in [13].

## 3. Experiment

### 3.1. Human evaluations

To evaluate MagicVideo-V2, we engaged human evaluators to conduct comparative analyses with contemporary state-of-the-art T2V systems. A panel of 61 evaluators rated 500 side-by-side comparisons between MagicVideo-V2 and an alternative T2V method. Each voter is presented with a random pair of videos, including one of ours vs one of the competitors, based on the same text prompt, for each



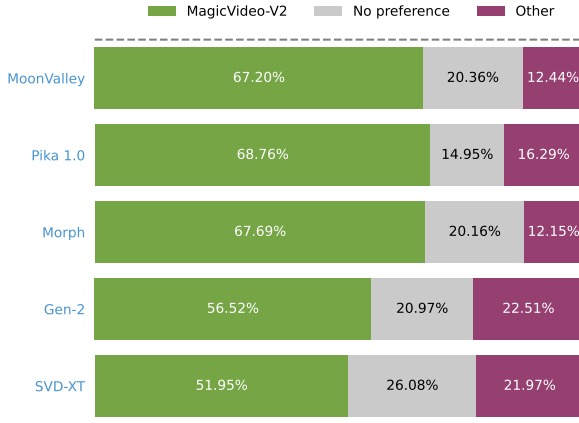


Figure 2. The distribution of human evaluators’ preferences, showing a dominant inclination towards MagicVideo-V2 over other state-of-the-art T2V methods. Green, gray, and pink bars represent trials where MagicVideo-V2 was judged better, equivalent, or inferior, respectively.

round of comparison. They were presented with three assessment options—Good, Same, or Bad—indicating a preference for MagicVideo-V2, no preference, or a preference for the competing T2V method, respectively. The voters are requested to cast their vote based on their overall preference on three criteria: 1) which video has higher frame quality and overall visual appealing. 2) which video is more temporal consistent, with better motion range and motion validity. 3) which video has fewer structure errors, or bad cases. The compiled statistics of these trials can be found in Table 1, with the proportions of preferences depicted in Figure 2. The results demonstrate a clear preference for MagicVideo-V2, evidencing its superior performance from the standpoint of human visual perception.

Method	Good (G)	Same (S)	Bad (B)	(G+S)/(B+S)
MoonValley [2]	4099	1242	759	2.67
Pika 1.0 [4]	4263	927	1010	2.68
Morph [3]	4129	1230	741	2.72
Gen-2 [1]	3448	1279	1373	1.78
SVD-XT [5]	3169	1591	1340	1.62

Table 1. Human side-by-side evaluations comparing MagicVideo-V2 with other state-of-the-art text-to-video generation methods, indicating a strong preference for MagicVideo-V2.

### 3.2. Qualitative examples

Selected qualitative examples of MagicVideo-V2 are presented in Figure 3. For a better-viewed experience, we invite readers to watch the accompanying videos on our



Prompt: A large blob of exploding splashing rainbow paint, with an apple emerging, 8k.



Prompt: An old-fashioned windmill surrounded by flowers, 3D design.



Prompt: A girl with a hairband performing a song with her guitar on a warm evening at a local market, children’s story book.



Prompt: A young, beautiful girl in a pink dress is playing piano gracefully.

Figure 3. Examples of MagicVideo-V2 generated videos via a text prompt.

project website<sup>1</sup>. As mentioned in Section 2, the I2V and V2V modules of MagicVideo-V2 excel at rectifying and refining imperfections from the T2I module, producing smooth and aesthetically pleasing videos. Select examples are showcased in Figure 4.

## 4. Conclusion

MagicVideo-V2 presents a new text-to-video generation pipeline. Our comprehensive evaluations, underscored by human judgment, affirm that MagicVideo-V2 surpasses SOTA methods. The modular design of MagicVideo-V2, integrating text-to-image, image-to-video, video-to-video and video frame interpolation, provides a new strategy for generating smooth and high-aesthetic videos.

## References

- [1] Gen-2. <https://research.runwayml.com/gen2>. Accessed: 2023-11-16. 3
- [2] MoonValley. <https://moonvalley.ai/>. Accessed: 2023-11-16. 3

<sup>1</sup><https://magicvideov2.github.io/>



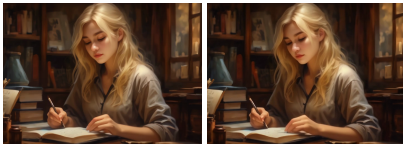
Prompt: “A gray British **Shorthair** skateboarding in Times Square, in cubist painting style.” The wrong dog generated from the T2I module is fixed by the I2V and V2V module.



Prompt: “Ironman flying over a burning city, very detailed surroundings, cities are blazing, shiny iron man suit, realistic, 4k ultra high defi” The ironman’s redundant arm is removed by the I2V and V2V module.



Prompt: “A lone traveller walks in a misty forest.” Left: low resolution video. Right: high resolution video. The tree details and scene brightness are refined by the V2V module.



Prompt: “A girl is writing something on a book. Oil painting style.” Left: low resolution video. Right: high resolution video. The background and aesthetic sense are improved by the V2V module.

Figure 4. Demonstrations of the I2V and V2V modules’ capabilities to correct and refine outputs, leading to polished and visually appealing videos.

- [3] Morph. <https://www.morphstudio.com/>. Accessed: 2023-11-16. 3
- [4] Pika 1.0. <https://pika.art/>. Accessed: 2023-12-26. 3
- [5] SVD-XT. <https://huggingface.co/stabilityai/stable-video-diffusion-img2vid-xt>. Accessed: 2023-11-27. 3
- [6] 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, 2023. 1
- [7] Xianhang Cheng and Zhenzhong Chen. Multiple video frame interpolation via enhanced deformable separable convolution, 2021. 2
- [8] Duolikun Danier, Fan Zhang, and David Bull. Ldmvfi: Video frame interpolation with latent diffusion models, 2023. 2
- [9] Rohit Girdhar, Mannat Singh, Andrew Brown, Quentin Duval, Samaneh Azadi, Sai Saketh Rambhatla, Akbar Shah, Xi Yin, Devi Parikh, and Ishan Misra. Emu video: Factorizing text-to-video generation by explicit image conditioning, 2023. 1
- [10] Yuwei Guo, Ceyuan Yang, Anyi Rao, Yaohui Wang, Yu Qiao, Dahua Lin, and Bo Dai. Animatediff: Animate your personalized text-to-image diffusion models without specific tuning, 2023. 1
- [11] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Rachel Hornung, Hartwig Adam, Hassan Akbari, Yair Alon, Vighnesh Birodkar, Yong Cheng, Ming-Chang Chiu, Josh Dillon, Irfan Essa, Agrim Gupta, Meera Hahn, Anja Hauth, David Hendon, Alonso Martinez, David Minnen, David Ross, Grant Schindler, Mikhail Sirotenko, Kihyuk Sohn, Krishna Somandepalli, Huisheng Wang, Jimmy Yan, Ming-Hsuan Yang, Xuan Yang, Bryan Seybold, and Lu Jiang. Videopoet: A large language model for zero-shot video generation, 2023. 1
- [12] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 1
- [13] Guozhen Zhang, Yuhao Zhu, Haonan Wang, Youxin Chen, Gangshan Wu, and Limin Wang. Extracting motion and appearance via inter-frame attention for efficient video frame interpolation. In *CVPR*, 2023. 2
- [14] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models, 2023. 2
- [15] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models, 2023. 1