

FlowZero: Zero-Shot Text-to-Video Synthesis with LLM-Driven Dynamic Scene Syntax

Yu Lu¹ Linchao Zhu² Hehe Fan² Yi Yang²
¹ReLER Lab, University of Technology Sydney
²CCAI, Zhejiang University
aniki.yulu@gmail.com

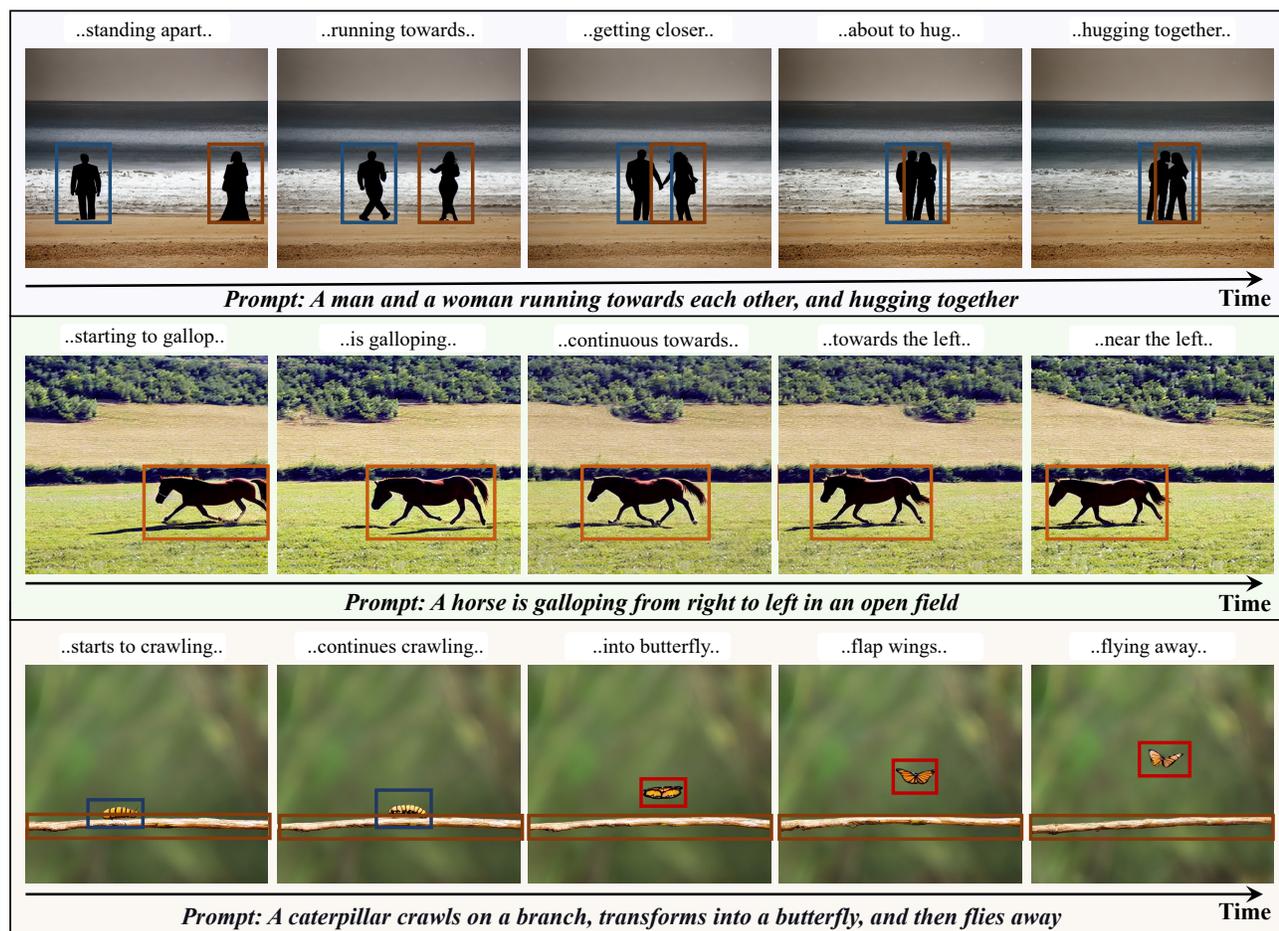


Figure 1. **Zero-shot text-to-video generation.** We present a new framework for text-to-video generation with exceptional temporal coherence, featuring **realistic object movements, transformations, and background motion** within the generated videos.

Abstract

Text-to-video (T2V) generation is a rapidly growing research area that aims to translate the scenes, objects, and actions within complex video text into a sequence of coherent visual frames. We present FlowZero, a novel framework that combines Large Language Models (LLMs) with

image diffusion models to generate temporally-coherent videos. FlowZero uses LLMs to understand complex spatio-temporal dynamics from text, where LLMs can generate a comprehensive dynamic scene syntax (DSS) containing scene descriptions, object layouts, and background motion patterns. These elements in DSS are then used to guide the image diffusion model for video generation with smooth

object motions and frame-to-frame coherence. Moreover, *FlowZero* incorporates an iterative self-refinement process, enhancing the alignment between the spatio-temporal layouts and the textual prompts for the videos. To enhance global coherence, we propose enriching the initial noise of each frame with motion dynamics to control the background movement and camera motion adaptively. By using spatio-temporal syntaxes to guide the diffusion process, *FlowZero* achieves improvement in zero-shot video synthesis, generating coherent videos with vivid motion. Project page: <https://flowzero-video.github.io/>

1. Introduction

In the field of AI-generated content, there has been growing interest in expanding the generative capabilities of pre-trained text-to-image (T2I) models to text-to-video (T2V) generation [5, 9–12, 14, 20, 27, 33]. Recent studies have introduced zero-shot T2V [10, 12, 14], which aims to adapt image diffusion models for video generation without additional training. These methods utilize the ability of image diffusion models, originally trained on static images, to generate frame sequences from video text prompts. However, generating coherent dynamic visual scenes in videos remains challenging due to the succinct and abstract nature of video text prompts.

Meanwhile, Large Language Models (LLMs) demonstrated their capability to generate layouts to control visual modules, especially image generation models [3, 19, 32]. These capabilities indicate a potential for LLMs to understand complex video prompts and generate fine-grained spatio-temporal layouts to guide video synthesis. However, generating spatio-temporal layouts for videos is more intricate, necessitating the LLMs to comprehend and illustrate how objects move and transform over time.

Furthermore, recent research [10, 12] in zero-shot T2V proposes utilizing LLMs to break down video text into frame-level descriptions. These descriptions are crafted to represent each moment or event within the video, guiding image diffusion models to generate semantic-coherent videos. However, these frame-level descriptions only capture the basic temporal semantics of video prompts, lacking detailed spatio-temporal information necessary for ensuring smooth object motion and consistent frame-to-frame coherence in videos. Additionally, representing global background movement to depict camera motion is crucial for immersive video generation [8, 30], which further complicates video generation.

In this paper, we introduce *FlowZero*, a novel framework that integrates LLMs with image diffusion models to generate temporally-coherent videos from text prompts. *FlowZero* utilizes LLMs for comprehensive analysis and translating the video text prompt into a proposed structured

Dynamic Scene Syntax (DSS). Unlike previous methods that only provide basic semantic descriptions, the DSS contains scene descriptions, layouts for foreground objects, and background motion patterns. Foreground layouts contain a series of bounding boxes that define each frame’s spatial arrangement and track changes in the positions and sizes of objects. This ensures that the coherent object motion and transformation align with the textual prompt. Additionally, *FlowZero* incorporates an iterative self-refinement process. This process effectively enhances the alignment between the generated layouts and the textual descriptions, specifically addressing inaccuracies such as spatial and temporal errors. In the self-refinement process, the generated layouts are iteratively compared and adjusted against the text through a feedback loop, ensuring a high fidelity and coherence of the spatio-temporal layouts.

FlowZero prompts LLMs to predict background motion patterns to enhance temporal coherence and consistency, which can be used to control global scenes and camera motion in video frames. For instance, consider a text that describes a horse running from right to left, as shown in the middle example of Figure 1. The LLMs predict a corresponding camera motion, making the background move from left to right, enhancing the video’s immersiveness [8, 30]. The background motion pattern includes specific directions and speeds. We introduce a motion-guided noise shifting (MNS) technique, shifting the initial noise of each frame according to the predicted background motion direction and speed, leading to smoother video synthesis.

FlowZero achieves a significant advancement in zero-shot text-to-video synthesis, utilizing the spatio-temporal planning ability of LLMs to generate detailed frame-by-frame syntax to enhance text-to-video generation. The fusion of these technologies within the *FlowZero* framework enables the generation of temporally-coherent, visually appealing videos directly from textual prompts.

Our contributions are summarised as follows:

- We introduce *FlowZero*, which uses LLMs to convert text into Dynamic Scene Syntax, leading to accurate frame-by-frame video instructions. The framework’s iterative self-refinement process ensures better alignment of spatio-temporal layouts with text prompts, enhancing video synthesis coherence and fidelity.
- The framework improves the global coherence of videos with adaptively controlled background motion through motion-guided noise shifting, increasing the realism of scene and camera motion.
- Through extensive experiments and evaluations, we demonstrate *FlowZero*’s capability to generate temporally-coherent videos that accurately depict complex motions and transformations as described in textual prompts.

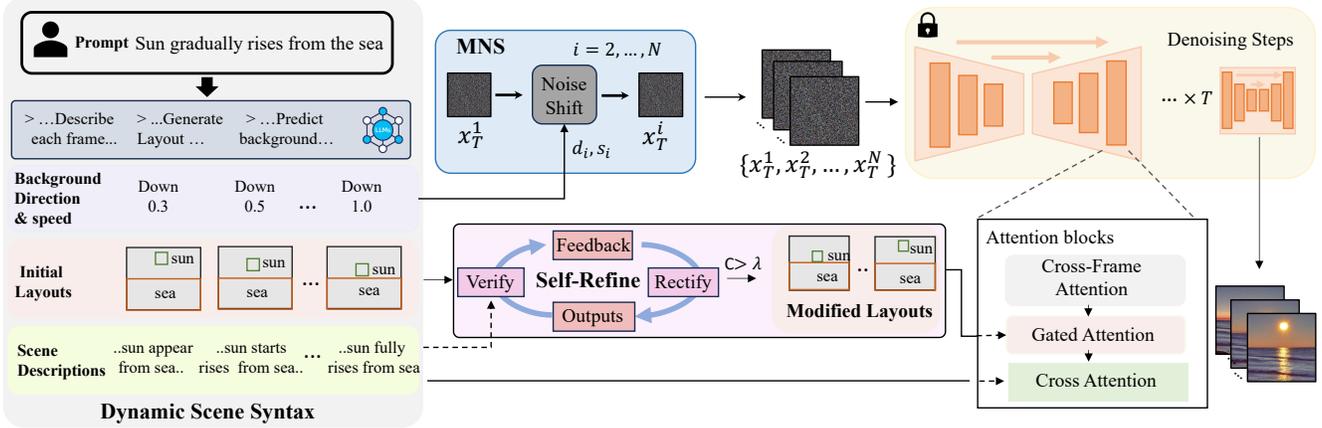


Figure 2. **Overview of FlowZero.** Starting from a video prompt, we first instruct the LLMs (*i.e.*, GPT4) to generate serial frame-by-frame syntax, including scene descriptions, foreground layouts, and background motion patterns. We employ an iterative self-refinement process to improve the generated spatio-temporal layouts. This process includes implementing a feedback loop where the LLM autonomously verifies and rectifies the spatial and temporal errors of the initial layouts. The loop continues until the confidence score C for the modified layouts exceeds a predefined threshold λ . Next, we perform motion-guided noise shifting (MNS) to obtain the initial noise for each frame i by shifting the first noise with predicted background motion direction d_i and speed s_i . Then, a U-Net with cross-attention, gated attention, and cross-frame attention is used to obtain N coherent video frames.

2. Related Work

Text-to-Video Generation. Text-to-video (T2V) generation has evolved from initial variational autoencoders [16, 17] and GANs [4] to advanced diffusion-based techniques [5, 7, 10–12, 14, 27, 33], signifying a major advancement in synthesis methods. Although video diffusion models create high-quality visuals, training T2V models is often computationally expensive. This has led to the exploration of alternative approaches that balance efficiency and quality. Recent advancements [10, 12, 14] have explored leveraging image diffusion models pre-trained on static images [19] to sidestep the demanding training process for T2V. For example, Text-to-Video Zero [14] uses linear transformations and attention mechanisms to maintain video coherence. FreeBloom [12] and Direct2V [10], use Large Language Models (LLMs) to guide image diffusion models with descriptive scene prompts for sequential frames. However, these approaches struggle to capture the intricate object dynamics and background motion of videos, often leading to less expressive and coherent video generation. In contrast to previous methods that only provide semantic descriptions for each frame, FlowZero utilizes LLMs to reason a more comprehensive Dynamic Scene Syntax, delivering detailed, frame-by-frame guidance to enhance the temporal coherence and realism of T2V outputs.

Visual Planning with Large Language Models. Advancements in text-to-image synthesis show that using an intermediate representation, such as a layout or segmentation map, greatly improves the alignment between gener-

ated images and their text descriptions [18, 31]. Various methodologies [3, 19, 29, 32] harness the vast world knowledge embedded in LLMs to craft spatial layouts to guide the image generation process. This has resulted in the creation of images with a reasonable spatial arrangement that closely matches the given textual prompts. For instance, LMD [19] introduces a novel, training-free approach, guiding a diffusion model with a unique controller to generate images based on layouts from LLMs. Similarly, LayoutGPT [3] employs a program-guided strategy, adapting LLMs to cater to layout-driven visual planning across diverse fields. Differing from these methods, FlowZero explores the spatio-temporal planning ability of LLMs for temporally-coherent video generation.

3. Method

As shown in Figure 2, FlowZero initially leverages LLMs (*e.g.*, GPT-4) to process a video prompt \mathcal{T} , generating frame-to-frame scene descriptions, foreground layouts, and background motion patterns. Subsequently, a self-refinement step corrects inconsistencies between the layouts and prompts, such as misaligned movement directions. The frame synthesis begins with a x_T^1 noise sampled from a Gaussian distribution. Then, we perform a motion-guided noise shifting (MNS) to shift the noise to obtain initial noises $\{x_T^1, x_T^2, \dots, x_T^N\}$, encoding background motion direction and speed into each frame. A modified U-Net with various attention mechanisms is employed to synthesize video frames. Finally, through DDIM sampling and a decoder, the final N video frames $\{\mathcal{D}(x_0^i)_{i=1}^N\} \in$

$\mathbb{R}^{N \times H \times W \times 3}$ are generated.

3.1. Dynamic Scene Syntax Generation

In this stage, we aim to use the LLMs, *i.e.*, GPT-4 [21] to convert textual prompts into structured syntaxes for guiding the generation of temporally-coherent videos. These syntaxes include frame-by-frame descriptions, foreground object layouts, and background motion patterns.

- **Scene Descriptions:** Videos often depict a series of continuous events, such as the sunrise, beginning with the “lighting in the edge” and gradually “rising from the horizon”. We propose using LLMs to break down the video text prompt into detailed frame descriptions to depict these events. Given a video text prompt \mathcal{T} , we instruct the LLMs to segment this prompt into detailed scene descriptions $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_N\}$. These descriptions maintain consistent linguistic structures, ensuring that each prompt accurately conveys the visual content in a detailed manner. By providing a description for each frame, we can capture the temporal semantics of the video prompt.
- **Foreground Layout:** While scene descriptions provide semantic details for each frame, these high-level constraints are not sufficient to accurately depict specific object motion and transformations. To achieve coherent object motion, we prompt LLMs to generate a sequence of frame-specific layouts $\{L_1, L_2, \dots, L_N\}$ that outline the spatial arrangement of foreground entities in each frame. These layouts are comprised of bounding boxes that define the position and size of the prompt-referenced objects, using the format: $object : \{x_1, y_1, x_2, y_2\}$. Here, $object$ represents the category of the object along with any relevant attributes (for example, “red car”), (x_1, y_1) and (x_2, y_2) denote the coordinates for the top-left and bottom-right vertices of the bounding box. These layouts provide more fine-grained conditions to ensure the foreground objects adhere to the visual and spatio-temporal cues the text provides.
- **Background Motion:** Background motion plays a crucial role in enhancing the global coherence of videos, especially when dynamic foreground objects are involved. For example, in a video showing a horse running to the left, synchronizing the camera motion with the horse’s direction can create a visually smooth effect, making the video more immersive and engaging [8, 30]. To effectively simulate this, we first categorize potential background motion into eight moving directions: {left, right, up, down, left_up, left_down, right_up, right_down}, and include a “random” option for non-directional movement. We also define a motion speed that ranges from 0 (no movement) to 1.0 (rapid movement). We use LLMs to determine the most appropriate background motion direction and speed for each frame. This helps us align it with the foreground movements as described in the scene. By

integrating background motions, we ensure global coherence and consistency in video sequences.

Based on previous studies [3, 10, 12, 29, 32], we instruct LLMs to generate these syntaxes through direct commands. For example, we use prompts like “describe each frame” to create descriptions and “generate layouts for each scene” to generate foreground layouts. We provide an example in context to enhance the stability and effectiveness of LLMs.

Iterative Self-Refinement. Due to the complex nature of reasoning in spatio-temporal dynamics, there may be discrepancies between the generated spatio-temporal layouts and the textual prompts. As illustrated in Figure 2, the sun initially moves downward over time, which contradicts the video prompt “sun gradually rises”. Previous research has shown that LLMs can verify and correct generated texts or codes [2, 15, 26, 28]. Inspired by this, we propose an iterative self-refinement process to address potential misalignments between the initial spatio-temporal layouts $\{L_1, L_2, \dots, L_N\}$ and the textual prompts $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_N\}$. The initial step of self-refinement involves prompting LLMs to verify spatial and temporal consistency between scene descriptions and layouts and provide detailed feedback. This feedback includes an analysis of problems, specific suggestions (*e.g.*, the sun should rise up instead of going down), and a confidence score c from 1 to 5 to measure the alignment of layouts with descriptions. We found that providing numerically supported analysis and suggestions enhances the effectiveness of the self-refinement process. For example, examining and comparing particular coordinates of the bounding boxes can be particularly helpful. We include in-context examples that clearly demonstrate the type of feedback most helpful for generating specific suggestions. Then, we prompt LLMs again to correct the layouts in the rectification step to improve spatial and temporal alignment with the textual prompts. This refinement process consists of multiple iterations, with the LLMs verifying and rectifying the layouts based on the previous iteration, leading to convergence toward an optimal layout representation. The iterations continue until the confidence score c is higher than a predefined alignment threshold λ .

3.2. Video Synthesis from Dynamic Scene Syntax

In this section, we seek to generate coherent video frames based on the generated DSS. As shown in the right part of Figure 2, beginning with a noise x_T^1 from standard Gaussian distribution, we first conduct Motion-guided noise shifting to obtain initial noises $\{x_T^1, x_T^2, x_T^3, \dots, x_T^N\}$ for each frame. These noises are obtained by shifting noise x_T^1 to match the background motion direction and speed predicted in the DSS generation stage. We will provide detailed information on Motion-guided noise shifting below. We employ a

modified U-Net with cross attention, gated attention, and cross-frame attention mechanisms [10, 12, 14]. The cross-attention mechanism within the U-Net is designed to input scene descriptions, enabling the capture of diverse semantics for each frame. Simultaneously, the gated attention [18] inputs foreground layouts into the U-Net, managing the arrangement of objects across different frames. We then convert the self-attention layer in the U-Net of the image diffusion model into cross-frame attention [10, 12, 14], which performs attention between the query frame and previous frames.

Motion-guided Noise Shifting. Previous method [14] performs a linear transformation on initial noises with fixed direction and speed to model global motion dynamics in video frames. In contrast, our approach allows LLMs to predict the background motion direction and speed adaptively for transforming noises, thereby significantly enhancing the global temporal coherence of videos.

Given the predicted background motion d and speed s , a straightforward method is directly shifting the noise spatially for each frame. However, this often results in abrupt changes in low-level visual effects, such as color and lighting alterations in the video frames. To address this problem, we propose a technique to shift the phase of noises in the frequency domain [13]. This method preserves the amplitude component to maintain low-level visual effects while modulating the phase component to simulate spatial noise shifting [6, 22–24], achieves smoother video frames.

Specifically, for each frame i , we use the predicted background motion direction (d_i) and speed (s_i) to guide the spatial shift of noise. This is achieved by modulating the phase component of noise x_T^1 in the frequency domain, T means the total diffusion step. The mathematical formulation is as follows:

$$x_T^i = \mathcal{F}^{-1} \left(\mathcal{F}(x_T^1) \cdot e^{-j \cdot 2\pi \cdot (i \cdot s_i) \cdot (d_y f_y + d_x f_x)} \right),$$

where \mathcal{F} and \mathcal{F}^{-1} denote the Discrete Fourier Transform (DFT) and its inverse, respectively. The frequencies in the y and x dimensions are represented by f_y and f_x . The direction multipliers d_y and d_x are derived based on the motion direction d_i . For instance, if the direction is “left,” d_x and d_y would be set to $\{0, 1\}$. When the background motion direction remains the same across frames, the index will increase linearly, resulting in smooth motion effects.

In scenarios with non-directional movements ($d = \text{random}$), such as “a goldfish swimming in a fish bowl,” setting a static background for all frames can be unrealistic. Our method addresses this by adding random disturbances at the phase components of all frequencies, simulating natural scene variability, and enhancing video realism.

We perform the noise-shifting technique in the frequency domain has several advantages:

1. Our method allows for easy modification of moving directions by adjusting the direction multipliers, offering greater flexibility than direct space shifting.
2. Since our technique operates in the frequency domain, it is more efficient and computationally less intensive than spatial domain transformations, particularly for handling high-resolution and long videos.
3. We can simulate realistic motionless scenes by adding random disturbances to the noises.

4. Experiments

4.1. Implementation Details

We utilize the GLIGEN [18] as the base image diffusion model, which is pre-trained to generate images adhering to a layout. We employ GPT-4 [21] to reason Dynamic Scene Syntax (DSS). In our tests, we generate $N = 8$ frames per video, each with a resolution of 512×512 . However, our framework allows for generating any desired number of frames by instructing LLMs and increasing N . We set a threshold λ of self-refinement as 3 and the maximum iteration as 5. All experiments are conducted on a single NVIDIA V100 GPU.

4.2. Comparisons with Baseline Methods

Qualitative Comparison In our qualitative comparative analysis, we compare videos generated using our FlowZero with several benchmark methods: zero-shot based methods, T2V-Z [14], DirecT2V [10], and training-based methods AnimateDiff [5], VideoFusion [20]. We assess the performance in three scenarios: basic object motion rendering, multiple object motion depiction, and complex object transformations.

In our initial assessment, shown in Figure 3, we analyze videos generated from prompts featuring basic object motion. FlowZero effectively demonstrates the ability to depict smooth object motion, particularly showcasing a butterfly’s departure from a flower. However, other zero-shot techniques, such as T2V-Z and DirecT2V, only capture temporal semantics and struggle to model the coherent object motion. AnimateDiff and VideoFusion were trained on extensive video-text data [1] and exhibit temporal frame coherence. However, they fall short in rendering nuanced motion details, resulting in slightly stilted animations. In scenarios involving multiple objects with designated movements, also presented in the right of Figure 3, FlowZero continues to excel, accurately animating specific objects and motion defined in the text prompts. Other methods struggle to precisely replicate the specified objects and movements, often resulting in a less accurate portrayal. In Figure 4, we compare all methods of generating videos from prompts that describe complex object transformations. FlowZero distinguishes itself by vividly rendering transformations, such as



Figure 3. **Qualitative comparison.** Our method can capture detailed object motion to generate temporally coherent frame sequences.

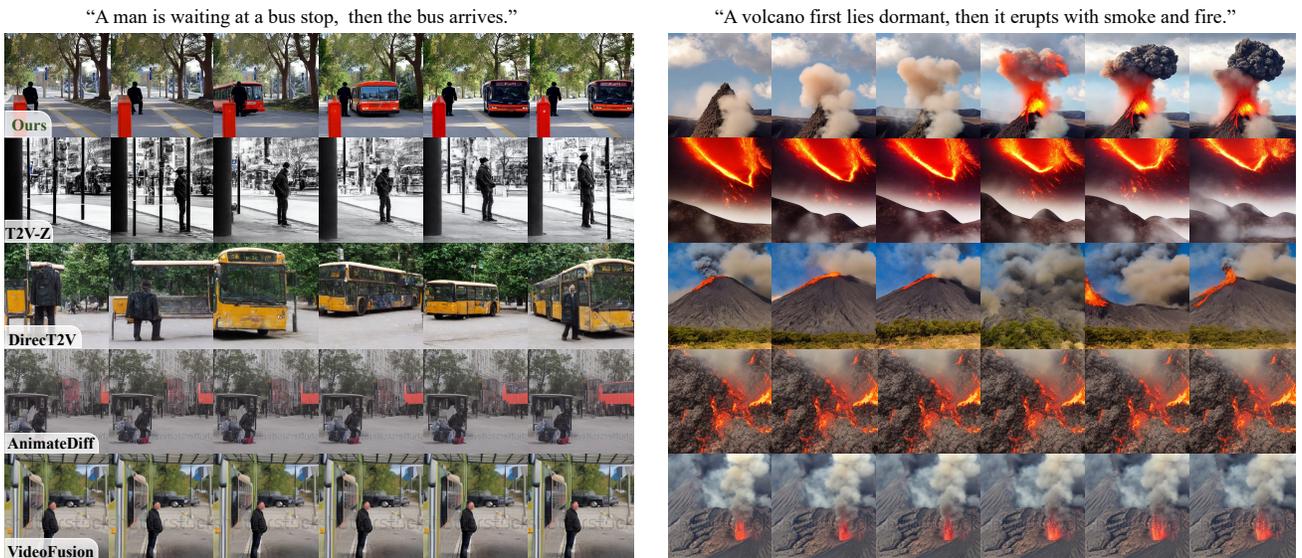


Figure 4. **Qualitative comparison.** Our method can model intricate object transformations representing narrative structures in the video prompt.

a tranquil volcano erupting. Other methods do not effectively translate these temporal dynamics, leading to less coherent visual transformations.

By utilizing LLMs to plan the spatio-temporal syntax as guidance for the diffusion model, FlowZero surpasses other methods in text-prompted video generation. Its superior performance is particularly noticeable in the accurate motion of multiple objects and intricate object transformations.

Quantitative Comparison As shown in Table 1, we first compare our methods with other four baseline methods,

i.e., AnimateDiff [5], VideoFusion [20], T2V-Z [14], DirecT2V [10] using CLIP score metrics [25]. The CLIP metrics measure the semantic similarity between the text and video frames. Our method achieves the highest performance by prompting LLMs to deduce more semantics from both spatial and temporal dimensions.

Due to the complexity of quantitatively evaluating videos with intricate temporal dynamics, we conducted a user study to validate the effectiveness of our method. We recruited 20 people from academia and industry to conduct this survey. We ask users to provide feedback on the se-

Table 1. **Quantitative Results.** We perform automatic metrics, i.e., CLIP score and user study, to validate the effectiveness.

Method	Training-Free	Automatic Metric		User Study		
		CLIP Score \uparrow	Semantic \uparrow	Temporal \uparrow	Quality \uparrow	Rank \downarrow
AnimateDiff [5]		0.244	3.15	2.75	2.97	3.42
VideoFusion [20]		0.264	3.38	2.92	3.11	3.17
T2V-Z [14]	✓	0.245	3.29	2.99	3.03	3.19
DirectT2V [10]	✓	0.244	3.39	3.29	2.52	2.97
Ours	✓	0.267	4.57	4.58	4.40	2.00

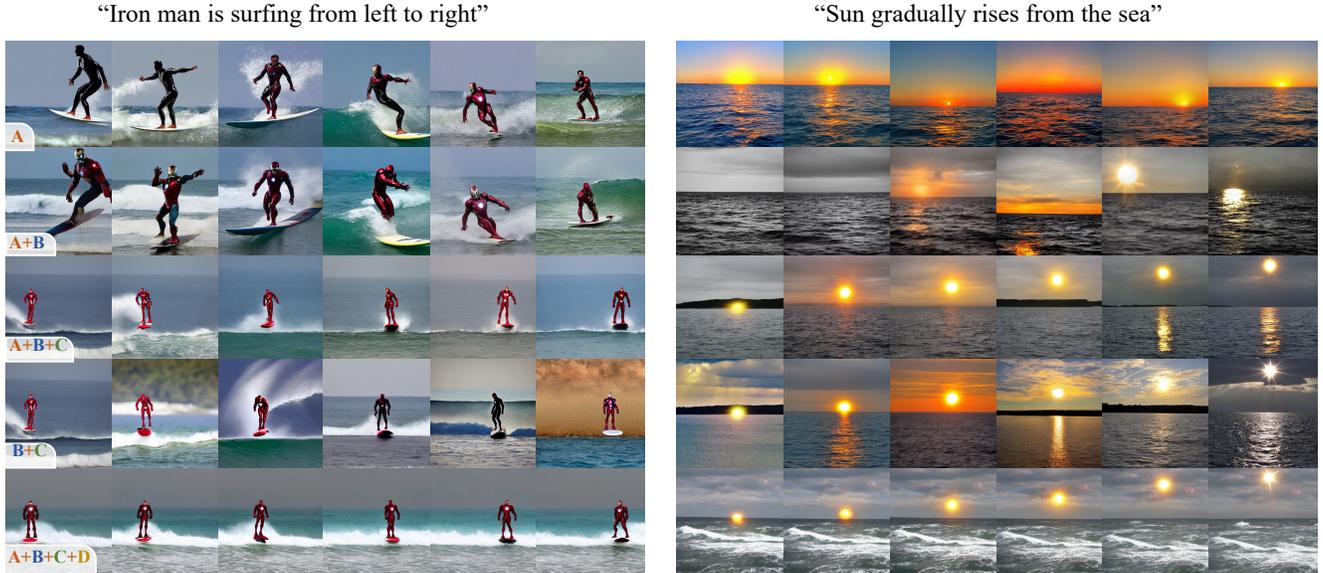


Figure 5. Ablation studies of the effectiveness of FlowZero. (A) cross-frame attention, (B) scene descriptions, (C) foreground layouts, (D) motion-guided noise shifting.

semantic accuracy, temporal coherence, and video quality of the videos generated by five different methods. It is evident that users prefer our method and achieve better results, surpassing even training-based methods, *e.g.*, AnimateDiff [5] and VideoFusion [20]. Furthermore, our methods achieve significant improvement over other zero-shot methods, *e.g.*, T2V-Zero [14] and DirectT2V [10] on temporal coherence, which validates the effectiveness of our approach.

4.3. Ablation Study

Effectiveness of FlowZero In Figure 5, we conduct a comprehensive ablation study to validate the effectiveness of key components in FlowZero. This includes cross-frame attention, scene descriptions, foreground object layouts, and background motion for noise shifting. We begin with a baseline model that employs cross-frame attention to adapt U-Net, feeding it original video prompts alongside independent random noise for each frame. Row #1 of Figure 5 demonstrates that the baseline generates videos with basic semantics like Ironman and surfing but fails to capture de-

tailed object motion. In row#2, we replace the video prompt with generated scene descriptions for each frame, similar to previous methods Free-Bloom [12] and DirectT2V [10]. However, we found that merely using temporal semantics resulted in a lack of coherent object motion, resulting in inconsistencies across frames. Instead, in row#3, by adding the layout to constrain the arrangement of foreground objects, we can clearly capture the coherent motion of the main object, such as “from left to right” and “rises”. However, the video frames still display temporal inconsistency, such as in color, lighting, and global scene. In row#4, we experimented with removing the cross-frame attention from U-Net, which means relying solely on a pure image diffusion model [18]. This modification resulted in a lack of inconsistencies in the representation of objects and backgrounds across frames, even though the layouts guide the object motion. Finally, by utilizing our motion-guided noise shifting technique in row#5, we can smoothly control the motion direction and speed in the background, resulting in a coherent global scene.

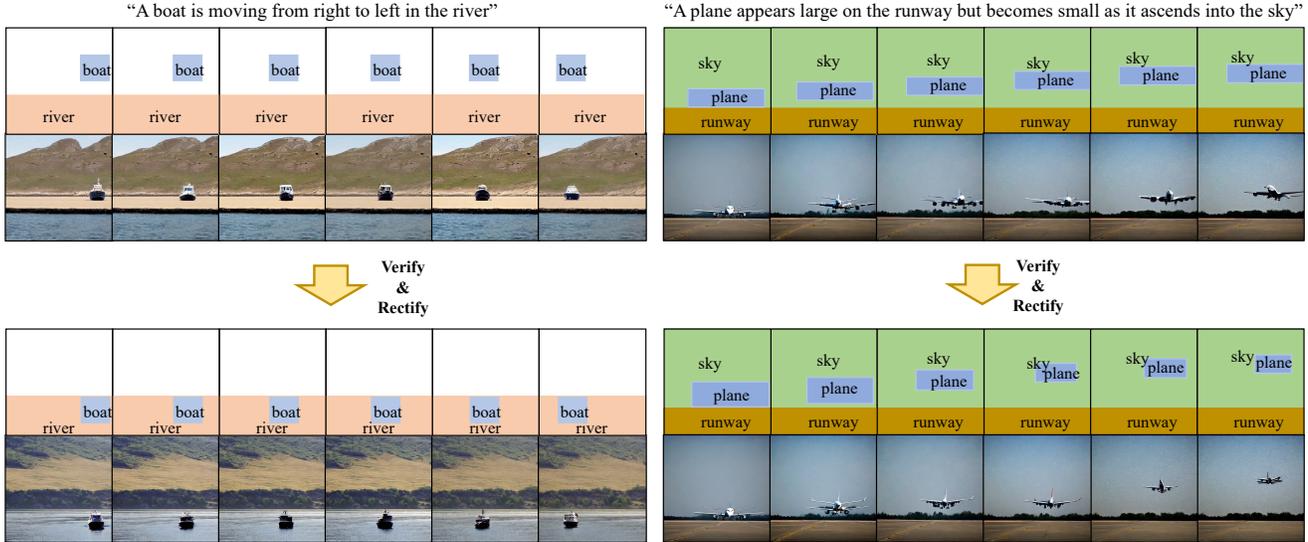


Figure 6. Analysis of the self-refinement process. The self-refinement mechanism verifies and rectifies spatial and temporal alignment between layouts and video prompts.

Effectiveness of Self-Refinement Process We present two examples in Figure 6 to illustrate the effectiveness of our self-refinement process in correcting spatial and temporal errors of initial layouts. In the first example, the video prompt describes “a boat moving in the river.” However, the initial generated layout incorrectly places the boat above the river. The boat is correctly positioned within the river through self-refinement, aligning with the prompt’s spatial arrangement. The second example describes a plane ascending into the sky. Initially, the size of the plane remains constant across all layouts over time, contradicting the expectation that it should appear smaller as it ascends. After refinement, the size of the plane decreases in later frames, accurately reflecting the prompt’s temporal dynamics.

To quantitatively evaluate the effectiveness of the self-refinement process, we propose a benchmark comprising four spatio-temporal layout generation tasks. These tasks include multiple objects, object movements (left, right, up, down), size changes (big to small or small to big), and visibility variations (half or quarter visibility). Each task includes 20 programmatically generated prompts, assessed using a rule-based metric. For instance, we calculate the change in object area across frames to evaluate size changes. The results are displayed in Table 2. We observe LLMs initially struggle to generate precise results that accurately reflect specific temporal changes, including object movement, size variation, and visibility. Moreover, through our self-refinement process, we noted a notable improvement in accuracy, particularly in tasks temporal visibility (from 61% to 78%). The self-refinement mechanism consistently enhances spatial-temporal layout genera-

Table 2. Quantitative analysis of the self-refinement process.

Method	Objects↑	Movement↑	Size↑	Visibility↑
w/o self-refine	90%	83%	80%	61%
w/ self-refine	96%	93%	93%	78%

tion, effectively aligning the generated content with specific temporal requirements. These experiments confirm the effectiveness of our self-refinement process in improving the spatial-temporal coherence of the generated scenes.

5. Conclusion

In this paper, we have investigated leveraging the spatial-temporal planning ability of Large Language Models to guide temporally-coherent text-to-video generation with image diffusion models. We prompt LLMs to generate comprehensive Dynamic Scene Syntax, including scene descriptions, layouts for foreground objects, and background motion patterns. The foreground layouts ensure coherent object motions and object transformations described in the prompt. Furthermore, the introduced iterative self-refinement can enhance the alignment between the generated spatio-temporal layouts and the textual descriptions, specifically addressing inaccuracies such as spatial and temporal errors. The background motion can be controlled by motion-guided noise shifting, leading to smoother video synthesis and a coherent global scene. We have performed extensive qualitative and quantitative experiments along with ablation studies to validate the effectiveness of our FlowZero framework. These experiments validate that FlowZero can generate temporally-coherent videos from complex video prompts.

References

- [1] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*, pages 1708–1718. IEEE, 2021. 5
- [2] Shehzaad Dhuliawala, Mojtaba Komeili, Jing Xu, Roberta Raileanu, Xian Li, Asli Celikyilmaz, and Jason Weston. Chain-of-verification reduces hallucination in large language models. *arXiv preprint arXiv:2309.11495*, 2023. 4
- [3] Weixi Feng, Wanrong Zhu, Tsu-Jui Fu, Varun Jampani, Arjun R. Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and William Yang Wang. Layoutgpt: Compositional visual planning and generation with large language models. *CoRR*, abs/2305.15393, 2023. 2, 3, 4
- [4] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron C. Courville, and Yoshua Bengio. Generative adversarial networks. *CoRR*, abs/1406.2661, 2014. 3
- [5] 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. *CoRR*, abs/2307.04725, 2023. 2, 3, 5, 6, 7
- [6] Bruce C Hansen and Robert F Hess. Structural sparseness and spatial phase alignment in natural scenes. *JOSA A*, 24(7):1873–1885, 2007. 5
- [7] Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. Latent video diffusion models for high-fidelity video generation with arbitrary lengths. *arXiv preprint arXiv:2211.13221*, 2022. 3
- [8] Katrin Heimann, Sebo Uithol, Marta Calbi, Maria Alessandra Umiltà, Michele Guerra, Joerg Fingerhut, and Vittorio Gallese. Embodying the camera: An eeg study on the effect of camera movements on film spectators sensorimotor cortex activation. *PLoS one*, 14(3):e0211026, 2019. 2, 4
- [9] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey A. Gritsenko, Diederik P. Kingma, Ben Poole, Mohammad Norouzi, David J. Fleet, and Tim Salimans. Imagen video: High definition video generation with diffusion models. *CoRR*, abs/2210.02303, 2022. 2
- [10] Susung Hong, Junyoung Seo, Sunghwan Hong, Heeseong Shin, and Seungryong Kim. Large language models are frame-level directors for zero-shot text-to-video generation. *CoRR*, abs/2305.14330, 2023. 2, 3, 4, 5, 6, 7
- [11] Wenyi Hong, Ming Ding, Wendi Zheng, Xinghan Liu, and Jie Tang. Cogvideo: Large-scale pretraining for text-to-video generation via transformers. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023.
- [12] Hanzhuo Huang, Yufan Feng, Cheng Shi, Lan Xu, Jingyi Yu, and Sibeiyang. Free-bloom: Zero-shot text-to-video generator with LLM director and LDM animator. *CoRR*, abs/2309.14494, 2023. 2, 3, 4, 5, 7
- [13] Bernd Jähne. *Digital image processing*. Springer Science & Business Media, 2005. 5
- [14] Levon Khachatryan, Andranik Movsisyan, Vahram Tadevosyan, Roberto Henschel, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. Text2video-zero: Text-to-image diffusion models are zero-shot video generators. *CoRR*, abs/2303.13439, 2023. 2, 3, 5, 6, 7
- [15] Geunwoo Kim, Pierre Baldi, and Stephen McAleer. Language models can solve computer tasks. *CoRR*, abs/2303.17491, 2023. 4
- [16] Yitong Li, Martin Renqiang Min, Dinghan Shen, David E. Carlson, and Lawrence Carin. Video generation from text. *CoRR*, abs/1710.00421, 2017. 3
- [17] Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David E. Carlson, and Jianfeng Gao. Storygan: A sequential conditional GAN for story visualization. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 6329–6338. Computer Vision Foundation / IEEE, 2019. 3
- [18] Yuheng Li, Haotian Liu, Qingyang Wu, Fangzhou Mu, Jianwei Yang, Jianfeng Gao, Chunyuan Li, and Yong Jae Lee. GLIGEN: open-set grounded text-to-image generation. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 22511–22521. IEEE, 2023. 3, 5, 7
- [19] Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models. *CoRR*, abs/2305.13655, 2023. 2, 3
- [20] 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. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2023, Vancouver, BC, Canada, June 17-24, 2023*, pages 10209–10218. IEEE, 2023. 2, 5, 6, 7
- [21] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. 4, 5
- [22] A Oppenheim, Jae Lim, Gary Kopec, and SC Pohlig. Phase in speech and pictures. In *ICASSP'79. IEEE International Conference on Acoustics, Speech, and Signal Processing*, pages 632–637. IEEE, 1979. 5
- [23] Alan V Oppenheim and Jae S Lim. The importance of phase in signals. *Proceedings of the IEEE*, 69(5):529–541, 1981.
- [24] Leon N Piotrowski and Fergus W Campbell. A demonstration of the visual importance and flexibility of spatial-frequency amplitude and phase. *Perception*, 11(3):337–346, 1982. 5
- [25] 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. 6
- [26] Noah Shinn, Federico Cassano, Ashwin Gopinath, Karthik R Narasimhan, and Shunyu Yao. Reflexion: Language agents

- with verbal reinforcement learning. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 4
- [27] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, Devi Parikh, Sonal Gupta, and Yaniv Taigman. Make-a-video: Text-to-video generation without text-video data. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net, 2023. 2, 3
- [28] Yixuan Weng, Minjun Zhu, Fei Xia, Bin Li, Shizhu He, Kang Liu, and Jun Zhao. Large language models are better reasoners with self-verification, 2023. 4
- [29] Jinheng Xie, Kai Ye, Yudong Li, Yuexiang Li, Kevin Qinghong Lin, Yefeng Zheng, Linlin Shen, and Mike Zheng Shou. Visorgpt: Learning visual prior via generative pre-training. *CoRR*, abs/2305.13777, 2023. 3, 4
- [30] Mehmet Burak Yilmaz, Elen Lotman, Andres Karjus, and Pia Tikka. An embodiment of the cinematographer: emotional and perceptual responses to different camera movement techniques. *Frontiers in Neuroscience*, 17, 2023. 2, 4
- [31] Lvmin Zhang and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. *CoRR*, abs/2302.05543, 2023. 3
- [32] Tianjun Zhang, Yi Zhang, Vibhav Vineet, Neel Joshi, and Xin Wang. Controllable text-to-image generation with GPT-4. *CoRR*, abs/2305.18583, 2023. 2, 3, 4
- [33] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *CoRR*, abs/2211.11018, 2022. 2, 3