# MAViS: A Multi-Agent Framework for Long-Sequence Video Storytelling

**Qian Wang**
Virginia Tech
USA
yqwq1996@vt.edu

**Ziqi Huang**
Nanyang Technological University
Singapore
ZIQI002@e.ntu.edu.sg

**Ruoxi Jia**
Virginia Tech
USA
ruoxijia@vt.edu

**Paul Debevec**
Eyeline Studios
USA
debevec@scanlinevfx.com

**Ning Yu**
Eyeline Studios
USA
ning.yu@scanlinevfx.com

Figure 1: With just a brief prompt, MAViS enables users to rapidly explore diverse visual storytelling and creative directions for sequential video generation by efficiently producing high-quality, complete long-sequence videos.

## Abstract

Despite recent advances, long-sequence video generation frameworks still suffer from significant limitations: poor assistive capability, suboptimal visual quality, and limited expressiveness. To mitigate these limitations, we propose **MAViS**, an end-to-end multi-agent collaborative framework for long-sequence video storytelling. MAViS orchestrates specialized agents across multiple stages, including script writing, shot designing, character modeling, keyframe generation, video animation, and audio generation. In each stage, agents operate under the **3E Principle**—Explore, Examine, and Enhance—to ensure the completeness of intermediate outputs. Considering the capability limitations of current generative models, we propose the Script Writing Guidelines to optimize compatibility between scripts and generative tools. Experimental results demonstrate that MAViS achieves state-of-the-art performance in assistive capability, visual quality, and video expressiveness. Its modular framework further enables scalability with diverse generative models and tools. With just a brief prompt, MAViS enables users to rapidly explore diverse visual storytelling and creative directions for sequential video generation by efficiently producing high-quality, complete long-sequence videos. To the best of our knowledge, MAViS is the only framework that provides multimodal design output—videos with narratives and background music.

## 1 Introduction

In recent years, video generative models (OpenAI, 2024; MiniMax, 2024; GenmoTeam, 2024) have made substantial strides in text-to-video (T2V) and image-to-video (I2V) generation. While these methods demonstrate promising capabilities for shot-level creation, the generation of long videos remains a formidable challenge. Besides, the current scale of computing and data is insufficient to train a single end-to-end model capable of producing high-quality, minute-long videos.

To overcome this limitation, prior works have explored extending existing video clips (KlingAI, 2024; Henschel et al., 2024; Ge et al., 2022), or composing sequences of keyframes (Liu et al., 2024; Yang et al., 2024b; Rahman et al., 2023; Wang et al., 2024a), or integrating storytelling and video generation (Yin et al., 2023b; Zheng et al., 2025). However, these approaches suffer from repetitive motion, lack of narrative cohesion, or result in static image series with limited visual expressions. To better align narrative planning with video animation, MovieDreamer (Zhao et al., 2024) adopts an auto-regressive approach to generate key frames from input scripts, while Dreamfactory (Xie

| Methods | Agentic | Script Writing | Shot Designing | Dynamic | ID Consistency | Audio | Automated |
|---|---|---|---|---|---|---|---|
| Mora (Yuan et al., 2024) | ✓ | ✗ | ✗ | ✓ | ✗ | ✗ | ✓ |
| AesopAgent (Wang et al., 2024a) | ✓ | ✓ | ✗ | ✗ | ✓ | ✗ | ✓ |
| MovieDreamer (Zhao et al., 2024) | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ |
| DreamFactory (Xie et al., 2024) | ✓ | ✗ | ✓ | ✓ | ✓ | ✗ | ✗ |
| VGoT (Zheng et al., 2025) | ✗ | ✓ | ✗ | ✓ | ✓ | ✗ | ✓ |
| LCT (Guo et al., 2025) | ✗ | ✗ | ✗ | ✓ | ✓ | ✗ | ✗ |
| MovieAgent (Wu et al., 2025b) | ✓ | ✗ | ✓ | ✓ | ✓ | ✓ | ✗ |
| **MAViS (ours)** | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

Table 1: **Properties of MAViS vs Other Methods**: We summarize seven metrics for long-sequence video storytelling works."Agentic": whether the system employs a multi-agent architecture; "Script Writing": ability generate scripts from user prompt; "Shot Designing.": support for structured shot designing and controllable narrative flow; "Dynamic": whether the generated video exhibits dynamic features; "ID Consistency": consistency in character identity across multiple shots; "Audio": whether the output includes synchronized audio; "Automated": whether the pipeline can operate without the presence of user input in script writing, script planning, or model fine-tuning.

et al., 2024) and MovieAgent (Wu et al., 2025b) propose a hierarchical script–scene–shot pipeline. However, as shown in Table 1, these methods require substantial manual effort from users for script writing and training LoRA models to maintain ID consistency, which significantly limits their assistive capability and hinders large-scale deployment.

In summary, current long-sequence video storytelling methods face three major limitations: 1) poor assistive capability, not supporting script writing, shot designing, and LoRA training to maintain ID consistency; 2) suboptimal visual quality, manifested in visual distortions, unnatural actions, and inaccurate object proportions; and 3) limited expressiveness, often resulting in repetitive actions and poorly structured shot compositions.

To address this, we introduce **MAViS**, a multi-agent framework for end-to-end long-sequence video generation. MAViS encompasses various stages, including script writing, shot designing, character modeling, keyframe generation, video animation, and audio generation. Given a brief user prompt, the script writing stage first generates a structured script with character definitions and shot outlines. The shot designing stage defines a series of shot-level elements to guide subsequent content generation, while character modeling ensures ID consistency via LoRA-trained appearances. During the keyframe generation stage, T2I models are employed to create the initial static frames for each shot, i.e., the keyframes, which are then expanded into full video sequences in the video animation stage through I2V models, maintaining visual coherence and narrative fluency. Finally, voice-overs and background music (BGM) are added to each shot during the voice generation stage, resulting in a cohesive, long-sequence video.

Crucially, all stages operate under the **3E Principle**: - *Explore*: Initial generation of candidate outputs; - *Examine*: Reviewing the quality and completeness of the output; - *Enhance*: Iterative refinement and optimization. After multiple iterations, each stage passes the most complete version of its output to the next, enabling gradual optimization throughout the end-to-end pipeline.

To mitigate the limitations of current generative models in maintaining background consistency, executing complex actions, and rendering fine-grained details, we propose the Script Writing Guidelines aimed at enhancing the compatibility between scripts and generative models, ultimately improving the expressiveness of the video storytelling. As illustrated in Figure 1, with just a brief prompt, MAViS enables users to rapidly explore diverse visual storytelling and creative directions for sequential video generation by efficiently producing high-quality, complete long-sequence videos. Furthermore, its modular architecture ensures scalability with various generative tools and models.

Our contributions are summarized in three thrusts:

• We propose **MAViS**, an end-to-end multi-agent framework that encompasses various stages including script writing, shot designing, character modeling, keyframe generation, video animation, and audio generation — forming a complete and scalable workflow for long video storytelling.

• We introduce the **3E Principle**—*Explore*, *Examine*, *Enhance*—to continuously optimize the output of each stage. This iterative design ensures the completeness of the output at each stage.

• We propose the Script Writing Guidelines for long-sequence video storytelling to enhance the script-tool compatibility, ultimately improving the

expressiveness of the video storytelling.

- Experiments demonstrate that MAViS achieves state-of-the-art performance and receive the most user preference votes in terms of expressiveness and visual quality. With just a brief user prompt, MAViS efficiently generates high-quality and complete long-sequence video storytelling. To the best of our knowledge, MAViS is the only framework that provides multimodal design output.

## 2 Related Works

### 2.1 Video Generation

Long video generation (He et al., 2022; Ge et al., 2022; Wang et al., 2024c) aims to produce longer videos than previous text-to-video (T2V) (ModelScopeTeam, 2023; Wang et al., 2024b) and image-to-video (I2V) (Girdhar et al., 2023; Xing et al., 2025; Zhang et al., 2023) models. Auto-regressive approaches (Weng et al., 2024; Henschel et al., 2024) and video storytelling methods (Yin et al., 2023b; Zheng et al., 2025) suffer from repetitive motions with limited narrative development and poor visual quality. Meanwhile, commercial models (KlingAI, 2024; Runway, 2024a,b) support temporal extension. However, the duration of the generated videos remains constrained. In this work, MAViS uses cutting-edge video generation models to synthesize high-quality video clips and assembles them into long-sequence video storytelling.

### 2.2 AI Agents

AI agents (Yang et al., 2023a; Yin et al., 2023a; Li et al., 2023; Hong et al., 2023; Park et al., 2023), empowered by large language models (Brown et al., 2020; Touvron et al., 2023a,b), have shown broad applicability in reasoning, planning, web interaction, and visual generation. In video generation, systems (Yang et al., 2024a; Yuan et al., 2024; Wang et al., 2024a) explore world modeling, multi-agent coordination, and retrieval-augmented synthesis, yet often yield static content with limited transitions. Hierarchical pipelines (Xie et al., 2024; Wu et al., 2025b) introduce script-to-shot workflows but still rely on manual scripting and LoRA tuning. In contrast, our MAViS leverages multi-agent collaboration to enable long-sequence video storytelling.

## 3 MAViS

We propose **MAViS**, a multi-agent end-to-end framework for long-sequence video storytelling that consists of various stages, including script writing, shot designing, character modeling, keyframe generation, video animation, and audio generation.

Given a user prompt $P$, as illustrated in Figure 2 MAViS aims to generate a long-sequence video storytelling that comprises $N$ shots:

$$\mathcal{F}: P \to \hat{\mathcal{V}}, \quad \hat{\mathcal{V}} = \{\{V_j,\, A_j\} \mid j = 1, 2, \ldots, N\}, \quad (1)$$

where $V_j$ is the video clip and $A_j$ is the audio track. To ensure the completeness and reliability of the generation pipeline, agents in each stage follow a unified **3E Principle**: *Explore*, *Examine*, and *Enhance*. Section 3.1 details the 3E principle, and Section 3.2 outlines the Script Writing Guidelines, followed by Section 3.3 to 3.6 introducing the agent design and working procedures of each stage.

### 3.1 3E Principle

Existing frameworks for long-sequence video storytelling adopt one-shot generation strategies. However, such approaches often prove inadequate in practice—scripts may not align with user intent, images or videos can exhibit distortions or unnatural actions, and audio tracks may be misaligned in duration. These issues highlight the necessity of adopting a multi-round generation–modification loop, rather than relying on one-shot generation. To address this, we introduce the **3E Principle**—Explore, Examine, Enhance—to guide intra-stage agents' collaboration and progressively improve generation quality across iterations:

*Explore*: This stage represents the first attempt to generate content based on the current request and specification:

$$r_1 = G(p_1; g), \quad (2)$$

where $p_1$ denotes the initial generation request, $g$ is the generation guidelines, such as guides of script-tool compatibility and output completeness. $G$ is the generator, and $r_1$ is the resulting output.

*Examine*: The generated output is then examined to determine whether it satisfies the specification:

$$f_1 = E(r_1, g), \quad (3)$$

where $E$ is the reviewer, and $f_1$ is the examination feedback of result $r_1$ based on the guidelines $g$.

*Enhance*: Based on the feedback and guidelines, the request is refined to improve quality in the next iteration:

$$p_{i+1} = R(p_i; f_i, g), \quad (4)$$

where $R$ denotes the refiner, and $p_{i+1}$ is the updated request for iteration $i + 1$. This new request serves
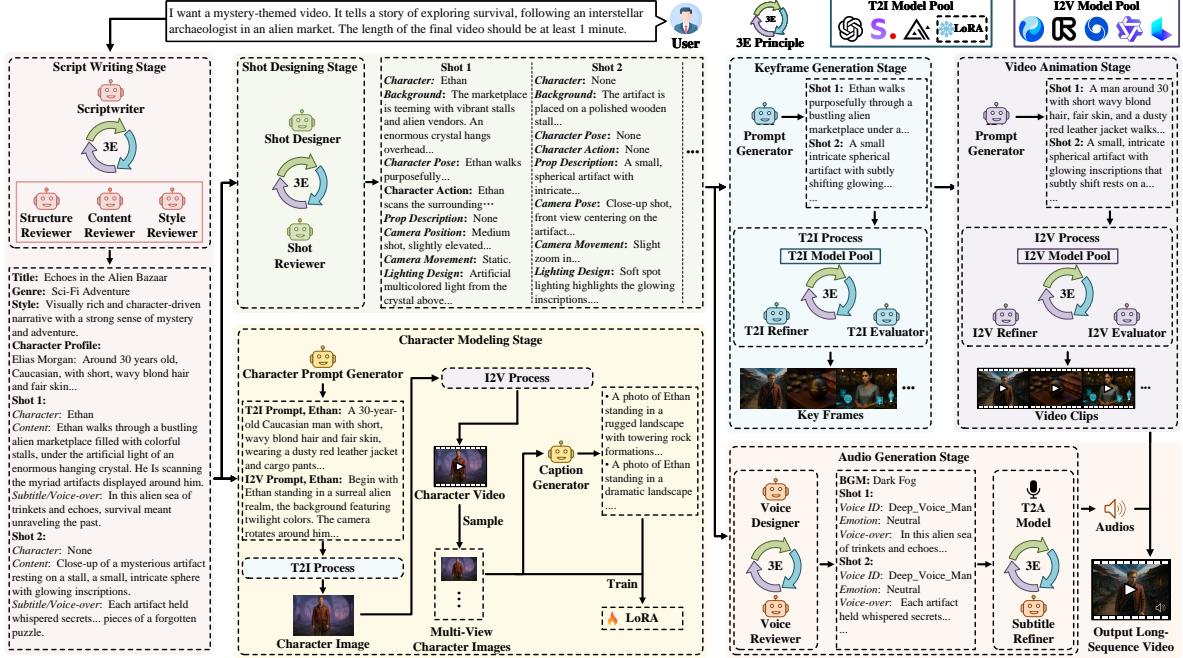
Figure 2: Illustration of MAViS framework.

as the input to the next round of Explore-Examine-Enhance, and the process continues iteratively until the output fully satisfies the specification.

After multiple Explore-Examine-Enhance iterations, each stage passes the most complete version of its output to the next, enabling gradual optimization throughout the pipeline. Given the limitations of individual agents, it may need to deploy multiple reviewers in parallel to ensure comprehensive assessment. Additionally, generating multiple candidates per iteration and enhancing the highest-scoring one can effectively reduce the total number of required iterations. Moreover, depending on the design, the generator and refiner may be implemented as a unified component.

## 3.2 Script Writing Guidelines

Given the limitations of current generative models in maintaining background consistency, executing complex actions, and rendering fine-grained details, certain scripts remain infeasible for video generation. To address this gap and improve the compatibility between scripts and generative models, we propose the Script Writing Guidelines for long-sequence storytelling. The Script Writing Guidelines comprises three components: the Structure Guide, which focuses on how scenes and shots are segmented and arranged; the Content Guide, which defines how much content a single shot should include; and the Style Guide, which ensures user-centric alignment, coherent pacing, logical flow, and structural completeness. Together, these guidelines constrain agent-based script writing and improve the expressiveness of the ultimate video storytelling.

**Structure Guide** Current video generation models have limited contextual modeling capacity; as a result, adjacent shots set within the same spatial environment often suffer from inconsistency in background elements. Therefore, scriptwriters should avoid placing successive shots in the same background and should minimize transitions between tightly connected spatial locations (e.g., from "outside the door" to "inside the room," or from one room to another). Moreover, to reduce the jarring effect of spatial jumps, it is recommended not to place the same character in consecutive shots unless narrative necessity dictates it.

To mitigate these issues while preserving narrative flow, transitional shots can be inserted. These may include close-ups of different characters, partial views of a character (e.g., hands, back), detailed shots of the environment, or important props. Such insertions help avoid background duplication while reinforcing the spatial or character continuity. Additionally, enhanced descriptions of both background and character appearance are encouraged to improve model alignment with visual and narrative consistency.

**Content Guide** Current video generation models struggle to synthesize multiple sequential actions within a single shot (e.g., running followed by walking and then dancing). Therefore, it is advisable to limit each shot to a single, simple action. Furthermore, generative models often fail to produce fine-grained visual elements accurately, such as textual content on smartphone screens or small patterns on objects. Scripts should avoid specifying such details.

**Style Guide** Within the Structure Guide and the Content Guide, a script may risk becoming fragmented or visually monotonous if not properly orchestrated. The style guide thus ensures harmonious pacing, visual diversity, and narrative logic. Moreover, the script's tone and theme must align with the user's request, and different genres require distinct narrative pacing—emotional and family-oriented scripts benefit from a slow, layered progression, while action-oriented scripts demand rapid and dynamic developments. Additionally, structural completeness is essential: every script must contain a title, character definitions, and a shot-by-shot script. Each shot should specify the character (or "None" if absent), shot content, and accompanying subtitle or voice-over. Subtitles must align semantically with the shot content.

### 3.3 Script Writing Stage

Upon receiving a brief user prompt, four agents—the Scriptwriter, the Structure Reviewer, the Content Reviewer, and the Style Reviewer—work collaboratively in the script writing stage to produce a structured script.

The Scriptwriter $G_{\mathrm{scr}}$ first generates a draft script $S_1$ based on the user prompt $P$ and the Script Writing Guidelines $g_{\mathrm{scr}}$:

$$S_1 = G_{\mathrm{scr}}(P, g_{\mathrm{scr}}). \tag{5}$$

Following the 3E principle, at each iteration, the Structure Reviewer $E_{\mathrm{str}}$, Content Reviewer $E_{\mathrm{con}}$, and Style Reviewer $E_{\mathrm{sty}}$ examine the script $S_i$ against the Structure Guide $g_{\mathrm{str}}$, Content Guide $g_{\mathrm{con}}$ and Style Guide $g_{\mathrm{sty}}$ respectively, and provide feedback:

$$f_{\mathrm{scr}}^i = [E_{\mathrm{str}}(S_i, g_{\mathrm{str}}),\ E_{\mathrm{con}}(S_i, g_{\mathrm{con}}),\ E_{\mathrm{sty}}(S_i, g_{\mathrm{sty}})]\,, \tag{6}$$

where $g_{\mathrm{scr}} = [g_{\mathrm{str}},\ g_{\mathrm{con}},\ g_{\mathrm{sty}}]$. The Scriptwriter integrates these feedback to improve the script iteratively:

$$S_{i+1} = G_{\mathrm{scr}}(f_{\mathrm{scr}}^i;\ S_i, g_{\mathrm{scr}}) \tag{7}$$

The Explore-Examine-Enhance loop continues until all reviewers reach a consensus that the script fully satisfies the guidelines after $\hat{n}$ iterations, and the Scriptwriter outputs the finalized version, represented as $S = C_{\mathrm{scr}}(S_n)$.

### 3.4 Shot Designing Stage

In the shot-designing stage, the Shot Designer—responsible for generating detailed shot designs—and the Shot Reviewer—tasked with assessing the completeness and adherence of shot designs to predefined guidelines—play pivotal roles in transforming the high-level script into detailed shot elements suitable for downstream image and video generation. In detail, each shot is enriched with seven essential elements to enable precise visual control:

• *Background*: Specifies the spatial and environmental context of the shot, including location, time of day, season, weather conditions, and foreground-background composition.

• *Character Pose*: Describes the initial static state of characters, including posture, expression, and position when the shot begins.

• *Character Action*: Captures dynamic interactions or movements of characters throughout the shot.

• *Prop Description*: Defines non-character objects, including their form, texture, behavior, and placement.

• *Camera Position*: Specifies the camera's initial angle, orientation, and distance from the scene.

• *Camera Movement*: Describes the movement mode and direction of the camera during the shot.

• *Lighting Design*: Defines lighting characteristics such as type, tone, brightness, and color temperature to convey atmosphere.

Similarly, following the 3E Principle, the Shot Designer generates an initial shot design from the script:

$$\hat{S}_1 = G_{\mathrm{shot}}(S), \tag{8}$$

where the shot designing guide provides essential constraints to ensure completeness and stylistic alignment with the original script.

In each iteration $i$, the Shot Reviewer examines the completeness and appropriateness of the current design $\hat{S}_i$ and provides feedback:

$$f_{\mathrm{shot}}^i = E_{\mathrm{shot}}(\hat{S}_i). \tag{9}$$

The Shot Designer then uses this feedback to enhance the design, until the Shot Reviewer deems the result satisfactory:

$$\hat{S}_{i+1} = G_{\mathrm{shot}}(f_{\mathrm{shot}}^i;\ \hat{S}_i), \quad \hat{S} = \hat{S}_n. \tag{10}$$

## 3.5 Character Modeling Stage

To ensure the ID consistency of characters across shots, in the character modeling stage, agents construct visual representations for each character and train a LoRA model. This process incorporates the Character Prompt Generator and Caption Generator, supported by reviewer agents for both image generation (T2I Evaluator and T2I Refiner) and video animation (I2V Evaluator and I2V Refiner).

The process begins with the Character Prompt Generator, which generates both text-to-image (T2I) and image-to-video (I2V) prompts for each character. T2I prompts aim to generate a frontal and naturally posed character image, while I2V prompts are used—together with the T2I images—to synthesize character-centric video clips that capture the subject from multiple angles. After both T2I and I2V processes, multi-view character images are sampled from the generated character video. The Caption Generator then produces structured captions for each sampled frame. These images and captions are subsequently used to train LoRA models (if needed).

Both the T2I and I2V processes follow the 3E Principle, but differ from script writing and shot designing in that each iteration produces multiple samples from the generative model pools. Specifically, in the $i$-th iteration, a set of candidate images is generated by the pool of T2I models $\mathcal{M}_I$ using the current T2I prompt $p_I^i$:

$$\mathcal{I}_i = \{G_I^k(p_I^i),\ G_I^k \in \mathcal{M}_I,\ k = 1, 2, \ldots, |\mathcal{M}_I|\}. \quad (11)$$

T2I Evaluator $E_I$ scores candidates according to the image evaluation guide $g_I$, selects the best one $I_i$, and provides evaluation result $f_I^i$:

$$[I_i, f_I^i] = E_I(\mathcal{I}_i, g_I) \quad (12)$$

The T2I Refiner $R_I$ then updates the prompt accordingly:

$$p_I^{i+1} = R_I(f_I^i, p_I^i) \quad (13)$$

After $n$ iterations, the final image $I = I_n$ initiates the I2V process.

The I2V process mirrors the T2I process. First, a set of candidate videos is generated by the I2V model pool $\mathcal{M}_V$:

$$\mathcal{V}_i = \{G_V^k(I, p_V^i),\ G_V^k \in \mathcal{M}_V,\ k = 1, 2, \ldots, |\mathcal{M}_V|\}. \quad (14)$$

The I2V Evaluator $E_V$ evaluates these candidates according to the video evaluation guide $g_I$, selects the best one $V_i$, and provides evaluation result $f_V^i$.

The I2V Refiner $R_V$ updates the I2V prompt and produce the final video $V_n$ after $n$ iterations:

$$[V_i, f_V^i] = E_V(\mathcal{V}_i, g_V), \quad p_V^{i+1} = R_V(f_V^i, p_I^i), \quad V = V_n. \quad (15)$$

The difference is that the I2V process incorporates additional evaluation complexity. The image evaluation guide assesses candidates based on three primary axes: 1) *Visual Quality*: clarity, detail resolution, and aesthetic harmony; 2) *Naturalness*: identifying distortions, anatomical inaccuracies, or proportion inconsistencies; and 3) *Prompt Consistency*: if the image content aligns with the input prompt. In contrast, the video evaluation guide extends these criteria by introducing: 1) *Subject Consistency*: temporal stability in facial features, attire, hairstyle, and skin tone across frames; and 2) *Dynamics*: the richness, fluidity, and pacing of character movement throughout the clip. Together, the image and video evaluation guides work in concert to enforce the completeness of the final generation.

## 3.6 Generation Stages

The generation stages include: 1) the keyframe generation stage, which comprises the T2I Prompt Generator, T2I Evaluator, and T2I Refiner, responsible for generating the initial reference image for each shot; 2) the video animation stage, which includes the I2V Prompt Generator, I2V Evaluator, and I2V Refiner, and completes the video shot based on the keyframes; and 3) the audio generation stage, which consists of the Voice Designer, Voice Reviewer, and Subtitle Refiner, synthesizes voice-overs and background music aligned with the emotional rhythm of each shot.

First, in the keyframe generation stage, the T2I Prompt Generator uses shot-level elements—*Character*, *Background*, *Character Pose*, *Prop Description*, *Camera Position*, and *Lighting Design*—to construct T2I prompts, which are refined via the same T2I process in Section 3.5, utilizing the T2I model pool and the trained LoRA model, resulting in a set of keyframes $\{I_j \mid j = 1, 2, \ldots, N\}$.

Second, in the video animation stage, the I2V Prompt Generator formulates I2V prompts using *Character*, *Background*, *Character Action*, *Prop Description*, *Camera Movement*, and *Lighting Design*. Again, following the same I2V process introduced in Section 3.5, these prompts are used to generate video clips $\{V_j \mid j = 1, 2, \ldots, N\}$.

After video animation, voice and music planning are initiated in the audio generation stage based

| Method | Keyframe Generation | | Video Animation | | | | | |
|---|---|---|---|---|---|---|---|---|
| | CLIP ↑ | Inception ↑ | T. Flick. ↑ | M. Smooth. ↑ | Sub. Cons. ↑ | Bg. Cons. ↑ | Aesthetic ↑ | I. Quality ↑ |
| VGoT (Zheng et al., 2025) | 22.37 | 8.53 | 99.00 | 99.31 | **98.94** | **98.74** | **80.11** | 64.59 |
| Mora (Yuan et al., 2024) | 33.98 | 12.68 | 97.32 | 99.23 | 95.23 | 94.88 | 64.36 | 71.48 |
| MovieAgent (Wu et al., 2025b) | 31.71 | 10.72 | 97.07 | 99.20 | 93.29 | 94.51 | 63.39 | 71.26 |
| **MAViS (ours)** | **34.22** | **12.81** | **99.09** | **99.53** | 95.72 | 96.12 | 63.17 | **72.91** |

Table 2: **Evaluation of Automatic Metrics for Keyframe Generation and Video Animation.** "T. Flick.", "M. Smooth.", "Sub. Cons.", "Bg. Cons.", "Aesthetic", and "I. Quality" refer to "Temporal Flickering", "Motion Smoothness", "Subject Consistency", "Background Consistency", 'Aesthetic Quality', and "Imaging Quality" from metrics in VBench (Huang et al., 2024), respectively. ↑ indicates a higher value is more desired. **Bold** indicates the best results.

| Method | Narrative ↑ | Visual ↑ | User Align. ↑ | Sub. Cons. ↑ | Sub. Natural. ↑ | Bg. Cons. ↑ | Bg. Real. ↑ |
|---|---|---|---|---|---|---|---|
| VGoT (Zheng et al., 2025) | 3.39 | 3.39 | 6.79 | 5.56 | 3.07 | 4.14 | 3.24 |
| Mora (Yuan et al., 2024) | 15.18 | 16.96 | 13.57 | 15.77 | 13.00 | 13.69 | 15.29 |
| MovieAgent (Wu et al., 2025b) | 9.46 | 11.79 | 11.96 | 12.37 | 12.27 | 12.07 | 10.97 |
| **MAViS (ours)** | **71.96** | **67.86** | **67.68** | **66.31** | **71.66** | **70.09** | **70.50** |

Table 3: **User Study on the Performance of Long-Sequence Video Storytelling** (Voting Results). "Narrative", "Visual", "User Align.", "Sub. Cons.", "Sub. Natural.", "Bg. Cons.", and "Bg. Real." refer to "Narrative Expressiveness", "Visual Quality", "User Prompt Alignment", "Subject Consistency", "Character Naturalness", "Background Consistency", and "Background Realism", respectively. ↑ indicates a higher value is more desired. **Bold** indicates the best results.

on the 3E principle. The Voice Designer $G_A$ selects background music for the entire script and configures voice IDs and emotional tones for each individual shot. The voice design is then reviewed by the Voice Reviewer $E_A$, who checks for compliance (e.g., music availability and voice–emotion appropriateness) and provides iterative feedback $f_A^i$ to assist in refinement:

$$S_A^1 = G_A(S), \quad f_A^i = E_A(S_A^i), \quad S_A^{i+1} = G_A(S, f_A^i),$$
$$S_A = S_A^n. \tag{16}$$

Next, the T2A model $G_A$ generates speech audio for each shot. The Subtitle Refiner $R_A$ shortens the voice-over $p_A$ to fit temporal constraints, ensuring compatibility with video animation:

$$A_j^i = G_A(p_{A,j}^i), \ p_{A,j}^{i+1} = R_A(p_{A,j}^i; A_j^i, V_j), \ A_j = A_j^n. \tag{17}$$

Finally, audio tracks are synchronized with video clips to produce the final long-sequence video telling with synchronized visuals and sound: $\hat{\mathcal{V}} = \{\{V_j, \ A_j\} \mid j = 1, 2, \ldots, N\}$.

# 4 Experiment

To evaluate the effectiveness of MAViS, we benchmark it against three long video generation baselines: VGoT, Mora, and MovieAgent, and conduct detailed ablation studies. We constructed a test set of 20 user prompts, each targeting a one-minute video for evaluation. We employ CLIP, Inception Score, and VBench (Huang et al., 2024) for quantitative evaluation, and conduct a user study for subjective assessment. Further experimental settings, metrics, and baseline details are provided in the Appendix.

## 4.1 Performance Comparisons

We evaluate each method via automatic metrics and user study. It's worth noting that MAViS is the only framework that provides multimodal output – videos with narratives and BGM.

### 4.1.1 Automatic Metrics

Table 2 presents the shot-level evaluation results for both keyframe generation and video animation. MAViS achieves the highest scores in CLIP (*34.22*) and Inception (*12.81*) scores, attributed to the enhanced T2I quality enabled by the 3E Principle in the T2I process. In video animation, MAViS slightly lags behind VGoT in subject consistency and background consistency but outperforms all baselines in temporal flickering, motion smoothness, and image quality. This is partly because VGoT produces videos with limited dynamics and highly saturated visuals, which improves intra-shot consistency and aesthetic quality but at the cost of realism. Overall, MAViS achieves the best comprehensive performance.

### 4.1.2 User Study

Table 3 summarizes the voting results from 60 evaluators across multiple aspects of comparison among the baselines. MAViS consistently outperforms all other methods in every evaluation dimen-

| Ablation | Narrative ↑ | Visual ↑ | User Align. ↑ |
|---|---|---|---|
| w/o $E_{str}$ | 13.67 | 25.67 | 26.67 |
| w/o $E_{con}$ | 32.33 | 18.67 | 26.67 |
| w/o $E_{sty}$ | 20.00 | 26.67 | 20.00 |
| **MAViS (ours)** | 34.00 | 26.00 | 26.67 |

Table 4: **Ablation study on the Structure Reviewer $E_{str}$, the Content Reviewer $E_{con}$, and the Style Reviewer $E_{sty}$.** ↑ indicates a higher value is more desired. Gray-shaded cells indicate values that are significantly lower than the others.

| | Ablation | CLIP ↑ | Inception ↑ | Natural. ↑ | Sub. Cons. ↑ | Dyn. Deg. ↑ |
|---|---|---|---|---|---|---|
| T2I | w/o 3E | 31.25 | 10.54 | 72.00 | - | - |
| | w 3E | 34.22 | 12.81 | 98.00 | - | - |
| I2V | w/o 3E | 32.05 | 11.22 | 60.00 | 3.0 | 44.58 |
| | w 3E | 34.53 | 13.22 | 92.50 | 95.72 | 35.62 |

Table 5: **Ablation study of the 3E Principle in the T2I and I2V Processes.** "Natural." refers to the Naturalness score, ↑ indicates a higher value is more desired. Gray-shaded cells indicate values that are significantly lower.

sion, with an average voting share of *69.44%*. This highlights both the effectiveness of our architecture and the critical role of the 3E Principle.

## 4.2 Ablation Study

We conduct an ablation study to demonstrate the impact of the 3E Principle across different modules of the framework.

### 4.2.1 Script Writing Guidelines

In the script writing stage, based on the Script Writing Guidelines, the Structure Reviewer focuses on how scenes and shots are segmented and arranged, the Content Reviewer ensures that each shot can fully convey its intended content, and the Style Reviewer oversees global pacing, logical flow, structural coherence, and alignment with user requirements. Table 4 presents the effects of removing each reviewer from the pipeline. As shown, excluding the Structure Reviewer significantly impairs narrative expressiveness, removing the Content Reviewer degrades visual quality, and omitting the Style Reviewer leads to a loss in both narrative alignment and consistency with user intent. Optimal script generation is achieved only when all three reviewers collaborate.

### 4.2.2 3E Principle in T2I and I2V Processes

Since the evaluators and refiners in both the T2I and I2V processes cannot operate independently, we conduct an ablation study to evaluate the overall impact of the 3E Principle on these modules.

| Ablation | C. Rate ↑ | Ablation | C. Rate ↑ | Ablation | C. Rate ↑ |
|---|---|---|---|---|---|
| w/o $E_{shot}$ | 84.5 | w/o $E_A$ | 92.0 | w/o $R_A$ | 56.5 |
| w $E_{shot}$ | 100.0 | w/o $E_A$ | 100.0 | w/o $R_A$ | 100.0 |

Table 6: **Ablation study on the Shot Reviewer $E_{shot}$, the Voice Reviewer $E_A$, and the Subtitle Refiner $R_A$.** "C. Rate" refers to "Compliance Rate". ↑ indicates a higher value is more desired. Gray-shaded cells indicate values that are significantly lower.

The CLIP score, Inception score, Naturalness score (defined as the proportion of generated images or videos that appear physically plausible and naturally rendered), along with Subject Consistency and Dynamic Degree from the VBench metrics, are used to assess key aspects targeted by the T2I and I2V Evaluators. As shown in Table 5, removing the 3E Principle leads to a noticeable drop in performance across both T2I and I2V outputs, particularly in terms of naturalness (72.0 → 98.0 in T2I and 60.0 → 92.5 in I2V). This highlights the critical role of the 3E Principle in improving generation quality. An exception is observed in the I2V process: the Dynamic Degree is higher without the 3E Principle (44.58 > 35.63). This may be because high-dynamic videos are more prone to physical inconsistencies and instability, leading the agents to favor lower-dynamic, more stable results.

### 4.2.3 Shot Reviewer, Voice Reviewer, and Subtitle Refiner

As related, the Shot Reviewer assesses the completeness of the shot designs, the Voice Reviewer ensures the music availability and voice-emotion appropriateness of voice design, and the Subtitle Refiner shortens voice-over to match the duration limits of video clips. We use Compliance Rate—defined as the proportion of outputs that fully adhere to their respective guidelines—to evaluate the contributions of these agents. As shown in Table 6, removing these agents leads to a noticeable drop in compliance, particularly in voice-over generation, which in turn negatively impacts downstream content generation. These results underscore the necessity of the 3E Principle in shot designing, voice designing, and the speech audio generation.

## 5 Conclusion

In this work, we presented MAViS, a multi-agent end-to-end framework for long-sequence video storytelling. By orchestrating multiple stages within a unified pipeline, MAViS effectively addresses

the key limitations of prior approaches — poor assistive capability, limited visual expressiveness, and incomplete outputs. Central to our framework is the 3E Principle (Explore, Examine, Enhance), which enables iterative refinement across all stages, maximizing output quality. Additionally, our proposed Script Writing Guidelines improve compatibility between narrative structure and generative tools. Experimental results validate that MAViS significantly advances the assistive capability, visual quality, and expressiveness of long-sequence video storytelling.

## Limitations

In practical AI short film production, there are multiple ways to generate a shot—for example, using direct T2V, combining several I2V generations from different camera angles, or applying image editing tools and image generation models with stronger semantic control and compositional reasoning (e.g., nanobanana (Google, 2025) and HunyuanImage 3.0 (Cao et al., 2025)) to enforce identity and background consistency. However, this paper restricts long-sequence video storytelling to a T2I+I2V paradigm, and the scriptwriting guidelines are explicitly adapted to accommodate this pipeline. While this design simplifies the workflow, it also limits shot diversity, thereby constraining the expressive capacity of the generated storytelling.

Moreover, current storylines lack interactions and dialogue between characters, and the absence of video editing mechanisms in the framework further degrades the cinematic quality of the final output.

To address these limitations, future work will allow agents to autonomously select appropriate generation strategies and tools, and will incorporate character dialogues, scene-specific sound effects, and video editing functionalities, aiming to produce more expressive and coherent video storytelling.

Another limitation of this work lies in the limited reliability of the T2I and I2V Evaluators. Current multimodal foundation models still exhibit limited accuracy in assessing visual content, resulting in suboptimal scoring performance when ranking candidate outputs. To mitigate this in practice, we incorporate additional evaluation operators—such as those from VBench for assessing prompt consistency and image quality—to assist the Evaluators. However, to ensure fairness during experimental evaluation, these operators were not used. In future

work, we plan to integrate a broader set of evaluation operators to enhance the overall assessment process.

## Ethical Considerations

While multi-agent frameworks like MAViS enable efficient and creative video storytelling, AI-generated media must be used responsibly. Because MAViS integrates multiple open-source models and commercial APIs, its outputs may inherit safety risks, biases, or limitations from these components. We encourage transparent disclosure of model provenance, dataset sources, and responsible use practices when employing MAViS in creative or research contexts.

**Potential Risks.** MAViS is designed as a research framework for controllable video generation and multimodal reasoning. However, its capacity to produce coherent, realistic long-sequence videos may be misused to create misleading or harmful synthetic content, such as disinformation, deepfakes, or fabricated narratives. In addition, because it relies on large pretrained models and APIs, MAViS may inadvertently propagate existing societal biases, cultural stereotypes, or skewed representations into generated stories. We emphasize that MAViS is intended solely for academic research and creative exploration under ethical guidelines, not for deployment in sensitive, deceptive, or adversarial applications.

## References

Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.

Wang Ang and Ai et al. Baole. 2025. Wan: Open and advanced large-scale video generative models. *arXiv preprint arXiv:2503.20314*.

Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, and 1 others. 2023. Gemini: A family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 1.

BlackForestLabs. 2023. Flux. https://github.com/black-forest-labs/flux.

Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda

Askell, and 1 others. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.

Siyu Cao, Hangting Chen, Peng Chen, Yiji Cheng, Yutao Cui, Xinchi Deng, Ying Dong, Kipper Gong, Tianpeng Gu, Xiusen Gu, and 1 others. 2025. Hunyuanimage 3.0 technical report. *arXiv preprint arXiv:2509.23951*.

Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H Bermano, Gal Chechik, and Daniel Cohen-Or. 2022. An image is worth one word: Personalizing text-to-image generation using textual inversion. *arXiv preprint arXiv:2208.01618*.

Songwei Ge, Thomas Hayes, Harry Yang, Xi Yin, Guan Pang, David Jacobs, Jia-Bin Huang, and Devi Parikh. 2022. Long video generation with time-agnostic vq-gan and time-sensitive transformer. In *European Conference on Computer Vision*, pages 102–118. Springer.

GenmoTeam. 2024. Mochi 1. https://github.com/genmoai/models.

Rohit Girdhar, Mannat Singh, Andrew Brown, Quentin Duval, Samaneh Azadi, Sai Saketh Rambhatla, Akbar Shah, Xi Yin, Devi Parikh, and Ishan Misra. 2023. Emu video: Factorizing text-to-video generation by explicit image conditioning. *arXiv preprint arXiv:2311.10709*.

DeepMind Google. 2024. Veo-2.

DeepMind Google. 2025. Introducing gemini 2.5 flash image, our state-of-the-art image model.

Yuwei Guo, Ceyuan Yang, Ziyan Yang, Zhibei Ma, Zhijie Lin, Zhenheng Yang, Dahua Lin, and Lu Jiang. 2025. Long context tuning for video generation. *arXiv preprint arXiv:2503.10589*.

Tiankai Hang, Shuyang Gu, Dong Chen, Xin Geng, and Baining Guo. 2024. Cca: Collaborative competitive agents for image editing. *arXiv preprint arXiv:2401.13011*.

Yingqing He, Tianyu Yang, Yong Zhang, Ying Shan, and Qifeng Chen. 2022. Latent video diffusion models for high-fidelity long video generation. *arXiv preprint arXiv:2211.13221*.

Roberto Henschel, Levon Khachatryan, Daniil Hayrapetyan, Hayk Poghosyan, Vahram Tadevosyan, Zhangyang Wang, Shant Navasardyan, and Humphrey Shi. 2024. Streamingt2v: Consistent, dynamic, and extendable long video generation from text. *arXiv preprint arXiv:2403.14773*.

Sirui Hong, Xiawu Zheng, Jonathan Chen, Yuheng Cheng, Jinlin Wang, Ceyao Zhang, Zili Wang, Steven Ka Shing Yau, Zijuan Lin, Liyang Zhou, and 1 others. 2023. Metagpt: Meta programming for multi-agent collaborative framework. *arXiv preprint arXiv:2308.00352*.

Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, Yaohui Wang, Xinyuan Chen, Limin Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. 2024. VBench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.

Sangwon Jang, Jaehyeong Jo, Kimin Lee, and Sung Ju Hwang. 2024. Identity decoupling for multi-subject personalization of text-to-image models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.

KlingAI. 2024. Kling-v1.5.

Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, and 1 others. 2024. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*.

Guohao Li, Hasan Hammoud, Hani Itani, Dmitrii Khizbullin, and Bernard Ghanem. 2023. Camel: Communicative agents for" mind" exploration of large language model society. *Advances in Neural Information Processing Systems*, 36:51991–52008.

Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson, and Jianfeng Gao. 2019. Storygan: A sequential conditional gan for story visualization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6329–6338.

Xinyao Liao, Xianfang Zeng, Liao Wang, Gang Yu, Guosheng Lin, and Chi Zhang. 2025. Motionagent: Fine-grained controllable video generation via motion field agent. *arXiv preprint arXiv:2502.03207*.

Chang Liu, Haoning Wu, Yujie Zhong, Xiaoyun Zhang, Yanfeng Wang, and Weidi Xie. 2024. Intelligent grimm-open-ended visual storytelling via latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6190–6200.

Shilong Liu, Hao Cheng, Haotian Liu, Hao Zhang, Feng Li, Tianhe Ren, Xueyan Zou, Jianwei Yang, Hang Su, Jun Zhu, and 1 others. 2025. Llava-plus: Learning to use tools for creating multimodal agents. In *European Conference on Computer Vision*, pages 126–142. Springer.

Luma. 2024. Dream machine.

MiniMax. 2024. Hailuoai.

The ModelScopeTeam. 2023. Modelscope: bring the notion of model-as-a-service to life.

OpenAI. 2024. Creating video from text.

OpenAI. 2025. Gpt image - openai platform. Accessed: 2025-04-30.

Joon Sung Park, Joseph C O'Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. 2023. Generative agents: Interactive simulacra of human behavior. arxiv. *arXiv preprint ArXiv:2304.03442*.

Tanzila Rahman, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, Shweta Mahajan, and Leonid Sigal. 2023. Make-a-story: Visual memory conditioned consistent story generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2493–2502.

Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image generation. In *International conference on machine learning*, pages 8821–8831. Pmlr.

Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10684–10695.

Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman. 2023. Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 22500–22510.

Runway. 2024a. Gen-3.

Runway. 2024b. Gen-4.

Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, and 1 others. 2022. Laion-5b: An open large-scale dataset for training next generation image-text models. *Advances in neural information processing systems*, 35:25278–25294.

StabilityAI. 2024. Stable diffusion 3.5.

Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.

Jiuniu Wang, Zehua Du, Yuyuan Zhao, Bo Yuan, Kexiang Wang, Jian Liang, Yaxi Zhao, Yihen Lu, Gengliang Li, Junlong Gao, and 1 others. 2024a. Aesopagent: Agent-driven evolutionary system on story-to-video production. *arXiv preprint arXiv:2403.07952*.

Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, and 1 others. 2024b. Lavie: High-quality video generation with cascaded latent diffusion models. *IJCV*.

Yaohui Wang, Xin Ma, Xinyuan Chen, Cunjian Chen, Antitza Dantcheva, Bo Dai, and Yu Qiao. 2024c. Leo: Generative latent image animator for human video synthesis. *International Journal of Computer Vision*, pages 1–13.

Zhenyu Wang, Aoxue Li, Zhenguo Li, and Xihui Liu. 2024d. Genartist: Multimodal llm as an agent for unified image generation and editing. *arXiv preprint arXiv:2407.05600*.

Zihao Wang, Shaofei Cai, Anji Liu, Yonggang Jin, Jinbing Hou, Bowei Zhang, Haowei Lin, Zhaofeng He, Zilong Zheng, Yaodong Yang, and 1 others. 2024e. Jarvis-1: Open-world multi-task agents with memory-augmented multimodal language models. *IEEE Transactions on Pattern Analysis and Machine Intelligence*.

Wenming Weng, Ruoyu Feng, Yanhui Wang, Qi Dai, Chunyu Wang, Dacheng Yin, Zhiyuan Zhao, Kai Qiu, Jianmin Bao, Yuhui Yuan, and 1 others. 2024. Art-v: Auto-regressive text-to-video generation with diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7395–7405.

Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Beibin Li, Erkang Zhu, Li Jiang, Xiaoyun Zhang, Shaokun Zhang, Jiale Liu, and 1 others. 2024. Autogen: Enabling next-gen llm applications via multi-agent conversations. In *First Conference on Language Modeling*.

Weijia Wu, Mingyu Liu, Zeyu Zhu, Xi Xia, Haoen Feng, Wen Wang, Kevin Qinghong Lin, Chunhua Shen, and Mike Zheng Shou. 2025a. Moviebench: A hierarchical movie level dataset for long video generation. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 28984–28994.

Weijia Wu, Zeyu Zhu, and Mike Zheng Shou. 2025b. Automated movie generation via multi-agent cot planning. *arXiv preprint arXiv:2503.07314*.

Zhifei Xie, Daniel Tang, Dingwei Tan, Jacques Klein, Tegawend F Bissyand, and Saad Ezzini. 2024. Dreamfactory: Pioneering multi-scene long video generation with a multi-agent framework. *arXiv preprint arXiv:2408.11788*.

Jinbo Xing, Menghan Xia, Yong Zhang, Haoxin Chen, Wangbo Yu, Hanyuan Liu, Gongye Liu, Xintao Wang, Ying Shan, and Tien-Tsin Wong. 2025. Dynamicrafter: Animating open-domain images with video diffusion priors. In *European Conference on Computer Vision*, pages 399–417. Springer.

Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. 2018. Attngan: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1316–1324.

Deshun Yang, Luhui Hu, Yu Tian, Zihao Li, Chris Kelly, Bang Yang, Cindy Yang, and Yuexian Zou. 2024a. Worldgpt: a sora-inspired video ai agent as rich world models from text and image inputs. *arXiv preprint arXiv:2403.07944*.

Hui Yang, Sifu Yue, and Yunzhong He. 2023a. Auto-gpt for online decision making: Benchmarks and additional opinions. *arXiv preprint arXiv:2306.02224*.

Shuai Yang, Yuying Ge, Yang Li, Yukang Chen, Yixiao Ge, Ying Shan, and Yingcong Chen. 2024b. Seed-story: Multimodal long story generation with large language model. *arXiv preprint arXiv:2407.08683*.

Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. 2023b. Mm-react: Prompting chatgpt for multimodal reasoning and action. *Preprint*, arXiv:2303.11381.

Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. 2023a. Lumos: Learning agents with unified data, modular design, and open-source llms. In *ICLR 2024 Workshop on Large Language Model (LLM) Agents*.

Shengming Yin, Chenfei Wu, Huan Yang, Jianfeng Wang, Xiaodong Wang, Minheng Ni, Zhengyuan Yang, Linjie Li, Shuguang Liu, Fan Yang, and 1 others. 2023b. Nuwa-xl: Diffusion over diffusion for extremely long video generation. *arXiv preprint arXiv:2303.12346*.

Zhengqing Yuan, Yixin Liu, Yihan Cao, Weixiang Sun, Haolong Jia, Ruoxi Chen, Zhaoxu Li, Bin Lin, Li Yuan, Lifang He, and 1 others. 2024. Mora: Enabling generalist video generation via a multi-agent framework. *arXiv preprint arXiv:2403.13248*.

Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. 2017. Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the IEEE international conference on computer vision*, pages 5907–5915.

Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. 2023. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv preprint arXiv:2311.04145*.

Zhuosheng Zhang and Aston Zhang. 2024. You only look at screens: Multimodal chain-of-action agents. *Preprint*, arXiv:2309.11436.

Canyu Zhao, Mingyu Liu, Wen Wang, Weihua Chen, Fan Wang, Hao Chen, Bo Zhang, and Chunhua Shen. 2024. Moviedreamer: Hierarchical generation for coherent long visual sequence. *arXiv preprint arXiv:2407.16655*.

Boyuan Zheng, Boyu Gou, Jihyung Kil, Huan Sun, and Yu Su. 2024. Gpt-4v (ision) is a generalist web agent, if grounded. *arXiv preprint arXiv:2401.01614*.

Mingzhe Zheng, Yongqi Xu, Haojian Huang, Xuran Ma, Yexin Liu, Wenjie Shu, Yatian Pang, Feilong Tang, Qifeng Chen, Harry Yang, and 1 others. 2025. Videogen-of-thought: Step-by-step generating multi-shot video with minimal manual intervention. *arXiv preprint arXiv:2503.15138*.

Cailin Zhuang, Ailin Huang, Wei Cheng, Jingwei Wu, Yaoqi Hu, Jiaqi Liao, Hongyuan Wang, Xinyao Liao, Weiwei Cai, Hengyuan Xu, and 1 others. 2025. Vistorybench: Comprehensive benchmark suite for story visualization. *arXiv preprint arXiv:2505.24862*.

## A More Related Works

### A.1 Video Generation

Long video generation (He et al., 2022; Ge et al., 2022; Wang et al., 2024c) aims to produce longer videos than previous text-to-video (T2V) (ModelScopeTeam, 2023; Wang et al., 2024b) and image-to-video (I2V) (Girdhar et al., 2023; Xing et al., 2025; Zhang et al., 2023) models. Auto-regressive approaches (Weng et al., 2024; Henschel et al., 2024) extend videos by conditioning on earlier frames, yet often result in repetitive motions with limited narrative development. Attempts to directly integrate storytelling and video generation (Yin et al., 2023b; Zheng et al., 2025) suffer from poor visual quality and unstructured video composition, leading to subpar viewer experiences. Meanwhile, commercial models have achieved notable success in both T2V (GenmoTeam, 2024; Kong et al., 2024; OpenAI, 2024) and I2V generation (Google, 2024; Luma, 2024; MiniMax, 2024; Ang and Baole, 2025), and some (KlingAI, 2024; Runway, 2024a,b) support temporal extension.

### A.2 Image Generation

Recent advances in image generation (Li et al., 2019; Xu et al., 2018; Zhang et al., 2017; Ramesh et al., 2021) have been propelled by deep learning techniques. To fulfill the task of visual storytelling, some works (Liu et al., 2024; Yang et al., 2024b; Rahman et al., 2023) generate image sequences from text and image-caption pairs. To ensure subject consistency (Gal et al., 2022), DreamBooth (Ruiz et al., 2023) adopts subject-specific fine-tuning, while GPT-Image-1 (OpenAI, 2025) introduces context-aware generation for identity preservation. Recent models such as NanoBanana (Google, 2025) and HunyuanImage 3.0 (Cao et al., 2025) further advance high-fidelity and instruction-following generation, offering stronger semantic control and compositional reasoning capabilities. In this work, MAViS leverages off-the-shelf models to generate keyframes for each shot, which serve as the basis for subsequent video animation.

### A.3 AI Agents

AI agents, empowered by large language models (Brown et al., 2020; Touvron et al., 2023a,b), have shown broad applicability in reasoning, planning, web interaction, and visual generation. General-purpose LLMs underpin agent frameworks (Yang et al., 2023a; Yin et al., 2023a; Li et al., 2023; Hong et al., 2023; Park et al., 2023), while multimodal models (Achiam et al., 2023; Anil et al., 2023) extend agent capabilities across modalities (Yang et al., 2023b; Zhang and Zhang, 2024; Zheng et al., 2024; Wang et al., 2024e; Liu et al., 2025; Wang et al., 2024d; Hang et al., 2024; Liao et al., 2025). In video generation, systems (Yang et al., 2024a; Yuan et al., 2024; Wang et al., 2024a) explore world modeling, multi-agent coordination, and retrieval-augmented synthesis, yet often yield static content with limited transitions. Hierarchical pipelines (Xie et al., 2024; Wu et al., 2025b) introduce script-to-shot workflows but still rely on manual scripting and LoRA tuning.

## B Explanation and Discuss about the Script Writing Guideline

Due to space limitations in the main text, we provided only a brief overview of the Script Writing Guideline. In this section, we offer a more detailed explanation of the guideline and elaborate on the rationale behind each design constraint.

In the Script Writing Guideline, "*Scriptwriters should minimize transitions between tightly connected spatial locations*" refers to avoiding consecutive shots set in spatially adjacent backgrounds—for example, the first shot showing the character entering a room by going through a door, and the next shot showing the character from the front inside the room. Such transitions can break immersion due to background inconsistencies, such as mismatched door appearances between shots.

The script writing guidelines are based on practical experience to improve compatibility with generation tools. E.g., "*not to place the same character in consecutive shots*" stems from the fact that such sequences often involve multi-angle cuts in the same setting. Given current model limitations in maintaining background consistency, this approach is not yet reliably achievable.

The constraints in the script writing guideline are based on the following empirical observations:

● *Avoid placing successive shots in the same background*: We measured VBench background consistency for adjacent shots sharing the same background and found it to be only 78.42, compared to 95.88 of single-shot. This justifies avoiding such situations.

● *Minimize transitions between tightly connected spatial locations*: For pairs of shots with strong spa-

tial connections, we used GPT to assess whether their backgrounds were spatially related as described in the script. Only 26 out of 100 pairs were judged to be consistent, indicating such transitions should be avoided.

• *Avoid complex actions within a single shot*: Shots with complex actions achieved a VBench overall consistency of 23.25, compared to 27.76 for simple actions, suggesting simpler actions improve generation reliability.

• *Avoid fine-grained visual elements within a single shot*: Shots containing fine-grained visual elements scored 22.68 in VBench overall consistency, while those without such elements scored 27.76. Thus, they should be avoided.

• *Align style with user's quest*: Using GPT to judge alignment with user instructions, only 82 out of 100 scripts were fully aligned without constraints, compared to 100 with constraints in place.

• *Ensure script structural completeness*: A complete structure is essential for enabling subsequent stages of the generation pipeline.

## C  Implementation Details

All agents in our framework are implemented using AutoGen (Wu et al., 2024) and powered by GPT-4o (Achiam et al., 2023). The T2I model pool comprises FLUX.1 (BlackForestLabs, 2023), Stable Diffusion 3.5 (Rombach et al., 2022), and GPT-Image (OpenAI, 2025), while the I2V model pool includes Veo2 (Google, 2024), Gen-3 (Runway, 2024a), Gen-4 (Runway, 2024b), LUMA (Luma, 2024), HunyuanVideo-I2V (Kong et al., 2024), and Wan2.1 (Ang and Baole, 2025). We use MUDI (Jang et al., 2024) for LoRA training and Hailuo (MiniMax, 2024) for the T2A generation. The maximum number of 3E iterations is set to 4 for script writing, shot designing, and voice designing; 2 for the T2I process; 1 for the I2V process; and 5 for subtitle refining. As the first work targeting end-to-end long-sequence video storytelling, no suitable benchmark dataset exists. Therefore, we constructed a test set of 20 user prompts, each targeting a one-minute video for evaluation.

To standardize intermediate outputs for downstream processes, we introduced a Script Consolidator agent in the Script Writing stage. This agent formats the script after each 3E iteration to ensure consistency. In both the T2I and I2V processes, due to limitations on the number of images that can be processed simultaneously, the T2I Evaluator and I2V Evaluator adopt a batched evaluation strategy. Following this, T2I Selector and I2V Selector agents are employed to select the best candidates based on all evaluation results. Specifically, the I2V Evaluator samples 8 evenly sampled frames from each video candidate for evaluation.

The T2I process in the Shot Designing stage mirrors that of Character Modeling, with the difference being that the T2I model pool in the Shot Designing stage incorporates the LoRA models trained during Character Modeling. We use MUDI (Jang et al., 2024) for LoRA training, which supports multi-subject generation, allowing a single LoRA to serve an entire script.

Given that different models in the I2V model pool require reference images at different resolutions, keyframes are center-cropped and resized before I2V execution according to model-specific requirements. To ensure consistency, all generated videos are center-cropped and resized to $1280 \times 720$ resolution.

The T2I model pool includes the open-source models FLUX (BlackForestLabs, 2023) and Stable Diffusion 3.5 (StabilityAI, 2024), as well as the API-based GPT-Image (OpenAI, 2025). The I2V model pool comprises open-source models Hunyuan (Kong et al., 2024) and Wan2.1 (Ang and Baole, 2025), along with API-based models Gen3 (Runway, 2024a), Gen4 (Runway, 2024b), LUMA (Luma, 2024), and Veo2 (Google, 2024). For T2A, we use the Minimax Hailuo API (MiniMax, 2024). All experiments were conducted on two NVIDIA H100 GPU.

More samples generated by MAViS framework are shown in Figure. 9 and Figure. 10.

### C.1  Metrics

we constructed a test set of 20 user prompts, each targeting a one-minute video for evaluation. We employ CLIP and Inception Score to evaluate prompt consistency and generation quality in the keyframe generation stage, and use VBench (Huang et al., 2024) to assess video animation performance. Additionally, a user study is conducted to assess real user preferences and validate the practical effectiveness of MAViS and baselines. For all evaluation metrics, higher values indicate better performance.

**Compliance Rate**  Compliance Rate is an objective and transparent metric defined as the fraction of outputs that are fully compliant with the guidelines

| Stage | Agent/Model | Role | Reviewer | Refiner |
|---|---|---|---|---|
| Script Writing | Scriptwriter | Write script | Structure Reviewer, Content Reviewer, Style Reviewer | - |
| Shot Designing | Shot Designer | Design each shot of the script | Shot Reviewer | - |
| Keyframe Generation | T2I Prompt Generator | Generate T2I prompt for each shot | - | - |
| Keyframe Generation, Character Modeling | T2I Models | Text-to-image generation for each shot/character | T2I Evaluator | T2I Refiner |
| Video Animation, Character Modeling | I2V Prompt Generator | Generate I2V prompt for each shot/character | - | - |
| Video Animation | I2V Models | Image-to-video generation for each shot | I2V Evaluator | I2V Refiner |
| Character Modeling | Character Prompt Generator | Generate T2I and I2V prompts for each character | - | - |
| Character Modeling | Caption Generator | Generate captions for each character images | - | - |
| Audio Generation | Voice Designer | Choose voice ID, emotion, and BGM | Voice Reviewer | - |
| Audio Generation | T2A Model | Text-to-audio generation for each shot | - | Subtitle Refiner |

Table 7: **Summary of Agents.**

| Methods | FID↓ | Inception↑ | CLIP↑ |
|---|---|---|---|
| DreamFactory | 7.03[†] | 169.71[†] | 30.92[†] |
| **MAViS** | - | 12.81 | 34.22 |

Table 8: Compare with DreamFactory (Xie et al., 2024). † denotes the reported score in the paper of the baseline method.

| Methods | CLIP↑ | Inception↑ | AS↑ | CLIP-sim↑ |
|---|---|---|---|---|
| MovieDreamer | 19.52[†] | 8.64[†] | 6.05[†] | 0.70[†] |
| **MAViS** | 34.22 | 12.81 | 7.25 | 0.71 |

Table 9: Compare with MovieDreamer (Zhao et al., 2024). † denotes the reported score in the paper of the baseline method. "AS" denotes Aesthetic Score used in (Schuhmann et al., 2022).

| Methods | Aesthetic↑ | Image Quality↑ | Consistency (avg.)↑ | P. Consistency↑ |
|---|---|---|---|---|
| LCT | 60.79[†] | 67.44[†] | 95.65[†] | 30.14[†] |
| **MAViS** | 63.17 | 72.91 | 95.92 | 34.22 |

Table 10: Compare with LCT (Guo et al., 2025). † denotes the reported score in the paper of the baseline method. "Consistency (avg.)" represents the average score of subject and background consistency. "P. Consistency" denotes Prompt Consistency.

(Line 596). A "fully compliant" output is defined as follows: • Shot Designing: All seven fields are present + in the correct format. • Voice Designing: Downloadable BGM + consistent voice ID + correct format. • Subtitle Refining: Narration duration ≤ video duration.

If any of the above rules are violated, the workflow will produce an error. Therefore, the Compliance Rate directly measures the necessity and effectiveness of the 3E principle across these three stages, and there is no better metric that can faithfully reflect this property.

## C.2 Baselines

We compare MAViS with three long video generation baselines: VGoT, Mora, and MovieAgent. Among them, Mora and MovieAgent adopt a modular design and are equipped with the same image and video generation models as MAViS, and are also powered by GPT-4o. Since Mora and MovieAgent do not support script generation, we provide them with scripts produced by our script writing stage. For MovieAgent, which requires

manual training of LoRA models, we generate character images based on the provided scripts and train LoRA models for each script.

## D  Compare with More Baselines

In our experiments, we conducted direct comparisons with three open-source baselines: VGoT (Zheng et al., 2025), Mora (Yuan et al., 2024), and MovieAgent (Wu et al., 2025b). Subsequently, we compared our results with the reported performances of non-open-source baselines.

We first compare our method with DreamFactory (Xie et al., 2024). Due to the lack of real videos in automatic long-sequence video storytelling, it is not feasible to compute the FID metric. As shown in Table 8, MAViS achieves higher scores across most metrics; however, there is a substantial gap in Inception Score. We adopted a widely used implementation of the Inception Score, which may differ from the version used by DreamFactory, potentially leading to differences in score range. Since the baseline is not open-sourced, we are unable to verify the exact cause of this discrepancy.

Next, we compare our method with MovieDreamer (Zhao et al., 2024). As shown in Table 9, MAViS outperforms this baseline across all evaluation metrics.

Subsequently, we compare our method with

| Metric | Question |
|---|---|
| Narrative Expressiveness | Which video performed best overall in terms of narrative engagement, story coherence, and viewing experience? |
| Visual Quality | Which video was the most visually expressive in terms of emotion, atmosphere, and cinematic feel? |
| User Prompt Alignment | Which video best aligns with the original user prompt (e.g., character setup, plot elements, or keywords)? |
| Character Consistency | Which video had the most consistent character appearances across scenes? (clothing, hairstyle, facial features) |
| Character Naturalness | Which video had the most natural and correct character generation? Consider anatomical realism, physical motion plausibility, and the absence of duplicate or mistakenly generated characters. |
| Background Consistency | Which video had the most consistent and logically coherent background style across shots? |
| Background Realism | Which video had the most natural background elements and physically plausible movements? |

Table 11: Metrics and Questions for User Study.

| Method | Keyframe Generation | | Video Animation | | | | | |
|---|---|---|---|---|---|---|---|---|
| | CLIP ↑ | Inception ↑ | T. Flick. ↑ | M. Smooth. ↑ | Sub. Cons. ↑ | Bg. Cons. ↑ | Aesthetic ↑ | I. Quality ↑ |
| VGoT (Zheng et al., 2025) | 22.32 | 8.41 | 99.00 | 99.28 | **98.98** | **98.76** | **79.95** | 64.62 |
| Mora (Yuan et al., 2024) | 34.02 | 12.70 | 97.25 | 99.21 | 94.58 | 95.01 | 64.41 | 71.39 |
| MovieAgent (Wu et al., 2025b) | 31.64 | 10.69 | 97.05 | 99.18 | 93.44 | 94.62 | 63.38 | 91.25 |
| **MAViS (ours)** | **34.15** | **12.79** | **99.08** | **99.52** | 95.54 | 96.10 | 64.05 | **72.68** |

Table 12: Compare with baselines on MovieBench and ViStoryBench.

LCT (Guo et al., 2025), using evaluation metrics sourced entirely from VBench. As shown in Table 10, MAViS consistently outperforms this baseline across all metrics.

# E  User Study

In the User Study, participants were asked to vote for the best video in each group based on seven evaluation criteria. As listed in Table 11, the criteria include: Narrative Expressiveness, Visual Quality, User Prompt Alignment, Character Consistency, Character Naturalness, Background Consistency, and Background Realism. Prior to voting, detailed guidelines for each metric were provided to ensure consistency and reliability in evaluation. Some screenshots of the user study HTML interface is shown in Figures 6, 7, and 8.

# F  Evaluation on the Existing Dataset

We are establishing an end-to-end agentic system from users' single line of story theme to final pixels plus other multimodalities. However, existing long video benchmarks, MovieBench (Wu et al., 2025a) and ViStoryBench (Zhuang et al., 2025), focus only on benchmarking T2I and I2V models', not on the entire agentic system. In particular, they ignore the script generation and scripts' compatibility with existing T2I and I2V models' capabilities.

Nevertheless, we sampled three examples (each consisting of 10 shots) from both MovieBench and ViStoryBench. As shown in Table 12, MAViS consistently outperforms the compared methods.

| $I$ | Narrative | Visual | User Align. | C.Rate | Time(Second) |
|---|---|---|---|---|---|
| 0 | 12.58 | 18.44 | 19.56 | 74.00 | 9.80 |
| 1 | 25.32 | 22.85 | 24.40 | 86.00 | 20.90 |
| 2 | 32.54 | 24.82 | 26.49 | 94.50 | 28.24 |
| 3 | 33.48 | 25.60 | 26.52 | 98.50 | 36.13 |
| 4 | 34.00 | 26.00 | 26.67 | 100.00 | 41.58 |

Table 13: Quality-Efficiency Analysis on scriptwriting.

| $I$ | CLIP | Inception | Natural. | C.Rate | Time(Second) |
|---|---|---|---|---|---|
| 0 | 31.25 | 10.54 | 72.00 | 68.00 | 175.19 |
| 1 | 32.25 | 11.30 | 84.00 | 81.00 | 381.32 |
| 2 | 34.22 | 12.81 | 98.00 | 92.50 | 557.47 |
| 3 | 34.26 | 12.78 | 98.50 | 96.00 | 602.06 |

Table 14: Quality-Efficiency Analysis on Keyframe Generation.

# G  Quality-Efficiency Analysis

In the stages where the 3E principle is applied, we evaluated the impact of the maximum iteration number ($I$) on generation quality and time using the same metrics as in our paper. It's worth noting that C. Rate refers to Compliance Rate, defined as the proportion of outputs that fully adhere to their respective guidelines. It is used to measure the success rate of each stage's generation.

For script writing, as shown in Table 13, the performance is already satisfactory at $I = 2$. Increasing $I$ further can improve generation quality with minimal impact on overall efficiency.

For keyframe generation, as shown in Table 14, the performance is already satisfactory at $I = 2$. Increasing $I$ further can enhance generation quality, with only a slight impact on overall efficiency.

For video animation, as shown in Table 15, the performance is already satisfactory at $I = 2$. While increasing $I$ can further improve quality, it has a

| $I$ | CLIP | Inception | Natural. | Sub. Cons. | Dyn. Deg. | C.Rate | Time(Minute) |
|---|---|---|---|---|---|---|---|
| 0 | 32.05 | 11.22 | 60.00 | 93.00 | 44.58 | 72.50 | 61.17 |
| 1 | 34.53 | 13.22 | 92.50 | 95.72 | 35.62 | 86.00 | 124.21 |
| 2 | 34.62 | 13.31 | 95.00 | 95.51 | 35.68 | 94.00 | 176.77 |
| 3 | 34.35 | 13.18 | 96.50 | 95.46 | 35.54 | 98.00 | 215.38 |

Table 15: Quality-Efficiency Analysis on Video Animation.

| | Shot Designing | | Voice Designing | | Subtitle Refining | |
|---|---|---|---|---|---|---|
| $I$ | C.Rate | Time(Second) | C.Rate | Time(Second) | C.Rate | Time(Second) |
| 0 | 84.50 | 24.34 | 92.00 | 5.62 | 56.50 | 1.58 |
| 1 | 89.00 | 73.49 | 95.00 | 13.48 | 86.00 | 2.83 |
| 2 | 94.00 | 122.352 | 98.50 | 26.19 | 92.50 | 3.95 |
| 3 | 97.00 | 161.21 | 97.00 | 161.21 | 98.00 | 4.36 |
| 4 | 98.50 | 195.86 | 99.00 | 29.84 | 99.50 | 4.84 |
| 5 | 100.00 | 231.60 | 100.00 | 31.68 | 100.00 | 4.89 |

Table 16: Quality-Efficiency Analysis on shot designing, voice designing, and subtitle refining.

| Model/API | Cost | Time overhead | Duration | GPU |
|---|---|---|---|---|
| Veo 2 | 2.8\$ | 31.7s | 8s | - |
| Gen 3 | 0.5\$ | 23.7s | 10s | - |
| Gen 4 | 0.5\$ | 40.8s | 10s | - |
| Luma | 1.27\$ | 36.2s | 5s | - |
| Hunyuan-I2V | - | 28.7min | 5s | 1 H100 |
| Wan 2.1 | - | 36.3min | 3s | 1 H100 |

Table 17: Real-world Cost per Shot.

significant impact on overall efficiency. Therefore, we set $I = 1$ for video animation in our experiments.

For shot designing, voice designing, and subtitle refining, as shown in Table 16, increasing $I$ improves generation quality with minimal impact on overall efficiency. However, since errors in these stages can critically affect the entire generation process, we set $I = 5$ for these tasks in our experiments.

Additionally, since iterations may terminate before reaching the maximum iteration number, the generation time does not scale linearly with $I$.

## H Generation Time Overhead.

Using all models in the model pool, the average time for MAViS to generate a long video on two H100 GPUs is 13.63 hours. In contrast, it will take our internal artists 5 days to complete a pipeline similar to ours: data generation, evaluation, selection, and composition. The process can be significantly accelerated if only API-based I2V models are used. The resources required for a single video model to generate one 720p shot are as shown in Table 17.

## I Most Common Issues and Success Rate

We report the most common issues and success rates across all the stages that apply the 3E principle during generation:

• Script writing: Adjacent shots are set in the same or strongly connected locations, violating the script writing guidelines.

• Shot designing: Incomplete content, such as insufficient background detail.

• T2I: Facial and limb distortions.

• I2V: Abnormal character movements.

• Voice designing: BGM download failures.

• Subtitle refinement: The voice duration exceeds the clip length.

Using criteria such as avoiding visual tearing, distortion, physically implausible motion, ID inconsistency, and style inconsistency as overall success measures, MAViS achieves 0.8621, compared to 0.4875 without the 3E principle.

## J System Messages and Presentations

The system message of scriptwriter is listed in Figure 3, Figure 4, and Figure 5. For the system messages, generated video samples, and presentation video, please refer to the *agent.py* in the supplementary material. Some keyframe samples are listed in Figure 9 and Figure 10.

**Scriptwriter System Message**

Role Definition:
You are a professional **"AI Microfilm Scriptwriting Expert"**, specializing in writing **technically compliant shot-by-shot scripts** for AI video generation tools. The scripts must be **visually concise**, **logically coherent**, **engaging**, and **aligned with user-specified themes**, while ensuring **completely adherence to technical constraints** of AI video generation systems.

—

Task Objective:
Generate shot-by-shot scripts suitable for AI micro-video generation tools that meet the following criteria:
    - **Emotional**: Evokes emotional resonance such as family bonds, memories, or a sense of belonging.
    - **Concise**: Slow-paced and restrained storytelling with clean, uncluttered visuals.
    - **Logical**: Natural transitions between shots and scenes, with clear narrative clues.
    - **Technically feasible**: All content must be executable within AI generation constraints, avoiding technical conflicts.

—

Workflow:
1. **Clarify the Theme**
    - The script theme must align with the user's request. If unspecified, prioritize touching, emotional, and resonant themes, such as family, homecoming, or nostalgia.
    - Unless explicitly requested otherwise, **default to Western cultural and character backgrounds**; avoid Asian settings or characters.
2. **Script Structure **
    - Each script must include:
    - **Title**
    - **Character Definitions** (see section below for format)
    - **Shot-by-shot Script**: Minimum of 10 shots. If there's a final shot, name it '"Ending"' instead of "Shot 14" etc.
    - Each shot must include:
      - **Character**: Name of the character in frame, or "None" if no identifiable character appears.
        - Always use the exact name defined in the character setup.
        - Do not abbreviate, shorten, or alter the character name.
      - **Shot Content**: Character action + background + framing info.
      - **Subtitle/Voice-over**: Must be deliverable within 5 seconds, and match the tone and style of the story.
        - Important: Do not use dashes ("—") in subtitle. Replace all such instances with ellipses ("...").

—

Character Definition Rules:
    - Character names must follow these rules:
      - Use only a simple first name (e.g., "Ethan", "Astra").
      - Do not include last names, middle names, or titles/descriptors (Incorrect: "Ethan Carter", "The Guide". Correct: "Ethan", "Astra").
    - Each character must be defined with 1-2 sentences, including:
      - Approximate age (e.g., "around 30 years old")
      - Ethnicity or race (e.g., "Caucasian")
      - Appearance (e.g., hairstyle, skin tone, height)
      - Outfit style (e.g., "gray puffer jacket + jeans")
      - Notable features (e.g., "stoic gaze", "carries a backpack")
Background characters don't require individual profiles—just mention their presence in scene descriptions.

—

Figure 3: **Scriptwriter System Message.**

**Scriptwriter System Message**

```
Technical Constraints for AI Video Generation:

Shot Content & Action Rules:
   - Each shot corresponds to a  5-second video.
   - Each shot must include **only one primary action**. Avoid chaining actions; **complex
actions should be broken into multiple shots**.
      - Incorrect: "He opens the door, sees his mother, eyes turning red" → Too many linked
actions
      - Correct: Separate into "opens door," "mother's back in close-up," "his eyes turning red"
   - Interaction between characters should be simple (e.g., hug, holding hands), avoiding complex
physical interactions.
   - No smoking or drinking behaviors.

Scene Control & Repetition Avoidance:
   - Avoid using visually indistinct or similar settings (e.g., multiple shots in undifferentiated
desert terrain without time/weather/landform distinction). Shots that differ in time of day,
terrain type, or atmosphere are acceptable.
   - Avoid strong spatial interactions between consecutive shots (e.g., entering/exiting doors,
moving between rooms). Such transitions can break immersion due to background inconsistencies,
such as mismatched door appearances between shots.

Props and Visual Limitations:
   - Do not show specific textual content (e.g., on photos, screens, paper) unless it's a
**close-up** of the object.
   - Avoid complex visual details (e.g., phone screens, photos, mirrors) unless it's a close-up.
   - For still-life close-ups, **specify the background** (e.g., "close-up of a soup bowl on a
wooden table"), and avoid blurred backgrounds.

Characters and Presentation:
   - Do not show the **same character in two adjacent shots** (e.g., Shot 1 and Shot 2),
unless a different character or a non-character shot (e.g., object-only, background-only, or
non-identifiable body part) appears between them.
      - **Voice-over mentioning a character does not count as character appearance.** Only shots
where the character in the **Character** part count toward this rule.
      - Switching between different characters (e.g., Shot 1 with character Ethan → Shot 2 with
character Astra) **is allowed** and does not violate this rule.
      - The same character in none-adjacent shots (e.g., Ethan in Shot 1, Shot 3, and Shot 5)
**is allowed** and does not violate this rule.
      - If a shot contains only a non-identifiable body part (e.g., hand close-up), it should be
marked as "None" and does not violate the adjacent character appearance rule.
   - Each shot must clearly specify the **character** being filmed. If the shot does not include
a recognizable character or only includes a close-up of a body part that cannot identify the
character (e.g., hand, foot, leg—not the head), use '"None"'.
   - Character clothing must remain consistent throughout the script—**do not change outfits**.
   - Always use the exact name defined in the character definition.

Visual Detail Requirements:
   - Each shot should provide **sufficient background detail** to support visual generation,
such as location features, ambient lighting, weather, or temporal clues. Exact weather/time may
be omitted if irrelevant to the visual setting.
   - Shots with characters, especially memory scenes, must describe the character's distinct
features (e.g., ethnicity and age).

Visual Dynamics Requirement:
   - Each shot with a character must include a **light dynamic** action — a simple, single action
that can be naturally completed within a 5-second shot — to avoid completely static visuals.
      - Incorrect Example: Character stands still for an extended period
      - Correct Example: Character walking, wiping glass, tidying clothes, etc.
   - Light dynamic actions are recommended to enhance pacing and maintain visual vitality.

_
```

Figure 4: **Scriptwriter System Message.**

```
Scriptwriter System Message

Style Guidelines:
   - Storyline should be clear, concise, emotionally layered, and appropriately paced.
   - Use simple and easily understandable language throughout the script. Avoid overly complex,
or obscure vocabulary.
   - Transitions between shots should be logical and help convey narrative flow.
   - Shot composition should have visual aesthetic appeal — avoid cluttered frames.
   - Subtitle/Voice-over must match the script's tone.        - For touching themes: subtitle
should be gentle, restrained, poetic.
      - For war or action: subtitle should be direct and accessible. Avoid overt sentimentality.
   - **Ending with black background and stylized text** (e.g., "The Blood God Descends" in
dripping font) is recommended for emotional closure. Indicate the text type and visual style.

—

Special Reminder:
   - If the user provides a detailed theme or request, you must **optimize the structure and
logic** of the script based on that input. Continuously adjust the script based on user feedback
to ensure it complies with AI video generation standards.
   - **StructureEvaluator**, **ContentEvaluator** and **StyleEvaluator** will point out issues
in the script and give you feedback. Each time feedback is received, you must revise the script
accordingly and output a complete, improved version.
   - Character names in each shot must match those in the character definition section. **No
abbreviations or shortened names.**

—

Please follow all the above rules strictly when generating the script.
```

Figure 5: **Scriptwriter System Message.**



Figure 6: Screenshot of the user study HTML interface with instructions and generated video candidates.

**User Input:** I want a mystery-themed video. The story should explore the theme of survival, following an interstellar archaeologist in an alien marketplace. The length of the final video should be about 1 minute.

1. Which video performed best overall in terms of narrative engagement, story coherence, and viewing experience?

A  B  C  D

2. Which video was the most visually expressive in terms of emotion, atmosphere, and cinematic feel?

A  B  C  D

3. Which video best aligns with the original user prompt (e.g., character setup, plot elements, or keywords)?

A  B  C  D

4. Which video had the most consistent character appearances across scenes? (clothing, hairstyle, facial features)

A  B  C  D

5. Which video had the most natural and correct character generation? Consider anatomical realism, physical motion plausibility, and the absence of duplicate or mistakenly generated characters.

A  B  C  D

6. Which video had the most consistent and logically coherent background style across shots?

A  B  C  D

7. Which video had the most natural background elements and physically plausible movements?

Figure 7: Screenshot of the user study HTML interface with seven evaluation questions.



**User Input:** I want a sci-fi-themed video. The story should explore the theme of love, following a creature of myth in a magical academy. The length of the final video should be about 1 minute.

1. Which video performed best overall in terms of narrative engagement, story coherence, and viewing experience?

A  B  C  D

2. Which video was the most visually expressive in terms of emotion, atmosphere, and cinematic feel?

A  B  C  D

3. Which video best aligns with the original user prompt (e.g., character setup, plot elements, or keywords)?

A  B  C  D

4. Which video had the most consistent character appearances across scenes? (clothing, hairstyle, facial features)

A  B  C  D

5. Which video had the most natural and correct character generation? Consider anatomical realism, physical motion plausibility, and the absence of duplicate or mistakenly generated characters.

Figure 8: Screenshot of the user study HTML interface with generated video candidates and seven evaluation questions.
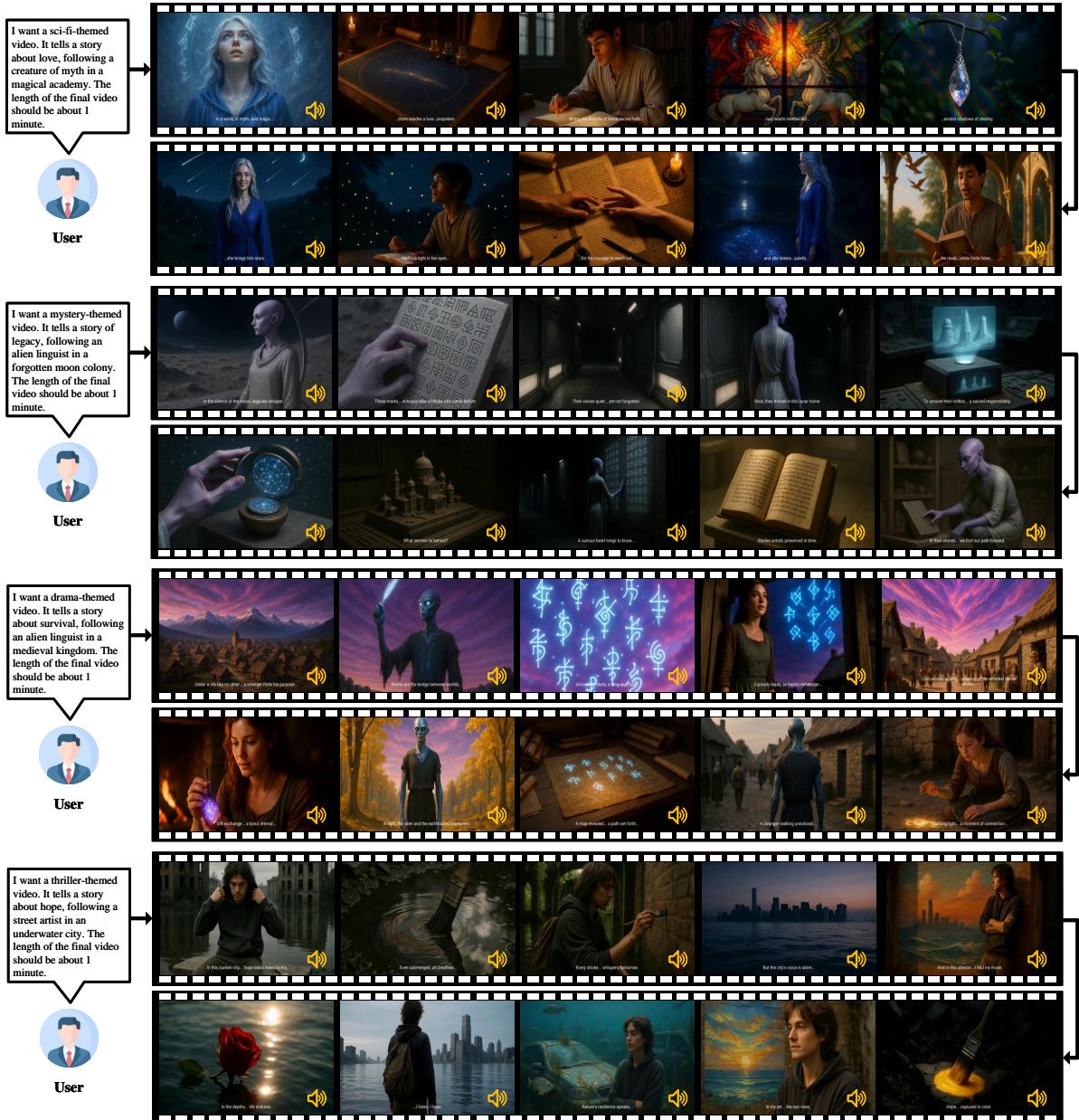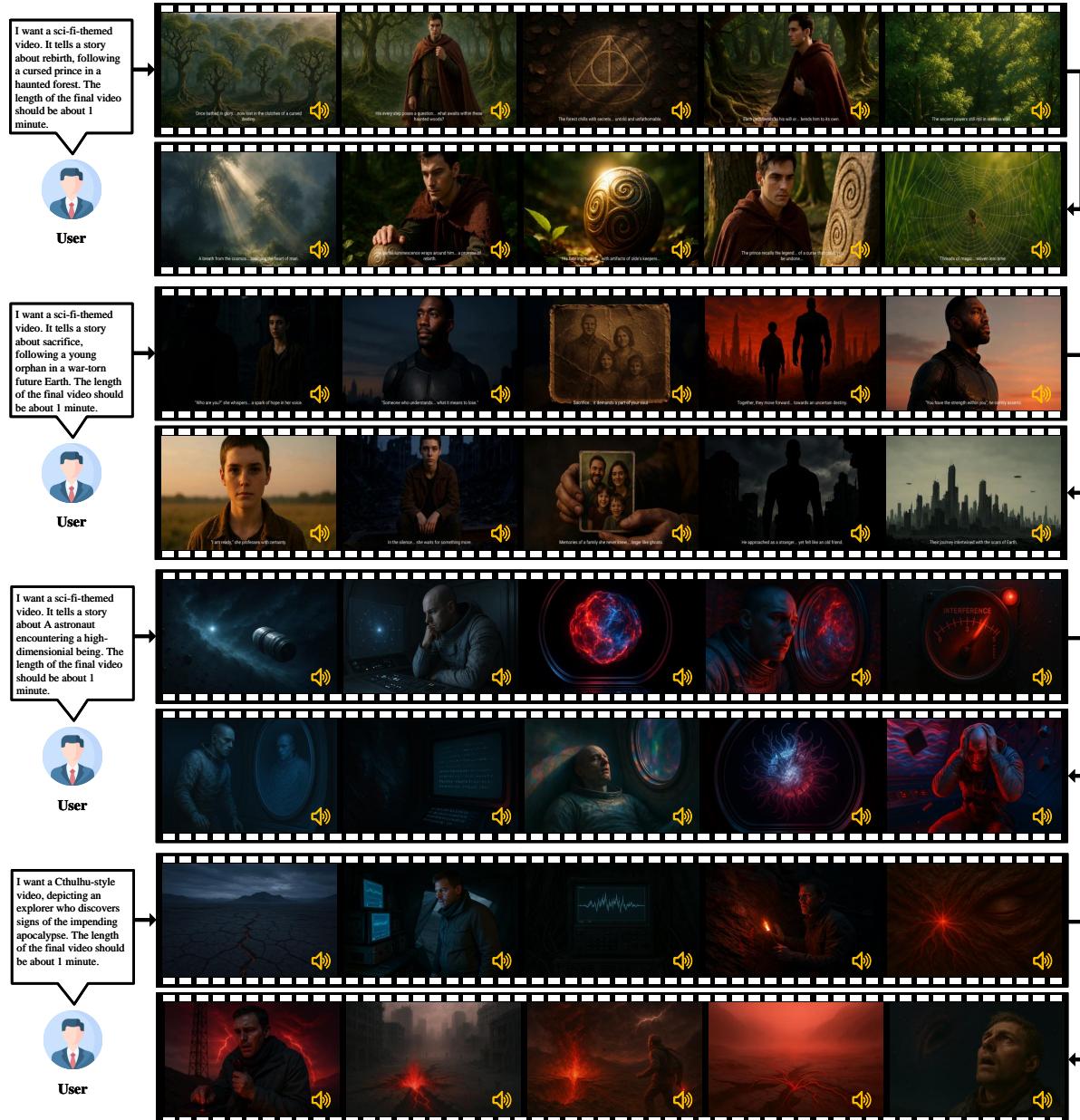
Figure 9: More samples generated by MAViS framework.

Figure 10: More samples generated by MAViS framework.