

# Automated Movie Generation via Multi-Agent CoT Planning

Weijia Wu , Zeyu Zhu , Mike Zheng Shou<sup>(✉)</sup>

Show Lab, National University of Singapore

## Abstract

*Existing long-form video generation frameworks lack automated planning, requiring manual input for storylines, scenes, cinematography, and character interactions, resulting in high costs and inefficiencies. To address these challenges, we present MovieAgent, an automated movie generation via multi-agent Chain of Thought (CoT) planning. MovieAgent offers two key advantages: 1) We firstly explore and define the paradigm of automated movie/long-video generation. Given a script and character bank, our MovieAgent can generate multi-scene, multi-shot long-form videos with a coherent narrative, while ensuring character consistency, synchronized subtitles, and stable audio throughout the film. 2) MovieAgent introduces a hierarchical CoT-based reasoning process to automatically structure scenes, camera settings, and cinematography, significantly reducing human effort. By employing multiple LLM agents to simulate the roles of a director, screenwriter, storyboard artist, and location manager, MovieAgent streamlines the production pipeline. Experiments demonstrate that MovieAgent achieves new state-of-the-art results in script faithfulness, character consistency, and narrative coherence. Our hierarchical framework takes a step forward and provides new insights into fully automated movie generation. The code and project website are available at: [Code](#) and [Website](#).*

## 1. Introduction

*“Every great movie should seem new every time you see it.”*

— Roger Ebert

Automated movie generation creates long-form videos with consistent characters, synchronized subtitles, and audio, given a script synopsis and character bank. It involves automating narrative planning, scene structuring,

and shot composition, replicating the hierarchical reasoning of real-world filmmaking. Most existing video generation research [1, 11, 32, 49] still focuses on short video generation without structured narratives, such as diffusion-based models like Stable Video Diffusion[1], Video LDM[2], and I2VGen-XL[44]. More recently, spatiotemporal transformer models, including Sora[3] and HunyuanVideo [18], have demonstrated superior performance in generating high-quality short videos (within 10 seconds) with realistic visuals and smoother motion dynamics. Compared to short-video generation, the development of long-form video generation [13, 26, 37, 42] has been relatively slow and still faces many challenges, such as maintaining narrative coherence, character consistency, structured scene transitions, and synchronized audio. DreamFactory[38] uses multi-agent systems and video generation models to synthesize keyframes, later expanded into long-form videos. Similarly, StoryAgent[13] employs multiple agents for customized storytelling video generation. However, these approaches are limited to basic long-video synthesis, lacking high-level planning and logically structured multi-scene narratives. They also fail to handle multi-object interactions, customization, and audio consistency, making them unsuitable for real-world applications. Thus, automated movie-level long-form generation remains an open challenge in the field.

Let us delve deeper into understanding what is essential and indispensable in the real-world movie production process, as shown in Figure 1 (a). In reality, real-world movie production is a **hierarchical** and **collaborative** process, involving multiple specialized roles: directors, screenwriters, storyboard artists, and cinematographers, who work together to maintain narrative coherence, character consistency, and structured scene transitions. Therefore, unlike short-video generation, movie level video generation is a complex process, including high-level cinematic themes and low-level cinematographic parameters, making it difficult to solve with a single model like an LLM or video generation framework.

Inspired by the real-world movie production process, we introduce multi-agent systems to simulate the roles of different filmmaking professionals and implement a hierarchical reasoning framework. As shown in Figure 1 (b), we propose

<sup>✉</sup> Corresponding author.

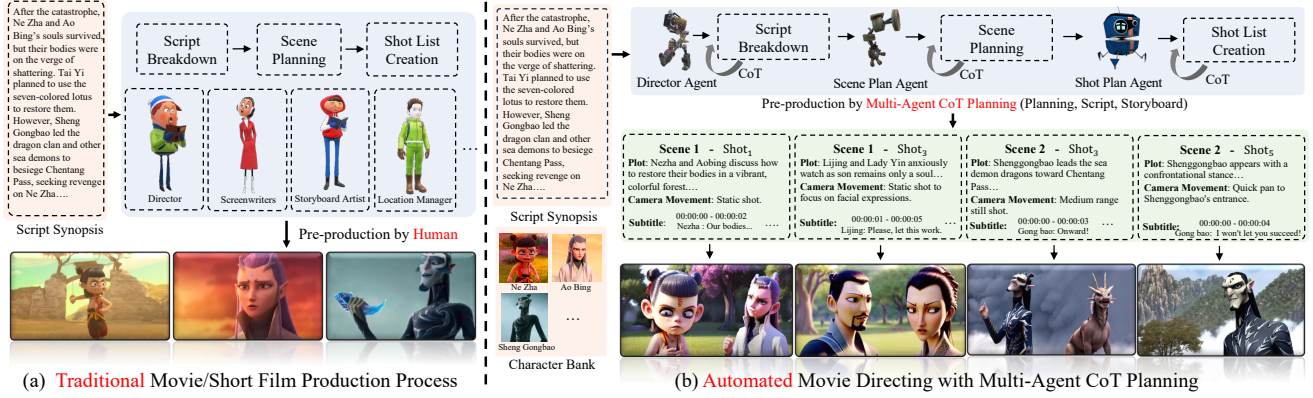


Figure 1. **Comparison of Traditional and Automated Movie Production.** Traditional filmmaking requires manual planning, while our MovieAgent automates script breakdown, scene planning, and shot design, enhancing efficiency and narrative coherence.

Table 1. **Comparison of Movie/Long-Video Production Costs.** ‘m’, ‘k’, ‘yrs’, and ‘mins’ denotes ‘million’, ‘thousand’, ‘years’, and ‘minutes’, respectively. Platform Videos refer to video content for platforms such as YouTube, TikTok. Compared to the high-cost, labor-intensive traditional long-video production, MovieAgent automatically generate video in under 10 minutes at almost no cost.

	Traditional Movies		Platform Videos	MovieAgent
	Frozen II	Inside Out 2		
Cost	150m	200m	1 to 100k	free
Time	6 to 7 yrs	5 to 6 yrs	1 to 100 days	2 to 10 mins

MovieAgent, a Multi-Agent CoT Planning framework that automatically structures and generates multi-scene, multi-shot videos with logical storylines, synchronized subtitles, and consistent character appearances. The key advantages of MovieAgent include: 1) **CoT-based Hierarchical Reasoning.** Unlike direct inference, which lacks structured and in-depth planning, CoT-based reasoning enables step-by-step, interpretable decision-making while recording the rationale behind each decision for use in subsequent steps. 2) **Near-zero Cost.** Compared to real-world movie production, which requires millions of dollars and over five years to complete, AI-driven movie generation (MovieAgent) is virtually cost-free, as shown in Table 1. 3) **Multi-Agent for Automated Filmmaking.** MovieAgent incorporates multiple specialized AI agents that simulate the roles of a director, screenwriter, and storyboard Artist. Therefore, the multi-agent framework can automatically decomposes a movie synopsis into structured acts, scenes, and shots, ensuring coherent plot development and seamless transitions, as shown in Figure 1 (b). These agents collaboratively enables precise control over both high-level cinematic themes and low-level cinematographic parameters simultaneously.

To summarize, the contributions of this paper are:

- We firstly explore and define the paradigm of automated movie/long-video generation. Given a script and character

bank, our MovieAgent can generates multi-scene, multi-shot long-form videos with a coherent narrative, ensuring character consistency, synchronized subtitles.

- MovieAgent employs a hierarchical CoT-based multi-agent reasoning framework to automate scene structuring, camera settings, and cinematography, reducing human effort. With internal CoT reasoning, MovieAgent effectively decouples and designs cinematic elements, including narrative structure, shot/scene composition, emotional tone, and subtitles.
- Experiments demonstrate that MovieAgent achieves state-of-the-art performance in automated storytelling and movie generation. Specifically, it excels in character consistency and narrative coherence, providing new insights into fully automated movie generation.

## 2. Related Works

### 2.1. Video Generation

Recent advancements in video generation have significantly improved quality and consistency, with approaches spanning diffusion models [5, 12, 36, 44, 46], and transformer frameworks [18, 39, 40]. Diffusion models have demonstrated remarkable success in image and video synthesis by gradually refining noise into realistic samples. VDM [12] pioneered the use of diffusion for video generation, introducing a spatiotemporal architecture to model frame dependencies. SVD [1] further advanced this by leveraging pre-trained text-to-image models for video generation, significantly improving quality. Lavie [30] introduces a high-quality video generation framework using cascaded latent diffusion models. More recently, SORA [3] and Hunyuan-video [18] showcased a highly coherent video generation system using advanced latent diffusion. Transformers have emerged as powerful architectures for modeling long-range dependencies. VideoPoet [17] introduced a VQ-VAE-based approach, tokenizing video frames and modeling them with

an autoregressive transformer. CogVideo [40] extended this paradigm with pre-trained text-to-video capabilities, leveraging hierarchical attention mechanisms to improve generation efficiency. VideoPoet [17] used a multimodal transformer to improve video-text understanding, enabling more controllable and expressive video synthesis. Despite advancements, existing frameworks still rely on manual input for narrative planning, cinematography, and scene composition. Our *MovieAgent* addresses these limitations by introducing a multi-agent, where agents simulate key filmmaking roles, enabling fully automated movie generation.

## 2.2. Story Visualization

Story visualization, which generates coherent visual sequences from text, is crucial for automated movie generation. Early GAN-based methods, such as StoryGAN [19], focused on maintaining narrative consistency in image sequences. With the rise of diffusion models, approaches like StoryDiffusion [50], Magic-Me [23] and DreamVideo [31] improved temporal coherence and motion dynamics in story-driven videos. Adapter-based techniques, including IP-Adapter [41], ROIctrl [8] and In-context LoRA [14], enabled efficient fine-tuning for personalized and character-consistent generation. Meanwhile, structured story-to-video frameworks like AutoStory [29] and Make-a-story [27] enhanced scene composition and transition planning. However, existing methods still lack automated high-level planning, often requiring manual intervention for cinematography, scene structuring. We introduce a multi-agent CoT-driven framework, enabling fully automated and coherent long-form movie generation.

## 2.3. LLM for Video Generation

Recent advancements in LLM-driven video generation [37, 51] have improved narrative structuring and interactive storytelling. VideoDirectorGPT [20] and VideoStudio [22] explored LLM-powered frameworks for scene composition, while Mora [43] enhanced video conceptualization for long-form coherence. For storyboarding and cinematic planning, DreamFactory [38] and StoryAgent [13] introduced LLM-based adaptive shot planning, reducing manual effort in camera control and character interactions. VideoGen-of-Thought [47] leveraged CoT reasoning to improve multi-shot video consistency. Although these methods enhance narrative structuring and storytelling with LLMs, they still require manual intervention or lack character and audio customization. In this paper, we firstly propose automated movie/long-video generation with a hierarchical CoT reasoning framework, which, given a script, character photos, and audio samples, automates planning, scene structuring, and cinematography for a more coherent and customizable filmmaking process.

## 3. Method

### 3.1. Task Formulation

Given a script synopsis  $S$  and a character bank  $C$ , the goal of automated movie generation is to generate a long-form video  $\hat{V}$  consisting of multiple scenes and shots while ensuring narrative coherence, character consistency, and audiovisual synchronization. Formally, the objective is to find an optimal mapping function:

$$\mathcal{F} : (S, C) \rightarrow \hat{V} \quad (1)$$

where character bank  $C = \{[\text{char}_k, I_k, A_k]\}_{k=1}^L$ ,  $L$  denotes the number of characters, and  $\text{char}_k$  is the  $k$ -th character names in the character list.  $I_k$  and  $A_k$  denotes the portrait images and audio samples of the character. The function  $\mathcal{F}(\cdot)$  refers to the automated movie generation function that systematically plans sub-scripts, scenes, and shots, along with various shot parameters, camera movements, and cinematographic settings. Ultimately, it generates a sequence of shots that collectively form the final movie output  $\hat{V} = \{\hat{V}_j^i \mid i = 1, 2, \dots, N, j = 1, 2, \dots, M\}$ , where  $\hat{V}_j^i$  denotes the  $j$ -th shot video in the  $i$ -th scene.

### 3.2. Automated Movie Generation

*MovieAgent* leverages a multi-agent Chain of Thought reasoning process (§3.3) to achieve structured and automated movie generation, as shown in Figure 2. The system decomposes the filmmaking process into a **hierarchical workflow**, simulating key roles in traditional movie production. Specifically, we introduce three specialized agents: Director Agent (§3.2.1), Scene Plan Agent (§3.2.2), and Shot Plan Agent (§3.2.3), which collaboratively structure narratives, plan scenes, and generate detailed cinematographic elements. Then, customized shot and audio generation (§3.2.4) is utilized to produce the final audio and video.

#### 3.2.1. Director Agent

The Director Agent is responsible for high-level narrative structuring. Given a script synopsis  $S$  and a character bank  $C$ , it systematically decomposes the storyline into sub-scripts  $\mathcal{S} = \{S_1, S_2, \dots, S_K\}$ , where each  $S_p$  represents  $p$ -th key narrative unit that contributes to the overall plot development. The segmentation function can be formulated as:

$$\mathcal{S} = \mathcal{F}_{\text{Director}}(S, C, p), \quad p \in \{1, \dots, K\} \quad (2)$$

where  $\mathcal{F}_{\text{Director}}(\cdot)$  is the decomposition function that segments the script  $S$  into meaningful sub-units  $\mathcal{S}$  based on character interactions, thematic continuity, and narrative flow. Specifically, the director agent follows a structured reasoning process: 1) *Identify Core Narrative Structure*: The director agent first analyzes the synopsis to identify main acts, key plot points, and turning points. 2) *Define Script Segmentation*: Based on these core narrative elements, the script

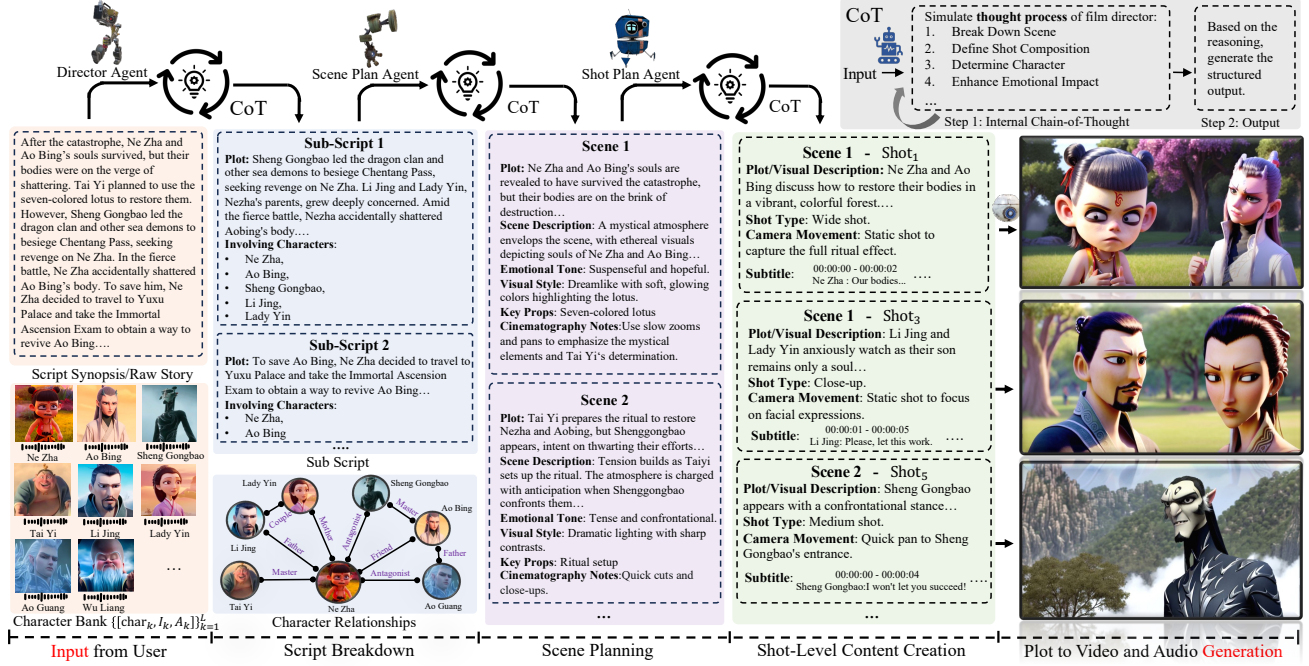


Figure 2. **The Overall Pipeline for MovieAgent.** The proposed framework employs a hierarchical CoT reasoning process with director, scene plan, and shot plan agents to automate long-form movie generation.

synopsis is divided into discrete, self-contained sub-scripts ( $S_p$ ). 3) *Ensure Logical Story Progression*: Each sub-script  $S_p$  maintains temporal and thematic coherence across  $S$  for a cohesive plot. 4) *Maintain Character Consistency*: The segmentation preserves the roles and relationships of characters from set  $C$ , ensuring their presence and interactions remain accurate throughout  $S$ . 5) *Justify the Division*: For each sub-script, a clear rationale for its segmentation (e.g., major event shift, emotional climax, new setting introduction) must be provided, serving as a reference for the subsequent step-by-step reasoning process.

### 3.2.2. Scene Plan Agent

With sub-scripts  $S$ , the next step involves determining the movie scenes, key scene elements, and scene boundaries. The Scene Plan Agent is designed to refine sub-scripts  $S$  into scene compositions  $\mathcal{P} = \{P_1, P_2, \dots, P_N\}$ , where each  $P_i$  represents a detailed scene with enriched descriptions for the  $i$ -th scene. The scene planning process is formalized as:

$$\mathcal{P} = \sum_{p=1}^K \mathcal{F}_{\text{Scene}}(S_p, C, i), \quad i \in \{1, \dots, N\} \quad (3)$$

where  $\mathcal{F}_{\text{Scene}}(\cdot)$  denotes the scene plan agent, which outputs a refined scene list  $\mathcal{P}$ . For the  $i$ -th scene  $P_i$ , the scene plan agent comprehensively summarizes factors such as involved characters, plot, emotional tone, visual style, and cinematography notes to thoroughly define the scene variables and expressive elements. Similar to the director agent, the scene

plan agent follow a structured reasoning process: 1) *Analyze the Narrative Structure*: The agent identifies key turning points and transitions, ensuring each scene forms a complete narrative with a clear start and end. 2) *Extract Key Scene Elements*: The model identifies all characters, their roles, interactions, and key events in each major scene. 3) *Define Scene Boundaries*: Finally, identify natural story breaks (e.g., location shifts, time jumps, emotional climaxes), ensuring each scene has a clear purpose, and justify each division (e.g., tone shift, new conflict). 4) *Justify the Division*: Preserve the internal Chain-of-Thought behind scene segmentation and reasoning to ensure traceability and analyzability.

### 3.2.3. Shot Plan Agent

Given the structured scenes  $\mathcal{P}$ , the shot plan agent is responsible for defining shot-level details, including character-aware plot, cinematographic parameters, and visual dynamics. Specifically, each scene  $P_i$  is further decomposed into detailed shot compositions  $\mathcal{V}^i = \{V_1^i, V_2^i, \dots, V_M^i\}$ , where each shot  $V_j^i$  captures distinct visual perspectives and cinematographic intentions of the  $j$ -th shot video in the  $i$ -th scene. Therefore, with scene list  $\mathcal{P} = \{P_1, P_2, \dots, P_N\}$ , the formalized function for shot-level decomposition is expressed as:

$$\mathcal{V} = \sum_{i=1}^N \mathcal{F}_{\text{Shot}}(P^i, C, j), \quad j \in \{1, \dots, M\} \quad (4)$$

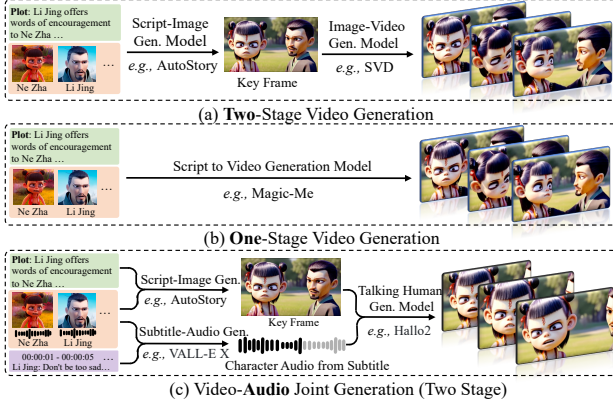


Figure 3. **Customized Shot-Level Video Generation for MovieAgent.** Current shot-level character-aware video generation approaches can be divided into three categories: (a) Keyframe-based two-stage video generation; (b) One-stage end-to-end video generation; (c) Keyframe-based joint video and audio generation.

where  $\mathcal{F}_{\text{Shot}}(\cdot)$  is the Shot Plan Agent function responsible for generating structured shots  $\mathcal{V} = \{V_j^i \mid i = 1, 2, \dots, N, j = 1, 2, \dots, M\}$ . Each shot level script  $V_j^i$  includes rich, structured shot script annotation, such as the involved characters, plot, camera movements, shot type, and character subtitles.

The Shot Plan Agent follows a structured reasoning workflow: 1) *Determine Shot Composition and Framing*: Identify appropriate shot types (e.g., wide, medium, close-up) and camera angles based on scene content and emotional impact. 2) *Specify Cinematographic Techniques*: Clearly define camera movements (static, pan, tilt, zoom, tracking), lighting styles, and visual effects necessary. 3) *Coordinate Visual Continuity*: Ensure visual coherence and consistency across shots within each scene, avoiding abrupt transitions or inconsistent visual styles. 4) *Align with Scene Narrative*: Each shot should advance the narrative or enhance emotions, with clear reasoning for traceability and analysis.

### 3.2.4. Customized Video and Audio Generation

At this stage, given the shot level script annotation  $V_j^i$  for the  $j$ -th shot in the  $i$ -th scene, the model invokes various customized image and video generation models (e.g., AutoStory [29], StoryDiffusion [50], and Magic-Me [23]) are used to produce the final shot-level video  $\hat{V}_j^i$ .

Current video generation models, such as SVD [1], DreamVideo [31] are unable to simultaneously support subtitle-to-audio generation. For the talking human generation task, some priors, such as Hallo2 [6] and Edtalk [28] primarily focuses on single human generation. Therefore, our current technology cannot fully address the simultaneous audio-video generation in a single model. Based on whether audio generation is required, we categorize the movie generation setting into two task:

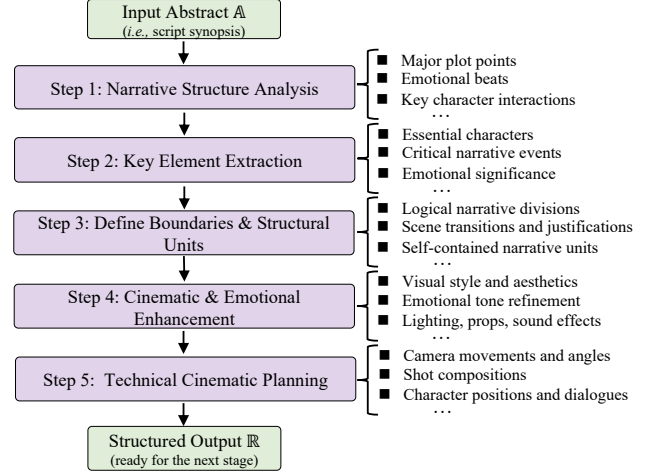


Figure 4. **Flowchart of the Internal Chain-of-Thought reasoning process.** Through Internal CoT, various agents can process and manage narrative elements more efficiently.

- **Pure Shot-level Video Generation in Figure 3 (a)-(b).** In this setting, we do not consider the audio generation of subtitle for the characters. Instead, we focus solely on generating pure video by modeling:  $\hat{V}_j^i = \mathcal{F}_{\text{Video}}(V_j^i, C)$ , where  $\mathcal{F}_{\text{Video}}(\cdot)$  can either be a two-stage video generation model (e.g., the combination of StoryDiffusion [50] and CogVideoX [40]) or a one-stage customized end-to-end video generation model (e.g., Magic-Me [23]), as illustrated in Figure 3 (a)-(b).
- **Video and Audio Joint Generation in Figure 3 (c).** In this setting, the character bank  $C$  must include the voice sample of each character, represented as  $\{[\text{char}_k, I_k, A_k]\}_{k=1}^L$ . Since no current model can simultaneously generate both audio and video, we adopt a two-stage video-audio joint generation strategy, as shown in Figure 3 (c). Formally, we express the formulation as:  $\hat{V}_j^i = \mathcal{F}_{\text{Talking}}(\mathcal{F}_{\text{Image}}(V_j^i, C), \mathcal{F}_{\text{Audio}}(V_j^i, C))$ , where  $V_j^i$  includes subtitles for all characters at the shot level. And  $\mathcal{F}_{\text{Talking}}(\cdot)$ ,  $\mathcal{F}_{\text{Image}}(\cdot)$ , and  $\mathcal{F}_{\text{Audio}}(\cdot)$  denote the talking-human generation model (e.g., Hallo2 [6]), customized image generation model (e.g., StoryDiffusion [50]), and customized audio generation model (e.g., VALL-E X [45]), respectively.

### 3.3. Internal Chain of Thought

The Internal Chain-of-Thought provides a general structured reasoning framework employed by various planning agents (e.g., Scene Plan Agent, Shot Plan Agent, Director Agent) to methodically translate abstract narrative and cinematic requirements into detailed, actionable plans. Formally, given an input abstract  $\mathbb{A}$  (e.g., scene synopsis, script synopsis), the Internal CoT generates structured reasoning:  $\mathbb{R} = \mathcal{F}_{\text{CoT}}(\mathbb{A})$ , where  $\mathcal{F}_{\text{CoT}}(\cdot)$  denotes the general internal reasoning func-

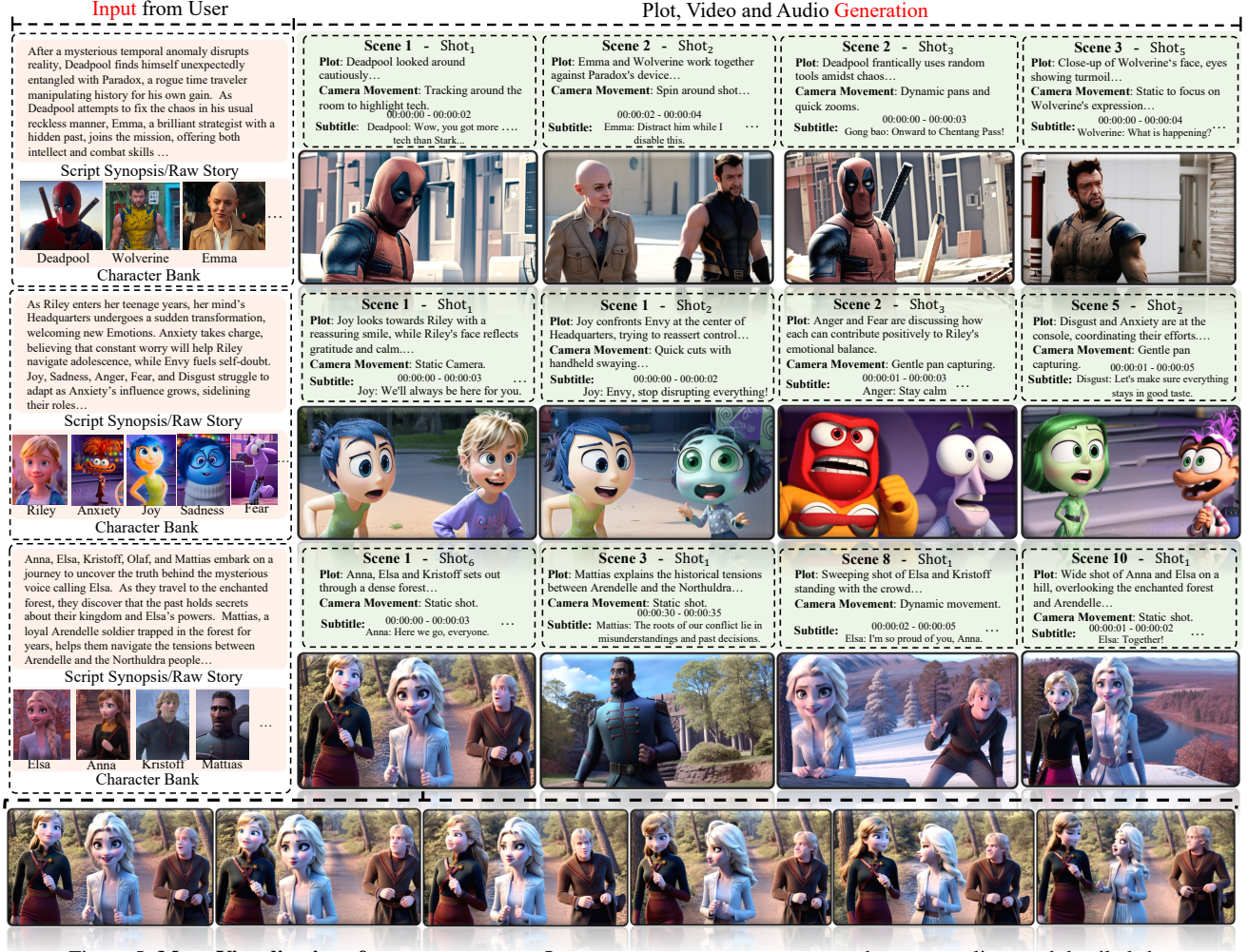


Figure 5. **More Visualizations for MovieAgent.** Our MovieAgent can generate coherent storylines and detailed shots.

tion utilized by planning agents, involving explicit, systematic reasoning steps prior to the final generation of detailed cinematic output.

As shown in Figure 4, the Internal CoT generally encompasses the five stages: (1) Narrative Structure Analysis, which identifies major plot points, emotional beats, and key character interactions; (2) Key Element Extraction, focusing on essential characters, critical narrative events, and emotional significance; (3) Define Boundaries & Structural Units, establishing logical narrative divisions, scene transitions, and self-contained units; (4) Cinematic & Emotional Enhancement, refining visual style, emotional tone, lighting, props, and sound effects; and (5) Technical Cinematic Planning, specifying camera movements, shot compositions, character positions, and dialogues. Upon completion of this structured internal reasoning, agents produce a comprehensive, structured output that encapsulates all detailed planning elements necessary for subsequent execution by next stage.

## 4. Experiments

### 4.1. Experiment Setting

**Metric.** Following prior works [4, 47], we evaluate the model using both automated metrics (*e.g.*, VBench [15]) and human voting. Automated metrics offer objective analysis but may not fully match human preferences, while user studies capture real preferences but can be biased. Since automated movie generation varies in shot count and lacks a ground truth video, metrics like FID [10] cannot be computed. Two A6000 GPUs were used for all experiments.

**Baseline.** Existing approaches, such as DreamFactory [38], StoryAgent [13], and StoryDiffusion [50], fail to address the automated movie generation task (Section 3.1), struggling with multi-characters consistency and automated script planning. Therefore, we decompose the task into three components: LLM-based script processing (GPT4-o [24], Deepseek-R1 [9], Llama3.3 [7]), image generation (AutoStory [29], StoryDiffusion [50]), and video generation (DreamVideo [31], Magic-Me [23]). For each com-

Table 2. **Performance of Automatic metric for Script to Keyframe/Video Generation on MovieAgent.** Models without character consistency (e.g., Open-Sora [48]) are excluded. ‘Subject Cons.’, ‘Bg Cons.’, ‘Motion Smth.’, and ‘Dyn. Degree’ refer to ‘Subject Consistency’, ‘Background Consistency’, ‘Motion Smoothness’, and ‘Dynamic Degree’ from the advanced VBench Metrics [15], respectively. GPT-4o serves as the LLM. And MovieAgent adopts multi-agent and internal CoT reasoning, while others rely on single-step generation.

Method	CLIP↑	Inception↑	VBench Metrics [16]/% ↑				
			Subject Cons.	Bg Cons.	Motion Smth.	Dyn. Degree	Aesthetic
<i>Script Synopsis to Keyframe/Storyboard Generation</i>							
StoryGen [21]	19.73	6.21	-	-	-	-	-
StoryDiffusion [50]	20.46	6.24	-	-	-	-	-
AutoStory [29]	20.21	6.01	-	-	-	-	-
MovieAgent	<b>22.12</b>	<b>7.23</b>	-	-	-	-	-
<i>Script Synopsis to Movie Generation</i>							
StoryDiffusion [50] + SVD [1]	21.39	8.36	93.64	93.78	96.30	74.48	56.69
StoryDiffusion [50] + CogVideoX [40]	21.83	9.01	93.45	94.56	96.60	27.89	56.05
AutoStory [29] + CogVideoX [40]	20.27	7.21	91.45	93.32	95.87	70.32	52.34
DreamVideo [31]	21.37	8.11	93.17	93.77	96.40	26.97	42.16
Magic-Me [23]	21.72	8.34	94.01	<b>94.68</b>	96.41	14.86	55.89
MovieAgent	<b>22.25</b>	<b>9.39</b>	<b>94.72</b>	93.52	<b>97.84</b>	<b>76.27</b>	<b>58.63</b>

ponent, we incorporate baseline models for evaluation and comparison, as detailed in Table 2.

**Evaluation Dataset.** Since the automated movie generation task (Section 3.1) is formally defined for the first time, we need to construct a new evaluation dataset. This evaluation dataset takes as input a script summary, character names, photos, and audio samples, and outputs a series of shot videos. To achieve this, we propose a test set, namely MoviePrompts, consisting of 10 script prompts: 8 prompts are derived from well-known movies (e.g., Ne Zha 2, Frozen II, and Inside Out 2), while the remaining 2 prompts (e.g., Fictional stories and characters) are privately designed by two annotators.

## 4.2. Performance Comparisons and Analysis

### 4.2.1. Automatic Metric on MoviePrompts

Table 2 experimental results for script-to-keyframe and script-to-movie generation. In keyframe generation, MovieAgent achieves the highest CLIP score 22.12 and Inception scores 7.23, indicating superior visual-semantic alignment and image quality. For movie generation, MovieAgent consistently outperforms nearly all baselines across VBench metrics [15], achieving the highest motion smoothness 97.84, dynamic degree 76.27, and aesthetic quality 58.63. These results highlight MovieAgent as a new state-of-the-art for automatic story-based video generation.

### 4.2.2. Human Rating

Figure 7 presents the user study for MovieAgent from two expert evaluators on the MoviePrompts dataset (10 movies). For a fairer and more fine-grained comparison, evaluators need rate each shot video on a scale of 1 to 5 based on the corresponding evaluation rules (detailed rules see supplementary materials). Due to the high cost of human



Figure 6. **Visualization Comparison for Different Methods.**

ratings, we limited assessments to key baselines (GPT-4o with DreamVideo and Magic Me). MovieAgent present a promising performance, outperforming the best baseline by up to 2 points on a five-point scale. Notably, it excels in Narrative Coherence (3.49), Visual Appeal (4.01), Script Faithfulness (3.89), Character Consistency (4.04), and Physical Law (3.42). Figure 6 provides a relevant visual comparison, while Figure 5 presents additional visualizations, including coherent storylines, keyframes, and shot videos. These results highlight the effectiveness of our multi-agent and CoT reasoning approach.

## 4.3. Ablation Study

We ablated three key aspects of MovieAgent: Internal CoT, LLM, and Multi-Agent. Currently, there are no automated metrics for evaluating script faithfulness and narrative coherence, which are key aspects of movie script and quality

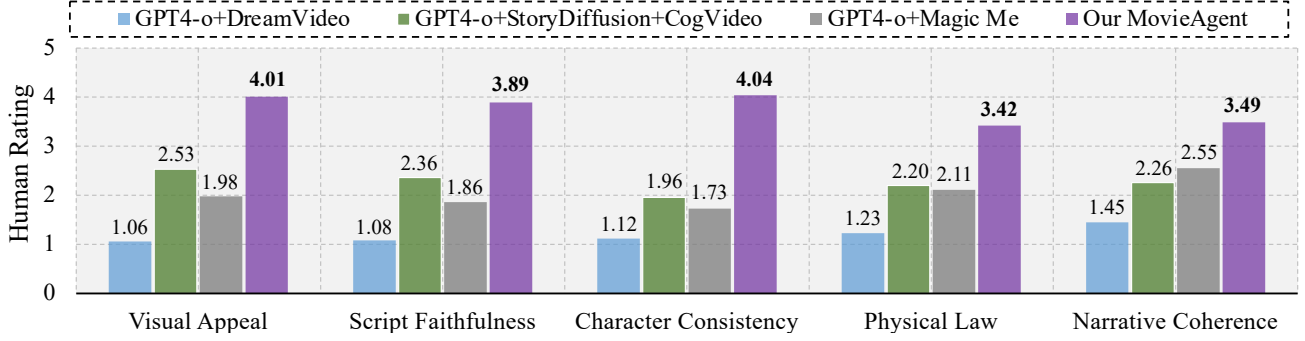


Figure 7. **Performance of User Study for Automated Movie Generation on MovieAgent.** Models without character consistency (e.g., Open-Sora [48]) are excluded. Under the 0-5 score rating system, MovieAgent demonstrates outstanding performance across multiple aspects, particularly in script faithfulness and character consistency.

Table 3. **Ablation Study for LLM, CoT, and Multi-Agent.** In this experiment, ROICtrl [8] and CogVideoX [40] are used as the image and video generation methods. Due to the high cost of human evaluation (up to 100 shots per movie), our ablation study focuses on three movies: Ne Zha 2, Frozen II, and Inside Out 2.

LLM Model	Internal CoT	Multi-Agent	Vis. Appeal	Script Faith.	Char. Consist.	Phys. Law	Narr. Coher.	Average
Llama3.3-70b	✓	✓	3.86	3.50	4.00	2.70	3.13	3.44
Deepseek-V3	✓	✓	3.74	3.66	3.78	3.07	3.45	3.54
Deepseek-R1	✓	✓	4.02	3.59	4.09	3.38	<b>3.79</b>	3.78
GPT4-o			3.89	3.36	4.08	3.37	3.09	3.55
GPT4-o		✓	4.02	3.69	4.13	3.46	3.31	3.72
GPT4-o	✓		3.92	3.38	4.08	3.37	3.29	3.61
GPT4-o	✓	✓	<b>4.04</b>	<b>3.92</b>	<b>4.11</b>	<b>3.49</b>	3.55	<b>3.82</b>

generation. Moreover, priors [4, 25, 33] confirm that human evaluation is more reliable than automated metrics for generation tasks. Therefore, our study focuses on human assessment.

#### 4.3.1. Effect of Internal Chain of Thought

Table 3 presents the ablation study for the internal Chain of Thought. Results show that GPT-4o with Internal CoT achieves a slight average score improvement (3.61 vs. 3.55 without CoT), with notable gains in Narrative Coherence (3.29 vs. 3.09). This is expected, as incorporating CoT enables step-by-step reasoning during story script generation, enhancing logical flow and coherence. By breaking down the reasoning process, CoT helps maintain narrative consistency and structure, as illustrated in Figure 4.

#### 4.3.2. Effect of Large Language Model

Table 3 compares the performance of various LLMs: Llama3.3-70b, Deepseek-V3, Deepseek-R1, and GPT4-o across human evaluation metrics including visual appeal, script faithfulness, character consistency, physical law, narrative coherence, and average. The results reveal that GPT4-o, particularly with multi-agent collaboration, achieves the highest average score of 3.82, outperforming Deepseek-V3 3.54, and Deepseek-R1 3.78. However, GPT-4o underperforms Deepseek-R1 in Narrative Coherence (3.55 vs. 3.79).

This is expected, as Deepseek-R1 is optimized for reasoning tasks with a built-in Internal CoT process, enabling more extensive reasoning and generating smoother, richer movie narratives.

#### 4.3.3. Effect of Multi-Agent

Table 3 presents the evaluation for the impact of multi-agent collaboration. The results show that multi-agent collaboration significantly enhances performance, with GPT4-o achieving an average score of 3.72, compared to 3.55 without multi-agent collaboration. The improvements in script faithfulness and narrative coherence are particularly significant, with increases of 0.33 and 0.22 on a five-point scale, respectively. This is reasonable because multi-agent system is a hierarchical structure. Multi-agent collaboration is more efficient than single-step script generation, as it better translates a script synopsis into a full movie script with detailed plot logic and parameters at the scene and shot levels.

## 5. Conclusion

In this paper, we firstly explore and define the paradigm of automated movie/long-video generation and propose MovieAgent, a multi-agent CoT-based framework for automated filmmaking. By integrating hierarchical reasoning and specialized AI agents, MovieAgent automates story structuring, scene planning, and shot composition, re-

Table 4. **Metric Summary.** We evaluate the generated long video on a 0-5 scale across five key aspects.

Metric	Description
Visual Appeal	Evaluates the overall visual quality, realism, and aesthetic consistency of the generated video.
Script Faithfulness	Measures how accurately the generated shot-level video content follows the provided shot level script and storyline.
Narrative Coherence	Assesses whether the narrative flows logically, maintaining consistent plot development.
Character Consistency	Checks whether characters maintain a stable appearance, behavior, and role throughout the movie.
Physical Law	Evaluates whether the generated video adheres to basic physical laws (e.g., gravity, motion realism).



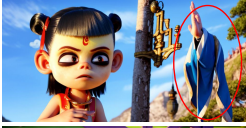
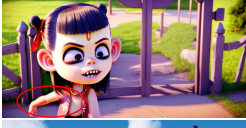

Score	Brief Explanation	Example
1	Completely <b>unrealistic</b> visuals, heavy distortions, severe artifacts, or glitches dominate the video. Scenes are barely recognizable.	
2	Poor quality with significant artifacts, <b>unstable</b> rendering, flickering, and unnatural textures. Characters and environments may be unrecognizable at times.	
3	Noticeable <b>inconsistencies</b> in textures, lighting, and transitions. Faces may distort occasionally, and the video lacks smoothness in motion.	
4	Generally <b>acceptable</b> quality, but some flickering, minor distortions, or occasional unnatural details remain.	
5	<b>Excellent</b> realism and aesthetic consistency. High resolution, smooth motion, and professional-grade visual appeal. Virtually indistinguishable from a high-quality animation.	

Figure 8. **Metric Criteria for Visual Appeal.**

ducing human intervention while ensuring narrative coherence and cinematographic quality. Experiments show that MovieAgent improves story consistency, character preservation, and audiovisual synchronization, addressing key challenges in AI-driven filmmaking. Our approach provides a scalable solution for automated storytelling, offering new insights into the future of AI-assisted movie production.


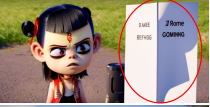



Score	Brief Explanation	Shot Level Plot	Example
1	The video is entirely <b>unrelated</b> to the script, failing to depict the intended scenes, characters, or events.	Tai yi offers his advice and encouragement, his expression a blend of wisdom and support. He gestures towards the <b>map</b> .	
2	Some scenes reflect the script, but there are noticeable <b>inconsistencies</b> , such as misplaced settings, incorrect character actions, or missing key moments.	Ne zha faces the <b>Exam Masters</b> , who stand in a semi-circle, their expressions stern and evaluative.	
3	The video mostly follows the script, but some <b>details</b> are incorrect or scenes are not fully aligned with the described narrative.	The mystical atmosphere envelops the scene, with a gentle, colorful <b>aura</b> surrounding Ne Zha and Ao Bing's souls.	
4	Well-aligned with the script, capturing most key moments correctly. Some <b>minor</b> deviations exist but do not significantly impact the storytelling.	Ne Zha sets out from the edge of a dense forest, <b>travel gear</b> in hand. The path ahead is lined with ethereal lights.	
5	<b>Perfectly</b> faithful to the script, with accurate depictions of all specified elements, actions, and narrative flow.	Nezha climbs a mountainous path, the Yuxu Palace visible on the horizon.	

Figure 9. **Metric Criteria for Script Faithfulness.**

## 6. Appendix

### 7. Metrics of Human Evaluation

#### 7.1. Definition for Metrics

To systematically evaluate the quality of generated long videos, we assess five key aspects on a 0-5 scale: visual appeal, script faithfulness, narrative coherence, character consistency, and adherence to physical laws. These metrics collectively measure both the aesthetic and structural integrity of the video. Specifically, they ensure that the visual quality remains realistic and consistent, the generated content accurately follows the script and storyline, and the narrative flows logically without abrupt transitions. Additionally, character consistency is crucial for maintaining a coherent identity and behavior throughout the video, while adherence to physical laws enhances motion realism and natural interactions within the scene. Together, these criteria provide a comprehensive framework for evaluating the effectiveness and realism of long video generation.

#### 7.2. Rules for User Study

In this user study, annotators are required to rate each generated long video on a 0-5 scale across five key evaluation aspects: Visual Appeal, Script Faithfulness, Narrative Coherence, Character Consistency, and Physical Law Adherence. Before scoring, the detailed criteria for each metric are established and provided as guidelines to ensure consistent evaluation.

**Visual Appeal.** The metric evaluates key factors such as realism, aesthetic consistency, artifact presence, motion smoothness, and rendering quality. Figure 8 presents the detailed criteria and examples corresponding to each score.

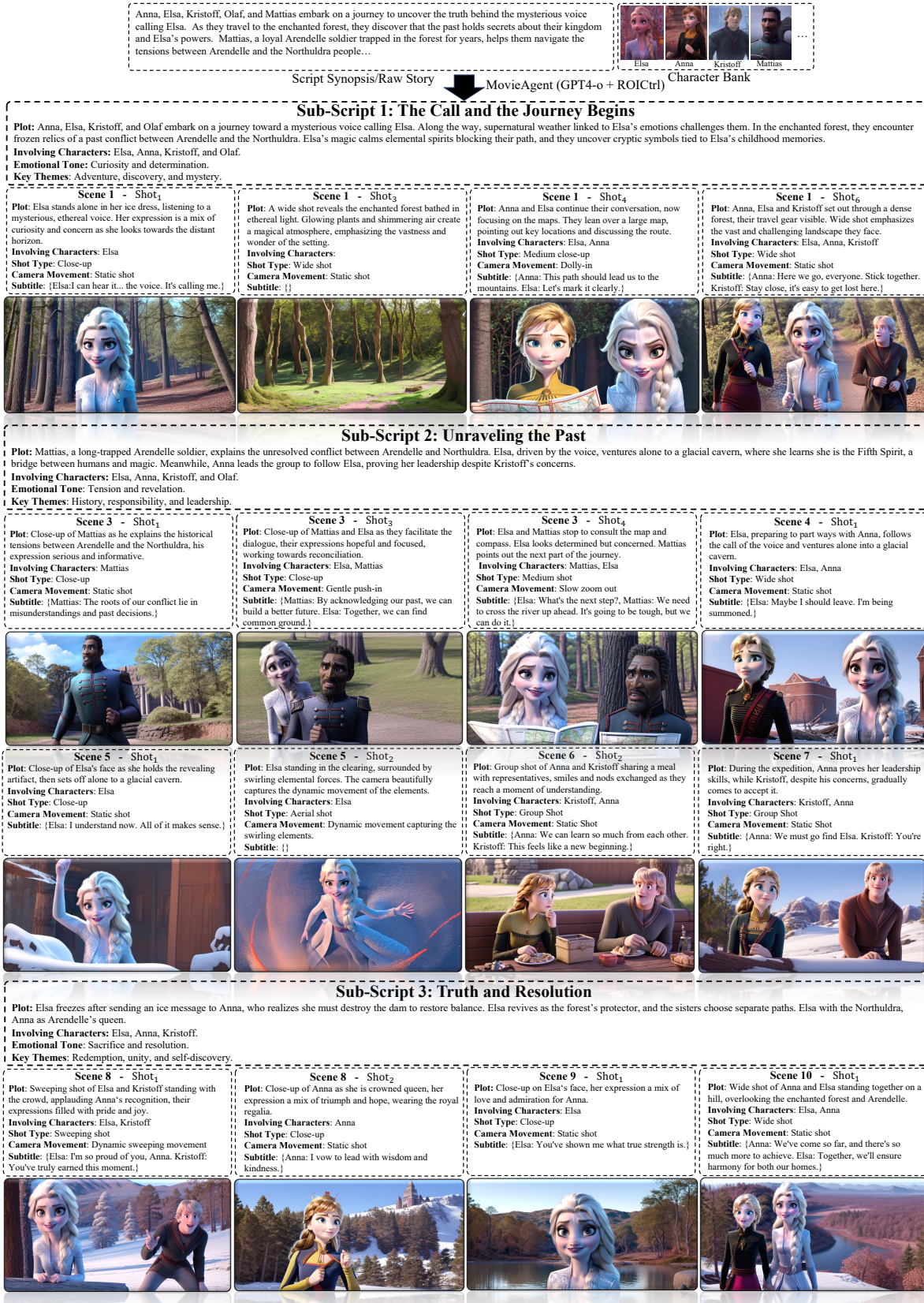


Figure 10. Detailed Visualization for Movie (e.g., FrozenII) Generation from MovieAgent.

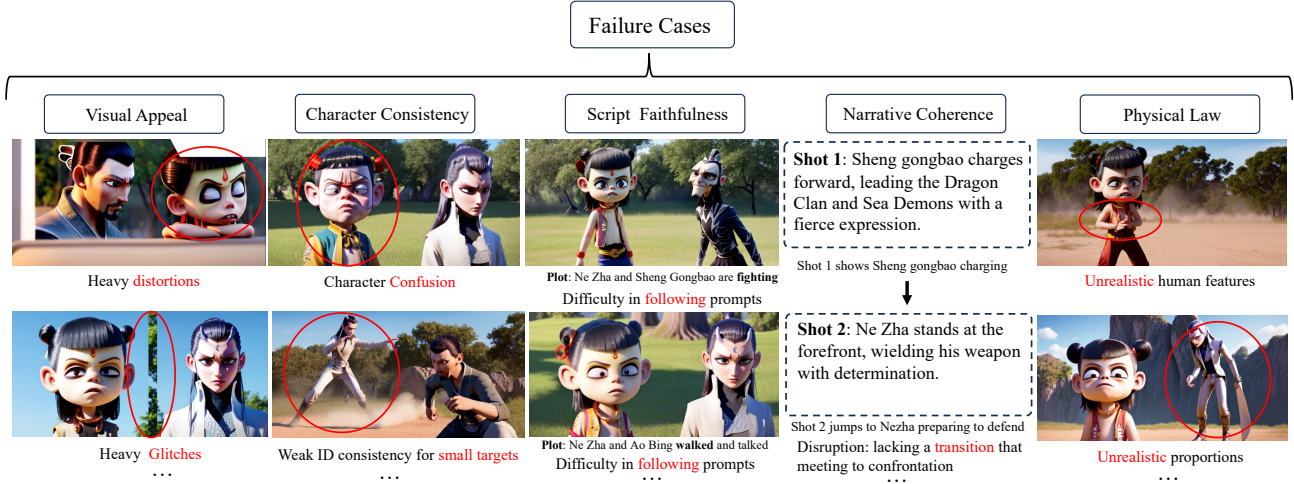


Figure 11. Failure Cases for Movie Generation from MovieAgent.

- *Score 1 (Severe Artifacts)*: Heavy distortions, glitches, and unrecognizable scenes dominate the video.
- *Score 2 (Poor Quality)*: Strong flickering, unstable textures, and character inconsistencies.
- *Score 3 (Inconsistent Rendering)*: Noticeable issues in textures, lighting, or motion transitions, though the content remains understandable.
- *Score 4 (Acceptable Quality)*: Good visual appeal with minor distortions or flickering but no major distractions.
- *Score 5 (High-Quality Rendering)*: Smooth, realistic, and professionally rendered visuals with high aesthetic consistency.

By applying this scoring framework, human evaluators can systematically assess the visual quality of generated videos, ensuring a consistent and objective evaluation standard.

**Script Faithfulness.** The metric assesses whether the generated content correctly follows the described plot, character actions, scene settings, and overall narrative structure. Figure 9 presents the detailed criteria.

- *Score 1 (Unrelated to Script)*: The video fails to depict the intended scenes, characters, or events, making it entirely disconnected from the script.
- *Score 2 (Major Inconsistencies)*: Some script elements are present, but there are misplaced settings, incorrect character actions, or missing key moments.
- *Score 3 (Partial Alignment)*: The video follows the script’s general structure but contains minor inaccuracies in details, scene transitions, or character behavior.
- *Score 4 (Good Faithfulness)*: Most key moments are accurately captured with only minor deviations that do not significantly affect storytelling.
- *Score 5 (Perfect Script Adherence)*: The video precisely reflects all specified elements, actions, and plot progression, achieving full alignment with the script.

This structured scoring system ensures an objective and con-

sistent evaluation of script adherence in generated videos, helping assess their narrative integrity.

**Narrative Coherence.** Narrative coherence assesses whether the storyline flows logically, maintaining consistent plot progression and smooth scene transitions. The detailed scoring criteria are as follows:

- *Score 1 (Completely Incoherent)*: The video lacks any meaningful structure, consisting of random and disconnected scenes that fail to form a logical storyline.
- *Score 2 (Frequent Disruptions)*: Contains abrupt cuts, illogical progressions, and a lack of continuity, making the plot difficult to follow.
- *Score 3 (Disjointed Flow)*: Some attempt at a storyline exists, but inconsistencies in sequencing and unnatural pacing disrupt the narrative flow.
- *Score 4 (Well-Structured Narrative)*: The storyline is logically structured with smooth transitions, though minor pacing issues may still be present.
- *Score 5 (Fully Coherent and Engaging)*: The video exhibits strong storytelling, seamless scene transitions, and consistent plot development, ensuring an immersive narrative experience.

**Character Consistency.** Character consistency evaluates whether characters maintain a stable appearance, behavior, and role throughout the video. The detailed scoring criteria are as follows:

- *Score 1 (Completely Inconsistent)*: Character appearances change frequently, making them unrecognizable.
- *Score 2 (Severe Inconsistencies)*: Major discrepancies in facial features, outfits, and personality.
- *Score 3 (Noticeable Variations)*: Some variations in facial structure, clothing, or expressions exist, and behaviors may not always align with character roles.
- *Score 4 (Good Consistency)*: Appearance and behavior remain stable, with only minor variations that do not sig-

nificantly impact immersion.

- *Score 5 (Perfect Consistency)*: Characters maintain a stable identity, facial features, and clothing consistently throughout the entire video.

**Physical Law.** Physical law adherence assesses whether the video follows basic physics principles, including motion realism, object interactions, and environmental consistency. The detailed scoring criteria are as follows:

- *Score 1 (Completely Unrealistic Physics)*: Objects float, characters phase through walls, and movements defy gravity without reason.
- *Score 2 (Frequent Violations)*: Major inconsistencies, such as erratic character movements, unnatural collisions, or unrealistic gravity interactions.
- *Score 3 (Limited Realism)*: Some aspects follow physics, but there are noticeable errors in object interactions, weight distribution, or motion continuity.
- *Score 4 (Good Physical Consistency)*: Characters and objects move naturally, with only minor deviations from real-world physics.
- *Score 5 (Flawless Adherence to Physics)*: The video demonstrates perfect motion realism, where all movements, interactions, and environmental effects behave.

## 8. More Visualization and Analysis

### 8.1. Advantages of Hierarchical Generation

Figure 10 presents more visualization for MovieAgent. From the image, the advantages of hierarchical generation are evident in its structured storytelling and controlled content organization. The image unfolds the narrative through Script Synopsis → Sub-Script → Scene → Shot, ensuring a clear and logical progression. Each level incorporates emotional themes, key characters, and cinematographic details, enhancing coherence and readability. This structured approach not only improves narrative consistency but also strengthens visual representation. Particularly in AI-generated content, hierarchical generation enables precise control over story details, ensuring seamless coordination between different levels and facilitating high-quality, film-grade script generation.

### 8.2. Discussion on Failure Cases and Improvements

Figure 11 presents some failure cases for MovieAgent, reveal several key challenges in automated movie generation. These issues span multiple aspects, including visual quality, character consistency, script faithfulness, narrative coherence, and adherence to physical laws. Below, we analyze the identified failures and propose potential improvements:

**Visual Appeal: Heavy Distortions and Glitches.** The generated images exhibit severe distortions and artifacts, particularly in facial structures and animations. These issues can break immersion and make the characters appear un-

natural. There are some potential improvements, such as further enhancing image or video generation models by using higher-quality data or incorporating a reward model to optimize outputs. The reward model can penalize low-quality image generations, driving the system toward producing more refined and visually appealing results.

**Character Consistency: Confusion and Weak ID Persistence.** In several frames, character identities are inconsistent, particularly for small targets. The same character may appear with different facial expressions, styles, or misaligned features, leading to continuity errors. At the data level, a possible improvement strategy is to train the model with larger and higher-quality datasets that ensure better character consistency in images and videos. At the algorithmic level, more efficient temporal consistency mechanisms can be explored, such as tracking embeddings across frames, enforcing ID-aware latent space regularization, and implementing strict feature-matching constraints to maintain coherence across sequences.

**Script Faithfulness: Prompt Following Issues.** Current video generation models struggle to produce complex human interactions, especially when prompts include actions like walking and talking. These models often fail to accurately capture and synchronize such interactions, making it challenging to generate realistic and coordinated human movements. A potential solution is to use higher-quality data and design more efficient strategies to enhance the prompt-following ability, enabling a stronger perception [34, 35] of various objects and interactions in the physical world.

**Narrative Coherence: Disjointed Scene Transitions.** The battle sequence between Shenggongbao and Nezha highlights an abrupt transition. Shot 1 shows Shenggongbao charging, but Shot 2 skips directly to Ne Zha preparing to defend, with no transition that meeting to confrontation. Some potential improvements include refining the LLM agent by utilizing a more precise internal Chain-of-Thought, along with hierarchical storyboarding to ensure smoother transitions between actions and enhance narrative coherence.

**Physical Law Violations: Unrealistic Human Features and Proportions.** Several frames showcase unnatural anatomical proportions and incorrect physics, such as distorted hands, limbs, and body movements that do not align with real-world constraints. Some potential improvements include integrating physics-based rendering and leveraging biomechanical constraints to ensure that character movements and proportions align with realistic human kinematics.

## References

- [1] Andreas Blattmann, Tim Dockhorn, Sumith Kulal, Daniel Mendelevitch, Maciej Kilian, Dominik Lorenz, Yam Levi, Zion English, Vikram Voleti, Adam Letts, et al. Stable video diffusion: Scaling latent video diffusion models to large

- datasets. *arXiv preprint arXiv:2311.15127*, 2023. 1, 2, 5, 7
- [2] Andreas Blattmann, Robin Rombach, Huan Ling, Tim Dockhorn, Seung Wook Kim, Sanja Fidler, and Karsten Kreis. Align your latents: High-resolution video synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 22563–22575, 2023. 1
  - [3] Tim Brooks, Bill Peebles, Connor Holmes, Will DePue, Yufei Guo, Li Jing, David Schnurr, Joe Taylor, Troy Luhman, Eric Luhman, Clarence Ng, Ricky Wang, and Aditya Ramesh. Video generation models as world simulators. 2024. 1, 2
  - [4] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, et al. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis. *arXiv preprint arXiv:2310.00426*, 2023. 6, 8
  - [5] Xinyuan Chen, Yaohui Wang, Lingjun Zhang, Shaobin Zhuang, Xin Ma, Jiashuo Yu, Yali Wang, Dahua Lin, Yu Qiao, and Ziwei Liu. Seine: Short-to-long video diffusion model for generative transition and prediction. In *The Twelfth International Conference on Learning Representations*, 2023. 2
  - [6] Jiahao Cui, Hui Li, Yao Yao, Hao Zhu, Hanlin Shang, Kaihui Cheng, Hang Zhou, Siyu Zhu, and Jingdong Wang. Hallo2: Long-duration and high-resolution audio-driven portrait image animation. *arXiv preprint arXiv:2410.07718*, 2024. 5
  - [7] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024. 6
  - [8] Yuchao Gu, Yipin Zhou, Yunfan Ye, Yixin Nie, Licheng Yu, Pingchuan Ma, Kevin Qinghong Lin, and Mike Zheng Shou. Roictrl: Boosting instance control for visual generation. *arXiv preprint arXiv:2411.17949*, 2024. 3, 8
  - [9] Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, et al. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*, 2025. 6
  - [10] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Advances in neural information processing systems*, 30, 2017. 6
  - [11] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 1
  - [12] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. *Advances in Neural Information Processing Systems*, 35:8633–8646, 2022. 2
  - [13] Panwen Hu, Jin Jiang, Jianqi Chen, Mingfei Han, Shengcai Liao, Xiaojun Chang, and Xiaodan Liang. Storyagent: Customized storytelling video generation via multi-agent collaboration. *arXiv preprint arXiv:2411.04925*, 2024. 1, 3, 6
  - [14] Lianghua Huang, Wei Wang, Zhi-Fan Wu, Yupeng Shi, Huanzhang Dou, Chen Liang, Yutong Feng, Yu Liu, and Jingren Zhou. In-context lora for diffusion transformers. *arXiv preprint arXiv:2410.23775*, 2024. 3
  - [15] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. *arXiv preprint arXiv:2311.17982*, 2023. 6, 7
  - [16] Ziqi Huang, Yinan He, Jiashuo Yu, Fan Zhang, Chenyang Si, Yuming Jiang, Yuanhan Zhang, Tianxing Wu, Qingyang Jin, Nattapol Chanpaisit, et al. Vbench: Comprehensive benchmark suite for video generative models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 21807–21818, 2024. 7
  - [17] Dan Kondratyuk, Lijun Yu, Xiuye Gu, José Lezama, Jonathan Huang, Grant Schindler, Rachel Hornung, Vighnesh Birodkar, Jimmy Yan, Ming-Chang Chiu, et al. Videopoet: A large language model for zero-shot video generation. *arXiv preprint arXiv:2312.14125*, 2023. 2, 3
  - [18] Weijie Kong, Qi Tian, Zijian Zhang, Rox Min, Zuozhuo Dai, Jin Zhou, Jiangfeng Xiong, Xin Li, Bo Wu, Jianwei Zhang, et al. Hunyuanvideo: A systematic framework for large video generative models. *arXiv preprint arXiv:2412.03603*, 2024. 1, 2
  - [19] Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson, and Jianfeng Gao. Storygan: A sequential conditional gan for story visualization. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 6329–6338, 2019. 3
  - [20] Han Lin, Abhay Zala, Jaemin Cho, and Mohit Bansal. Videodirectorgpt: Consistent multi-scene video generation via llm-guided planning. *arXiv preprint arXiv:2309.15091*, 2023. 3
  - [21] Chang Liu, Haoning Wu, Yujie Zhong, Xiaoyun Zhang, Yanfeng Wang, and Weidi Xie. 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, 2024. 7
  - [22] Fuchen Long, Zhaofan Qiu, Ting Yao, and Tao Mei. Videostudio: Generating consistent-content and multi-scene videos. In *European Conference on Computer Vision*, pages 468–485. Springer, 2024. 3
  - [23] Ze Ma, Daquan Zhou, Chun-Hsiao Yeh, Xue-She Wang, Xiuyu Li, Huanrui Yang, Zhen Dong, Kurt Keutzer, and Jiashi Feng. Magic-me: Identity-specific video customized diffusion. *arXiv preprint arXiv:2402.09368*, 2024. 3, 5, 6, 7
  - [24] OpenAI. Gpt-4o: Multimodal large language model. <https://openai.com/research>, 2025. Accessed: 2025-03-02. 6
  - [25] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 8

- [26] Adam Polyak, Amit Zohar, Andrew Brown, Andros Tjandra, Animesh Sinha, Ann Lee, Apoorv Vyas, Bowen Shi, Chih-Yao Ma, Ching-Yao Chuang, et al. Movie gen: A cast of media foundation models. *arXiv preprint arXiv:2410.13720*, 2024. 1
- [27] Tanzila Rahman, Hsin-Ying Lee, Jian Ren, Sergey Tulyakov, Shweta Mahajan, and Leonid Sigal. 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, 2023. 3
- [28] Shuai Tan, Bin Ji, Mengxiao Bi, and Ye Pan. Edtalk: Efficient disentanglement for emotional talking head synthesis. In *European Conference on Computer Vision*, pages 398–416. Springer, 2024. 5
- [29] Wen Wang, Canyu Zhao, Hao Chen, Zhekai Chen, Kecheng Zheng, and Chunhua Shen. Autostory: Generating diverse storytelling images with minimal human efforts. *International Journal of Computer Vision*, pages 1–22, 2024. 3, 5, 6, 7
- [30] Yaohui Wang, Xinyuan Chen, Xin Ma, Shangchen Zhou, Ziqi Huang, Yi Wang, Ceyuan Yang, Yinan He, Jiashuo Yu, Peiqing Yang, et al. Lavie: High-quality video generation with cascaded latent diffusion models. *International Journal of Computer Vision*, pages 1–20, 2024. 2
- [31] Yujie Wei, Shiwei Zhang, Zhiwu Qing, Hangjie Yuan, Zhiheng Liu, Yu Liu, Yingya Zhang, Jingren Zhou, and Hongming Shan. Dreamvideo: Composing your dream videos with customized subject and motion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6537–6549, 2024. 3, 5, 6, 7
- [32] Jay Zhangjie Wu, Yixiao Ge, Xintao Wang, Stan Weixian Lei, Yuchao Gu, Yufei Shi, Wynne Hsu, Ying Shan, Xiaohu Qie, and Mike Zheng Shou. Tune-a-video: One-shot tuning of image diffusion models for text-to-video generation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 7623–7633, 2023. 1
- [33] Weijia Wu, Zhuang Li, Yefei He, Mike Zheng Shou, Chunhua Shen, Lele Cheng, Yan Li, Tingting Gao, Di Zhang, and Zhongyuan Wang. Paragraph-to-image generation with information-enriched diffusion model. *arXiv preprint arXiv:2311.14284*, 2023. 8
- [34] Weijia Wu, Yuzhong Zhao, Hao Chen, Yuchao Gu, Rui Zhao, Yefei He, Hong Zhou, Mike Zheng Shou, and Chunhua Shen. Datasetdm: Synthesizing data with perception annotations using diffusion models. *Advances in Neural Information Processing Systems*, 36:54683–54695, 2023. 12
- [35] Weijia Wu, Yuzhong Zhao, Mike Zheng Shou, Hong Zhou, and Chunhua Shen. Diffumask: Synthesizing images with pixel-level annotations for semantic segmentation using diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1206–1217, 2023. 12
- [36] Weijia Wu, Zhuang Li, Yuchao Gu, Rui Zhao, Yefei He, David Junhao Zhang, Mike Zheng Shou, Yan Li, Tingting Gao, and Di Zhang. Draganything: Motion control for anything using entity representation. In *European Conference on Computer Vision*, pages 331–348. Springer, 2024. 2
- [37] Weijia Wu, Mingyu Liu, Zeyu Zhu, Xi Xia, Haoen Feng, Wen Wang, Kevin Qinghong Lin, Chunhua Shen, and Mike Zheng Shou. Moviebench: A hierarchical movie level dataset for long video generation. *arXiv preprint arXiv:2411.15262*, 2024. 1, 3
- [38] Zhifei Xie, Daniel Tang, Dingwei Tan, Jacques Klein, Tegawend F Bissyand, and Saad Ezzini. Dreamfactory: Pioneering multi-scene long video generation with a multi-agent framework. *arXiv preprint arXiv:2408.11788*, 2024. 1, 3, 6
- [39] Wilson Yan, Yunzhi Zhang, Pieter Abbeel, and Aravind Srinivas. Videogpt: Video generation using vq-vae and transformers. *arXiv preprint arXiv:2104.10157*, 2021. 2
- [40] Zhuoyi Yang, Jiayan Teng, Wendi Zheng, Ming Ding, Shiyu Huang, Jiazheng Xu, Yuanming Yang, Wenyi Hong, Xiaohan Zhang, Guanyu Feng, et al. Cogvideox: Text-to-video diffusion models with an expert transformer. *arXiv preprint arXiv:2408.06072*, 2024. 2, 3, 5, 7, 8
- [41] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. Ip-adapt: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arXiv:2308.06721*, 2023. 3
- [42] Shengming Yin, Chenfei Wu, Huan Yang, Jianfeng Wang, Xiaodong Wang, Minheng Ni, Zhengyuan Yang, Linjie Li, Shuguang Liu, Fan Yang, et al. NUWA-XL: Diffusion over diffusion for extremely long video generation. *arXiv preprint arXiv:2303.12346*, 2023. 1
- [43] Zhengqing Yuan, Yixin Liu, Yihan Cao, Weixiang Sun, Hao-long Jia, Ruoxi Chen, Zhaoxu Li, Bin Lin, Li Yuan, Lifang He, et al. Mora: Enabling generalist video generation via a multi-agent framework. *arXiv preprint arXiv:2403.13248*, 2024. 3
- [44] Shiwei Zhang, Jiayu Wang, Yingya Zhang, Kang Zhao, Hangjie Yuan, Zhiwu Qin, Xiang Wang, Deli Zhao, and Jingren Zhou. I2vgen-xl: High-quality image-to-video synthesis via cascaded diffusion models. *arXiv preprint arXiv:2311.04145*, 2023. 1, 2
- [45] Ziqiang Zhang, Long Zhou, Chengyi Wang, Sanyuan Chen, Yu Wu, Shujie Liu, Zhuo Chen, Yanqing Liu, Huaming Wang, Jinyu Li, et al. Speak foreign languages with your own voice: Cross-lingual neural codec language modeling. *arXiv preprint arXiv:2303.03926*, 2023. 5
- [46] Rui Zhao, Yuchao Gu, Jay Zhangjie Wu, David Junhao Zhang, Jia-Wei Liu, Weijia Wu, Jussi Keppo, and Mike Zheng Shou. Motiondirector: Motion customization of text-to-video diffusion models. In *European Conference on Computer Vision*, pages 273–290. Springer, 2024. 2
- [47] Mingzhe Zheng, Yongqi Xu, Haojian Huang, Xuran Ma, Yexin Liu, Wenjie Shu, Yatian Pang, Feilong Tang, Qifeng Chen, Harry Yang, et al. Videogen-of-thought: A collaborative framework for multi-shot video generation. *arXiv preprint arXiv:2412.02259*, 2024. 3, 6
- [48] Zangwei Zheng, Xiangyu Peng, Tianji Yang, Chenhui Shen, Shenggui Li, Hongxin Liu, Yukun Zhou, Tianyi Li, and Yang You. Open-sora: Democratizing efficient video production for all. <https://github.com/hpcaitech/Open-Sora>, 2024. 7, 8
- [49] Daquan Zhou, Weimin Wang, Hanshu Yan, Weiwei Lv, Yizhe Zhu, and Jiashi Feng. Magicvideo: Efficient video generation with latent diffusion models. *arXiv preprint arXiv:2211.11018*, 2022. 1

- [50] Yupeng Zhou, Daquan Zhou, Ming-Ming Cheng, Jiashi Feng, and Qibin Hou. Storydiffusion: Consistent self-attention for long-range image and video generation. *arXiv preprint arXiv:2405.01434*, 2024. [3](#), [5](#), [6](#), [7](#)
- [51] Junchen Zhu, Huan Yang, Huiguo He, Wenjing Wang, Zixi Tuo, Wen-Huang Cheng, Lianli Gao, Jingkuan Song, and Jianlong Fu. Moviefactory: Automatic movie creation from text using large generative models for language and images. In *Proceedings of the 31st ACM International Conference on Multimedia*, pages 9313–9319, 2023. [3](#)