

VisAgent: Narrative-Preserving Story Visualization Framework

Seungkwon Kim*
NAVER WEBTOON AI
Republic of Korea

GyuTae Park*
NAVER WEBTOON AI
Seoul National University
Republic of Korea

Sangyeon Kim
NAVER WEBTOON AI
Republic of Korea

Seung-Hun Nam†
NAVER WEBTOON AI
Republic of Korea

Abstract—Story visualization is the transformation of narrative elements into image sequences. While existing research has primarily focused on visual contextual coherence, the deeper narrative essence of stories often remains overlooked. This limitation hinders the practical application of these approaches, as generated images frequently fail to capture the intended meaning and nuances of the narrative fully. To address these challenges, we propose VisAgent, a training-free multi-agent framework designed to comprehend and visualize pivotal scenes within a given story. By considering story distillation, semantic consistency, and contextual coherence, VisAgent employs an agentic workflow. In this workflow, multiple specialized agents collaborate to: (i) refine layered prompts based on the narrative structure and (ii) seamlessly integrate generated elements, including refined prompts, scene elements, and subject placement, into the final image. The empirically validated effectiveness confirms the framework’s suitability for practical story visualization applications.

Index Terms—Narrative-preserving story visualization, Multi-agent framework, Story distillation, Semantic consistency

I. INTRODUCTION

Story visualization, the process of translating narrative content into visual representations, presents a significant challenge due to the need to convey the story’s essence to a diverse audience [1], [2]. Recent advancements in diffusion models (DMs), such as Stable Diffusion (SD) [3], and large language models (LLMs) like GPT-4 [4], have significantly enhanced story visualization capabilities. DMs facilitate the generation of coherent and high-quality scenes by effectively capturing complex conditional distributions. Meanwhile, LLMs complement this process by providing accurate text-to-visual alignment—termed *semantic consistency*, ensuring that generated images closely reflect the narrative context and character properties, such as appearance, actions, and interactions [5]–[7].

In the developing field of story visualization, early efforts were centered around generating isolated scenes utilizing DMs and manual prompts for text-to-image (T2I) generation [3], [8]. As the field progressed, the focus shifted towards more sophisticated multi-scene visualizations with training emphasizing character identity preservation—*contextual consistency* [9]–[13]. For effective multi-scene generation, precise foreground (FG) element placement within each background (BG) is crucial. Recent advancements in layout suggestion techniques, driven by language models with vision capabilities [14]–[16], have introduced innovative training-free approaches that enhance flexibility, implementation, and computational efficiency in visual storytelling [17]–[19]. Beyond fixed sequential prompts and story completion concepts [20], [21], emerging approaches focus on *story distillation*, a process that refines a given plain narrative into well-crafted prompts for multi-scene visualization [12], [13].

Through the analysis of existing literature, we have identified and separated two core components essential for narrative-preserving story visualization (Fig. 1). The primary objectives of these components are story distillation, semantic consistency, and contextual

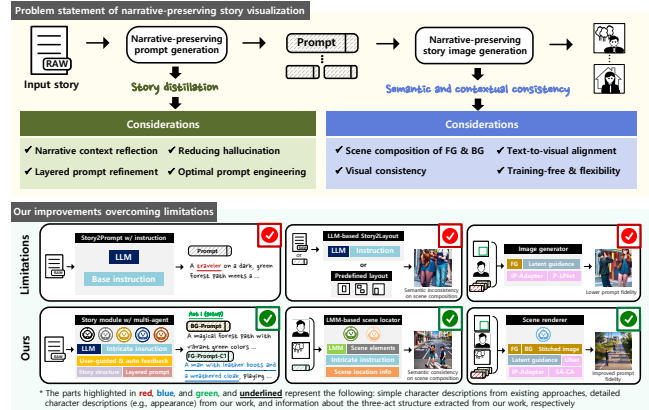


Fig. 1. Problem statement of narrative-preserving story visualization and our improvements for overcoming limitations.

consistency. Consideration in story distillation involves refining the detailed information of the scenes, including background (BG) and foreground (FG), into prompts suitable for diffusion models (DMs) to interpret, while maintaining the narrative context of the original story without inducing hallucinations. When generating multiple scenes with prompts, it is crucial to consider the compositional and stylistic harmony between FG and BG elements, while maintaining the benefits of a training-free approach and ensuring coherence across multiple scenes. Existing approaches often face limitations, as illustrated in the lower part of Fig. 1.

To overcome these limitations, we propose a multi-agent framework named **VisAgent** to enable narrative-preserving story visualization with exquisite perceptual quality. VisAgent comprises the following two core modules: *story module* and *image module*. The story module leverages LLM-based agents to comprehend the input story in a three-act structure format, recognize the narrative context, and refine the extracted context information into well-crafted layered prompts (e.g., FG and BG prompts). The resultant prompts are optimized for seamless interpretation by DMs, considering prompt engineering principles. In addition, to minimize hallucinations and ensure rationale verification, the framework incorporates user-guided and automated agents for feedback and reflection on intermediate and final outputs. The image module, which is composed of three agents, is designed to generate narrative-preserving story images from layered prompts. To achieve this, it requires high fidelity in the FG and BG images, along with the appropriate placement of FG subjects within the BG to align with the narrative context. In the initial stage, the module prepares scene elements (e.g., character and BG images) in the scene element generator agent. Subsequently, a scene locator agent, powered by a large multi-modal model (LMM), uses these elements to generate a layout for subject placement. Finally, the scene renderer agent integrates all element images and the location information to generate the final rendered image using the proposed

*: Equal contribution, †: Corresponding author

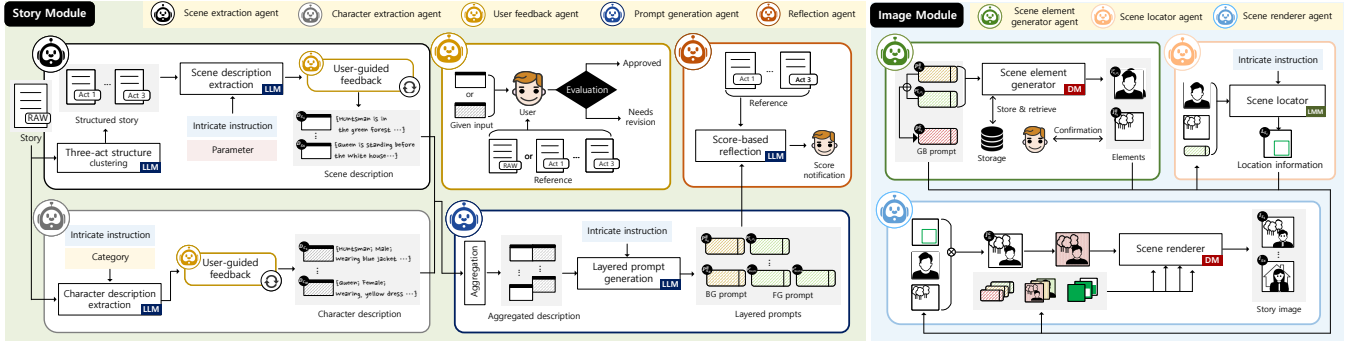


Fig. 2. Architectural overview of VisAgent, a multi-agent framework designed for story visualization. The recycle symbol (♻️) represents a process that repeats until approval. ⊕ refers to the process of generating a GB prompt by concatenating BG and FG prompts to each scene while ⊗ specifically denotes a segmentation-based image stitching process.

semantic-aware cross-attention layer. As emphasized in Fig. 1, our contributions are as follows.

- We propose a story module, a specialized multi-agent framework that refines a given story into layered prompts by analyzing its narrative structure and distilling key events and character attributes, achieving effective story distillation.
- We propose an image module comprising multiple agents that collaboratively generate narrative-consistent story images by separately generating FG and BG, determining subject placement, and rendering the final image while ensuring semantic coherence with a novel semantic-aware cross-attention layer.
- Our framework, designed with the interplay of multi-agents and user feedback, provides an application that enables users to efficiently visualize their own stories, realizing their creative idea.

II. PROPOSED FRAMEWORK

Our framework, VisAgent, comprises story and image modules. All agents within the framework are meticulously devised to collaborate, optimizing a multi-modal setting for story visualization.

A. Preliminaries

Let R denote the input story in plain text format, with N representing the number of scenes for story distillation and M denoting the number of characters featured in R . In this study, the elements corresponding to each scene and character are denoted by S_i and C_j , where $i \in \{0, 1, \dots, N-1\}$ and $j \in \{0, 1, \dots, M-1\}$. We define D_{S_i} as the scene description, D_{C_j} as the character description, $P_{S_i}^B$ as the BG prompt, P_{S_i, C_j}^F as the FG prompt for a specific character, and $P_{S_i}^G$ as the global (GB) prompt. The BG and FG character images denoted as $I_{S_i}^B$ and I_{S_i, C_j}^F respectively, are generated from $P_{S_i}^B$ and P_{S_i, C_j}^F . The location of the subject, defined by the coordinates $(x_{\min}, y_{\min}, x_{\max}, y_{\max})$, is represented as L_{S_i} , while the generated story image is denoted as I_{S_i} .

B. Story Module

1) *Narrative-Preserving Prompt Generation*: While LLMs have been extensively studied for their effectiveness in generating prompts tailored for DMs [13], [15], [17], [22], existing approaches in story visualization often prioritize image quality over the preservation of the narrative structure [15], [16]. This can result in generated images that fail to effectively convey the story’s intentions, potentially leading to incomplete or insufficient visualizations.

2) *Methodology*: To address this limitation, we develop the story module, a multi-agent framework that incorporates the human-in-the-loop concept and introduces a narrative-preserving concept for the first time in story visualization tasks. Our module analyzes narrative structures from plain text input and refines them into layered prompts for BG and FG, enabling story distillation. In this way, the refined output can preserve the key scenes of the story, while incorporating detailed descriptions for both BG and FG (see Table I)

In detail, the story module is defined as $SM(R, \epsilon_S, \delta, N)$, where ϵ_S and δ represent the devised intricate instruction and predefined category, respectively. In detail, the input story R is thoroughly analyzed and transformed into effective prompts for a generative model while preserving its essential narratives. Our approach is grounded in the assumption that all stories, to varying degrees, rely on specific narrative structures to enhance their appeal. Central to this approach is the *scene extraction agent*, which deconstructs and distills the story utilizing the storytelling structure (i.e., three-act structure [23]), a widely recognized and long-standing framework in storytelling. In addition, the *character extraction agent* identifies and extracts descriptions about all characters, including their attire, gender, and other attributes defined by the predefined category. If the attire is unspecified, the agent infers details from context, essential for defining character style in DM prompts.

Although these agents effectively identify crucial scenes and characters automatically, the *user feedback agent* allows users to confirm or modify parts of the outputs (D_{S_i} and D_{C_j}), ensuring the preservation of all critical information and reducing hallucinations throughout the distillation process. Based on the user-confirmed results, the *prompt generation agent* then generates separate BG and FG prompts ($P_{S_i}^B$ and P_{S_i, C_j}^F), each covering mutually exclusive aspects of the scene and having a format tailored to the targeted DM by applying prompt engineering principles. Finally, in addition to the LLM’s capability for textual similarity assessment, the *reflection agent* conducts a thorough review by comparing the prompts with the corresponding story segments, informing users of any potential deviations caused by the story module. The resulting layered prompts offer two advantages: (i) preserving the essential narrative elements of the original story, and (ii) aligning with the standard graphic narrative synthesis process for FG and BG composition, which are crucial components for the image module.

C. Image Module

1) *Narrative-Preserving Story Image Generation*: To generate a narrative-preserving story image, it is essential to create an image that achieves high prompt fidelity for both FG and BG elements and

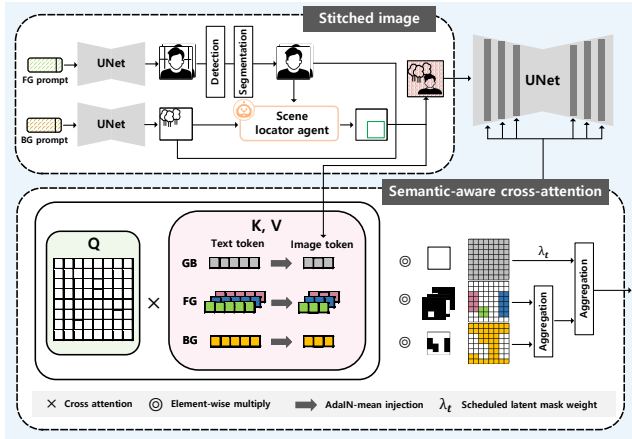


Fig. 3. Schematic of semantic-aware cross-attention layer.

ensures that FG subjects are appropriately positioned and integrated within the BG. Inspired by a common process of creating cartoons, we generate FG and BG elements separately, determine the layout for placing FG elements on the BG, and finally produce the complete story image. This approach preserves the quality of each scene element while considering its global context, resulting in a narrative-preserving story image.

2) *Methodology*: From this perspective, we propose the image module, composed of three specifically designed agents to generate story scene images (see Fig 2). The image module is defined as $IM(P_{S_i}^B, P_{S_i, C_j}^F, \epsilon_I)$, where ϵ_I represents the devised intricate instruction.

First, the *scene element generator agent* generates each element image (i.e., FG character images and a BG image) from refined layered prompts, which serve as key components in the composition of resultant images. Note that we adopt a strategy of utilizing an IP-Adapter [24] and a subject storage, similar to [18], [19], to manage consistency between element images (i.e., contextual consistency). Users may optionally repeat and confirm this process multiple times until the scene elements are satisfactorily produced.

Subsequently, the *scene locator agent*, powered by an LMM model, suggests suitable FG placement within the BG, considering the semantic context of layered prompts and the BG image and ensuring a harmonious scene composition. Unlike existing approaches [18], [19] that solely rely on prompts for scene layout determination, our agent leverages visual context from the scene element images to ensure accurate scene composition, thereby significantly enhancing semantic consistency and achieving narrative-preserving scene synthesis.

Finally, the *scene renderer agent* generates a story scene image using a stitched image as latent guidance input, along with the scene renderer, a DM that incorporates a modified version of the cross-attention layer. Initially, the agent generates a stitched image by meticulously integrating the FG with the BG, employing an open-vocabulary detection model [25] and a segmentation model [26] in a layer-stacking manner. Subsequently, the stitched image is encoded and passed through the forward diffusion process, and then used as a latent guidance for the input of the scene renderer. In particular, to augment semantic consistency, we propose a novel Semantic-Aware Cross-Attention (SA-CA) layer that comprehensively utilizes all layered prompts, scene elements (FG and BG), and latent guidance inputs (see Fig. 3). Although it is inspired by the parallel text and cross-attention mechanisms of P-UNet [19], the following key differences are introduced to further enhance semantic consistency.

TABLE I

RESULTS ON NARRATIVE-PRESERVING PROMPT GENERATION

Method	Scene	Refined prompt
GPT-4o w/ base ins. [13]	#1	A young boy, Jack, standing at the base of a towering beanstalk reaching into a mystical sky, surrounded by lush, vibrant greenery
	#2	A massive, golden giant's castle in the clouds, with Jack peering in through a grand window, his face full of awe and trepidation
	#3	Jack climbing a colossal, twisting beanstalk against rolling clouds and a shimmering sunset, filled with adventure and wonder
	#4	An enchanted giant's treasure room filled with golden coins, shimmering jewels, and magical artifacts, with Jack cautiously exploring
	#5	A dramatic moment of Jack slaying a fierce giant, with intense action and dynamic lighting highlighting the epic confrontation
Story module	BG:	A humble rustic and weathered traditional market, no building, surrounded by a few sparse trees and a patchy garden, highres detailed, soft lighting, daytime
	Act 1	FG-C1: A small boy with worn-out blue medieval clothing, standing, handing over a cow
	Act 2	FG-C2: A old man with worn-out medieval merchant clothing, standing, holding a basket of magical beans
	Act 1	FG-C1: A towering, fantastically gigantic beanstalk spiraling up into the sky, disappearing into the clouds, a blue sky with wisps of clouds, highres, detailed, soft lighting, daytime
	Act 2	FG-C2: A small boy with worn-out blue medieval clothing, climbing, holding onto the gigantic beanstalk
	Act 3	FG-C3: An old medieval cottage surrounded by trees, clear sky, highres, detailed, soft lighting, daytime
	Act 2	FG-C1: A small boy with worn-out blue medieval clothing, standing, discovering a mysterious cottage
	Act 2	FG-C1: An antique interior of mysterious medieval cottage, furniture, beams, highres, detailed, soft lighting, daytime
	Act 3	FG-C2: A small boy with worn-out blue medieval clothing, holding, standing next to a gigantic beanstalk
	FG-C3: A giant human monster with muscle, beard, big nose and with black village clothing, sitting, closed eyes	
	FG-C1: A towering, fantastically gigantic chunk of beanstalk, meadow, grass, clouds, trees, highres, detailed, soft lighting, daytime	
	FG-C2: A small boy with worn-out blue medieval clothing, holding, standing next to a gigantic beanstalk	
	FG-C3: A giant human monster with muscle, beard, big nose and with black village clothing, falling from sky, floating in the air	

* Here, Acts 1, 2, and 3 correspond to a three-act structure's setup, conflict, and resolution. C1 to C3 each represent the index of a character. In the story module, scenes such as exchanging the cow for beans (#1) and stealing treasures (#3) are highlighted because of their critical narrative significance. In contrast, the baseline focuses merely on major events, often disconnected from the overall narrative.

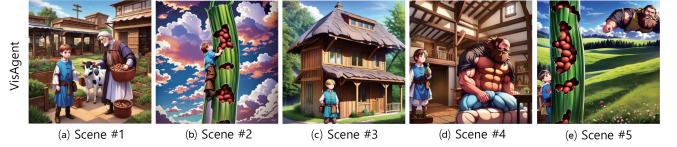


Fig. 4. Results of qualitative evaluation: narrative-preserving story visualization results of VisAgent using refined layered prompts as listed in Table I.

- *Using BG and stitched image as reference and guidance input*: Our baseline [19] employs FG images as reference image tokens in the modified cross-attention layer. The SA-CA layer extends this by leveraging FG, BG, and a stitched image as reference tokens for their respective regions, to generate a semantically coherent scene. Besides, we use the stitched image as latent guidance instead of a black BG with the pasted FG image. As FG and BG are generated independently earlier, the SA-CA layer is expected to enhance prompt fidelity over the baseline, particularly for the BG.
- *Global latent aggregation for global region*: Cross-attention is applied separately to the FG and BG regions using their respective layered text and image prompt tokens, as well as across the entire region using a global text (combining FG and BG) prompt token and a global image prompt token derived from the stitched image. To aggregate all latents, we introduce λ_t adjusting the influence of a global latent during scene generation. We adopt a stepwise strategy where λ_t progressively increases over the timesteps.
- *Token alignment between text and image prompt token*: Inspired by [27], we apply an adaptive mean normalization (AdaIN-mean) operation between text and image prompt tokens to enhance semantic consistency while maintaining identity fidelity, as shown by the thick gray arrow in Fig. 3. For FG, BG and GB latent, the key and value from the image tokens are aligned with those from the text tokens before applying cross-attention.

III. EXPERIMENTS

A. Settings

1) *Implementation Details*: We adopt GPT-4o, a variant based on GPT-4 [4], as LLM and LMM. The multi-agent system in the story module is implemented based on the LangGraph framework [28]. For the story module, plain text about *<Jack and the Beanstalk>* is employed as input story¹, featuring three characters and narrative key moments, and the N for story distillation is set to 5. For the image module, we adopt a model based on SD v1.5, as specified in [19], using a DDIM sampler with 30 steps. λ_t is set to 0.1 for Steps

¹<https://americanliterature.com/childrens-stories/jack-and-the-beanstalk>



Fig. 5. Results of qualitative evaluation: performance analysis of five example prompts-based story visualizations compared to the baseline.

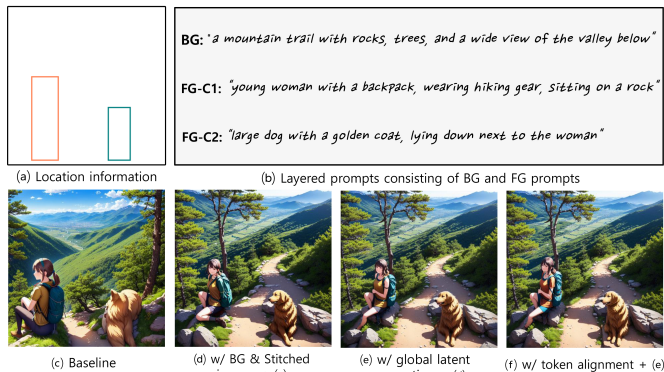


Fig. 6. Results of ablation study on SA-CA of image module.

1–10, 0.3 for Steps 11–20, and 0.5 for Steps 21–30, incrementing at intervals of ten steps. The evaluations are completed with one NVIDIA A100 GPU with 80 GB of GPU memory.

2) *Evaluations*: The story and image modules constituting the VisAgent are evaluated via quantitative and qualitative evaluations. Referencing [18], [29], we select the quantitative metrics of Fréchet inception distance (FID) and character–character similarity (CCS) to evaluate contextual consistency and text–image similarity (TIS) to assess semantic consistency. To quantitatively evaluate the scene renderer agent, a new VisAgent benchmark is introduced, comprising 400 narrative stories generated by an LLM. Based on a format similar to CMIGBench [18], it includes richer FG and BG prompts to assess the ability to generate detailed and narrative scene content. In addition, we exhibit the results for qualitative evaluation of the distilled story and the visualized scene image, and conduct a human preference study via an A/B test with 20 volunteers.

3) *Baselines*: To evaluate the story distillation capabilities of our story module, we compared its performance to the baseline approach (i.e., GPT-4o with base instruction) specified in [13]. We compare our image module with the state-of-the-art training-free method, AutoStudio [19], from the perspective of visualization performance.

B. Experimental Results

1) *Results of Story Module*: Our story module effectively captures key narrative elements, including build-up scenes leading to the climax and connecting scenes between major events (see Table I). In contrast, the baseline approach often produces an irregular distribution of scenes, neglecting narrative flow and focusing solely on major events. Our module also ensures consistent character styles (e.g., appearance) across prompts and provides refined layered prompts, enhancing fidelity and narrative representation in the visualization

TABLE II
RESULTS OF QUANTITATIVE EVALUATION

Method	VisAgent Benchmark			CMIGBench		
	TIS (%) ↑	FID ↓	CCS (%) ↑	TIS (%) ↑	FID ↓	CCS (%) ↑
AutoStudio	25.04	267.48	83.42	31.90 [†]	246.85 [†]	79.02 [†]
VisAgent	25.43	263.74	83.76	32.58	243.11	80.14

[†] denotes results reproduced by our implementation using the official code.

TABLE III
RESULTS OF HUMAN EVALUATION (%)

Metric	Story distillation		Semantic and contextual consistency	
	GPT-4o w/ [13]	VisAgent	AutoStudio [19]	VisAgent
User score	17.78	82.22	15.56	84.44

process. User evaluations, assessing the distillation quality of the prompts listed in Table I and their suitability for DMs, further confirm the effectiveness of our module (see Table III).

2) *Results of Image Module*: First, we perform a qualitative evaluation on visualized resultant images. As shown in Fig. 4, the image module achieves narrative-preserving visualization with layered prompts refined from the plain story. In a comparison with [19], our module demonstrated superior quality across five example scenes as illustrated in Fig. 5). This achievement stems from leveraging our scene locator agent, which considers semantic composition and other agents enhancing semantic consistency and contextual consistency with the layered prompts.

Second, Table II presents the quantitative results of the scene renderer agent which has a SA-CA layer, compared to the baseline, using FID, CCS, and TIS metrics, demonstrating improved performance across both benchmarks. Note that the scene locator is not used in this experiment to ensure a fair comparison, and the metrics of the baseline were reproduced using the official code in our setting.

Moreover, Fig. 6 presents output images illustrating the application of a several strategies to the baseline, resulting in a SA-CA layer. Specifically, Fig. 6(a)–(b) depict the input components, while Fig. 6(d)–(f) show output images generated by the SA-CA with submodules specified in the subcaptions. Compared to the baseline, Fig. 6(d) demonstrates improved prompt fidelity for the BG prompt, “a mountain trail,” and produces a more natural image of FG-C2. By integrating the global latent with λ_t , Fig. 6(e) shows contextually improved image, such as FG-C1’s “sitting on a rock,” with a more natural pose. As a final step, Fig. 6(f) demonstrates slightly improved quality through the incorporation of token alignment.

Finally, we conduct a human preference study to provide an intuitive evaluation of our work compared to the baseline [19]. To this end, we generated story scene images with prompts and exposed them with guidelines (i.e., the criteria for measuring contextual consistency include prompt fidelity and the scene composition between the BG and FG), requesting volunteers to perform an A/B test. The results presented in Table III indicate the superiority of the proposed module.

IV. CONCLUSION

This study introduced VisAgent, a multi-agent framework designed for narrative-preserving story visualization. The framework’s story module effectively refines the input story through a sophisticated multi-agentic workflow, producing layered prompts that accurately capture the narrative context while minimizing hallucinations. The image module, composed of scene element generator, scene locator and scene renderer agents, generate narrative-preserving story image by effectively integrating scene elements. Experimental results validate the effectiveness of VisAgent.

REFERENCES

- [1] E. Segel and J. Heer, "Narrative visualization: Telling stories with data," *IEEE transactions on visualization and computer graphics*, vol. 16, no. 6, pp. 1139–1148, 2010.
- [2] Y.-Z. Song, Z. Rui Tam, H.-J. Chen, H.-H. Lu, and H.-H. Shuai, "Character-preserving coherent story visualization," in *European Conference on Computer Vision*. Springer, 2020, pp. 18–33.
- [3] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2022, pp. 10 684–10 695.
- [4] OpenAI, "Gpt-4 technical report," 2023, accessed: 2023-03-31. [Online]. Available: <https://cdn.openai.com/papers/gpt-4.pdf>
- [5] S. Kim, S. Kim, and S.-H. Nam, "A framework for portrait stylization with skin-tone awareness and nudity identification," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 3660–3664.
- [6] T. Guo, X. Chen, Y. Wang, R. Chang, S. Pei, N. V. Chawla, O. Wiest, and X. Zhang, "Large language model based multi-agents: A survey of progress and challenges," *arXiv preprint arXiv:2402.01680*, 2024.
- [7] L. Zhang and M. Agrawala, "Adding conditional control to text-to-image diffusion models," *arXiv preprint arXiv:2302.05543*, 2023.
- [8] D. Podell, Z. English, K. Lacey, A. Blattmann, T. Dockhorn, J. Müller, J. Penna, and R. Rombach, "Sdxl: Improving latent diffusion models for high-resolution image synthesis," *arXiv preprint arXiv:2307.01952*, 2023.
- [9] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen, "Lora: Low-rank adaptation of large language models," *arXiv preprint arXiv:2106.09685*, 2021.
- [10] Q. Wang, X. Bai, H. Wang, Z. Qin, and A. Chen, "Instantid: Zero-shot identity-preserving generation in seconds," *arXiv preprint arXiv:2401.07519*, 2024.
- [11] Z. Li, M. Cao, X. Wang, Z. Qi, M.-M. Cheng, and Y. Shan, "Photomaker: Customizing realistic human photos via stacked id embedding," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 8640–8650.
- [12] W. Wang, C. Zhao, H. Chen, Z. Chen, K. Zheng, and C. Shen, "Autostory: Generating diverse storytelling images with minimal human effort," *arXiv preprint arXiv:2311.11243*, 2023.
- [13] Y. Gong, Y. Pang, X. Cun, M. Xia, Y. He, H. Chen, L. Wang, Y. Zhang, X. Wang, Y. Shan *et al.*, "Talecrafter: Interactive story visualization with multiple characters," *arXiv preprint arXiv:2305.18247*, 2023.
- [14] Y. Li, H. Liu, Q. Wu, F. Mu, J. Yang, J. Gao, C. Li, and Y. J. Lee, "Gligen: Open-set grounded text-to-image generation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 22 511–22 521.
- [15] L. Lian, B. Li, A. Yala, and T. Darrell, "Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models," *arXiv preprint arXiv:2305.13655*, 2023.
- [16] C. Liu, H. Wu, Y. Zhong, X. Zhang, Y. Wang, and W. Xie, "Intelligent grimm-open-ended visual storytelling via latent diffusion models," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2024, pp. 6190–6200.
- [17] Y. Zhou, D. Zhou, M.-M. Cheng, J. Feng, and Q. Hou, "Storydiffusion: Consistent self-attention for long-range image and video generation," *arXiv preprint arXiv:2405.01434*, 2024.
- [18] J. Cheng, B. Yin, K. Cai, M. Huang, H. Li, Y. He, X. Lu, Y. Li, Y. Li, Y. Cheng *et al.*, "Theatergen: Character management with llm for consistent multi-turn image generation," *arXiv preprint arXiv:2404.18919*, 2024.
- [19] J. Cheng, X. Lu, H. Li, K. L. Zai, B. Yin, Y. Cheng, Y. Yan, and X. Liang, "Autostudio: Crafting consistent subjects in multi-turn interactive image generation," *arXiv preprint arXiv:2406.01388*, 2024.
- [20] X. Pan, P. Qin, Y. Li, H. Xue, and W. Chen, "Synthesizing coherent story with auto-regressive latent diffusion models," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2024, pp. 2920–2930.
- [21] S. Yang, Y. Ge, Y. Li, Y. Chen, Y. Ge, Y. Shan, and Y. Chen, "Seed-story: Multimodal long story generation with large language model," *arXiv preprint arXiv:2407.08683*, 2024.
- [22] J. Qin, J. Wu, W. Chen, Y. Ren, H. Li, H. Wu, X. Xiao, R. Wang, and S. Wen, "Diffusiongpt: Llm-driven text-to-image generation system," *arXiv preprint arXiv:2401.10061*, 2024.
- [23] P. Papalampidi, F. Keller, L. Frermann, and M. Lapata, "Screen-play summarization using latent narrative structure," *arXiv preprint arXiv:2004.12727*, 2020.
- [24] H. Ye, J. Zhang, S. Liu, X. Han, and W. Yang, "Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models," 2023.
- [25] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu *et al.*, "Grounding dino: Marrying dino with grounded pre-training for open-set object detection," *arXiv preprint arXiv:2303.05499*, 2023.
- [26] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick, "Segment anything," *arXiv:2304.02643*, 2023.
- [27] Y. Wu, Z. Li, H. Zheng, C. Wang, and B. Li, "Infinite-id: Identity-preserved personalization via id-semantics decoupling paradigm," *arXiv preprint arXiv:2403.11781*, 2024.
- [28] LangGraph, "Langgraph," <https://python.langchain.com/docs/langgraph/>, 2023, accessed: 2023-04-25.
- [29] J. Hessel, A. Holtzman, M. Forbes, R. L. Bras, and Y. Choi, "Clipscore: A reference-free evaluation metric for image captioning," *arXiv preprint arXiv:2104.08718*, 2021.