

LLMs Behind the Scenes: Enabling Narrative Scene Illustration

Anonymous ACL submission

Abstract

Generative AI has established the opportunity to readily transform content from one medium to another. This capability is especially powerful for storytelling, where visual illustrations can illuminate a story originally expressed in text. In this paper, we focus on the task of narrative scene illustration, which involves automatically generating an image depicting a scene in a story. Motivated by recent progress on text-to-image models, we consider a pipeline that uses LLMs as an interface for prompting text-to-image models to generate scene illustrations given raw story text. We apply variations of this pipeline to a prominent story corpus in order to synthesize illustrations for scenes in these stories. We conduct a human annotation task to obtain pairwise quality judgments for these illustrations. The outcome of this process is the SCENEILLUSTRATIONS dataset, which we release as a new resource for future work on cross-modal narrative transformation. Through our analysis of this dataset and experiments modeling illustration quality, we demonstrate that LLMs can effectively verbalize scene knowledge implicitly evoked by story text. Moreover, this capability is impactful for generating and evaluating illustrations.

1 Introduction

Observing the transformation of a story from one modality to another (e.g. from text to visual form) can make the story more compelling to its audience. Recent advances in generative AI have enabled this kind of cross-modal transformation to be performed automatically. In particular, text-to-image models allow people to create visual material using natural language alone. Current interaction with these models typically involves users envisioning a particular visual target and then crafting language that realizes that target. Many stories that currently only exist in text form would be well-suited for transfer to an image modality, but the text itself of

these stories may not be naturally optimal for directly applying text-to-image models. Given their demonstrated success at meta-prompting (e.g. Zhou et al., 2023), large language models (LLMs) may be able to interface with story text to synthesize suitable prompts for text-to-image models towards this end. The cooperation between these AI models would make it possible to automatically generate illustrations for any given text-based story.

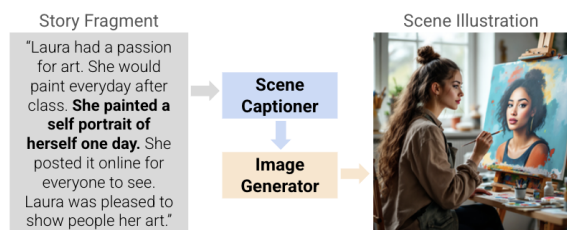


Figure 1: Overview of scene illustration pipeline

In this paper, we exemplify this approach to visual transfer of story text. Generating illustrations for stories, a task that has been termed *story visualization*, encompasses a myriad of challenges. Some of these challenges pertain to modeling the relation between the story text and illustrations (text-image alignment), while others pertain to the relation between illustrations for different scenes in the story (image-image alignment). Existing story visualization research (e.g. Li et al., 2019) has largely focused on image-image alignment, in particular the problem of ensuring visual consistency between depictions of story elements like characters and settings. We aim to bring more research attention to issues of text-image alignment in this domain. Thus, our work is scoped to focus on individual *scene illustrations*. In particular, we consider scene-level units of stories (*fragments*). We present a pipeline (outlined in Figure 1) that generates a scene illustration given a fragment in its story context. Through systematic variation and ablation of the components of this pipeline, we pro-

duce a novel set of scene illustrations for fragments in a notable story corpus. We then conduct a human annotation task to obtain relative quality judgments for pairs of illustrations. We refer to the resulting quality-annotated items as the SCENEILLUSTRATIONS dataset.

We leverage the SCENEILLUSTRATIONS dataset to demonstrate that LLMs can explicate visual knowledge of narrative scenes by inferring this knowledge directly from story text, without any visual input. We establish this capability through two findings. First, we show that LLMs are an effective interface for transforming story text into prompts that facilitate text-to-image models to produce illustrations. Second, we show that LLMs can verbalize scene characteristics in a way that is useful for evaluating the quality of illustrations. In particular, we demonstrate an approach to predicting human-favored illustrations among pairs in our dataset, through which we apply LLM-specified scene characteristics as evaluation criteria for scoring illustrations. The success of this approach relative to a criteria-ablated baseline further suggests the utility of LLMs for explicating scene knowledge that is implicitly conveyed by story text.

Contributions This paper makes the following contributions¹:

- We define and motivate the task of narrative scene illustration in relation to existing research on visually aligned storytelling.
- We demonstrate a pipeline for producing scene illustrations for any given story text. The pipeline components are fully interchangeable and can be used with any LLM and text-to-image models.
- We apply our pipeline to synthesize scene illustrations for existing stories and elicit human quality annotations for pairs of these illustrations, resulting in the newly created SCENEILLUSTRATIONS dataset.
- Through analysis of the quality annotations in SCENEILLUSTRATIONS, we show that LLMs are an effective interface between story text and text-to-image models in facilitating scene illustration.
- We assess an approach to predicting these quality annotations that involves applying LLM verbalizations of scene characteristics as evaluation criteria. We discuss the evaluation results as additional evidence that LLMs can explicate visual scene knowledge inferred from story text.

¹The SCENEILLUSTRATIONS dataset and all code for our experiments is available at: [withheld/during/review](#)

2 Background and Related Work

Image-Aligned Story Data Datasets that pair story text with images have emerged from research on visually grounded story generation, which involves writing a story given a sequence of images. Human authors have performed this task for existing media-sourced images (Halperin and Lukin, 2023; Huang et al., 2016; Hong et al., 2023). For the reverse-direction task of story visualization, which involves generating a sequence of images given story text, some research has leveraged videos for data creation (Li et al., 2019; Tao et al., 2024). Distinct frames of video are sampled as static images, while crowdsourced descriptions of frames are designated as the story text (Li et al., 2019; Maharana and Bansal, 2021; Maharana et al., 2022). A key design factor of all the above datasets is that the story text is authored specifically in response to the images, rather than originating in text form. We explore an alternative process for visually aligning narratives by synthesizing images for existing text-based stories.

Multimodal Storytelling Systems In addition to datasets, there are increasing demonstrations of story visualization systems, as well as systems that generate story text and images in parallel, i.e. multimodal story generation (An et al., 2024; Koh et al., 2023; Singh et al., 2023; Wan et al., 2024; Yang et al., 2024). While some models applied to these use cases have been trained end-to-end on the specialized datasets described above (Feng et al., 2023; Maharana and Bansal, 2021; Tao et al., 2024), researchers have also begun to leverage generically pretrained models to expand the scope of these systems to open-domain storytelling (de Lima et al., 2024; Gong et al., 2023; Soumik Rakshit, 2024). We follow suit in leveraging a plug-and-play pipeline for scene illustration.

Meta-Prompting for Text-to-Image Models

One challenge with using generic models for story visualization is that the story text itself is not necessarily an optimal prompt for text-to-image models. In particular, this text tends to lack detailed visual descriptions (e.g. the physical appearance of story elements like entities and locations), which are considered essential when providing instructions to text-to-image models (Maharana et al., 2022). Users of these models who have become skilled in writing prompts have done so largely through an iterative process of observing what prompt language

yields desirable images (Don-Yehiya et al., 2023). Even with this skill, significant effort is required to manually compose a prompt that captures the intended visual features of the scene corresponding to a story fragment. Following the paradigm of meta-prompting (e.g. Zhou et al., 2023), there is a variety of research on automated prompt optimization for text-to-image models (Brade et al., 2023; Feng et al., 2024; Hao et al., 2023; Wang et al., 2024), some of which establishes the effectiveness of LLMs in facilitating this process (Lian et al., 2024). Accordingly, recent story visualization work has used LLMs as an interface for deriving text-to-image prompts from story text. In particular, Gong et al. (2023) and He et al. (2024) instructed GPT-4 to transform a story into a series of scene-level prompts intended as input to text-to-image models. It is presumed that these synthesized prompts are more visually descriptive than the story text and thus produce better images, but this has not been empirically validated. Thus, we address this opportunity in our work.

LLMs for Image Evaluation Assessing the degree of semantic alignment between images and text is a prominent research endeavor, which has primarily involved measuring their similarity when projected into a shared embedding space (e.g. Hessel et al., 2021). Because of their capacity for visually descriptive language, even unimodal (text-only) LLMs can contribute to this endeavor. Several works have demonstrated the utility of LLMs for zero-shot visual recognition tasks (Li et al., 2023; Maniparambil et al., 2023; Menon and Vondrick, 2023; Pratt et al., 2023). This line of research has recently extended to eliciting visual knowledge from LLMs as a strategy for text-to-image evaluation (Lin et al., 2025; Lu et al., 2023; Hu et al., 2023). Encouraged by recent demonstrations of LLM-based evaluation in multimodal story generation (An et al., 2024), we pursue this method for evaluating scene illustrations.

Criteria-based Evaluation with LLMs In NLP, criteria is a means of anchoring evaluation to certain objectives (Yuan et al., 2024). With the rapidly expanding LLM-as-a-judge paradigm, this has evolved to the point where LLMs are not just applying human-authored criteria to assess text, but are also generating their own criteria (Cook et al., 2024). We examine LLMs’ capacity to generate evaluation criteria for the scene illustration task.

3 Scene Illustration Pipeline

We first outline the high-level components² of the illustration pipeline in this section, before describing their application in the next section.

Story Fragmentation In our work, we consider a *scene* to be an abstract unit of a story that can be distinctly illustrated by a single image. The story text that aligns to a scene is referred to as a *fragment*. Thus, the first step of producing a scene illustration is to identify its source fragment. Recent work has validated the use of LLMs for the related task of segmenting events in narrative text (Michelmann et al., 2025). Accordingly, we utilize an LLM for this fragmentation task, by instructing it to explicitly annotate the boundaries of all fragments in a given story. Table A.13 shows the prompt we provide to the LLM to facilitate this, where the input contains the story text and the LLM is expected to generate the same text with brackets demarcating the left and right boundaries of each fragment, as demonstrated by the exemplars. We parse this output with a simple regular expression to gather the list of fragments.

Scene Descriptions Once a fragment is identified, the fragment with its story context can then be mapped to a *scene description*. A scene description is a verbalization of what should be illustrated in the image corresponding to the fragment. This text serves as the input to the text-to-image model used to produce the scene illustration. As described below in §4, we consider different types of scene descriptions in order to evaluate the capability of LLMs to generate these descriptions.

Image Generation As mentioned, the scene descriptions are the inputs to a text-to-image model, referred to here as an *image generator*. While we use the term ‘illustration’ to describe the end-to-end process that yields an image depicting a scene, the output of this process (i.e. the image generator output) is also called an *illustration*.

4 SCENEILLUSTRATIONS Dataset

Each item in the SCENEILLUSTRATIONS dataset consists of a fragment with its story context, along with two illustrations depicting the fragment. The illustrations vary based on their scene description

²We ran all model components using APIs, which we specify here for each model. Unless otherwise indicated, we used the default inference parameters defined by the model’s API.

and/or the image generator used to produce them. The dataset consists of 2995 items in total, which were created in two phases: Phase 1 yielded 1777 items and Phase 2 yielded 1218 items. In this section, we detail the pipeline for synthesizing and annotating the items in both phases, then we present analyses of the annotation results.

4.1 Story Text

Seeking out a story corpus suitable for the scene illustration task, we ultimately selected the well-studied ROCStories corpus (Mostafazadeh et al., 2016) based on some key considerations. In particular, these English-language stories were authored to adhere to basic narrative structure in a tightly length-constrained format. In particular, each story consists of five sentences conveying “a causally (logically) linked set of events involving some shared characters”. Thus, we can expect that stories are composed of distinct fragments that are each appropriately visualized as a scene illustration. Moreover, the stories are narrations of everyday experiences that can be interpreted according to commonsense knowledge. This knowledge is general enough it is likely to be familiar to the model components of our illustration pipeline.

4.2 Phase 1 Pipeline Details

We applied the pipeline outlined in §3 to produce an initial set of scene illustrations, which we refer to as *Phase 1* data. As inputs to the pipeline, we used the first 50 stories in the ROCStories dev set.³

Fragmentation We divided these stories into fragments as described in §3, using CLAUDE-3.5⁴ as the LLM, which has displayed notable storytelling-related capabilities (e.g. Mazur, 2025). As shown in Table A.6, this resulted in 206 total fragments across all 50 stories, an average of 4.12 per story. §A.1.1 presents some additional analysis.

Scene Descriptions We applied an LLM to transform a fragment alongside its story context into a scene description, using the prompt in Table A.15 with CLAUDE-3.5 as the LLM. We employ the term *scene captioner* to refer to an LLM’s role when

running this prompt, and we refer to the outputs as CAPTION scene descriptions. As outlined in Table 1, CAPTION is one of four scene description types we consider for Phase 1. We compare CAPTION with baseline scene descriptions composed of the raw story text. In the first baseline case, NC-FRAGMENT (i.e. no context fragment), we use the original fragment isolated from its story context as a scene description. The obvious limitation of NC-FRAGMENT is that the ablated context may be necessary for understanding certain information in the fragment (for example, a fragment might use a pronoun whose referent is only specified in the context). Thus, we considered two additional baseline scene descriptions that account for the story context, referred to as VC-FRAGMENT and SC-FRAGMENT. As Table 1 shows, VC-FRAGMENT (i.e. verbose context) inserts the full story text into the scene description, which is formatted as an instruction to consider this context when illustrating the fragment. Alternatively, SC-FRAGMENT (i.e. succinct context) is a rewritten version of the fragment where references to information in the story context are made explicit, enabling the fragment to be understood independently of the context. We prompt an LLM (also CLAUDE-3.5) to do this rewriting task, using the prompt in Table A.14. Table 1 gives examples of these different scene descriptions, with additional examples in Table A.16.

Image Generation We then applied two image generators⁵ to generate images using the scene descriptions as prompts. In particular, we used Midjourney v6.1, denoted here as MJ-6.1 (Midjourney, 2024), and FLUX-1[pro], denoted here as FLUX-1-PRO (Black Forest Labs, 2024a). We selected these image generators because they topped the *Artificial Analysis Image Arena Leaderboard* at the time of Phase 1 in August 2024. This leaderboard captures the relative ELO score (Boubdir et al., 2023) of text-to-image models based on pairwise human judgments regarding how well images from different models reflect the input prompt. Table A.16 includes examples of generated illustrations.

4.3 Phase 1 Annotation Task

Illustration Pairs Our primary objective for Phase 1 was to assess the effectiveness of the LLM scene captioner (i.e. CAPTION scene descriptions) in generating illustrations relative to generating

³The dev and test items in ROCStories are actually designated as the Story Cloze Test, where items have a specific format: each story consists of four sentences plus two alternative fifth sentences, where one is the ‘correct’ story ending and the other is the ‘incorrect’ ending. For each item, we discarded the incorrect ending and appended the correct ending after the initial four sentences to form a single five-sentence story.

⁴claude-3-5-sonnet-20240620, ran via the Anthropic API

⁵With exception to Midjourney, we ran all image generation models via the Replicate API.

Type	Format	Example
NC-FRAGMENT	{{fragment}}	"Alice called her mother and apologized profusely."
VC-FRAGMENT	"Consider this story: [{{story}}] Based on this context, illustrate this fragment of the story: [{{fragment}}]"	"Consider this story: [Alice was getting married in a few weeks. One night, her mother called and she forgot to call her back. Her mother left an angry message on her phone. She threatened not to come to the wedding. Alice called her mother and apologized profusely.] Based on this context, illustrate this fragment of the story: [Alice called her mother and apologized profusely.]"
SC-FRAGMENT	LLM output of fragment rewriting prompt (Table A.14)	"The bride-to-be called her mother and apologized profusely for forgetting to return her call and for the resulting angry message threatening not to attend the wedding."
CAPTION	LLM output of scene captioning prompt (Table A.15)	"A young woman with a worried expression sits on a couch, holding a phone to her ear. She's gesticulating with her free hand, appearing to speak emphatically. In the background, a wedding dress can be seen hanging on a closet door. The room is dimly lit, suggesting it's evening, and there's a notepad with wedding plans visible on a nearby coffee table."

Table 1: Types of scene descriptions for Phase 1

them directly from the raw story text. To address this, we randomly sampled pairs of illustrations each belonging to the same story fragment (across 206 possible fragments), where one illustration used CAPTION as the scene description, while the other used one of the baseline scene descriptions, NC-FRAGMENT, VC-FRAGMENT, or SC-FRAGMENT. This sampling resulted in some pairs where the illustrations used the same image generator and others that used different image generators. Ultimately there were 1777 illustration pairs. The exact distribution of types for the illustrations and resulting pairs is specified in A.1.2 and A.1.3.

Task Design We designed an annotation task to assess the relative quality of the two illustrations in each pair. In judging a pair, human annotators were shown the full story with the target fragment for that scene underlined, along with the two alternative images. As shown in Figure A.3, annotators were instructed to select the image that was "the better visualization of the underlined fragment". Annotators could express uncertainty by selecting "I can't decide which is better". We implemented the UI for this task using POTATO (Pei et al., 2022).

Procedure We deployed the task on Prolific to obtain annotators. English proficiency was the only requirement for participation. We sought 2 annotators to judge each illustration pair. Each participant judged between 33 and 74 pairs (median=47), plus 3 "attention check" items where one illustration in the pair was replaced with one for a different story, making it trivially easy which image to select. Participants were paid \$6 for an expected completion time of 30 minutes. We filtered out participants who did not pass all of the attention check items. Ultimately, 75 (out of 80) participants

passed the attention checks. This resulted in a total of 3554 responses for the 1777 pairs, where each item received a response from 2 annotators.

4.4 Phase 1 Annotation Results

Inter-annotator Agreement Given the annotated pairs resulting from §4.3, we computed the inter-annotator agreement of which illustration was selected as the better one in each pair. We did this using an *uncertainty-weighted* variation of Cohen's Kappa score (Cohen, 1960), which we abbreviate here as κ_u . This variation considers that response disagreements arising from one annotator expressing uncertainty (i.e. selecting "I can't decide") should be weighted half as much as disagreements where the two annotators each select a different illustration as better. As indicated in Table A.8, the overall κ_u for all 1777 items was 0.436, which can be classified as moderate agreement (Landis and Koch, 1977). Annotators agreed in their responses for 62.3% of items. §A.1.3 provides a finer-grained analysis of agreement for different types of pairs.

Win Rates for Scene Description Types To determine whether using an LLM as a scene captioner helps illustration quality, we counted how often the favored illustration was associated with each scene description type, i.e. each type's *win rate*. Table 2 shows the win rate for CAPTION illustrations when alternatively paired with NC-FRAGMENT, VC-FRAGMENT, and SC-FRAGMENT illustrations. This win rate is represented as the percentage of responses in which annotators selected the CAPTION illustration as better among all responses for each respective set of pairs. In all three cases, the

CAPTION is significantly⁶ better: it has an overall win rate of $\approx 78\%$ against NC-FRAGMENT, $\approx 75\%$ against VC-FRAGMENT, and $\approx 73\%$ against SC-FRAGMENT. Table A.17 further examines the win rates for pairs that used the same image generator, verifying that CAPTION is equally favorable regardless of which image generator was used. This validates the importance of the scene captioner in the pipeline: the resulting verbalization enables the image generator to better depict how a story fragment should be visually illustrated as a scene.

Scene Description Pair	# Pairs	CAPTION %
CAPTION vs. NC-FRAGMENT	680	78.1
CAPTION vs. VC-FRAGMENT	384	74.7
CAPTION vs. SC-FRAGMENT	393	72.5

Table 2: Win rates of CAPTION over the baseline scene descriptions in Phase 1

Win Rates for Image Generators While we focus primarily on how scene descriptions affect illustration quality, we also considered whether there were quality differences based on which image generator was used. These results are given in §A.1.4.

4.5 Phase 2 Motivation and Design

After observing that the CAPTION scene descriptions significantly contribute to illustration quality, we wanted to compare the impact of different LLMs as scene captioners. Phase 1 only considered CLAUDE-3.5. In Phase 2, we included other LLMs with storytelling-relevant capabilities (e.g. Tian et al., 2024): GPT-4o⁷ (OpenAI et al., 2024) and LLAMA-3.1⁸ (Grattafiori et al., 2024). We used the same captioning prompt from §4 (Table A.15).

We expanded the Phase 2 data to include a larger set of fragments compared with those of Phase 1. We randomly sampled 1000 stories from the ROCStories dev set, split them into fragments using the same method from Phase 1 (CLAUDE-3.5 with the Table A.13 prompt), then randomly selected one fragment per story for inclusion in the dataset.

We also considered a larger set of image generators in Phase 2. Based on the state of the Artificial Analysis Leaderboard in November 2024, we selected five image generators. This included MJ-6.1 from Phase 1, as well as FLUX1.1[pro]

(referred to here as FLUX-1.1-PRO) (Black Forest Labs, 2024b), Ideogram 2.0 (IDEOGRAM-2.0) (Ideogram, 2024), Recraft V3 (RECREAF-V3) (Recraft, 2024), and Stable Diffusion 3.5 Large (SD-3.5-LARGE) (Stability AI, 2024).

We applied the scene illustration pipeline to produce illustrations for all 1000 fragments, varying runs of the pipeline between the 3 scene captioners and 5 image generators. We sampled a roughly equal ratio of pairs where the illustrations varied by scene captioner, image generator, or both scene captioner and image generator. The exact distribution is specified in A.1.2 and A.1.3. We repeated the same §4.3 procedure to obtain selections from two annotators for the better illustration in each pair. There were 48 (out of 49 total) annotators on Prolific who passed the attention checks, each annotating between 46 and 109 pairs (median=50), resulting in a total of 2436 responses for 1218 pairs.

4.6 Phase 2 Annotation Results

Inter-annotator Agreement As shown in Table A.8, the overall κ_u for all 1218 items in Phase 2 was 0.228, and annotators agreed in their responses for 52.4% of these items. This is lower than the agreement observed for Phase 1. §A.1.3 analyzes agreement across different pair subsets.

Win Rates for Scene Captioners Table 3 shows the win rates for each LLM scene captioner against each of the others. In particular, each value is the percentage of responses where the illustration associated with the scene captioner in the row label was selected as better than the illustration associated with the scene captioner in the column label. Statistically significant win rates are denoted with an asterisk. These results show that CLAUDE-3.5 yields the highest win rates, followed by GPT-4o, with LLAMA-3.1 having lowest rates. The win rate for CLAUDE-3.5 against LLAMA-3.1 is statistically significant, suggesting that the former generates more descriptive captions compared with the latter.

	CLAUDE-3.5	GPT-4o	LLAMA-3.1
CLAUDE-3.5	-	46.1	49.6*
GPT-4o	41.2	-	47.8
LLAMA-3.1	39.7	42.9	-

Table 3: Win rates (%) by scene captioner for Phase 2

Win Rates for Image Generators While not the focus of our analysis, we observed some significant differences in the win rates of different image generators. These results appear in §A.1.4.

⁶Statistical significance was computed using a one-sample binomial test at $\alpha = 0.05$ to determine if the win rate was higher than that expected by chance, where chance is defined as $(1 - \#ties/\#responses)/2$

⁷gpt-4o-2024-05-13, ran via the OpenAI API

⁸llama-3.1-405b-instruct, ran via the Replicate API

Fragment (within story)

Illustration 1

Illustration 2

Sophie’s nana was terminally ill. Sophie visited her in the hospital to say goodbye. **Her nana gave Sophie her prized gold locket. She told Sophie to keep it to remember her by.** Sophie cried.



Criteria	Response	Response
1. The image shows two people: an elderly woman (nana) and a younger woman (Sophie)	✓	✗
2. The setting appears to be a hospital room or medical facility	✓	✓
3. The elderly woman is in a hospital bed or medical chair	✓	✗
4. The image shows a gold locket	✓	✓
5. The locket is clearly visible and recognizable as a piece of jewelry	✓	✓
6. The elderly woman is holding or presenting the locket to the younger woman	✗	✓
7. The younger woman’s hand is positioned to receive or touch the locket	✓	✗
8. The facial expressions of both women convey emotional significance	✓	✓
9. The elderly woman’s expression shows love, tenderness, or sadness	✓	✓
10. The younger woman’s expression shows a mix of emotions (sadness, gratitude, love)	✓	✗
11. The body language of both women suggests intimacy and connection	✓	✓
12. The composition focuses on the moment of giving/receiving the locket	✓	✓
13. The lighting adequately illuminates the locket and the faces of both women	✓	✓
14. The locket appears to be in good condition, suggesting its value as a keepsake	✓	✓
15. The elderly woman’s appearance suggests illness or frailty	✓	✗
16. The younger woman’s appearance and demeanor suggest she is visiting	✓	✗
17. The overall atmosphere of the image conveys a solemn and meaningful moment	✓	✓
18. The spatial relationship between the two women suggests closeness and care	✓	✓
19. Any medical equipment or hospital elements are present but not dominating the scene	✓	✓
20. The perspective allows viewers to see both the locket and the emotional exchange between the women	✓	✓
	Score=19.0	Score=14.0

Table 4: Demonstration of criterial rating approach applied to both illustrations in a given pair. In this particular example, the criteria writer is CLAUDE-3.5, and the rater providing each response is GPT-4O.

5 Predicting Illustration Quality

5.1 Perfect-Agreement Data Subset

The SCENEILLUSTRATIONS dataset provides an opportunity to understand what characterizes a successful transformation of a narrative scene from text to image form. To initiate this line of work, we explored a particular approach to modeling annotators’ judgments of relative illustration quality. For this experiment, we combined the items from Phase 1 and Phase 2, and disregarded items involving annotator disagreement. The resulting *Perfect-Agreement* subset consists of 1745 items ($\approx 58\%$ of the full dataset) where both annotators agreed in their selection of the better illustration in the pair.

5.2 Criteria Generation

Our approach leverages the finding from §4 that LLMs can effectively verbalize visual descriptions of scenes based on the story text. We consider whether these descriptions can be used as *criteria* for predicting illustration quality. For each fragment, we ran the prompt in Table A.18 to produce criteria articulating the expected visual characteristics of the scene illustration. We use the term *criteria writer* to refer to an LLM’s role when running this prompt, and we refer to its output as a *criteria set*. An example of a criteria set is included in Table 4. Note that a criteria writer model does not require vision capabilities, since it observes

only the story text as input. §A.2.1 discusses some design considerations for generating criteria.

Criteria Writer Details We examined three criteria writers, the same LLMs that operated as scene captioners in §4.5: CLAUDE-3.5, GPT-4O, and LLAMA-3.1. Applying the Table A.18 prompt with temperature=0 to facilitate deterministic output, each criteria writer generated one criteria set per fragment. We post-processed this output to identify each individual criterion according to its expected numerical label in the set. §A.2.2 gives some descriptive analysis of the criteria sets.

5.3 Criteria-based Ratings

After obtaining the criteria sets, we then enlisted visually-enabled models to assess illustrations based on this criteria. In our scheme, when applying a criteria set to score a given illustration, each criterion receives a response indicating whether or not it is satisfied by the image. The overall illustration quality is quantified by the total number of satisfied criteria. Our scoring protocol is as follows: a response conveying that the criterion is satisfied is assigned 1.0 points; a response conveying “maybe” or partial satisfaction is assigned 0.5 points; and a response conveying the criterion is not satisfied is assigned 0.0 points. The total score for an illustration is the sum of these point values.

We implemented this by prompting a visually-

Criteria Writer	VLM Rater							
	CLAUDE-3.5		GPT-4O		PIXTRAL		Average	
	Criterial	Base	Criterial	Base	Criterial	Base	Criterial	Base
CLAUDE-3.5	0.717	0.606	0.709	0.567	0.712	0.589	0.713	0.587
GPT-4O	0.701	0.602	0.687	0.583	0.699	0.589	0.695	0.592
LLAMA-3.1	0.684	0.597	0.678	0.589	0.677	0.581	0.679	0.589
Average	0.700	0.602	0.691	0.580	0.696	0.586	0.696	0.589

Table 5: Accuracy of criterial and baseline (Base) raters grouped by criteria writer and VLM

enabled LLM (i.e. VLM) to assign responses to each criterion for a given illustration. We use the term *criterial rater* to refer to a VLM’s role when running this prompt, which appears in Table A.19. As shown, the rater observes an illustration and the criteria set for the corresponding fragment. The rater is asked to respond to each criterion (where a response of ‘✓’ means the criterion is satisfied, ‘✗’ means not satisfied, and ‘?’ means “maybe”). As post-processing, we parsed these response tokens and mapped them to the point values defined above to obtain the illustration score. Table 4 exemplifies this approach applied to both illustrations in a pair.

Rater Details For raters, we utilized three VLMs that have obtained notable performance on visual understanding benchmarks: CLAUDE-3.5, GPT-4O, and PIXTRAL⁹ (Mistral AI, 2024). Each rater ran the prompt in Table A.19 with temperature=0. All images were resized to a height of 240 pixels with proportional width. We briefly assessed the correctness of raters’ responses, which appears in §A.2.3.

Comparative Baseline To determine the impact of criteria in assessing quality, we designed a comparable rating approach that scores illustrations on the same scale as the criterial rater but without observing the criteria itself. We use the term *baseline rater* to refer to a VLM’s application of the prompt for this approach, which is shown in Table A.20. The prompt presents the fragment and illustration, and instructs the VLM to assign a rating in half-point increments between 0 and a maximum that is dynamically set to the length of the given criteria set. For each criteria writer, we compare the result obtained by a particular criterial rater to the analogous result obtained by the baseline rater.

5.4 Selection Performance Results

We applied all raters to score the illustrations in the *Perfect-Agreement* subset of SCENEILLUSTRATIONS. For a given pair, a rater’s selection was the image it assigned a higher score. We measured

each rater’s performance in terms of proportion of pairs where the rater’s selection matched the human selection. We refer to this metric as *accuracy*.

Table 5 shows the accuracy for all raters on these pairs, with the respective averages for each criteria writer and rater. For reference, always selecting the second illustration in each pair yields 49.4% accuracy. We observe that the criterial raters all considerably outperform the baseline raters (an average accuracy of $\approx 70\%$ vs. 59%). Criteria from different writers yields comparable results, with CLAUDE-3.5 averaging the highest accuracy across raters ($\approx 71\%$). The raters obtain similar accuracies when applied to the same criteria. Overall this outcome suggests that criteria are an effective strategy for modeling illustration quality, which in turn provides further evidence of LLMs’ capacity to verbalize visual characteristics of narrative scenes. This leaves room for further accuracy gains, motivating future exploration of this dataset for understanding what makes a compelling scene illustration.

6 Conclusion and Future Work

This paper details a pipeline for generating illustrations of narrative scenes, which we apply to produce SCENEILLUSTRATIONS, a quality-annotated dataset of illustrations for a popular story corpus. We identify that LLMs can facilitate this illustration task by distilling scene descriptions from story text. We show that this capacity to verbalize implicit scene knowledge is also useful for modeling illustration quality.

The scene illustration task isolates text-image alignment challenges in story visualization from issues of image-image alignment. In future work, we plan to consider recent approaches addressing the latter, such as ensuring visual consistency between story elements (e.g. Liu et al., 2025) and progressive story development across images (e.g. Maharana et al., 2022), in order to extend our illustration pipeline to generate multi-scene image sequences that depict complete stories.

⁹pixtral-large-2411, ran via the MistralAI API

Limitations

We consider the following limitations:

Proprietary Models Our scene illustration pipeline has a plug-and-play design, enabling any LLM to be used for fragmentation and scene captioning and any text-to-image model to be used for image generation. However, most of the models we assessed in this paper are proprietary (i.e. closed-weight), with exception to LLAMA-3.1 and SD-3.5-LARGE. While the gap between closed and open-weight models is narrowing (Cottier et al., 2024), currently most models with capabilities relevant to the illustration task are closed-weight. This poses a general disadvantage in accessibility and reproducibility, which applies likewise to this work.

Prompt Design Currently there is no tractable way to ensure that a particular prompt is optimal for the task it is intended to perform. Prompt optimization is fundamentally a process of iterative trial-and-error, even when automation is used to increase the number of trials. For our experiments, we primarily employed a principled approach to writing prompts, which involved adhering to general guidance on effective prompt design such as explaining instructions clearly and including representative exemplars (e.g. DAIR.AI, 2025). We iterated on this design according to qualitative subjective assessment of model outputs for inputs not included in our scene illustration dataset (i.e. “vibe-based” prompt engineering), rather than employing a quantitative optimization approach (e.g. Khattab et al., 2024) based on targets in a designated development set. There are tradeoffs to this technique: while it avoids overfitting to our presented dataset, it leaves open the possibility of further prompt optimization, which could yield a different view of model behavior compared with our observations.

Story Corpus The story corpus we use, ROCStories, is popular in NLP research for some of the same reasons discussed in §4: the constrained language and structure of the text makes the narrative elements more accessible to computational modeling techniques. The stories were authored specifically for the benefit of this research. However, this corpus is distinct from “naturally” authored stories whose complexity is what makes them compelling to readers. We have not yet fully assessed whether our scene illustration pipeline generalizes to more complex narratives.

Ethical Considerations

Generative AI models, and in particular text-to-image models, pose various ethical risks (Bird et al., 2023). In this work, we were primarily concerned with the risk of exposing Prolific annotators to harmful content. We attempted to mitigate this risk by manually reviewing stories sampled for inclusion in our dataset. We flagged stories that we anticipated could yield objectionable illustrations, and re-sampled a different story to replace each of these. Ultimately, this re-sampling was triggered for 10 stories. Of course, this procedure did not eliminate the risk, so we also utilized the content warning feature on the Prolific platform, which indicated to potential annotators that the task could expose them to offensive and/or biased content.

References

- Jie An, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Lijuan Wang, and Jiebo Luo. 2024. [Openleaf: A novel benchmark for open-domain interleaved image-text generation](#). In *Proceedings of the 32nd ACM International Conference on Multimedia*, MM '24, page 11137–11145, New York, NY, USA. Association for Computing Machinery.
- Charlotte Bird, Eddie Ungless, and Atoosa Kasirzadeh. 2023. [Typology of risks of generative text-to-image models](#). In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '23, page 396–410, New York, NY, USA. Association for Computing Machinery.
- Black Forest Labs. 2024a. [Announcing black forest labs: Flux.1 model family](#).
- Black Forest Labs. 2024b. [Announcing flux1.1 \[pro\] and the bfl api](#).
- Meriem Boubdir, Edward Kim, Beyza Ermis, Sara Hooker, and Marzieh Fadaee. 2023. [Elo uncovered: Robustness and best practices in language model evaluation](#). In *Proceedings of the Third Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)*, pages 339–352, Singapore. Association for Computational Linguistics.
- Stephen Brade, Bryan Wang, Mauricio Sousa, Sageev Oore, and Tovi Grossman. 2023. [Promptify: Text-to-image generation through interactive prompt exploration with large language models](#). In *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*, UIST '23, New York, NY, USA. Association for Computing Machinery.
- Jacob Cohen. 1960. A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.

860	Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu,	Sebastian Michelmann, Manoj Kumar, Kenneth A Nor-	915
861	Yu Cheng, Yuexin Wu, Lawrence Carin, David Carl-	man, and Mariya Toneva. 2025. Large language mod-	916
862	son, and Jianfeng Gao. 2019. StoryGAN: A Sequen-	els can segment narrative events similarly to humans.	917
863	tial Conditional GAN for Story Visualization . In	<i>Behavior Research Methods</i> , 57(1):1–13.	918
864	<i>2019 IEEE/CVF Conference on Computer Vision and</i>		
865	<i>Pattern Recognition (CVPR)</i> , pages 6322–6331, Los	Midjourney. 2024. Version 6.1 .	919
866	Alamitos, CA, USA. IEEE Computer Society.		
867	Long Lian, Boyi Li, Adam Yala, and Trevor Darrell.	Mistral AI. 2024. Pixtral large .	920
868	2024. LLM-grounded diffusion: Enhancing prompt	Nasrin Mostafazadeh, Nathanael Chambers, Xiaodong	921
869	understanding of text-to-image diffusion models with	He, Devi Parikh, Dhruv Batra, Lucy Vanderwende,	922
870	large language models . <i>Transactions on Machine</i>	Pushmeet Kohli, and James Allen. 2016. A corpus	923
871	<i>Learning Research</i> .	and cloze evaluation for deeper understanding of	924
872	Zhiqiu Lin, Deepak Pathak, Baiqi Li, Jiayao Li, Xide	commonsense stories . In <i>Proceedings of the 2016</i>	925
873	Xia, Graham Neubig, Pengchuan Zhang, and Deva	<i>Conference of the North American Chapter of the</i>	926
874	Ramanan. 2025. Evaluating text-to-visual generation	<i>Association for Computational Linguistics: Human</i>	927
875	with image-to-text generation. In <i>Computer Vision –</i>	<i>Language Technologies</i> , pages 839–849, San Diego,	928
876	<i>ECCV 2024</i> , pages 366–384, Cham. Springer Nature	California. Association for Computational Linguis-	929
877	Switzerland.	tics.	930
878	Tao Liu, Kai Wang, Senmao Li, Joost van de Weijer,	OpenAI, :, Aaron Hurst, Adam Lerer, Adam P. Goucher,	931
879	Fahad Shahbaz Khan, Shiqi Yang, Yaxing Wang, Jian	Adam Perelman, Aditya Ramesh, Aidan Clark,	932
880	Yang, and Ming-Ming Cheng. 2025. One-prompt-	AJ Ostrow, Akila Welihinda, Alan Hayes, Alec	933
881	one-story: Free-lunch consistent text-to-image gen-	Radford, Aleksander Mądry, Alex Baker-Whitcomb,	934
882	eration using a single prompt . In <i>The Thirteenth</i>	Alex Beutel, Alex Borzunov, Alex Carney, Alex	935
883	<i>International Conference on Learning Representa-</i>	Chow, Alex Kirillov, Alex Nichol, Alex Paino, Alex	936
884	<i>tions</i> .	Renzin, et al. 2024. Gpt-4o system card . <i>Preprint</i> ,	937
885	Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and	arXiv:2410.21276.	938
886	William Yang Wang. 2023. LLMScore: Unveiling	Jiaxin Pei, Aparna Ananthasubramaniam, Xingyao	939
887	the power of large language models in text-to-image	Wang, Naitian Zhou, Apostolos Dedeloudis, Jack-	940
888	synthesis evaluation . In <i>Thirty-seventh Conference</i>	son Sargent, and David Jurgens. 2022. POTATO:	941
889	<i>on Neural Information Processing Systems</i> .	The portable text annotation tool . In <i>Proceedings of</i>	942
890	Adyasha Maharana and Mohit Bansal. 2021. Integrat-	<i>the 2022 Conference on Empirical Methods in Nat-</i>	943
891	ing visuospatial, linguistic, and commonsense struc-	<i>tural Language Processing: System Demonstrations</i> ,	944
892	ture into story visualization . In <i>Proceedings of the</i>	pages 327–337, Abu Dhabi, UAE. Association for	945
893	<i>2021 Conference on Empirical Methods in Natural</i>	Computational Linguistics.	946
894	<i>Language Processing</i> , pages 6772–6786, Online and	Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi.	947
895	Punta Cana, Dominican Republic. Association for	2023. What does a platypus look like? generating	948
896	Computational Linguistics.	customized prompts for zero-shot image classifica-	949
897	Adyasha Maharana, Darryl Hannan, and Mohit Bansal.	tion. In <i>2023 IEEE/CVF International Conference</i>	950
898	2022. Storydall-e: Adapting pretrained text-to-image	<i>on Computer Vision (ICCV)</i> , pages 15645–15655.	951
899	transformers for story continuation. In <i>Computer</i>	Recraft. 2024. Recraft introduces a revolutionary ai	952
900	<i>Vision – ECCV 2022</i> , pages 70–87, Cham. Springer	model that thinks in design language .	953
901	Nature Switzerland.	Nikhil Singh, Guillermo Bernal, Daria Savchenko, and	954
902	Mayug Maniparambil, Chris Vorster, Derek Molloy,	Elena L. Glassman. 2023. Where to hide a stolen	955
903	Noel Murphy, Kevin McGuinness, and Noel E.	elephant: Leaps in creative writing with multimodal	956
904	O’Connor. 2023. Enhancing CLIP with GPT-4: Har-	machine intelligence . <i>ACM Trans. Comput.-Hum.</i>	957
905	nessing Visual Descriptions as Prompts . In <i>2023</i>	<i>Interact.</i> , 30(5).	958
906	<i>IEEE/CVF International Conference on Computer</i>	Soumik Rakshit. 2024. Building a genai-assisted auto-	959
907	<i>Vision Workshops (ICCVW)</i> , pages 262–271, Los	matic story illustrator .	960
908	Alamitos, CA, USA. IEEE Computer Society.	Stability AI. 2024. Introducing stable diffusion 3.5 .	961
909	Lech Mazur. 2025. Llm creative story-writing bench-	Ming Tao, Bing-Kun Bao, Hao Tang, Yaowei Wang,	962
910	mark .	and Changsheng Xu. 2024. Coin: A lightweight	963
911	Sachit Menon and Carl Vondrick. 2023. Visual classifi-	and effective framework for story visualization and	964
912	cation via description from large language models . In	continuation . In <i>Proceedings of the 32nd ACM Inter-</i>	965
913	<i>The Eleventh International Conference on Learning</i>	<i>national Conference on Multimedia</i> , MM ’24, page	966
914	<i>Representations</i> .	10659–10668, New York, NY, USA. Association for	967
		Computing Machinery.	968

- Yufei Tian, Tenghao Huang, Miri Liu, Derek Jiang, Alexander Spangher, Muhao Chen, Jonathan May, and Nanyun Peng. 2024. [Are large language models capable of generating human-level narratives?](#) In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 17659–17681, Miami, Florida, USA. Association for Computational Linguistics.
- Qian Wan, Xin Feng, Yining Bei, Zhiqi Gao, and Zhicong Lu. 2024. [Metamorpheus: Interactive, affective, and creative dream narration through metaphorical visual storytelling](#). In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI ’24, New York, NY, USA. Association for Computing Machinery.
- Zhijie Wang, Yuheng Huang, Da Song, Lei Ma, and Tianyi Zhang. 2024. [Promptcharm: Text-to-image generation through multi-modal prompting and refinement](#). In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, CHI ’24, New York, NY, USA. Association for Computing Machinery.
- Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, Nathan Cooper, Griffin Adams, Jeremy Howard, and Iacopo Poli. 2024. [Smarter, better, faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetuning and inference](#). *Preprint*, arXiv:2412.13663.
- Shuai Yang, Yuying Ge, Yang Li, Yukang Chen, Yixiao Ge, Ying Shan, and Yingcong Chen. 2024. [Seed-story: Multimodal long story generation with large language model](#). *Preprint*, arXiv:2407.08683.
- Weizhe Yuan, Pengfei Liu, and Matthias Gallé. 2024. [LLMCrit: Teaching large language models to use criteria](#). In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 7929–7960, Bangkok, Thailand. Association for Computational Linguistics.
- Yongchao Zhou, Andrei Ioan Muresanu, Ziwen Han, Keiran Paster, Silviu Pitis, Harris Chan, and Jimmy Ba. 2023. [Large language models are human-level prompt engineers](#). In *The Eleventh International Conference on Learning Representations*.

A Appendix

A.1 Additional Statistics for SCENEILLUSTRATIONS Dataset

A.1.1 Analysis of Phase 1 Fragments

As mentioned in §4.2, there were 206 total fragments derived from the 50 stories in Phase 1, based on applying CLAUDE-3.5 to the prompt in Table 13. As shown in Table 6, the majority consist of a single sentence, with some consisting of 2 sentences and a few having 3 sentences. An internal annotator assessed each fragment to determine if it was the correctly-sized unit for a scene illustration. A fragment was considered incorrectly-sized if it either did not include all the text in the story relevant to a single scene (i.e. the fragment was too short) or if it included text pertaining to more than one scene (i.e. the fragment was too long). The annotator considered the vast majority of fragments to be correctly-sized ($\approx 96\%$).

# Total Fragments	206
# 1-Sentence Fragments	164
# 2-Sentence Fragments	40
# 3-Sentence Fragments	2
Mean # Sentences Per Fragment	1.21
Mean # Fragments Per Story	4.12
% of Correctly-Sized Fragments	96.1%

Table 6: Fragmentation statistics for stories in Phase 1

A.1.2 Illustration Types

Table 7 characterizes the illustrations contained in the SCENEILLUSTRATIONS dataset. In particular, we display the number of unique illustrations derived from each scene description and image generator considered for Phase 1 and Phase 2.

A.1.3 Illustration Pair Types and Inter-annotator Agreement

Table 8 lists the number of illustration pairs associated with each combination of scene description types and image generators. This table also presents the specific inter-annotator agreement of different subsets of pairs. As presented in §4.4, the inter-annotator agreement in terms of uncertainty-weighted kappa (κ_u) for all Phase 1 pairs was 0.436. The first section of Table 8 shows that agreement was higher among the 1457 Phase 1 pairs where the illustrations used different scene descriptions ($\kappa_u=0.483$), while agreement was lower among the 661 pairs where the illustrations used different image generators ($\kappa_u=0.364$). This indicates that

Illustration Type	# Illustrations
<i>Phase 1</i>	
All	1576
By Scene Description	
NC-FRAGMENT	395
VC-FRAGMENT	384
SC-FRAGMENT	393
CAPTION	404
By Image Generator	
FLUX-1-PRO	791
MJ-6.1	785
<i>Phase 2</i>	
All	1582
By Scene Captioner	
CLAUDE-3.5	496
GPT-4O	532
LLAMA-3.1	554
By Image Generator	
FLUX-1.1-PRO	307
IDEOGRAM-2.0	300
MJ-6.1	321
RECRAFT-V3	323
SD-3.5-LARGE	331

Table 7: Number of unique illustrations associated with each scene description type and image generator in Phase 1 (top) and Phase 2 (bottom)

Illustration Pair Type	# Pairs	κ_u
<i>Phase 1</i>		
All	1777	0.436
Different Scene Descriptions	1457	0.483
NC-FRAGMENT vs. CAPTION	680	0.520
VC-FRAGMENT vs. CAPTION	384	0.504
SC-FRAGMENT vs. CAPTION	393	0.398
Different Image Generators		
FLUX-1-PRO vs. MJ-6.1	661	0.364
<i>Phase 2</i>		
All	1218	0.228
Different Scene Captioners	810	0.236
CLAUDE-3.5 vs. GPT-4O	266	0.198
CLAUDE-3.5 vs. LLAMA-3.1	268	0.234
GPT-4O vs. LLAMA-3.1	276	0.273
Different Image Generators	813	0.228
FLUX-1.1-PRO vs. IDEOGRAM-2.0	72	0.079
FLUX-1.1-PRO vs. MJ-6.1	74	0.183
FLUX-1.1-PRO vs. RECRAFT-V3	75	0.089
FLUX-1.1-PRO vs. SD-3.5-LARGE	71	0.164
IDEOGRAM-2.0 vs. MJ-6.1	99	0.312
IDEOGRAM-2.0 vs. RECRAFT-V3	81	0.159
IDEOGRAM-2.0 vs. SD-3.5-LARGE	94	0.339
MJ-6.1 vs. RECRAFT-V3	73	0.419
MJ-6.1 vs. SD-3.5-LARGE	88	0.184
RECRAFT-V3 vs. SD-3.5-LARGE	86	0.271
<i>Phase 1 & 2</i>		
All	2995	0.351

Table 8: Illustration pair statistics for Phase 1 and Phase 2 of the SCENEILLUSTRATIONS dataset, including inter-annotator agreement (κ_u) for different pair types

the scene description type was particularly influential to annotators’ judgments of which illustration was better. Considered along with Table 2, we can specifically conclude that ablating the scene captioner (i.e. using the baseline NC-FRAGMENT/VC-

FRAGMENT/SC-FRAGMENT scene descriptions) yielded illustrations that annotators consistently judged as lower quality relative to those that used the scene captioner.

For Phase 2, as reported in §4.6, the overall κ_u was 0.228 among all 1218 pairs. The second section of Table 8 shows that the agreement level was similar between the 810 pairs where illustrations involved different scene captioners ($\kappa_u=0.236$) and the 813 pairs that involved different image generators ($\kappa_u=0.228$). Agreement varied especially widely based on which particular image generators were paired together (ranging from 0.079 for FLUX-1.1-PRO vs. IDEOGRAM-2.0, up to 0.419 for MJ-6.1 vs. RECRAFT-V3). This indicates that in contrast to Phase 1 where there was a significant variable (the presence/absence of the scene captioner) that made the relative quality of illustrations more consistently distinguishable to annotators, the Phase 2 pairs were less reliably distinct.

A.1.4 Win Rates for Image Generators

To determine whether the choice of image generator influenced illustration quality in both Phase 1 and Phase 2, we computed the win rates for each image generator against each other among the pairs that used different image generators.

For Phase 1, there were only two image generators used to produce illustrations, FLUX-1-PRO vs. MJ-6.1. We did not find any significant difference in the win rates of these image generators. Table 9 shows these results.

FLUX-1-PRO	MJ-6.1
42.6%	41.0%

Table 9: Win rates (percentages) of FLUX-1-PRO vs MJ-6.1 for Phase 1 pairs

The Phase 2 data utilized a larger set of image generators. Table 10 shows the win rates of these image generators, presented comparably to Table 3 where each value is the percentage of selections for the image generator in the row against the image generator in the column. According to these results, IDEOGRAM-2.0 obtains the highest win rates against the other image generators, with significant success against FLUX-1.1-PRO, MJ-6.1, and SD-3.5-LARGE. Additionally, RECRAFT-V3 is significantly favored over MJ-6.1. Further analysis of these model differences for this task is an opportunity for future work.

A.2 Criteria-based Evaluation Details

A.2.1 Criteria Design Considerations

As referenced in §5.2, two design considerations for the criteria generation prompt (Table 18) were *flexibility* and *atomicity*. Flexibility emphasizes that a scene characteristic referenced by a criterion may be depicted with multiple alternative visual details that all align equally with the story text. For example, if a criterion conveys that the scene should take place at a particular location, it should be flexible about how the location is portrayed. Regarding atomicity, we aimed for each criterion to be as atomic as possible, meaning that it should refer to only a single characteristic of the scene. This promotes concise and easy-to-parse responses when judging whether the criterion is satisfied by an image, as opposed to a criterion that conflates multiple characteristics, some of which are satisfied and others that are not. Concerning the length of the generated criteria, our prompt did not specify a particular number of criteria to return, but the exemplar and instructions indicated that the criteria should comprehensively refer to as many scene characteristics as possible without redundancy.

A.2.2 Descriptive Analysis of Criteria Sets

Regarding the generated criteria sets (§5.2), Table 11 compares the average number of criteria in the sets generated by each criteria writer, revealing that CLAUDE-3.5 generated the longest criteria sets, followed by GPT-4O, and LLAMA-3.1.

CLAUDE-3.5	GPT-4O	LLAMA-3.1
19.3	17.3	15.8

Table 11: Mean number of criteria per set for each writer

Additionally, Figure 2 visualizes all criteria, based on encoding each criterion with the ModernBERT embedding model (Warner et al., 2024), then running PCA + t-SNE to yield a 2D embedding. While there are no distinct clusters associated with each criteria writer, some separation can be observed between the criteria generated by CLAUDE-3.5 and GPT-4O, while those generated by LLAMA-3.1 are more distributed alongside both other writers.

A.2.3 Criterial Rater Assessment

As referenced in §5.3, we conducted a small assessment of the correctness of the VLM raters’ responses to criteria. To do this, we randomly sampled 100 items, each with a unique image and crite-

	FLUX-1.1-PRO	IDEOGRAM-2.0	MJ-6.1	RECRAFT-V3	SD-3.5-LARGE
FLUX-1.1-PRO	-	35.4	43.2	45.3	48.6
IDEOGRAM-2.0	53.5*	-	61.6*	46.9	58.5*
MJ-6.1	39.9	31.8	-	28.8	44.9
RECRAFT-V3	44.0	40.1	61.6*	-	50.6
SD-3.5-LARGE	43.7	30.3	43.8	37.8	-

Table 10: Win rates (percentages) by image generator for Phase 2. Statistically significant win rates are denoted with an asterisk.

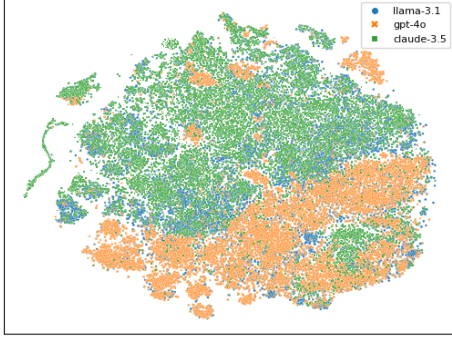


Figure 2: Visualization of criteria generated by each writer. Each point is a single criterion represented by its ModernBERT embedding. We applied PCA followed by t-SNE to plot the embedding in 2D space.

ria set. We then enlisted an expert human annotator to assign a response to each criterion, which we treated as the gold standard criterion response for the sampled image. We measured rater correctness in terms of linear-weighted κ agreement with the gold standard, where responses of ‘**X**’ were mapped to -1, ‘?’ to 0, and ‘✓’ to 1; this results in less weight assigned to disagreements involving ‘?’ (“maybe”) responses. Table 12 shows the κ on these 1699 criterion responses. It indicates that raters are all substantially aligned with the human annotator, though GPT-4O appears to have the highest human agreement, followed by CLAUDE-3.5, and then PIXTRAL.

Rater	κ
CLAUDE-3.5	0.676
GPT-4O	0.710
PIXTRAL	0.622

Table 12: Correctness of criterial rater responses (κ)

You are performing the task of story fragmentation. The task is to split a story into fragments where each fragment consists of a distinct part of the story. A fragment contains enough information to yield a visualization that is unique to that part of the story. In this version of the task, you will insert brackets (i.e. [and]) into the given story text to annotate the beginning and end of each fragment. Write the fragments without preamble. Here are some examples:

Story: Mia sat at home in her living room watching sports. Her favorite soccer team was playing their rival. To encourage her team, she began chanting positive phrases. During her chant, her favorite team scored a goal. Mia cheered loudly and thought that she helped score that goal.

Fragmented Story: [Mia sat at home in her living room watching sports. Her favorite soccer team was playing their rival.] [To encourage her team, she began chanting positive phrases.] [During her chant, her favorite team scored a goal.] [Mia cheered loudly and thought that she helped score that goal.]

[...2 more exemplars...]

Story: {{story}}

Fragmented Story:

Table 13: **Fragmentation prompt.** LLM prompt for annotating fragment boundaries in a story, which consists of a task instruction and exemplars demonstrating the task. The stories in the exemplars are from various corpora (ROCStories, TinyStories, and the ARL Creative Visual Storytelling Anthology).

You will be shown a story fragment along with its story context. Your task is to rewrite the fragment so that its meaning can be fully understood if read independently of the story context. For instance, you should replace names of characters with generic nouns. You should replace pronouns with the nouns they refer to (if the reference is a character, replace it with the appropriate generic noun). For first-person pronouns, replace the pronoun with a generic identifier (e.g. "I" -> "A person", "my" -> "the person's"). If the fragment implicitly refers to any other information in the story context, this information should be made explicit in the revised fragment. Write the revised fragment without preamble. Here are some examples:

Story Context: Anna was filling her bird feeders. But a chunk of suet fell onto the ground. Her dog rushed over and lapped it up! Anna was astonished. She had no idea dogs loved bird food!

Story Fragment: Her dog rushed over and lapped it up!

Revised Story Fragment: The woman's dog rushed over and lapped up the chunk of suet that had fallen onto the ground.

[...2 more exemplars...]

Story Context: {{story}}

Story Fragment: {{fragment}}

Revised Story Fragment:

Table 14: **Fragment rewriting prompt.** LLM prompt for generating SC-FRAGMENT scene descriptions. The prompt consists of a task instruction and exemplars demonstrating the task. The stories in the exemplars are from the ROCStories corpus.

Imagine an AI system will be used to generate illustrations for story fragments. This AI illustrator generates a single image given a caption describing what is contained in the image. Your task is to read a story fragment along with its story context and write a caption that describes how to illustrate the fragment. The caption should elaborately describe the most salient way to visualize the fragment. It should completely specify all the information the illustrator needs to generate the image. Write the caption without preamble. Here are some examples:

Story Context: Carrie had just learned how to ride a bike. She didn't have a bike of her own. Carrie would sneak rides on her sister's bike. She got nervous on a hill and crashed into a wall. The bike frame bent and Carrie got a deep gash on her leg.

Story Fragment: Carrie would sneak rides on her sister's bike.

Caption for Story Fragment: A young girl with a mischievous expression carefully wheels a bicycle that's slightly too big for her out of a garage, glancing over her shoulder as if making sure no one sees her.

[...2 more exemplars...]

Story Context: {{story}}

Story Fragment: {{fragment}}

Caption for Story Fragment:

Table 15: **Scene captioning prompt.** LLM prompt for generating CAPTION scene descriptions. The prompt consists of a task instruction and exemplars demonstrating the task. The stories in the exemplars are from various corpora (ROCStories, TinyStories, and the ARL Creative Visual Storytelling Anthology).









Fragment (within story)	CAPTION	CAPTION Illustration	VC-FRAGMENT Illustration
<p>Joe went to the junkyard to find treasures. He dug through a stack of old clothes. At the bottom was a vintage shirt with a rock band on it. Joe put the shirt on to try it out. Joe was happy the shirt fit.</p>	<p>A young man with an excited expression stands at the entrance of a sprawling junkyard. He's wearing casual clothes and holding a small backpack. The junkyard is filled with piles of discarded items, including old appliances, furniture, and scrap metal. The sky is bright and sunny, creating an adventurous atmosphere. In the foreground, a rusty sign reads "Junkyard" with operating hours listed below.</p>		
<p>Trish hated the outdoors. Her friends convinced her to go camping. She wasn't having a good time. They showed her how to fish and showed her the stars. Eventually Trish grudgingly came to accept nature.</p>	<p>A young woman with a scowling expression sits on a log by a campfire, arms crossed and looking miserable. She's surrounded by cheerful friends setting up tents and unpacking camping gear in a forest clearing. Her clean, urban clothing contrasts with the rugged outdoor setting, emphasizing her discomfort with nature.</p>		
<p>Sammy's coffee grinder was broken. He needed something to crush up his coffee beans. He put his coffee beans in a plastic bag. He tried crushing them with a hammer. It worked for Sammy.</p>	<p>A man in casual clothing stands at a kitchen counter, holding a hammer above a clear plastic bag filled with whole coffee beans. The hammer is poised mid-swing, about to strike the bag. The man's face shows a mix of determination and uncertainty. Scattered around the counter are a few stray coffee beans and an unplugged, visibly broken coffee grinder.</p>		
<p>I decided to go on a bike ride with my brother. We both headed out in the morning. We were having a lot of fun. Suddenly, he hit a rock and broke his wheel! I felt very badly for my brother.</p>	<p>A concerned young person stands next to their brother, who sits dejectedly on the ground next to a fallen bicycle with a visibly bent front wheel. The scene takes place on a sunny morning on a bike path, with trees and nature in the background. The standing sibling has a sympathetic expression, reaching out to comfort their brother, who looks disappointed and upset about the broken bike.</p>		

Table 16: Examples of scene illustrations in Phase 1. For each story fragment, we show an illustration resulting from the LLM-generated CAPTION scene description and one resulting from the baseline VC-FRAGMENT scene description. The image generator for all illustrations is FLUX-1-PRO.

Scene Description Pair	CAPTION Win %		
	MJ-6.1 & FLUX-1-PRO	MJ-6.1 Only	FLUX-1-PRO Only
CAPTION vs. NC-FRAGMENT	78.1	79.2	77.7
CAPTION vs. VC-FRAGMENT	74.7	74.5	75.0
CAPTION vs. SC-FRAGMENT	72.5	76.1	68.9

Table 17: Extended view of Table 2. Here, the win rates (percentages) for CAPTION vs. baseline scene descriptions in Phase 1 are split out by pairs where both illustrations used the same image generator (the MJ-6.1 Only and FLUX-1-PRO Only columns). This shows that the CAPTION win rate is similar regardless of which image generator is used.

Ellen dreamed of winning a prize for her roses. She planned to enter her special purple rose at the fair. She fertilized the rose bush and covered it each night. The roses grew more beautiful every day. Ellen ended up winning the prize.

Read the entire story above. Which image is a better visualization of the underlined fragment?



☐ I can't decide which image is better

Figure 3: Example of a item shown to participants in the annotation task described in §4.3

Imagine an AI system will be used to judge the quality of images intended to illustrate story fragments. This AI judge scores the images given some criteria about what should be depicted in the images. Your task involves writing the criteria for this AI judge. In particular, you will read a story and focus on a fragment at a specific position in the story. You will write the criteria defining the characteristics the image for that fragment should satisfy in order to be considered a good illustration of the fragment. There are a few things to consider when writing the criteria. First, while the criteria should define the fundamental characteristics depicted in the image, the visual details of these characteristics may vary across images, and alternative details may be similarly effective in illustrating the fragment. Each criterion should be written in a way that accommodates these potential variations in detail, without assuming specific information that is not defined explicitly in the story. Additionally, each criterion should refer to only a single atomic characteristic of the image. If a criterion references multiple characteristics such that an image might satisfy some but not others, it should be further split into multiple separate criteria. For example, instead of writing "the image shows a sapphire ring on the bathroom floor" as one criterion, you should write "the image shows a ring", "the ring contains a sapphire", and "the ring is on the bathroom floor" as separate criteria. The criteria should not only consider the presence of certain elements in the image, but also the visual quality of their depiction. Write the criteria without preamble, with a number header (e.g. '1.') for each criterion. Try to write as many criteria as possible, but avoid specifying extraneous or redundant criteria. Here is an example:

Story Context: Lisa has a beautiful sapphire ring. She always takes it off to wash her hands. One afternoon, she noticed it was missing from her finger! Lisa searched everywhere she had been that day. She was elated when she found it on the bathroom floor!

Story Fragment: She was elated when she found it on the bathroom floor!

Image Criteria for Story Fragment:

1. The image shows a clearly visible ring
2. The image portrays a bathroom setting recognizable through typical bathroom elements (tiles, fixtures, etc.)
3. The ring contains a blue gemstone recognizable as a sapphire
4. The ring is on the bathroom floor
5. The ring appears to be positioned naturally as if it had fallen or been dropped
6. A female figure (Lisa) is present in the image
7. The woman's facial expression clearly conveys joy or elation
8. The woman's body language demonstrates excitement or relief
9. The woman's positioning suggests she has just discovered or is reaching for the ring
10. The lighting adequately illuminates the ring to make it visible as the focal point
11. The perspective of the image allows viewers to see both the ring and the woman's emotional reaction
12. The composition draws attention to the moment of discovery
13. The spatial relationship between the woman and ring suggests imminent retrieval
14. The overall scene composition captures the spontaneous nature of the discovery
15. The woman's appearance suggests this is taking place during daytime/afternoon
16. The ring appears intact and undamaged, justifying the woman's relief
17. The bathroom setting appears residential rather than public

Story Context: {{story}}

Story Fragment: {{fragment}}

Image Criteria for Story Fragment:

Table 18: **Criteria generation prompt.** LLM prompt used to generate evaluation criteria for assessing the quality of scene illustrations.

You will observe an image along with a list of criteria, where each criterion describes a characteristic or quality that may or may not be depicted in the image. Your task is to determine whether or not each criterion is satisfied by the image. For each criterion, if the image fully satisfies that criterion, write a checkmark ('✓') after it. If the image only partially satisfies the criterion but not completely, write a question mark ('?') after it. Otherwise, if the image does not satisfy that criterion, write an X mark ('X') after it. Reiterate each criterion before giving your assessment for it, but do not provide additional preamble in your response. Here is an example:

Criteria:

1. The image shows a young woman (Laura) in an apartment setting
2. The woman's facial expression conveys happiness or contentment
3. The apartment appears to be newly moved into, with some visible unpacked items
4. There are visible windows in the apartment
5. The view through the windows shows recognizable California scenery (palm trees, ocean, mountains, or urban landscape)
6. The lighting suggests natural daylight entering the apartment
7. The apartment appears residential and suitable for a recent college graduate

Image: <IMAGE WILL APPEAR HERE>

Criteria Responses:

1. The image shows a young woman (Laura) in an apartment setting ✓
2. The woman's facial expression conveys happiness or contentment X
3. The apartment appears to be newly moved into, with some visible unpacked items ?
4. There are visible windows in the apartment ✓
5. The view through the windows shows recognizable California scenery (palm trees, ocean, mountains, or urban landscape) X
6. The lighting suggests natural daylight entering the apartment ✓
7. The apartment appears residential and suitable for a recent college graduate ✓

Criteria:

{{criteria}}

Image: {{image}}

Criteria Responses:

Table 19: **Criteria-based rating prompt.** VLM prompt used to score the quality of a scene illustration by assigning responses to each criterion in a provided criteria set

Your task is to rate how well a particular image illustrates a fragment of a story. You will observe the fragment with its story context, alongside the image depicting the fragment. Provide a rating on a scale ranging from 0.0 to {{len(criteria)}} in half-point increments, where 0.0 indicates the image is unrelated to the fragment and {{len(criteria)}} indicates the image is a perfect illustration of the fragment. Do not provide additional preamble before the rating. Here is an example:

Story: Laura had just graduated college. She was planning on moving on California. She packed all her belongings in her car and drove 18 hours. When she arrived at her new apartment she unpacked all her things. Laura loved the new change of scenery at her new place.

Fragment: Laura loved the new change of scenery at her new place.

Image: <IMAGE WILL APPEAR HERE>

Rating: 4.5

Story: {{story}}

Fragment: {{fragment}}

Image: {{image}}

Rating:

Table 20: **Baseline rating prompt.** VLM prompt used to score the quality of a scene illustration by directly assigning a rating between 0 and a maximum. This maximum is dynamically set to the total number of criteria in a provided criteria set ({{len(criteria)}}), even though the criteria themselves are not referenced in the prompt.