

# Intelligent Grimm - Open-ended Visual Storytelling via Latent Diffusion Models

Chang Liu<sup>1,3\*</sup>, Haoning Wu<sup>1\*</sup>, Yujie Zhong<sup>2</sup>, Xiaoyun Zhang<sup>1</sup>, Yanfeng Wang<sup>1,3</sup>, Weidi Xie<sup>1,3</sup>

<sup>1</sup>Coop. Medianet Innovation Center, Shanghai Jiao Tong University, China

<sup>2</sup>Meituan Inc., China

<sup>3</sup>Shanghai AI Laboratory, China

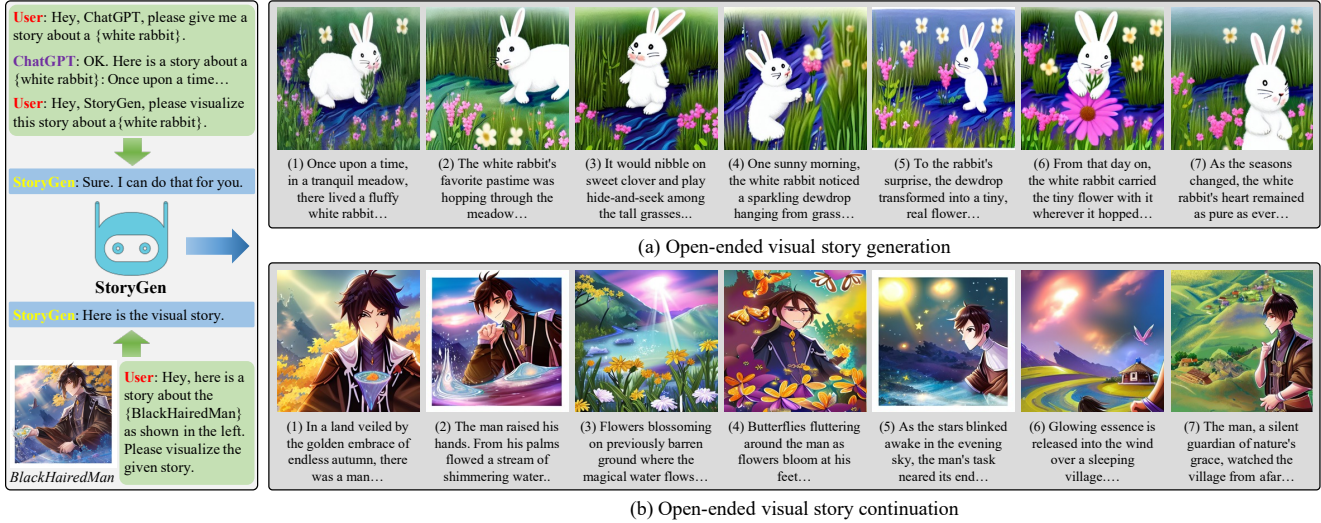


Figure 1. **An illustration of open-ended visual storytelling.** In practice, users can feed a unique and engaging story synthesized by a large language model into our proposed **StoryGen** model to generate a sequence of images coherently, denoted as *open-ended visual story generation*. And they can also provide a pre-defined character with its corresponding storyline, to perform *open-ended visual story continuation*. We recommend the reader to zoom in and read the story.

## Abstract

Generative models have recently exhibited exceptional capabilities in text-to-image generation, but still struggle to generate image sequences coherently. In this work, we focus on a novel, yet challenging task of generating a coherent image sequence based on a given storyline, denoted as **open-ended visual storytelling**. We make the following three contributions: (i) to fulfill the task of visual storytelling, we propose a learning-based auto-regressive image generation model, termed as **StoryGen**, with a novel vision-language context module, that enables to generate the current frame by conditioning on the corresponding text prompt and preceding image-caption pairs; (ii) to address the data shortage of visual storytelling, we collect paired image-text sequences by sourcing from online videos and open-source E-books, establishing processing pipeline for constructing a large-scale dataset with diverse char-

acters, storylines, and artistic styles, named **StorySalon**; (iii) Quantitative experiments and human evaluations have validated the superiority of our StoryGen, where we show StoryGen can generalize to unseen characters without any optimization, and generate image sequences with coherent content and consistent character. Code, dataset, and models are available at [https://haoningwu3639.github.io/StoryGen\\_Webpage/](https://haoningwu3639.github.io/StoryGen_Webpage/).

“Mirror mirror on the wall, who’s the fairest of them all?”

— Grimms’ Fairy Tales

## 1. Introduction

This paper considers an interesting, yet challenging task, namely, *open-ended visual storytelling*. The goal is to train a generative model that effectively captures the relation between visual elements and corresponding text descriptions, to generate a sequence of images that tell a visually coherent story, as shown in Figure 1. The outcome of this task has significant potential for education, as it provides chil-

\*: These authors contribute equally to this work.

dren with an engaging and interactive way to learn complex visual concepts and develop imagination, creativity, emotional intelligence, and language skills, as evidenced by research in psychology [5, 48].

The recent literature has witnessed tremendous progress in image generation, particularly with the guidance of text as prompt, such as stable diffusion [42], DALL-E [40] and Imagen [14]. However, to generalize the models for open-ended visual storytelling, we are facing three challenges: (i) previous models are designed to only generate images independently, without considering context, for example, preceding frames or overall narrative, resulting in a lack of visual consistency; (ii) most methods generate images by only conditioning on text, which potentially leads to ambiguities or requires unnecessarily long descriptions to maintain character appearances; (iii) existing datasets are limited to a few animations, covering a closed set of vocabulary or characters [25, 31, 36]. Training on such datasets suffers from severe overfitting on seen characters, leading to unsatisfactory generalization capability for open-ended generation.

This paper describes a learning-based model for open-ended visual storytelling, termed as **StoryGen**, that enables to generate unseen characters without any further optimization, while having character consistency. At inference, StoryGen can synthesize frames either by taking text prompts, or along with preceding image-text pairs as conditions, *i.e.*, iteratively creating visual sequences that are aligned with language description, while being consistent with preceding frames in both style and character perspectives. Specifically, to achieve consistency within the generated image sequence, we incorporate a novel **vision-language context module** into the pre-trained stable diffusion model, which provides visual context by conditioning the generation process on extracted diffusion denoising feature of previous frames under the guidance of corresponding captions.

As for training, we construct a dataset called **StorySalon**, that features a rich source of coherent images and stories, primarily comprising children’s storybooks collected from videos and E-books. As a result, our dataset includes a diverse vocabulary with different characters, storylines, and artistic styles. The scale and diversity of our collected dataset enable the model for open-vocabulary visual storytelling, *i.e.*, generating new image sequences that are not limited to pre-defined storylines, characters, or scenes. For example, we can prompt a large language model to create unique and engaging stories, then feed them into StoryGen for generation, as shown in Figure 1.

To summarize, we make the following contributions in this paper: (i) we initiate a fun yet challenging task, namely, *open-ended visual storytelling*, that involves generating engaging image sequences aligned to a given storyline; (ii) we propose a learning-based open-ended visual storytelling model, termed as **StoryGen**, which can generalize to un-

seen characters without any further optimization and generate coherent visual stories, utilizing a novel vision-language context module; (iii) we establish a data processing pipeline and collect a large-scale dataset of storybooks, called **StorySalon**, from online videos and open-source E-books, resulting in a diverse vocabulary with various characters, storylines, and artistic styles; (iv) we conduct quantitative experiments and human evaluations to validate the effectiveness of our proposed modules, demonstrating the superiority of our model, in terms of image quality, consistency, and visual-language alignment of generated contents.

## 2. Related Works

**Text-to-image Generation** has been tackled using various generative models, with GAN [8] as the first widely used model. Several GAN-based methods [53, 56, 57] have achieved notable success, and auto-regressive transformers [49], such as DALL-E [40], have also demonstrated the ability to generate high-quality images based on text prompts. Recently, diffusion models, such as Imagen [44] and DALL-E 2 [41], have emerged as a popular approach. Stable Diffusion Models [42] performs diffusion process in latent space, and can generate impressive images after pre-training on a large-scale text-image dataset.

**Diffusion Models** learn to model a data distribution via iterative denoising and are trained with denoising score matching. Notably, DDPM [13] demonstrates improved performance over other generative models, while DDIM [46] significantly boosts efficiency. In view of their superior generative capabilities, diffusion models have found extensive utility in various downstream applications besides image generation, such as video generation [6, 14, 15, 45], image manipulation [2, 10, 18, 33], grounded generation [26], image restoration [4], and image inpainting [1, 28, 35, 51].

**Story Synthesis** is first introduced as the task of story visualization by StoryGAN [25], which presents a GAN-based framework and the PororoSV dataset, derived from cartoons. Some works [29, 30] follow the GAN-based framework, whereas others [3, 21] emphasize more on text representation. StoryDALL-E [31] extends story synthesis to story continuation with the initial image given, and exploits a pre-trained DALL-E model [40] to produce coherent images. AR-LDM [36] introduces an auto-regressive latent diffusion model to generate image sequences, but only consistent within a limited character vocabulary. NUWA-XL [55] exploits hierarchical diffusion models to synthesize long videos, but still achieve character consistency by memorizing. TaleCrafter [7] proposes a story visualization system and utilizes LoRA [16] to achieve character consistency. However, large-scale applications will be constrained due to its optimization-based nature. In this paper, we target more ambitious applications, to develop an open-ended visual storytelling model, that can synthesize coherent image

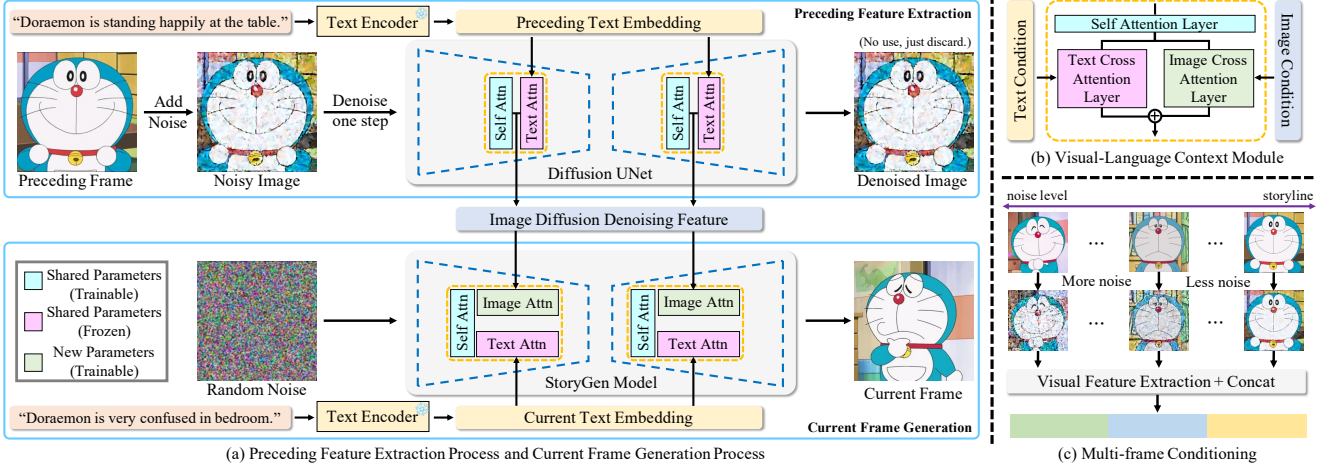


Figure 2. **Architecture Overview.** (a) Our StoryGen model utilizes current text prompt and previous visual-language contexts as conditions to generate an image, iteratively synthesizing a coherent image sequence. Note the parameters of the corresponding attention layers are shared between Diffusion UNet and StoryGen. To avoid potential ambiguity, the parameters are not shared across UNet blocks in a single model. (b) The proposed Visual-Language Context Module can effectively combine the information from current text prompt and contexts from preceding image-caption pairs. (c) We add more noise to reference frames with longer temporal distances to the current frame as positional encoding to distinguish the temporal order. The multiple features can then be directly concatenated to serve as context conditions.

sequences based on storylines of diverse topics.

### 3. Method

In this section, we start by formulating the problem of open-ended visual storytelling in Section 3.1; then we elaborate on the proposed StoryGen architecture in Section 3.2; lastly, we present details for model training in Section 3.3.

#### 3.1. Problem Formulation

In this paper, we focus on a challenging task, termed as *open-ended visual storytelling*, the goal is to generate continuous image sequence from a given story in the form of natural language. Specifically, we propose a learning-based auto-regressive image generation model, called **StoryGen**, that generates the current frame  $\hat{\mathcal{I}}_k$  by conditioning on the current text prompt  $\mathcal{T}_k$ , and image-text pairs  $(\hat{\mathcal{I}}_{<k}, \mathcal{T}_{<k})$  of previous frames, as illustrated in Figure 2 (a). The model is formulated as follows:

$$\{\hat{\mathcal{I}}_1, \hat{\mathcal{I}}_2, \dots, \hat{\mathcal{I}}_L\} = \Phi_{\text{StoryGen}}(\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_L\}; \Theta)$$

$$\hat{\mathcal{I}}_k := \Phi_{\text{StoryGen}}(\hat{\mathcal{I}}_k | \mathcal{T}_k, (\hat{\mathcal{I}}_{<k}, \mathcal{T}_{<k}))$$

Here,  $\{\mathcal{T}_1, \mathcal{T}_2, \dots, \mathcal{T}_L\}$  refer to the given storylines, and  $\{\hat{\mathcal{I}}_1, \hat{\mathcal{I}}_2, \dots, \hat{\mathcal{I}}_L\}$  denote the generated image sequence.  $\Phi_{\text{StoryGen}}(\cdot)$  represents our proposed StoryGen model. In one-step generation, StoryGen takes the current text prompt, and preceding image-caption pairs as conditions, and generates an image consistent with both the story’s narrative and previous frames. The whole image sequence can then be synthesized with step-by-step inference.

**Relation to Existing Tasks.** In contrast to existing story visualization works, this paper makes improvements from two aspects: (i) conventional generation/continuation tasks are limited to training on specific characters/stories, for example, [25, 31, 36] only exploits datasets from animation *The Flintstones* and *Pororo*, while our model enables to generate visual stories based on any given storyline, such as a brand-new one generated by ChatGPT; and any pre-defined character, for example, ‘Doraemon’ from the Internet; (ii) unlike existing work that requires costly character-specific optimization, for example, [7, 36] rely on LoRA-based [16] optimization to adapt to new characters, our model is learning-based and expected to generalize to any unseen character without any further optimization.

#### 3.2. Architecture

To tackle the problem of open-ended visual storytelling, we expect the model to not only condition on the current text prompt, but also preceding image-text pairs. In this section, we describe the procedure for one-step generation, i.e., generating the  $k$ -th frame ( $k > 1$ ) by conditioning on  $\{(\hat{\mathcal{I}}_1, \mathcal{T}_1), \dots, (\hat{\mathcal{I}}_{k-1}, \mathcal{T}_{k-1}), \mathcal{T}_k\}$ . Generally speaking, our proposed **StoryGen** model comprises four components: (i) Input Initialization, (ii) Context Encoding, (iii) Visual-Language Contextual Fusion, (iv) Conditional Generation.

**Input Initialization.** Our model is built upon the foundation of a pre-trained stable diffusion model (SDM), which randomly samples a noisy latent  $\mathbf{x}$  from the latent space of the VAE [19] encoder. Moreover, for a given text prompt  $\mathcal{T}_k$ , the text condition will be extracted by a pre-trained CLIP [38] text encoder  $\phi_{\text{CLIP}}$  via  $\mathcal{C}^T = \phi_{\text{CLIP}}(\mathcal{T}_k)$ .



**Context Encoding.** In standard SDM, the noisy latent is recursively denoised with a UNet, conditioning on the text prompt. However, in our case, it is crucial for the generation procedure to also condition on context features of preceding frames, to maintain consistency in characters and storyline.

In practice, to extract the contextual features, we add noise to the preceding frames and exploit the pre-trained SDM to denoise for one diffusion step under the guidance of their corresponding captions. The diffusion features after every self-attention layer in the UNet blocks can be directly selected to serve as the conditioning visual context features, thus constituting a pyramid of visual context features. The visual condition features for  $\hat{\mathcal{I}}_k$  can be expressed as:

$$\mathcal{C}^V = [\phi_{\text{SDM}}(\hat{\mathcal{I}}_1, \phi_{\text{CLIP}}(\mathcal{T}_1)), \dots, \phi_{\text{SDM}}(\hat{\mathcal{I}}_{k-1}, \phi_{\text{CLIP}}(\mathcal{T}_{k-1}))]$$

Experimentally, we notice that, the magnitude of noise added to the preceding frames can greatly affect the conditional generation quality, *i.e.*, large-scale noise on preceding frames incurs severe information loss. Thus, we propose to use a much smaller diffusion timestep  $t'$  for preceding frames compared with the diffusion timestep  $t$  of the current image  $\hat{\mathcal{I}}_k$ , and follow a  $t' = t/10$  rule. As depicted in Figure 2 (c), in case of multiple preceding image-caption pairs, we use larger  $t'$  for frames with longer temporal distances to  $\hat{\mathcal{I}}_k$ . Therefore, the extracted multi-frame visual context features can be directly concatenated, and their different noise level will serve as temporal positional embedding. Such design reflects the intuition that frames with longer distances will incur less effect on generating the current frame.

**Vision-Language Contextual Fusion.** Here, our vision-language context module is designed to fuse information from current text prompt and contextual information from preceding image-caption pairs. This is achieved by augmenting the transformer decoder in SDM with an additional image cross-attention layer. Note that, the math expression in this section is not strict, we omit the footnote of diffusion timestep  $t$  and UNet block level  $l$  for simplicity.

Specifically, on visual context conditioning, the noisy latent  $\mathbf{x}$  is projected into query, and cross-attends to the visual context features from the corresponding-level UNet block that act as key and value, denoted as:

$$\mathbf{Q}_I = \mathbf{x}\mathbf{W}_I^Q, \quad \mathbf{K}_I = \mathcal{C}^V\mathbf{W}_I^K, \quad \mathbf{V}_I = \mathcal{C}^V\mathbf{W}_I^V$$

where  $\mathbf{W}_I^Q$ ,  $\mathbf{W}_I^K$ , and  $\mathbf{W}_I^V$  represent different projection matrices, respectively.

On text conditioning, the noisy latent  $\mathbf{x}$  is again projected to query, and cross-attends to the text features of the current prompt encoded by CLIP text encoder, *i.e.*,

$$\mathbf{Q}_T = \mathbf{x}\mathbf{W}_T^Q, \quad \mathbf{K}_T = \mathcal{C}^T\mathbf{W}_T^K, \quad \mathbf{V}_T = \mathcal{C}^T\mathbf{W}_T^V$$

where  $\mathbf{W}_T^Q$ ,  $\mathbf{W}_T^K$ , and  $\mathbf{W}_T^V$  also represent corresponding projection matrices.

As depicted in Figure 2 (b), the image cross-attention layer is inserted in parallel to the text cross-attention layer in the transformer decoder of UNet blocks. Drawing inspiration from ControlNet [58], the results from these two cross-attention layers are simply summed up as the final output  $\mathbf{O}$ . The final output can thus be expressed as:

$$\mathbf{O} = \text{Softmax}\left(\frac{\mathbf{Q}_I(\mathbf{K}_I)^\top}{\sqrt{d}}\right)\mathbf{V}_I + \text{Softmax}\left(\frac{\mathbf{Q}_T(\mathbf{K}_T)^\top}{\sqrt{d}}\right)\mathbf{V}_T$$

**Conditional Generation.** With the fused vision-language condition features from above, our StoryGen can now generate visual stories that achieve both content coherence and character consistency. Here, our conditional generation procedure can be represented as:

$$\hat{\mathcal{I}}_k = \Phi_{\text{StoryGen}}(\hat{\mathcal{I}}_k | \mathcal{T}_k, (\hat{\mathcal{I}}_{<k}, \mathcal{T}_{<k})) = \Phi_{\text{StoryGen}}(\mathbf{x}, \mathcal{C}^T, \mathcal{C}^V)$$

With the new conditioning modality introduced, we also adopt another classifier-free guidance term [12], as has been done in [2]. Concretely, we exploit two different guidance scales,  $w_v$  and  $w_t$  for the visual condition and the text condition. The relation between the final noise for inference  $\bar{\epsilon}_\theta$  and UNet-predicted noise  $\epsilon_\theta$  is now expressed as:

$$\begin{aligned} \bar{\epsilon}_\theta(\mathbf{x}_t, t, \mathcal{C}^V, \mathcal{C}^T) &= \epsilon_\theta(\mathbf{x}_t, t, \emptyset, \emptyset) \\ &\quad + w_v(\epsilon_\theta(\mathbf{x}_t, t, \mathcal{C}^V, \emptyset) - \epsilon_\theta(\mathbf{x}_t, t, \emptyset, \emptyset)) \\ &\quad + w_t(\epsilon_\theta(\mathbf{x}_t, t, \mathcal{C}^V, \mathcal{C}^T) - \epsilon_\theta(\mathbf{x}_t, t, \mathcal{C}^V, \emptyset)) \end{aligned}$$

**Discussion.** Our work differs from previous ones from two aspects. First, our StoryGen is a learning-based method, which can directly generalize to unseen characters by attending to reference images. Second, we propose to condition the generation process on diffusion features of preceding image-text pairs from the same SDM, which preserves more visual details, greatly differing from existing works [22, 52, 54] using CLIP, BLIP [24], or VAE features.

### 3.3. Model Training

**Training Objective.** At training stage, we randomly sample a triplet each time, *i.e.*,  $\{\mathcal{I}_k, \mathcal{T}_k, (\mathcal{I}_{<k}, \mathcal{T}_{<k})\}$ . The objective function can be expressed as:

$$\mathcal{L}_t = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t, \mathcal{C}^V, \mathcal{C}^T} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t, \mathcal{C}^V, \mathcal{C}^T)\|^2 \right]$$

**Two-stage Training Strategy.** Our two-stage training strategy includes single-frame pre-training and multiple-frame fine-tuning. To be specific, at the first stage, we do not introduce additional image cross-attention layers, and only train self-attention layers in standard SDM to ensure the single-frame generation ability. In multiple-frame fine-tuning, we train additional image cross-attention layers in vision-language context module on our dataset, with all other parameters frozen. This enables the generation procedure to





Figure 3. **Dataset Pipeline and Visualization.** **Left:** Meta-data sourced from the Internet undergoes a three-step pipeline including frame extraction, visual-language alignment and post-processing, resulting in properly aligned image-text pairs. **Right:** Our StorySalon dataset contains diverse styles and characters.

utilize information from not only current prompt, but also preceding image-caption pairs.

**Inference.** As shown in Figure 1, at inference time, we can prompt ChatGPT to generate novel storylines, and synthesize the first image directly or attending to a pre-defined character. Then the previously synthesized frames, along with the story descriptions, are treated as conditions to synthesize the image sequence in an auto-regressive manner. Experimentally, our proposed StoryGen is shown to generate images that align with the storyline, as well as maintain consistency with previously generated frames.

#### 4. StorySalon Dataset

In order to train our proposed *open-ended visual storytelling* model, we construct a large-scale dataset, termed as **StorySalon**. The dataset contains videos and E-books with diverse characters, storylines, and artistic styles. Specifically, we download a large number of videos and subtitles from YouTube, by querying keywords related to story-telling for children, for instance, *storytime*. Additionally, we collect E-books (partially with corresponding audios available) from six open-source libraries which are all registered under the Creative Commons 4.0 International Attribution (CC BY 4.0) license. In the following, we elaborate on the data processing pipeline and statistics of our collected dataset.

**Visual Frame Extraction.** We extract keyframes from the videos, along with the corresponding subtitles and their timestamps. To remove duplicate frames, we extract ViT features for each frame using pre-trained DINO [32]. For the image groups with high similarity scores, we only keep one of each. Then, we use YOLOv7 [50] to segment and remove real-person frames and headshots, as they often correspond to the story-teller and are unrelated to the content of the storybook. Similarly, we extract images from the

| Dataset           | Style     | #Frames        | Avg.Length | #Categories |
|-------------------|-----------|----------------|------------|-------------|
| PororoSV [25]     | Animation | 73,665         | 5          | 9           |
| FlintstonesSV [9] | Animation | 122,560        | 5          | 7           |
| DiDeMoSV [31]     | Real      | 52,905         | 3          | -           |
| VIST [17]         | Real      | 145,950        | 5          | -           |
| <b>StorySalon</b> | Animation | <b>159,778</b> | <b>14</b>  | <b>446</b>  |

Table 1. **Dataset Statistics.** Our StorySalon dataset far exceeds previous story generation datasets in terms of the total number of images, average length, and categories of characters included.

downloaded E-books, except for those with extraneous information, for example, the authorship page. We acquire the corresponding text description with Whisper [39] from the audio file, and for E-books that do not have corresponding audio files, but with available storyline text, we use OCR algorithms, to directly recognize the text on each page.

**Visual-Language Alignment.** As shown in Figure 3, for each of the image, we can collect two types of text descriptions, *e.g.*, story-level narration, and descriptive captions. This is based on our observation that there actually exists a semantic gap between narrative storyline and descriptive text, for example, the same image can be well described as “*The cat is isolated by others, sitting alone in front of a village.*” in the story, or “*A black cat sits in front of a number of houses.*” as descriptive caption, therefore, directly fine-tuning stable diffusion models with story narration may be detrimental to its pre-trained text-image alignment. In practice, to get story-level paired image-text samples, we align the subtitles with visual frames by using Dynamic Time Warping (DTW) algorithm [34]. To get visual descriptions, we use TextBind [23] to generate captions for each image, with both the image and the corresponding narrative text as inputs. At training time, this allows us to substitute the original story with more accurate and descriptive captions.

**Visual Frame Post-processing.** In practice, we observe that book pages and borders in images can potentially interfere with our generative model by having story texts printed on them. To tackle this, we use an OCR detector to identify text regions in images and an image inpainting model [42] to fill in the text and headshot regions, resulting in more precise image-text pairs that are suitable for model training.

**Discussion.** After the three-step pipeline above, we obtain our StorySalon dataset. As shown in Table 1, our dataset has nearly 160K animation-style images in total with an average length of 14 frames per story, which is conducive to building long-range semantic correspondence. Finally, we query MiniGPT-4 [59] about the main character category of each image in our dataset, like *Dog* and *Cat*, then count the categories and filter out those appear less than 3 times. Compared with previous datasets with less than 10 characters, our dataset comprises hundreds of character categories, and even more character instances, which provides a data basis

| Model           | FID ↓        | CLIP-I ↑      | CLIP-T ↑      |
|-----------------|--------------|---------------|---------------|
| GT              | -            | 1.0           | 0.2668        |
| SDM             | 73.50        | 0.6155        | 0.3218        |
| Prompt-SDM      | 67.35        | 0.6272        | <b>0.3225</b> |
| Finetuned-SDM   | 42.01        | 0.6970        | 0.3005        |
| StoryDALL-E     | 38.34        | 0.6823        | 0.2366        |
| AR-LDM          | 39.55        | 0.6864        | 0.2614        |
| <b>StoryGen</b> | <b>33.90</b> | <b>0.7467</b> | 0.2875        |

Table 2. **Comparison of automatic metrics** on StorySalon test set. Prompt-SDM denotes Stable Diffusion model with cartoon-style-directed prompts and Finetuned-SDM represents a Stable Diffusion model with all parameters fine-tuned on our StorySalon dataset.

for training open-ended visual storytelling models, showing a significantly broader range of visual styles and character appearances over existing datasets.

## 5. Experiments

In this section, we start by describing our experimental settings, then compare with other models from three different perspectives: image-text alignment, consistency and image quality with subjective human evaluation and quantitative metrics. Additionally, we present results for ablation experiments to prove the effectiveness of our proposed modules.

### 5.1. Experimental Settings

**Training Details.** Our model is built on the stable diffusion v1.5 model, and trained with a learning rate of  $1 \times 10^{-5}$  and a batch size of 256. We begin with a single-frame self-attention pre-training stage, which involves 3,000 iterations on 8 NVIDIA RTX3090. Next, we incorporate our proposed vision-language context module, and train it for 5,000 iterations using a single preceding image-caption pair as context condition, then continue to train it for another 5,000 iterations with multiple image-caption pairs for multi-frame conditioning. To maintain our model’s unconditional denoising ability for classifier-free guidance, we randomly drop the current text and the context image-caption pairs with a probability of 5% and 15%, respectively. During inference, we utilize DDIM [46] with 40 steps of sampling and select the guidance weight  $w_v = 7.0$  and  $w_t = 3.5$ .

**Baselines.** We consider two scenarios of our proposed open-ended storytelling task, namely, story generation and story continuation. For **story generation**, we need the model to be able to generate a complete visual story only based on a given storyline. So we present a comparison with Stable Diffusion Model (SDM) and **Prompt-SDM**, which conditions on an additional cartoon-style-directed prompt “A cartoon style image”. For **story continuation**, the first

| Story Generation   |             |             |             |             |             |               |
|--------------------|-------------|-------------|-------------|-------------|-------------|---------------|
| Model              | Align. ↑    | Style ↑     | Cont. ↑     | Char. ↑     | Qual. ↑     | Pref. ↑       |
| GT                 | 4.04        | 4.66        | 4.41        | 4.54        | 4.29        | —             |
| SDM                | 3.61        | 2.88        | 2.90        | 2.51        | 3.74        | 14.05%        |
| Prompt-SDM         | 3.39        | 2.56        | 2.68        | 2.10        | 3.44        | 8.57%         |
| StoryGen-S         | 3.50        | 2.73        | 2.81        | 2.21        | 3.19        | 10.24%        |
| <b>StoryGen</b>    | <b>3.78</b> | <b>4.79</b> | <b>4.26</b> | <b>4.64</b> | <b>3.76</b> | <b>67.14%</b> |
| Story Continuation |             |             |             |             |             |               |
| StoryDALL-E        | 1.18        | 1.55        | 1.20        | 1.14        | 1.19        | 0.63%         |
| AR-LDM             | 2.47        | 2.82        | 2.40        | 1.87        | 2.54        | 2.50%         |
| <b>StoryGen</b>    | <b>4.23</b> | <b>4.70</b> | <b>4.35</b> | <b>4.38</b> | <b>4.18</b> | <b>96.87%</b> |

Table 3. **Comparison results of human evaluation.** GT stands for ground truth from the test set. StoryGen-S represents StoryGen without context conditions. The abbreviated metrics are Text-image alignment, Style consistency, Content consistency, Character consistency, image quality, and Preference, respectively.

frame or the main character is given, and the model is expected to generate coherent images based on the storyline. In this scenario, we compare our model with two closed-set story continuation models: namely, **StoryDALL-E** [31] and **AR-LDM** [36] re-trained on our StorySalon dataset.

**Automatic Metrics.** To evaluate the quality of generated image sequences, we adopt three widely-used metrics, including Fréchet Inception Distance score (FID) [11], CLIP image-image similarity (CLIP-I), and CLIP text-image similarity (CLIP-T). Notably, in order to avoid the impact of randomness in synthesis quality, we utilize a CLIP-based scoring function trained exclusively on text-to-image generated images, namely, PickScore [20], to automatically select the generated images with better quality. Each chosen image is selected from a pool of 10 candidates.

### 5.2. Quantitative Evaluation Results

We compare our StoryGen model with other baselines on StorySalon test set, which contains 5% of total data (nearly 7K pairs). Each contains a current prompt and the image-text context of the previous frame. The models are expected to generate the current frame based on given conditions.

The quantitative results in Table 2 demonstrate that our StoryGen model exhibits significant performance improvement in terms of FID score and CLIP-I similarity compared to existing models, while maintaining comparable CLIP-T similarity. This confirms that our model can effectively exploit contextual information, thus generating animation-style visual stories based on the given storyline. Notably, CLIP trained on natural images tends to have an understanding bias towards animation-style images, and the slight decline in CLIP-T is an inevitable result of the conflict between text condition and newly introduced image condition.

### 5.3. Human Evaluation Results

Considering that the above metrics may not reflect the quality of the generated stories accurately, and there is no stan-



(a) *Open-ended story generation for: a story of a {white dog}*: (1) Once upon a time, in a peaceful countryside, there lived a white dog... (2) The white dog had an adventurous spirit, always eager to discover... (3) One afternoon, the white dog was staring at a sunflower... (4) The white dog ventured into a sunflower field... (5) The white dog discovered a bird's nest in the field... (6) From that day on, the white dog became a guardian... (7) The white dog's spirit remained steadfast and bright, as the seasons changed and the leaves fell...



(b) *Open-ended story continuation for: a story of a {a white-haired man}*: (1) In a perpetual twilight, the white-haired man reached towards the twilight sky, stars appearing at his touch... (2) One twilight, the white-haired man looked concerned at a dark void in the sky... (3) The white-haired man drew stars in the sky with a silver quill... (4) The white-haired man was observing new constellations shining where the void once was... (5) The white-haired man with a serene expression was watching the peaceful starry sky... (6) The white-haired man walked towards a tower observatory under the starry sky... (7) Alone but content, the white-haired man's gaze traversed the depths of space...

Figure 4. **Qualitative Comparison with other methods.** The image sequences in orange, green, and blue boxes are generated by Prompt-SDM, AR-LDM and StoryGen respectively. Our synthesis results exhibit impressive performance superiority in terms of style, content and character consistency, text-image alignment, and image quality. Please refer to the Appendix for more qualitative results.

standardized metric for evaluating the consistency within the visual story, we further include human evaluation for comparison of image-text alignment, image style, story consistency, character consistency and synthesis quality.

For the two scenarios mentioned above, we respectively conduct two types of human evaluation to assess the quality of generated visual stories. To mitigate bias, participants are unaware of the type of storybooks they are evaluating. Concretely, we prompt GPT-4 to produce multiple storylines for both test modes, and for story continuation, we search the Internet for multiple characters that have never appeared in our dataset. Then we utilize our StoryGen along with other baselines to generate corresponding sequences of images.

**Protocol-I.** We randomly select an equal number of samples from the generated results of our StoryGen and other baselines. Each time we randomly sample a visual story from these sources, and participants are then invited to rate the sample with a score ranging from 1 to 5, taking into account text-image alignment, style consistency, content con-

sistency, character consistency and image quality. Higher scores indicate better samples. We also evaluate the same number of samples from StorySalon test set as a reference.

**Protocol-II.** Each time we randomly sample a storyline and its corresponding visual storybooks generated by StoryGen and other methods. Participants are invited to select their preferred generated result among these different image sequences of the same storyline.

**Results.** The results of human evaluation presented in Table 3 illustrate that our StoryGen model demonstrates excellent performance in overall score, especially in terms of consistency and quality. This indicates that our model can generate coherent image sequences that are highly consistent with given text prompts and visual-language contexts.

## 5.4. Qualitative Results

In Figure 4, we present visualization results of both open-ended visual story generation and visual story continuation,



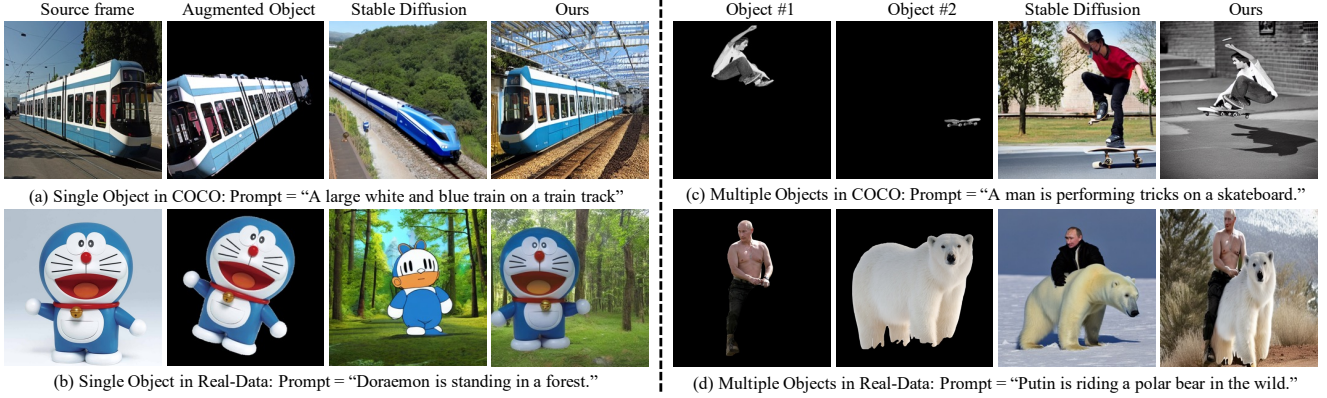


Figure 5. **Ablation studies on consistency.** We incorporate our proposed Visual-Language Context Module into a pre-trained SDM, and train it on MS-COCO [27] with other parameters frozen. The content consistency of single-object and multi-object generation on COCO and real data has demonstrated the effectiveness of our module. Please refer to the Appendix for experiment details and quantitative results.

showing that our StoryGen can generate visual stories with a broad vocabulary, while maintaining content coherence and character consistency throughout the narrative, whereas other methods fail to do so. Moreover, our model can stably maintain the animation style of generated images, which satisfies the requirements of visual storytelling for children. More results can be found in the supplementary material.

## 5.5. Ablation Studies

In order to demonstrate the effectiveness of our proposed modules, we conduct ablation studies from both quantitative metrics and qualitative visualization.

**On Variants of StoryGen.** We evaluate the performance of multiple model variants on the StorySalon test set, including (i) our model without the context module, marked as **StoryGen-Single**, which solely fine-tunes the self-attention layers on our dataset. (ii) our model with context features encoded by the VAE of SDM as context condition, without text-guided diffusion process, denoted as **StoryGen-VAE**; (iii) our model with CLIP image embedding as context condition (**StoryGen-CLIP**); (iv) our model with context features extracted by BLIP image encoder (**StoryGen-BLIP**); (v) our model with naive denoising features at Large-scale diffusion Timestep, satisfying  $t' = t$ , as condition (**StoryGen-LT**); and (vi) our full model (**StoryGen**). We also employ PickScore to filter generation results of all these models. The findings presented in Table 4 illustrate the inclusion of our context module can significantly improve the model performance, in terms of CLIP-I and FID. As for the slight inferiority in CLIP-T, we have claimed above that this is due to the understanding bias towards animation-style images for CLIP trained on natural images. **Qualitative Visualization.** As mentioned above, consistency is a crucial factor in visual story generation. We hope to more intuitively demonstrate that our proposed context module can accurately capture the image content of

| Model           | FID ↓        | CLIP-I ↑      | CLIP-T ↑      |
|-----------------|--------------|---------------|---------------|
| StoryGen-Single | 38.81        | 0.6869        | <b>0.3140</b> |
| StoryGen-VAE    | 36.98        | 0.6846        | 0.3061        |
| StoryGen-CLIP   | 36.66        | 0.6934        | <b>0.3140</b> |
| StoryGen-BLIP   | 34.78        | 0.7026        | 0.2838        |
| StoryGen-LT     | 36.41        | 0.7141        | 0.3025        |
| <b>StoryGen</b> | <b>33.90</b> | <b>0.7467</b> | 0.2875        |

Table 4. **Ablation studies** on Visual-Language Context Module.

the previous frame. To this end, we incorporate our context module into SDM and train it from scratch on the MS-COCO [27] with other parameters frozen. Specifically, we crop the object and perform data augmentations such as translation and rotation to use it as image condition. The category of the cropped object is used as its corresponding text, and the caption of the original image serves as the text prompt. We expect the model to reconstruct the original image relying on the conditions above, which enables the context module to learn how to leverage the previous image. As shown in Figure 5, our model can make full use of the objects in the reference frame and generate new images that are consistent with them, while SDM fails to do so. In addition, this can also be transferred to any real-world reference image, which strongly illustrates the robustness and capability of our context module to assist diffusion models in generating images based on any given object.

## 6. Conclusion

In this paper, we consider an interesting, yet challenging task, termed as *open-ended visual storytelling*, which involves generating a sequence of images that tell a coherent visual story based on the given storyline. Our proposed learning-based **StoryGen** model can take input from the preceding image-caption context along with the text prompt to generate coherent image sequences in an auto-regressive

manner, *i.e.*, without test-time optimization. On the data side, we establish a data processing pipeline to collect a large-scale dataset named **StorySalon** that comprises storybooks with diverse characters, storylines, and artistic styles sourced from videos and E-books. Extensive human evaluation and quantitative comparison have illustrated that our proposed model substantially outperforms existing models, from the perspective of image quality, content coherence, character consistency, and visual-language alignment.

## Acknowledgments

This work is supported by National Key R&D Program of China (No. 2022ZD0161400), National Natural Science Foundation of China (62271308), STCSM (22511105700, 22DZ2229005), 111 plan (BP0719010), and State Key Laboratory of UHD Video and Audio Production and Presentation.

## References

- [1] Omri Avrahami, Dani Lischinski, and Ohad Fried. Blended diffusion for text-driven editing of natural images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [2] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 2, 4
- [3] Hong Chen, Rujun Han, Te-Lin Wu, Hideki Nakayama, and Nanyun Peng. Character-centric story visualization via visual planning and token alignment. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2022. 2
- [4] Zheng Chen, Yulun Zhang\*, Ding Liu, Bin Xia, Jinjin Gu, Linghe Kong\*, and Xin Yuan. Hierarchical integration diffusion model for realistic image deblurring. In *Advances in Neural Information Processing Systems*, 2023. 2
- [5] K. Dickinson David, A. Griffith Julie, Golinkoff Roberta, Michnick, and Hirsh-Pasek Kathy. How reading books fosters language development around the world. *Child Development Research*, 2012. 2
- [6] Patrick Esser, Johnathan Chiu, Parmida Atighehchian, Jonathan Granskog, and Anastasis Germanidis. Structure and content-guided video synthesis with diffusion models. In *Proceedings of the International Conference on Computer Vision*, 2023. 2
- [7] Yuan Gong, Youxin Pang, Xiaodong Cun, Menghan Xia, Yingqing He, Haoxin Chen, Longyue Wang, Yong Zhang, Xintao Wang, Ying Shan, and Yujiu Yang. Talecrafter: Interactive story visualization with multiple characters. *SIGGRAPH Asia*, 2023. 2, 3
- [8] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the ACM*, 2020. 2
- [9] Tanmay Gupta, Dustin Schwenk, Ali Farhadi, Derek Hoiem, and Aniruddha Kembhavi. Imagine this! scripts to compositions to videos. In *Proceedings of the European Conference on Computer Vision*, 2018. 5
- [10] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. In *Proceedings of the International Conference on Learning Representations*, 2023. 2
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. In *Advances in Neural Information Processing Systems*, 2017. 6
- [12] Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. In *NeurIPS 2021 Workshop on Deep Generative Models and Downstream Applications*, 2021. 4
- [13] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *Advances in Neural Information Processing Systems*, 2020. 2
- [14] Jonathan Ho, William Chan, Chitwan Saharia, Jay Whang, Ruiqi Gao, Alexey Gritsenko, Diederik P Kingma, Ben Poole, Mohammad Norouzi, David J Fleet, et al. Imagen video: High definition video generation with diffusion models. *arXiv preprint arXiv:2210.02303*, 2022. 2
- [15] Jonathan Ho, Tim Salimans, Alexey Gritsenko, William Chan, Mohammad Norouzi, and David J Fleet. Video diffusion models. In *Advances in Neural Information Processing Systems*, 2022. 2
- [16] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *Proceedings of the International Conference on Learning Representations*, 2022. 2, 3
- [17] Ting-Hao Huang, Francis Ferraro, Nasrin Mostafazadeh, Ishan Misra, Aishwarya Agrawal, Jacob Devlin, Ross Girshick, Xiaodong He, Pushmeet Kohli, Dhruv Batra, et al. Visual storytelling. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2016. 5
- [18] Bahjat Kavar, Shiran Zada, Oran Lang, Omer Tov, Huiwen Chang, Tali Dekel, Inbar Mosseri, and Michal Irani. Imagic: Text-based real image editing with diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 2
- [19] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations*, 2014. 3
- [20] Yuval Kirstain, Adam Polyak, Uriel Singer, Shahbuland Matiana, Joe Penna, and Omer Levy. Pick-a-pic: An open dataset of user preferences for text-to-image generation. In *Advances in Neural Information Processing Systems*, 2023. 6, 5
- [21] Bowen Li. Word-level fine-grained story visualization. In *Proceedings of the European Conference on Computer Vision*, 2022. 2

- [22] Dongxu Li, Junnan Li, and Steven CH Hoi. Blip-diffusion: Pre-trained subject representation for controllable text-to-image generation and editing. In *Advances in Neural Information Processing Systems*, 2023. 4
- [23] Huayang Li, Siheng Li, Deng Cai, Longyue Wang, Lemao Liu, Taro Watanabe, Yujiu Yang, and Shuming Shi. Textbind: Multi-turn interleaved multimodal instruction-following. *arXiv preprint arXiv:2309.08637*, 2023. 5
- [24] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. In *Proceedings of the International Conference on Machine Learning*, 2022. 4
- [25] Yitong Li, Zhe Gan, Yelong Shen, Jingjing Liu, Yu Cheng, Yuexin Wu, Lawrence Carin, David Carlson, and Jianfeng Gao. Storygan: A sequential conditional gan for story visualization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019. 2, 3, 5
- [26] Ziyi Li, Qinye Zhou, Xiaoyun Zhang, Ya Zhang, Yanfeng Wang, and Weidi Xie. Guiding text-to-image diffusion model towards grounded generation. In *Proceedings of the International Conference on Computer Vision*, 2023. 2
- [27] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision*, 2014. 8, 4
- [28] Andreas Lugmayr, Martin Danelljan, Andres Romero, Fisher Yu, Radu Timofte, and Luc Van Gool. Repaint: Inpainting using denoising diffusion probabilistic models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2
- [29] Adyasha Maharana and Mohit Bansal. Integrating visuospatial, linguistic, and commonsense structure into story visualization. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, 2021. 2
- [30] Adyasha Maharana, Darryl Hannan, and Mohit Bansal. Improving generation and evaluation of visual stories via semantic consistency. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, 2021. 2
- [31] Adyasha Maharana, Darryl Hannan, and Mohit Bansal. Storydall-e: Adapting pretrained text-to-image transformers for story continuation. In *Proceedings of the European Conference on Computer Vision*, 2022. 2, 3, 5, 6
- [32] Caron Mathilde, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision*, 2021. 5
- [33] Chenlin Meng, Yutong He, Yang Song, Jiaming Song, Jiajun Wu, Jun-Yan Zhu, and Stefano Ermon. Sdedit: Guided image synthesis and editing with stochastic differential equations. In *Proceedings of the International Conference on Learning Representations*, 2021. 2
- [34] Meinard Müller. Dynamic time warping. *Information retrieval for music and motion*, 2007. 5
- [35] Alex Nichol, Prafulla Dhariwal, Aditya Ramesh, Pranav Shyam, Pamela Mishkin, Bob McGrew, Ilya Sutskever, and Mark Chen. Glide: Towards photorealistic image generation and editing with text-guided diffusion models. In *Proceedings of the International Conference on Machine Learning*, 2022. 2
- [36] Xichen Pan, Pengda Qin, Yuhong Li, Hui Xue, and Wenhui Chen. Synthesizing coherent story with auto-regressive latent diffusion models. In *Winter Conference on Applications of Computer Vision*, 2024. 2, 3, 6
- [37] Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe Penna, and Robin Rombach. Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*, 2023. 6
- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *Proceedings of the International Conference on Machine Learning*, 2021. 3
- [39] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision. In *Proceedings of the International Conference on Machine Learning*, 2023. 5
- [40] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation. In *Proceedings of the International Conference on Machine Learning*, 2021. 2
- [41] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022. 2
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Bjorn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2022. 2, 5, 3, 6
- [43] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *Medical Image Computing and Computer-Assisted Intervention—MICCAI 2015: 18th International Conference, Munich, Germany, October 5-9, 2015, Proceedings, Part III* 18, 2015. 2
- [44] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L Denton, Kamyar Ghasemipour, Raphael Gontijo Lopes, Burcu Karagol Ayan, Tim Salimans, et al. Photorealistic text-to-image diffusion models with deep language understanding. In *Advances in Neural Information Processing Systems*, 2022. 2
- [45] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiuyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. Make-a-video: Text-to-video generation without text-video data. In *Proceedings of the International Conference on Learning Representations*, 2023. 2
- [46] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *Proceedings of the Inter-*



*national Conference on Learning Representations*, 2020. 2, 6

- [47] Yang Song, Prafulla Dhariwal, Mark Chen, and Ilya Sutskever. Consistency models. In *Proceedings of the International Conference on Machine Learning*, 2023. 6
- [48] Gabrielle A. Strouse, Angela Nyhout, and Patricia A. Ganea. The role of book features in young children’s transfer of information from picture books to real-world contexts. *Frontiers in Psychology*, 2018. 2
- [49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in Neural Information Processing Systems*, 2017. 2
- [50] Chien-Yao Wang, Alexey Bochkovskiy, and Hong-Yuan Mark Liao. YOLOv7: Trainable bag-of-freebies sets new state-of-the-art for real-time object detectors. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 5
- [51] Shaoan Xie, Zhifei Zhang, Zhe Lin, Tobias Hinz, and Kun Zhang. Smartbrush: Text and shape guided object inpainting with diffusion model. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2023. 2
- [52] Li Xin, Chu Wenqing, Wu Ye, Yuan Weihang, Liu Fanglong, Zhang Qi, Li Fu, Feng Haocheng, Ding Errui, and Wang Jingdong. Videogen: A reference-guided latent diffusion approach for high definition text-to-video generation. *arXiv preprint arXiv:2309.00398*, 2023. 4
- [53] Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. AttnGAN: Fine-grained text to image generation with attentional generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018. 2
- [54] Hu Ye, Jun Zhang, Sibao Liu, Xiao Han, and Wei Yang. Ip-adapter: Text compatible image prompt adapter for text-to-image diffusion models. *arXiv preprint arxiv:2308.06721*, 2023. 4
- [55] Shengming Yin, Chenfei Wu, Huan Yang, Jianfeng Wang, Xiaodong Wang, Minheng Ni, Zhengyuan Yang, Linjie Li, Shuguang Liu, Fan Yang, et al. Nuwa-xl: Diffusion over diffusion for extremely long video generation. In *Association for Computational Linguistics*, 2023. 2
- [56] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. StackGAN: Text to photo-realistic image synthesis with stacked generative adversarial networks. In *Proceedings of the International Conference on Computer Vision*, 2017. 2
- [57] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris N Metaxas. StackGAN++: Realistic image synthesis with stacked generative adversarial networks. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 2
- [58] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image diffusion models. In *Proceedings of the International Conference on Computer Vision*, 2023. 4
- [59] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language

understanding with advanced large language models. *arXiv preprint arXiv:2304.10592*, 2023. 5, 4

# Intelligent Grimm - Open-ended Visual Storytelling via Latent Diffusion Models

## Supplementary Material

### Contents

|  |          |
|--|----------|
| <b>1. Introduction</b>                                     | <b>1</b> |
| <b>2. Related Works</b>                                    | <b>2</b> |
| <b>3. Method</b>   | <b>3</b> |
| 3.1. Problem Formulation . . . . .                         | 3        |
| 3.2. Architecture . . . . .                                | 3        |
| 3.3. Model Training . . . . .                              | 4        |
| <b>4. StorySalon Dataset</b>                               | <b>5</b> |
| <b>5. Experiments</b>                                      | <b>6</b> |
| 5.1. Experimental Settings . . . . .                       | 6        |
| 5.2. Quantitative Evaluation Results . . . . .             | 6        |
| 5.3. Human Evaluation Results . . . . .                    | 6        |
| 5.4. Qualitative Results . . . . .                         | 7        |
| 5.5. Ablation Studies . . . . .                            | 8        |
| <b>6. Conclusion</b>                                       | <b>8</b> |
| <b>A Preliminaries on Diffusion Models</b>                 | <b>2</b> |
| <b>B Further Architecture Details</b>                      | <b>2</b> |
| B.1. Parameter Sharing Strategy . . . . .                  | 2        |
| B.2. Multi-frame Condition Strategy . . . . .              | 2        |
| B.3. Two-stage Training Strategy . . . . .                 | 3        |
| <b>C Dataset Details</b>                                   | <b>3</b> |
| C.1. Data Sources . . . . .                                | 3        |
| C.2. Dataset Statistics . . . . .                          | 4        |
| <b>D Consistency Ablation on COCO</b>                      | <b>4</b> |
| D.1. Experiment Details . . . . .                          | 4        |
| D.2. Quantitative Results . . . . .                        | 5        |
| D.3. Qualitative Results . . . . .                         | 5        |
| <b>E Broader Impacts</b>                                   | <b>5</b> |
| <b>F. Limitations</b>                                      | <b>6</b> |
| <b>G More Experiments</b>                                  | <b>6</b> |
| G.1. Analysis on multi-frame conditioning . . . . .        | 6        |
| G.2. Multi-object conditioned Story Continuation . . . . . | 6        |
| G.3. Story Generation Visualization . . . . .              | 6        |
| G.4. Story Continuation Visualization . . . . .            | 6        |
| G.5. Failure Case Visualization . . . . .                  | 15       |

## A. Preliminaries on Diffusion Models

Diffusion models are a type of generative models that undergo a denoising process, converting input noise into meaningful data samples. Diffusion models comprise a forward diffusion process that incorporates Gaussian noise into an image sample  $\mathbf{x}_0$ , accomplished via a Markov process over  $T$  steps. If we denote the noisy image at step  $t$  as  $\mathbf{x}_t$ , the transition function  $q(\mathbf{x}_t|\mathbf{x}_{t-1})$  connecting  $\mathbf{x}_{t-1}$  and  $\mathbf{x}_t$  can be expressed as follows:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \quad q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^T q(\mathbf{x}_t|\mathbf{x}_{t-1})$$

where  $\beta_t \in (0, 1)$  is the variance schedule controlling the step size.

Using Gaussian distribution property and reparameterization, if we define  $\alpha_t = 1 - \beta_t$  and  $\bar{\alpha}_t = \prod_{i=1}^t \alpha_i$ , we can write the equation above as follows:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$$

Diffusion models also comprise a reverse diffusion process that learns to restore the initial image sample from noise. A UNet-based model [43] is utilized in the diffusion model to learn the reverse diffusion process  $p_\theta$ . The process  $p_\theta$  can be expressed using the following equation.

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^T p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

where  $\boldsymbol{\mu}_\theta$  is the predicted Gaussian distribution mean value.

As we compute the loss function by taking the mean absolute error of the noise term  $\epsilon_\theta$  into account, we can express the mean value  $\boldsymbol{\mu}_\theta$  in terms of the noise term  $\epsilon_\theta$  as follows:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(\mathbf{x}_t, t) \right)$$

Therefore, the objective can be written as:

$$\mathcal{L}_t = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right]$$

## B. Further Architecture Details

In this section, we will provide a more comprehensive illustration about more design details of our model.

### B.1. Parameter Sharing Strategy

Initially, we will discuss the strategy of parameter sharing between the standard diffusion UNet and our StoryGen. As illustrated in Figure 2 of the main paper, the standard diffusion UNet is exploited in *Preceding Feature Extraction* to extract diffusion context features, and StoryGen is utilized in *Current Frame Generation* to generate new frames with consistency and coherence.

Specifically, the parameters of the corresponding attention layers are shared between the standard diffusion UNet and our StoryGen, including self-attention layers and text cross-attention layers. In practise, the standard diffusion UNet here is a modified version of our StoryGen, without the image cross-attention layers, and all other parameters are shared. This design allows the UNet to extract contextual diffusion features within the same latent space as StoryGen.

### B.2. Multi-frame Condition Strategy

As illustrated in Figure 2 and Section 3.2 of the main paper, when dealing with multiple preceding image-caption pairs, we add more noise (corresponding to a larger diffusion timestep  $t'$ ) to reference frames with longer temporal distances to the current frame. Such design effectively serves two purposes, *first*, it is based on the observation that frames with longer temporal distances will incur less effect on the generation of the current frame; *second*, the different noise level also serves as positional encoding, allowing for the differentiation of temporal order, which enables us to directly concatenate these diffusion context features.

Specifically, we use  $t$  to represent the diffusion timestep of the current image  $\hat{\mathcal{I}}_k$ , and use  $t'_j$  to represent the diffusion timestep of the preceding image-text pair  $(\hat{\mathcal{I}}_j, \mathcal{T}_j)$ . When generating image  $\hat{\mathcal{I}}_k$ , we use  $t'_{k-1} = t/10$  for  $(\hat{\mathcal{I}}_{k-1}, \mathcal{T}_{k-1})$  pair,  $t'_{k-2} = 2t/10$  for  $(\hat{\mathcal{I}}_{k-2}, \mathcal{T}_{k-2})$  pair, and so on. In summary, the diffusion timestep  $t'_{k-i}$  for  $(\hat{\mathcal{I}}_{k-i}, \mathcal{T}_{k-i})$  pair will follow a  $t'_{k-i} = i * t/10$  rule.



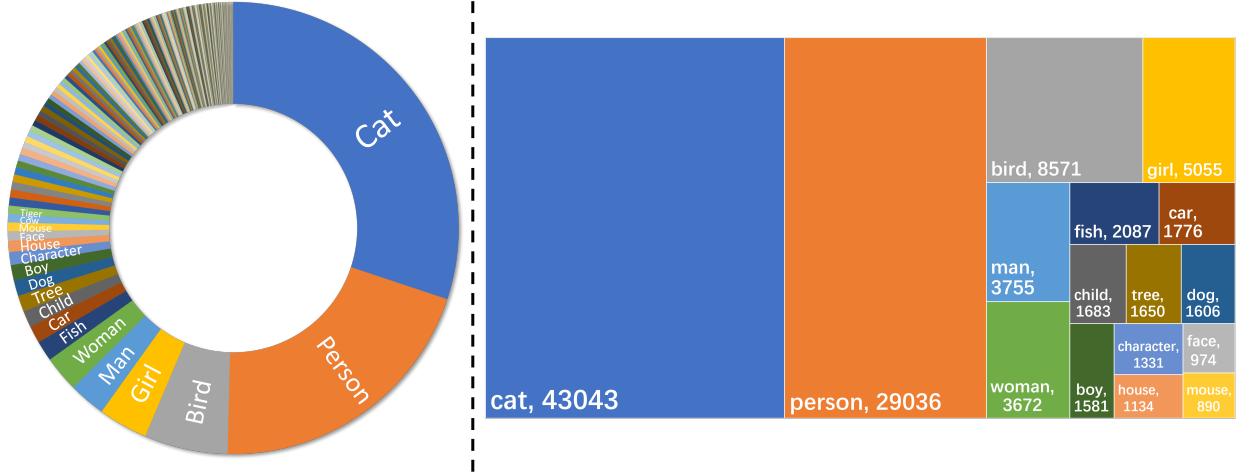


Figure 6. **Dataset Statistics Results.** **Left:** Distribution of text-image pairs classified by the main character categories in our collected StorySalon dataset. **Right:** The top 16 character categories and corresponding numbers in StorySalon, cover a wide range of character types.

### B.3. Two-stage Training Strategy

As illustrated in Section 3.3 of the main paper, we exploit a two-stage training strategy, including single-frame pre-training and multiple-frame fine-tuning. In **single-frame pre-training stage**, we do not introduce additional image cross-attention layers, and only train the self-attention layers in standard SDM [42] on our dataset. The goal of this stage is to train the model for single-frame generation in the style of storybooks. In **multiple-frame fine-tuning stage**, we train the additional image cross-attention layers in vision-language context module on our dataset, with all other parameters frozen. We first train image cross-attention layers with a single preceding image-caption pair, and then continue with multiple image-caption pairs. Consequently, StoryGen acquires the capability to condition on multiple preceding image-caption pairs and generate image sequences in an auto-regressive manner. Throughout the entire two-stage training, the text cross-attention layers remain frozen, preserving the vision-language alignment inherited from the pre-trained stable diffusion models.

## C. Dataset Details

In Section C.1, we present additional details about the data sources of our StorySalon dataset. Subsequently, we show the detailed statistics of our dataset in Section C.2.

### C.1. Data Sources

Our StorySalon dataset mainly comprises of two components, *e.g.*, online videos collected from YouTube, and open-source E-books collected from six online libraries. For online videos, we download a large number of videos and corresponding subtitles from YouTube, by querying keywords related to story-telling for children, for instance, *storytime*. For open-source E-books, we collect E-books (partially with corresponding audios available) from six open-source online libraries which are all registered under the Creative Commons 4.0 International Attribution (CC BY 4.0) license. These online libraries have consistently dedicated themselves to assisting children in underdeveloped regions. We extend our appreciation for their ongoing endeavors and contributions. Specifically, these open-source online libraries include:

- **African Storybook.** <https://africanstorybook.org/>;
- **Bloom Library.** <https://bloomlibrary.org/>;
- **Book Dash.** <https://bookdash.org/>;
- **Global Digital Library.** <https://digitallibrary.io/topic/library-books/>;
- **Room to Read.** <https://literacycloud.org/>;
- **Digital Library of Illustrated Storybooks.** <https://storyweaver.org.in/en>.

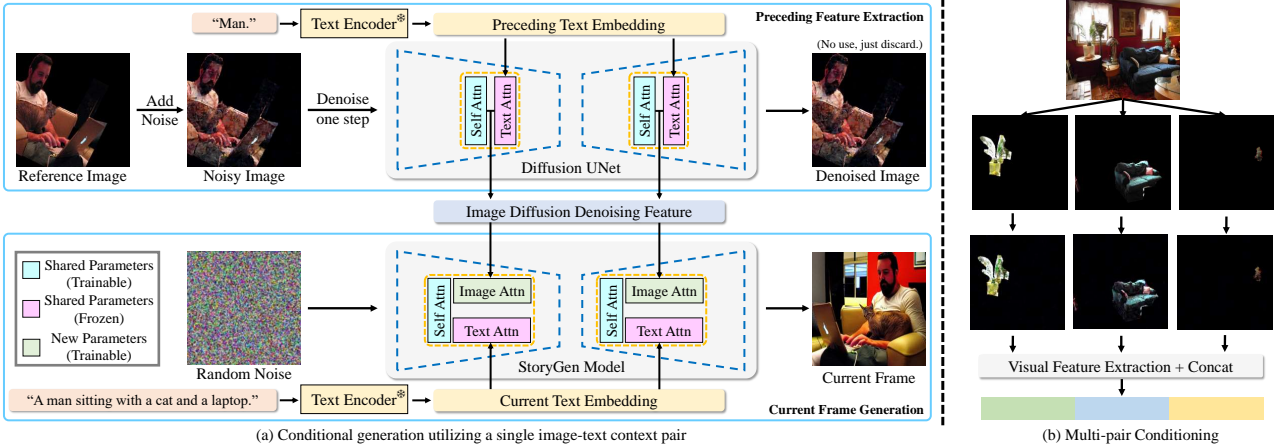


Figure 7. **Architecture Overview.** (a) Our StoryGen model utilizes current text prompt and previous visual-language contexts as conditions to generate an image. (b) In case of multiple image-text context pairs, the multiple features can be directly concatenated to serve as context conditions.

## C.2. Dataset Statistics

Our StorySalon dataset comprises a total of 11,280 storybooks and 159,778 text-image pairs, with approximately 160K animation-style images, averaging 14 frames per narrative, as demonstrated in Table 1 of the main paper. We divide the dataset into train and test sets following a 9 : 1 ratio. Both the video and E-book components are randomly split into train and test sets according to this proportion.

To categorize the characters in storybooks, we use MiniGPT-4 [59] to infer the predominant character category in each image of our dataset, such as *Dog* and *Cat*. Subsequently, we count these categories and exclude those occurring fewer than three times. The distribution of these character categories is depicted in Figure 6. In contrast to preceding datasets featuring fewer than ten characters, our dataset encompasses hundreds of character categories, with rich appearances. Consequently, StorySalon offers the data foundation for training open-ended visual storytelling models.

## D. Consistency Ablation on COCO

In Section D.1, we will give a brief illustration on further details about our qualitative consistency ablation experiment on COCO [27]. Subsequently, in Section D.2, we will design a new quantitative ablation experiment on COCO, and present additional results. Finally, in Section D.3, we will provide more visualization results of this consistency ablation on COCO.

### D.1. Experiment Details

**Motivation.** This experiment serves as an ablation study, designed to show StoryGen’s ability in utilizing image conditions, and preserving visual details. Specifically, at training time, we train StoryGen on images from COCO datasets, through a self-supervised learning, by reconstructing the input image, with the image caption as text prompt, and cropped objects as reference. At inference time, the model enables to directly generate images with cropped objects as reference.

**Experiment Settings.** Our experiments on COCO include two scenarios: conditional generation utilizing a single image-text context pair, and alternatively, employing multiple image-text context pairs, as shown in Figure 7.

**Training on single image-text context pair.** In this case, we randomly select an image-text pair from the COCO dataset. Initially, we extract all objects with their respective masks, collage them together, and then apply data augmentation, such as translation and rotation to create a composite reference image. The categories of the extracted objects serve as the reference text prompt candidates, while the captions of the initial image are employed as the text prompt candidates. If multiple candidates for the text prompt or reference text prompt exist, we will randomly choose one. We expect the model to reconstruct the original image based on these conditions, which enables the context module to learn how to leverage the given reference objects. We fine-tune our StoryGen on the train set of COCO2017 for 5,000 iterations. We only fine-tune the additional image cross-attention layers, and keep self-attention layers and text cross-attention layers frozen.

**Training on multiple image-text context pairs.** In this case, instead of collaging the objects from an image together, we

individually apply data augmentation to these objects, thereby generating several reference images for a single image-text pair. The categories of these objects are consistently selected as the corresponding reference text prompts. Note that, we employ same diffusion noise scales across these multiple reference images, given that they lack a temporal sequence and hold equal significance. Taking the model pre-trained for single image-text pair, we continue to fine-tune the image cross-attention layers for another 5,000 iterations, with all other parameters frozen.

## D.2. Quantitative Results

To measure consistency, we compute the similarity between the generated image and reference image with a pre-trained DINO [32] model. Specifically, we compare our StoryGen and variants with the original Stable Diffusion model on the validation set of COCO2017. Notably, the only difference between original SDM and our StoryGen here, is that StoryGen is augmented with additional image cross-attention layers trained on COCO train set, and all other parameters remain identical between the models. Thus, we utilize StoryGen to synthesize the original images with the reference objects, reference text prompts, and image captions, in both single and multiple image-text context pair scenarios; and we also exploit the original SDM to generate images with the image captions of current frames. For every image in COCO validation set, ten candidate images are generated. We use PickScore [20] to automatically identify those with better visual quality, then calculate the average DINO feature similarity between the ground truth images and the generated images for both StoryGen and original Stable Diffusion models.

The quantitative results in Table 5 demonstrate that our StoryGen model exhibits significant performance improvement in terms of consistency between the generated image and given references. Both StoryGen and StoryGen (Multiple) outperform the standard SDM. Moreover, compared to utilizing CLIP or BLIP features as visual conditions, our StoryGen model using diffusion-denoising features as condition demonstrates significant performance advantage, showing its effectiveness for retaining visual details from the reference image.

| Model                        | SDM    | StoryGen-CLIP | StoryGen-BLIP | StoryGen | StoryGen (Multiple) |
|------------------------------|--------|---------------|---------------|----------|---------------------|
| Consistency Score $\uparrow$ | 0.4804 | 0.5103        | 0.5475        | 0.7076   | <b>0.7317</b>       |

Table 5. Quantitative results on measuring consistency. SDM stands for standard stable diffusion models. StoryGen-CLIP and StoryGen-BLIP represent StoryGen with features extracted by CLIP and BLIP image encoder as context condition, as stated in our ablation study. StoryGen and StoryGen (Multiple) stand for StoryGen in single image-text context pair scenario and multiple image-text context pairs scenario, respectively. Consistency Score stands for the average DINO feature similarity, and the higher score yields better results.

## D.3. Qualitative Results

We provide more visualization samples of our consistency ablation on COCO in this section. The results of StoryGen are synthesized with the reference objects, reference text prompts, and image captions, in both single and multiple image-text context pair scenarios. As the original SDM cannot utilize reference images to exploit contextual visual information, the results of SDM are generated with the image captions only. The visualization results are all selected from the generated samples on the validation set of COCO2017.

**Single-pair COCO Visualization.** The qualitative results of our StoryGen on COCO, in single image-text context pair scenario, are depicted in Figure 9. Compared with the results of SDM, our model demonstrates obvious consistency in its generation results. Despite SDM can generate images that satisfy the text prompts, *i.e.*, good image-text alignment, the generated images are unable to maintain consistency with the reference image.

**Multi-pair COCO Visualization.** The qualitative results of our StoryGen on COCO, in multiple image-text context pairs scenario, are depicted in Figure 10. Despite multiple context pairs lead to more potential compositions, StoryGen is still able to generate high-quality images with the given reference objects, maintaining strong object consistency and semantic coherence, showing the ability of our architecture to exploit reference images for generation.

## E. Broader Impacts

Our storytelling model also has some positive impacts on the industry of creation and education: The widespread application of our visual storytelling model has the potential to inspire creators and artists to create a large number of visual storybooks rich in basic knowledge, which will have a profound impact on children’s early education, as demonstrated by related work in psychology. Our work draws inspiration from those open-source online libraries assisting children in underdeveloped regions, can potentially assist artists in creating educational storybooks tailored to these young learners.



## F. Limitations

The principal constraint of our storytelling model lies in the selection of stable diffusion models as its foundational architecture. Stable diffusion models are known to grapple with significant issues, notably the generation of images with inaccuracies in limb counts (such as legs, arms, or fingers) and decreased quality in the synthesis of images with multiple objects. Regrettably, our storytelling model inherits these limitations from the stable diffusion model. We anticipate addressing these shortcomings in future research endeavours by considering the adoption of more robust architectures, such as DALL-E 3, SD-XL [37], or consistency models [47].

## G. More Experiments

We will provide more quantitative evaluations and visualization samples of our *open-ended visual storytelling* in this section.

### G.1. Analysis on multi-frame conditioning

We conduct the multi-frame condition experiments on a test subset with 5,400 samples. As shown in Table 6, we find that: (i) Conditioning on previous frame is critical, (ii) the number of conditioning frames gives similar results. However, we do find differences in qualitative results, so we use the 3 closest frames as conditions in auto-regressive generation. As for frames with less than 3 previous frames, we use all previous frames instead.

|          | 0-frame | 1-frame      | 2-frame | 3-frame       |
|----------|---------|--------------|---------|---------------|
| FID ↓    | 40.29   | <b>32.34</b> | 33.27   | 33.65         |
| CLIP-I ↑ | 0.6841  | 0.7368       | 0.7419  | <b>0.7435</b> |

Table 6. Comparison of multi-frame conditioning.

### G.2. Multi-object conditioned Story Continuation

As illustrated in Figure 8, benefiting from our diverse StorySalon dataset and training strategy, our StoryGen model also demonstrates excellent performance on multi-object story continuation.

### G.3. Story Generation Visualization

We conduct a comparative analysis of results from StoryGen against those generated by SDM [42], Prompt-SDM [42], and StoryGen-Single. As illustrated in Figure 11, Figure 12, and Figure 13, the story generation results of our proposed models demonstrate significantly better consistency in style and character, as well as improved alignment between text and image.

While the results from SDM and Prompt-SDM are visually appealing, they exhibit a lack of stylistic and character consistency. The results of StoryGen-Single also display inconsistency, which proves that our StoryGen effectively utilizes the additional image condition to achieve consistency, rather than naive memorization.

### G.4. Story Continuation Visualization

We undertake another comparative analysis between the results from StoryGen and those from StoryDALL-E [31] and AR-LDM [36]. As depicted in Figure 14, Figure 15, and Figure 16, the story continuation results of our proposed models exhibit superior proficiency in maintaining style and character consistency, achieving stronger alignment between story and image, and enhancing image quality. Note that, all characters in the given reference image are unseen in our StorySalon datasets.

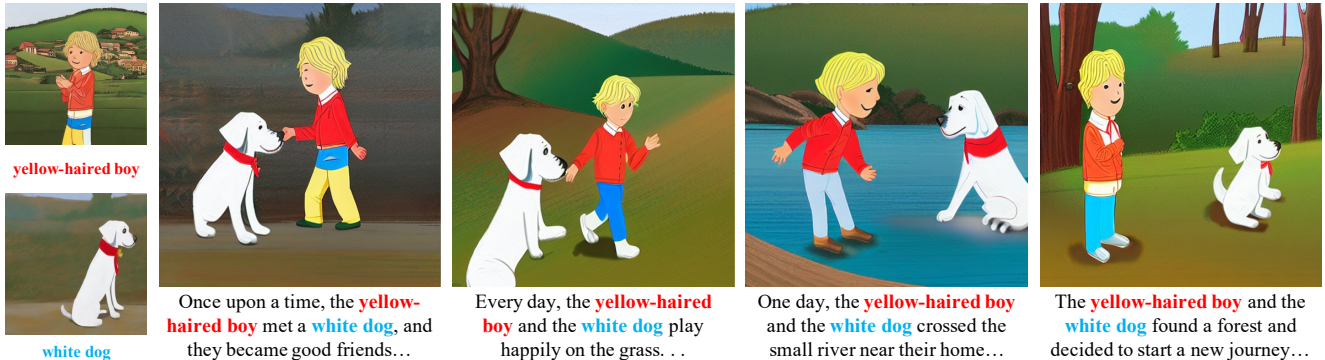
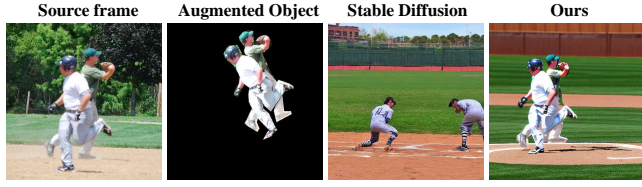


Figure 8. Example of Multi-object Story Continuation.



Single Object in COCO: Prompt = "Two baseball players are playing baseball on a field."



Single Object in COCO: Prompt = "A big purple public bus called south tyne."



Single Object in COCO: Prompt = "Two men shake hands at a formal dinner gathering."



Single Object in COCO: Prompt = "A large golden airplane is on the runway."



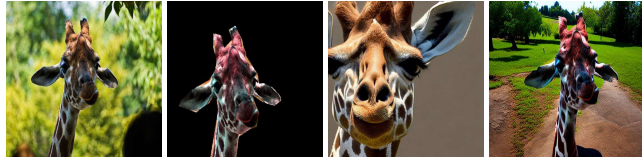
Single Object in COCO: Prompt = "Stylish man with blurred background."



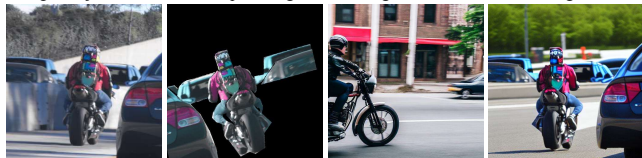
Single Object in COCO: Prompt = "A man letting his baby pet a horse."



Single Object in COCO: Prompt = "The fire hydrant in the green grass is red."



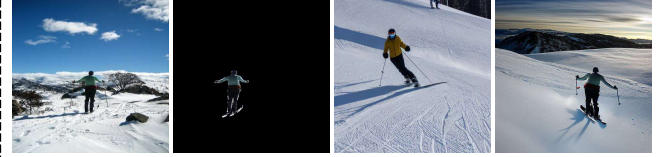
Single Object in COCO: Prompt = "A giraffe looking at the camera and making a face."



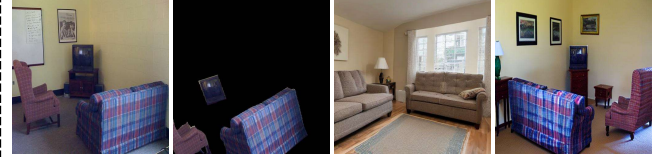
Single Object in COCO: Prompt = "A person riding a motorcycle down a street."



Single Object in COCO: Prompt = "Two apple computers are on a desk."



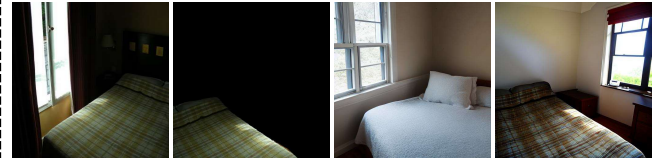
Single Object in COCO: Prompt = "A person is skiing on a snowy hill top."



Single Object in COCO: Prompt = "A living room with a couch and chair."



Single Object in COCO: Prompt = "A small group of sheep standing next to a building."



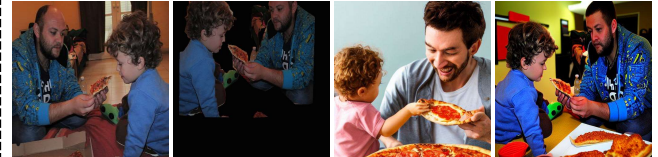
Single Object in COCO: Prompt = "The neatly made bed is beside an open window."



Single Object in COCO: Prompt = "A four compartment tray holding various food."



Single Object in COCO: Prompt = "A sheep standing on top of a rock."



Single Object in COCO: Prompt = "A man feeding pizza to a reluctant toddler."



Single Object in COCO: Prompt = "A little boy sticking his hand in a bowl of water."

Figure 9. Visualization results of StoryGen on COCO in single image-text context pair scenario.



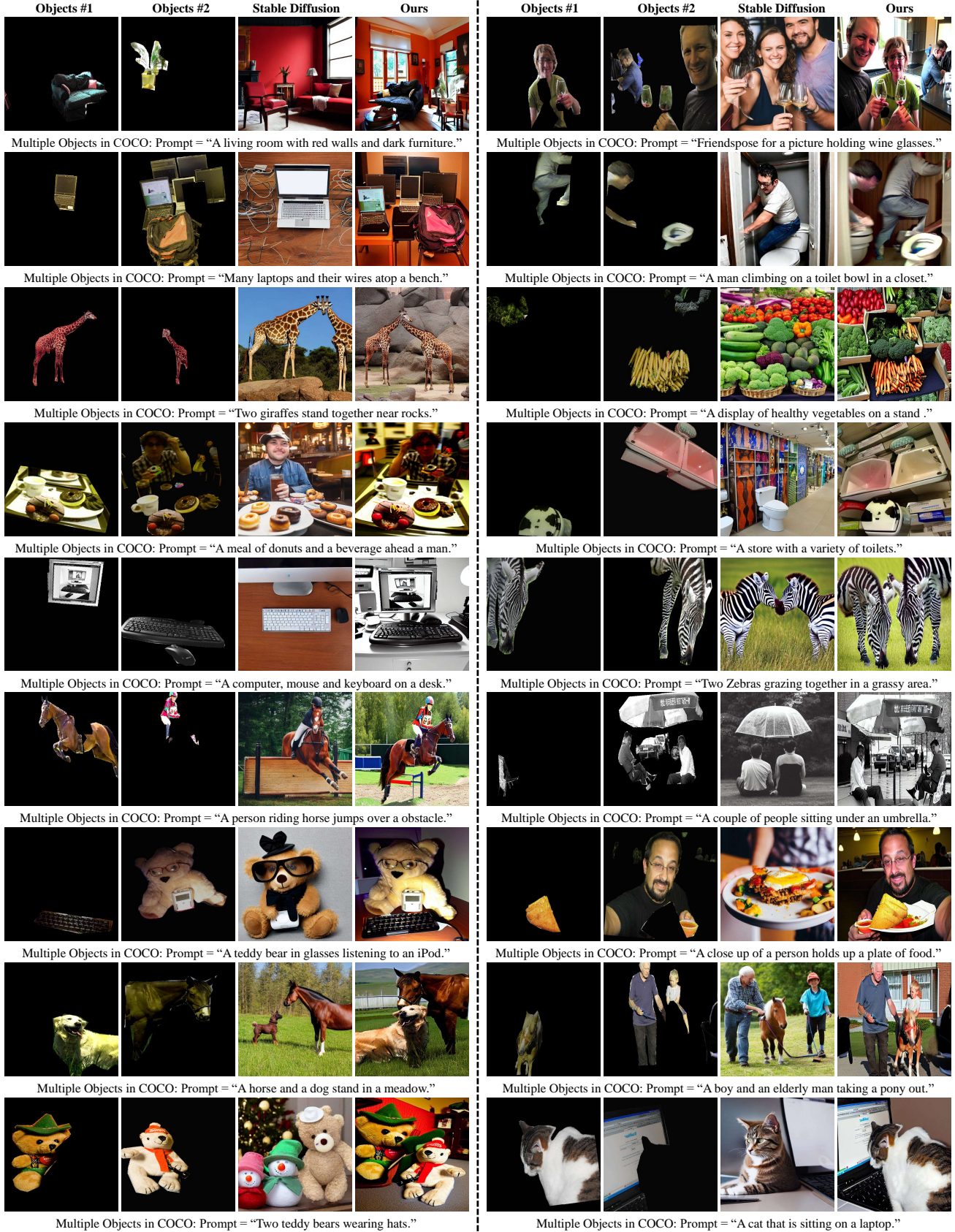


Figure 10. Visualization results of StoryGen on COCO in multiple image-text context pairs scenario.



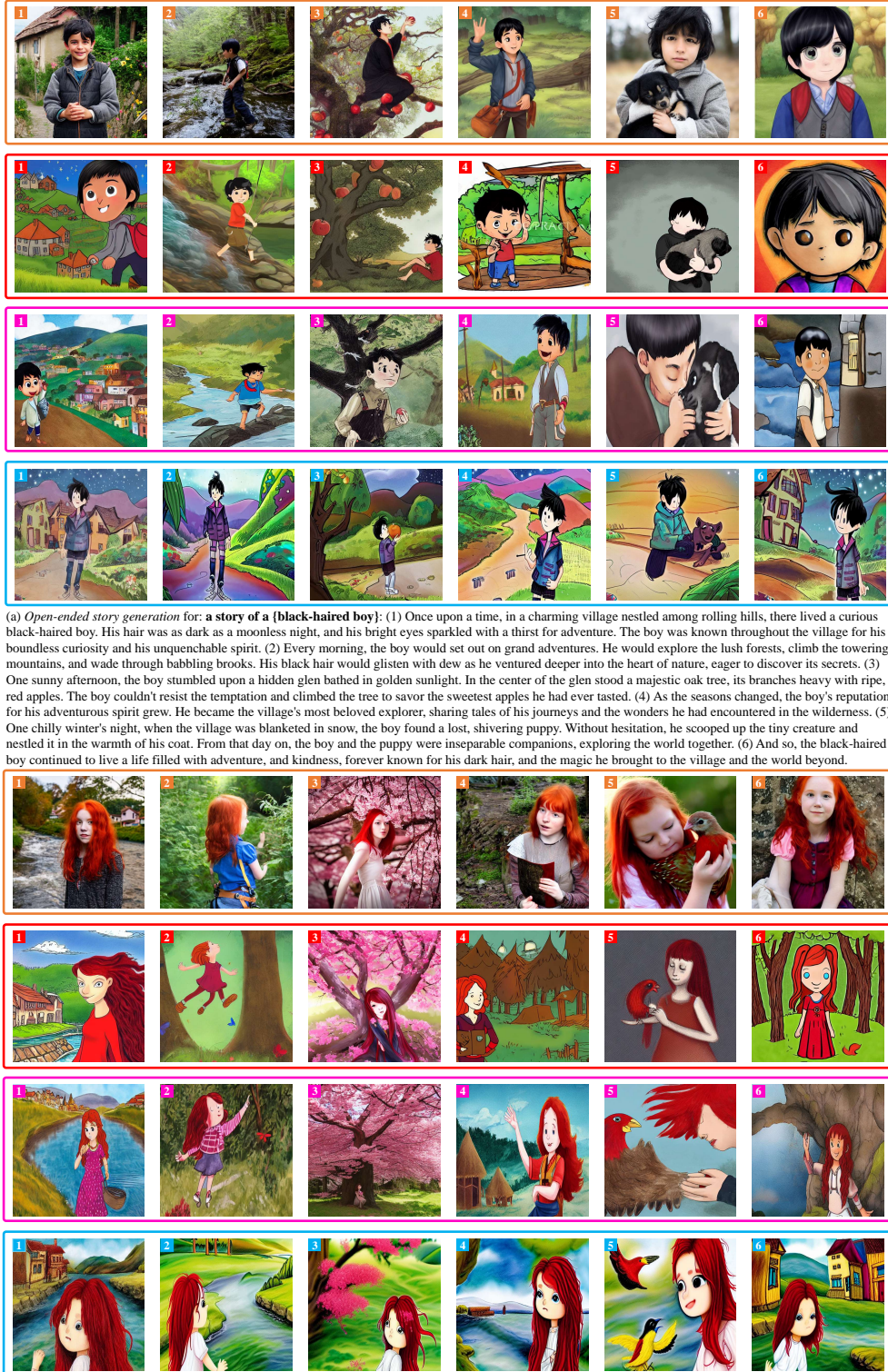


Figure 11. **Visualization results of Story Generation.** The images in orange, red, pink, and blue boxes are generated by SDM, Prompt-SDM, StoryGen-Single, StoryGen, respectively.





(a) *Open-ended story generation for: a story of a (black girl):* (1) Once upon a time, in a vibrant forest brimming with emerald leaves and silver streams, there lived a young black girl with hair as curly as the ferns and a smile as bright as the sun. She wore a dress spun from the purple petals of the wildflowers that grew in abundance around her tiny, ivy-covered cottage. (2) Every morning, the girl would step outside, her basket in hand, to gather the most colorful fruits the forest had to offer. She loved the way the morning dew made the world look like it was sprinkled with diamonds, and she'd often dance, twirling amidst the shimmering mist. (3) One afternoon, while exploring deeper into the forest, she discovered a clear, tranquil pond that mirrored the sky so perfectly it seemed as if it held the clouds and the sun within its depths. Here, she would sit and daydream, tossing pebbles to create ripples that carried her thoughts to the stars. (4) As the seasons changed, so did the forest, and the girl witnessed the leaves painting themselves in oranges and reds. She collected these leaves, pressing them into a book, creating a mosaic of memories, each leaf a reminder of the day's joy and wonder. (5) When winter whispered in, the girl found beauty in the silence of the forest covered in snow. She built sculptures of snow, each one more fantastical than the last, giving life to the winter's quiet. (6) Spring brought a chorus of blooms, and the girl, now with a crown of flowers in her hair, joined in the celebration, planting seeds that she had gathered, ensuring the cycle of growth and beauty continued. (7) The girl grew, and with each passing year, she learned the secrets of the trees, the whispers of the wind, and the dance of the seasons. In harmony with the forest, she became its guardian, a symbol of the enduring dance of life.



(b) *Open-ended story generation for: a story of a (black wolf):* (1) In the depths of a snowy wilderness, there was a solitary black wolf whose coat shimmered against the stark white of the frozen landscape. The wolf had a majestic presence, with eyes that glinted like the first stars of the evening sky. His powerful paws left a trail of footprints as the only evidence of his passage through the thick blanket of snow. (2) Each day, the black wolf would climb to the peak of a great mountain, letting out a deep, resounding howl that echoed through the valleys. The sound would carry for miles, a song of strength and solitude that resonated with the whispering pines and the crisp, winter air. (3) With the arrival of spring, the black wolf watched as the snow melted, revealing a carpet of wildflowers. He roamed through the blossoming terrain, his black fur juxtaposed against the riot of colors, a guardian of the waking world. (4) Summer brought with it an abundance of life, and the black wolf would spend his nights chasing the golden orb of the moon, racing through forests where the fireflies lit his path, a living embodiment of the night's spirit. (5) When autumn arrived, the wolf found joy in the crunch of the leaves beneath his paws. The forest was a cascade of oranges, reds, and yellows, and he moved through it like a shadow, part of the tapestry of the changing seasons. (6) The wolf, now older and wiser, took to resting by a tranquil lake during the quiet afternoons, reflecting on the cycles of nature. The calm waters mirrored his noble figure and the fiery sunsets, a scene of perfect peace and solitude. (7) As years passed, the black wolf became a legend of the wilderness, a solitary figure that moved with the grace of the seasons. His story was written in the earth, a tale of harmony with the world, a spirit both wild and free.

Figure 12. **Visualization results of Story Generation.** The images in orange, red, pink, and blue boxes are generated by SDM, Prompt-SDM, StoryGen-Single, StoryGen, respectively.





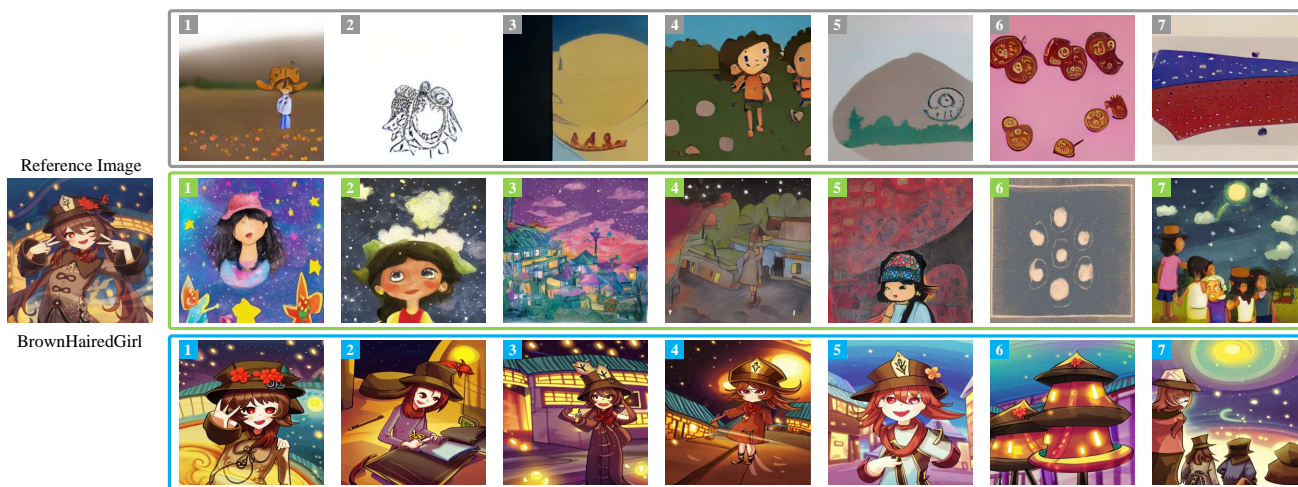
(a) *Open-ended story generation for: a story of a (princess):* (1) Once upon a time, in a land ruled by the rhythm of seasons, there lived a beautiful princess who was in love with the stars. Her castle was adorned with towers that reached towards the heavens, and every night, she would climb to the highest balcony to converse with the twinkling dots scattered across the night sky. (2) One evening, as a shooting star sliced through the darkness, the princess made a wish. She longed to visit the stars and dance among them, to learn their ancient secrets and to see the world from their eternal vantage point. (3) The next morning, the princess discovered a mysterious, silvery seed had fallen from the sky and landed in her royal garden. She planted the seed in the earth, watering it with water from the enchanted spring that ran through the castle grounds. (4) To the princess's wonder, the seed grew rapidly, unfurling into a magnificent vine with leaves that shimmered like stars and flowers that glowed with the luminescence of the moon. The vine spiraled up one of the castle towers, beckoning the princess to climb. (5) With her heart pounding with excitement, the princess began to ascend the vine. As she climbed higher, the air grew thinner, and the sky seemed to embrace her. The vine ended at the threshold of the stars. (6) The princess stepped off the vine and found herself walking on a path of stardust. Each step she took was lighter than the last, and she danced among the stars, just as she had wished. They whispered to her in the language of light, sharing stories of distant worlds and the dance of the cosmos. (7) As dawn approached, the princess knew it was time to return to her own world. She glided back down the stardust path, down the vine, and stepped onto her balcony as the first rays of sunlight kissed the horizon. Her heart was full of starlight, and her eyes shone with the reflection of her night among the stars.



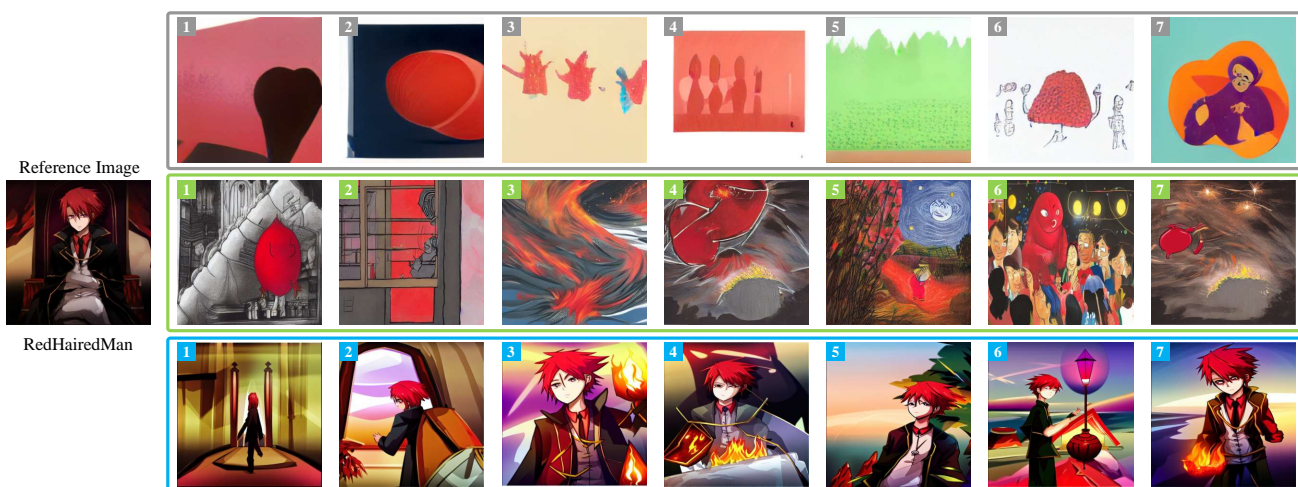
(b) *Open-ended story generation for: a story of a (prince):* (1) Once upon a time, in a kingdom draped in emerald valleys and crystal lakes, there lived a young prince who was fascinated by the secrets of nature. The prince spent his days wandering the expansive gardens of the palace. He was particularly enchanted by a single rosebush that grew at the edge of the garden. (2) One day, as the prince watched the rosebush, he noticed that it began to wilt despite his careful attention. Worried, he consulted the ancient library in his castle, scouring old texts and botany books for a cure. He learned of a rare water source, hidden deep in the forest, whose waters were said to rejuvenate any plant. (3) Determined to save his cherished rosebush, the prince set out alone into the forest. He traversed through thickets and over streams, guided by the chirps of birds and the rustling of leaves. His journey was long, and the forest seemed to whisper secrets as he passed. (4) Finally, after several days, he arrived at a glade where the sunlight shimmered down like warm gold. There, at the center, was a spring that sparkled with water so clear it looked like liquid diamonds. The prince filled his flask with the water, feeling its coolness and vitality. (5) On his way back to the castle, the prince encountered a variety of creatures. A wise owl nodded at him from a tree branch, a family of rabbits watched curiously from the bushes, and a graceful deer bowed its head as he passed. The prince realized he was not alone in his quest; the forest itself was guiding and protecting him. (6) Upon his return, the prince immediately watered the ailing rosebush with the magical spring water. Overnight, the rosebush regained its vigor, its petals unfolding with colors so vivid and fragrant that they seemed to glow in the moonlight. (7) The prince's dedication to his rosebush became a legend in the kingdom. He went on to create the most magnificent garden, filled with plants and flowers from all over the world, each thriving under his care. The prince became known not just as a ruler, but as a guardian of nature, with a garden that was a testament to his love for all living things.

Figure 13. Visualization results of Story Generation. The images in orange, red, purple, and blue boxes are generated by SDM, Prompt-SDM, StoryGen-Single, StoryGen, respectively.



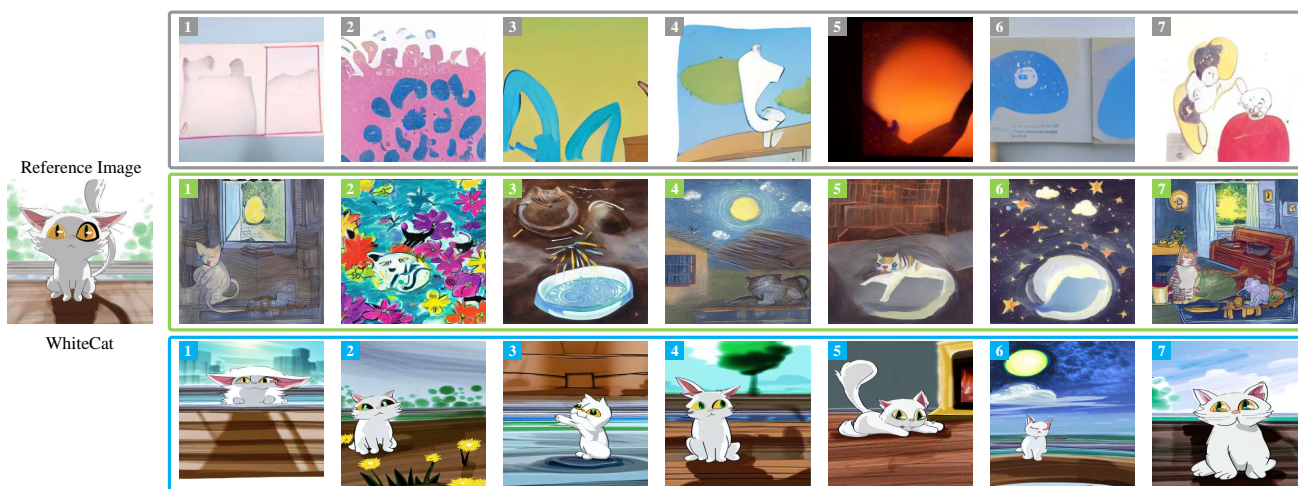


(a) *Open-ended story continuation for: a story of a {a brown-haired girl}*: (1) Once upon a time, in a village where the night sky always shimmered with stars, there lived a brown-haired girl known for her enchanting hats. She wore a unique hat every day, each adorned with symbols of dreams and hopes, like flowers and stars. This girl had a special gift; she could weave the essence of the night sky into her hats, making them glow with a soft, comforting light. (2) One evening, as the crescent moon hung low, the girl decided to craft a hat that would capture the beauty of the night forever. She worked under the starlight, her fingers dancing with threads that glimmered like constellations. She poured her joy and the whispers of the night wind into her creation, a hat that seemed to hold the entire galaxy within its folds. (3) When she placed the hat upon her head, the village was bathed in a celestial glow. The villagers looked out from their homes in awe, as gentle light cascaded down the cobblestone streets. They felt peace and wonder, and the night seemed less dark, less fearful, as if the girl's hat had captured the heart of the night itself. (4) The girl realized that her gift was more than just crafting hats; she could bring hope to her people. She began to wander the streets each night, her hat glowing ever so brightly, a beacon for those who were lost or in need of comfort. She didn't need words; her presence was enough to lift spirits. (5) As the seasons changed, so did her hats. The girl and her hats became a symbol of the village's identity, a reminder that beauty and magic could be found in the simplest of things. (6) Years passed, and the girl grew older, but her legacy remained timeless. The hats she had crafted were passed down through generations, each one a treasured heirloom that continued to glow with a piece of the night. (7) And so, the brown-haired girl with her magical hats lived on in the hearts of the people, a legend woven into the fabric of the village. They would look up at the stars and remember the girl who walked with the night, who showed them that even in darkness, there is light to be found.

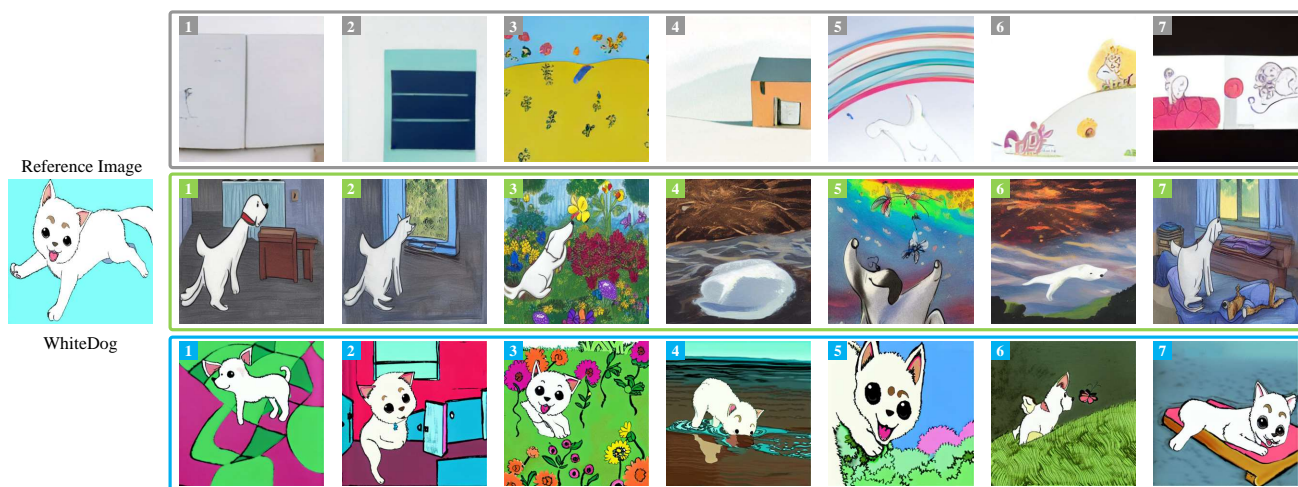


(b) *Open-ended story continuation for: a story of a {a red-haired man}*: (1) In a land where colors held magic, there lived a red-haired man who was the guardian of the Flame of Creation. His hair, the color of burning embers, was a symbol of the fire that he protected — a fire that had the power to ignite inspiration and passion in the hearts of the people. (2) Despite the grandeur of his task, the red-haired man felt a growing sense of solitude. His life was an endless cycle of tending to the flame. He longed to see the world beyond his hall, to witness the wonders his flame brought to life. Yet, he dared not leave, for the flame required constant care, and there was no one else to shoulder the burden. (3) One evening, as shadows danced along the walls of his hall, the red-haired man noticed that the flame flickered unusually. It whispered to him of a possibility he had never considered. The flame could divide, sharing a spark that could be carried out into the world without letting the original fire die. (4) With a mixture of trepidation and excitement, the red-haired man fashioned a lantern from the hall's curtains and an old chair. He captured the wayward spark in the lantern, ensuring that the main flame continued to burn strong. Now he held a piece of the Flame of Creation, a portable spark that would allow him to venture into the world. (5) As the red-haired man stepped outside, the lantern's glow seemed to brighten the world in hues he had only imagined. Wherever he walked, life sprang forth: flowers bloomed, trees bore fruit, and the night sky shimmered with new stars. (6) In time, the red-haired man realized that while the flame's magic was powerful, it was the actions of the people that truly created change. The man's solitude was replaced by a sense of connection to the world, fulfilled by the knowledge that the flame's inspiration was at the heart of all creativity. (7) The guardian of the Flame of Creation returned to his hall, understanding now that his duty was not only to protect the flame but also to share its gift. He would continue to venture out into the world, carrying the spark that brought dreams to life.

Figure 14. **Visualization results of Story Continuation.** The images in gray, green, and blue boxes are generated by StoryDALL-E, AR-LDM, StoryGen, respectively.



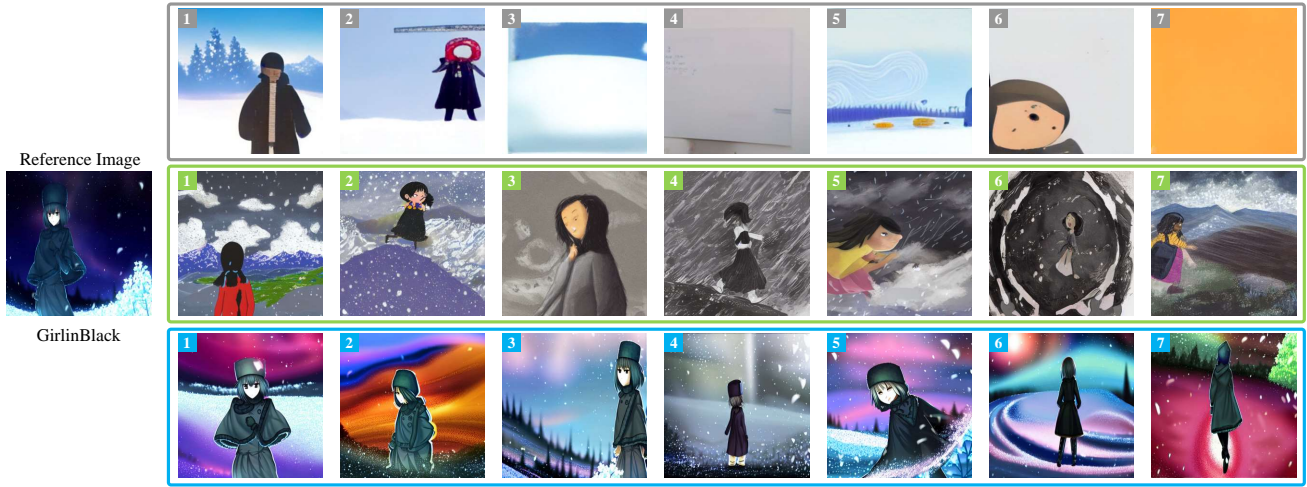
(a) *Open-ended story continuation for: a story of a {a white cat}*: (1) Once upon a time, in a cozy little house bathed in sunshine, there lived a small white cat with big, bright eyes that sparkled like gems. Every morning, the cat would find a sunbeam streaming through the window, jump onto the sill, and bask in the golden glow, feeling the warmth on its soft fur. (2) One afternoon, while the cat was lounging in its favorite spot, a gentle breeze carried in the sweet scent of flowers from the garden. The cat, curious and playful, leapt off the sill and followed the fragrance, finding itself in the midst of a colorful flower bed, with petals dancing in the wind. (3) In the garden, the cat discovered a little fountain, its water sparkling in the sunlight. With a paw, the cat gently tapped the surface, causing ripples to spread across the water. Fascinated, the cat watched the ripples, feeling as if the water was playing a game with it. (4) As the day turned to dusk, and the sky painted itself in hues of orange and pink, the cat climbed atop a wooden fence, gazing at the setting sun. There, it sat, silhouette outlined against the sky, feeling the cool evening breeze ruffling its fur. (5) When the stars began to twinkle in the night sky, the white cat returned inside, finding a cozy spot in front of the fireplace. The flickering flames cast warm, dancing shadows, and the cat curled up, closing its eyes, the warmth lulling it into a peaceful sleep. (6) Throughout the night, in its dreams, the cat journeyed to fantastical places, flying among the clouds and walking on the moon, each dream more vivid and wondrous than the last, until the soft rays of dawn signaled a new day. (7) And so the days went by, with the little white cat finding joy in the simplest things: a ray of sunshine, the scent of flowers, the playfulness of water, the serenity of sunsets, the warmth of fire, and the boundless adventures found in dreams.



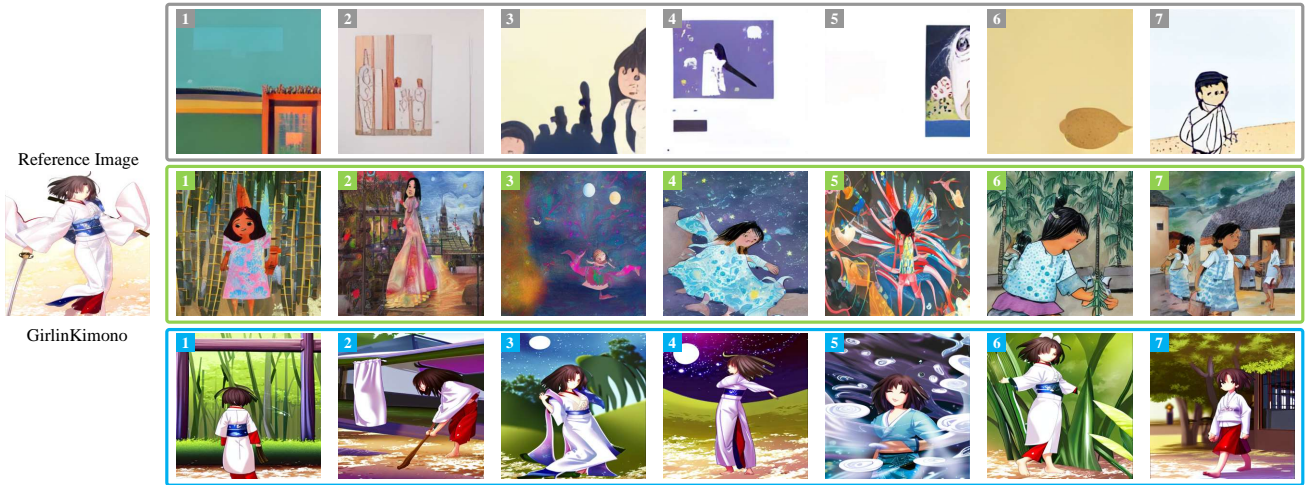
(b) *Open-ended story continuation for: a story of a {a white dog}*: (1) Once upon a time in a colorful world, there was a small white dog with playful spots over its fur. This curious little dog loved to explore every corner of its bright and cheerful home. It would scamper around with boundless energy, its tongue lolling out in a happy pant, and its tail wagging like a fluffy pendulum. (2) One sunny morning, the white dog found a mysterious blue butterfly fluttering near the window. It was unlike any butterfly it had seen before. The dog tilted its head, eyes wide with wonder, and decided to follow the butterfly wherever it might lead. With a joyful bounce, the dog leapt towards the fluttering creature, embarking on a new adventure. (3) The butterfly led the dog through the garden where flowers bloomed in every shade imaginable. The dog marveled at the sights, sniffing the fragrant air filled with the scent of fresh blooms. It chased the butterfly around the garden, over the green grass, under the flowering bushes, and around the stone path. (4) As they journeyed together, the white dog and the butterfly came across a clear, bubbling stream. The dog had never seen water so clear, and it watched in amazement as the sunlight danced upon the water's surface. Feeling adventurous, the dog dipped its paws into the cool stream, sending ripples across the water. (5) Suddenly, the butterfly soared up high, with the white dog gazing after it. The dog noticed a rainbow arching across the sky, its colors reflecting the vibrant world below. The dog felt a surge of joy and decided to race along the stream, as if it were racing the colors of the rainbow. (6) The day turned to evening, and the sky painted itself with the hues of sunset. The white dog found itself on a hill, watching the sun dip below the horizon. The butterfly landed gently on the dog's nose, as if to say goodbye. The dog sat peacefully, feeling grateful for the day's journey and the beauty it had seen. (7) As the stars began to twinkle in the night sky, the white dog returned home, its heart full of the day's wonders. It curled up in its cozy bed, dreaming of the gardens, the stream, the rainbow, and the butterfly. The white dog knew that tomorrow was another day for adventure, but for now, it rested, wrapped in the warmth of its memories.

Figure 15. **Visualization results of Story Continuation.** The images in gray, green, and blue boxes are generated by StoryDALL-E, AR-LDM, StoryGen, respectively.





(a) *Open-ended story continuation for: a story of a {a girl in black}*: (1) In a realm where winter reigned eternal, the girl in the black coat walked alone, her presence the only warmth in the icy world. She was the Whisperer of the Wind, a gentle spirit who could speak to the cold breezes and soothe their icy fury. (2) Each morning, as the sun struggled to pierce the wintry gloom, she would climb the highest hill and listen to the stories the wind told. Tales of distant lands, of sun-soaked shores, and of children playing under the warmth of a softer sun. (3) One day, the wind spoke of a coming storm, a tempest that could bury her world in snow and silence forever. The Whisperer knew she had to calm the storm's heart, or springtime's hope would never return to her frozen home. (4) With courage in her step, she walked into the heart of the storm, her black coat fluttering like a banner of night. She spoke to the blizzard, her voice a melody that rivaled the storm's howl, a plea for peace and tranquility. (5) The storm, taken aback by the girl's bravery and the sweetness of her voice, began to lessen its wrath. Snowflakes slowed their dance, and the icy gales held their breath, listening to the Whisperer's song. (6) As the storm's heart calmed, the snow ceased, and the winds carried the girl's song across the land. Wherever her voice reached, ice melted, revealing the first glimpses of the soil beneath—a promise of the spring to come. (7) The girl in the black coat became the legend of the winter world, the one who conversed with the wind and turned the fiercest of storms into a peaceful slumber. And though she wandered alone, her song of warmth and the hope of spring lived on in the hearts of all those who yearned for the thaw.



(b) *Open-ended story continuation for: a story of a {a girl in kimono}*: (1) In the heart of a dense bamboo forest, there stood a solitary shrine, its red torii gate a stark contrast against the sea of green. It was here that the girl in the white kimono found solace and purpose. She was the shrine's keeper, tasked with ensuring that the balance between the human realm and the spirits was maintained. (2) Each morning, with the first light of dawn casting a soft glow over the land, the girl would sweep the shrine's grounds with a handmade broom, her white kimono glimmering in the sun's gentle rays. She took great care in her work, for she knew that cleanliness was a gesture of respect to the spirits. (3) It was during the night of the full moon that the girl's responsibilities took on a magical turn. The air would thrum with energy, and the border between worlds grew thin. On these nights, the girl would perform a sacred dance, a ritual to honor the spirits and ensure their goodwill towards the villagers. (4) With each precise step and wave of her sleeve, the girl's dance would draw luminous orbs from the moonlit sky, each one a spirit coming to witness her devotion. The orbs hovered around her, pulsating with the serene energy of the unseen world. (5) As the dance reached its crescendo, the spirits would begin to swirl around the girl, creating a vortex of otherworldly light. The girl's connection to the spirits was strongest at this moment, and she would whisper her wishes for the village's safety and prosperity. (6) When the dance ended and the first light of dawn approached, the spirits would depart, leaving behind a trail of sparkling dew on the bamboo leaves. This dew was said to have healing properties, and the girl would collect it carefully, a gift from the spirits to the villagers. (7) The villagers rarely saw the mysterious events that took place at the shrine, but they felt the peace and prosperity that the girl's rituals brought to their lives. The girl in the white kimono remained ever vigilant, a silent guardian whose dance with the spirits kept their world in harmony.

Figure 16. **Visualization results of Story Continuation.** The images in gray, green, and blue boxes are generated by StoryDALL-E, AR-LDM, StoryGen, respectively.

### G.5. Failure Case Visualization

Figure 17 presents some instances where StoryGen did not perform optimally. These failure cases primarily stem from the inherent limitations of SDM. Figures (a), (b), and (c) illustrate occurrences where StoryGen is prone to generating images with limb count inaccuracies, such as incorrect numbers of legs. Figures (d) and (e) show scenarios where the generation of multiple objects results in each object being of subpar quality. Figures (e), (f), and (g) depict instances of StoryGen producing low-quality human faces. Regarding Figure (h), despite the visual prompt being "A *black wolf walking through a forest with autumn leaves falling*", the generated image erroneously includes snowfall, due to the winter setting of the reference image. This discrepancy arises from the conflict between the image and text conditions.

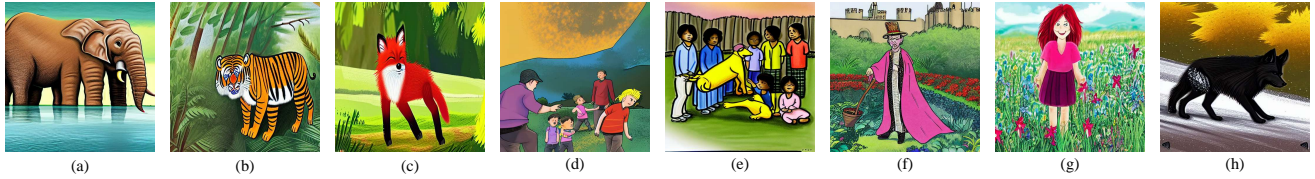


Figure 17. Some failure cases of StoryGen.