

# GraPE 🍇: A Generate-Plan-Edit Framework for Compositional T2I Synthesis

Ashish Goswami<sup>1</sup>, Satyam Kumar Modi<sup>1\*</sup>, Santhosh Rishi Deshineni<sup>1\*</sup>,  
Harman Singh<sup>1</sup>, Prathosh A. P<sup>2</sup>, Parag Singla<sup>1</sup>  
<sup>1</sup>IIT-Delhi, <sup>2</sup>IISc Bangalore

ashish.goswami@scai.iitd.ac.in, parags@cse.iitd.ac.in

## Abstract

Text-to-image (T2I) generation has seen significant progress with diffusion models, enabling generation of photo-realistic images from text prompts. Despite this progress, existing methods still face challenges in following complex text prompts, especially those requiring compositional and multi-step reasoning. Given such complex instructions, SOTA models often make mistakes in faithfully modeling object attributes, and relationships among them.

In this work, we present an alternate paradigm for T2I synthesis, decomposing the task of complex multi-step generation into three steps, (a) Generate: we first generate an image using existing diffusion models (b) Plan: we make use of Multi-Modal LLMs (MLLMs) to identify the mistakes in the generated image expressed in terms of individual objects and their properties, and produce a sequence of corrective steps required in the form of an edit-plan. (c) Edit: we make use of an existing text-guided image editing models to sequentially execute our edit-plan over the generated image to get the desired image which is faithful to the original instruction. Our approach derives its strength from the fact that it is modular in nature, is training free, and can be applied over any combination of image generation and editing models. As an added contribution, we also develop a model capable of compositional editing, which further helps improve the overall accuracy of our proposed approach. Our method flexibly trades inference time compute with performance on compositional text prompts. We perform extensive experimental evaluation across 3 benchmarks and 10 T2I models including DALLE-3 and the latest – SD-3.5-Large. Our approach not only improves the performance of the SOTA models, by upto 3 points, it also reduces the performance gap between weaker and stronger models. The code is available at <https://dair-iitd.github.io/GraPE/>

## 1. Introduction

Diffusion models [45], have revolutionized the space of generating photo-realistic images with multiple works

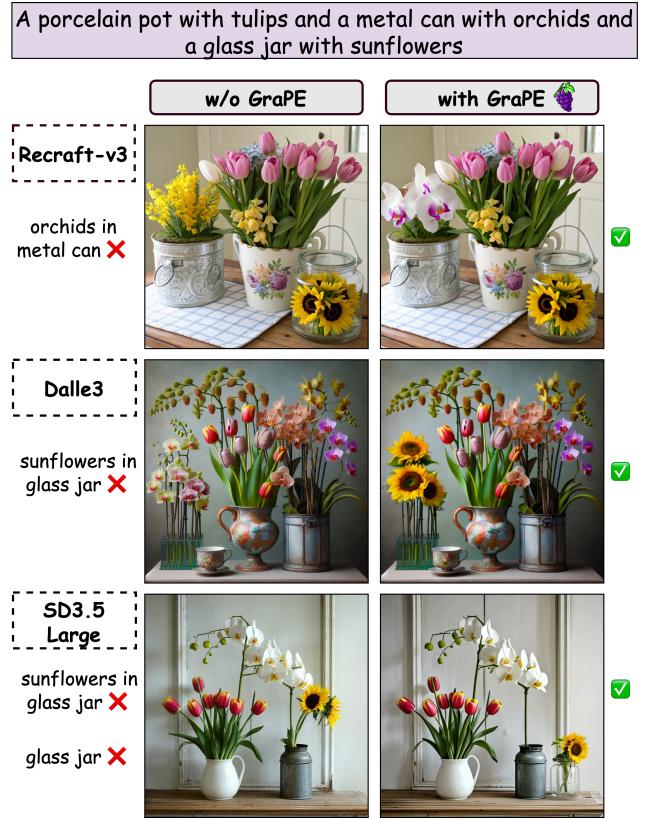


Figure 1. Illustration of GraPE’s capability to align the image with the input prompt. *Left*: Generations from various state-of-the-art diffusion models [5, 41, 46], along with their inaccuracies. *Right*: Images produced by GraPE, in a completely training free manner.

showing their superior performance compared to their older competitors such as GANs, VAEs and Flow-based models [14, 34, 43]. Of specific interest has been their ability to generate images from text. Despite the success, it has been observed that SOTA diffusion models still struggle to generate accurate images for instructions which involve multi-step compositional reasoning, often resulting in errors over object attributes as well as their interactions with other ob-

jects in the image [9, 26, 40]. Figure 1 shows an example of a text-prompt and the corresponding images generated by 3 of the SOTA text-to-image generation models. Several reasons have been pointed for the relative performance on such complex instructions, including lack of appropriate training data, and use of weak text-encoders to name a few. This has naturally limited their applicability for the task of automated and reliable image generation in real life scenarios. To make matters worse, image-editing has been left even further behind, with existing models hardly capable of correctly handling complex edit tasks.

To overcome these challenges, recent work has seen several credible attempts to address the challenge of generating images faithful to complex instructions. Broadly, these can be divided in two categories (1) those requiring fine-tuning of existing generation models, examples include [8, 24, 27, 51, 54] (2) and those which are based only on adapting the inference procedure, examples include [1, 9, 16, 31]. While these approaches have been able to show some improvements over the base diffusion models, the problem remains far from solved. We survey the existing related works in detail in Section 2.

In this work, we present a novel approach which is based on the following observations: Instead of performing the image synthesis for a complex compositional instruction in one go, we can redefine the task into iterative refinement of a partially faithful generated image. Our approach can be described as a pipeline of the following 3 steps: (a) *Generate*: Initial generation using existing T2I models (b) *Plan*: Identification of errors in the form of a series of simple edits (largely over individual objects and their relations) required to fix the original generation; the edits are expressed in the form of a *sequential editing plan* generated via the use of a MLLM (c) *Edit* carrying out the edit plan in a sequential manner over the original generation to get the desired image. We refer to our approach as *GeneRAte-Plan-Edit* (GraPE). To the best of our knowledge, we are the first ones to decompose the task of fine-grained generation as a series of simpler (atomic) edits which are preceded by a partially correct generation.

We perform a series of experiments to evaluate the efficacy of our approach over three benchmark datasets. We compare with 10 different T2I models, and show improvement in each case, with improvements up-to a significant 20+ points in some cases. We also note that we see a more significant improvement in weaker (typically smaller) generation models, resulting in narrowing of performance gap between smaller and larger SoTA models. We also perform ablations showing the value of our object-centric approach for generating edit plans. Along with this, significant boost in performance is shown when we use the compositional image editor that we develop, as part of our pipeline. We also present detailed error analysis and point to limitations

of our model in our experimental analysis.

The contributions of our paper can be summarized as follows: (a) We propose a simple yet effective approach for improving the performance of T2I models via a generate-plan-edit pipeline. (b) Our approach is guided by object-centric edit plans, and exploits the power of editing models which can handle compositions well. (c) Our approach is modular, and can be used with any existing base generation or editing model (d) We perform a series of experiments demonstrating the efficacy of our approach over a large number of T2I models, and also present detailed insights into the working of model.

## 2. Related Work

**Compositional Image Generation:** T2I models struggle with compositional image generation, such as generating objects with correct attributes, adhering to object relations, etc. [4, 17, 21]. Several works aim to improve the compositional image generation capabilities of T2I diffusion models. Some prior works intervene on the cross-attention maps in diffusion models which are responsible for incorporating information from the text into the image. [1, 9] strengthen the attention activations of neglected concepts in the image during the generative process using inference time losses, primarily targeting the reduction of catastrophic neglect and incorrect attribute binding. [27, 51] proposes training models focused on alignment of Cross-attention maps for each token with pseudo-ground truth segmentation maps. [16] uses linguistic structures and structured representations of text such as constituency tree, for manipulating cross-attention representations. However, all these works often focus on narrow compositional concepts, such as only attribute binding or reducing catastrophic neglect of objects, than being a general purpose method for aligning images with text. [50] uses MLLMs as an agent with an extensive set of tools for decomposing and planning the generation process into multiple steps. [57] uses an LLM for planning the generation process and suggesting sub-regions for region-wise diffusion, however, it may struggle with either with decomposing into or merging of these regions, since their re-weighting and sampling method is different from what diffusion generative models are trained for. [12, 18, 31, 35] utilize LLMs for generating intermediate representations such as layouts, panels or blobs followed by generating using specialized models that require such complex annotations to be trained. [11, 19, 32, 56, 64] utilize LLMs such as Gemma, LLaMA, T5-XXL [39, 47, 49] as text-encoders for improved text representations compared to contrastively trained models like CLIP [38] that are largely agnostic to word order [28, 59]. Orthogonally, we use MLLMs for planning the sequential editing process based on the initial generated image, and use general purpose generator and editor models in their native fashion.

**Instruction Guided Image Editing:** Editing images with human written instructions have seen significant interest in the community [6, 15, 29, 44, 60–62]. Among recent work, InstructPix2Pix [6] leverages a training free editing model, Prompt-to-Prompt [22] along with GPT3 [7] to produce paired image data with editing instructions. MagicBrush [60] is a manually annotated real world image editing dataset that improves InstructPix2Pix. [29] release Aurora, for action and reasoning centric image editing. In our work we explore the use of general-purpose image editors for aligning images generated via T2I models to their complex text prompts. We also develop a compositional image editor leading to further improved performance.

### 3. Method

In this section, we outline our *Generate-Plan-Edit* (GraPE) framework. Let us introduce some notation. Let  $T$  be a textual instruction. In the task of T2I synthesis, given an instruction  $T$ , our goal is to be able to generate an image  $I_o$  which satisfies various requirements expressed via the instruction  $T$ . While most existing techniques take the approach of directly generating  $I_o$  via  $T$ , they often result in various kinds of inaccuracies, due to the complexity of the instruction. We are motivated by the observation that the task of T2I synthesis can be broken down into simpler steps of first generation, followed by identification of errors, and a sequence of corrective edits, each of which is simple and object specific in nature. Accordingly, we propose the following generation pipeline.

We first generate an initial image  $I_g$  using a SOTA generative model,  $\mathcal{G}$ . This  $I_g$  may have mistakes, or inaccuracies with respect to the intent expressed in  $T$ . We then make use of an existing MLLM ( $\mathcal{P}$ ) to identify the mistakes in  $I_g$ , as a difference between image and textual descriptions in  $T$  and  $I_g$ , respectively, for each object of interest, and corrective steps suggested, in the form of an edit plan expressed as  $(T_{e_1}, T_{e_2}, \dots, T_{e_n})$ , where each  $T_{e_k}$  is an edit instruction fixing some aspect of the  $I_g$  so that it can be aligned with the original prompt  $T$ . Note that the number of edits is instance specific, and is a part of the output produced by  $\mathcal{P}$ . Finally, we make use of an editing model  $\mathcal{E}$  which inputs the current image  $I_{e_k}$  along with an edit instruction,  $T_{e_{k+1}}$  and produces next image in the sequence  $I_{e_{k+1}} = \mathcal{E}(I_{e_k}, T_{e_{k+1}})$  with  $k \in \{1, 2, \dots, n-1\}$ . Note that  $I_g = I_{e_0}$  and  $I_{e_n} = I_o$ . This algorithm comprises of 3 broad steps: (a) *Generate*: Generate the image for textual instruction  $\mathcal{T}$  (b) *Plan*: Identify mistakes and propose a correction plan (c) *Edit*: Perform the sequence of edits based on the corrective plan. We now describe each step in detail.

#### 3.1. Image Generation

For our framework’s initial step, we create a base image  $I_g$ , which serves as a foundational input for the subsequent

processes involved in planning and editing. This image  $I_g$  and its corresponding text-prompt  $T$  are considered as primary inputs for the next steps, which allows our framework to operate in a plug-and-play manner, regardless of the T2I model used in generation.

#### 3.2. Multi-Modal Planner

The planner is a key component of GraPE, with its role being executed by a Multi-Modal Large Language Model (MLLM) that assesses the initial image,  $I_g$ , for potential object misalignment relative to the textual description provided in the prompt,  $T$ . This discrepancy analysis acts similarly to the Chain-of-Thought prompting technique [52], designed to enhance reasoning by encouraging the planner to break down the assessment process into object level steps.

**Prompting Style:** The prompting style used with the MLLM play a crucial role in its ability to generate clear and concise plans across images from diverse domains. Unlike existing approaches [31, 57], which typically rely on LLMs to create abstract or structured representations—such as bounding boxes, layouts or blobs—either in a few-shot manner or through explicit training on annotated data [63], our approach takes a streamlined alternative that bypasses the complexity of these strategies.

Our method focuses on leveraging the inherent strengths of MLLMs, focusing on tasks where they already excel, such as image captioning and language comprehension, which reduces the need for intricate representations enabling our framework to maintain simplicity and generalization across domains. By combining these strengths with carefully designed few-shot examples, we guide the planner to produce object-centric, structured outputs without necessitating additional, specialized training.

The planner’s output is organized into four key steps, explained below using the example in Fig 2.

- **Analyzing Textual Elements:** The goal here is to extract high level object-attribute pairs from the text-prompt,  $T$  focusing on relationships among the described objects. As shown in Fig. 2, the key extracted entities are *Tiny Dog*, *White Car* and *plate of sushi*.
- **Analyzing Image Elements:** Next, the planner generates a detailed object-level analysis of  $I_g$ , which is akin to creating a comprehensive image caption limited to object-level information. As seen in our example, MLLM effectively finds the objects present in the image in context of the text-prompt,  $T$  i.e *Tiny Dog* and *White Car* and non-existence of *sushi*, *scattered oranges*.
- **Error Identification:** The extracted entities from both modalities extracted in the prior steps are compared and a descriptive summary of the identified errors is generated, this step not only grounds the mistakes but also provides

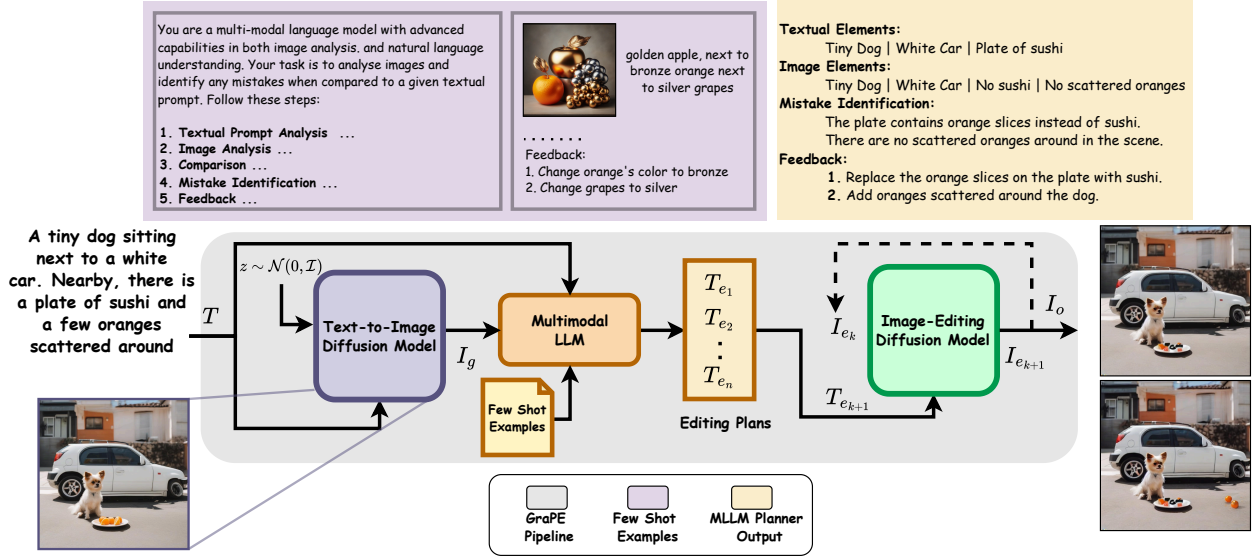


Figure 2. Proposed GraPE framework, a given text prompt is used to generate an initial image from T2I model,  $I_g$  which is then fed into a MLLM based planner along with the text prompt which identifies the objects that are misaligned in the image and outputs a set of edit plans guided by few-shot prompting. The plans are executed as a series of edits over the initial image to produce the final image

an interpretable way into planner’s reasoning.

- **Feedback Generation:** Finally, leveraging its extensive knowledge, the MLLM generates actionable feedback in the form of editing instructions to re-align image elements with the textual description. This step ensures that output aligns with the prompt’s intent. As in Figure 2, the MLLM plans to first replaces the objects (orange slices) on the plate with sushi and adds oranges around the dog in two distinct steps.

This structured approach not only enhances the interpretability and accuracy of the planner’s outputs but also provides a user-friendly, transparent process for refining image alignments in GraPE. The plans are extracted from MLLM output using regex matching of Feedback header.

### 3.3. Text Guided Image Editing

In our framework, the editing model plays a crucial yet straightforward role: it iteratively refines the generated image ( $I_g$ ) based on the plans extracted from the structured output produced by the MLLM in the planning phase. This design ensures flexibility, allowing for the use of any pre-trained model as a plug-and-play component.

**Compositional Image Editing** T2I models based on CLIP [38], and by extension, editing models built on it are limited by CLIP’s compositional reasoning capabilities including understanding word order in text and object relations in images [28, 48, 59]. We hypothesize and show for the first time that enhancing word order understanding in

the text space and training the model on high-quality object-centric and reasoning-oriented data is key to induce compositionality in image editors. For this, we introduce PixEdit, a text-guided image editing model based on PixArt-Sigma [11] and trained on the reasoning-centric dataset used by AURORA [29]. With the T5-XXL [39] language model as its text encoder, PixEdit demonstrates enhanced performance on compositional and spatially complex edits, benefiting from both robust language comprehension and reasoning-centric training data. Providing gains over all existing image-editing models for the generation task in our GraPE framework. Note that PixEdit is developed to be a *general-purpose* image editor, with enhanced compositional and object centric editing capabilities, and is not specific to the generation task we tackle. For additional details about PixEdit refer Supplementary Section 9.

### 3.4. Implementation Details

We utilize a frontier multi-modal model, GPT-4o as a multi-modal planner for its exceptional image-understanding and instruction following capabilities, both of which are essential for generating high-quality plans. While GraPE works with any strong image editing model, for best results we use our developed PixEdit and the recently proposed AURORA model [29] as our editors due to their ability to perform object-centric and reasoning-oriented edits on real images, supported by training on an editing dataset derived from videos and simulators. See Supplementary Section 9 for more details on implementation.









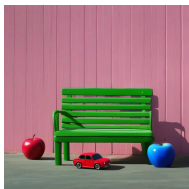
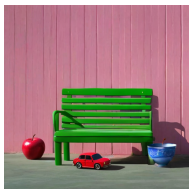
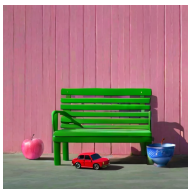
Text Prompt	Generated Image	Edit Step 1	Edit Step 2	Edit Step 3
A porcelain pot with tulips and a metal can with orchids and a glass jar with sunflowers				
		Remove sunflowers from image.	Add a glass jar on table.	Put sunflowers in the glass jar.
Three chairs are arranged in a row. A gray duck is positioned on the right side of the rightmost chair. A skyline in background.				
		Remove the duck from the leftmost chair.	Remove the duck from the middle chair.	Position a gray duck on the right side of the rightmost chair.
A green bench, a red car, a blue bowl, and a pink apple.				
		Change the blue apple to a blue bowl.	Change the red apple to a pink apple.	

Table 1. Iterative results by applying GraPE. on images generated by SD3.5 and SDXL, these images are edited via the proposed PixEdit editing model. Please see 8 for additional details on the plans generated for these images.

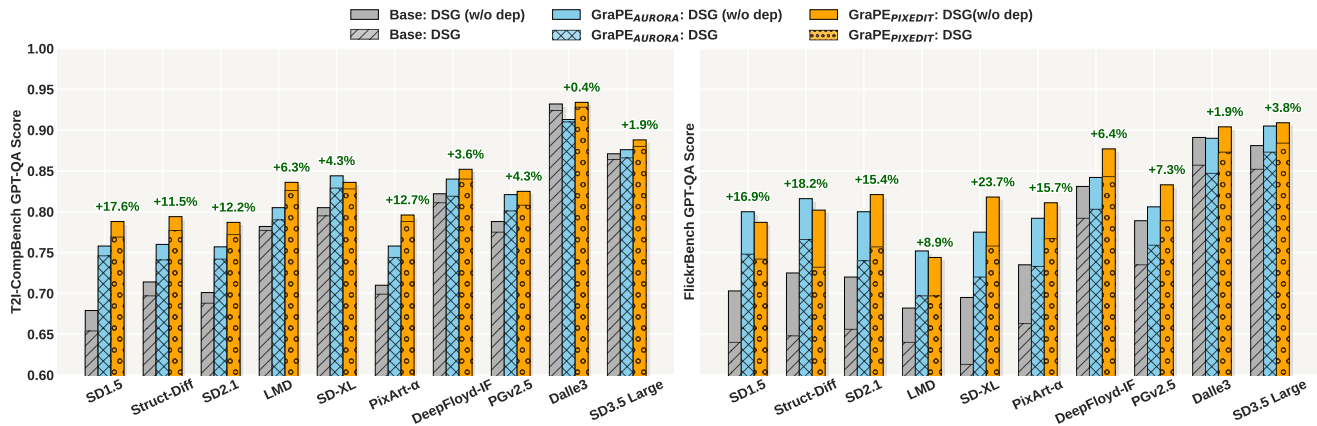


Figure 3. Experimental results showcasing the maximum gain in DSG score by GraPE with both AURORA and PixEdit as editing models. The figure presents both DSG and DSG (w/o dependency) scores. The percentage gain is measured over DSG scores.

## 4. Experiments

### 4.1. Experimental Settings

**Models:** To evaluate the effectiveness of GraPE, we conduct a comprehensive assessment across 10 state-of-the-art

text-to-image (T2I) models of varying sizes and capabilities. These models include: Stable Diffusion v1.5 [42], Structure Diffusion [16], Stable Diffusion V2.1 [42], LMD [31], SD-XL [37], PixArt- $\alpha$  [10], DeepFloyd IF [2], PlaygroundV2.5 [30], Dalle 3 [5], and the latest Stable Diffu-

Method	Concept K=1		Concept K=3		Concept K=5		Concept K=7	
	Base	+GraPE <sub>PixEdit</sub>	Base	+GraPE <sub>PixEdit</sub>	Base	+GraPE <sub>PixEdit</sub>	Base	+GraPE <sub>PixEdit</sub>
Stable-Diffusion v1.5 [42]	0.808 ±0.009	<b>0.892</b> ±0.002	0.606 ±0.018	<b>0.752</b> ±0.002	0.497 ±0.010	<b>0.689</b> ±0.002	0.450 ±0.005	<b>0.610</b> ±0.005
Structure Diffusion [16]	0.823 ±0.002	<b>0.885</b> ±0.004	0.606 ±0.002	<b>0.723</b> ±0.005	0.542 ±0.014	<b>0.703</b> ±0.006	0.447 ±0.001	<b>0.571</b> ±0.002
Stable Diffusion v2.1 [42]	0.833 ±0.002	<b>0.868</b> ±0.005	0.639 ±0.014	<b>0.737</b> ±0.002	0.579 ±0.012	<b>0.687</b> ±0.005	0.466 ±0.002	<b>0.626</b> ±0.002
LMD [31]	0.855 ±0.004	<b>0.873</b> ±0.002	0.711 ±0.008	<b>0.773</b> ±0.006	0.643 ±0.011	<b>0.725</b> ±0.009	0.591 ±0.002	<b>0.668</b> ±0.009
SD-XL [37]	0.848 ±0.010	<b>0.877</b> ±0.005	0.708 ±0.018	<b>0.780</b> ±0.007	0.635 ±0.014	<b>0.729</b> ±0.004	0.520 ±0.003	<b>0.628</b> ±0.002
PixArt-α [10]	0.813 ±0.010	<b>0.872</b> ±0.010	0.668 ±0.011	<b>0.722</b> ±0.002	0.649 ±0.011	<b>0.742</b> ±0.002	0.507 ±0.001	<b>0.625</b> ±0.002
DeepFloyd IF [2]	0.883 ±0.009	<b>0.915</b> ±0.000	0.680 ±0.016	<b>0.765</b> ±0.007	0.663 ±0.014	<b>0.745</b> ±0.002	0.583 ±0.002	<b>0.662</b> ±0.006
PlaygroundV2.5 [30]	0.908 ±0.010	<b>0.955</b> ±0.004	0.737 ±0.023	<b>0.792</b> ±0.009	0.658 ±0.015	<b>0.721</b> ±0.002	0.540 ±0.003	<b>0.640</b> ±0.005
Dalle3 [5]	0.947 ±0.002	<b>0.953</b> ±0.002	0.832 ±0.012	<b>0.861</b> ±0.003	0.812 ±0.014	<b>0.832</b> ±0.004	0.728 ±0.006	<b>0.737</b> ±0.004
Stable Diffusion v3.5 Large [46]	0.927 ±0.005	<b>0.948</b> ±0.002	0.815 ±0.002	<b>0.817</b> ±0.002	0.803 ±0.003	<b>0.831</b> ±0.004	0.759 ±0.004	<b>0.784</b> ±0.005

Table 2. Results on Concept-mix benchmark GraPE<sub>PixEdit</sub>

sion v3.5 Large [46]. This selection encompasses a diverse range of recent T2I models, allowing for a thorough evaluation of GraPE’s robustness across different architectures and configurations.

### Benchmarks:

- T2I Compbench** [26]: We evaluate the compositional generation capabilities of GraPE using this well-established benchmark, specifically designed to assess performance in this area. This includes compositional challenges across categories such as shape, color, texture, spatial, non-spatial, and complex. We utilize a subset of 100 prompts, sampled uniformly across all six categories—shape, color, texture, spatial, non-spatial, and complex—to ensure a comprehensive assessment of the model’s performance across diverse compositional challenges.
- ConceptMix** [55]: This benchmark evaluates models across varying levels of controllable compositionality. Using the official ConceptMix codebase, we generate 100 prompts for each  $K \in \{1, 3, 5, 7\}$ . Each prompt includes at least one object paired with  $K$  additional visual concepts—such as color, number, shape, size, style, spatial arrangement, and texture—offering a diverse and rigorous assessment of the model’s ability to handle increasingly complex visual compositions.
- Flicker-Bench**: Flicker-30k [36] is a widely used benchmark for evaluating the generation quality of T2I models using metrics such as FID [23]. The human-generated prompts in this dataset offer a well-balanced mix of compositionality and realism, making it suitable for assessing our method’s performance on general-purpose tasks beyond compositionality. To this end, we sample 100 prompts from the dataset’s test split for evaluation.

**Evaluation Metrics:** Many of the SOTA automated evaluation metrics belong to the QA pair generation and Visual Question-Answering (VQA) framework. Prior works

[13, 21, 25, 53, 58] have demonstrated its effectiveness and correlation with human-predictions when judging semantic accuracies. The ability to generate object/attribute centric ‘yes’/‘no’ questions over a compositional text prompt allows for fine-grained assessment of image elements using MLLMs. We therefore follow [13] to generate grounded binary questions for T2I-Compbench and Flickr-Bench and utilize ConceptMix’s existing set of questions that are generated simultaneously with the prompts and utilize GPT-4o as the choice VQA model following [55].

## 4.2. Results

**GraPE 🍷 Improves Compositional T2I Synthesis** GraPE can be used with any given back-bone in a plug-n-play manner, and substantially improves performance across the board for all of the 10 T2I models tested in this study, including SoTA models like SD3.5 [46]. Results on T2I Compbench and Flickr-Bench are present in Figure 3. Results on ConceptMix are present in Table 2. On T2I Compbench we see improvements ranging from 17.6% in the case of SD1.5 to nearly 2% for the latest SD3.5 Large model. On ConceptMix as well, we see large performance gains across prompt complexities  $K=\{1, 3, 5, 7\}$ . For prompts with higher complexity ( $K=7$ ), we see gains ranging from 3.1% to 35.5% across models.

**GraPE 🍷 is more effective with complex T2I prompts.** As seen in Table 2, the absolute performance difference between generation from base model vs our method increases as we go towards more complex prompts (from  $K=1$  to  $K=7$ ). For e.g., for SD XL, the gains compared to the base model increase from 2.9% for  $K=1$  to 10.8% for  $K=7$ , and this observation remains constant across T2I models.

**GraPE 🍷 exceeds performance on general T2I tasks** Gains achieved by GraPE are not limited to compositional or complex prompts as is the case in T2I Compbench and ConceptMix datasets, rather it also generalizes to more widely used general T2I benchmarks such as Flickr-Bench. Fig. 3 shows that GraPE leads to large performance gains

ranging from 2 to 23% across different models.

**GraPE** reduces the gap between models of varied T2I capabilities GraPE leads to a larger absolute performance increase for relatively less capable base diffusion models. This closes the gap between performance of such models with more capable ones. For e.g., for ConceptMix dataset ( $K=3$ ), the difference between SD-1.5 and SD-3.5 reduces from 20.9% to 6.5% when comparing the base generations vs generations from GraPE respectively.

**GraPE** scales well with more compute Our approach can be viewed from a lens of using extra inference time compute (calls to the editing model) for creating images better aligned with complex prompts, achieving better performance on the T2I Synthesis task. GraPE flexibly trades off performance with the amount of compute (or the number of edits at inference time). See Fig. 4 for results on four representative models. This also shows that each planning step of GraPE is important on average for improving on the T2I synthesis task gradually.

For additional results, including comparisons with other recent approaches, see Supplementary Section 8.

### 4.3. Ablations

**Naive Planner** This section highlights the importance of object and attribute centric decomposition of plans to generate effective editing instructions. We modify the few-shot structure of the MLLM planner to skip the decomposition into image and text elements and prompt it to identify errors and suggest editing plans given the text-prompt and generated image. We refer to the pipeline thus constructed as  $\text{GraPE}_{naive}$ . Table 3 presents the results of both GraPE and  $\text{GraPE}_{naive}$  on  $K=7$  subset of ConceptMix benchmark for a subset of models. See Supplementary Section 7 for more ablations.

Method	Concept K=7	
	$\text{GraPE}_{naive}$	GraPE
Stable Diffusion v2.1 [42]	0.618	<b>0.626</b>
LMD [31]	0.622	<b>0.668</b>
SD-XL [37]	0.598	<b>0.628</b>
PlaygroundV2.5 [30]	0.622	<b>0.640</b>
Stable Diffusion v3.5 Large [46]	0.724	<b>0.784</b>

Table 3. Results on Concept-mix benchmark ( $K=7$  subset). Comparing both  $\text{GraPE}_{naive}$  and GraPE (Using PixEdit)

### 4.4. Analysis and Insights

We conduct a detailed examination of each component in GraPE, focusing specifically on identifying and analyzing failure cases in both the MLLM planner and the editing model. For this purpose, we curate a dataset of 120 image-prompt pairs, uniformly sampled across all 10 models and

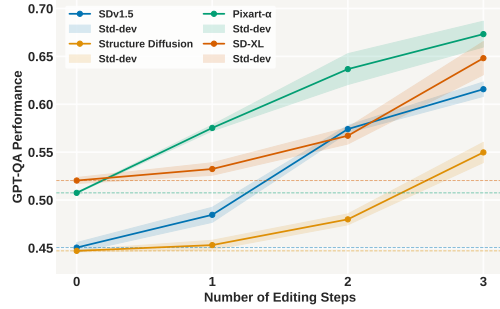


Figure 4. Trend of GPT-QA score with increasing number of editing steps.

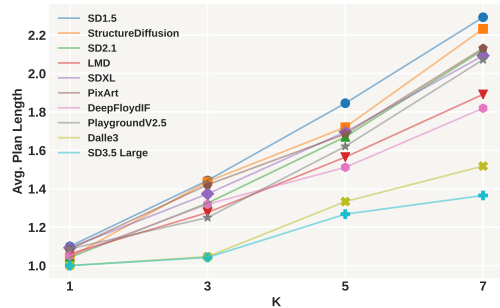


Figure 5. Average Steps per Plan for ConceptMix Benchmark across the models.

6 benchmarks. This dataset is carefully evaluated by six human annotators, compensated well above the national average hourly wage.

Type	Q1		Q2			Q3	
	Yes	No	No	Partial	Full	Yes	No
GPT-4o	25.4	74.5	8.7	35.7	55.5	89.4	10.5
Qwen-VL-72B	25.4	74.5	26.9	30.0	43.0	85.9	14.0

Table 4. Results from human study on MLLM planners. We present results with GPT-4o based planner used in GraPE and an open-source alternative Qwen-VL-72B [3]

**MLLM planner Accuracy** For each image-prompt pair, along with the *plan* generated by GraPE, evaluators were asked to answer three multiple-choice questions:

- Q1: Does the generated image accurately and completely match the content described in the text prompt?
- Q2: To what extent does executing the proposed plan improve the image’s alignment with the text prompt?
- Q3: Are the instructions in the plan indivisible, meaning each instruction is a fundamental action that cannot be further subdivided?

Questions Q1 and Q3 require binary answers, with 'yes' (1) or 'no' (0) responses. Q2 offers three response options: 'no improvement', 'partial improvement' and 'full alignment'. This structured evaluation allows us to assess the accuracy and granularity of the MLLM planner in relation to image-prompt alignment.

Our analysis reveals that Qwen-VL-72B as planner has around 18% more cases of 'no improvement' compared to GPT-4o, largely due to generating a considerable number of empty plans. This suggests that while Qwen-VL-72B may trail GPT-4o in visual capabilities it often fails to identify discrepancies, but tends to produce a solid plan when it does recognize them, as evident from Table 4. We also carry out the annotation process on a smaller model, Qwen-VL-7B which in contrast to the above planners performed poorly, generating approximately 100 empty plans out of 120 images, indicating significant limitations in both identifying discrepancies and generating actionable plans.

**Error Analysis of Editing Model** We qualitatively take a look at the failure cases of the editing model (PixEdit) specifically utilizing the above dataset to gather a subset of plans which are deemed to fully align the image with the text prompt by the majority of evaluators (based on Q2 earlier). This subset consists of about 46 plans in which any errors in the final image solely result from the shortcomings of the image-editing model. We specifically look at images in which correct, partial or incorrect edits were performed by the model. Out of these 14 were completely correct, 17 partially correct, and 15 were incorrect. Within the correct ones, the distribution of correct add, remove and modify was 12, 0 and 2, respectively. Within the partially correct ones, the distribution of correct add, remove and modify was 11, 2 and 4, respectively. Figure 6 provides some representative examples. Fixing these issues, possibly through an RL based mechanism for incorporating feedback is a direction for future work. This also points to the fact that the primary reason for any errors originates from the editing model, rather than because of the planner.

**Multi-Modal Planner as Evaluation Metric** Can planning help evaluating T2I generative models? For this, we look at the analysis and mistake identification capabilities of the planner. Specifically, we post-process the planner-generated output filtering out only the textual-image element comparison and pass this back to the planner which is provided a zero-shot system prompt. This prompt directs the planner to generate a score on a 1-100 scale identifying the alignment based on the objects, their relationships and attributes. When compared with LLMscore [33], We observed a moderate correlation of 0.423 and 0.404 when tested using Spearman and Pearson's test respectively. This shows the future potential of our MLLM planner to be used



Figure 6. Results illustrating failure cases of Image-Editing model

in conjunction with SoTA evaluation metrics to further improve them.

**Number of planning steps** In Fig. 5, we see a gradual increase in the average number of plans as we increase the prompt complexity from  $K=1$  to  $K=7$  using the Concept-Mix benchmark. This shows our planner is able to recognize the complexity of the prompt and generate a longer plan in response. We also see that on average, more capable T2I models like SD3.5 require fewer editing steps compared to weaker models such as SD1.5.

Additional qualitative examples of our approach are shown in Supplementary Section 10.

## 5. Discussion and Limitations

**Summary:** In this paper, we have looked at the task of fine-grained image generation for compositional prompts. Unlike most existing methods, which try to achieve this task in one-go, we take a different route, and decompose the problem in 3-steps: (a) Generate an initial image based on the prompts (b) Identify the mistakes and plan out corrective steps by analyzing the difference between textual and visual elements for each object (c) Execute the corrective steps via an editing model in a sequential manner. The ability of our model to break the corrective steps into multiple simple edits is key to the success of our approach. Extensive experimentation on multiple datasets shows that our method can significantly improve the generation quality over 10 different models, including SoTA approaches.

**Limitations and future work:** Limitations of our approach include its inability to handle certain kinds of complex prompts, and its dependence on a pre-trained MLLM in the planning step. Future work includes working with more complex prompts, integrating our approach with models which explicitly change the architecture to better align with text, and incorporating corrective feedback using an RL based framework and extending to image editing tasks.



## References

- [1] Aishwarya Agarwal, Srikrishna Karanam, K J Joseph, Apoorv Saxena, Koustava Goswami, and Balaji Vasan Srinivasan. A-star: Test-time attention segregation and retention for text-to-image synthesis, 2023. 2
- [2] DeepFloyd Lab at StabilityAI. DeepFloyd IF: a novel state-of-the-art open-source text-to-image model with a high degree of photorealism and language understanding. <https://www.deepfloyd.ai/deepfloyd-if>, 2023. Retrieved on 2023-11-08. 5, 6, 4
- [3] Jinze Bai, Shuai Bai, Shusheng Yang, Shijie Wang, Sinan Tan, Peng Wang, Junyang Lin, Chang Zhou, and Jingren Zhou. Qwen-vl: A versatile vision-language model for understanding, localization, text reading, and beyond, 2023. 7
- [4] Eslam Mohamed Bakr, Pengzhan Sun, Xiaoqian Shen, Faizan Farooq Khan, Li Erran Li, and Mohamed Elhoseiny. Hrs-bench: Holistic, reliable and scalable benchmark for text-to-image models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 20041–20053, 2023. 2
- [5] James Betker, Gabriel Goh, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, et al. Improving image generation with better captions. *Computer Science*, 2(3):8, 2023. 1, 5, 6, 4
- [6] Tim Brooks, Aleksander Holynski, and Alexei A. Efros. Instructpix2pix: Learning to follow image editing instructions, 2023. 3, 2
- [7] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Advances in Neural Information Processing Systems*, 2020. 3
- [8] Agneet Chatterjee, Gabriela Ben Melech Stan, Estelle Aflalo, Sayak Paul, Dhruva Ghosh, Tejas Gokhale, Ludwig Schmidt, Hannaneh Hajishirzi, Vasudev Lal, Chitta Baral, and Yezhou Yang. Getting it right: Improving spatial consistency in text-to-image models, 2024. 2
- [9] Hila Chefer, Yuval Alaluf, Yael Vinker, Lior Wolf, and Daniel Cohen-Or. Attend-and-excite: Attention-based semantic guidance for text-to-image diffusion models, 2023. 2
- [10] Junsong Chen, Jincheng Yu, Chongjian Ge, Lewei Yao, Enze Xie, Yue Wu, Zhongdao Wang, James Kwok, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\alpha$ : Fast training of diffusion transformer for photorealistic text-to-image synthesis, 2023. 5, 6, 4
- [11] Junsong Chen, Chongjian Ge, Enze Xie, Yue Wu, Lewei Yao, Xiaozhe Ren, Zhongdao Wang, Ping Luo, Huchuan Lu, and Zhenguo Li. Pixart- $\sigma$ : Weak-to-strong training of diffusion transformer for 4k text-to-image generation, 2024. 2, 4
- [12] Jaemin Cho, Abhay Zala, and Mohit Bansal. Visual programming for text-to-image generation and evaluation. In *NeurIPS*, 2023. 2
- [13] Jaemin Cho, Yushi Hu, Roopal Garg, Peter Anderson, Ranjay Krishna, Jason Baldridge, Mohit Bansal, Jordi Pont-Tuset, and Su Wang. Davidsonian scene graph: Improving reliability in fine-grained evaluation for text-to-image generation, 2024. 6
- [14] Prafulla Dhariwal and Alex Nichol. Diffusion models beat gans on image synthesis. *ArXiv*, abs/2105.05233, 2021. 1
- [15] Alaaeldin El-Nouby, Shikhar Sharma, Hannes Schulz, Devon Hjelm, Layla El Asri, Samira Ebrahimi Kahou, Yoshua Bengio, and Graham W. Taylor. Tell, draw, and repeat: Generating and modifying images based on continual linguistic instruction. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 3
- [16] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis, 2023. 2, 5, 6, 4
- [17] Weixi Feng, Xuehai He, Tsu-Jui Fu, Varun Jampani, Arjun Reddy Akula, Pradyumna Narayana, Sugato Basu, Xin Eric Wang, and William Yang Wang. Training-free structured diffusion guidance for compositional text-to-image synthesis. In *The Eleventh International Conference on Learning Representations*, 2023. 2
- [18] Yutong Feng, Biao Gong, Di Chen, Yujun Shen, Yu Liu, and Jingren Zhou. Ranni: Taming text-to-image diffusion for accurate instruction following, 2024. 2
- [19] Peng Gao, Le Zhuo, Dongyang Liu, Ruoyi Du, Xu Luo, Longtian Qiu, Yuhang Zhang, Chen Lin, Rongjie Huang, Shijie Geng, Renrui Zhang, Junlin Xi, Wenqi Shao, Zhengkai Jiang, Tianshuo Yang, Weicai Ye, He Tong, Jingwen He, Yu Qiao, and Hongsheng Li. Lumina-t2x: Transforming text into any modality, resolution, and duration via flow-based large diffusion transformers, 2024. 2
- [20] Yuying Ge, Sijie Zhao, Chen Li, Yixiao Ge, and Ying Shan. Seed-data-edit technical report: A hybrid dataset for instructional image editing, 2024. 2
- [21] Tejas Gokhale, Hamid Palangi, Besmira Nushi, Vibhav Vineet, Eric Horvitz, Ece Kamar, Chitta Baral, and Yezhou Yang. Benchmarking spatial relationships in text-to-image generation, 2023. 2, 6
- [22] Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-Or. Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*, 2022. 3
- [23] 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, 2018. 6
- [24] Xiwei Hu, Rui Wang, Yixiao Fang, Bin Fu, Pei Cheng, and Gang Yu. Ella: Equip diffusion models with llm for enhanced semantic alignment, 2024. 2
- [25] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate

- and interpretable text-to-image faithfulness evaluation with question answering, 2023. 6
- [26] Kaiyi Huang, Kaiyue Sun, Enze Xie, Zhenguo Li, and Xihui Liu. T2i-compbench: A comprehensive benchmark for open-world compositional text-to-image generation, 2023. 2, 6
- [27] Dongzhi Jiang, Guanglu Song, Xiaoshi Wu, Renrui Zhang, Dazhong Shen, Zhuofan Zong, Yu Liu, and Hongsheng Li. Comat: Aligning text-to-image diffusion model with image-to-text concept matching, 2024. 2
- [28] Amita Kamath, Jack Hessel, and Kai-Wei Chang. Text encoders bottleneck compositionality in contrastive vision-language models, 2023. 2, 4
- [29] Benno Krojer, Dheeraj Vattikonda, Luis Lara, Varun Jampani, Eva Portelance, Christopher Pal, and Siva Reddy. Learning action and reasoning-centric image editing from videos and simulations, 2024. 3, 4, 1, 2
- [30] Daiqing Li, Aleks Kamko, Ehsan Akhgari, Ali Sabet, Linmiao Xu, and Suhail Doshi. Playground v2.5: Three insights towards enhancing aesthetic quality in text-to-image generation, 2024. 5, 6, 7, 4
- [31] Long Lian, Boyi Li, Adam Yala, and Trevor Darrell. Llm-grounded diffusion: Enhancing prompt understanding of text-to-image diffusion models with large language models, 2024. 2, 3, 5, 6, 7, 4
- [32] Bingchen Liu, Ehsan Akhgari, Alexander Visheratin, Aleks Kamko, Linmiao Xu, Shivam Shrirao, Joao Souza, Suhail Doshi, and Daiqing Li. Playground v3: Improving text-to-image alignment with deep-fusion large language models, 2024. 2
- [33] Yujie Lu, Xianjun Yang, Xiujun Li, Xin Eric Wang, and William Yang Wang. LLMscore: Unveiling the power of large language models in text-to-image synthesis evaluation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. 8
- [34] Alex Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. *ArXiv*, abs/2102.09672, 2021. 1
- [35] Weili Nie, Sifei Liu, Morteza Mardani, Chao Liu, Benjamin Eckart, and Arash Vahdat. Compositional text-to-image generation with dense blob representations. In *International Conference on Machine Learning (ICML)*, 2024. 2
- [36] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models, 2016. 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, 2023. 5, 6, 7, 1, 4
- [38] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision, 2021. 2, 4
- [39] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer, 2023. 2, 4
- [40] Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-conditional image generation with clip latents, 2022. 2
- [41] Recraft AI. Recraft-V3. <https://www.recraft.ai/ai-image-generator>, 2024. 1
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10684–10695, 2022. 5, 6, 7, 4
- [43] Chitwan Saharia, William Chan, Saurabh Saxena, Lala Li, Jay Whang, Emily L. Denton, Seyed Kamyar Seyed Ghasemipour, Burcu Karagol Ayan, Seyedeh Sara Mahdavi, Raphael Gontijo Lopes, Tim Salimans, Jonathan Ho, David J. Fleet, and Mohammad Norouzi. Photorealistic text-to-image diffusion models with deep language understanding. *ArXiv*, abs/2205.11487, 2022. 1
- [44] Harman Singh, Poorva Garg, Mohit Gupta, Kevin Shah, Ashish Goswami, Satyam Modi, Arnab Mondal, Dinesh Khandelwal, Dinesh Garg, and Parag Singla. Image manipulation via multi-hop instructions - a new dataset and weakly-supervised neuro-symbolic approach. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 2975–3007. Association for Computational Linguistics, 2023. 3
- [45] Jascha Sohl-Dickstein, Eric A. Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics, 2015. 1
- [46] stability.ai. Stable diffusion 3.5. <https://stability.ai/news/introducing-stable-diffusion-3-5>, 2024. 1, 6, 7, 4
- [47] Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, Pouya Tafti, Léonard Hussenot, Pier Giuseppe Sessa, Aakanksha Chowdhery, Adam Roberts, Aditya Barua, Alex Botev, Alex Castro-Ros, Ambrose Slone, Amélie Héliou, Andrea Tacchetti, Anna Bulanova, Antonia Paterson, Beth Tsai, Bobak Shahriari, Charline Le Lan, Christopher A. Choquette-Choo, Clément Crepy, Daniel Cer, Daphne Ippolito, David Reid, Elena Buchatskaya, Eric Ni, Eric Noland, Geng Yan, George Tucker, George-Christian Muraru, Grigory Rozhdestvenskiy, Henryk Michalewski, Ian Tenney, Ivan Grishchenko, Jacob Austin, James Keeling, Jane Labanowski, Jean-Baptiste Lespiau, Jeff Stanway, Jenny Brennan, Jeremy Chen, Johan Ferret, Justin Chiu, Justin Mao-Jones, Katherine Lee, Kathy Yu, Katie Millican, Lars Lowe Sjoesund, Lisa Lee, Lucas Dixon, Machel Reid, Maciej Mikula, Mateo Wirth, Michael Sharman, Nikolai Chinaev, Nithum Thain, Olivier Bachem, Oscar Chang, Oscar Wahltinez, Paige Bailey, Paul Michel, Petko Yotov, Rahma Chaabouni, Ramona Comanescu, Reena Jana, Rohan Anil, Ross McIlroy, Ruibo Liu, Ryan Mullins, Samuel L Smith, Sebastian Borgeaud, Sertan Girgin, Sholto Douglas, Shree Pandya, Siamak Shakeri, Soham De, Ted Klimenko, Tom Hennigan, Vlad Feinberg, Wojciech Stokowiec, Yu hui Chen, Zafarali Ahmed, Zhitao

- Gong, Tris Warkentin, Ludovic Peran, Minh Giang, Clément Farabet, Oriol Vinyals, Jeff Dean, Koray Kavukcuoglu, Demis Hassabis, Zoubin Ghahramani, Douglas Eck, Joelle Barral, Fernando Pereira, Eli Collins, Armand Joulin, Noah Fiedel, Evan Senter, Alek Andreev, and Kathleen Kenealy. Gemma: Open models based on gemini research and technology, 2024. [2](#)
- [48] Tristan Thrush, Ryan Jiang, Max Bartolo, Amanpreet Singh, Adina Williams, Douwe Kiela, and Candace Ross. Winoground: Probing vision and language models for visiolinguistic compositionality. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5228–5238, 2022. [4](#)
- [49] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023. [2](#)
- [50] Zhenyu Wang, Aoxue Li, Zhenguo Li, and Xihui Liu. Genartist: Multimodal llm as an agent for unified image generation and editing. *arXiv preprint arXiv:2407.05600*, 2024. [2](#)
- [51] Zirui Wang, Zhizhou Sha, Zheng Ding, Yilin Wang, and Zhuowen Tu. Tokencompose: Text-to-image diffusion with token-level supervision, 2024. [2](#)
- [52] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. [3](#)
- [53] Olivia Wiles, Chuhan Zhang, Isabela Albuquerque, Ivana Kajić, Su Wang, Emanuele Bugliarello, Yasumasa Onoe, Chris Knutsen, Cyrus Rashtchian, Jordi Pont-Tuset, and Aida Nematzadeh. Revisiting text-to-image evaluation with gecko: On metrics, prompts, and human ratings, 2024. [6](#)
- [54] Weijia Wu, Zhuang Li, Yefei He, Mike Zheng Shou, Chunhua Shen, Lele Cheng, Yan Li, Tingting Gao, Di Zhang, and Zhongyuan Wang. Paragraph-to-image generation with information-enriched diffusion model, 2023. [2](#)
- [55] Xindi Wu, Dingli Yu, Yangsibo Huang, Olga Russakovsky, and Sanjeev Arora. Conceptmix: A compositional image generation benchmark with controllable difficulty, 2024. [6](#)
- [56] Enze Xie, Junsong Chen, Junyu Chen, Han Cai, Haotian Tang, Yujun Lin, Zhekai Zhang, Muyang Li, Ligeng Zhu, Yao Lu, and Song Han. Sana: Efficient high-resolution image synthesis with linear diffusion transformer, 2024. [2](#)
- [57] Ling Yang, Zhaochen Yu, Chenlin Meng, Minkai Xu, Stefano Ermon, and Bin Cui. Mastering text-to-image diffusion: Recaptioning, planning, and generating with multimodal llms. In *International Conference on Machine Learning*, 2024. [2](#), [3](#), [1](#)
- [58] Michal Yarom, Yonatan Bitton, Soravit Changpinyo, Roei Aharoni, Jonathan Herzig, Oran Lang, Eran Ofek, and Idan Szpektor. What you see is what you read? improving text-image alignment evaluation, 2023. [6](#)
- [59] Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and why vision-language models behave like bags-of-words, and what to do about it?, 2023. [2](#), [4](#)
- [60] Kai Zhang, Lingbo Mo, Wenhui Chen, Huan Sun, and Yu Su. Magicbrush: A manually annotated dataset for instruction-guided image editing, 2024. [3](#)
- [61] Shu Zhang, Xinyi Yang, Yihao Feng, Can Qin, Chia-Chih Chen, Ning Yu, Zeyuan Chen, Huan Wang, Silvio Savarese, Stefano Ermon, Caiming Xiong, and Ran Xu. Hive: Harnessing human feedback for instructional visual editing. *arXiv preprint arXiv:2303.09618*, 2023.
- [62] Tianhao Zhang, Hung-Yu Tseng, Lu Jiang, Weilong Yang, Honglak Lee, and Irfan Essa. Text as neural operator: Image manipulation by text instruction. In *Proceedings of the 29th ACM International Conference on Multimedia*, pages 1893–1902, 2021. [3](#)
- [63] Xinyu Zhang, Mengxue Kang, Fei Wei, Shuang Xu, Yuhe Liu, and Lin Ma. Tie: Revolutionizing text-based image editing for complex-prompt following and high-fidelity editing, 2024. [3](#)
- [64] Le Zhuo, Ruoyi Du, Han Xiao, Yangguang Li, Dongyang Liu, Rongjie Huang, Wenzhe Liu, Lirui Zhao, Fu-Yun Wang, Zhanyu Ma, Xu Luo, Zehan Wang, Kaipeng Zhang, Xi-angyang Zhu, Si Liu, Xiangyu Yue, Dingning Liu, Wanli Ouyang, Ziwei Liu, Yu Qiao, Hongsheng Li, and Peng Gao. Lumina-next: Making lumina-t2x stronger and faster with next-dit, 2024. [2](#)

# GraPE 🍇: A Generate-Plan-Edit Framework for Compositional T2I Synthesis

## Supplementary Material

### 6. Algorithm:

We present the algorithm showcasing the sequence of steps in GraPE below.

---

**Algorithm 1** Generate-Plan-Edit : GraPE 🍇

---

**Require:** Text Prompt:  $T$ , T2I Model:  $G$ , MLLM:  $\mathcal{P}$ , Editing Model:  $\mathcal{E}$  and Few-Shot Examples:  $[E_1 \dots E_p]$

*#Generate*

$I_g \leftarrow \mathcal{G}(T)$  ▷ Initial Generated Image

*#Plan*

$\{T_{e_1}, T_{e_2}, \dots, T_{e_n}\} \leftarrow \mathcal{P}(I_g, T, [E_1 \dots E_p])$

*#Edit*

$I_{e_0} \leftarrow I_g$

**for**  $k = 1, \dots, n-1$  **do**

$I_{e_{k+1}} \leftarrow \mathcal{E}(I_{e_k}, T_{e_{k+1}})$  ▷ Intermediate Edited Image

**end for**

$I_{e_o} \leftarrow I_{e_n}$  ▷ Final Edited Image

---

### 7. Ablations:

We provide the system prompt and an in-context prompting example in Fig. 9 for the GraPE<sub>naive</sub> pipeline used in section 4.3. The planner follows the system prompt to generate reasoning and editing instructions from the given image and text prompt, without the proposed decomposition into image and text based elements.

### 8. Additional and Detailed Results:

**Comparison with SOTA LLM-Based Approaches** We compare the performance of GraPE against RPG [57], a recent LLM-based image generation method. RPG leverages an LLM to decompose complex image generation tasks into simpler subtasks within localized subregions and employs complementary-regional-diffusion to coherently merge the generated subregions into a unified image. The comparison is conducted on the ConceptMix benchmark, with RPG using SD-XL [37] as its base model, while GraPE employs the PixEdit editing model. Results are summarized in Table 5.

Our evaluation reveals that RPG fails to generate valid plans for over 50% of the samples in the comparison, primarily due to frequent parsing errors. For the remaining samples, GraPE demonstrates significant improvements in GPT-QA accuracy over images generated by RPG. Specif-

ically, GraPE achieves gains of 21.5%, 3.0%, 18.0%, and 20.3% for  $K=1$ ,  $K=3$ ,  $K=5$ , and  $K=7$ , respectively.

**Table for T2I-Bench and Flickr-Bench** Table 10, 11 presents the absolute values of GraPE with PixEdit and AURORA [29] as editing models, for the graphs in Fig. 3.

**Results on ConceptMix Benchmark using GraPE with AURORA** Table 9 presents the numbers on ConceptMix Dataset using GraPE with AURORA editing model, this table complements Table 2, where PixEdit was used as the editing model.

**Performance on Editing Benchmarks** We employ two established benchmarks to evaluate and compare the editing performance of PixEdit with other models. The first benchmark, AURORA-BENCH, consists of image-edit instruction pairs sourced from eight distinct datasets. This benchmark enables the assessment of the discriminative editing capabilities of models across a wide variety of images. Performance is measured using an automated metric called **DiscEdit** [29]. The second benchmark we evaluate is the MagicBrush Test-set, which comprises 1K editing turns, including both single-turn and multi-turn edits. For quantitative analysis on this benchmark, we utilize standard metrics such as L1, L2, and CLIP/DINO similarity scores, providing a comprehensive evaluation of the editing performance. The results in Table 6 indicates strong discriminative capability in PixEdit compared to baselines. Table 7 shows that PixEdit is great at single turn edits and follows closely with baselines in multi-turn edits.

### 9. Implementation Details:

#### 9.1. Prompts and Few-shot Examples:

Figure 7 and Figure 8 illustrate the full system prompt and a selection of few-shot examples employed in the MLLM-based planner for GraPE.

#### 9.2. Hyperparameters

**Generation** To generate images using various diffusion models, we adhere to the default hyperparameters specific to each model, such as the number of inference steps and the sampling method. All images are generated with a fixed random seed of 0 to ensure reproducibility.

You are a multi-modal language model with advanced capabilities in both image analysis and natural language understanding. Your task is to analyse images and identify any mistakes, inconsistencies, or discrepancies when compared to a given textual prompt. Follow these steps:

**Textual Prompt Analysis :**

- Thoroughly read and understand the textual prompt.
- Extract key elements, objects, and attributes described in the text.

**Image Analysis:**

- Examine the image in detail.
- Identify and describe the key elements, objects, and attributes present in the image.

**Comparison:**

- Compare the elements, objects, and attributes found in the image with those described in the textual prompt.
- Note any discrepancies, such as missing elements, additional elements not described in the text, or any incorrect attributes (e.g., color, size, position).

**Mistake Identification:**

- Clearly identify and list any mistakes or discrepancies found.
- Provide a detailed explanation for each identified mistake.

**Feedback:**

- Offer suggestions for correcting the identified mistakes to ensure the image accurately reflects the textual prompt.
- Keep the feedback minimal and to the point. It should be object centric.

The following rules should also be noted:

- In case there are no mistakes, provide no feedback.
- If there is extreme misalignment between the textual prompt and the image only return <REGENERATE> tag.
- Do Not use phrases like 'in order to match the prompt' and similar ones the feedback
- The feedback must not dictate what will happen after the change. It is not a suggestion but an instruction.
- The feedback instructions must be atomic or single step. Any instruction requiring multiple steps must be divided into two instructions

Figure 7. System-prompt used with GraPE’s MLLM Planner

Method	Concept K=1		Concept K=3		Concept K=5		Concept K=7	
	Base	+GraPE <sub>PixEdit</sub>	Base	+GraPE <sub>PixEdit</sub>	Base	+GraPE <sub>PixEdit</sub>	Base	+GraPE <sub>PixEdit</sub>
RPG [57]	0.696 ±0.015	<b>0.845</b> ±0.008	0.694 ±0.005	<b>0.715</b> ±0.007	0.583 ±0.002	<b>0.688</b> ±0.005	0.388 ±0.007	<b>0.467</b> ±0.005

Table 5. Comparison of RPG [57] and GraPE on ConceptMix

Model	WhatsUp	Something	AG	Kubric	CLEVR
MagicBrush	0.472	0.371	0.477	0.392	0.400
AURORA	0.565	0.548	0.583	<b>0.592</b>	0.450
PixEdit	<b>0.566</b>	<b>0.613</b>	<b>0.606</b>	0.400	<b>0.600</b>

Table 6. Performance comparison of PixEdit with other models on AURORA-BENCH using *DiscEdit* computed as specified in [29]. Note: The baseline scores are taken from [29]

**Planning** We utilize GPT-4o with a sampling temperature set to 0, ensuring deterministic outputs for both planning and Visual Question Answering (VQA) tasks.

**Editing** For image editing, we use AURORA with 50 inference steps and its default sampler as described in its paper. Additionally, we employ 14 inference steps with DPM-

Solver++ in combination with the PixEdit model to generate edited images, ensuring reproducibility by using a fixed random seed of 0.

**9.3. PixEdit:**

This section details the implementation of PixEdit model. We chose PixArt-sigma [11] as our backbone T2I model which we convert into an editing checkpoint. We follow similar training strategy as used in [6, 29] i.e adding additional channels in the initial convolution layer of the diffusion model and use the randomly sampled noise concatenated with input image as input to de-noise into edited-image. We started with the SEED-EDIT [20] dataset for initial pre-training stage. We skip the automatically generated data and use the 147K sample subset of this dataset that are either procured from the web or manually created to en-

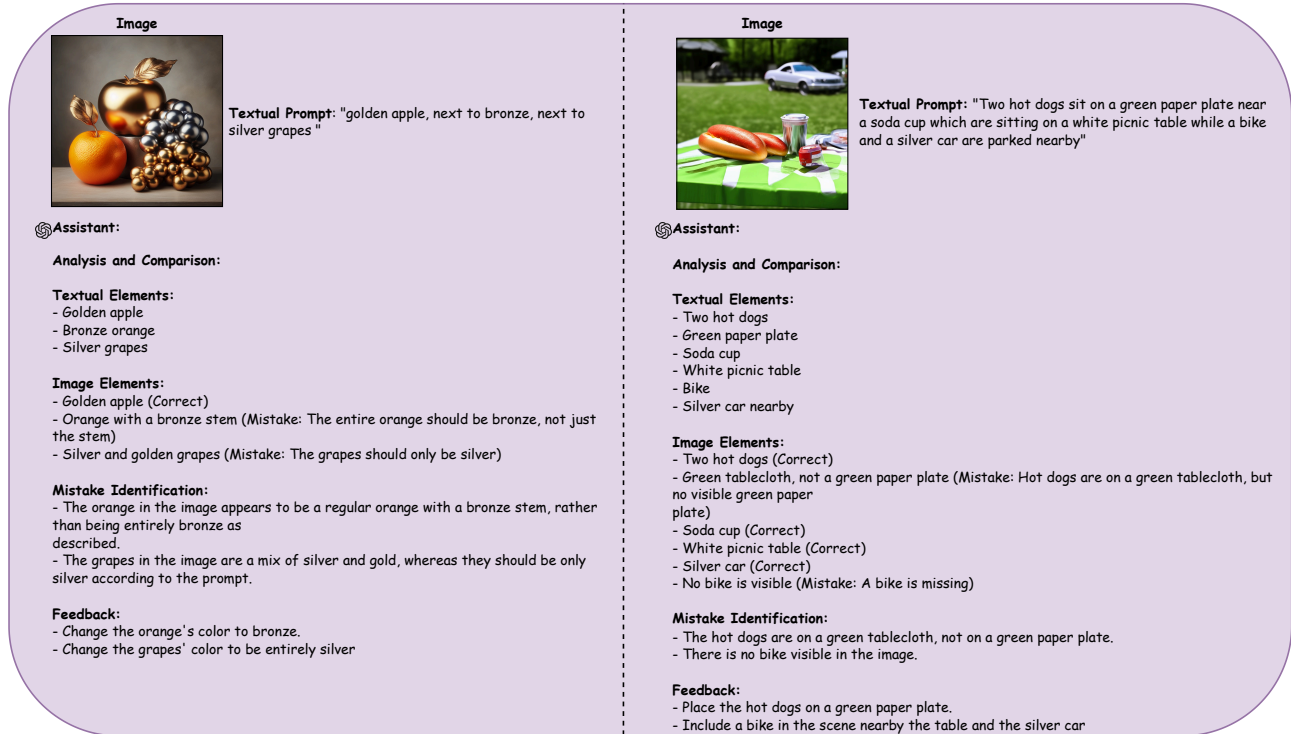


Figure 8. Selective Few shot examples used in GraPE’s MLLM Planner

Settings	Methods	L1↓	L2↓	CLIP-I↑	DINO↑	CLIP-T↑
Single-turn	MagicBrush	0.0788	0.0274	0.8978	0.8313	0.2973
	AURORA	0.0754	0.0270	0.9105	0.8594	<b>0.2981</b>
	PixEdit	<b>0.0719</b>	0.0278	0.9082	<b>0.8661</b>	0.2946
Multi-Turn	MagicBrush	0.0921	0.0327	0.8777	0.7996	<b>0.3020</b>
	AURORA	0.0904	0.0334	<b>0.8887</b>	<b>0.8272</b>	0.3005
	PixEdit	<b>0.0892</b>	0.0357	0.8729	0.8068	0.2947

Table 7. Comparison of methods across Single-turn and Multi-turn settings on Magic-Brush Test set. **Note:** The numbers were reproduced using the publicly available codebases of the respective baselines with the same random seed.

sure high quality pre-training data. The model is trained for 32K steps with an effective batch size of 256 at this stage. We further fine-tune this model for 42K steps on the mixture of editing dataset proposed in [29] to create PixEdit.

### 9.3.1. Analysis of Improvement in PixEdit

PixEdit introduces two significant advancements over AURORA:

**Upgraded Backbone Diffusion Model:** Transitioning from SD1.5 to PixArt-Sigma.

**Refined Pre-Training Data:** Replacing the Instruct-Pix2Pix dataset with the real-only subset of the Seed-Edit dataset for pre-training.

To isolate the impact of change (2), we compare PixEdit

with SD1.5 adapted for editing using the same Seed-Edit dataset and evaluate their performance on the  $K = 7$  subset of the ConceptMix benchmark. Results in Table 8 show that SD1.5 adapted with Seed-Edit data, and subsequently fine-tuned using the Aurora data-mixture, performs slightly better than AURORA alone on average. This highlights the advantage of incorporating real-world, complex pre-training data.

PixEdit, leveraging its more powerful PixArt-Sigma backbone and a compositional text-encoder, achieves the best performance on average. This underscores the benefits of combining an advanced diffusion model with improved text-image alignment for editing tasks.

Given a seed image that captures the fundamental structure of a prompt given to a text-to-image model. Generate a set of precise and object centric edit instructions for sequentially modifying the generated seed image. Adhere to the following guidelines:

1. Create the edit instructions with clarity , focussing on object-centric details.
2. Avoid generating vague or ambiguous edit instructions.

Refer to the provided samples to guide your response:

**INPUT:** A woman wearing a white shirt and grey shorts using a shovel to dig in snow.  
**ASSISTANT:** Let's think step by step:  
 Alright, let's break down your example into a detailed plan following the rules.

**IMAGE:**



**Original Prompt:**  
 "A woman wearing a white shirt and gray shorts using a shovel to dig in snow ."

**Image Re-Caption:**  
 A woman clad in a rich brown coat and stylish cap is energetically excavating through the pristine snow, wielding a sturdy shovel with determination.

**Reasoning:**  
 Based on the current image caption the action being performed by the woman matches the input, however these are the following errors that this image has :  
 1. The woman is wearing a coat while the prompt requires her to wear white shirt with grey shorts. Therefore we can use the following editing instructions to correct these.

**Edits:**  
 1. Make the woman wear shirt and shorts.  
 2. Color the shirt white.  
 3. Color the shorts grey.

Figure 9. System-prompt and few-shot example format for GraPE<sub>naive</sub>

Method	Concept K=7		
	GraPE <sub>AURORA</sub>	GraPE <sub>SD1.5*</sub>	GraPE <sub>PixEdit</sub>
Stable-Diffusion v1.5 [42]	0.598 ±0.003	<b>0.618</b> ±0.000	0.610 ±0.005
Structure Diffusion [16]	<b>0.612</b> ±0.003	0.596 ±0.005	0.571 ±0.002
Stable Diffusion v2.1 [42]	0.613 ±0.009	0.588 ±0.003	<b>0.626</b> ±0.002
LMD [31]	0.625 ±0.005	0.649 ±0.004	<b>0.668</b> ±0.009
SD-XL [37]	0.624 ±0.010	0.615 ±0.008	<b>0.628</b> ±0.002
PixArt-α [10]	0.624 ±0.006	0.614 ±0.006	<b>0.625</b> ±0.002
DeepFloyd IF [2]	0.638 ±0.006	<b>0.672</b> ±0.007	0.662 ±0.006
PlaygroundV2.5 [30]	0.611 ±0.002	<b>0.722</b> ±0.004	0.640 ±0.005
Dalle3 [5]	0.718 ±0.006	0.719 ±0.007	<b>0.737</b> ±0.004
Stable Diffusion v3.5 Large [46]	0.765 ±0.003	0.761 ±0.007	<b>0.784</b> ±0.005

Table 8. GPT-QA scores using GraPE with different editing models SD1.5\* vs AURORA vs PixEdit. Note GraPE<sub>SD1.5\*</sub> represents the editing model based on SD1.5 trained with refined pre-training data and fine-tuned on Aurora data-mixture.

## 10. Additional Qualitative Examples:

**Errors due to Planner vs Editing Model** We try to understand the broader context, consider a hypothetical set of 100 plans (randomly sampled from the human study on MLLM planners). We look at the breakup of these 100 plans and how the planner and editing-model fare on them. ~ 74.5% of the plans are deemed to have generation errors,

~ 55.5% (41.3 plans in absolute terms) of which are completely corrected at the planner level (considering GPT-4o) indicating a strong alignment between the generated plan and the text prompt. We then note that about one-third of these (~ 30.5% or 12.5 images) are successfully edited to produce completely correct final images, showcasing the editing model’s ability to translate plans into precise vi-

sual outputs. The remaining two-thirds of the aligned plans ( $\sim 69.5\%$  or 28.8 images) fall short, with errors entirely attributable to the editing model. This underscores that while the planner demonstrates robust performance in generating accurate and contextually aligned plans, the editing model remains the predominant source of errors.

**Editing Errors in Perfect Plans** The subset of 46 edit plans showing full improvement as per the majority of human annotators chosen for qualitative error analysis predominantly involved plans requiring "add" actions, with relatively fewer plans focusing on "remove" or "modify" actions. This can also be seen from the breakdown of the corresponding final edited images on the basis of mentioned actions. This distribution reflects the natural bias of generative models which often fail to identify certain textual elements and the image generated then has errors which can be fixed by "add" edit actions.

**Plans showing no or partial improvement** Finally, considering the plans which show no/partial improvement we note that distinct patterns emerge providing further insights into the limitations of the planner. The plans with no improvement are primarily empty plans generated by the planner, where it fails to propose any actionable edits to address the misalignment between the text prompt and the image. On the other hand, the plans with partial improvement typically exhibit incomplete details, where some, but not all, elements of the prompt are addressed. A small subset of these also includes hallucinations, where the planner mistakenly assumes the presence of an object in the image that does not actually exist. These observations highlight areas where the planner could be enhanced, particularly in terms of improving its ability to identify actionable edits and reducing cases of over- or under-specification in the generated plans.



Method	Concept K=1		Concept K=3		Concept K=5		Concept K=7	
	Base	GraPE <sub>AURORA</sub>	Base	GraPE <sub>AURORA</sub>	Base	GraPE <sub>AURORA</sub>	Base	GraPE <sub>AURORA</sub>
Stable-Diffusion v1.5 [42]	0.808 ±0.009	<b>0.863</b> ±0.005	0.606 ±0.018	<b>0.717</b> ±0.024	0.497 ±0.010	<b>0.688</b> ±0.010	0.450 ±0.005	<b>0.598</b> ±0.003
Structure Diffusion [16]	0.823 ±0.002	<b>0.895</b> ±0.004	0.606 ±0.002	<b>0.734</b> ±0.021	0.542 ±0.014	<b>0.698</b> ±0.015	0.447 ±0.001	<b>0.612</b> ±0.003
Stable Diffusion v2.1 [42]	0.833 ±0.002	<b>0.875</b> ±0.014	0.639 ±0.014	<b>0.719</b> ±0.012	0.579 ±0.012	<b>0.694</b> ±0.009	0.466 ±0.002	<b>0.613</b> ±0.009
LMD [31]	0.855 ±0.004	<b>0.907</b> ±0.008	0.711 ±0.008	<b>0.765</b> ±0.005	0.643 ±0.011	<b>0.721</b> ±0.012	0.591 ±0.002	<b>0.625</b> ±0.005
SD-XL [37]	0.848 ±0.010	<b>0.893</b> ±0.006	0.708 ±0.018	<b>0.740</b> ±0.013	0.635 ±0.014	<b>0.692</b> ±0.018	0.520 ±0.003	<b>0.624</b> ±0.010
PixArt- $\alpha$ [10]	0.813 ±0.010	<b>0.883</b> ±0.008	0.668 ±0.011	<b>0.717</b> ±0.012	0.649 ±0.011	<b>0.711</b> ±0.022	0.507 ±0.001	<b>0.624</b> ±0.006
DeepFloyd IF [2]	0.883 ±0.009	<b>0.898</b> ±0.005	0.680 ±0.016	<b>0.735</b> ±0.016	0.663 ±0.014	<b>0.713</b> ±0.007	0.583 ±0.002	<b>0.638</b> ±0.006
PlaygroundV2.5 [30]	0.908 ±0.010	<b>0.945</b> ±0.004	0.737 ±0.023	<b>0.783</b> ±0.012	0.658 ±0.015	<b>0.709</b> ±0.006	0.540 ±0.003	<b>0.611</b> ±0.002
Dalle3 [5]	0.947 ±0.002	<b>0.950</b> ±0.004	<b>0.832</b> ±0.012	0.809 ±0.007	<b>0.812</b> ±0.014	0.783 ±0.015	<b>0.728</b> ±0.006	0.718 ±0.006
Stable Diffusion v3.5 Large [46]	<b>0.927</b> ±0.005	<b>0.927</b> ±0.002	<b>0.815</b> ±0.002	0.812 ±0.004	<b>0.803</b> ±0.003	0.792 ±0.002	0.759 ±0.004	<b>0.765</b> ±0.003

Table 9. Results on Concept-mix benchmark - GraPE<sub>AURORA</sub>

Method	Base		GraPE <sub>AURORA</sub>		GraPE <sub>PixEdit</sub>	
	DSG (w/o dep)	DSG	DSG (w/o dep)	DSG	DSG (w/o dep)	DSG
Stable-Diffusion v1.5 [42]	0.679 ±0.006	0.654 ±0.007	0.758 ±0.008	0.746 ±0.009	<b>0.788</b> ±0.006	<b>0.769</b> ±0.006
Structure Diffusion [16]	0.714 ±0.004	0.697 ±0.004	0.760 ±0.017	0.741 ±0.014	<b>0.794</b> ±0.006	<b>0.777</b> ±0.008
Stable Diffusion v2.1 [42]	0.701 ±0.001	0.688 ±0.000	0.757 ±0.012	0.742 ±0.014	<b>0.787</b> ±0.002	<b>0.772</b> ±0.004
LMD [31]	0.782 ±0.001	0.777 ±0.002	0.805 ±0.005	0.790 ±0.006	<b>0.836</b> ±0.008	<b>0.826</b> ±0.008
SD-XL [37]	0.805 ±0.002	0.795 ±0.001	<b>0.844</b> ±0.004	<b>0.829</b> ±0.005	0.839 ±0.003	0.828 ±0.002
PixArt- $\alpha$ [10]	0.710 ±0.004	0.699 ±0.006	0.758 ±0.006	0.744 ±0.009	<b>0.796</b> ±0.006	<b>0.788</b> ±0.006
DeepFloyd IF [2]	0.822 ±0.002	0.811 ±0.004	0.840 ±0.008	0.819 ±0.011	<b>0.852</b> ±0.001	<b>0.840</b> ±0.001
PlaygroundV2.5 [30]	0.788 ±0.002	0.775 ±0.004	0.821 ±0.011	0.801 ±0.011	<b>0.825</b> ±0.006	<b>0.808</b> ±0.006
Dalle3 [5]	0.932 ±0.003	0.924 ±0.003	0.913 ±0.014	0.910 ±0.015	<b>0.934</b> ±0.003	<b>0.928</b> ±0.002
Stable Diffusion v3.5 Large [46]	0.871 ±0.003	0.864 ±0.004	0.876 ±0.006	0.866 ±0.007	<b>0.888</b> ±0.003	<b>0.880</b> ±0.003

Table 10. GPT-QA scores on T2I Comp-Bench, GraPE<sub>X</sub> refers to using GraPE with X editing model

Method	Base		GraPE <sub>AURORA</sub>		GraPE <sub>PixEdit</sub>	
	DSG (w/o dep)	DSG	DSG (w/o dep)	DSG	DSG (w/o dep)	DSG
Stable-Diffusion v1.5 [42]	0.703 ±0.004	0.640 ±0.005	<b>0.800</b> ±0.016	<b>0.748</b> ±0.017	0.787 ±0.005	0.742 ±0.006
Structure Diffusion [16]	0.725 ±0.012	0.648 ±0.014	<b>0.816</b> ±0.003	<b>0.766</b> ±0.002	0.802 ±0.005	0.732 ±0.006
Stable Diffusion v2.1 [42]	0.720 ±0.009	0.656 ±0.007	0.800 ±0.010	0.740 ±0.010	<b>0.821</b> ±0.001	<b>0.757</b> ±0.001
LMD [31]	0.682 ±0.006	0.640 ±0.012	<b>0.752</b> ±0.008	<b>0.697</b> ±0.008	0.744 ±0.002	<b>0.697</b> ±0.004
SD-XL [37]	0.695 ±0.004	0.613 ±0.003	0.775 ±0.010	0.720 ±0.011	<b>0.818</b> ±0.005	<b>0.758</b> ±0.007
PixArt- $\alpha$ [10]	0.735 ±0.006	0.663 ±0.005	0.792 ±0.005	0.733 ±0.005	<b>0.811</b> ±0.011	<b>0.767</b> ±0.011
DeepFloyd IF [2]	0.831 ±0.003	0.792 ±0.003	0.842 ±0.006	0.803 ±0.007	<b>0.877</b> ±0.003	<b>0.843</b> ±0.005
PlaygroundV2.5 [30]	0.789 ±0.005	0.735 ±0.007	0.806 ±0.010	0.759 ±0.012	<b>0.833</b> ±0.005	<b>0.789</b> ±0.004
Dalle3 [5]	0.891 ±0.004	0.857 ±0.005	0.890 ±0.002	0.847 ±0.001	<b>0.904</b> ±0.004	<b>0.873</b> ±0.007
Stable Diffusion v3.5 Large [46]	0.881 ±0.004	0.852 ±0.004	0.905 ±0.002	0.873 ±0.004	<b>0.909</b> ±0.003	<b>0.884</b> ±0.003

Table 11. GPT-QA scores on Flick-Bench, GraPE<sub>X</sub> refers to using GraPE with X editing model

### Text Prompt

The image is an oil painting in the style of impressionism. It shows two tiny corgis in front of a giraffe. The dogs are positioned nearer to the observer than the giraffe.

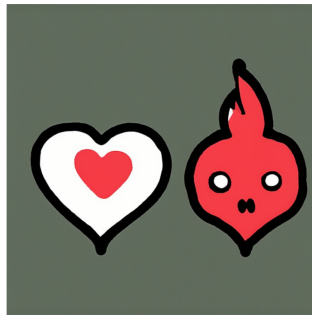
### Generated Image (Before Edit)



### Edit Plan

1. Ensure the two animals are clearly depicted as tiny corgis ✓
2. Reposition the corgis so they are nearer to the observer than the giraffe (giraffe not present after 1) ✗

A tiny, heart-shaped chicken is depicted in a pop art style.



1. Change the red character to resemble a chicken ✓  
(Heart Ignored) ✗

A purple glass-textured rose is positioned in the foreground of the image. In front of the rose, there are exactly four green spiders. The image also contains a skyscraper located behind the rose.



1. Add three more green spiders in front of the rose (1 spider hallucinated) ✗
2. Add a skyscraper behind the rose ✓

A huge dog stands with three tiny, fluffy blue cactuses touching its lowest point.



1. Change the cactuses to be fluffy ✓
2. Change the cactuses' color to blue ✓  
(cactus count ignored) ✗

Figure 10. Results illustrating failure cases of generated Edit Plans








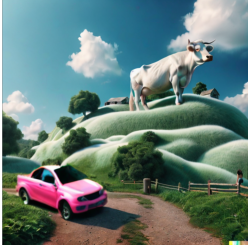
Text Prompt	Generated Image	Edit Plan	Final Edited Image
<p>The image contains a square-shaped doll.</p>		<p>1. Change the doll's body to be square-shaped (distortion) ❌</p>	
<p>A tiny brown fork is positioned to the left of a knife.</p>		<p>1. Change the fork to be tiny (no change) ❌  2. Change the fork to be entirely brown (partial) ✅</p>	
<p>A painting in the impressionist style shows four forks placed at different angles in the foreground. In the background, a hill with a fluffy texture rises gently, while a volcano with a slightly jagged peak stands further behind.</p>		<p>1. Reduce the number of forks in the foreground to four (done with partial visibility) ❌</p>	
<p>On a fluffy-textured hill, there is a white cow. A tiny, metallic, pink car is parked next to the hill. Nearby, a human is standing.</p>		<p>1. Change the cow's color to white ✅  2. Change the car to be tiny and entirely pink ✅  3. Add a human standing nearby (not added) ❌</p>	

Figure 11. More results illustrating failure cases of Editing Models