

VISUALLY GUIDED DECODING: GRADIENT-FREE HARD PROMPT INVERSION WITH LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

Text-to-image generative models like DALL-E and Stable Diffusion have revolutionized visual content creation across various applications, including advertising, personalized media, and design prototyping. However, crafting effective textual prompts to guide these models remains challenging, often requiring extensive trial and error. Existing prompt inversion methods, such as soft and hard prompt techniques, suffer from issues like limited interpretability and incoherent prompt generation. To address these limitations, we introduce Visually Guided Decoding (VGD), a gradient-free approach that leverages large language models (LLMs) and CLIP-based guidance to generate coherent and semantically aligned prompts. VGD utilizes the robust text generation capabilities of LLMs to produce human-readable prompts while employing CLIP scores to ensure alignment with user-specified visual concepts. This method enhances the interpretability, generalization, and flexibility of prompt generation without the need for additional training. Our experiments demonstrate that VGD outperforms existing prompt inversion techniques in generating understandable and contextually relevant prompts, facilitating more intuitive and controllable interactions with text-to-image models.

1 INTRODUCTION

In recent years, image generative models such as DALL-E and Stable Diffusion have shown remarkable success in generating high-fidelity images (Ramesh et al., 2022; Rombach et al., 2022; Podell et al., 2024). These models have been widely used in a variety of applications, including visual content generation (e.g., advertisement, movie, game), personalized content generation (e.g., caricature, photo editing), and also prototyping (e.g., architecture and product design). Previous studies have shown that, just as humans can draw an image of an object solely based on detailed descriptions (e.g., criminal composite sketch), generative models can generate images of objects using a well-crafted prompt, even if they have not been trained on those specific objects (Gal et al., 2023; Everaert et al., 2023).

A well-known drawback of these approaches is the difficulty in finding a textual description (a.k.a., prompt) that effectively guides the generation of the desired visual content (Hao et al., 2024). For example, to create a culinary masterpiece (see Figure 1), one should have good knowledge of exquisite cooking styles such as French cuisine or molecular gastronomy. Without the expertise to craft precise prompts, users have to rely on laborious trial and error to generate the desired images. Automation of the prompt generation process will reduce the time and effort required to generate the desired image.

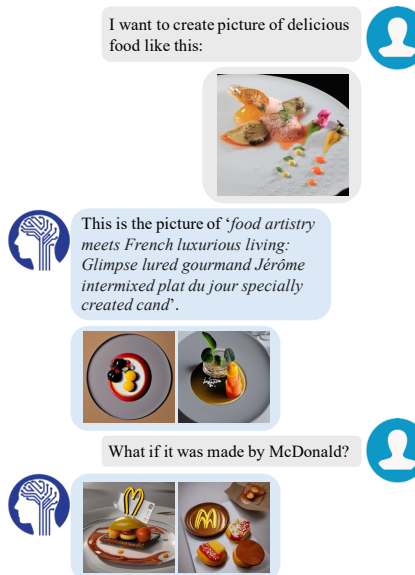


Figure 1: Visually Guided Decoding (VGD) works with any LLM without extra training, making it easy to integrate into a chat-based interface that offers interpretable and controllable text-to-image generation.



Figure 2: VGD generates fully interpretable prompts that enhance generalizability across tasks and models in text-to-image generation.

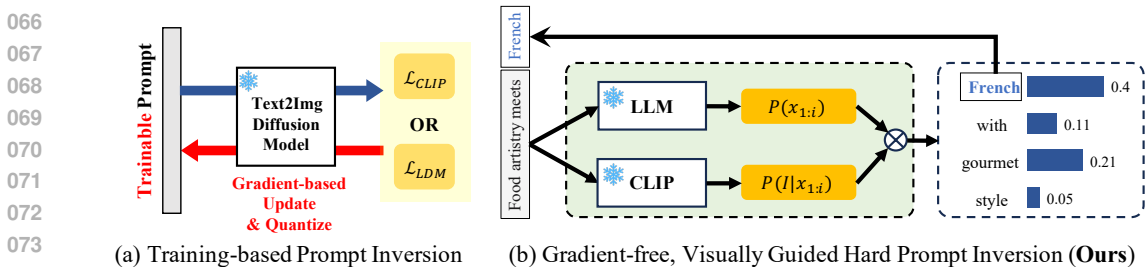


Figure 3: While conventional prompt inversion techniques update prompt embeddings through gradient-based optimization and quantization, VGD is a gradient-free technique that utilizes large language models and CLIP to generate relevant sentences.

To address the difficulty in finding the prompt, soft prompt inversion techniques have been proposed (Gal et al., 2023; Kumari et al., 2023; Ruiz et al., 2023; Voynov et al., 2023). In these approaches, following the prompt learning works in natural language processing (Li & Liang, 2021; Lester et al., 2021; Liu et al., 2022), an image is transformed into the embedding vector. While this so-called soft prompt contains useful information about the image, since it is a string of numbers, it is not human-readable and hence users cannot handle it easily (see Figure 2).

Recently, to mitigate the lack of interpretability and generalization capability, hard prompt inversion has been proposed (Mahajan et al., 2024; Wen et al., 2024). In this approach, soft prompts (embedding vectors) are projected to the nearest neighbors in a predefined vocabulary (see Figure 3 (a)). Specifically, the learnable input vectors for the conditioning network (Radford et al., 2021) in Stable Diffusion are trained to reconstruct user-provided images through an iterative de-noising process (*a.k.a.*, reverse diffusion process) and then mapped to the nearest word embedding in the model’s vocabulary. During this training process, the long gradient path in the reverse diffusion process to the input prompt often results in a vanishing gradient problem. Further, since the hard prompt is a simple collection of words (see Figure 2), it is difficult to identify which words are essential in generating the desired images.

An aim of this paper is to propose a technique that generates a fully interpretable prompt. Essence of our technique, henceforth referred to as visually guided decoding (VGD), is to generate contextually meaningful sentence for hard prompt through the token generation process of large language models (LLMs) (Ouyang et al., 2022; Touvron et al., 2023). To this end, during the token generation process, we employ CLIP to align the generated token with the provided images (see Figure 3 (b)). The benefit of the proposed approach is that one can fully understand the meaning of the prompt (*interpretability*) and synthesize image for various applications (*generalizability*) (*e.g.*, style transfer, image editing, and image variation) with greater control over the output image (*flexibility*). Since VGD seamlessly stitches LLMs and text-to-image models without the need for training (*i.e.*, gradient-free), VGD provides a flexible and efficient solution for chat-interfaced image generation services using LLM (*e.g.*, DALL-E extension on ChatGPT) (see Figure 1).

In our experiments, VGD qualitatively and quantitatively achieves state-of-the-art (SOTA) performance, demonstrating superior interpretability, generalizability, and flexibility in text-to-image generation compared to previous soft and hard prompt inversion methods. We also show that VGD is compatible with various LLMs including LLaMA2, LLaMA3 and Mistral.

2 BACKGROUND AND RELATED WORKS

2.1 CLIP MODEL

CLIP model aligns semantically related visual and textual content within a shared representation space (Radford et al., 2021). It comprises an image encoder $\text{CLIP}_{\text{img}}(\cdot)$ and a text encoder $\text{CLIP}_{\text{txt}}(\cdot)$. The image encoder encodes an input image I into a visual embedding $\mathbf{f}_{\text{img}} = \text{CLIP}_{\text{img}}(I)$. Similarly, The text encoder processes the input text $T = [x_1^{\text{txt}}, x_2^{\text{txt}}, \dots, x_N^{\text{txt}}]$, yielding a textual embedding $\mathbf{f}_{\text{txt}} = \text{CLIP}_{\text{txt}}(T)$. CLIP is trained using a contrastive learning approach that maximizes the similarity between \mathbf{f}_{txt} and \mathbf{f}_{img} for matched sentence-image pair, while minimizing the similarity for unmatched pairs. The similarity is defined as:

$$\text{Sim}(\mathbf{f}_{\text{txt}}, \mathbf{f}_{\text{img}}) := s \cdot \frac{\mathbf{f}_{\text{txt}} \cdot \mathbf{f}_{\text{img}}}{\|\mathbf{f}_{\text{txt}}\|_2 \|\mathbf{f}_{\text{img}}\|_2} \quad (1)$$

where s is a scalar scaling factor. Then, the training objective of CLIP is formulated as:

$$\underset{\text{CLIP}_{\text{img}}, \text{CLIP}_{\text{txt}}}{\text{Maximize}} \quad \log P(T|I) + \log P(I|T), \quad (2)$$

$$P(T|I) = \frac{\exp(\text{Sim}(\mathbf{f}_{\text{txt}} \cdot \mathbf{f}_{\text{img}}))}{\sum_{j=1}^B \exp(\text{Sim}(\mathbf{f}_{\text{txt}_j} \cdot \mathbf{f}_{\text{img}}))}, \quad P(I|T) = \frac{\exp(\text{Sim}(\mathbf{f}_{\text{txt}} \cdot \mathbf{f}_{\text{img}}))}{\sum_{j=1}^B \exp(\text{Sim}(\mathbf{f}_{\text{txt}} \cdot \mathbf{f}_{\text{img}_j}))}, \quad (3)$$

where B indicates the number of text-image pairs in the mini-batch.

2.2 TEXT-TO-IMAGE DIFFUSION MODEL

Text-to-image models generate an image I that maximizes $P(I|T)$ given textual description T . Modern text-to-image models, such as Stable Diffusion, are based on Latent Diffusion Model (LDM), which usually takes CLIP text embedding \mathbf{f}_{txt} as a condition to generate an image. The output of LDM, denoted as $\epsilon_\theta(\mathbf{z}_k, \mathbf{f}_{\text{txt}}, k)$, represents the predicted denoising result after k diffusion steps, starting from a random Gaussian noise \mathbf{z}_k conditioned on \mathbf{f}_{txt} . The training objective of LDM is as follows:

$$L_{\text{LDM}} = \mathbb{E}_{\epsilon, \mathbf{z}, \mathbf{f}_{\text{txt}}, k} [\|\epsilon - \epsilon_\theta(\mathbf{z}_k, \mathbf{f}_{\text{txt}}, k)\|_2^2], \quad (4)$$

where $\epsilon \sim N(0, \mathbf{I})$ indicates the noise used to corrupt clean latent variables. Our goal is to find the optimal text condition T that yields the image containing desired visual concepts.

2.3 PROMPT INVERSION

Soft Prompt Inversion Following prompt-tuning approach Li & Liang (2021), Soft Prompt Inversion techniques (Gal et al., 2023; Kumari et al., 2023) extend the vocabulary of the model with a new special token S_* , which embeds the objects from the provided images. During the generation process, the continuous embedding vector of S_* is prepended to the embeddings of the input text tokens (e.g., “A photo of S_* ”, “A rendition of S_* ”) and then used as an input to the LDM model. While effective in generating specific visual concepts, the soft prompt is in continuous vector form which is not human-readable and difficult for users to modify to generate desired images.

Hard Prompt Inversion To mitigate the limitations of Soft Prompt Inversion, PEZ (Wen et al., 2024) and PH2P (Mahajan et al., 2024) have recently introduced hard prompt (discrete prompt) inversion approaches. On top of the soft prompt inversion methods, hard prompt inversion methods leverage projected gradient descent as a quantization technique to produce the learned soft prompt within the word embedding space. These hard prompts, once generated, often appear as a disjointed collection of words (commonly referred to as the incoherence problem), making it difficult to determine which words need to be modified to generate the desired images. Moreover, hard prompt inversion involves a complex training process due to discrete optimization and multi-stage pipelining. VGD does not need for gradient-based optimization, reducing computational requirements and streamlining the process.

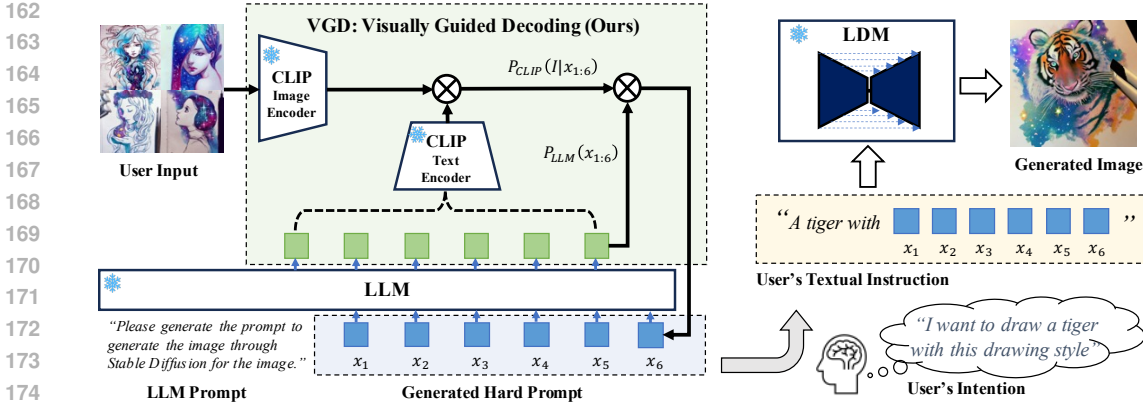


Figure 4: Overview of the proposed hard prompt inversion method, VGD, which seamlessly integrates LLM, CLIP, and LDM for user-friendly image generation process.

Image Captioning One might think that generating an image directly from captions produced by large multi-modal models (LMMs) (Li et al., 2022; 2023; Alayrac et al., 2022; Liu et al., 2024) would be an alternative option to the hard prompt inversion. However, LMM-generated captions often lack the fine details necessary for detailed image synthesis control. To complement the missing information, prompt generation services like CLIP-Interrogator¹ have been introduced to complement image captioning models such as LLaVA (see Appendix A.4 for more details).

3 GRADIENT FREE PROMPT INVERSION WITH LANGUAGE MODELS

3.1 INTERPRETABILITY DEGRADATION OF PREVIOUS HARD PROMPT INVERSION

The goal of hard prompt inversion is to determine the text prompt $T = [x_1^{\text{txt}}, x_2^{\text{txt}}, \dots, x_N^{\text{txt}}]$ that maximizes the LDM’s probability $P(I|T)$ for a target image I (where typically $N \leq 77$), with each x_i^{txt} selected from the model’s vocabulary. By applying Bayes’ theorem, this objective can be reformulated as:

$$P(I|T) = \frac{P(I)P(T|I)}{P(T)}. \tag{5}$$

Since the prior probability of the image $P(I)$ is independent of T , finding the optimal \hat{T} is equivalent to finding T that maximizes $P(T|I)/P(T)$, expressed as:

$$\hat{T} = \operatorname{argmax}_T \frac{P(T|I)}{P(T)}. \tag{6}$$

As indicated in Eq. 6, the inversion process maximizes $P(T|I)$ (a.k.a., image captioning objective) while simultaneously minimizing the prior probability of text, $P(T)$. Minimizing $P(T)$ can cause the hard prompt T to contain uncommon or awkward phrasing, leading to reduced interpretability. We suggest that this is the primary reason why existing hard prompt inversion techniques exhibit lower interpretability compared to our approach (see Section 4.3 for qualitative comparison).

3.2 VISUALLY GUIDED DECODING

Step 1 - Problem Formulation Our goal is to find an *interpretable* text prompt without training prompt embeddings, using the *gradient-free* approach. Inspired by the noisy channel model (Jurafsky, 2000; Brown et al., 1993), widely used in machine translation and speech recognition, we model the prompt discovery process by integrating an LDM objective with a regularization term for language modeling. The objective is formulated as:

$$\hat{T} = \operatorname{argmax}_T P(I|T)P(T)^\alpha, \tag{7}$$

where α is a hyperparameter that balances the influence of $P(I|T)$ and $P(T)$.

¹<https://github.com/pharmapsychotic/clip-interrogator>

Step 2 - Approximation with CLIP Score Computing $P(I|T)$ by forward and reverse diffusion processes at each token generation step is computationally intractable. To alleviate this computational burden, we approximate $P(I|T)$ using CLIP, such that $P(I|T) \approx P_{\text{CLIP}}(I|T)$. We assert that this approximation is justified by the fact that Stable Diffusion utilizes the frozen CLIP text encoder. We empirically show that when a CLIP model different from the one aligned with Stable Diffusion is used, the approximation no longer holds, leading to a decline in performance (see Section 4.4 for the ablation study on CLIP image encoders).

In addition, we leverage an external language model such as LLaMA (Touvron et al., 2023) to model the $P(T) \approx P_{\text{LLM}}(T)$. This leads to our main objective:

$$\hat{T} = \underset{T}{\operatorname{argmax}} P_{\text{CLIP}}(I|T)P_{\text{LLM}}(T)^\alpha. \quad (8)$$

Step 3 - Token-by-Token Generation We can further express Eq. 8 using a left-to-right decomposition of the text probability $P_{\text{LLM}}(T)$:

$$P_{\text{CLIP}}(I|x_{1:N}^{\text{txt}}) \prod_{i=1}^N P_{\text{LLM}}(x_i^{\text{txt}}|x_{1:i-1}^{\text{txt}})^\alpha. \quad (9)$$

This decomposition facilitates token-by-token text generation using a beam search decoding strategy (Anderson et al., 2017; Post & Vilar, 2018; Hu et al., 2019; Holtzman et al., 2020). Specifically, we iteratively select the i -th token x_i^{txt} that maximizes the objective, Eq. 9, given previous tokens $x_{1:i-1}^{\text{txt}}$. In doing so, VGD generates prompts that are semantically aligned with the target image (*via* CLIP), linguistically coherent, and interpretable (*via* LLM).

3.3 IMPLEMENTATION OF VISUALLY GUIDED DECODING

LLM Prompting Instead of fine-tuning the LLM for each image generation task, we design specific system and user prompts tailored to different tasks. The LLM is then queried to generate tokens as if it were creating text prompts for text-to-image models (see Appendix A.1 for details).

Beam Initialization Providing the LLM with image-related words at the start of the generation process increases the probability of generating tokens that align with the image in subsequent steps. To achieve this, we first compare all text embeddings \mathbf{f}_{txt} in the CLIP vocabulary with the target image embedding \mathbf{f}_{img} , selecting the top- M tokens $x_{1:M}^{\text{txt}}$ as the initial input prefix for the LLM. Empirically, we find that even $M = 1$ is sufficient for VGD. This process is computationally efficient since all text embeddings can be precomputed.

Beam Expansion and Pruning Since CLIP does not provide next-token probabilities, we select beam search candidates using the LLM’s next-token prediction results. Specifically, we expand each beam by appending K next-token candidates with the highest probabilities $P_{\text{LLM}}(x_i^{\text{txt}}|x_{1:i-1}^{\text{txt}})$. We then prune the expanded beam candidates (K^2 candidates from K beams) to retain only K beams based on the combined score. Unless otherwise specified, we set $K = 10$ for the experiments.

Beam Search Termination The beam search terminates when either 1) beam expansion fails to improve the score (Eq. 9), or 2) a pre-determined maximum prompt length is reached. The final result \hat{T} is selected among the candidates based on the score.

4 EXPERIMENTS

4.1 SETUP

Datasets We conduct experiments on four datasets with diverse distributions: LAION-400M (Schuhmann et al., 2021; 2022), MS COCO (Lin et al., 2014), Celeb-A (Liu et al., 2015), and Lexica.art². Following PEZ (Wen et al., 2024), we randomly sample 100 images from each dataset and evaluate prompt inversion methods across 5 runs using different random seeds. See Appendix A.3 for more details.

²<https://huggingface.co/datasets/Gustavosta/Stable-Diffusion-Prompts>

Table 1: Image quality (CLIP-I score) and prompt quality (BERTScore) comparison.

Method	#Tokens	LAION	MS COCO	Celeb-A	Lexica.art	MS COCO (BERTScore)			Lexica.art (BERTScore)		
						Precision	Recall	F1	Precision	Recall	F1
CLIP-I Score											
Textual Inversion	1	0.388	0.569	0.522	0.697	-	-	-	-	-	-
LLaVA-1.5	32	0.513	0.642	0.463	0.580	0.914	0.918	0.916	0.858	0.780	0.817
+ CLIP Interrogator	~77	0.540	0.685	0.511	0.762	0.794	0.898	0.843	0.818	0.819	0.819
LLaVA-1.5 + VGD	32	0.573	0.718	0.546	0.769	0.831	0.879	0.854	0.835	0.793	0.813
	~77	0.569	0.724	0.544	0.785	0.800	0.874	0.835	0.810	0.794	0.802
PEZ	16	0.538	0.687	0.622	0.743	0.760	0.834	0.795	0.772	0.783	0.777
	32	0.530	0.685	0.619	0.745	0.736	0.830	0.780	0.752	0.784	0.768
	64	0.507	0.670	0.584	0.728	0.715	0.825	0.766	0.734	0.783	0.758
VGD (Ours)	16	0.484	0.650	0.482	0.700	0.833	0.862	0.847	0.827	0.779	0.802
	32	0.493	0.670	0.506	0.735	0.818	0.868	0.842	0.816	0.786	0.801
	64	0.511	0.678	0.514	0.753	0.787	0.863	0.823	0.799	0.789	0.794
	~77	0.510	0.678	0.513	0.754	0.788	0.864	0.824	0.801	0.791	0.795



Figure 5: Generated images using hard prompts produced by the proposed VGD.

Baselines We compare the proposed method with the hard prompt inversion (PEZ (Wen et al., 2024)), soft prompt inversion (Textual Inversion (Gal et al., 2023)), LLaVA-1.5 (Liu et al., 2024) generated caption, and LLaVA-1.5 combined with CLIP-Interrogator. Images are generated with the Stable Diffusion 2.1-768 model across all comparisons (Podell et al., 2024).

Evaluation Metric We evaluate the quality of the prompt via image embedding similarity between the target image and an image generated using the generated hard prompt. We call this similarity-based metric CLIP-I score. To ensure fairness, we utilize different CLIP models for generation and evaluation: CLIP-ViT-H-14 for generation and larger CLIP-ViT-G-14 for similarity evaluation. Following (Mahajan et al., 2024), we also compare the contextual similarity between the generated prompt and the ground-truth annotations (captions) of the target image, using BERTScore (Zhang et al., 2020) as metric.

4.2 IMAGE GENERATION WITH VGD PROMPTS

In Figure 5, we qualitatively demonstrate that VGD generates realistic and diverse images without overfitting to the original target image.

Qualitative Evaluation As shown in Table 1, the images generated by VGD are the most similar to the original images, as measured by the CLIP-I score, even without using a gradient-based optimization process like PEZ. When using LLaVA-1.5 as the language model, VGD achieves SOTA CLIP-I scores across all four datasets, surpassing the strong baseline that exploits an external database

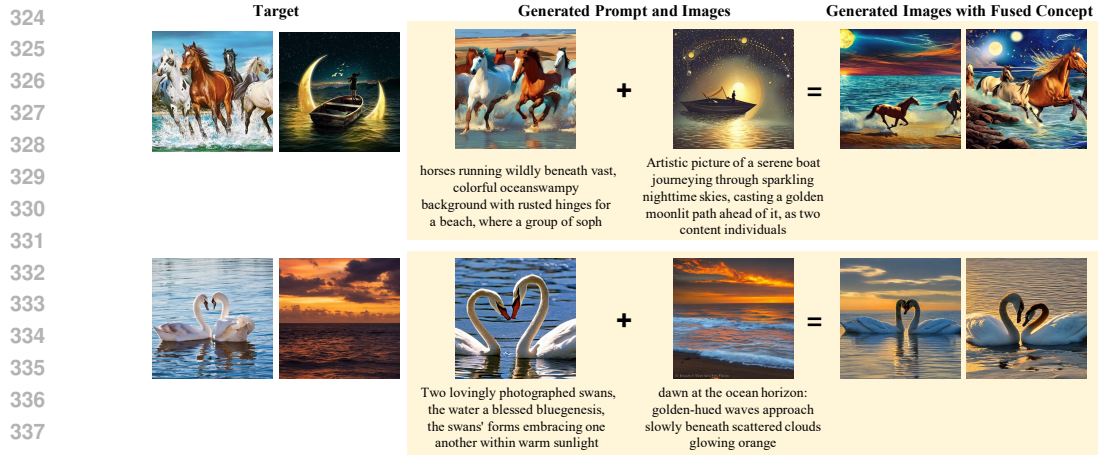


Figure 6: Multi-concept image generation with VGD. We show that the hard prompts obtained from two different images can be concatenated to fuse the semantic concepts.

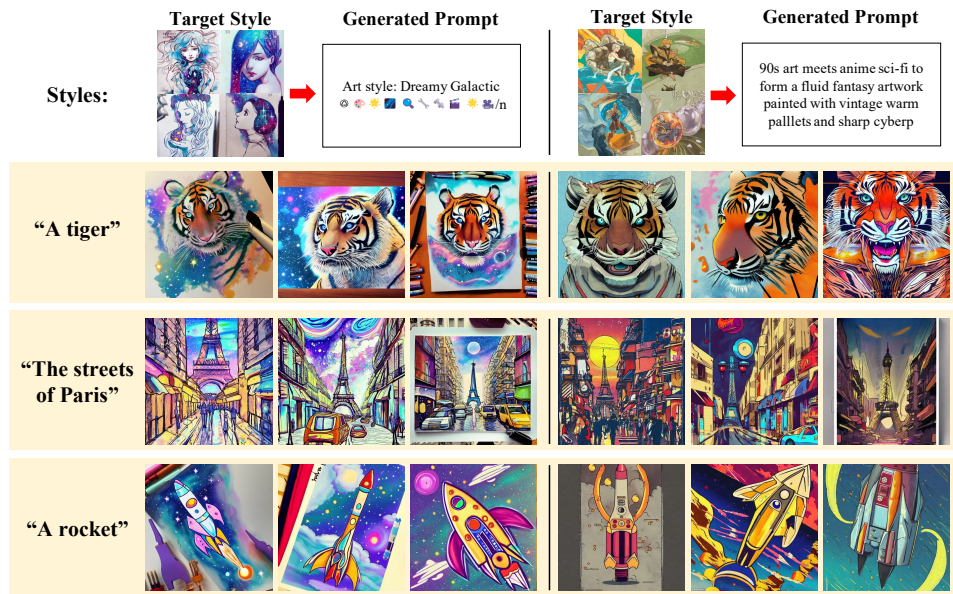


Figure 7: Decoded hard prompt for style transfer. Given several sample images with the same style, we can extract the style with a hard prompt and transfer it to other objects or scenes.

(LLaVA 1.5 + CLIP Interrogator) by a considerable margin. Furthermore, unlike PEZ, VGD exhibits a positive correlation between the CLIP-I score and the number of tokens, highlighting the effectiveness of VGD’s token-by-token generation process.

Multi-concept Generation VGD is capable of fusing multiple concepts from different images without additional effort. For each image with its own distinct concept, we individually extract prompts and concatenate them to generate a single image that combines their concepts. Figure 6 presents two examples of multi-concept generation, illustrating that VGD accurately extracts the key elements from each image so that the generated image can contain multiple concepts without omitting any. We show the example of multi-concept generation of 5 images with VGD and PEZ in Figure 21 of the Appendix.

Style Transfer VGD can also be easily adapted to style transfer. Given several examples (typically 4 images) that share the same style, we extract their shared style characteristics as a single hard prompt and use this prompt to apply the style to new objects or scenes. Figure 7 shows two examples of style transfer, demonstrating that VGD correctly embeds the shared style elements in the prompt, which can be easily combined with new objects.

4.3 QUALITY OF VGD PROMPTS

Quantitative Evaluation To assess the coherence of generated prompts, we measure the BERTScore between the prompts generated by each method and the ground-truth captions. As shown in Table 1, VGD achieves a significant improvement in BERTScore compared to PEZ, especially when using more than 64 tokens. Although LLaVA-based baselines show higher BERTScore, their CLIP-I scores are worse than those of VGD. This implies that the LLaVA captions fail to describe the target image in detail, even though the text itself may be semantically similar to ground-truth captions.

Interpretability When a prompt closely aligns with the image content and is interpretable by humans, it becomes easier for users to modify the image in the desired direction by making slight edits to the prompt. To evaluate the interpretability of the generated prompts, we make minimal modifications to the prompt and observe the effects. For example, we change “winter” to “spring” and see if the generated image successfully reflects the change. As illustrated in Figure 9, we compare VGD with previous gradient-based prompt inversion techniques, namely PEZ (hard) and Textual Inversion (soft).

The results indicate that VGD exhibits superior interpretability. When we adjust the number of objects, background elements, atmosphere (seasons), and species in the prompt, VGD successfully incorporates these changes. The relevant terms are already present in the original prompt, enabling accurate image generation with the modified prompts. In contrast, the competitors struggle to produce the desired outcomes. For instance, when we try to generate four wolves, both PEZ and Textual Inversion fail to change the number of wolves included.

Prompt Distillation We can also utilize prompt inversion to reduce the length of prompts while preserving their capability, *a.k.a.* prompt distillation. Long prompts may contain redundant and unimportant information, especially when hand-crafted. This becomes problematic, particularly when there exists a limit on the number of tokens that the model can process. Using VGD, we can generate a shorter prompt that preserves the key aspects of the original prompt. As illustrated in Figure 10, the images generated with distilled prompts are very similar to the original image. In addition, images generated with distilled prompts do not show significant performance degradation on CLIP-I score compared to the original prompt, as shown in Figure 8.

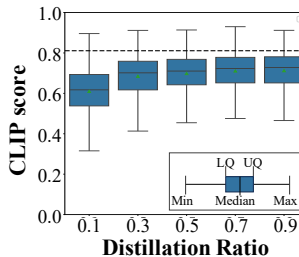


Figure 8: Prompt distillation on Lexica.art dataset.

4.4 EFFECT OF CLIP AND LANGUAGE MODEL

Table 2: Comparison of VGD variants.

Model	CLIP Score	BERTScore		
		Precision	Recall	F1
VGD	0.735	0.816	0.786	0.801
LLM only	0.487	0.811	0.773	0.791
CLIP only	0.768	0.727	0.779	0.751

Table 3: Ablation on CLIP models.

CLIP Model	#Tokens	CLIP Score
laion/ViT-H-14 (original)	32	0.735
openai/ViT-B-32	32	0.655
openai/ViT-L-14	32	0.626

Role of CLIP and LLM To investigate the effect of using both CLIP and LLM, we modify VGD in two ways and compare these variants on the Lexica.art dataset. In the ‘LLM only’ variant, VGD selects the next token based solely on $P_{\text{LLM}}(x_i^{\text{txt}}|x_{1:i-1}^{\text{txt}})$ without considering CLIP score, operating like typical LLM text decoding. In the ‘CLIP only’ variant, VGD determines the next token using only the CLIP model. Table 2 shows that the ‘LLM only’ variant suffers from a drastic performance degradation in the CLIP score because it does not consider the target image at all. On the other hand, the ‘CLIP only’ variant achieves the highest CLIP score but the lowest BERTScore, suggesting that the generated prompt is less coherent and harder to interpret. In comparison, the proposed VGD demonstrates strong and balanced performance on both the CLIP score and BERTScore.

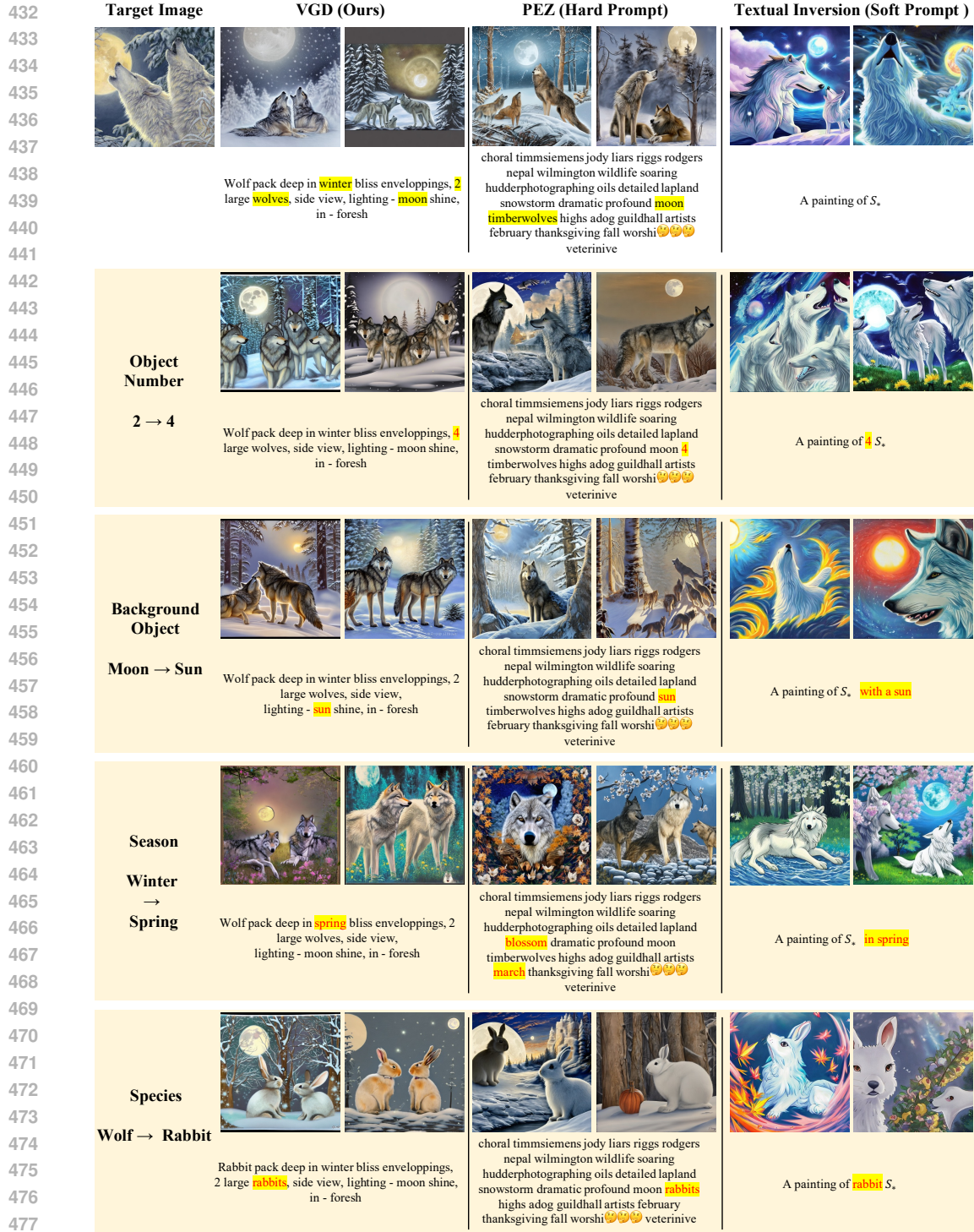


Figure 9: Prompt interpretability. With the good interpretability of the prompt, we can modify them to edit the original image to the desired direction.

CLIP Model Ablation VGD leverages the same CLIP model that is used in the text-to-image model, Stable Diffusion 2.1. To assess the importance of this alignment, we perform an ablation study using two alternative CLIP models: `openai/ViT-B-32` and `openai/ViT-L-14`. Table 3 shows that using a misaligned CLIP model leads to significant performance degradation. This highlights the importance of the proper approximation $P(T|I) \approx P_{\text{CLIP}}(T|I)$ in the decoding pro-

REFERENCES

- 540
541
542 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel
543 Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language
544 model for few-shot learning. *Advances in neural information processing systems*, 35:23716–
545 23736, 2022.
- 546 Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. Guided open vocabulary
547 image captioning with constrained beam search. In *Proceedings of the 2017 Conference on Em-
548 pirical Methods in Natural Language Processing*, pp. 936–945, 2017.
- 549 Tim Brooks, Aleksander Holynski, and Alexei A Efros. Instructpix2pix: Learning to follow image
550 editing instructions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
551 Recognition*, pp. 18392–18402, 2023.
- 552
553 Peter F Brown, Stephen A Della Pietra, Vincent J Della Pietra, and Robert L Mercer. The mathe-
554 matics of statistical machine translation: Parameter estimation. *Computational linguistics*, 19(2):
555 263–311, 1993.
- 556 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha
557 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models.
558 *arXiv preprint arXiv:2407.21783*, 2024.
- 559 Martin Nicolas Everaert, Marco Bocchio, Sami Arpa, Sabine Süssstrunk, and Radhakrishna Achanta.
560 Vetim: Expanding the vocabulary of text-to-image models only with text. In *The 34th British
561 Machine Vision Conference (BMVC 2023)*. BMVA, 2023.
- 562
563 Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit Haim Bermano, Gal Chechik, and
564 Daniel Cohen-or. An image is worth one word: Personalizing text-to-image generation using
565 textual inversion. In *The Eleventh International Conference on Learning Representations*, 2023.
- 566 Songwei Ge, Taesung Park, Jun-Yan Zhu, and Jia-Bin Huang. Expressive text-to-image generation
567 with rich text. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
568 pp. 7545–7556, 2023.
- 569 Qin Guo and Tianwei Lin. Focus on your instruction: Fine-grained and multi-instruction image
570 editing by attention modulation. In *Proceedings of the IEEE/CVF Conference on Computer Vision
571 and Pattern Recognition*, pp. 6986–6996, 2024.
- 572
573 Yaru Hao, Zewen Chi, Li Dong, and Furu Wei. Optimizing prompts for text-to-image generation.
574 *Advances in Neural Information Processing Systems*, 36, 2024.
- 575 Amir Hertz, Ron Mokady, Jay Tenenbaum, Kfir Aberman, Yael Pritch, and Daniel Cohen-or.
576 Prompt-to-prompt image editing with cross-attention control. In *The Eleventh International Con-
577 ference on Learning Representations*, 2023.
- 578
579 Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes, and Yejin Choi. The curious case of neural text
580 degeneration. In *International Conference on Learning Representations*, 2020.
- 581
582 J Edward Hu, Huda Khayrallah, Ryan Culkin, Patrick Xia, Tongfei Chen, Matt Post, and Benjamin
583 Van Durme. Improved lexically constrained decoding for translation and monolingual rewriting.
584 In *Proceedings of the 2019 Conference of the North American Chapter of the Association for
585 Computational Linguistics: Human Language Technologies*, pp. 839–850, 2019.
- 586
587 Albert Q Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot,
588 Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, et al.
589 Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- 590 Dong Jing, Xiaolong He, Yutian Luo, Nanyi Fei, Guoxing Yang, Wei Wei, Huiwen Zhao, and Zhiwu
591 Lu. Fineclip: Self-distilled region-based clip for better fine-grained understanding. In *The Thirty-
592 eighth Annual Conference on Neural Information Processing Systems*.
- 593 Daniel Jurafsky. *Speech and language processing*, 2000.

- 594 Xianghao Kong, Ollie Liu, Han Li, Dani Yogatama, and Greg Ver Steeg. Interpretable diffusion via
595 information decomposition. In *The Twelfth International Conference on Learning Representations*,
596 2024.
- 597 Nupur Kumari, Bingliang Zhang, Richard Zhang, Eli Shechtman, and Jun-Yan Zhu. Multi-concept
598 customization of text-to-image diffusion. In *Proceedings of the IEEE/CVF Conference on Com-*
599 *puter Vision and Pattern Recognition*, pp. 1931–1941, 2023.
- 600 Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt
601 tuning. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih (eds.),
602 *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp.
603 3045–3059, November 2021.
- 604 Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-
605 training for unified vision-language understanding and generation. In *International conference on*
606 *machine learning*, pp. 12888–12900. PMLR, 2022.
- 607 Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image
608 pre-training with frozen image encoders and large language models. In *International conference*
609 *on machine learning*, pp. 19730–19742. PMLR, 2023.
- 610 Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
611 *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the*
612 *11th International Joint Conference on Natural Language Processing*, 2021.
- 613 Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr
614 Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *Computer*
615 *Vision–ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014,*
616 *Proceedings, Part V 13*, pp. 740–755. Springer, 2014.
- 617 Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *Advances*
618 *in neural information processing systems*, 36, 2024.
- 619 Xiao Liu, Kaixuan Ji, Yicheng Fu, Weng Tam, Zhengxiao Du, Zhilin Yang, and Jie Tang. P-tuning:
620 Prompt tuning can be comparable to fine-tuning across scales and tasks. In *Proceedings of the*
621 *60th Annual Meeting of the Association for Computational Linguistics*, 2022.
- 622 Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang. Deep learning face attributes in the wild.
623 In *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- 624 Shweta Mahajan, Tanzila Rahman, Kwang Moo Yi, and Leonid Sigal. Prompting hard or hardly
625 prompting: Prompt inversion for text-to-image diffusion models. In *Proceedings of the IEEE/CVF*
626 *Conference on Computer Vision and Pattern Recognition*, pp. 6808–6817, 2024.
- 627 Hadas Orgad, Bahjat Kawar, and Yonatan Belinkov. Editing implicit assumptions in text-to-image
628 diffusion models. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*,
629 pp. 7053–7061, 2023.
- 630 Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong
631 Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow
632 instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:
633 27730–27744, 2022.
- 634 Roni Paiss, Ariel Ephrat, Omer Tov, Shiran Zada, Inbar Mosseri, Michal Irani, and Tali Dekel.
635 Teaching clip to count to ten. In *Proceedings of the IEEE/CVF International Conference on*
636 *Computer Vision*, pp. 3170–3180, 2023.
- 637 Dustin Podell, Zion English, Kyle Lacey, Andreas Blattmann, Tim Dockhorn, Jonas Müller, Joe
638 Penna, and Robin Rombach. SDXL: Improving latent diffusion models for high-resolution image
639 synthesis. In *The Twelfth International Conference on Learning Representations*, 2024.
- 640 Matt Post and David Vilar. Fast lexically constrained decoding with dynamic beam allocation for
641 neural machine translation. In *Proceedings of the 2018 Conference of the North American Chapter*
642 *of the Association for Computational Linguistics: Human Language Technologies*, 2018.

- 648 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
649 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
650 models from natural language supervision. In *International conference on machine learning*, pp.
651 8748–8763. PMLR, 2021.
- 652 Aditya Ramesh, Prafulla Dhariwal, Alex Nichol, Casey Chu, and Mark Chen. Hierarchical text-
653 conditional image generation with clip latents. *arXiv preprint arXiv:2204.06125*, 2022.
- 654 Kanchana Ranasinghe, Brandon McKinzie, Sachin Ravi, Yinfei Yang, Alexander Toshev, and
655 Jonathon Shlens. Perceptual grouping in contrastive vision-language models. In *Proceedings*
656 *of the IEEE/CVF International Conference on Computer Vision*, pp. 5571–5584, 2023.
- 657 Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-
658 resolution image synthesis with latent diffusion models. In *Proceedings of the IEEE/CVF confer-*
659 *ence on computer vision and pattern recognition*, pp. 10684–10695, 2022.
- 660 Nataniel Ruiz, Yuanzhen Li, Varun Jampani, Yael Pritch, Michael Rubinstein, and Kfir Aberman.
661 Dreambooth: Fine tuning text-to-image diffusion models for subject-driven generation. In *Pro-*
662 *ceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 22500–
663 22510, 2023.
- 664 Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis,
665 Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of
666 clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021.
- 667 Christoph Schuhmann, Romain Beaumont, Richard Vencu, Cade Gordon, Ross Wightman, Mehdi
668 Cherti, Theo Coombes, Aarush Katta, Clayton Mullis, Mitchell Wortsman, et al. Laion-5b: An
669 open large-scale dataset for training next generation image-text models. *Advances in Neural*
670 *Information Processing Systems*, 35:25278–25294, 2022.
- 671 Raphael Tang, Linqing Liu, Akshat Pandey, Zhiying Jiang, Gefei Yang, Karun Kumar, Pontus Stene-
672 torp, Jimmy Lin, and Ferhan Ture. What the DAAM: Interpreting stable diffusion using cross
673 attention. In *Proceedings of the 61st Annual Meeting of the Association for Computational Lin-*
674 *guistics*, 2023.
- 675 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Niko-
676 lay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open founda-
677 tion and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 678 Andrey Voynov, Qinghao Chu, Daniel Cohen-Or, and Kfir Aberman. p+: Extended textual condi-
679 tioning in text-to-image generation. *arXiv preprint arXiv:2303.09522*, 2023.
- 680 Zhijie Wang, Yuheng Huang, Da Song, Lei Ma, and Tianyi Zhang. Promptcharm: Text-to-image
681 generation through multi-modal prompting and refinement. In *Proceedings of the CHI Conference*
682 *on Human Factors in Computing Systems*, pp. 1–21, 2024.
- 683 Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein.
684 Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery.
685 *Advances in Neural Information Processing Systems*, 36, 2024.
- 686 Mert Yuksekgonul, Federico Bianchi, Pratyusha Kalluri, Dan Jurafsky, and James Zou. When and
687 why vision-language models behave like bags-of-words, and what to do about it? *arXiv preprint*
688 *arXiv:2210.01936*, 2022.
- 689 Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. Bertscore: Evalu-
690 ating text generation with bert. In *International Conference on Learning Representations*, 2020.
- 691 Yuechen Zhang, Shengju Qian, Bohao Peng, Shu Liu, and Jiaya Jia. Prompt highlighter: Interactive
692 control for multi-modal llms. In *Proceedings of the IEEE/CVF Conference on Computer Vision*
693 *and Pattern Recognition*, pp. 13215–13224, 2024.
- 694 Yiwu Zhong, Jianwei Yang, Pengchuan Zhang, Chunyuan Li, Noel Codella, Liunian Harold Li,
695 Luowei Zhou, Xiyang Dai, Lu Yuan, Yin Li, et al. Regionclip: Region-based language-image
696 pretraining. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recog-*
697 *nition*, pp. 16793–16803, 2022.

A DETAILS

A.1 PROMPT SETUP FOR LARGE LANGUAGE MODELS

We use different system and user prompts for different image generation tasks. Tables 5 and 6 show the prompts used for the experiments.

Table 5: LLM prompt setting for different image generation tasks.

System Prompt
You are a respectful and honest visual description generator for Stable Diffusion text prompt. Answer in 1 sentence and do not mention anything other than the prompt. Do not mention ‘description’.
User Prompt (Prompt Inversion)
Please generate the diffusion prompt on the given condition containing the objects, people, background, and the style of the image:
User Prompt (Style Transfer)
Please generate the diffusion prompt of the image style based on the given condition containing the painting style, color, and shapes of the image:
User Prompt (Prompt Distillation)
Please generate the diffusion prompt within $\{\text{max_length}\}$ tokens so that you can generate same images with a given prompt: $\{\text{target_prompt}\}$
Model prompt
Answer: Sure, here is a prompt for stable diffusion within $\{\text{model.max_length}\}$ tokens:

Table 6: LLaVA prompt setting for different image generation tasks.

Prompt (Image Captioning)
USER: < image > Describe the scene in this image with one sentence. ASSISTANT:
Prompt (VGD)
USER: < image > Please generate the diffusion prompt containing the objects, people, background and the style of the image. ASSISTANT: Sure, here is a prompt for stable diffusion within $\{\text{model.max_length}\}$ tokens:

A.2 HYPERPARAMETERS AND MODELS

The beam width K is set to 10. The balancing hyperparameter α is set to 0.67. For VGD generation, we used `laion/CLIP-ViT-H-14-laion2B-s32B-b79K`³. For CLIP-I score evaluation, we used `laion/CLIP-ViT-g-14-laion2B-s12B-b42K`⁴. We used Stable Diffusion `stabilityai/stable-diffusion-2-1`⁵ for text-to-image model.

A.3 DATASET CHARACTERISTICS

LAION-400M (Schuhmann et al., 2021; 2022) comprises 400 million diverse CLIP-filtered images. MS COCO (Lin et al., 2014) mainly contains real-life photographs with multiple common objects,

³<https://huggingface.co/laion/CLIP-ViT-H-14-laion2B-s32B-b79K>

⁴<https://huggingface.co/laion/CLIP-ViT-g-14-laion2B-s12B-b42K>

⁵<https://huggingface.co/stabilityai/stable-diffusion-2-1>

756 whereas Celeb-A (Liu et al., 2015) consists of celebrity portraits. Lexica.art contains AI-generated
757 paintings with their prompts. While different datasets exhibit different characteristics, VGD works
758 effectively regardless of the data type.

760 A.4 CLIP-INTERROGATOR’S APPROACH

761
762 CLIP-Interrogator supplements missing details in generated image captions by extensive trial-and-
763 error approach. Specifically, CLIP-Interrogator augments the caption with additional tokens selected
764 from a pre-collected bank of 100K keywords (*e.g.*, artist names, styles, mediums, movements), and
765 compare each with the target image through CLIP model to filter out image-irrelevant keywords.
766 While effective in generating the hard prompt, CLIP-Interrogator has several drawbacks: 1) an
767 inability to generate images for styles, artists, or objects that are not registered in the keyword
768 bank, and 2) the computational burden of processing 100K phrases through a CLIP encoder. These
769 limitations highlight the challenges of relying solely on image captions for image synthesis, and
770 underscore the need for efficient and flexible hard prompt generation techniques.

771 B ADDITIONAL RELATED WORKS

772 B.1 PROMPT EDITING FOR TEXT-TO-IMAGE MODELS

773
774 To fully harness the capabilities of text-to-image models, various approaches have been proposed
775 to tailor text prompts to the user’s specific intentions (Hertz et al., 2023; Brooks et al., 2023; Wang
776 et al., 2024) Notably, a rich text editor has been introduced, enabling users to design complex textual
777 guidance to image generation models (Ge et al., 2023). Similarly, Prompt Highlighter, an interactive
778 text-to-image interface utilizing multi-modal LLMs, has been proposed (Zhang et al., 2024). These
779 studies underscore the necessity for user-friendly and controllable text editing methods. Our VGD
780 can be seamlessly integrated into LLM-based user interfaces, as illustrated in Figure 1. VGD lever-
781 ages the same token generation process as LLMs, benefiting from the highly optimized inference
782 process without requiring additional training.

783 B.2 INTERPRETABILITY OF TEXT-TO-IMAGE MODELS

784
785 As text-to-image applications become standard in the industry, interpretability studies for these mod-
786 els have gained significant attention. For instance, research has focused on the cross-attention mech-
787 anism to identify the model’s attention for each word (Tang et al., 2023). In advance, the diffusion
788 objective has been investigated to understand the internal behavior of text-to-image models (Kong
789 et al., 2024). Building on these findings, recent works modified the model’s architecture to more
790 faithfully reflect user’s instructions (Orgad et al., 2023; Guo & Lin, 2024). We emphasize that VGD
791 aims to enhance the human interpretability without architectural modifications. We believe that
792 previous works can also be also applied to VGD to further improve usability.

793 C ADDITIONAL EXPERIMENTAL RESULTS

794 C.1 GENERALIZATION OF GENERATED PROMPTS TO DIFFERENT T2I MODELS

795
796 We demonstrate how prompts generated by our method and other baseline methods behave when
797 applied to different text-to-image models (see Fig. ??, 15, and 16). In this experiment, we did not
798 generate any new prompt for the newly tested T2I model (*i.e.*, MidJourney, DALL-E 2), meaning
799 that we simply reused the prompts generated for SD in evaluating the effectiveness of producing
800 high-quality images with various T2I models. In doing so, we can ensure a fair comparison and
801 also highlight the generalizability of the prompts across models. Finally, we would like to point out
802 that this type of generalization cannot be done with soft prompt methods like Textual Inversion, as
803 they are inherently tied to the specific model they are trained on. When compared to hard prompt
804 inversion methods (such as PH2P and PEZ), VGD generates consistently better results in various
805 text-to-image models, which underscores the robustness and flexibility of the proposed VGD.
806
807
808
809

810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863

C.2 PROMPT GENERATION TIME

We further investigate the efficiency of VGD in comparison with other baseline methods, measured on a single A100 80GB GPU. As shown in Table 7, the proposed VGD needs significantly less time to generate each image; 7x faster than Textual Inversion, 12x faster than PEZ, and more than 2200x faster than PH2P. This is because VGD does not require any training steps and only the decoding process is slightly modified.

Table 7: Prompt generation time of VGD, PEZ, PH2P, Textual Inversion, and CLIP Interrogator.

Method	PEZ	PH2P	Textual Inversion	CLIP Interrogator	VGD (ours)
Time (sec)	193.14	34264.74	103.45	24.85	15.33

C.3 MORE QUALITATIVE COMPARISON

As shown in Fig. 17 and 19, we conduct a qualitative comparison of image variation and style transfer performance across different methods, including PEZ, Textual Inversion, and CLIP Interrogator. The results demonstrate that, similar to the findings in Table 1, our method outperforms prior works, with CLIP Interrogator coming second.

While CLIP Interrogator captures objects effectively by generating captions using LLaVA (Liu et al., 2024) and predefined rich keyword bank, there are some critical weaknesses in prompts generated by the CLIP Interrogator (see Fig. 17). First, the subsequent keywords extracted from the caption’s keyword bank using CLIP similarity often suffer from redundancy and lack of interpretability. Second, CLIP Interrogator tends to generate prompts of full length (77 words), which limits its applicability to diverse tasks like multi-concept image generation and style transfer. Finally, as seen in the prompt example for the tiger image in Fig. 17, the quality of the generated images often decreases when using the extracted keywords, compared to images generated without attaching the keywords. In some cases, this leads to all extracted keywords being rejected entirely.

C.4 HUMAN EVALUATION

To evaluate whether the images generated by VGD are semantically aligned with the original image and to assess image quality, we conduct human preference evaluation. For this experiment, participants were shown images generated for image variation and style transfer by five different methods (VGD, Textual Inversion, PEZ, PH2P, and CLIP Interrogator) and were asked to select the image most similar to the target image and style (Kumari et al., 2023) (see Fig. 12). The evaluation was conducted using 65 images randomly sampled from Lexica.art dataset. Fig. 11 demonstrates that our method consistently performed best in both image variation and style transfer tasks according to user preferences. Specifically, 60.5% and 52.6% of the participants selected VGD’s results as the most similarly generated ones for image variation and most instruction-following generations for style transfer tasks, respectively.

C.5 BAD CASES

While VGD generally shows decent results, our method sometimes struggles to capture fine-grained details of regional or background objects in complex images. These cases occur when the image contains multiple objects (see Fig. 20). This is because VGD has difficulty in generating prompts for regional or local objects in the background. We believe this phenomenon is mainly due to the CLIP vision encoder’s tendency to prioritize main objects in an image (Jing et al.; Paiss et al., 2023; Ranasinghe et al., 2023; Yuksekgonul et al., 2022; Zhong et al., 2022).

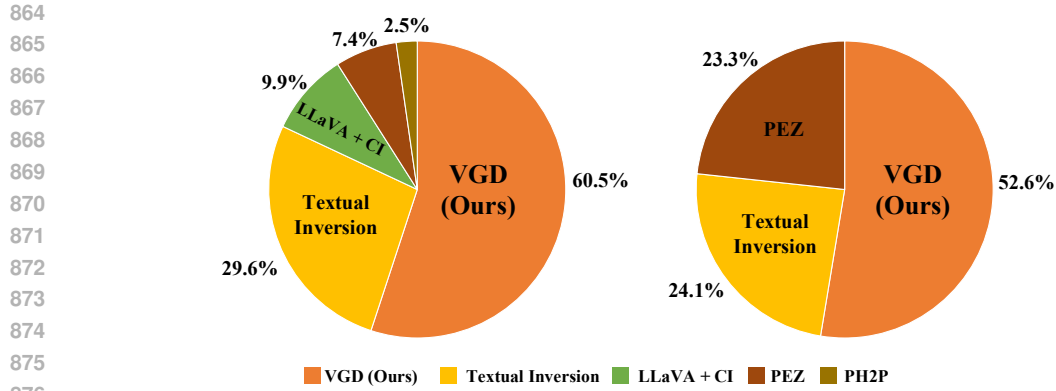


Figure 11: Human preference evaluation for image variation (left) and style transfer (right) tasks using different prompt generation methods. CI indicates CLIP Interrogator. Best viewed in color.

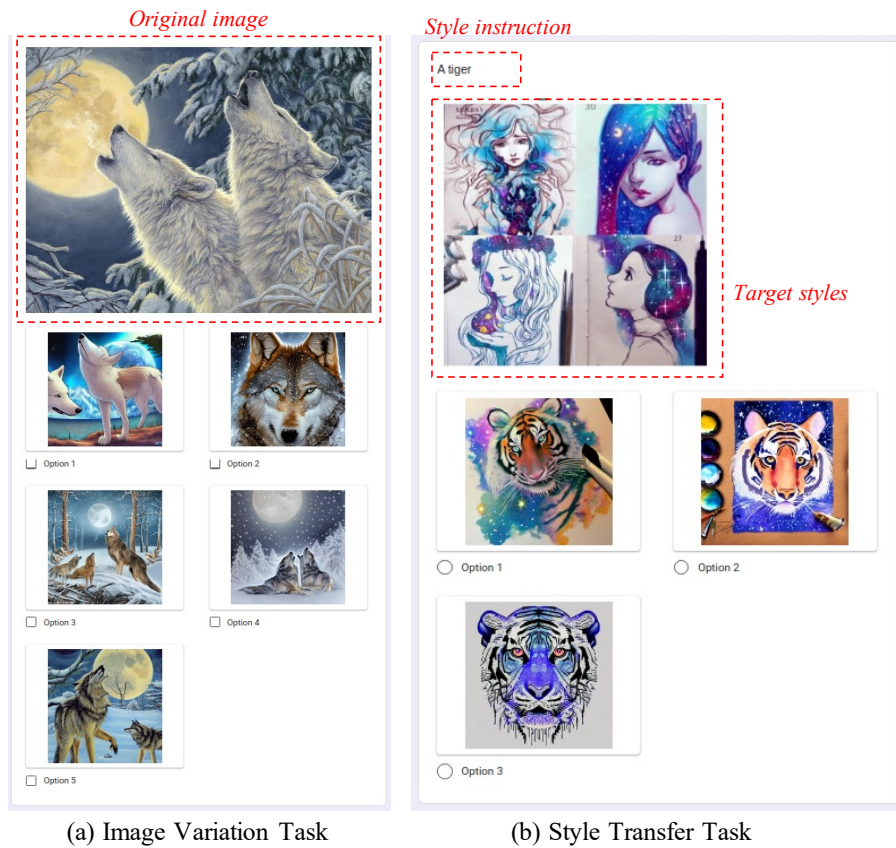


Figure 12: Webpage used for human preference evaluation.

C.6 MULTI-CONCEPT GENERATION

We further explore the potential of VGD by combining five different concepts into a single image. Figure 21 demonstrates that the prompts generated by VGD can be concatenated to generate a complex image without missing concepts. However, generated prompts from PEZ suffers omits important details, such as the parliamentary, wolves, and the style of *The Starry Night*. This shows that VGD generates hard prompts that are compact and meaningful.

Table 8: Prompt interpretability evaluation with Perplexity metric. Lower scores indicate better performance.

Method	#Tokens	LAION-400M	MS COCO	Celeb-A	Lexica.art
PEZ	32	813.27	766.81	971.35	659.39
LLaVA-1.5 + CI	~77	102.62	80.97	67.59	76.04
PH2P	16	-	-	-	1126.11
VGD	32	94.96	89.48	68.69	91.12
VGD+LLaVA	32	190.37	106.84	108.80	121.11
VGD+LLaVA	77	42.00	36.11	34.39	47.07

D LIMITATION AND ETHICAL STATEMENT

D.1 LIMITATIONS

Although we showcase the generalizability of VGD-generated prompts across multiple text-to-image models, our evaluations primarily focus on popular models such as Stable Diffusion and MidJourney. However, since many other models are trained on similar datasets and follow comparable methods and architectures, we anticipate that the results would be consistent across these models as well.

D.2 ETHICS STATEMENT

The development of VGD raises important ethical considerations, particularly regarding potential misuse. While our method improves the interpretability and generalizability of text-to-image generation, it could also be exploited to generate harmful or inappropriate content, such as deepfakes or offensive imagery. We strongly discourage such uses and emphasize that this method is intended for responsible applications in creative, educational, and professional contexts.

972
 973
 974
 975
 976
 977
 978
 979
 980
 981
 982
 983
 984
 985
 986
 987
 988
 989
 990
 991
 992
 993
 994
 995
 996
 997
 998
 999
 1000
 1001
 1002
 1003
 1004
 1005
 1006
 1007
 1008
 1009
 1010
 1011
 1012
 1013
 1014
 1015
 1016
 1017
 1018
 1019
 1020
 1021
 1022
 1023
 1024
 1025



Figure 13: Qualitative examples generated by Stable Diffusion, Midjourney and Dall-E 2 using prompts generated from VGD, PEZ, and LLaVA+CLIP-Interrogator.

1026
 1027
 1028
 1029
 1030
 1031
 1032
 1033
 1034
 1035
 1036
 1037
 1038
 1039
 1040
 1041
 1042
 1043
 1044
 1045
 1046
 1047
 1048
 1049
 1050
 1051
 1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064
 1065
 1066
 1067
 1068
 1069
 1070
 1071
 1072
 1073
 1074
 1075
 1076
 1077
 1078
 1079

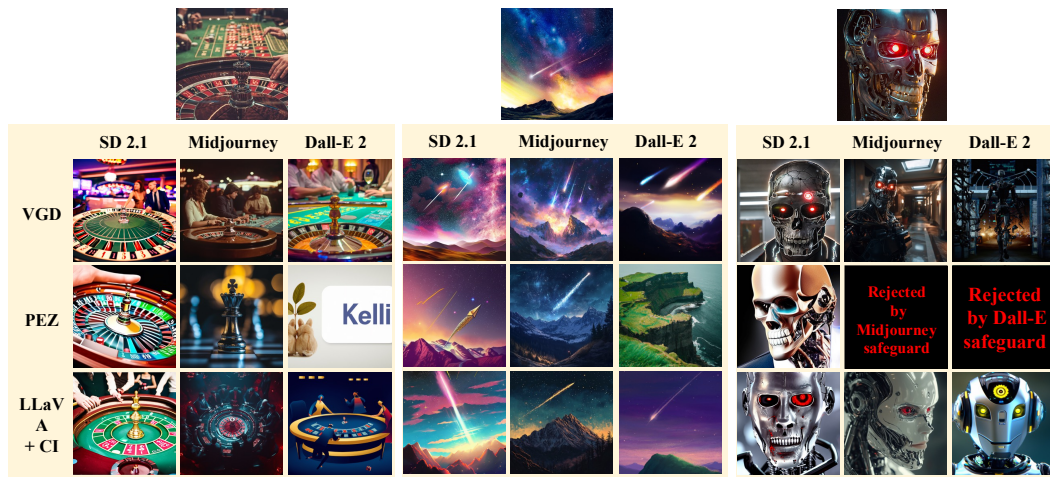


Figure 14: More qualitative examples generated by Stable Diffusion, Midjourney and Dall-E 2 using prompts generated from VGD, PEZ, and LLaVA+CLIP-Interrogator.



Figure 15: More examples generated by MidJourney.

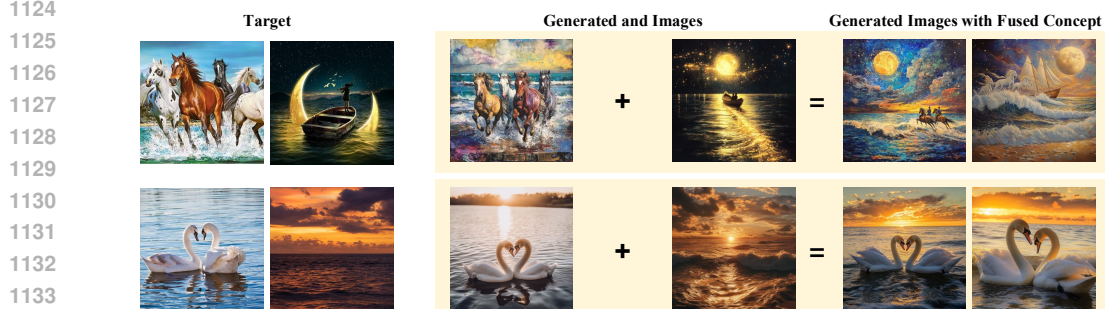


Figure 16: Multi-concept image generation by MidJourney.

1134
1135
1136
1137
1138
1139
1140
1141
1142
1143
1144
1145
1146
1147
1148
1149
1150
1151
1152
1153
1154
1155
1156
1157
1158
1159
1160
1161
1162
1163
1164
1165
1166
1167
1168
1169
1170
1171
1172
1173
1174
1175
1176
1177
1178
1179
1180
1181
1182
1183
1184
1185
1186
1187



Figure 17: Image variation comparison between VGD and baselines (generated by Stable Diffusion).

1188
1189
1190
1191
1192
1193
1194
1195
1196
1197
1198
1199
1200
1201
1202
1203
1204
1205
1206
1207
1208
1209
1210
1211
1212
1213
1214
1215
1216
1217
1218
1219
1220
1221
1222
1223
1224
1225
1226
1227
1228
1229
1230
1231
1232
1233
1234
1235
1236
1237
1238
1239
1240
1241



Figure 18: More image variation comparison between VGD and baselines (generated by Stable Diffusion).



Figure 19: Style transfer comparison between VGD and baselines.

1242
1243
1244
1245
1246
1247
1248
1249
1250
1251
1252
1253
1254
1255
1256
1257
1258
1259
1260
1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

Target Image	Generated Prompt and Images		
			
	<p>Cardinals winter wonderland scenes with birds nestled snug against the firs background as they prepare for their big day in their brightly -> missing birdhouse</p>		
			
	<p>Bears explore misty mountains en route a rusty railroad snaking through emerald pine bush; their footsteps echoing through time. -> missing train</p>		
			
	<p>Nature photographer captured a moment: squirrel twitch squirrel twitch background image is surrounded by various objects 100% res -> missing bird, mushroom</p>		

Figure 20: Bad case examples of VGD.

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Target Images to Combine

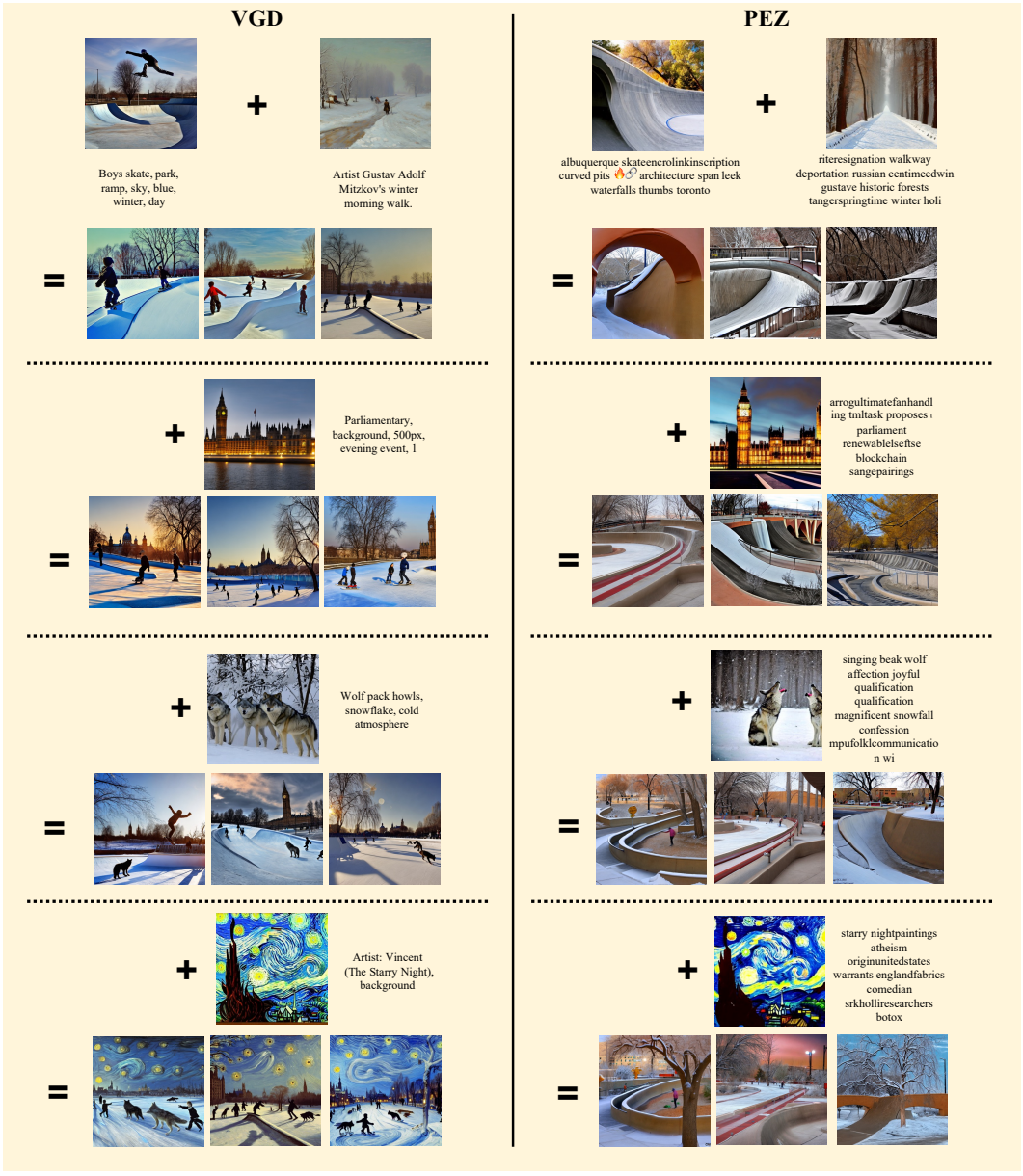


Figure 21: Five-concept image generation comparison between VGD and PEZ.