# Stylus: Automatic Adapter Selection for Diffusion Models

**Michael Luo**
UC Berkeley
michael.luo@berkeley.edu

**Justin Wong**
UC Berkeley
wong.justin@berkeley.edu

**Brandon Trabucco**
CMU MLD
brandon@btrabucco.com

**Yanping Huang**
Google Deepmind
huangyp@google.com

**Joseph E. Gonzalez**
UC Berkeley
jegonzal@berkeley.edu

**Zhifeng Chen**
Google Deepmind
zhifengc@google.com

**Ruslan Salakhutdinov**
CMU MLD
rsalakhu@cs.cmu.edu

**Ion Stoica**
UC Berkeley
istoica@berkeley.edu

Wooden dish rack on a counter holding plates, saucers, a bowl, mugs and glasses.

A boy holding an umbrella on the edge of a cliff.

A open field with large elephants standing in it.

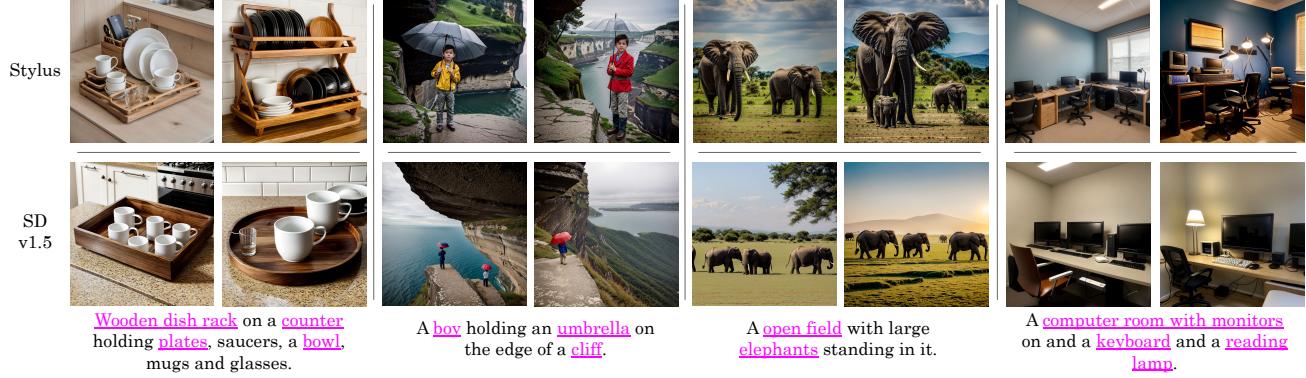A computer room with monitors on and a keyboard and a reading lamp.

Figure 1. **Adapter Selection.** Given a user-provided prompt, our method identifies highly relevant adapters (e.g. Low-Rank Adaptation, LoRA) that are closely aligned with the prompt's context and at least one of the prompt's keywords. Composing relevant adapters into Stable Diffusion improves visual fidelity, image diversity, and textual alignment. Note that these prompts are sampled from MS-COCO [19].

## Abstract

*Beyond scaling base models with more data or parameters, fine-tuned adapters provide an alternative way to generate high fidelity, custom images at reduced costs. As such, adapters have been widely adopted by open-source communities, accumulating a database of over 100K adapters—most of which are highly customized with insufficient descriptions. To generate high quality images, This paper explores the problem of matching the prompt to a set of relevant adapters, built on recent work that highlight the performance gains of composing adapters. We introduce Stylus, which efficiently selects and automatically composes task-specific adapters based on a prompt's keywords. Stylus outlines a three-stage approach that first summarizes adapters with improved descriptions and embeddings, retrieves relevant adapters, and then further assembles adapters based on prompts' keywords by checking how well they fit the prompt. To evaluate Stylus, we developed* `StylusDocs`, *a curated dataset featuring 75K adapters with pre-computed adapter embeddings. In our evaluation on popular Stable Diffusion checkpoints, Stylus achieves* greater CLIP/FID Pareto efficiency and is twice as preferred, with humans and multimodal models as evaluators, over the base model. See *stylus-diffusion.github.io* for more.

## 1. Introduction

In the evolving field of generative image models, finetuned adapters [7, 9] have become the standard, enabling custom image creation with reduced storage requirements. This shift has spurred the growth of extensive open-source platforms that encourage communities to develop and share different adapters and model checkpoints, fueling the proliferation of creative AI art [24, 43]. As the ecosystem expands, the number of adapters has grown to over 100K, with Low-Rank Adaptation (LoRA) [12] emerging as the dominant finetuning approach (see Fig. 3). A new paradigm has emerged where users manually select and creatively compose multiple adapters, on top of existing checkpoints, to generate high-fidelity images, moving beyond the standard approach of improving model class or scale.

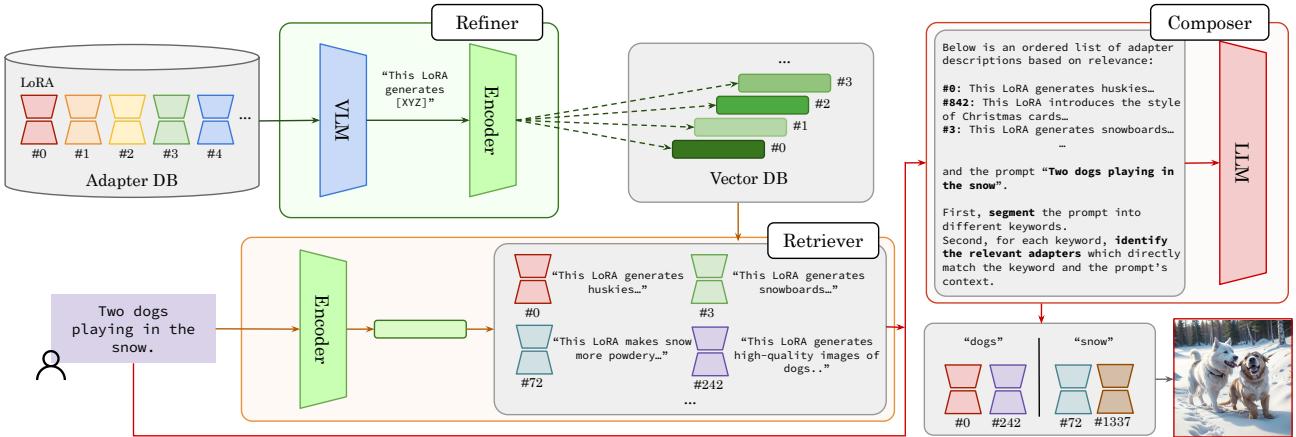In light of performance gains, our paper explores the automatic selection of adapters based on user-provided prompts

1

Figure 2. **Stylus algorithm.** Stylus consists of three stages. The *refiner* plugs an adapter's model card through a VLM to generate textual descriptions of an adapter's task and then through an encoder to produce the corresponding text embedding. The *retriever* fetches candidate adapters that are relevant to the entire user prompt. Finally, the *composer* prunes and jointly categorizes the remaining adapters based on the prompt's tasks, which correspond to a set of keywords.

(see Fig. 1). However, selecting relevant adapters presents unique challenges compared to existing retrieval-based systems, which rank relevant texts via lookup embeddings [16]. Specifically, efficiently retrieving adapters requires converting adapters into lookup embeddings, a step made difficult with low-quality documentation or no direct access to training data—a common issue on open-source platforms. Furthermore, in the context of image generation, user prompts often imply multiple highly-specific tasks. For instance, the prompt "two dogs playing the snow" suggests that there are two tasks: generating images of "dogs" and "snow". This necessitates segmenting the prompt into various tasks (i.e. keywords) and selecting relevant adapters for each task, a requirement beyond the scope of existing retrieval-based systems [8]. Finally, composing multiple adapters can degrade image quality, override existing concepts, and introduce unwanted biases into the model (see App. A.4).

We propose Stylus, a system that efficiently assesses user prompts to retrieve and compose sets of highly-relevant adapters, automatically augmenting generative models to produce diverse sets of high quality images. Stylus employs a three-stage framework to address the above challenges. As shown in Fig. 2, the *refiner* plugs in an adapter's model card, including generated images and prompts, through a multi-modal vision-language model (VLM) and a text encoder to pre-compute concise adapter descriptions as lookup embeddings. Similar to prior retrieval methods [16], the *retriever* scores the relevance of each embedding against the user's entire prompt to retrieve a set of candidate adapters. Finally, the *composer* segments the prompt into disjoint tasks, further prunes irrelevant candidate adapters, and assigns the remaining adapters to each task. We show that the composer identifies highly-relevant adapters and avoids conceptually-similar adapters that
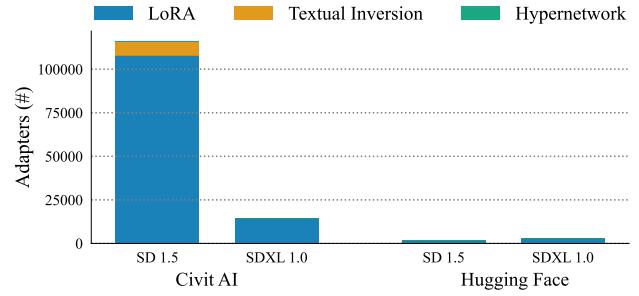


Figure 3. **Number of Adapters.** Civit AI boasts 100K+ adapters for Stable Diffusion, outpacing that of Hugging Face. Low-Rank Adaptation (LoRA) is the dominant approach for finetuning.

introduce biases detrimental to image generation (§ 4.3). Finally, Stylus applies a binary mask to control the number of adapters per task, ensuring high image diversity by using different adapters for each image and mitigating challenges with composing many adapters.

To evaluate our system, we introduce StylusDocs, an adapter dataset consisting of 75K LoRAs[1], that contains pre-computed adapter documentations and embeddings from Stylus's *refiner*. Our results demonstrate that Stylus improves visual fidelity, textual alignment, and image diversity over popular Stable Diffusion (SD 1.5) checkpoints—shifting the CLIP-FID Pareto curve towards greater efficiency and achieving up to 2x higher preference scores with humans and vision-language models (VLMs) as evaluators. As a system, Stylus is practical and does not present large overheads to the image generation process. Finally, Stylus can extend to different image-to-image application domains, such as image inpainting and translation.

---

[1]Sourced from `https://civitai.com/` [24].

## 2. Related Works

**Adapters.** Adapters efficiently fine-tune models on specific tasks with minimal parameter changes, reducing computational and storage requirements while maintaining similar performance to full fine-tuning [7, 9, 12]. Our study focuses on retrieving and merging multiple Low-Rank adapters (LoRA), the popular approach within existing open-source communities [24, 25, 43].

Adapter composition has emerged as a crucial mechanism for enhancing the capabilities of foundational models across various applications [17, 30, 34, 38, 39]. For large language models (LLM), the linear combination of multiple adapters improves in-domain performance and cross-task generalization [3, 13, 14, 40, 41, 46]. In the image domain, merging LoRAs effectively composes different tasks—concepts, characters, poses, actions, and styles—together, yielding images of high fidelity that closely align with user specifications [21, 47]. Our approach advances this further by actively segmenting user prompts into distinct tasks and merging the appropriate adapters for each task.

**Retrieval-based Methods.** Retrieval-based methods, such as retrieval-augmented generation (RAG), significantly improve model responses by adding semantically similar texts from a vast external database [16]. These methods convert text to vector embeddings using text encoders, which are then ranked against a user prompt based on similarity metrics [4, 8, 18, 23, 31, 33]. Similarly, our work draws inspiration from RAG to encode adapters as vector embedings: leveraging visual-language foundational models (VLM) to generate semantic descriptions of adapters, which are then translated into embeddings.

A core limitation to RAG is limited precision, retrieving distracting irrelevant documents. This leads to a "needle-in-the-haystack" problem, where more relevant documents are buried further down the list [8]. Recent work introduce *reranking* step; this technique uses cross-encoders to assess both the raw user prompt and the ranked set of raw texts individually, thereby discovering texts based on actual relevance [23, 32]. Rerankers have been successfully integrated with various LLM-application frameworks [2, 20, 29].

## 3. Our Method: Stylus

Adapter selection presents distinct challenges compared to existing methods for retrieving text-based documents, as outlined in Section 2. First, computing embeddings for adapters is a novel task, made more difficult without access to training datasets. Furthermore, in the context of image generation, user prompts often specify multiple highly fine-grained tasks. This challenge extends beyond retrieving relevant adapters relative to the entire user prompt, but also matching them with specific tasks within the

prompt. Finally, composing multiple adapters can degrade image quality and inject foreign biases into the model. Our three-stage framework below—**R**efine, **R**etrieve, and **C**ompose—addresses the above challenges (Fig. 2).

### 3.1. Refiner

The *refiner* is a two-stage pipeline designed to generate textual descriptions of an adapter's task and the corresponding text embeddings for retrieval purposes. This approach mirrors retrieval-based methods [16], which pre-compute embeddings over an external database of texts.

Given an adapter $A_i$, the first stage is a vision-language model (VLM) that takes in the adapter's model card—a set of randomly sampled example images from the model card $\mathcal{I}_i \in \{I_{i1}, I_{i2}, ...\}$, the corresponding prompts $\mathcal{P}_i \in \{p_{i1}, p_{i2}, ...\}$, and an author-provided description,[2] $D_i$—and returns an improved description $D_i^*$. Optionally, the VLM also recommends the weight for LoRA-based adapters, as the adapter weight is usually specified either in the author's description $D_i$ or the set of prompts $P_i$, a feature present in popular image generation software [1]. If information cannot be found, the LoRA's weight is set to 0.8. In our experiments, these improved descriptions were generated by Gemini Ultra [37] (see § A.1 for prompt).

The second stage uses an embedding model ($\mathcal{E}$) to generate embeddings $e = \mathcal{E}(D^*)$ for all adapters. In our experiments, we create embeddings from OpenAI's `text-embedding-3-large` model [18, 26]. We store pre-computed embeddings in a vector database.

### 3.2. Retriever

The *retriever* fetches the most relevant adapters over the entirety of the user's prompt using similarity metrics. Mathematically, the retriever employs the same embedding model ($\mathcal{E}$) to process the user prompt, $s$, generating embedding $e_s = \mathcal{E}(s)$. It then calculates cosine similarity scores between the prompt's embedding $e_s$ and the embedding of each adapter in the matrix $\mathcal{M}$. The top $K$ adapters $\mathcal{A}_K$ ($K = 150$, in our experiments) are selected based on the highest similarity scores: $\mathcal{A}_K = \arg\text{top-K}_i \left( \frac{e_s \cdot \mathcal{M}_i}{\|e_s\| \|\mathcal{M}_i\|} \right)$, where $\mathcal{M}_i$ is the $i^{th}$ row of the embedding matrix, representing the $i^{th}$ adapter's embedding.

### 3.3. Composer

The *composer* serves a dual purpose: segmenting the prompt into tasks from a prompt's keywords and assigning retrieved adapters to tasks. This implicitly filters out adapters that are not semantically aligned with the prompt and detects those likely to introduce foreign bias to the

---

[2]We note that a large set of author descriptions are inaccurate, misleading, or absent. The *refiner* helped correct for human errors by using generated images as the ground truth, significantly improving our system.

Figure 4. **Qualitative comparison between Stylus over realistic (left) and cartoon (right) style Stable Diffusion checkpoints.** Stylus produces highly detailed images that correctly depicts keywords in the context of the prompt. For the prompt "A graffiti of a corgi on the wall", our method correctly depicts a spray-painted corgi, whereas the checkpoint generates a realistic dog.
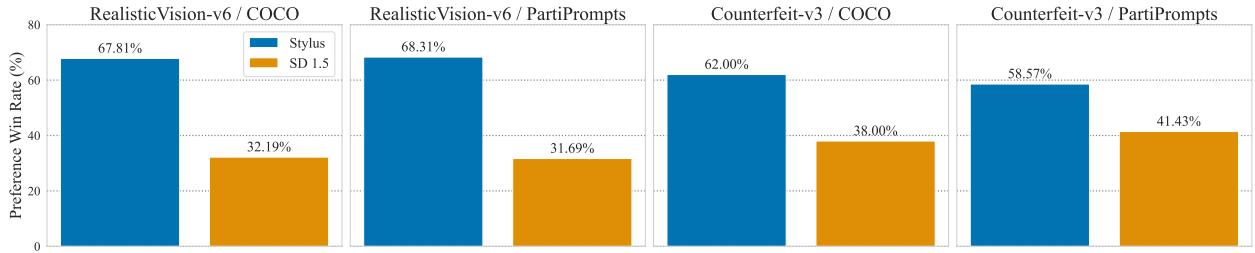


Figure 5. **Human Evaluation.** Stylus achieves a higher preference scores (2:1) over different datasets and Stable Diffusion checkpoints.

prompt through keyword grounding. For example, if the prompt is "pandas eating bamboo", the composer may discard an irrelevant "grizzly bears" adapter and a biased "panda mascots" adapter. Mathematically, the composer ($\mathcal{C}$) takes in the prompt ($s$) and the top $K$ adapters ($\mathcal{A}_K$) from the retriever, classifying them over different tasks, $\mathcal{T}(s) = \{t_1, t_2, \ldots, t_n\}$. This can be expressed as:

$$\mathcal{C}(s, \mathcal{A}_K) = \{(t_i, \mathcal{A}_{k_i}) \,|\, t_i \in \mathcal{T}(s), \mathcal{A}_{k_i} \subseteq \mathcal{A}_K,$$
$$\forall A_j \in \mathcal{A}_{k_i}, Sim(A_j, t_i) \geq \tau\} \qquad (1)$$

, where $\mathcal{A}_{k_i}$ is the subset of adapters per task $t_i$, $Sim(A_j, t_i)$ measures the similarity score between an adapter and a task, and $\tau$ is an arbitary threshold.

While the composer can be trained with human-labeled data [28], we opt for a simpler approach that requires no training—prompting a long-context Large Language Model (LLM). The LLM accepts the adapter descriptions and the prompt as part of its context and returns a mapping of tasks to a curated set of adapters. In our implementation, we choose Gemini 1.5, with a 128K context window, as the composer's LLM (see App. A.2 for the full prompt).

Most importantly, Stylus's composer parallels *reranking*, an advanced RAG technique. Rerankers employ cross encoders ($\mathcal{F}$) that compare the retriever's individual adapter descriptions, generated from the refiner, against the user prompt to determine better similarity scores: $\mathcal{F}(s, D^*)$. This prunes for adapters based on semantic relevance, thereby improving search quality, but not over keyword alignment. Our experimental ablations (§ 4.3) show that

our composer outperforms existing rerankers (Cohere, `rerank-english-v2.0`) [32].

### 3.4. Masking

The composer identifies tasks $t_i$ from the prompt $s$ and assigns each task a set of relevant adapters $\mathcal{A}_{k_i}$, formalized as: $C(\mathcal{A}_k, s) = \{(t_1, \mathcal{A}_{k_1}), (t_2, \mathcal{A}_{k_2}), \ldots\}$. Our masking strategy applies a binary mask, $M_i$, for each task $t_i$. Each mask, $M_i$, can either be an one hot encoding, all ones, or all zeroes vector. Across all tasks, we perform a cross-product across masks, $M_1 \times M_2 \times M_3 \times \ldots$, generating an exponential number of masking schemes.
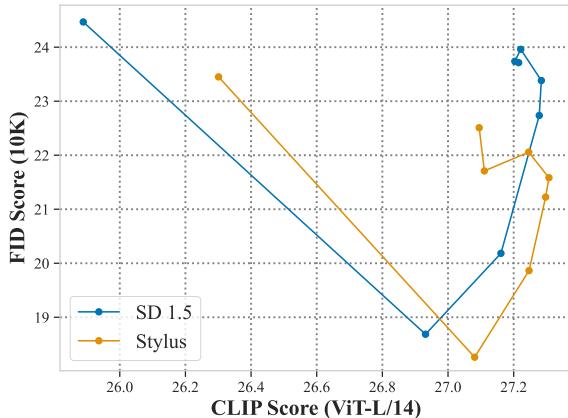
The combinatorial sets of masking schemes enable diverse linear combinations of adapters for a single prompt, leading to highly-diverse images (§ 4.2.3). This approach also curtails the number of final adapters merged into the base model, minimizing the risk of composed adapters introducing undesirable effects to the image [47].

Finally, an adapter's weight (i.e. LoRA), which is extracted from the refiner (§ 3.1), is divided by the total number of adapters (after masking) in its corresponding task. This solves the problem of image saturation, where assigning high adapter weights to an individual task (or concept) leads to sharp decreases in image quality (see App. A.4).

## 4. Results

### 4.1. Experimental Setup

**Adapter Testbed.** Adapter selection requires a large database of adapters to properly evaluate its performance.

4

Figure 6. **Clip/FID Pareto Curve for COCO.** We observe Stylus can improve visual fidelity (FID) and/or textual alignment (CLIP) over a range of guidance values (CFG): [1, 1.5, 2, 3, 4, 6, 9, 12].

| | CLIP ($\triangle$) | FID ($\triangle$) |
|---|---|---|
| Stylus | **27.25** (+0.03) | **22.05** (-1.91) |
| Reranker | 25.48 (-1.74) | 22.81 (-1.15) |
| Retriever-only | 24.93 (-2.29) | 24.68 (+0.72) |
| Random | 26.34 (-0.88) | 24.39 (+0.43) |
| SD v1.5 | 27.22 | 23.96 |

Table 1. **Evaluation over different retrieval methods (CFG=6).** Stylus outperforms existing retrieval-based methods, attains the best FID score, and similar CLIP score to Stable Diffusion.

However, existing methods [13, 46] only evaluate against 50-350 adapters for language-based tasks, which is insufficient for our use case, since image generation relies on highly fine grained tasks that span across many concepts, poses, styles, and characters. To bridge this gap, we introduce StylusDocs, a comprehensive dataset that pulls 75K LoRAs from popular model repositories, Civit AI and HuggingFace [24, 43]. This dataset contains precomputed OpenAI embeddings [18] and improved adapter descriptions from Gemini Ultra-Vision [37], the output of Stylus's refiner component (§ 3.1). We further characterize the distribution of adapters in App. A.3.

**Generation Details.** We assess Stylus against Stable-Diffusion-v1.5 [34] as the baseline model. Across experiments, we employ two well-known checkpoints: Realistic-Vision-v6, which excels in producing realistic images, and Counterfeit-v3, which generates cartoon and anime-style images. Our image generation process integrates directly with Stable-Diffusion WebUI [1] and defaults to 35 denoising steps using the default DPM Solver++ scheduler [22]. To replicate high-quality images from existing users, we enable high-resolution upscaling to generate 1024x1024 from 512x512 images, with the default latent upscaler [15] and denoising strength set to 0.7. For images generated by Stylus, we discovered adapters could

shift the image style away from the checkpoint's original style. To counteract this, we introduce a *debias prompt* injected at the end of a user prompt to steer images back to the checkpoint's style[3].

### 4.2. Main Experiments

#### 4.2.1 Human Evaluation.

To demonstrate our method's general applicability, we evaluate Stylus over a cross product of two datasets, Microsoft COCO [19] and PartiPrompts [45], and two checkpoints, which generate realistic and anime-style images respectively. Examples of images generated in these styles are displayed in Figure 4; Stylus generates highly detailed images that better focus on specific elements in the prompt.

To conduct human evaluation, we enlisted four users to assess 150 images from both Stylus and Stable Diffusion v1.5 for each dataset-checkpoint combination. These raters were asked to indicate their preference for Stylus or Stable-Diffusion-v1.5. In Fig. 5, users generally showed a preference for Stylus over existing model checkpoints. Although preference rates were consistent across datasets, they varied significantly between different checkpoints. Adapters generally improve details to their corresponding tasks (e.g. generate detailed elephants); however, for anime-style checkpoints, detail is less important, lowering preference scores.

#### 4.2.2 Automatic Benchmarks.

We assess Stylus using two automatic benchmarks: CLIP [10], which measures the correlation between a generated images' caption and users' prompts, and FID [11], which evaluates the diversity and aesthetic quality of image sets. We evaluate COCO 2014 validation dataset, with 10K sampled prompts, and the Realistic-Vision-v6 checkpoint. Fig. 6 shows that Stylus shifts the Pareto curve towards greater efficiency, achieving better visual fidelity and textual alignment. This improvement aligns with our human evaluations, which suggest a correlation between human preferences and the FID scores.
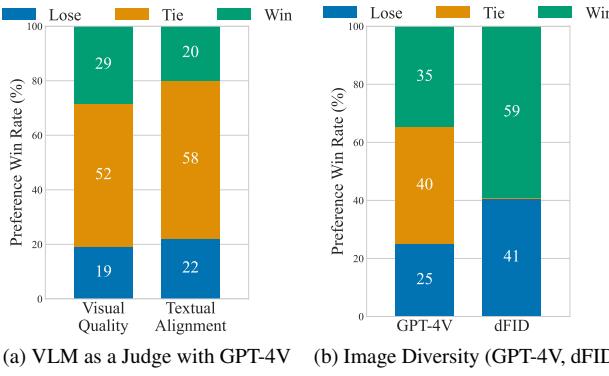
#### 4.2.3 VLM as a Judge

We use *VLM as a Judge* to assess two key metrics: textual alignment and visual fidelity, simulating subjective assessments [5]. For visual fidelity, the VLM scores based on disfigured limbs and unrealistic composition of objects. When asked to make subjective judgements, autoregressive models tend to exhibit bias towards the first option presented. To combat this, we evaluate Stylus under both orderings and only consider judgements that are consistent across reorderings; otherwise, we label it a tie. In Fig. 8a, we assess eval-

---

[3]The debias prompts are "realistic, high quality" for Realistic-Vision-v6 and "anime style, high quality" for Counterfeit-v3, respectively.

5

Figure 7. **Image Diversity.** Given the same prompt, our method (left) generates more diverse and comprehensive sets of images than that of existing Stable Diffusion checkpoints (right). Stylus's diversity comes from its masking scheme and the composer LLM's temperature.



(a) VLM as a Judge with GPT-4V    (b) Image Diversity (GPT-4V, dFID)

Figure 8. **Preference Win Rate over GPT-4V as a judge.** Stylus achieves higher preference scores over GPT-4V for visual quality and image diversity.
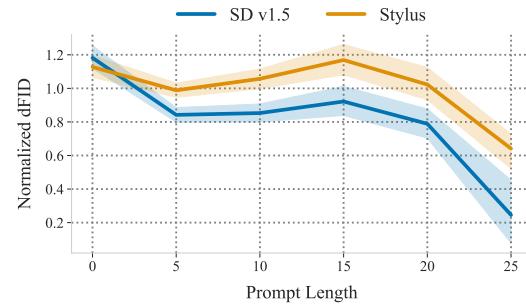


Figure 9. **Image diversity ($d$FID) across prompt length.** Stylus achieves higher diversity score than Stable Diffusion when prompt length increases.

uate 100 randomly sampled prompts from the PartiPrompts dataset [45]. Barring ties, we find visual fidelity achieves 60% win rate between Stylus and the Stable Diffusion realistic checkpoint, which is conclusively consistent with the 68% win rate from our human evaluation. For textual alignment, we find negligible differences between Stylus and the Stable Diffusion checkpoint. As most prompts lead to a tie, this indicates Stylus does not introduce additional artifacts. We provide the full prompt in Appendix A.5.

#### 4.2.4 Diversity per Prompt

Given identical prompts, Stylus generates highly diverse images due to different composer outputs and masking schemes. Qualitatively, Fig. 7 shows that Stylus generates dragons, maidens, and kitchens in diverse positions, concepts, and styles. To quantitatively assess this diversity, we use two metrics:

$d$**FID:** Previous evaluations with FID [11] show that Stylus improves image quality and diversity *across prompts*[4]. We define $d$FID specifically to evaluate diversity per prompt,

---

[4]FID fails to disentangle image fidelity from diversity [27, 35].

calculated as the variance of latent embeddings from InceptionV3 [36]. Mathematically, $d$FID involves fitting a Normal distribution $\mathcal{N}(\mu, \Sigma)$ to the latent features of InceptionV3, with the metric given by the trace of the covariance matrix, $d\text{FID} = \text{Tr}\,\Sigma$.

**GPT-4V:** We use *VLM as a Judge* to assess image diversity between images generated using Stylus and the Stable Diffusion checkpoint over PartiPrompts. Five images are sampled per group, Stylus and SD v1.5, with group positions randomly swapped across runs to avoid GPT-4V's positional bias [47]. Similar to VisDiff, we ask GPT-4V to rate on a scale from 0-2, where 0 indicates no diversity and 2 indicates high diversity [6]. Full prompt and additional details are provided in App A.5.

Fig. 8b displays preference rates and defines a win when Stylus achieves higher $d$FID or receives a higher score from GPT-4V for a given prompt. Across 200 prompts, Stylus prevails in approximately 60% and 58% cases for $d$FID and GPT-4V respectively, excluding ties. Figure 9 compares Stylus with base Stable Diffusion 1.5 across prompt lengths, revealing that Stylus consistently produces more diverse images. Additional results measuring diversity per keyword are presented in Appendix A.6.

Figure 10. **Different Retrieval Methods.** Stylus outperforms all other retrieval methods, which choose adapters than either introduce foreign concepts to the image or override other concepts in the prompt, reducing textual alignment.
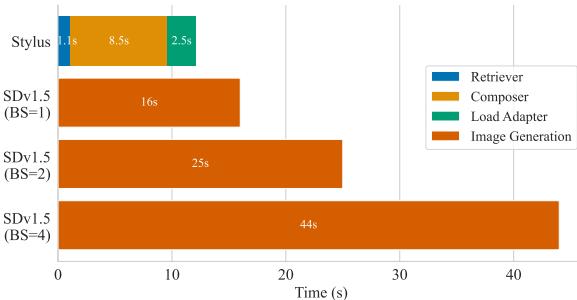


Figure 11. **Comparison of Stylus's inference overheads with Stable Diffusion's inference time by batch size (BS).** At BS=1, Stylus accounts for 75% of the image generation time, primarily due to the composer processing long context prompts from adapter descriptions. However, Stylus's overhead decreases when batch size increases.

## 4.3. Ablations

### 4.3.1 Alternative Retrieval-based Methods

We benchmark Stylus's performance relative to different retrieval methods. For all baselines below, we select the top three adapters and merge them into the base model.

**Random**: Adapters are randomly sampled without replacement from `StylusDocs`.

**Retriever**: The retriever emulates standard RAG pipelines [16, 46], functionally equivalent to Stylus without the composer stage. Top adapters are fetched via cosine similarity over adapter embeddings.

**Reranker**: An alternative to Stylus's composer, the reranker fetches the retriever's adapters and plugs a cross-encoder that outputs the semantic similarity between adapters' descriptions and the prompt. We evaluate with Cohere's reranker endpoint [32].

As shown in Tab. 1, Stylus achieves the highest CLIP and FID scores, outperforming all other baselines which fall behind the base Stable Diffusion model. First, both the retriever and reranker significantly underperform compared to Stable Diffusion. Each method selects adapters that are *similar* to the prompt but potentially introduce unrelated biases. In Fig. 10, both methods choose adapters related to elephant movie characters, which biases the concept of elephants and results in depictions of unrealistic elephants. Furthermore, both methods incorrectly assign weights to adapters, causing adapters' tasks to overshadow other tasks within the same prompt. In Fig. 10, both the reranker and retriever generate images solely focused on singular items—beds, chairs, suitcases, or trains—while ignoring other elements specified in the prompt. We provide an analysis of failure modes in A.4.

Conversely, the random policy exhibits performance comparable, but slightly worse, to Stable Diffusion. The random baseline chooses adapters that are orthogonal to the user prompt. Thus, these adapters alter unrelated concepts, which does not affect image generation. In fact, we observed that the distribution of random policy's images in Fig. 10 were nearly identical to Stable Diffusion.
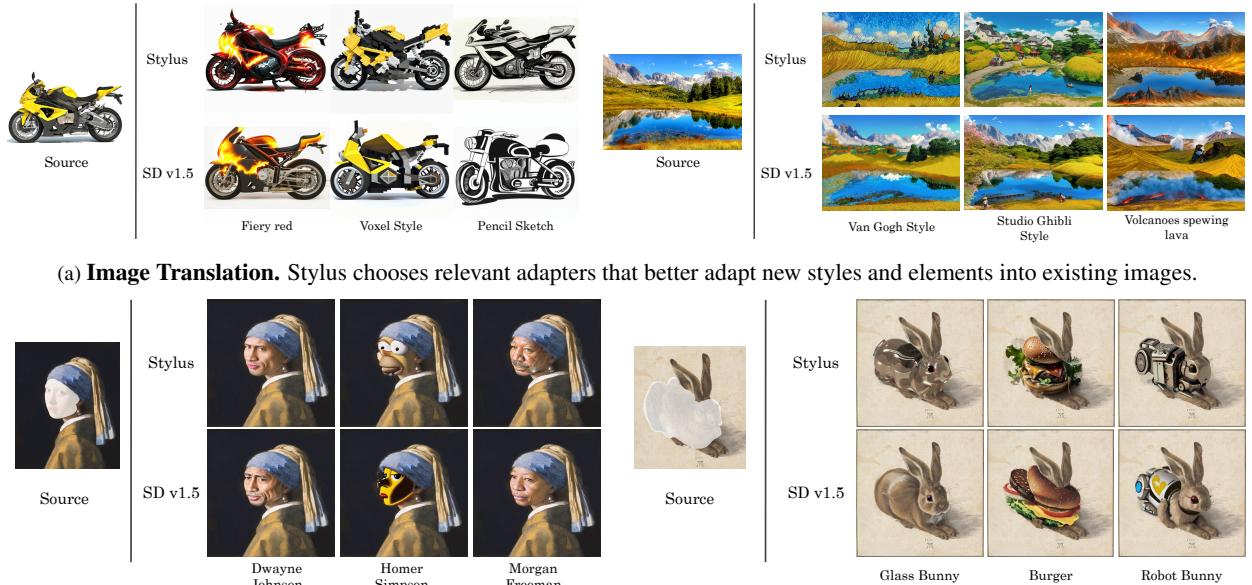
(a) **Image Translation.** Stylus chooses relevant adapters that better adapt new styles and elements into existing images.



(b) **Inpainting.** Stylus chooses adapters than can better introduce new characters or concepts into the inpainted mask.

Figure 12. **Stylus over different image-to-image tasks.**

### 4.3.2 Breakdown of Stylus's Inference Time

This section breaks down the latency introduced by various components of Stylus. We note that image generation time is independent of Stylus, as adapter weights are merged into the base model [12].

Figure 11 demonstrates the additional time Stylus contributes to the image generation process across different batch sizes (BS), averaged over 100 randomly selected prompts. Specifically, Stylus adds *12.1* seconds to the image generation time, with the composer accounting for 8.5 seconds. The composer's large overhead is due to long-context prompts, which include adapter descriptions for the top 150 adapters and can reach up to 20K+ tokens. Finally, when the BS is 1, Stylus presents a 75% increase in overhead to the image generation process. However, Stylus's latency remains consistent across all batch sizes, as the composer and retriever run only once. Hence, for batch inference workloads, Stylus incurs smaller overheads as batch size increases.

### 4.3.3 Image-Domain Tasks

Beyond text-to-image, Stylus applies across various image-to-image tasks. Fig. 12 demonstrates Stylus applied to two different image-to-image tasks: image translation and inpainting, described as follows:

**Image translation:** Image translation involves transforming a source image into a variant image where the content remains unchanged, but the style is adapted to match the prompt's definition. Stylus effectively converts images into their target domains by selecting the appropriate LoRA, which provides a higher fidelity description of the style. We present examples in Fig 12a. For a yellow motorcycle, Stylus identifies a voxel LoRA that more effectively decomposes the motorcycle into discrete 3D bits. For a natural landscape, Stylus successfully incorporates more volcanic elements, covering the landscape in magma.

**Inpainting:** Inpainting involves filling in missing data within a designated region of an image, typically outlined by a binary mask. Stylus excels in accurately filling the masked regions with specific characters and themes, enhancing visual fidelity. We provide further examples in Fig. 12b, demonstrating how Stylus can precisely inpaint various celebrities and characters (left), as well as effectively introduce new styles to a rabbit (right).

## 5. Conclusion

We propose Stylus, a new algorithm that automatically selects and composes adapters to generate better images. Our method leverages a three-stage framework that precomputes adapters as lookup embeddings and retrieves most relevant adapters based on prompts' keywords. To evaluate Stylus, we develop StylusDocs, a processed dataset featuring 75K adapters and pre-computed adapter embeddings. Our evaluation of Stylus, across automatic metrics, humans, and vision-language models, demonstrate that Stylus achieves better visual fidelity, textual alignment, and image diversity than existing Stable Diffusion checkpoints.

## References

[1] AUTOMATIC1111. Stable Diffusion Web UI, Aug. 2022. 3, 5

[2] Harrison Chase. LangChain, Oct. 2022. 3

[3] Alexandra Chronopoulou, Matthew E. Peters, Alexander Fraser, and Jesse Dodge. Adaptersoup: Weight averaging to improve generalization of pretrained language models, 2023. 3

[4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019. 3

[5] Yann Dubois, Chen Xuechen Li, Rohan Taori, Tianyi Zhang, Ishaan Gulrajani, Jimmy Ba, Carlos Guestrin, Percy S Liang, and Tatsunori B Hashimoto. Alpacafarm: A simulation framework for methods that learn from human feedback. In A. Oh, T. Neumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in Neural Information Processing Systems*, volume 36, pages 30039–30069. Curran Associates, Inc., 2023. 5

[6] Lisa Dunlap, Yuhui Zhang, Xiaohan Wang, Ruiqi Zhong, Trevor Darrell, Jacob Steinhardt, Joseph E. Gonzalez, and Serena Yeung-Levy. Describing differences in image sets with natural language. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2024. 6, 14

[7] Rinon Gal, Yuval Alaluf, Yuval Atzmon, Or Patashnik, Amit H. Bermano, Gal Chechik, and Daniel Cohen-Or. An image is worth one word: Personalizing text-to-image generation using textual inversion, 2022. 1, 3

[8] Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Meng Wang, and Haofen Wang. Retrieval-augmented generation for large language models: A survey, 2024. 2, 3

[9] David Ha, Andrew Dai, and Quoc V. Le. Hypernetworks, 2016. 1, 3

[10] Jack Hessel, Ari Holtzman, Maxwell Forbes, Ronan Le Bras, and Yejin Choi. Clipscore: A reference-free evaluation metric for image captioning, 2022. 5

[11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium, 2018. 5, 6

[12] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021. 1, 3, 8

[13] Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Efficient cross-task generalization via dynamic lora composition, 2024. 3, 5

[14] Joel Jang, Seungone Kim, Bill Yuchen Lin, Yizhong Wang, Jack Hessel, Luke Zettlemoyer, Hannaneh Hajishirzi, Yejin Choi, and Prithviraj Ammanabrolu. Personalized soups: Personalized large language model alignment via post-hoc parameter merging, 2023. 3

[15] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of diffusion-based generative models, 2022. 5

[16] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021. 2, 3, 7

[17] 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. 3

[18] Jimmy Lin, Ronak Pradeep, Tommaso Teofili, and Jasper Xian. Vector search with openai embeddings: Lucene is all you need, 2023. 3, 5

[19] Tsung-Yi Lin, Michael Maire, Serge Belongie, Lubomir Bourdev, Ross Girshick, James Hays, Pietro Perona, Deva Ramanan, C. Lawrence Zitnick, and Piotr Dollár. Microsoft coco: Common objects in context, 2015. 1, 5, 11

[20] Jerry Liu. LlamaIndex, 11 2022. 3

[21] Nan Liu, Yilun Du, Shuang Li, Joshua B. Tenenbaum, and Antonio Torralba. Unsupervised compositional concepts discovery with text-to-image generative models, 2023. 3

[22] Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. DPM-solver++: Fast solver for guided sampling of diffusion probabilistic models, 2023. 5

[23] Tengyu Ma. Vectorize your data to gear up your ai stack., 2023. 3

[24] Justin Maier. The home of open-source generative ai, 2022. 1, 2, 3, 5, 11

[25] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan. Peft: State-of-the-art parameter-efficient fine-tuning methods. https://github.com/huggingface/peft, 2022. 3

[26] Rui Meng, Ye Liu, Shafiq Rayhan Joty, Caiming Xiong, Yingbo Zhou, and Semih Yavuz. Sfr-embedding-mistral:enhance text retrieval with transfer learning. Salesforce AI Research Blog, 2024. 3

[27] Muhammad Ferjad Naeem, Seong Joon Oh, Youngjung Uh, Yunjey Choi, and Jaejun Yoo. Reliable fidelity and diversity metrics for generative models. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 7176–7185. PMLR, 2020. 6

[28] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback, 2022. 4

[29] Malte Pietsch, Timo Möller, Bogdan Kostic, Julian Risch, Massimiliano Pippi, Mayank Jobanputra, Sara Zanzottera, Silvano Cerza, Vladimir Blagojevic, Thomas Stadelmann, Tanay Soni, and Sebastian Lee. Haystack: the end-to-end NLP framework for pragmatic builders, Nov. 2019. 3

[30] 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. 3

[31] 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. 3

[32] Nils Reimers. Say goodbye to irrelevant search results: Cohere rerank is here, 2023. 3, 4, 7

[33] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. 3

[34] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models, 2022. 3, 5

[35] Mehdi S. M. Sajjadi, Olivier Bachem, Mario Lucic, Olivier Bousquet, and Sylvain Gelly. Assessing generative models via precision and recall. In Samy Bengio, Hanna M. Wallach, Hugo Larochelle, Kristen Grauman, Nicolò Cesa-Bianchi, and Roman Garnett, editors, *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, December 3-8, 2018, Montréal, Canada*, pages 5234–5243, 2018. 6

[36] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2818–2826. IEEE Computer Society, 2016. 6

[37] Gemini Team. Gemini: A family of highly capable multimodal models, 2023. 3, 5, 11

[38] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models, 2023. 3

[39] 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. 3

[40] Hanqing Wang, Bowen Ping, Shuo Wang, Xu Han, Yun Chen, Zhiyuan Liu, and Maosong Sun. Lora-flow: Dynamic lora fusion for large language models in generative tasks, 2024. 3

[41] Xinyi Wang, Yulia Tsvetkov, Sebastian Ruder, and Graham Neubig. Efficient test time adapter ensembling for low-resource language varieties. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 730–737, Punta Cana, Dominican Republic, Nov. 2021. Association for Computational Linguistics. 3

[42] 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. 11, 14

[43] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Huggingface's transformers: State-of-the-art natural language processing, 2020. 1, 3, 5, 11

[44] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Thomas L. Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models, 2023. 11, 14

[45] Jiahui Yu, Yuanzhong Xu, Jing Yu Koh, Thang Luong, Gunjan Baid, Zirui Wang, Vijay Vasudevan, Alexander Ku, Yinfei Yang, Burcu Karagol Ayan, Ben Hutchinson, Wei Han, Zarana Parekh, Xin Li, Han Zhang, Jason Baldridge, and Yonghui Wu. Scaling autoregressive models for content-rich text-to-image generation, 2022. 5, 6

[46] Ziyu Zhao, Leilei Gan, Guoyin Wang, Wangchunshu Zhou, Hongxia Yang, Kun Kuang, and Fei Wu. Loraretriever: Input-aware lora retrieval and composition for mixed tasks in the wild, 2024. 3, 5, 7

[47] Ming Zhong, Yelong Shen, Shuohang Wang, Yadong Lu, Yizhu Jiao, Siru Ouyang, Donghan Yu, Jiawei Han, and Weizhu Chen. Multi-lora composition for image generation, 2024. 3, 4, 6

# A. Appendix

## A.1. Details of the Refiner VLM

We provide a complete example input to the refiner's VLM in Tab. 2. The prompt utilizes Chain-of-Thought (CoT) prompting, which decomposes the VLM's goal of producing better adapter descriptions into two steps [42, 44]. Initially, the VLM categorizes the adapter's task into one of several topics—such as concepts, styles, characters, or poses. Subsequently, the VLM is prompted to elaborate on why the adapter is associated with a particular topic and how it modifies images within that context. We found that this two step logical process significantly improved the structure and quality of model responses.

## A.2. Details of the Composer LLM

We provide a full example prompt of the composer's LLM component in Tab. 3, which is plugged through the Gemini 1.5 endpoint [37]. Our experiments feed in descriptions of the top 150 adapters into the LLM's context. Using a Chain-of-Thought (CoT) approach, the prompt is structured to first identify keywords or tasks, then allocate appropriate adapters to these tasks. If necessary, it merges keywords for adapters that span multiple tasks [42, 44].

## A.3. Stylus-Bench Characterization

This section describes `StylusDocs`, which comprises of 76K Low Rank Adapters (LoRAs) from public repositories, including Civit AI and Hugging Face [24, 43]. We excluded NSFW-labeled adapters from the Civit AI dataset, which originally contained over 100K LoRAs. Figure 13 illustrates the distribution of adapters across various semantic categories and their popularity, measured by download counts. A significant majority, 70%, of adapters belong to the character and celebrity category, primarily consisting of anime or game characters. Another 13% of adapters modify image style, 8% adjust clothing, and 4% represent various concepts (Fig. 13a). These statistics indicate that our experiments consider a minor proportion of adapters, as the COCO dataset does not feature characters or celebrities [19]. Despite this, Stylus outperforms base Stable Diffusion. Furthermore, the popularity of adapters follows a Pareto distribution, where the top adapters receive exponentially more downloads than the others (Fig. 13a). However, the top adapter accounts for only 0.5% of total downloads, which suggests that the distribution is long-tailed.
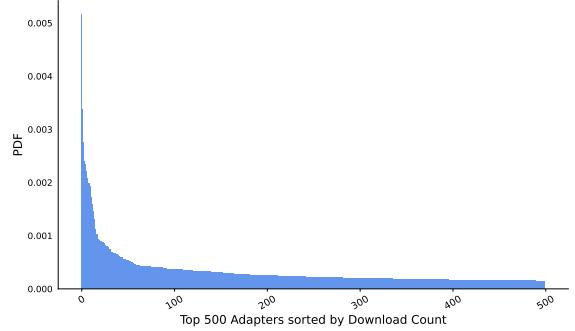
## A.4. Failure Modes

We detail different failure modes that were discovered while developing Stylus.

**Image saturation.** The quality of image generation is highly depend on adapters' weights. If the assigned weight



(a) Distribution of adapters across *categories*.



(b) Top 500 adapters ranked by *percentage of downloads*.

Figure 13. **Workload Characterization of `StylusDocs`.** (a) Most adapters are categorized as characters or celebrities. (b) Adapter popularity exhibits a power-law distribution, with the top adapters receiving exponentially more downloads than the others.

is above the recommended value, the adapter negatively impacts image generation, leading to a growing number of visual inconsistencies and artifacts. In Fig. 14a, assigning a high weight to a "James Bond" LoRA increases images exposure and introducing significant visual tearing. Stylus mitigates over-saturation with its refiner component, which extract the right adapter weights from the adapter's model card. Lastly, Stylus uniformly weights adapters based on their associated tasks, ensuring that similar adapters do not significantly impact their corresponding tasks.

**Task Blocking.** Composing adapters presents the risk of overwriting existing concepts or tasks specified in the prompt and other selected adapters. We illustrate several examples in Figure 2—a train LoRA overrides the toy train concept (left), a park bench LoRA masks a person in an orange blanket (middle), and a fancy cake LoRA erases the image of a man eating the cake (right). Task blocking typically arises from two main issues: the adapter weight set too high or too many adapters merged into the base model. Stylus addresses this by reducing an adapter's weight with uniform weighting per task, while the masking scheme reduces the number of selected adapters. Although Stylus does not completely solve task blocking, it offers simple heuristics to mitigate the issue.

**Image Prompts**

Prompt 1: Photo of Dwayne Johnson, wearing military clothes and cap, dramatic lighting, `<lora:TheRockV3:0.9>`.
Prompt 2: Photo of Dwayne Johnson, wearing a Superman suit, high quality, `<lora:TheRockV3:1>`.
Prompt 3: Photo of Dwayne Johnson, wearing an Armani tuxedo, `<lora:TheRockV3:0.9>`

**Model Card Description**

- Title: Dwayne "The Rock" Johnson (LoRA)
- Tags: Celebrity, Photorealistic, Hollywood, Celeb
- Trigger Words: Th3R0ck
- Description: Had to make this one, due to Kevin Hart Lora. Recommended lora strength: 0.9. *% Author descriptions may be misleading or incomplete.*

Your goal is improve the description of a model adapter's task for Stable Diffusion, with images, prompts, and descriptions pulled from popular model repositories. Above, we have provided the following information and the associated constraints:

1. Examples of generated images (from left to right) from the adapter and the corresponding user-provided prompts.
- Some prompts may specify the adapter weight (i.e. `<lora:NAME:WEIGHT>`). If provided, you will need to infer the adapter's name and weight. Prioritize this weight over the author's recommended weight.
2. The adapter's model card from the original author. This includes the title, tags, trigger words, and description.
- The model card description may be incorrect, misleading, or incomplete.
- The model card may specify the weight of the model adapter, or the recommended range. Find the recommended weight of the adapter (default is 0.8).

*% Chain-of-Thought Prompting*

Again, your mission is to provide a clear description of the model's adapter purpose and its impact on the image. To do so, you should implicitly categorize the model adapter into only one of the following topics: [Concept, Style, Pose, Action, Celebrity/Character, Clothing, Background, Building, Vehicle, Animal, Action]. Do not associate an adapter with a topic that is vague or uninteresting.

First, describe the topic associated with the adapter and explain how this adapter alters the images, based on the common elements observed in the example images. Your requirements are:
- Do not describe any training or dataset-related details.
- Provide additional context from your prior knowledge if there is insufficient information.
- Do not hallucinate and repeat text. Output only english words and sentences.

Second, recommend an optimal weight for the adapter as a float. Do not specify a range, only give one value.

Please format your output as follows:

Example 1: [*Description of adapter and its weight*]

Example 2: [*Description of adapter and its weight*]

Table 2. Full prompt for the refiner VLM to generate better adapter descriptions.

**Task Diversity.** Merging adapters into the base model overwrites the base model's prior distribution over an adapter's corresponding tasks. If an adapter is not finetuned on a diverse set of images, diversity is significantly reduced among different instances of the same task. We present three examples in Fig. 14c, over different prompts that specify multiple instances of the same task (teddy bears, women, and apples). We observe that all instances of each task are highly identical with one another. Stylus offers no solution to address or mitigate this problem.

**Low quality adapters.** Low quality adapters can significantly degrade the quality of image generation, as shown by corrupted images in Fig. 14d. This issue typically arises from poor training data or from fine-tuning the adapter for too many epochs. Stylus attempts to blacklist such adapters.

| **Retrieved Adapter Descriptions** |
| --- |

**42**: This LoRA is for the concept of dragon, a mythical creature. It generates images of dragons with a variety of different appearances, including both Western and Eastern styles...

**120**: This LoRA steers the image generation towards a fantasy castle, with a focus on the building and its surroundings. The castle is depicted as a grand structure, often with towering spires, intricate architecture, and a sense of grandeur...

**3478**: This LoRA is designed to generate images of a Chinese dragon breathing fire. It generates images of a dragon with a long, serpentine body, covered in scales, with a large head and sharp teeth. The dragon is breathing fire, with flames coming out of its mouth...

**1337**: This LoRA is designed to generate images of animals breathing fire. It generates images of animals, such as rabbits, dragons, and frogs, breathing fire. The fire is shown as a bright, orange-yellow flame that is coming out of the animal's mouth...

**...**

Provided above are the IDs and descriptions for different model adapters (e.g. LoRA) for Stable Diffusion that may be related to the prompt. Your goal is to fetch adapters that can improve image fidelity. The prompt is:

*Dragon breathing fire on a castle.*

*% Chain-of-Thought Prompting*
First, segment the prompt into different tasks—concepts, styles, poses, celebrities, backgrounds, objects, actions, or adjectives—from the prompt's keywords.

Here are the requirements for <u>tasks</u>:
- Tasks should never introduce new information to the prompt. The topic must be selected from the prompt's keywords.
- Different tasks must be orthogonal from each other.
- All tasks combined must span the entirety of the prompt.
- Prioritize choosing narrower tasks. You may merge tasks if a relevant adapter spans several tasks.

Second, for each task, provide 0-5 of the most relevant model adapters to the task. For each adapter, infer an adapter's main function from its description. This function must directly match at least one task and the context of the prompt. If the adapter is indirectly relevant, do not include it.

Here are the requirements for <u>adapters</u>:
- Adapters should only be used at most once across all tasks. If an adapter is used in one task, it should not be used in another task.
- Adapters should not introduce novel concepts or biases to the topic or the prompt. Do not include such adapters.
- Adapters cannot encompass a broader scope relative to its assigned task. For example, if the task is about a "dog", the adapter cannot be about general "animals".
- Adapters cannot be too narrow in scope relative to its assigned task. For example, if a task is about pandas, do not choose highly specific pandas such as the character "Po" from Kung Fu Panda. However, it is acceptable to choose adapters that modify the style of the task, such as "Red Pandas".
- If an adapter spans multiple tasks, merge these tasks together. For example, if there is an adapter that is about "fluffy cats", merge the topics "fluffy" and "cats" together.
- Avoid choosing NSFW and anthropormorphic adapters.

Finally, for each selected adapter, provide a strong reason for why this adapter is relevant to the prompt, directly matches the keyword, and improves image quality.

Give me the answer only. Please format your output as follows:

Example 1: [*Dictionary of tasks to the associated adapter ids and reasons for their selection.*]

Example 2: [*Dictionary of tasks to the associated adapter ids and reasons for their selection.*]

Table 3. Full prompt for the composer LLM.

However, our blacklist is not comprehensive, and as a result, Stylus may still occasionally select low-quality adapters.

**Retrieval Errors.** Stylus's retrieval process involves three stages, each introducing potential errors that can compound in later stages. For instance, the refiner may return incor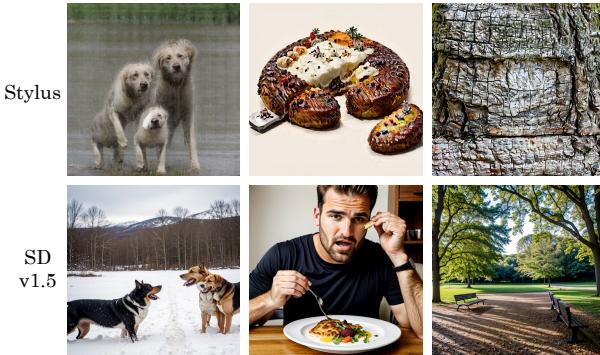rect descriptions of an adapter's task, while the composer may classify the adapter into an incorrect task. We detail three examples in Figure 4. Stylus selects an "okapi" (forest giraffe) LoRA, known for its distinctive zebra-like appearance, causing the generated giraffes to adopt the okapi's skin texture. In the middle, Stylus selects a flowery vase LoRA, a misinterpretation of the prompt "orange flowers placed in a vase." On the right, the composer

(a) **Image Saturation.** Assigning too high of a weight to a "James Bond" adapter leads to significant degradation in visual fidelity.
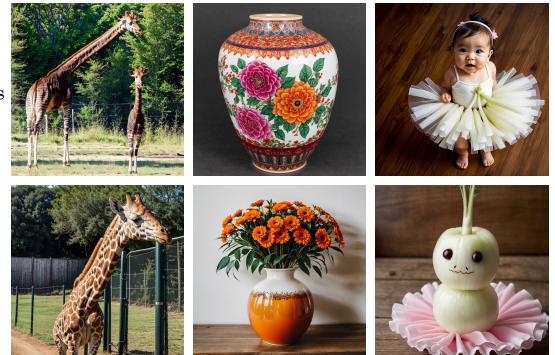


(b) **Task Blocking.** Adapters can block a prompt's or other adapter's tasks (i.e. toy trains, person in orange blanket, or man eating cake).



(c) **Task Diversity.** Adding an adapter reduces diversity of instances within a single task (i.e. teddy bears, woman, and apples).



(d) **Low quality adapters.** Low quality adapters can significantly impact visual fidelity. We blacklist such adapters.



(e) **Retrieval Errors.** Retrieval errors can lead to foreign biases in image generation and deliberate misinterpretations of the prompt.

Figure 14. **Categorization of Different Failure Modes.**

incorrectly chooses a human baby adapter for the prompt "a baby daikon radish in a tutu.", resulting in images of babies instead of daikons. Stylus includes an option to self-repair faulty composer outputs with multi-turn conversations, which can improve adapter selection.

## A.5. VLM as a Judge

The full prompts to GPT-4V as a judge for textual alignment, visual fidelity, and image diversity are specified in Tables 4 and 5.

To distinguish the two images (or groups of images), the VLM exploits multi-turn prompting: We provide each image (or group of images) labeled with IMAGE/GROUP A or IMAGE/GROUP B. Note that the ACK messages are not generated by the VLM; instead, it is part of VLM's context window. We provide the rubric, detailed instructions, reminders, and example model outputs in our prompt. For scoring, the VLM employs Chain-of-Thought (CoT) prompting to output scores 0-2, similar to VisDiff [6, 42, 44]. We observe that larger ranges (5-10) leads the model towards abstaining from making decisions, as it avoids out-

---

**System Prompt:**
You are a photoshop expert judging which image has better composition quality.

**Scoring:** Compositional quality scores can be 2 (very high quality), 1 (visually aesthetic but has elements with distortion/missing features/extra features), 0 (low visual quality, issues with texture/blur/visual artifacts).

Composition can be broken down into three main aspects:
- **Clarity**: If the image is blurry, poorly lit, or has poor composition (objects obstructing each other), it gets scores 0.
- **Disfigured Parts**: This applies to both body parts of humans and animals as well as objects like motorcycles. If the image has a hand that has 6 fingers it gets a 1 for having otherwise normal fingers, but the hand should not have two fingers. If the fingers themselves are disfigured showing lips and teeth warped in, it gets a 0.
- **Detail**: If the sail of a sailboat's sail shows dynamic ripples and ornate patterns, this shows detail and should get a score of 2. If it's monochrome and flat, it gets a score of 1. If it looks like a cartoon and is inconsistent with the environment, give a score of 0.

**Scoring:** Alignment scores can be 2 (fully aligned), 1 (incorporates part of the theme but not all), 0 (not aligned).

We provide several examples:
- If the prompt is 'shoes', and an image is a sock, this is not aligned and gets a score of 0.
- If the prompt is 'shoes without laces', but the shoes have laces, this is somewhat aligned and gets a score of 1.
- If the prompt is 'a concert without fans', but there's fans in the image, pick the images that show fewer fans.

---

**User:**
This is IMAGE A. Reply 'ACK'.
*% Generated Image from Group A*

---

**Assistant:** ACK

---

**User:**
This is IMAGE B. Reply 'ACK'.
*% Generated Images from Group B*

---

**Assistant:** ACK

---

**User:**
Rate the quality of the images in GROUP A and GROUP B. For each image, provide a score and explanation.

Image A Quality: <SCORE>(<EXPLANATION>)
Image B Quality: <SCORE>(<EXPLANATION>)
Preference: Group <CHOICE>(<EXPLANATION>)

*% Prevent VLM from returning neutral results.*
I'll make my own judgement using your results, your response is just an opinion as part of a rigorous process. I provide additional requirements below:
- You must pick a group for 'Better Quality' / 'Better Alignment', neither is not an option.
- If it's a close call, make a choice first then explain why in parenthesis.

---

Table 4. Full prompt judging compositional quality (left) or textual alignment (right) using VLM.

putting extreme scores. However, the score range 0-2 provides the VLM sufficient granularity to express preferences and prompt the model to summarize the key differences.

**Textual Alignment.** The VLM scores how well a generated image follows the prompt's specifications. We note that prompts with negations (e.g. "concert with no fans" or "harbor with no boats") fail for both Stylus and the Stable Diffusion checkpoint. Hence, we prompted the VLM to assign better scores for images that produced less fans or boats. Furthermore, as adapters can potentially block

existing concepts in the image (see Fig. 14b), the VLM allocates partial credit in scenarios where images partially capture the set of keywords in the prompt.

**Visual Quality.** Our evaluation assesses visual quality through three metrics: clarify, disfigurements, and detail. First, the VLM assigns low clarity scores if an image is blurry, poorly lighted, or exhibits poor compositional quality. We note that LoRAs are trained over specific tasks/concepts; the model determines how to compose different concepts. For instance, a rhinoceros LoRA combined

Table 5. Full prompt judging diversity using VLM.

with a motorcycle LoRA led to images of motorcycles draped with rhinoceros hide. As such, the VLM assigns partial credit when the model fails to combine concepts in a meaningful way. Second, the VLM assigns lower scores by judging if an image has disfigured parts. For instance, diffusion models have trouble accurately depicting a human hand, oftentimes generating extra fingers. Finally, the VLM's final score depends on the detail of image. We find that adapters are able to bring greater detail to certain concepts. For example, an elephant adapter generates elephants with much greater detail than that of the base model. However, we note that the VLM is not good at detecting subtleties in detail.

**Diversity.** For each prompt, we generate five images each for Stylus and the Stable Diffusion checkpoint. These images are then assessed with a VLM (Visual Language Model, GPT-4V) judge, which rates and ranks them based on diversity. In Tab. 5, we measure diversity through two metrics. The first metric, theme interpretation, measures di-versity based on the interpretation of the prompt, which is often under-specified. We find that different thematic interpretations improves model response due to non-ambiguity. The second metric measures diversity by the variance of focus across different subjects. We find that many prompts often under-specify which subject is the focus on the image.

**A.6. Additional Diversity Scores**

Fig. 15 decomposes $d$FID scores over the top 100 keywords in the PartiPrompts dataset. We highlight that the largest differences stem from concepts, appearances, attributes, or styles. For example, Stylus excels over concepts ranging from animals ("bears", "sloth", and 'squirrel') to objects ("microphone", "box", and "jacket"). Selected attributes can include but are not limited to: ("white", "blue", and "photographic"). Regardless of keyword, Stylus attains higher diversity scores across the board.
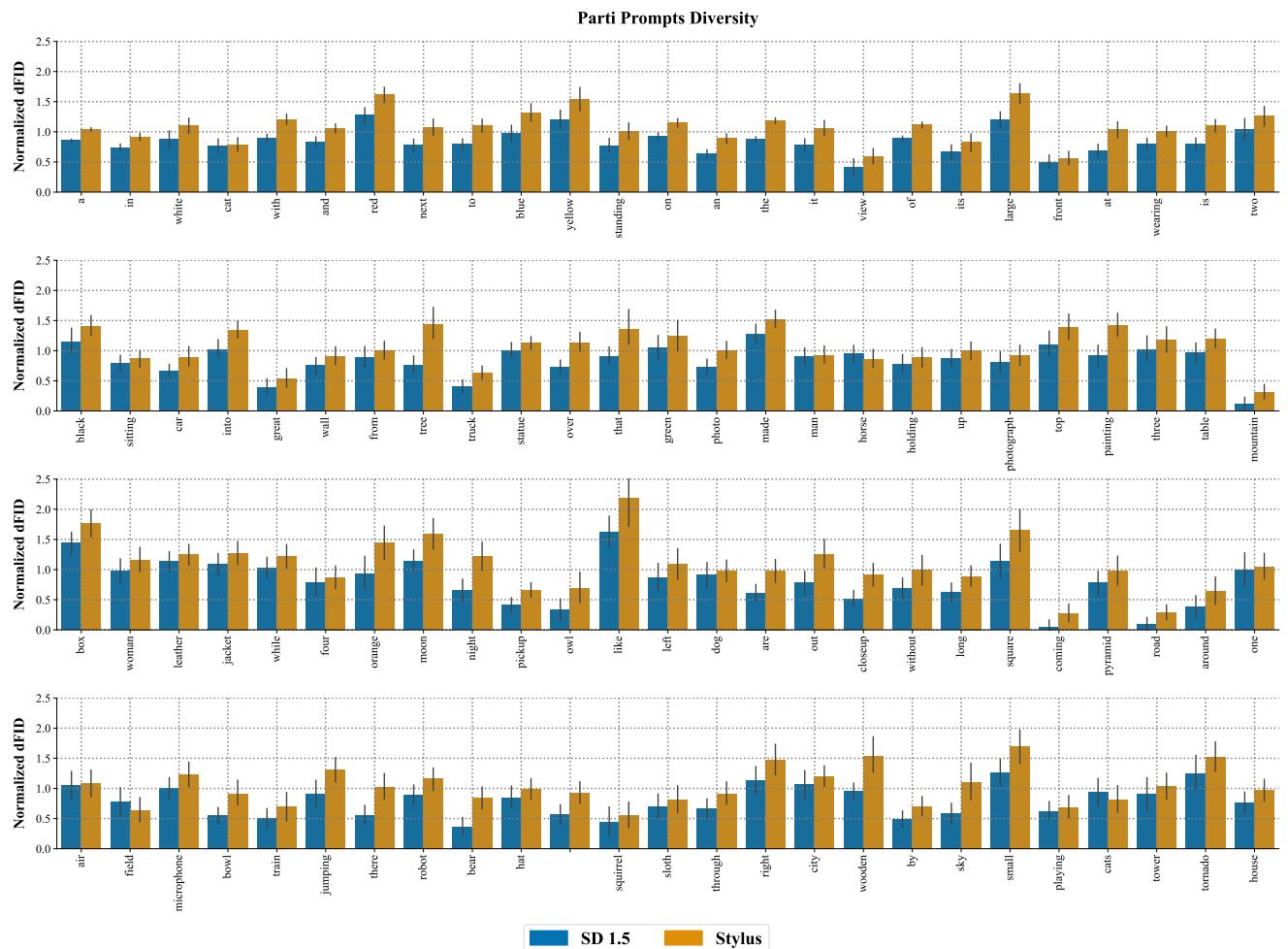
Figure 15. *d***FID for top 100 keywords in PartiPrompts dataset.** Stylus leads to consistently higher diversity when compared to Stable Diffusion checkpoints, especially for words describing concepts and attributes.