

Stylus: Automatic Adapter Selection for Diffusion Models

Michael Luo¹ Justin Wong¹ Brandon Trabucco² Yanping Huang³ Joseph E. Gonzalez¹ Zhifeng Chen³ Ruslan Salakhutdinov² Ion Stoica¹

¹UC Berkeley ²Carnegie Mellon University MLD ³Google Deepmind

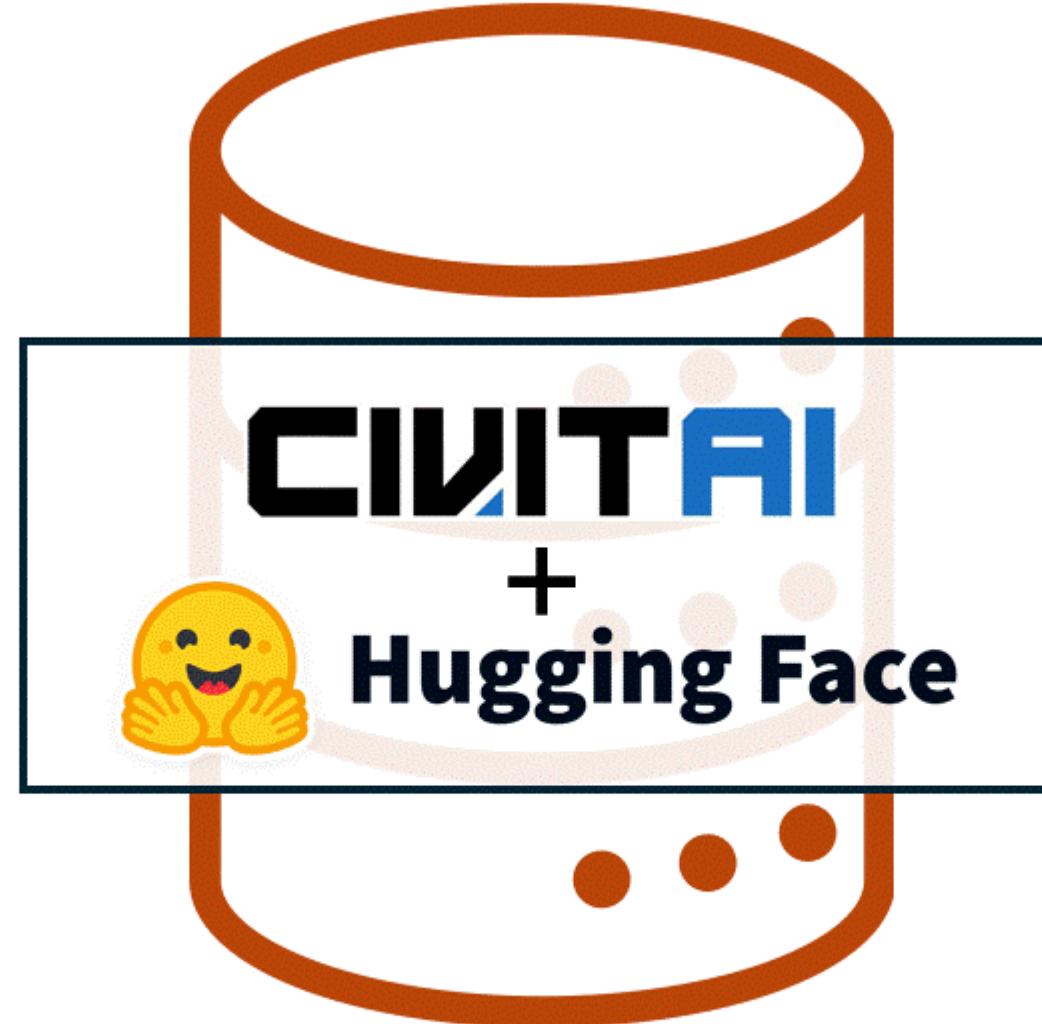


Stylus: Automatic Adapter
Selection For Diffusion Models



A castle that looks
like Studio Ghibli.

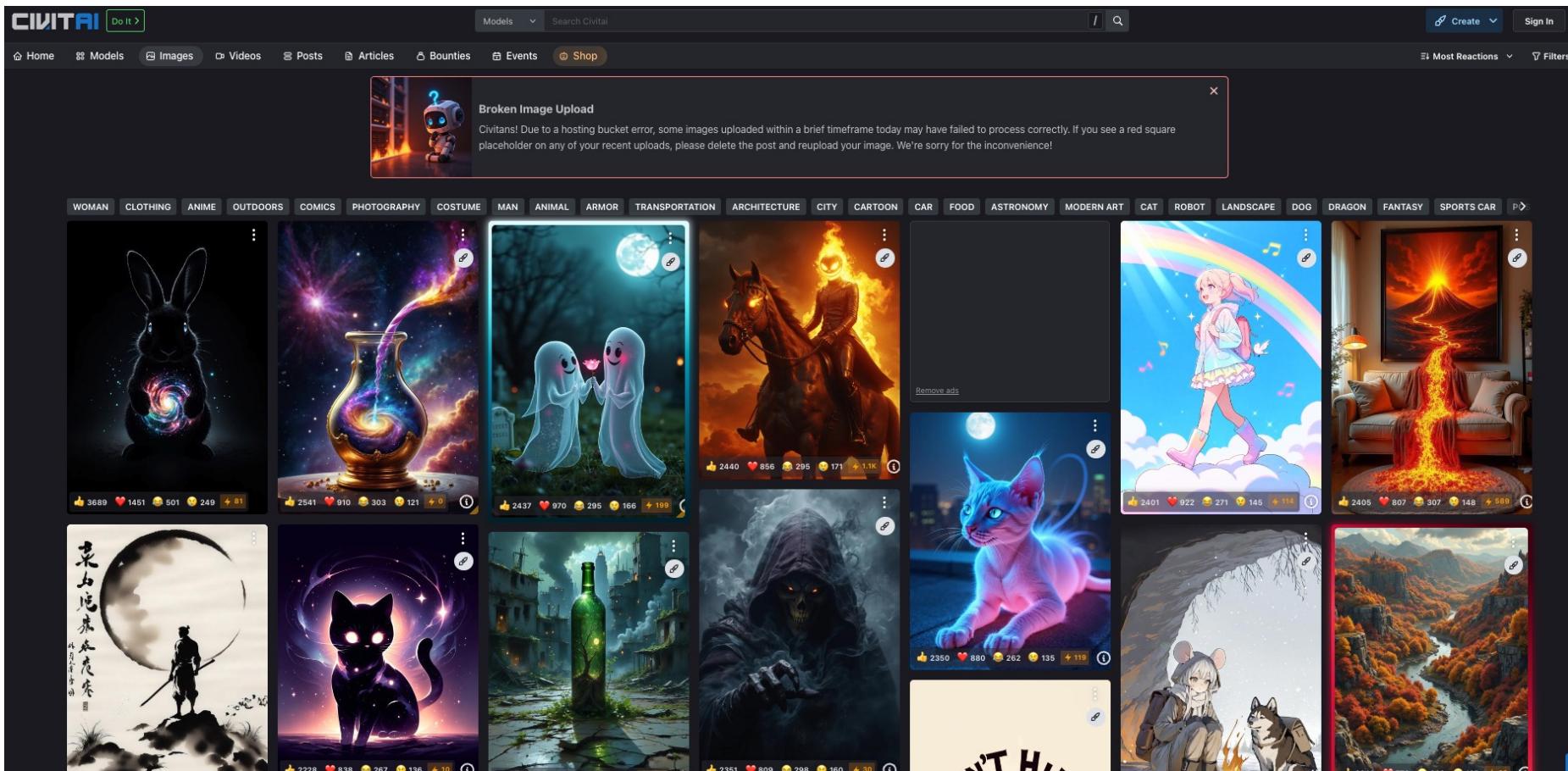




stylus-diffusion.github.io

CIVITAI

- Diffusion Model은 굉장히 다양한 이미지를 생성할 수 있다.
- CIVIT AI 커뮤니티를 Diffusion Model을 생성하는 방법을 갤러리로 저장



- Diffusion Model은 굉장히 다양한 이미지를 생성할 수 있다.
- CIVIT AI 커뮤니티를 Diffusion Model을 생성하는 방법을 갤러리로 저장

LoRA (Low-Rank Adaptation)

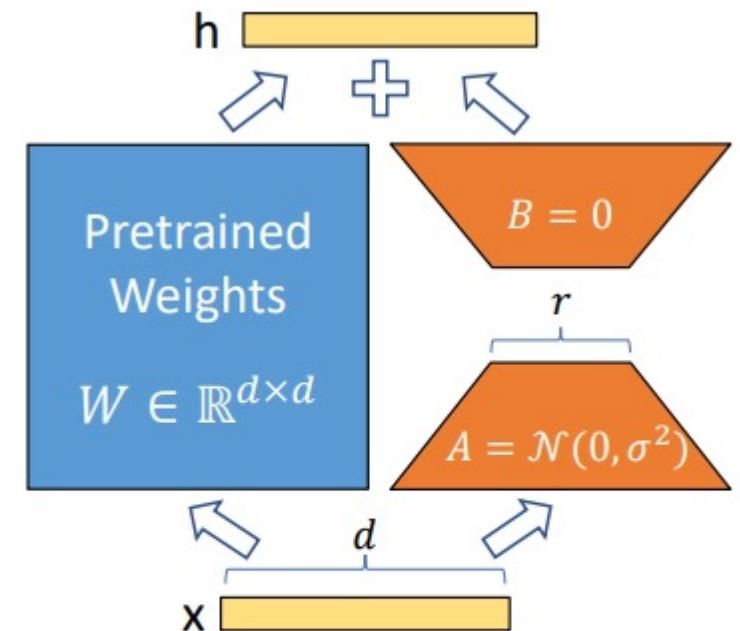
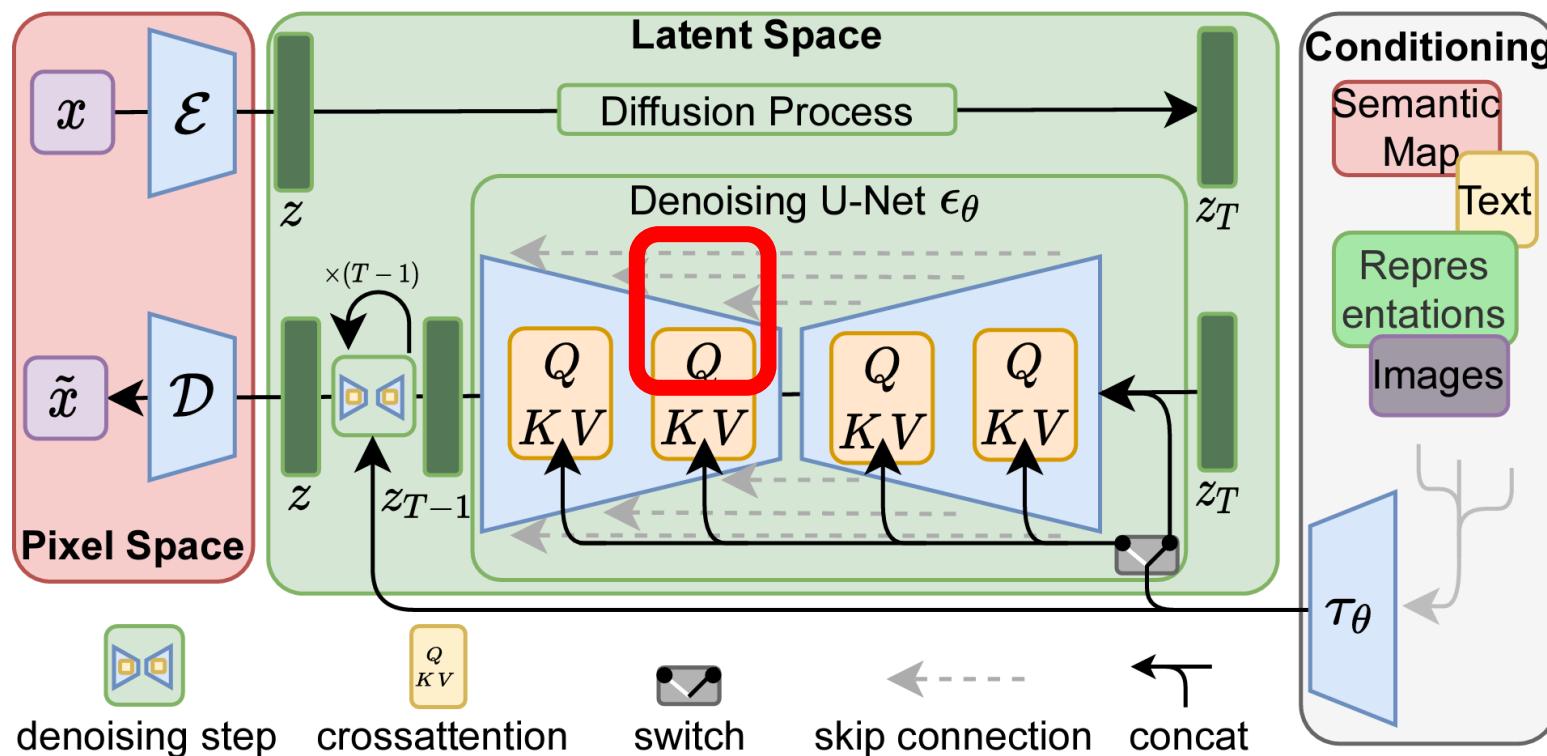
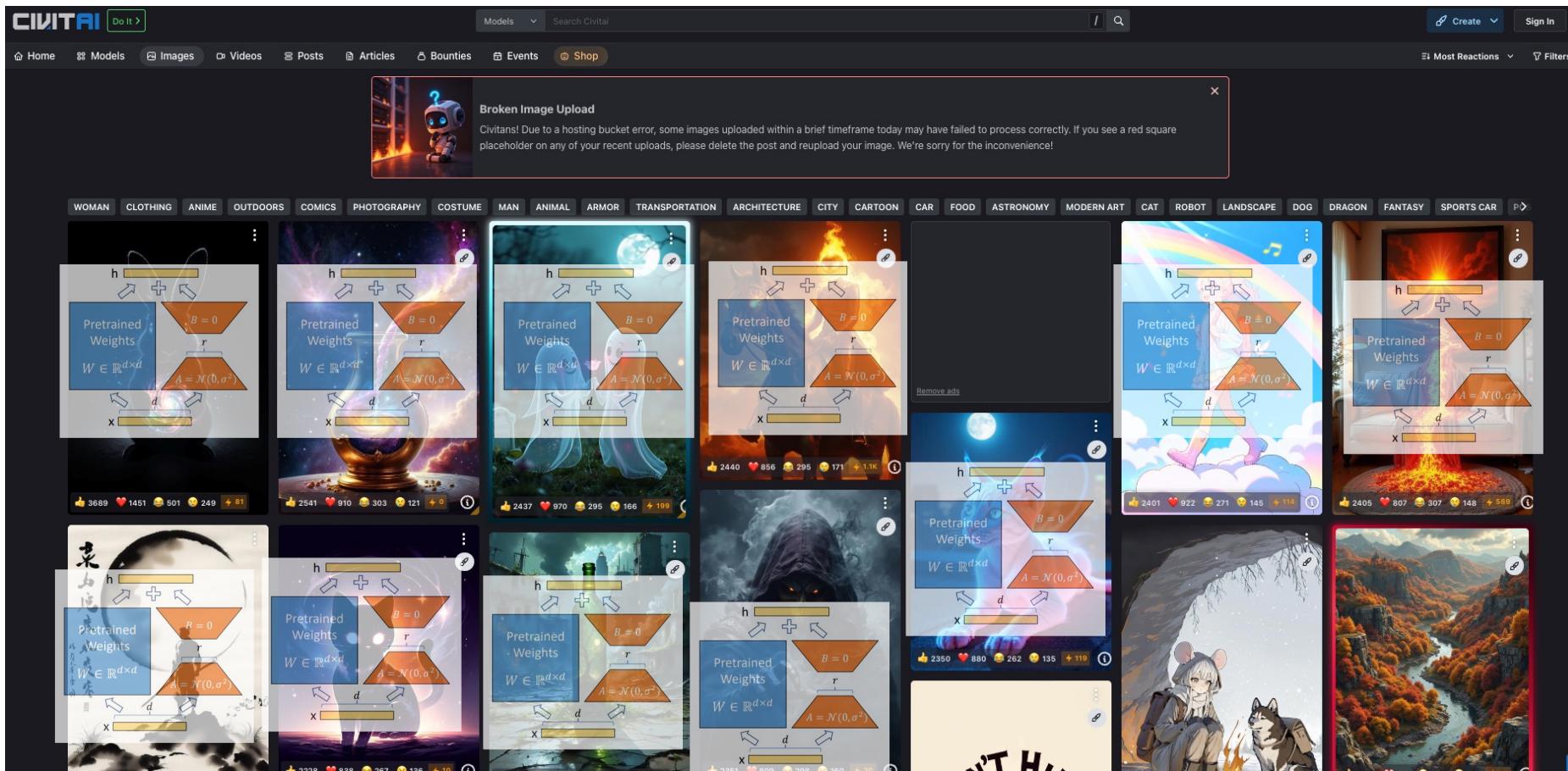


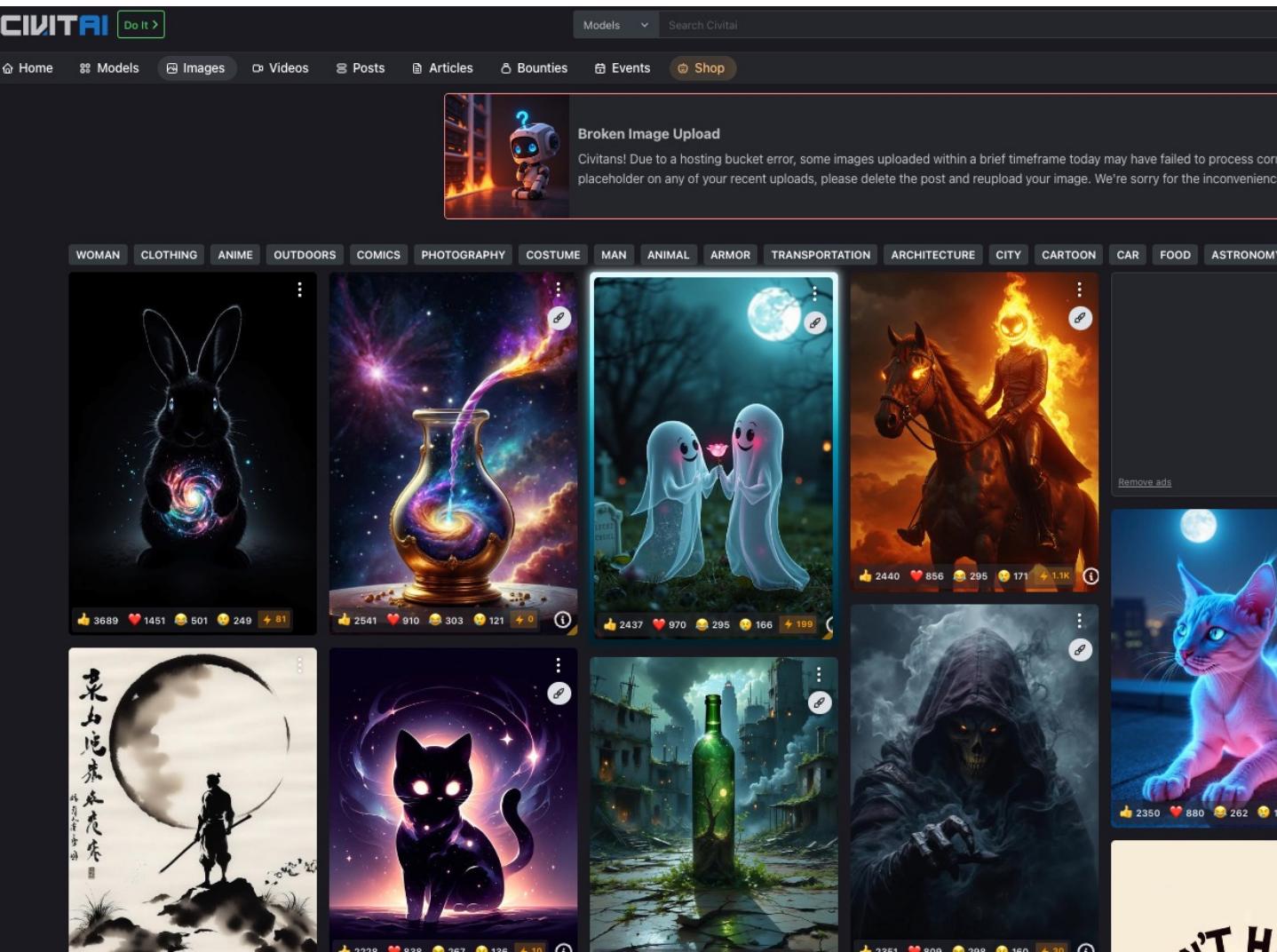
Figure 1: Our reparametrization. We only train A and B .

- Diffusion Model은 굉장히 다양한 이미지를 생성할 수 있다.
- CIVIT AI 커뮤니티를 Diffusion Model을 생성하는 방법을 갤러리로 저장 (각 이미지마다)

LoRA (Low-Rank Adaptation)



- 내가 원하는 정확한 이미지를 찾기 어려움.



Stylus: Retrieval을 편하게!

- **입력:** 텍스트 (Two dogs playing in the snow)
- **Retriever:** VectorDB에서 LoRA 파라미터를 찾아 줌

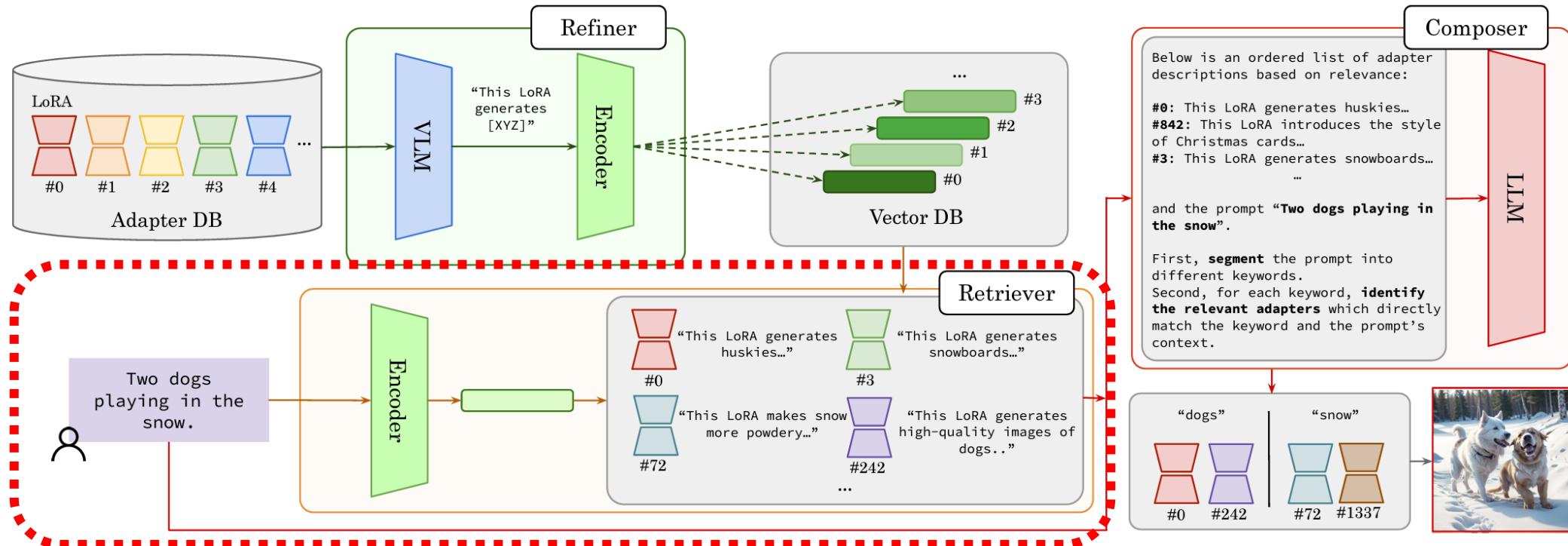


Figure 2. **Stylus algorithm.** Stylus consists of three stages. The *refiner* plugs an adapter’s model card through a VLM to generate textual descriptions of an adapter’s task and then through an encoder to produce the corresponding text embedding. The *retriever* fetches candidate adapters that are 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.

Stylus: Retrieval을 편하게!

- **입력:** 텍스트 (Two dogs playing in the snow)
- **Composer:** LoRA들의 설명과 입력을 매칭해서 어떤 것을 사용할지 결정

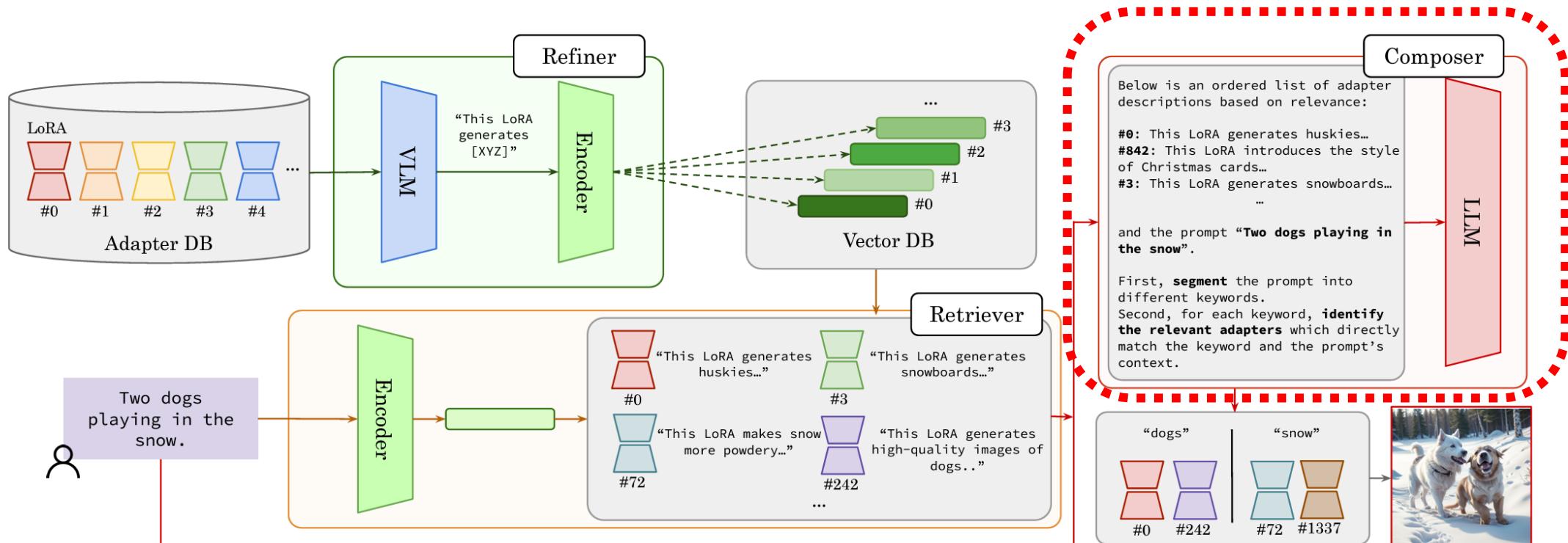


Figure 2. **Stylus algorithm.** Stylus consists of three stages. The *refiner* plugs an adapter’s model card through a VLM to generate textual descriptions of an adapter’s task and then through an encoder to produce the corresponding text embedding. The *retriever* fetches candidate adapters that are 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.

Stylus: Retrieval을 편하게!

- **입력:** 텍스트 (Two dogs playing in the snow)
 - **Composer:** LoRA들의 설명과 **입력을** 매칭해서 어떤 것을 사용할지 결정

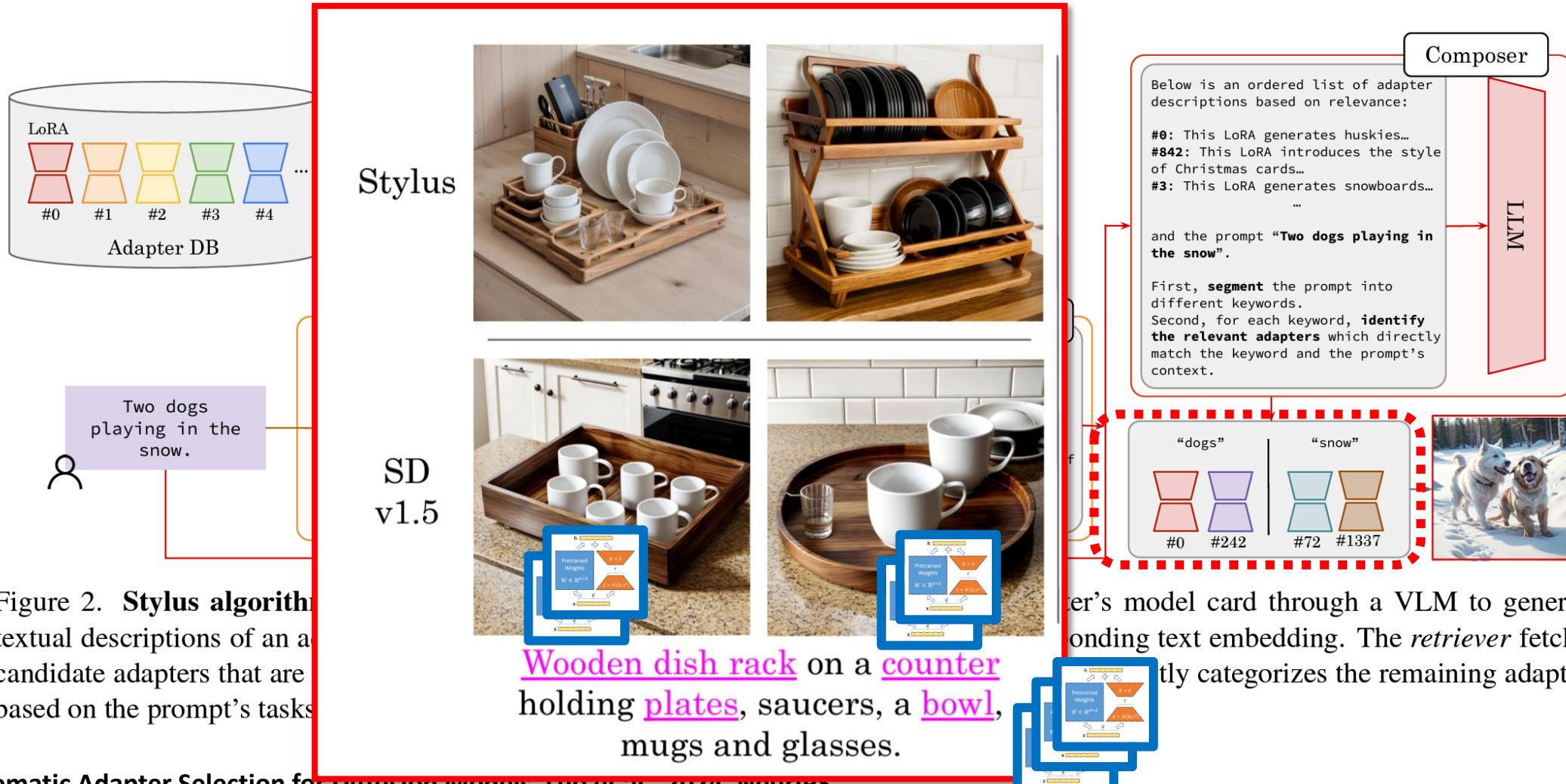


Figure 2. Stylus algorithm
textual descriptions of an array of candidate adapters that are based on the prompt’s tasks

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

ter's model card through a VLM to generate a bonding text embedding. The *retriever* fetches the adapter category and then the *adapter* itself categorizes the remaining adapters.

Stylus: Retrieval을 편하게!

- **입력:** 텍스트 (Two dogs playing in the snow)
- **Retriever:** VectorDB에서 LoRA 파라미터를 찾아 줌 (**Embedding 유사도 기반**)

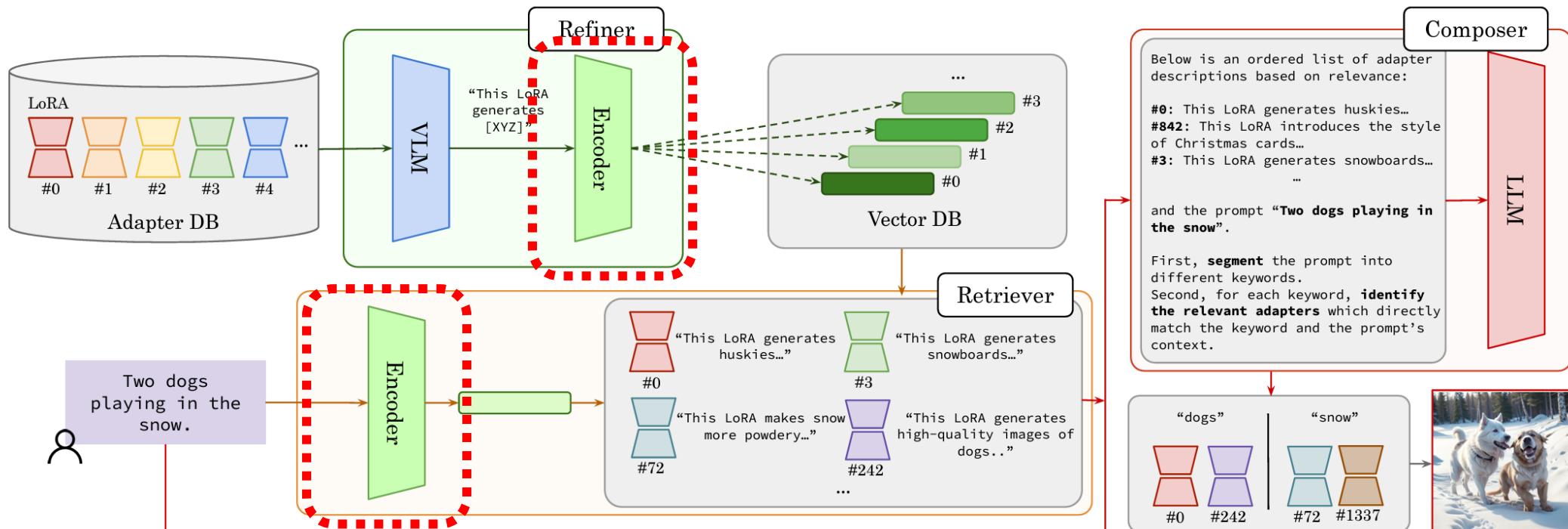
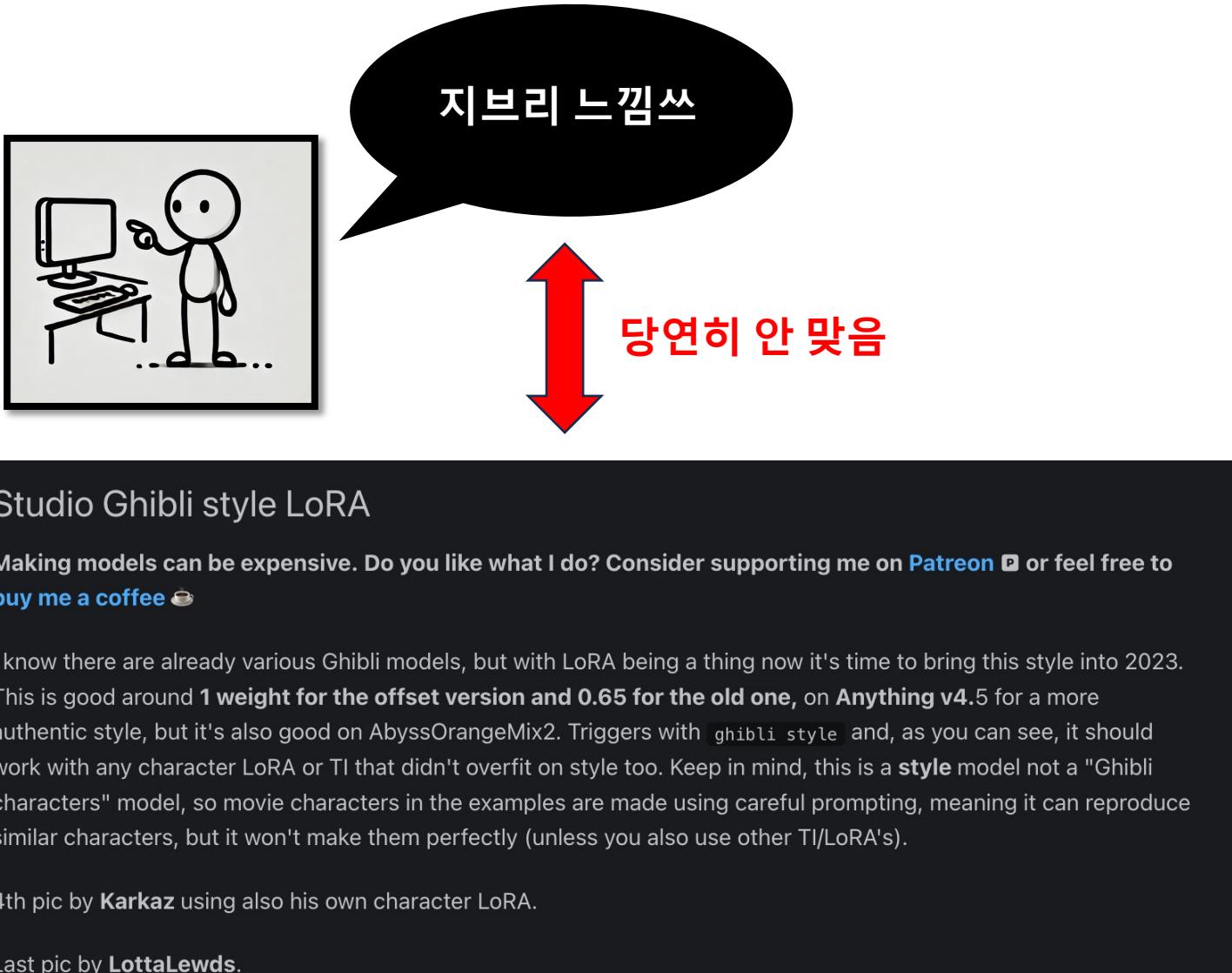
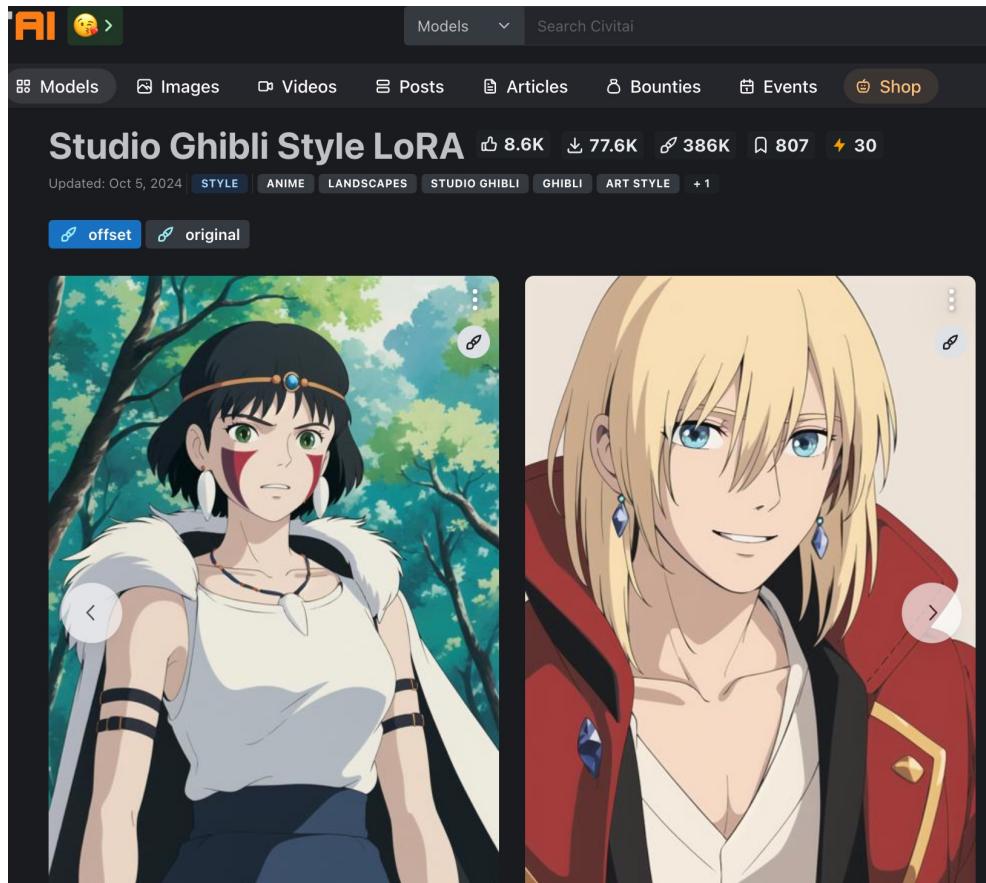
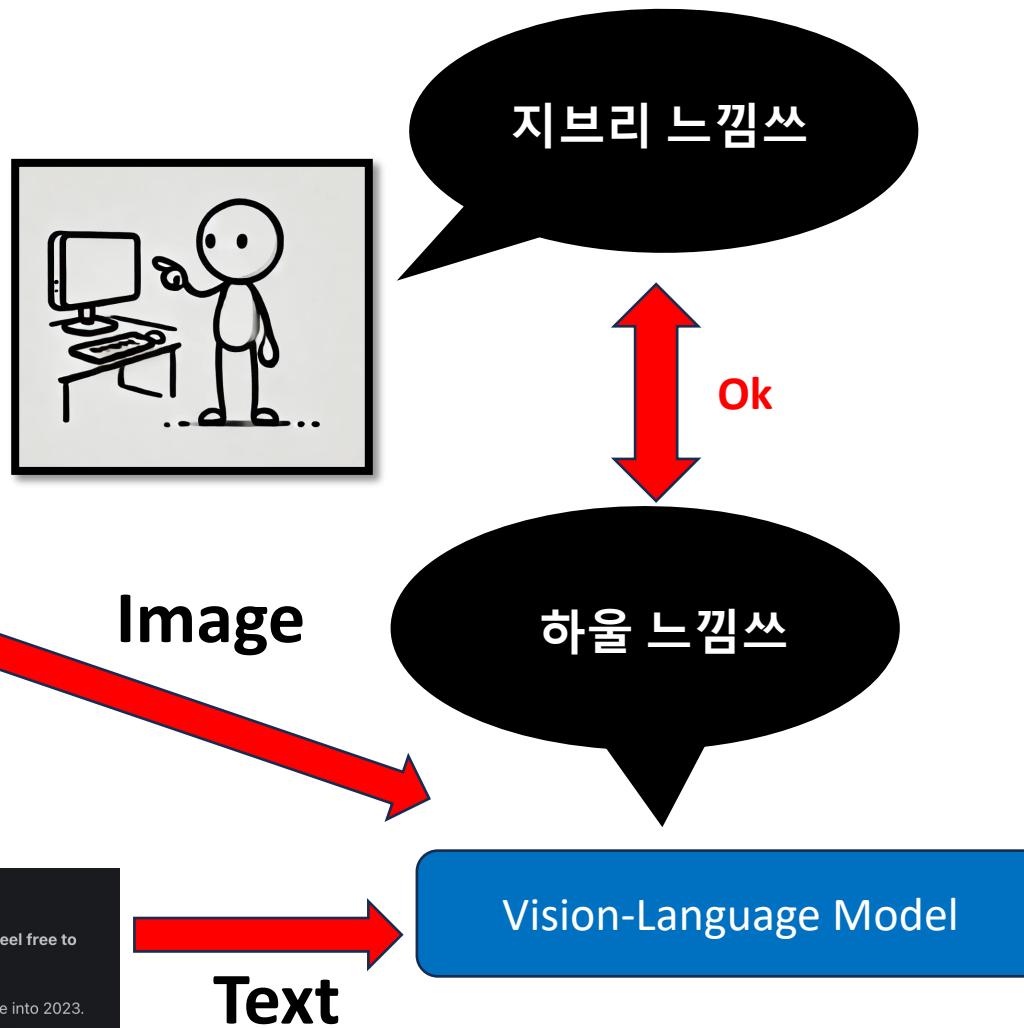
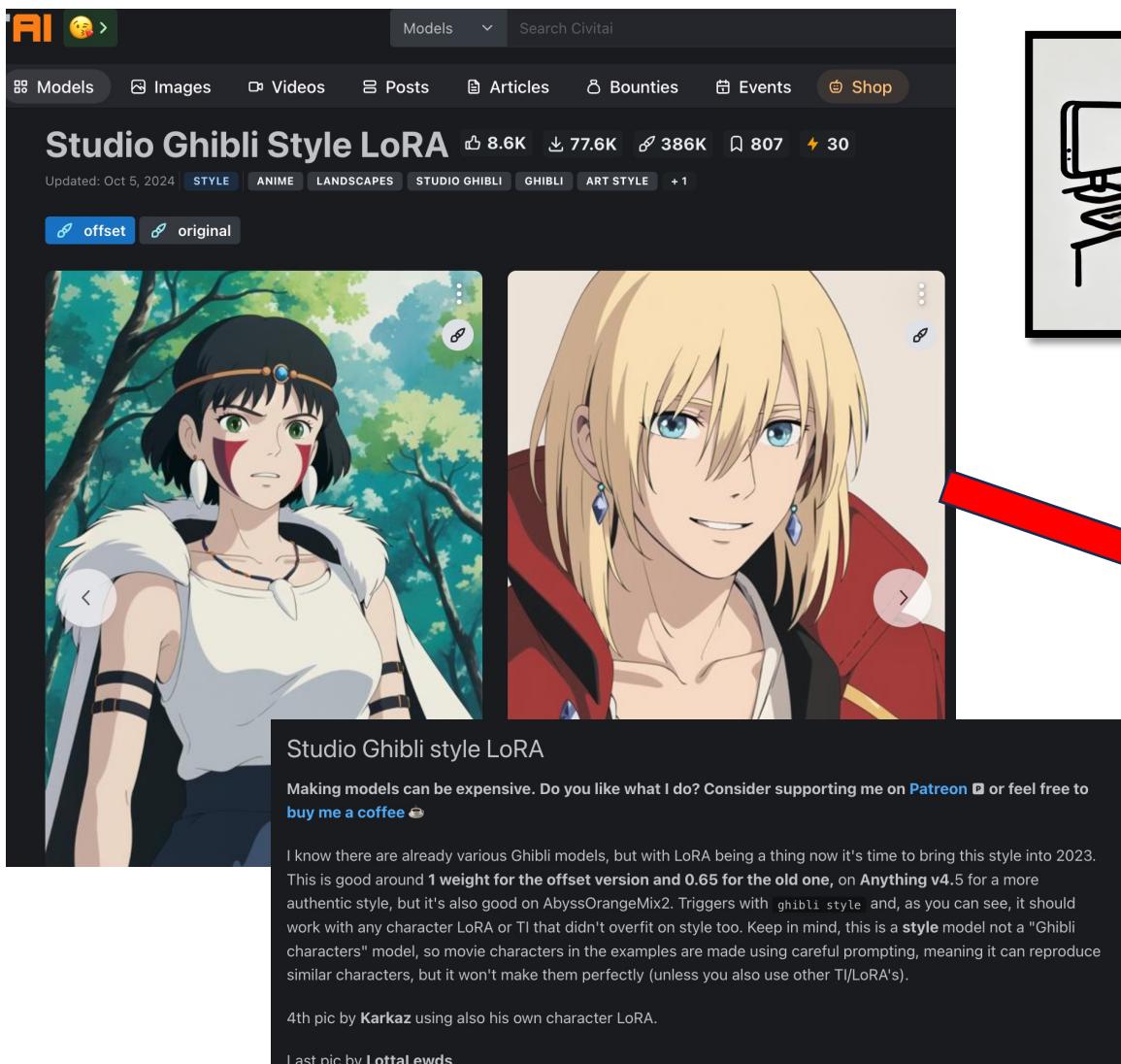


Figure 2. **Stylus algorithm.** Stylus consists of three stages. The *refiner* plugs an adapter’s model card through a VLM to generate textual descriptions of an adapter’s task and then through an encoder to produce the corresponding text embedding. The *retriever* fetches candidate adapters that are 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.

사용자 프롬프트 vs LoRA Description



사용자 프롬프트 vs LoRA Description



Stylus: Retrieval을 편하게!

- Refiner (VLM) : CIVIT-AI LoRA 정보 후처리
- Encoder (Text Embedding): 언어를 벡터로 변환

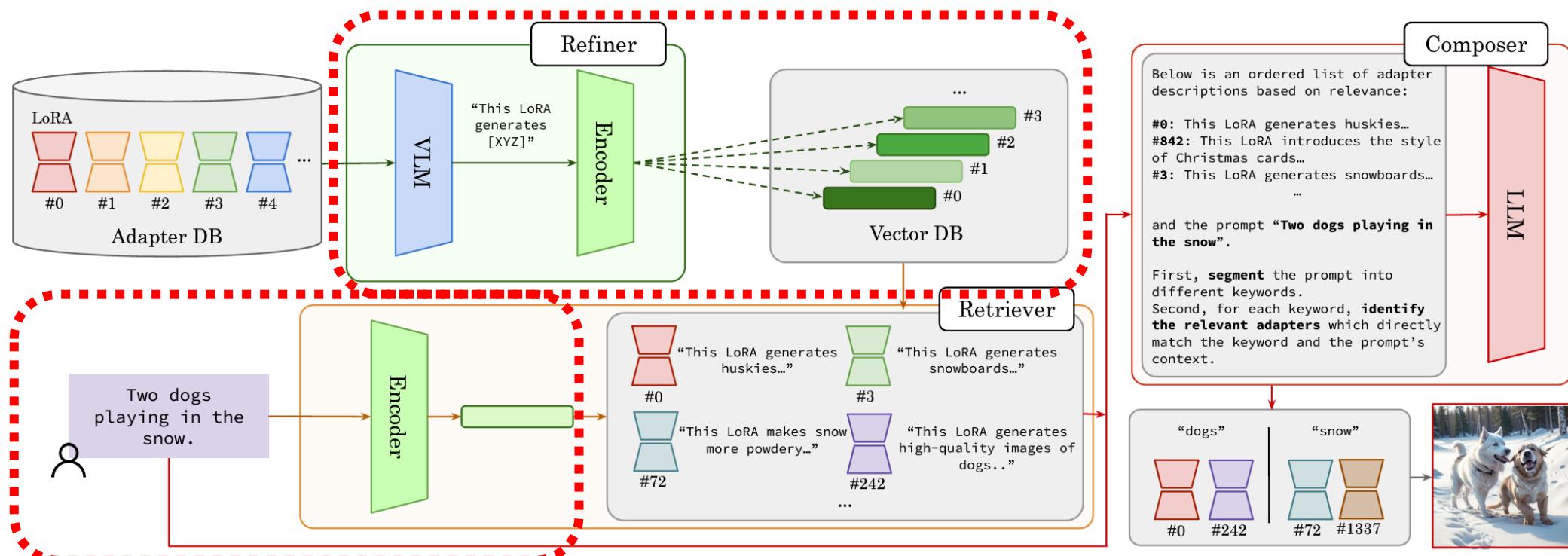


Figure 2. **Stylus algorithm.** Stylus consists of three stages. The *refiner* plugs an adapter’s model card through a VLM to generate textual descriptions of an adapter’s task and then through an encoder to produce the corresponding text embedding. The *retriever* fetches candidate adapters that are 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.

Stable Diffusion vs Stylus (SD + LoRA)

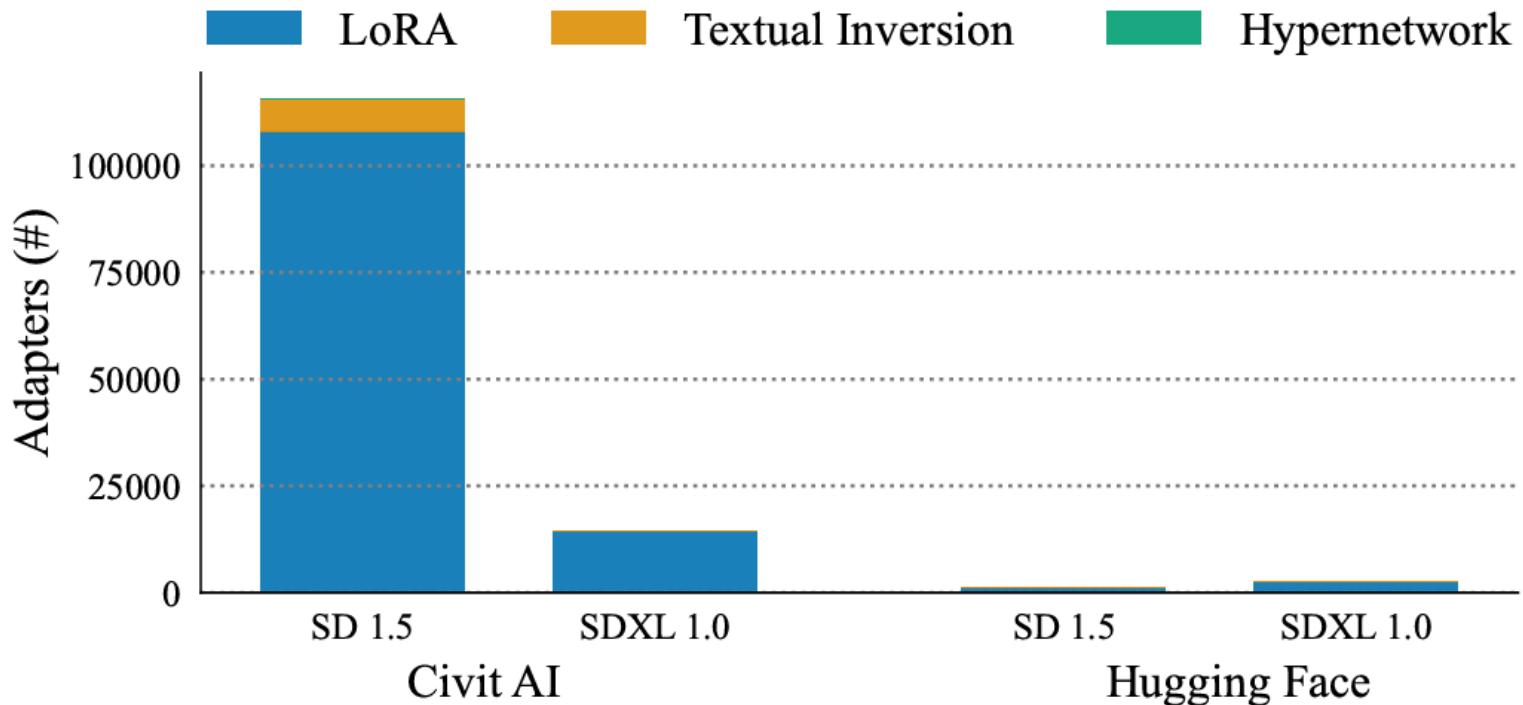


Figure 3. **Number of Adapters.** Civit AI boasts 100K+ adapters for Stable Diffusion, outpacing that of Hugging Face. Low-Rank Adaptation (LoRA) is the dominant approach for finetuning.

Stable Diffusion vs Stylus (SD + LoRA)



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.*

Stable Diffusion vs Stylus (SD + LoRA)

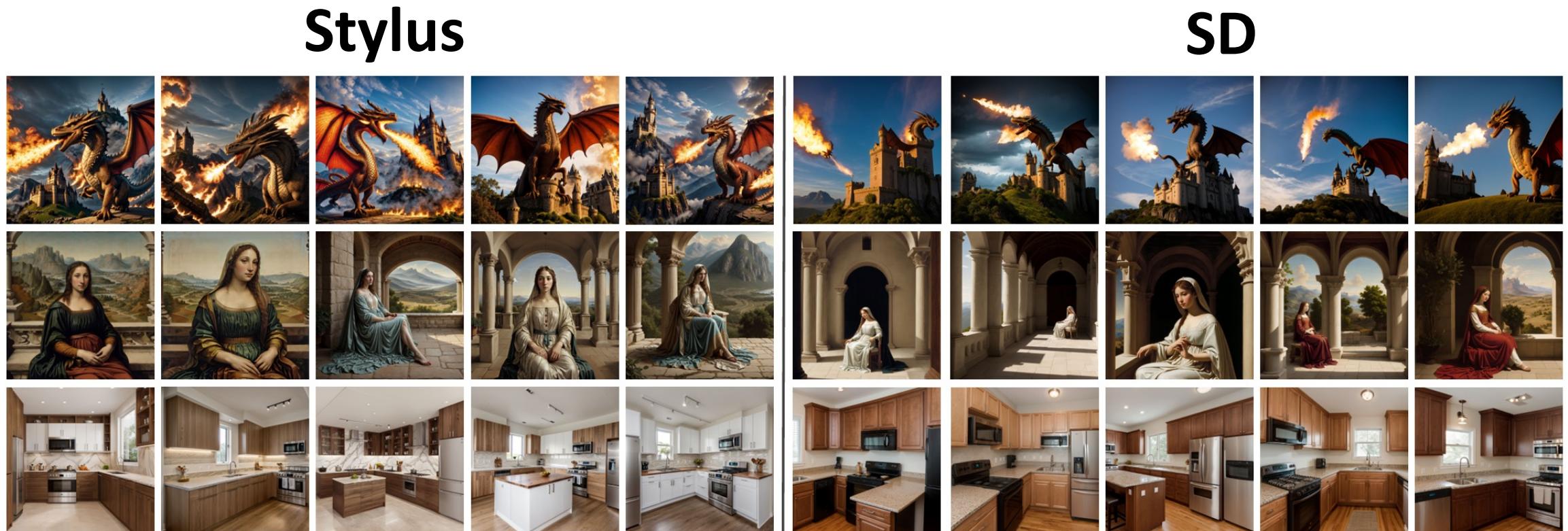


Figure 7. **Image Diversity.** Given the same prompt, our method (left) generates more diverse and comprehensive sets of images than that of existing Stable Diffusion checkpoints (right). Stylus's diversity comes from its masking scheme and the composer LLM's temperature.

Stable Diffusion vs Stylus (SD + LoRA)



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.

Stable Diffusion vs Stylus (SD + LoRA)

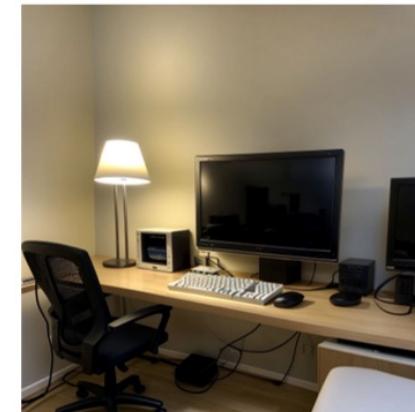
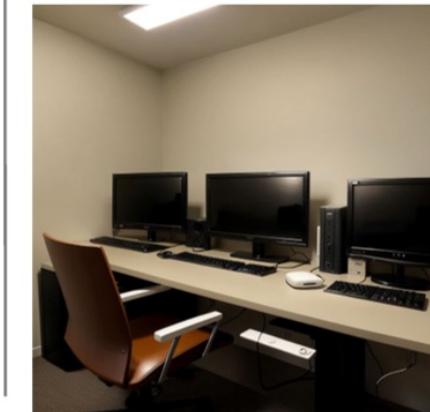
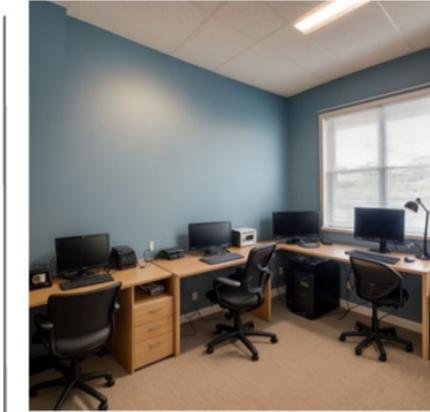
Stylus



SD
v1.5



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



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

Stable Diffusion vs Stylus (SD + LoRA)

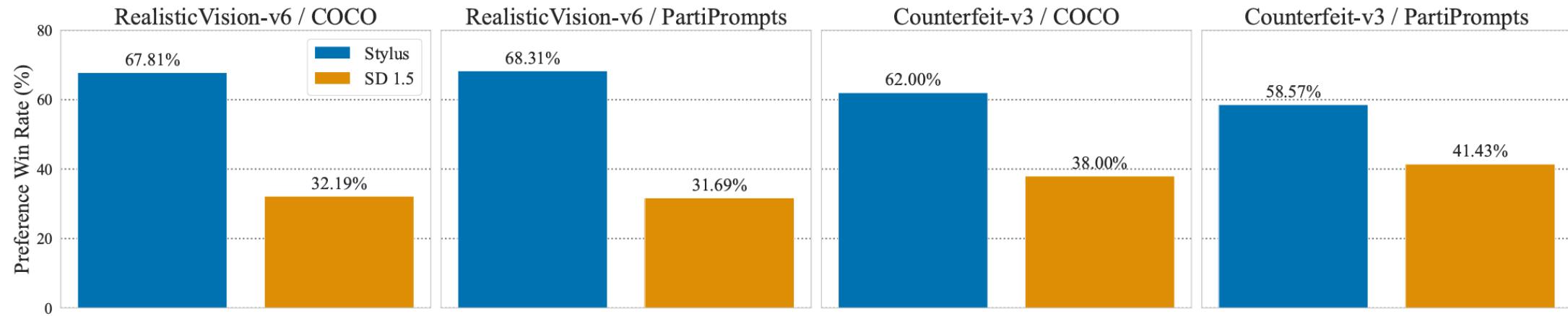
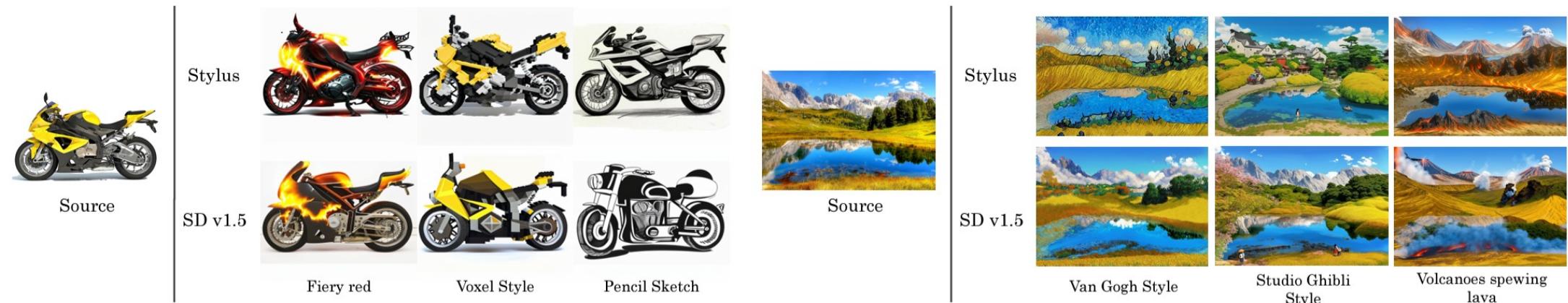
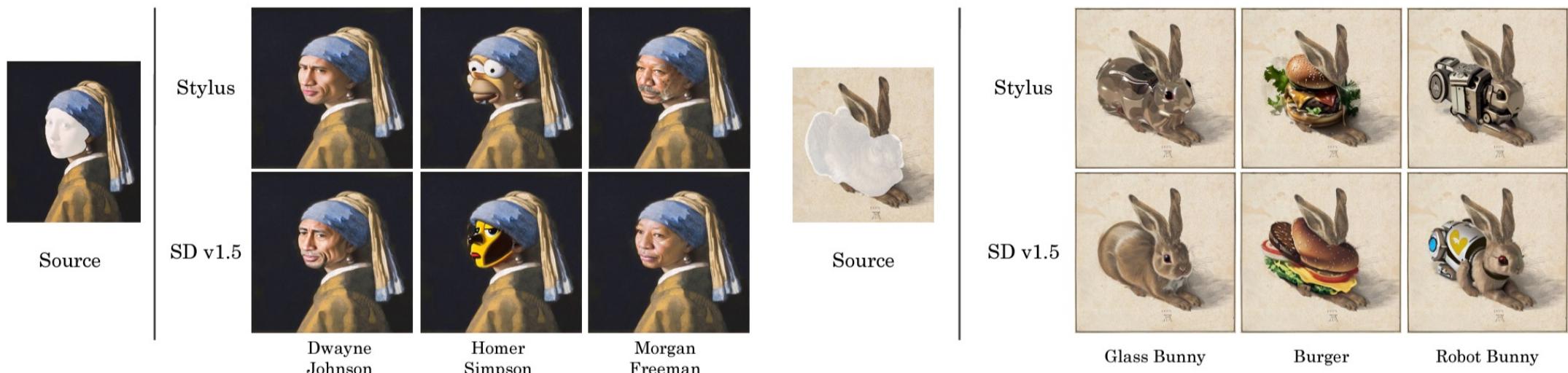


Figure 5. **Human Evaluation.** Stylus achieves a higher preference scores (2:1) over different datasets and Stable Diffusion checkpoints.

Stable Diffusion vs Stylus (SD + LoRA)



(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.

Stable Diffusion vs Stylus (SD + LoRA)

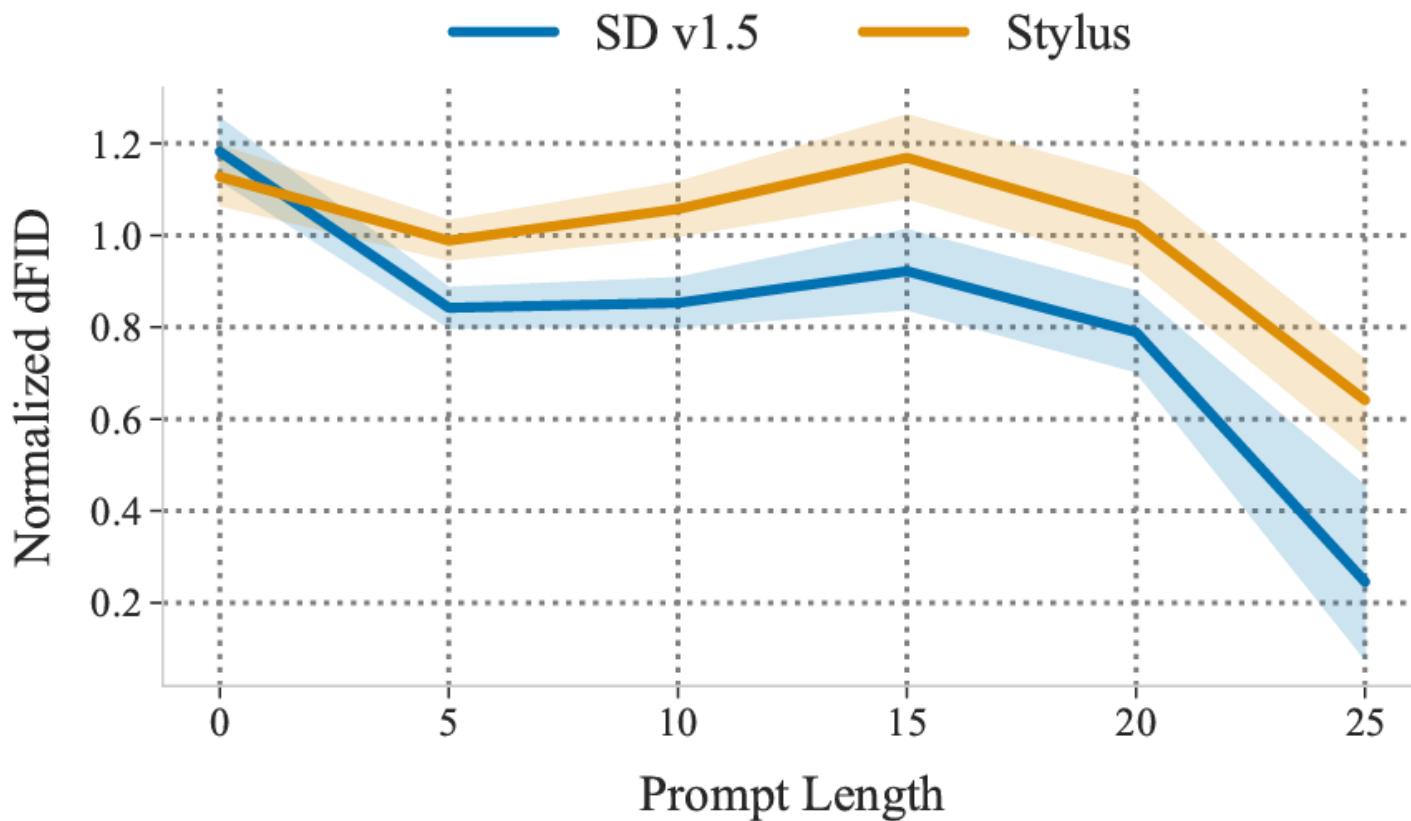


Figure 9. Image diversity ($dFID$) across prompt length. Stylus achieves higher diversity score than Stable Diffusion when prompt length increases.

Stable Diffusion vs Stylus (SD + LoRA)

	CLIP (Δ)	FID (Δ)
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.