# Stylus 🖌️

## Automatic Adapter Selection

**Main Presenter: Michael Luo**

Paper Authors: Michael Luo, Justin Wong, Brandon Trabucco, Yanping Huang, Joseph Gonzalez, Zhifeng Chen, Ruslan Salakhutdinov, Ion Stoica
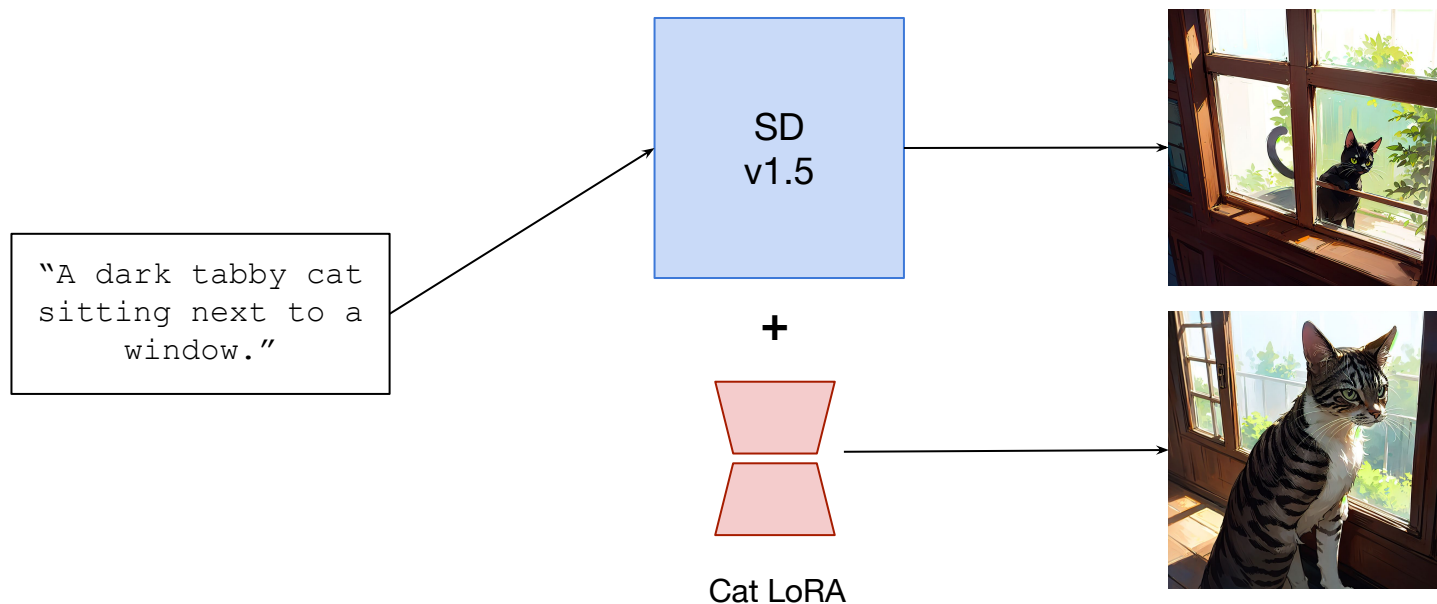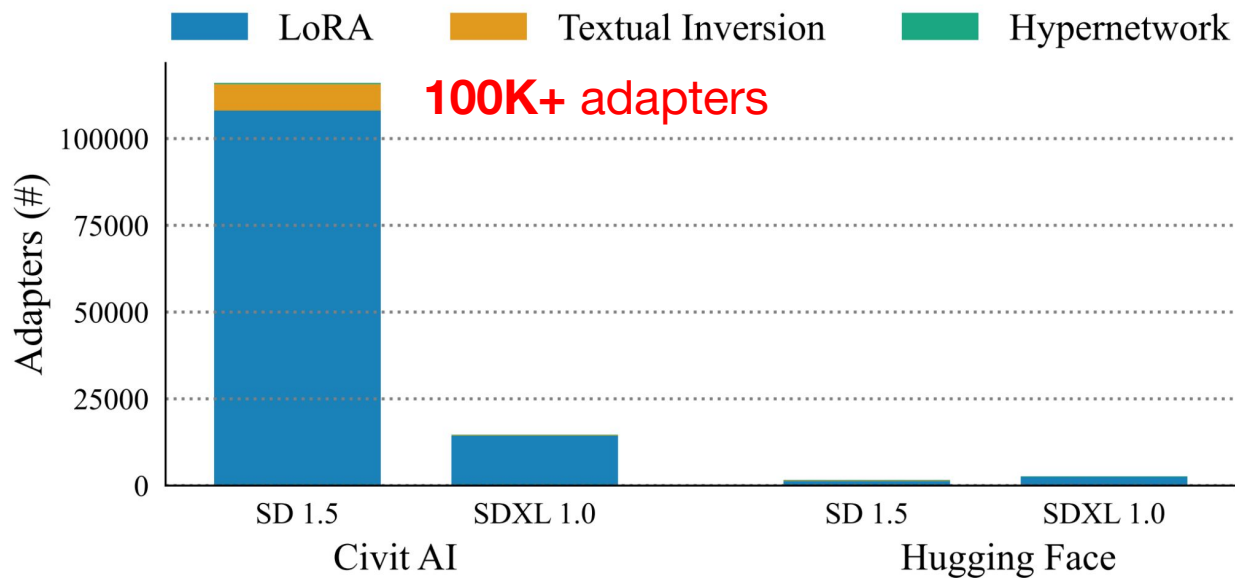
# Situation

Finetuned Adapters (i.e. LoRA, textual inversion) introduce novel concepts and styles to the base model, thereby improving image quality.

# Situation

Open-source contributors have created over **100K+ adapters**!

# An Emergent Problem

- Base models are no longer good enough for generating images.

- Users *manually mix and match* many checkpoints and adapters to generate the images they want.





*Image taken from Civit.ai*

# Stylus 🖌️

**Goal**: Automatically *select* and *compose* the right adapters given a user prompt.
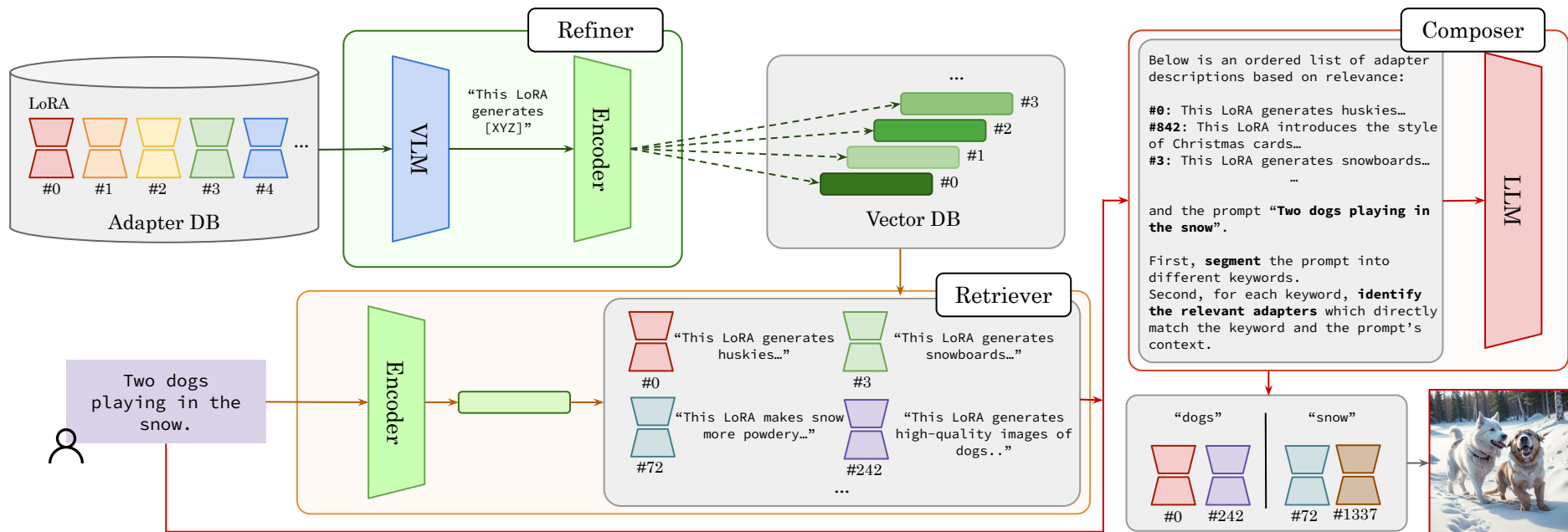


*stylus-diffusion.github.io*

# Stylus 🖌️

**Goal**: Automatically *select* and *compose* the right adapters given a user prompt.
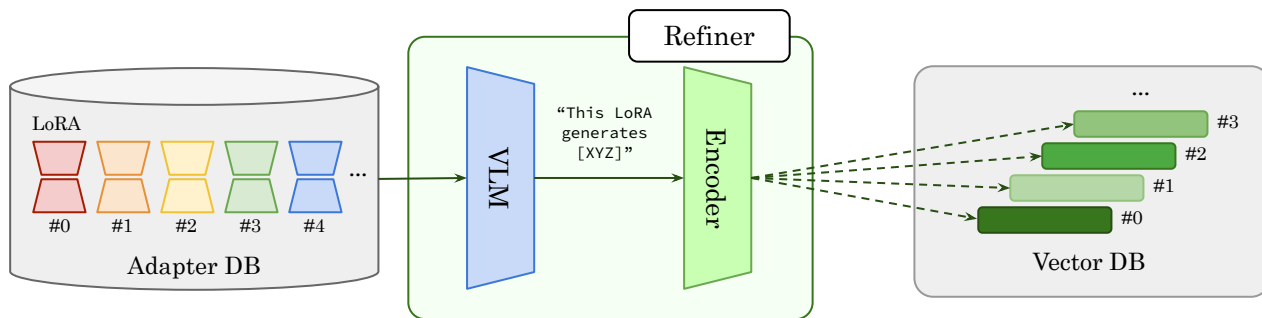
**Requirements**:
- No training. Scales to new adapters over time.
- Performant. Must do better than existing retrieval approaches.
- Low inference overhead. Identifies the right adapters quickly.

# Stylus Architecture

# Refiner



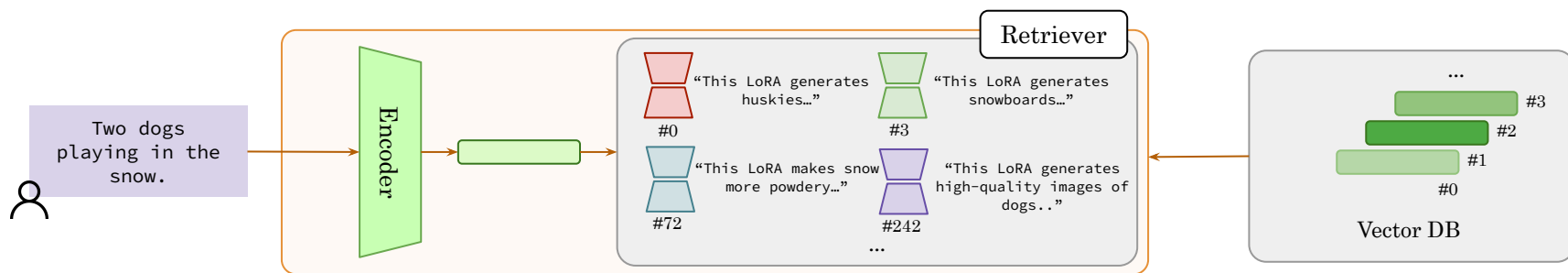**Goal:** Generate a vector embedding for each adapter.

**How?**
- Vision Language Model (VLM) to infer adapters' descriptions
  - With adapter's model card (author description + example images)
- Embedding Model to embed adapter description

# Retriever (RAG)

**Goal:** Retrieves the most relevant candidate adapters via similarity metric.

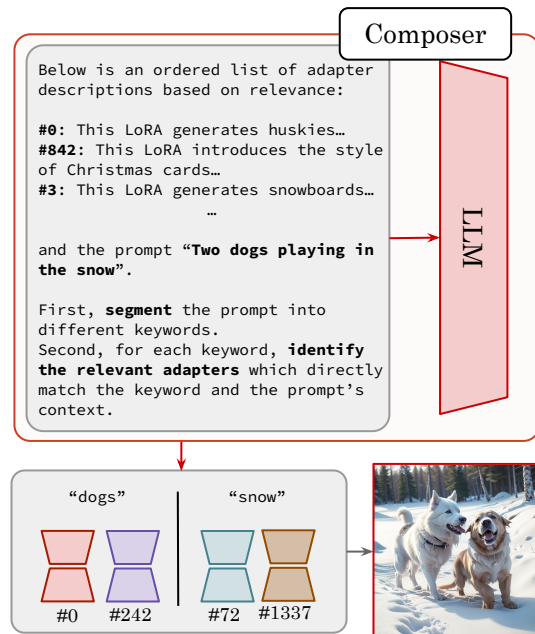**Problem**: RAG lacks precision. Easy to add *slightly relevant* adapters.

# Composer

**Goal:** Significantly improves precision for retrieving the *right* adapters.
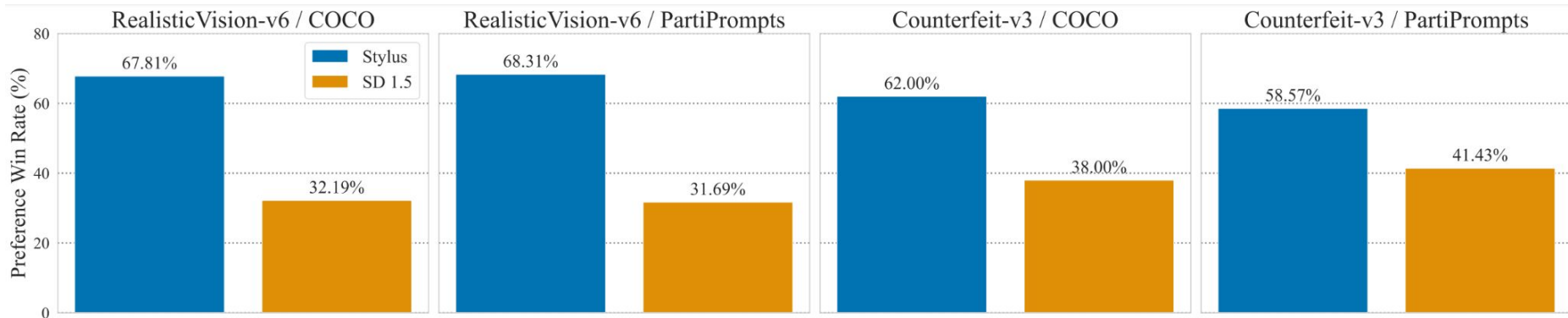
**How?**
- Evaluate adapters based on *semantic relevance*.
  - Maps relevant adapters to keywords in user prompt.
    - Enforces that adapters stay relevant.
  - Prunes irrelevant adapters.
- Efficient. Only requires one LLM call.

# Evaluation

- **Prompt Datasets**
  - COCO 2014
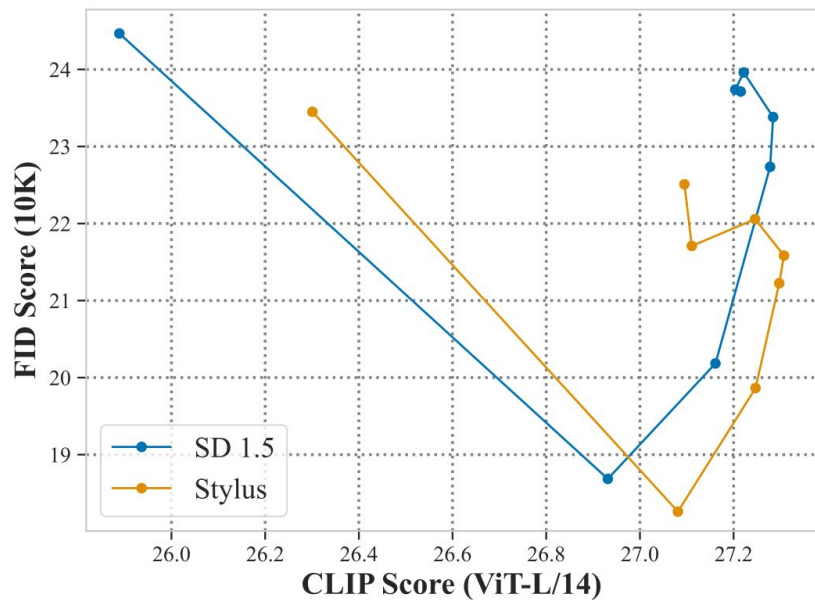  - PartiPrompts

- **Adapter Database**
  - StylusDocs v2
    - Adapters from Civit AI + Huggingface.
    - 75K generated adapter descriptions from GPT-4o.

- **Evaluation**
  - Human Evaluation
  - FID/CLIP Pareto Curve
  - VLM as a judge

- **Base Models**
  - RealisticVision (SD v1.5)
  - CounterFeit (SD v1.5)
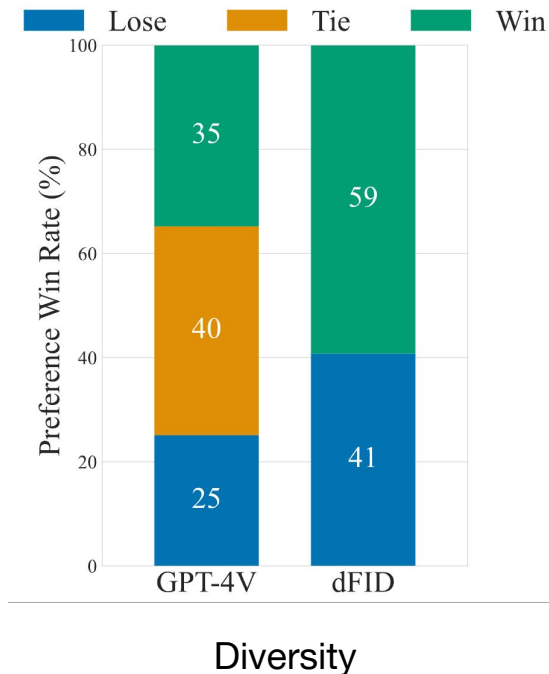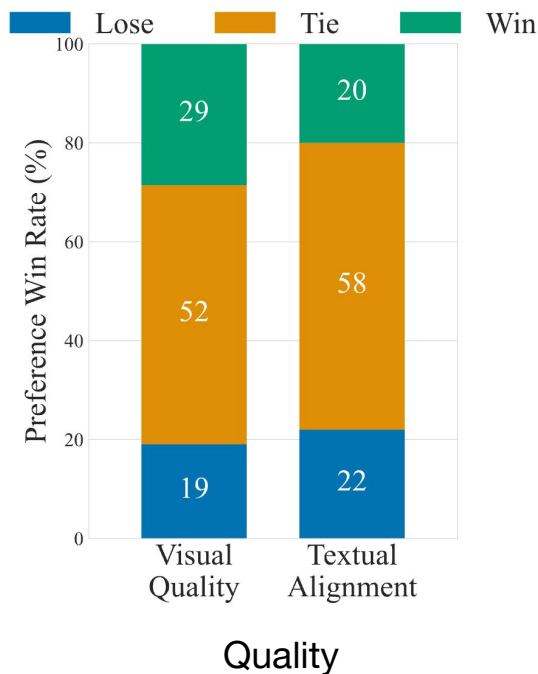
# Human Evaluation



Human prefer Stylus **2x** more than the base model.

# CLIP/FID Pareto Curve



Stylus improves **pareto efficiency** for the CLIP/FID curve.

# Vision Language Model (VLM) as a Judge



Quality

Diversity

VLM believes Stylus generates **higher quality** and **diverse** images.

# Ablations

|  | 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 |

|  | CLIP ($\triangle$) | FID ($\triangle$) |
|---|---|---|
| No-Refiner | 24.91 (-2.31) | 24.26 (+0.30) |
| Gemini-Ultra Refiner | 27.25 (+0.03) | 22.05 (-1.91) |
| GPT-4o Refiner | 28.04 (+0.82) | 21.96 (-2.00) |
| SD v1.5 | 27.22 | 23.96 |

**Retrieval Methods.**
Stylus (with composer) is necessary.
Retriever-only (RAG) hurts end2end
performance.

**Refiner.**
Better adapter descriptions lead to much
better performance.

# Conclusion

Stylus automatically **selects** and **composes** adapters given a user prompt, improving *image fidelity* and *textual alignment*.

Our Contributions
- **StylusDocsv2**: a **75K** entry dataset for adapter descriptions
- Stylus's composer is the first to use an LLM to improve retrieval methods, outperforming rerankers.
- Among the first to employ VLM as a judge for evaluation.