# Manipulating Embeddings of Stable Diffusion Prompts to Control Image Compositionality

Bachelor Thesis

Dinara Imambayeva

Supervised by Niklas Deckers

Leipzig, 18.01.2024

# OVERVIEW

- Motivation
  - Problem Description
- Theoretical Background
- Related Works
- Approach
- Experimental Setup
- Results
- Further Approaches
- Summary
- Outlook

# MOTIVATION

## Image generation with Stable Diffusion

- User expectation: the generated image corresponds to the prompt

A picture of a white cat and a black cat → Stable Diffusion →

# MOTIVATION

## Image generation with Stable Diffusion
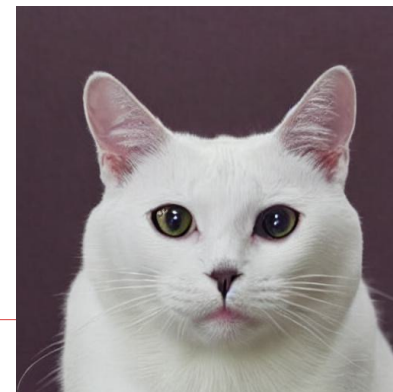
- Reality:

A picture of a white cat and a black cat → Stable Diffusion →

failed composition of multiple objects

failed attribute binding

missing objects

# MOTIVATION: PROBLEM DESCRIPTION

- Often failed image compositionality:

  o Color leakage

  o Incorrect number of objects

  o Missing objects

  o Failed attribute binding

  o Failed composition of multiple objects

- Problem: generated images do not satisfy the user

# MOTIVATION: PROBLEM DESCRIPTION

- Possible approaches:
  - Trying different seeds



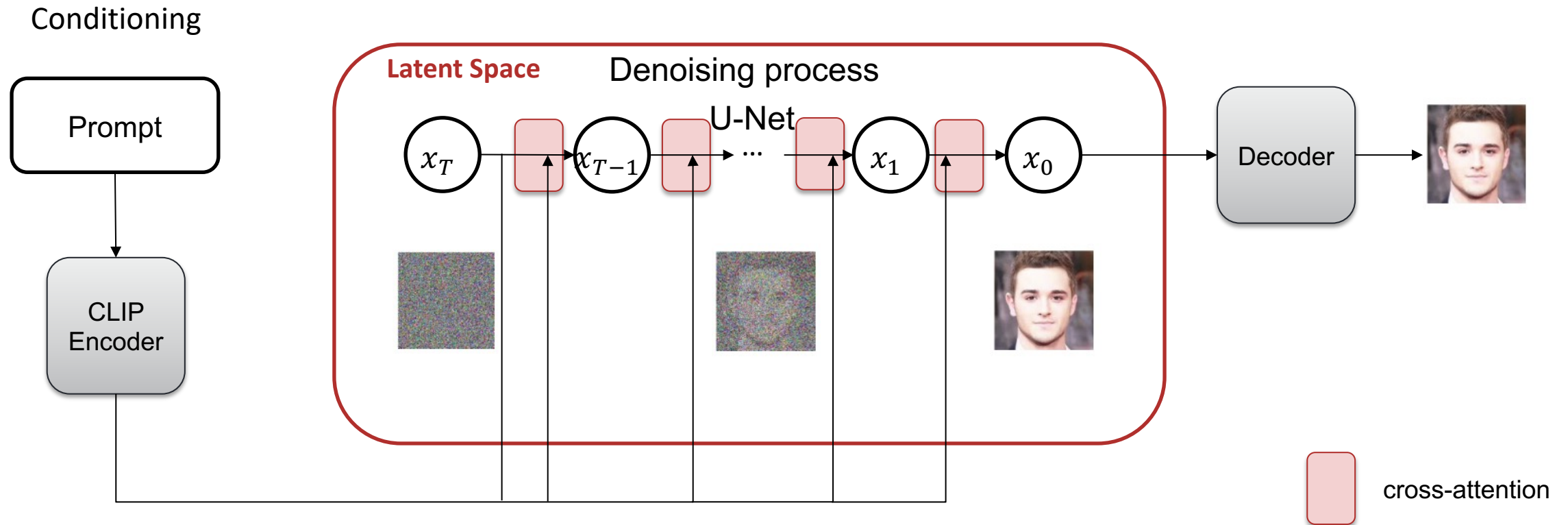seed = 154      seed = 510321      seed = 99

# MOTIVATION: PROBLEM DESCRIPTION

- Possible approaches:
  - Prompt engineering



- Requires sometimes many attempts
- Still not desired output
- No user control, user frustration
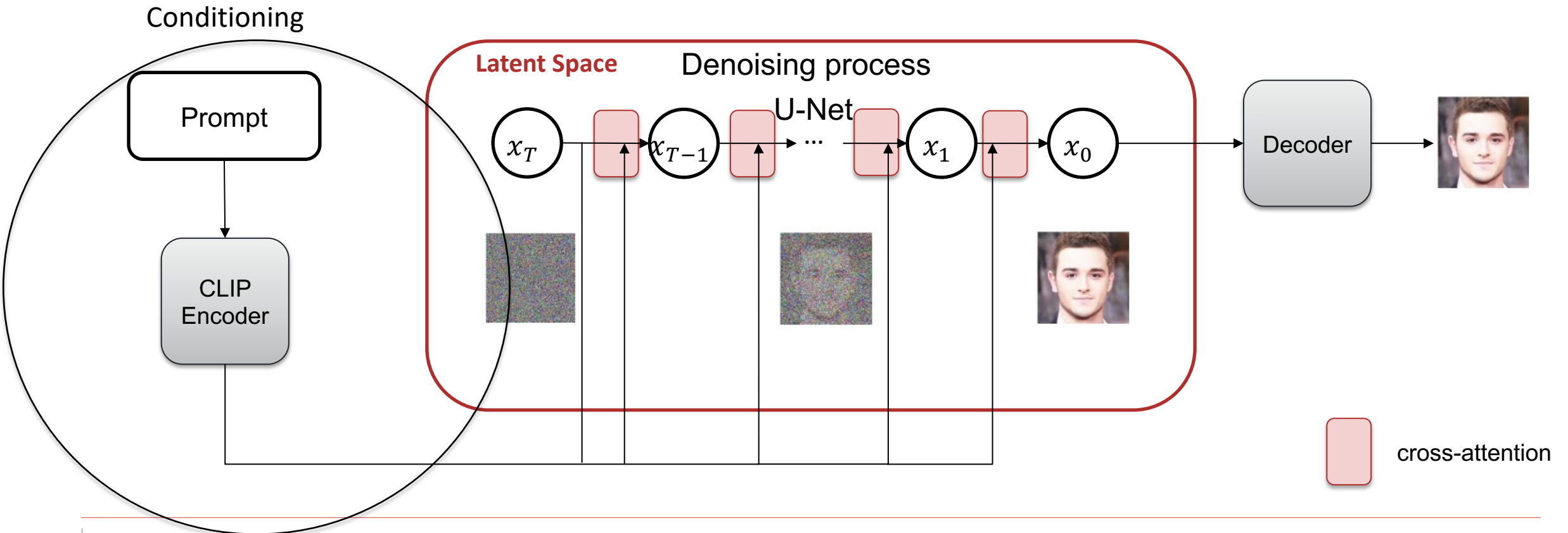
# THEORETICAL BACKGROUND
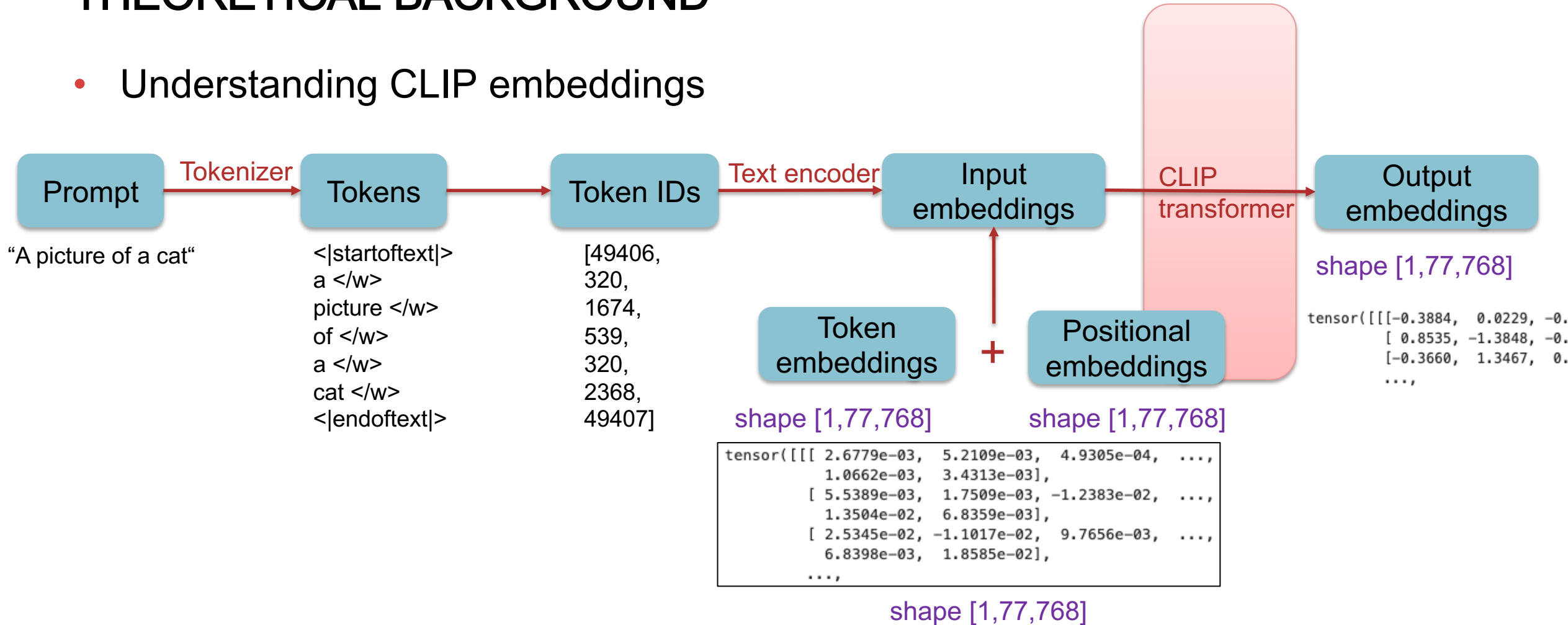
- Image generation with Stable Diffusion

# THEORETICAL BACKGROUND

- Image generation with Stable Diffusion
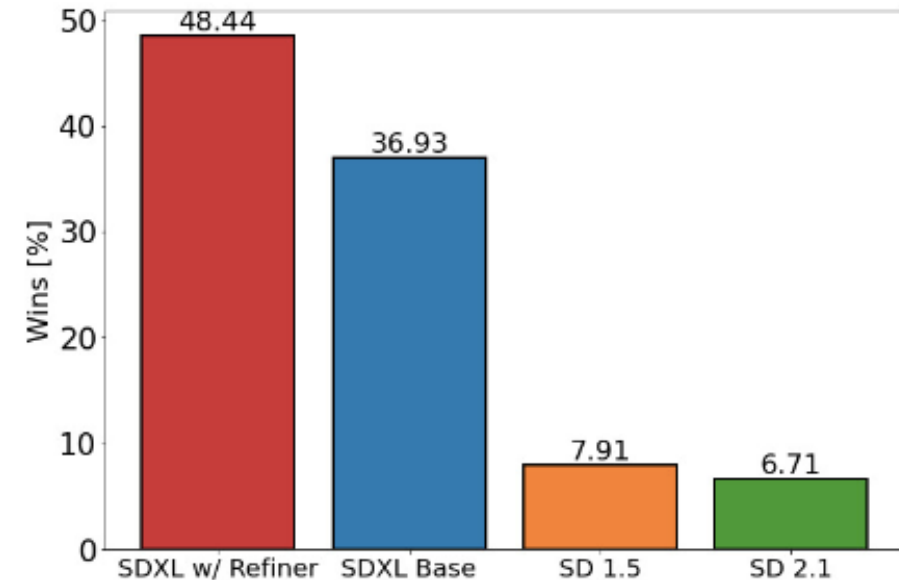
# THEORETICAL BACKGROUND

- Understanding CLIP embeddings

**Prompt** → (Tokenizer) → **Tokens** → **Token IDs** → (Text encoder) → **Input embeddings** → (CLIP transformer) → **Output embeddings**

"A picture of a cat"

**Tokens**
```
<|startoftext|>
a </w>
picture </w>
of </w>
a </w>
cat </w>
<|endoftext|>
```

**Token IDs**
```
[49406,
320,
1674,
539,
320,
2368,
49407]
```

**Output embeddings** shape [1,77,768]

```
tensor([[[−0.3884, 0.0229, −0.
       [ 0.8535, −1.3848, −0.
       [−0.3660, 1.3467, 0.
       ...,
```

**Token embeddings**     +     **Positional embeddings**

shape [1,77,768]          shape [1,77,768]

```
tensor([[[ 2.6779e−03, 5.2109e−03, 4.9305e−04, ...,
         1.0662e−03, 3.4313e−03],
        [ 5.5389e−03, 1.7509e−03, −1.2383e−02, ...,
         1.3504e−02, 6.8359e−03],
        [ 2.5345e−02, −1.1017e−02, 9.7656e−03, ...,
         6.8398e−03, 1.8585e−02],
        ...,
```

shape [1,77,768]

# THEORETICAL BACKGROUND

## Stable Diffusion XL

- Released in July 2023

- High-resolution image synthesis

- Higher fidelity

- Modified architecture:
  - Larger U-Net
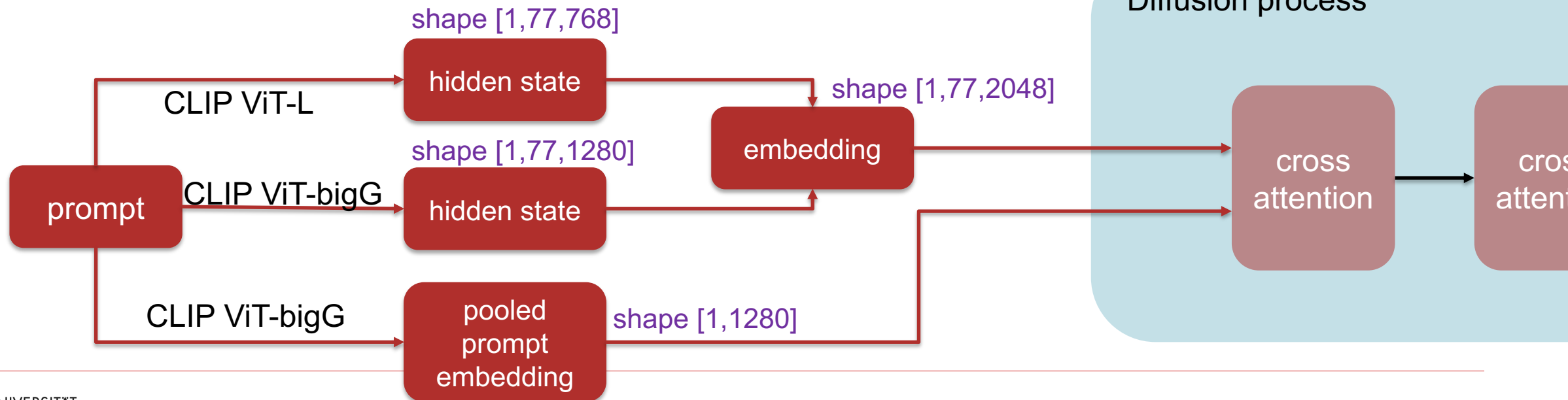  - 2 CLIP text encoders

- Additional refinement model



Comparison of user preferencies

Source: Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., ... & Rombach, R. (2023). Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952.*

# THEORETICAL BACKGROUND

## Stable Diffusion XL

- 2 CLIP text encoders for high-resolution image synthesis
  - CLIP ViT-L
  - CLIP ViT-bigG

# THEORETICAL BACKGROUND

## Stable Diffusion XL

A picture of a white
cat and a black cat

Stable
Diffusion XL

# THEORETICAL BACKGROUND

## Stable Diffusion XL



`A picture of a white cat and a black cat` → Stable Diffusion XL →

**Still failed compositionality!**

UNIVERSITÄT
LEIPZIG

# RELATED WORKS: COMPOSABLE DIFFUSION

- Compositional image generation
- Each concept is generated separately
- 2 compositional operators: Conjunction and Negation



(a) Composing Language Descriptions (Composed Stable Diffusion)

"A photo of cherry blossom trees" AND "Sun dog" AND "Green grass"

"A church" AND "Lightning in the background" AND "A beautiful pink sky"

"A stone castle surrounded by lakes and trees," AND "Black and white"

"A stone castle surrounded by lakes and trees," AND (NOT "Black and white")

Source:Liu, N., Li, S., Du, Y., Torralba, A., & Tenenbaum, J. B. (2022, October). Compositional visual generation with composable diffusion models. In *European Conference on Computer Vision*(pp. 423-439). Cham: Springer Nature Switzerland.

UNIVERSITÄT
LEIPZIG

# RELATED WORKS: STRUCTURED DIFFUSION

- Training-free guidance for compositional text-to-image synthesis

- **Idea**: separate encoding of noun phrases combined with manipulation in cross-attention layers

- 5-8% advantage compared to Stable Diffusion



A red car and a white sheep.

A brown bench sits in front of an old white building

A blue backpack and a brown elephant

Source: Feng, W., He, X., Fu, T. J., Jampani, V., Akula, A., Narayana, P., ... & Wang, W. Y. (2022). Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*.

# RELATED WORKS: STRUCTURED DIFFUSION

Source: Feng, W., He, X., Fu, T. J., Jampani, V., Akula, A., Narayana, P., ... & Wang, W. Y. (2022). Training-free structured diffusion guidance for compositional text-to-image synthesis. *arXiv preprint arXiv:2212.05032*.

# RELATED WORKS: POOLING

- Prompt manipulation by pooling the noun phrases
- Idea: additional attribute binding through pooling



```
green cubes in a pink bowl
```
without NP pooling



```
green cubes in a pink bowl
```
with NP pooling

Source: https://colab.research.google.com/drive/1izMKdvBMfThVSRp8Tg4Fc1eiGVP0NKZP?usp=sharing#scrollTo=wa3iaNL3Mq28

# RELATED WORKS: POOLING

- Common technique in NLP
- Conversion of the NP embedding vectors into a single vector
- Pooling:
  - Maximum pooling
  - Mean square pooling
  - Mean pooling

$$[a,b].mean/\sqrt{2}$$

a – vector for `green`
b – vector for `cubes`

```
  *0                    *0
green cubes in a pink bowl


    pool              pool


 green cubes        pink bowl
```

# APPROACH

- Combining Structured Diffusion approach with pooling
- Extract noun phrases from the constituency tree for each prompt
- Embedding manipulation:
  - Separate embedding of NPs
  - Additional pooling of NPs

- Adjusted approach for SDXL:
  - Generating embeddings with both CLIP encoders
  - No manipulation in pooled prompt embeddings

# EXPERIMENTAL SETUP



Dataset → Prompt preprocessing → Noun phrases extraction → Generating embeddings for each noun phrase → Pooling → Inserting noun phrases into a prompt → Generating images → Evaluation

UNIVERSITÄT LEIPZIG

# EXPERIMENTAL SETUP

Not every prompt can demonstrate the problem

## Datasets

- **CC-500 (**Concept Conjuction dataset)
- Prompt format: "a [colorA] [objectA] and a [colorB] [objectB]"

- **ABC-6K** (Attribute Binding Contrast dataset)
- 3200 prompt pairs
- „a kitchen with white appliances and brown cupboards"
- „a kitchen with brown appliances and white cupboards"

# EXPERIMENTAL SETUP

## Technical setup

- Stable Diffusion Model: "runwayml/stable-diffusion-v1-5"
- Stable Diffusion XL Base 0.9 (released July 2023)

- Fixed seed: 0

# RESULTS: CC500

| | SD 1.5 | Embedding NPs | Embedding NPs + Pooling |
|---|---|---|---|
| a yellow bowl and a blue cat |  |  |  |
| a red bowl and a blue cup |  |  |  |

24

# RESULTS: CC500

a yellow car and a red cat

SD XL

Embedding NPs

Embedding NPs + Pooling

# RESULTS: CC500

`a pink cow and a brown sheep`

SD XL

Embedding NPs

Embedding NPs + Pooling

# RESULTS

**Evaluation criteria for manual evaluation on CC-500**

- Adopted from Structured Diffusion

Categories:
1. Zero or one object is depicted
2. Two objects are depicted
3. Two object are depicted with correct colors
4. Image looks natural

# RESULTS

| SD 1.5 | Zero/One object | Two objects | Two objects with correct color | Natural looking image |
|---|---|---|---|---|
| SD | 82% | 12% | 6% | 48% |
| NP embedding | 82% | 10% | 8% | 40% |
| NP embedding + pooling | 94% | 6% | 0% | 34% |

| SD XL | Zero/One object | Two objects | Two objects with correct color | Natural looking image |
|---|---|---|---|---|
| SD | 40% | 42% | 18% | 44% |
| NP embedding | 42% | 44% | 14% | 30% |
| NP embedding + pooling | 40% | 46% | 14% | 66% |

# FURTHER APPROACHES: SECOND APPROACH

**Idea**: Can we subtract/add the color to improve attribute binding?



red cube          -          cube          +          car          =

Generated by Stable Diffusion 1.5, seed = 0

# FURTHER APPROACHES: SECOND APPROACH

## Experiment on SDXL

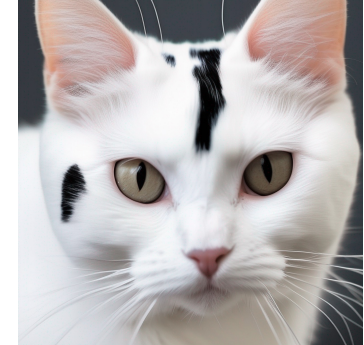`A white cat with black ears and markings`



No modification



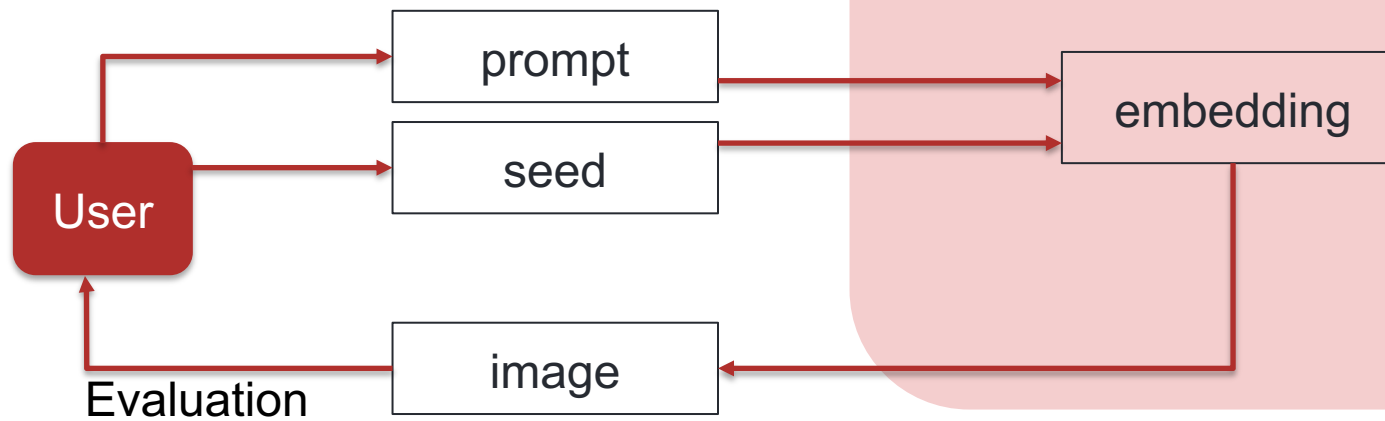NP embedding + pooling



Adding `black` vector to the according NP

# FURTHER APPROACHES: SECOND APPROACH



**Idea**: User interaction interface for image adjustment

# FURTHER APPROACHES: SECOND APPROACH

**Idea**: User interaction interface for image adjustment

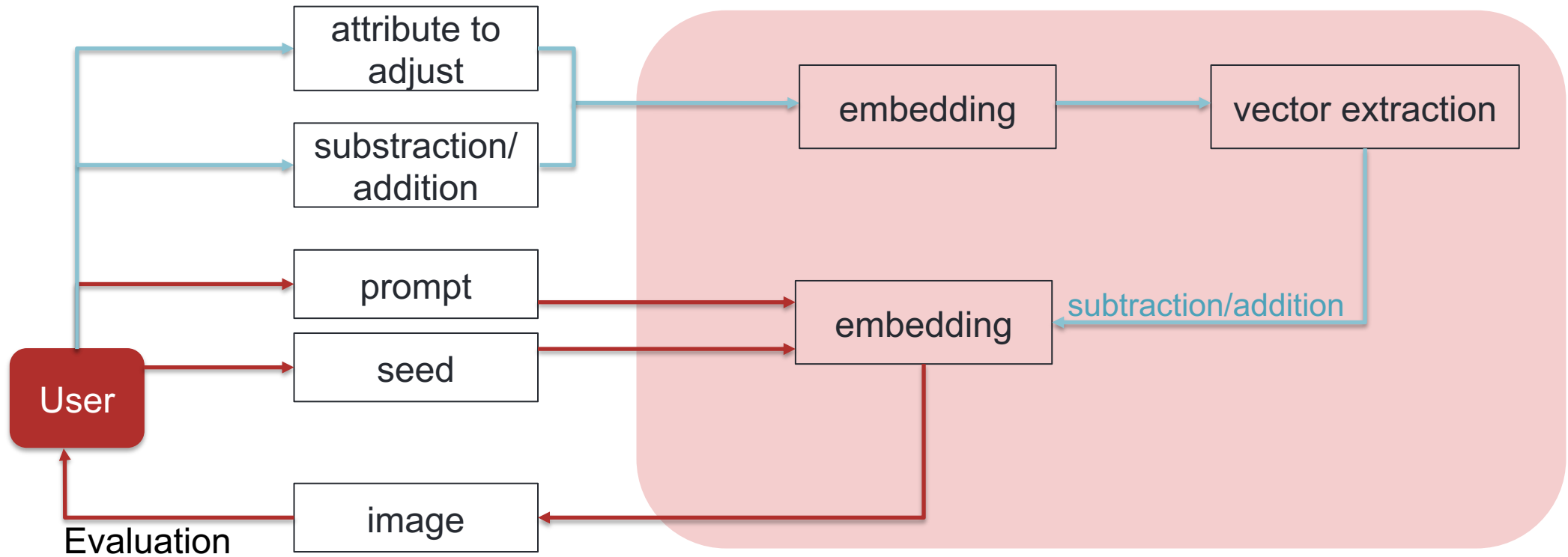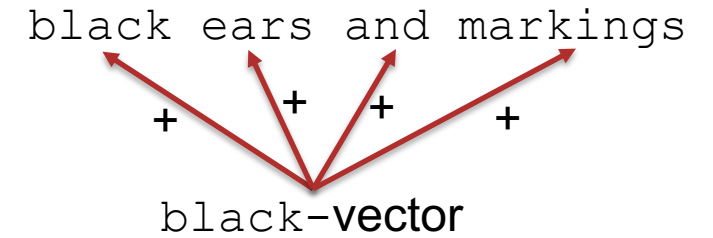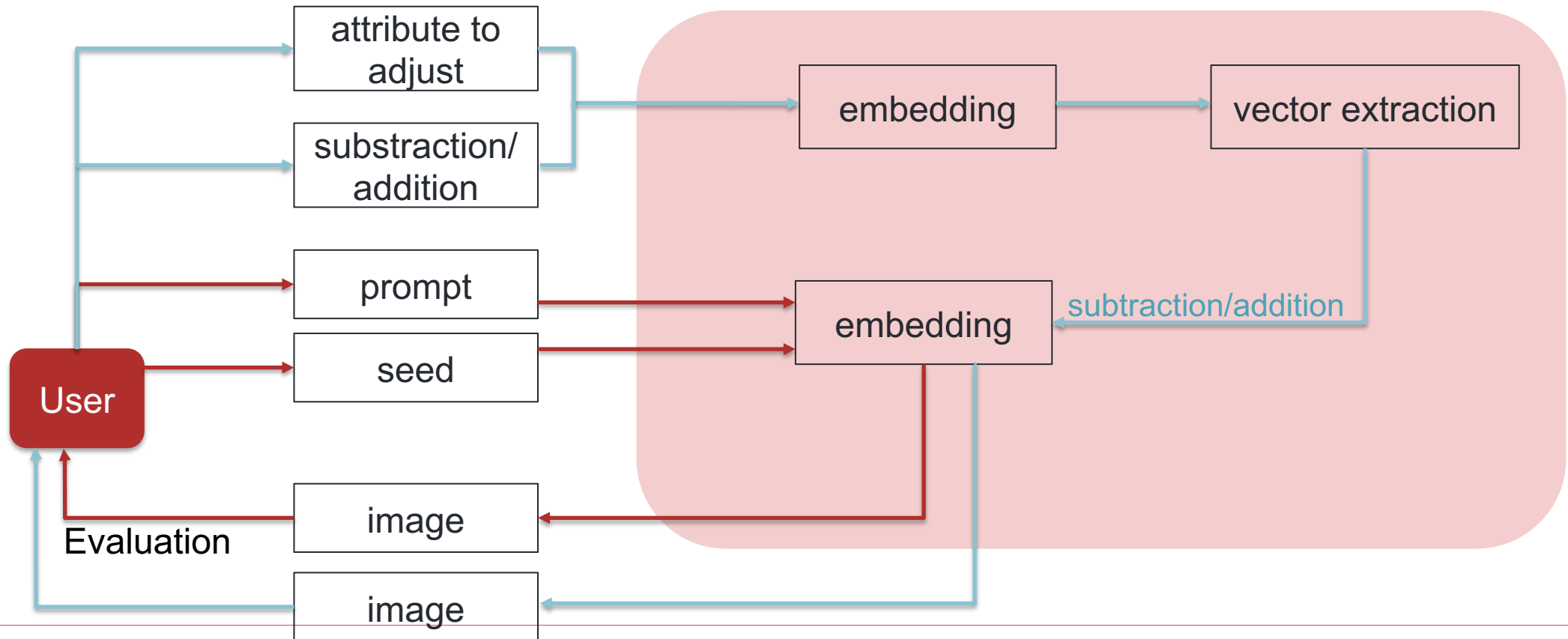# FURTHER APPROACHES: SECOND APPROACH

**Idea**: User interaction interface for image adjustment

`black` → CLIP → embedding

([[[-0.3884, 0.0229, -0.0522, ..
[0.8535, -1.3848, -0.4604, ..
[-0.3660, 1.3467, 0.8745, ...
...,



attribute to adjust

substraction/ addition

embedding

prompt

seed

User

image

Evaluation
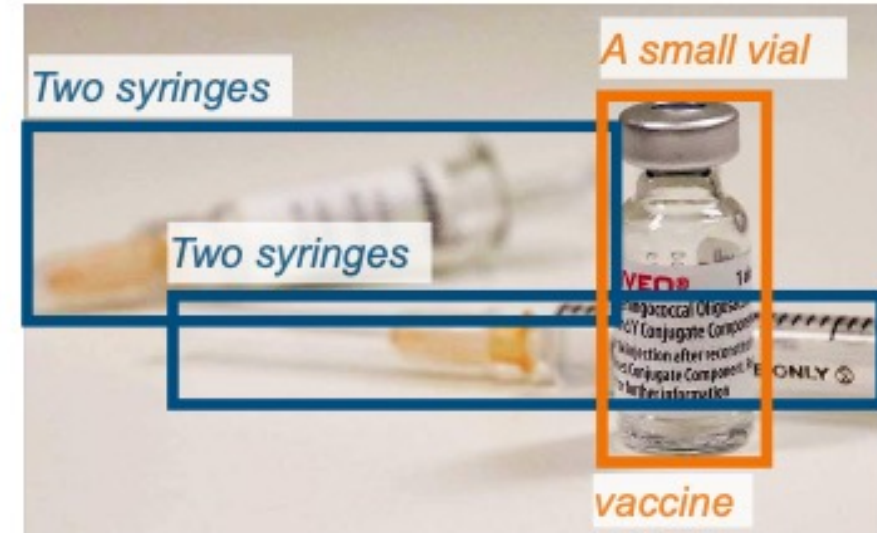
embedding

# FURTHER APPROACHES: SECOND APPROACH



**Idea**: User interaction interface for image adjustment

# FURTHER APPROACHES: EVALUATION

Automatic evaluation with GLIP

- Framework for object detection and phrase grounding

- Phrase grounding – identifying the correspondence between phrases in a prompt and objects in an image

- Object detection as a grounding task



Source: Li, L. H., Zhang, P., Zhang, H., Yang, J., Li, C., Zhong, Y., ... & Gao, J. (2022). Grounded language-image pre-training. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 10965-10975).

# SUMMARY

- Failed compositionality problem limits user control over output
- Improving image compositionality approaches needed
- My approach:
  - Separate embedding of NPs
  - Additional pooling of NPs

- Limitation: architecture-dependent approach

# OUTLOOK

**Current Objectives**:

- Further image generation and evaluation of results (CC-500 and ABC-6K)
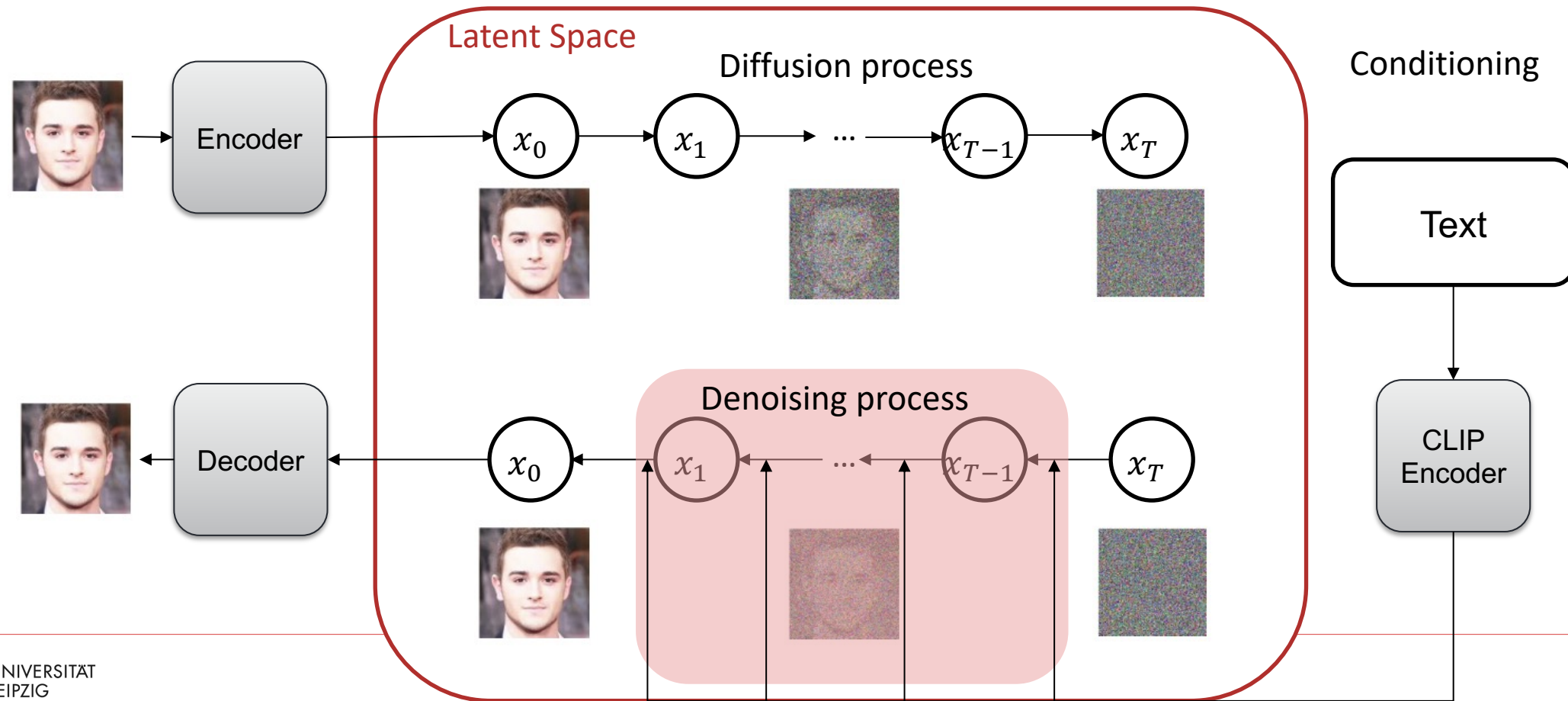- Implementation of automatic evaluation with GLIP

**Future Objectives (beyond thesis)**:

- Refinement of the second approach
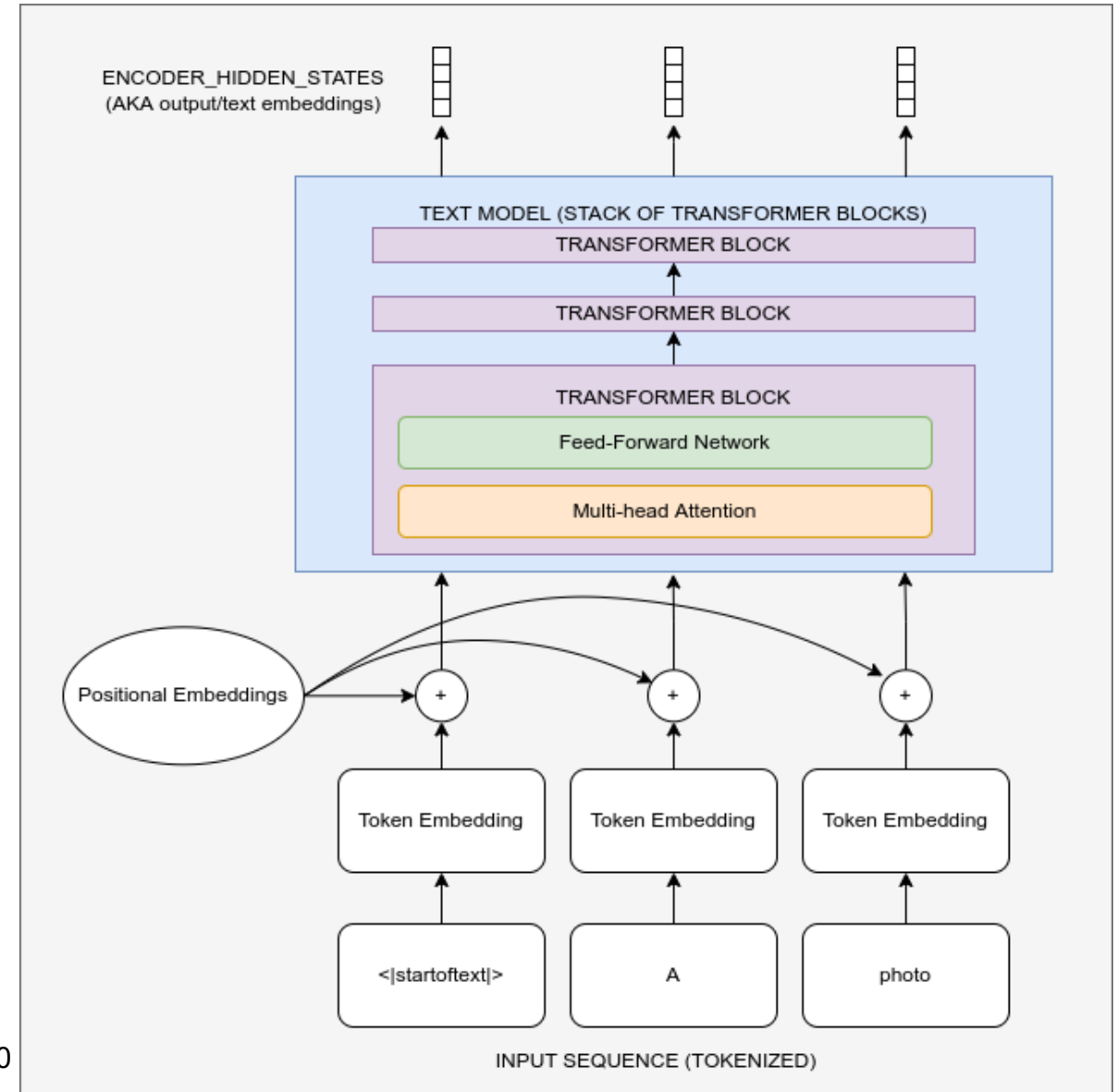- Implementation of a user interface

# BACKUP SLIDES

# TECHNICAL BACKGROUND

– Image generation with Stable Diffusion

# CLIP

- Prompt-embedding pipeline
- CLIP encoder has transformer architecture:
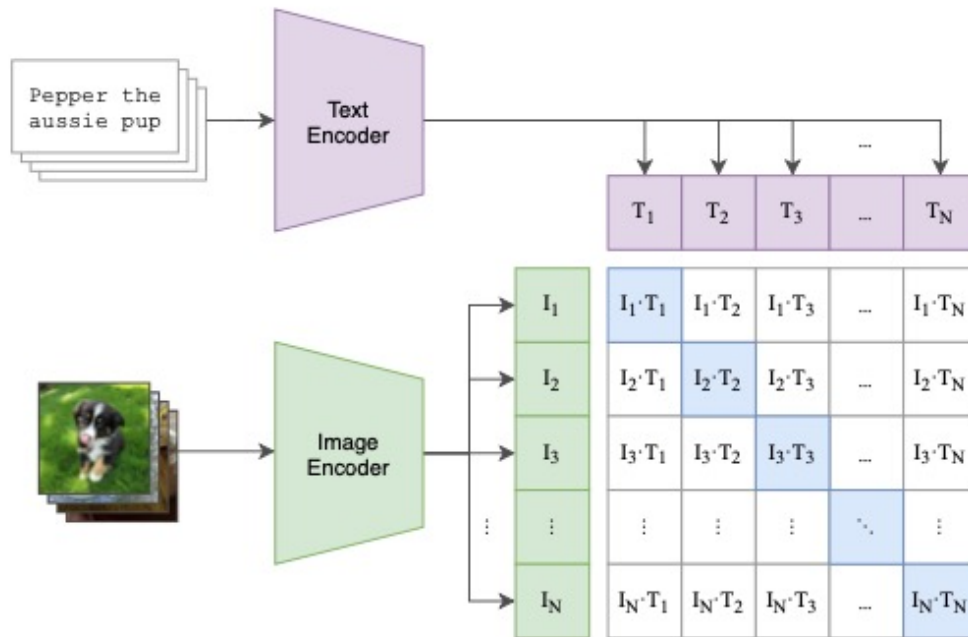  - 12 layers
  - 8 attention-heads

# CLIP

- CLIP text encoder = transformer
- Training in multimodal space

Source: Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021, July). Learning transferable visual models from natural language supervision. In *International conference on machine learning* (pp. 8748-8763). PMLR.

# STABLE DIFFUSION XL

## Two-stage-pipeline



## Comparison between SD 1.5 and SDXL

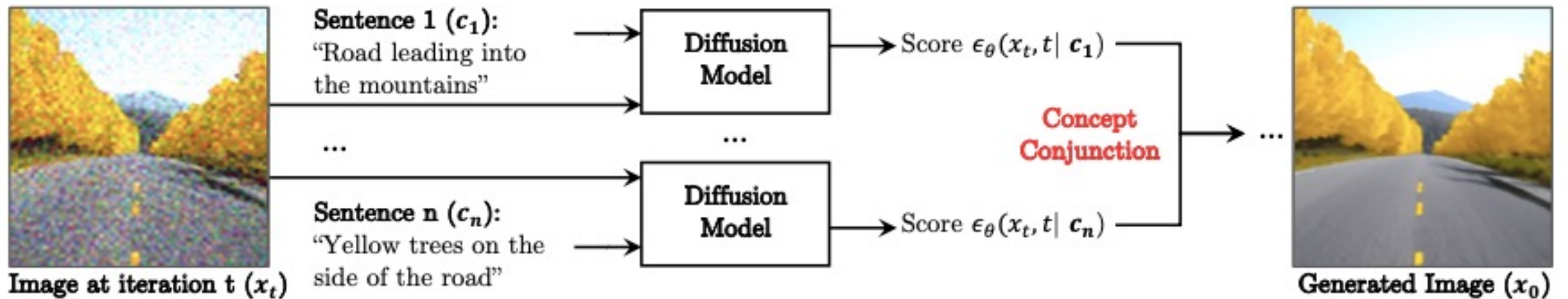| Model | SDXL | SD 1.4/1.5 |
|---|---|---|
| # of UNet params | 2.6B | 860M |
| Transformer blocks | [0, 2, 10] | [1, 1, 1, 1] |
| Channel mult. | [1, 2, 4] | [1, 2, 4, 4] |
| Text encoder | CLIP ViT-L & OpenCLIP ViT-bigG | CLIP ViT-L |
| Context dim. | 2048 | 768 |
| Pooled text emb. | OpenCLIP ViT-bigG | N/A |

## Comparison of user preferences



Source: Podell, D., English, Z., Lacey, K., Blattmann, A., Dockhorn, T., Müller, J., ... & Rombach, R. (2023). Sdxl: Improving latent diffusion models for high-resolution image synthesis. *arXiv preprint arXiv:2307.01952*.

# COMPOSABLE DIFFUSION

# STRUCTURED DIFFUSION: CROSS-ATTENTION

**Algorithm 1** StructureDiffusion Guidance.

**Require:**

    **Input:** Prompt $\mathcal{P}$, Parser $\xi$, decoder $\psi$, trained diffusion model $\phi$.

    **Output:** Generated image $x$.

1: Retrieve concept set $\mathcal{C} = [c_1, \ldots, c_k]$ by traversing $\xi(\mathcal{P})$;

2: $\mathcal{W}_p \leftarrow \text{CLIP}_{\text{text}}(\mathcal{P}), \; \mathcal{W}_i \leftarrow \text{CLIP}_{\text{text}}(c_i)$;          $i = 1, \ldots, k$

3: **for** $t = T, T-1, \ldots, 1$ **do**

4:     **for** each cross attention layer in $\phi$ **do**

5:         Obtain previous layer's output $\mathcal{X}^t$.

6:         $Q^t \leftarrow f_Q(\mathcal{X}^t), \; K_p \leftarrow f_K(\mathcal{W}_p), \; V_i \leftarrow f_V(\overline{\mathcal{W}_i})$;      $i = \text{p}, 1, \ldots, k$

7:         Obtain attention maps $M^t$ from $Q^t, K_p$;      {Eq. 1}

8:         Obtain $O^t$ from $M^t, \{V_i\}$, and feed to following layers;      {Eq. 4}

9:     **end for**

10: **end for**

11: Feed $z^0$ to decoder $\psi(\cdot)$ to generate x.

# STABLE DIFFUSION XL

Why not modifying projected pooled embedding?

→    No visual improvement

`Red cube in a blue bowl`



Not manipulated
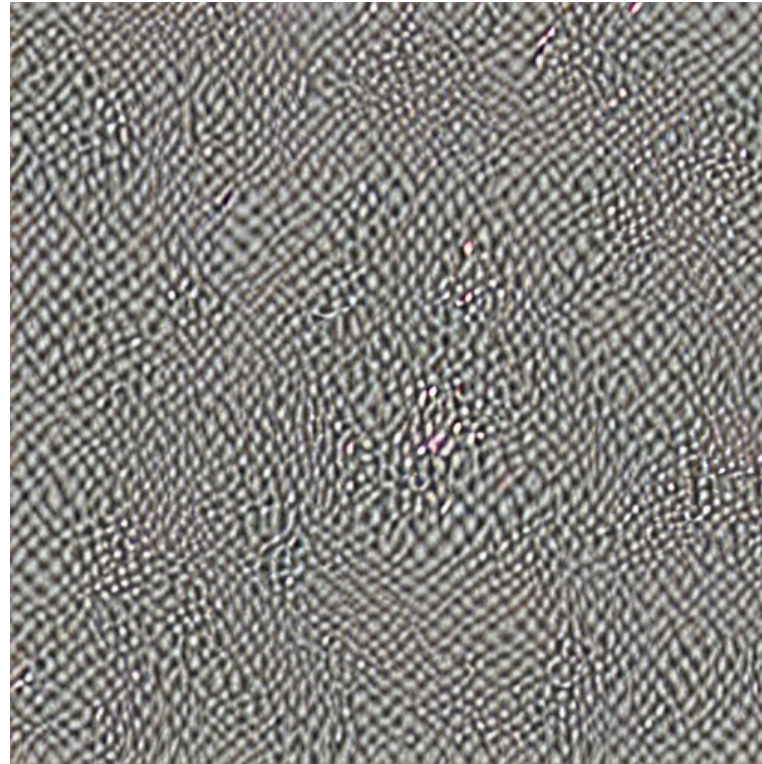


Manipulating both hidden states



Manipulating both hidden states + projected pooled layer

# STABLE DIFFUSION

Why not manipulating starttoken?

- `Red cube in a blue bowl`
- Impact of the start token



Pooling all tokens



Pooling all tokens without start token