



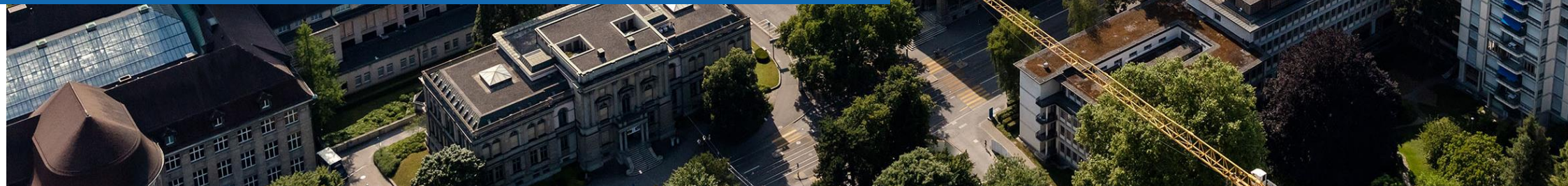
# InstructPix2Pix: Learning to Follow Image Editing Instructions

- University of California, Berkeley

Presenter: Qi Ma

Advisor: Grötschla Florian

Seminar in Deep Neural Networks



Motivation :

# Motivation:

- What kind of text-based image editing functionality do you expect to use?

# Motivation:

- What kind of text-based image editing functionality do you expect to use?



# Motivation:

- What kind of text-based image editing functionality do you expect to use?



Remove the people

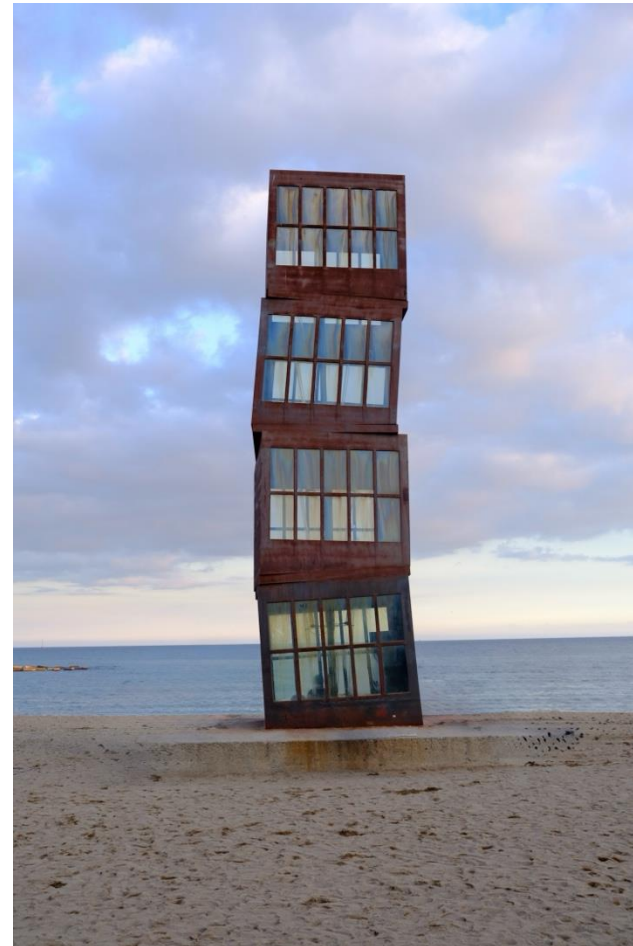


# Motivation:

- What kind of text-based image editing functionality do you expect to use?

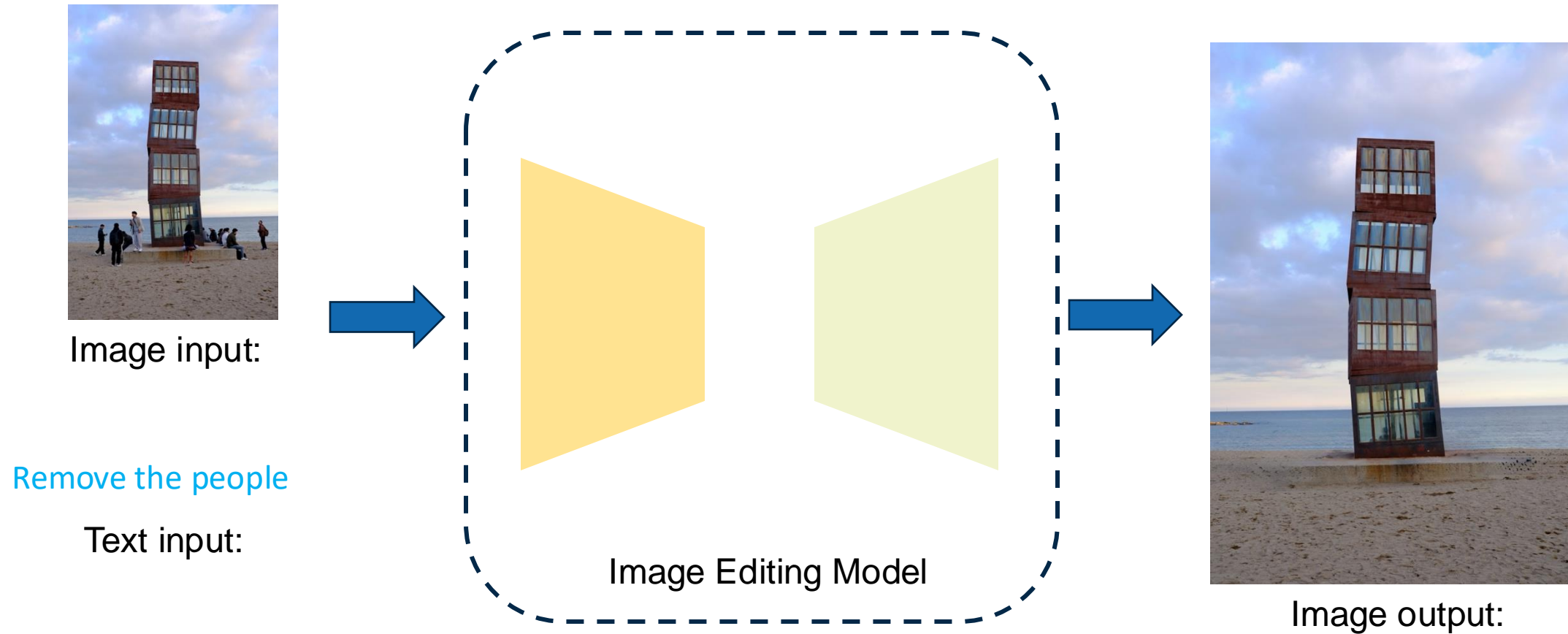


Remove the people



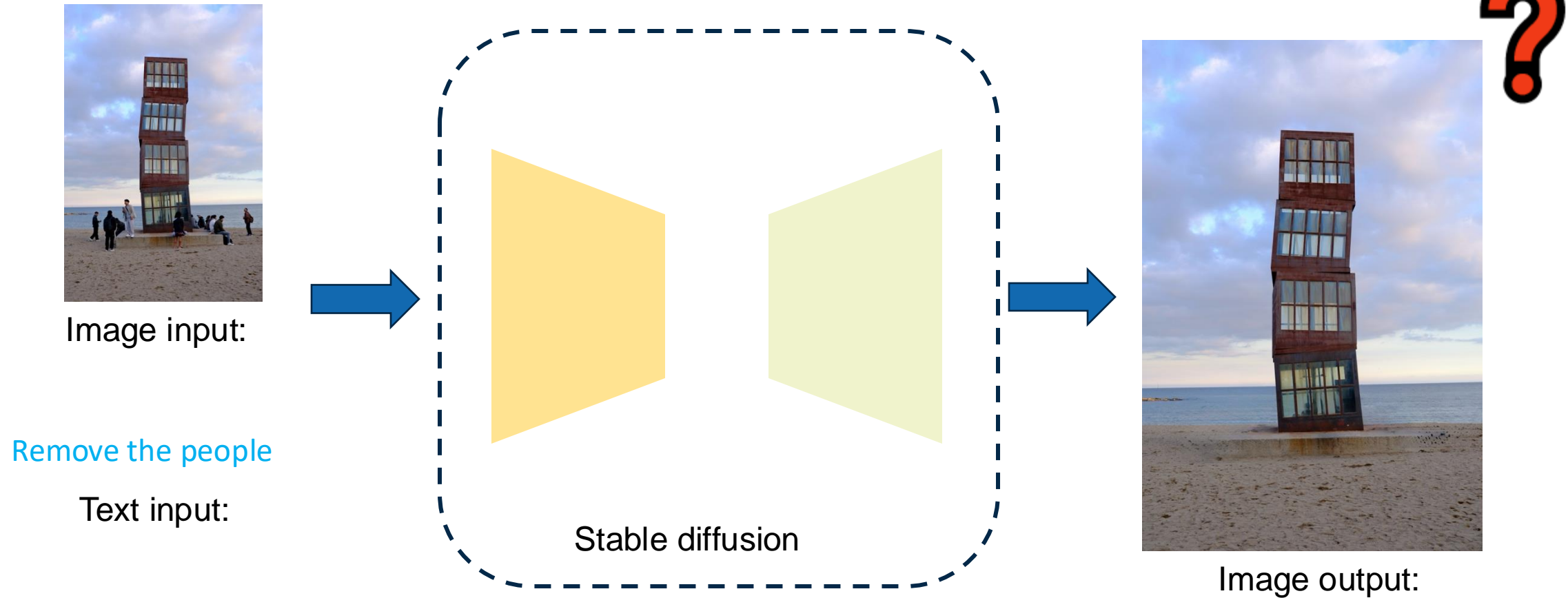
# Motivation:

- What kind of text-based image editing functionality do you expect to use?



# Motivation:

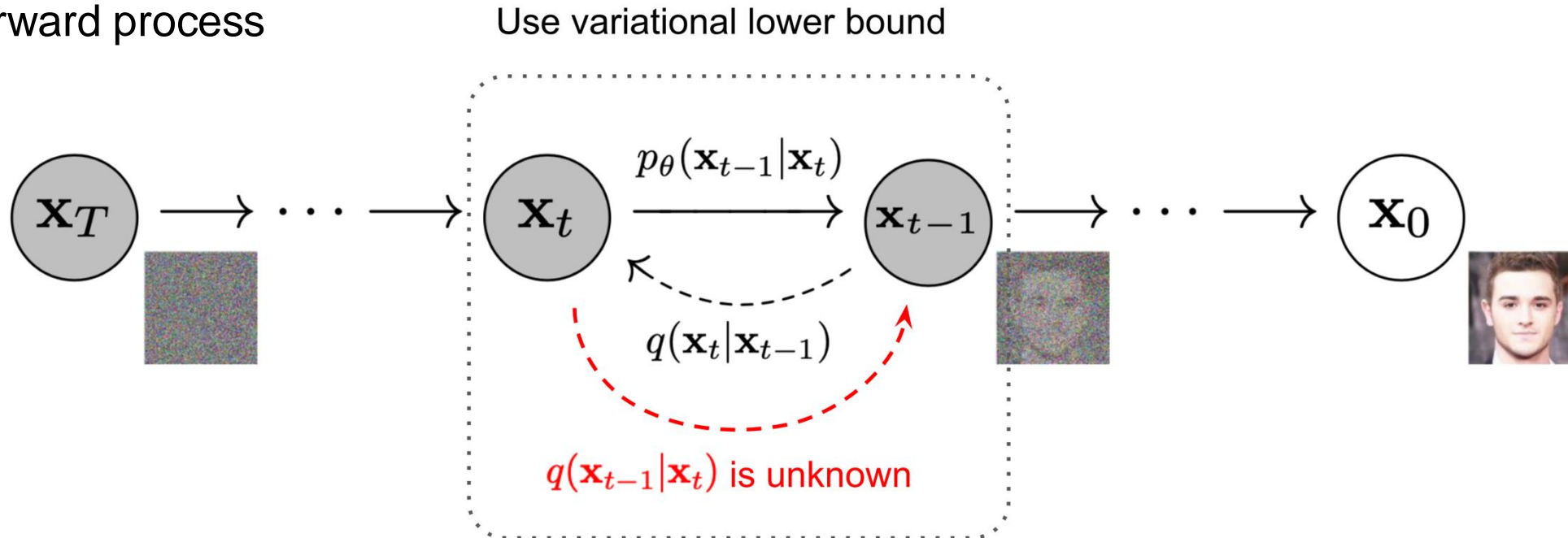
- Sound easy? Why not try stable diffusion<sup>[1]</sup>?





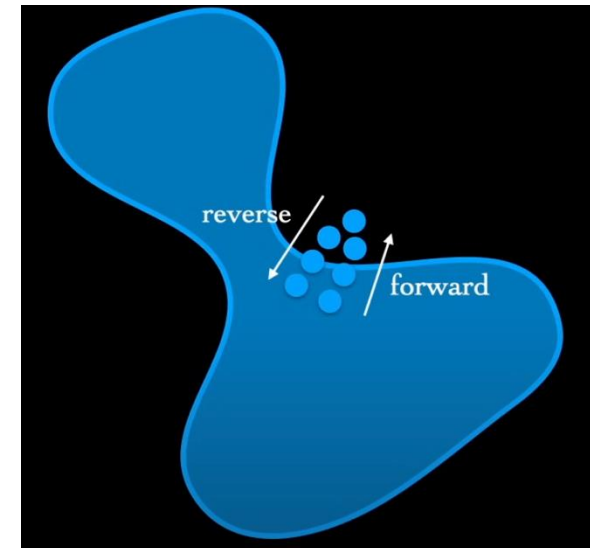
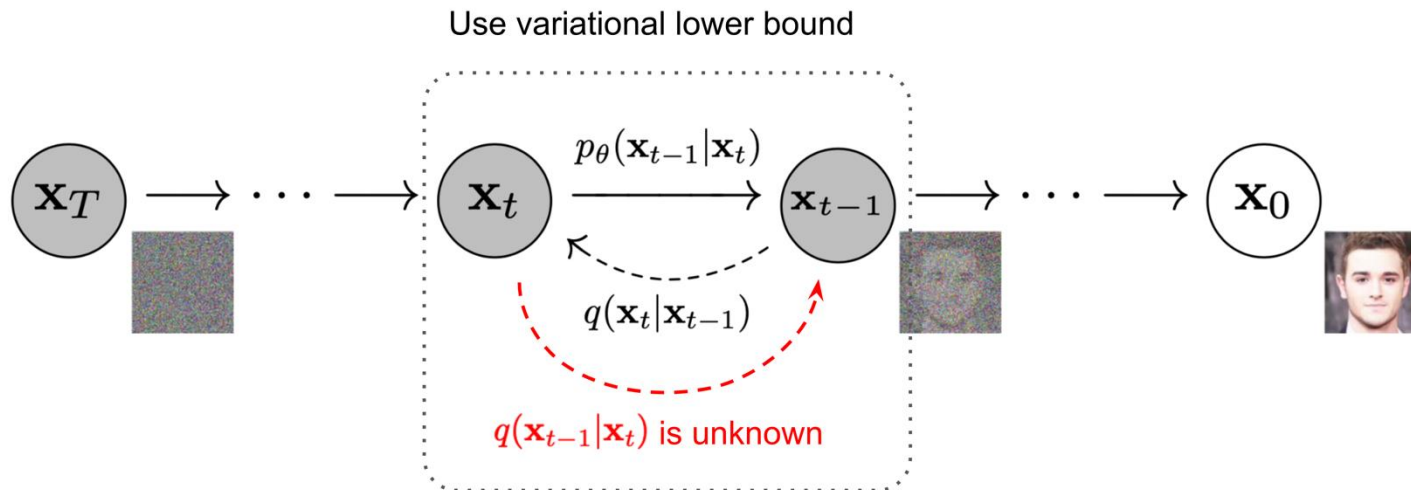
# Motivation:

- What is diffusion model?
  - Forward process



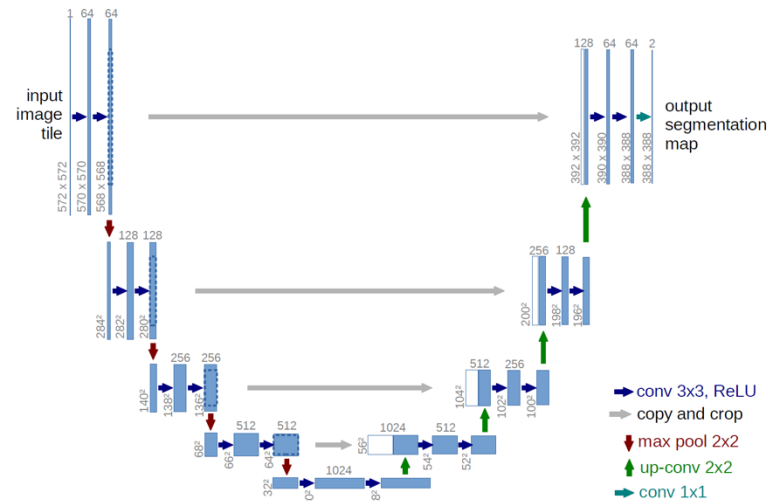
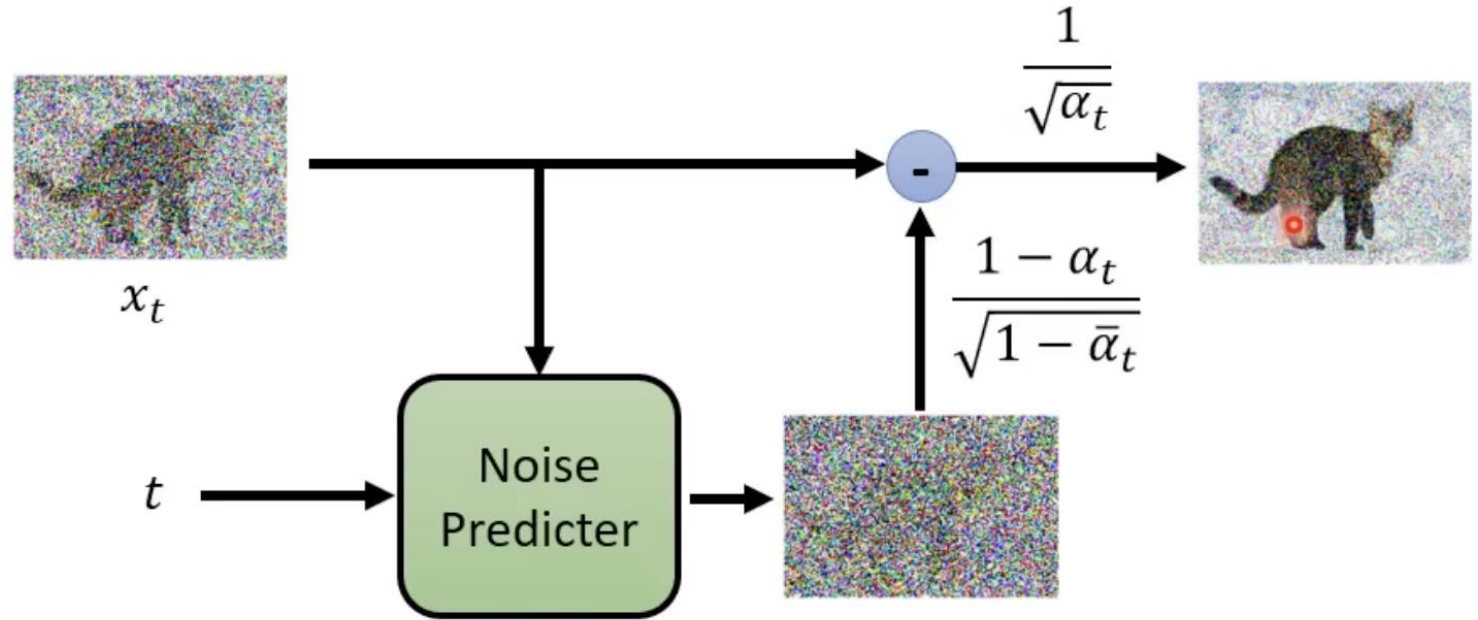
# Motivation:

- What is diffusion model?
  - Reverse process



# Motivation:

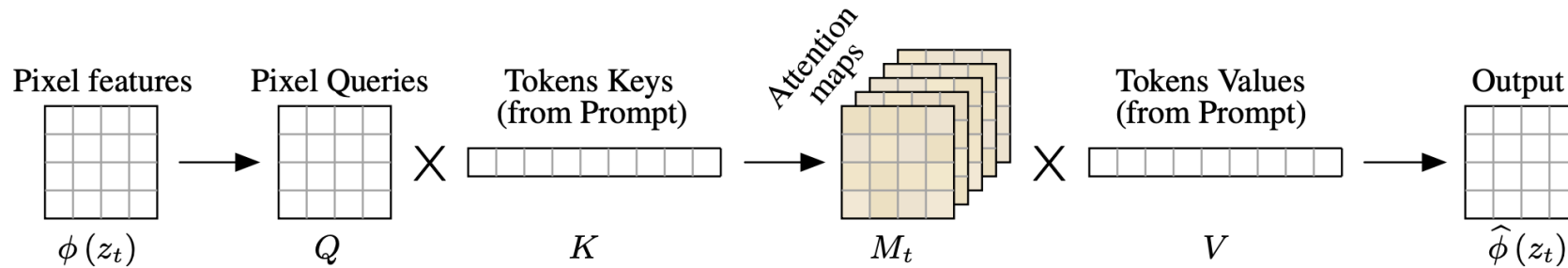
- What is diffusion model?



# Motivation:

- How to condition on text prompt?
  - Cross attention operation in noise predictor

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right) \cdot V,$$



Q: (Token num, Feature\_dim)

M: (Token num x Token num')

K: (Token num', Feature\_dim)

# Motivation:

- Sound easy? Why not try stable diffusion?



Image input:

Remove the people

Text input:

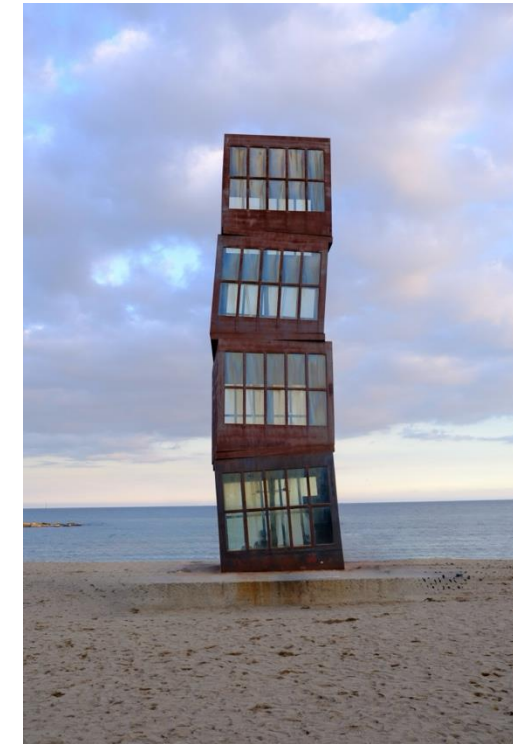
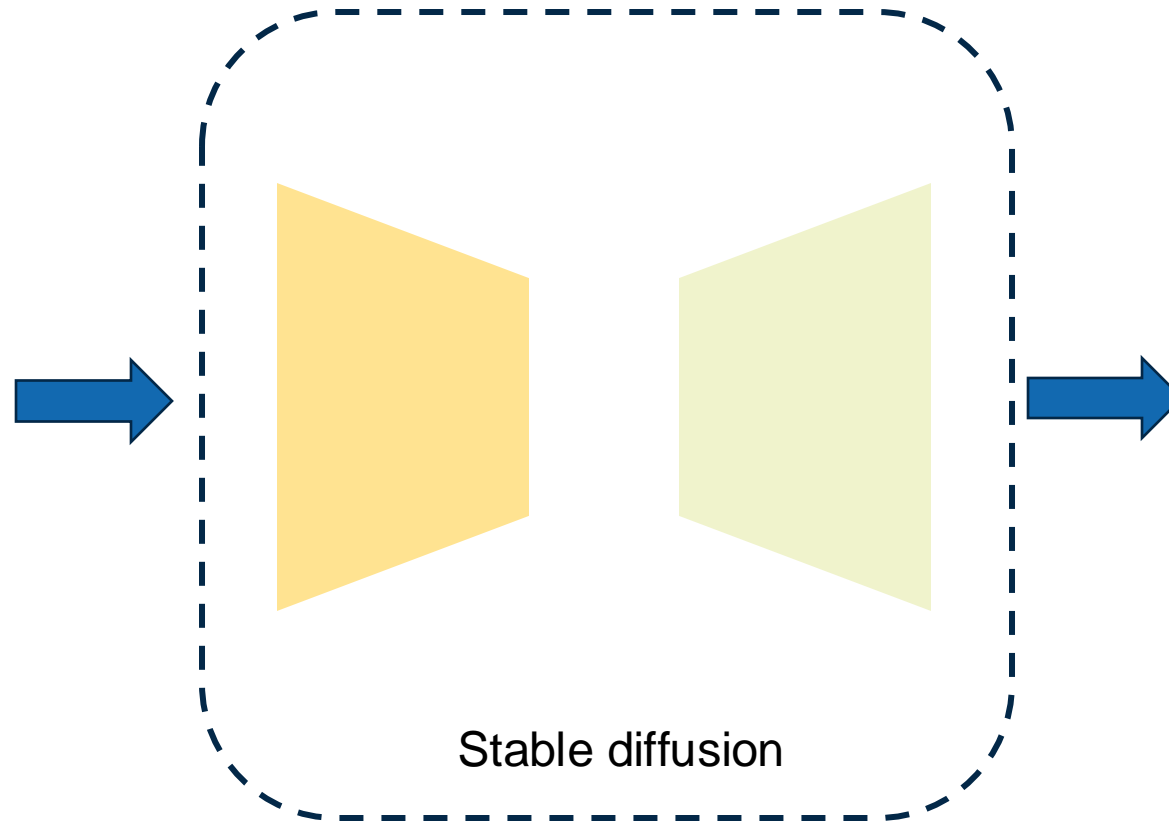
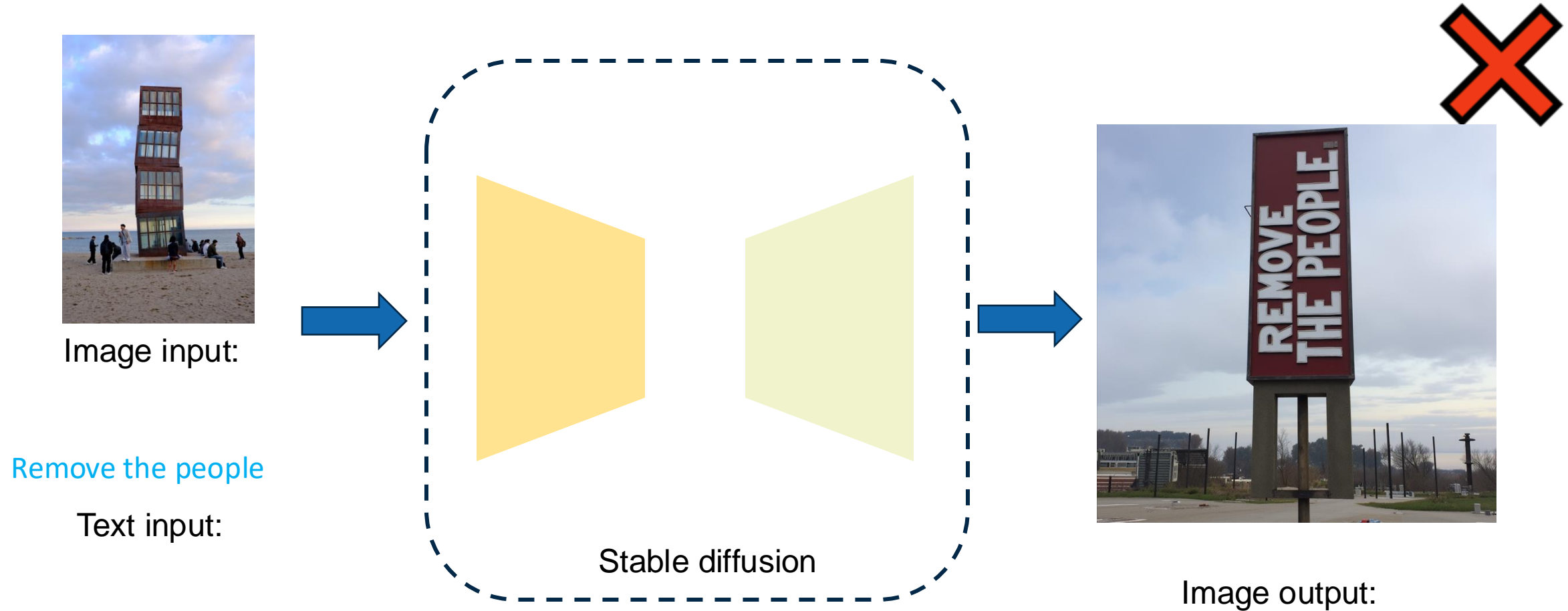


Image output:

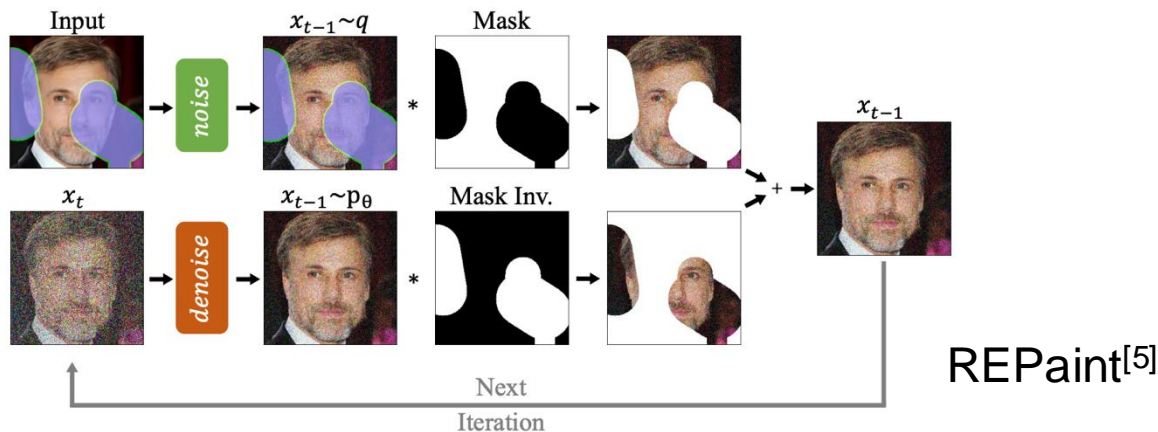
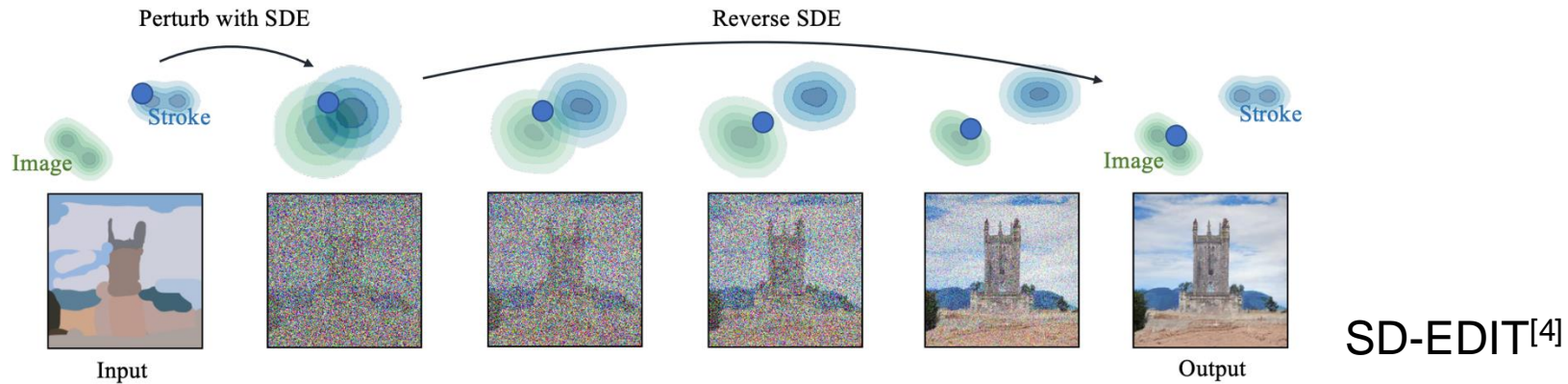
# Motivation:

- Problem: Stable diffusion will also change the content we interested



# Motivation:

- Stable Diffusion combined with other noise/modules to make it work

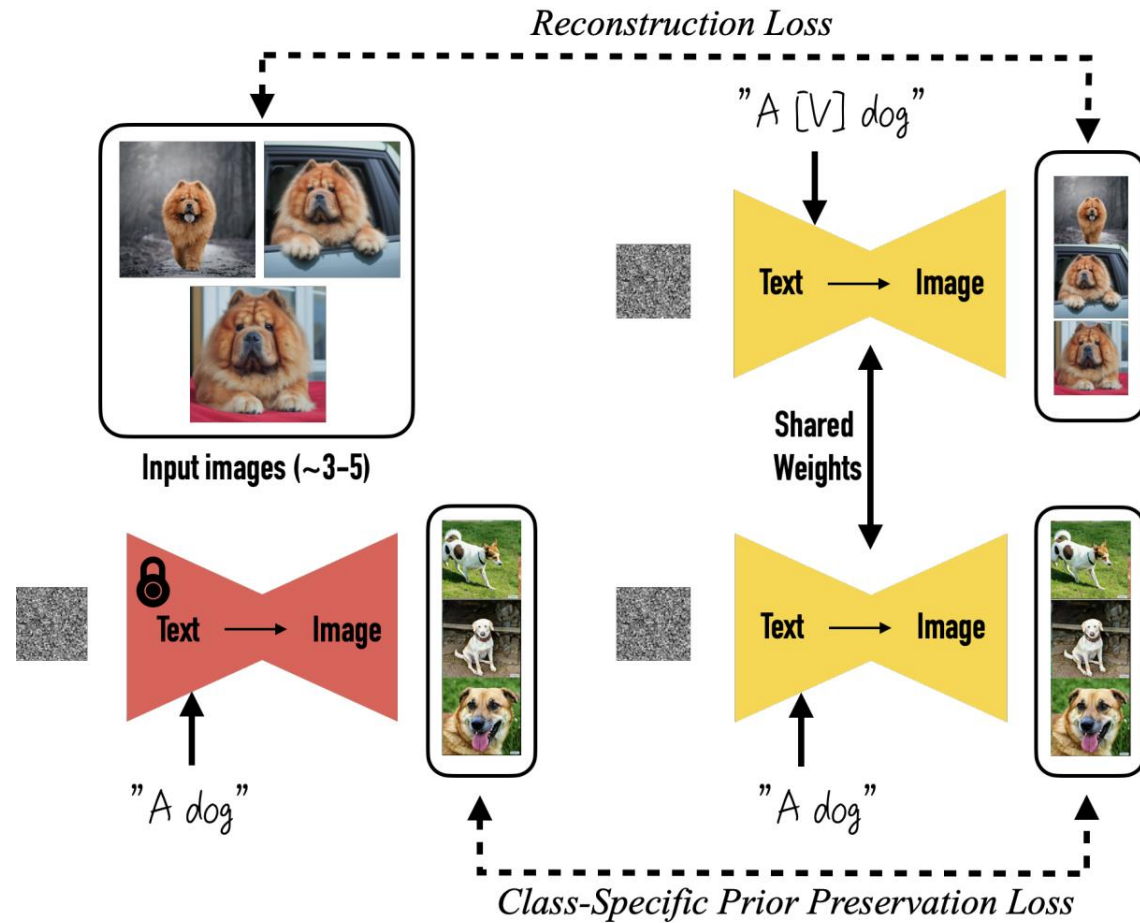


Limitation:

- Poor performance
- Extra work on masking

# Motivation:

- Stable Diffusion combined with other modules to make it work

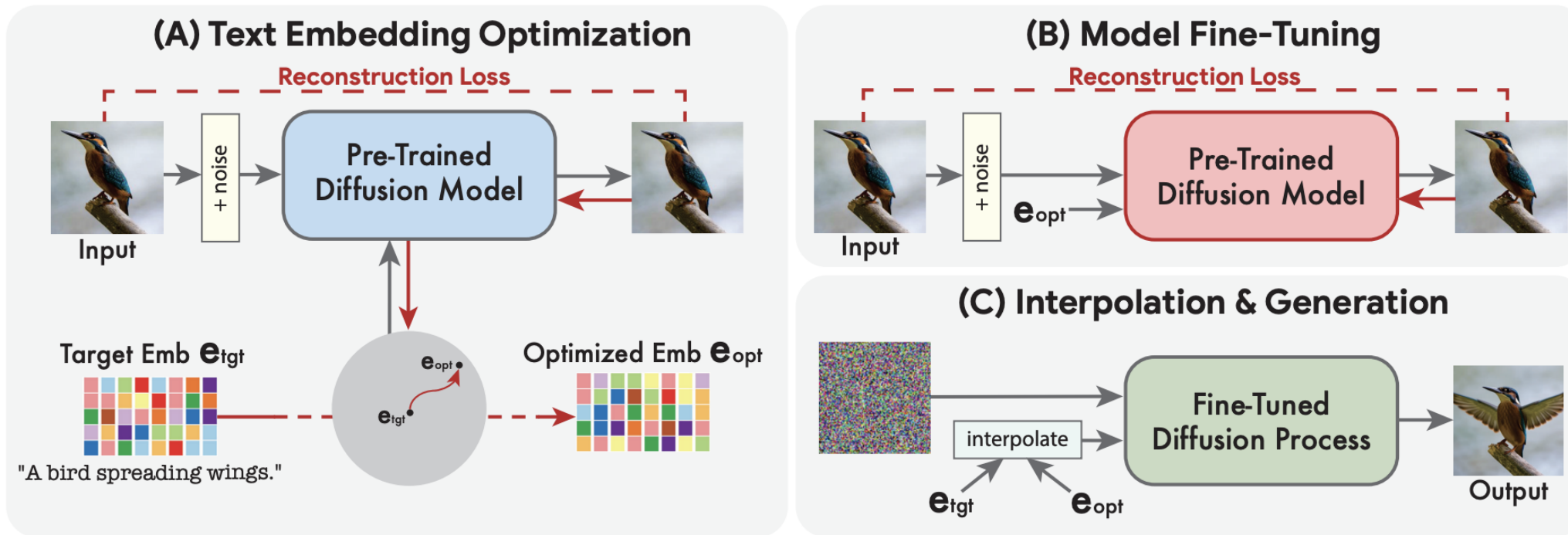


DreamBooth[6]



# Motivation:

- Stable Diffusion combined with other modules to make it work



Imagic<sup>[7]</sup>

Limitation:

- Need to finetune whole stable diffusion model

# Motivation:

- Target: Can we train a diffusion model to follow instructions simple as below?



Image input:

Remove the people

Text input:

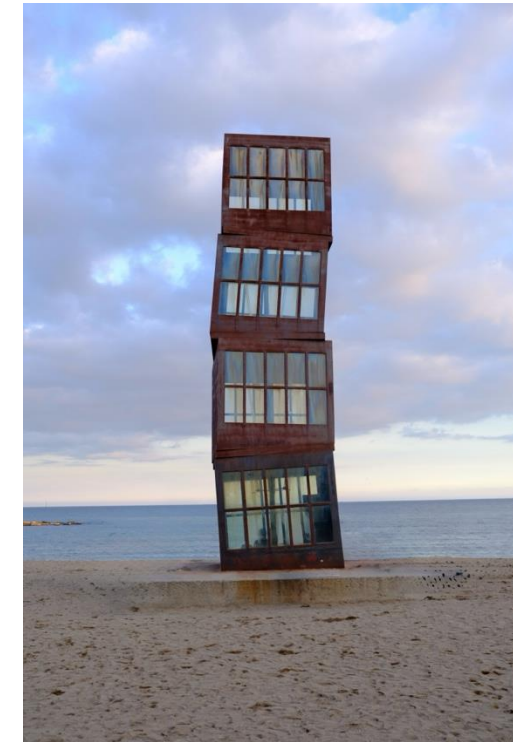
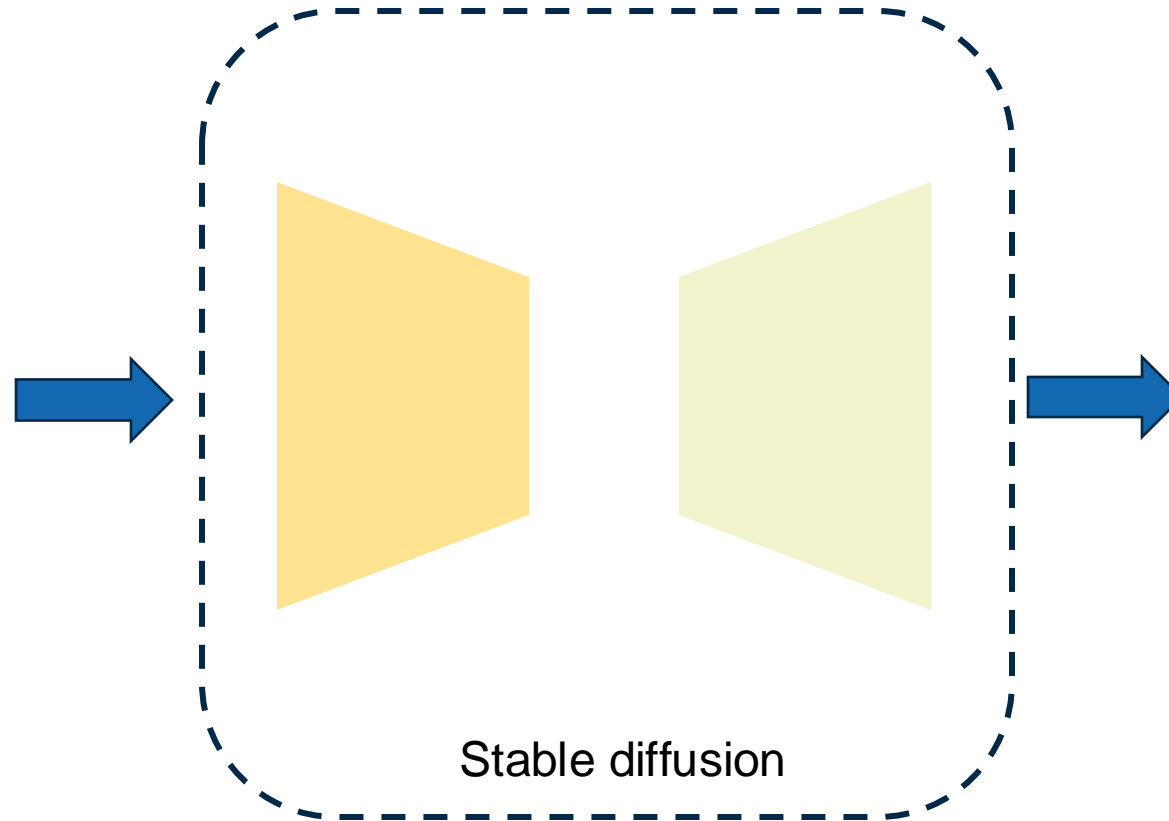
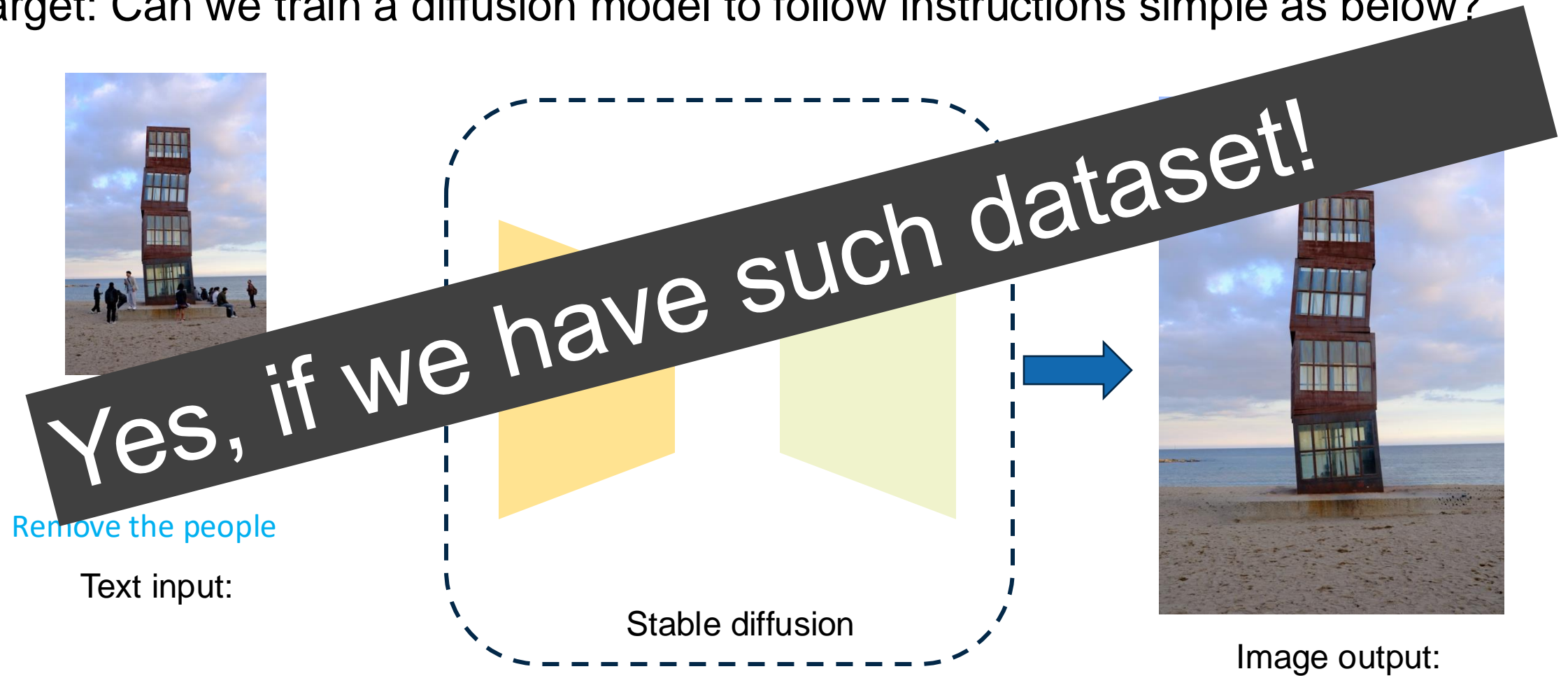


Image output:

# Motivation:

- Target: Can we train a diffusion model to follow instructions simple as below?



Method:

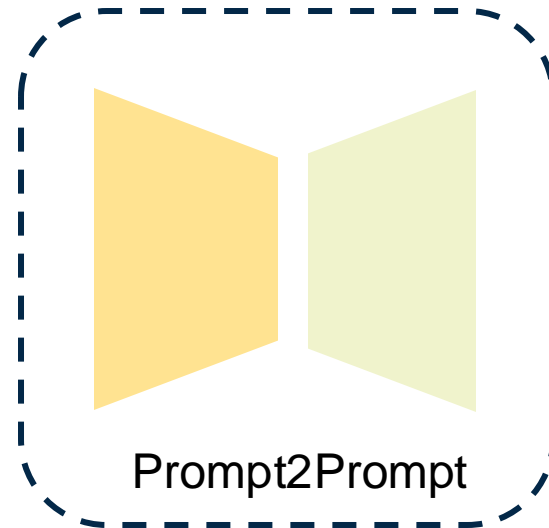
# Prompt-to-Prompt<sup>[8]</sup>

- Given two similar textual descriptions as input, can synthesize image pairs while keep original structure and content

1) Photo of a cat riding on a bicycle

2) Photo of a cat riding on a car

Input text pair:



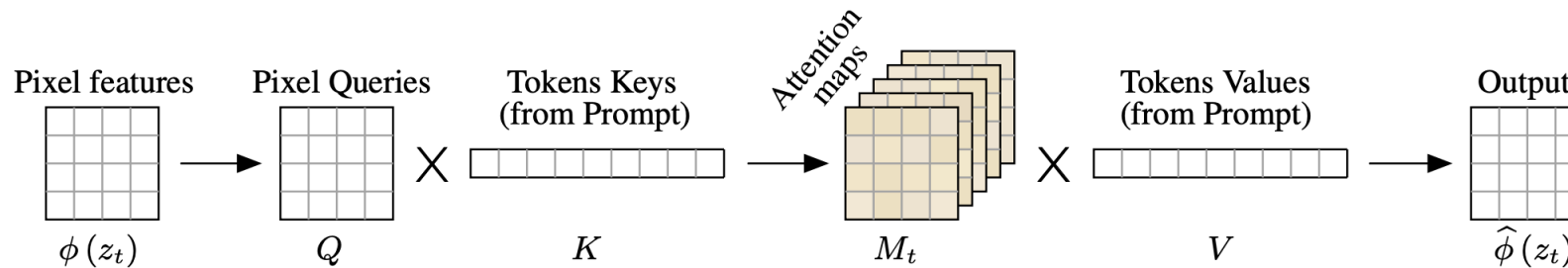
# Prompt-to-Prompt<sup>[8]</sup>

- How ?

Cross attention control!

$$Q = W_Q^{(i)} \cdot \varphi_i(z_t), \quad K = W_K^{(i)} \cdot \tau_\theta(y), \quad V = W_V^{(i)} \cdot \tau_\theta(y).$$

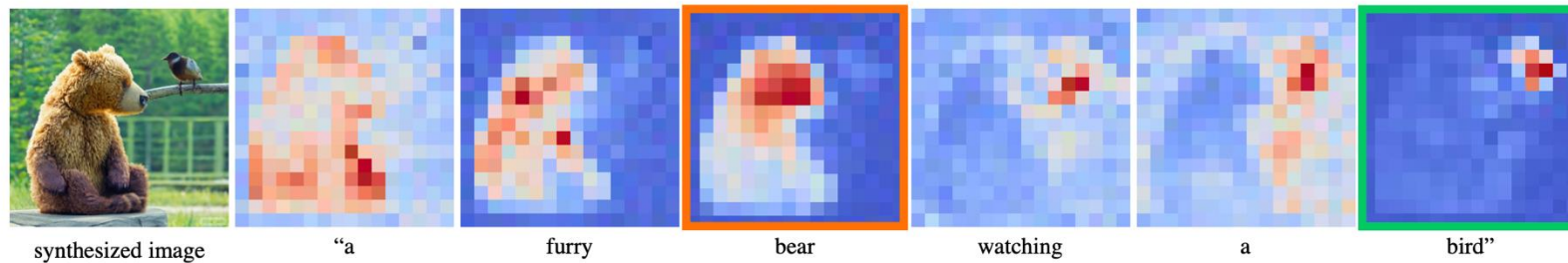
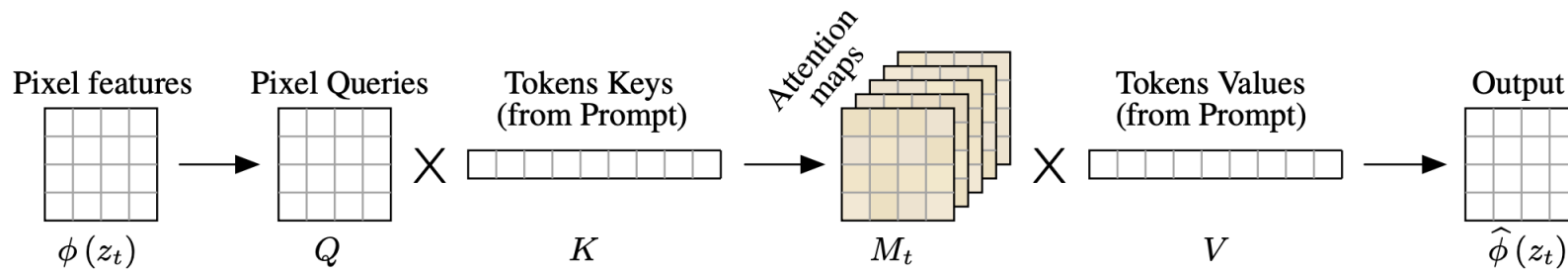
$$M = \text{Softmax} \left( \frac{QK^T}{\sqrt{d}} \right),$$



[8] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., & Cohen-Or, D. (2022). Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.

# Prompt-to-Prompt<sup>[8]</sup>

- Cross Attention: pixels are more attracted to the words that describe them

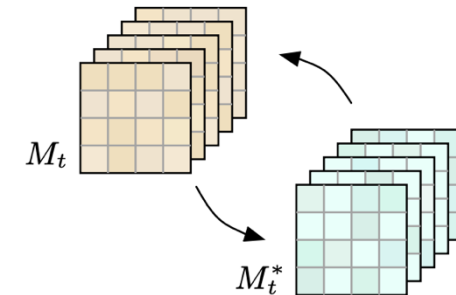


Average attention maps across all timestamps

[8] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., & Cohen-Or, D. (2022). Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.

# Prompt-to-Prompt<sup>[8]</sup>

- How to control cross attention?
  - Source prompt: Photo of a cat riding on a bicycle
  - Target prompt: Photo of a cat riding on a car
- The Attention map will be different!



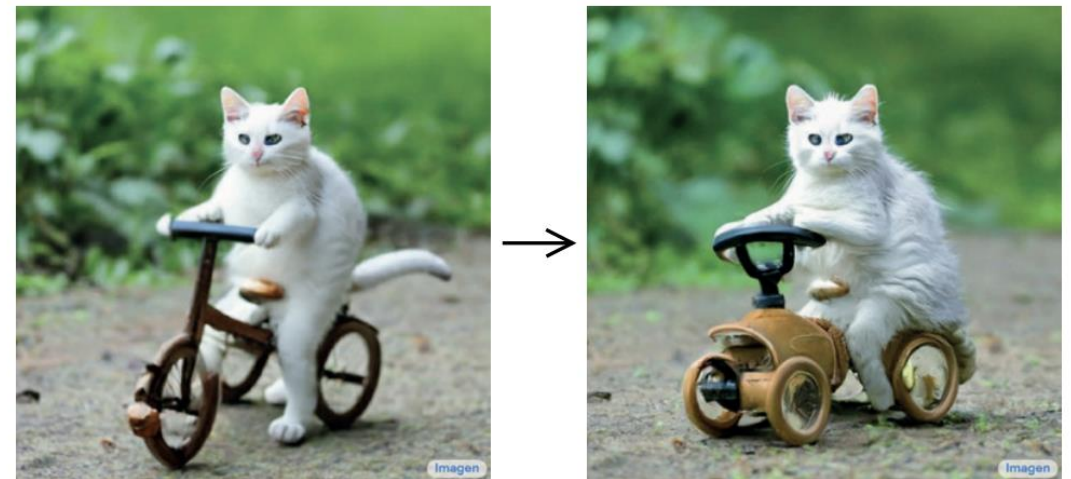
Word Swap

---

## Algorithm 1: Prompt-to-Prompt image editing

---

- 1 **Input:** A source prompt  $\mathcal{P}$ , a target prompt  $\mathcal{P}^*$ , and a random seed  $s$ .
  - 2 **Output:** A source image  $x_{src}$  and an edited image  $x_{dst}$ .
  - 3  $z_T \sim N(0, I)$  a unit Gaussian random variable with random seed  $s$ ;
  - 4  $z_T^* \leftarrow z_T$ ;
  - 5 **for**  $t = T, T - 1, \dots, 1$  **do**
  - 6      $z_{t-1}, M_t \leftarrow DM(z_t, \mathcal{P}, t, s)$ ;
  - 7      $M_t^* \leftarrow DM(z_t^*, \mathcal{P}^*, t, s)$ ;
  - 8      $\widehat{M}_t \leftarrow Edit(M_t, M_t^*, t)$ ;
  - 9      $z_{t-1}^* \leftarrow DM(z_t^*, \mathcal{P}^*, t, s) \{M \leftarrow \widehat{M}_t\}$ ;
  - 10 **end**
  - 11 **Return**  $(z_0, z_0^*)$
- 

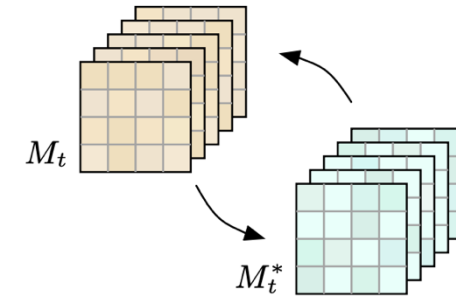


[8] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., & Cohen-Or, D. (2022). Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.



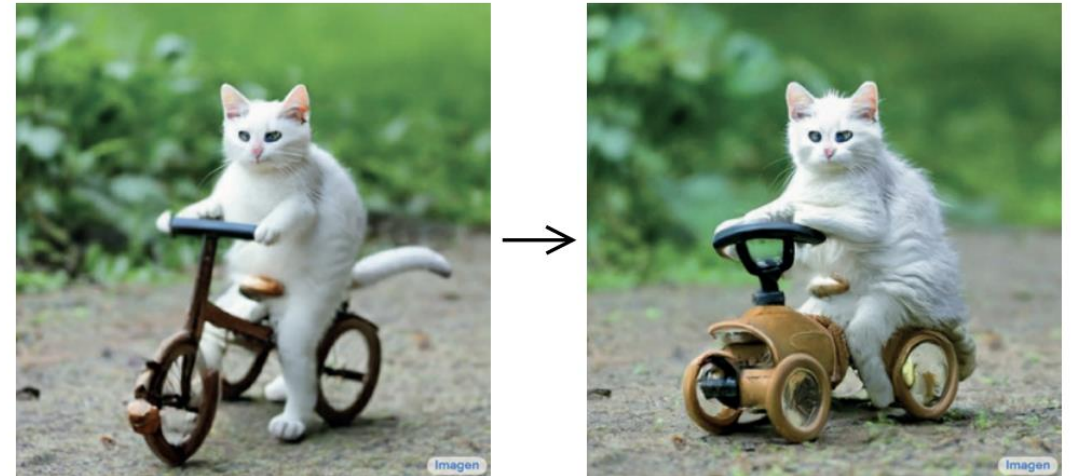
# Prompt-to-Prompt<sup>[8]</sup>

- How to control cross attention?
  - Source prompt: Photo of a cat riding on a bicycle
  - Target prompt: Photo of a cat riding on a car
- Scheduling is also important



Word Swap

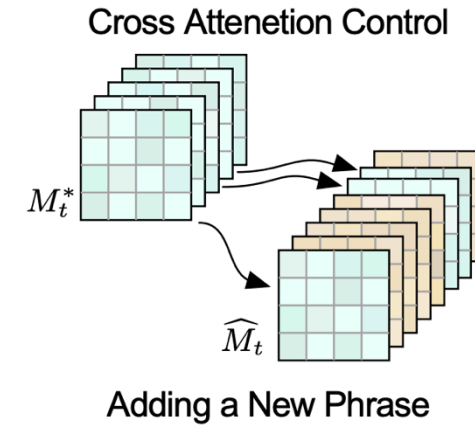
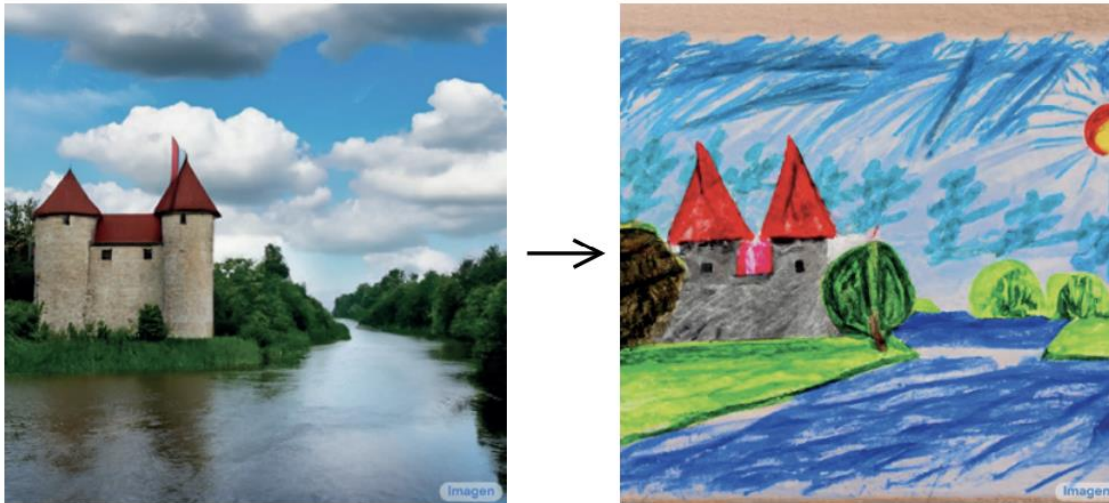
$$\text{Edit}(M_t, M_t^*, t) := \begin{cases} M_t^* & \text{if } t < \tau \\ M_t & \text{otherwise.} \end{cases}$$



[8] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., & Cohen-Or, D. (2022). Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.

# Prompt-to-Prompt<sup>[8]</sup>

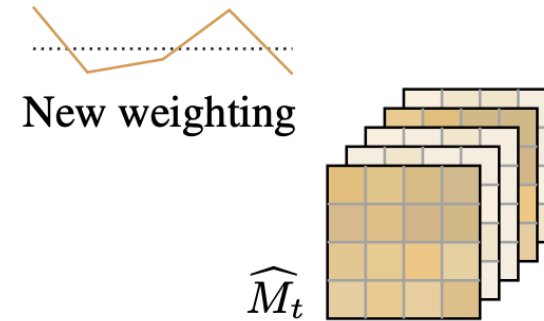
- How to control cross attention?
  - Source prompt: a castle next to a river
  - Target prompt: children drawing of a castle next to a river



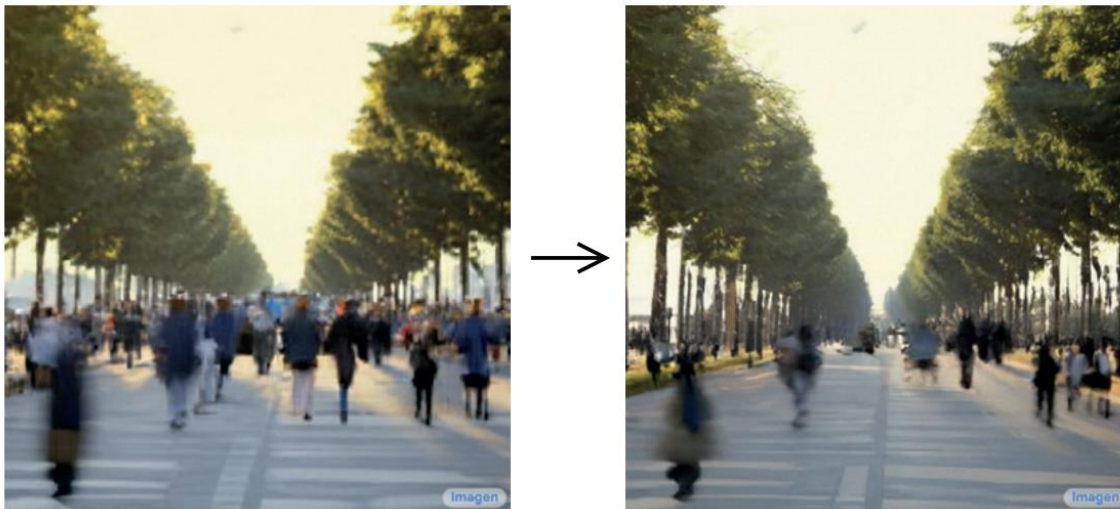
[8] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., & Cohen-Or, D. (2022). Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.

# Prompt-to-Prompt<sup>[8]</sup>

- How to control cross attention?
  - Source prompt: The boulevards are crowded today
  - Target prompt: The boulevards are crowded today
- The Attention map will be different!



Attention Re-weighting



[8] Hertz, A., Mokady, R., Tenenbaum, J., Aberman, K., Pritch, Y., & Cohen-Or, D. (2022). Prompt-to-prompt image editing with cross attention control. *arXiv preprint arXiv:2208.01626*.

# GPT-3<sup>[9]</sup>

- With Prompt to Prompt
  - We can generate images pair given captions pair
- How to automatically get captions pair?
  - Large pretrained language model GPT-3
  - Finetune GPT-3 on 700 captions and manually wrote instructions and output captions

## (a) Generate text edits:



[9] Brown, Tom, et al. "Language models are few-shot learners." Advances in neural information processing systems 33 (2020): 1877-1901.

# GPT-3<sup>[9]</sup>

- How to automatically get captions pair?
  - Large pretrained language model GPT-3
  - Finetune GPT-3 on 700 captions and manually wrote instructions and output captions

	<b>Input LAION caption</b>	<b>Edit instruction</b>	<b>Edited caption</b>
<b>Human-written (700 edits)</b>	<i>Yefim Volkov, Misty Morning</i>	<i>make it afternoon</i>	<i>Yefim Volkov, Misty Afternoon</i>
	<i>girl with horse at sunset</i>	<i>change the background to a city</i>	<i>girl with horse at sunset in front of city</i>
	<i>painting-of-forest-and-pond</i>	<i>Without the water.</i>	<i>painting-of-forest</i>
	...	...	...
<b>GPT-3 generated (&gt;450,000 edits)</b>	<i>Alex Hill, Original oil painting on canvas, Moonlight Bay</i>	<i>in the style of a coloring book</i>	<i>Alex Hill, Original coloring book illustration, Moonlight Bay</i>
	<i>The great elf city of Rivendell, sitting atop a waterfall as cascades of water spill around it</i>	<i>Add a giant red dragon</i>	<i>The great elf city of Rivendell, sitting atop a waterfall as cascades of water spill around it with a giant red dragon flying overhead</i>
	<i>Kate Hudson arriving at the Golden Globes 2015</i>	<i>make her look like a zombie</i>	<i>Zombie Kate Hudson arriving at the Golden Globes 2015</i>
	...	...	...

# Dataset Generate Flow

## Training Data Generation

### (a) Generate text edits:

Input Caption: *"photograph of a girl riding a horse"* →

GPT-3

Instruction: *"have her ride a dragon"*

Edited Caption: *"photograph of a girl riding a dragon"*

### (b) Generate paired images:

Input Caption: *"photograph of a girl riding a horse"* →

Stable Diffusion  
+ Prompt2Prompt



### (c) Generated training examples:

*"convert to brick"*



*"Color the cars pink"*



*"Make it lit by fireworks"*



*"have her ride a dragon"*



...



(a) Without Prompt-to-Prompt.

(b) With Prompt-to-Prompt.

Importance of using  
Prompt-to-Prompt

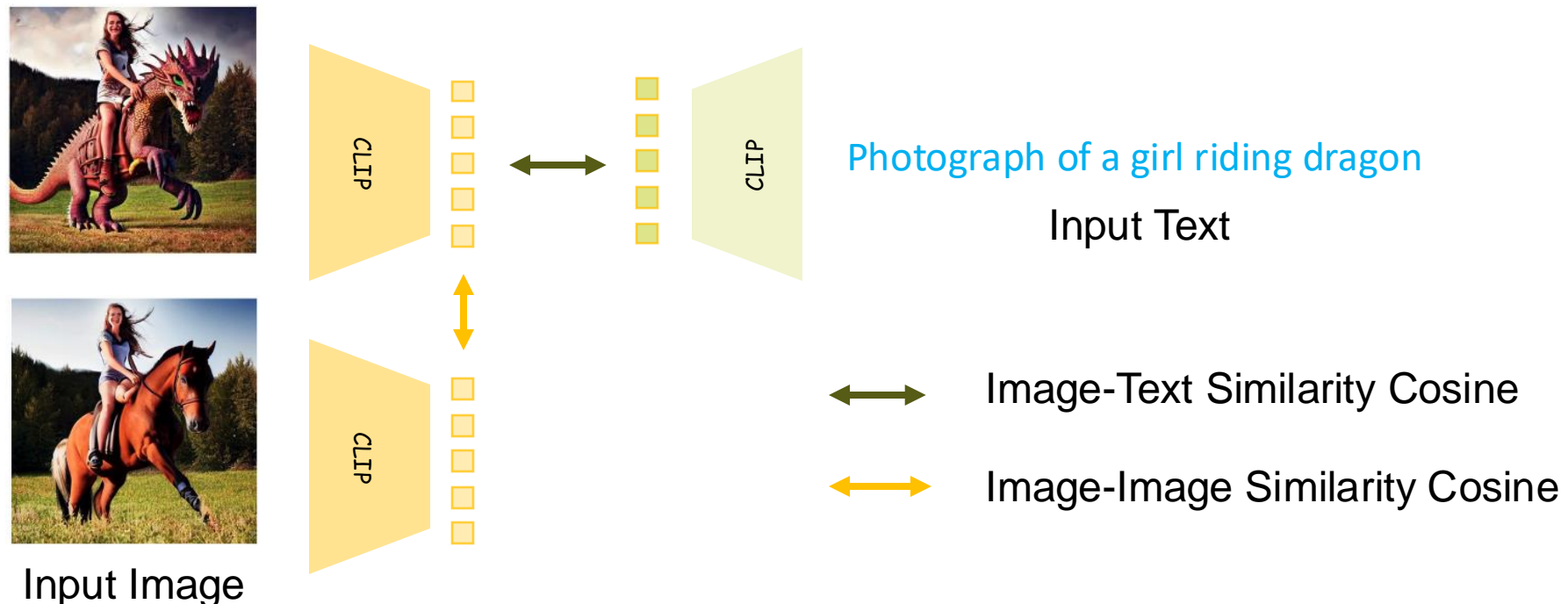
# Dataset Quality Control

- What if Prompt-to-Prompt fail?



# Dataset Quality Control

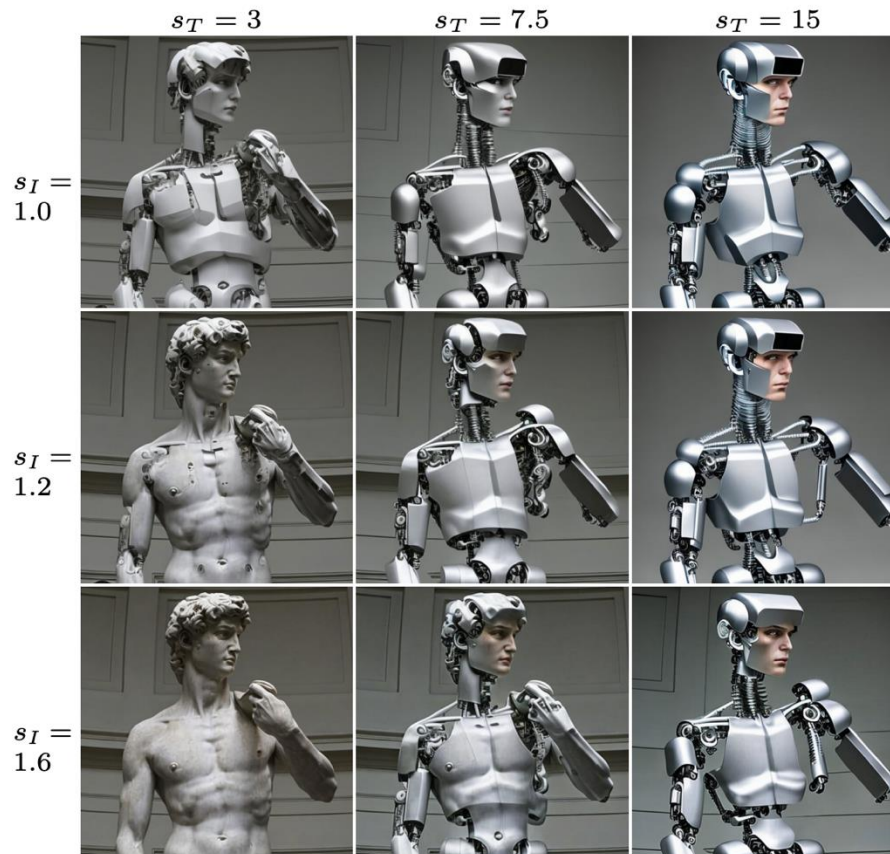
- How to control the data quality?
  - Try different prompt to prompt hyperparameters
    - One caption -> 100 images pair
  - Filter the pairs with CLIP<sup>[10]</sup>
    - CLIP can estimate image-text similarity





# Classifier Free Guidance

- Stable diffusion training in both conditional and unconditional
  - 5% randomly set to null during training
  - During inference:



Edit instruction: "Turn him into a cyborg!"

$c_I$  Image conditioning

$c_T$  Text instruction conditioning

$s_I$   $s_T$  Guidance Scales

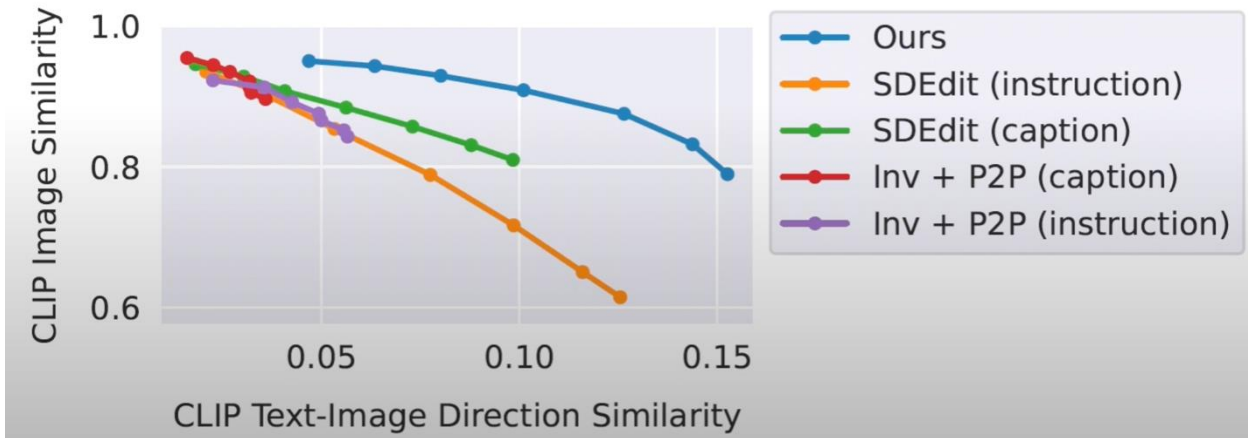
$$\begin{aligned}\tilde{e}_\theta(z_t, c_I, c_T) &= e_\theta(z_t, \emptyset, \emptyset) \\ &+ s_I \cdot (e_\theta(z_t, c_I, \emptyset) - e_\theta(z_t, \emptyset, \emptyset)) \\ &+ s_T \cdot (e_\theta(z_t, c_I, c_T) - e_\theta(z_t, c_I, \emptyset))\end{aligned}$$

# Experiments:

# Experiment:

- Quantitative Results:

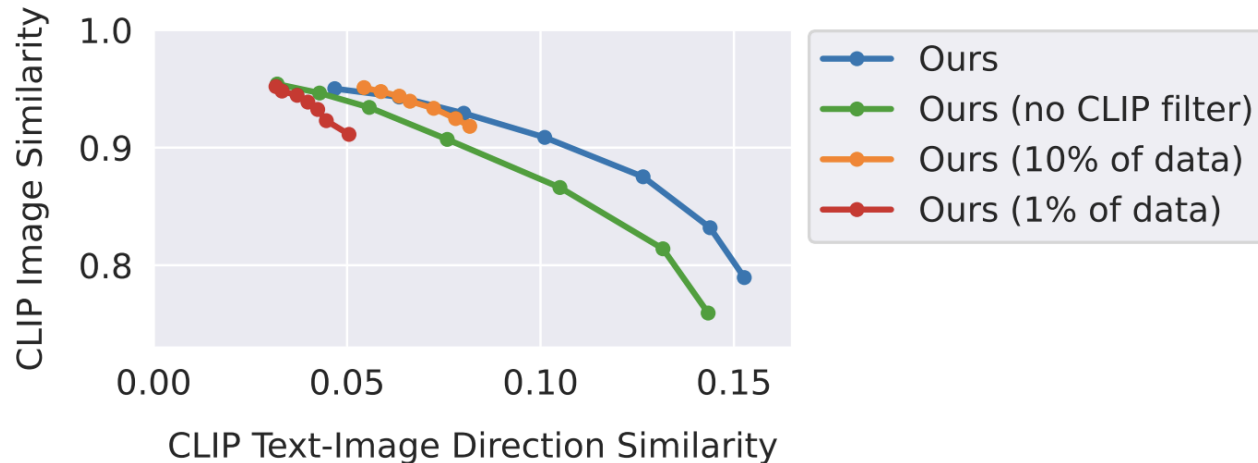
- Compare With baseline SDEdit, P2P
- CLIP Text-image Direction Similarity:
  - How much the change in text captions agrees with the change in the images
- Fix text guidance  $s_T = 7.5$  and vary the image guidance  $s_I = [1.0, 2.2]$



$$\begin{aligned}\Delta T &= E_T(t_{target}) - E_T(t_{source}) , \\ \Delta I &= E_I(G_{train}(w)) - E_I(G_{frozen}(w)) , \\ \mathcal{L}_{direction} &= 1 - \frac{\Delta I \cdot \Delta T}{|\Delta I| |\Delta T|} .\end{aligned}$$

# Experiment:

- Quantitative Results:
  - Importance of scale of dataset and filtering
  - Fix text guidance  $s_T = 7.5$  and vary the image guidance  $s_I = [1.0, 2.2]$



$$\begin{aligned}\Delta T &= E_T(t_{target}) - E_T(t_{source}) , \\ \Delta I &= E_I(G_{train}(w)) - E_I(G_{frozen}(w)) , \\ \mathcal{L}_{direction} &= 1 - \frac{\Delta I \cdot \Delta T}{|\Delta I| |\Delta T|} .\end{aligned}$$

# Experiment:

- Qualitative Results:
  - Trained on generated images, but generalize to real images and text instructions



Input



*"Make it a grocery store"*



Input



*"Add boats on the water"*



*"Replace the mountains with a city skyline"*

# Experiment:

- Qualitative Results:
  - Compared with baseline

Input



SDEdit-OC [39]



T2L [6]



SDEdit-E [39]

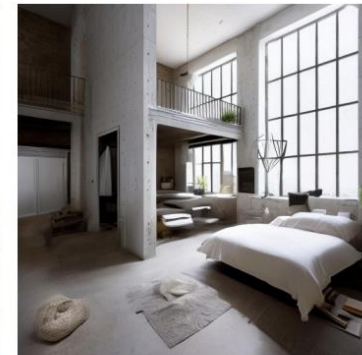


Ours



“Dali Painting of Nimbus Cloud...”

“make it look like a Dali Painting”



“Industrial design bedroom furniture...”

“add a bedroom”

# Custom demo:



Have him wear a hat.



Remove the people

Make it Minecraft



Turn the cat to dog



Some Thoughts:



# Contribution:

- Large Scale Dataset
  - 454,445 pairs after filtering
- Fast and simple text instruction based method
  - Classifier Free guidance on two conditionings
  - Impressive results and generalization to real / drawing
- Inspiring:
  - Take advantage of pretrained weights

# Limitation:

- Large Scale Dataset, large enough?
  - 0.4 million vs 3 billion (Imagen) vs 5.8 billion (stable diffusion)
- Large Scale Dataset, good enough?
  - Limited data quality control with CLIP,
  - 0.75 threshold image to image, 0.2 threshold image to text
  - User Study



↔  
CLIP Score:  
0.96



- 0.33 ↔ photograph of a girl riding a dragon
- 0.21 ↔ photograph of a girl riding a Horse
- 0.14 ↔ photograph of a girl riding a *bicycle*
- 0.34 ↔ photograph of a dragon *riding* a girl
- 0.32 ↔ photograph of a dragon *chasing* a girl

# Potential Improvement:

- Data Scalability

- Need to train on synthetic image, can it train on real image?
- Need to inverse the Gaussian process!

“A black bear is walking in the grass.”



real image

reconstructed

“...next to red flowers.”

“...when snow comes

“while another black bear

“Oil painting of...”



Real image

Reconstructed



Real image

Reconstructed



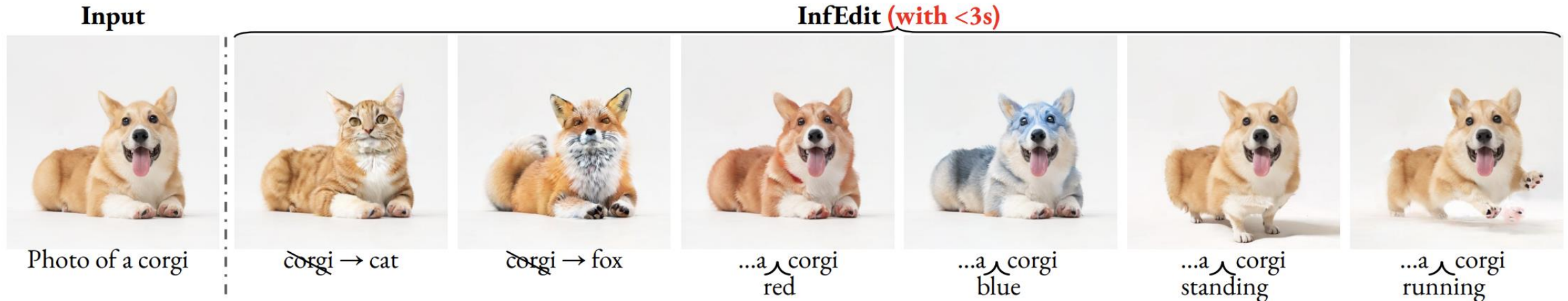
Real image

Reconstructed

# Potential Improvement:

- Data Scalability

- Improved work on this: InfEdit<sup>[11]</sup>



- Improved work on captions: ChatGPT, DeepSeek...

Q&A:

# Q&A: