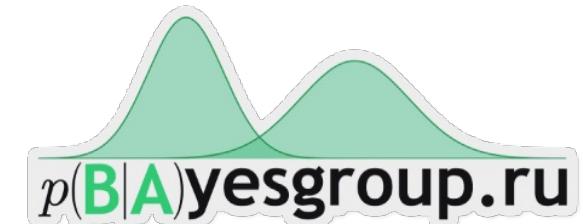


# Image Manipulation by Diffusion Models

Aibek Alanov

Research scientist at Artificial Intelligence Research Institute (AIRI) and  
Centre of Deep Learning and Bayesian Methods HSE University



# Text-to-Image Diffusion Models

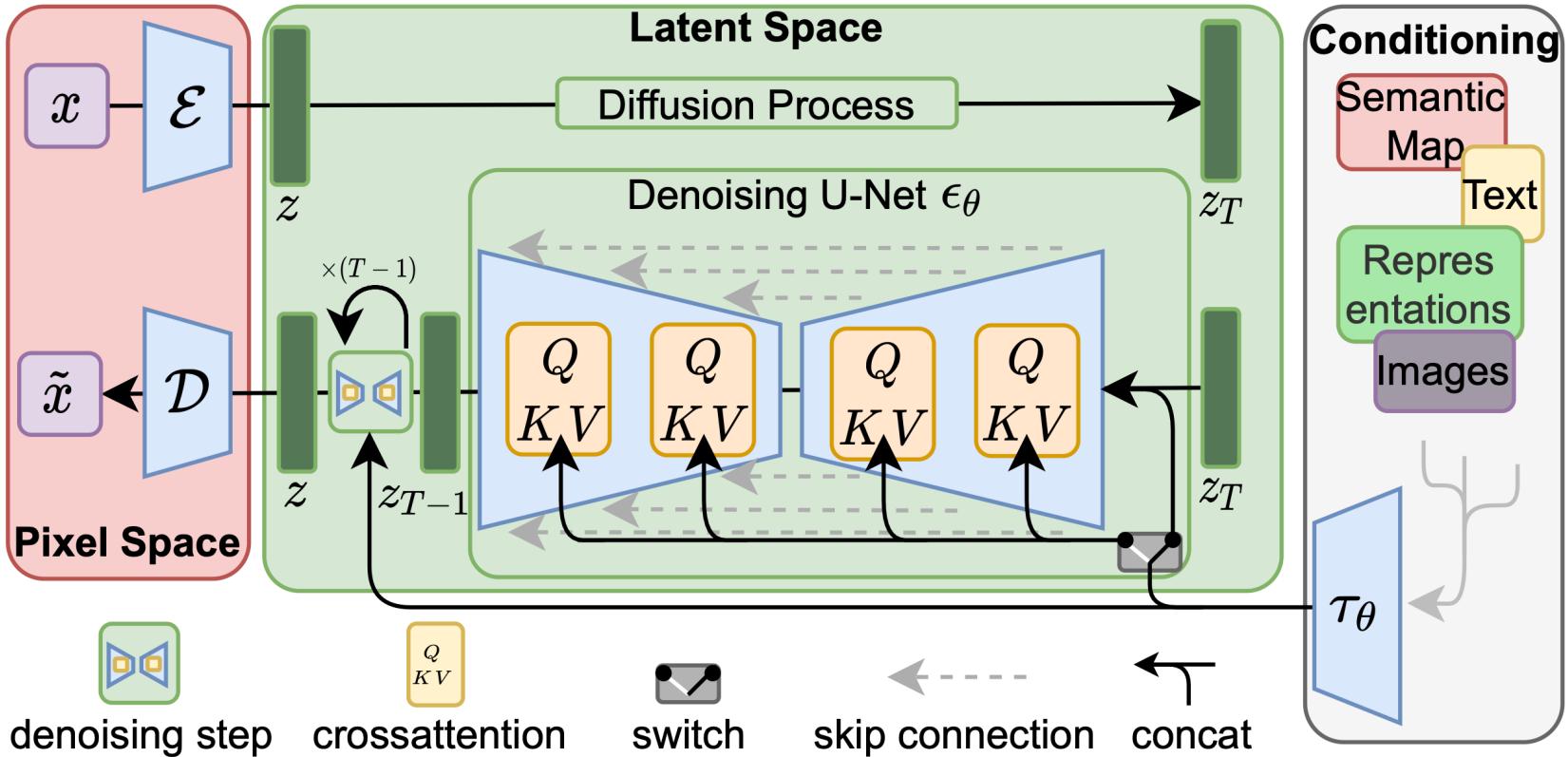


ai-forever/  
**Kandinsky-2.0**

Kandinsky 2.0 – multilingual text2image latent diffusion model



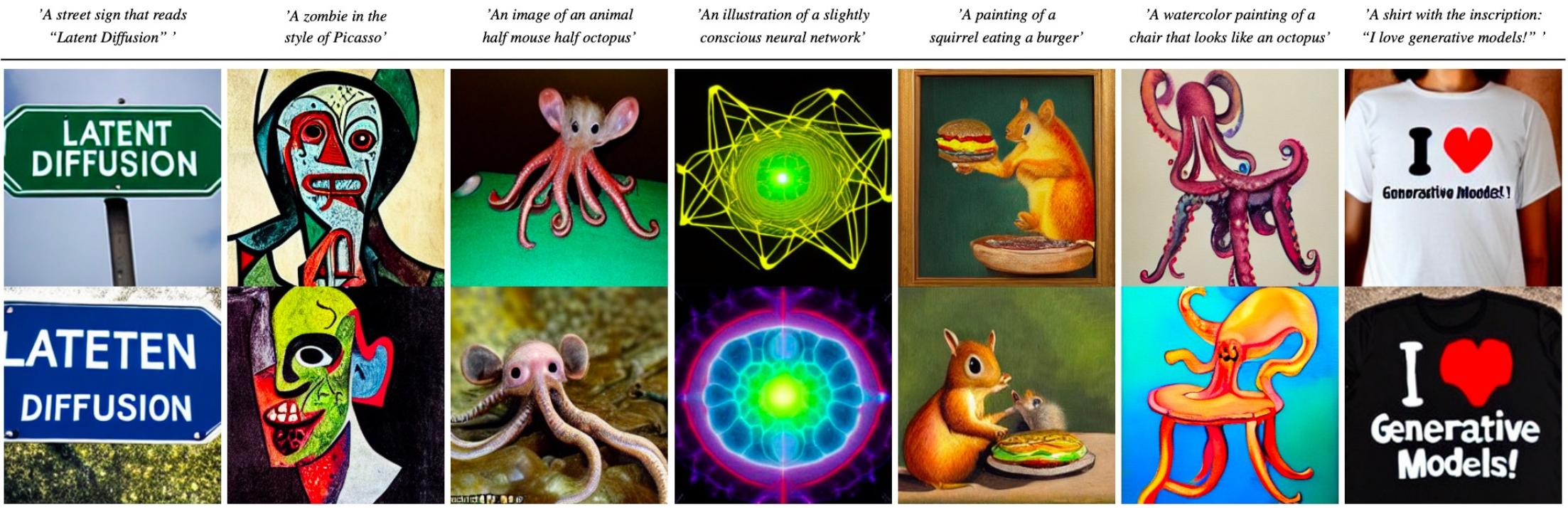
# Stable Diffusion



$$L_{LDM} := \mathbb{E}_{\mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0, 1), t} \left[ \|\epsilon - \epsilon_\theta(z_t, t, \tau_\theta(y))\|_2^2 \right]$$

# Stable Diffusion – Sample Examples

Text-to-Image Synthesis on LAION. 1.45B Model.



# Image Manipulation Problems

- Generating unique concepts in novel scenes
- Editing the given real image

# Concept Generation



Input images



in the Acropolis



swimming



sleeping



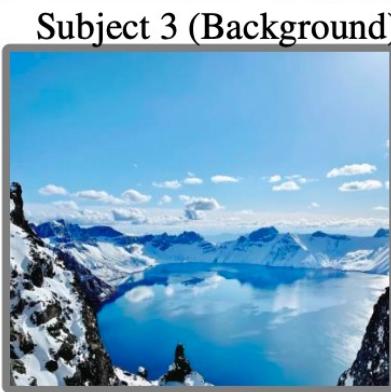
in a doghouse

in a bucket



getting a haircut

# Concept Generation



# Editing Real Image

**Input caption:** “A baby wearing a blue shirt lying on the sofa.”



**Input Image**



“... blond baby...”



“... floral shirt...”



“... golden shirt...”



“... sleeping baby...”



“baby” → “robot”

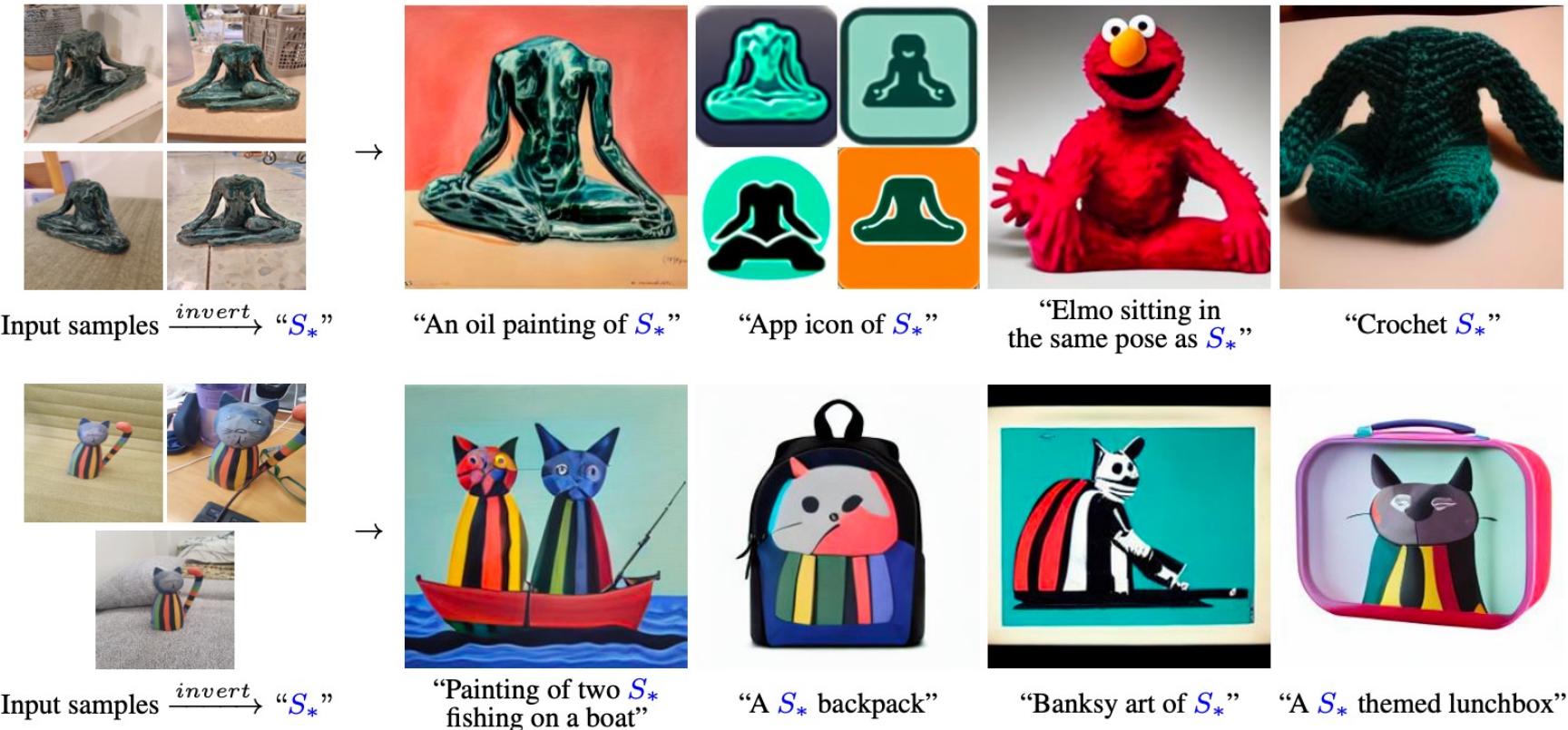


“sofa” → “grass”

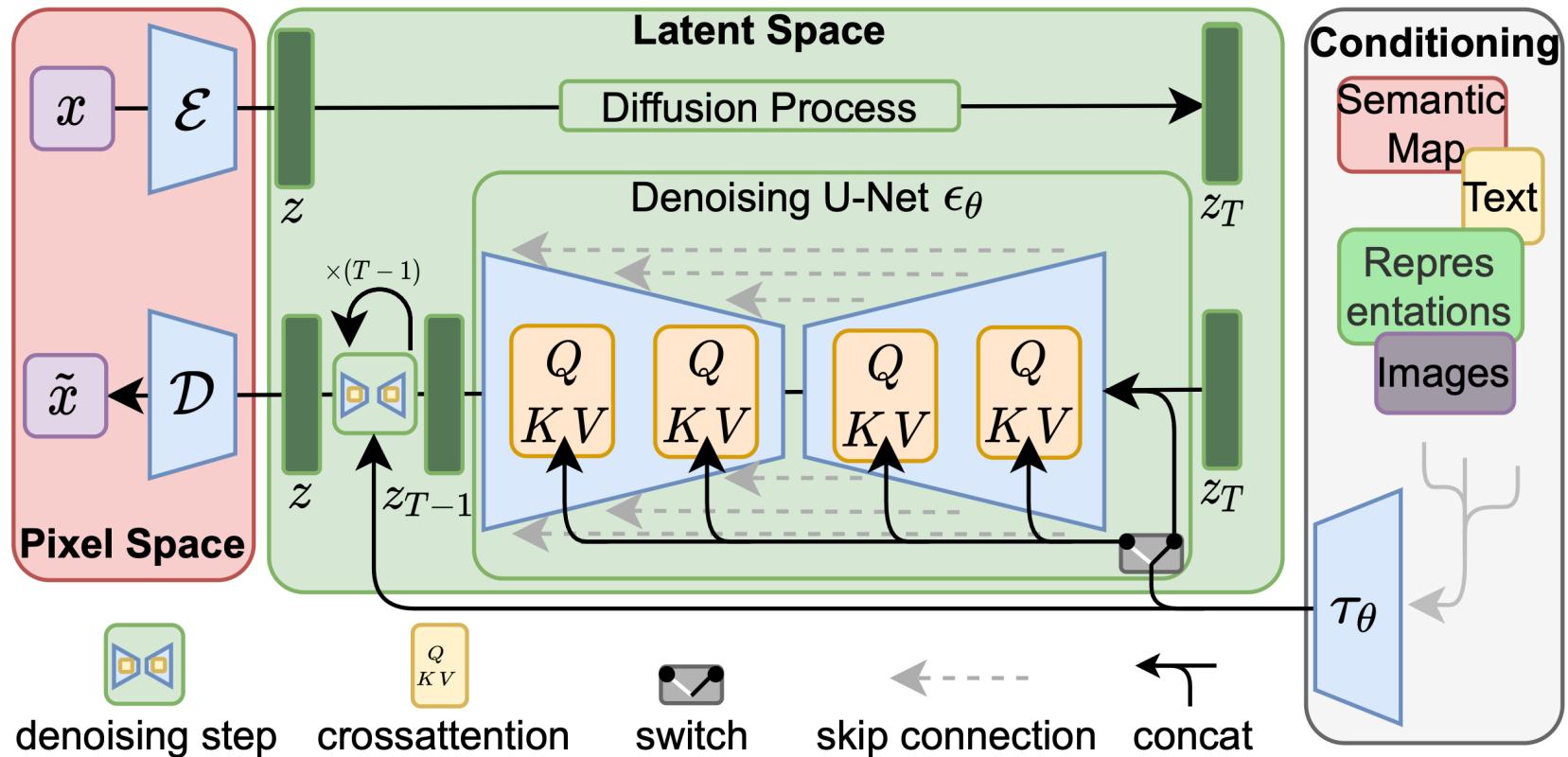


“sofa” → “ball pit”

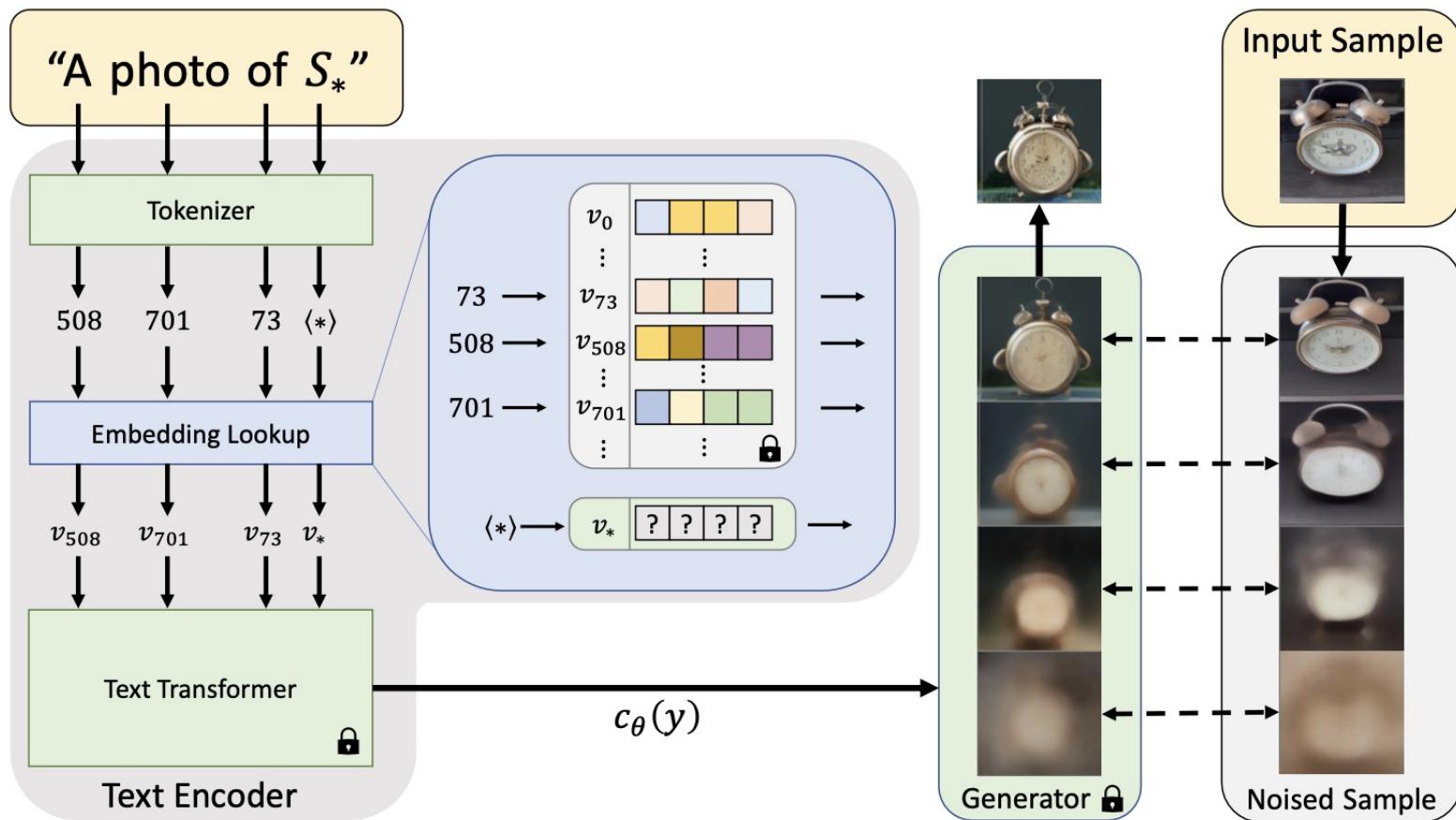
# Textual Inversion



# Textual Inversion



# Textual Inversion



# Textual Inversion

Our optimization goal can then be defined as:

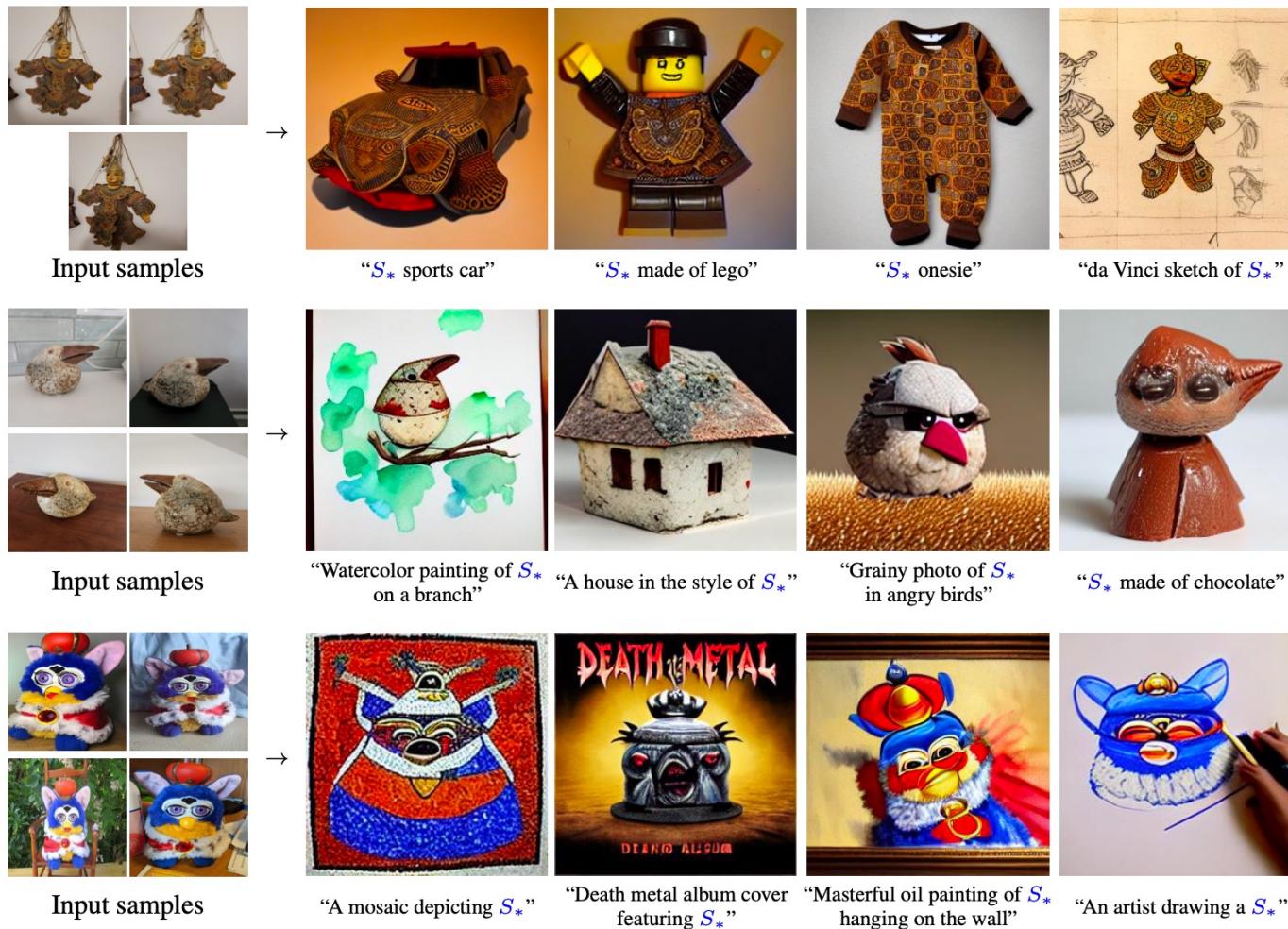
$$v_* = \arg \min_v \mathbb{E}_{z \sim \mathcal{E}(x), y, \epsilon \sim \mathcal{N}(0,1), t} \left[ \|\epsilon - \epsilon_\theta(z_t, t, c_\theta(y))\|_2^2 \right], \quad (2)$$

and is realized by re-using the same training scheme as the original LDM model, while keeping both  $c_\theta$  and  $\epsilon_\theta$  fixed. Notably, this is a reconstruction task. As such, we expect it to motivate the learned embedding to capture fine visual details unique to the concept.

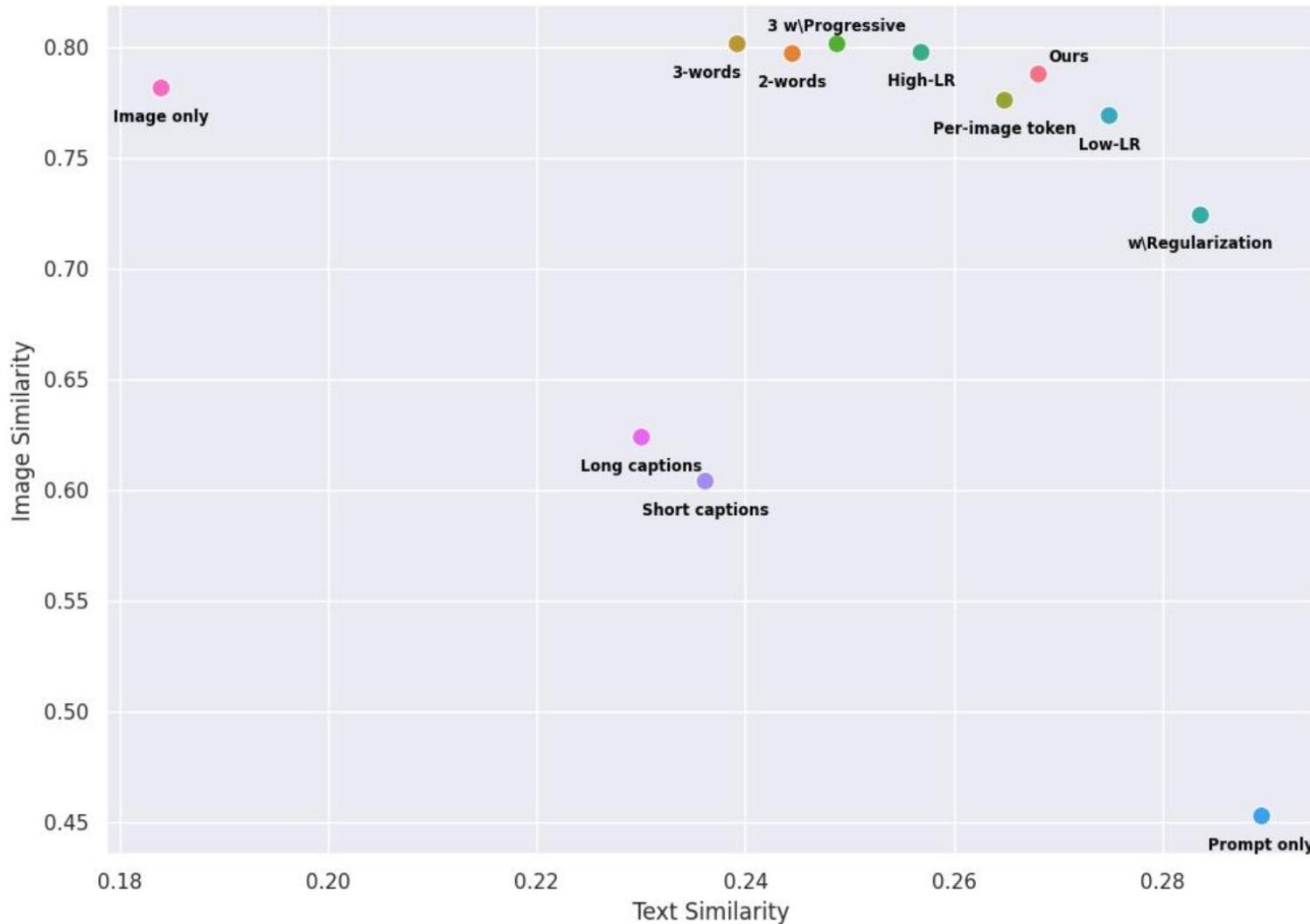
# Textual Inversion



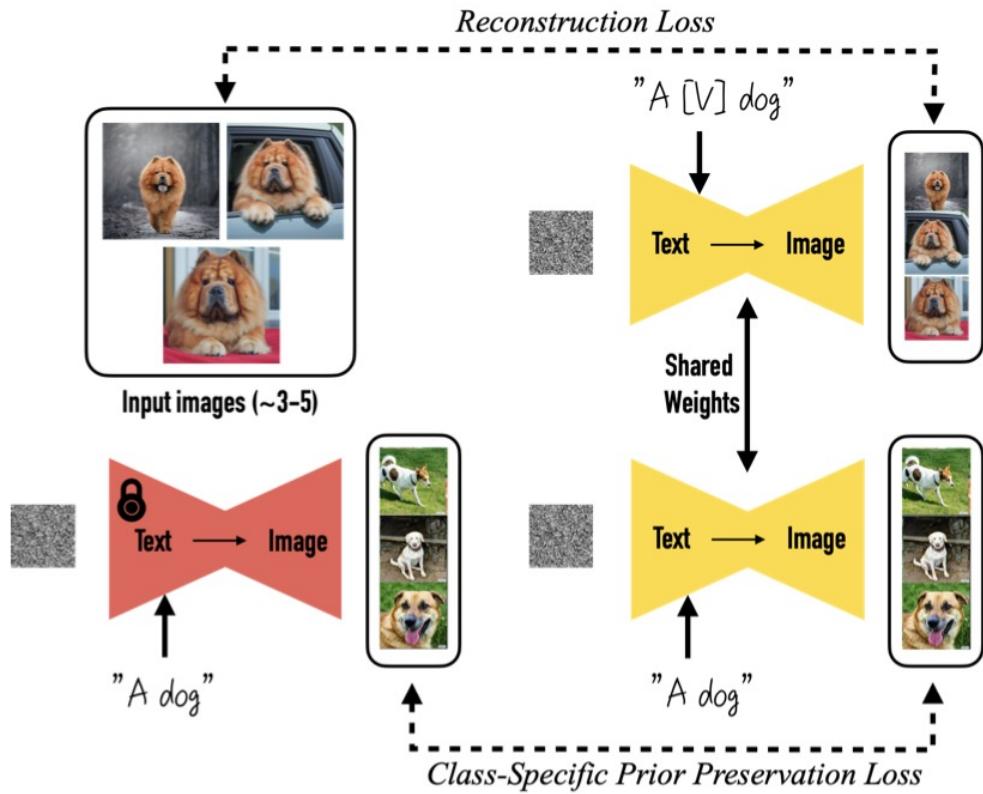
# Textual Inversion



# Textual Inversion



# DreamBooth



English words. Our approach is to find rare tokens in the vocabulary, and then invert these tokens into text space, in order to minimize the probability of the identifier having a strong prior. We perform a rare-token lookup in the vocab-

$\mathbf{c}_{\text{pr}} := \Gamma(f(\text{"a [class noun]}))$ . The loss becomes:

$$\begin{aligned} & \mathbb{E}_{\mathbf{x}, \mathbf{c}, \epsilon, \epsilon', t} [w_t \|\hat{\mathbf{x}}_\theta(\alpha_t \mathbf{x} + \sigma_t \epsilon, \mathbf{c}) - \mathbf{x}\|_2^2 + \\ & \quad \lambda w_{t'} \|\hat{\mathbf{x}}_\theta(\alpha_{t'} \mathbf{x}_{\text{pr}} + \sigma_{t'} \epsilon', \mathbf{c}_{\text{pr}}) - \mathbf{x}_{\text{pr}}\|_2^2], \end{aligned} \quad (2)$$

# DreamBooth



Input images



A [V] backpack in the  
Grand Canyon



A wet [V] backpack  
in water



A [V] backpack in Boston



A [V] backpack with the  
night sky



Input images



A [V] teapot floating  
in milk



A transparent [V] teapot  
with milk inside

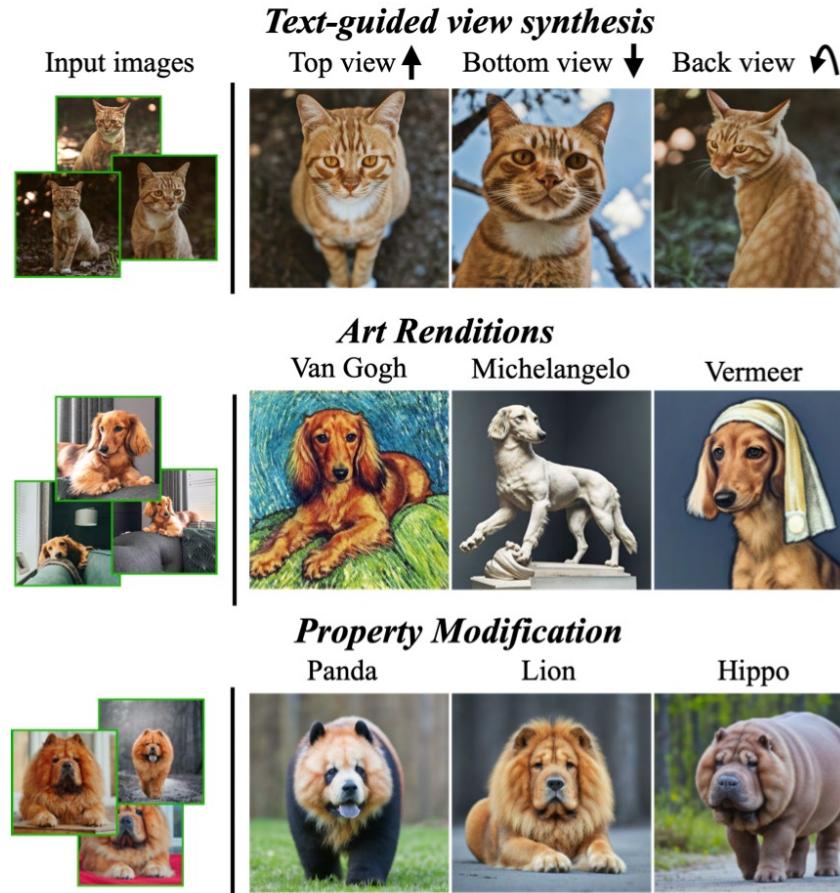


A [V] teapot  
pouring tea



A [V] teapot floating  
in the sea

# DreamBooth



# DreamBooth vs Textual Inversion



Figure 4. Comparisons with Textual Inversion [20] Given 4 input images (top row), we compare: DreamBooth Imagen (2nd row), DreamBooth Stable Diffusion (3rd row), Textual Inversion (bottom row). Output images were created with the following prompts (left to right): “a [V] vase in the snow”, “a [V] vase on the beach”, “a [V] vase in the jungle”, “a [V] vase with the Eiffel Tower in the background”. DreamBooth is stronger in both subject and prompt fidelity.

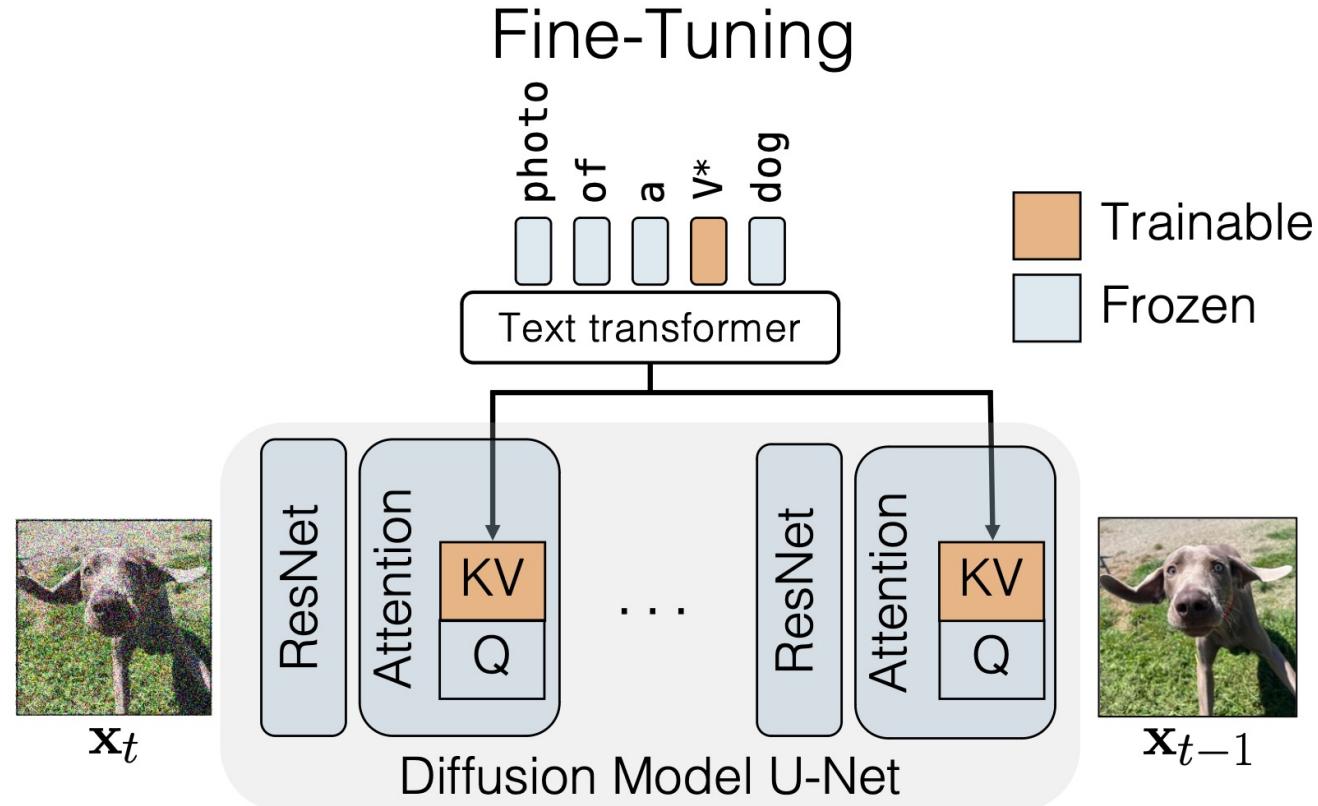
Method	DINO $\uparrow$	CLIP-I $\uparrow$	CLIP-T $\uparrow$
Real Images	0.774	0.885	N/A
DreamBooth (Imagen)	<b>0.696</b>	<b>0.812</b>	<b>0.306</b>
DreamBooth (Stable Diffusion)	0.668	0.803	0.305
Textual Inversion (Stable Diffusion)	0.569	0.780	0.255

Table 1. Subject fidelity (DINO, CLIP-I) and prompt fidelity (CLIP-T, CLIP-T-L) quantitative metric comparison.

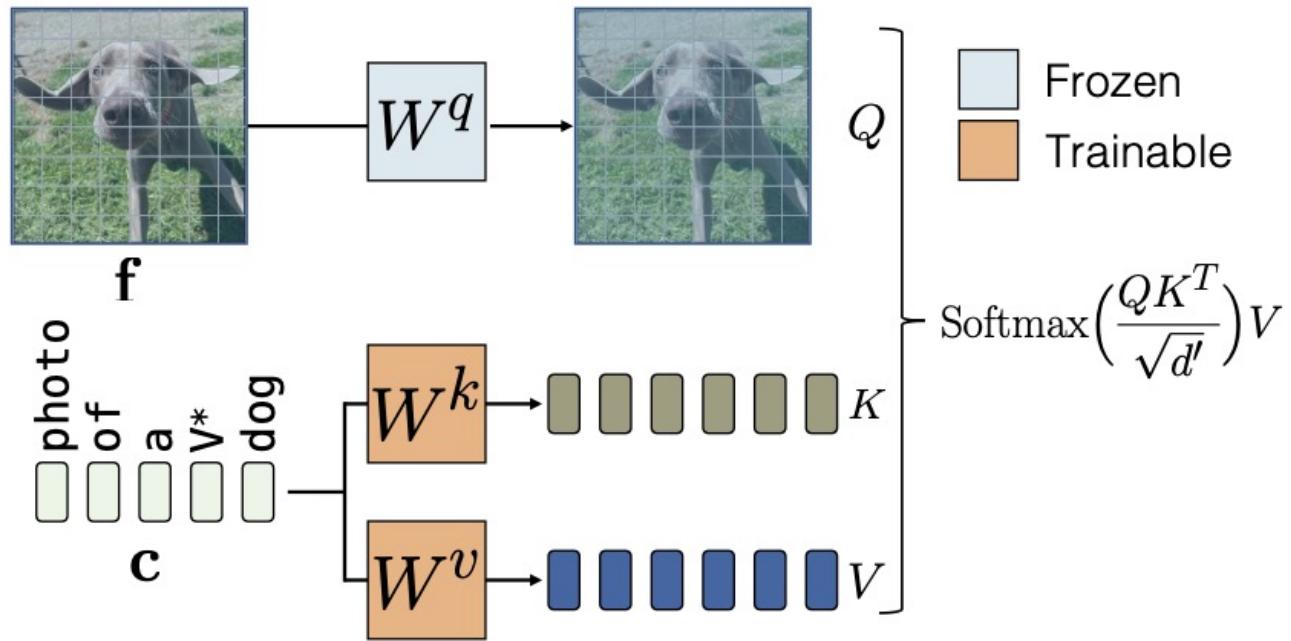
Method	Subject Fidelity $\uparrow$	Prompt Fidelity $\uparrow$
DreamBooth (Stable Diffusion)	<b>68%</b>	<b>81%</b>
Textual Inversion (Stable Diffusion)	22%	12%
Undecided	10%	7%

Table 2. Subject fidelity and prompt fidelity user preference.

# Custom Diffusion



# Custom Diffusion



# Custom Diffusion – Multiple Concept Training

- Joint training on multiple concepts
- Constrained optimization to merge concepts

$$\begin{aligned}\hat{W} = \arg \min_W & ||WC_{\text{reg}}^{\top} - W_0C_{\text{reg}}^{\top}||_F \\ \text{s.t. } & WC^{\top} = V, \text{ where } C = [\mathbf{c}_1 \cdots \mathbf{c}_N]^{\top} \\ & \text{and } V = [W_1\mathbf{c}_1^{\top} \cdots W_N\mathbf{c}_N^{\top}]^{\top}. \end{aligned}\tag{4}$$

# Custom Diffusion – Constrained Optimization

By using the method of Lagrange multipliers [8], we need to minimize the following objective:

$$L = \frac{1}{2} ||WC_{\text{reg}}^{\top} - W_0C_{\text{reg}}^{\top}|| - \text{trace}(\mathbf{v}(WC^{\top} - V)), \quad (7)$$

$\hat{W} = W_0 + \mathbf{v}^{\top} \mathbf{d}$ , where  $\mathbf{d} = C(C_{\text{reg}}^{\top} C_{\text{reg}})^{-1}$   
and  $\mathbf{v}^{\top} = (V - W_0 C^{\top})(\mathbf{d} C^{\top})^{-1}$ .

here  $\mathbf{v} \in \mathbb{R}^{s \times o}$  is the Lagrangian multiplier corresponding to the constraints. Differentiating the above objective and equating it to 0, we obtain:

$$\begin{aligned} WC_{\text{reg}}^{\top} C_{\text{reg}} - W_0 C_{\text{reg}}^{\top} C_{\text{reg}} - \mathbf{v}^{\top} C &= 0 \\ \implies W &= W_0 + \mathbf{v}^{\top} C (C_{\text{reg}}^{\top} C_{\text{reg}})^{-1}. \end{aligned} \quad (8)$$

We assume  $C_{\text{reg}}$  is non-degenerate. Using the above solution in Eqn. 6,  $WC^{\top} = V$ , we obtain:

$$\begin{aligned} (W_0 + \mathbf{v}^{\top} C (C_{\text{reg}}^{\top} C_{\text{reg}})^{-1}) C^{\top} &= V \\ \text{Let } \mathbf{d} &= C (C_{\text{reg}}^{\top} C_{\text{reg}})^{-1} \\ \mathbf{v}^{\top} &= (V - W_0 C^{\top}) (\mathbf{d} C^{\top})^{-1}. \end{aligned} \quad (9)$$

# Custom Diffusion – Results

Target Images



Custom Diffusion (Ours)



DreamBooth



Textual Inversion



Artistic variations: A watercolor painting of V\* tortoise plush on a mountain



Scene change: V\* teddybear in Times Square



Add object: V\* table and an orange sofa

# Custom Diffusion – Results

Target Images



Ours (joint training)



Ours (optimization)



DreamBooth



$V_1^*$  chair with the  $V_2^*$  cat sitting on it near a beach



The  $V_1^*$  cat is sitting inside a  $V_2^*$  wooden pot and looking up



$V_1^*$  flower in the  $V_2^*$  wooden pot on a table



# Custom Diffusion – Results

	<b>Method</b>	<b>Text-alignment</b>	<b>Image-alignment</b>	<b>KID (validation)</b>
<b>Single-Concept</b>	Textual Inversion	0.670	<b>0.827</b>	22.27
	DreamBooth	0.781	0.776	32.53
	Ours (w/ fine-tune all)	<b>0.795</b>	0.748	<b>19.27</b>
	Ours	<b>0.795</b>	0.775	20.96
<b>Multi-Concept</b>	Textual Inversion	0.544	0.630	
	DreamBooth	0.783	0.695	
	Ours (w/ fine-tune all)	0.787	0.691	
	Ours (Sequential)	0.797	0.700	—
	Ours (Optimization)	0.800	0.695	
	Ours (Joint)	<b>0.801</b>	<b>0.706</b>	

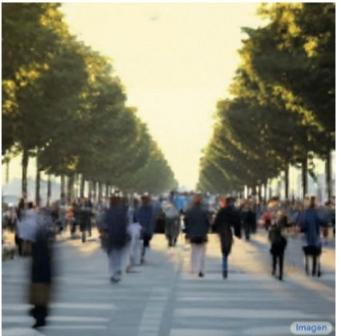
# Other Works on Concept Generation

**“ELITE: Encoding Visual Concepts into Textual Embeddings for Customized Text-to-Image Generation” (2023)**

**“Cones: Concept Neurons in Diffusion Models for Customized Generation” (2023)**

**“SVDiff: Compact Parameter Space for Diffusion Fine-Tuning” (2023)**

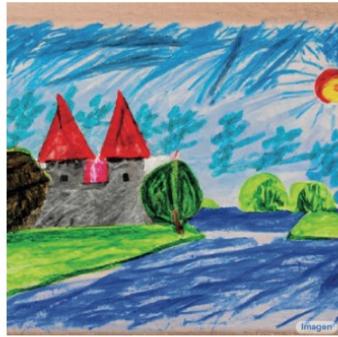
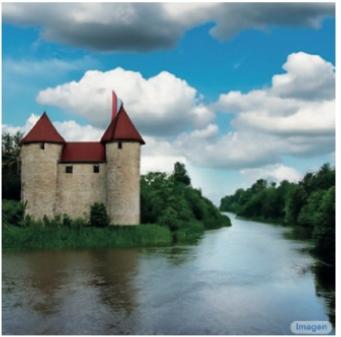
# Prompt-to-Prompt Image Editing



“The boulevards are crowded today.”



“Photo of a cat riding on a ~~bicycle~~ <sup>car</sup>.”

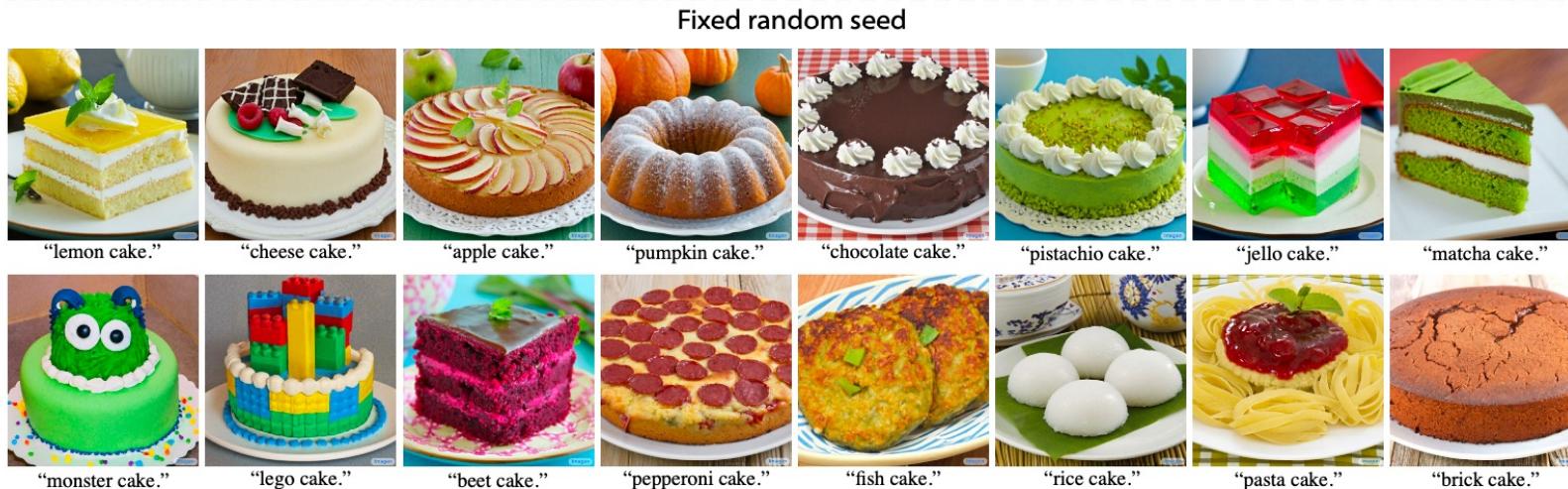
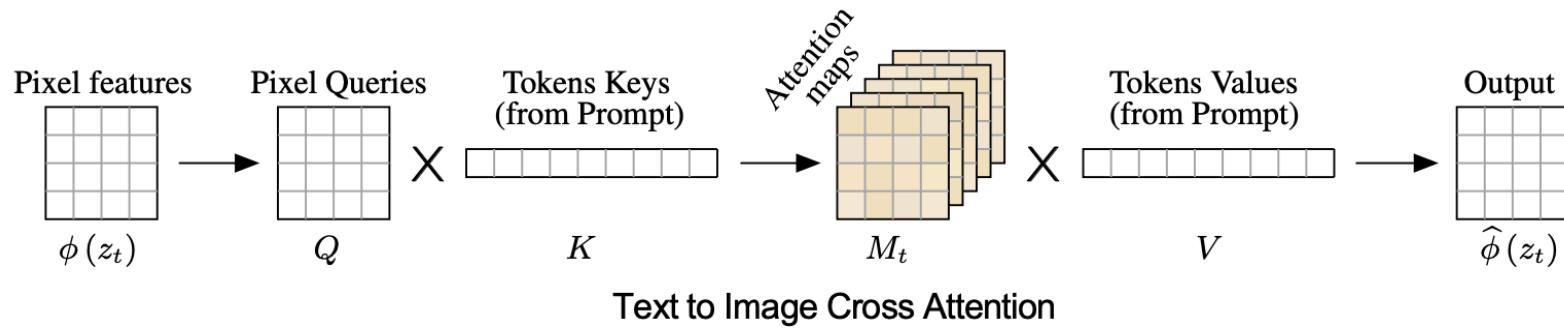


“~~Children drawing of~~ a castle next to a river.”

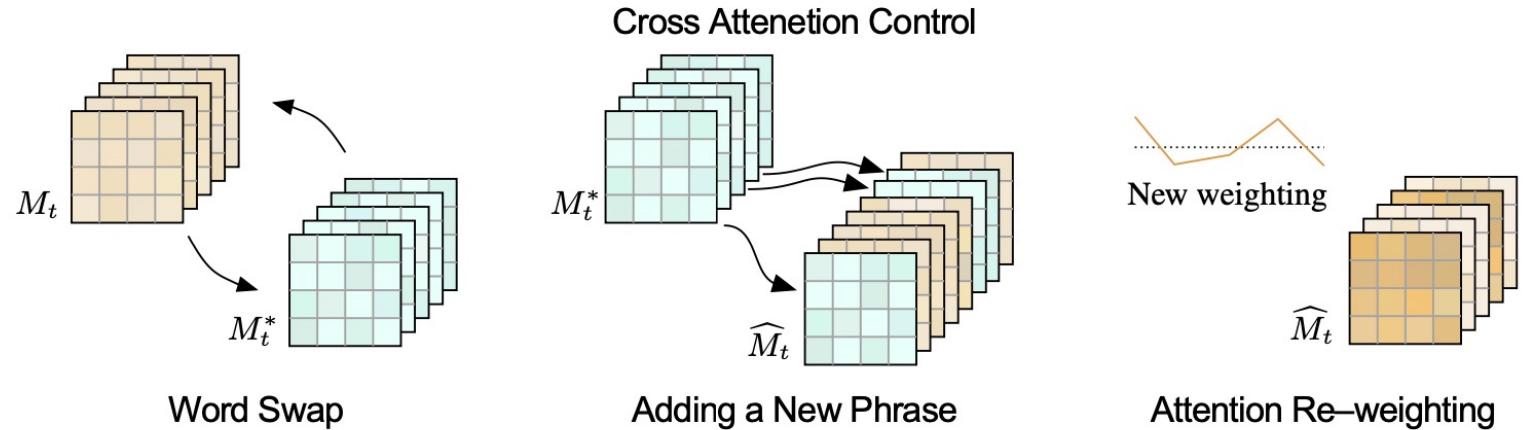


“a cake with ~~decorations~~ <sup>jelly beans</sup>.”

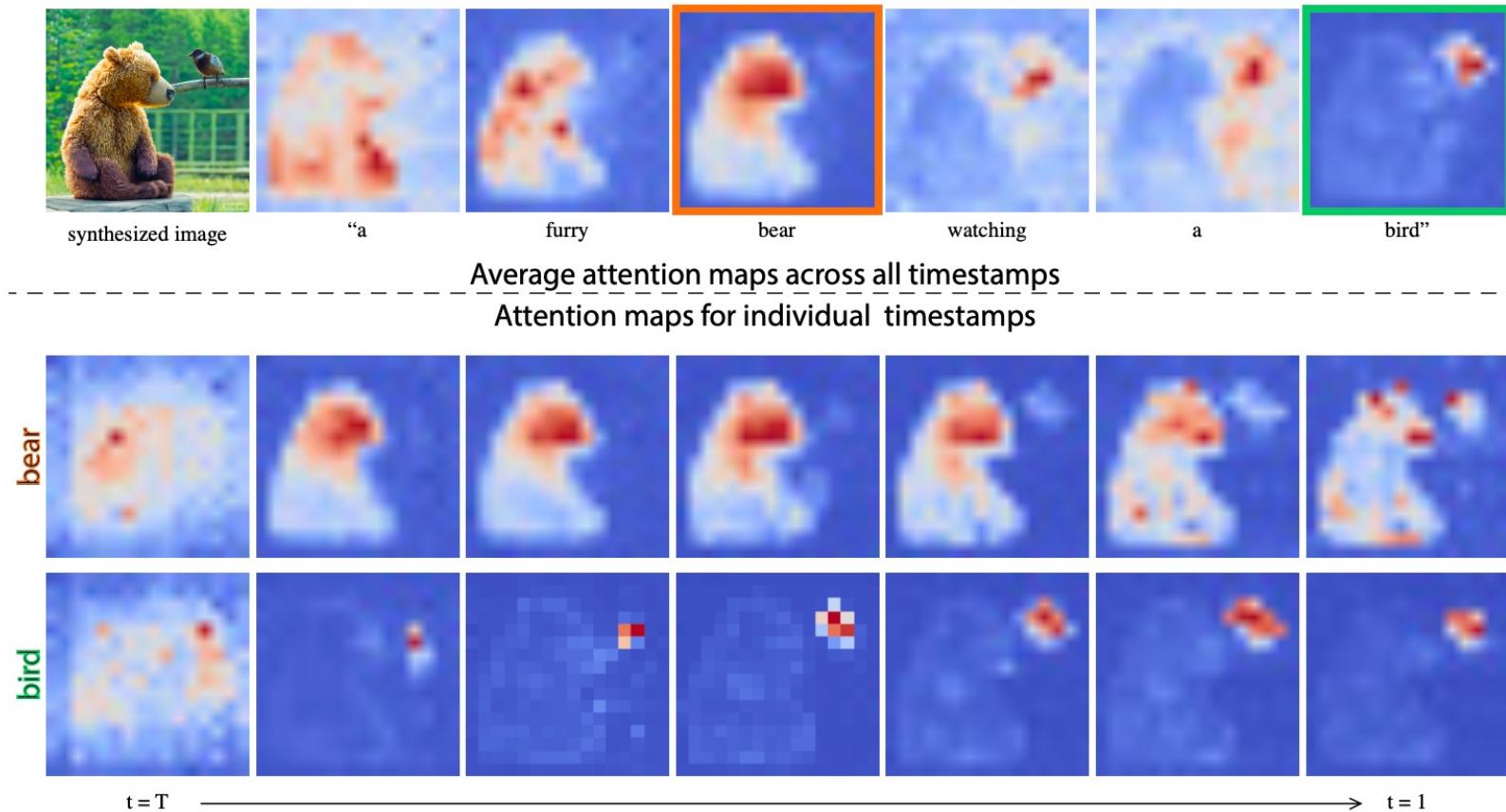
# Prompt-to-Prompt Image Editing



# Prompt-to-Prompt Image Editing



# Prompt-to-Prompt Image Editing



# Prompt-to-Prompt Image Editing

---

**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      $z_{t-1}^* \leftarrow DM(z_t^*, \mathcal{P}^*, t, s_t) \{M \leftarrow \widehat{M}_t\}$ ;
- 10 **end**
- 11 **Return**  $(z_0, z_0^*)$

---

# Null-Text Inversion for Editing Real Images

**Input Image**



**DDIM Inversion**



**Null-text Inversion**



**Prompt-to-Prompt image editing**



Input caption: "Zoom photo of flowers."

"...origami flowers..."

flowers → cupcakes

"...wither flowers..."

photo → sketch



Input caption: "A cat sitting next to a mirror."

"...silver cat sculpture..."

cat → tiger

"...sleeping cat..."

"Watercolor drawing of..."

# Null-Text Inversion for Editing Real Images

$$\min_{\theta} E_{z_0, \varepsilon \sim N(0, I), t \sim \text{Uniform}(1, T)} \|\varepsilon - \varepsilon_{\theta}(z_t, t, \mathcal{C})\|_2^2. \quad (1)$$

$$z_{t-1} = \sqrt{\frac{\alpha_{t-1}}{\alpha_t}} z_t + \left( \sqrt{\frac{1}{\alpha_{t-1}} - 1} - \sqrt{\frac{1}{\alpha_t} - 1} \right) \cdot \varepsilon_{\theta}(z_t, t, \mathcal{C}).$$

$$\tilde{\varepsilon}_{\theta}(z_t, t, \mathcal{C}, \emptyset) = w \cdot \varepsilon_{\theta}(z_t, t, \mathcal{C}) + (1 - w) \cdot \varepsilon_{\theta}(z_t, t, \emptyset).$$

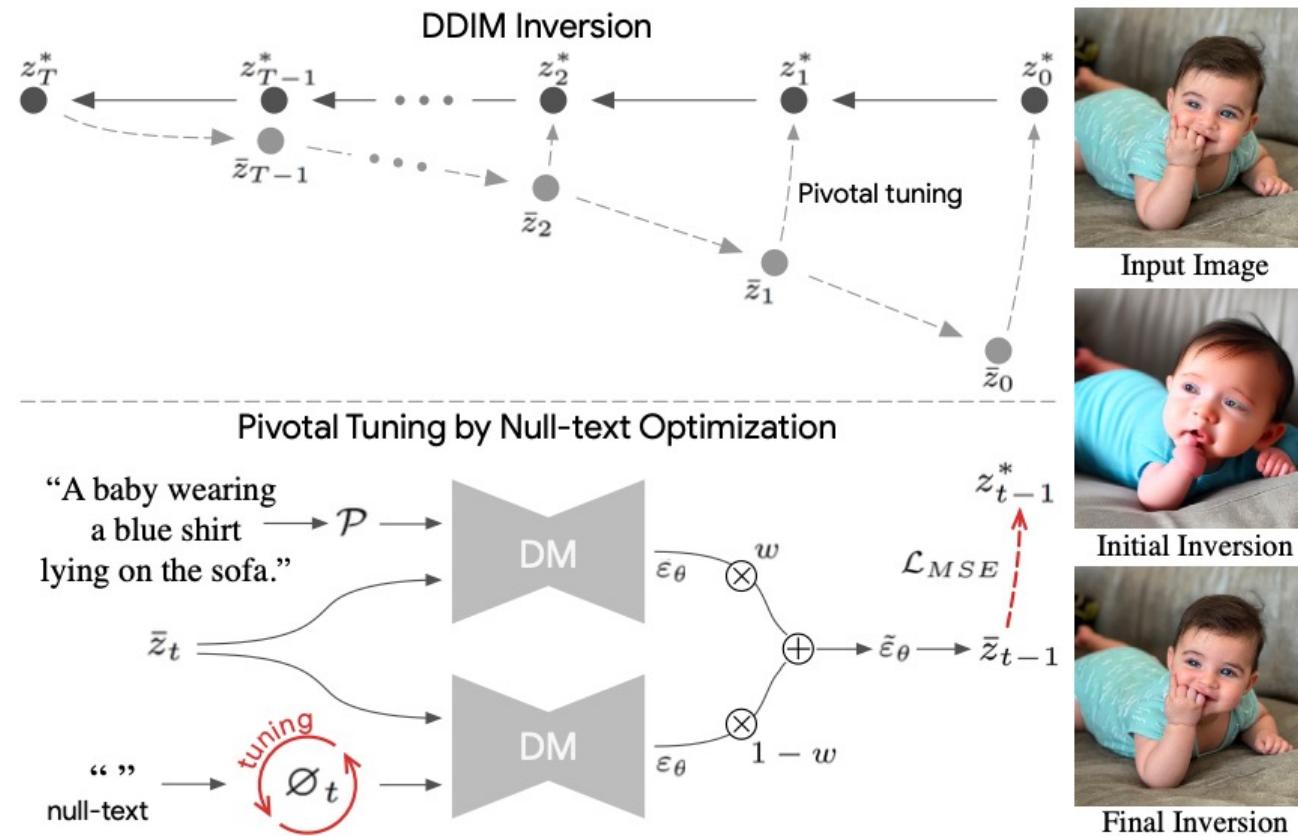
E.g.,  $w = 7.5$  is the default parameter for Stable Diffusion.

**DDIM inversion.** A simple inversion technique was suggested for the DDIM sampling [13, 35], based on the assumption that the ODE process can be reversed in the limit of small steps:

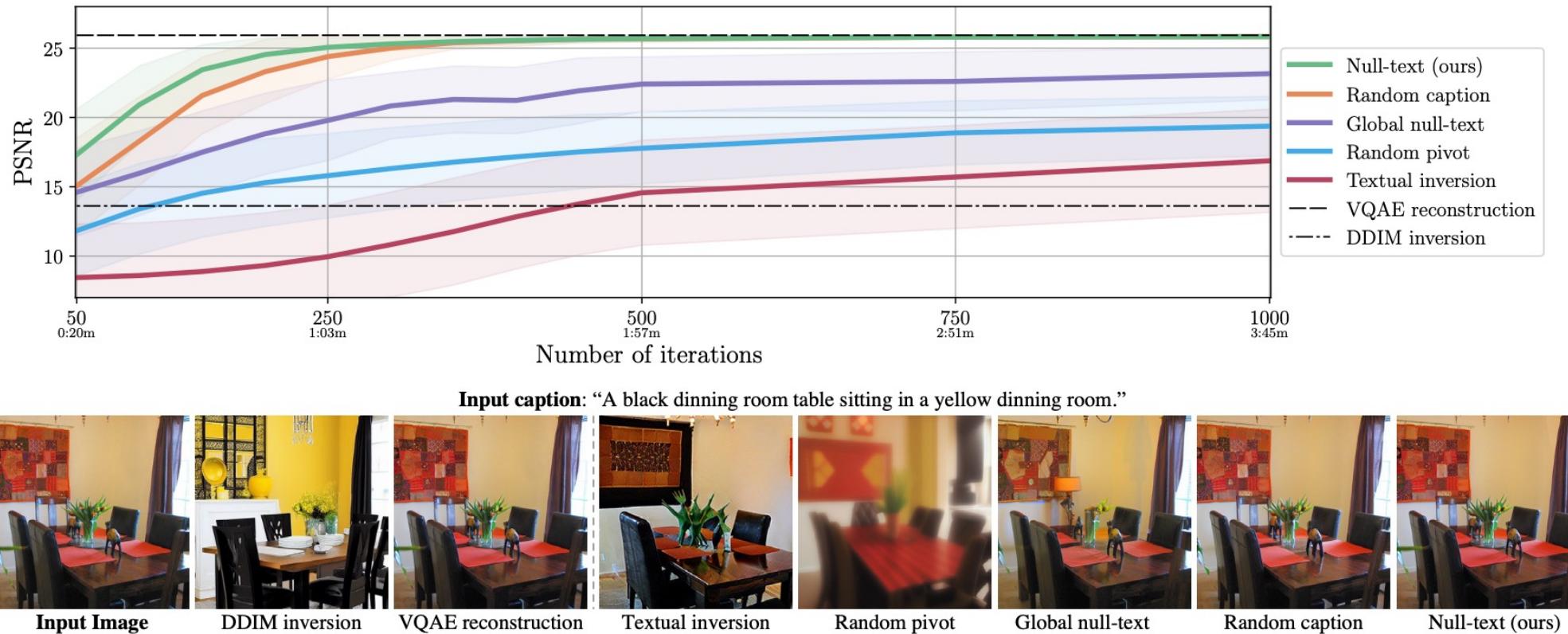
$$z_{t+1} = \sqrt{\frac{\alpha_{t+1}}{\alpha_t}} z_t + \left( \sqrt{\frac{1}{\alpha_{t+1}} - 1} - \sqrt{\frac{1}{\alpha_t} - 1} \right) \cdot \varepsilon_{\theta}(z_t, t, \mathcal{C}).$$

In other words, the diffusion process is performed in the reverse direction, that is  $z_0 \rightarrow z_T$  instead of  $z_T \rightarrow z_0$ , where  $z_0$  is set to be the encoding of the given real image.

# Null-Text Inversion for Editing Real Images



# Null-Text Inversion for Editing Real Images



# Null-Text Inversion for Editing Real Images

**Input caption:** “A baby wearing a blue shirt lying on the sofa.”



**Input Image**



“... blond baby...”



“... floral shirt...”



“... golden shirt...”



“... sleeping baby...”



“baby” $\rightarrow$ “robot”

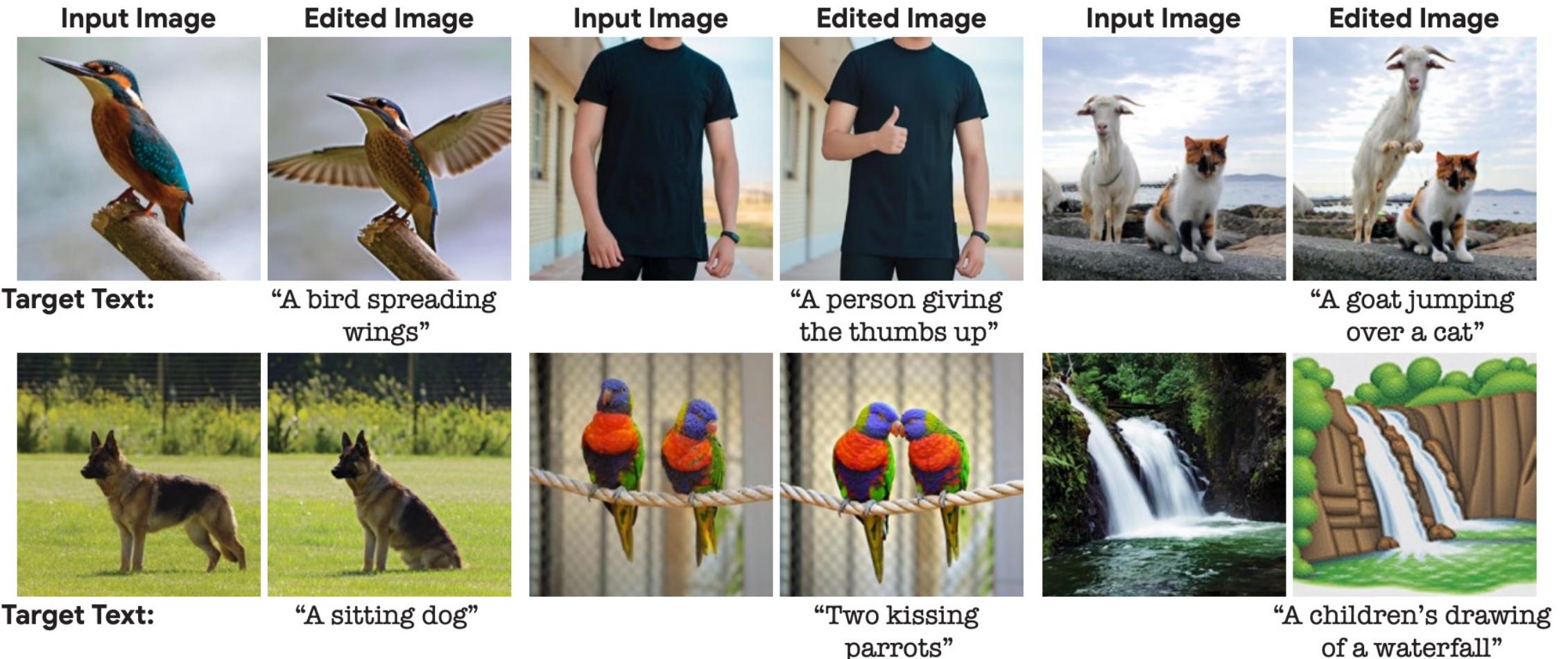


“sofa” $\rightarrow$ “grass”

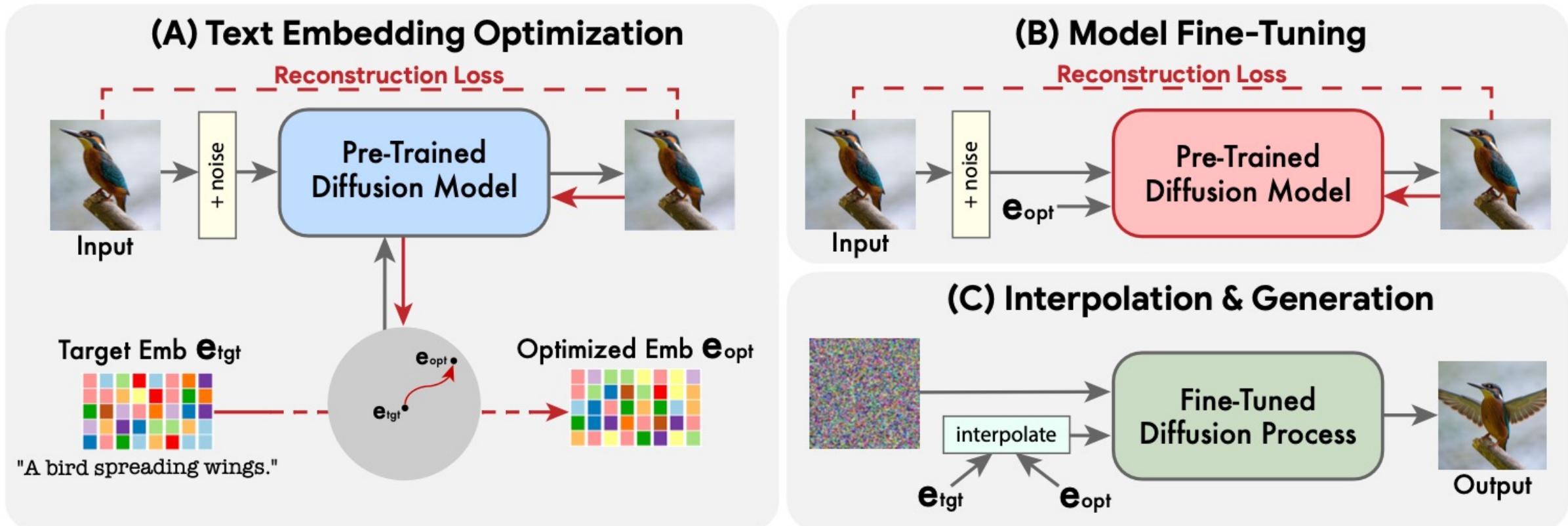


“sofa” $\rightarrow$ “ball pit”

# Imagic: Text-Based Real Image Editing



# Imagic: Text-Based Real Image Editing



# Imagic: Text-Based Real Image Editing

**Text embedding optimization** The target text is first passed through a text encoder [43], which outputs its corresponding text embedding  $\mathbf{e}_{tgt} \in \mathbb{R}^{T \times d}$ , where  $T$  is the number of tokens in the given target text, and  $d$  is the token embedding dimension. We then freeze the parameters of the generative diffusion model  $f_\theta$ , and optimize the target text embedding  $\mathbf{e}_{tgt}$  using the denoising diffusion objective [19]:

$$\mathcal{L}(\mathbf{x}, \mathbf{e}, \theta) = \mathbb{E}_{t, \epsilon} \left[ \|\epsilon - f_\theta(\mathbf{x}_t, t, \mathbf{e})\|_2^2 \right], \quad (2)$$

where  $t \sim Uniform[1, T]$ ,  $\mathbf{x}_t$  is a noisy version of  $\mathbf{x}$  (the input image) obtained using  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$  and [Equation 1](#), and  $\theta$  are the pre-trained diffusion model weights. This results in a text embedding that matches our input image as closely as possible. We run this process for relatively few steps, in order to remain close to the initial target text embedding, obtaining  $\mathbf{e}_{opt}$ . This proximity enables meaningful linear

**Text embedding interpolation** Since the generative diffusion model was trained to fully recreate the input image  $\mathbf{x}$  at the optimized embedding  $\mathbf{e}_{opt}$ , we use it to apply the desired edit by advancing in the direction of the target text embedding  $\mathbf{e}_{tgt}$ . More formally, our third stage is a simple linear interpolation between  $\mathbf{e}_{tgt}$  and  $\mathbf{e}_{opt}$ . For a given hyperparameter  $\eta \in [0, 1]$ , we obtain

$$\bar{\mathbf{e}} = \eta \cdot \mathbf{e}_{tgt} + (1 - \eta) \cdot \mathbf{e}_{opt}, \quad (3)$$

# Imagic: Text-Based Real Image Editing

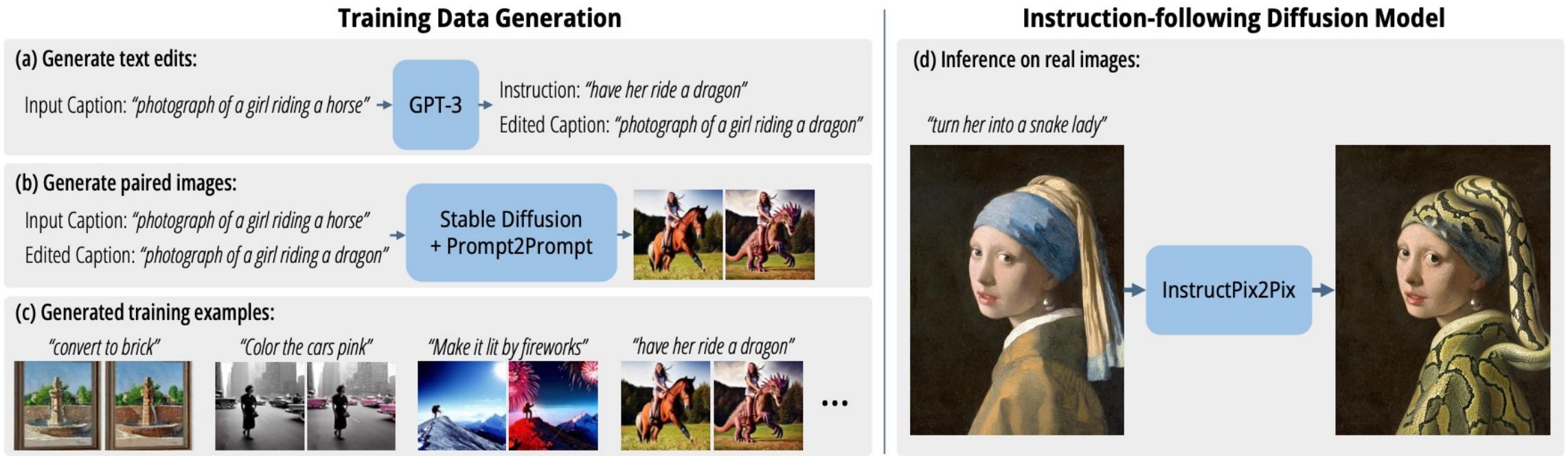


Figure 7. **Embedding interpolation.** Varying  $\eta$  with the same seed, using the pre-trained (top) and fine-tuned (bottom) models.

# InstructPix2Pix



# InstructPix2Pix



# InstructPix2Pix

	<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>
	...	...	...

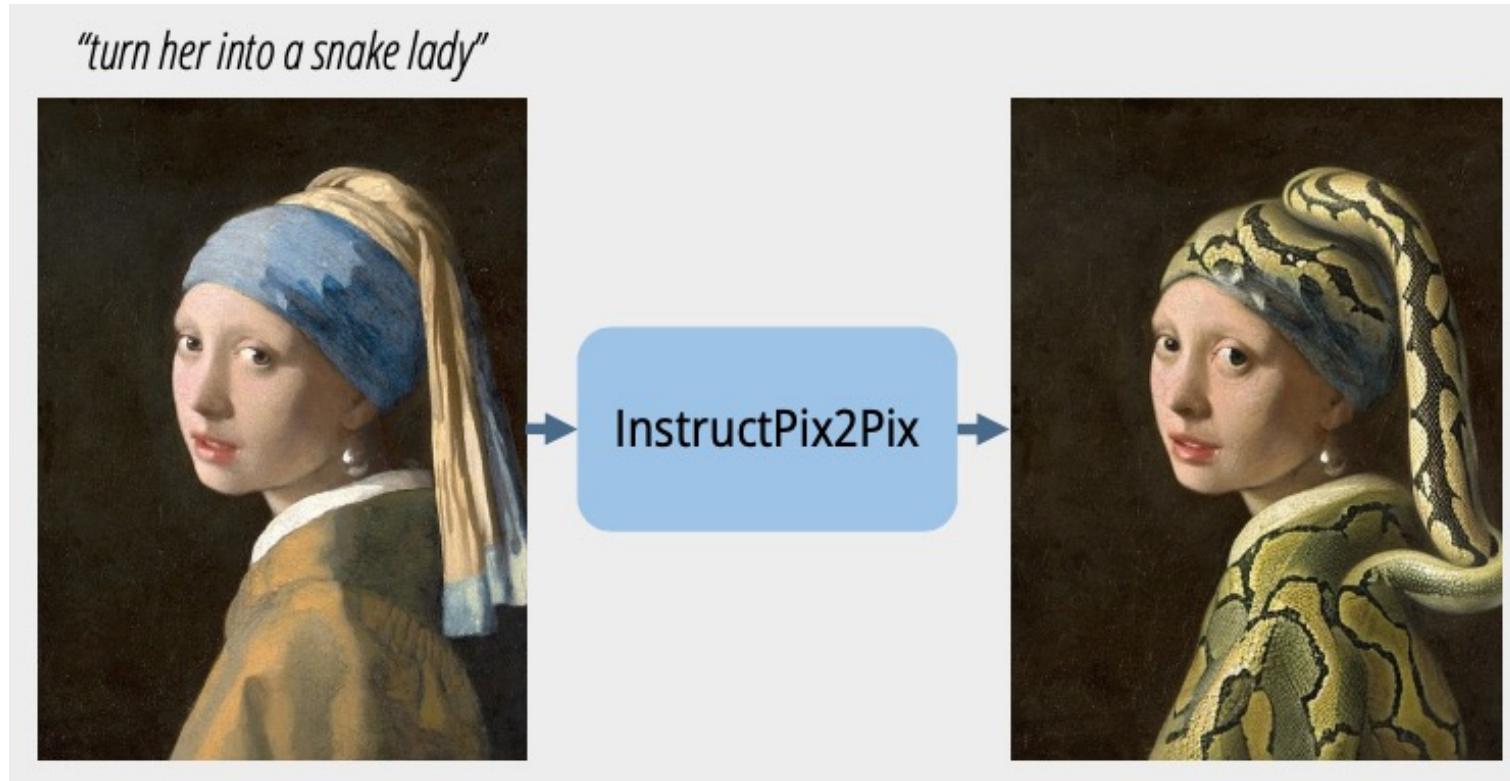
# InstructPix2Pix



**(a) Without Prompt-to-Prompt.**

**(b) With Prompt-to-Prompt.**

# InstructPix2Pix



# InstructPix2Pix



Figure 5. *Mona Lisa* transformed into various artistic mediums.



Figure 6. *The Creation of Adam* with new context and subjects (generated at 768 resolution).

# InstructPix2Pix



Figure 11. Applying our model recurrently with different instructions results in compounded edits.

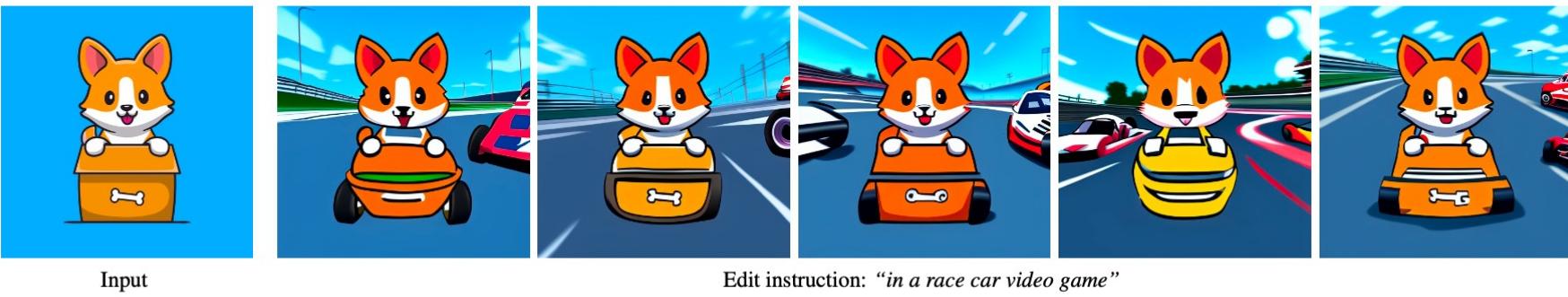
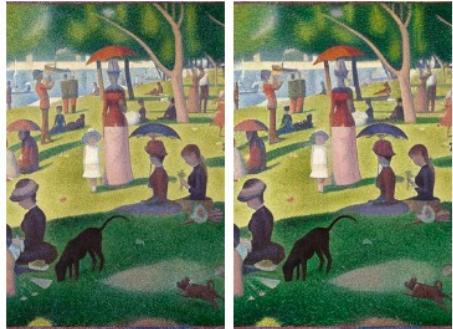


Figure 12. By varying the latent noise, our model can produce many possible image edits for the same input image and instruction.

# InstructPix2Pix



*"Zoom into the image"*



*"Move it to Mars"*



*"Color the tie blue"*



*"Have the people swap places"*

Figure 13. Failure cases. Left to right: our model is not capable of performing viewpoint changes, can make undesired excessive changes to the image, can sometimes fail to isolate the specified object, and has difficulty reorganizing or swapping objects with each other.

# Other Works on Real Image Editing

**“UniTune: Text-Driven Image Editing by Fine Tuning an Image Generation Model on a Single Image” (2022)**

**“SINE: SIngle Image Editing with Text-to-Image Diffusion Models” (2022)**

**“Zero-shot Image-to-Image Translation” (2023)**