

# Diffusion Models for Visual Computing

Niloy Mitra, Daniel Cohen-Or, Minhyuk Sung,  
Chun-Hao Huang, Duygu Ceylan, Paul Guerrero

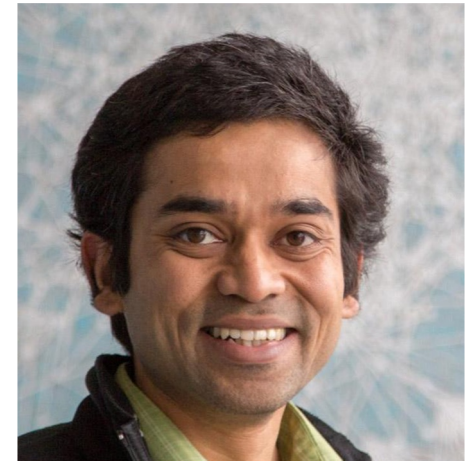


## Part 1: Introduction to Diffusion Models



[https://geometry.cs.ucl.ac.uk/courses/diffusion4VC\\_eg24/](https://geometry.cs.ucl.ac.uk/courses/diffusion4VC_eg24/)

# People



Niloy Mitra



Duygu Ceylan



Daniel Cohen-Or



ChunHao Huang



Paul Guerrero



Minhyuk Sung

# Why do we need this Tutorial?

What are **diffusion model**?

What are the **design choices**?

**Controls** and **adaptation** in the context of Visual Computing

*Learn together*

# Related Materials

- Survey papers
- Past tutorials/courses
- Blogs and recorded videos

# Presentation Schedule

Introduction to Diffusion Models

Guidance and Conditioning Sampling

Attention

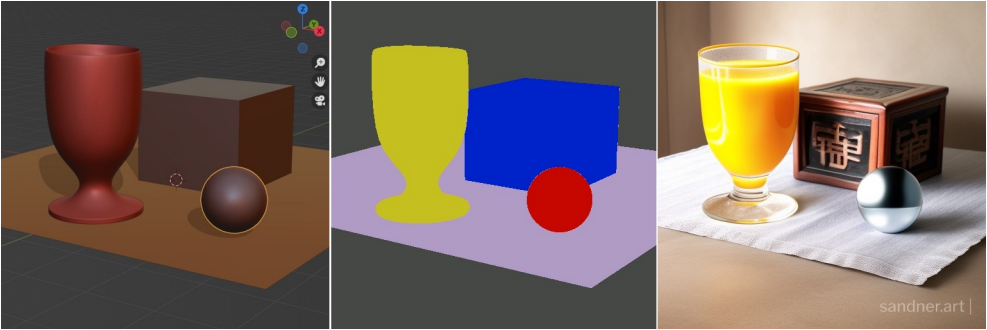
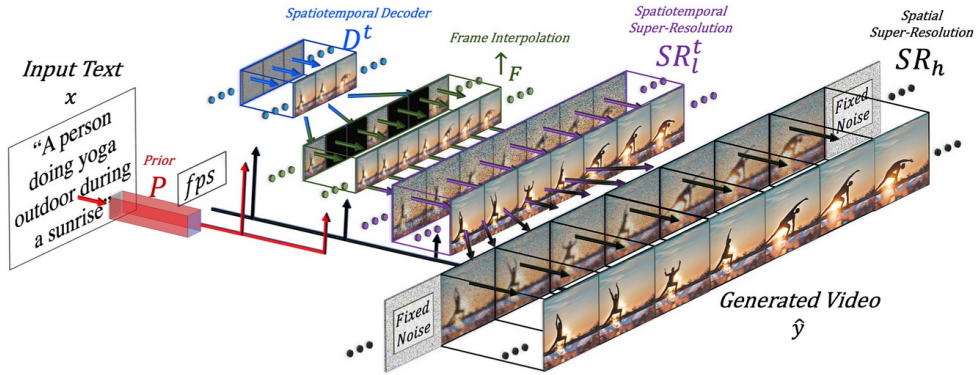
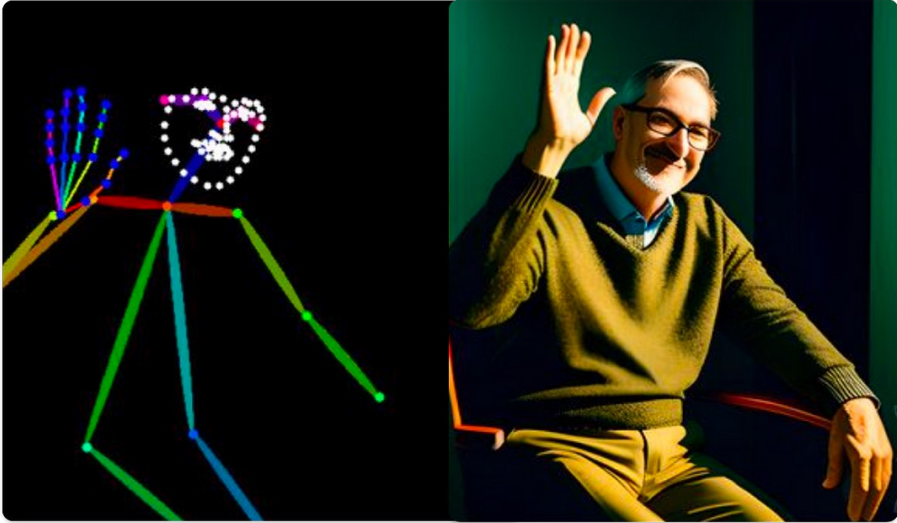
Break

Personalization and Editing

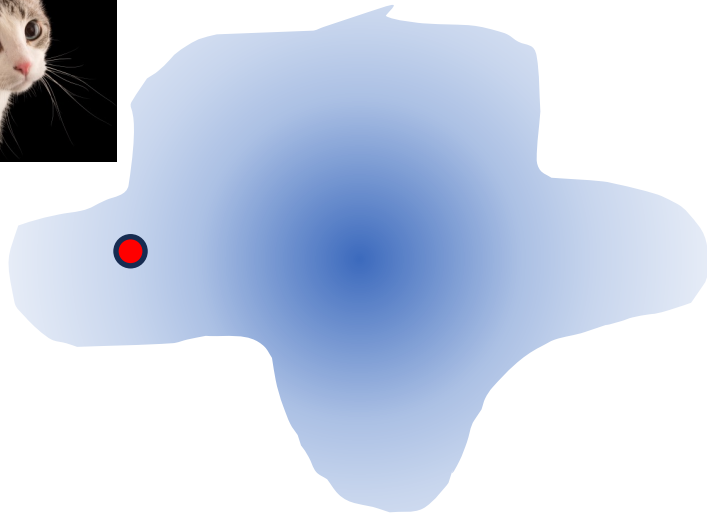
Beyond Single (RGB) Image Generation

Diffusion Models for 3D Generation

# Applications

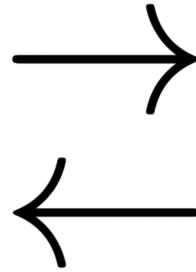


# What is a Diffusion Process?



(unknown) data distribution

unknown map



known distribution

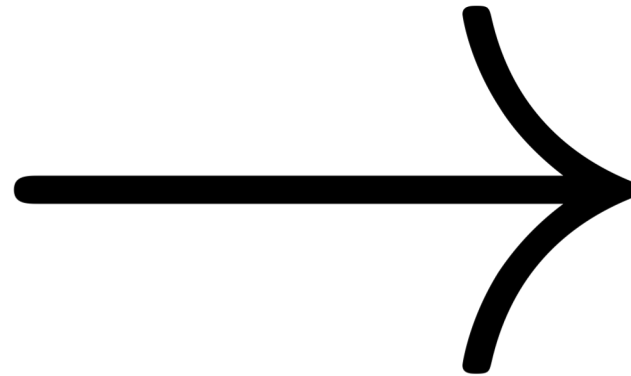
*Sampling  $\Leftrightarrow$  (Unconditional generation)*

# Mapping between Distributions



$\mathbf{X}_0$

data  
distribution



$\mathbf{X}_T$

$\mathcal{N}(\mathbf{0}, \mathbb{I})$

known  
distribution

# Gaussian (Normal) Distribution

- Uniquely defined by **Mean** and **Variance**

$$\mu, \Sigma$$

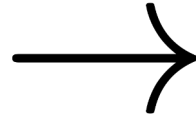
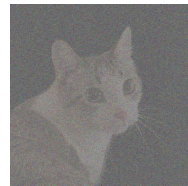
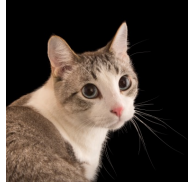
- Reparameterization 'trick'

$$\mathcal{N}(\mu, \Sigma)$$

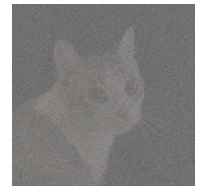
$$x \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

- Many results on combining Gaussian distributions

# Mapping in Many Steps



ward mapping



$\mathbf{X}_0 \rightarrow$

$\mathbf{X}_{t-1} \rightarrow \mathbf{X}_t$

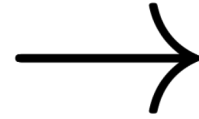
$\rightarrow$

$\rightarrow \mathbf{X}_T$

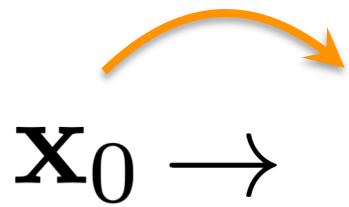
$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbb{I})$$

$$\mathcal{N}(\mathbf{0}, \mathbb{I})$$

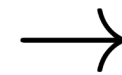
# Mapping in Many Steps



forward mapping



$\mathbf{X}_{t-1} \rightarrow \mathbf{X}_t$

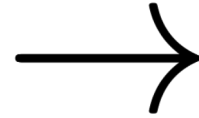


$\rightarrow \mathbf{X}_T$

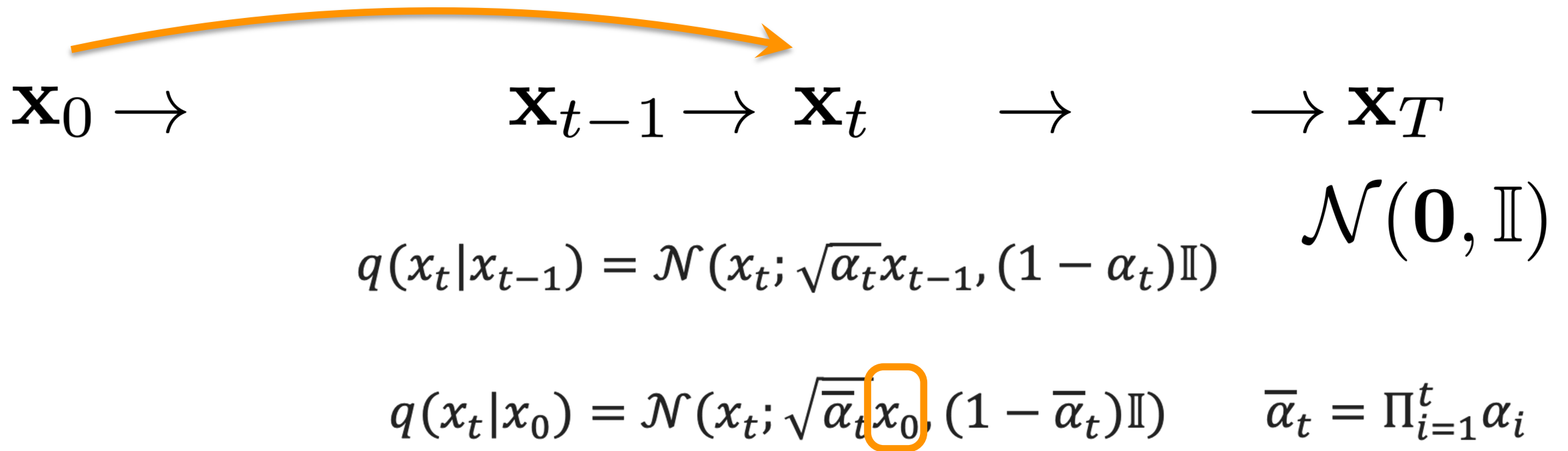
$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \sqrt{\alpha_t}x_{t-1}, (1 - \alpha_t)\mathbb{I})$$

$$\mathcal{N}(\mathbf{0}, \mathbb{I})$$

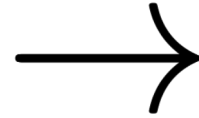
# Mapping in Many Steps



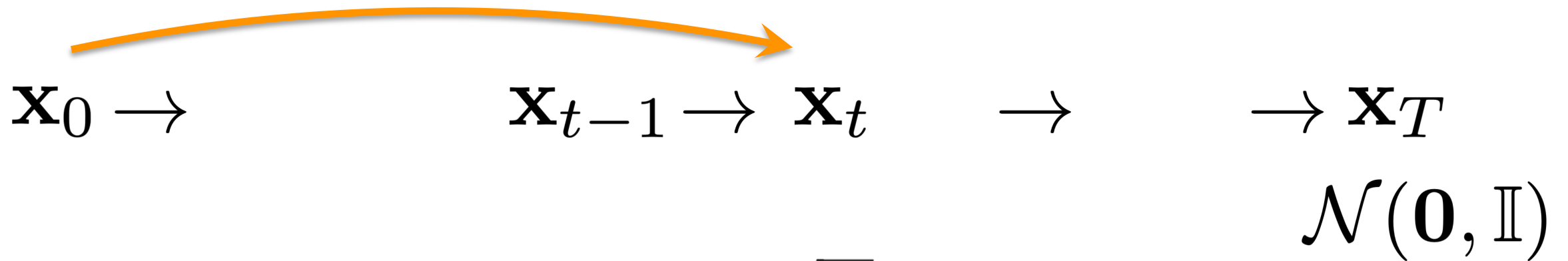
forward mapping



# Mapping in Many Steps



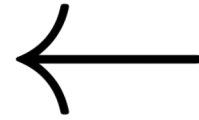
forward mapping



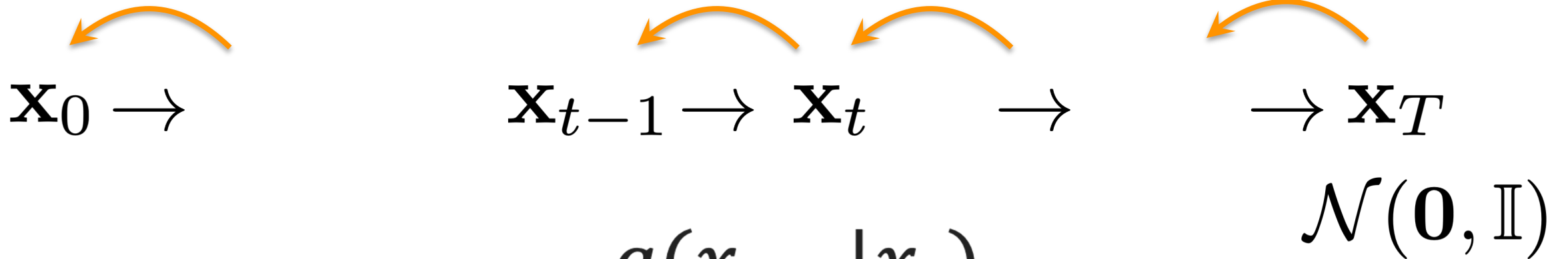
$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbb{I})$$

$$\boxed{x_t} = \hat{\epsilon}(x_0) = \sqrt{\bar{\alpha}_t}\boxed{x_0} + \sqrt{(1 - \bar{\alpha}_t)}\boxed{\epsilon_t} \quad \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

# Generative Modeling: Sampling



reverse mapping



$$q(x_{t-1} | x_t)$$

$\approx$

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{(1 - \bar{\alpha}_t)} \epsilon_t$$

$$p(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}} D_\theta(x_t, t), (1 - \bar{\alpha}_{t-1}) \mathbb{I})$$

# Loss Functions

$$\mathcal{L}_{simple}(\theta) = \mathbb{E}_{t, x_0, \epsilon} [C_t \| \epsilon_{\theta}(x_t, t) - \epsilon \|^2]$$

$$\mathcal{L}(\theta) = \mathbb{E}_{t, \epsilon, x_0} \left[ C_t \| D_{\theta}(\hat{\epsilon}_t(x_0), t) - x_0 \|^2 \right]$$

$$p(x_{t-1} | x_t) = \mathcal{N}(x_{t-1}; \sqrt{\bar{\alpha}_{t-1}} D_{\theta}(x_t, t), (1 - \bar{\alpha}_{t-1}) \mathbb{I})$$

# Algorithm (How to Train?)

---

## Algorithm 1 Training

---

1: **repeat**

2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$

3:  $t \sim \text{Uniform}(\{1, \dots, T\})$

4:  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

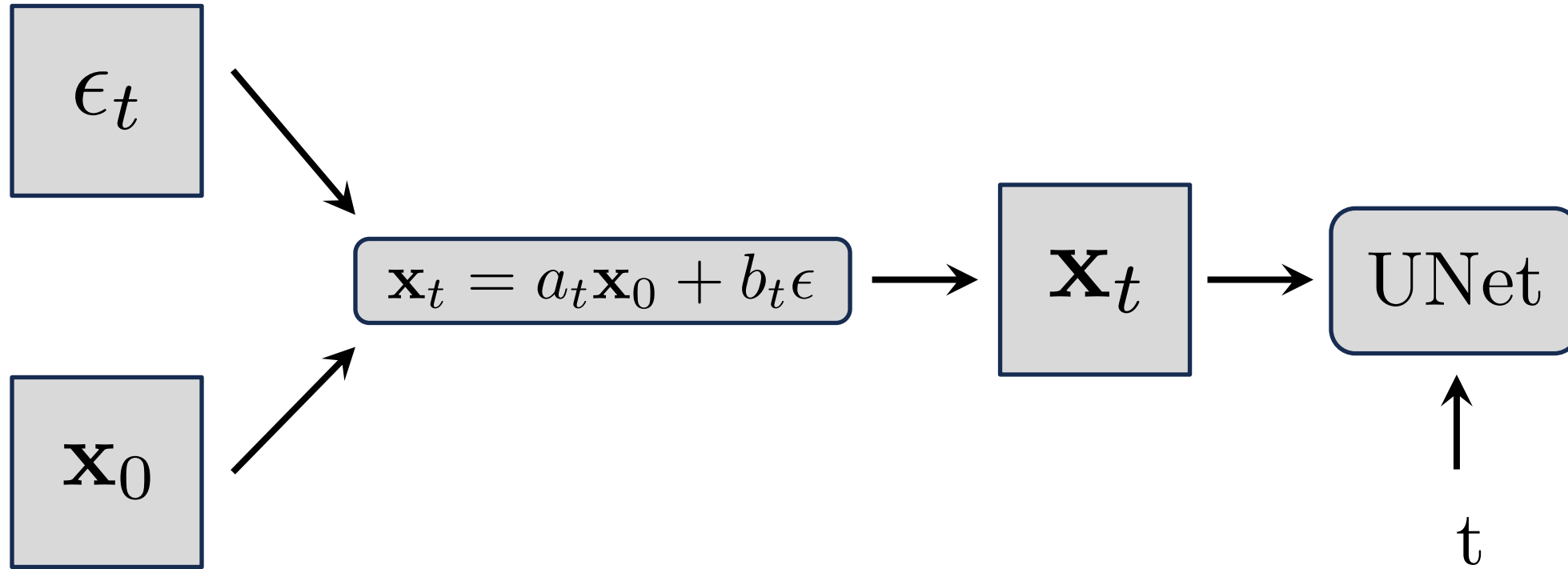
5: Take gradient descent step on

$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}, t)\|^2$$

6: **until** converged

---

# Training Loss



# Three Interpretations

- Predict Noise  $\epsilon_t$

- Predict clean image  $\mathbf{x}_0$

- Score-based optimization  $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_0) = -\frac{\epsilon}{\sqrt{1 - \bar{\alpha}_t}}$

*they are equivalent!!*

# Algorithm (How to Sample?)

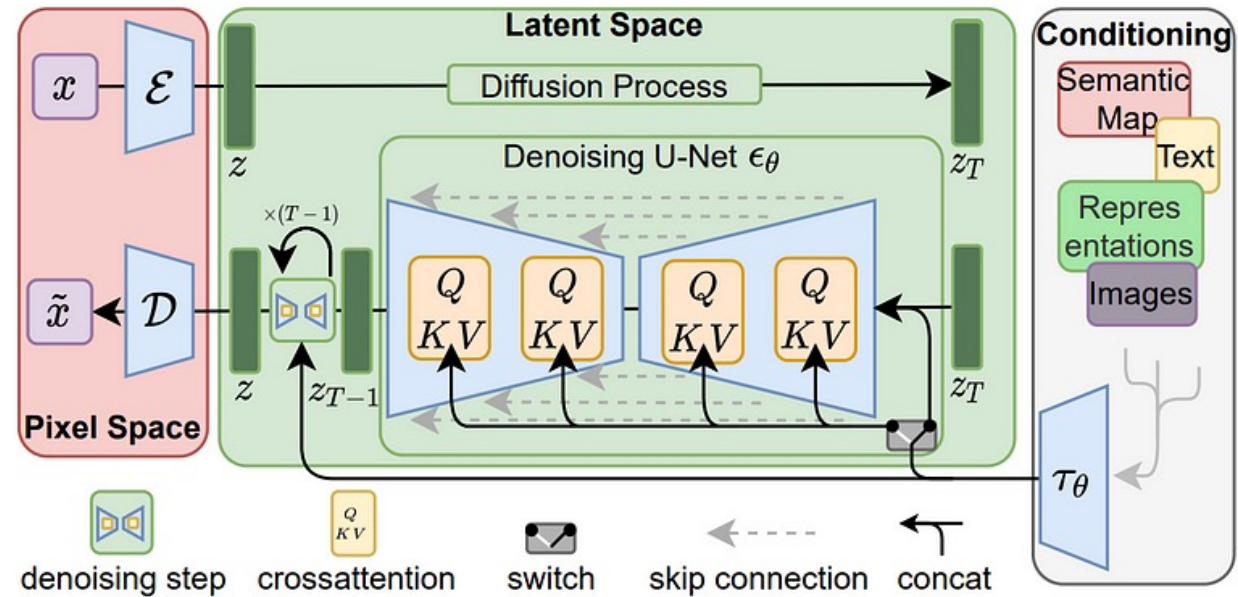
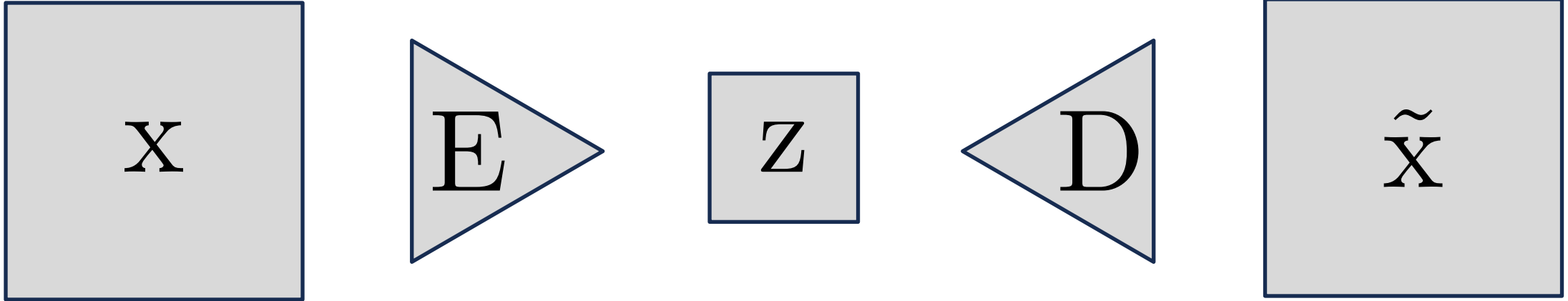
---

## Algorithm 2 Sampling

---

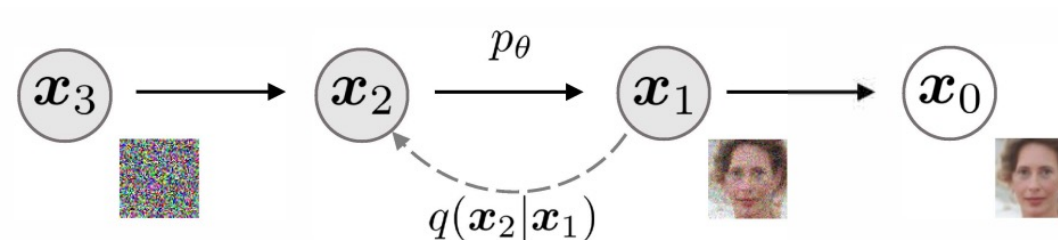
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2: **for**  $t = T, \dots, 1$  **do**
  - 3:    $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
  - 4:    $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
  - 5: **end for**
  - 6: **return**  $\mathbf{x}_0$
-

# Latent Diffusion Model

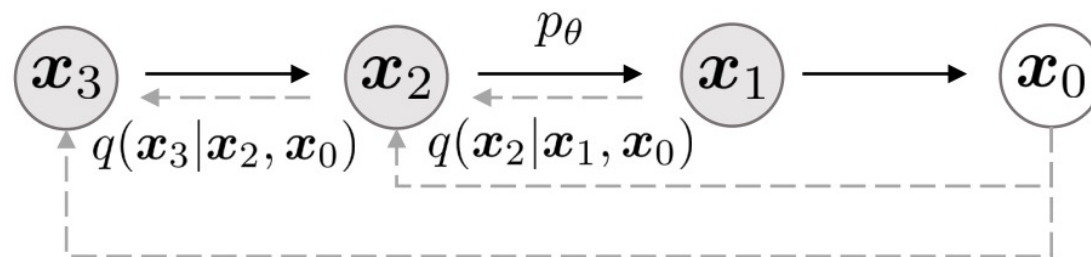


# DDPM vs DDIM

- DDPM: Markovian process



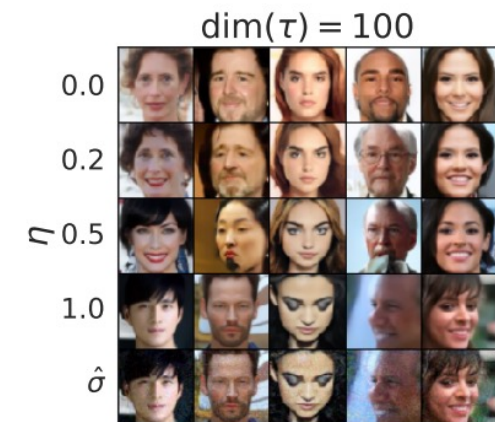
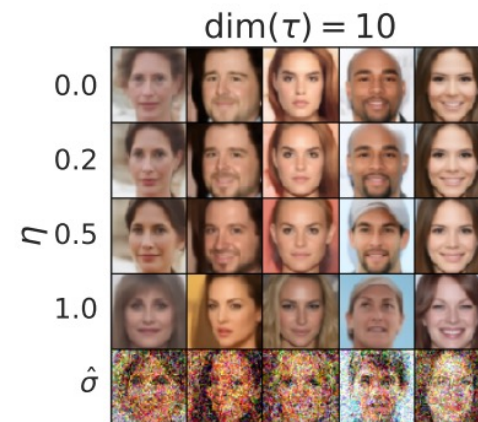
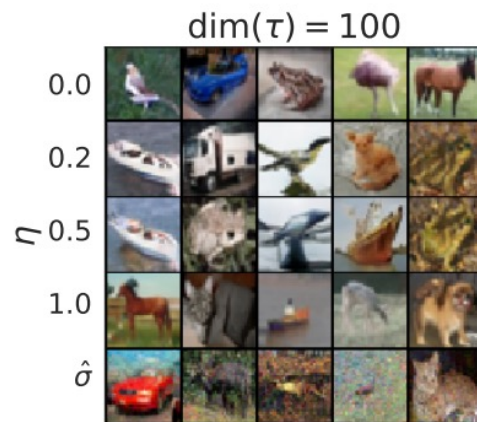
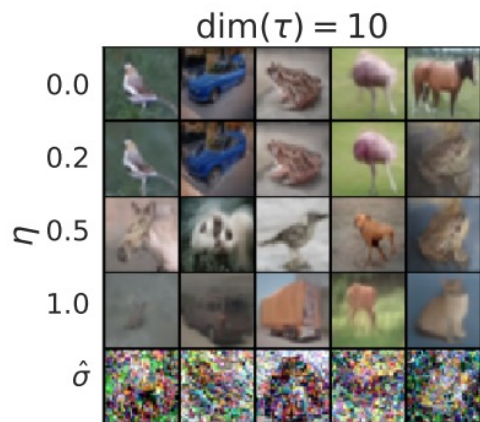
- DDIM: Non-Markovian process but 10-50x faster!!
  - Trained w/ pretrained DDPM diffusion



$$\mathbf{x}_{t-1} = \underbrace{\sqrt{\alpha_{t-1}} \left( \frac{\mathbf{x}_t - \sqrt{1 - \alpha_t} \epsilon_\theta^{(t)}(\mathbf{x}_t)}{\sqrt{\alpha_t}} \right)}_{\text{"predicted } \mathbf{x}_0"} + \underbrace{\sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta^{(t)}(\mathbf{x}_t)}_{\text{"direction pointing to } \mathbf{x}_t"} + \underbrace{\sigma_t \epsilon_t}_{\text{random noise}}$$

# DDPM vs DDIM

| $S$            | CIFAR10 ( $32 \times 32$ ) |             |             |             |             | CelebA ( $64 \times 64$ ) |              |             |             |             |
|----------------|----------------------------|-------------|-------------|-------------|-------------|---------------------------|--------------|-------------|-------------|-------------|
|                | 10                         | 20          | 50          | 100         | 1000        | 10                        | 20           | 50          | 100         | 1000        |
| $\eta$ 0.0     | <b>13.36</b>               | <b>6.84</b> | <b>4.67</b> | <b>4.16</b> | 4.04        | <b>17.33</b>              | <b>13.73</b> | <b>9.17</b> | <b>6.53</b> | 3.51        |
| $\eta$ 0.2     | 14.04                      | 7.11        | 4.77        | 4.25        | 4.09        | 17.66                     | 14.11        | 9.51        | 6.79        | 3.64        |
| $\eta$ 0.5     | 16.66                      | 8.35        | 5.25        | 4.46        | 4.29        | 19.86                     | 16.06        | 11.01       | 8.09        | 4.28        |
| $\eta$ 1.0     | 41.07                      | 18.36       | 8.01        | 5.78        | 4.73        | 33.12                     | 26.03        | 18.48       | 13.93       | 5.98        |
| $\hat{\sigma}$ | 367.43                     | 133.37      | 32.72       | 9.99        | <b>3.17</b> | 299.71                    | 183.83       | 71.71       | 45.20       | <b>3.26</b> |



# Summary so far

---

## Algorithm 1 Training

---

- 1: **repeat**
  - 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
  - 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
  - 4:  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 5: Take gradient descent step on  
$$\nabla_{\theta} \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t)\|^2$$
  - 6: **until** converged
- 

---

## Algorithm 2 Sampling

---

- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
  - 2: **for**  $t = T, \dots, 1$  **do**
  - 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
  - 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
  - 5: **end for**
  - 6: **return**  $\mathbf{x}_0$
-

# Presentation Schedule

Introduction to Diffusion Models

Guidance and Conditioning Sampling

Attention

Break

Personalization and Editing

Beyond Single (RGB) Image Generation

Diffusion Models for 3D Generation

# Diffusion Models for Visual Content Generation

## Guidance

**Presenter: Minhyuk Sung**

Eurographics 2024 Tutorial

# Guidance Using Additional Information

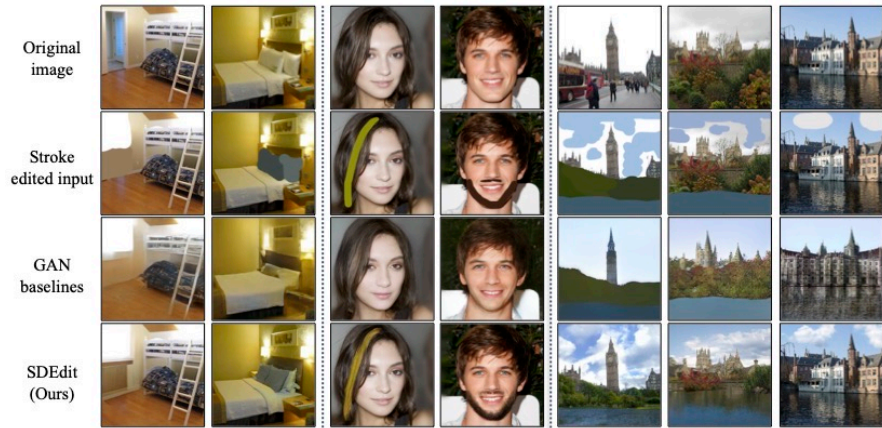
If we have **additional information** about the data, such as class labels for images, can we utilize this information to **improve the quality of generated outputs?**



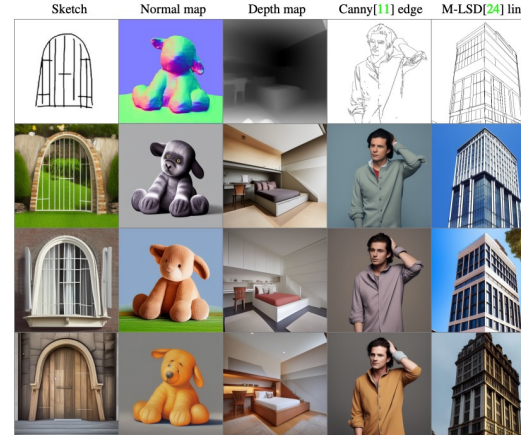
Ho and Salimans, Classifier-Free Diffusion Guidance, NeurIPS 2021 Workshop.

# Zero-Shot / Few-Shot Adaptations

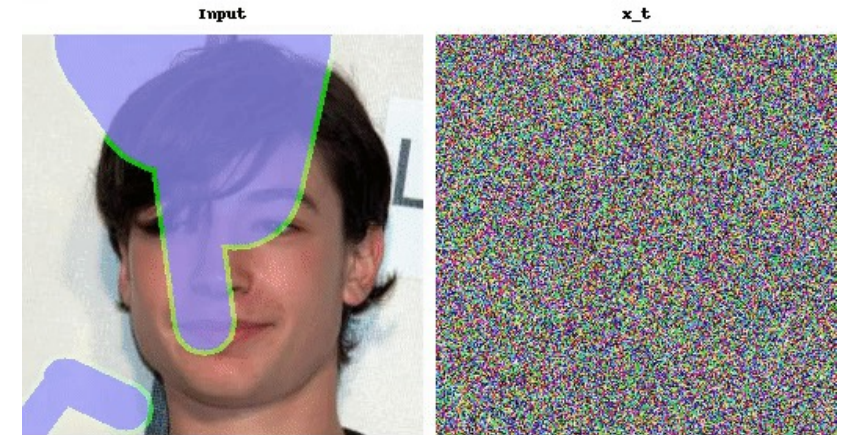
Can we apply a **pretrained** diffusion model to various **conditional generation** setups in a **zero-shot** or **few-shot** manner?



Meng et al., SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, ICLR 2022.



Zhang et al., Adding Conditional Control to Text-to-Image Diffusion Models, ICCV 2022.



Lugmayr et al., RePaint: Inpainting using Denoising Diffusion Probabilistic Models, CVPR 2022.

# Content

1. Classifier Guidance / Classifier-Free Guidance
2. ControlNet
3. SDEdit / RePaint

# **Classifier Guidance / Classifier-Free Guidance (CFG)**

# Guidance Using Additional Information

How to use **class labels** or some **additional information** about the data to **improve the quality of generated outputs**?



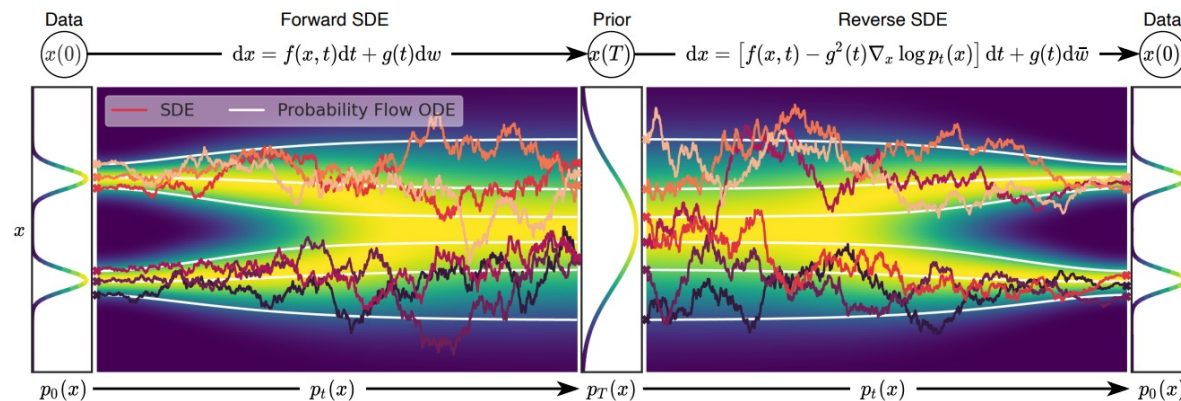
Ho and Salimans, Classifier-Free Diffusion Guidance, NeurIPS 2021 Workshop.

# Recap: Stochastic Differential Equations

In a **continuous** time domain, the **reverse** process is formulated as the following stochastic differential equation:

$$d\mathbf{x} = [\mathbf{f}(\mathbf{x}, t)dt - g^2(t)\nabla_{\mathbf{x}} \log p_t(\mathbf{x})]dt + g(t)d\mathbf{w}$$

**DDPM is a specific discretization of the SDE formulation.**



# Recap: Stochastic Differential Equations

The gradient of the logarithm of the PDF  $\nabla_{\mathbf{x}} \log q(\mathbf{x}_t)$  is referred to as a **score** function.

$$\begin{aligned} \nabla_{\mathbf{x}} \log q(\mathbf{x}_t) &= \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0)} [\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t | \mathbf{x}_0)] \\ &= -\frac{1}{\sqrt{1 - \bar{\alpha}_t}} \mathbb{E}_{\mathbf{x}_0 \sim q(\mathbf{x}_0)} [\boldsymbol{\varepsilon}(\mathbf{x}_t, \mathbf{x}_0)] = -\frac{\boldsymbol{\varepsilon}_{\theta}(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}} \end{aligned}$$

The **noise** predicted by the neural network  $\boldsymbol{\varepsilon}_{\theta}(\mathbf{x}_t, t)$  is the scaled **score**!

# Classifier Guidance

- Assume that data  $\mathbf{x}$  and **class label**  $y$  pairs are given:

$$p(\mathbf{x}_t, y) = p(\mathbf{x}_t)p(y|\mathbf{x}_t)$$

- At each timestep  $t$ , we are interested in the **score** function

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t, y):$$

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t, y) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t)$$

$$= -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) + \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t)$$

$$= -\frac{1}{\sqrt{1-\bar{\alpha}_t}} \left( \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) - \sqrt{1-\bar{\alpha}_t} \nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t) \right)$$

Noise  
predictor

How to  
compute this?

# Classifier Guidance

How to compute  $\nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t)$ ?

- Train a classifier  $p_\phi(y|\mathbf{x}_t)$ , taking noisy data  $\mathbf{x}_t$  (and timestep  $t$ ) as input and classifying it.
- Use  $\nabla_{\mathbf{x}_t} \log p_\phi(y|\mathbf{x}_t)$  as an approximation of  $\nabla_{\mathbf{x}_t} \log p(y|\mathbf{x}_t)$ .

# Classifier Guidance

- Update the noise predictor  $\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t)$  as follows:

$$\bar{\boldsymbol{\varepsilon}}_\theta(\mathbf{x}_t, t) = \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) - \sqrt{1 - \bar{\alpha}_t} \nabla_{\mathbf{x}_t} \log p_\phi(y|\mathbf{x}_t)$$

- The **strength** of the classifier guidance can be controlled by adding a **weight  $w$** :

$$\bar{\boldsymbol{\varepsilon}}_\theta(\mathbf{x}_t, t) = \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, t) - w \sqrt{1 - \bar{\alpha}_t} \nabla_{\mathbf{x}_t} \log p_\phi(y|\mathbf{x}_t)$$

- **Limitation:** Additional training of a classifier is required.

# Classifier-Free Guidance (CFG)

- Jointly train both the **conditional** and **unconditional** diffusion models.
  - For the **unconditional** case, simply feed a **null token  $\emptyset$**  as the condition.
  - In the reverse process, given the condition  $y$ , take  $\boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, y, t)$ .
- 
- Take a **linear combination** of unconditional and conditional noises as follows:

$$\hat{\boldsymbol{\varepsilon}}_\theta(\mathbf{x}_t, y, t) = \underbrace{(1 + w)}_{\text{Increase } \uparrow \text{ the conditional noise}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, y, t) - \underbrace{w}_{\text{Decrease } \downarrow \text{ the unconditional noise}} \boldsymbol{\varepsilon}_\theta(\mathbf{x}_t, \emptyset, t) \quad w \geq 0$$

# Classifier-Free Guidance (CFG)

|                     | Model                                    | FID ( $\downarrow$ )             |
|---------------------|--|----------------------------------|
|                     | BigGAN-deep, max IS (Brock et al., 2019) | 25                               |
|                     | BigGAN-deep (Brock et al., 2019)         | 5.7                              |
|                     | CDM (Ho et al., 2021)                    | 3.52                             |
|                     | LOGAN (Wu et al., 2019)                  | 3.36                             |
| Classifier Guidance | ADM-G (Dhariwal & Nichol, 2021)          | 2.97                             |
|                     | <b>Ours</b>                              | $T = 128 / 256 / 1024$           |
|                     | $w = 0.0$                                | 8.11 / <b>7.27</b> / 7.22        |
|                     | $w = 0.1$                                | 5.31 / 4.53 / 4.5                |
|                     | $w = 0.2$                                | 3.7 / 3.03 / 3                   |
|                     | $w = 0.3$                                | 3.04 / <b>2.43</b> / <b>2.43</b> |
|                     | $w = 0.4$                                | 3.02 / 2.49 / 2.48               |
|                     | $w = 0.5$                                | 3.43 / 2.98 / 2.96               |
|                     | $w = 0.6$                                | 4.09 / 3.76 / 3.73               |
|                     | $w = 0.7$                                | 4.96 / 4.67 / 4.69               |
|                     | $w = 0.8$                                | 5.93 / 5.74 / 5.71               |
|                     | $w = 0.9$                                | 6.89 / 6.8 / 6.81                |
|                     | $w = 1.0$                                | 7.88 / 7.86 / 7.8                |
|                     | $w = 2.0$                                | 15.9 / 15.93 / 15.75             |
|                     | $w = 3.0$                                | 19.77 / 19.77 / 19.56            |
|                     | $w = 4.0$                                | 21.55 / 21.53 / 21.45            |

$T = 256$   
ImageNet Results

# Classifier-Free Guidance (CFG)

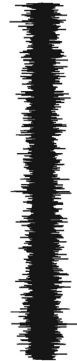
- (+) The classifier does not need to be trained.
- (+) It is more versatile and can be used not only for class labels but also for any additional information (e.g., text descriptions).
- (−) In the generation process, the noise predictor needs to be evaluated twice.
- (−) Determining the optimal weight  $w$  can be challenging.

# Example: Audio Conditioned Generation

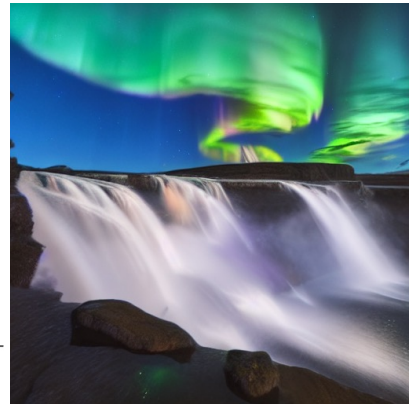


# Example: Audio Conditioned Generation

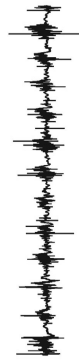
Waterfall Burbling



**Text Prompt:** "Aurora Borealis Lights"



Snow



**Text Prompt:** "Lego"



Forest



**Text Prompt:** "Surreal Dreamscapes"



# Example: Two-Hand Interaction Generation

One hand is generated based on the condition of the other hand.



# **ControlNet: Few-Shot Adaptation**

# LAION-5B: A NEW ERA OF OPEN LARGE-SCALE MULTI-MODAL DATASETS

by: Romain Beaumont, 31 Mar, 2022

We present a dataset of **5.85 billion** CLIP-filtered image-text pairs, 14x bigger than LAION-400M, previously the biggest openly accessible image-text dataset in the world - see also our [NeurIPS2022 paper](#)


Authors: Christoph Schuhmann, Richard Vencu, Romain Beaumont, Theo Coombes, Cade Gordon, Aarush Katta, Robert Kaczmarczyk, Jenia Jitsev

Backend url:  french cat

Index:

[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions   
Display full captions   
Display similarities   
Safe mode   
Hide duplicate urls   
Hide (near) duplicate images   
Search over   
Search with



french cat

french cat

How to tell if your feline is french. He wears a b...

イケメン猫モデル「トキ・ナンタケツト」がかっこいい - NAVER まとめ

Hilarious pics of funny cats! funnycatsgif.com

cat in a suit Georgian sells tomatoes

網友挑戰「加幾筆畫

# ControlNet [Zhang et al., 2023]

- Should we have 5 billion **input-output pairs** to train conditional image diffusion models?
- Should we have the new dataset for every type of condition?

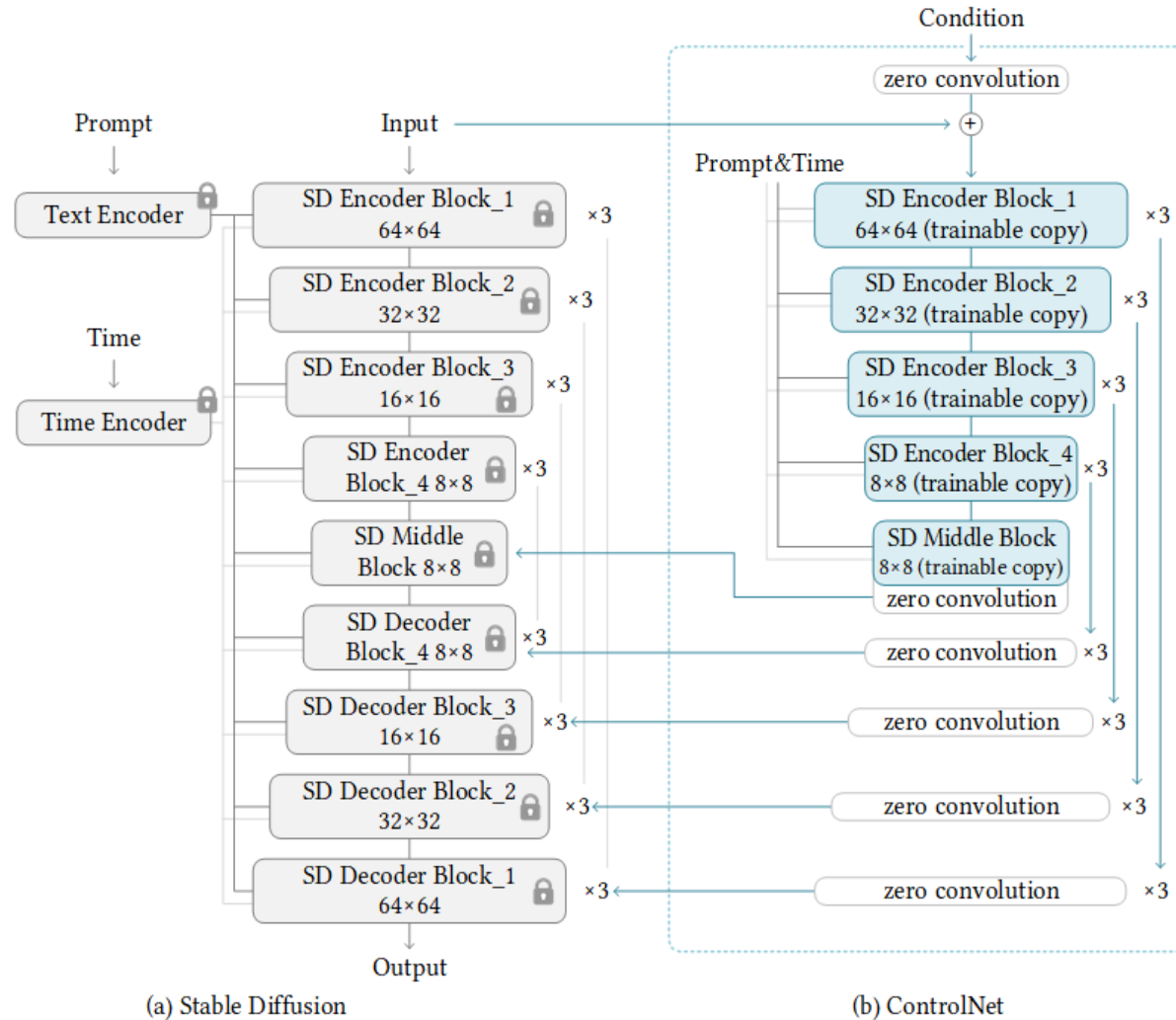


# ControlNet [Zhang et al., 2023]

How to **convert** a pretrained unconditional image diffusion model into an **image-conditioned generative model** using a relatively **smaller set of input-output pairs** ( $\sim 100k$ )?



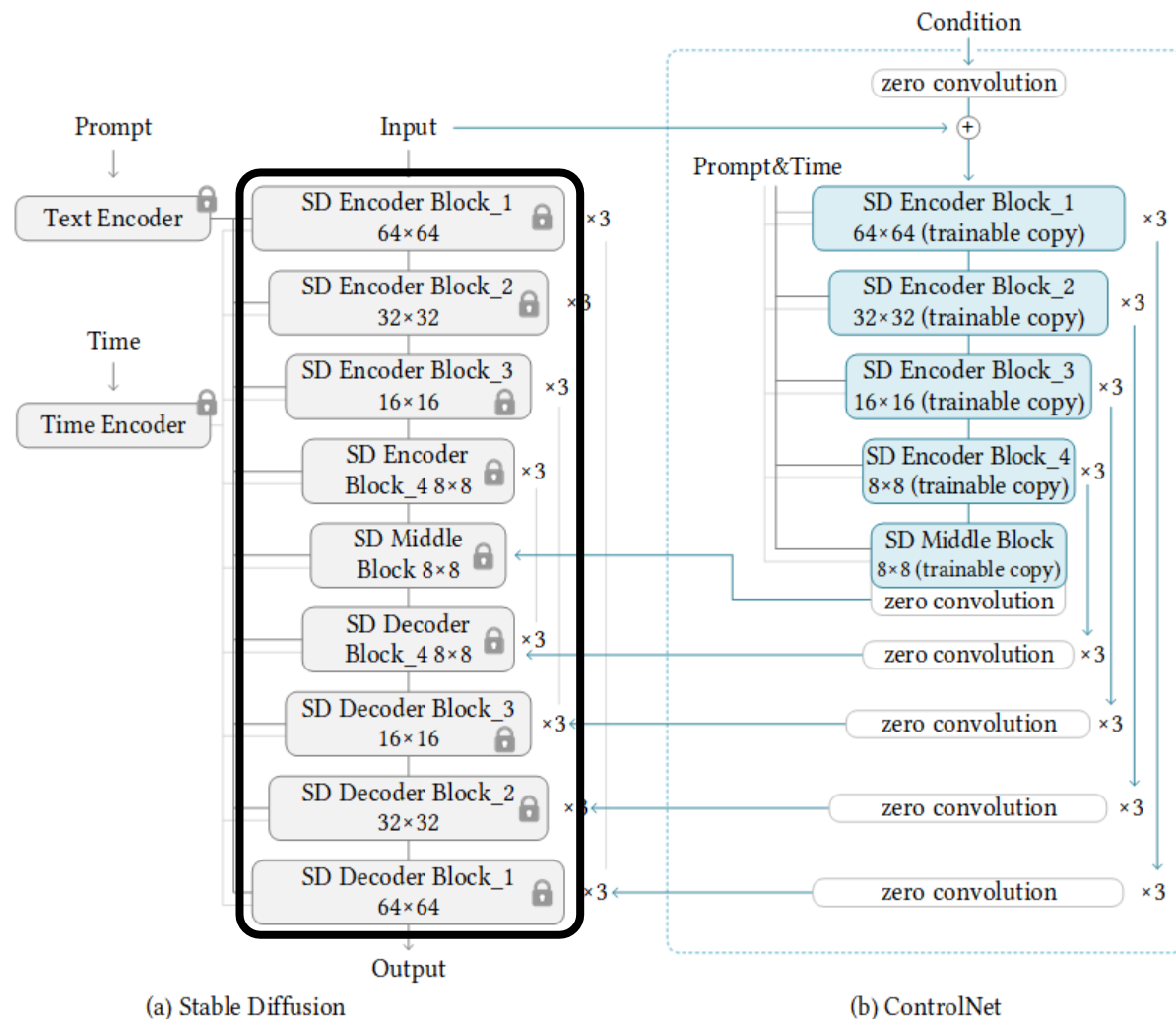
# ControlNet [Zhang et al., 2023]



## Key Idea:

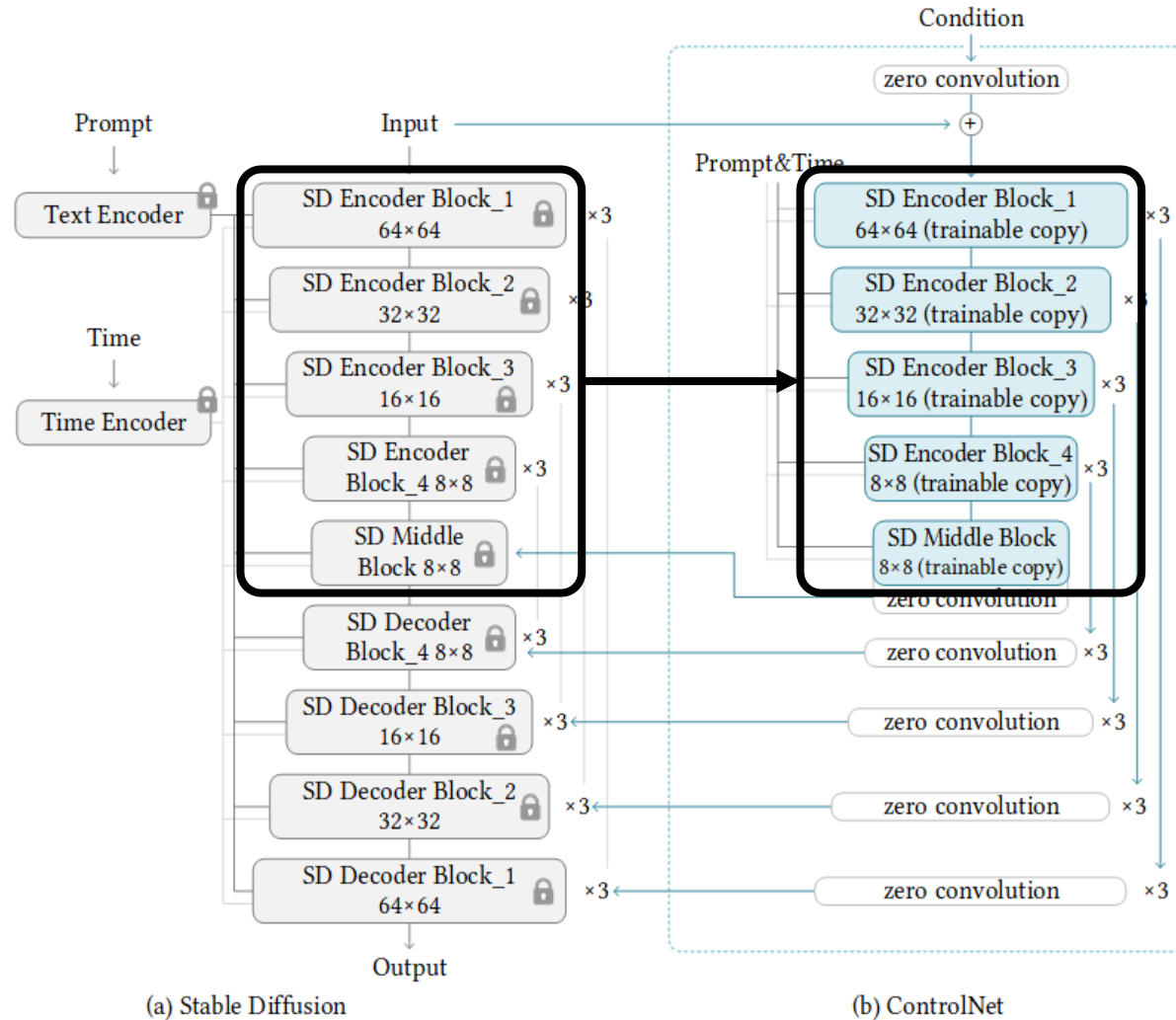
Fully leverage the pretrained noise prediction network to process the conditional image.

# ControlNet [Zhang et al., 2023]



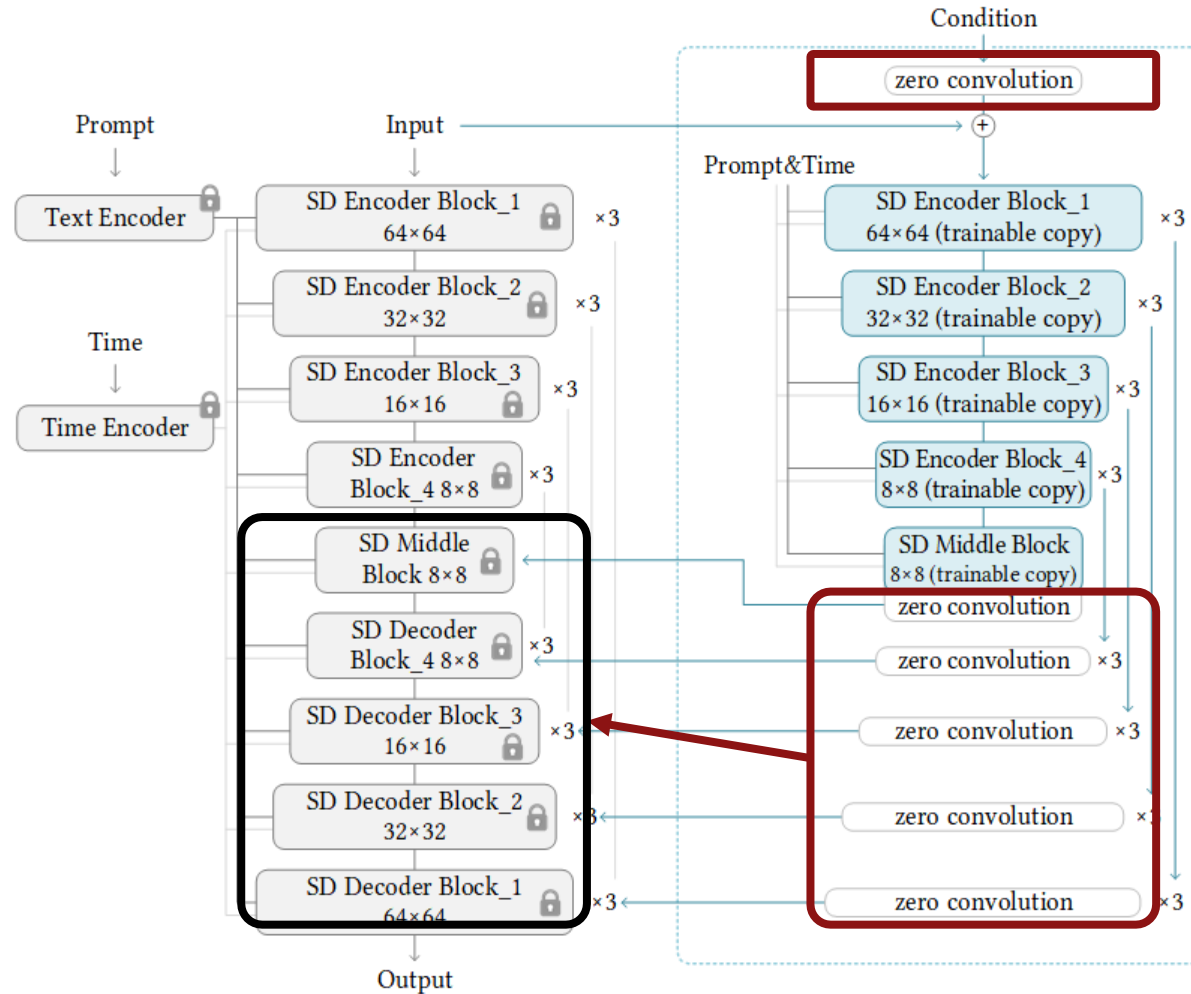
1. Freeze the noise prediction network.

# ControlNet [Zhang et al., 2023]



- For the encoding of the conditional image, copy the pretrained encoder parameters while allowing them to be updated during fine-tuning.

# ControlNet [Zhang et al., 2023]



(a) Stable Diffusion

(b) ControlNet

- Combine the encoded conditional image information with the noisy image using **zero convolution**.

# ControlNet [Zhang et al., 2023]

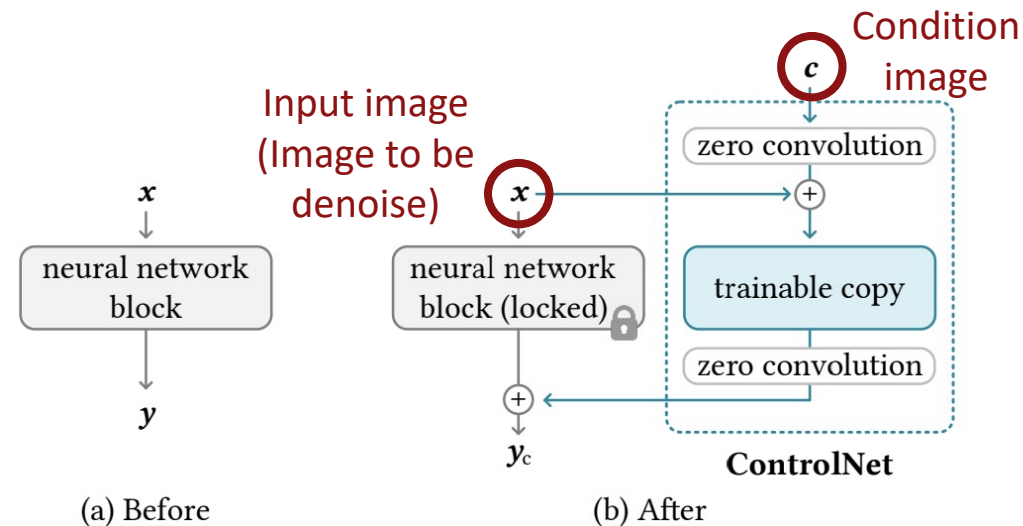
Zero Convolution  $Z$  is a  $1 \times 1$  convolution layer with learnable weight (scaling) parameters  $\mathbf{a}$  and bias (offset) parameters  $\mathbf{b}$ , both of which are **initialized with zero**:

$$Z(\mathbf{x}; \mathbf{a}, \mathbf{b}) = \mathbf{a} \cdot \mathbf{x} + \mathbf{b}$$

# ControlNet [Zhang et al., 2023]

ControlNet modifies each neural network block using zero convolution as follows:

$$y_c = F(x; \Theta) + Z(F(x + Z(c; a_1, b_1); \Theta_c); a_2, b_2)$$

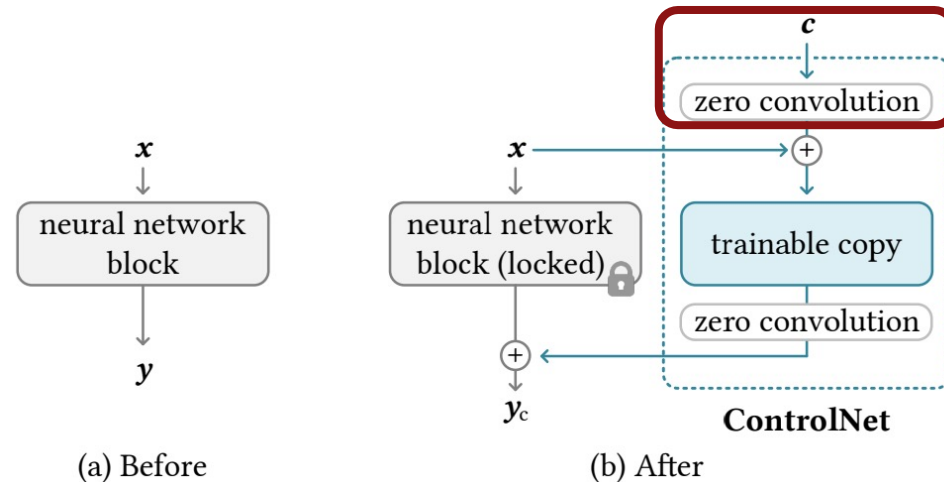


# ControlNet [Zhang et al., 2023]

ControlNet modifies each neural network block using zero convolution as follows:

Zero in the beginning since  $\mathbf{a}_1 = \mathbf{b}_1 = \mathbf{0}$ .

$$\mathbf{y}_c = F(\mathbf{x}; \Theta) + Z(F(\mathbf{x} + \boxed{Z(\mathbf{c}; \mathbf{a}_1, \mathbf{b}_1)}; \Theta_c); \mathbf{a}_2, \mathbf{b}_2)$$

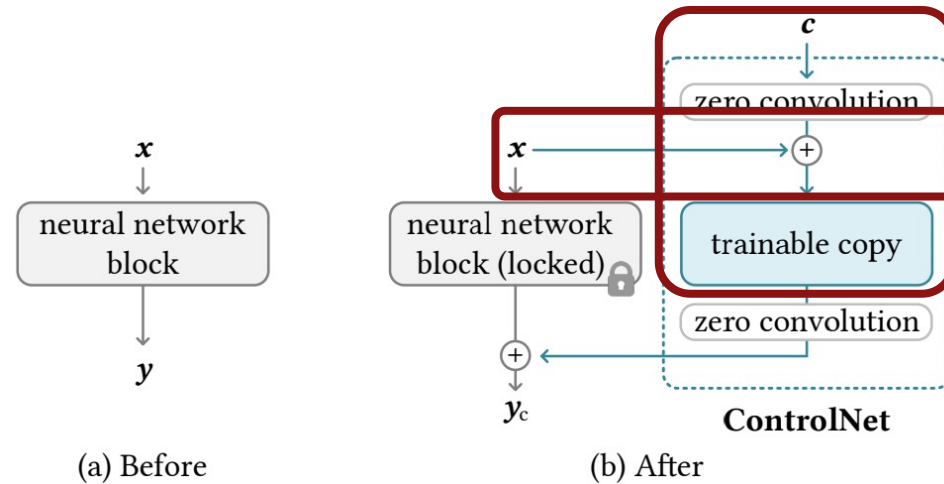


# ControlNet [Zhang et al., 2023]

ControlNet modifies each neural network block using zero convolution as follows:

Same as  $F(x; \Theta)$  in the beginning since  $\Theta_c = \Theta$ .

$$y_c = F(x; \Theta) + Z(F(x + Z(c; a_1, b_1); \Theta_c); a_2, b_2)$$

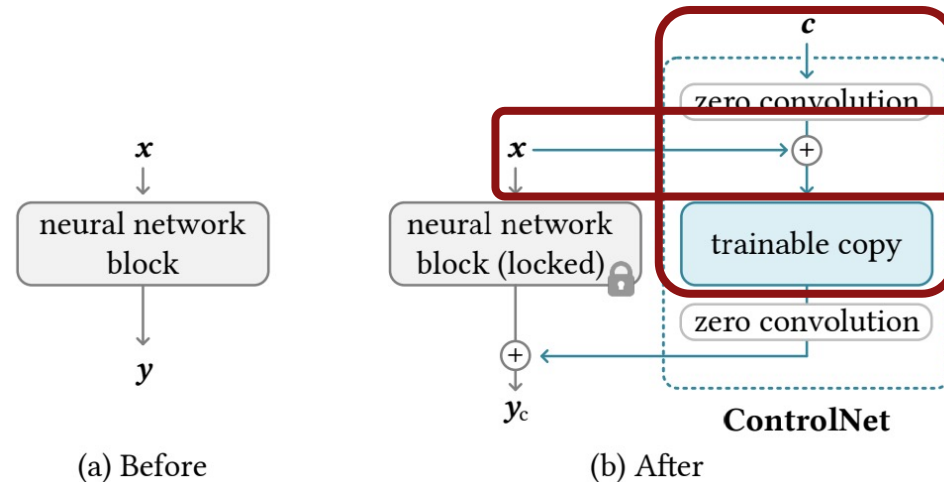


# ControlNet [Zhang et al., 2023]

For encoding of the conditional image  $c$ , **begin with the state where the input is  $x$ , while gradually incorporating  $c$ .**

Same as  $F(x; \Theta)$  in the beginning since  $\Theta_c = \Theta$ .

$$y_c = F(x; \Theta) + Z(F(x + Z(c; a_1, b_1); \Theta_c); a_2, b_2)$$

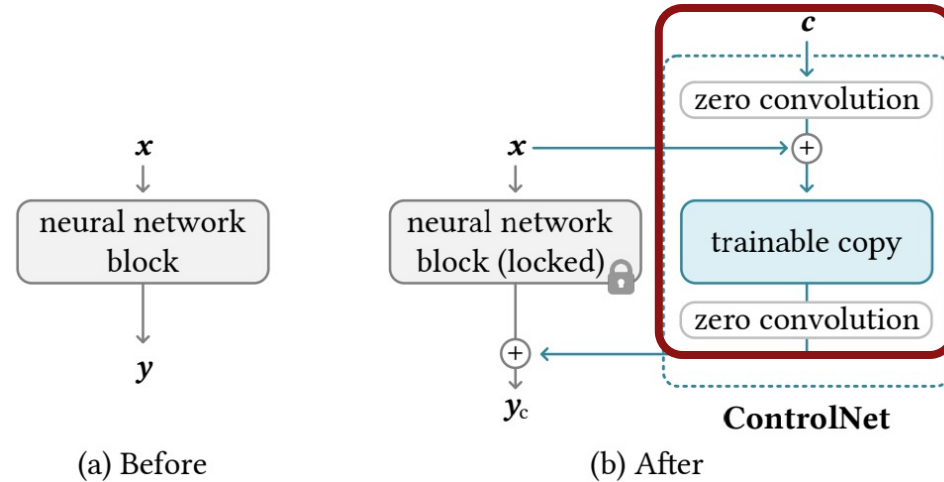


# ControlNet [Zhang et al., 2023]

ControlNet modifies each neural network block using zero convolution as follows:

Zero in the beginning since  $\mathbf{a}_2 = \mathbf{b}_2 = \mathbf{0}$ .

$$\mathbf{y}_c = F(\mathbf{x}; \Theta) + Z(F(\mathbf{x} + Z(\mathbf{c}; \mathbf{a}_1, \mathbf{b}_1); \Theta_c); \mathbf{a}_2, \mathbf{b}_2)$$

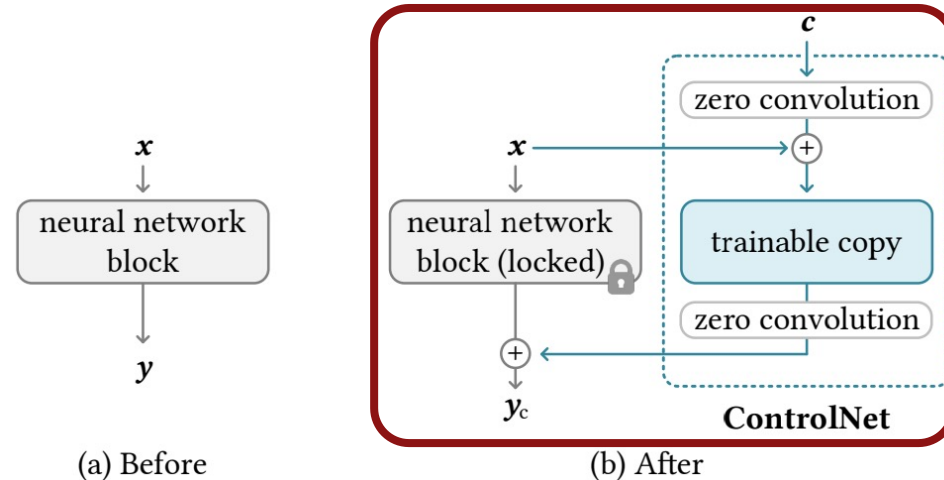


# ControlNet [Zhang et al., 2023]

ControlNet modifies each neural network block using zero convolution as follows:

Same as  $F(\mathbf{x}; \Theta)$  in the beginning.

$$y_c = F(\mathbf{x}; \Theta) + Z(F(\mathbf{x} + Z(\mathbf{c}; \mathbf{a}_1, \mathbf{b}_1); \Theta_c); \mathbf{a}_2, \mathbf{b}_2)$$

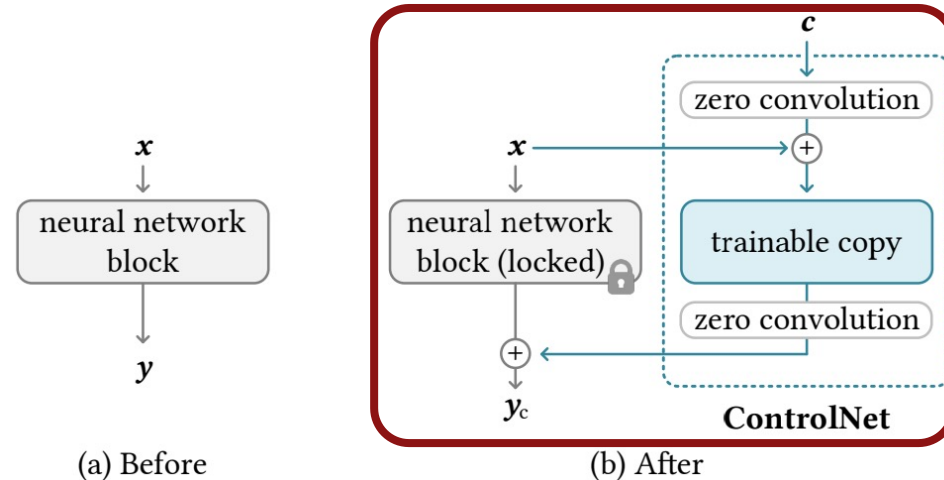


# ControlNet [Zhang et al., 2023]

Also for the noise prediction, gradually incorporate the output of conditional image encoding.

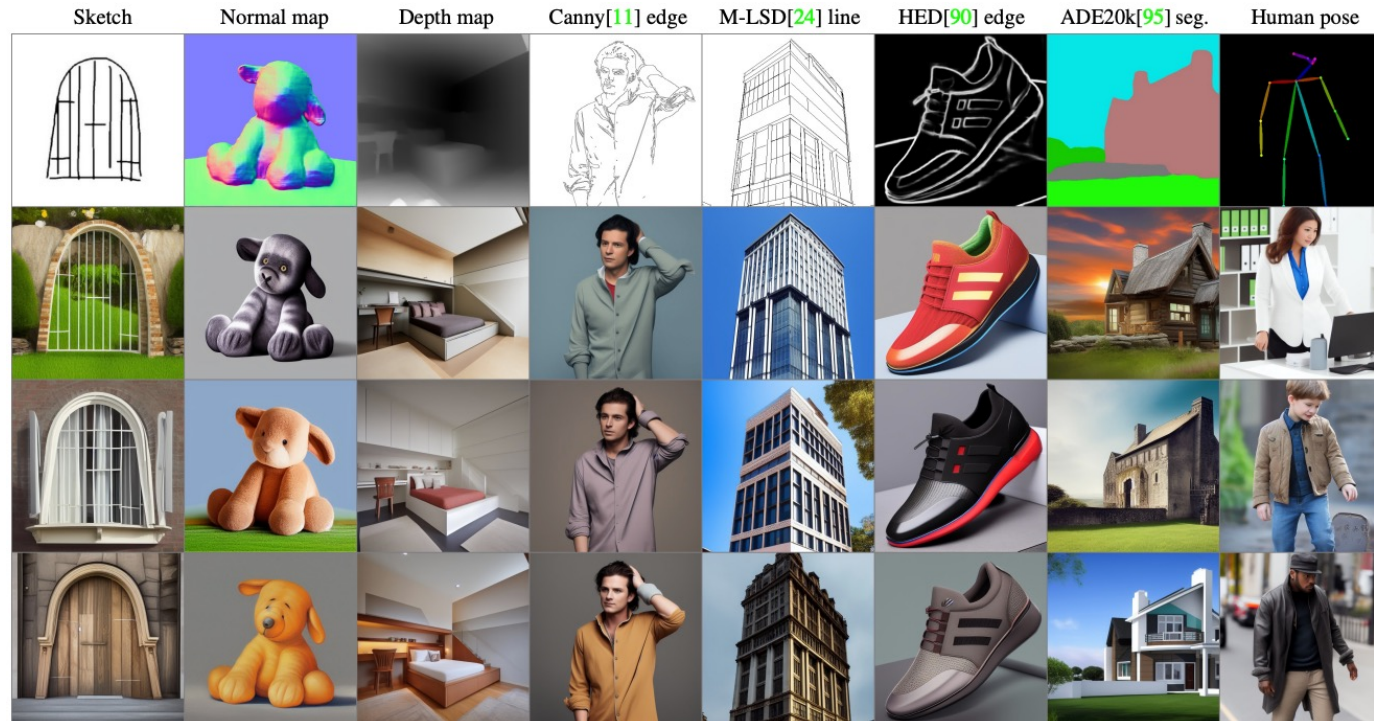
Same as  $F(\mathbf{x}; \Theta)$  in the beginning.

$$y_c = F(\mathbf{x}; \Theta) + Z(F(\mathbf{x} + Z(\mathbf{c}; \mathbf{a}_1, \mathbf{b}_1); \Theta_c); \mathbf{a}_2, \mathbf{b}_2)$$



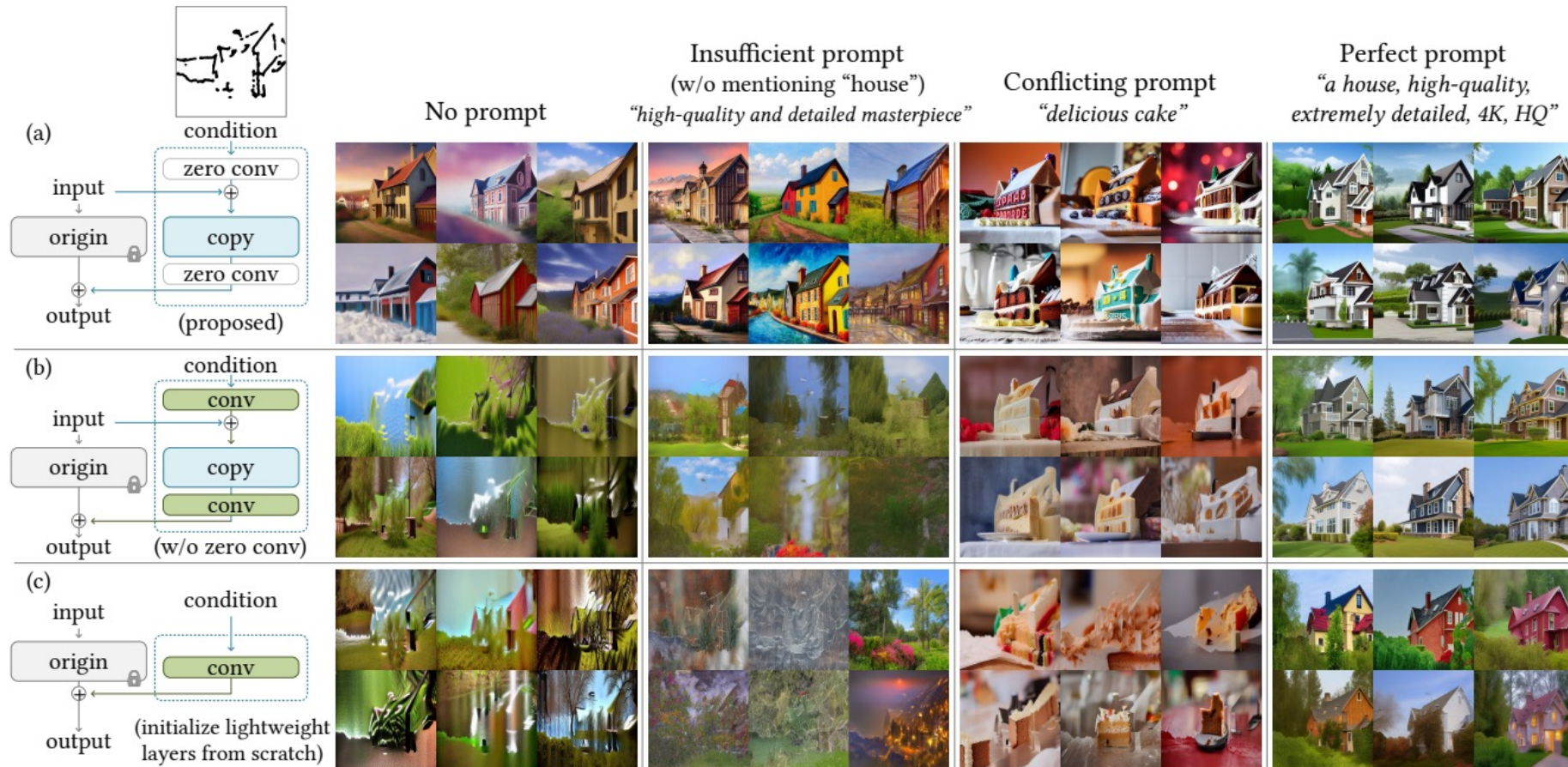
# ControlNet [Zhang et al., 2023]

- The idea can be used for any image conditioning.
- The training dataset is “ ~100k in size, which is 50,000 time smaller than the LAION-5B.”

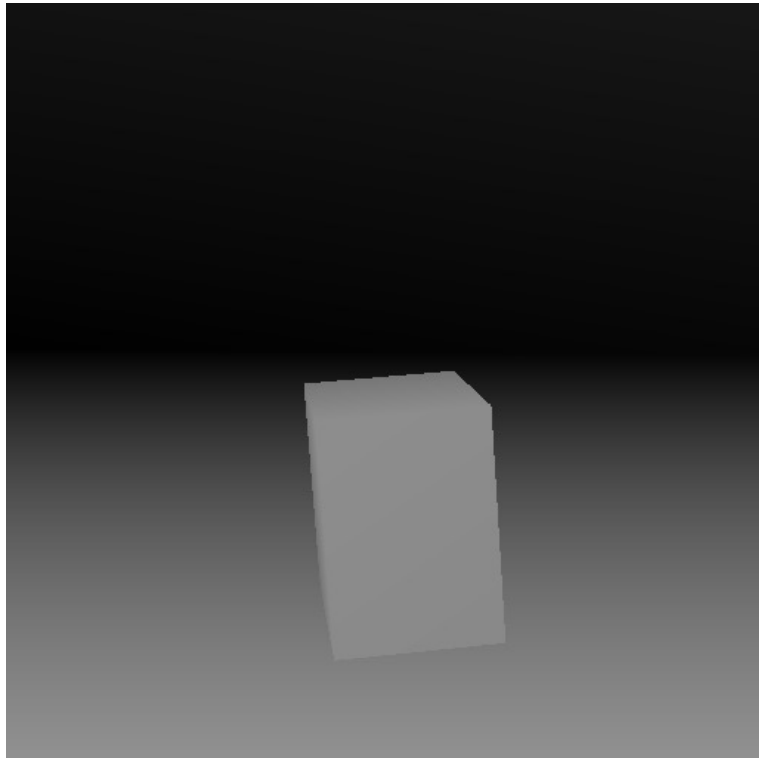


# ControlNet [Zhang et al., 2023]

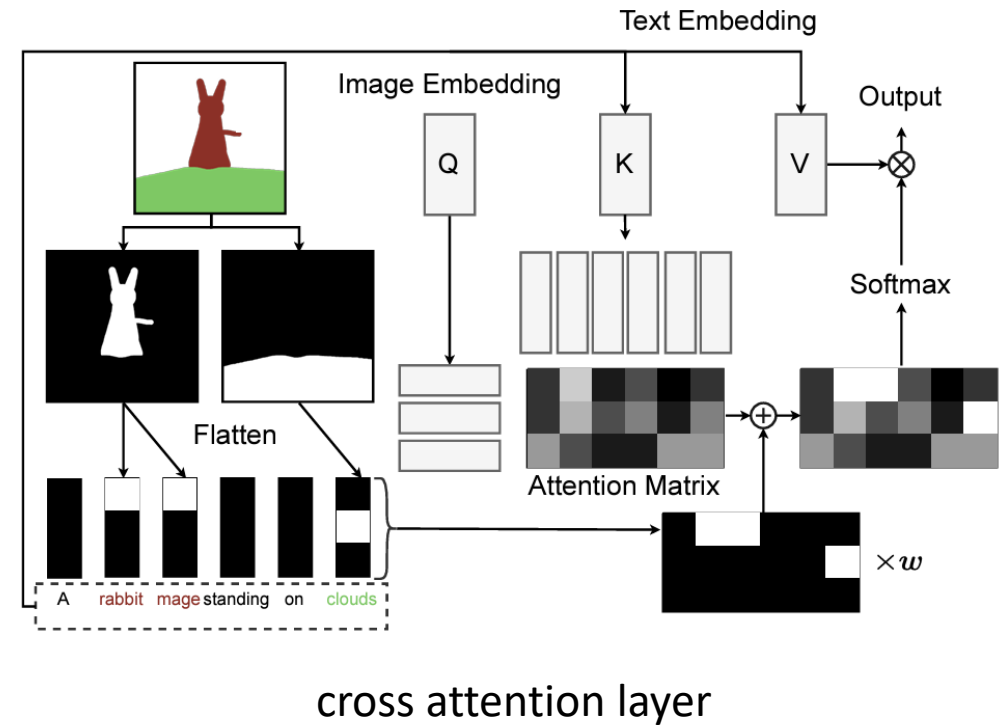
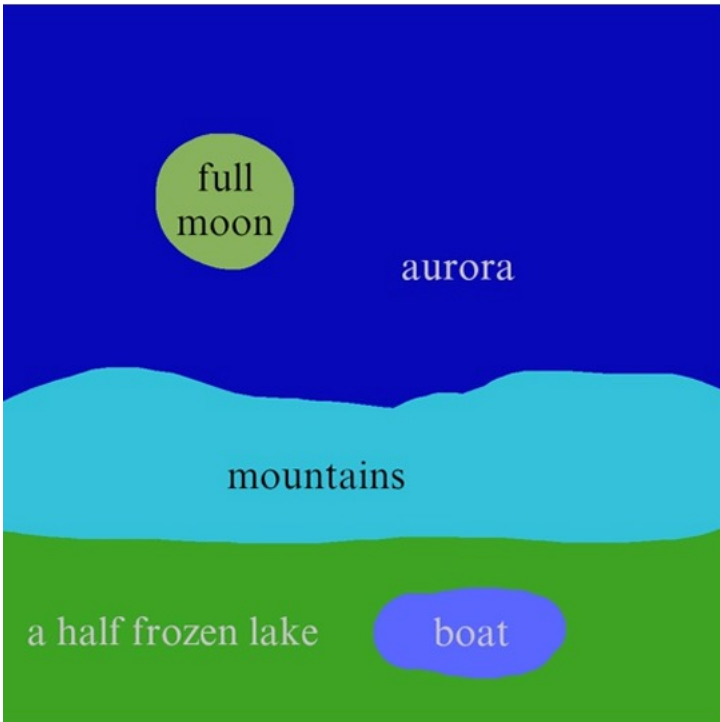
Ablative study of different architectures.



# Adherence of Conditioning Signal



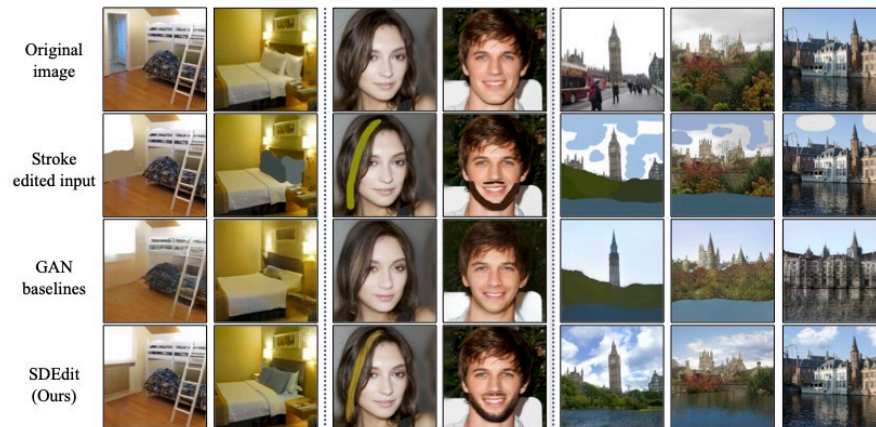
# Regional Control of Conditioning Signal



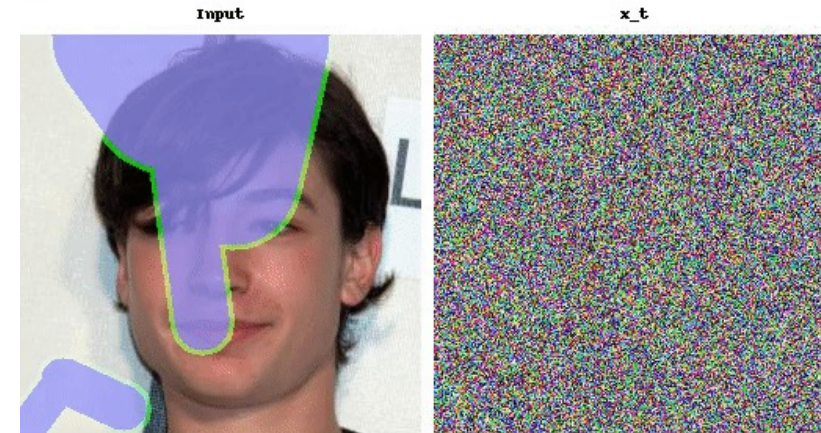
# Zero-Shot Applications

# Zero-Shot Adaptations

How to **edit** and **inpaint** images using a pretrained image diffusion model even without the need for fine-tuning?



Meng et al., SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations, ICLR 2022.

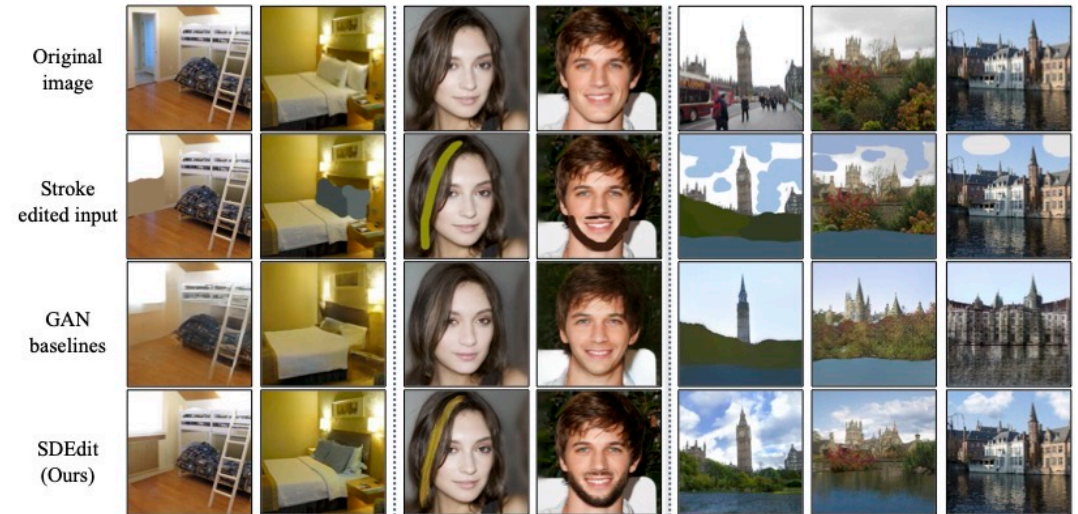
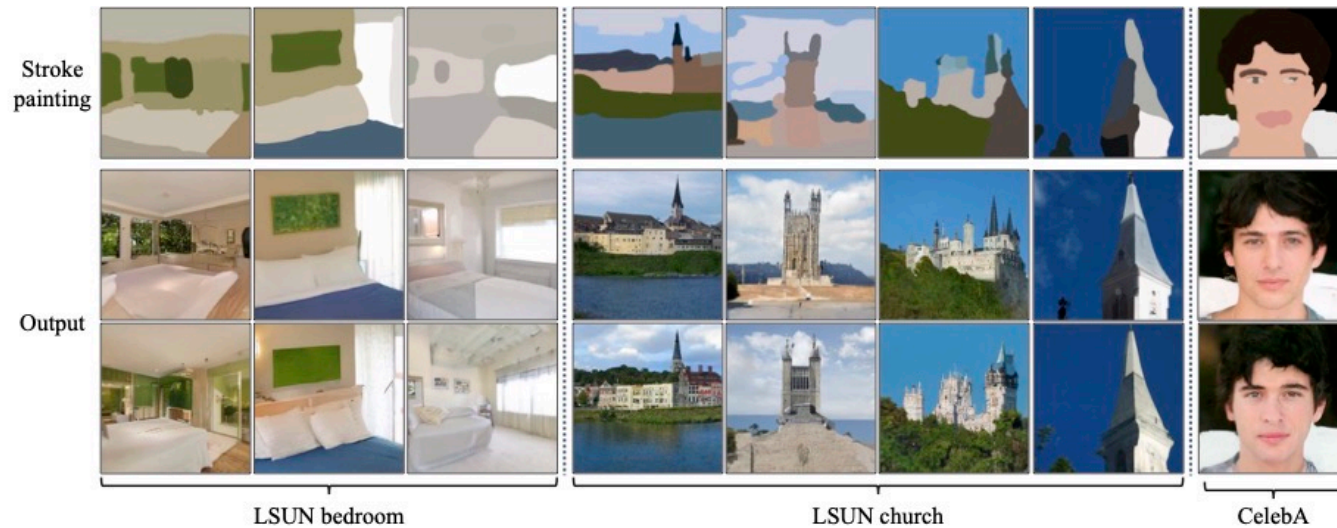


Lugmayr et al., RePaint: Inpainting using Denoising Diffusion Probabilistic Models, CVPR 2022.

# SDEdit [Meng et al., 2022]

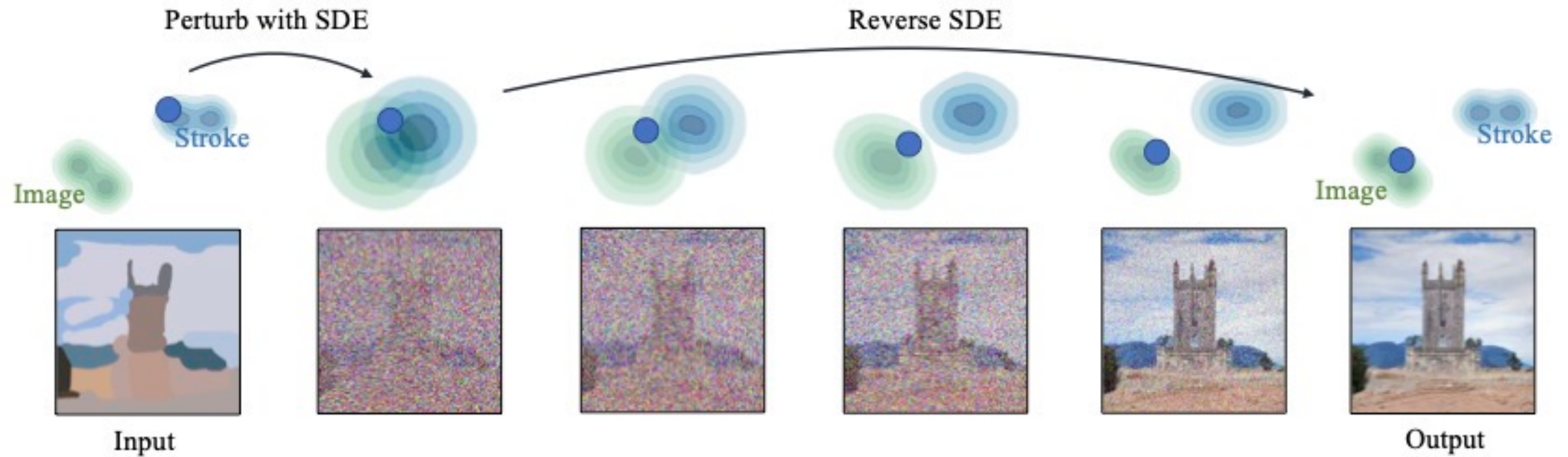
Image generation/editing through **user interaction**

- Image generation from sketches
- Image editing from scribbles



# SDEdit [Meng et al., 2022]

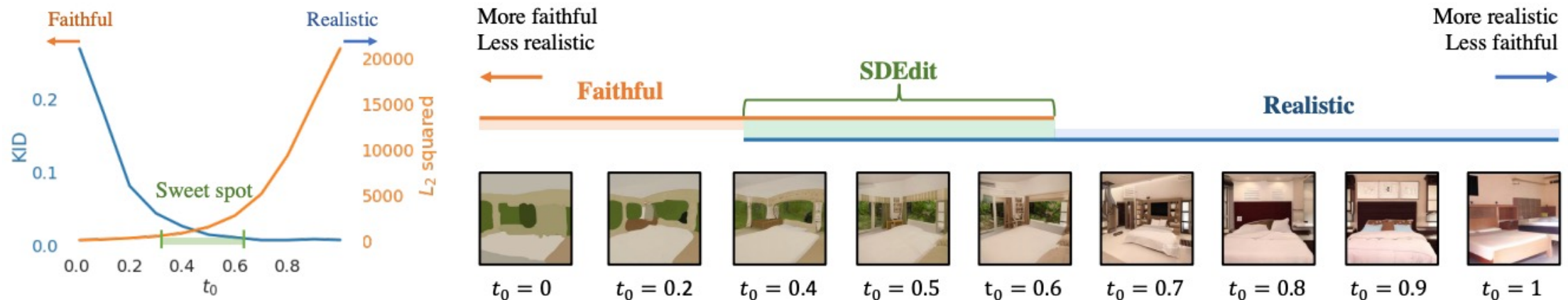
Perform the **forward** process for a bit and then **reverse** the process.



# SDEdit [Meng et al., 2022]

## Realism vs. faithfulness

As you perform the forward process with a **longer** timestep, the output image becomes **more realistic but less faithful** (deviating from the given condition).



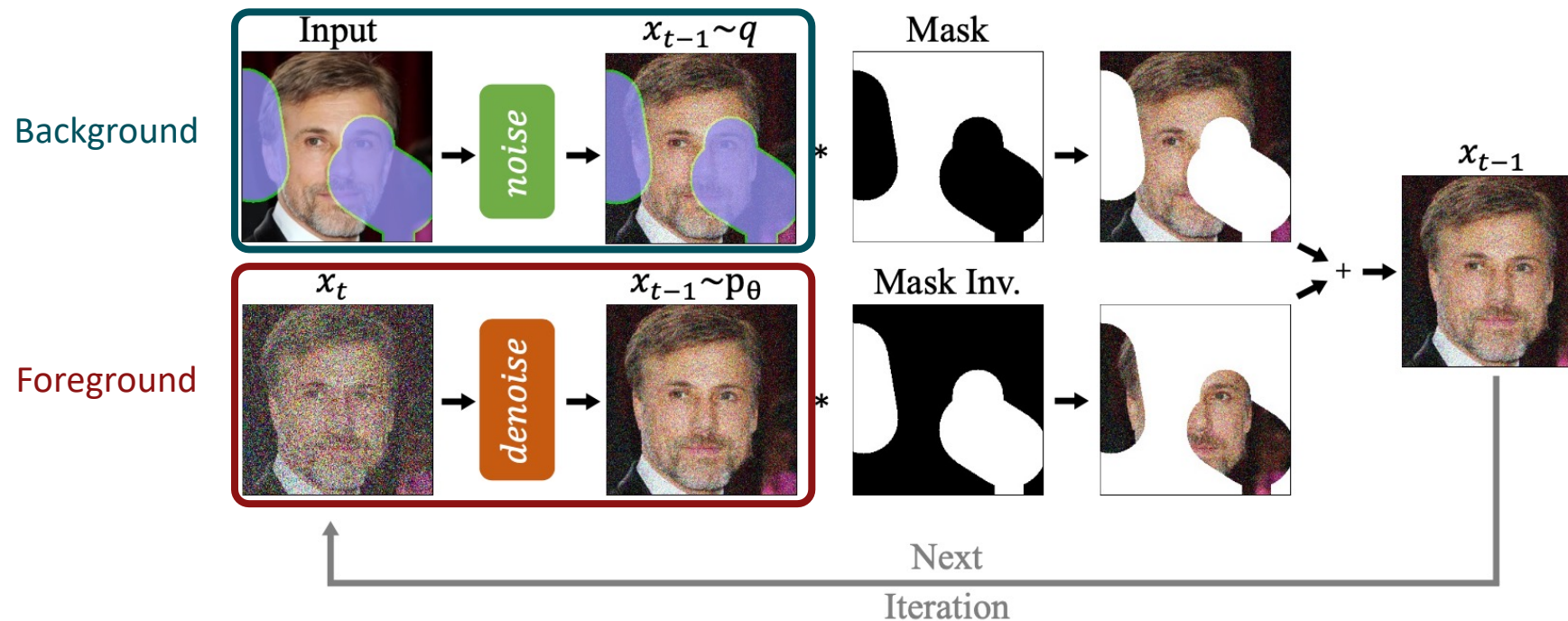
# RePaint [Lugmayr et al., 2022]

- Filling regions in an image using a pretrained image diffusion model.
- Assume that the missing **foreground** is filled while the **background** is fixed.



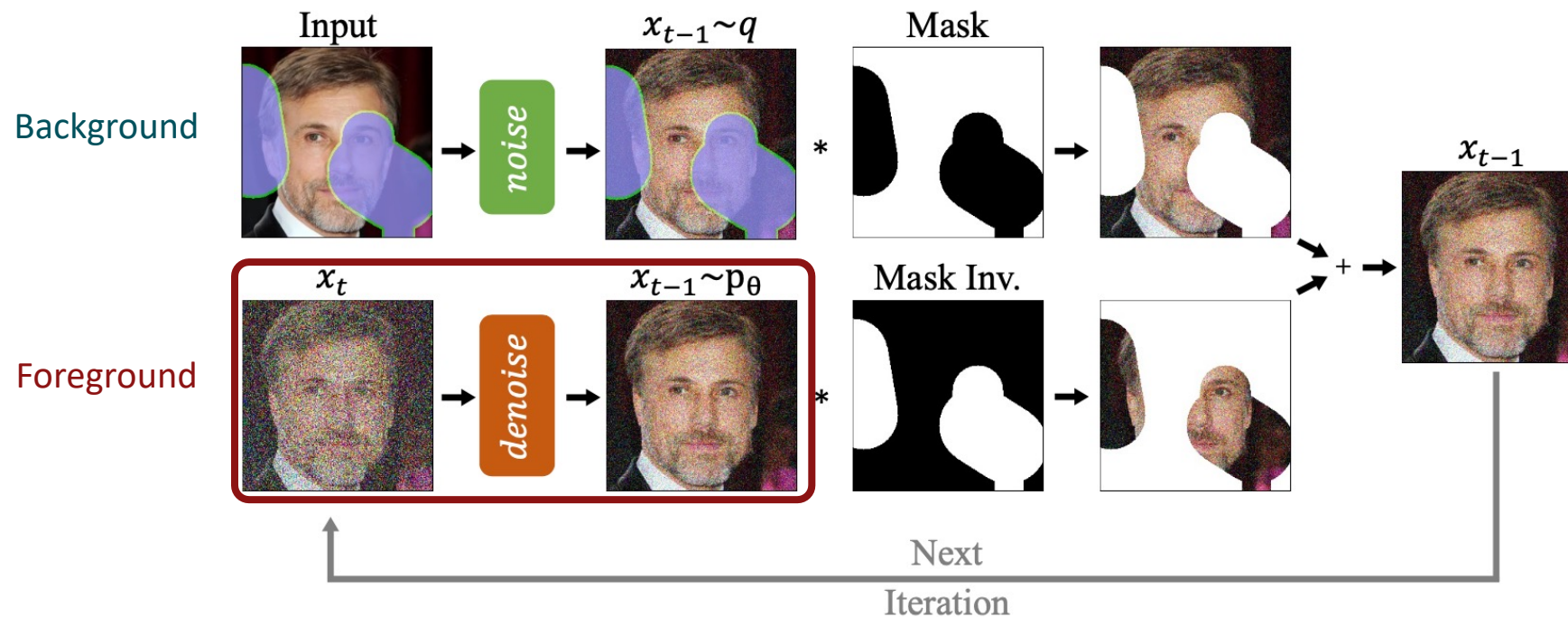
# RePaint [Lugmayr et al., 2022]

Combine denoised **foreground** images (to be filled) and noisy **background** (to be fixed) images at each iteration of the reverse process.



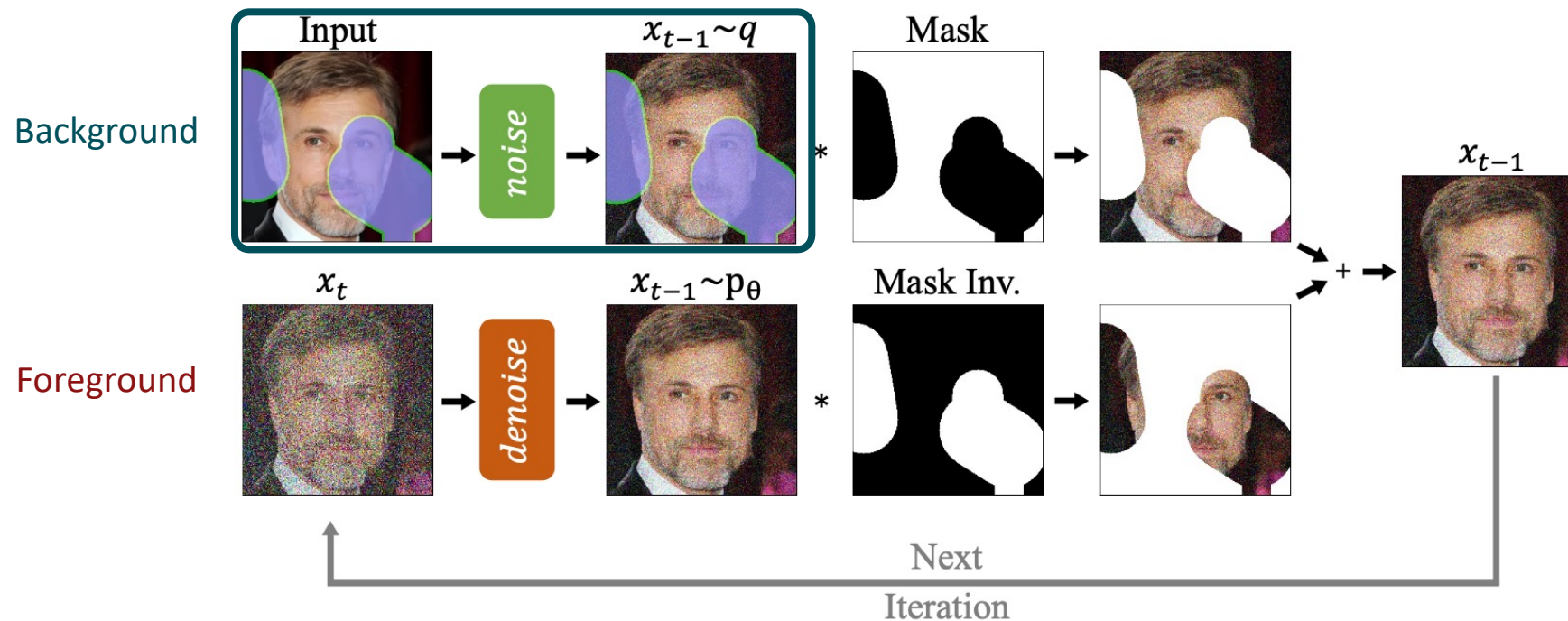
# RePaint [Lugmayr et al., 2022]

- Starting from  $x_T$ , at each timestep  $t$ , denoise  $x_t$  one step (reverse process), resulting  $x_{t-1}^f$ .



# RePaint [Lugmayr et al., 2022]

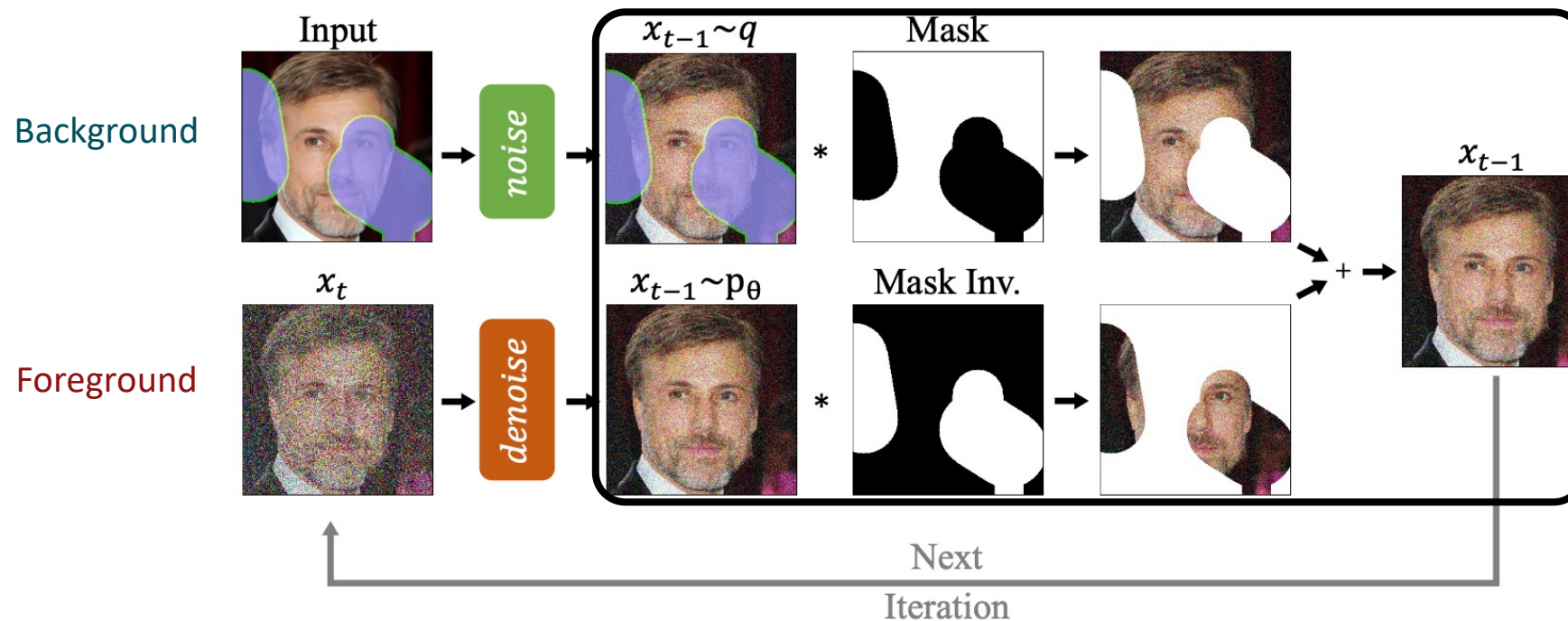
2. Perturb the input background image  $x_0^b$  (forward process) with a noise scale of the timestep  $t - 1$ , resulting  $x_{t-1}^b$ .



# RePaint [Lugmayr et al., 2022]

3. Combine  $x_{t-1}^b$  and  $x_{t-1}^f$  (where  $M$  is the background mask):

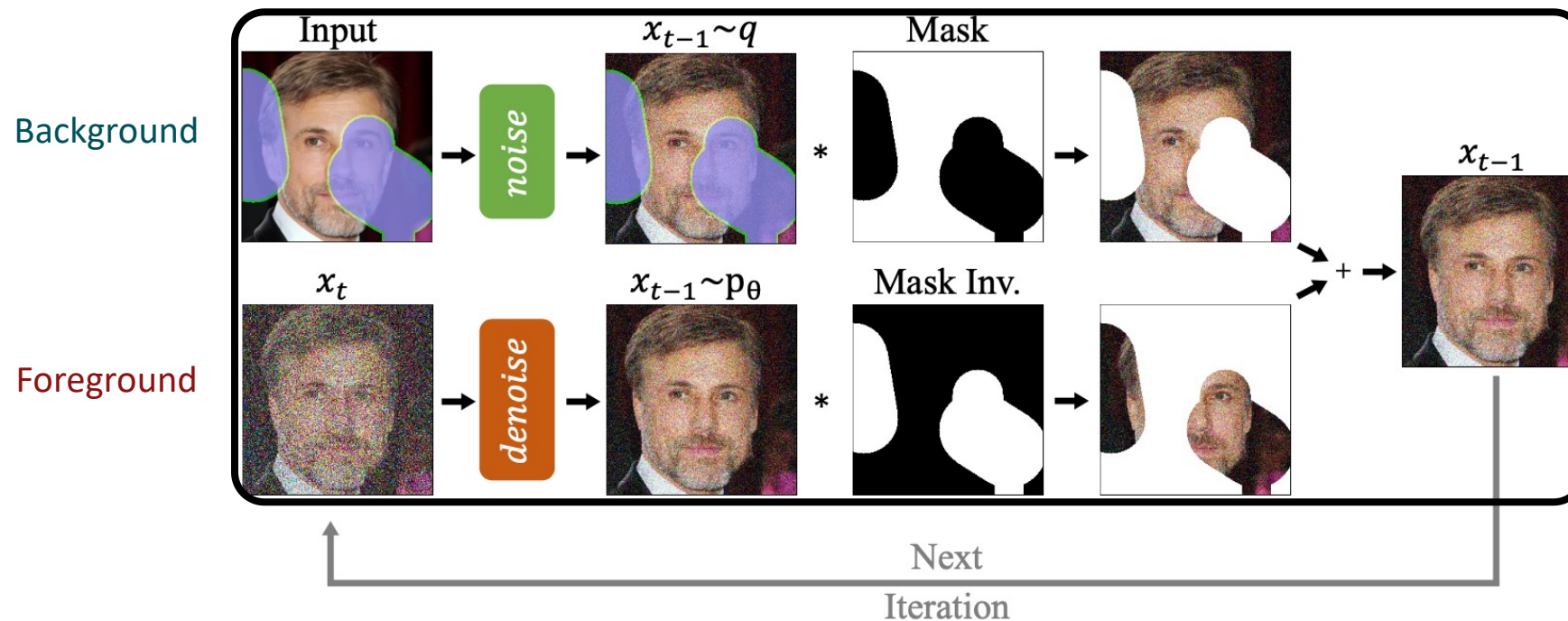
$$x_{t-1} = M \odot x_{t-1}^b + (1 - M) \odot x_{t-1}^f$$



# RePaint [Lugmayr et al., 2022]

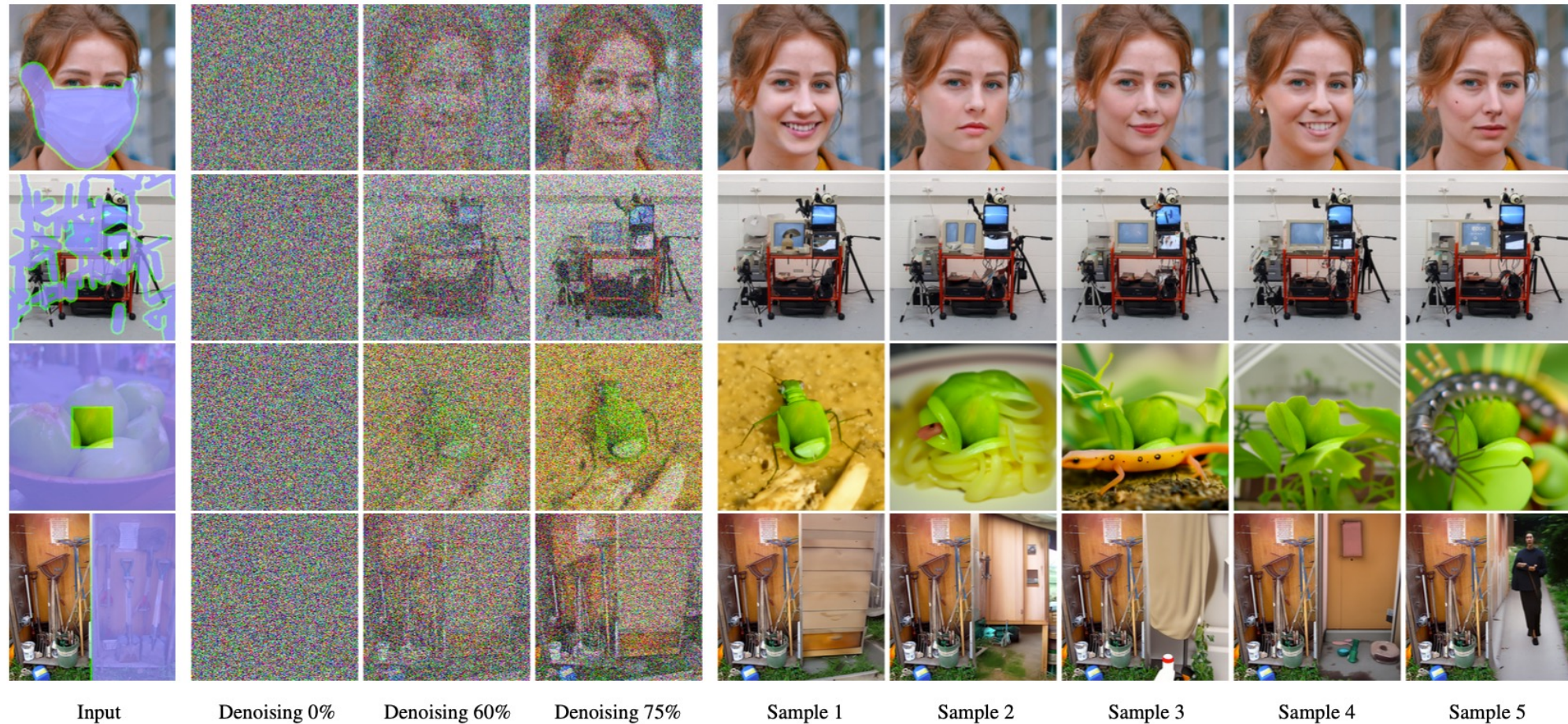
4. Repeat this process for  $t = T, \dots, 1$ .

(You may need to **replay** the forward/reverse processes in intermediate intervals.)



# RePaint [Lugmayr et al., 2022]

## Results with different masks



# Summary

## 1. Classifier Guidance / Classifier-Free Guidance

Enhancing the quality of generated outputs using additional information.

## 2. ControlNet

Adapting the pretrained diffusion model with relatively few conditional data using zero convolution.

## 3. SDEdit / RePaint

Utilizing the pretrained diffusion model for editing and image inpainting.

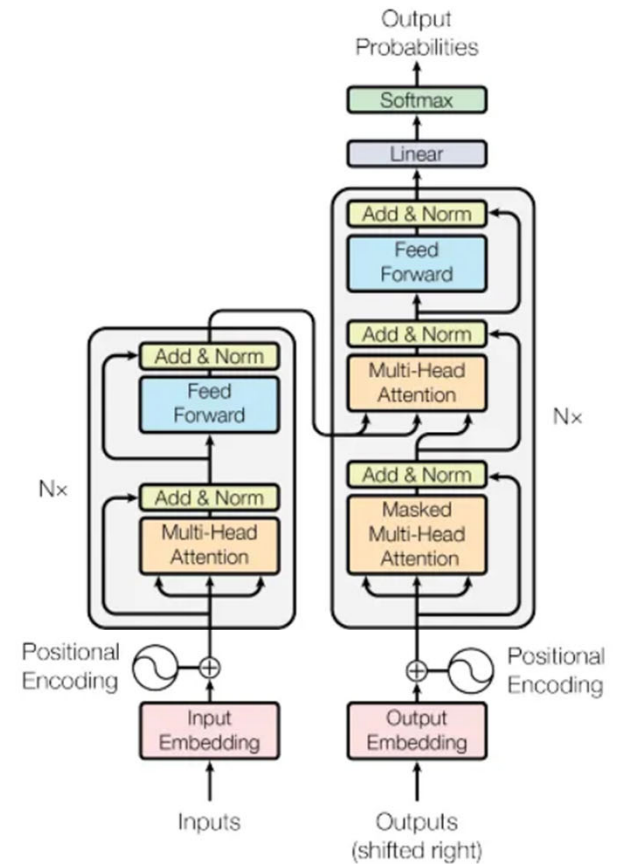
# Attention Please

Attention Mechanism in Generative Models

Daniel Cohen-Or

# “Attention is all you need” , Vaswani et al. 2017

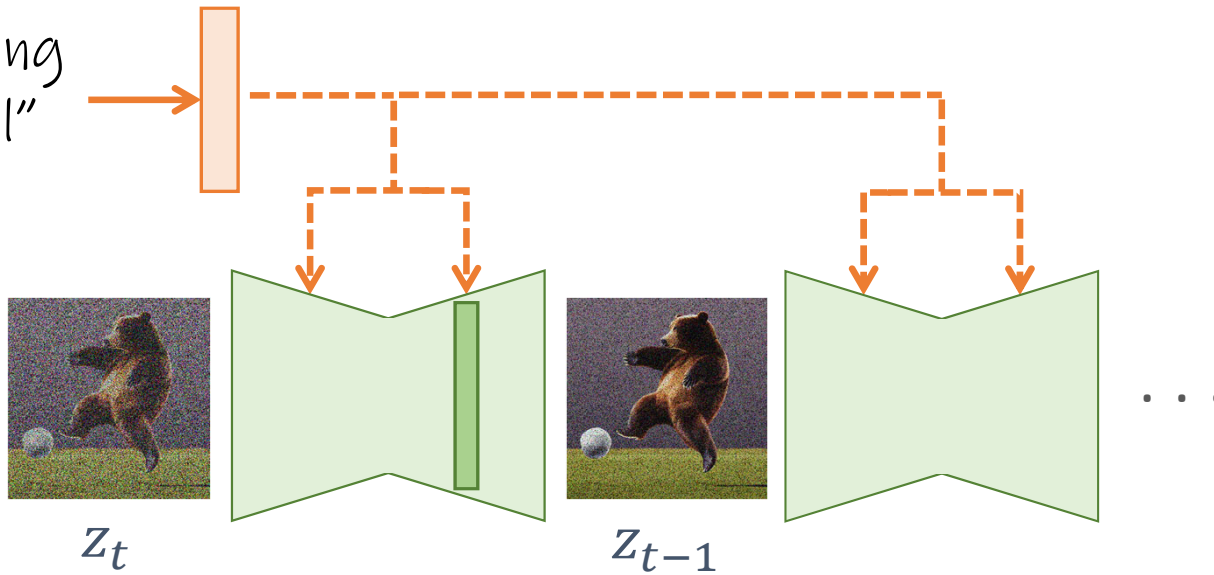
$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



From “Attention is all you need” paper by Vaswani, et al., 2017 [1]

# The Denoising Process with Attention Layers

"A bear kicking a soccer ball"



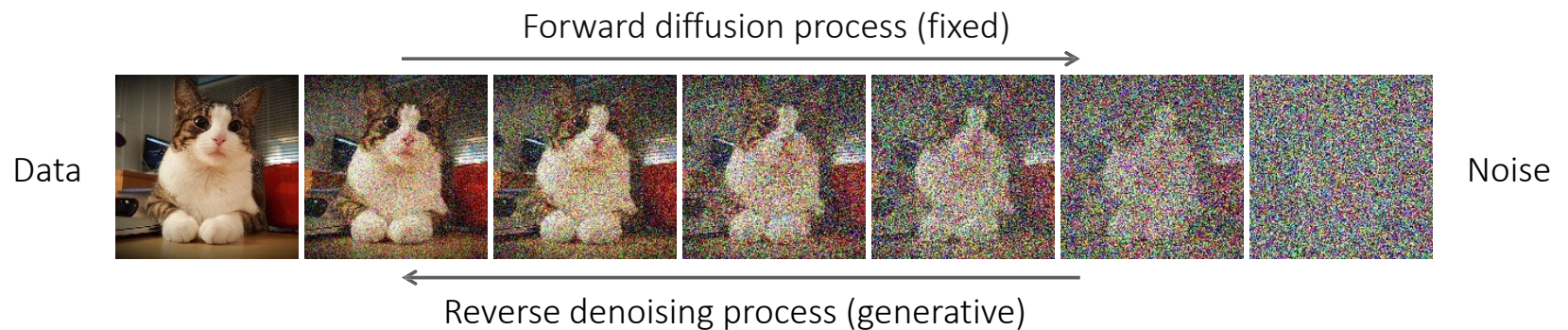
# Denoising network

# Denoising Diffusion Models

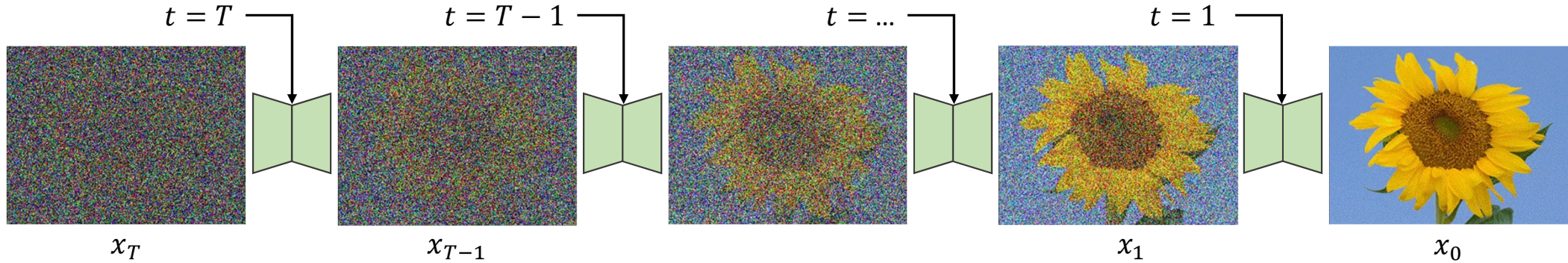
## Learning to generate by denoising

Denoising diffusion models consist of two processes:

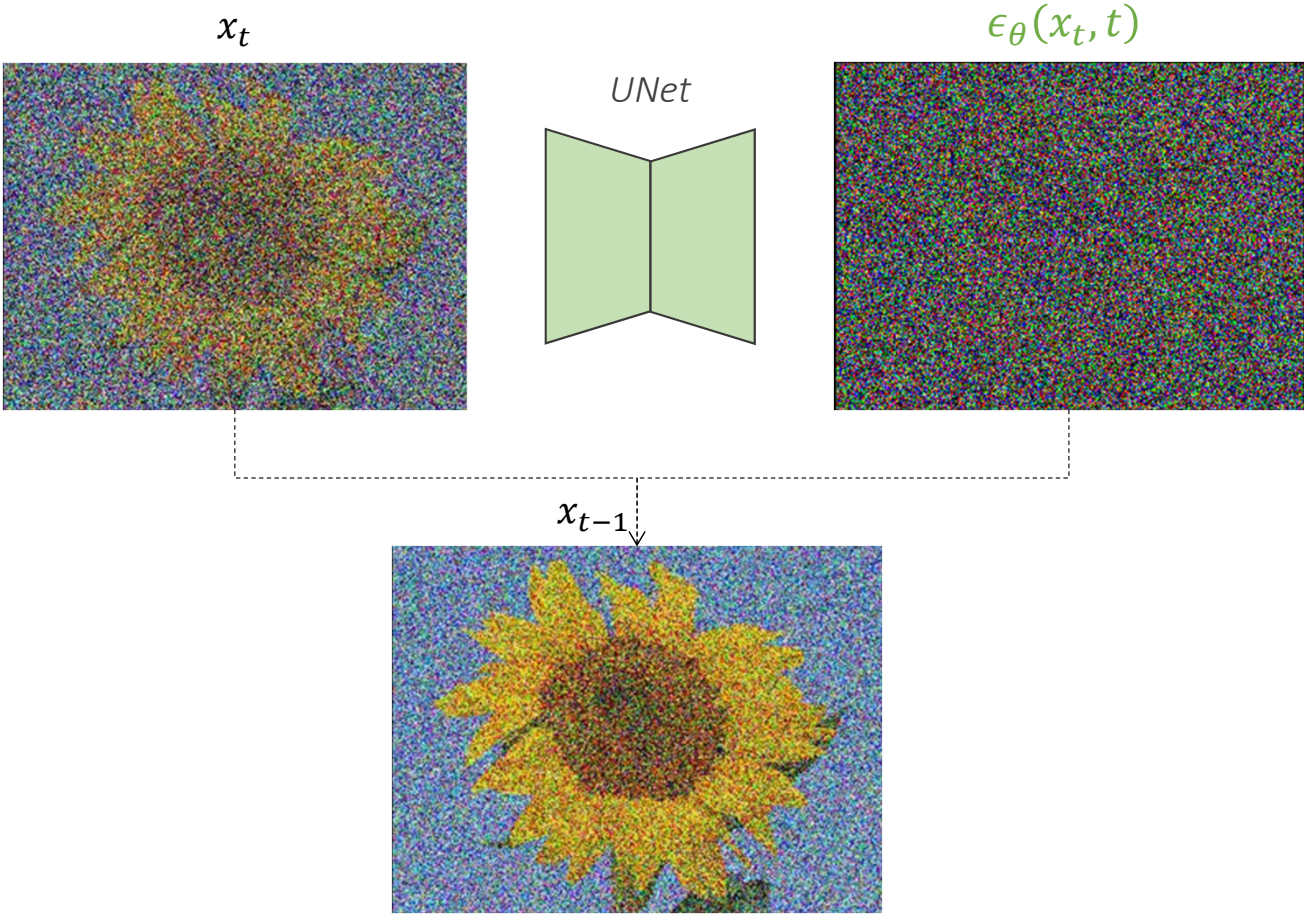
- Forward diffusion process that gradually adds noise to input
- Reverse denoising process that learns to generate data by denoising



# Diffusion Models

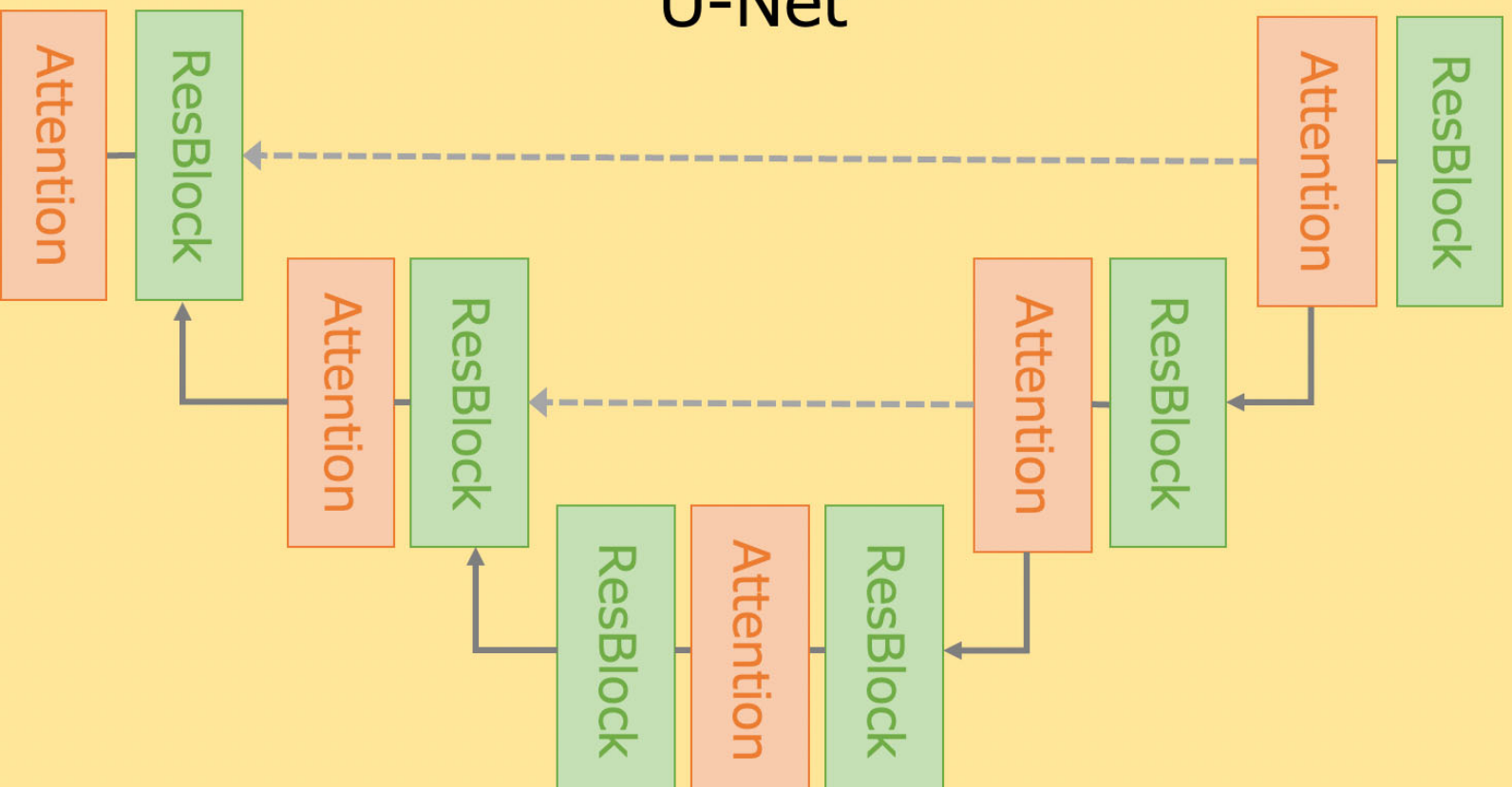


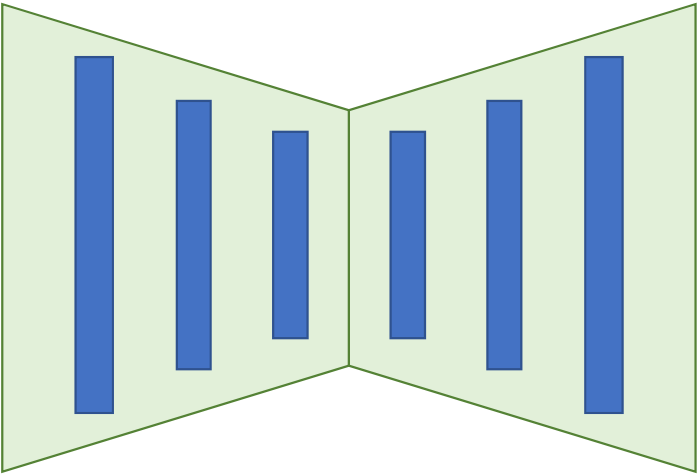
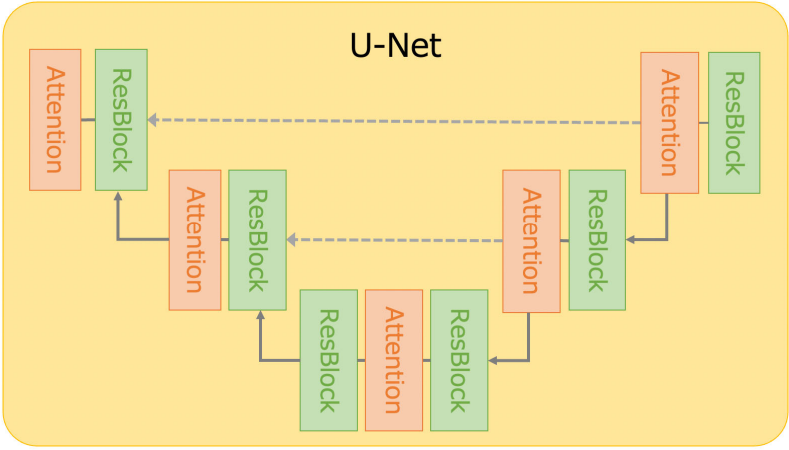
# Diffusion Models



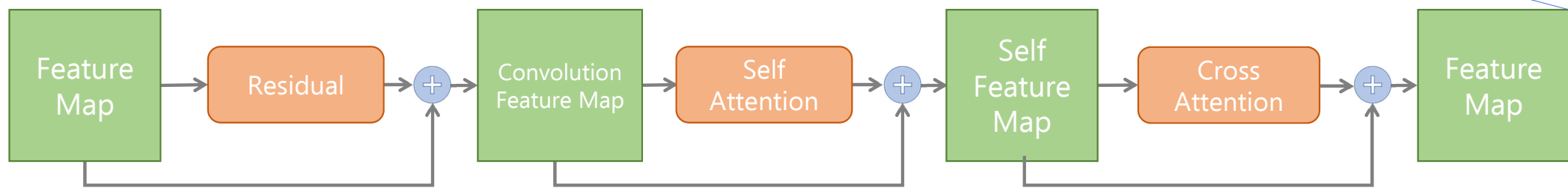
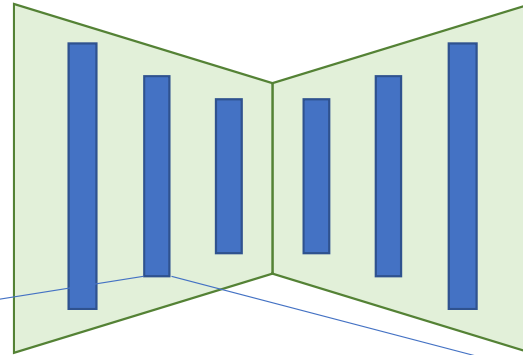
Unet

# U-Net

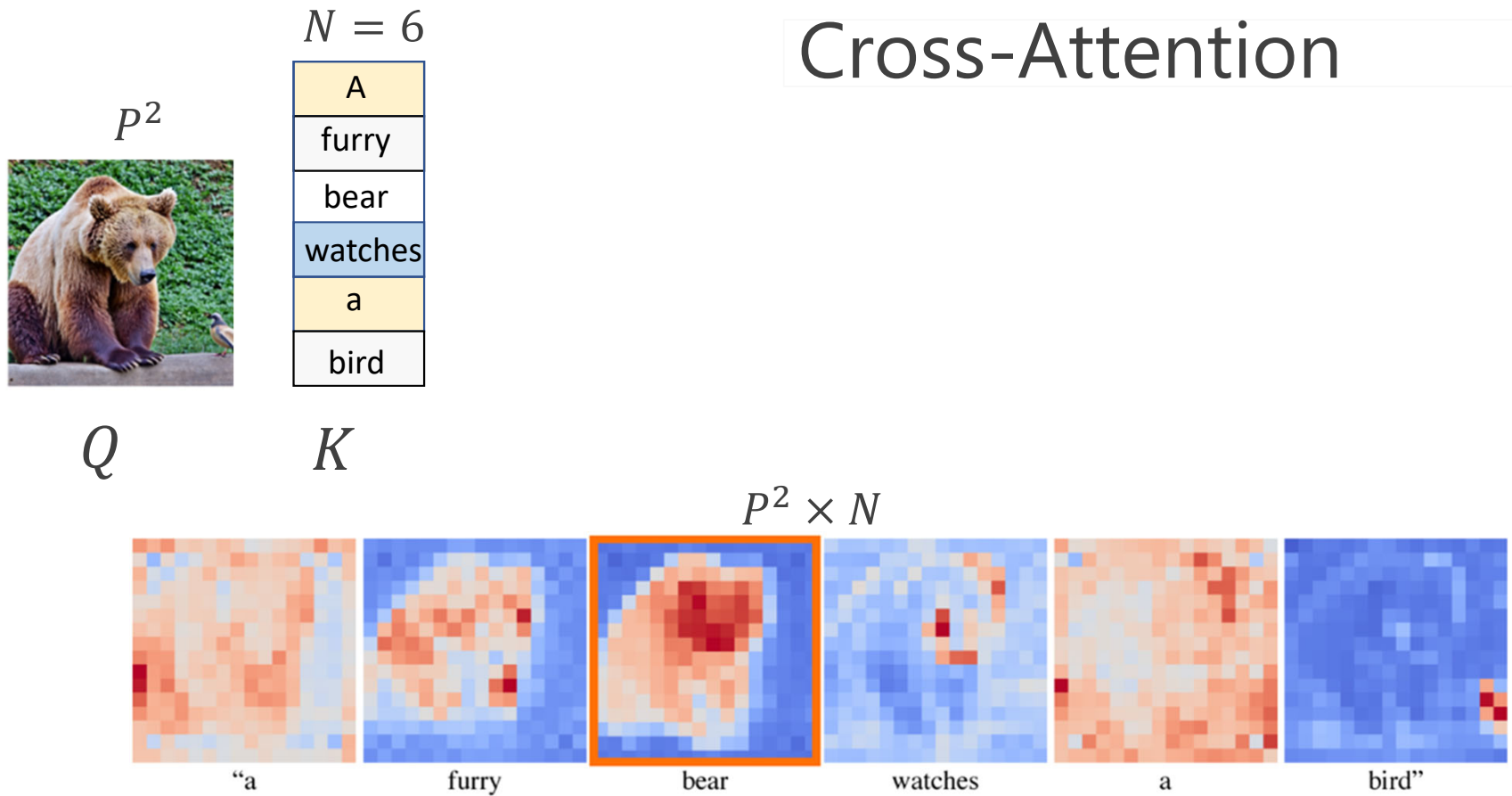




# The UNet Layer



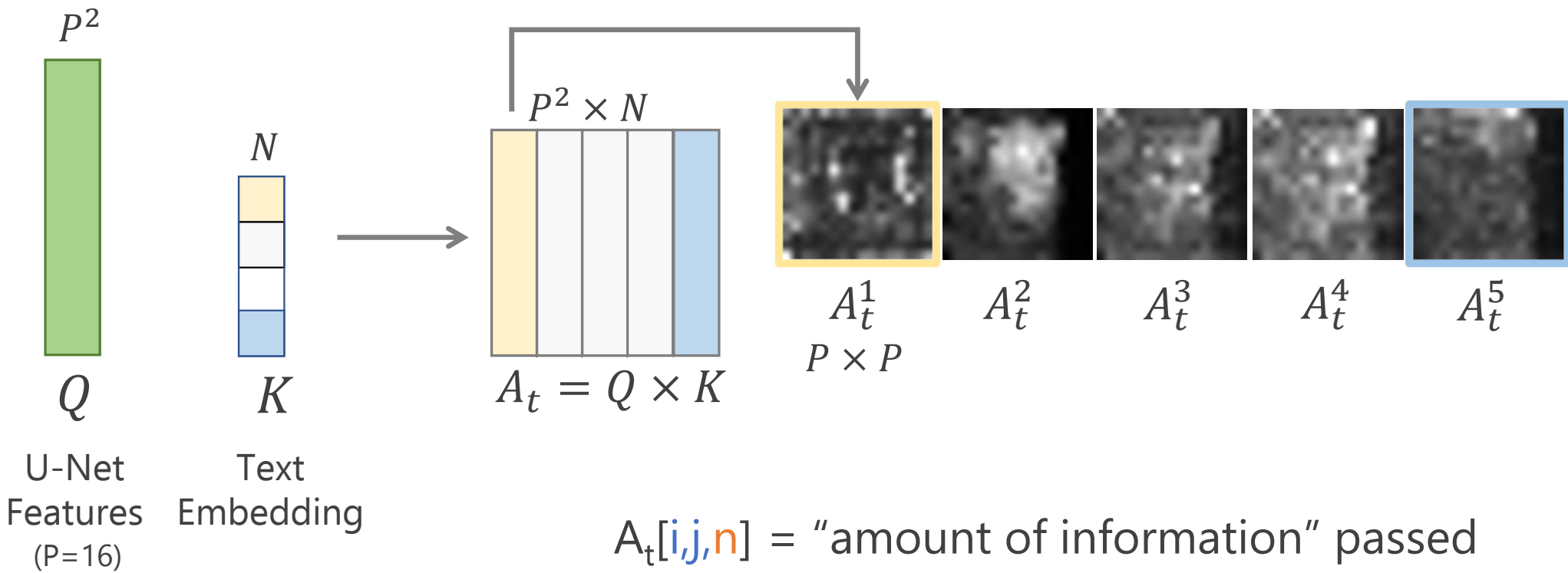
# Cross-Attention



$$A_t = Q \times K$$

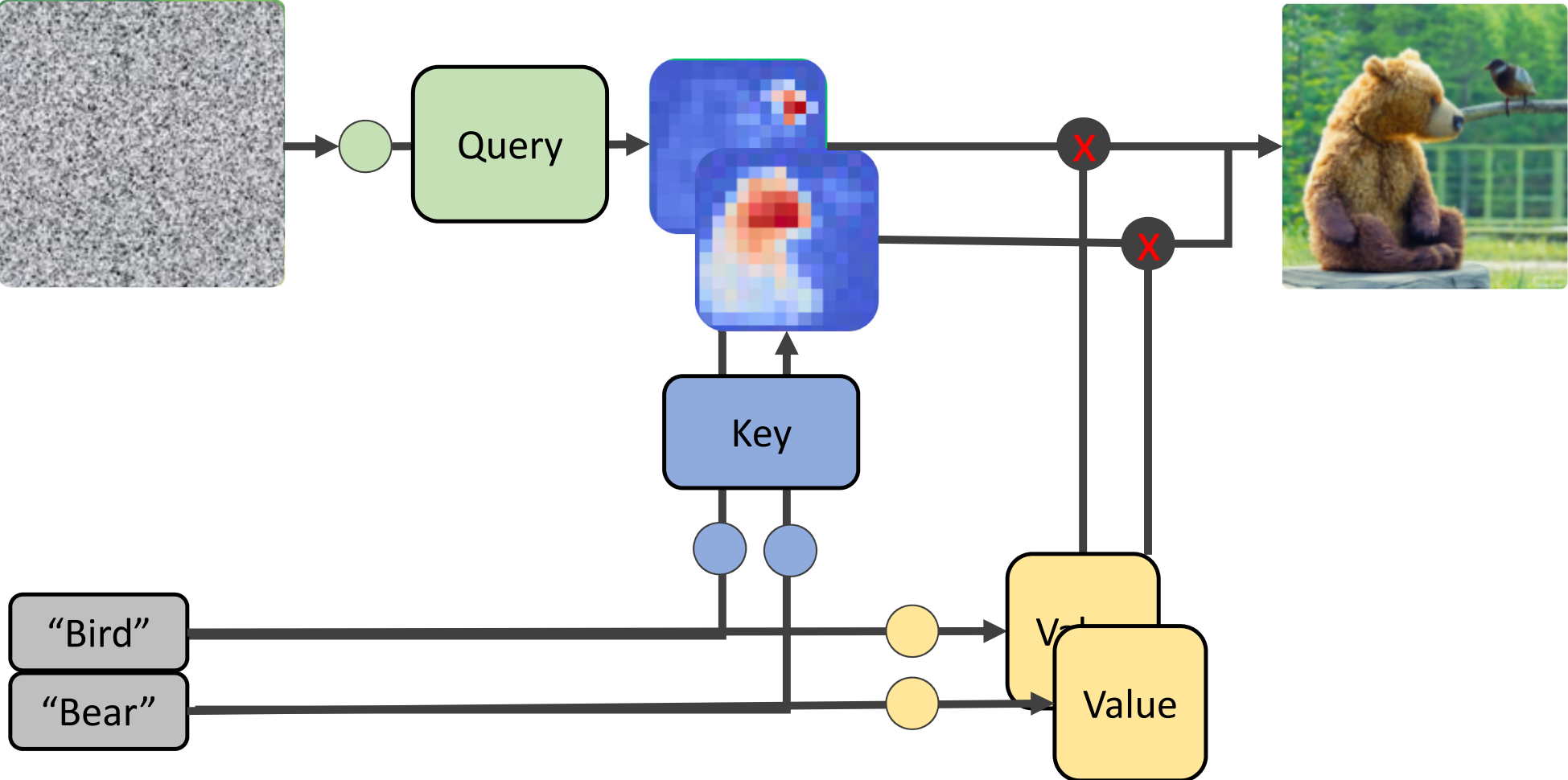
$A_t[i,j,n]$  = “amount of information” passed from token  $n$  to patch  $(i,j)$

# The Cross-Attention

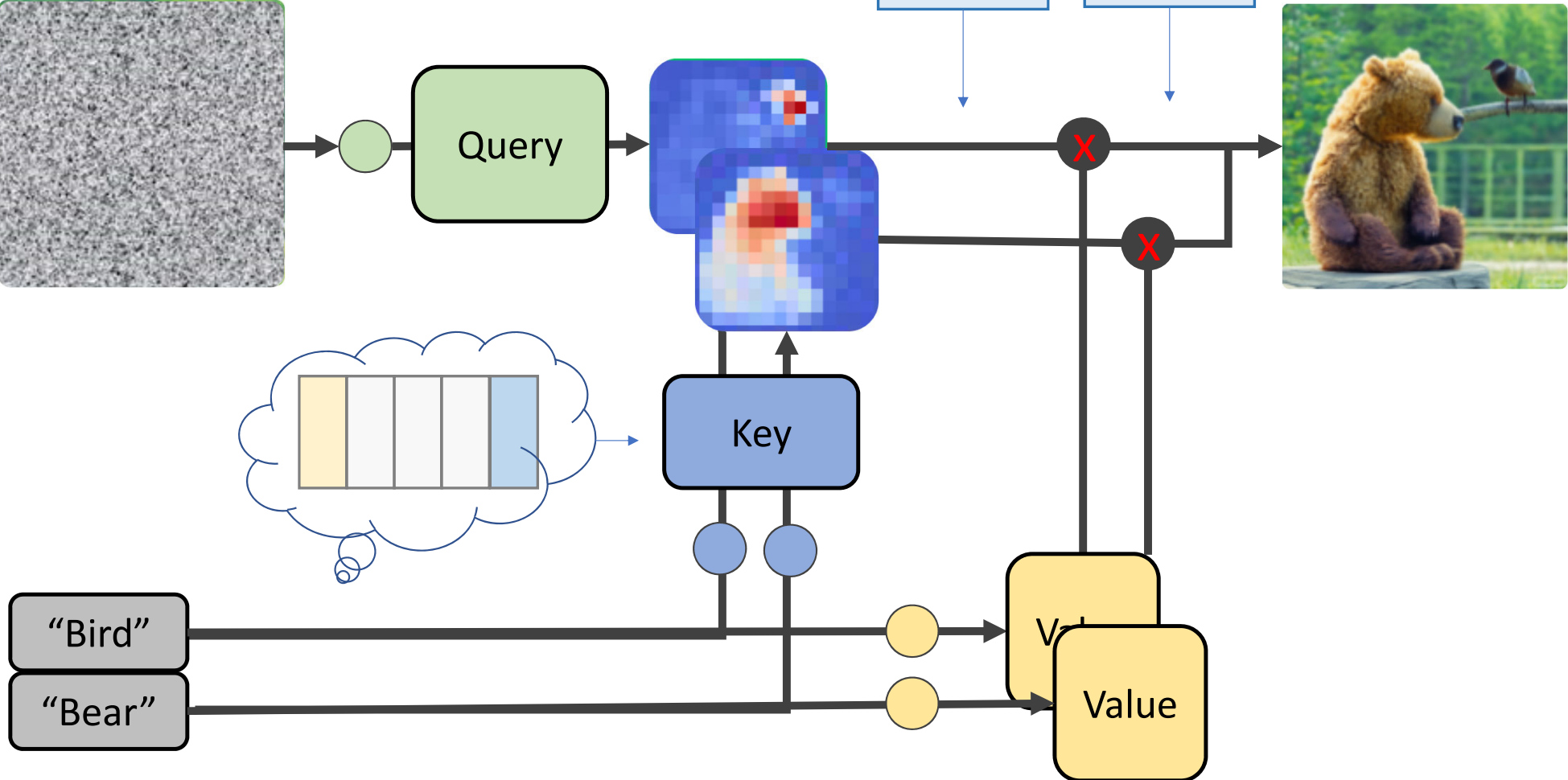


$A_t[i,j,n]$  = "amount of information" passed from token  $n$  to patch  $(i,j)$

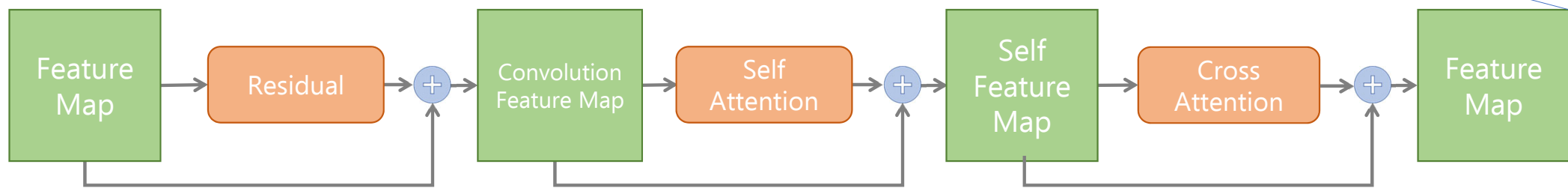
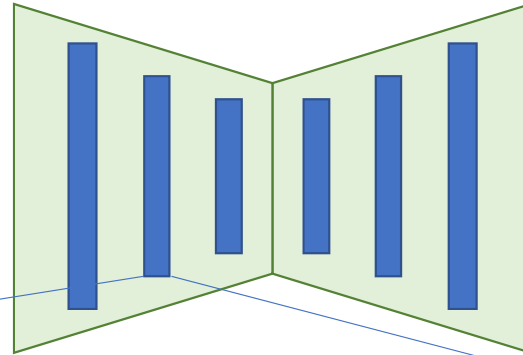
# Cross-attention



# Cross-attention

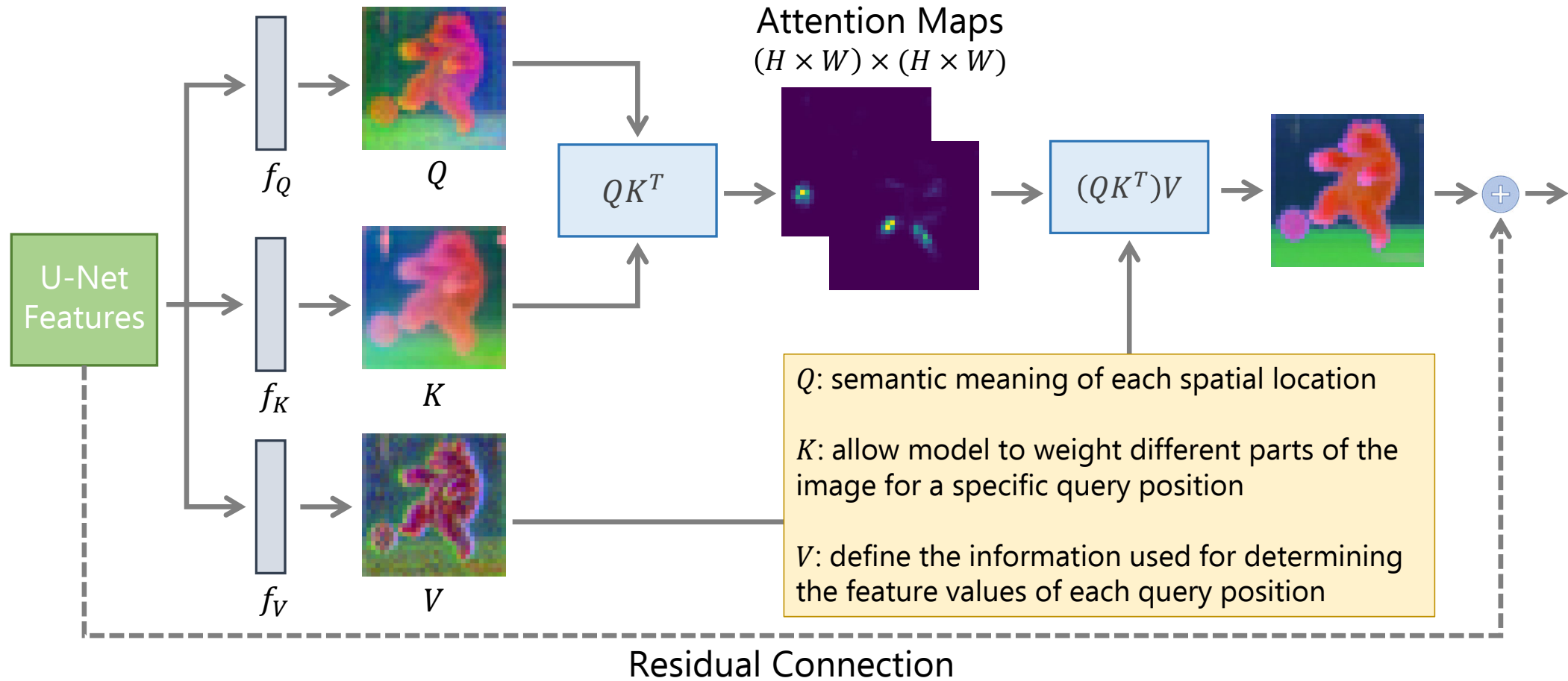


# The UNet Layer

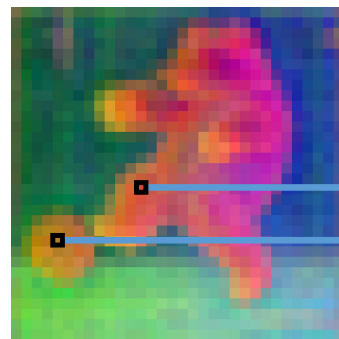


# The Self-Attention

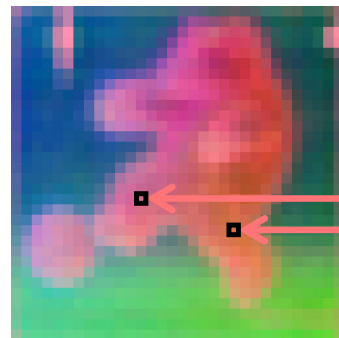
Represents where the model "looks" in the image for each spatial position in Q



# The Self-Attention

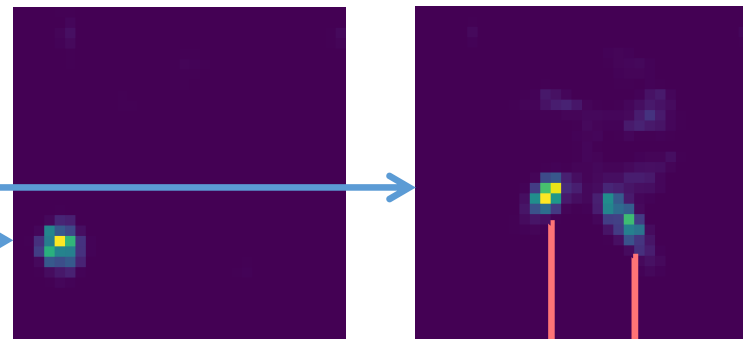


Queries  
 $H \times W \times C$



Keys  
 $H \times W \times C$

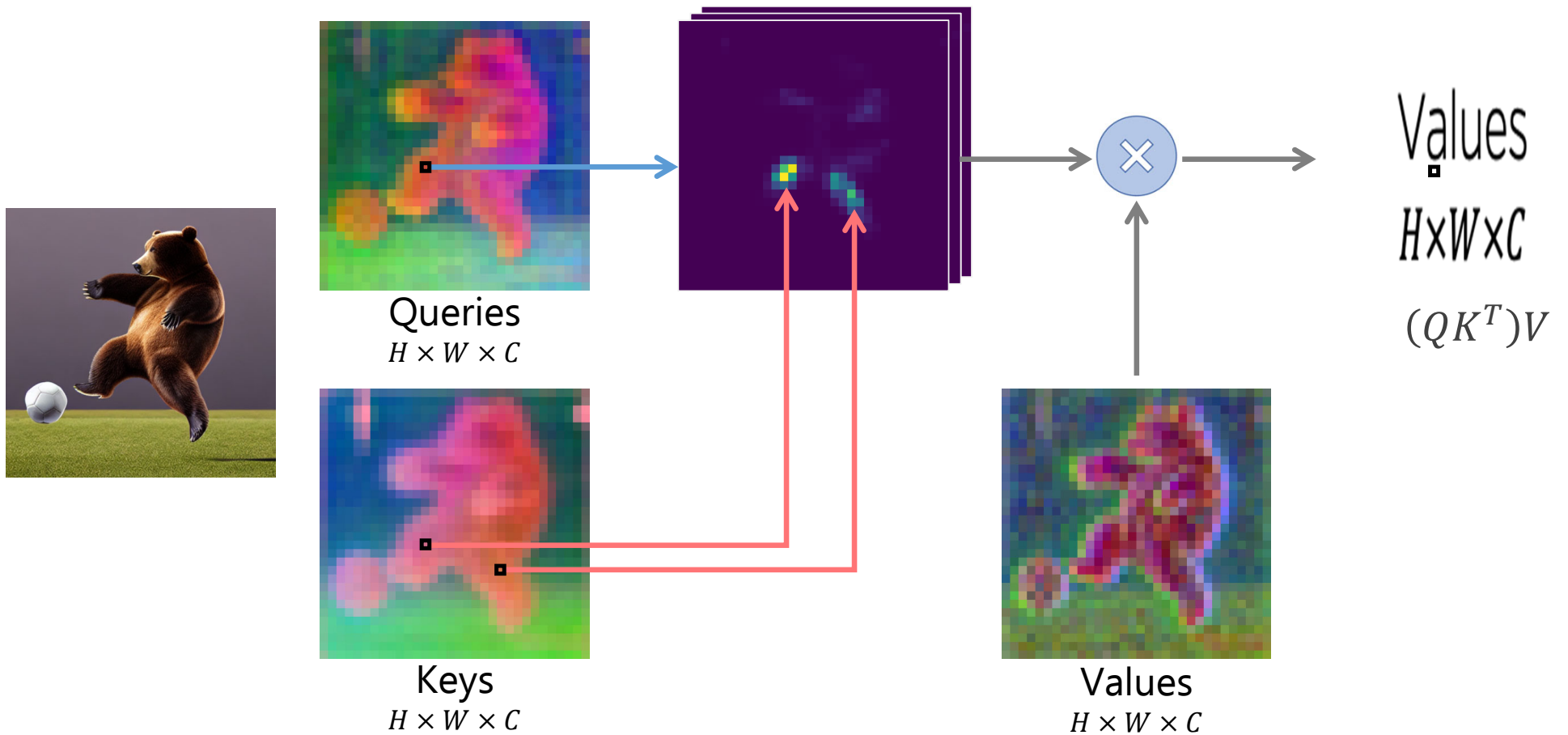
Attention Maps  
 $(H \times W) \times (H \times W)$



Each query defines to a  
 $H \times W$  attention map!

A query on the leg of the bear "attends"  
to keys located on the leg of the bear!

# The Self-Attention



t = 0.6, layer: 10 / 70





# Self-Segmentation

---

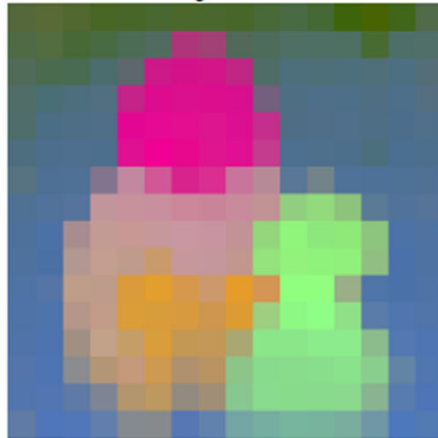
Localizing Object-level Shape Variations [Patashnik et al., ICCV 2023]

# Self-Attention Maps

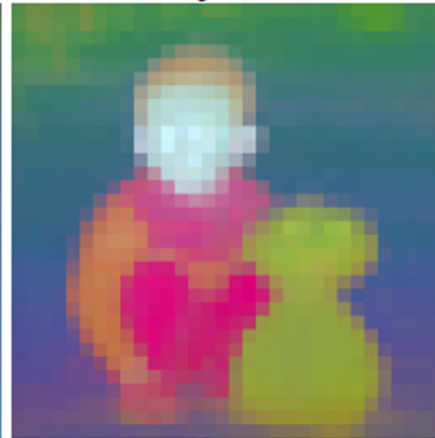
**Input image**



**layer=4**



**layer=8**



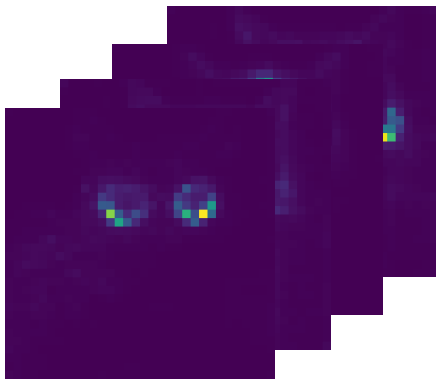
**layer=11**



Are these PCA on the self-attention ? On what exactly the QK maps?

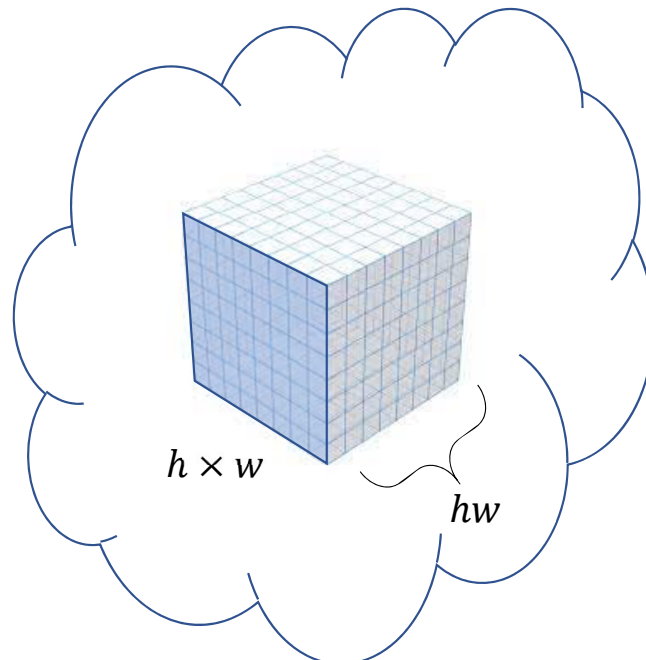
# Self-Segmentation

"a **cat** is wearing sunglasses"



$hw \times (h \times w)$

There is a lot of semantics in the self attention features!!!



cluster



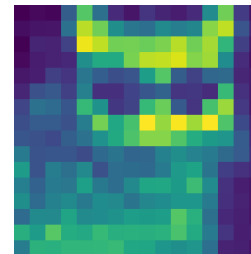
# Segments labeling



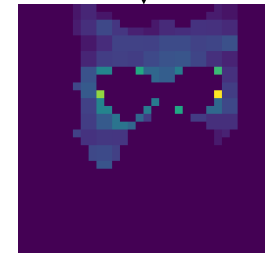
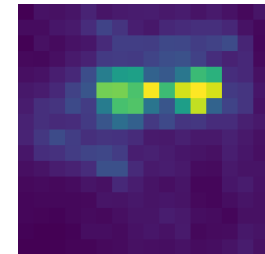
"a cat is wearing sunglasses"



cat



sunglasses



score: 0.65

score: 0.19

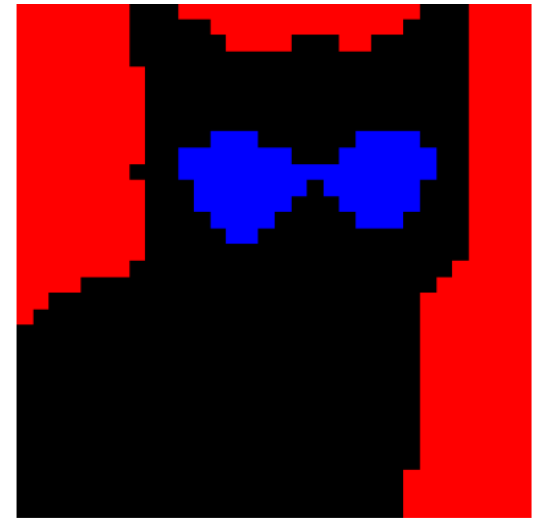
# Segments labeling



"a cat is wearing sunglasses"



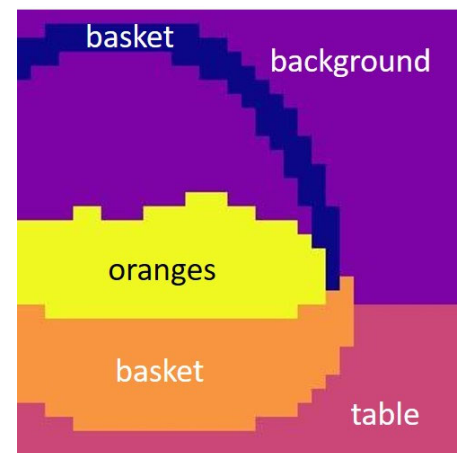
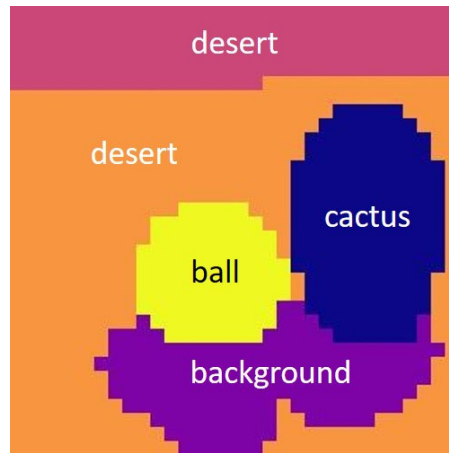
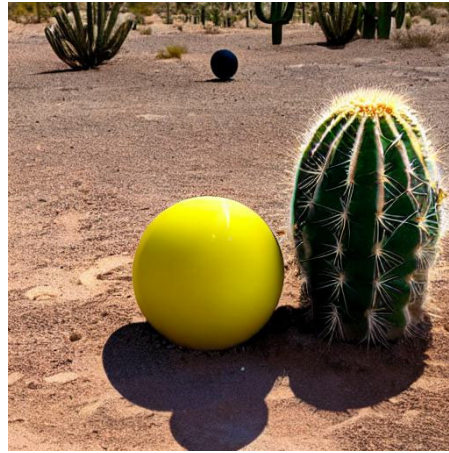
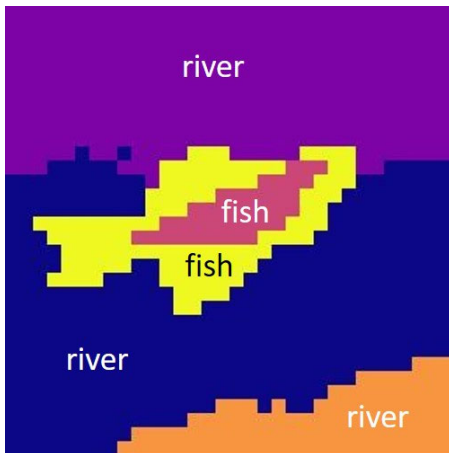
|              |              |     |
|--------------|--------------|-----|
| ■ 1          | ■ 1          | ■ 4 |
| ■ background | ■ background | ■ 1 |



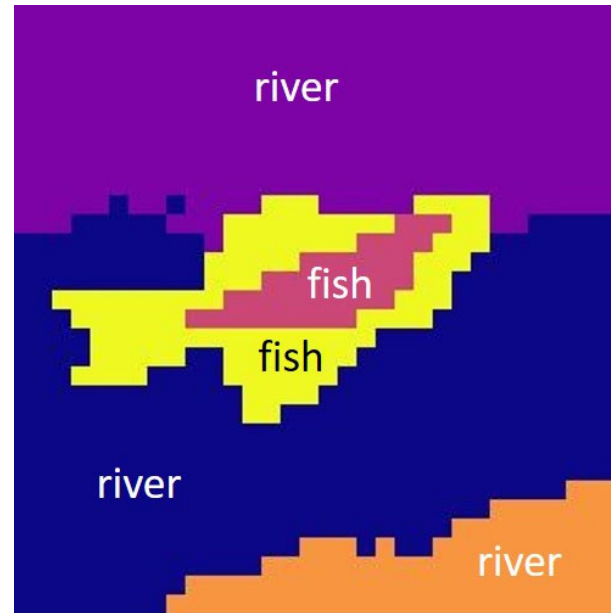
|     |              |     |
|-----|--------------|-----|
| ■ 1 | ■ background | ■ 4 |
|-----|--------------|-----|

1-cat, 4-sunglasses

# Self-Segmentation Results



# Self-Segmentation Results



# Cross-Image Attention

**Cross-Image Attention for Zero-Shot Appearance Transfer**

[Yuval Alaluf\\*](#), Daniel Garibi\*, Or Patashnik, [Hadar Averbuch Elor](#),  
[Daniel Cohen-Or](#)



# Motivation



Structure

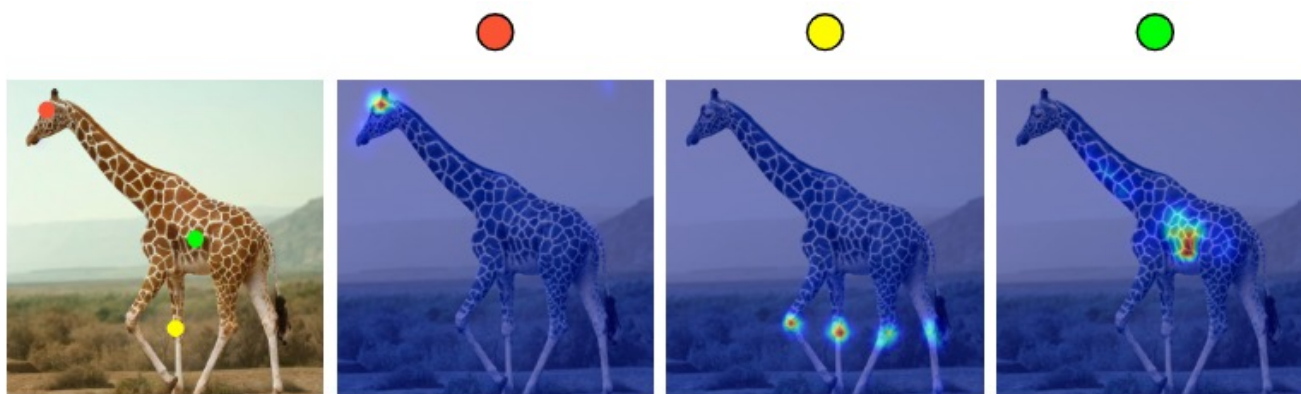


Appearance



Output

# The Roles of the Queries, Keys, and Values



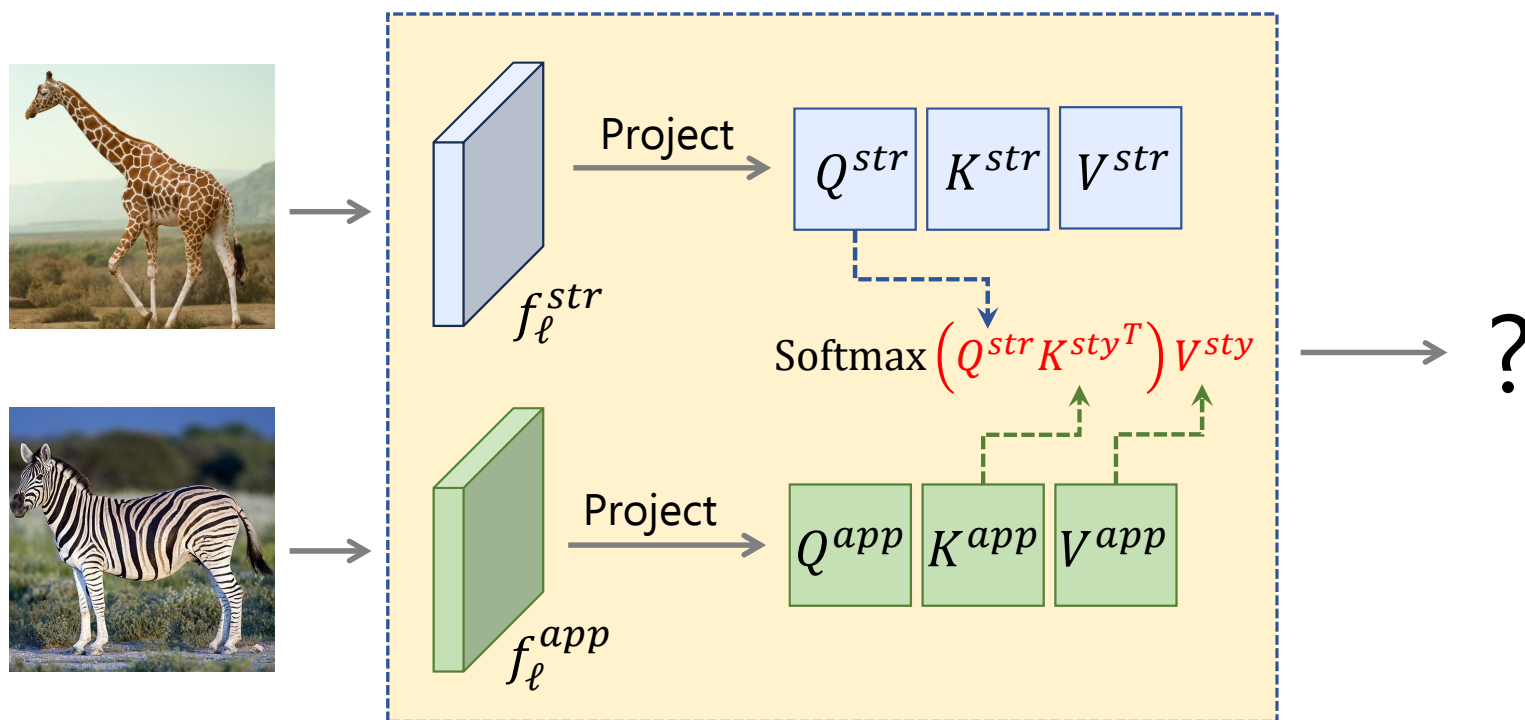
$$Q_{struct} \cdot K_{struct}^T$$



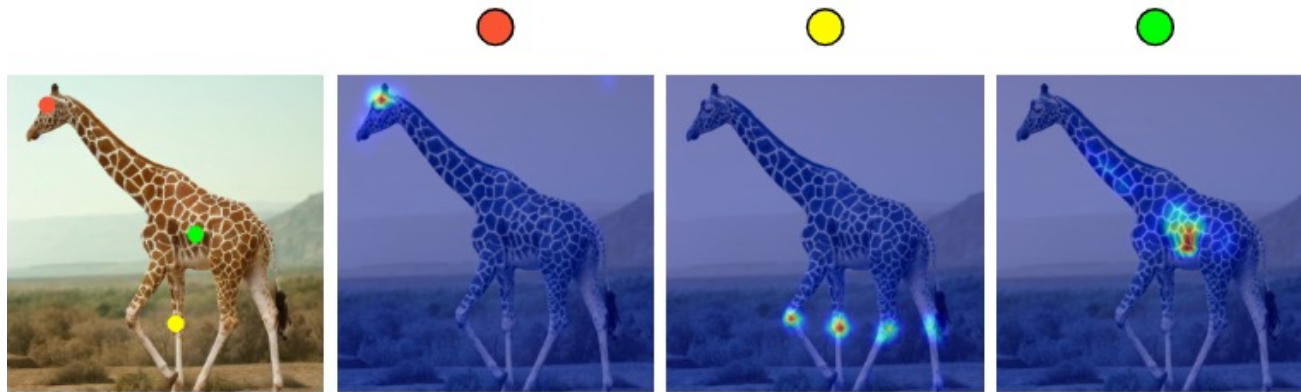
$$Q_{app} \cdot K_{app}^T$$

Self-attention maps, which focus on semantically similar regions in the image.

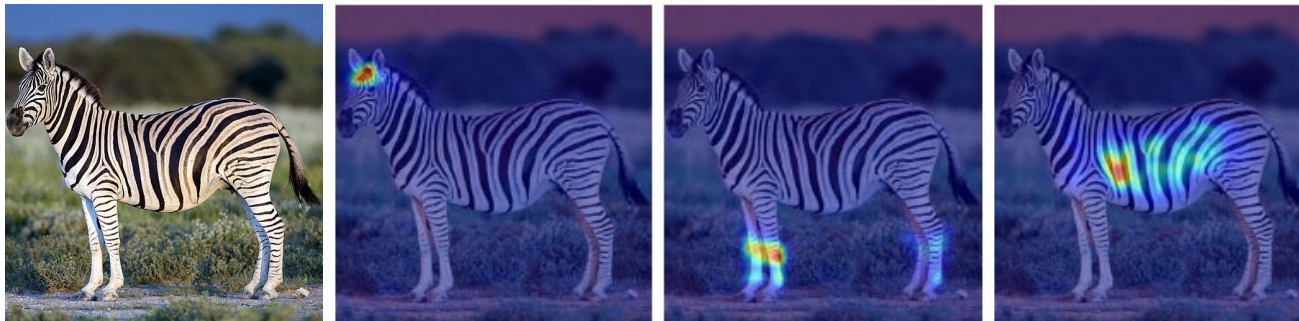
# What If we Swapped the Queries, Keys, and Values Between Different Images?



# The Roles of the Queries, Keys, and Values



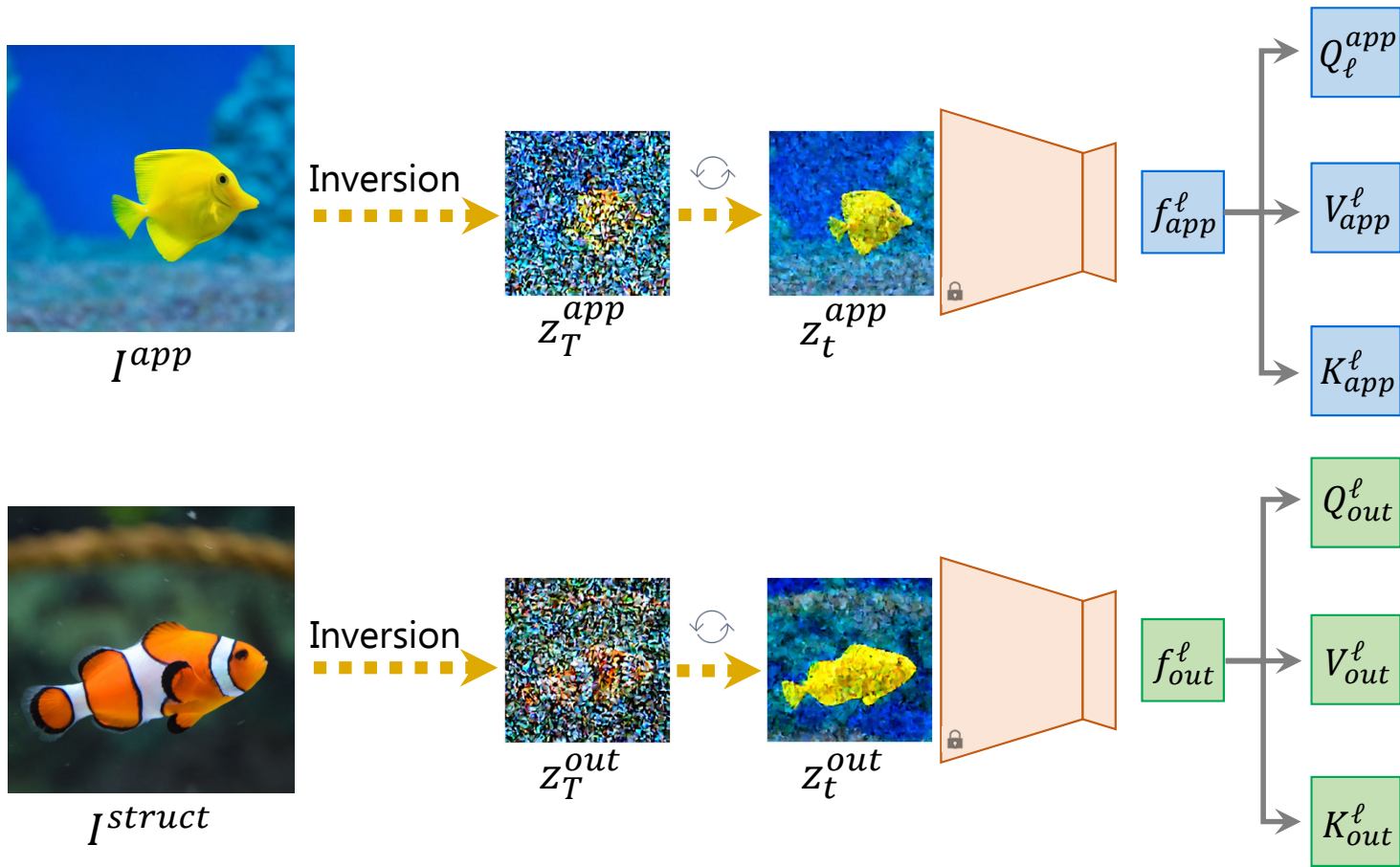
$$Q_{struct} \cdot K_{struct}^T$$



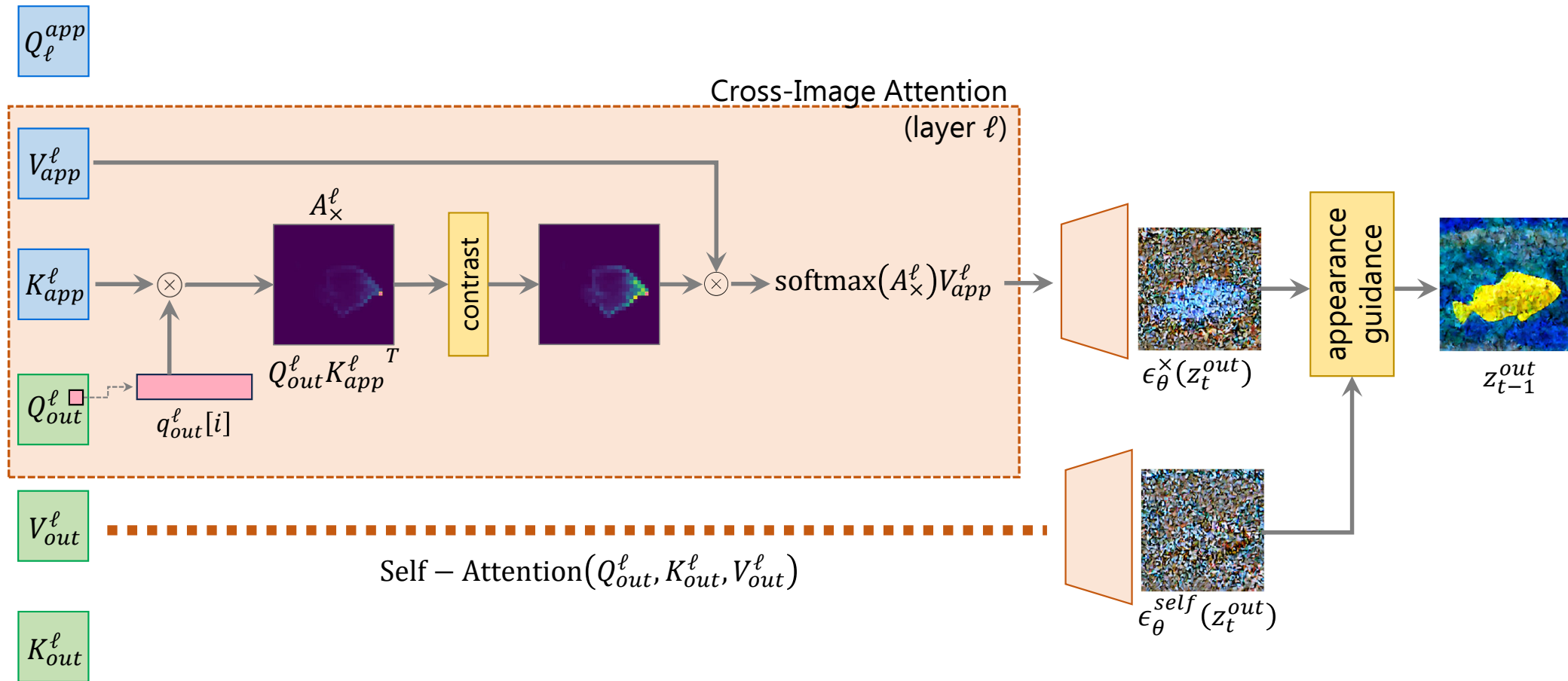
$$Q_{struct} \cdot K_{app}^T$$

Taking the **queries** from the structure image and the **keys** from the appearance image gives semantic correspondences between objects!

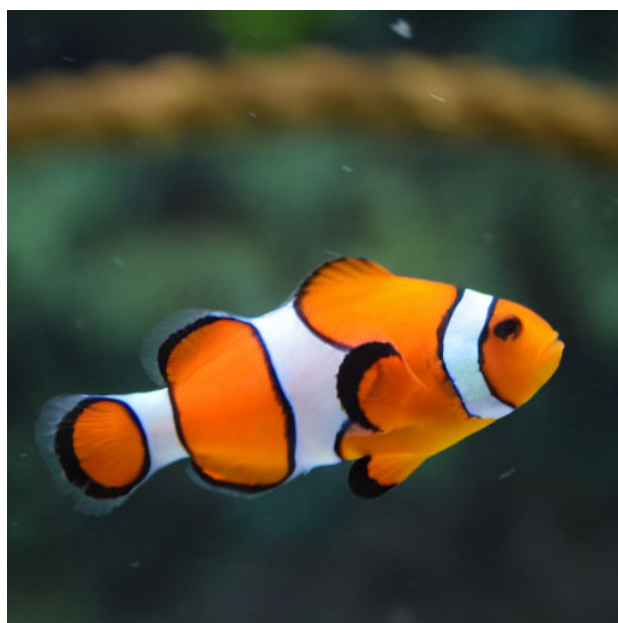
# The Cross-Image Attention



# The Cross-Image Attention



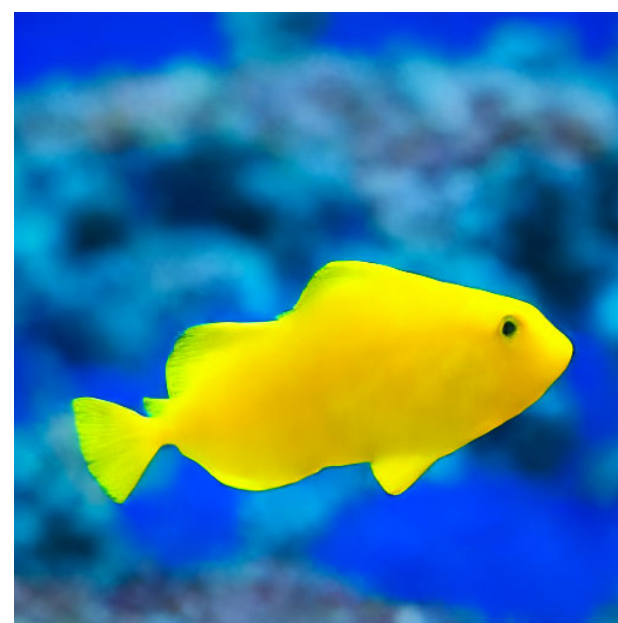
# The Cross-Image Attention



Structure



Appearance



Output

# Appearance Transfer Results



Structure



Appearance



Output

# Appearance Transfer Results



Structure



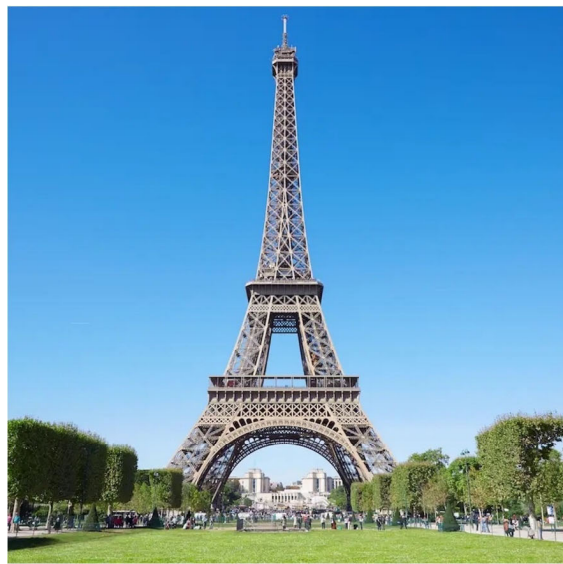
Appearance



Output

# Appearance Transfer Results

Eiffel Tower



Structure

Sagrada Família



Appearance



Output

# Appearance Transfer Results



Structure



Appearance



Output

# Appearance Transfer Results





Structure



Appearance



Output

# StyleAlign

[\(link\)](#)

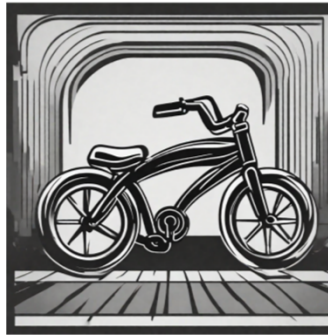
“Toy train...”



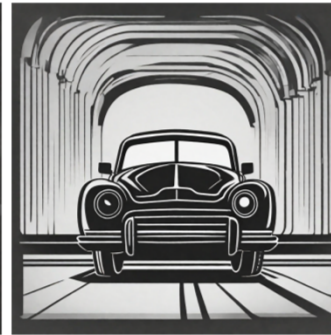
“Toy airplane...”



“Toy bicycle...”



“Toy car...”



“Toy boat...”



“...BW logo, high contrast.”

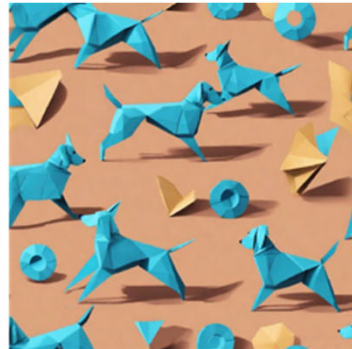


“...colorful, macro photo.”

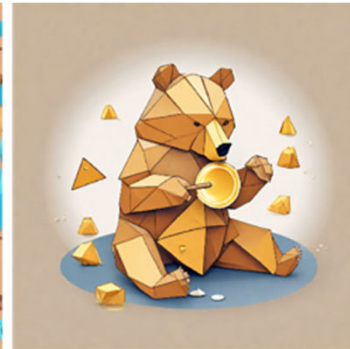
# Text-to-Image Generation



"A cat playing with a ball of wool..."



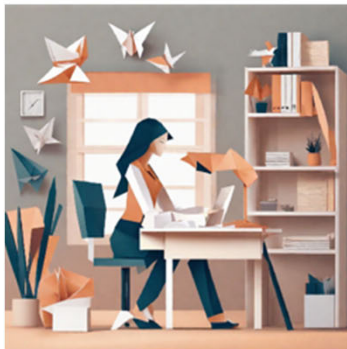
"A dog catching a frisbee..."



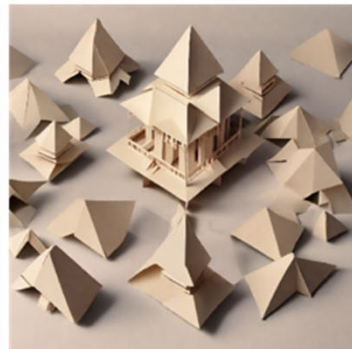
"A bear eating honey..."



"A whale playing with a ball..."



"A woman working in the office..."



"A temple..."



"A person riding a bike..."



"A cactus..."

"... in minimal origami style."

# Text-to-Image Generation with StyleAligned



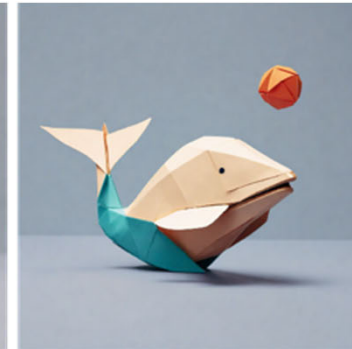
"A cat playing with a ball of wool..."



"A dog catching a frisbee..."



"A bear eating honey..."



"A whale playing with a ball..."



"A woman working in the office..."



"A temple..."



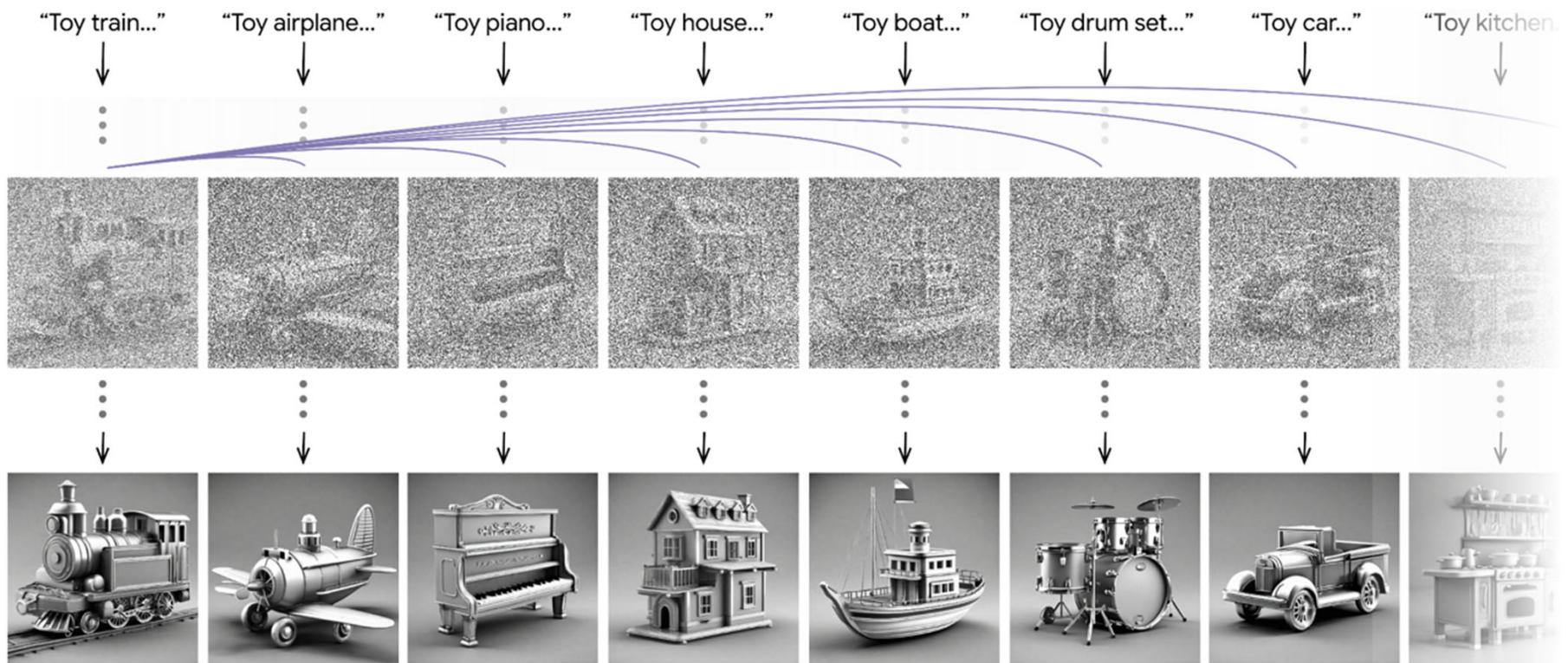
"A person riding a bike..."



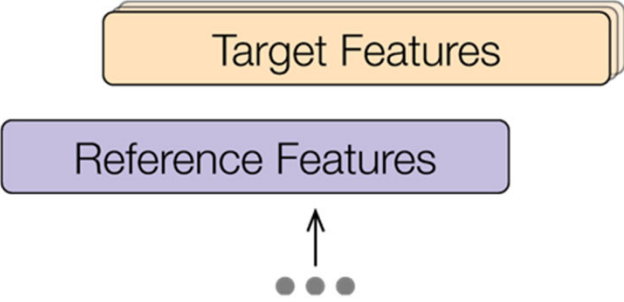
"A cactus..."

"... in minimal origami style."

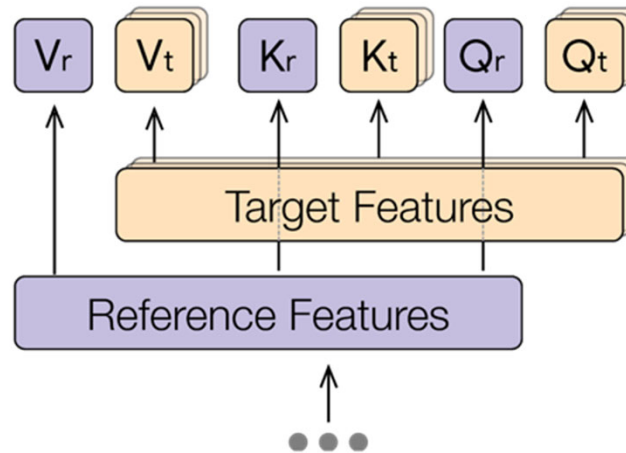
# Shared attention during the diffusion process



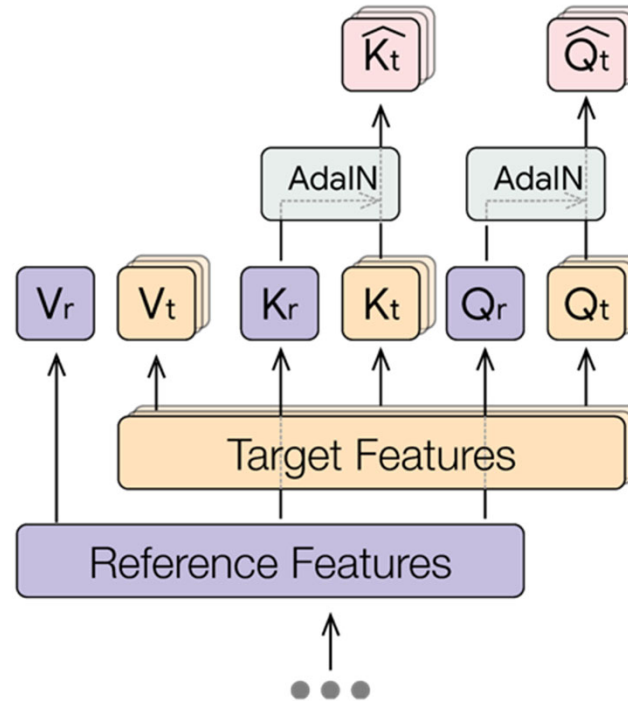
# Shared Attention Layer



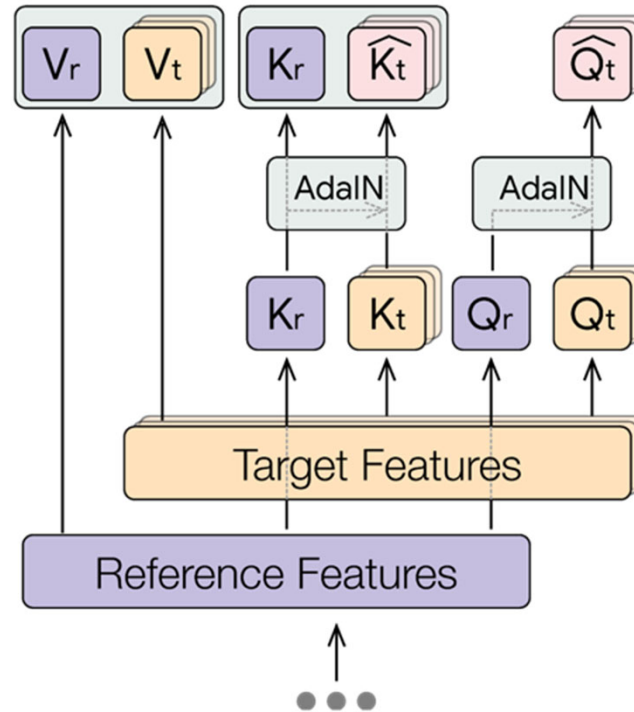
# Shared Attention Layer



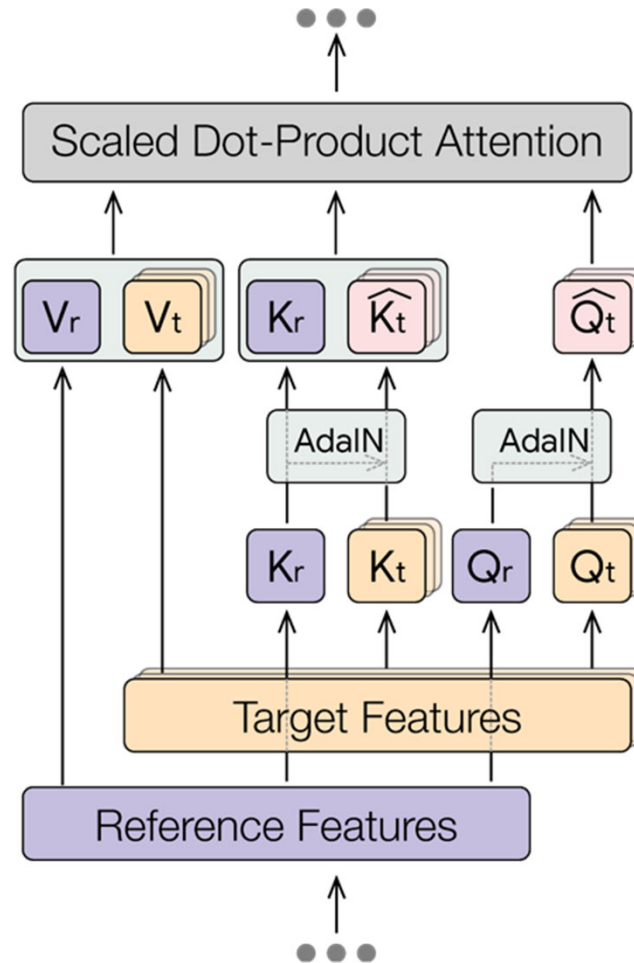
# Shared Attention Layer



# Shared Attention Layer



# Shared Attention Layer



# Style Aligned generation of Synthetic Images

“Firewoman...”



“Gardner...”



“Scientist ...”



“Police woman...”



“Saxophone player...”



“Painter...”



“Astronaut...”



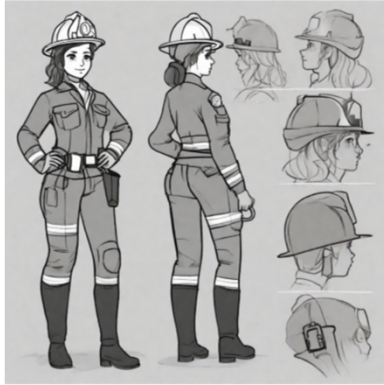
“Taxi Driver...”



“...made of claymation, stop motion animation.”

# Style Aligned generation of Synthetic Images

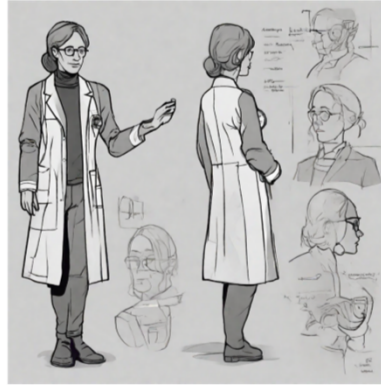
“Firewoman...”



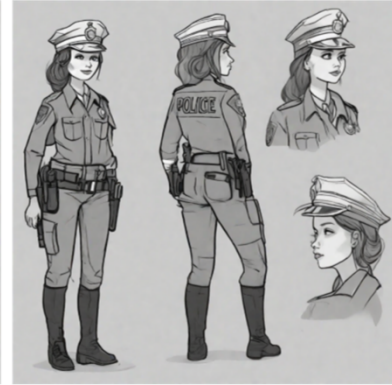
“Gardner...”



“Scientist ...”



“Police woman...”



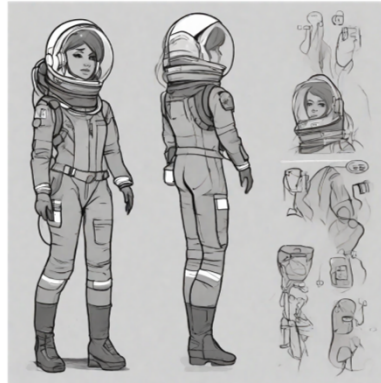
“Saxophone player...”



“Painter...”



“Astronaut...”



“Taxi Driver...”



“...sketch, character sheet.”

# Style Aligned generation of Synthetic Images

“Firewoman...”



“Gardner...”



“Scientist ...”



“Police woman...”



“Saxophone player...”



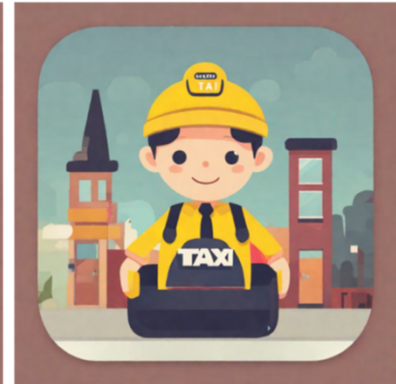
“Painter...”



“Astronaut...”



“Taxi Driver...”



“ ...in minimal flat design illustration.”

# Style Aligned generation from an Input Image

Reference image



Space rocket



Boy riding a bicycle



Matterhorn mountain



Mime artist



Seattle needle



# Style Aligned generation from an Input Image

Reference image



Space rocket



Boy riding a bicycle



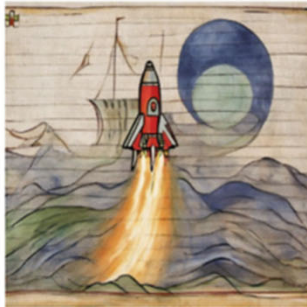
Matterhorn mountain



Mime artist

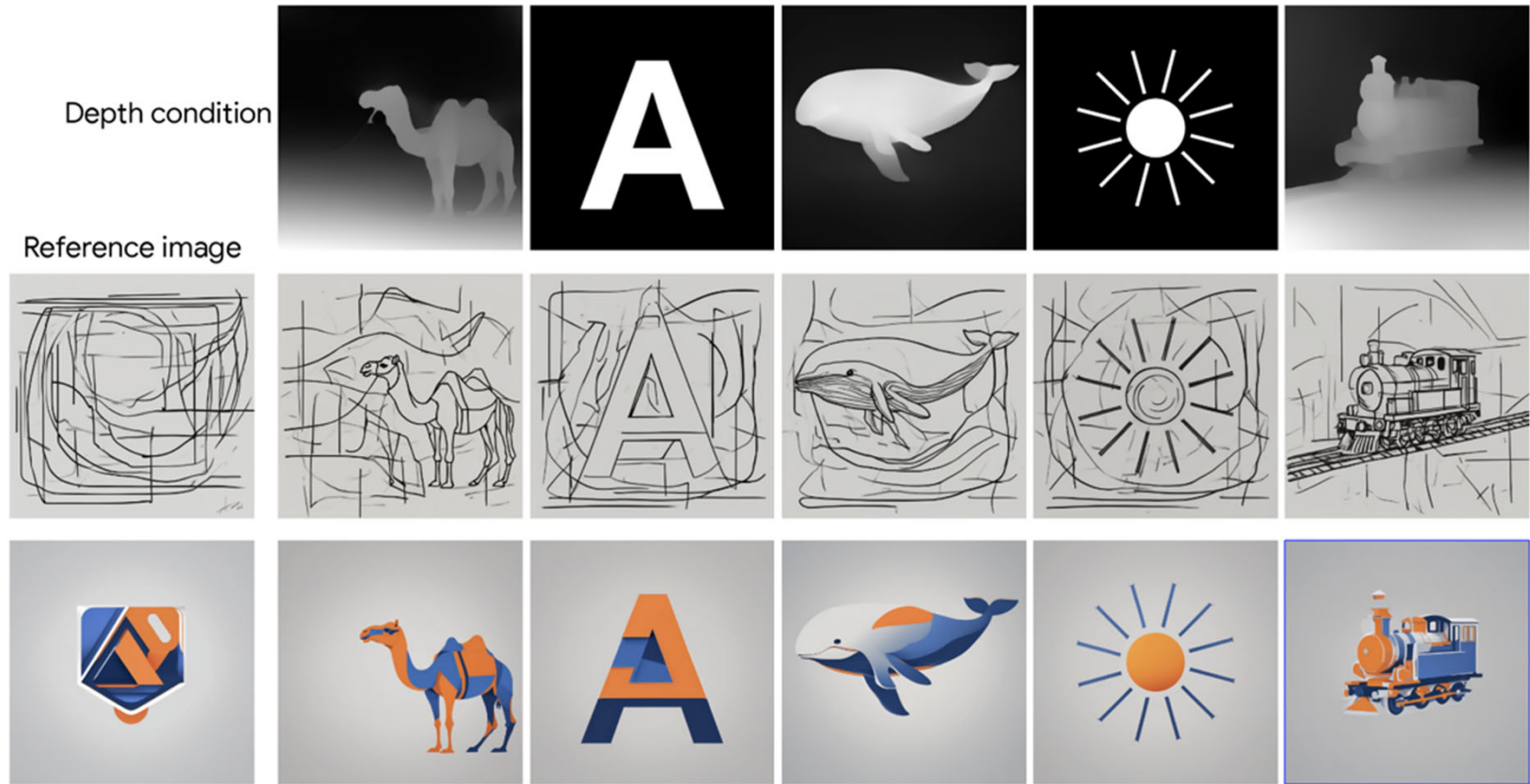


Seattle needle



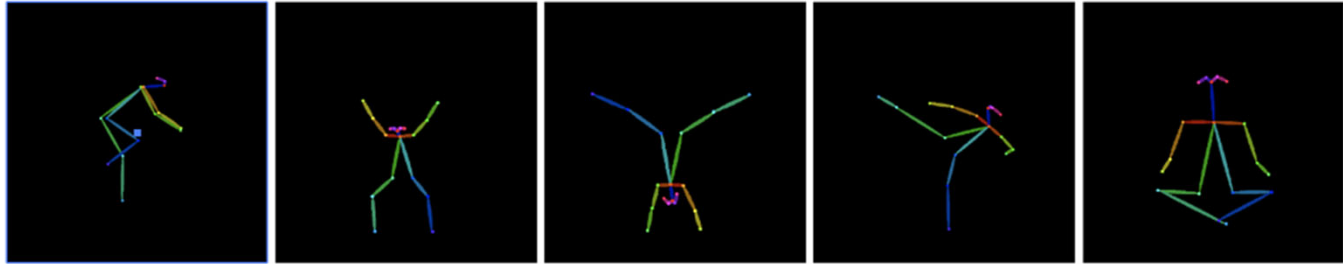
**StyleAligned with other methods**

# ControlNet + StyleAligned

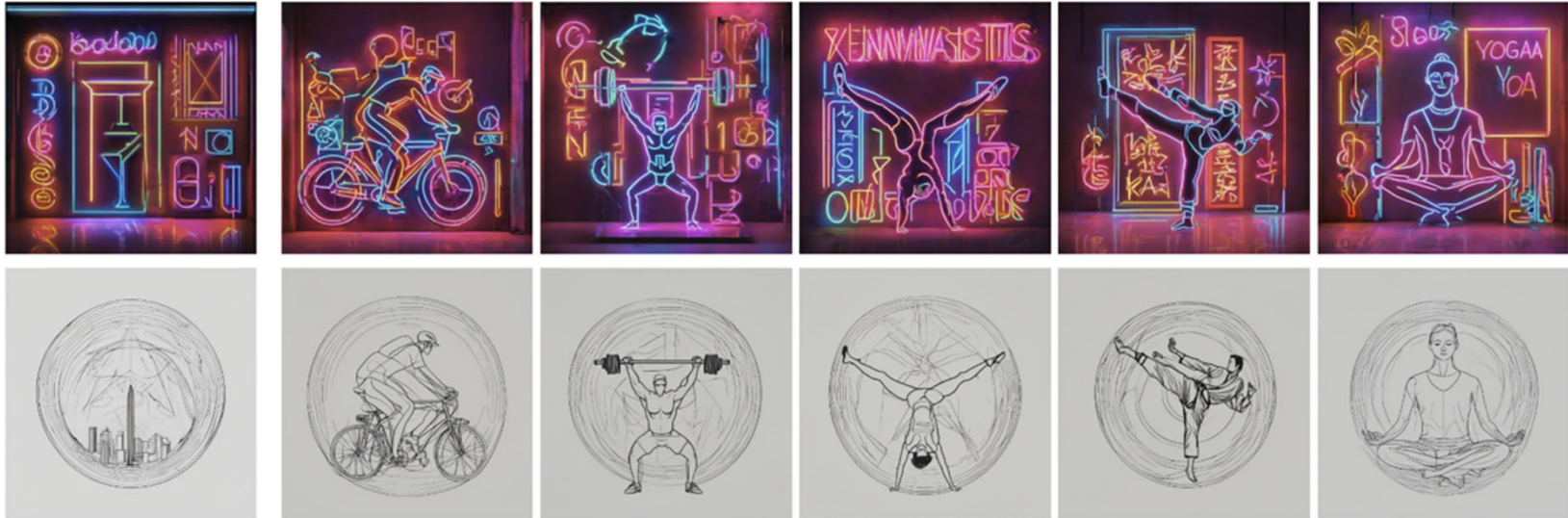


# ControlNet + StyleAligned

Pose condition

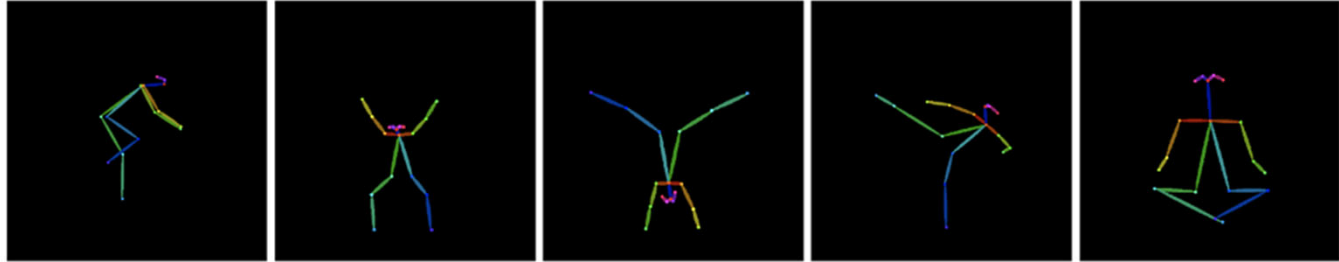


Reference image

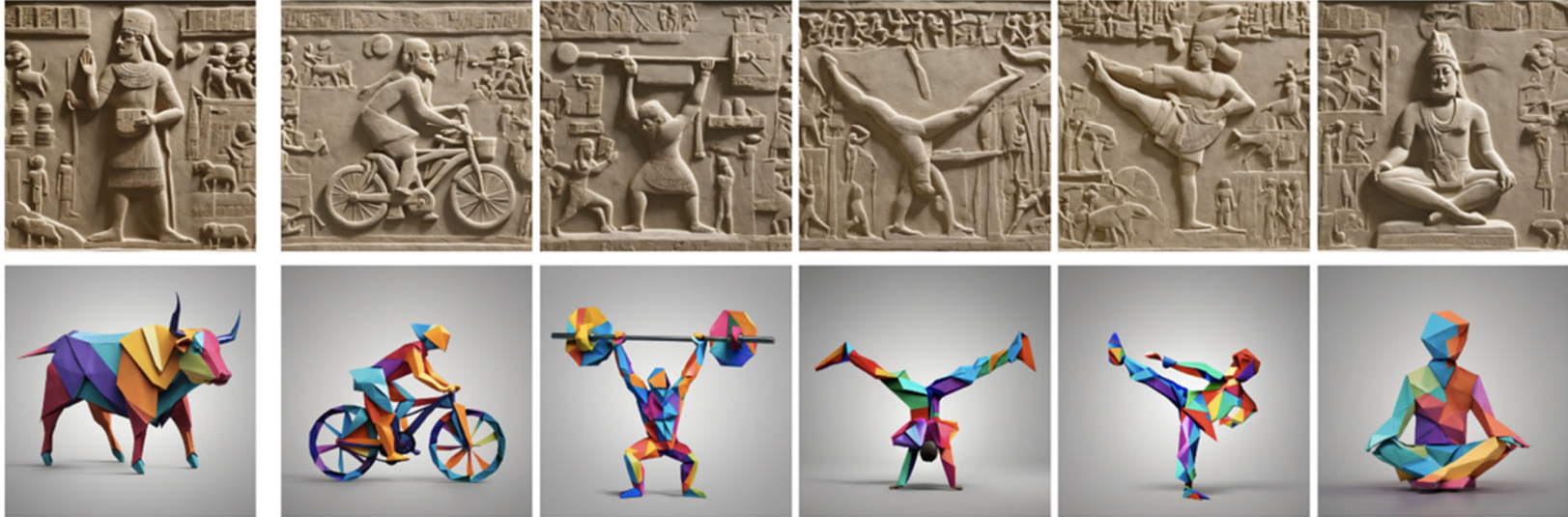


# ControlNet + StyleAligned

Pose condition



Reference image



# Textual Inversion+Dreambooth

Personalized content



"<V object> in the style of a beautiful paper-cut art."

# + StyleAligned

Personalized  
content



Reference image



"<V object> in the style of a beautiful paper-cut art."

# W.O AdalN

Personalized content



Reference image



"<V object> in the style of a beautiful paper-cut art."

# DreamBooth + StyleAligned

Personalized content



Reference image



# MultiDiffusion + StyleAligned

Reference image



"A poster in a flat design style."



"Houses in a flat design style."



"Mountains in a flat design style."



"Girrafes in a flat design style."

# MultiDiffusion + StyleAligned

Reference image



"A poster in a flat design style."



W,O shared attention



W,O Attention Adaln



StyleAligned full

# MultiDiffusion in Multi Styles

Left Reference



Right Reference



# MultiDiffusion in Multi Styles

Left reference



Right reference

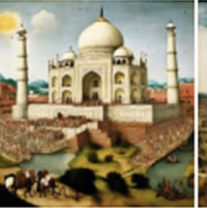
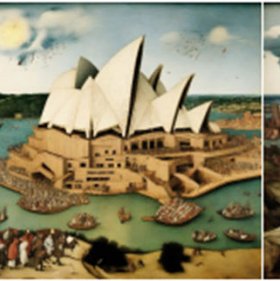




Left Reference

Right Reference







# Diffusion Models for Visual Computing

Paul Guerrero

Niloy Mitra, Daniel Cohen-Or, Minhyuk Sung, Chun-Hao Huang, Duygu Ceylan

## Part 4: Personalization & Editing



[https://geometry.cs.ucl.ac.uk/courses/diffusion4VC\\_eg24/](https://geometry.cs.ucl.ac.uk/courses/diffusion4VC_eg24/)

# Presentation Schedule

Introduction to Diffusion Models

Guidance and Conditioning Sampling

Attention

Break

**Personalization and Editing**

Beyond Single Images

Diffusion Models for 3D Generation

# Personalization

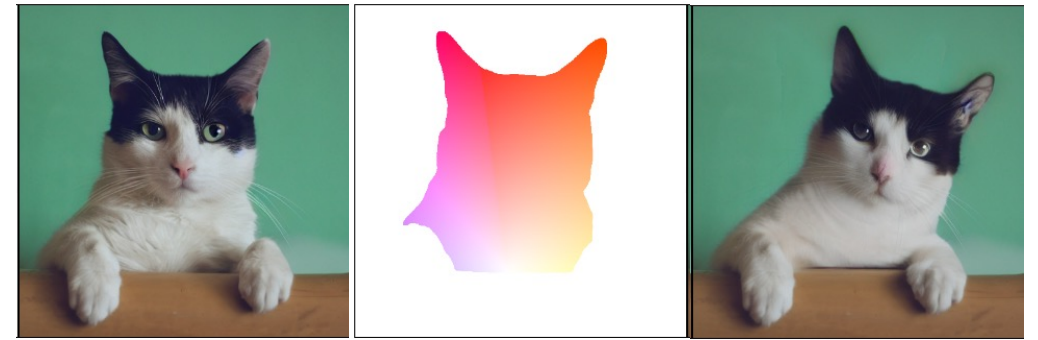


Input images

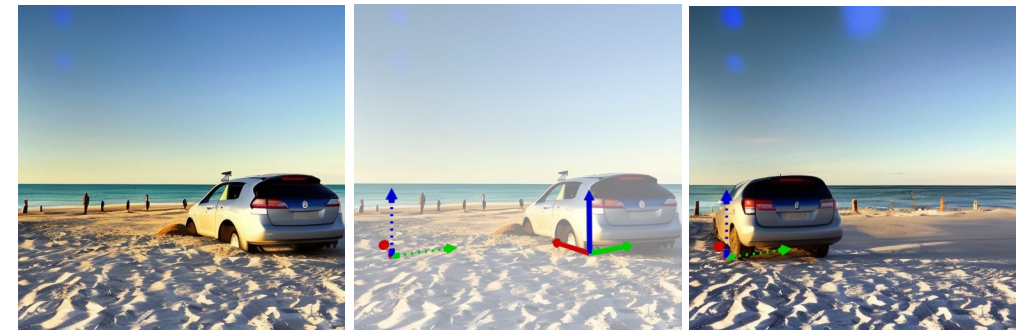


DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation, Ruiz et al, CVPR 2023

# Editing



Motion Guidance: Diffusion-Based Image Editing with Differentiable Motion Estimators, Geng and Owens, ICLR 2024



Diffusion Handles: Enabling 3D Edits for Diffusion Models by Lifting Activations to 3D Pandey et al, CVPR 2024

# Personalization

*“a hyper-realistic digital painting of a happy girl, with brown eyes”*

Without Personalization



Generated with StabelDiffusion 2.1

With Personalization



ConsiStory: Training-Free Consistent Text-to-Image Generation  
Tewel et al., ArXiv Feb. 2024

# Personalization

With Personalization



Same subject in different settings.

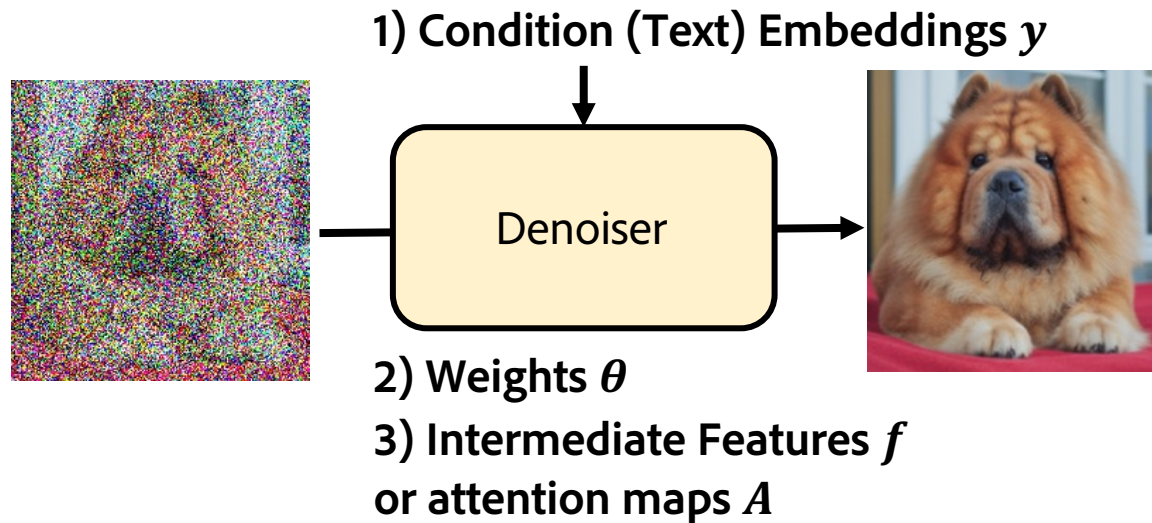
Personalization:

Generative Model  
+ **Identity Preservation**

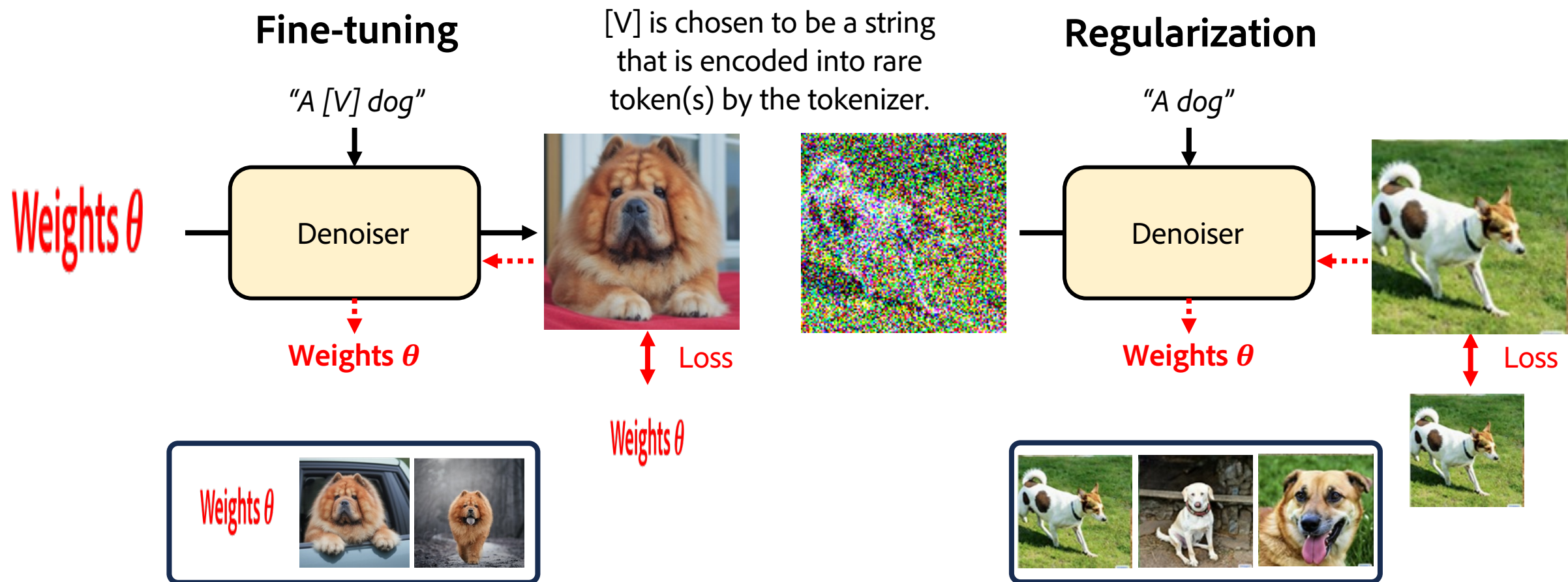
ConsiStory: Training-Free Consistent Text-to-Image Generation  
Tewel et al., ArXiv Feb. 2024

# Identity Preservation

How can we represent the identity of a subject?



# ID Preservation by Fine-Tuning Denoiser Weights

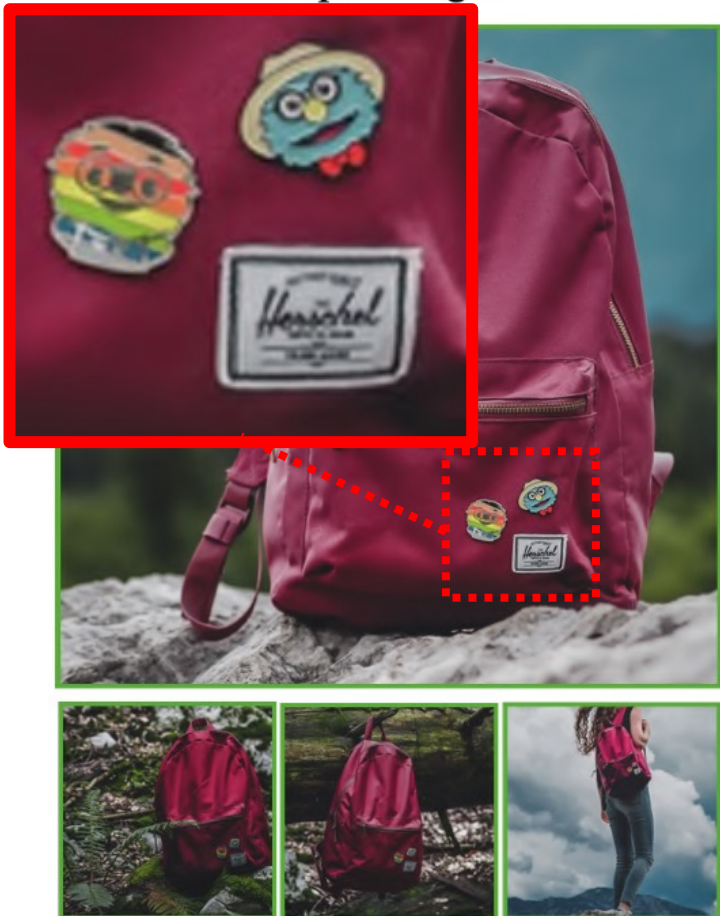


3-5 example images showing the identity.

Regularization data can be generated by the original, non-finetuned model, or can come from a large dataset.

# ID Preservation by Fine-Tuning Denoiser Weights

Input images



A [V] backpack in the Grand Canyon



A [V] backpack with the night sky



A [V] backpack in the city of Versailles



A wet [V] backpack in water



A [V] backpack in Boston

DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation, Ruiz et al., CVPR 2023

# ID Preservation by Fine-Tuning Denoiser Weights

Input images



Vincent Van Gogh



Michelangelo



Rembrandt



Johannes Vermeer



Pierre-Auguste Renoir



Leonardo da Vinci



Code



<https://github.com/google/dreambooth>



Code



<https://huggingface.co/docs/diffusers/en/training/dreambooth>

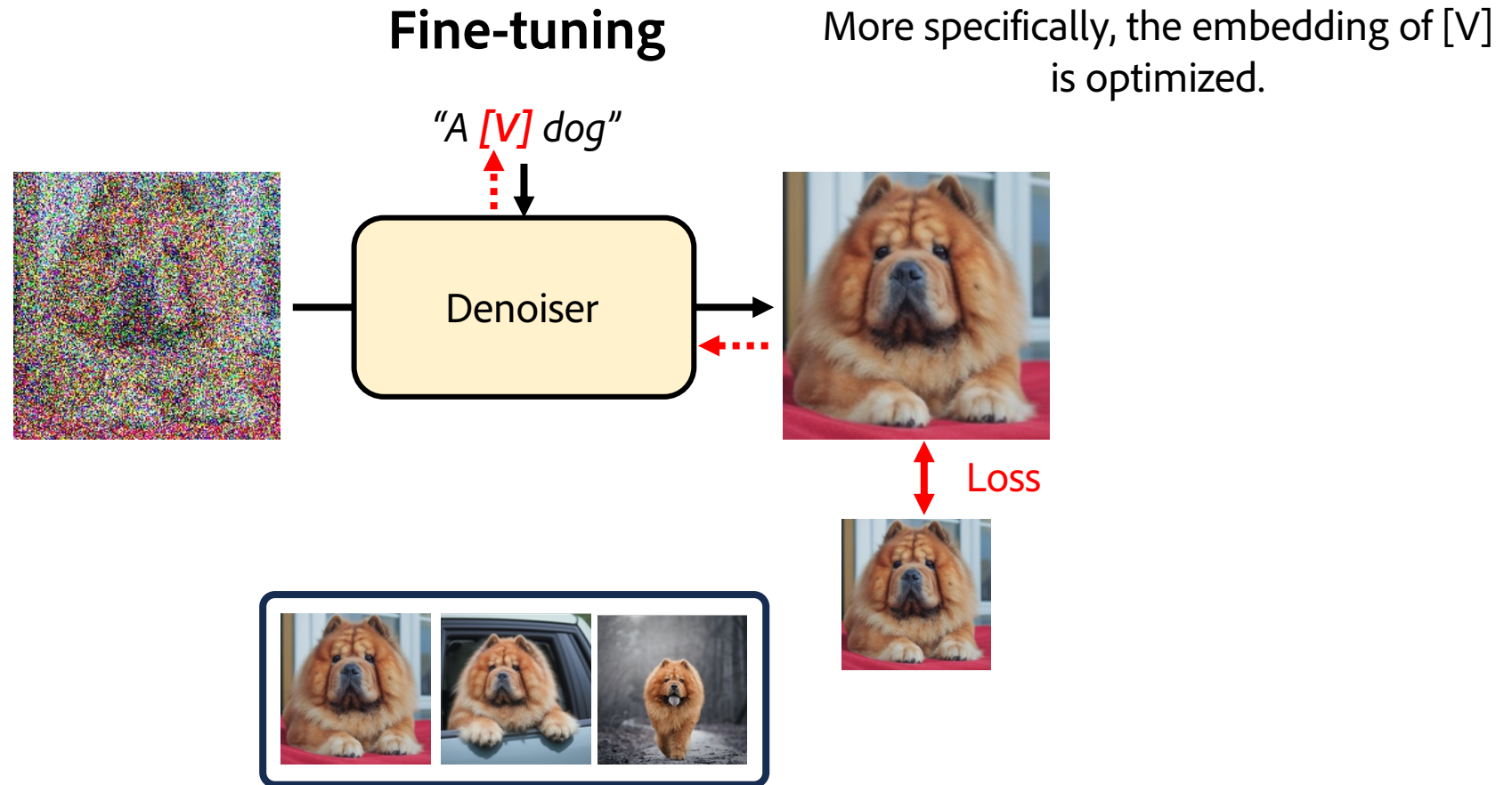


Demo



<https://huggingface.co/spaces/multimodalart/dreambooth-training>

# ID Preservation by Fine-Tuning Text Embeddings



An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion,  
Gal et al., ICLR 2023

# ID Preservation by Fine-Tuning Text Embeddings

Compared to fine-tuning denoiser weights:

- Only the optimized embedding needs to be stored, not the full denoiser weights.
- ID preservation is a bit weaker.

Input Images



“A photo of  $[V]$ ”



“A photo of  $S_*$ ”

“A photo of  $S_*$   
on the beach”

“A photo of  $S_*$   
on the moon”

“Elmo holding  
a  $S_*$ ”

An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion,  
Gal et al., ICLR 2023

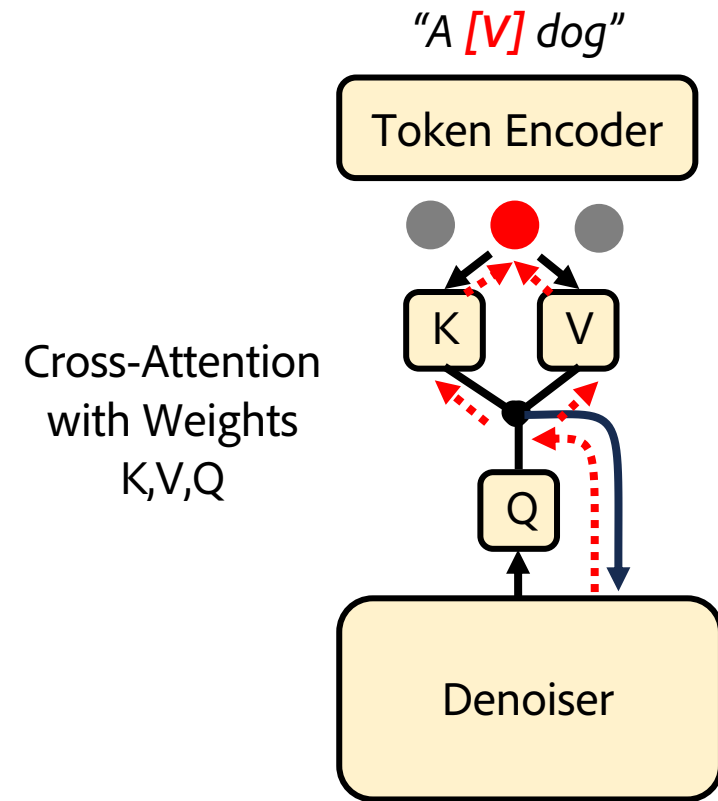
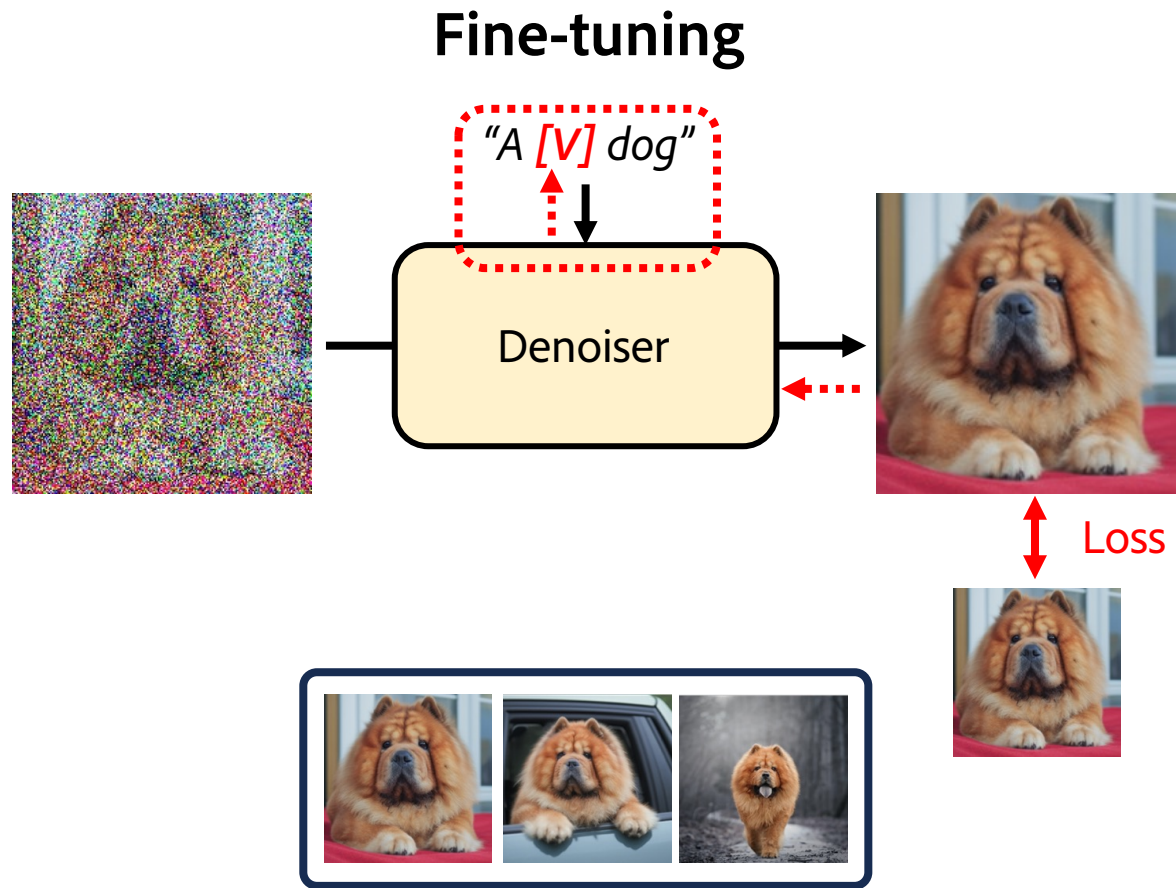


Code



[https://github.com/rinongal/textual\\_inversion](https://github.com/rinongal/textual_inversion)

# Fine-Tuning Text Embeddings & Cross-Att. Weights



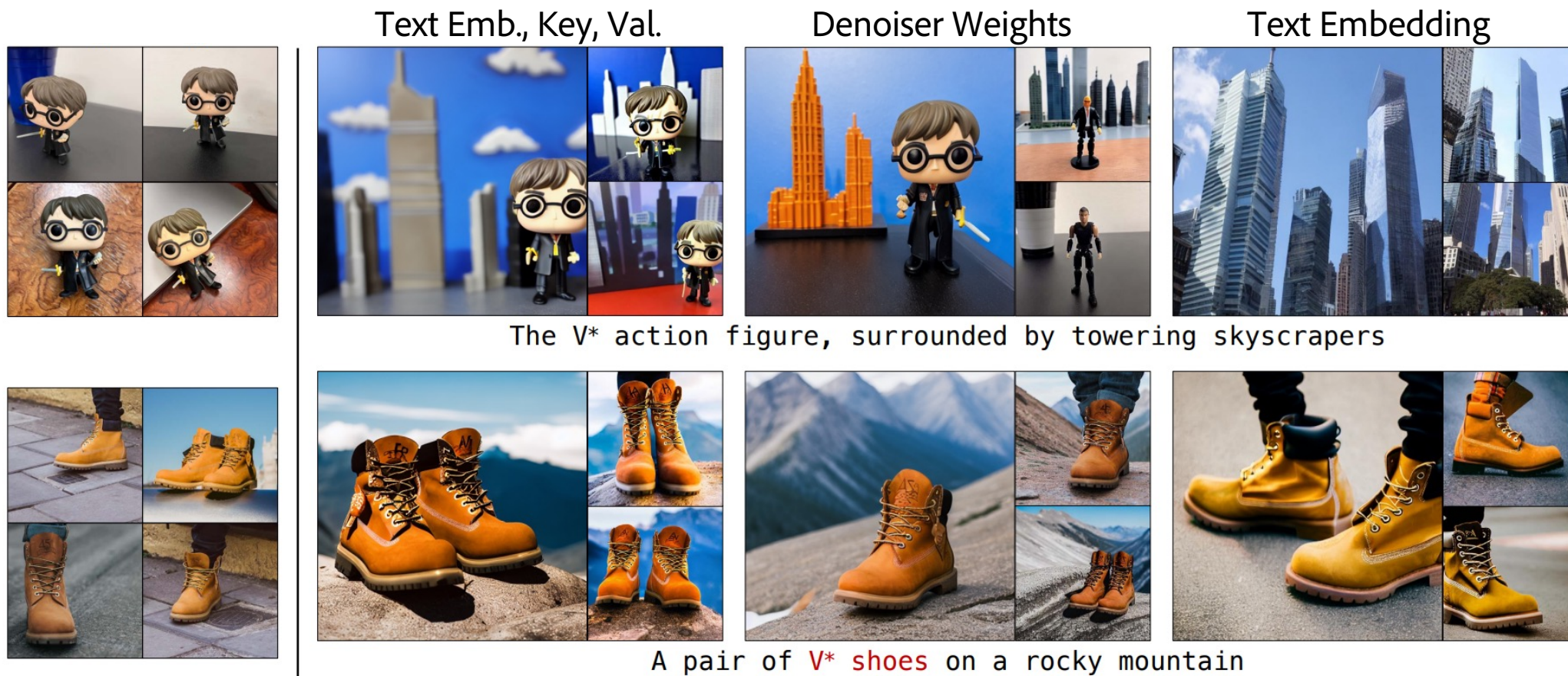
Multi-Concept Customization of Text-to-Image Diffusion, Kumari et al, CVPR 2023

Key-Locked Rank One Editing for Text-to-Image Personalization, Tewel et al, SIGGRAPH 2023

# Fine-Tuning Text Embeddings & Cross-Att. Weights

Fine-tuning text embeddings, keys and values.

ID preservation is close to tuning denoiser weights, while requiring less storage.



Multi-Concept Customization of Text-to-Image Diffusion, Kumari et al, CVPR 2023

# Fine-Tuning Text Embeddings & Cross-Att. Weights

Fine-tuning text embeddings and values only.

ID preservation is slightly worse than tuning denoiser weights but follows the prompt better.

Text Emb., Value



pot\*

Text Emb., Key, Val.



A pot\* with mountains and sunset in background

Denoiser Weights



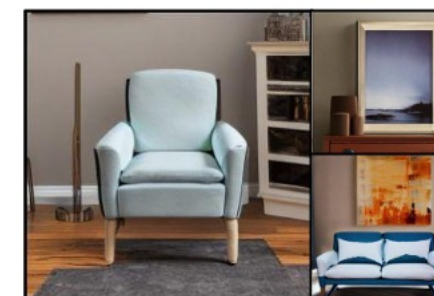
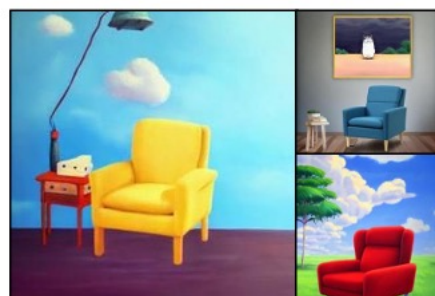
Text Embedding



chair\*



A chair\* oil painting ghibli inspired



## ***Custom Diffusion***

*Multi-Concept Customization of Text-to-Image  
Diffusion*

(fine-tuning keys, values and text embeddings)



Code



<https://github.com/adobe-research/custom-diffusion>



Code



[https://huggingface.co/docs/diffusers/en/training/custom\\_diffusion](https://huggingface.co/docs/diffusers/en/training/custom_diffusion)

## ***Perfusion***

*Key-Locked Rank One Editing for  
Text-to-Image Personalization*

(fine-tuning values and text embeddings only)



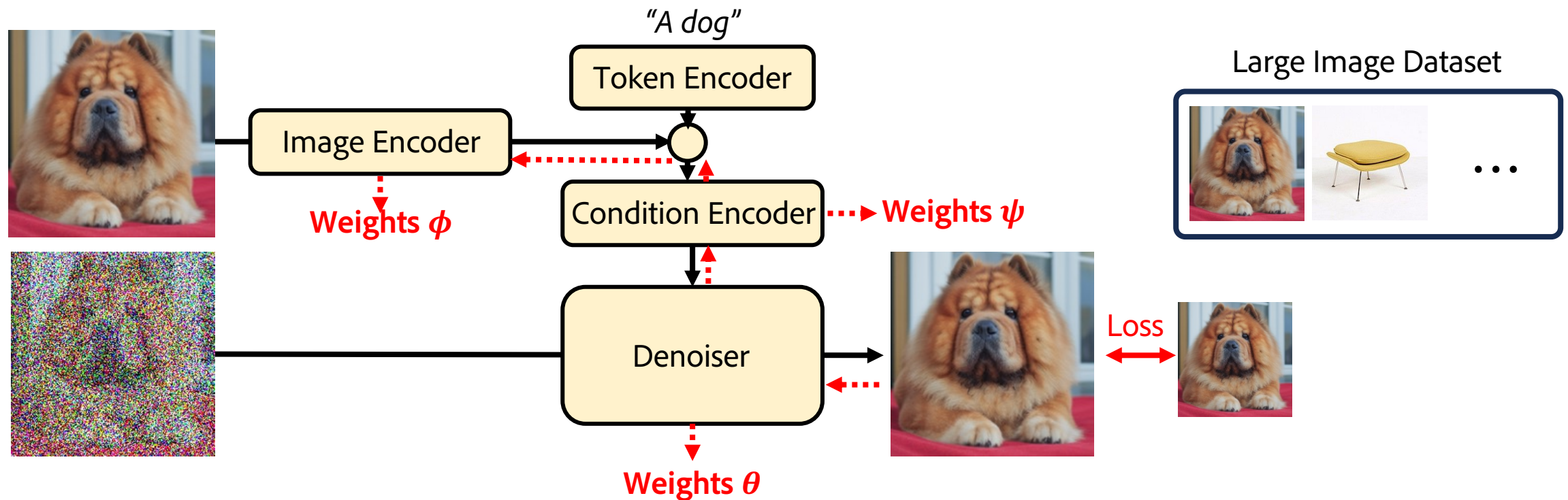
Code



<https://github.com/ChenDarYen/Key-Locked-Rank-One-Editing-for-Text-to-Image-Personalization>

# ID Preservation with Learned Condition Encoders

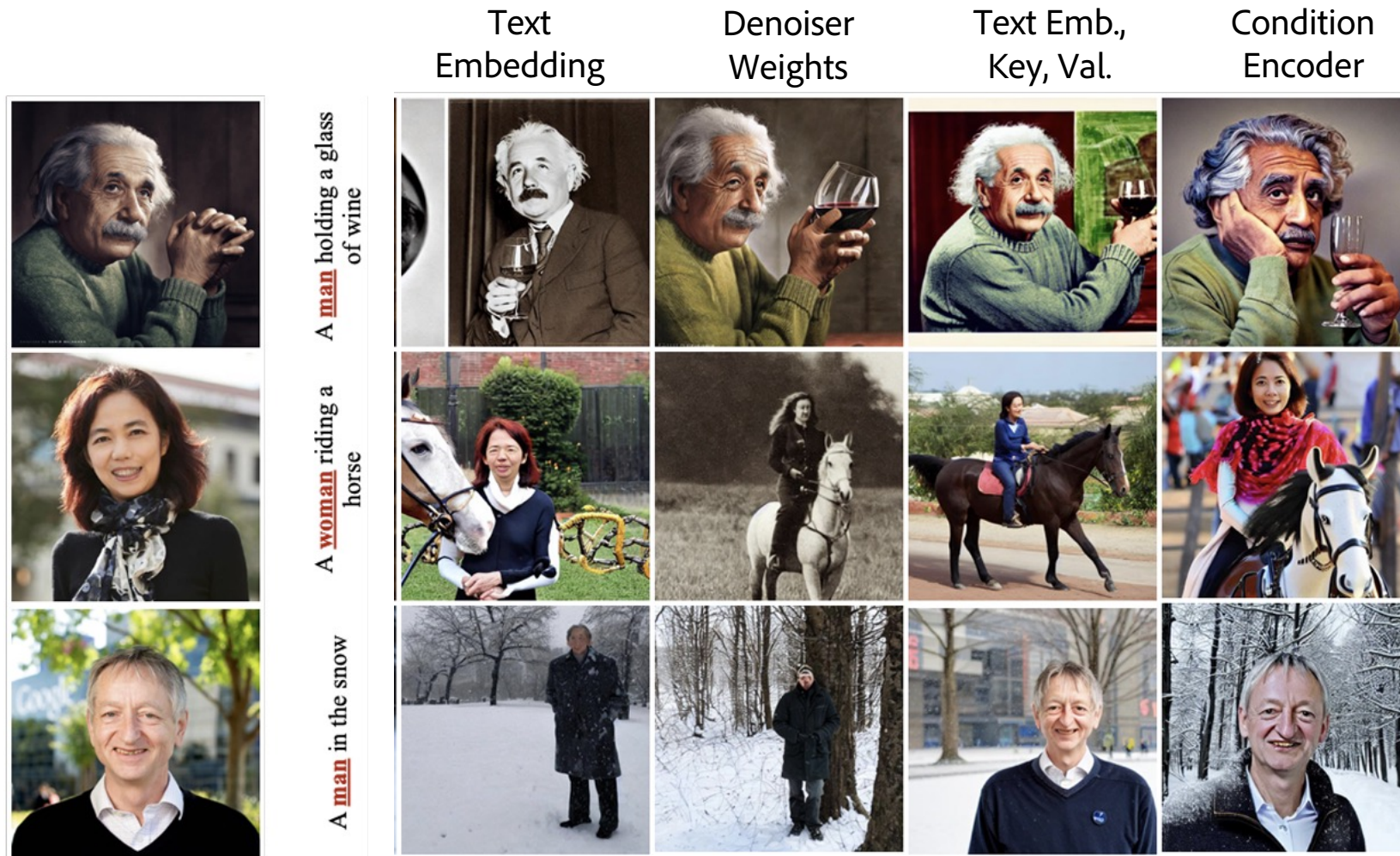
Motivation: avoid the need to fine-tune for each object identity.



FastComposer: Tuning-Free Multi-Subject Image Generation with Localized Attention, Gal et al., ArXiv May 2023

BLIP-Diffusion: Pre-trained Subject Representation for Controllable Text-to-Image Generation and Editing, Li et al., NeurIPS 2024

# ID Preservation with Learned Condition Encoders



FastComposer: Tuning-Free Multi-Subject Image Generation with Localized Attention, Gal et al, ArXiv May 2023

# ID Preservation with Learned Condition Encoders



BLIP-Diffusion: Pre-trained Subject Representation for Controllable Text-to-Image Generation and Editing, Li et al., NeurIPS 2024

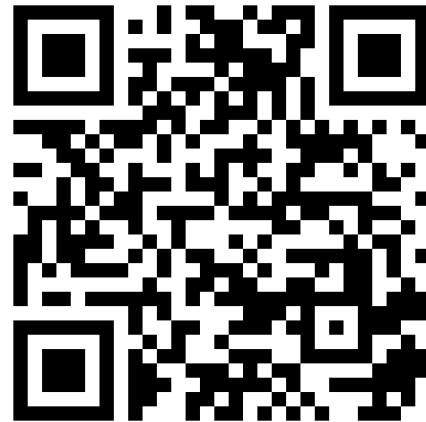
## ***FastComposer***

*uning-Free Multi-Subject Image Generation  
with Localized Attention  
(CLIP-based image encoder)*



Code

Demo



<https://github.com/mit-han-lab/fastcomposer/tree/main>

<https://replicate.com/cjwbw/fastcomposer>

## ***BLIP-Diffusion***

*Pre-trained Subject Representation for Controllable  
Text-to-Image Generation and Editing  
(BLIP-based image encoder)*



Code



Code

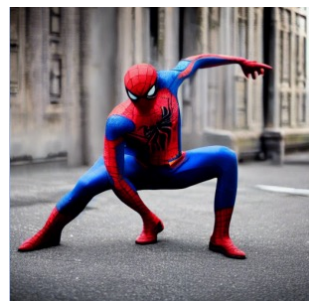


[https://huggingface.co/docs/diffusers/en/api/pipelines/blip\\_diffusion](https://huggingface.co/docs/diffusers/en/api/pipelines/blip_diffusion)

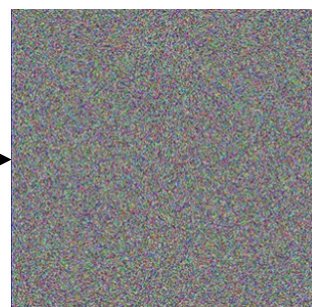
<https://github.com/salesforce/LAVIS/tree/main/projects/blip-diffusion>

# ID Preservation through Intermediate Features

Non-Generated  
Image

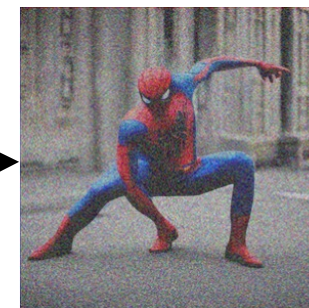


Inversion



"A photo of spiderman"

Denoiser



same noise

Intermediate Features  $f$   
or Attention Maps  $A$

Denoiser



"A photo of a statue  
in the snow"

## Feature injection approaches:

- 1) Directly overwrite denoiser features with target features  $f$ .
- 2) Guidance energy towards target features  $f$ .
- 3) Cross-Attention from denoiser features to target features  $f$ .

Plug-and-Play Diffusion Features for Text-Driven Image-to-Image Translation,  
Tumanyan et al., CVPR 2023

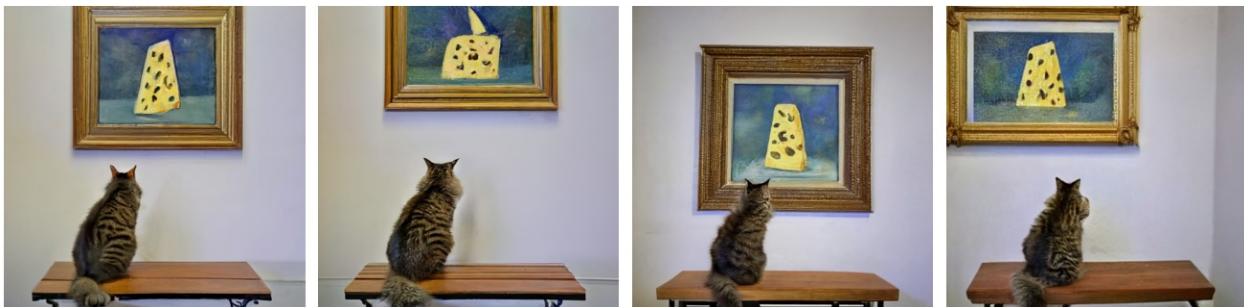
Diffusion Self-Guidance for Controllable Image Generation,  
Epstein et al., NeurIPS 2023

MasaCtrl: Tuning-Free Mutual Self-Attention Control for Consistent Image  
Synthesis and Editing, Cao et al., ICCV 2023

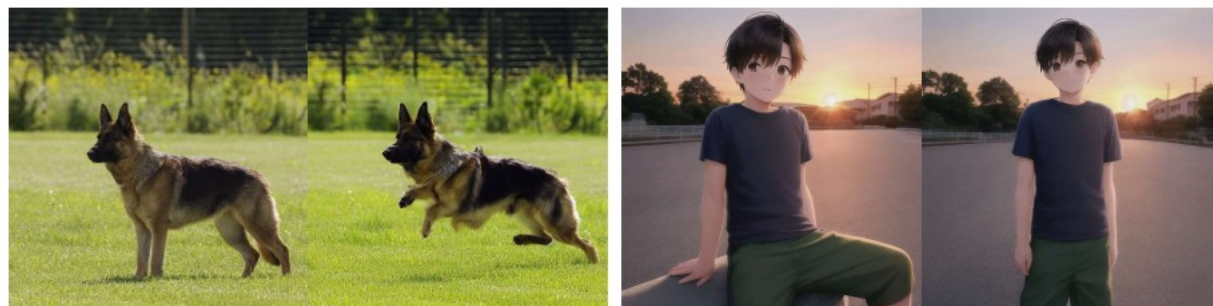
# ID Preservation through Intermediate Features

- No training or fine-tuning needed.
- Intermediate features entangle identity with location and context, thus it requires additional work to allow for large changes in locations & contexts.

Guidance energy towards target features  $f$ .



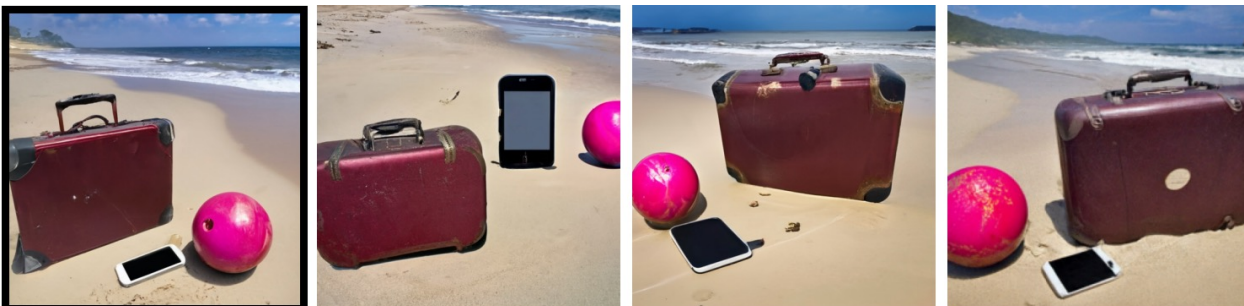
Cross-Attention from denoiser features to target features  $f$ .



Input real image

“... jumping ...”

“A sitting boy” → “... standing ...”



Diffusion Self-Guidance for Controllable Image Generation,  
Epstein et al, NeurIPS 2023

MasaCtrl: Tuning-Free Mutual Self-Attention Control for Consistent Image  
Synthesis and Editing, Cao et al, ICCV 2023

# ID Preservation through Intermediate Features

- No training or fine-tuning needed.
- Intermediate features entangle identity with location and context, thus it requires additional work to allow for large changes in locations & contexts.

Directly overwrite denoiser features with target features  $f$ .



Plug-and-Play Diffusion Features for Text-Driven Image-to-Image Translation,  
Tumanyan et al., CVPR 2023

## MasaCtrl

*Tuning-Free Mutual Self-Attention Control for  
Consistent Image Synthesis and Editing*  
(Cross-Attention-Based Feature Injection)



Code



<https://github.com/TencentARC/MasaCtrl>

## PnP-Diffusers

*Plug-and-Play Diffusion Features for Text-Driven  
Image-to-Image Translation*  
(Overwrite-Based Feature Injection)



Demo



<https://huggingface.co/spaces/hysts/PnP-diffusion-features>



Code



<https://github.com/MichalGeyer/pnp-diffusers>

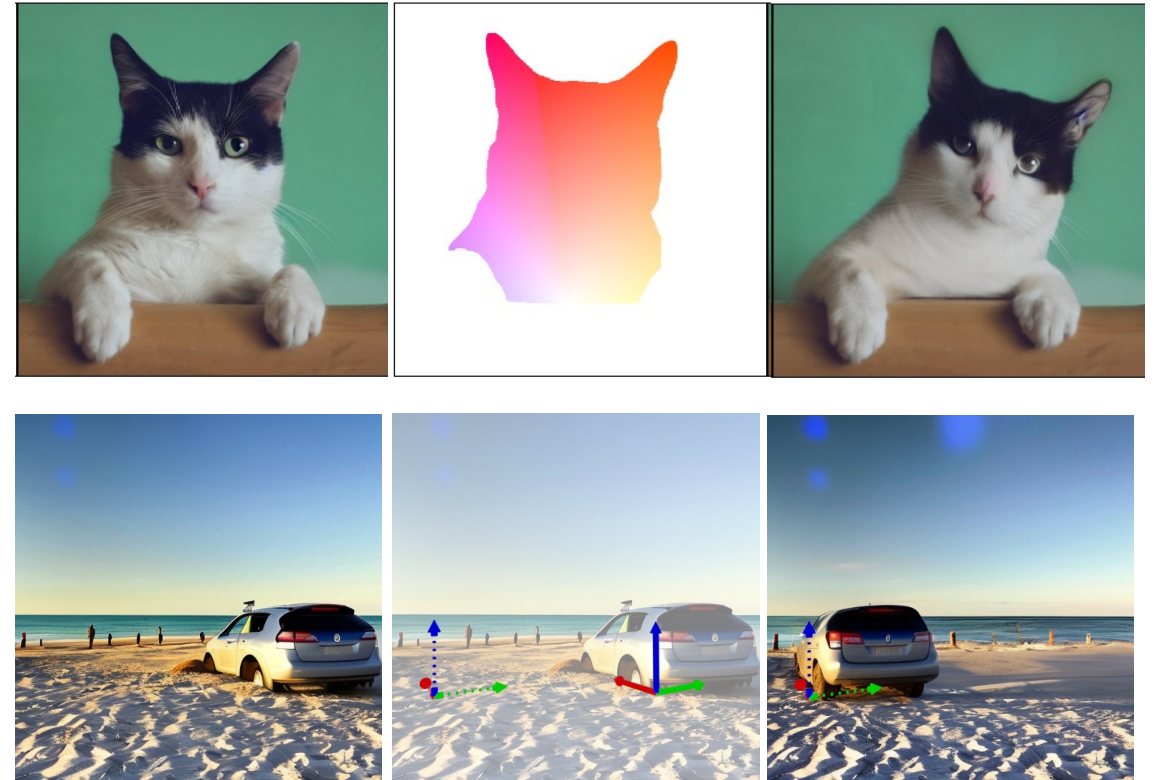
# Editing with Generative Models

## Personalization



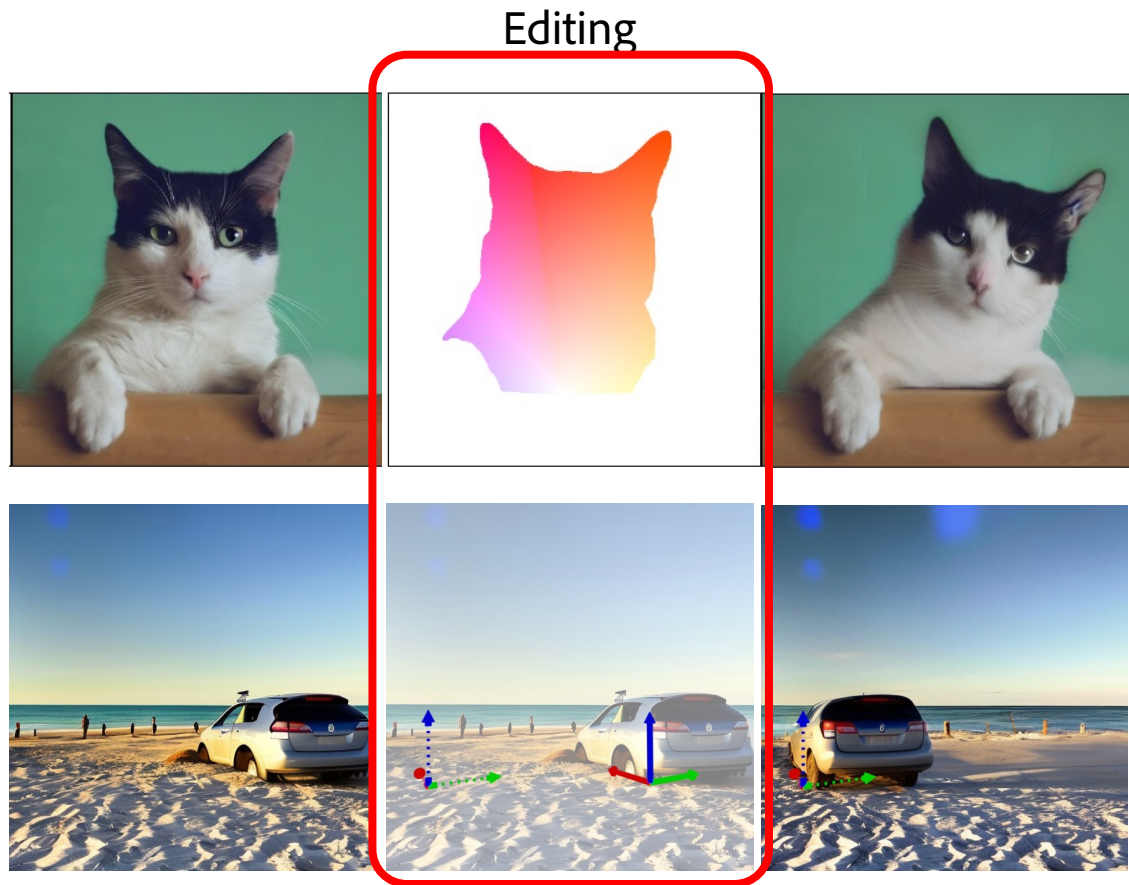
ConsiStory: Training-Free Consistent Text-to-Image Generation  
Tewel et al., ArXiv Feb. 2024

## Editing



Diffusion Handles: Enabling 3D Edits for Diffusion Models by Lifting Activations to 3D Pandey et al., CVPR 2024  
Motion Guidance: Diffusion-Based Image Editing with Differentiable Motion Estimators, Geng and Owens, ICLR 2024

# Editing with Generative Models



Same subject, same scene.  
Subject property changed by user **edit**.  
(Property such as position, pose, etc.)

Editing:

Generative Model  
+ Identity Preservation  
+ Edit Control

Diffusion Handles: Enabling 3D Edits for Diffusion Models by Lifting Activations to 3D Pandey et al., CVPR 2024

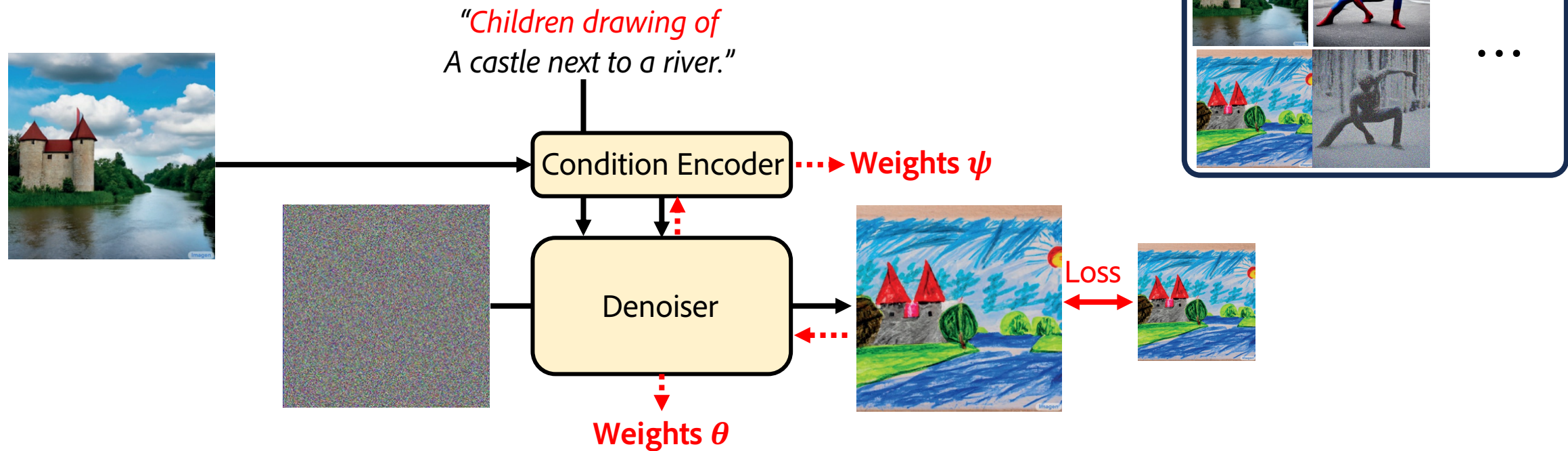
Motion Guidance: Diffusion-Based Image Editing with Differentiable Motion Estimators, Geng and Owens, ICLR 2024

# Image-to-Image Translation

Edit Control: Text Prompt

Identity Preservation: Condition Encoder

Large Paired Image Dataset

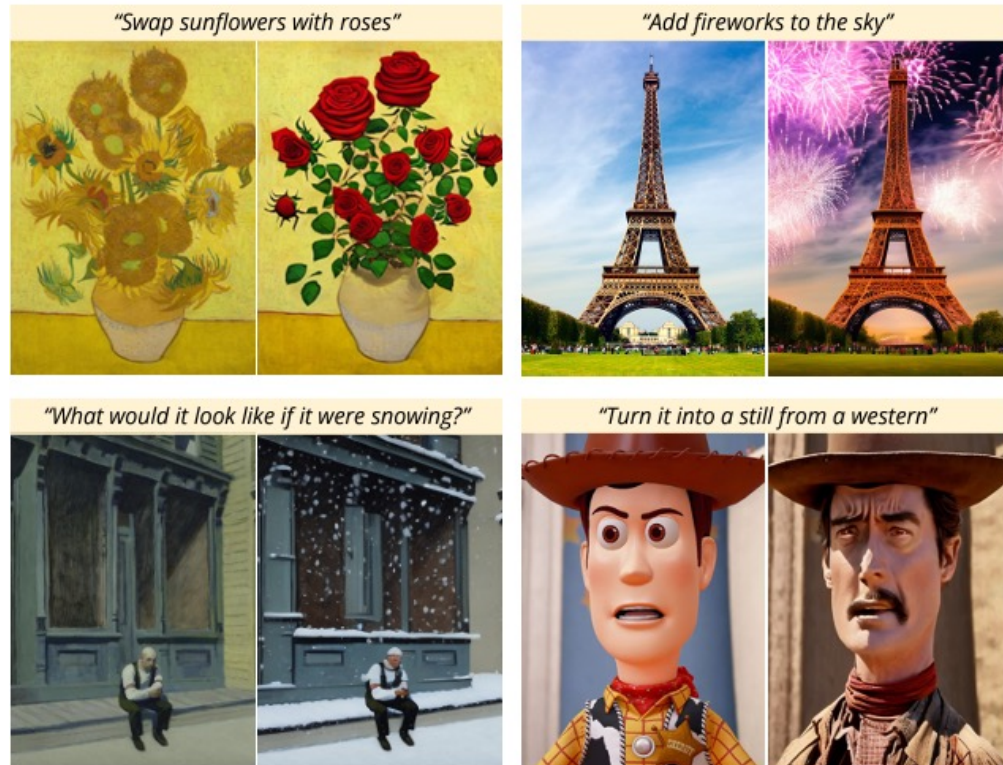


InstructPix2Pix: Learning to Follow Image Editing Instructions, Brooks et al, CVPR 2023

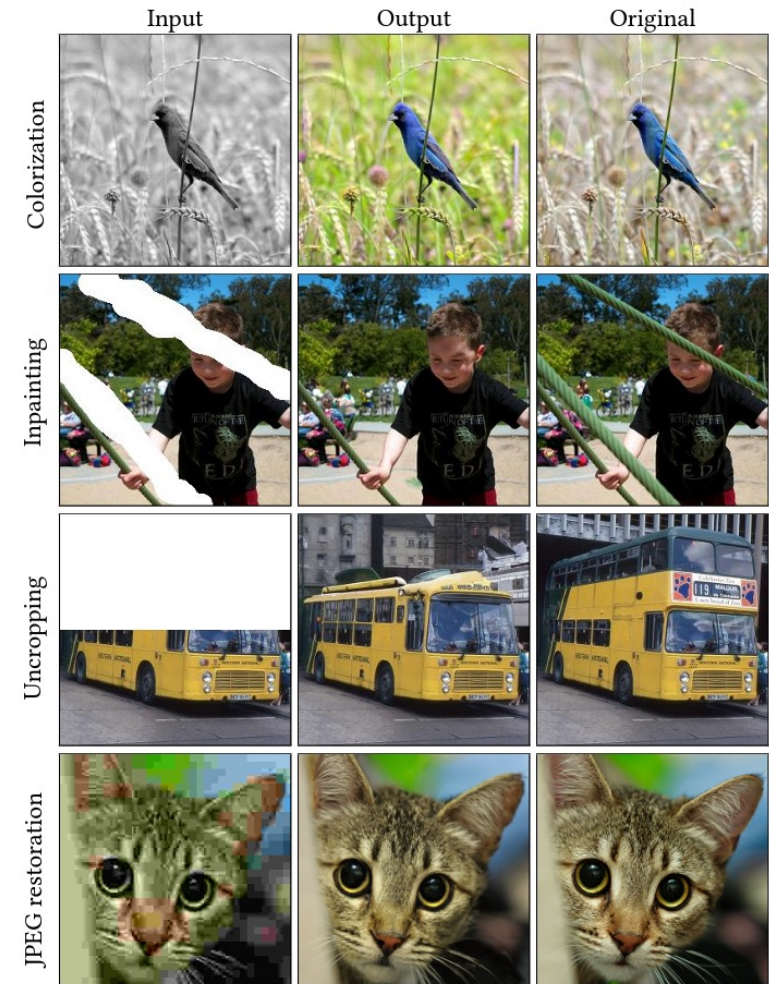
Palette: Image-to-Image Diffusion Models, Saharia et al., SIGGRAPH 2022

# Image-to-Image Translation

- Requires paired image dataset.
- Text prompt provides only coarse control.



InstructPix2Pix: Learning to Follow Image Editing Instructions,  
Brooks et al, CVPR 2023



Palette: Image-to-Image Diffusion Models,  
Saharia et al., SIGGRAPH 2022

## *InstructPix2Pix*

*Learning to Follow Image Editing Instructions*  
(Edit Instructions)



Code



<https://github.com/timothybrooks/instruct-pix2pix>



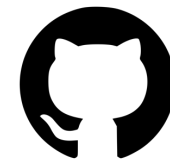
Demo



<https://huggingface.co/spaces/timbrooks/instruct-pix2pix>

## *Palette*

*Image-to-Image Diffusion Models*  
(Per-Task Fine-tuning)



Code



<https://github.com/Janspiry/Palette-Image-to-Image-Diffusion-Models>

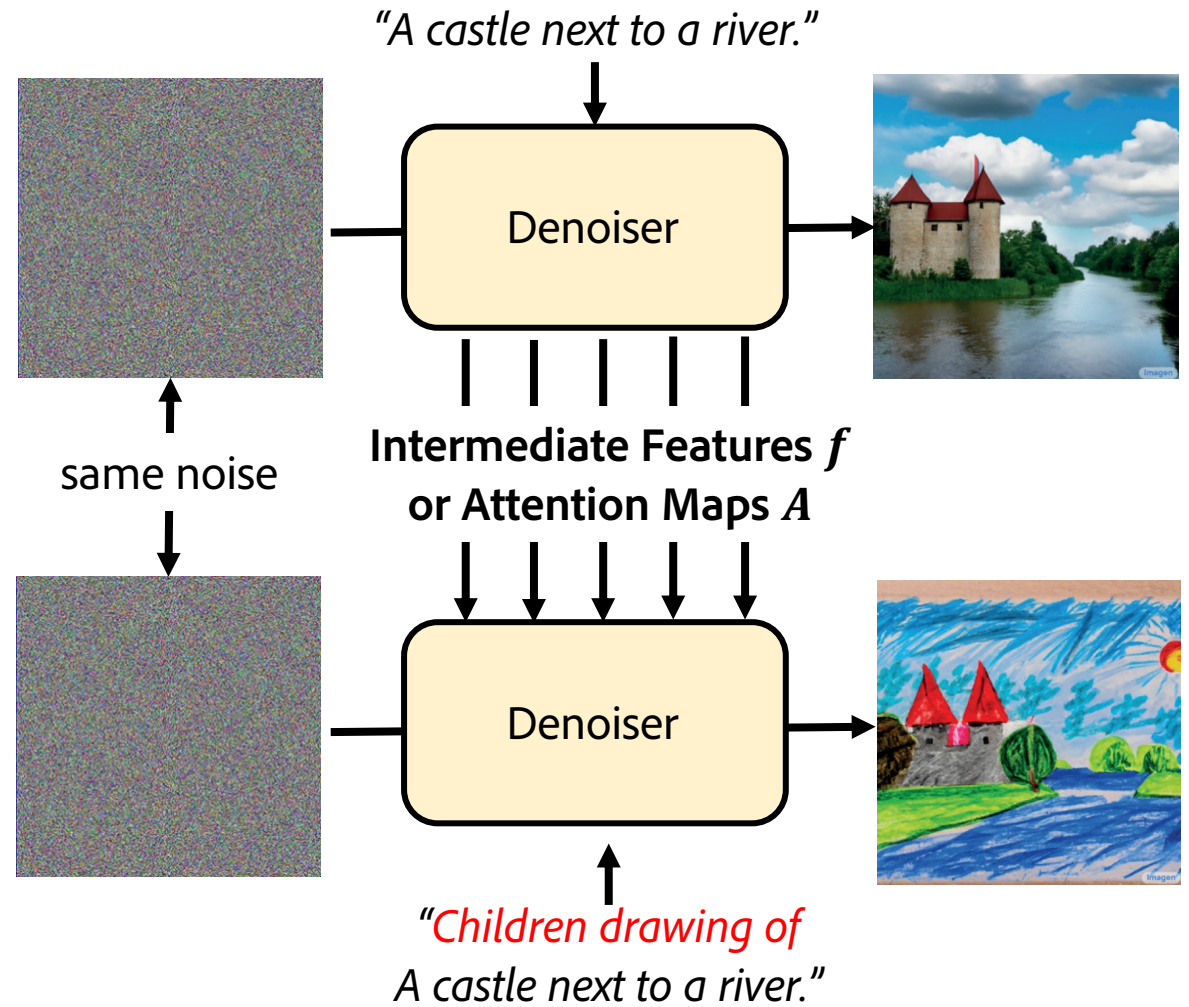
# Image-to-Image Translation

## Edit Control:

Text Prompt

## Identity Preservation:

Intermediate Features / Attention Maps



Plug-and-Play Diffusion Features for Text-Driven Image-to-Image Translation, Tumanyan et al., CVPR 2023

Prompt-to-Prompt Image Editing with Cross Attention Control, Hertz et al., ArXiv Aug. 2022

MasaCtrl: Tuning-Free Mutual Self-Attention Control for Consistent Image Synthesis and Editing, Cao et al., ArXiv Aug. 2022

# Image-to-Image Translation

- No training or fine-tuning needed.
- Cannot strongly change scene layout (object positions in the image remain roughly the same)
- Text prompt provides only coarse control.



Input Real Image



"a photo of a bronze horse in a museum"



"A photo of a pink horse on the beach"



"A photo of a robot horse"

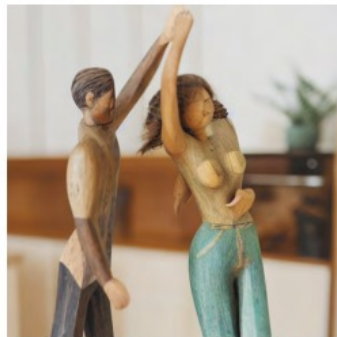


"a cake with decorations."

jelly beans



Input Real Image



"A wooden sculpture of a couple dancing"



"A cartoon of a couple dancing"



"a photo of robots dancing"



"Photo of a cat riding on a bicycle."

car



Plug-and-Play Diffusion Features for Text-Driven Image-to-Image Translation, Tumanyan et al., CVPR 2023

Prompt-to-Prompt Image Editing with Cross Attention Control, Hertz et al., ArXiv Aug. 2022

# Image-to-Image Translation

- No training or fine-tuning needed.
- Cannot strongly change scene layout (object positions in the image remain roughly the same)
- Text prompt provides only coarse control.



Input real image

“... jumping ...”



“A sitting boy” → “... standing ...”



Input real image

“...giving a thumbs up...”



“Elon Musk → ... side view ...”



“An apple” → “... two ...”



“A standing bird” → “... spreading wings ...”

MasaCtrl: Tuning-Free Mutual Self-Attention Control for Consistent Image Synthesis and Editing, Cao et al, ICCV 2023

## ***PnP-Diffusers***

*Plug-and-Play Diffusion Features for Text-Driven  
Image-to-Image Translation  
(Overwrite-Based Feature Injection)*



Demo



Code



<https://huggingface.co/spaces/hysts/PnP-diffusion-features>

## ***Prompt-to-Prompt***

*Image Editing with Cross Attention  
Control  
(Overwrite-Based Feature Injection)*



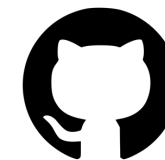
Code



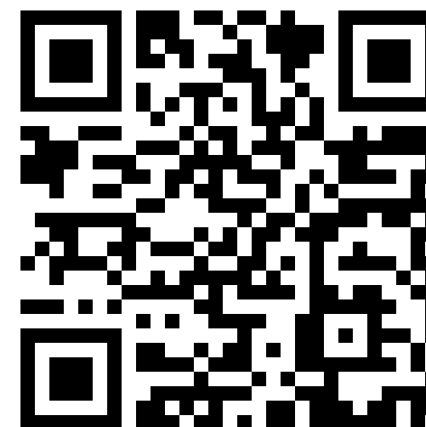
<https://github.com/google/prompt-to-prompt/>

## ***MasaCtrl***

*Tuning-Free Mutual Self-Attention Control  
for Consistent Image Synthesis and Editing  
(Cross-Attention-Based Feature Injection)*



Code

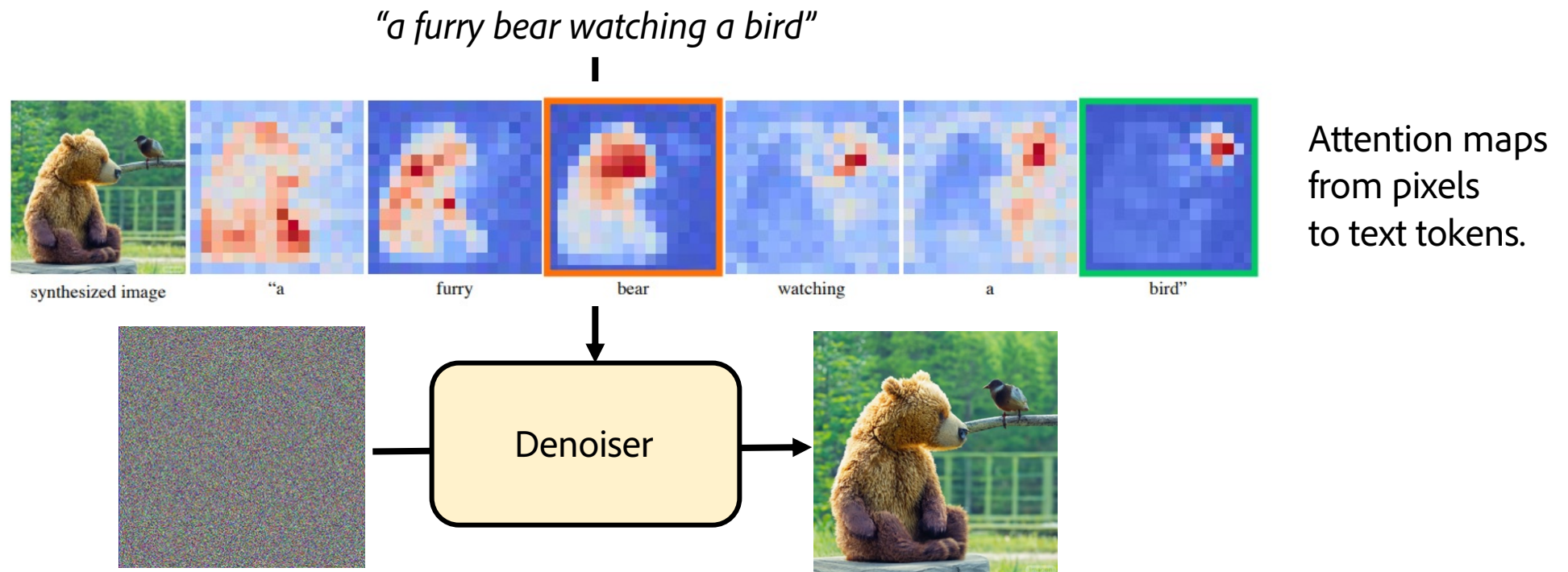


<https://github.com/TencentARC/MasaCtrl>

# Using Attention Maps & Intermediate Features

**Edit Control:** Transform Intermediate Features / Attention Maps

**Identity Preservation:** Intermediate Features / Attention Maps

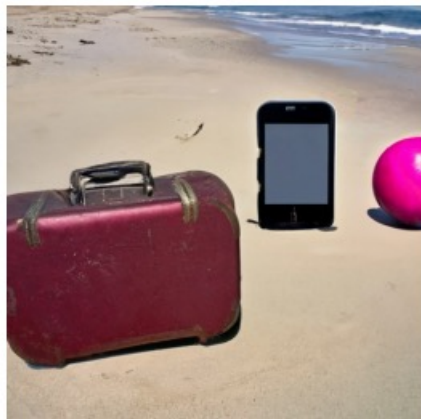


Prompt-to-Prompt Image Editing with Cross Attention Control, Hertz et al, ArXiv Aug, 2022

Diffusion Self-Guidance for Controllable Image Generation, Epstein et al, NeurIPS 2023

Diffusion Handles Enabling 3D Edits for Diffusion Models by Lifting Activations to 3D, CVPR 2024

# Using Attention Maps & Intermediate Features



(c) Swap w. fries

(d) Width ↓

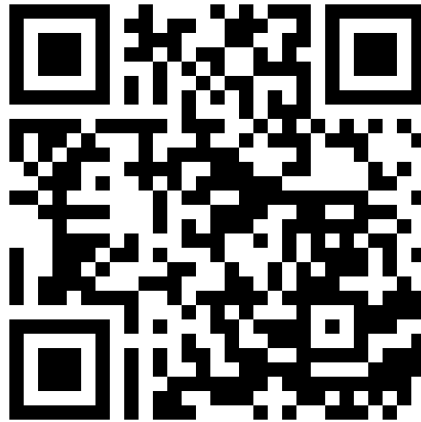
(e) Width ↓, height ↑

Diffusion Self-Guidance for Controllable Image Generation,  
Epstein et al., NeurIPS 2023

**Prompt-to-Prompt**  
*Image Editing with Cross Attention  
Control*  
(Overwrite-Based Feature Injection)



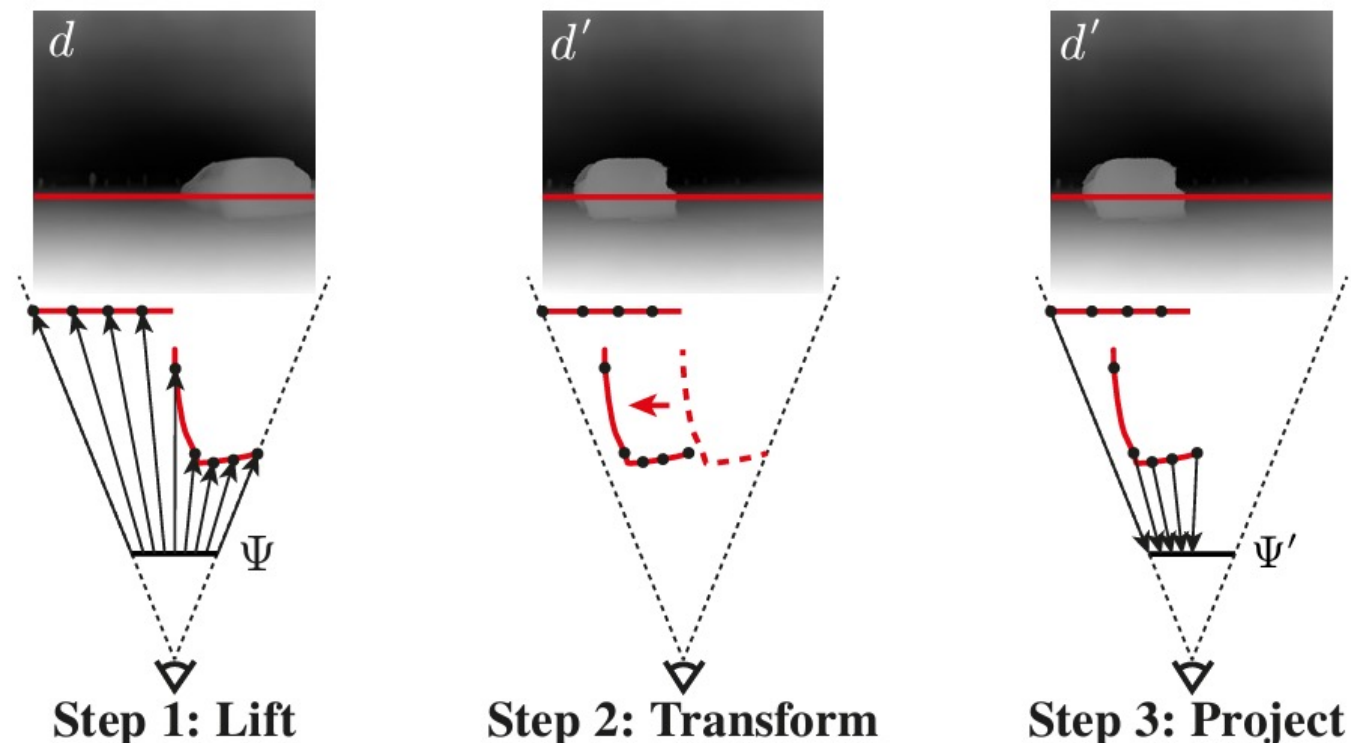
Code



<https://github.com/google/prompt-to-prompt/>

# Editing with Attention Maps and Intermediate Features

Attention maps / intermediate features can be 3D-transformed using monocular depth estimates.



Diffusion Handles Enabling 3D Edits for Diffusion Models by Lifting Activations to 3D, CVPR 2024

# Attention Maps & Intermediate Features

Attention maps / intermediate features can be 3D-transformed using monocular depth estimates.



Diffusion Handles Enabling 3D Edits for Diffusion Models by  
Lifting Activations to 3D, CVPR 2024

# Attention Maps & Intermediate Features

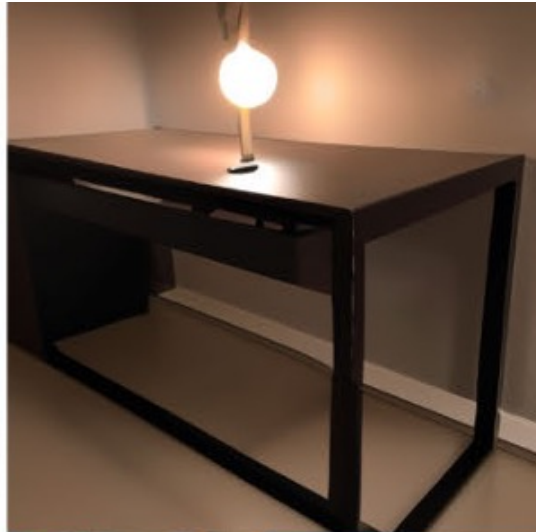
Attention maps / intermediate features can be 3D-transformed using monocular depth estimates.



Diffusion Handles Enabling 3D Edits for Diffusion Models by  
Lifting Activations to 3D, CVPR 2024

# Attention Maps & Intermediate Features

Attention maps / intermediate features can be 3D-transformed using monocular depth estimates.



Diffusion Handles Enabling 3D Edits for Diffusion Models by  
Lifting Activations to 3D, CVPR 2024

# Diffusion Models for Visual Computing

Chun-Hao Huang

Niloy Mitra, Daniel Cohen-Or, Minhyuk Sung, Duygu Ceylan, Paul Guerrero

## Part 5: Beyond Single Images



[https://geometry.cs.ucl.ac.uk/courses/diffusion4VC\\_eg24/](https://geometry.cs.ucl.ac.uk/courses/diffusion4VC_eg24/)

# Presentation Schedule

Introduction to Diffusion Models

Guidance and Conditioning Sampling

Attention

Break

Personalization and Editing

**Beyond Single Images**

Diffusion Models for 3D Generation

# So Far

- Theory and principles of diffusion models
- 2D image generation and editing
  - Editing one single image
  - Identity preservation

# This Section

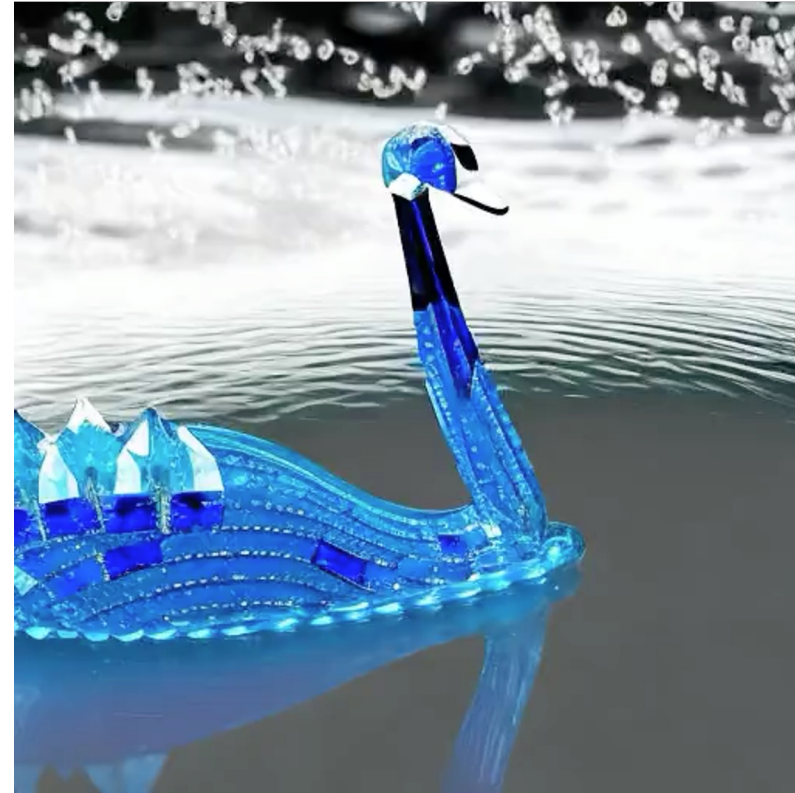
- Editing “multiple images”
  - in temporal dimension → videos
  - in 2D image domain → image montage
  - in 3D spatial domain → multi-views (*without* explicit 3D awareness)
- Training-free / zero-shot setup
  - Repurposing existing pretrained image diffusion model
- Training setup
  - Video diffusion model

# Individual Editing – Videos

input



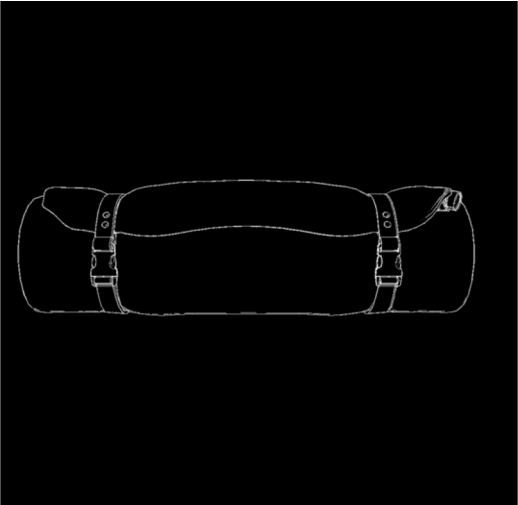
per-frame edit



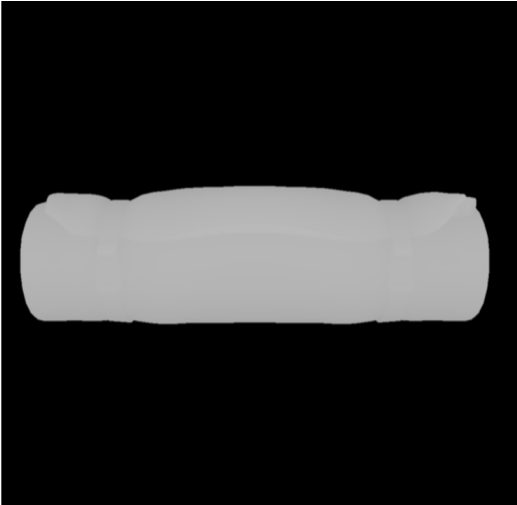
depth-conditioned SD

# Individual Editing – Multiple Views

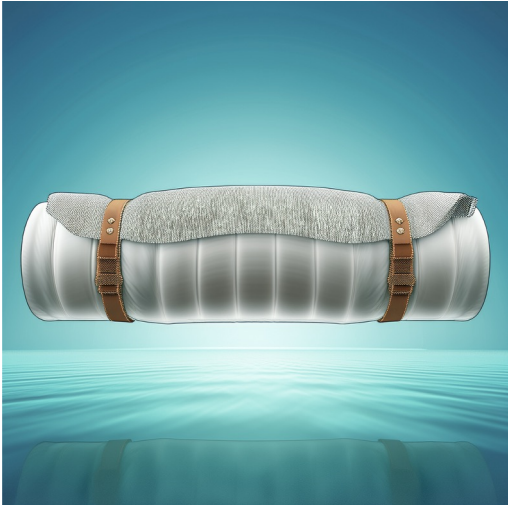
input



input



per-view edits



Firefly structure match

# Design Space for Consistency

- If *correspondences are known*, synchronizing:
  1. Initial noise and/or noisy latents
  2. Features
- Loss-guided denoising.
- Repurposing self attention for cross-view, cross-frame attention.

# Noise Model

- Denoting the initial noise of each frame as  $\epsilon^1, \epsilon^2, \dots, \epsilon^i, \dots$
- I.I.D.:  $\epsilon^1, \epsilon^2, \dots, \epsilon^i, \dots \sim \mathcal{N}(0, \mathbf{I})$
- Re-using the same noise:  $\epsilon^1 = \epsilon^2 = \epsilon^i = \dots \sim \mathcal{N}(0, \mathbf{I})$
- Mixed Noise Model [1]

$$\epsilon_{\text{shared}} \sim \mathcal{N}\left(0, \frac{\alpha^2}{1+\alpha^2} \mathbf{I}\right), \epsilon_{\text{ind}}^i \sim \mathcal{N}\left(0, \frac{\alpha^2}{1+\alpha^2} \mathbf{I}\right)$$
$$\epsilon^i = \epsilon_{\text{shared}} + \epsilon_{\text{ind}}^i$$

- Progressive Noise Model [1]

$$\epsilon^0 \sim \mathcal{N}(0, \mathbf{I}), \epsilon_{\text{ind}}^i \sim \mathcal{N}\left(0, \frac{1}{\sqrt{1+\alpha^2}} \mathbf{I}\right)$$
$$\epsilon^i = \frac{\alpha}{1+\alpha^2} \epsilon^{i-1} + \epsilon_{\text{ind}}^i$$

[1] Ge et al., Preserve Your Own Correlation: A Noise Prior for Video Diffusion Models, *ICCV'23*

# Correspondence-guided Noise



Kass and Pesare, Coherent noise for non-photorealistic rendering, *ToG'11*

# Correspondence-guided Noise

- $\epsilon^0 = \sqrt{\frac{1-\alpha}{1+\alpha}} \epsilon_{\text{ind}}^0$ , where  $\epsilon_{\text{ind}}^0 \sim \mathcal{N}(0, \mathbf{I})$

- $\epsilon^i(x, y) = \begin{cases} \sqrt{\frac{1-\alpha}{1+\alpha}} \epsilon_{\text{ind}}^i(x, y), & \epsilon_{\text{ind}}^i \sim \mathcal{N}(0, \mathbf{I}) \\ \alpha \epsilon^{i-1}(x', y') + (1 - \alpha) \epsilon_{\text{ind}}^i(x, y) \end{cases}$

disocclusion

if correspondence pairs  $(x', y')$  &  $(x, y)$  exist e.g., obtained by optical flow.

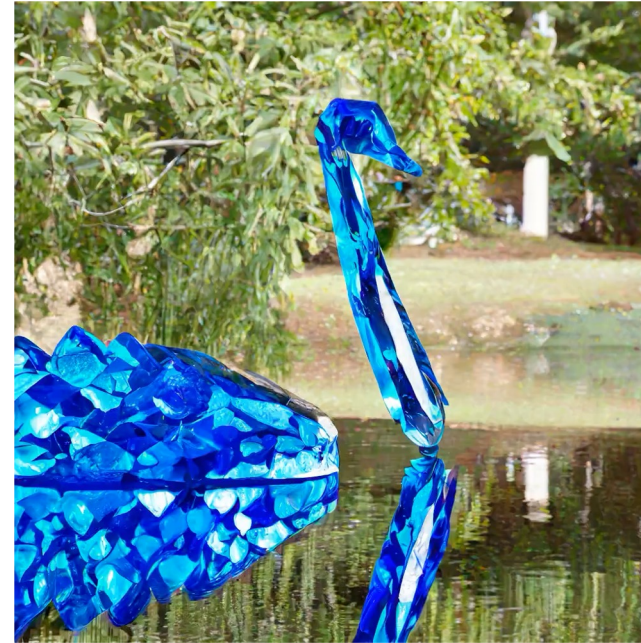
input



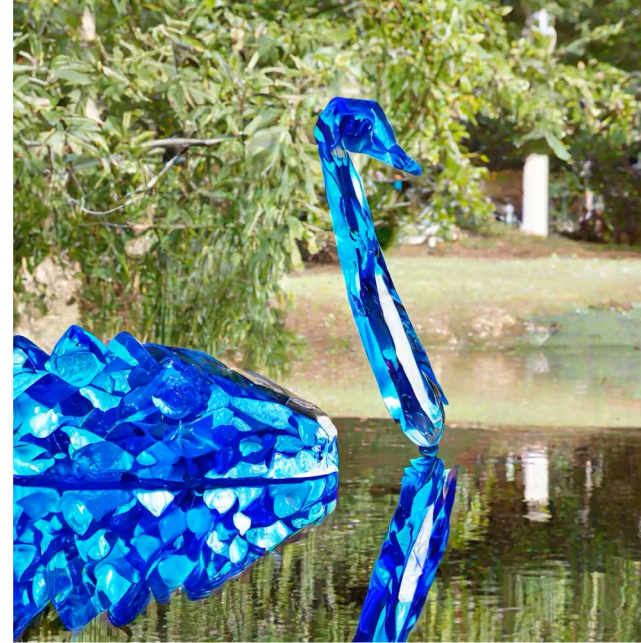
correspondence guided noise



results using same rand initialized noise

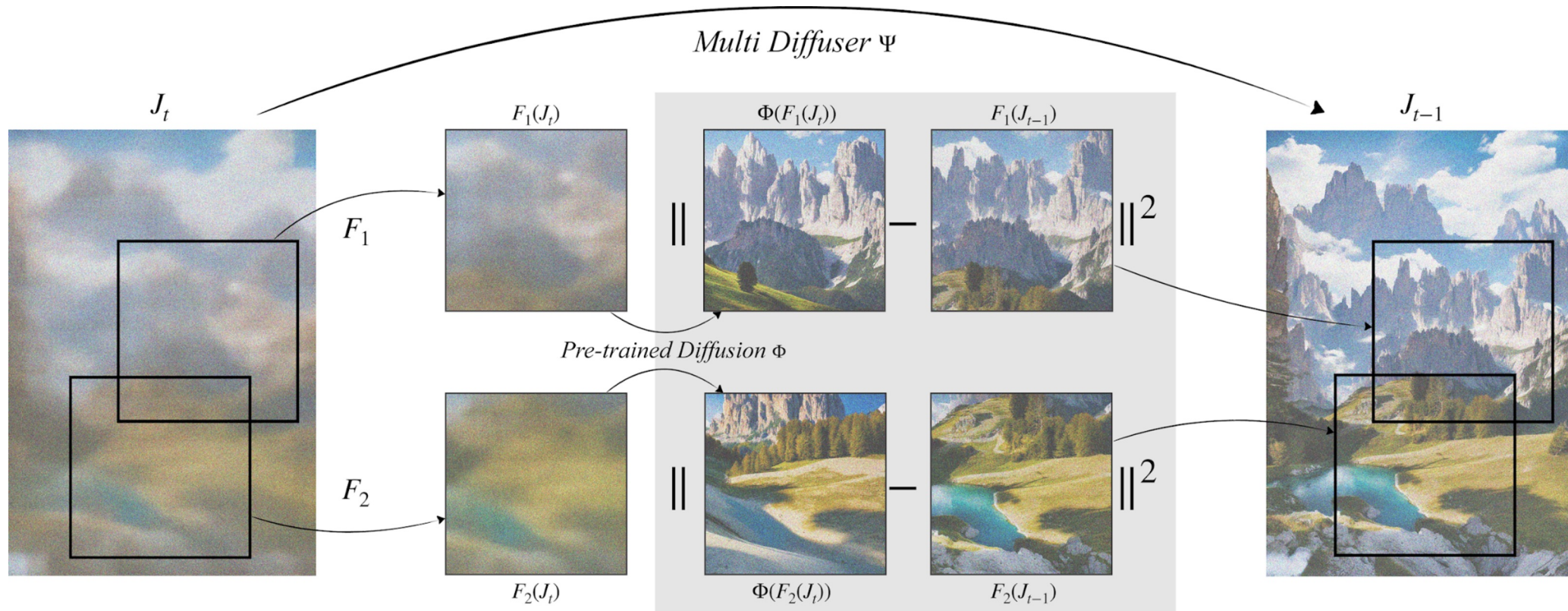


results using correspondence guided noise



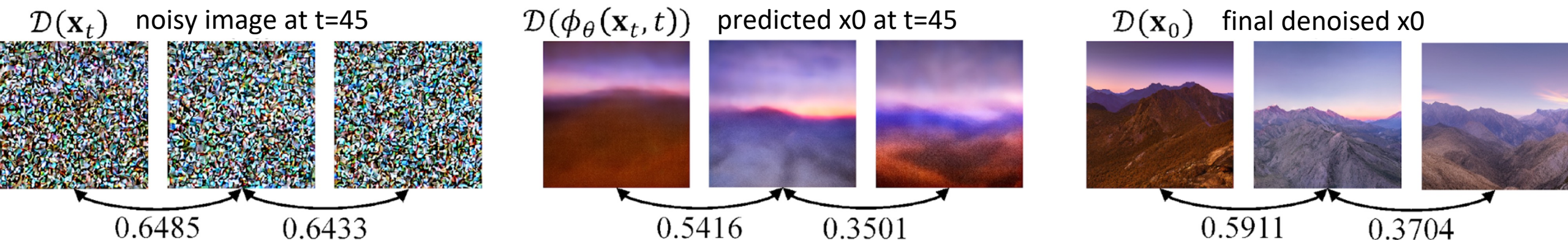
# Synchronizing Colors or Latent Features

- Averaging noisy latents of the same pixels



Bar-Tal et al., MultiDiffusion: Fusing Diffusion Paths for Controlled Image Generation, *ICML '23*

# Loss-guided Denoising



$$\hat{\mathbf{x}}_t^{(i)} = \mathbf{x}_t^{(i)} - w \nabla_{\mathbf{x}_t^{(i)}} \mathcal{L} \left( \mathcal{D}(\phi_\theta(\mathbf{x}_t^{(i)}, t)), \mathcal{D}(\phi_\theta(\mathbf{x}_t^{(0)}, t)) \right)$$

Loss  $\mathcal{L}$  : LPIPS score on the predicted denoised  $x_0$

Lee et al., SyncDiffusion: Coherent Montage via Synchronized Joint Diffusions, *NeurIPS'23*

without loss-guided denoising



with loss-guided denoising



## ***MultiDiffusion***

*Fusing Diffusion Paths for Controlled Image Generation  
(averaging noisy latents)*



Code



<https://github.com/omerbt/MultiDiffusion>

## ***SyncDiffusion***

*Coherent Montage via Synchronized Joint Diffusions  
(loss-guided denoising)*



Code

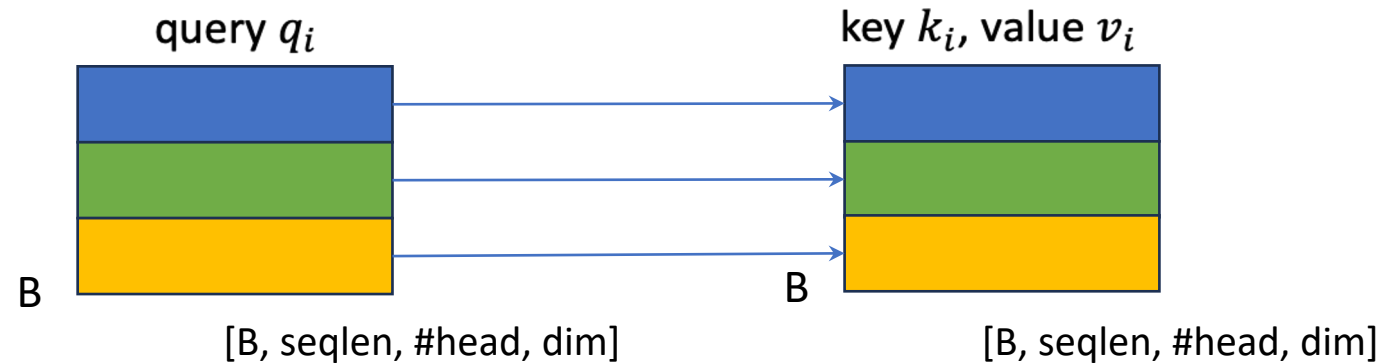


<https://github.com/KAI-ST-Visual-AI-Group/SyncDiffusion>

# Self-attention in UNet

$$\text{Att}(q_i, k_i, v_i) = \text{softmax}\left(\frac{q_i k_i^T}{\sqrt{d}}\right) v_i$$

- Independent generation: each sample in the batch attends to only itself.

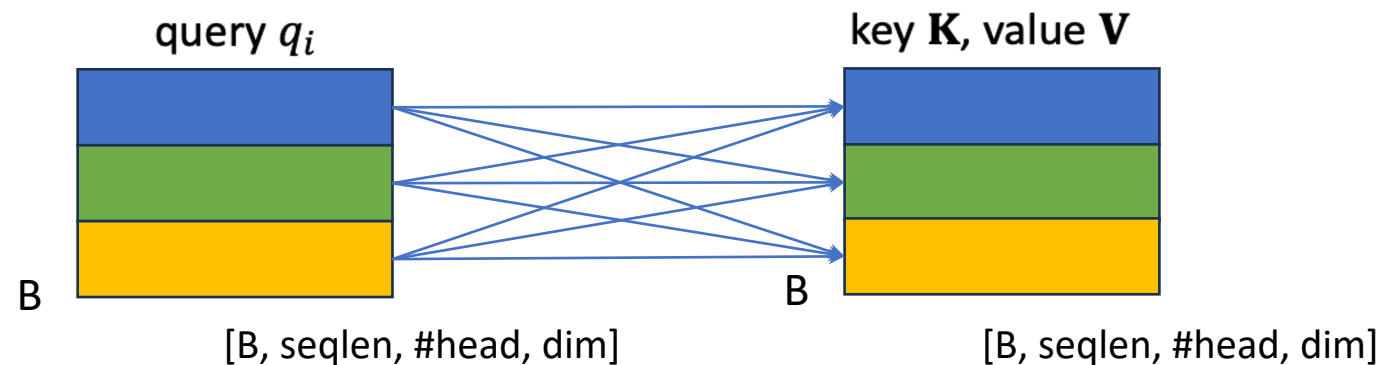


# Repurposing Self Attention

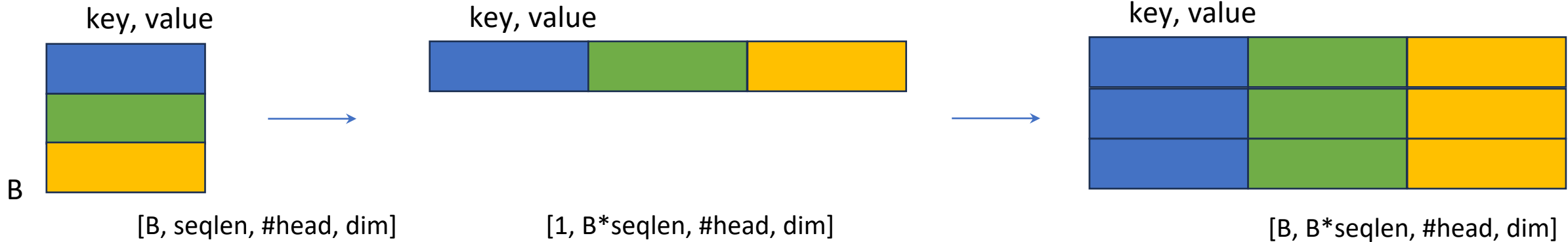
$$Att(q_i, \mathbf{K}, \mathbf{V}) = softmax\left(\frac{q_i \mathbf{K}^T}{\sqrt{d}}\right) \mathbf{V},$$

where  $\mathbf{K} = [k_1, \dots, k_i, \dots]$ ,  $\mathbf{V} = [v_1, \dots, v_i, \dots]$

- Batch generation: each sample attends to every samples.



# Pseudo Code



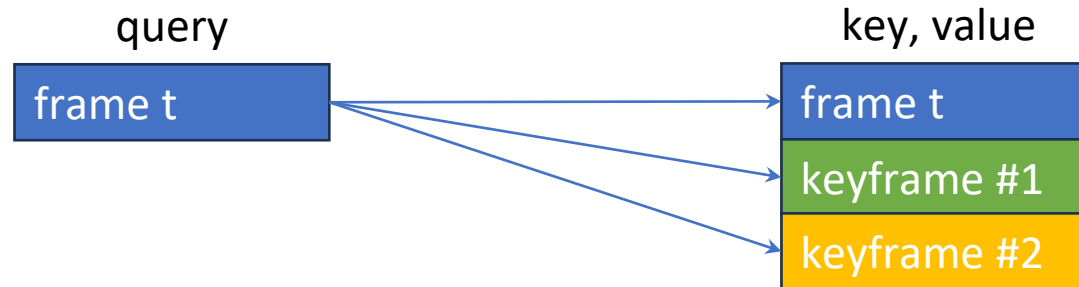
```
bs = k.shape[0]
```

```
k, v = (t.reshape(1, -1, attn.heads, attn.dim_head) for t in (k, v))
```

```
k = k.expand(bs, -1, -1, -1)
```

```
v = v.expand(bs, -1, -1, -1)
```

# Training-free/zero-shot Video Stylization



- Each frame attends to pre-defined “keyframes”
  - Pix2Video [1]: keyframes = [frame 1, frame t-1]
  - FateZero [2]: keyframes = [middle frame]
  - Text2Video-Zero [3]: keyframes = [frame 1]

[1] Ceylan et al., Pix2Video: Video Editing using Image Diffusion, *ICCV'23*

[2] Qi et al., FateZero: Fusing Attention for Zero-shot Text-based Video Editing, *ICCV'23*

[3] Khachatryan et al., Text2Video-Zero: Text-to-Image Diffusion Models are Zero-Shot Video Generators, *ICCV'23*

# Text-guided Video Stylization



a group of chocolate pigs looking for food



Ceylan et al., Pix2Video: Video Editing using Image Diffusion, *ICCV'23*

# Another Example

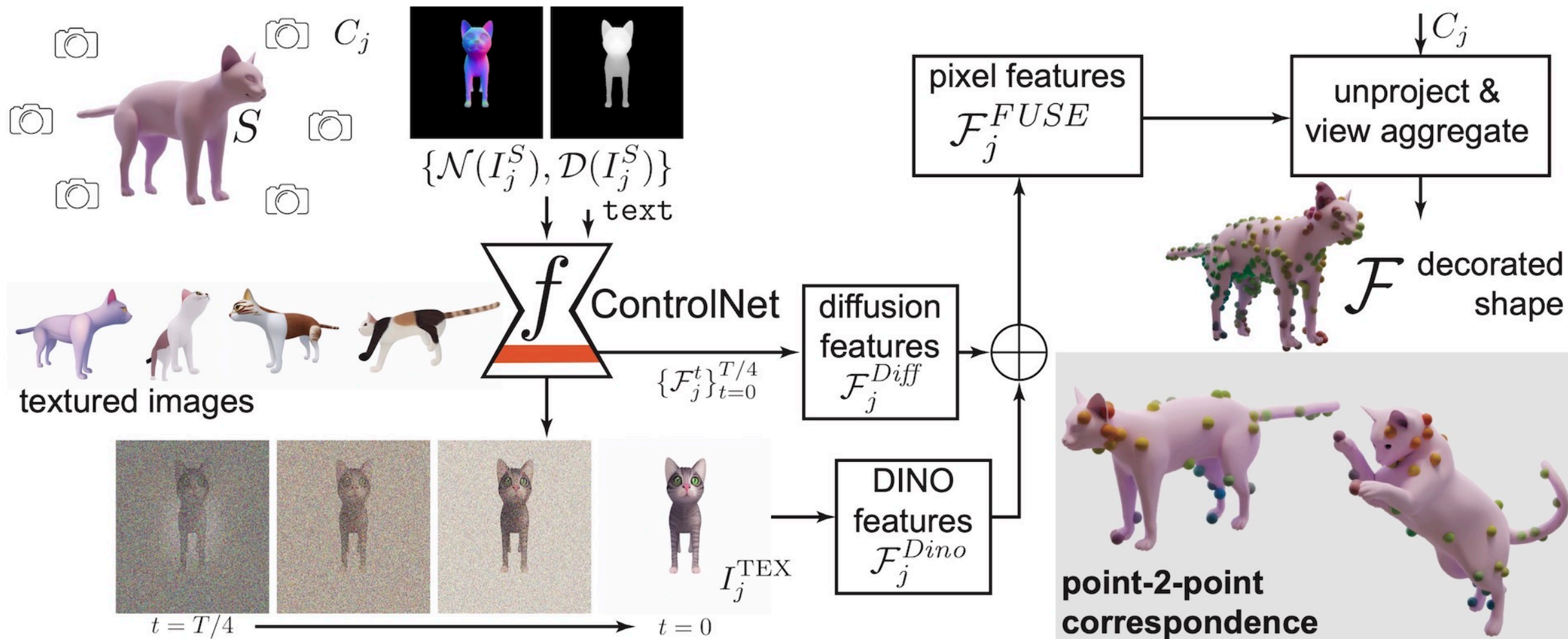


a Swarovski blue crystal swan on the lake



Ceylan et al., Pix2Video: Video Editing using Image Diffusion, *ICCV'23*

# Diffusion Features



Dutt et al., Diffusion 3D Features (Diff3F): Decorating Untextured Shapes with Distilled Semantic Features, *CVPR'24*

# Diffusion Features



Dutt et al., Diffusion 3D Features (Diff3F): Decorating Untextured Shapes with Distilled Semantic Features, *CVPR'24*

## ***Pix2Video***

*Video Editing using Image Diffusion*



Code



<https://github.com/duyguceylan/pix2video>

## ***FateZero***

*Fusing Attentions for Zero-shot Text-based Video Editing*



Code



<https://github.com/CheonyangQiQi/FateZero>

## ***Text2Video-Zero***

*Text-to-Image Diffusion Models are Zero-Shot Video Generators*



Code



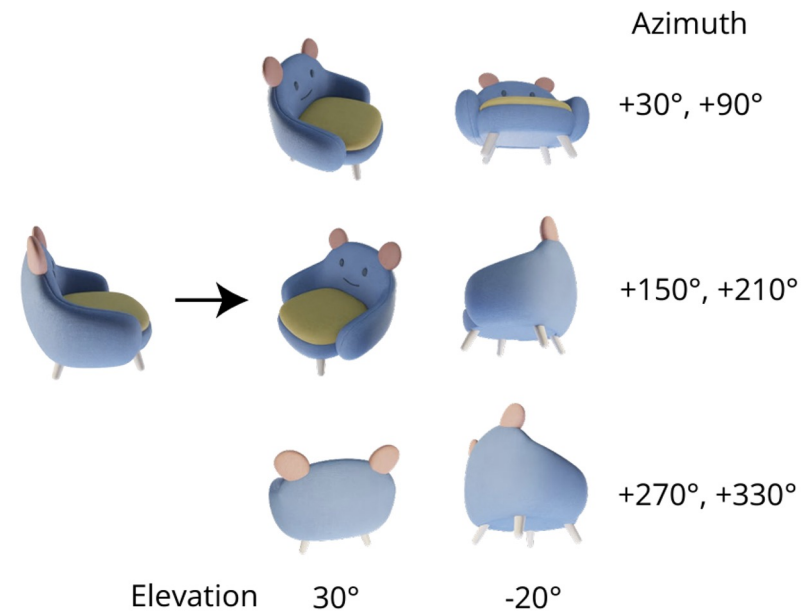
<https://github.com/Picsart-AI-Research/Text2Video-Zero>

# Applied on Multi-view Generation

- Stacking views into an **image grid**  
a simple yet effective way of repurposing self-attention as cross-view attention.



Tsalicoglou et al., Textmesh: Generation of realistic 3D meshes from text prompts, arXiv'23



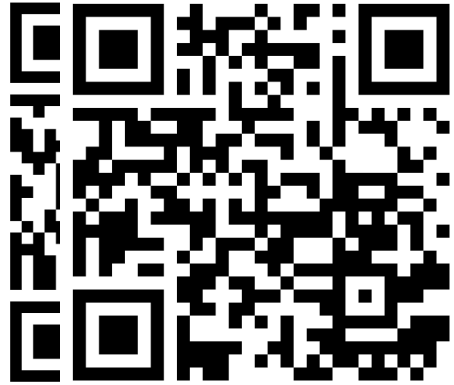
Shi et al., Zero123++: a Single Image to Consistent Multi-view Diffusion Base Model, arXiv'23

# Zero123++

*A Single Image to Consistent Multi-view Diffusion Base Model  
(stacking 6 predefined views)*



Code



<https://github.com/SUDO-AI-3D/zero123plus>



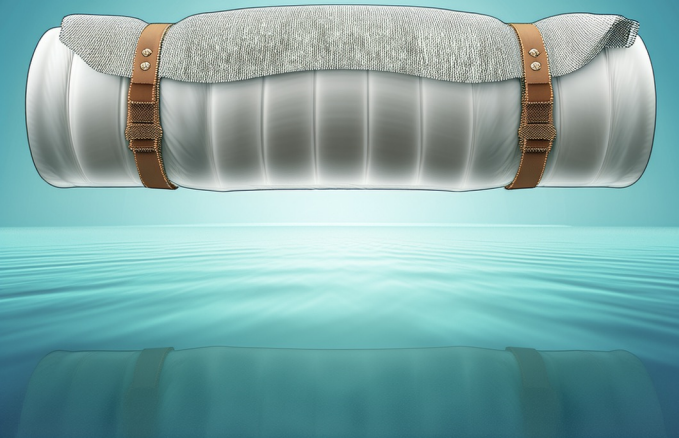
Demo



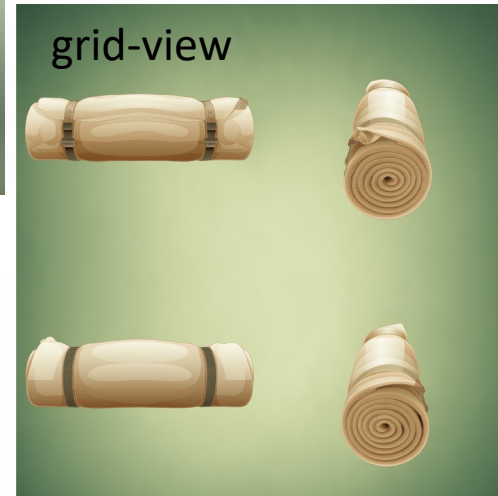
<https://github.com/SUDO-AI-3D/zero123plus>

independent

# Multi-view

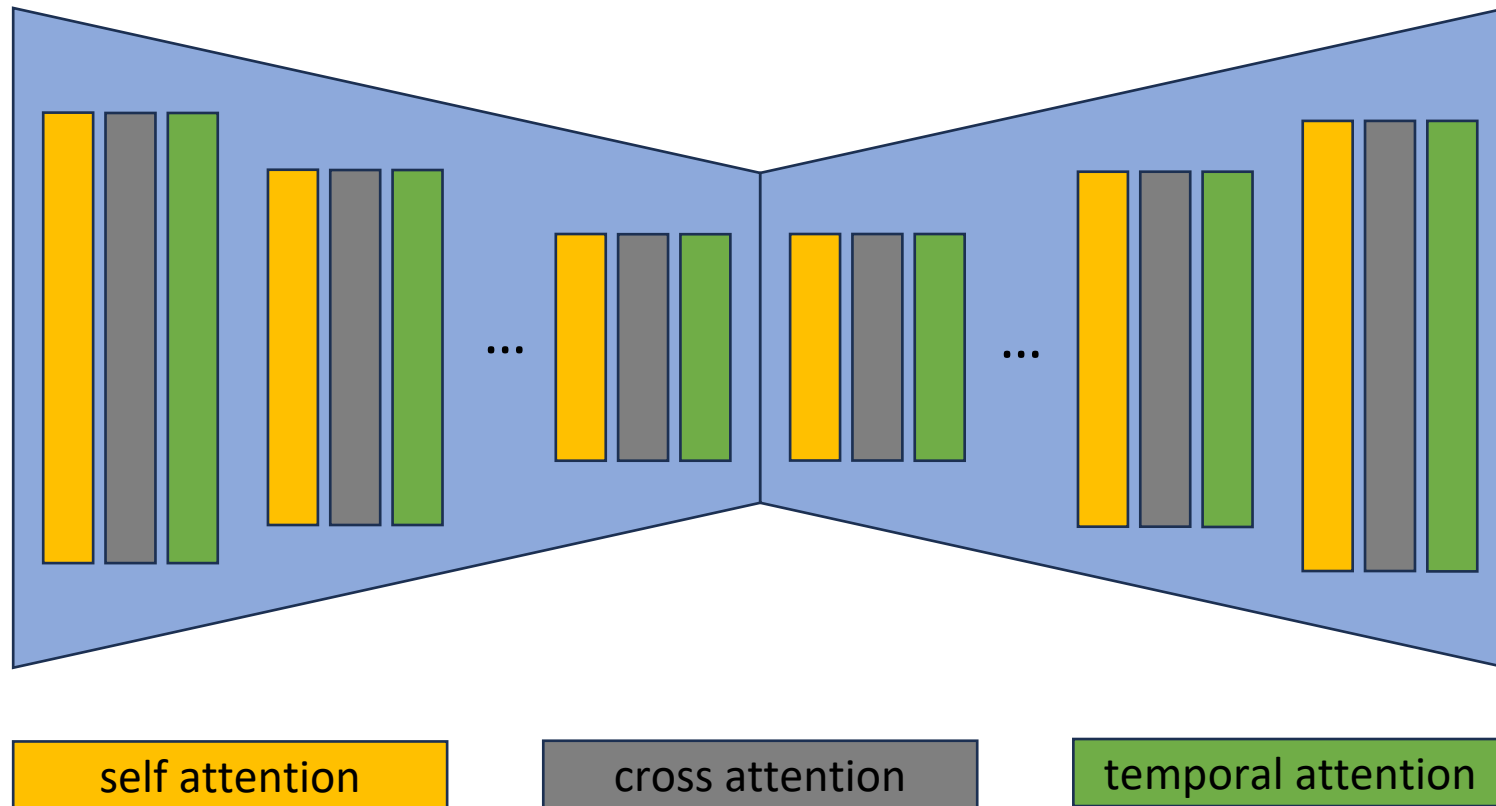


batch

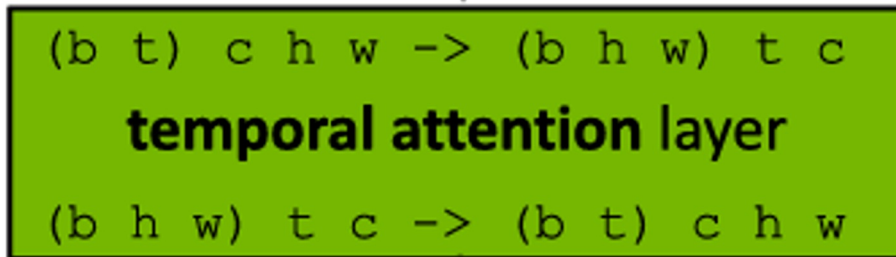


# Video Diffusion Model

- Adding a **temporal attention** layer after spatial cross attention



# Inflated Temporal Attention Layer



Blattmann et al., Align your Latents: High-Resolution Video Synthesis with Latent Diffusion Models, *CVPR'23*

Blattmann et al., Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets, *arXiv'23*

Guo et al., AnimateDiff: Animate Your Personalized Text-to-Image Diffusion Models without Specific Tuning, *ICLR'24*



Code



[https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/transformers/transformer\\_temporal.py](https://github.com/huggingface/diffusers/blob/main/src/diffusers/models/transformers/transformer_temporal.py)

# Applied on Multi-view Generation

Generating 360 turn-table videos of 3D objects

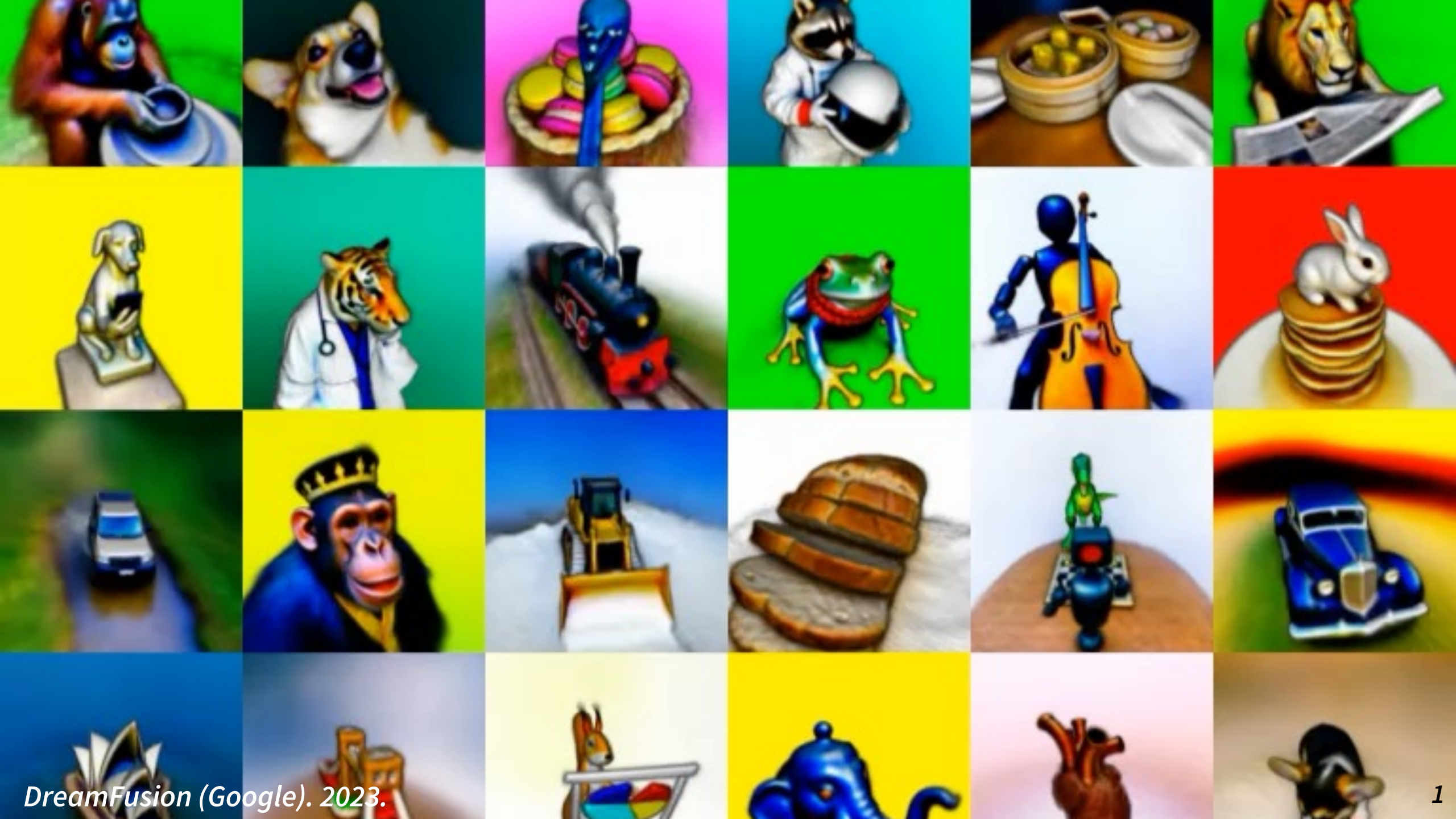


Blattmann et al., Stable Video Diffusion: Scaling Latent Video Diffusion Models to Large Datasets, arXiv'23

*Diffusion Models for  
Visual Content Generation  
3D Generation*

*Presenter: Minhyuk Sung*

*Eurographics 2024 Tutorial*



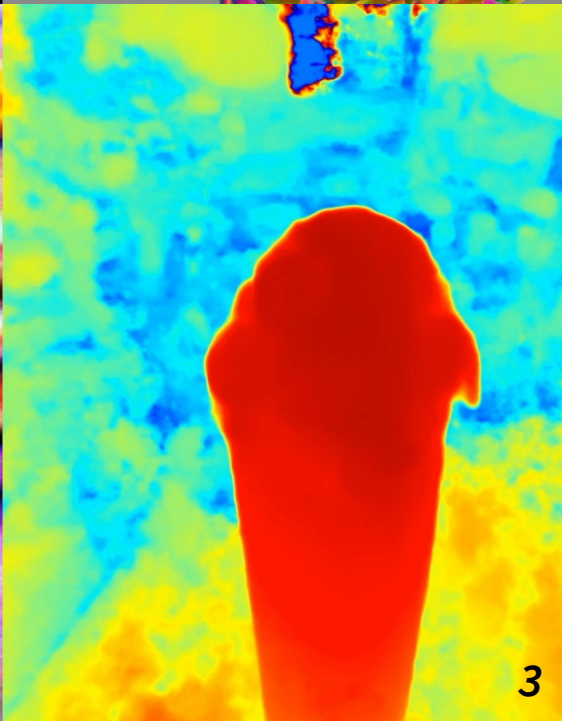
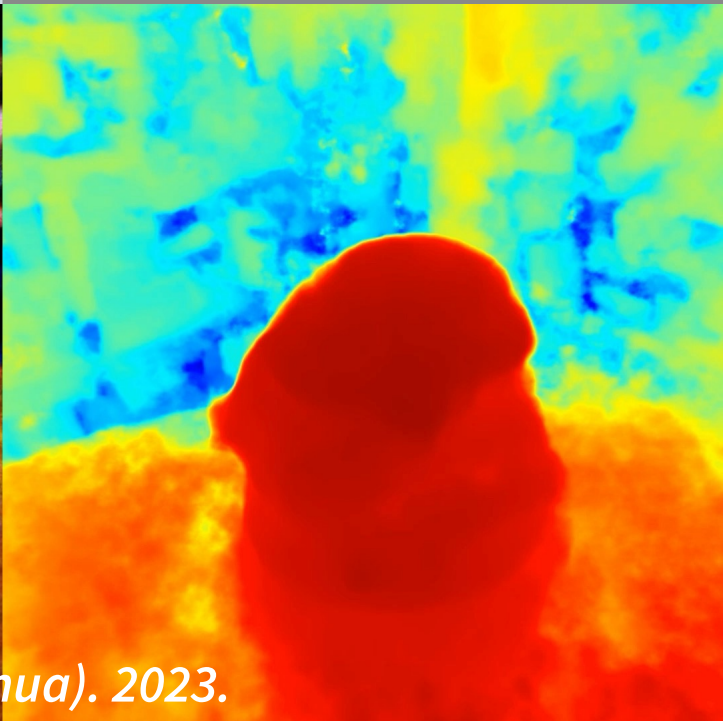
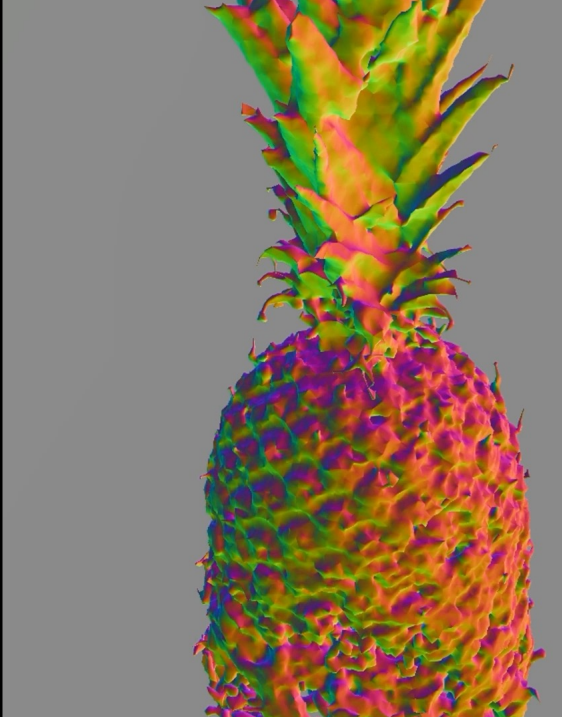
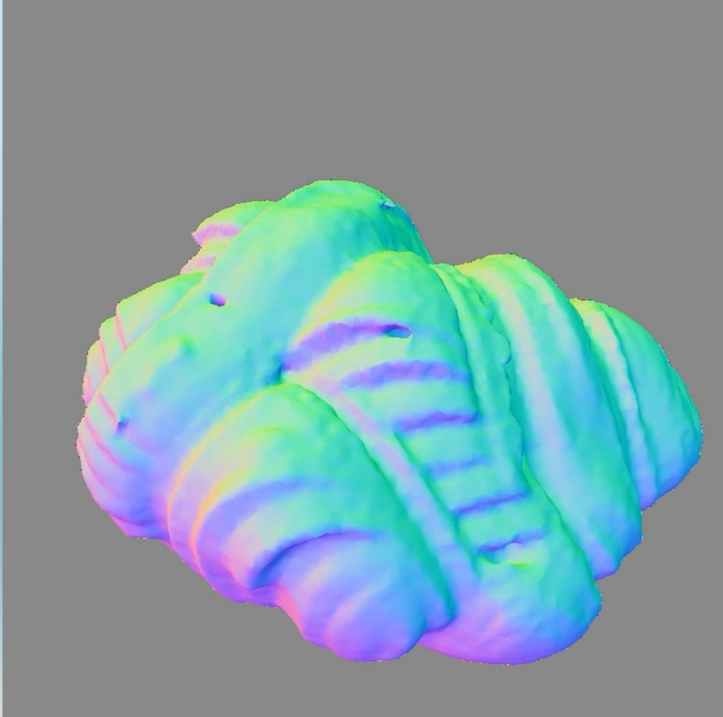


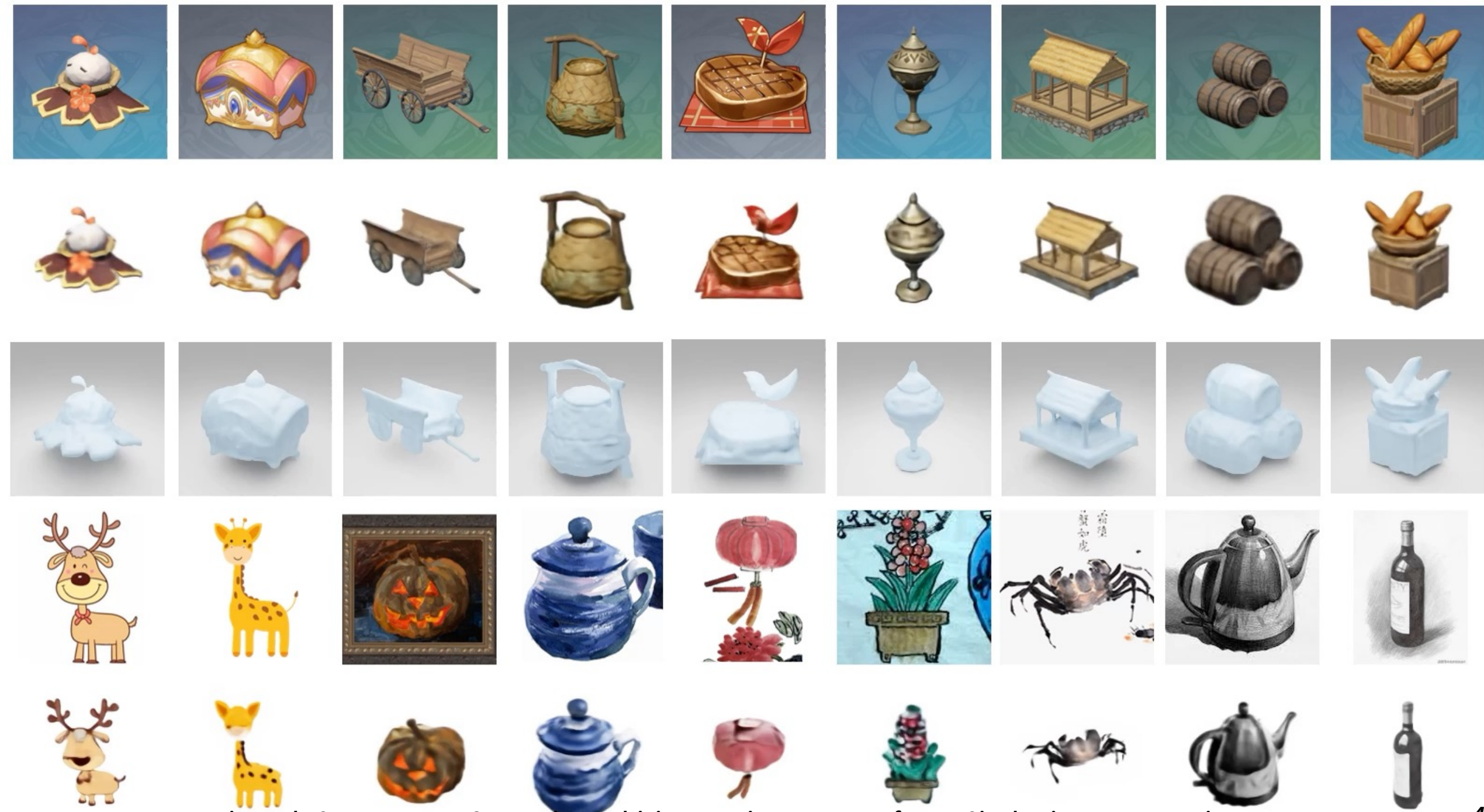
# High-Resolution Text-to-3D Content Creation

Chen-Hsuan Lin\* Jun Gao\* Luming Tang\* Towaki Takikawa\* Xiaohui Zeng\*  
Xun Huang Karsten Kreis Sanja Fidler# Ming-Yu Liu# Tsung-Yi Lin

\*# : equal contributions

**NVIDIA Corporation**



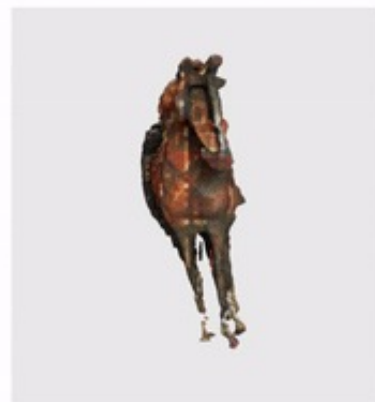
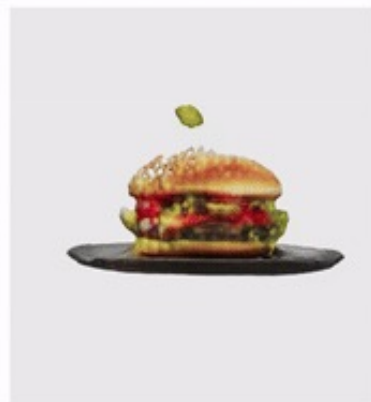
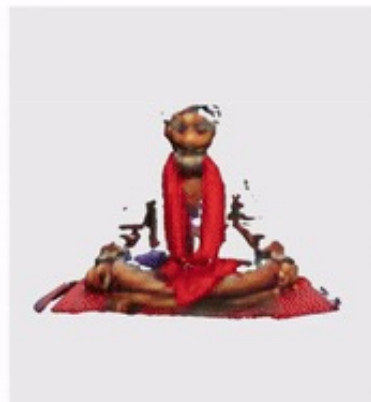


Liu et al., SyncDreamer: Generating Multiview-consistent Images from a Single-view Image, arXiv 2023.

Zero123-XL



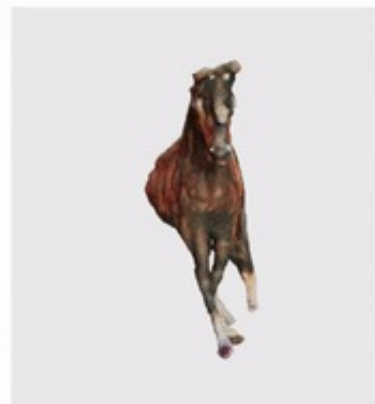
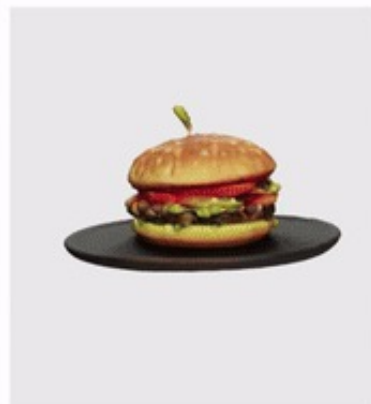
OpenLRM



Stable Zero123



TripoSr (ours)



# *3D Diffusion Models*

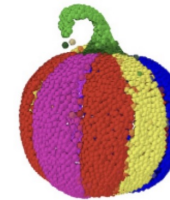
# 3D Diffusion Models

## Point clouds

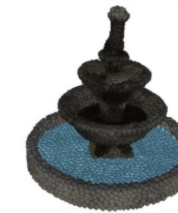
- *ShapeGF (Cai et al., 2020)*
- *DPM (Luo and Hu, 2021)*
- *LION (Nichol et al., 2022)*



"a corgi wearing a red santa hat"



"a multicolored rainbow pumpkin"



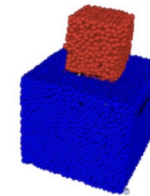
"an elaborate fountain"



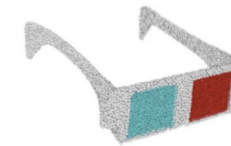
"a traffic cone"



"a vase of purple flowers"



"a small red cube is sitting on top of a large blue cube. red on top, blue on bottom"



"a pair of 3d glasses, left lens is red right is blue"



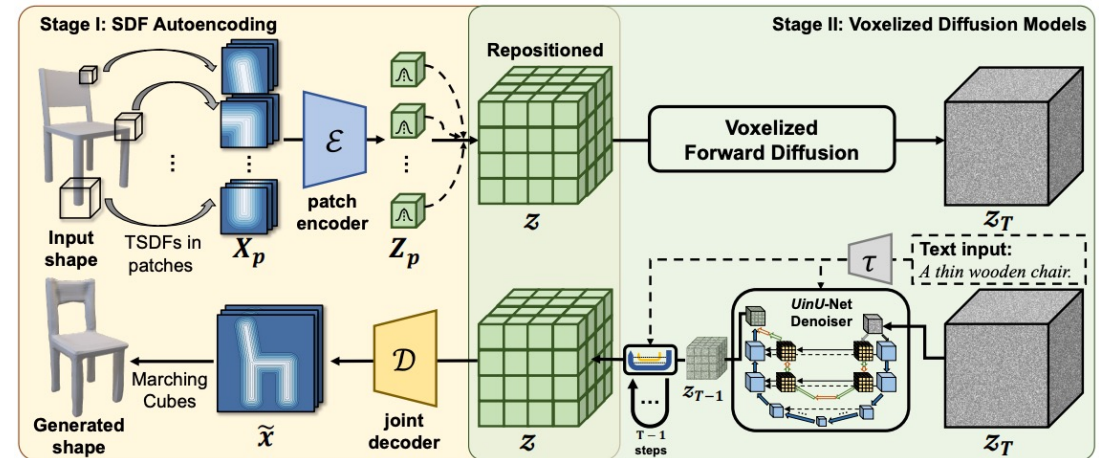
"an avocado chair, a chair imitating an avocado"

*LION, Nichol et al. 2022.*

# 3D Diffusion Models

## Voxel representation

- *PVD (Zhou et al., 2022)*
- *Diffusion-SDF (Li et al., 2022)*

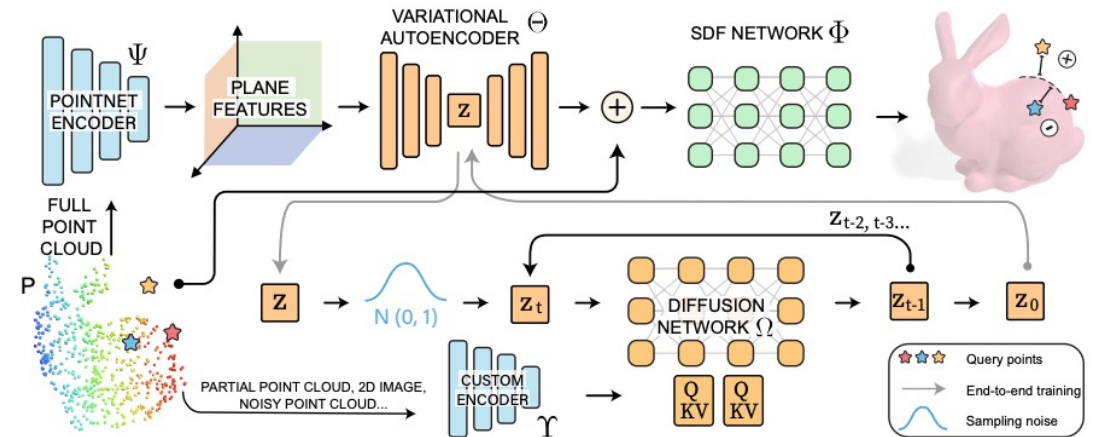


*Diffusion-SDF, Li et al., 2022.*

# 3D Diffusion Models

## Latent representation

- *SDFusion*  
(Cheng et al., 2022)
- *DiffusionSDF* (w/o hyphen,  
Chou et al., 2023)



DiffusionSDF, Chou et al., 2023

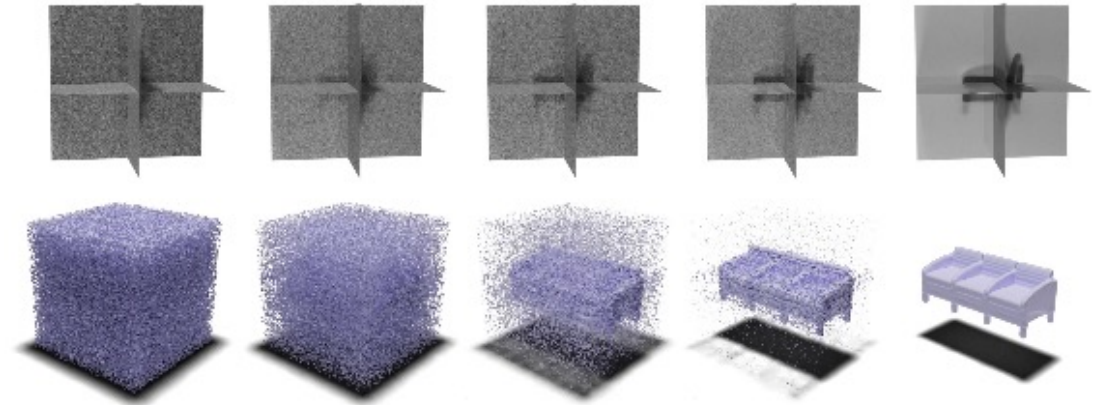
# 3D Diffusion Models

*Triplane representation*

- *NFD (Shue et al., 2022)*

*Diffusion in the spectral domain:*

- *NeuralWavelet (Hui et al., 2022)*



Shape Generation



Shape Interpolation

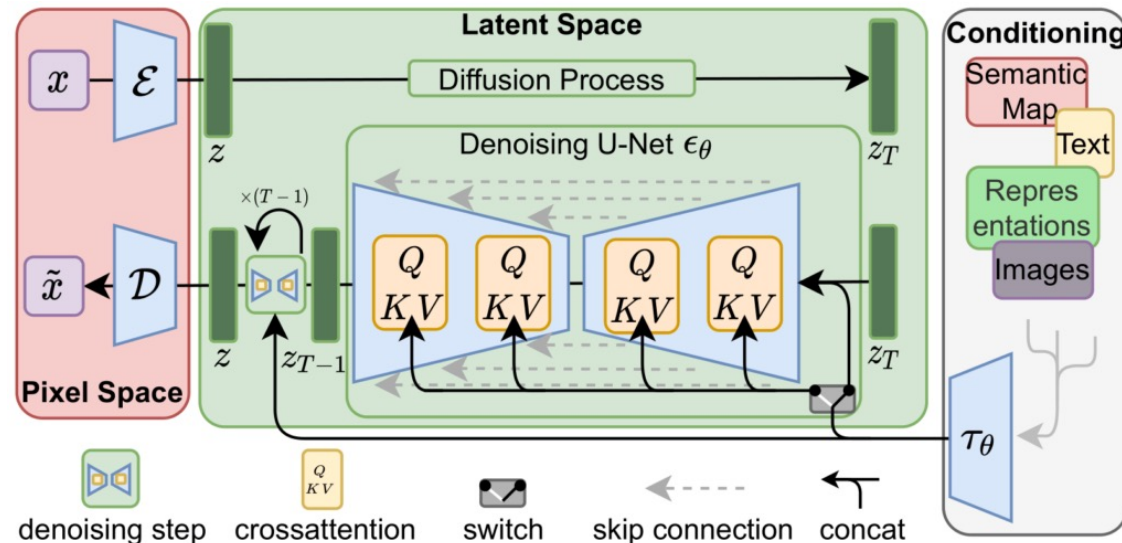
NFD, Shue et al., 2022.

# *Diffusion w/ Different Representations*

- *Implicit representation*
- *Explicit representation*

# Diffusion w/ Different Representations

- **Implicit** representation (i.e., **latent** features)
  - E.g., Latent diffusion (Rombach et al., 2022)
  - (+) **Best quality of the generated data.**
  - (-) **Requires retraining** for each conditional generation setup.



Latent Diffusion, Rombach et al., 2022.

# Diffusion w/ Different Representations

- **Explicit** representation (E.g., **pixels** in images)
  - E.g.: The original DDPM model (Ho et al., 2020) for images.
  - (–) **Suboptimal performance** due to the high dimensionality.
  - (–) Cannot change the **resolution** of the data.
  - (+) Can be directly leveraged in **conditional generation** setups in a **zero-shot** manner.

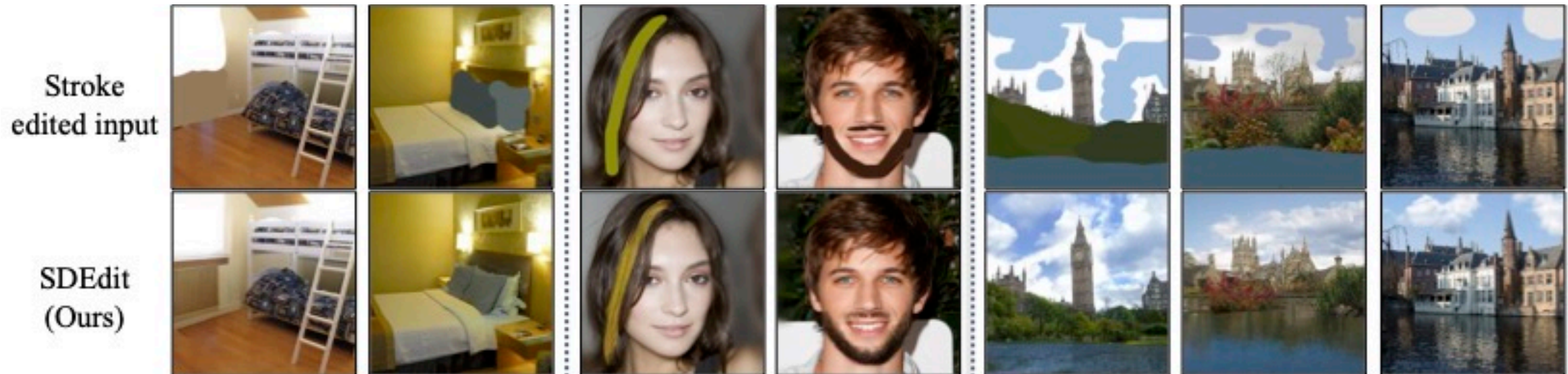
# SDEdit [Meng et al., 2022]

Image **editing** using a pretrained pixel-space diffusion model.

$$\mathbf{x}^{(0)} = m \odot \mathbf{x}_a^{(0)} + (1 - m) \odot \mathbf{x}_b^{(0)}$$

$$\mathbf{x}^{(t)} = \text{denoise}(\mathbf{x}^{(0)}, t)$$

$$\mathbf{x}'^{(0)} = \text{add\_noise}(\mathbf{x}^{(t)}, t)$$



# RePaint [Lugmayr et al., 2022]

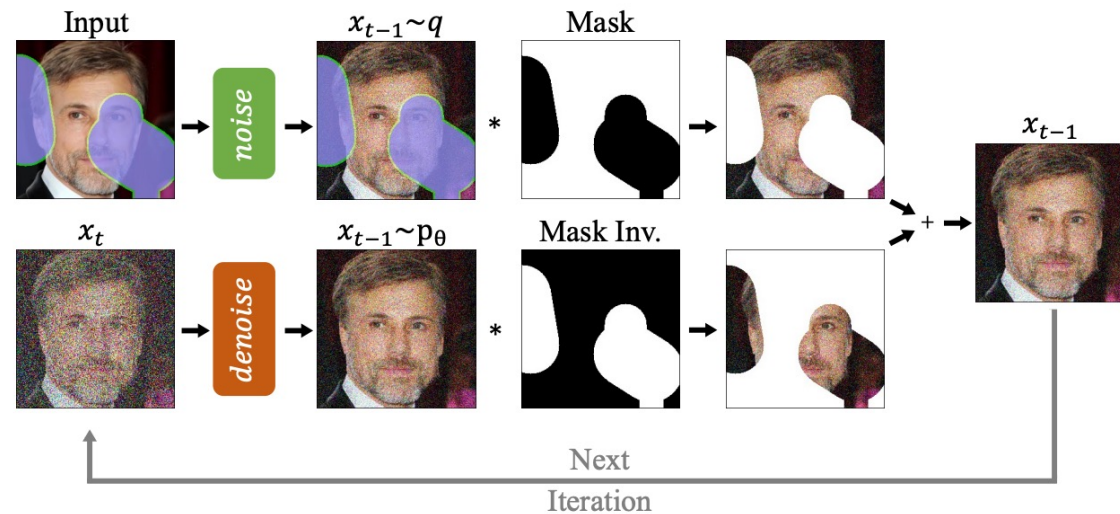
Image *inpainting* using a pretrained pixel-space diffusion model.

$$\mathbf{x}_f^{(t-1)} = \text{denoise}(\mathbf{x}^{(t)}, 1)$$

$$\mathbf{x}_b^{(t-1)} = \text{add\_noise}(\mathbf{x}^{(0)}, t)$$

$$\mathbf{x}^{(t-1)} = m \odot \mathbf{x}_f^{(t-1)} + (1 - m) \odot \mathbf{x}_b^{(t-1)}$$

Repeat for  $t = T, \dots, 1$ .



# Hybrid Representation

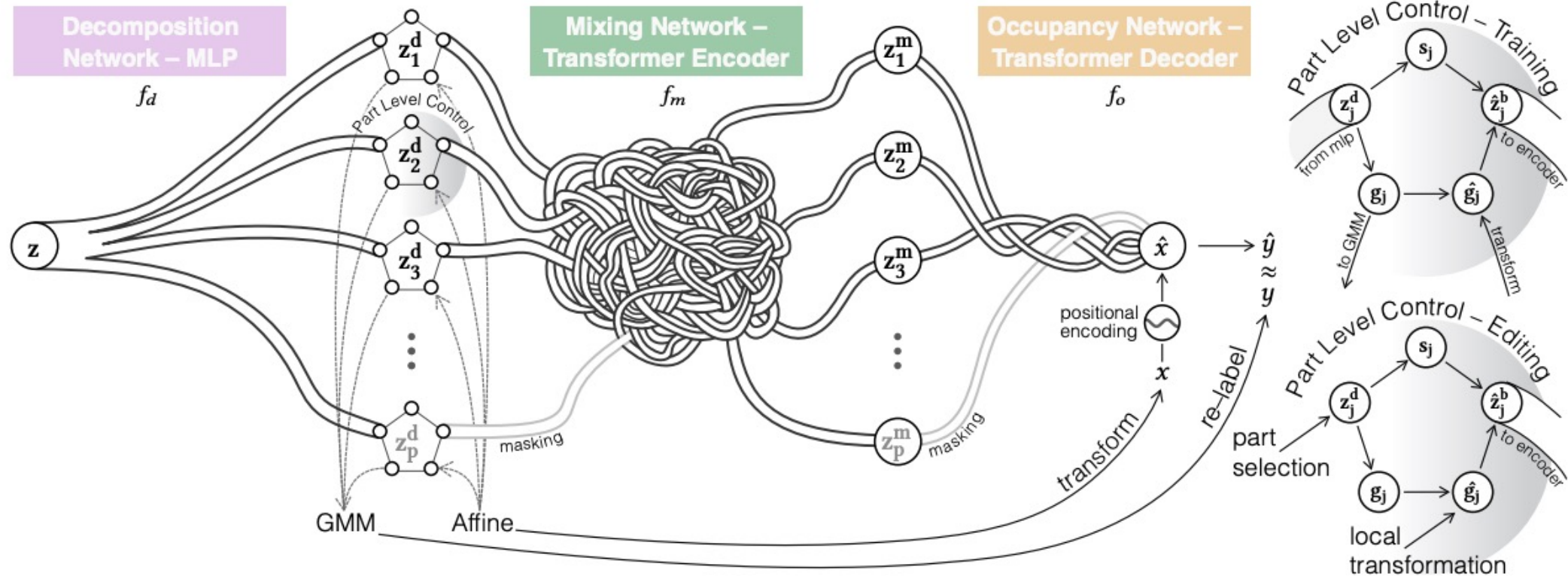
Leverages a novel *hybrid* representation describing

- *global* part-level structure *explicitly*, and
- *local* geometry *implicitly*.



# SPAGHETTI [Hertz et al., 2022]

The part decomposition and the explicit/implicit representations are learned in an **unsupervised** way.



# Part-Level Representation

For each part of an object learned in an *unsupervised* way,

- *Explicit* parameters of *Gaussian blubs* indicate position, scale, and rotation.



# Part-Level Representation

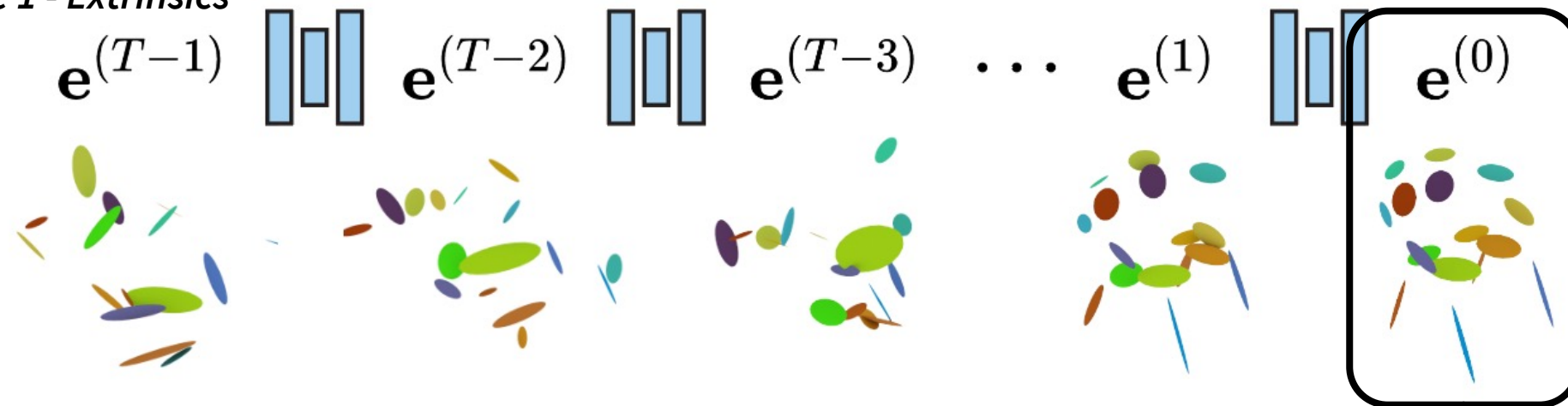
For each part of an object learned in an *unsupervised* way,

- *Implicit* latent feature is decoded into an *occupancy* function.

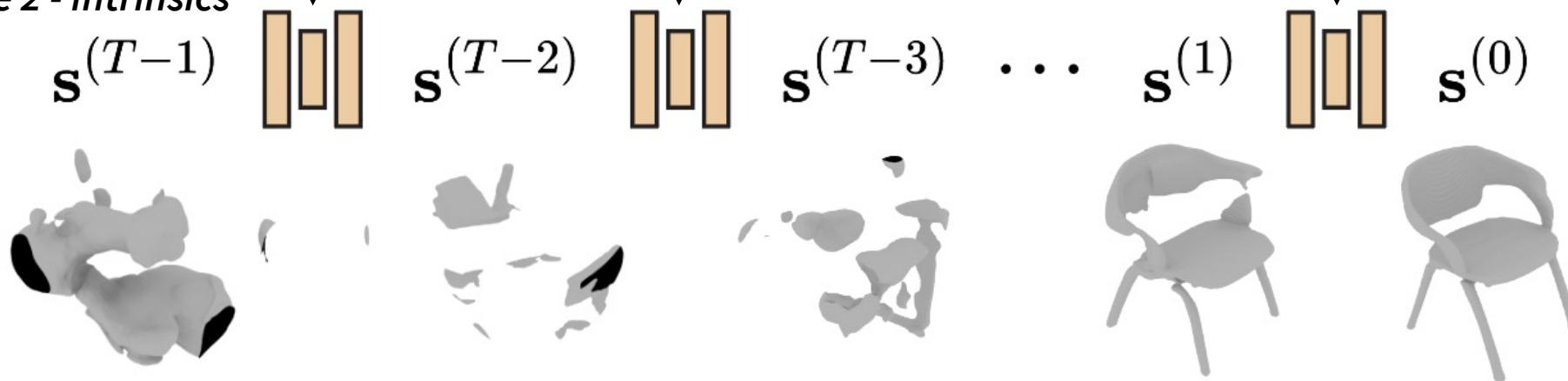


# Two-Phase Cascaded Diffusion

Phase 1 - Extrinsics



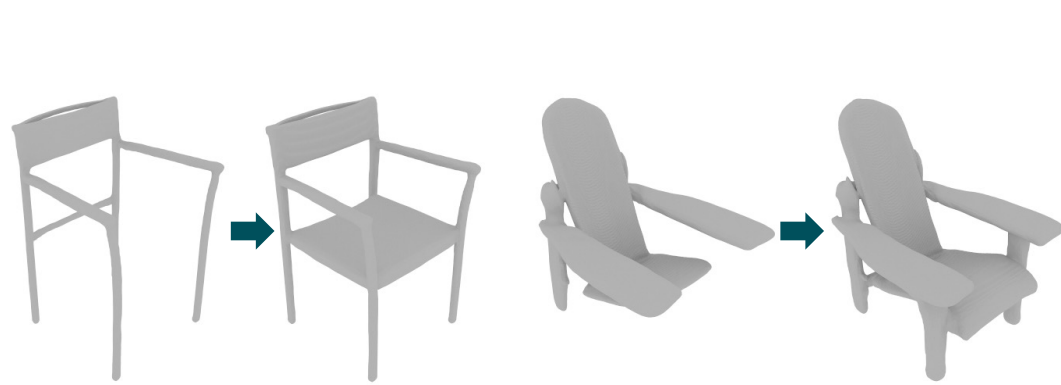
Phase 2 - Intrinsic



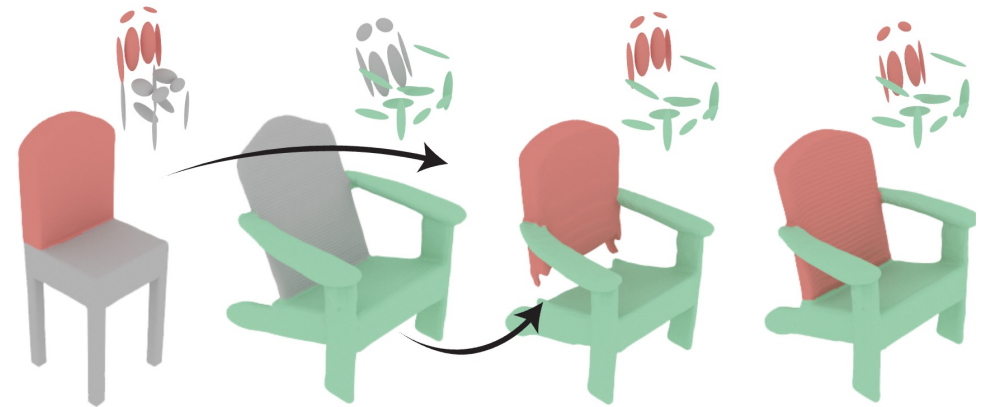
# Diversity of 3D Shapes



# Applications



**Part Completion**



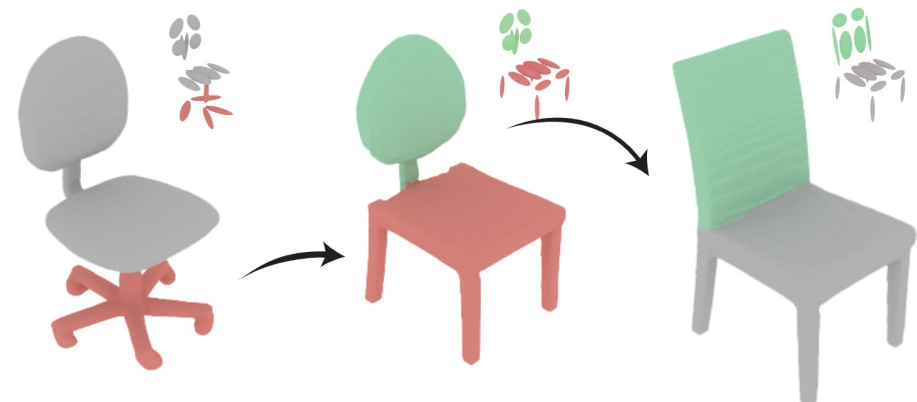
**Part Mixing**



"Chair has **round arms** and **wheels**."

"Its the one with **gaps** in the **back**."

**Text-to-3D Generation**



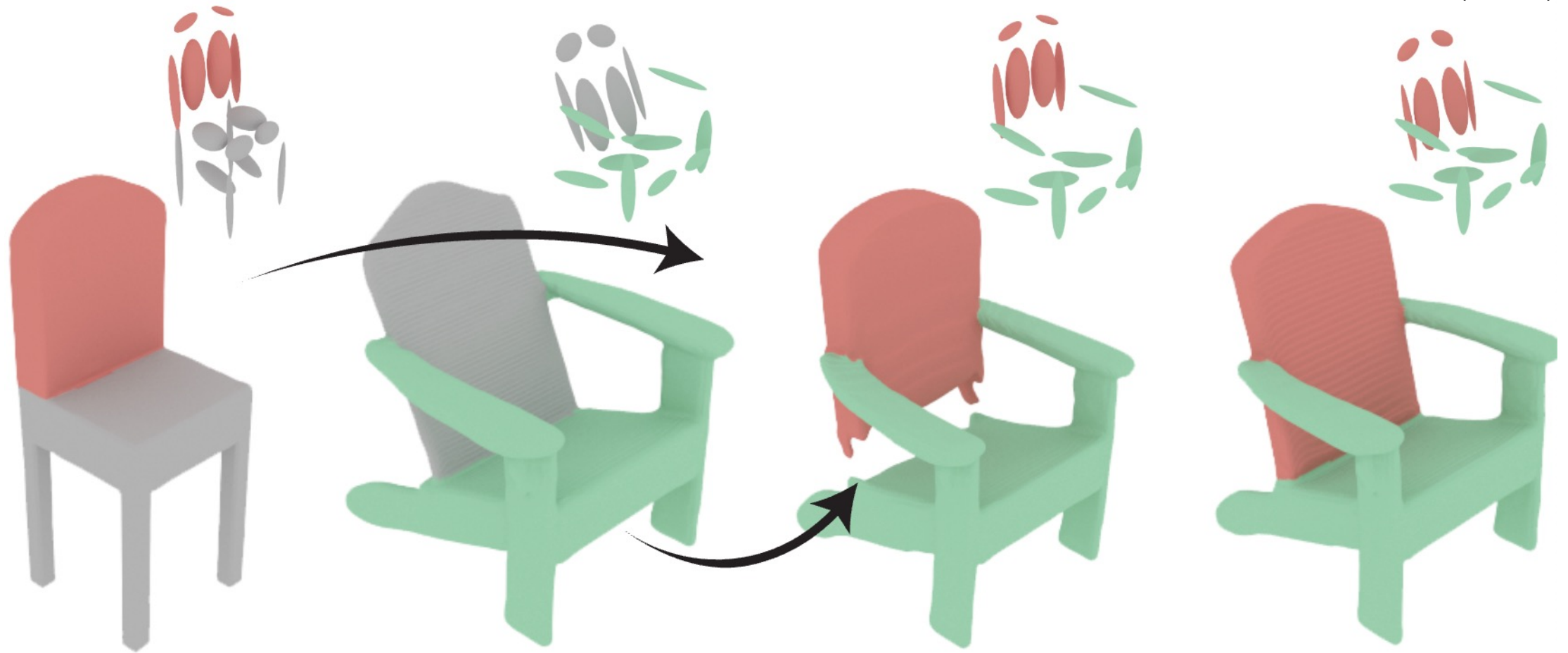
"A chair with **four legs**"

"**rectangle back** chair"

**Text-Guided Part Editing**

# Part Mixing

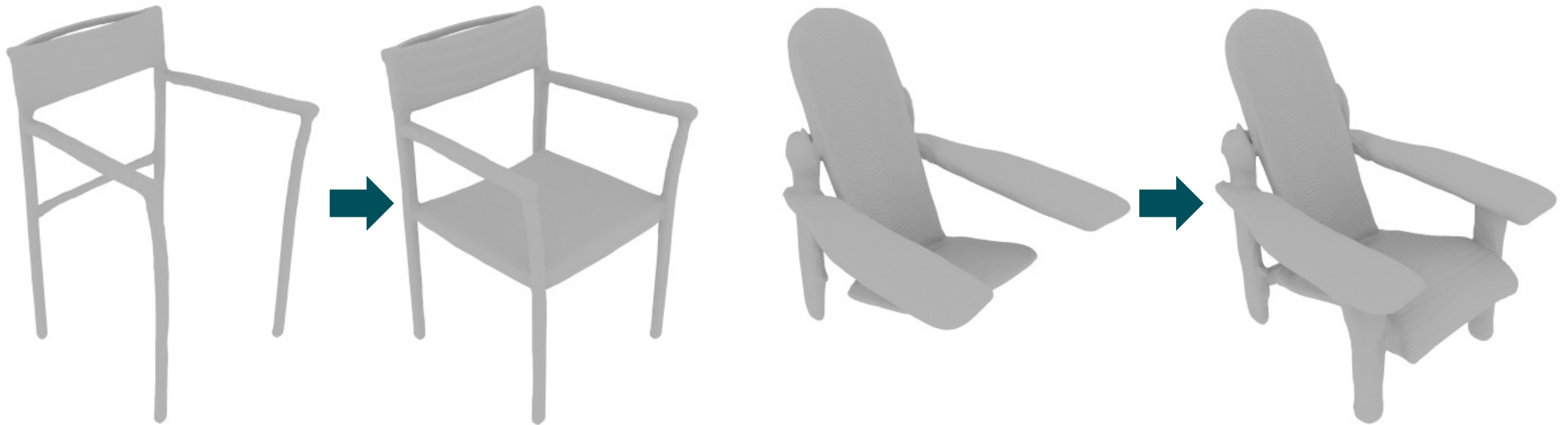
$$\mathbf{x}^{(0)} = m \odot \mathbf{x}_a^{(0)} + (1 - m) \odot \mathbf{x}_b^{(0)}$$
$$\mathbf{x}^{(t)} = \text{denoise}(\mathbf{x}^{(0)}, t)$$
$$\mathbf{x}'^{(0)} = \text{add\_noise}(\mathbf{x}^{(t)}, t)$$



# Part Completion

$$\begin{aligned} \mathbf{x}_f^{(t-1)} &= \text{denoise}(\mathbf{x}^{(t)}, 1) \\ \mathbf{x}_b^{(t-1)} &= \text{add\_noise}(\mathbf{x}^{(0)}, t) \\ \mathbf{x}^{(t-1)} &= m \odot \mathbf{x}_f^{(t-1)} + (1 - m) \odot \mathbf{x}_b^{(t-1)} \end{aligned}$$

Repeat for  $t = T, \dots, 1$ .



# Text-to-3D Generation

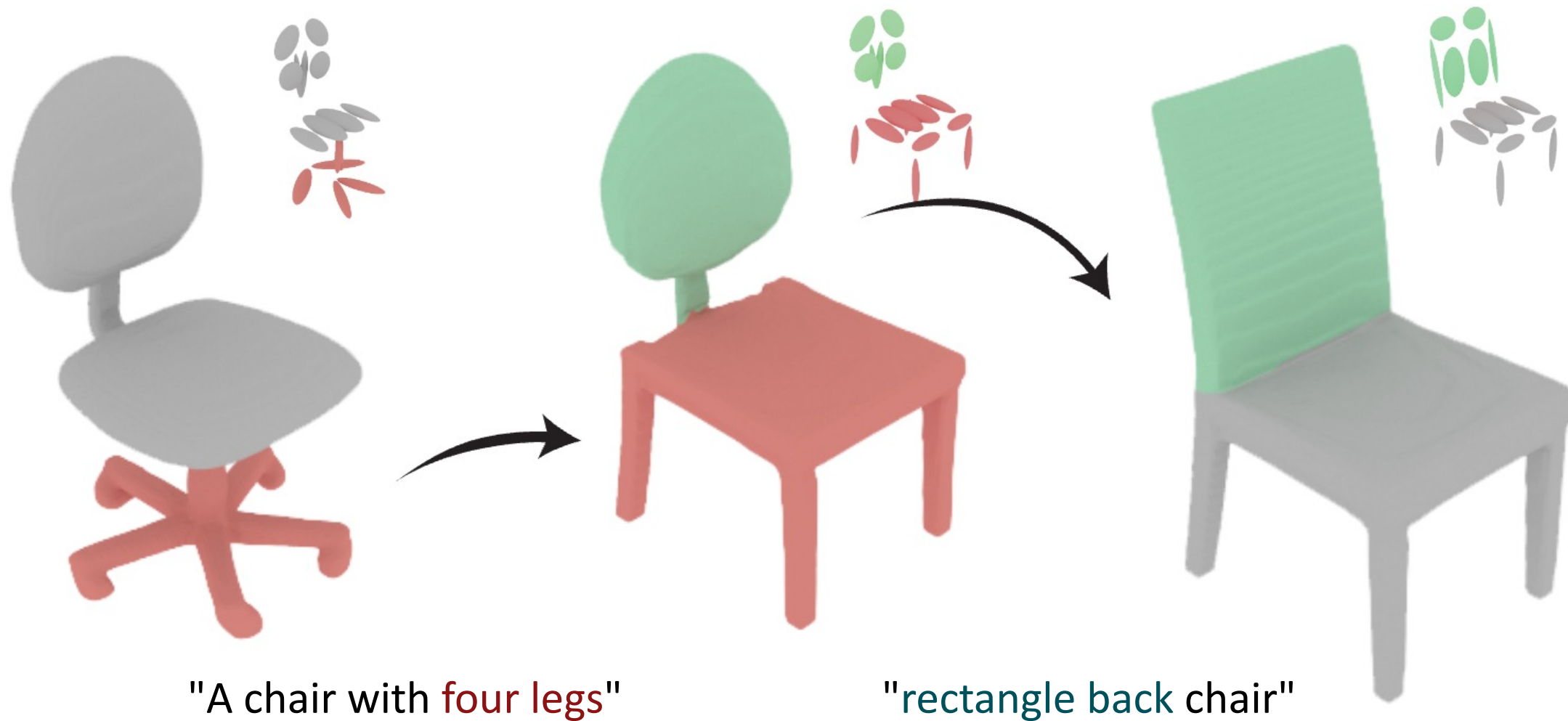


"Chair has **round arms** and **wheels**."



"Its the one with **gaps** in the **back**."

# Text-Guided Part Editing



# Limitation?

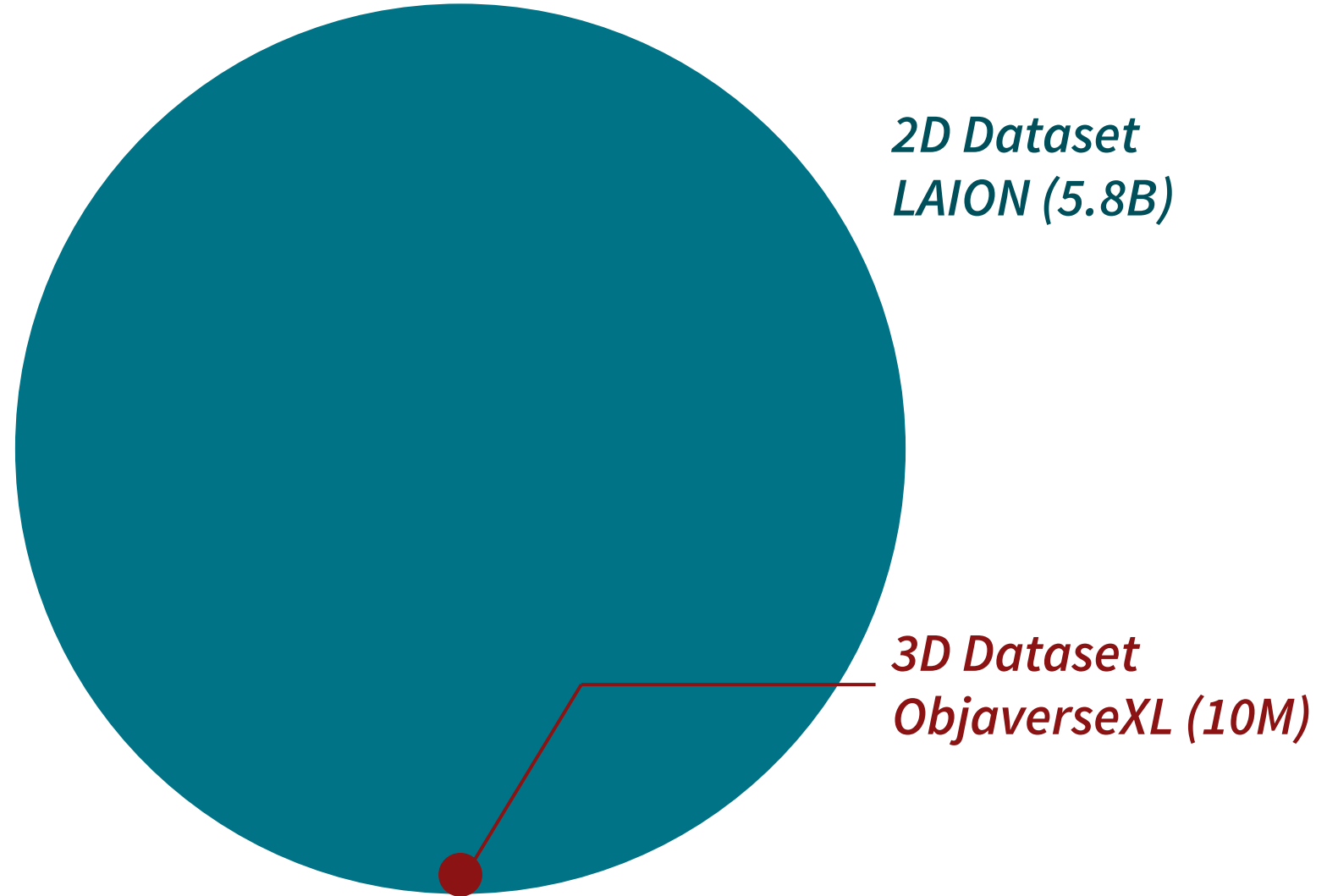


# Challenge: Lack of Large-Scale 3D Dataset



Objaverse-XL  
Allen Institute

# *Data Deficiency*



# Diversity of Imaginable 3D Shapes



“frog wearing a sweater”



“eggshell broken in two  
with an adorable chick  
standing next to it”



“ghost eating a hamburger”



“a pig wearing a backpack”

*We have a small-scale 3D dataset  
but images on an internet scale.*

*Images are projections of 3D  
from specific angles.*



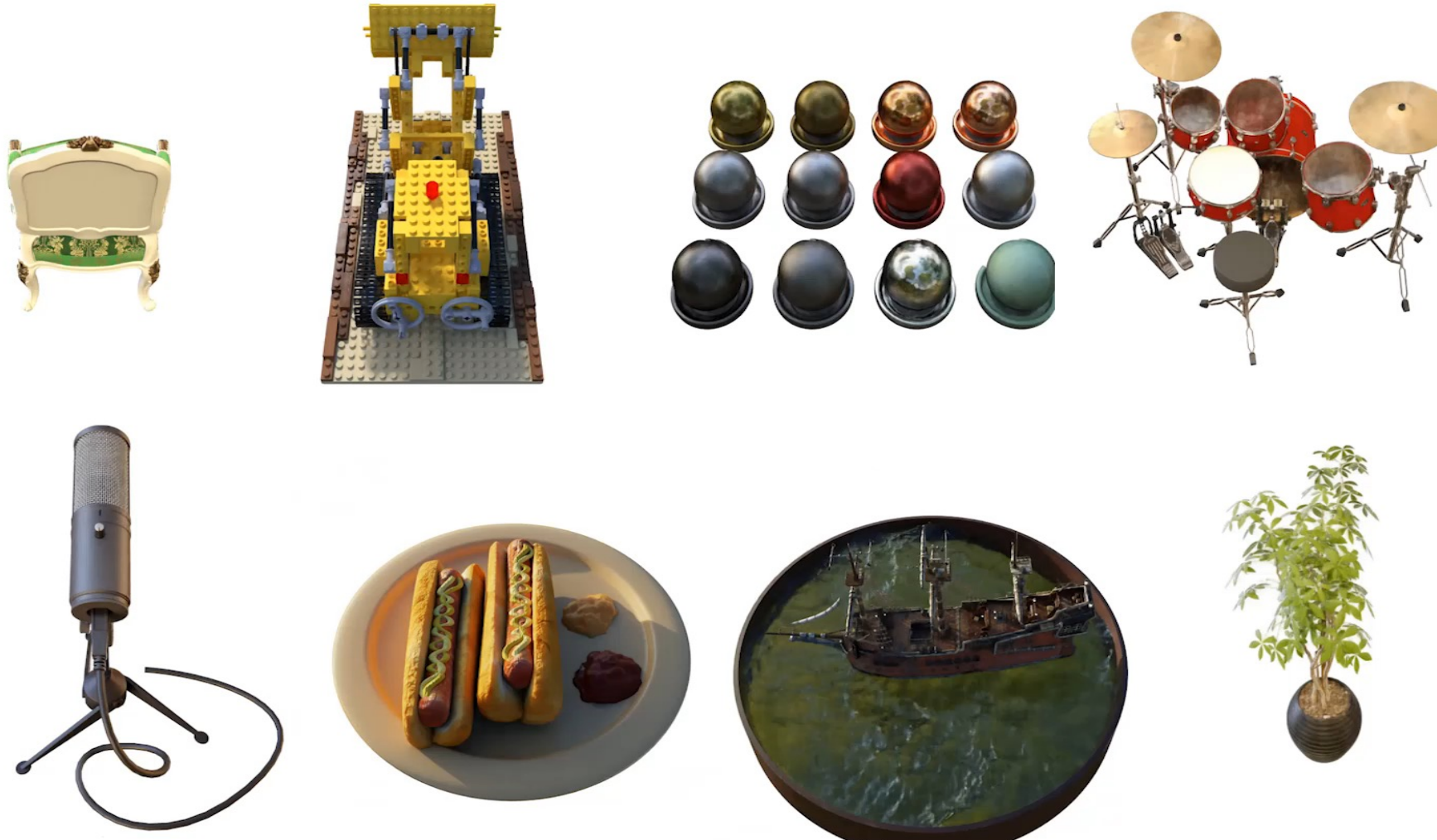
# *3D Generation*

*How to generate a 3D object  
from a collection of 2D images?*

# *3D Reconstruction*

*How to reconstruct a 3D object  
from a collection of 2D images  
of a specific object?*

# NeRF



# 3D Reconstruction

- *Input: A set of **images** with **camera poses**.*
- *Output: A representation of the **3D object**.*





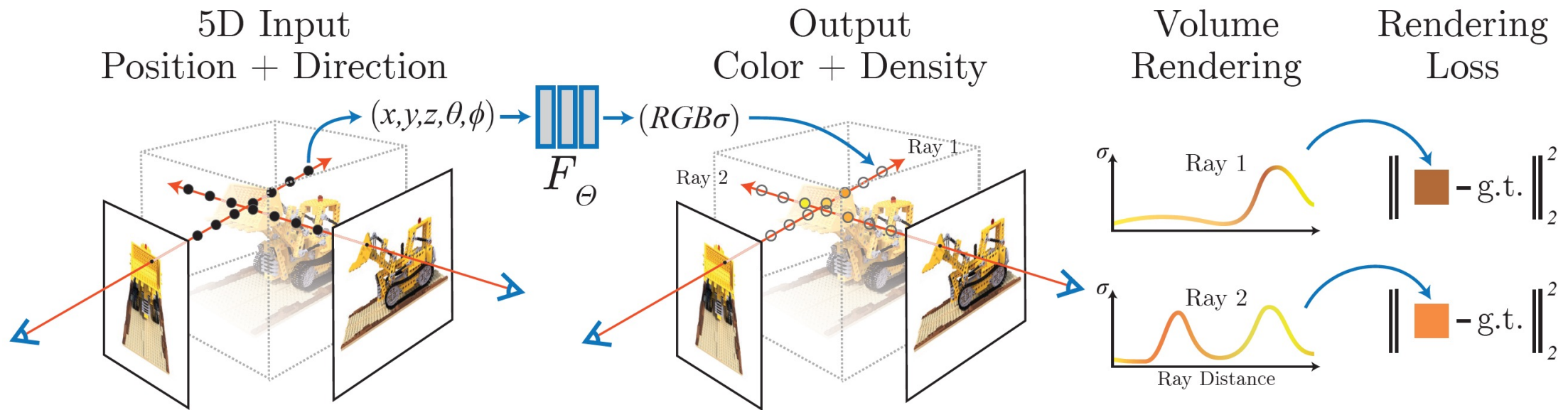
▼ Metrics

59.97 (16.67 ms)

**Gaussian Splatting**

# NeRF Optimization

1. **Render** a NeRF representation into a specific view.
2. Compute the **difference** with the **given image**.
3. Update the NeRF using **gradient descent**.



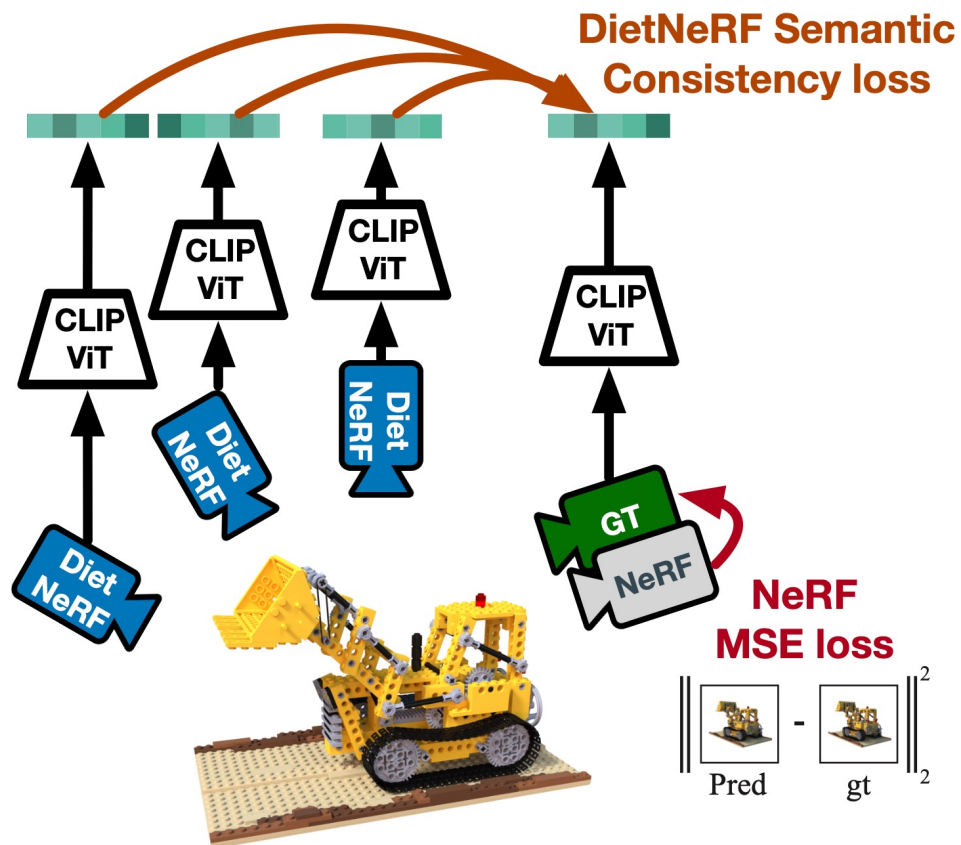
*Can we perform NeRF reconstruction  
with a few images?*

Then, we need additional information!

# Few-Shot NeRF

DietNeRF [Jain et al., ICCV 2021]

Use *priors* learned by **CLIP** [Radford et al., 2021], a *text-image model*.



“a bulldozer is a bulldozer from any perspective”

# CLIP [Radford et al., 2021]

CLIP takes a text-image pair as input and assesses the *alignment* between the text and the image.

“A beautiful painting of a dog riding a dolphin” +  = CLIP 70.7%

STEP 500

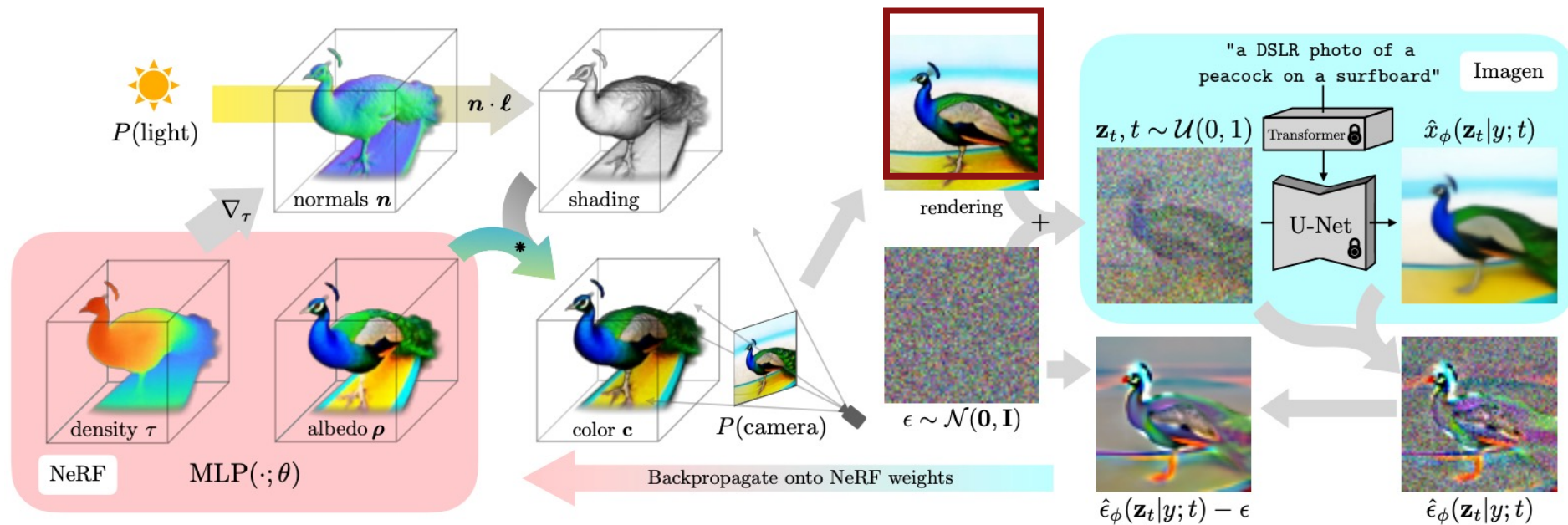
The image shows a text prompt on the left, a plus sign, a painting of a dog riding a dolphin in the center, an equals sign, and the CLIP score of 70.7% on the right. The painting is a colorful, detailed work showing a golden retriever sitting on the back of a grey dolphin. They are in a body of water with a town and a church in the background. The text prompt is in a monospace font. The CLIP score is in a bold, sans-serif font. The text 'STEP 500' is at the bottom of the image area.

# Knowledge Distillation in 3D Generation

1. *Render* the NeRF representation into a specific view

2. Compute the *alignment* to the given text.

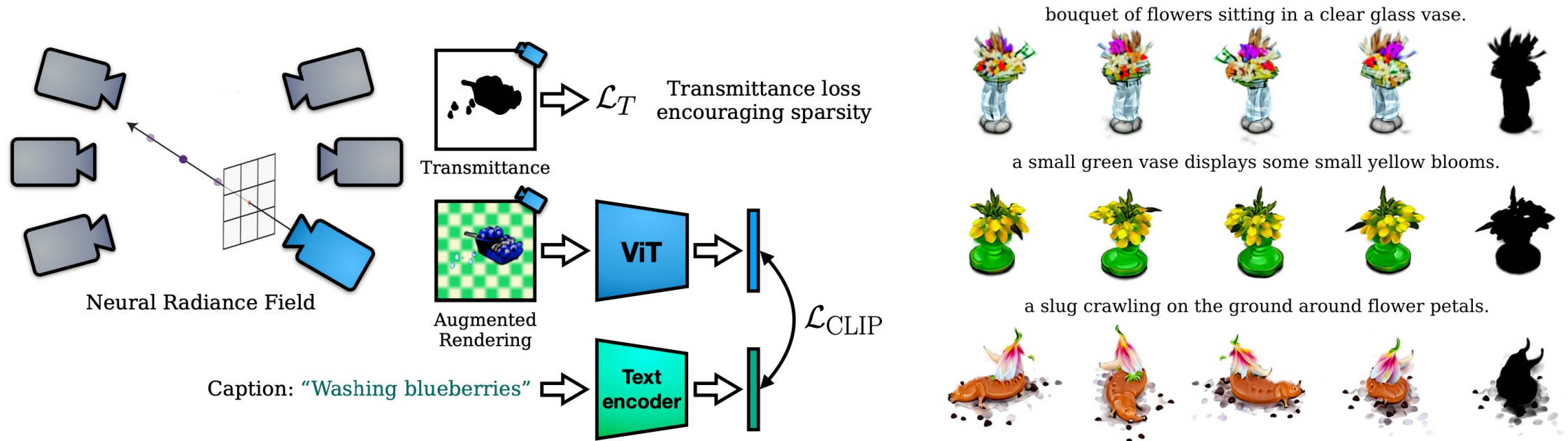
3. Update the NeRF using the *gradient descent*.



# Extrem Case: Zero-Shot NeRF

*DreamFields [Jain et al., CVPR 2022]*

Given *a text prompt but no images*, generate a 3D shape by maximizing similarity between a *rendered image* and the input prompt in the CLIP embedding space.



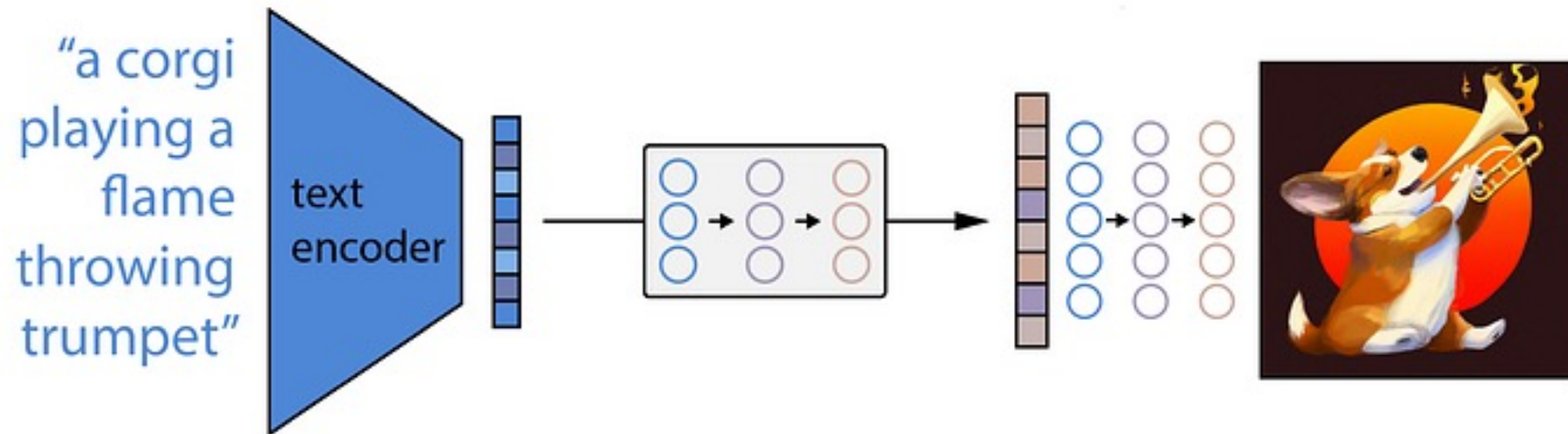
# DreamFields [Jain et al., CVPR 2022]



an archair in the shape of a \_\_\_\_.  
an archair imitating a \_\_\_\_.

a teapot in the shape of a \_\_\_\_.  
a teapot imitating a \_\_\_\_.

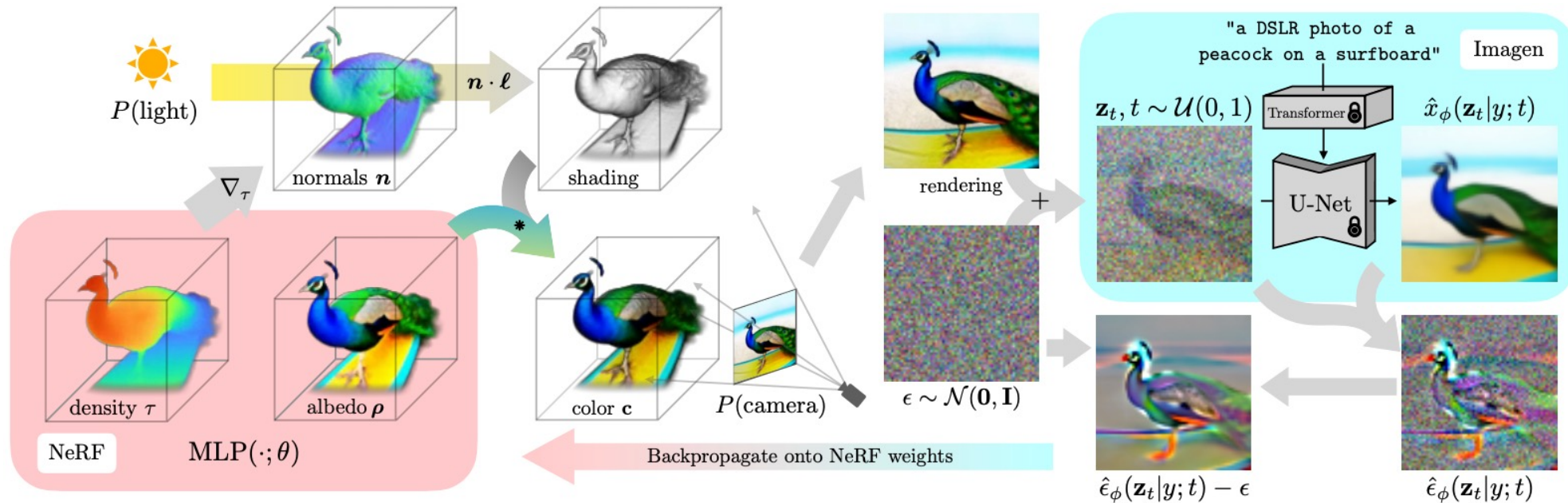
# Can we use an *image diffusion model* instead of CLIP?



# *Score Distillation Sampling (SDS)*

# DreamFusion [Poole et al., ICLR 2023]

Proposed the idea of **Score Distillation Sampling (SDS)**, leveraging a pretrained diffusion model to **measure the plausibility** of rendered images.



# Score Distillation Sampling (SDS)

- How can we utilize a pretrained diffusion model to *measure the plausibility* of rendered images?

- Review the *loss* function:

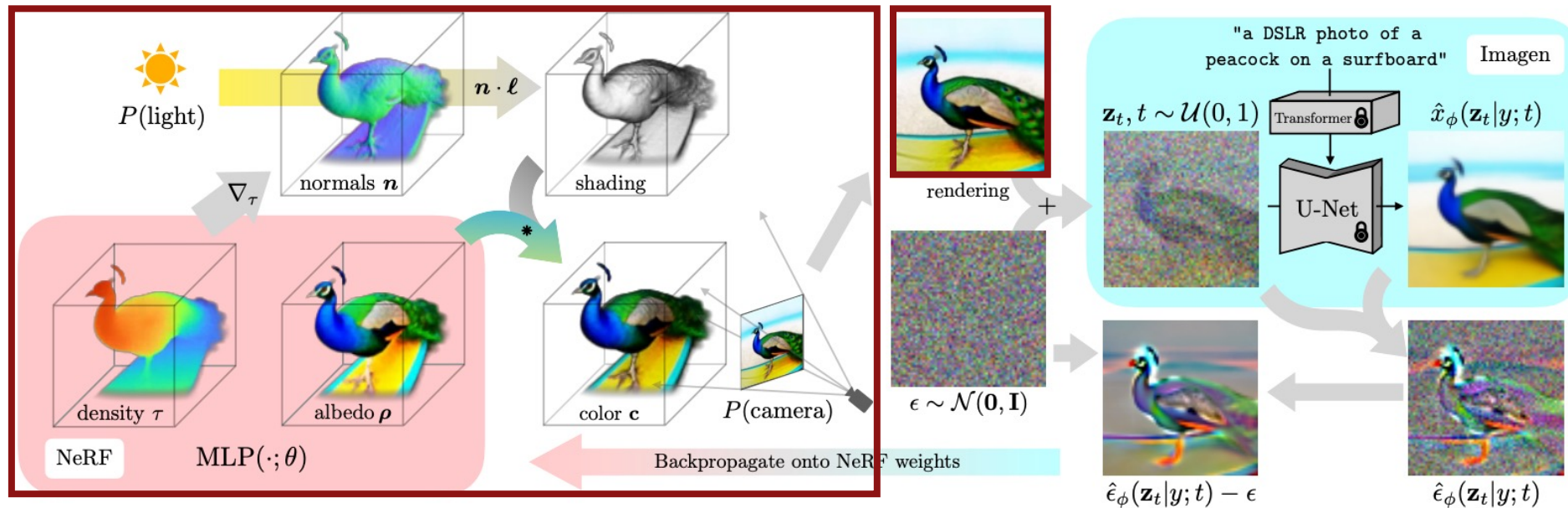
$$\mathcal{L} = \mathbb{E}_{t \sim [1, T], \mathbf{x}_0, \boldsymbol{\epsilon}_t} \left[ \left\| \boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta \left( \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t, t \right) \right\|^2 \right]$$

If the training of the diffusion model has converged, the loss for *real* data  $\mathbf{x}_0$  will be close to zero.

# Score Distillation Sampling (SDS)

1. **Render** the NeRF representation into a specific view.

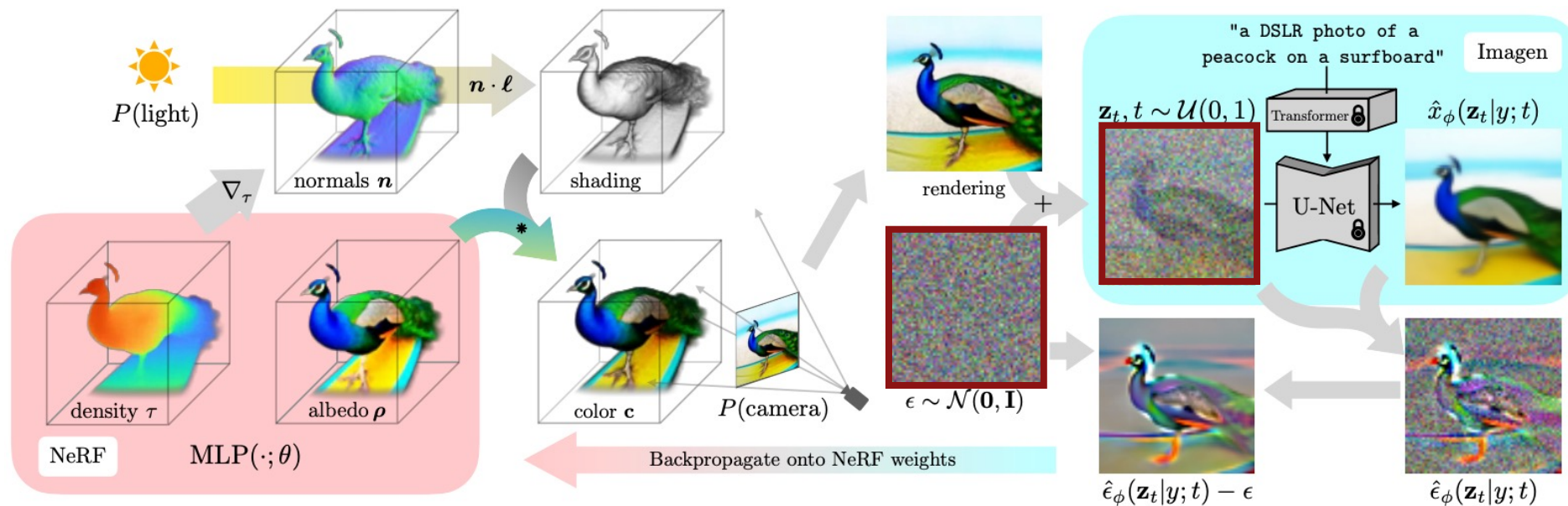
Let  $\phi$  denote the NeRF parameter, and  $\mathbf{x}_0 = g(\phi)$  denote the rendered image.



# Score Distillation Sampling (SDS)

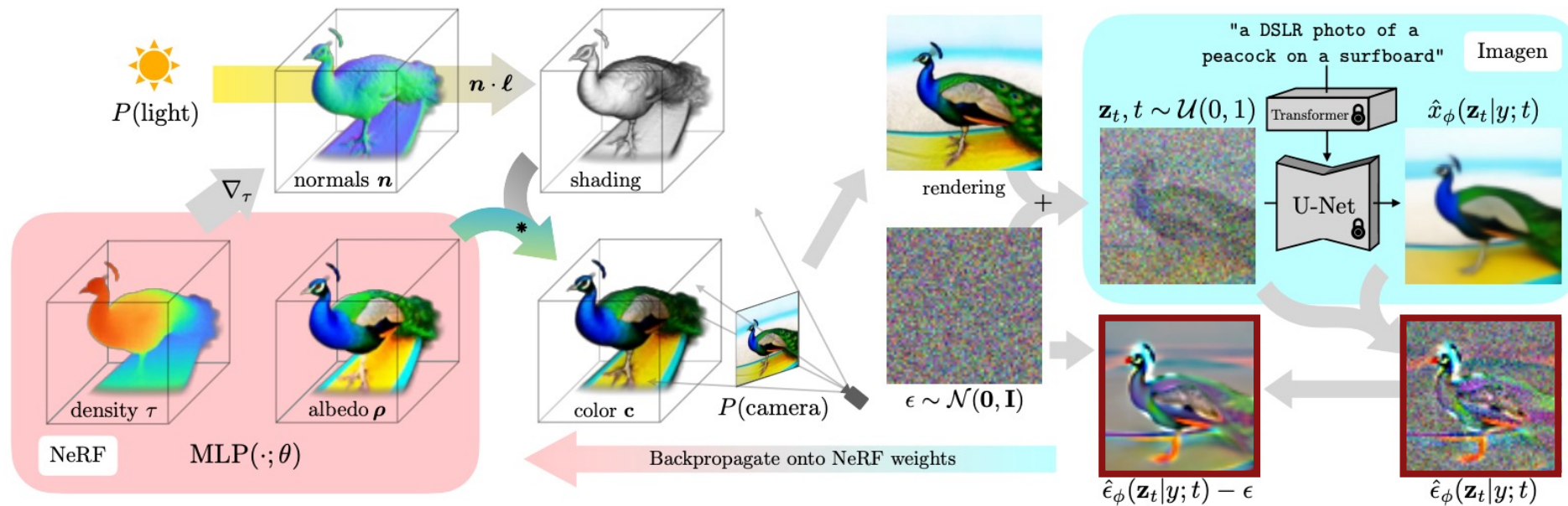
2. Add *noise* to the rendered image  $\mathbf{x}_0$ :

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_t.$$



# Score Distillation Sampling (SDS)

3. Perform *gradient descent* on  $\mathcal{L}$  with respect to the NeRF parameters  $\phi$ .



# DreamFusion Results



“frog wearing a sweater”



“eggshell broken in two with an adorable chick standing next to it”



“ghost eating a hamburger”



“a pig wearing a backpack”

# *Why SDS Instead of Reverse Diffusion?*

*This is a scenario where the images are **parameterized differently from how they were represented during the training of the diffusion model.***

- *Training: Per-pixel colors.*
- *Inference: NeRF rendering.*

# Example: Vector Images / Sketches

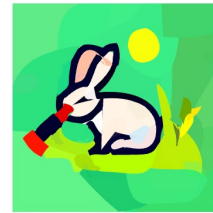
The same idea but with a *different parameterization* of images.



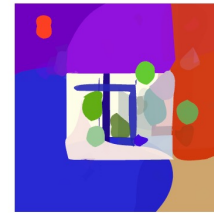
A blue poison-dart frog sitting on a water lily. [...]



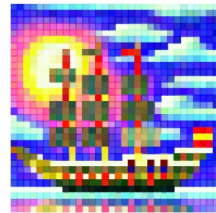
the Great Pyramid. [...]



A rabbit cutting grass with a lawnmower. [...]



Translation. [...]



a pirate ship landing on the moon. [...]



A snail on a leaf. [...]



Yeti taking a selfie. [...]



a hedgehog. [...]



The space between infinity. [...]



A realistic photograph of a cat. [...]



A watercolor painting of a cat. [...]



A Japanese woodblock print of one cat. [...]

# Example: Mesh Editing

The same idea but with a *different parameterization of images.*

Source



ARAP

[Sorkine and Alexa, 2007]



APAP

[Yoo et al., 2024]



# Stable-DreamFusion

## Stable-Dreamfusion

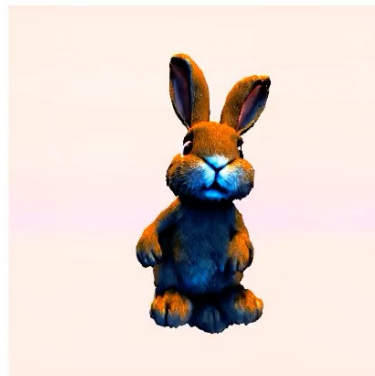
A pytorch implementation of the text-to-3D model Dreamfusion, powered by the [Stable Diffusion](#) text-to-2D model.

**ADVERTISEMENT:** Please check out [threestudio](#) for recent improvements and better implementation in 3D content generation!

**NEWS (2023.6.12):**

- Support of [Perp-Neg](#) to alleviate multi-head problem in Text-to-3D.
- Support of Perp-Neg for both [Stable Diffusion](#) and [DeepFloyd-IF](#).

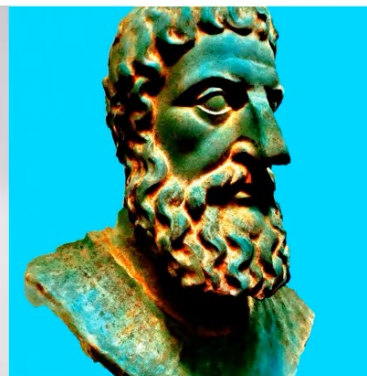
### Text-to-3D



a rabbit, animated movie character, high detail 3d model



a DSLR photo of a delicious hamburger



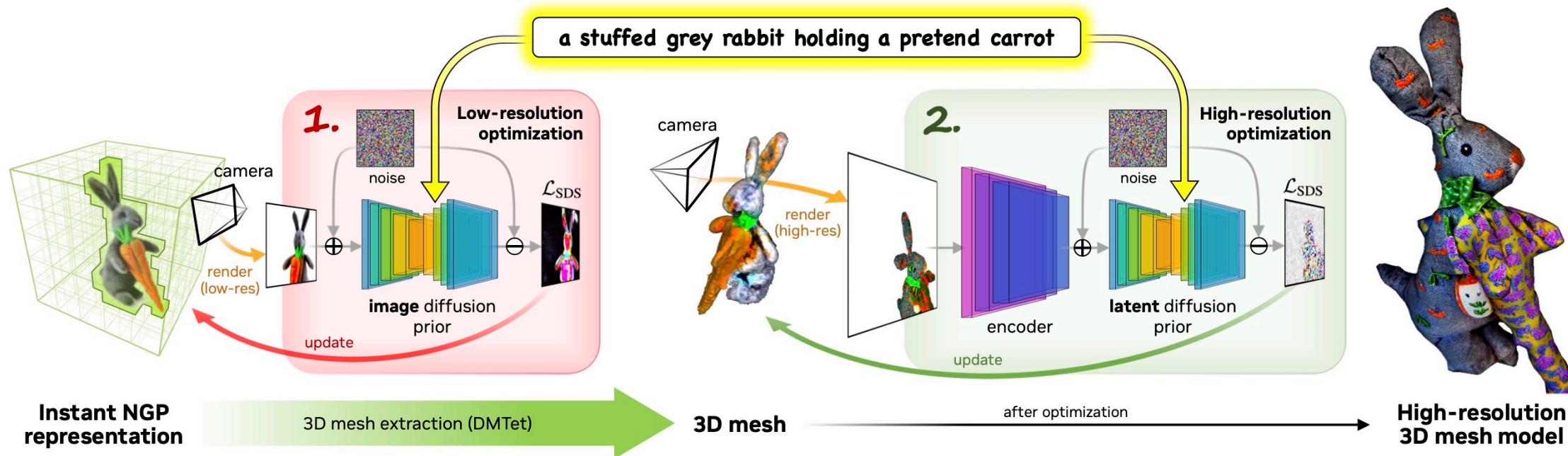
a highly detailed stone bust of Theodoros Kolokotronis



a small saguaro cactus planted in a clay pot

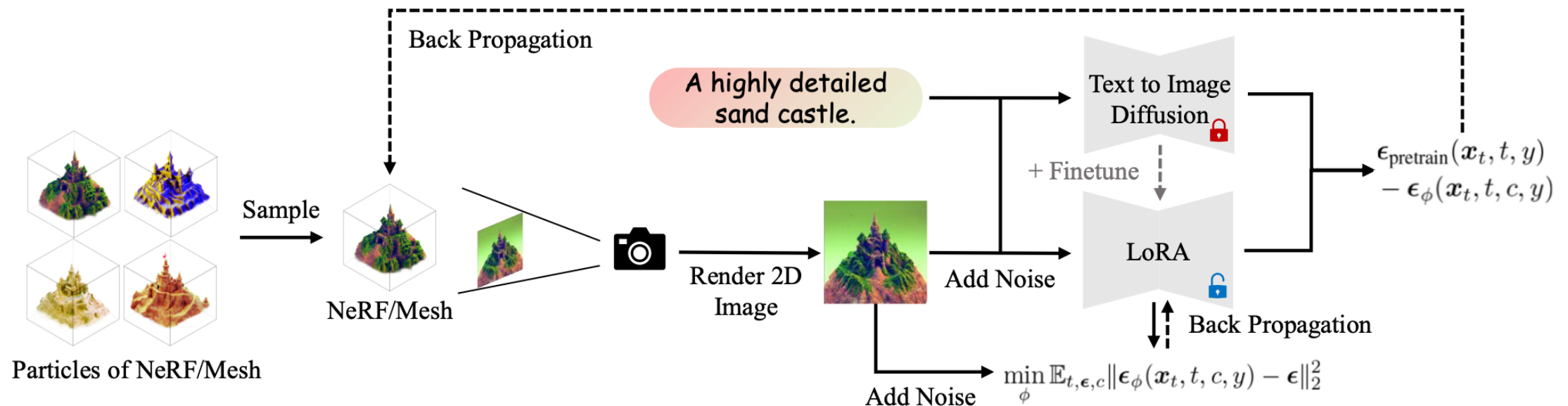
# Magic3D [Lin et al., CVPR 2023]

- *Two-stage approach.*
- *Extract a coarse mesh in the first stage, and then texture the mesh in the second stage.*



# ProlificDreamer [Wang et al., arXiv 2023]

- Minimize the SDS loss for the **multiple samples** of the NeRF parameters  $\phi$ .
- Finetune the diffusion model with the **Low Rank Adaptation (LoRA)** technique.

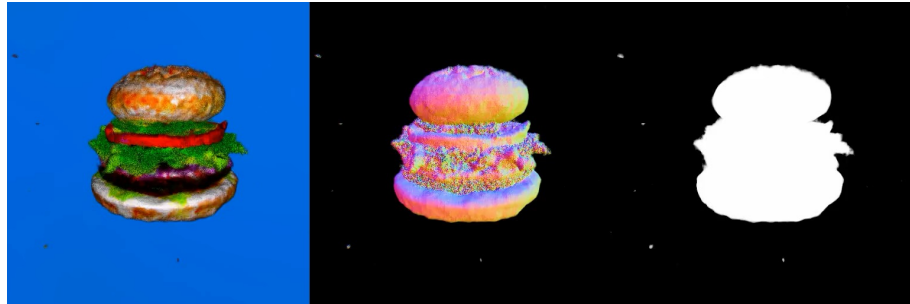


# Limitation of SDS

It does not converge well without a **high CFG weight** (e.g.,  $w = 400$ ) and thus suffers from **model collapse**.

“a delicious hamburger”

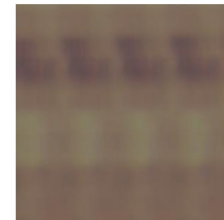
$w = 400$



$w = 7.5$



Credit: Jaihoon Kim



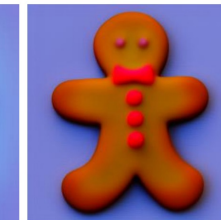
NeRF Initial.



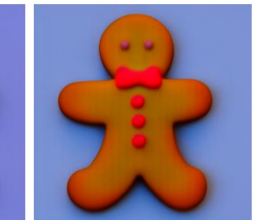
Seed=0



Seed=1



Seed=2

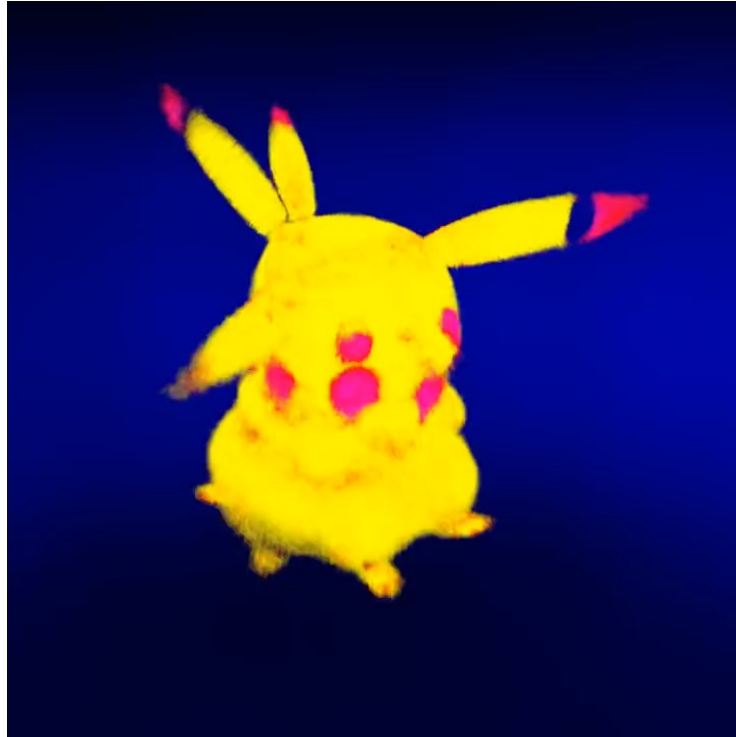
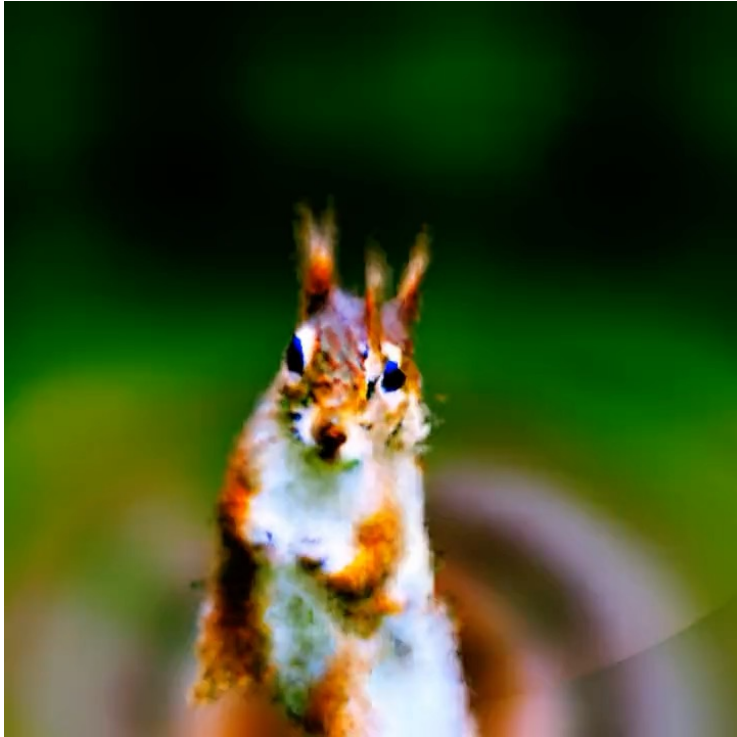


Seed=3

“gingerbread man”

Huang et al., *DreamTime: An Improved Optimization Strategy for Text-to-3D Content Creation*, arXiv 2023.

# *Limitation of 3D Generation from 2D Priors*



*Twitter @\_akhaliq*

# Limitation of 3D Generation from 2D Priors

Supervision for **geometry** is still needed!



*ProlificDreamer, Wang et al., 2023.*

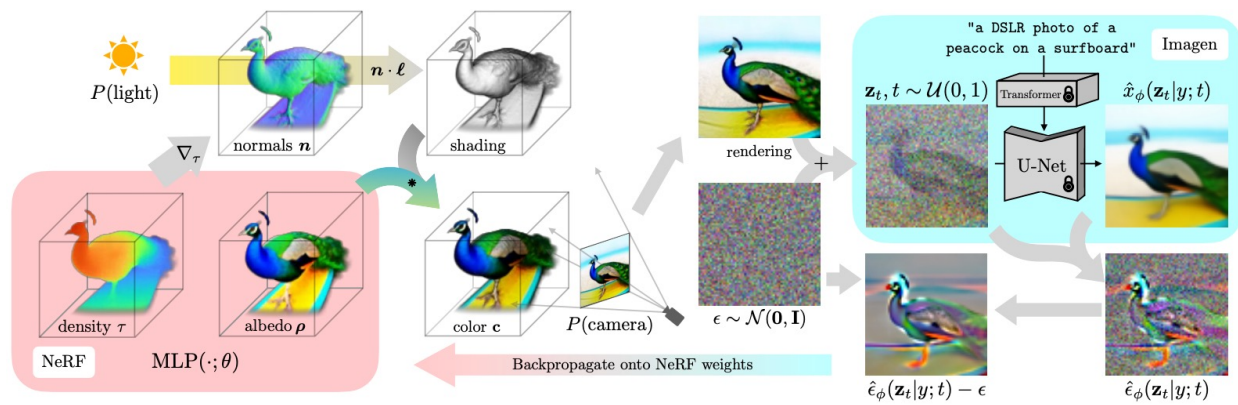


*StableDreamFusion*

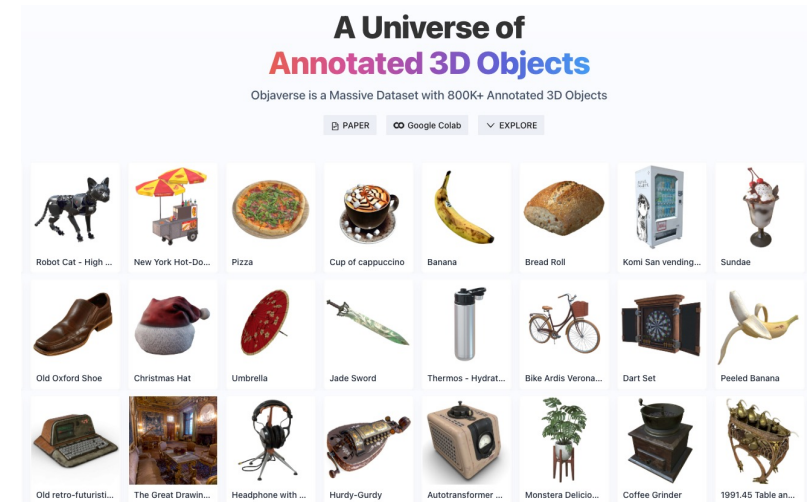
# *What's Next for 3D Generative Models?*

# 1. Combining 3D Supervision

While *pretrained image generative models* will continue to be valuable for 3D generation, the key to producing realistic and solid 3D shapes would lie in utilizing *small-scale 3D priors*.



DreamFusion  
Google

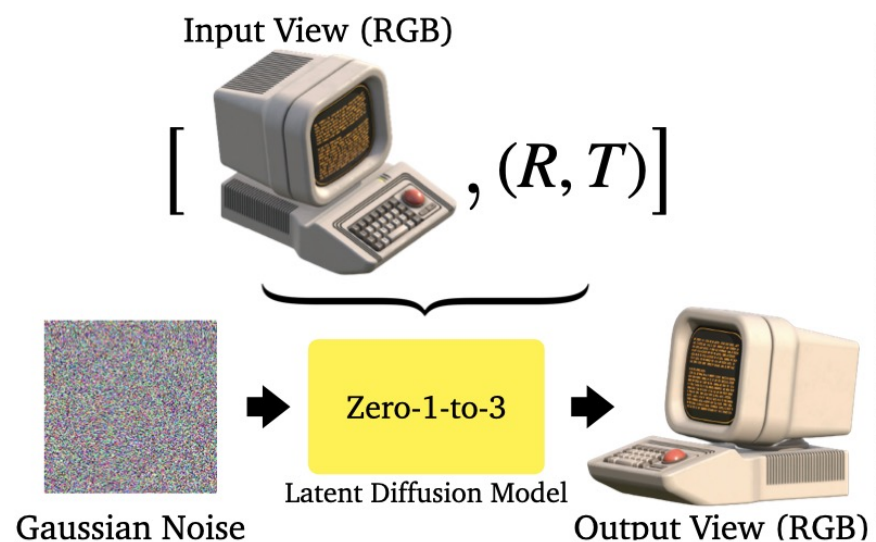


Objaverse  
Allen Institute

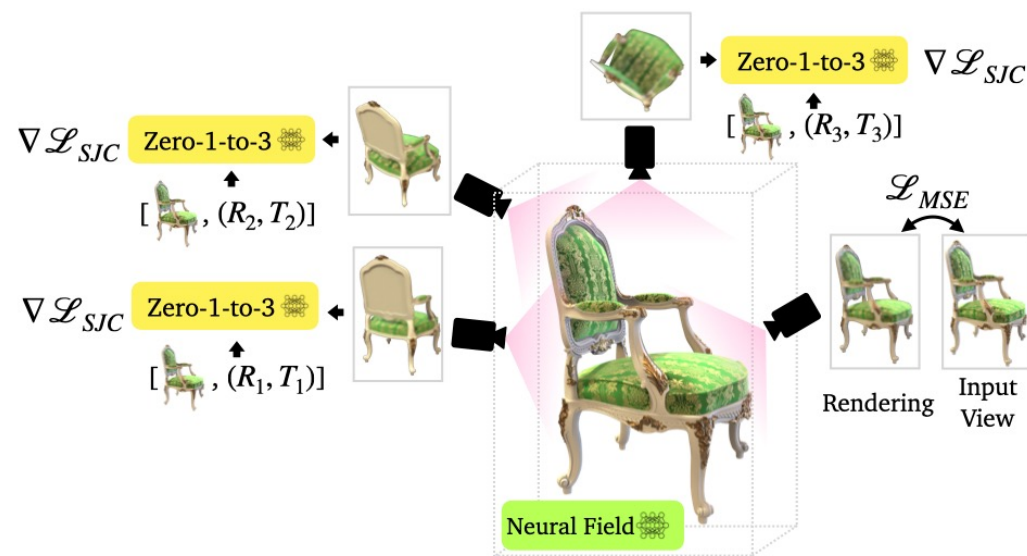
# Zero-1-to-3 [Liu et al., 2023]

## Novel view generation

*An image diffusion model generating a novel view image conditioned by another view image and camera pose.*



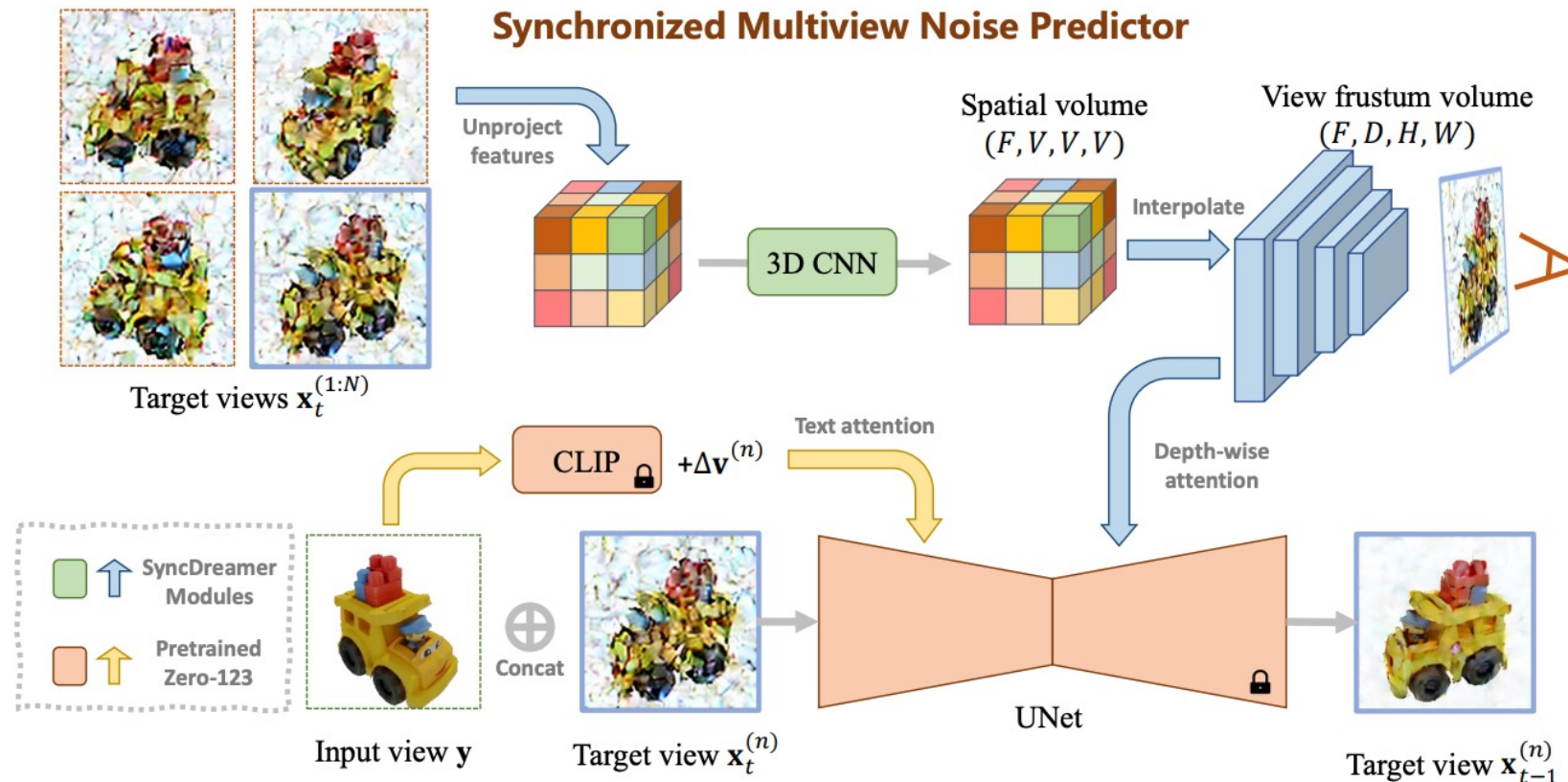
Novel View Synthesis



3D Reconstruction

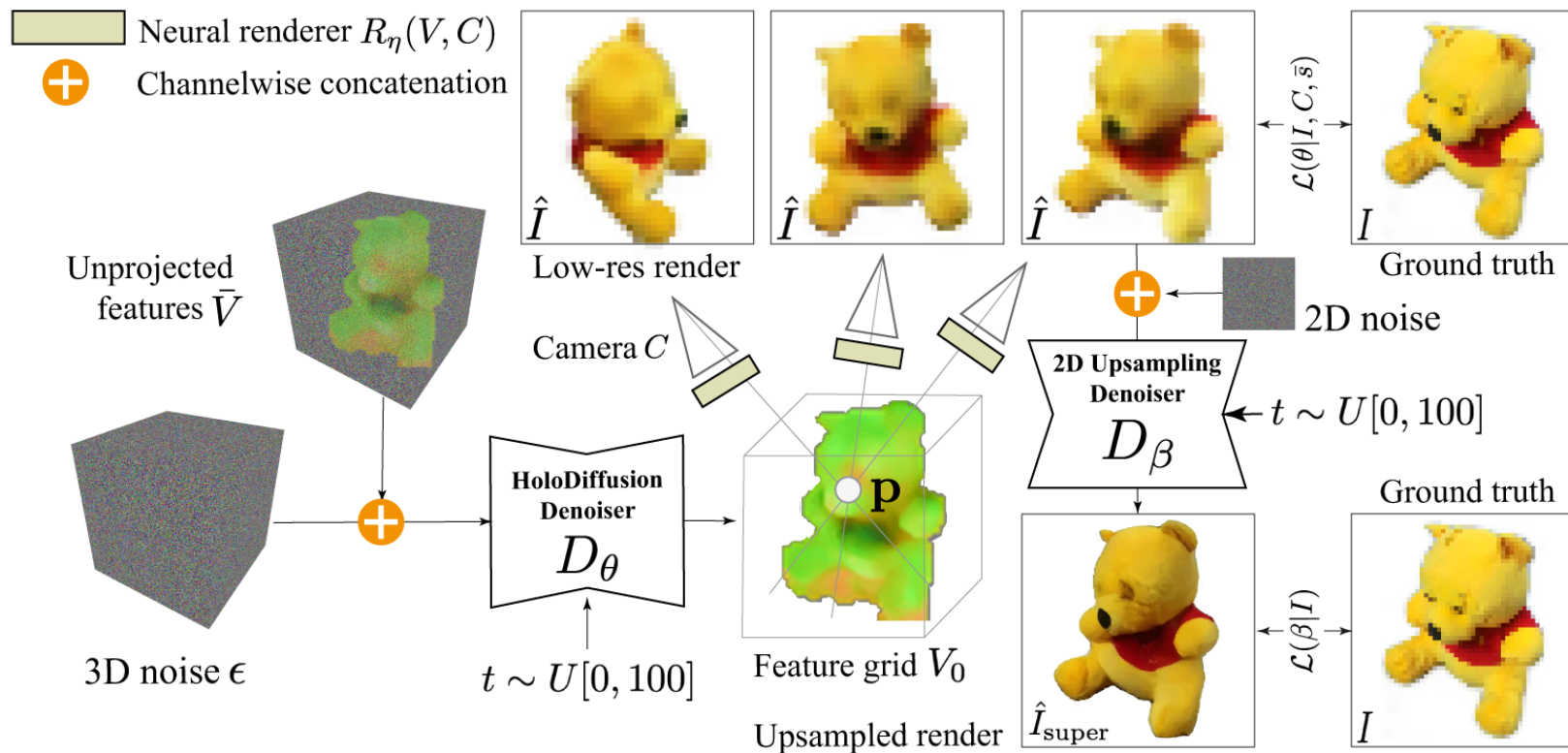
# SyncDreamer [Liu et al., 2023]

Utilize Zero 1-to-3 to learn the joint probability distribution of multi-view images.



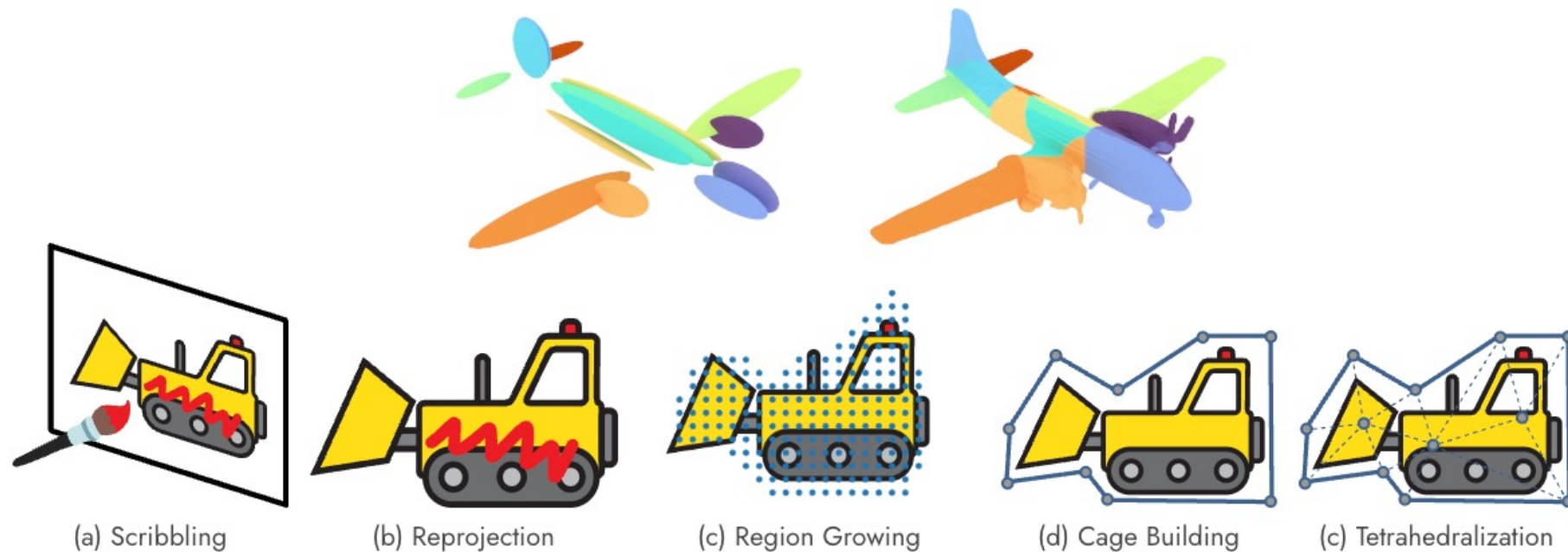
# HoloFusion [Karnewar et al., 2023]

- Train a 3D diffusion model using multi-view images only.
- Can be extended to integrate 2D priors.



## 2. Generation → Editing

The focus of 3D generative models will shift towards creating *versatile models* capable of not only generating but *editing* and *manipulating 3D shapes*.



*NeRFshop, Jambon et al., I3D 2023.*

# Delta Denoising Score [Hertz et al., 2023]

A new loss function for zero-shot *image* editing.



“Photo of sunset, winter, river.” → “Photo of sunset, winter, river with polar bears.”



“Amaryllis, North Window, oil on linen.” → “Cactus, North Window, oil on linen.”



“Forest painting, autumn trees.” → “Moon over forest painting, autumn trees.”

# Posterior Distillation Sampling [Koo et al., 2024]

A new loss function for zero-shot **NeRF** editing.

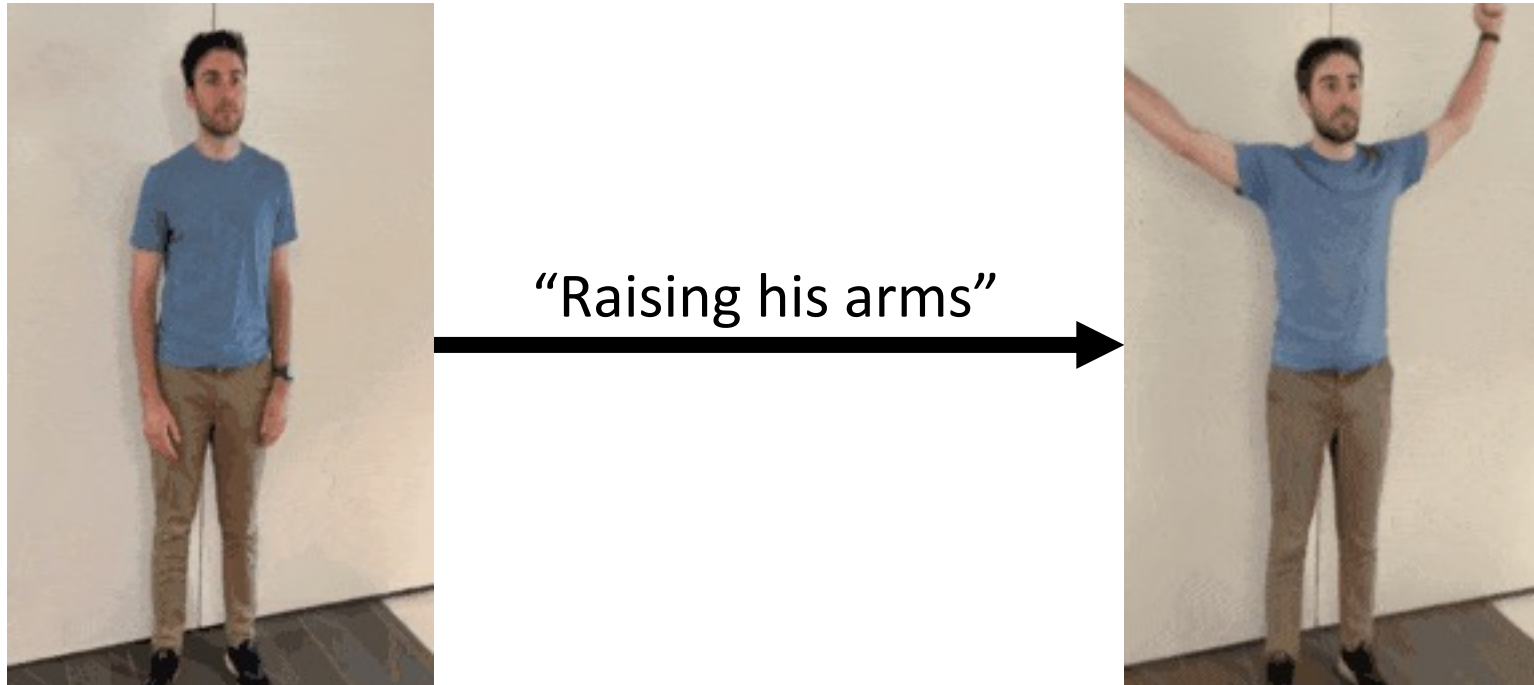


“Leonardo DiCaprio”



# Posterior Distillation Sampling [Koo et al., 2024]

A new loss function for zero-shot **NeRF** editing.



## 3. Texture Generation

How can we generate *photorealistic texture* given a 3D mesh or Gaussian splats?

