

StyleGuide

- ❖ Crafting visual style prompting with negative visual query guidance
-

NAVER AI Lab, Generation research

김준호

<https://github.com/taki0112>



NAVER AI LAB

Visual Style Prompting with Swapping Self-Attention

Jaeseok Jeong ^{*1,2} Junho Kim ^{*2} Yunjey Choi ² Gayoung Lee ² Youngjung Uh ¹

Abstract

In the evolving domain of text-to-image generation, diffusion models have emerged as powerful tools in content creation. Despite their remarkable capability, existing models still face challenges in achieving controlled generation with a consistent style, requiring costly fine-tuning or often inadequately transferring the visual elements due to content leakage. To address these challenges, we propose a novel approach, visual style prompting, to produce a diverse range of images while maintaining specific style elements and nuances. During the denoising process, we keep the query from original features while swapping the key and value with those from reference features in the late self-attention layers. This approach allows for the visual style prompting without any fine-tuning, en-

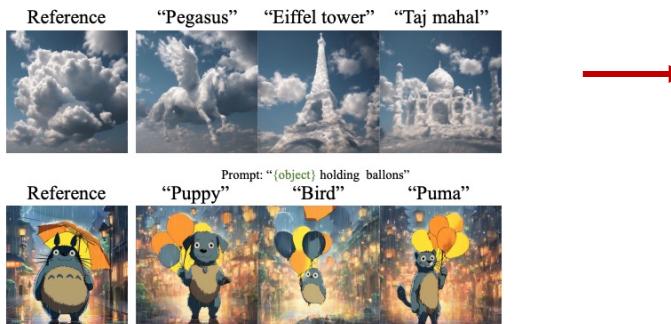


Figure 1. We tackle visual style prompting, reflecting style elements from reference images and contents from text prompts, in a training-free manner.

arxiv

STYLEGUIDE: PREVENT CONTENT LEAKAGE USING NEGATIVE QUERY GUIDANCE

Anonymous authors

Paper under double-blind review

ABSTRACT

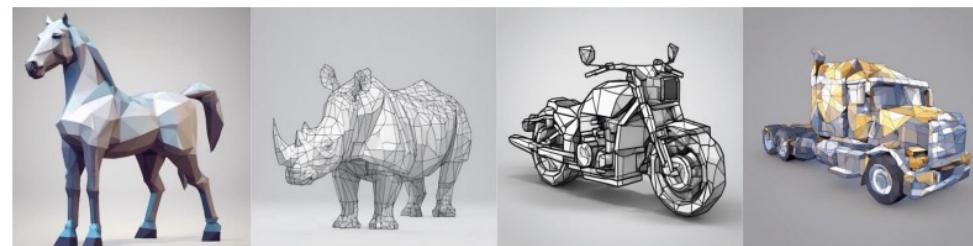
In the domain of text-to-image generation, diffusion models have emerged as powerful tools. Recently, studies on visual prompting, where images are used as prompts, have enabled more precise control over style and content. However, existing methods often suffer from content leakage, where undesired elements from the visual style prompt are transferred along with the intended style (content leakage). To address this issue, we 1) extend classifier-free guidance (CFG) to utilize swapping self-attention and propose 2) negative visual query guidance (NVQG) to reduce the transfer of unwanted contents. NVQG employs negative score by intentionally simulating content leakage scenarios which swaps queries instead of key and values of self-attention layers from visual style prompts. This simple yet effective method significantly reduces content leakage. Furthermore, we provide careful solutions for using a real image as a visual style prompts and for image-to-image (I2I) tasks. Through extensive evaluation across various styles and text prompts, our method demonstrates superiority over existing approaches, reflecting the style of the references and ensuring that resulting images match the text prompts.

ICLR 2025



“low-poly style horse”

(a) Various style results within a text description



“Horse” “Rhino” “Motorcycle” “Truck”
+ “low-poly style {object}. low-poly game art, polygon mesh, jagged, blocky,
wireframe edges, centered composition, simple ..”

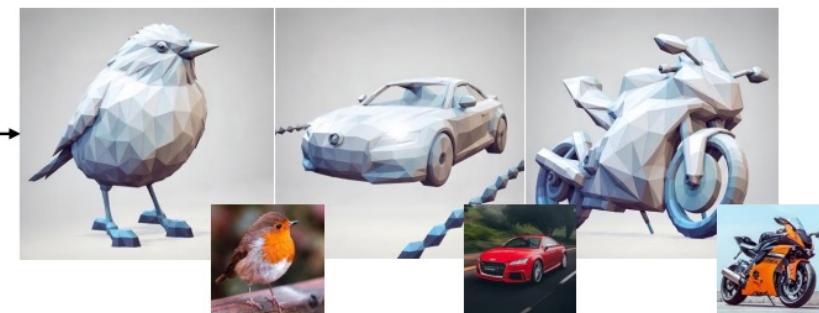
(b) Various style results with a highly detailed text description

Reference



“Rhino” “Motorcycle” “Truck”

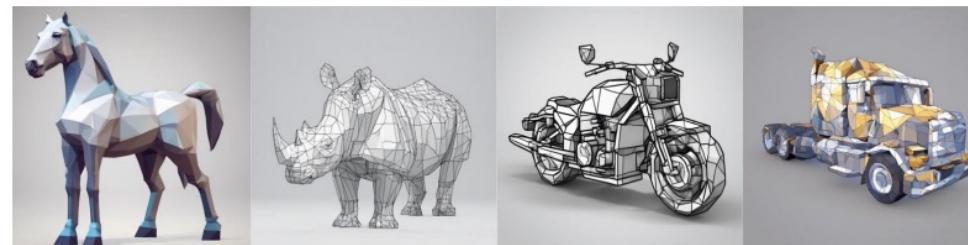
(c) T2I results specified by a style reference



(d) Style transfer results specified by a style reference



(a) Various style results within a text description



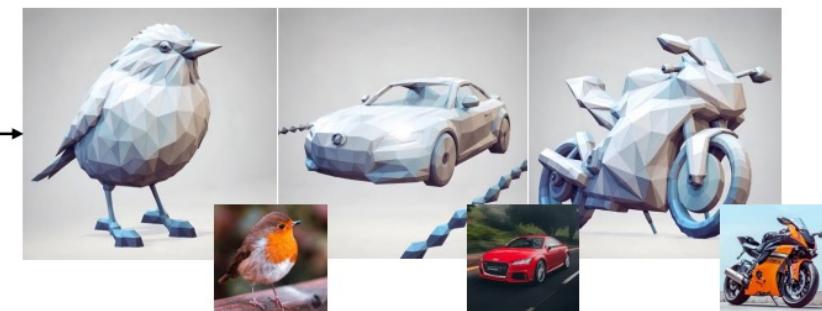
+ “low-poly style {object}. low-poly game art, polygon mesh, jagged, blocky, wireframe edges, centered composition, simple ..”

(b) Various style results with a highly detailed text description

Reference



(c) T2I results specified by a style reference

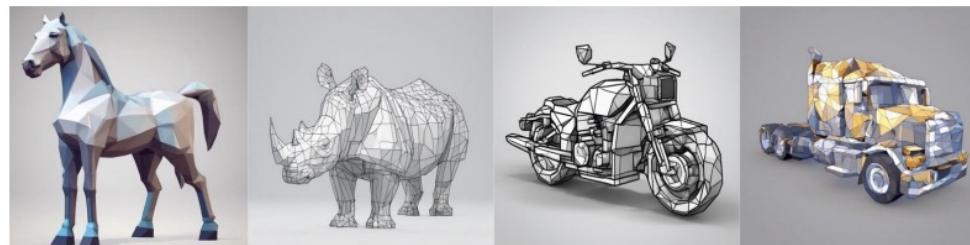


(d) Style transfer results specified by a style reference



“low-poly style horse”

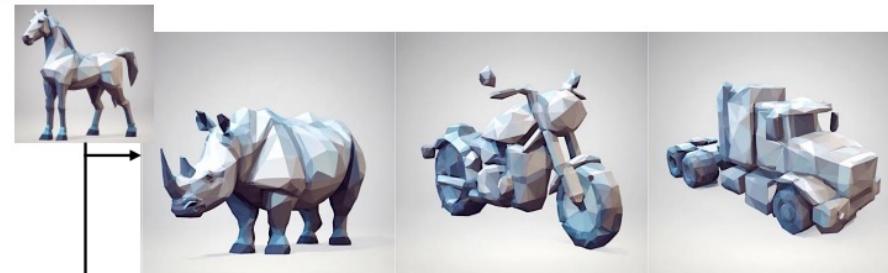
(a) Various style results within a text description



“Horse” “Rhino” “Motorcycle” “Truck”
+ “low-poly style {object}. low-poly game art, polygon mesh, jagged, blocky,
wireframe edges, centered composition, simple ..”

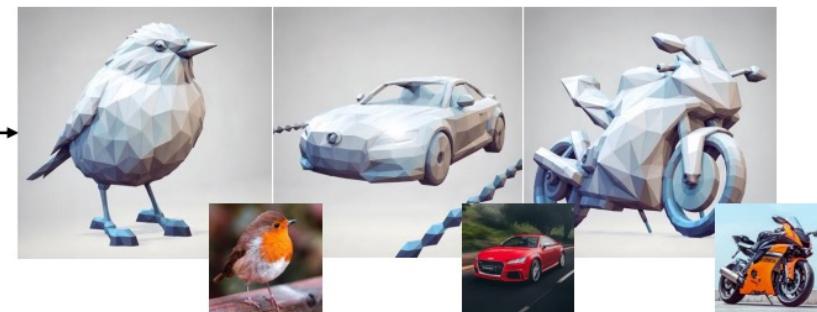
(b) Various style results with a highly detailed text description

Reference



“Rhino” “Motorcycle” “Truck”

(c) T2I results specified by a style reference

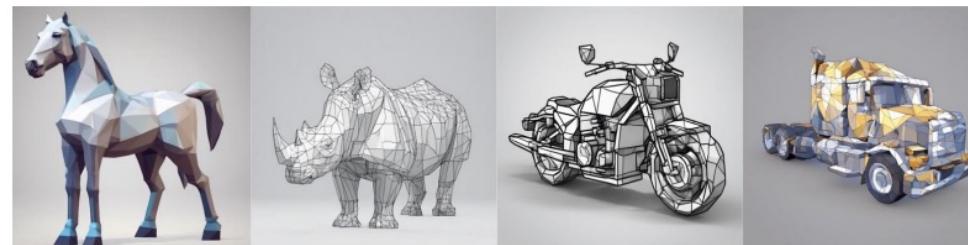


(d) Style transfer results specified by a style reference



“low-poly style horse”

(a) Various style results within a text description



“Horse” “Rhino” “Motorcycle” “Truck”
+ “low-poly style {object}. low-poly game art, polygon mesh, jagged, blocky,
wireframe edges, centered composition, simple ..”

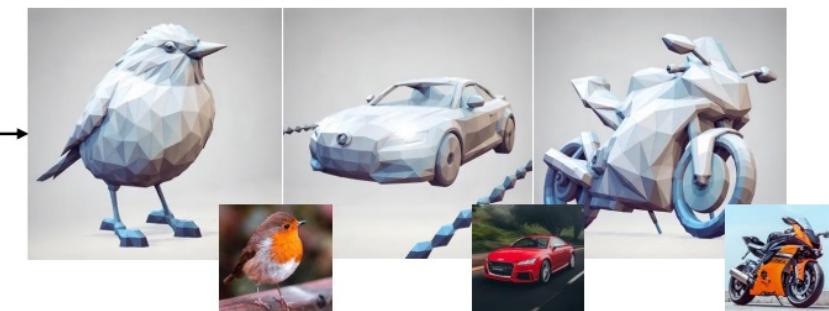
(b) Various style results with a highly detailed text description

Reference



“Rhino” “Motorcycle” “Truck”

(c) T2I results specified by a style reference

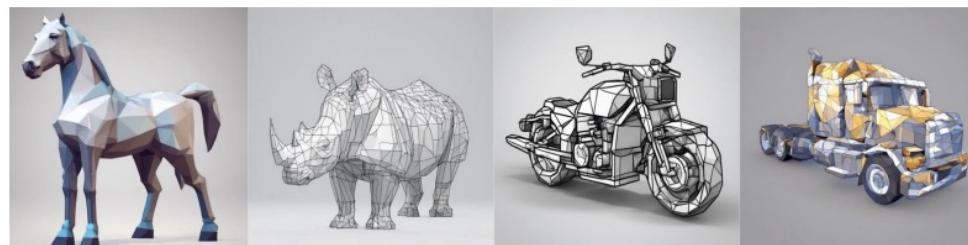


(d) Style transfer results specified by a style reference



“low-poly style horse”

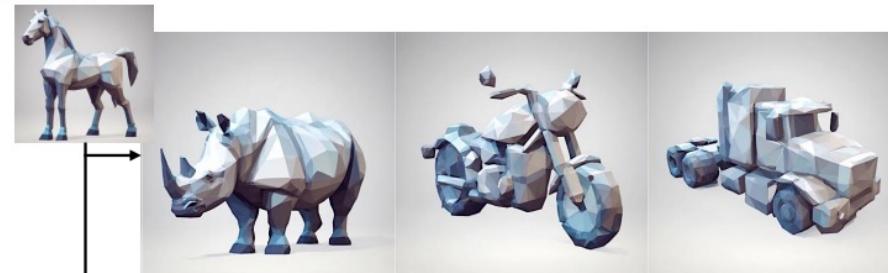
(a) Various style results within a text description



“Horse” “Rhino” “Motorcycle” “Truck”
+ “low-poly style {object}. low-poly game art, polygon mesh, jagged, blocky,
wireframe edges, centered composition, simple ..”

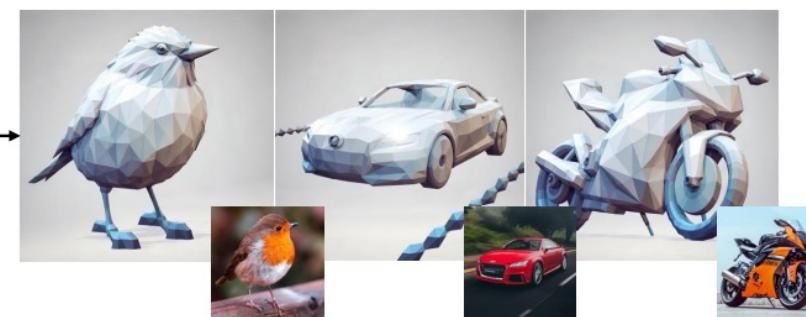
(b) Various style results with a highly detailed text description

Reference



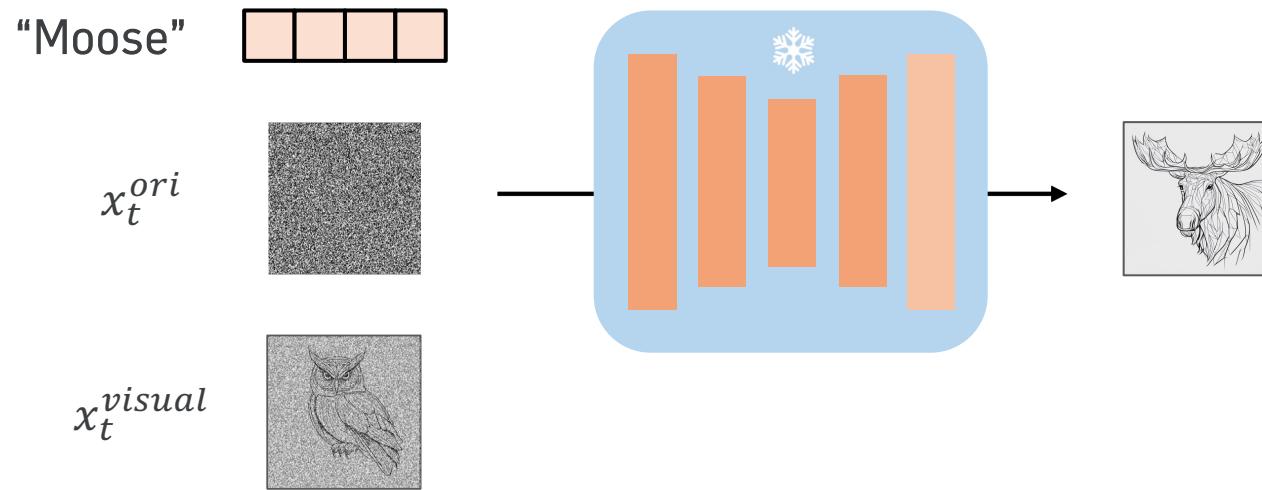
“Rhino” “Motorcycle” “Truck”

(c) T2I results specified by a style reference

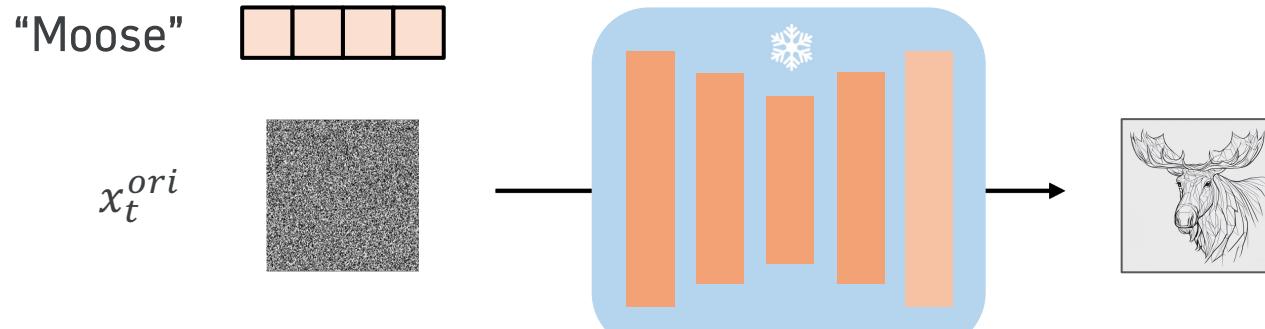


(d) Style transfer results specified by a style reference

Overview



Overview

 x_t^{visual} 

stochastic encoding

For mapping the latent.

Overview

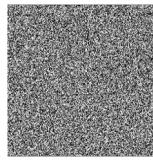
For reflecting the reference image's style.

swapping
self-attention

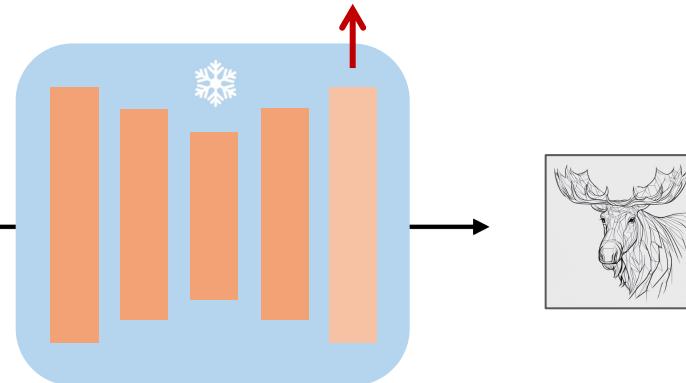
"Moose"



x_t^{ori}



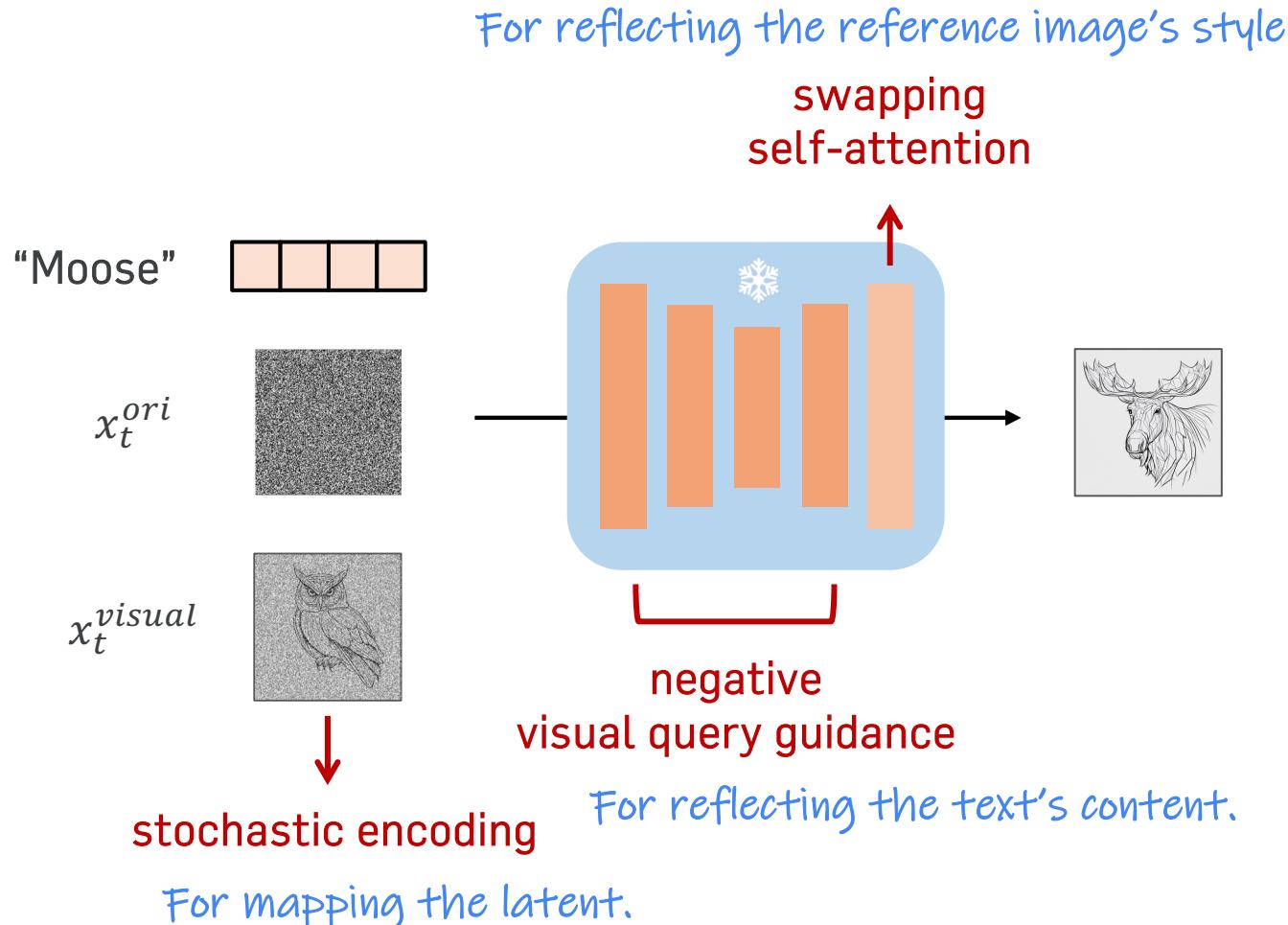
x_t^{visual}



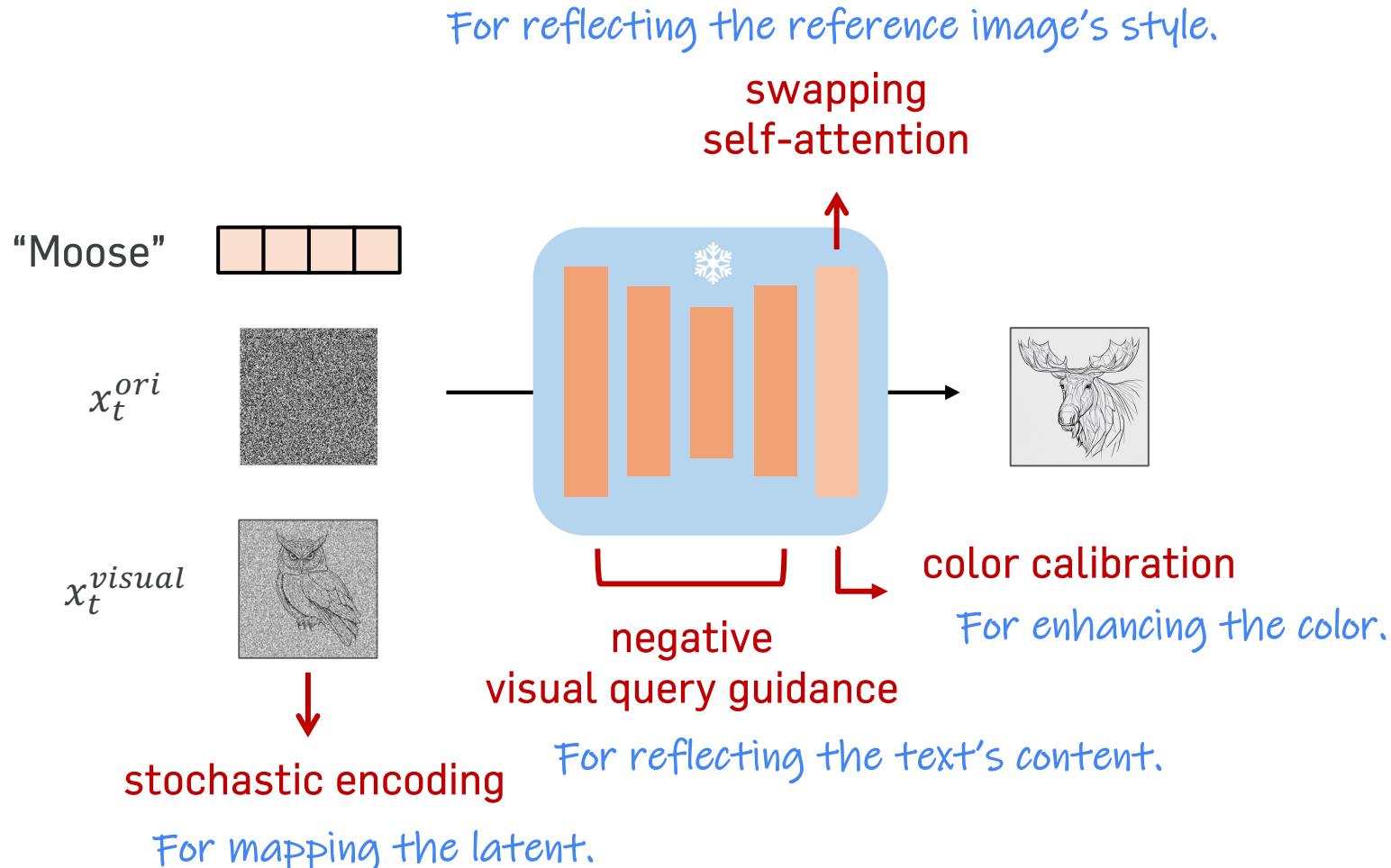
stochastic encoding

For mapping the latent.

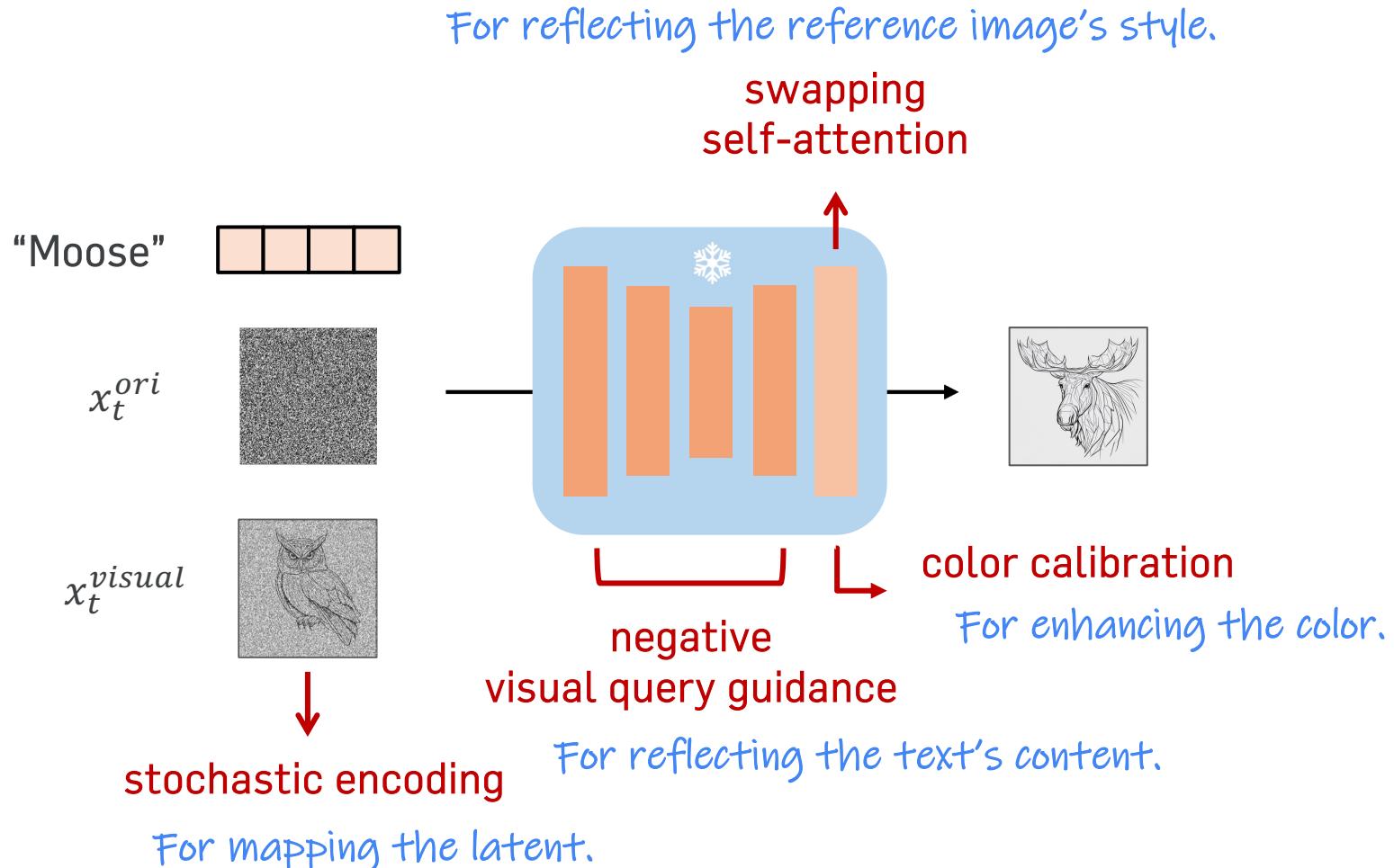
Overview



Overview



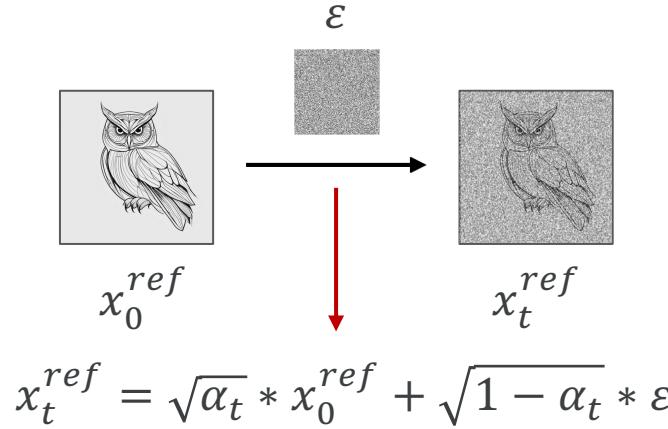
Overview



Summary

- Training-free
- Versatile
- SoTA
- Efficient & Reliable

Stochastic encoding



For mapping the latent.

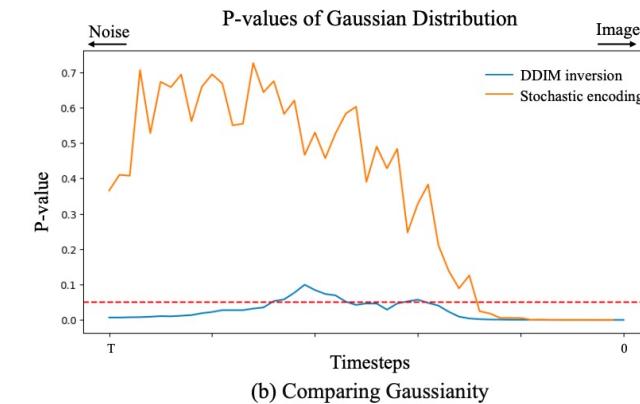
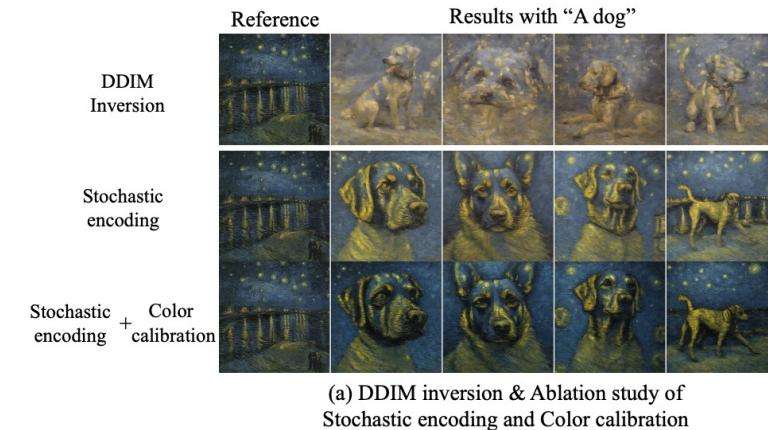
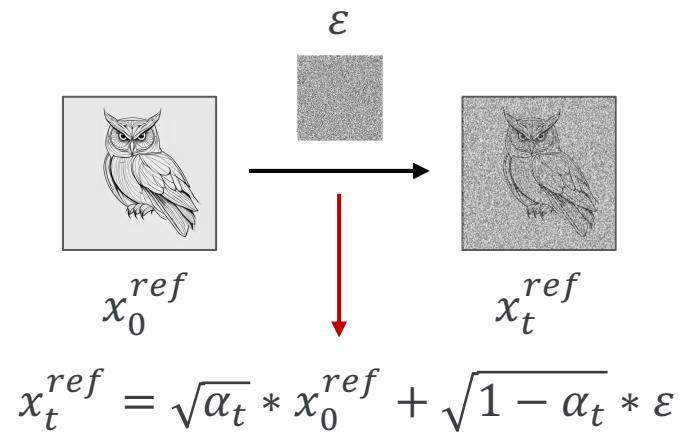


Figure 12: (a) Comparison of DDIM inversion vs. stochastic encoding and the effect of color calibration. Stochastic encoding reduces artifacts in the resulting images, while color calibration better reflects the colors of the reference image. (b) Stochastic encoding produces the latents closer to the standard Gaussian distribution compared to DDIM inversion. A P-value above 0.05 suggests that the data likely follows the standard Gaussian distribution.

Stochastic encoding



For mapping the latent.

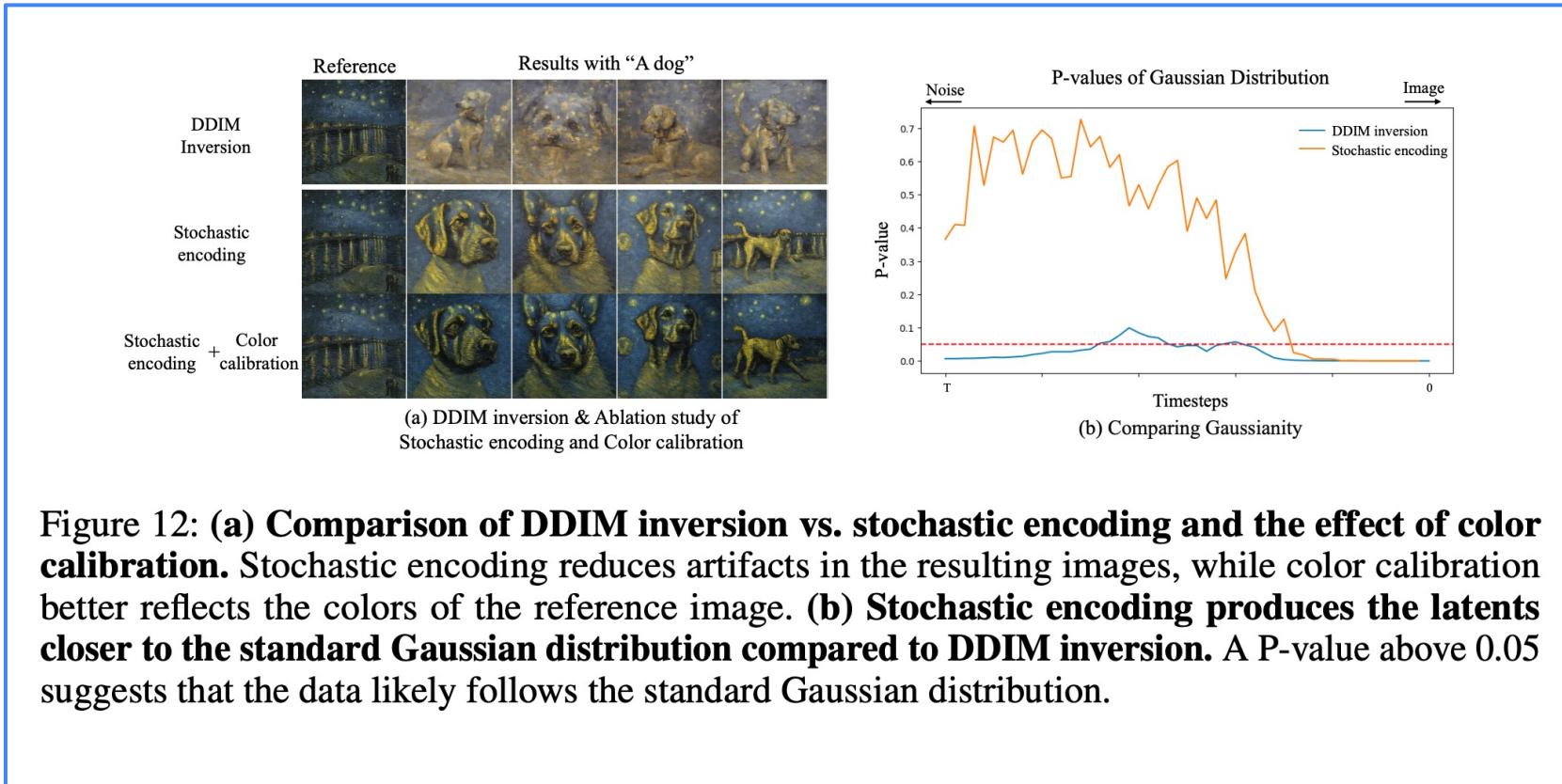
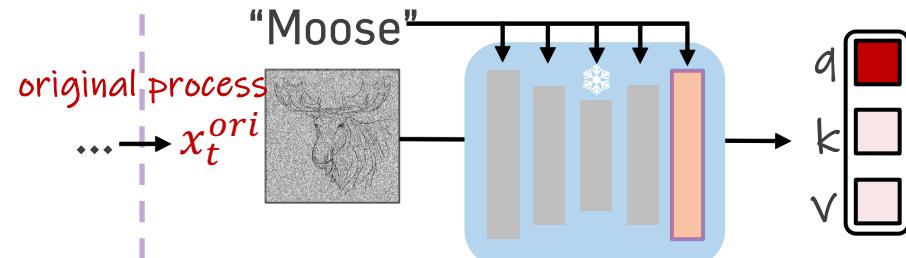


Figure 12: **(a) Comparison of DDIM inversion vs. stochastic encoding and the effect of color calibration.** Stochastic encoding reduces artifacts in the resulting images, while color calibration better reflects the colors of the reference image. **(b) Stochastic encoding produces the latents closer to the standard Gaussian distribution compared to DDIM inversion.** A P-value above 0.05 suggests that the data likely follows the standard Gaussian distribution.

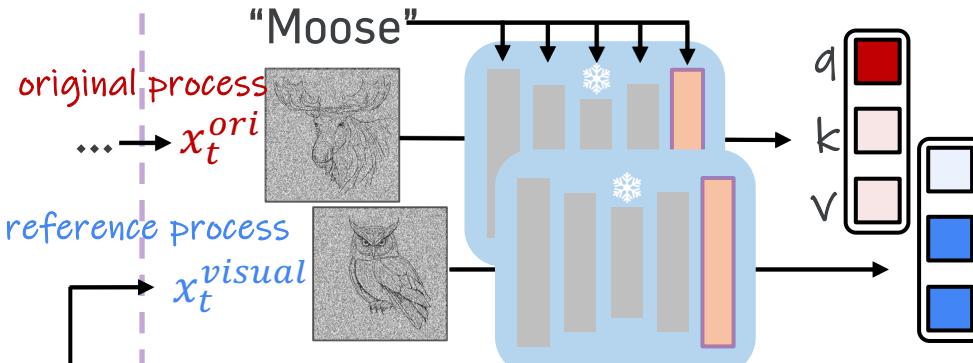
* q, k, v = query, key, value

For $\epsilon_\theta(x_t, c)$: key, value swapping in self-attn



* q, k, v = query, key, value

For $\epsilon_\theta(x_t, c)$: key, value swapping in self-attn



Sec 3
stochastic encoding



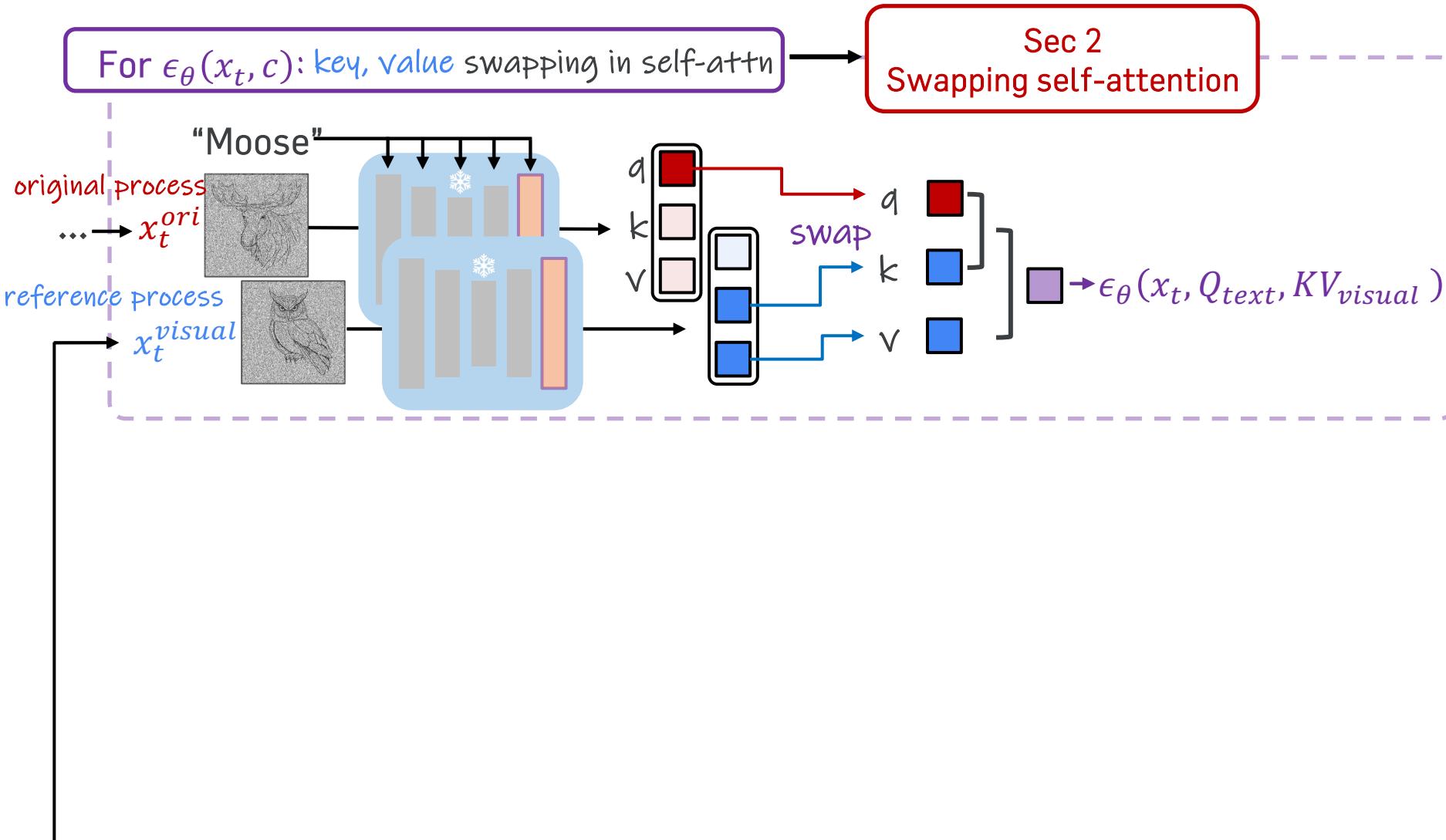
For mapping the latent

Method

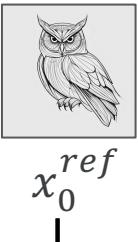
Guide in self-attention

* $q, k, v = \text{query, key, value}$

For reflecting the reference image's style



Sec 3
stochastic encoding



Motivation:

$$\tilde{\epsilon}_\theta(x_t, c) = \epsilon_\theta(x_t, c) + \omega * (\epsilon_\theta(x_t, c) - \epsilon_\theta(x_t, \phi))$$

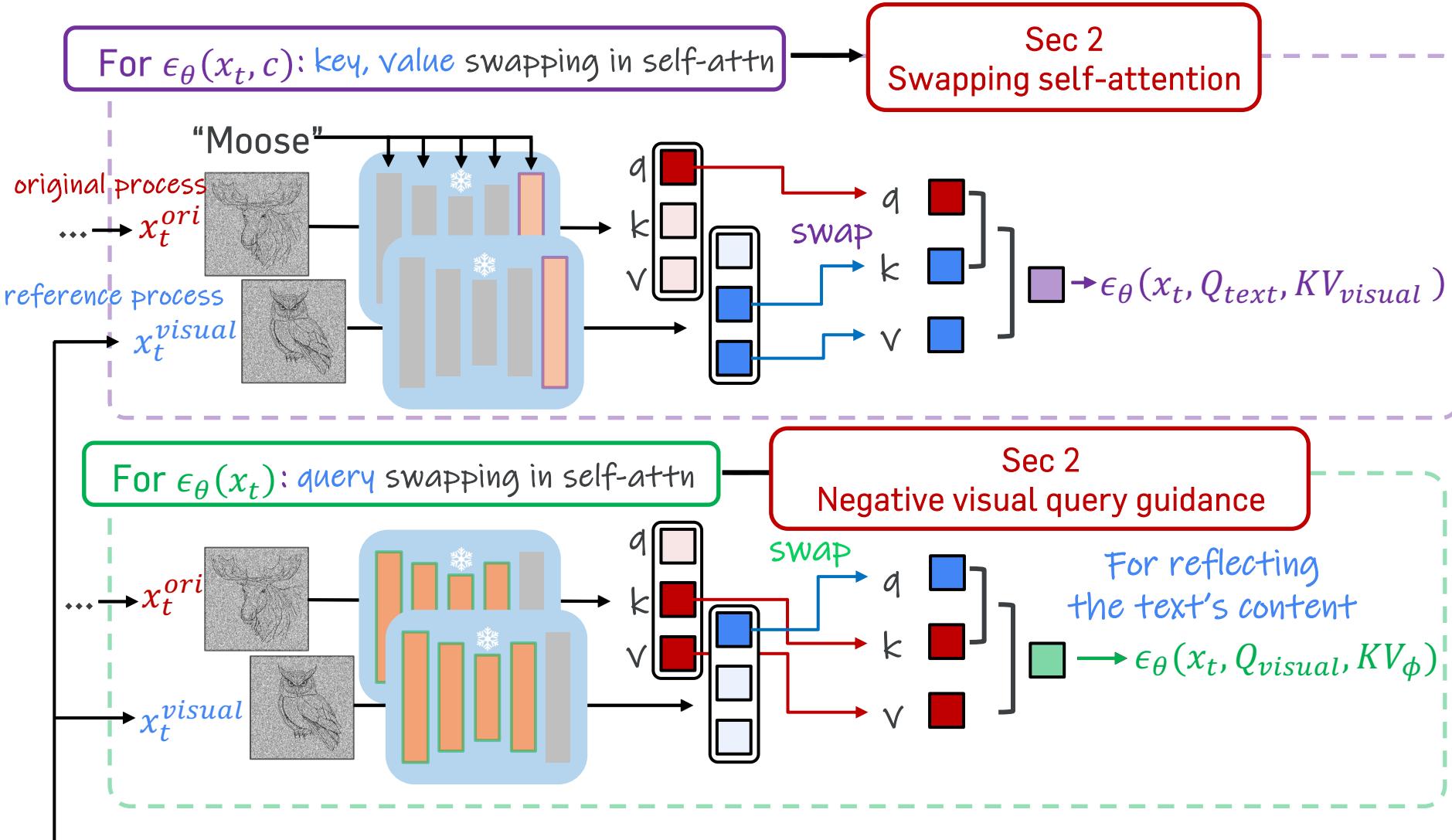
For mapping the latent

Method

Guide in self-attention

* $q, k, v = \text{query, key, value}$

For reflecting the reference image's style



Motivation:

$$\tilde{\epsilon}_\theta(x_t, c) = \epsilon_\theta(x_t, c) + \omega * (\epsilon_\theta(x_t, c) - \epsilon_\theta(x_t, \phi))$$

For mapping the latent

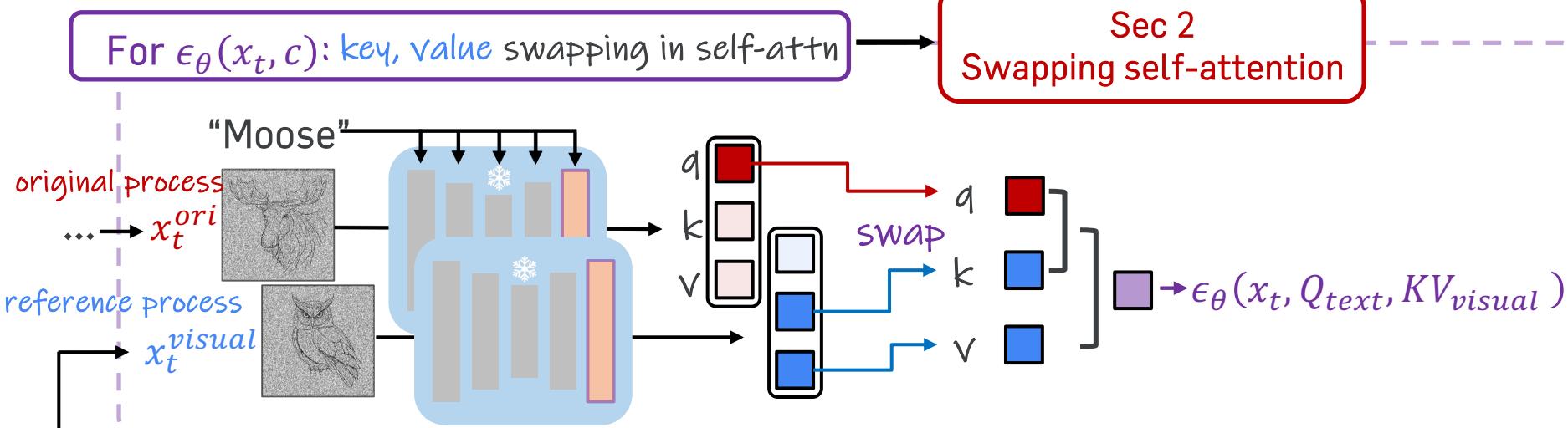


Method

Guide in self-attention

* q, k, v = query, key, value

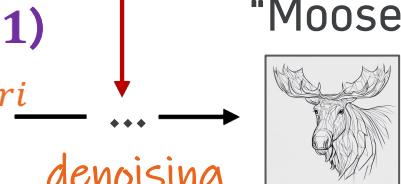
For reflecting the reference image's style



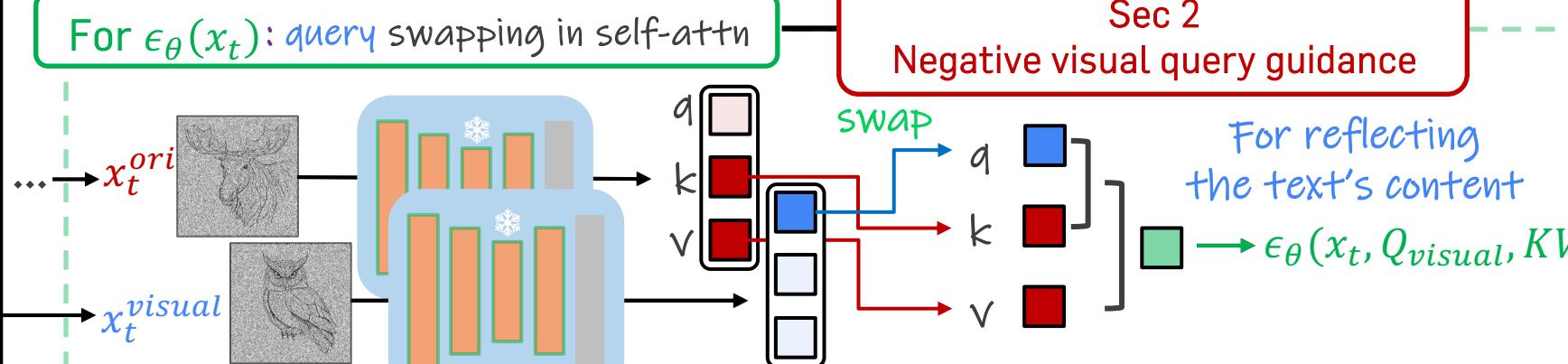
For enhancing the color



"Moose"



denoising



x_0^{visual}

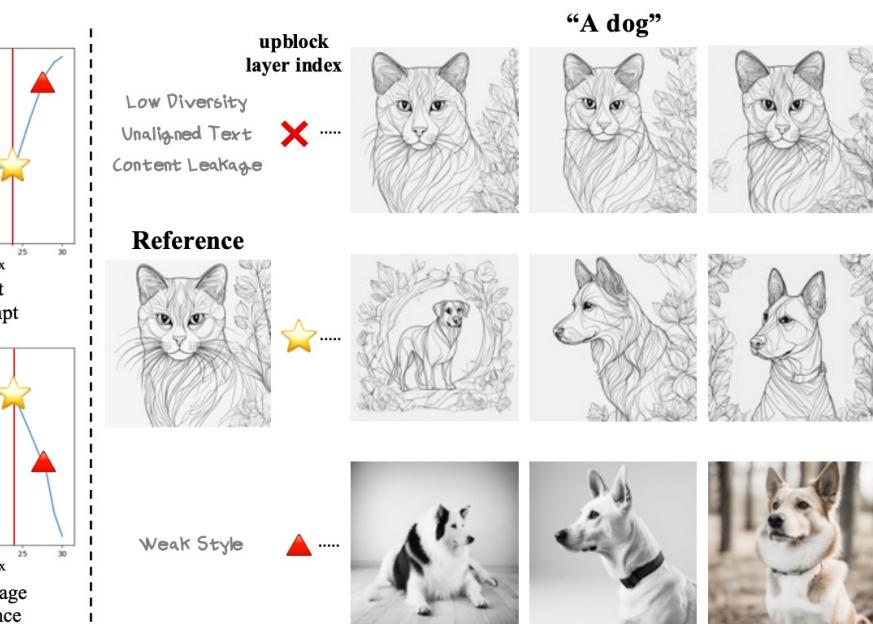
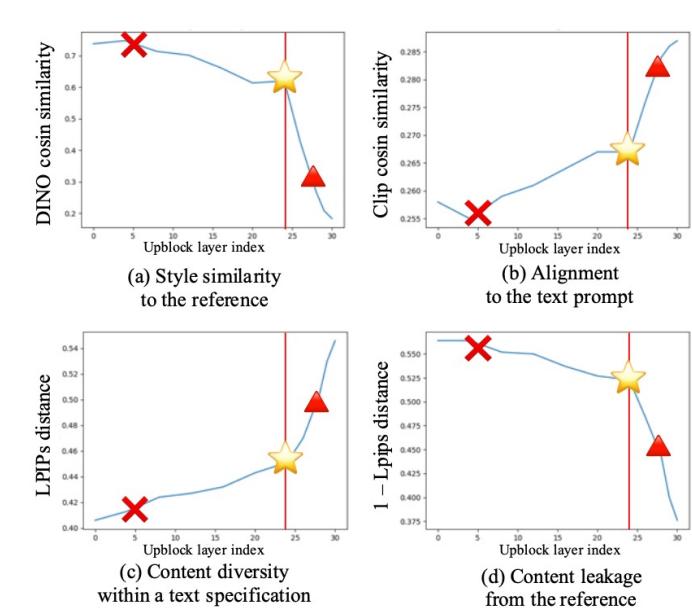
Motivation:

$$\tilde{\epsilon}_\theta(x_t, c) = \epsilon_\theta(x_t, c) + \omega * (\epsilon_\theta(x_t, c) - \epsilon_\theta(x_t, \phi))$$

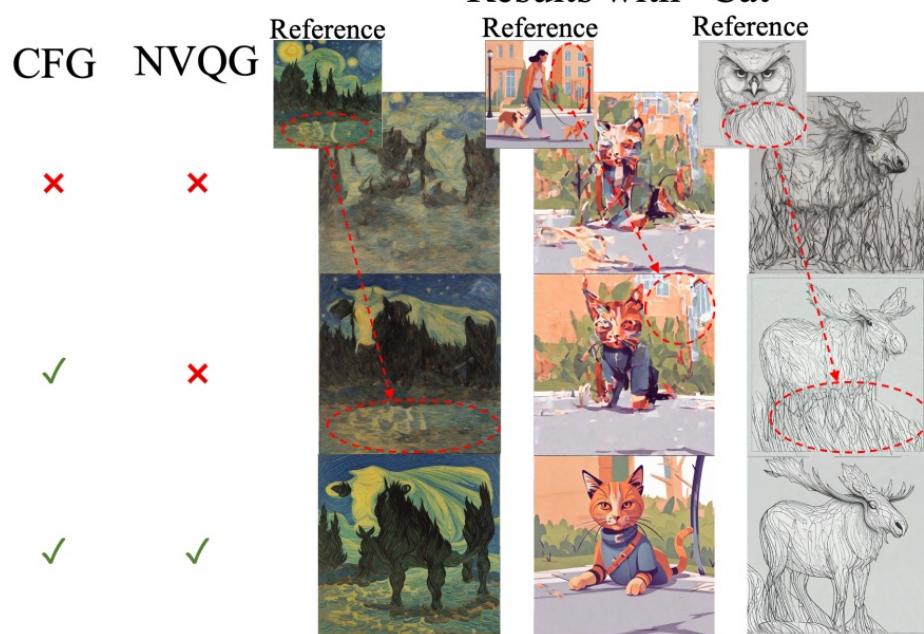
For mapping the latent

	CFG	NVQG
Reference		
“Cat”		
“Dog”		
“Cow”		
“Lion”		

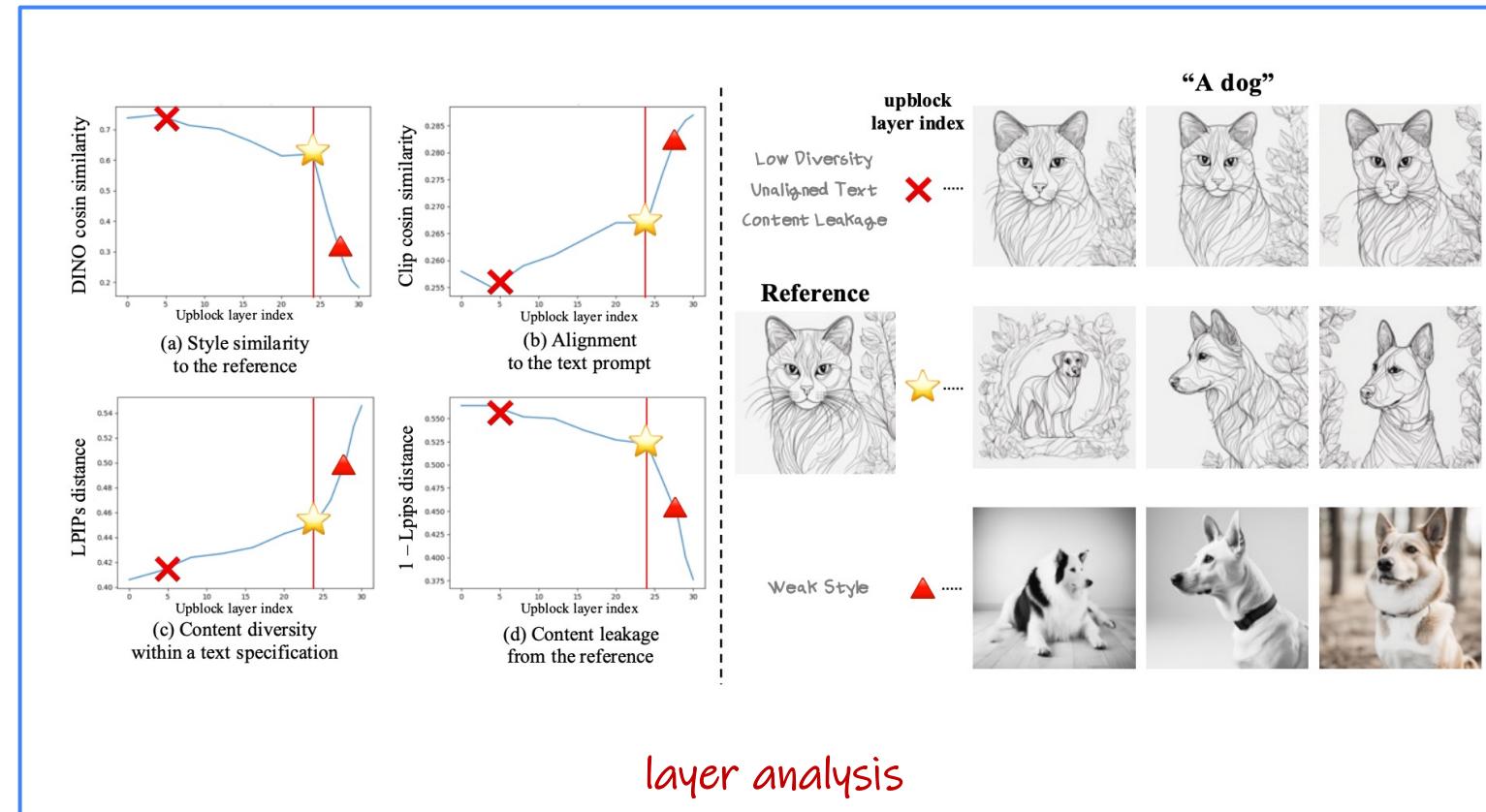
negative query visual guidance



layer analysis



negative query visual guidance



layer analysis

Function $\text{color_calibration}(x_t, \hat{x}_t, x_0^{\text{visual}})$:

$$x_{\text{pred}} \leftarrow \frac{x_t - \sqrt{1-\alpha_t} \cdot \epsilon_\theta(\hat{x}_t)}{\sqrt{\alpha_t}} ;$$
$$x_{\text{dir}} \leftarrow \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\hat{x}_t) ;$$
$$\epsilon \sim \mathcal{N}(0, I) ;$$
$$x_{\text{noise}} \leftarrow \sigma_t \cdot \epsilon ;$$
$$\hat{x}_{\text{pred}} \leftarrow \text{adain}(x_{\text{pred}}, x_0^{\text{visual}}) ;$$
$$x_{t-1} \leftarrow \sqrt{\alpha_{t-1}} \cdot \hat{x}_{\text{pred}} + x_{\text{dir}} + x_{\text{noise}} ;$$

return x_{t-1}

Color calibration

For enhancing the color

Function `color_calibration($x_t, \hat{x}_t, x_0^{visual}$)`:

$$x_{pred} \leftarrow \frac{x_t - \sqrt{1 - \alpha_t} \cdot \epsilon_\theta(\hat{x}_t)}{\sqrt{\alpha_t}} ;$$
$$x_{dir} \leftarrow \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\hat{x}_t) ;$$
$$\epsilon \sim \mathcal{N}(0, I) ;$$
$$x_{noise} \leftarrow \sigma_t \cdot \epsilon ;$$
$$\hat{x}_{pred} \leftarrow \text{adain}(x_{pred}, x_0^{visual}) ;$$
$$x_{t-1} \leftarrow \sqrt{\alpha_{t-1}} \cdot \hat{x}_{pred} + x_{dir} + x_{noise} ;$$

return x_{t-1}

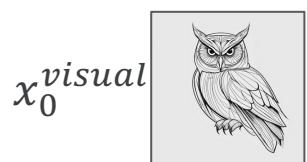
Color calibration

 x_0^{visual} 

For enhancing the color

Function `color_calibration($x_t, \hat{x}_t, x_0^{visual}$)`:
$$\begin{aligned}x_{pred} &\leftarrow \frac{x_t - \sqrt{1 - \alpha_t} \cdot \epsilon_\theta(\hat{x}_t)}{\sqrt{\alpha_t}} ; \\x_{dir} &\leftarrow \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\hat{x}_t) ; \\\epsilon &\sim \mathcal{N}(0, I) ; \\x_{noise} &\leftarrow \sigma_t \cdot \epsilon ; \\\hat{x}_{pred} &\leftarrow \text{adain}(x_{pred}, x_0^{visual}) ; \\x_{t-1} &\leftarrow \sqrt{\alpha_{t-1}} \cdot \hat{x}_{pred} + x_{dir} + x_{noise} ;\end{aligned}$$
return x_{t-1}

Color calibration



For enhancing the color

Function `color_calibration($x_t, \hat{x}_t, x_0^{\text{visual}}$)`:

```
 $x_{\text{pred}} \leftarrow \frac{x_t - \sqrt{1 - \alpha_t} \cdot \epsilon_\theta(\hat{x}_t)}{\sqrt{\alpha_t}} ;$ 
 $x_{\text{dir}} \leftarrow \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\hat{x}_t) ;$ 
 $\epsilon \sim \mathcal{N}(0, I) ;$ 
 $x_{\text{noise}} \leftarrow \sigma_t \cdot \epsilon ;$ 
 $\hat{x}_{\text{pred}} \leftarrow \text{adain}(x_{\text{pred}}, x_0^{\text{visual}}) ;$ 
 $x_{t-1} \leftarrow \sqrt{\alpha_{t-1}} \cdot \hat{x}_{\text{pred}} + x_{\text{dir}} + x_{\text{noise}} ;$ 
return  $x_{t-1}$ 
```

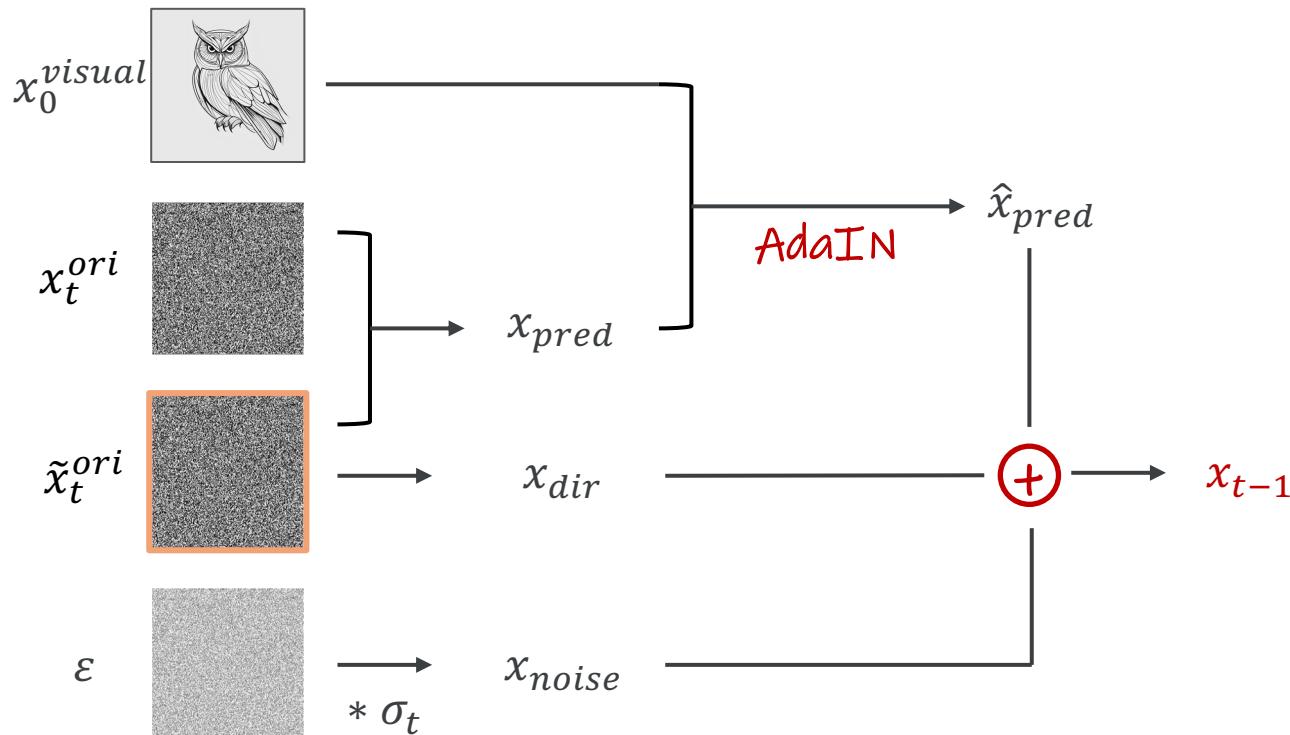
Color calibration



For enhancing the color

Function `color_calibration($x_t, \hat{x}_t, x_0^{visual}$)`:
$$\begin{aligned}x_{pred} &\leftarrow \frac{x_t - \sqrt{1-\alpha_t} \cdot \epsilon_\theta(\hat{x}_t)}{\sqrt{\alpha_t}} ; \\x_{dir} &\leftarrow \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\hat{x}_t) ; \\ \epsilon &\sim \mathcal{N}(0, I) ; \\x_{noise} &\leftarrow \sigma_t \cdot \epsilon ; \\ \hat{x}_{pred} &\leftarrow \text{adain}(x_{pred}, x_0^{visual}) ; \\x_{t-1} &\leftarrow \sqrt{\alpha_{t-1}} \cdot \hat{x}_{pred} + x_{dir} + x_{noise} ; \\ \text{return } &x_{t-1}\end{aligned}$$

Color calibration



For enhancing the color

Function `color_calibration($x_t, \hat{x}_t, x_0^{visual}$)`:

```

 $x_{pred} \leftarrow \frac{x_t - \sqrt{1-\alpha_t} \cdot \epsilon_\theta(\hat{x}_t)}{\sqrt{\alpha_t}} ;$ 
 $x_{dir} \leftarrow \sqrt{1 - \alpha_{t-1} - \sigma_t^2} \cdot \epsilon_\theta(\hat{x}_t) ;$ 
 $\epsilon \sim \mathcal{N}(0, I) ;$ 
 $x_{noise} \leftarrow \sigma_t \cdot \epsilon ;$ 
 $\hat{x}_{pred} \leftarrow \text{adain}(x_{pred}, x_0^{visual}) ;$ 
 $x_{t-1} \leftarrow \sqrt{\alpha_{t-1}} \cdot \hat{x}_{pred} + x_{dir} + x_{noise} ;$ 
return  $x_{t-1}$ 

```

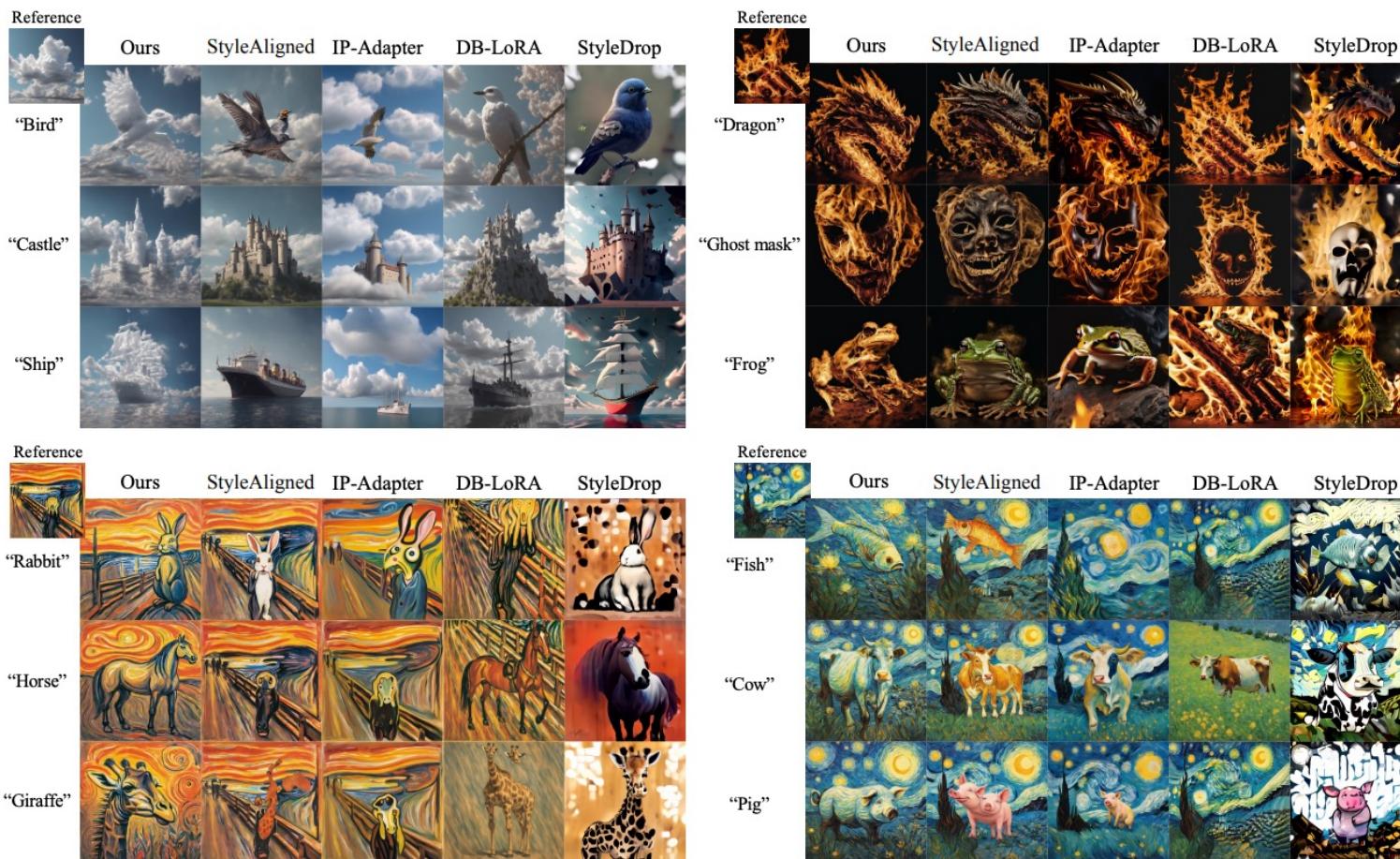
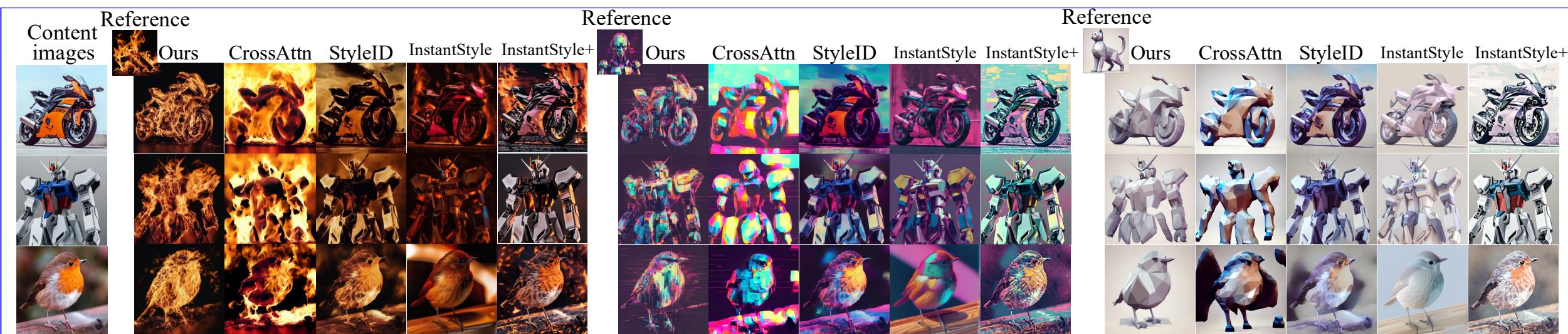


Figure 7: **Qualitative comparison across various styles and text prompts.** StyleGuide faithfully reflects style elements in reference images without causing content leakage from the reference images.

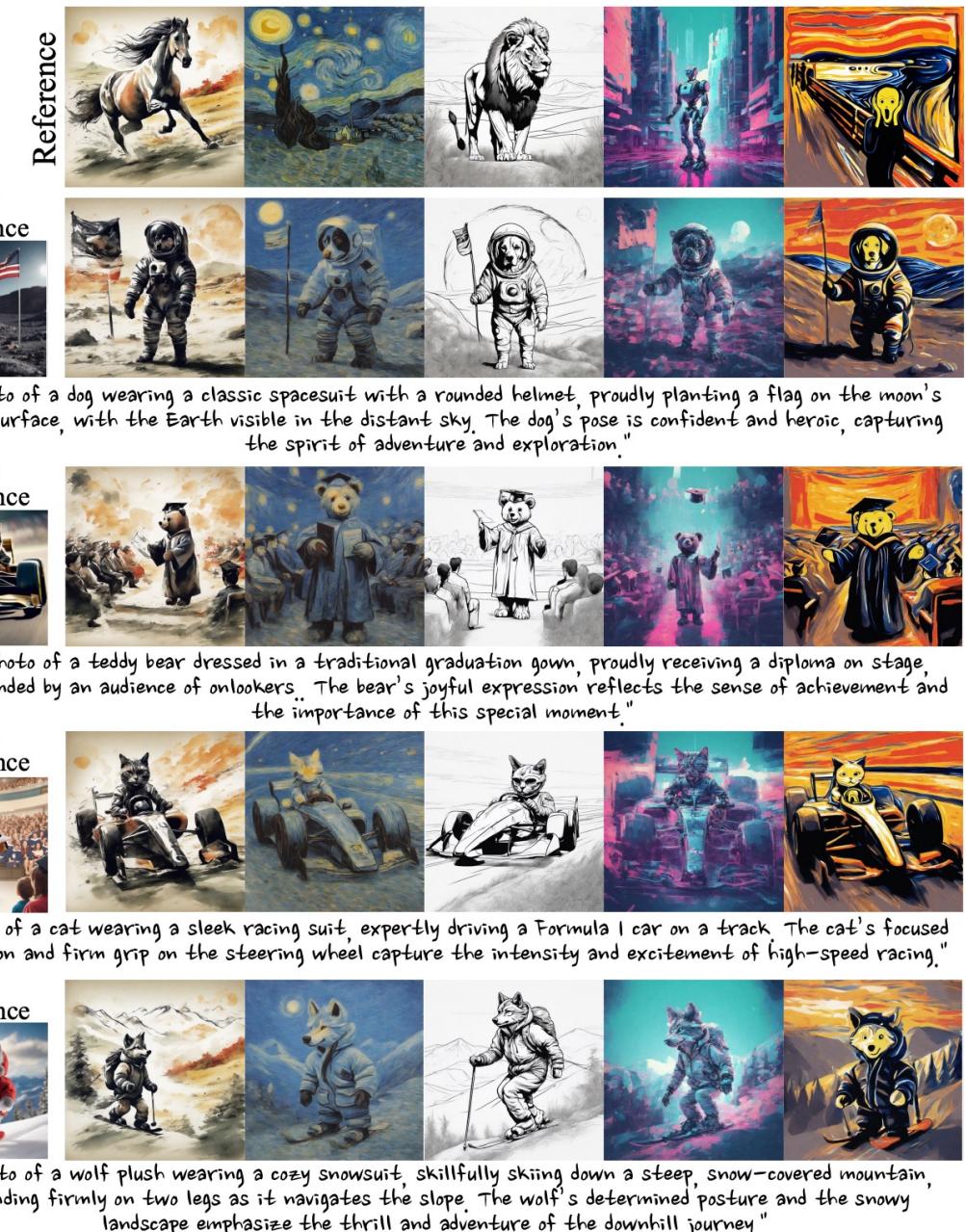
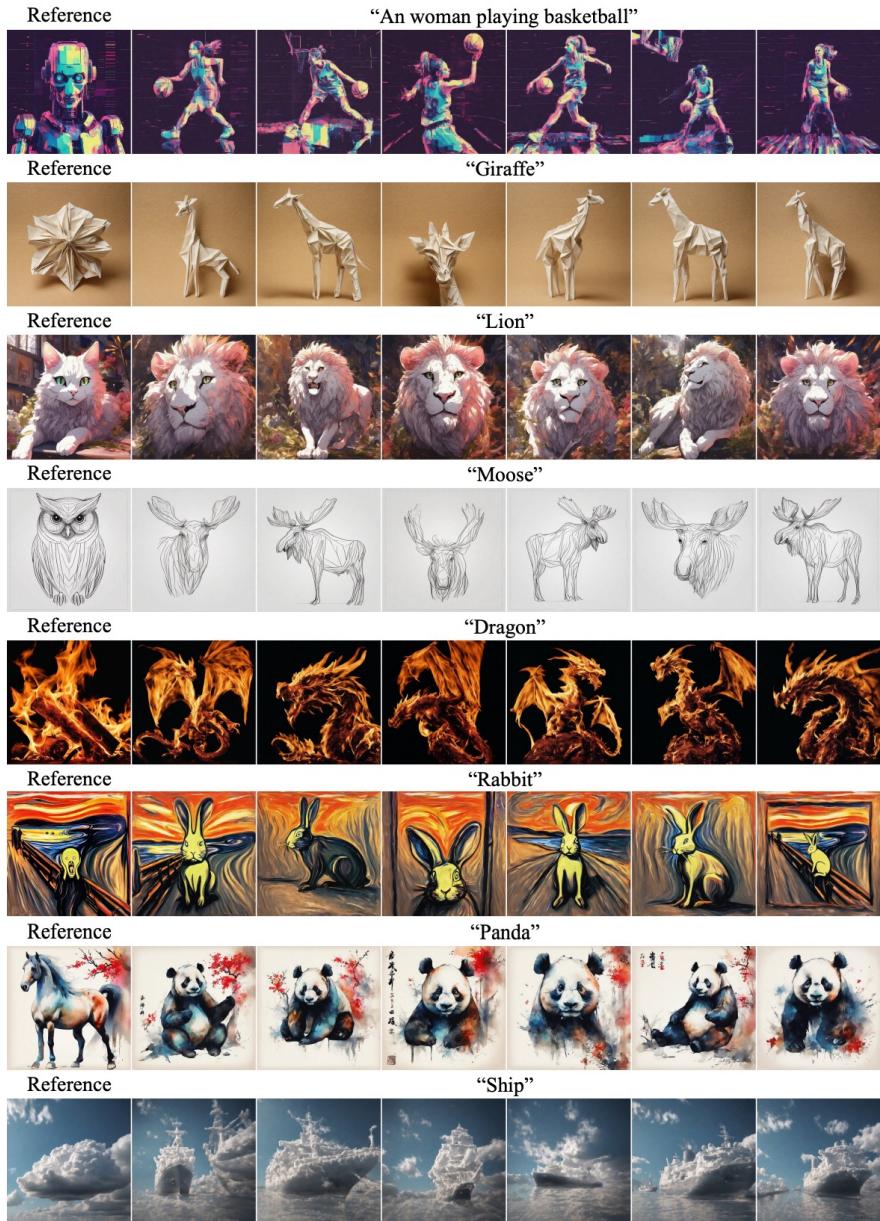
Results

Comparison (image-to-image)



Results

simple prompt & complex prompt





Junho Kim
(NAVER AI Lab)

Thank you for your attention !