

# High Fidelity Guided Image Synthesis using Latent Diffusion Models



Jaskirat Singh<sup>†</sup>



Stephen Gould<sup>†\*</sup>



Liang Zheng<sup>†\*</sup>



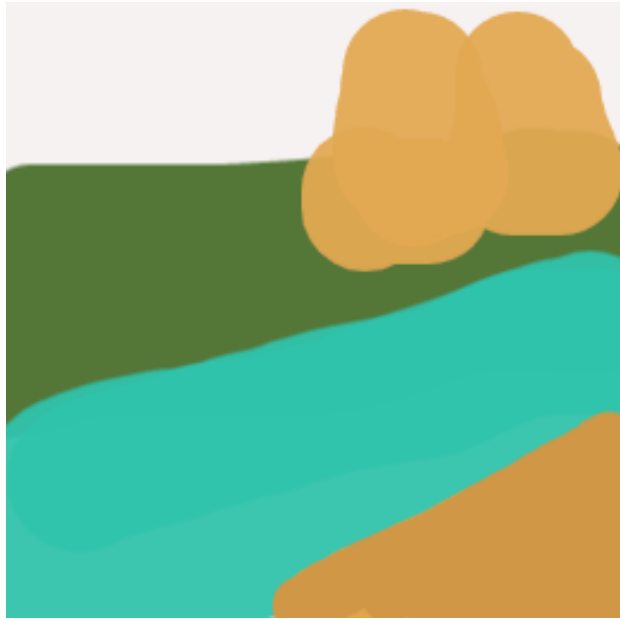
Australian  
National  
University



Poster Tag: TUE-PM-179

# Guided Image Synthesis with User Scribbles

Text Prompt: "a photo of a beautiful landscape"



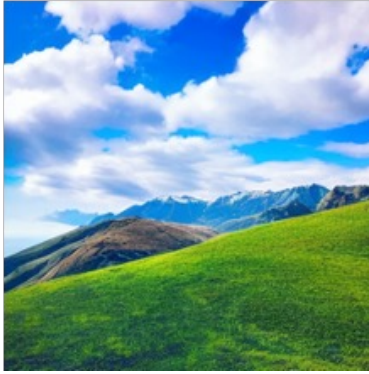
Reference Image  
Containing User Scribbles

Diffusion  
Model



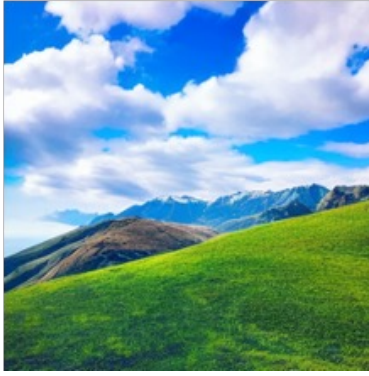
Guided Synthesis Output  
Conditioned on both text and reference input

Text Prompt: "a photo of a beautiful landscape"



*Text-Conditioned  
Image Outputs*

Text Prompt: "a photo of a beautiful landscape"

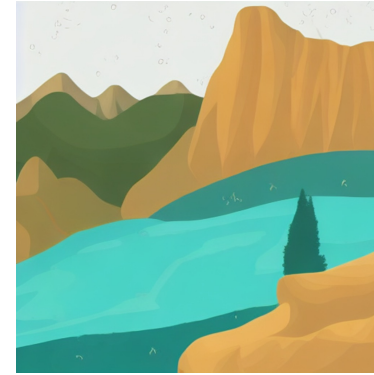
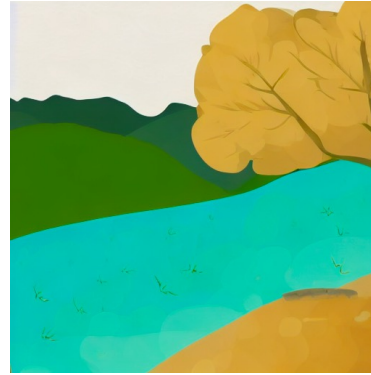


*Text-Conditioned  
Image Outputs*

Text Prompt:  
"a photo of a beautiful  
landscape"

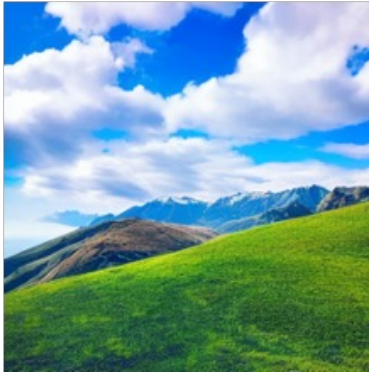


Reference Painting



*SDEdit*

Text Prompt: "a photo of a beautiful landscape"

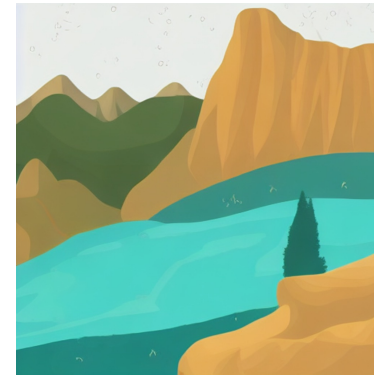
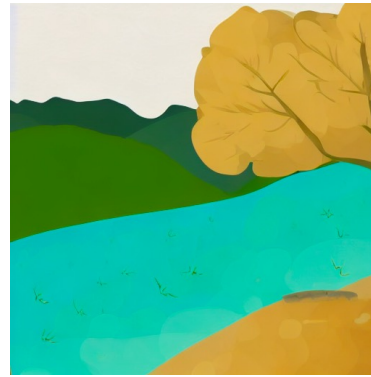


*Text-Conditioned  
Image Outputs*

Text Prompt:  
"a photo of a beautiful  
landscape"



Reference Painting



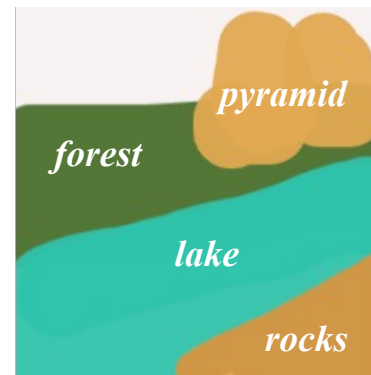
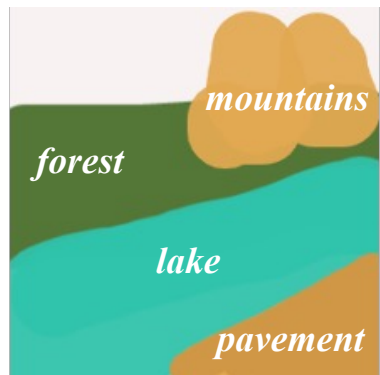
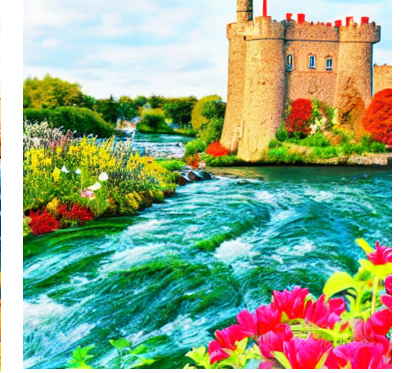
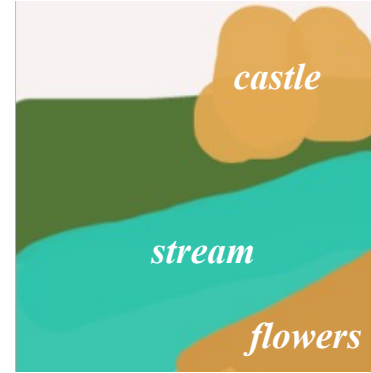
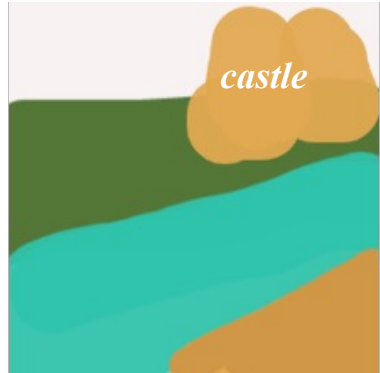
*SDEdit*



*GradOP  
(Ours)*

# Controlling Semantics of Different Painting Regions

Text Prompt: "a photo of a beautiful landscape"



# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"

Reference:  $y$



# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"

Reference:  $y$



Target Subspace  
*conditioned only on text*



$S_{\tau_{text}}$



# GradOP: Constrained Optimization Solution

Text Prompt: "a **photo** of a red tiger in a green field"

Reference:  $y$



$$\begin{aligned} x^* &= \operatorname{argmin}_x \mathcal{L}(f(x), y) \\ \text{s.t. } x &\in \mathcal{S}_{\tau_{\text{text}}} \end{aligned}$$

Constrained Optimization  
Formulation

Target Subspace  
*conditioned only on text*



# GradOP: Constrained Optimization Solution

Text Prompt: "a *photo* of a red tiger in a green field"

Reference:  $y$



$$x^* = \operatorname{argmin}_x \mathcal{L}(f(x), y)$$
$$s.t. \quad x \in \mathcal{S}_{\tau_{text}}$$

Constrained Optimization  
Formulation

Target Subspace  
*conditioned only on text*

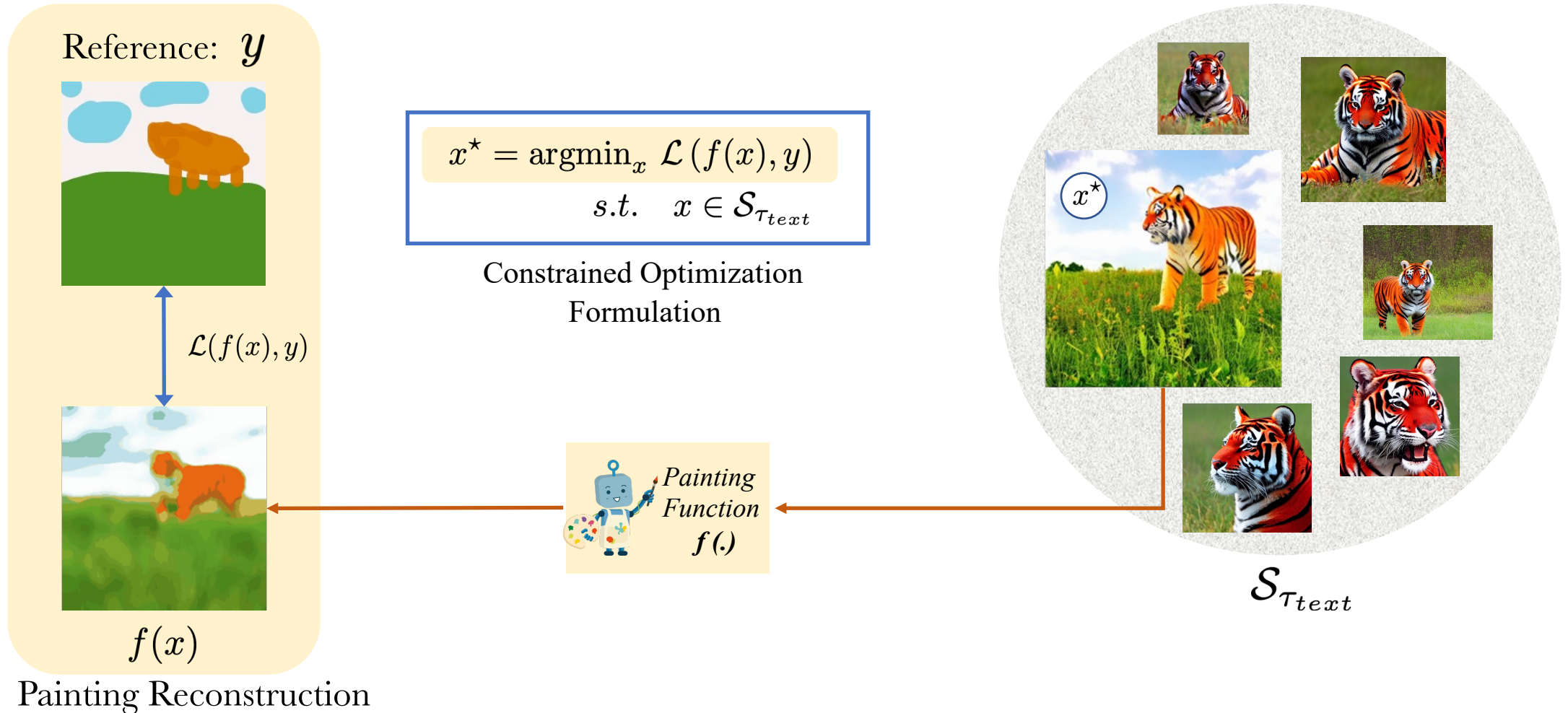


$\mathcal{S}_{\tau_{text}}$

# GradOP: Constrained Optimization Solution

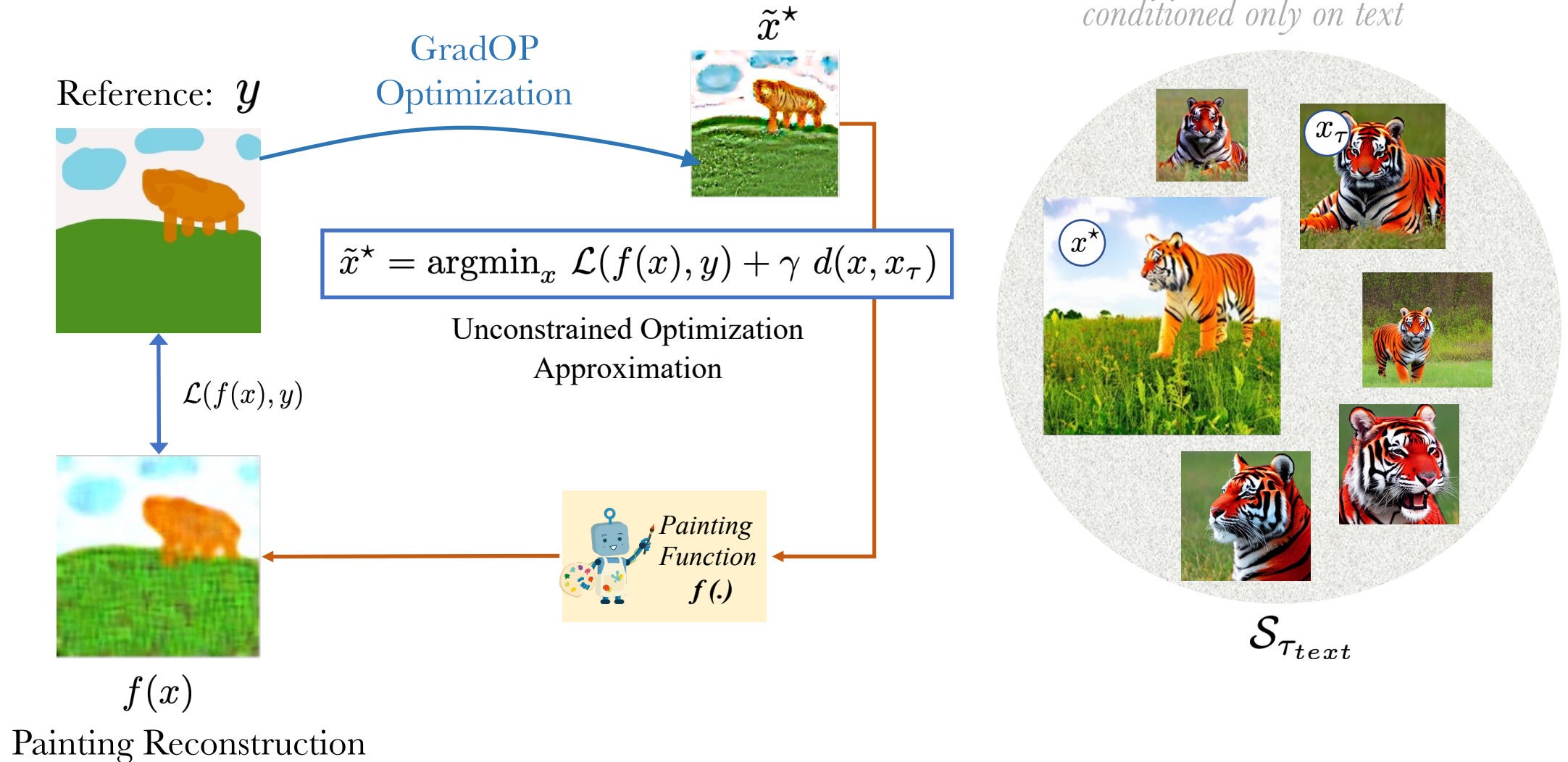
Text Prompt: "a photo of a red tiger in a green field"

Target Subspace  
*conditioned only on text*



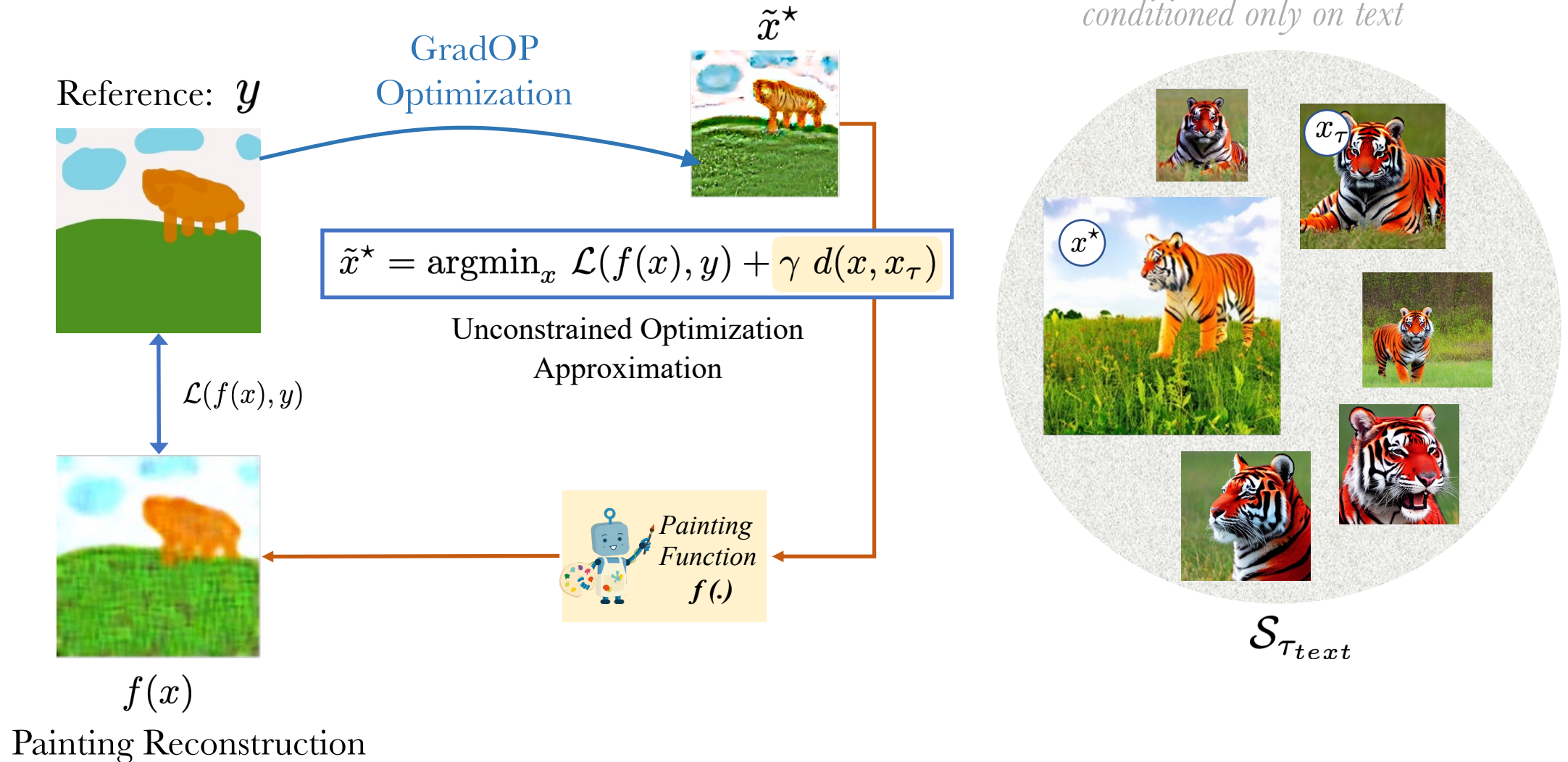
# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"



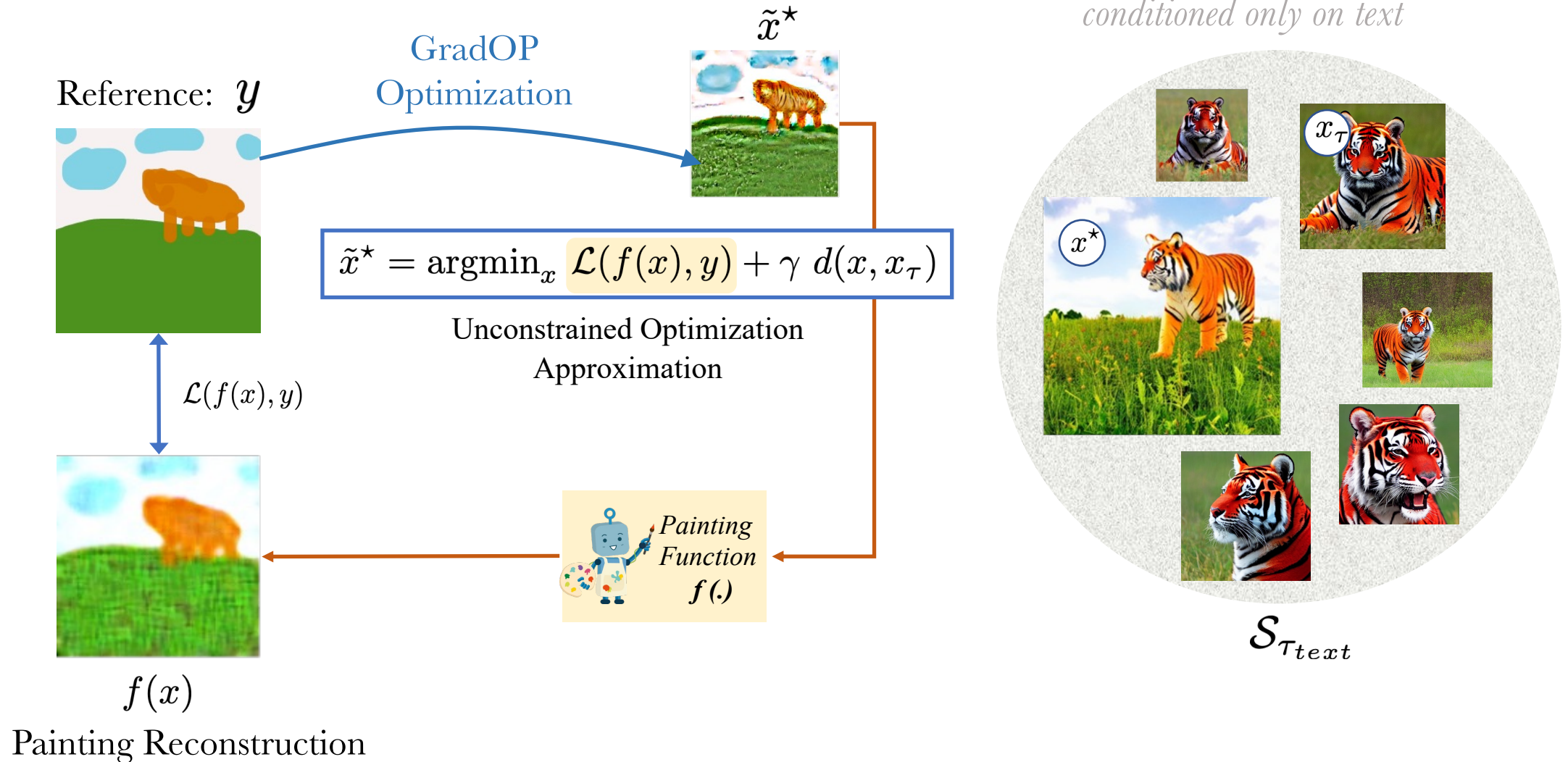
# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"



# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"



# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"

Reference:  $y$



GradOP  
Optimization

$\tilde{x}^*$



DDIM  
Inversion

$$\tilde{x}^* = \operatorname{argmin}_x \mathcal{L}(f(x), y) + \gamma d(x, x_\tau)$$

Unconstrained Optimization  
Approximation

Target Subspace  
*conditioned only on text*



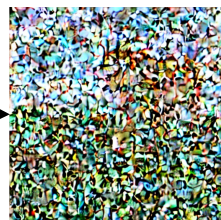
$S_{\tau_{text}}$

$\tilde{x}^*$

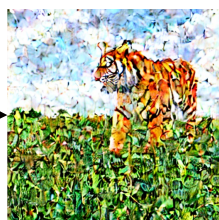


Forward  
Diffusion

$x_{t_0}$



$x_t$



...

$x^*$



Output Image

# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"

Reference:  $y$



Target Subspace  
*conditioned only on text*



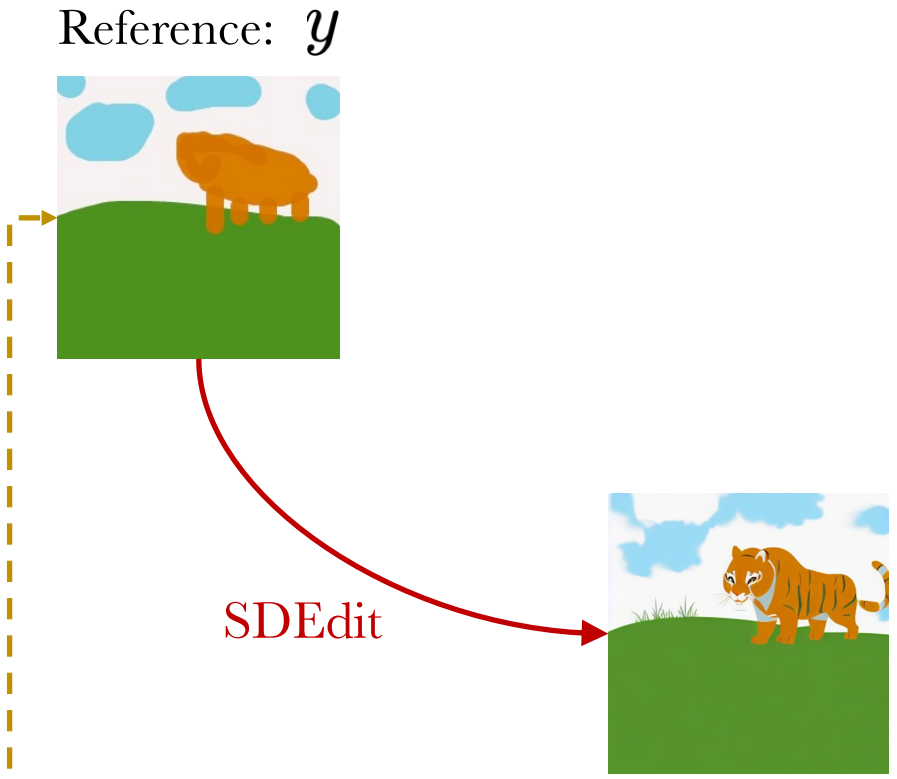
$S_{\tau_{text}}$

High Domain Gap



# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"



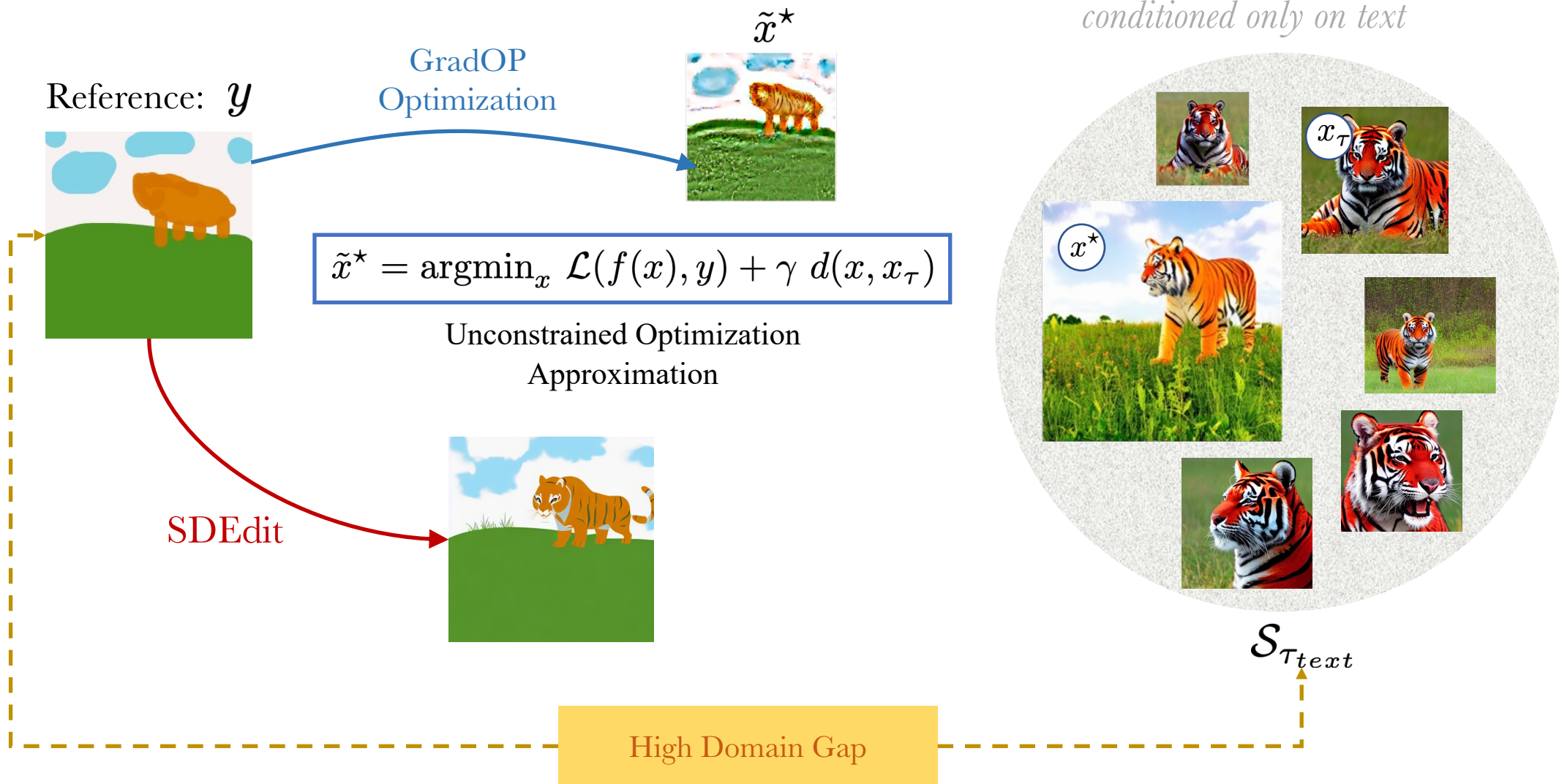
Target Subspace  
*conditioned only on text*



High Domain Gap

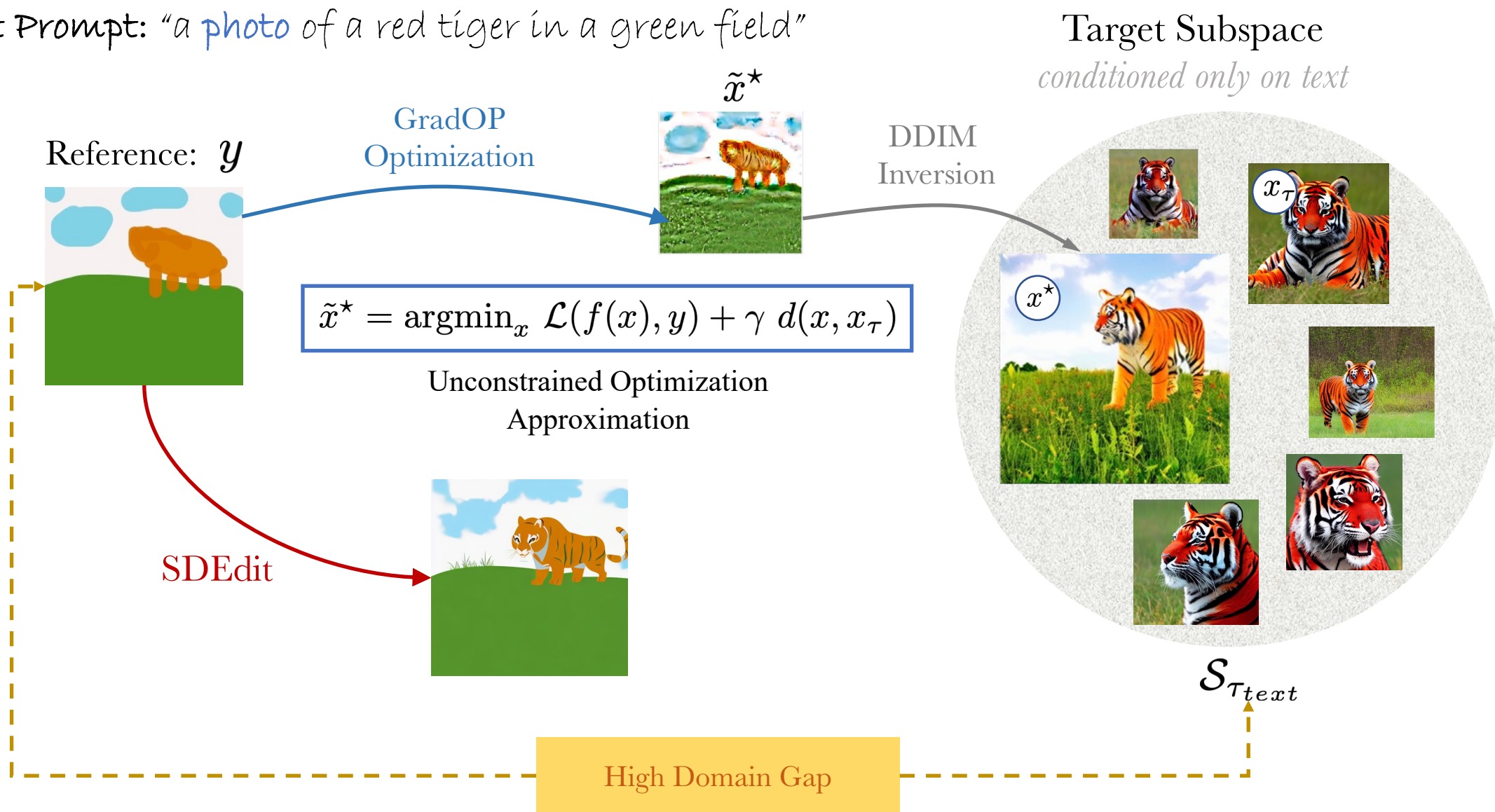
# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"



# GradOP: Constrained Optimization Solution

Text Prompt: "a photo of a red tiger in a green field"

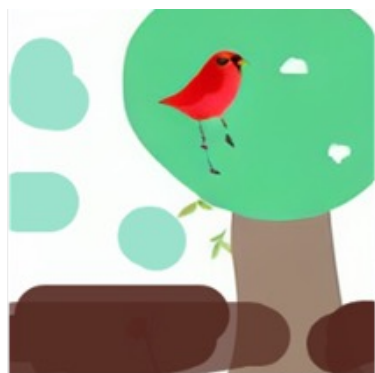


Text Prompt: "a photo of a red bird in a tree"

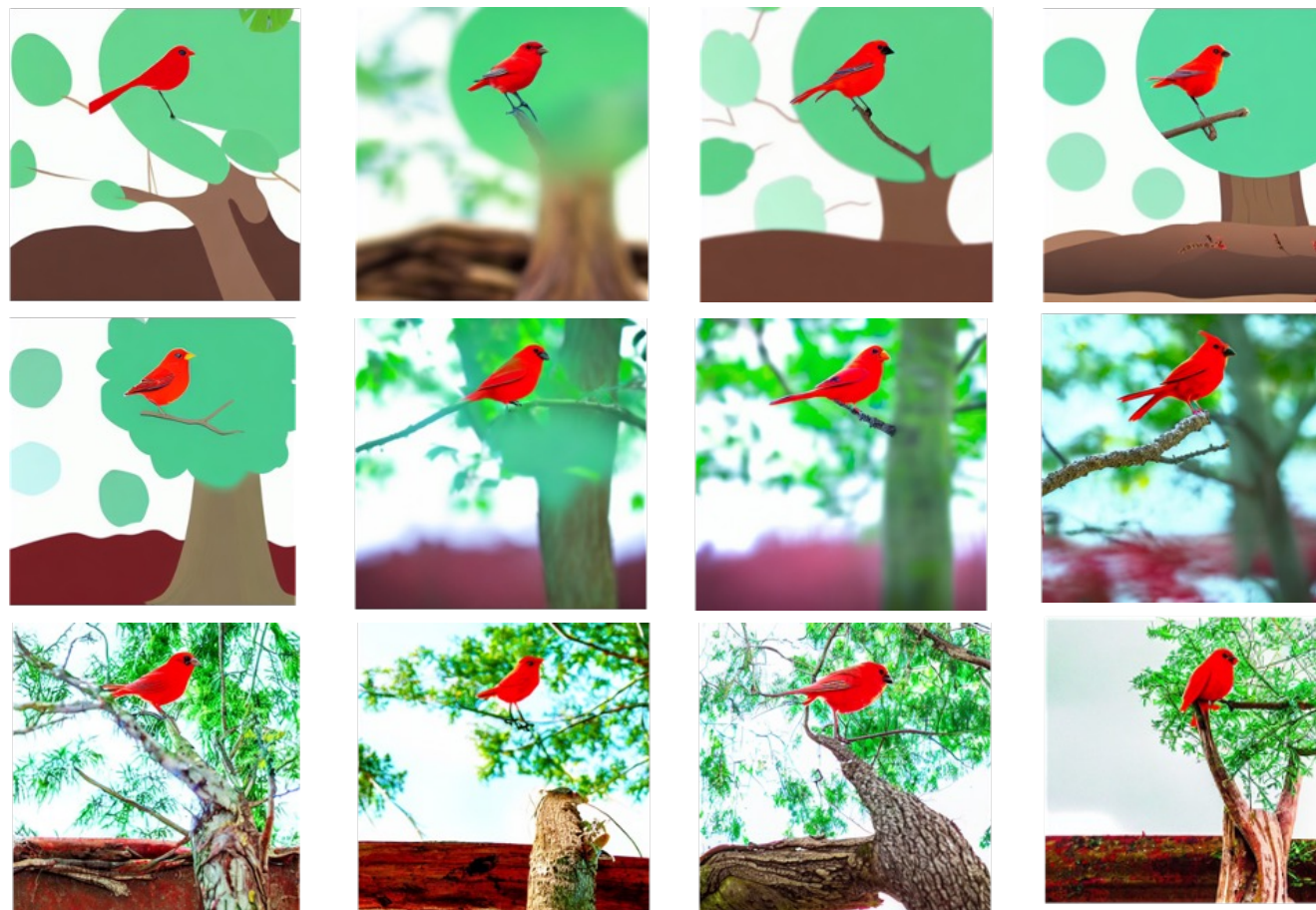


Target  
Subspace

Text Prompt:  
"a photo of a red bird in a tree"



Reference Painting



SDEdit

Iterative  
Loopback

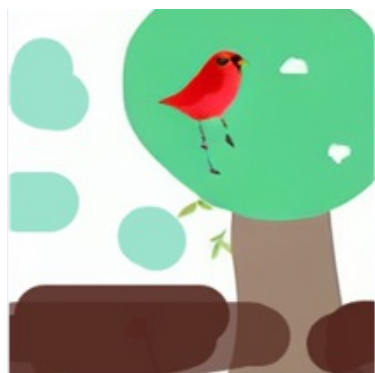
GradOP  
(Ours)

Text Prompt: "a photo of a red bird in a tree"

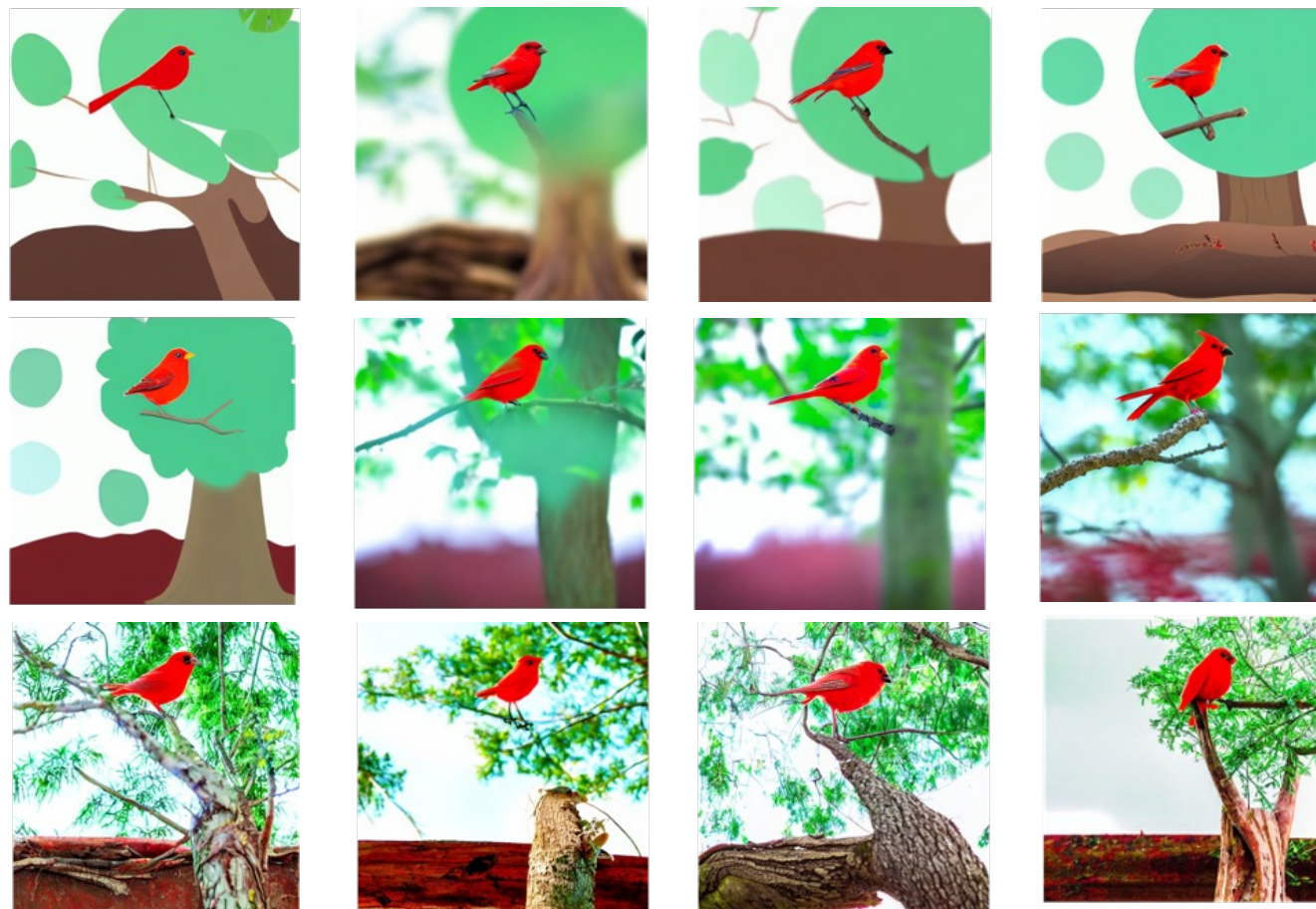


Target  
Subspace

Text Prompt:  
"a photo of a red bird in a tree"



Reference Painting



SDEdit

Iterative  
Loopback

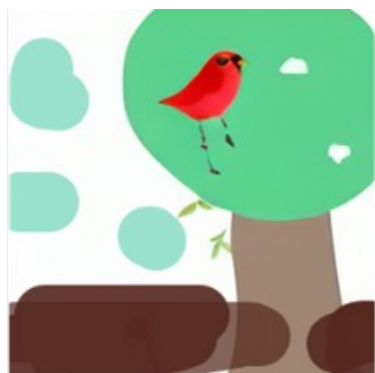
GradOP  
(Ours)

Text Prompt: "a photo of a red bird in a tree"

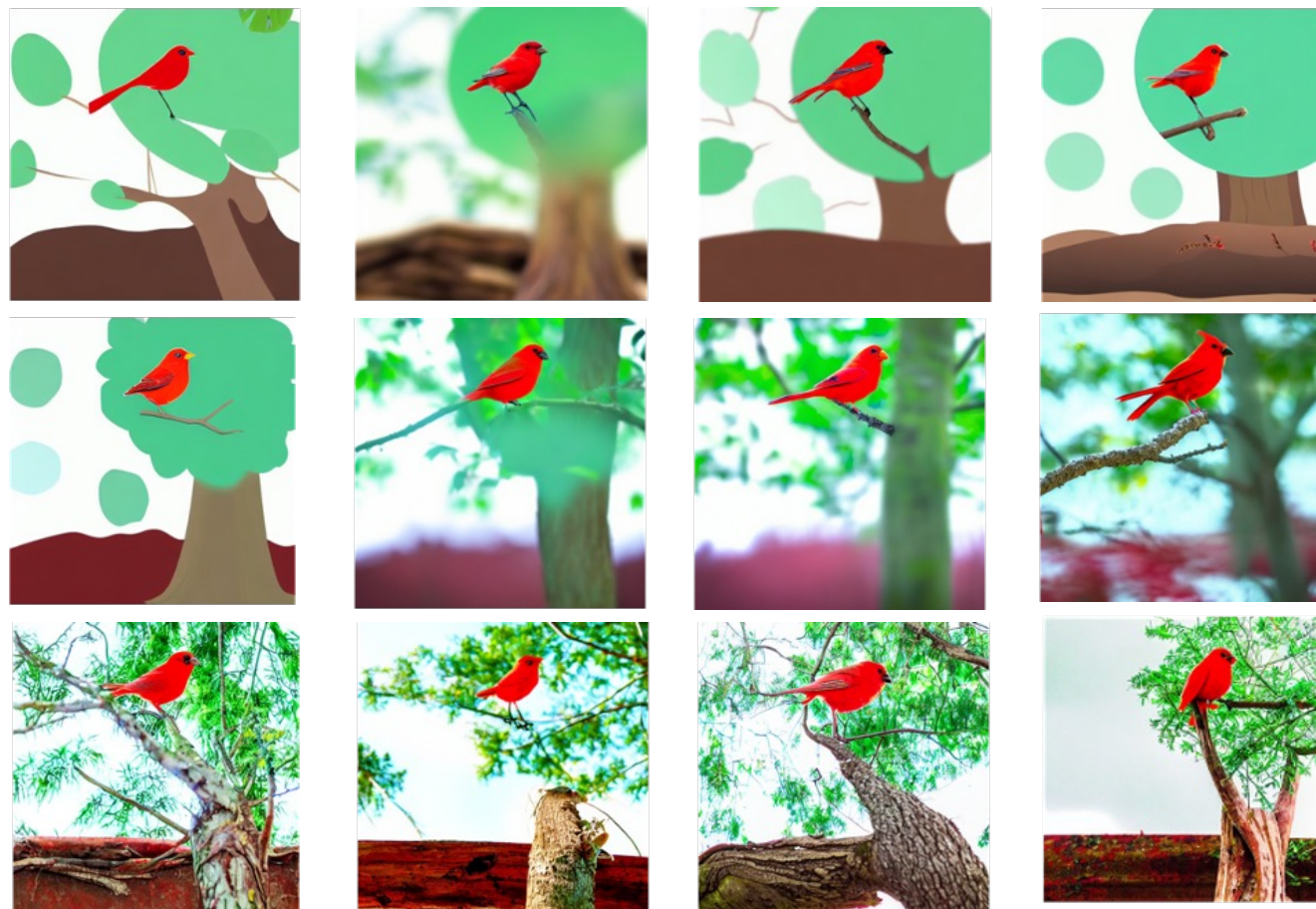


Target  
Subspace

Text Prompt:  
"a photo of a red bird in a tree"



Reference Painting



SDEdit

Iterative  
Loopback

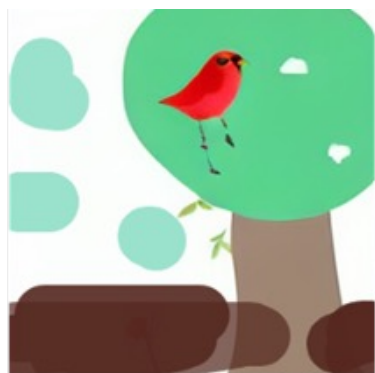
GradOP  
(Ours)

Text Prompt: "a photo of a red bird in a tree"

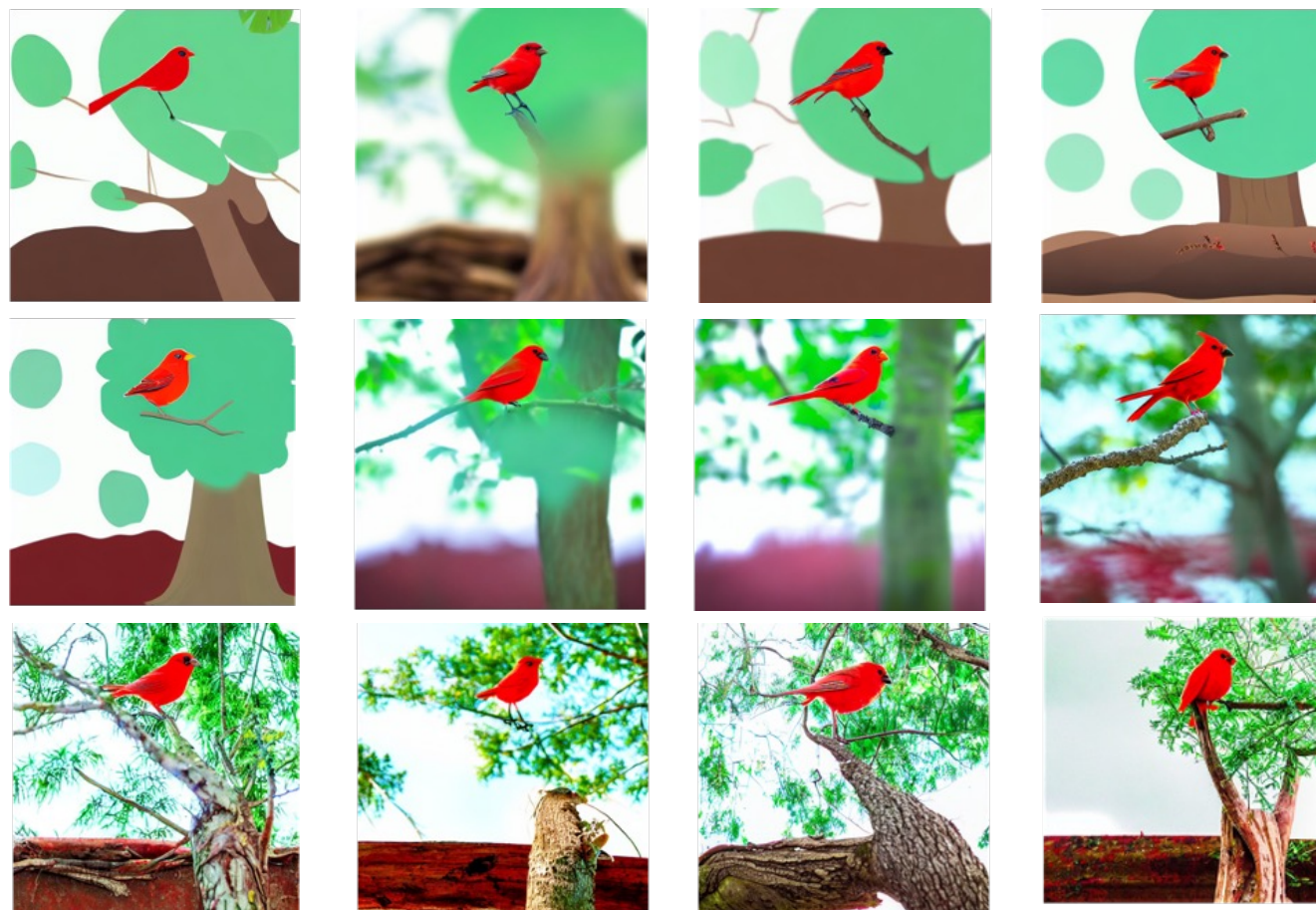


Target  
Subspace

Text Prompt:  
"a photo of a red bird in a tree"



Reference Painting

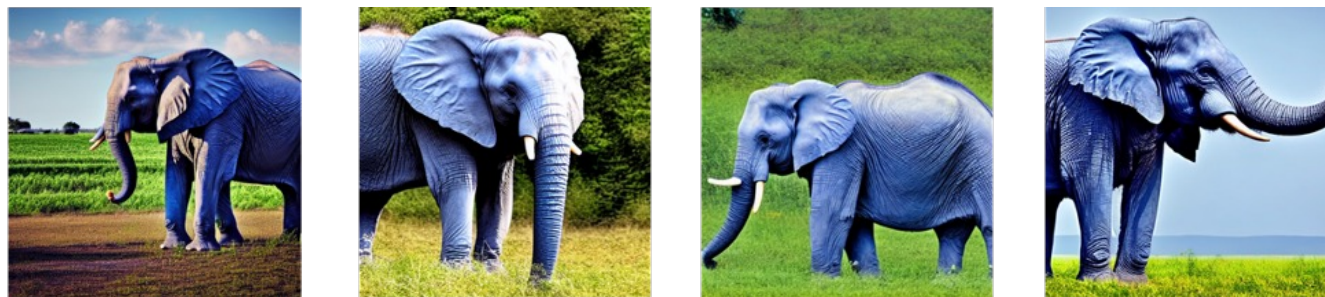


SDEdit

Iterative  
Loopback

GradOP  
(Ours)

Text Prompt: "a photo of a blue elephant in a field"

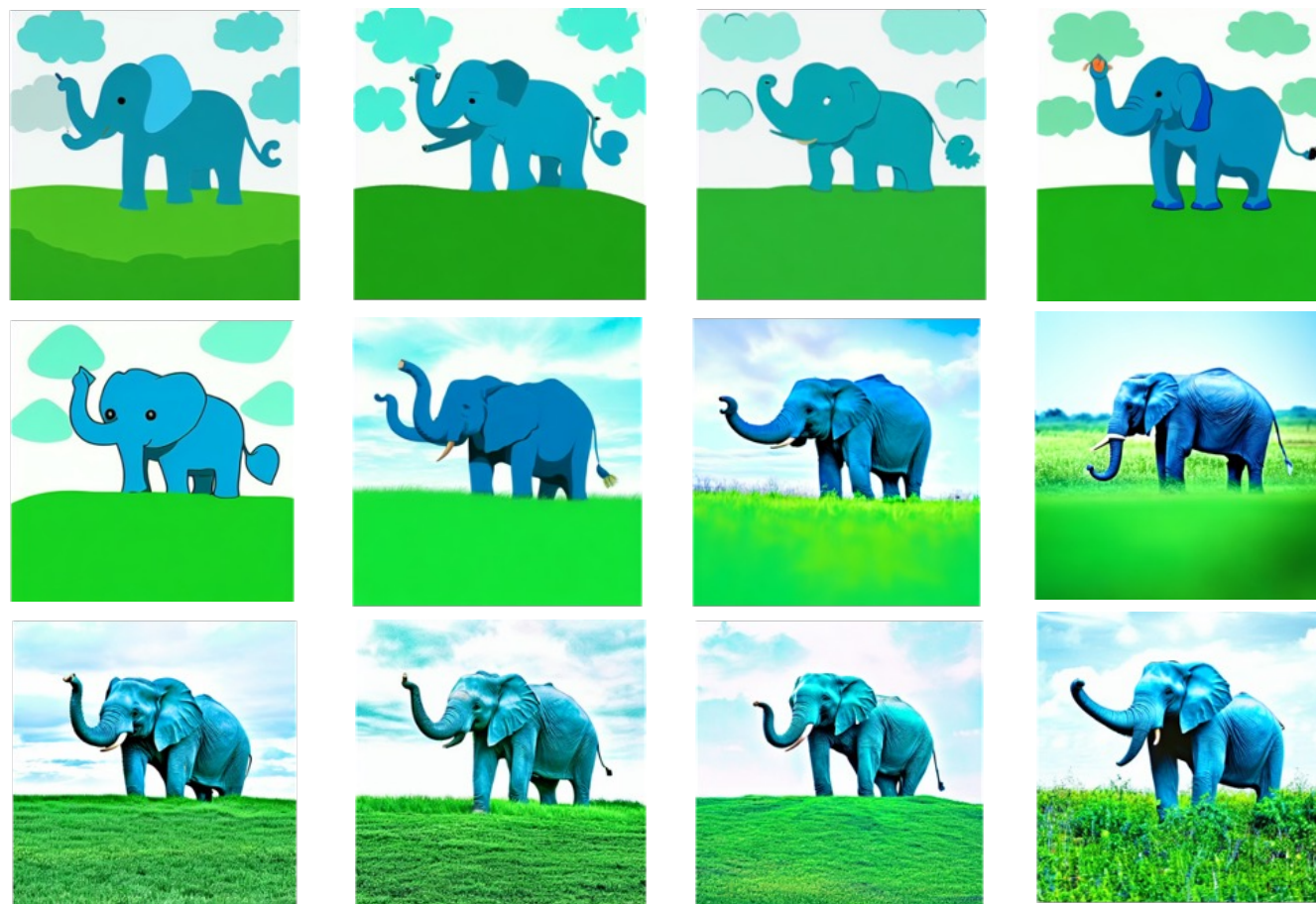


Target  
Subspace

Text Prompt:  
"a photo of a blue elephant  
in a field"



Reference Painting



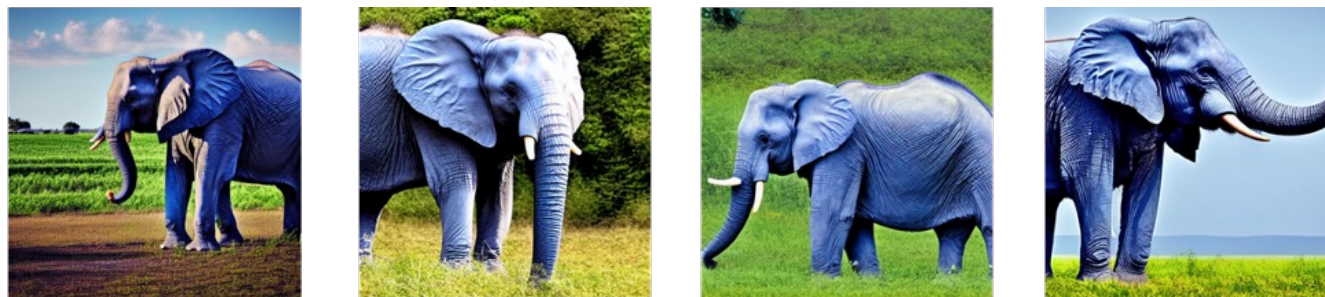
SDEdit

Iterative  
Loopback

GradOP  
(Ours)



Text Prompt: "a photo of a blue elephant in a field"

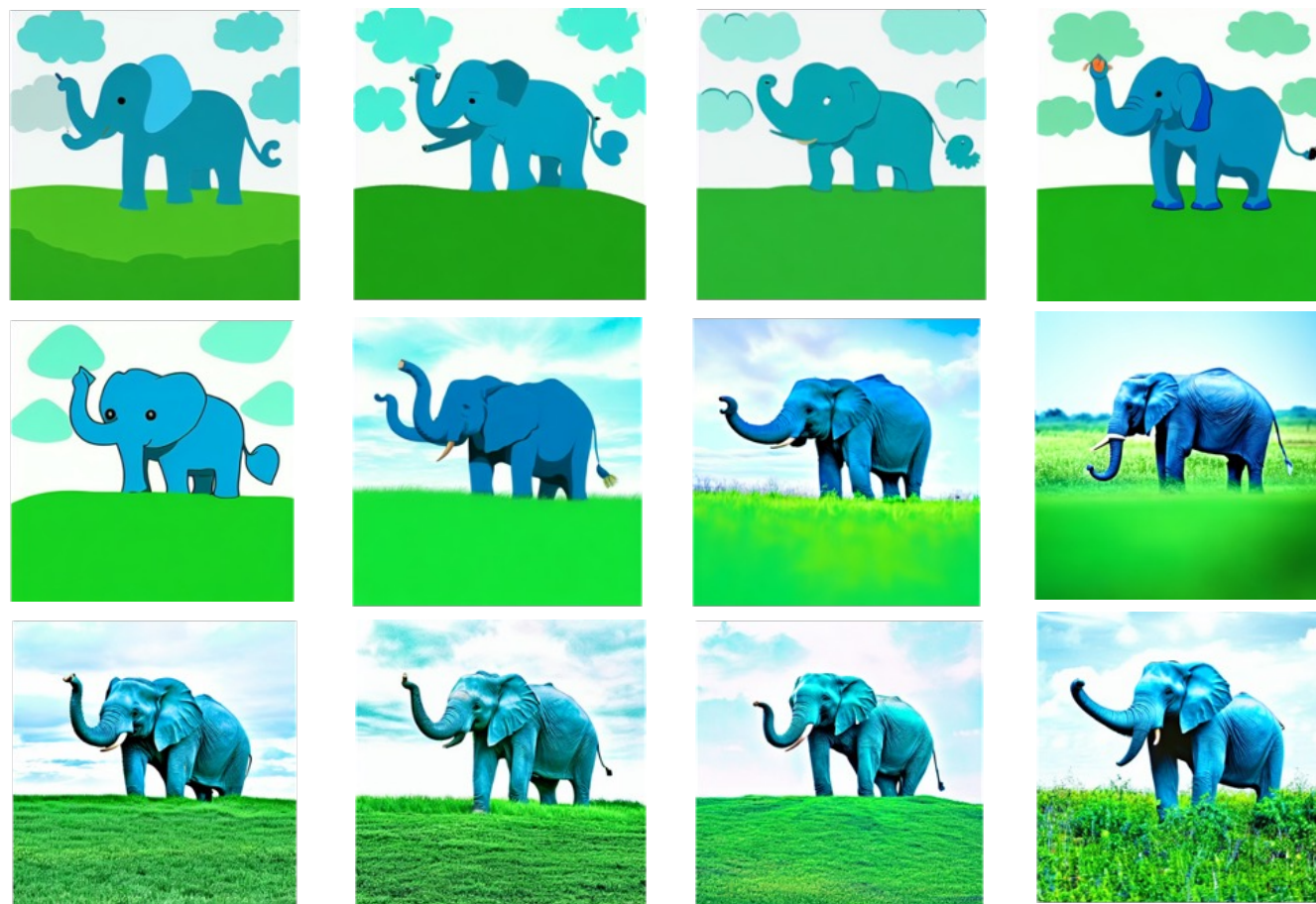


Target  
Subspace

Text Prompt:  
"a photo of a blue elephant  
in a field"



Reference Painting

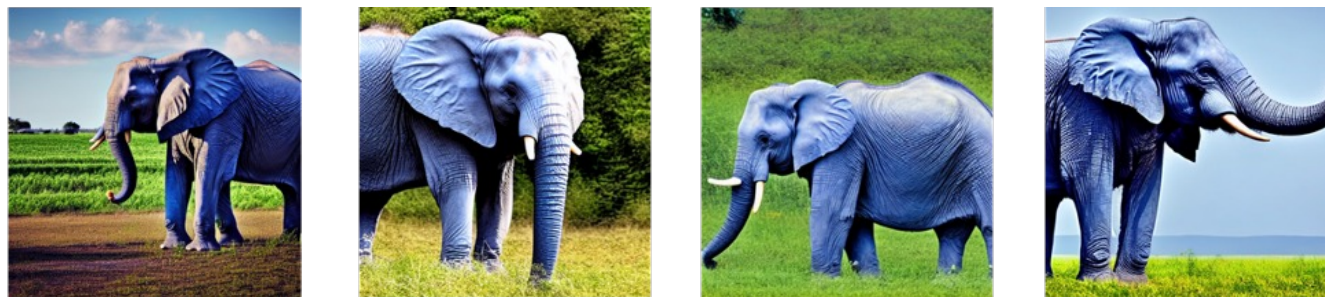


SDEdit

Iterative  
Loopback

GradOP  
(Ours)

Text Prompt: "a photo of a blue elephant in a field"

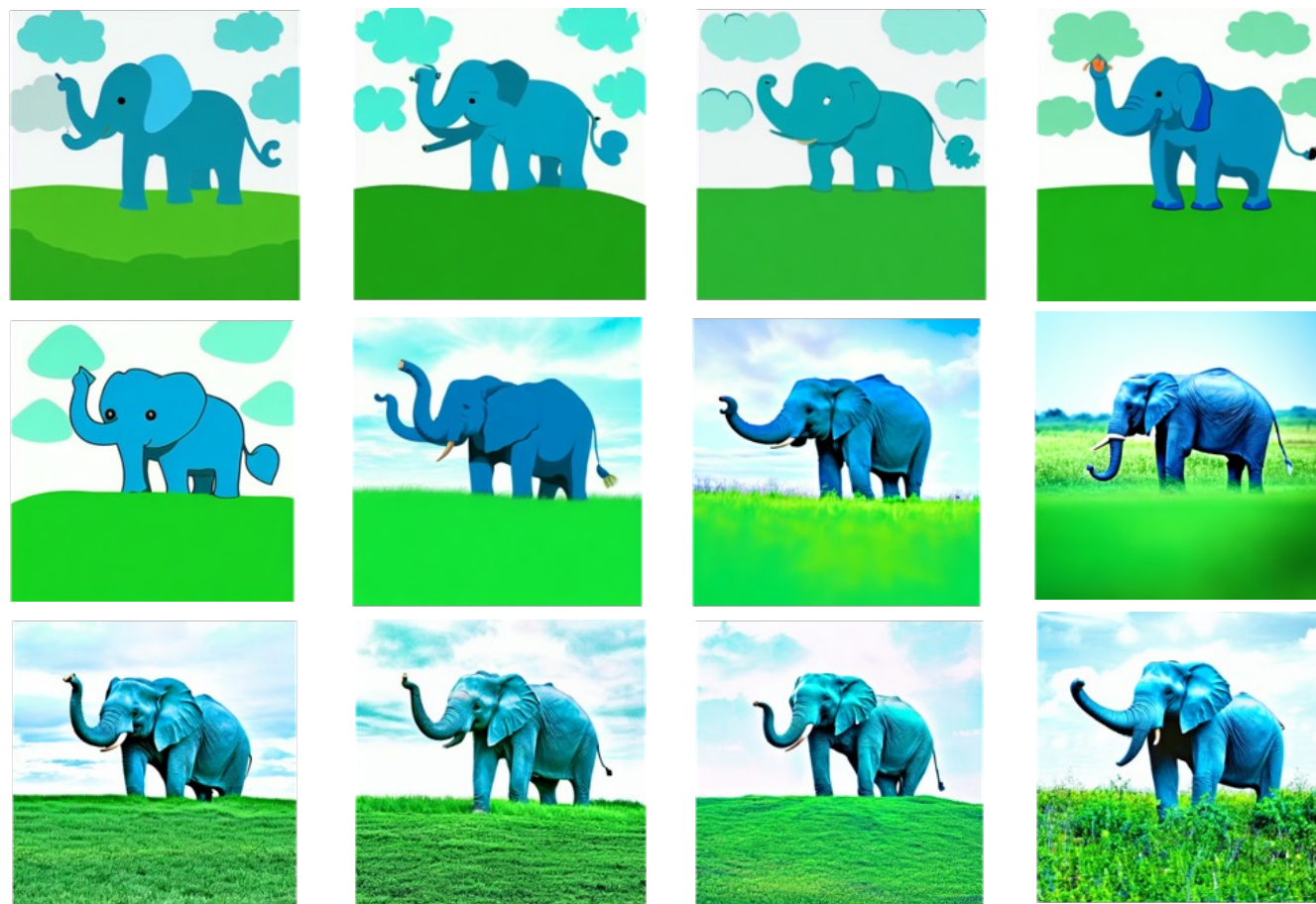


Target  
Subspace

Text Prompt:  
"a photo of a blue elephant  
in a field"



Reference Painting



SDEdit

Iterative  
Loopback

GradOP  
(Ours)

Text Prompt: "a photo of castle by the lake"



Target  
Subspace

$t_0 = 0.7$

$t_0 = 0.75$

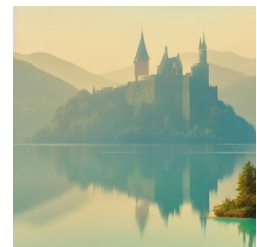
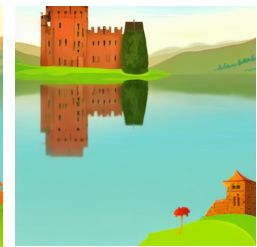
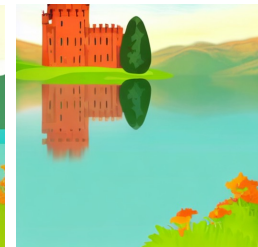
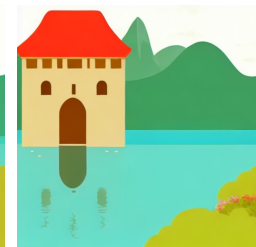
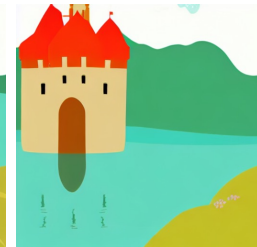
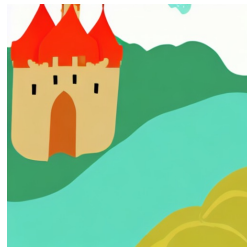
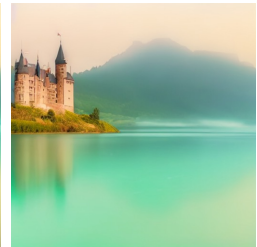
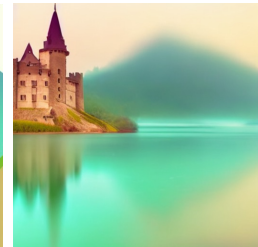
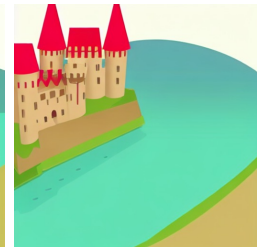
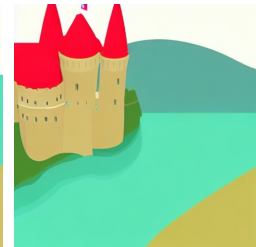
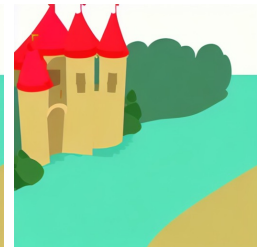
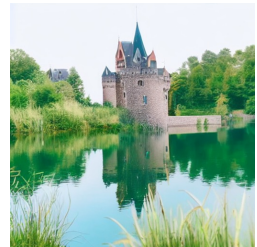
$t_0 = 0.8$

$t_0 = 0.85$

$t_0 = 0.9$

$t_0 = 0.95$

$t_0 = 0.99$



Reference Painting

SDEdit

Ours

Text Prompt: "a photo of castle by the lake"



Target  
Subspace

$t_0 = 0.7$

$t_0 = 0.75$

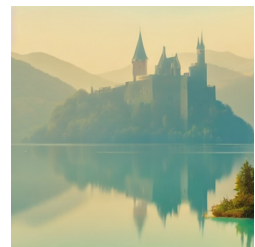
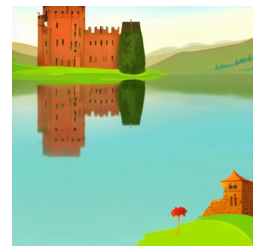
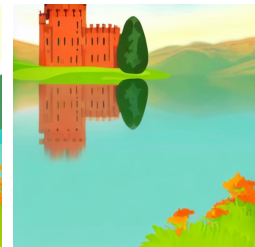
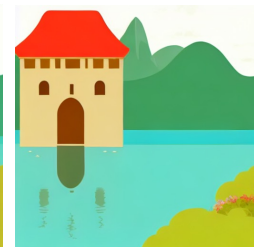
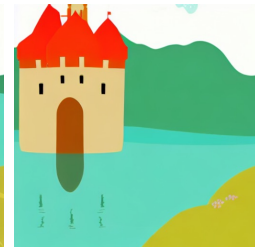
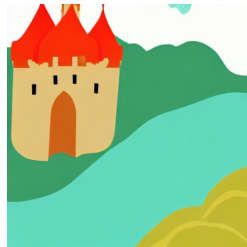
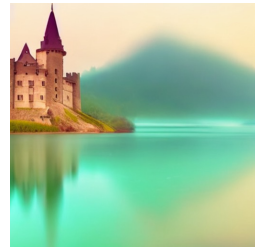
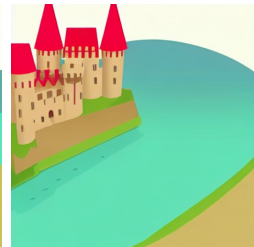
$t_0 = 0.8$

$t_0 = 0.85$

$t_0 = 0.9$

$t_0 = 0.95$

$t_0 = 0.99$



Reference Painting

Faithful but unrealistic

SDEdit

Ours

Text Prompt: "a photo of castle by the lake"



Target  
Subspace

$t_0 = 0.7$

$t_0 = 0.75$

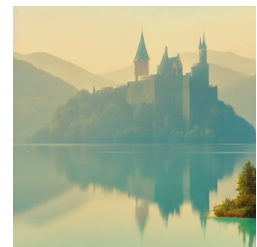
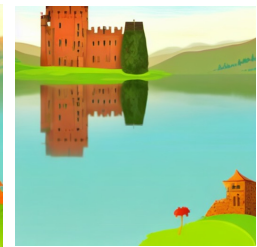
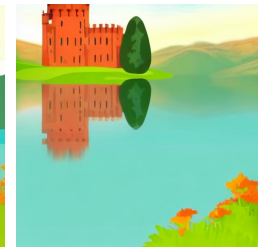
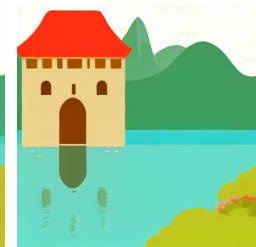
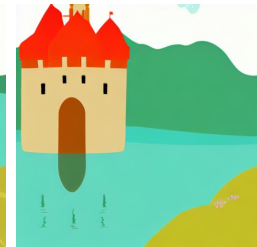
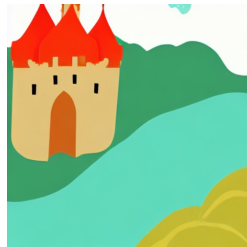
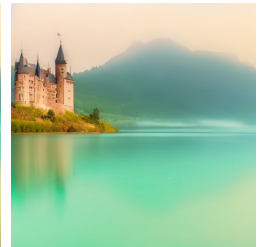
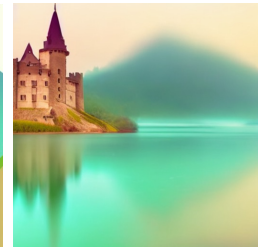
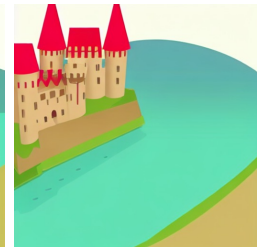
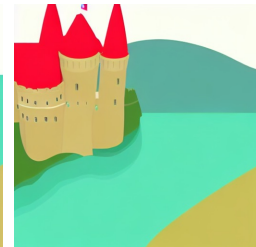
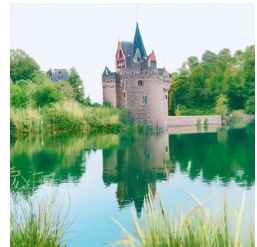
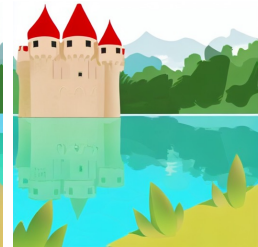
$t_0 = 0.8$

$t_0 = 0.85$

$t_0 = 0.9$

$t_0 = 0.95$

$t_0 = 0.99$



Reference Painting

SDEdit

Faithful but unrealistic

Realistic but unfaithful

Ours

Text Prompt: "a photo of castle by the lake"



Target  
Subspace

$t_0 = 0.7$

$t_0 = 0.75$

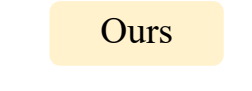
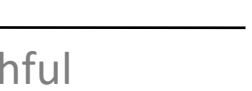
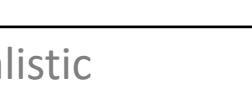
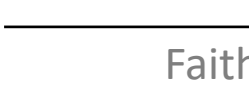
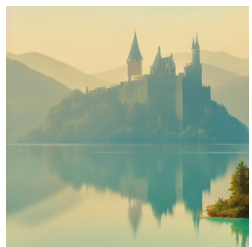
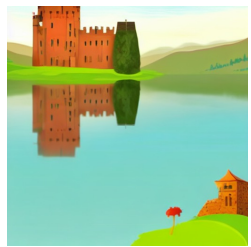
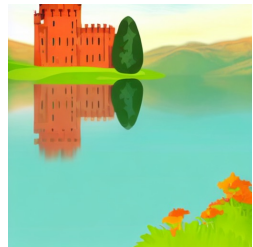
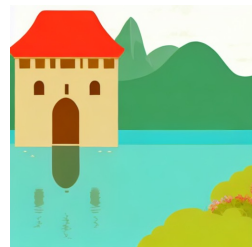
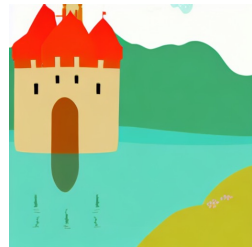
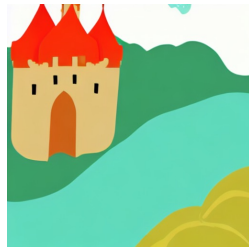
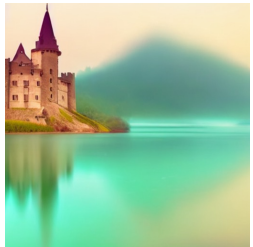
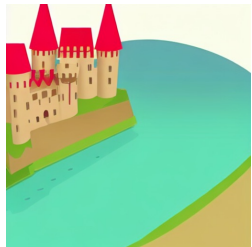
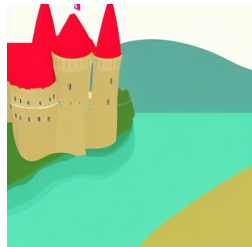
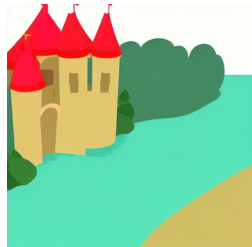
$t_0 = 0.8$

$t_0 = 0.85$

$t_0 = 0.9$

$t_0 = 0.95$

$t_0 = 0.99$



Reference Painting

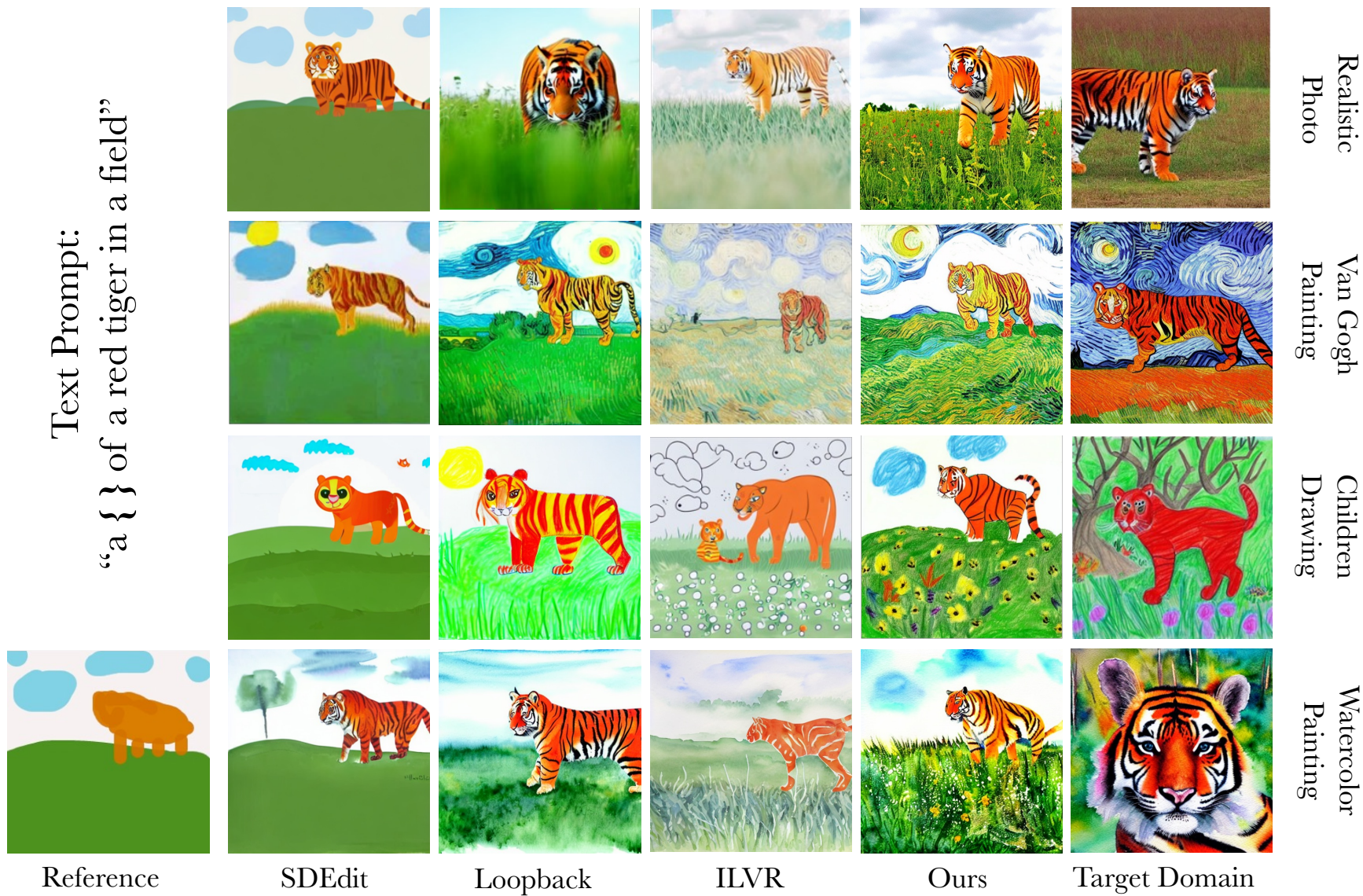
Faithful but unrealistic

SDEdit

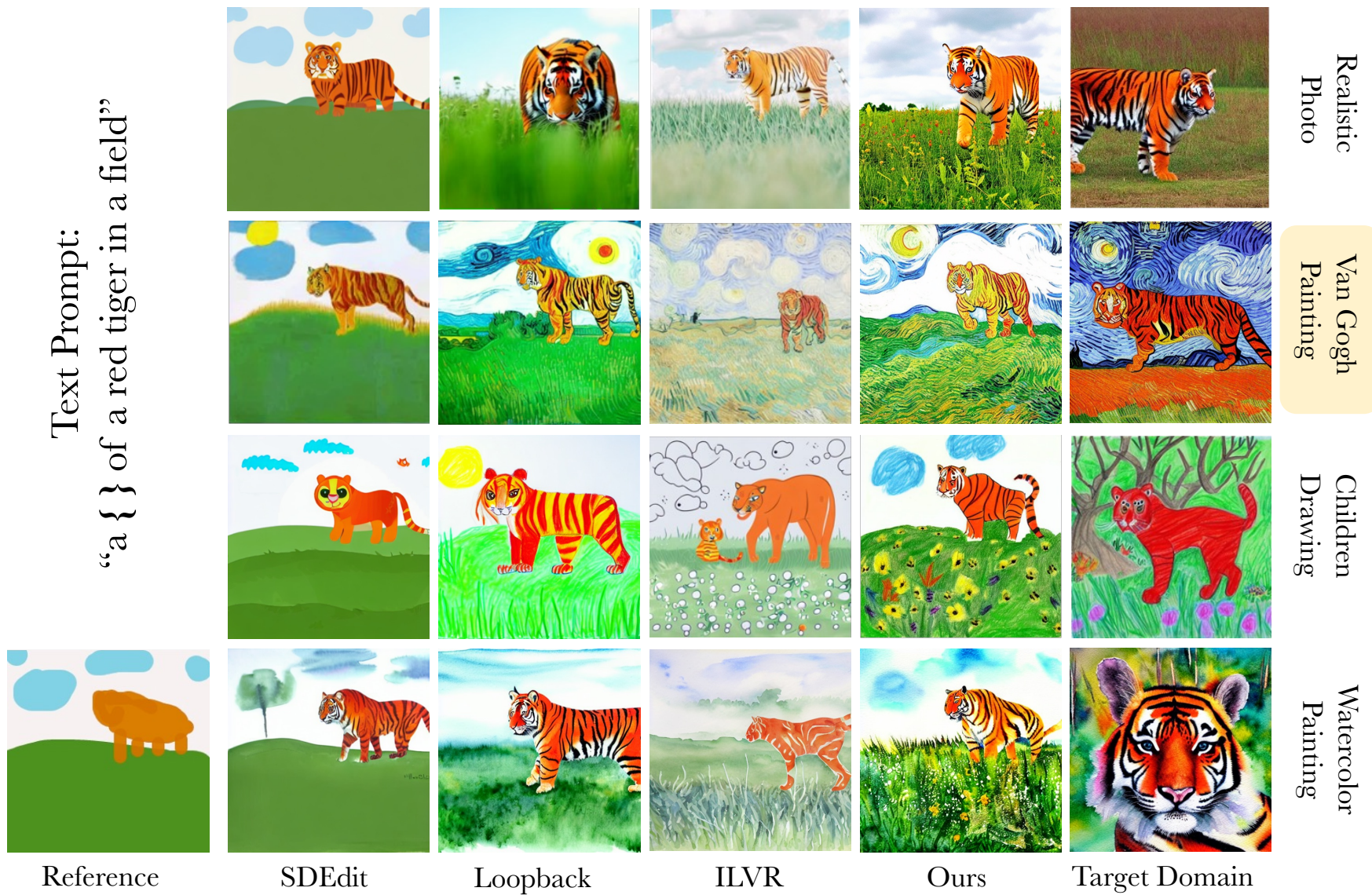
Realistic but unfaithful

Ours

# Generalizability across Target Domains

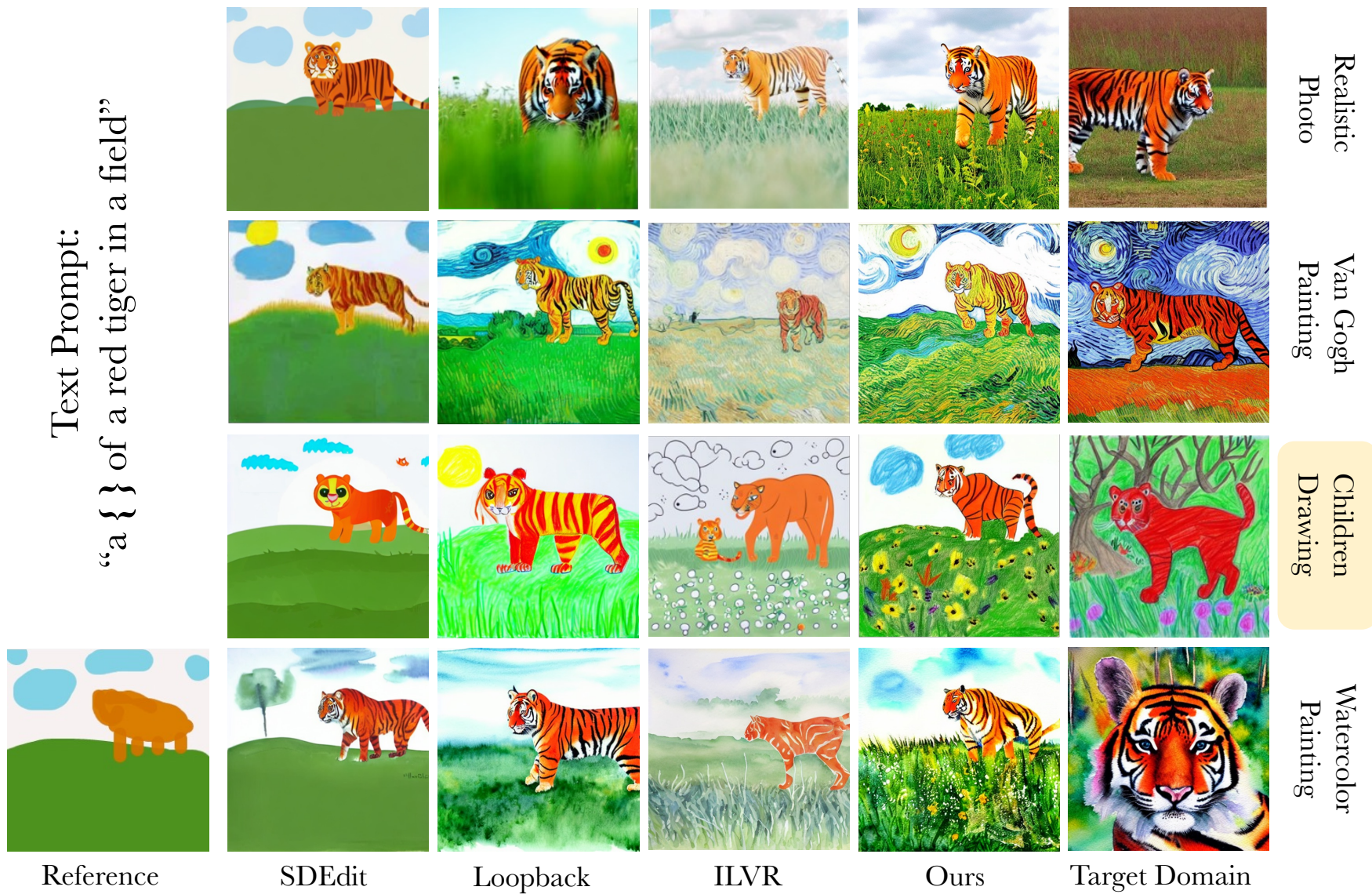


# Generalizability across Target Domains

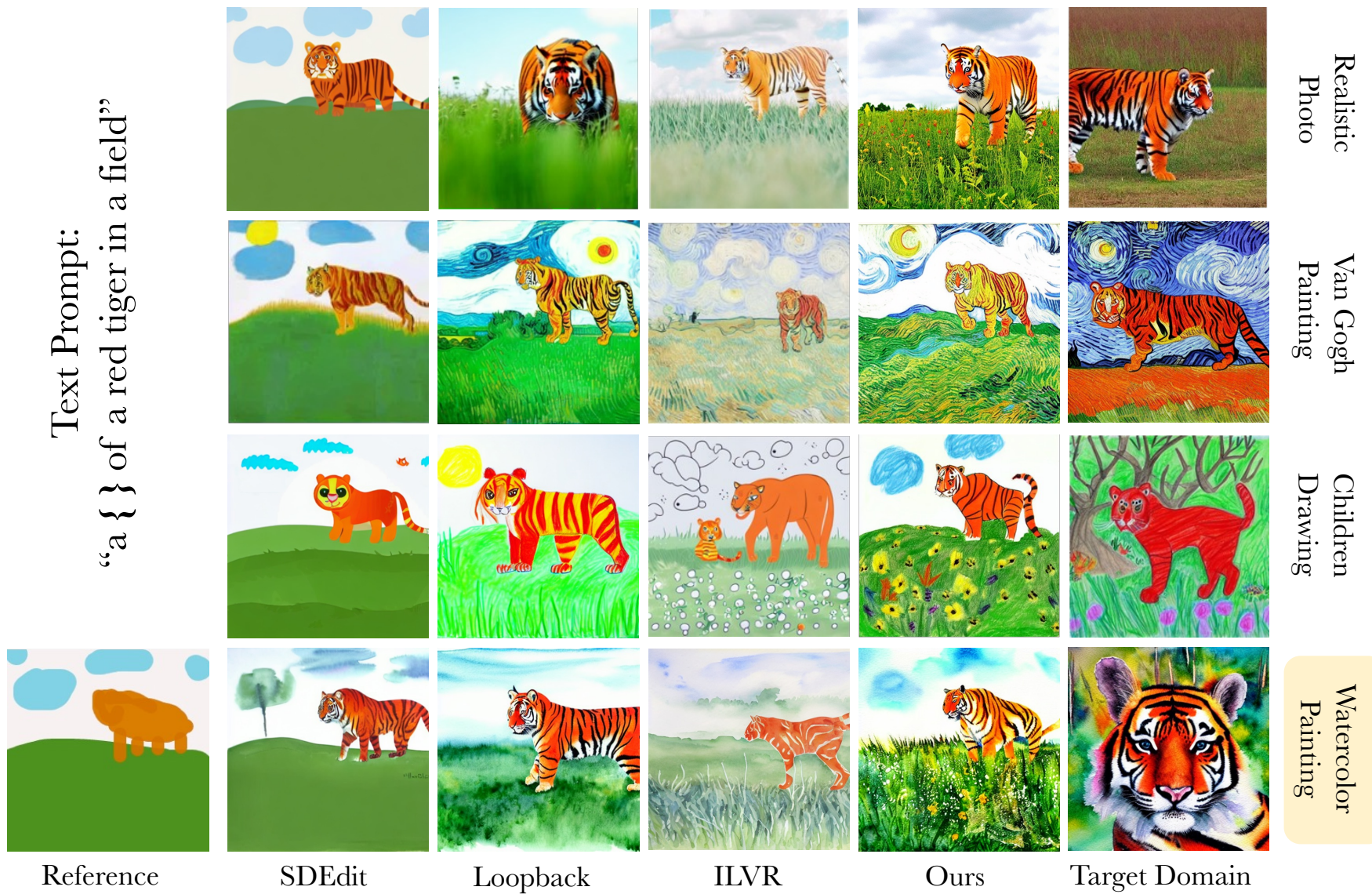




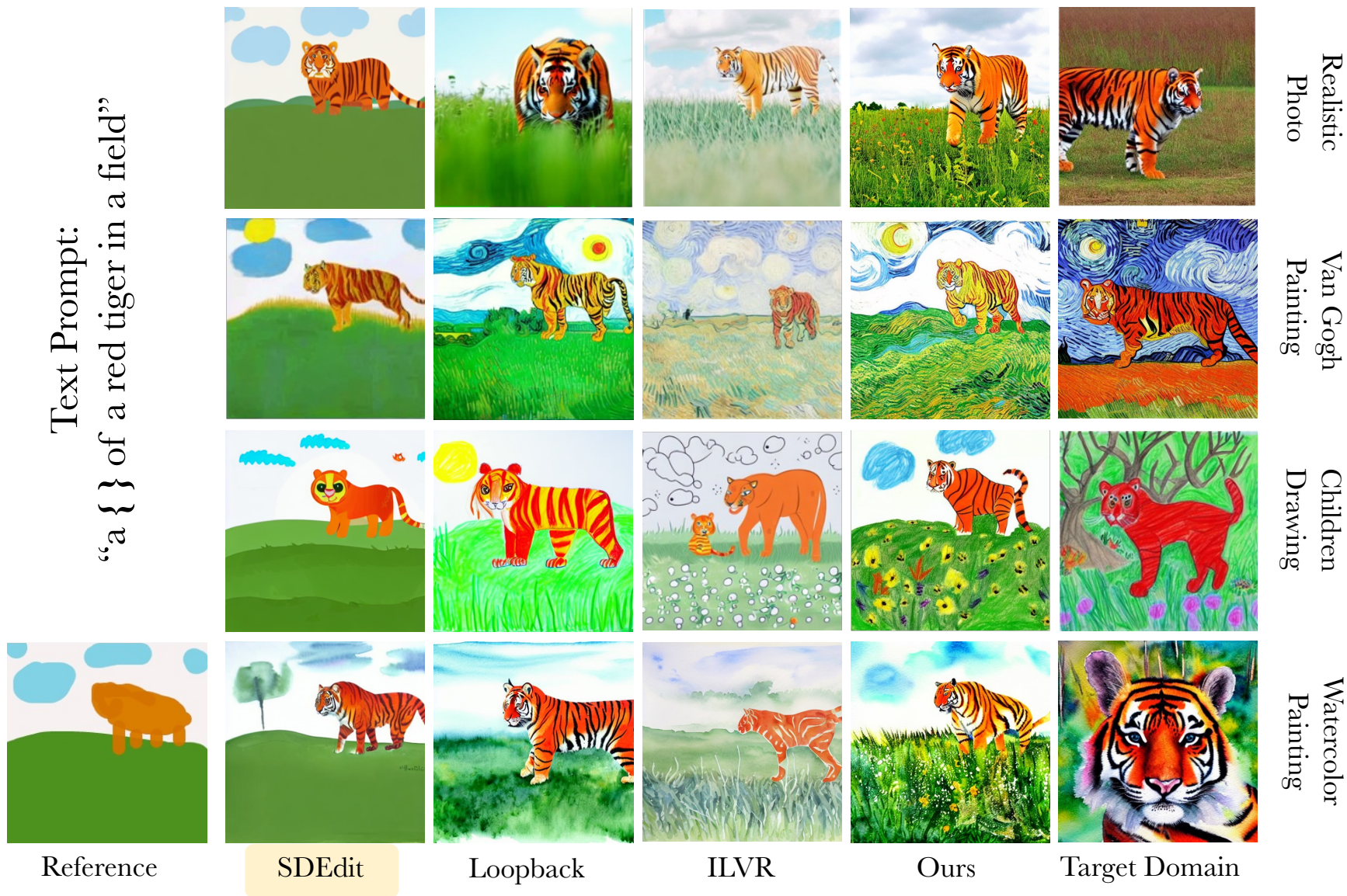
# Generalizability across Target Domains



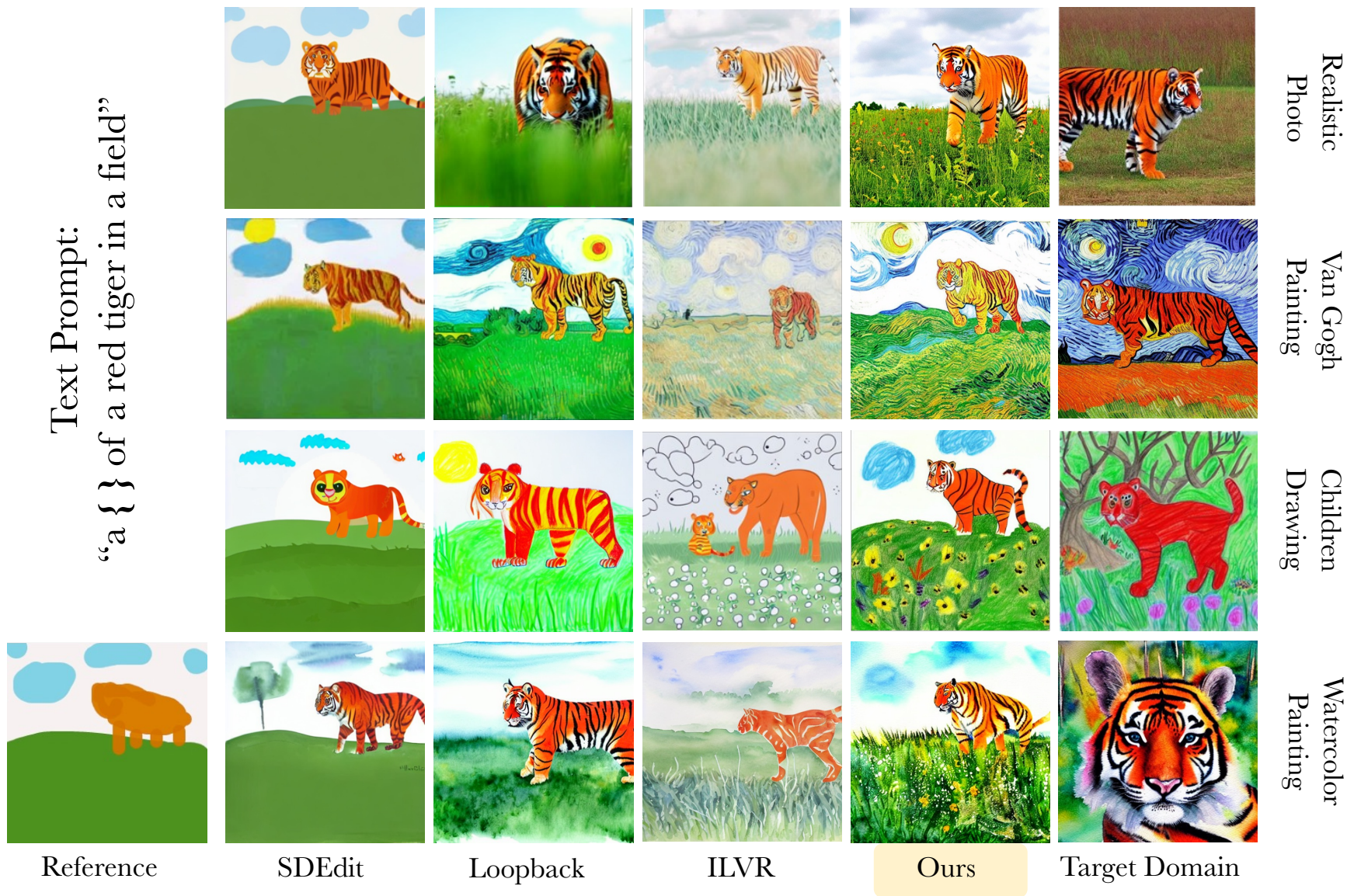
# Generalizability across Target Domains



# Generalizability across Target Domains



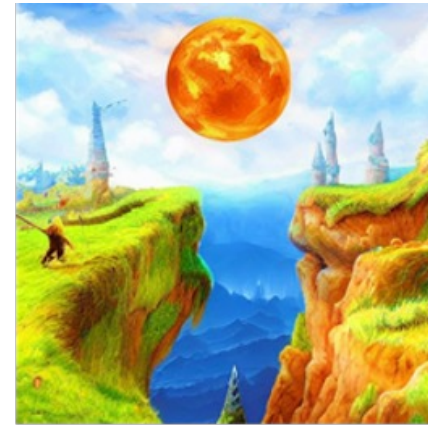
# Generalizability across Target Domains



# Controlling Semantics of Different Painting Regions

Text Prompt:

*"a fantasy landscape, trending on artstation"*



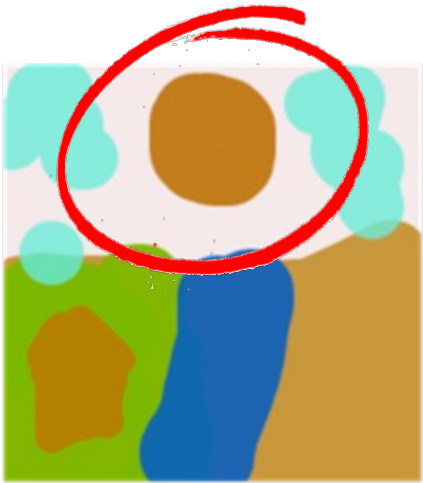
*Reference Painting*

Semantics of different painting regions might not accurately reflect user's intent.

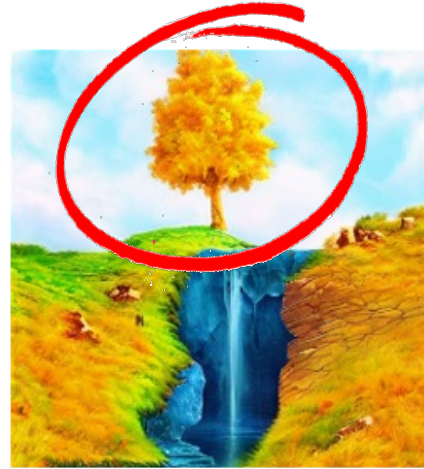
# Controlling Semantics of Different Painting Regions

Text Prompt:

*"a fantasy landscape, trending on artstation"*



*Reference Painting*



Semantics of different painting regions might not accurately reflect user's intent.

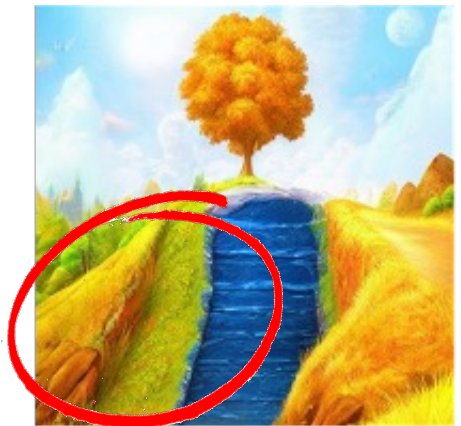
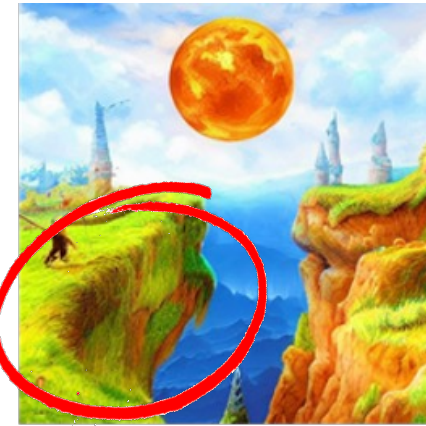
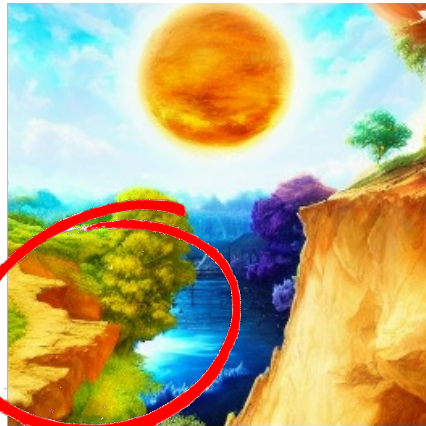
# Controlling Semantics of Different Painting Regions

Text Prompt:

*"a fantasy landscape, trending on artstation"*



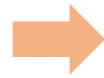
Reference Painting



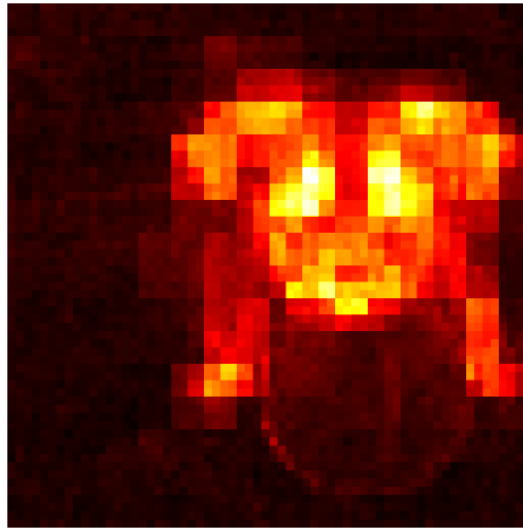
Some regions might be entirely omitted if the model does not understand that it corresponds to a separate semantic entity

# Controlling Semantics of Different Painting Regions

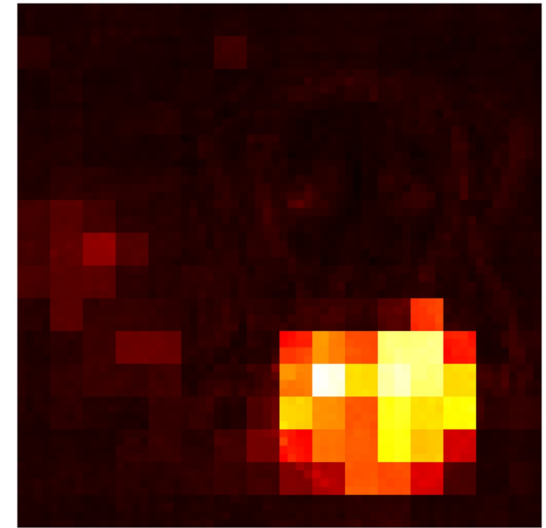
Text Prompt:  
"a *dog* with a *ball*"



"dog"



"ball"



*Average cross-attention maps during reverse diffusion process*

Cross-Attention maps show high overlap with semantic segmentation of different entities





Reference Painting



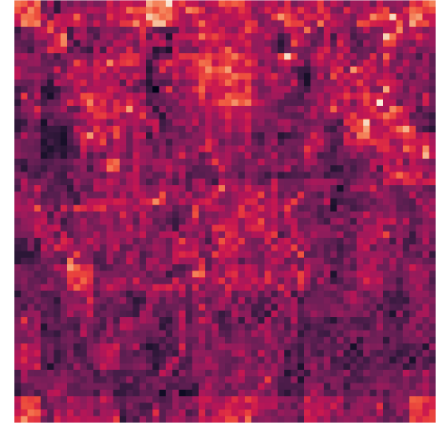
Reference Painting

Binary Mask for  
Painting Region



$B_i$  : "hut"

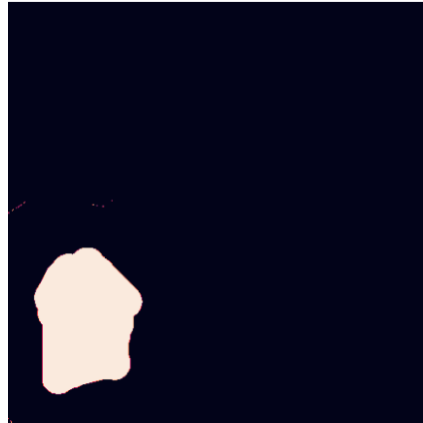
Original  
Cross-attention Map



$A_i^t$

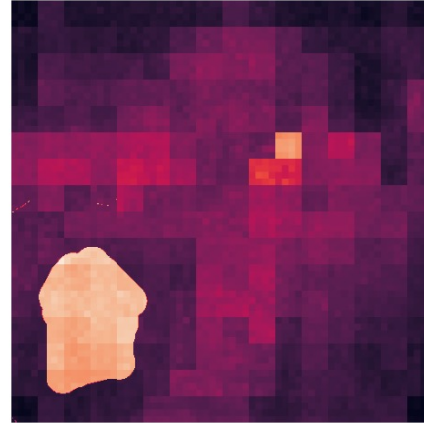


Reference Painting



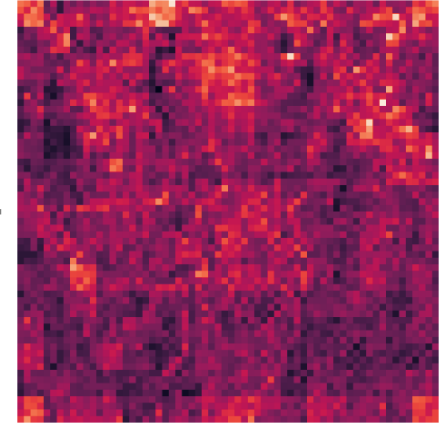
Binary Mask for  
Painting Region

$\mathcal{B}_i$  : “hut”



Modified  
Cross-attention Map

$\tilde{\mathcal{A}}_i^t$



Original  
Cross-attention Map

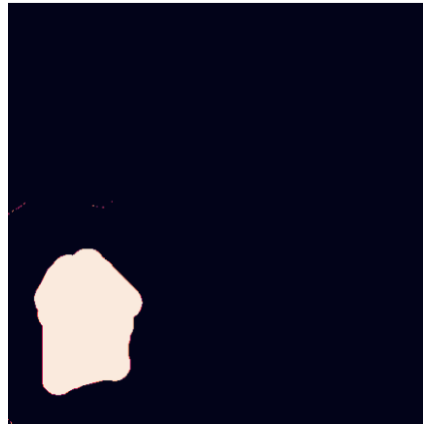
$\mathcal{A}_i^t$

Introducing Cross-attention Overlap  
with Desired Painting Regions

$$\tilde{\mathcal{A}}_i^t = w_i \left[ (1 - \kappa_t) \mathcal{A}_i^t + \kappa_t \frac{\mathcal{B}_i}{\|\mathcal{B}_i\|_F} \|\mathcal{A}_i^t\|_F \right]$$

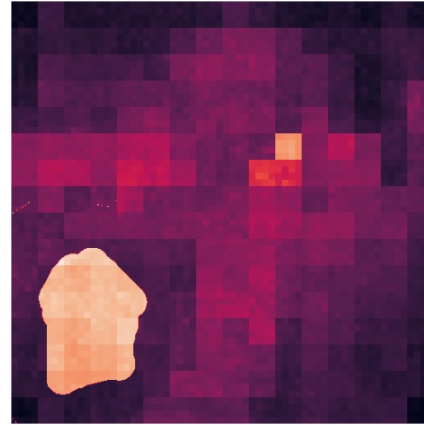


Reference Painting



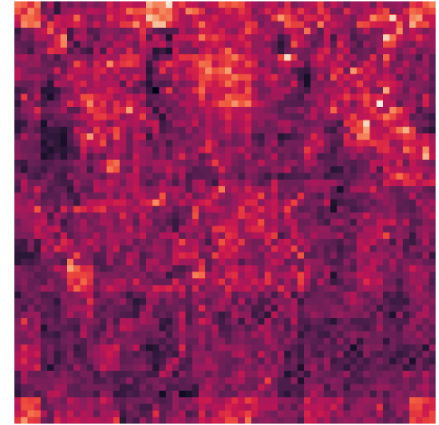
Binary Mask for Painting Region

$\mathcal{B}_i$  : "hut"



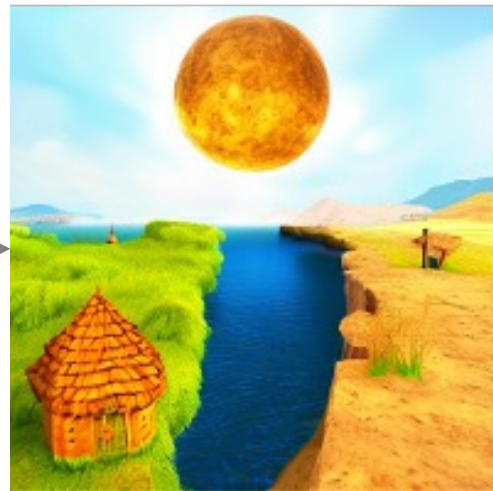
Modified Cross-attention Map

$\tilde{\mathcal{A}}_i^t$



Original Cross-attention Map

$\mathcal{A}_i^t$



Generated Image

Introducing Cross-attention Overlap with Desired Painting Regions

$$\tilde{\mathcal{A}}_i^t = w_i \left[ (1 - \kappa_t) \mathcal{A}_i^t + \kappa_t \frac{\mathcal{B}_i}{\|\mathcal{B}_i\|_F} \|\mathcal{A}_i^t\|_F \right]$$

Reference Painting



Generation Outputs - w/o semantic control

Reference Painting



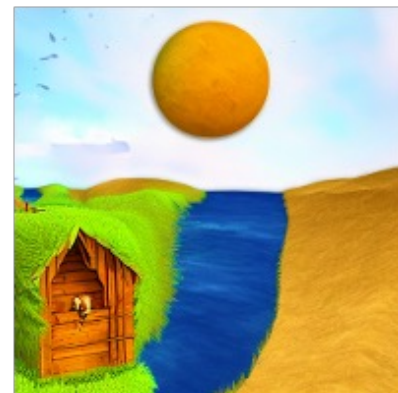
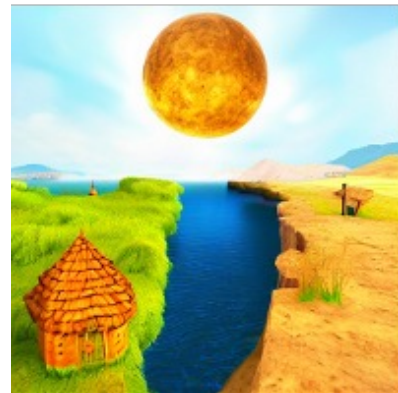
Generation Outputs - w/o semantic control



Semantic Guide



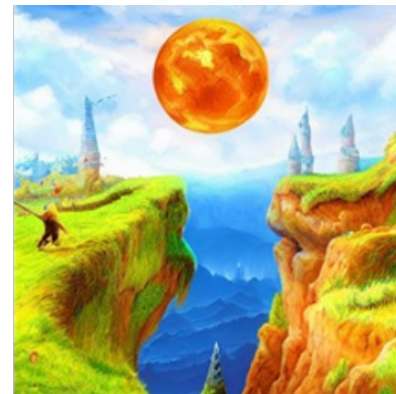
Generation Outputs - with semantic control



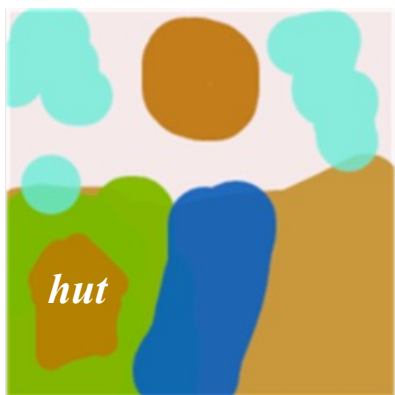
Reference Painting



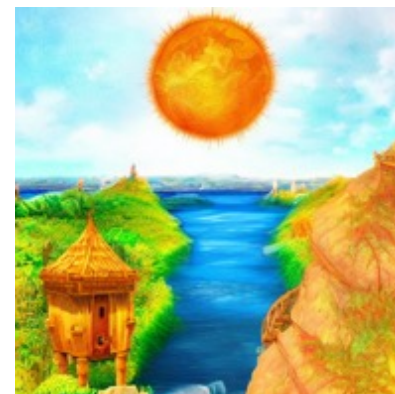
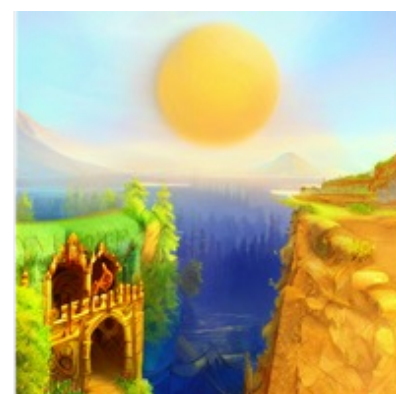
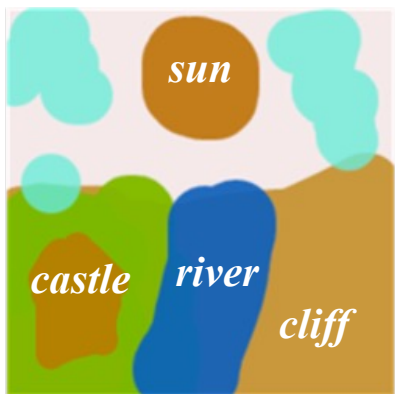
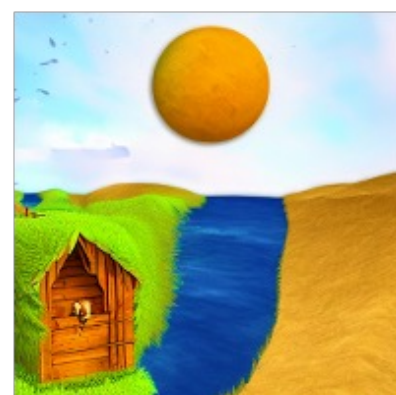
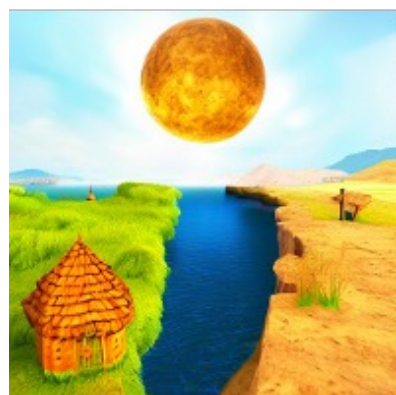
Generation Outputs - w/o semantic control



Semantic Guide



Generation Outputs - with semantic control





Thanks for listening!

Project Page and Online Demo   
<https://1jsingh.github.io/gradop>



Jaskirat Singh<sup>†</sup>



Stephen Gould<sup>†\*</sup>



Liang Zheng<sup>†\*</sup>

Poster Tag: TUE-PM-179