

# SINE: SINgLe Image EditIng with Text-to-Image Diffusion Models

Zhixing Zhang<sup>1</sup> Ligong Han<sup>1</sup> Arnab Ghosh<sup>2</sup> Dimitris Metaxas<sup>1</sup> Jian Ren<sup>2</sup>

<sup>1</sup> Rutgers University <sup>2</sup> Snap Inc.

# SINE: SINgLe Image EditIng with Text-to-Image Diffusion Models

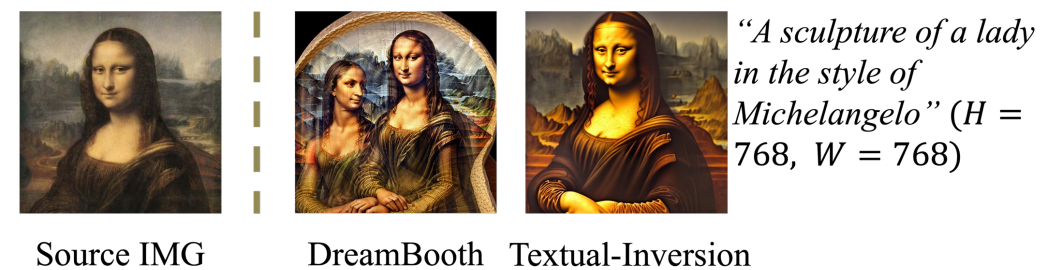
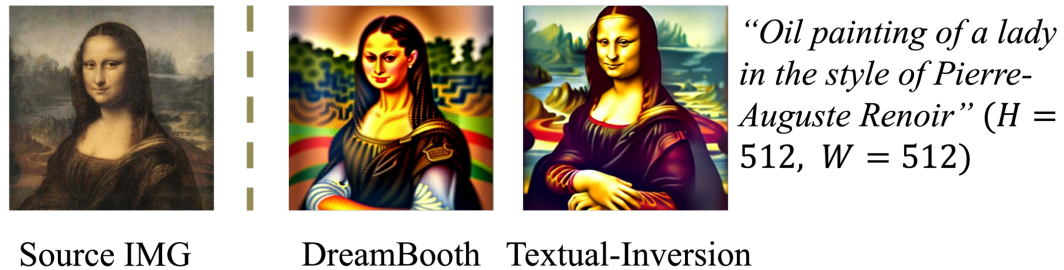
Zhixing Zhang<sup>1</sup> Ligong Han<sup>1</sup> Arnab Ghosh<sup>2</sup> Dimitris Metaxas<sup>1</sup> Jian Ren<sup>2</sup>

<sup>1</sup> Rutgers University <sup>2</sup> Snap Inc.



# Motivation

- Pre-trained diffusion model for real image editing



Overfitting  
 $\Rightarrow$  *Model-based guidance*

Arbitrary resolution Generation  
 $\Rightarrow$  Patch-based fine-tuning

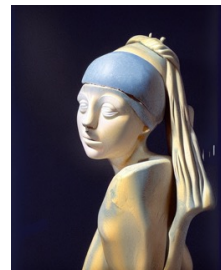


# Motivation

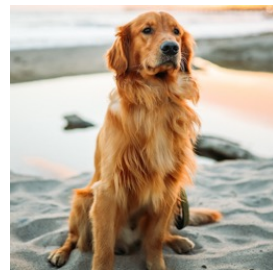
- Pre-trained diffusion model for real image editing:  
*e.g.* Textual-Inversion, DreamBooth
- Overfitting  $\Rightarrow$  *Model-based guidance*
- Arbitrary resolution  $\Rightarrow$  Patch-based fine-tuning



Source Image



*"... sculpture"*



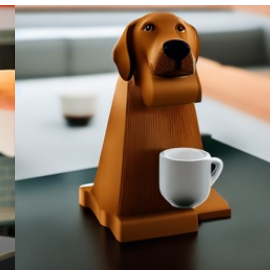
Source Image



*"...super-hero cape"*



*"coffee machine..."*



*"coffee maker..."*



Source Image



*"...covered by snow"*



Source Image



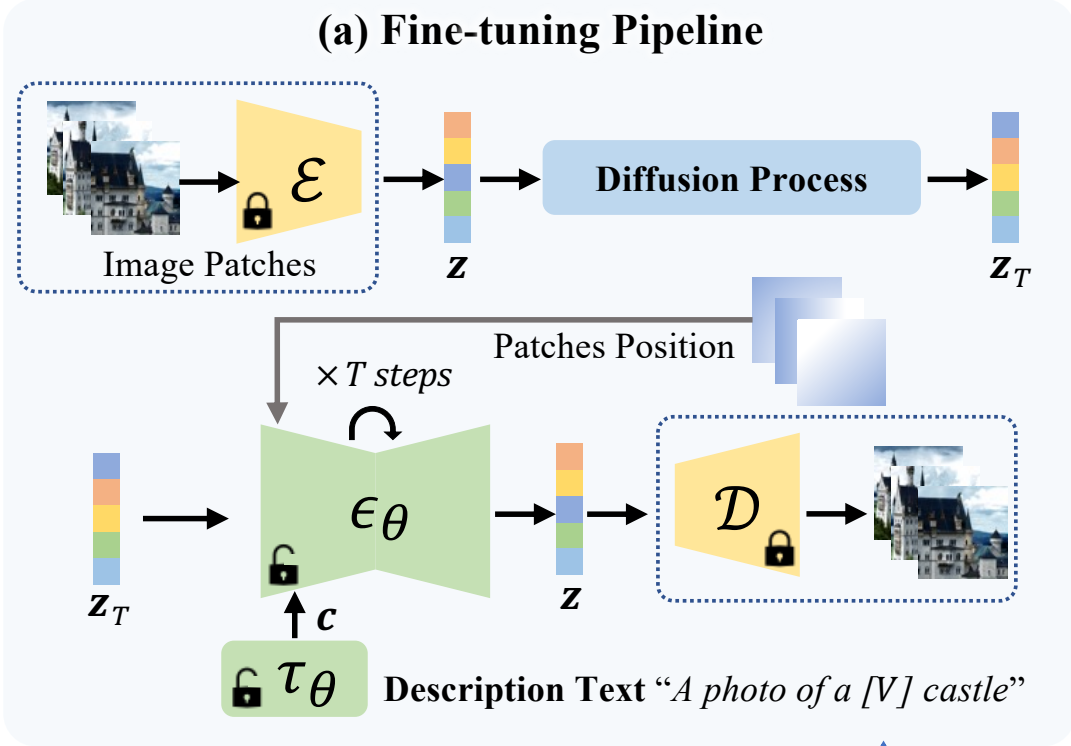
*"... in snow"*



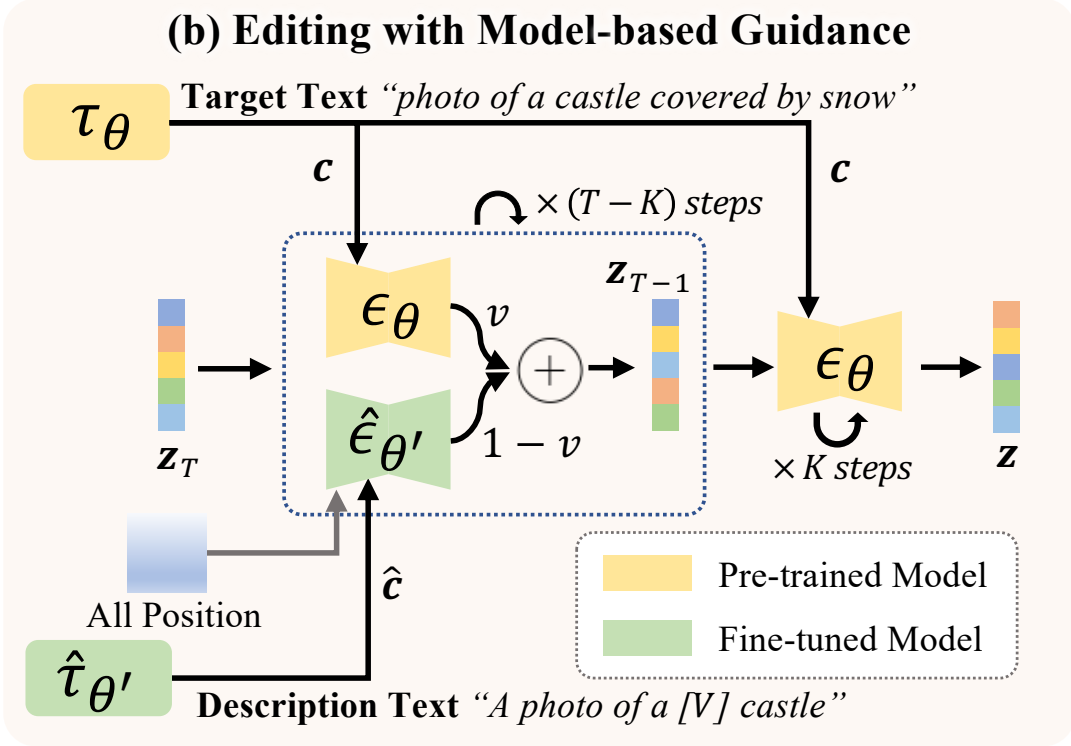
*"... blossom street"*



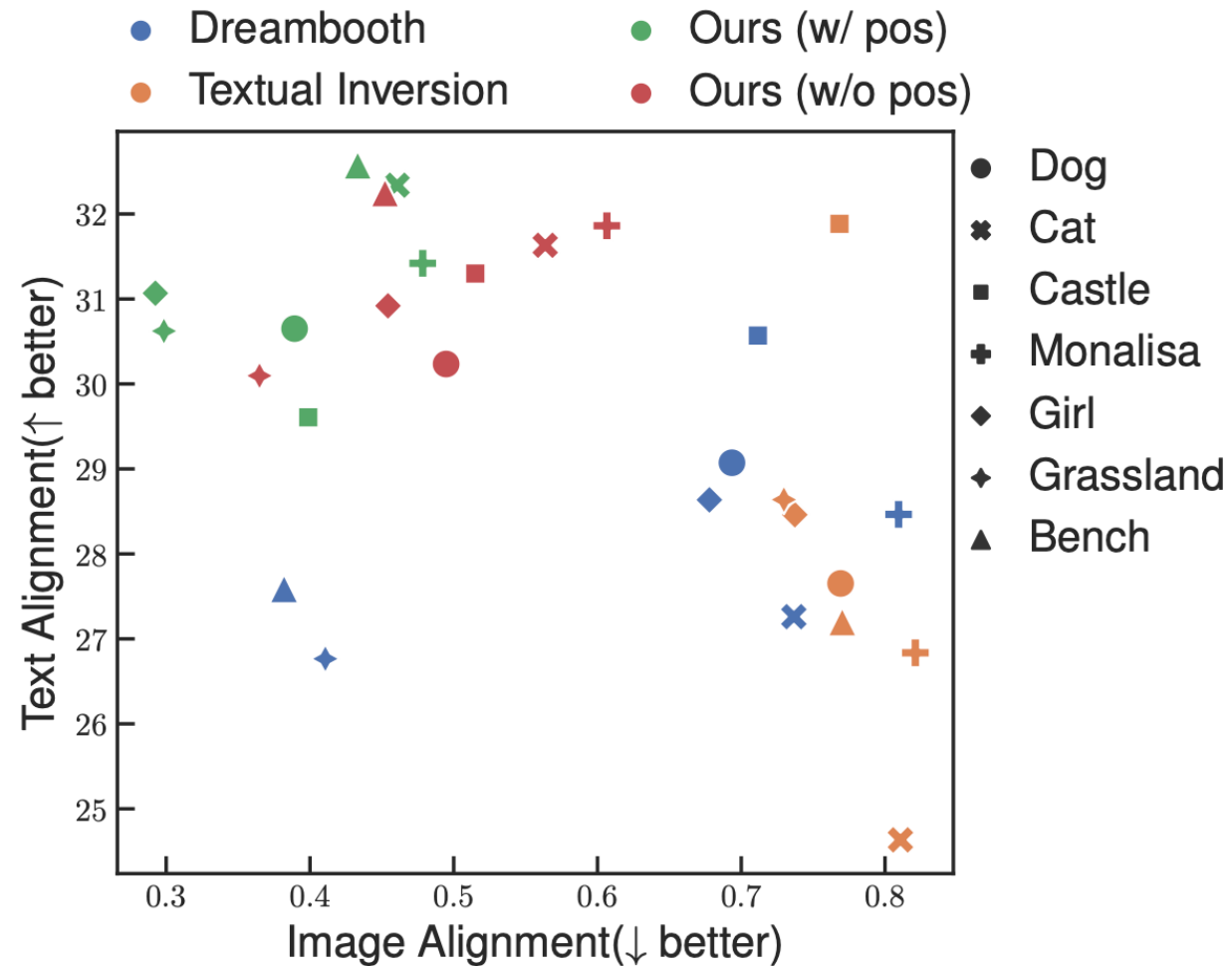
# Methods



special token



# Results



# Results

- $K = 400, v = 0.7$



Source Image



*"A painting of a grassland in the style of Vincent Van Gogh"*



*"A photo of a grassland in thunderstorm"*



Source Image



*"A photo of a bench on sand"*



*"A painting of a bench in the style of Vincent Van Gogh"*

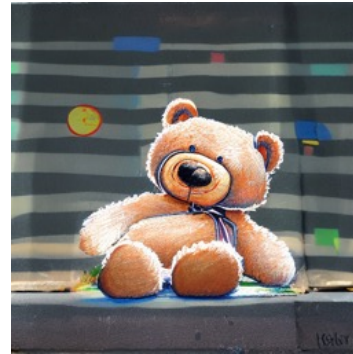


# Results

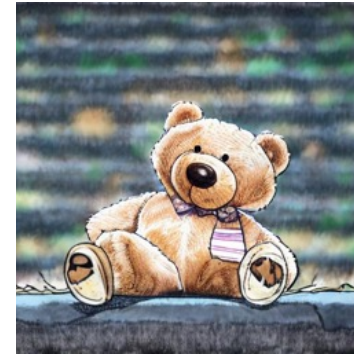
- $K = 400, v = 0.7$



**Source Image**



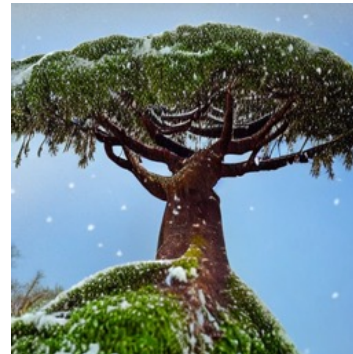
*“A painting of a teddy bear in the style of street painting”*



*“A cartoon of a teddy bear”*



**Source Image**



*“A photo of a tree covered by snow”*



*“A photo of a tree on fire”*



# Results

- $K = 400, v = 0.7$



**Source Image**



*“A sculpture of a horse”*



*“A cartoon of a horse”*



**Source Image**



*“modern art painting of a vase”*



*“painting of a vase in the style of Vincent Van Gogh”*

# Results

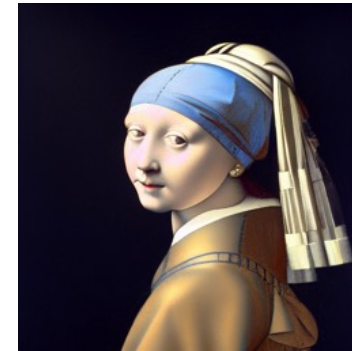
- $K = 400, v = 0.7$



**Source Image**



*"A painting of a girl in the style of Vincent Van Gogh"*



*"A painting of a girl in the style of Leonardo da Vinci"*



**Source Image**



*"A photo of a building covered by snow"*



*"A photo of a cyberpunk building"*

# Results

- $K = 400, v = 0.7$



**Source Image**



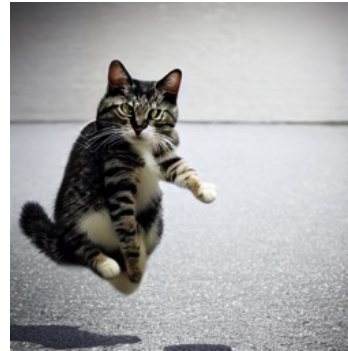
*"A photo of a Husky in ocean"*



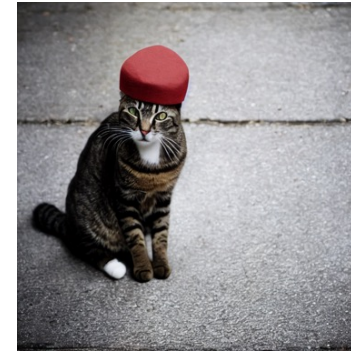
*"A cartoon of a dog on grass"*



**Source Image**



*"A photo of a jumping cat"*



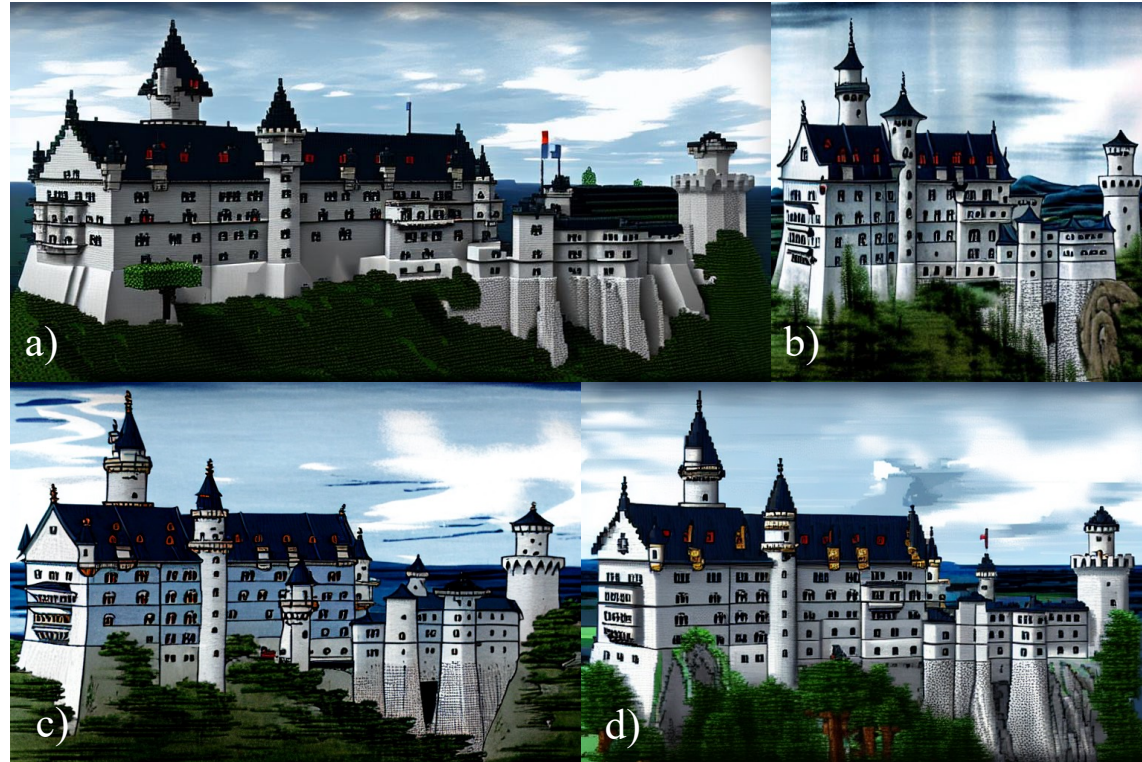
*"A photo of a cat wearing a red hat"*

# More editing tasks: Different resolutions

- Synthesize images w/ various resolutions



Source Image



a) “a castle in the style of Minecraft” (512×1024)

b) “a Chinese painting of a castle”(512×512)

c) “painting of a castle in the style of Hokusai”(512×768)

d) “pixel art of a castle”(512×768)

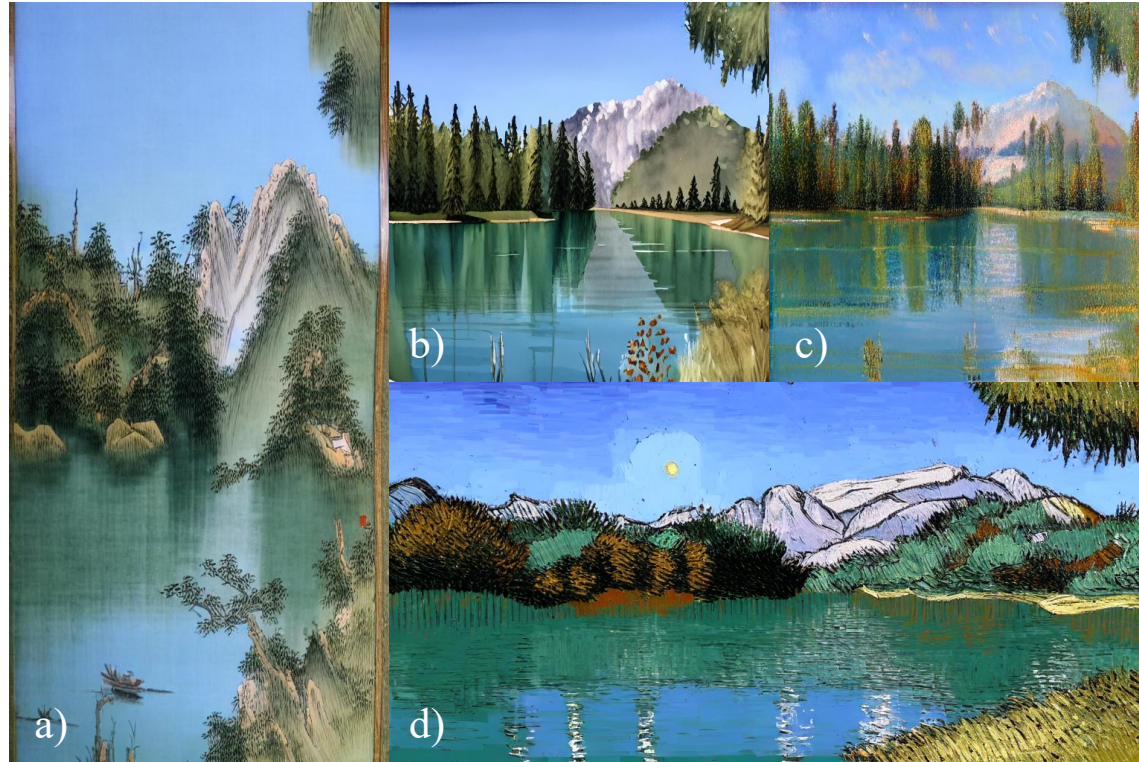


# More editing tasks: Different resolutions

- Synthesize images w/ various resolutions



**Source Image**



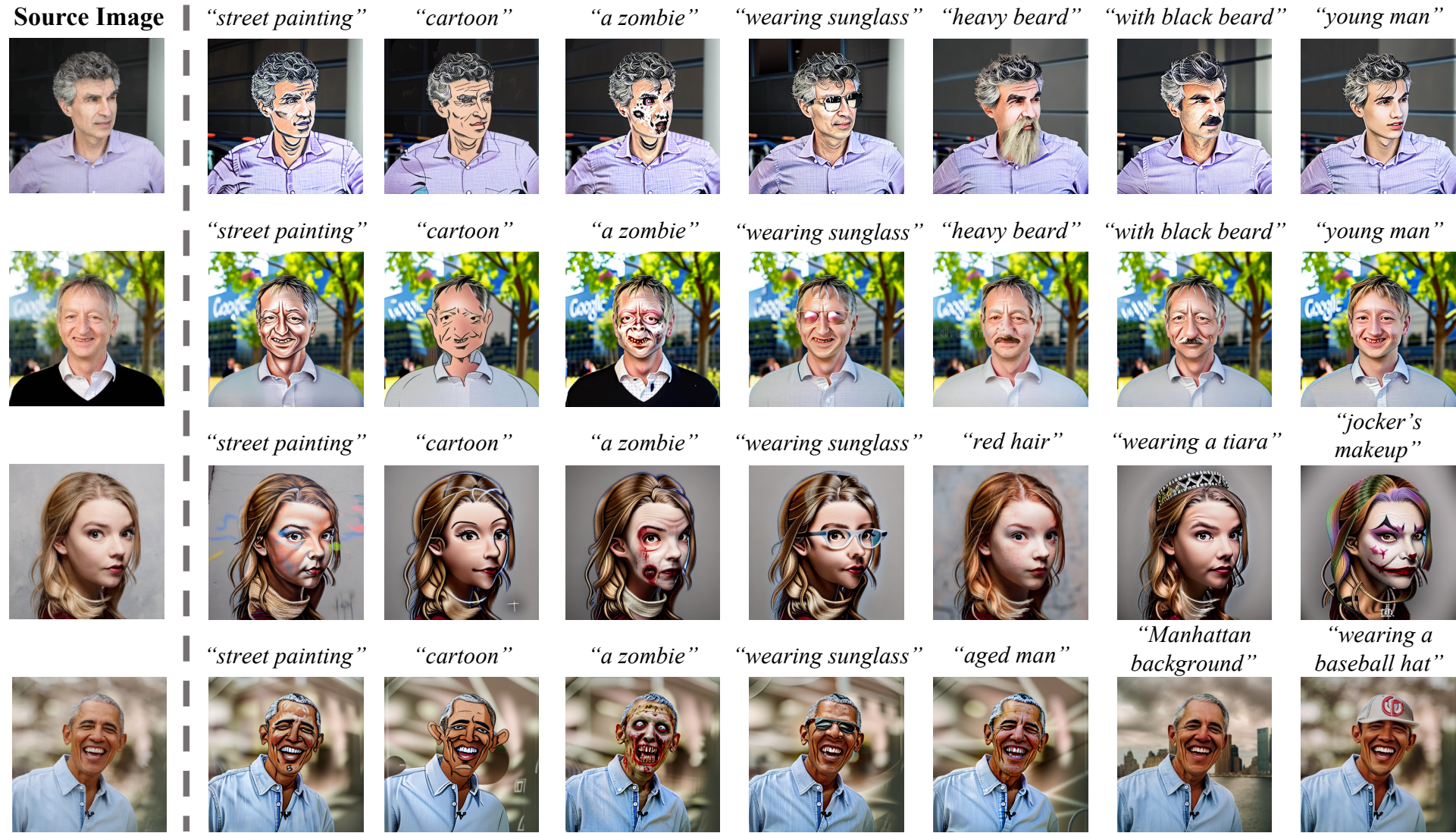
a) “*traditional Chinese painting of a lake*”(1024×512)

b) “*a watercolor painting of a lake of a lake*”(512×512)

c) “*an impressionism painting of a lake*”(512×512)

d) “*painting of a lake ... Vincent Van Gogh*”(512×1024)

# More editing tasks: Face editing



# More editing tasks: High-resolution

- “A children’s painting of a castle.”
- $H=768, W=1024$
- $K = 400, v = 0.7$



# More editing tasks: High-resolution

- “A painting of a castle in the style of Claude Monet.”
- $H=768, W=1024$
- $K = 400, v = 0.65$





# More editing tasks: High-resolution

- “A desert.”
- $H=768, W=1024$
- $K = 500, v = 0.8$



# More editing tasks: High-resolution

- “A desert.”
- $H=768, W=1024$
- $K = 500, v = 0.8$



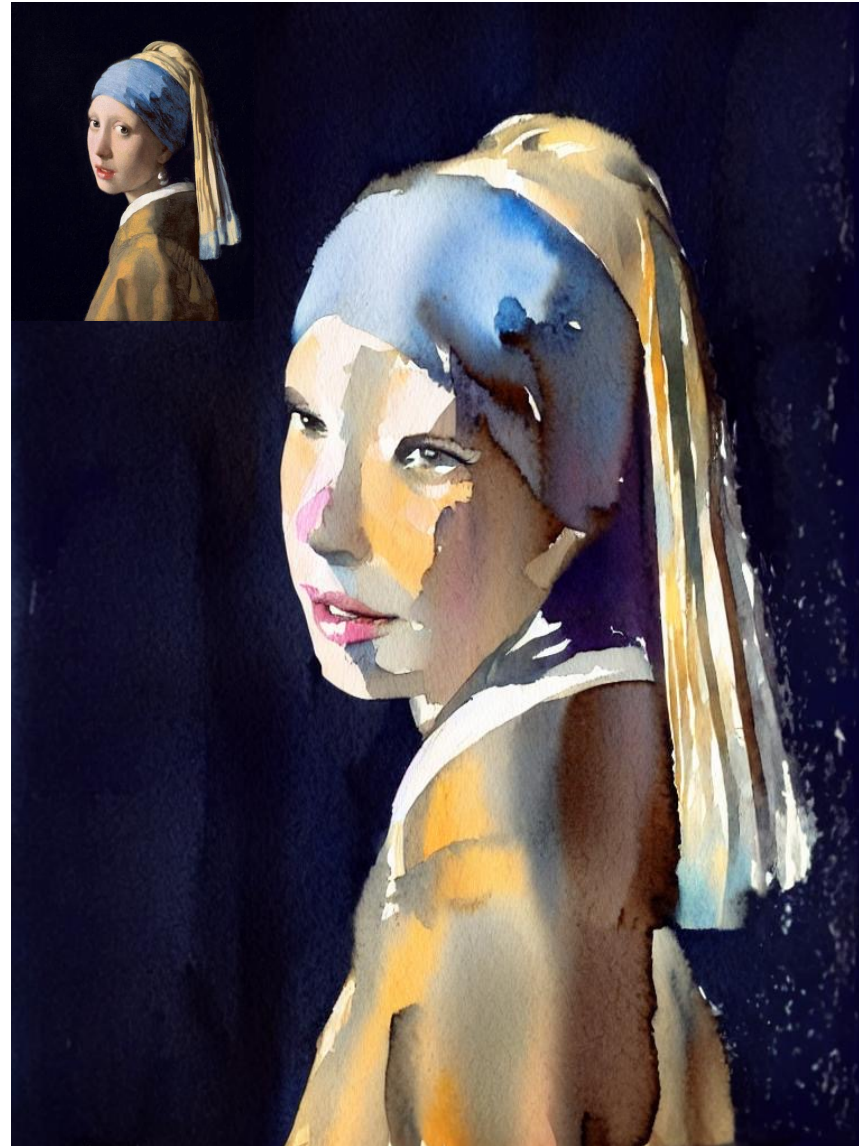
# More editing tasks: High-resolution

- “A photo of a lake with many sailboats.”
- $H=768, W=1024$
- $K = 400, v = 0.7$



# More editing tasks: High-resolution

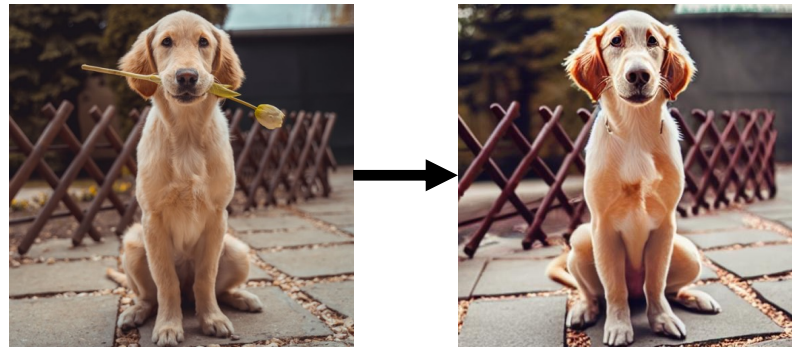
- “A watercolor painting of a girl.”
- $H=1024, W=768$
- $K = 400, v = 0.6$



# More editing tasks



Style Transfer



Content Removal



Style Generation



# Results: Comparison

Source Image



DreamBooth



Textual-Inversion



Ours (w/o pos)



Ours (w/ pos)



*"Painting of a castle in the style of Vincent Van Gogh" (H = 512, W = 512)*

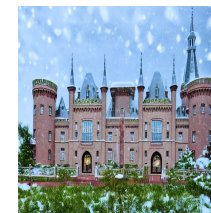
DreamBooth



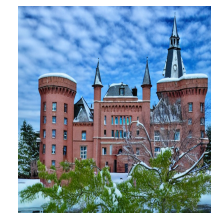
Textual-Inversion



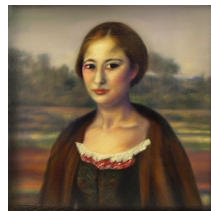
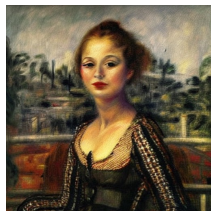
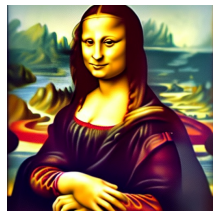
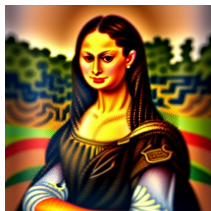
Ours (w/o pos)



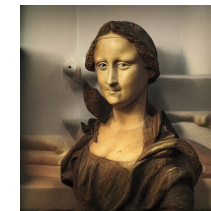
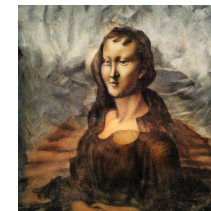
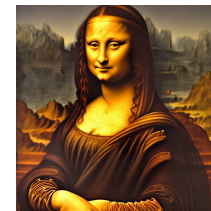
Ours (w/ pos)



*"A photo of a castle covered by snow" (H = 512, W = 768)*



*"Oil painting of a lady in the style of Pierre-Auguste Renoir" (H = 512, W = 512)*



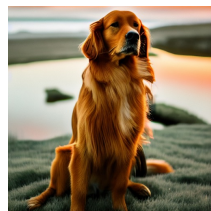
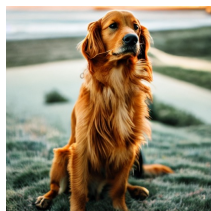
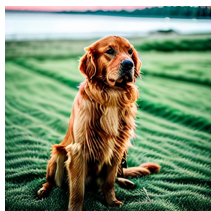
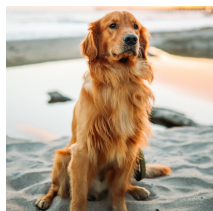
*"A sculpture of a lady in the style of Michelangelo" (H = 768, W = 768)*



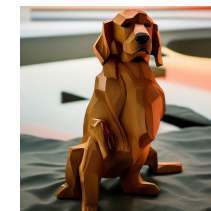
*"A children's drawing of a bench" (H = 512, W = 512)*



*"A bench in a desert" (H = 768, W = 768)*



*"A dog standing on grass" (H = 512, W = 512)*

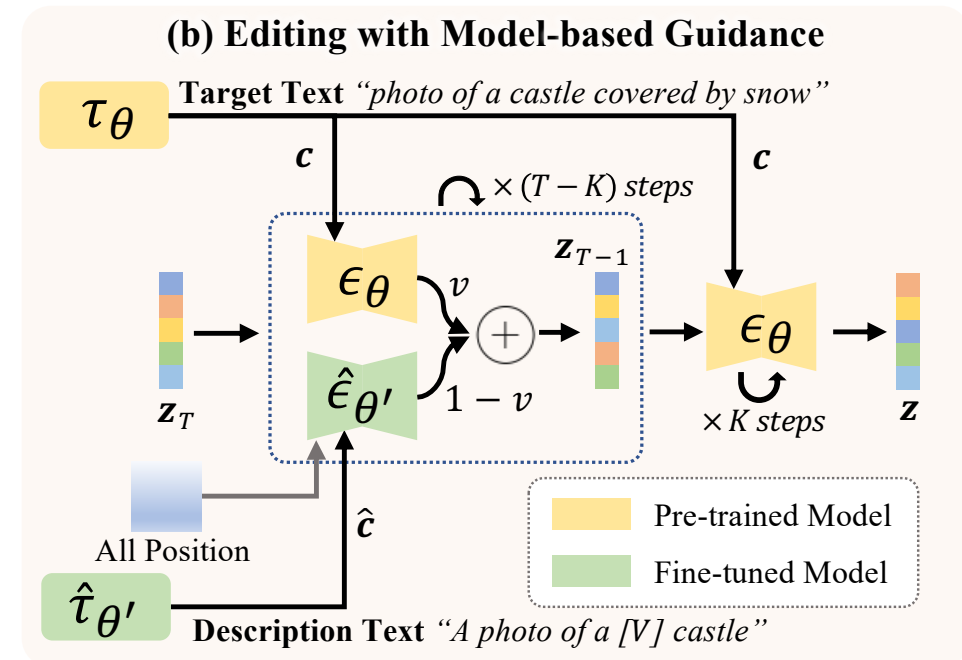
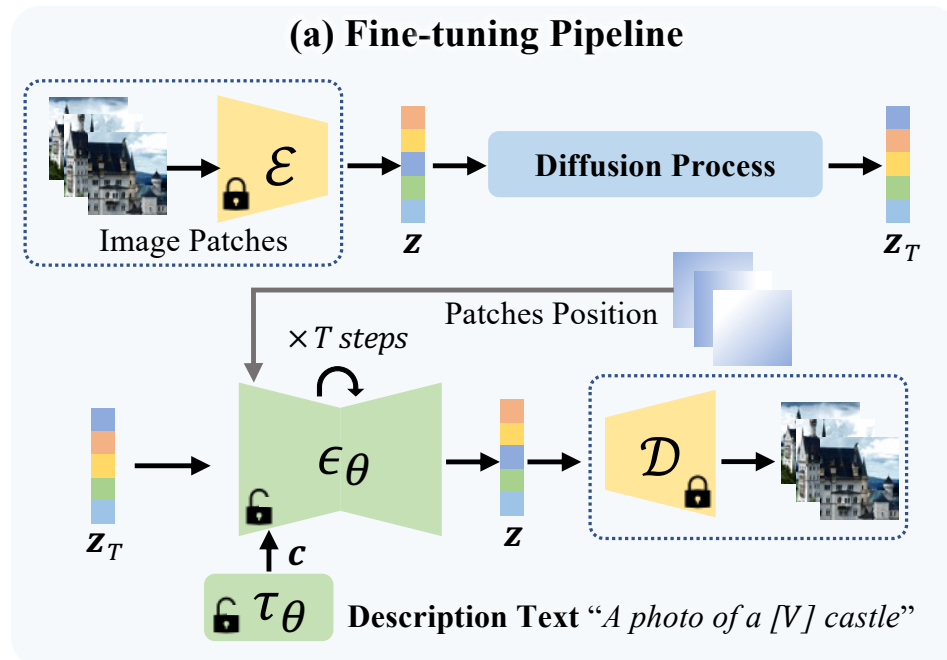


*"A sculpture of a dog in the style of Michelangelo" (H = 768, W = 768)*



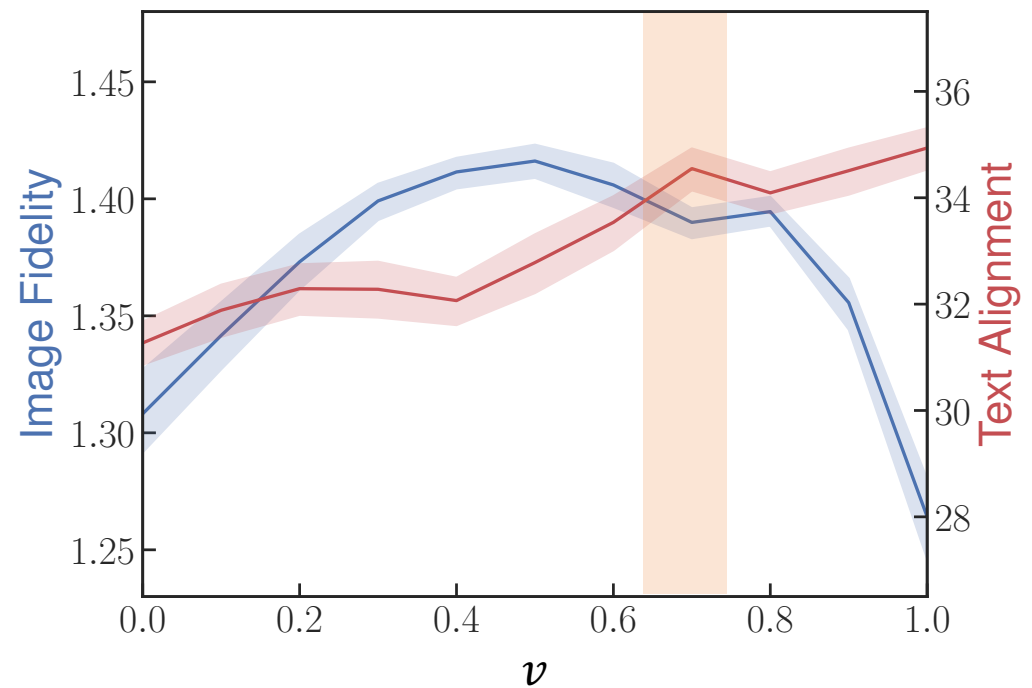
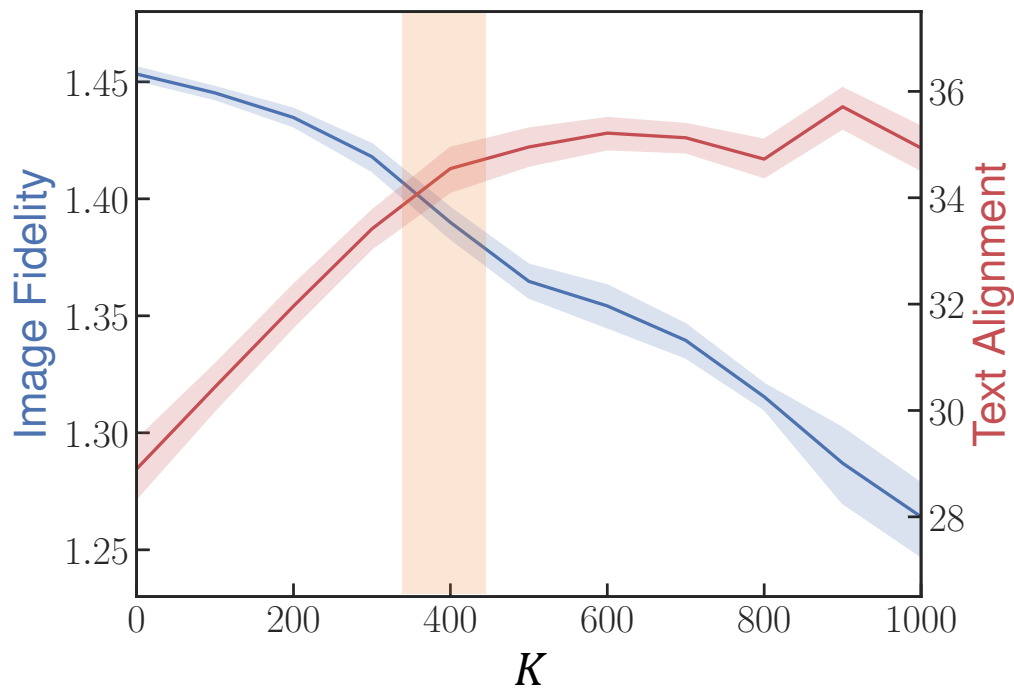
# Model-based Guidance

- $\tilde{\epsilon}_\theta(\mathbf{z}_t, \mathbf{c}) = w \left( v \epsilon_\theta(\mathbf{z}_t, \mathbf{c}) + (1 - v) \hat{\epsilon}_{\theta'}(\mathbf{z}_t, \hat{\mathbf{c}}) \right) + (1 - w) \epsilon_\theta(\mathbf{z}_t, \emptyset)$  for  $t > K$ 
  - $v$ : strength of editing
  - $K$ : perform guidance from  $t = 1000$  to  $t = K$
- decrease  $\downarrow v$  or  $\downarrow K \Rightarrow$  Increase  $\uparrow$  similarity between output sample & input image



# Analysis: Model-based guidance

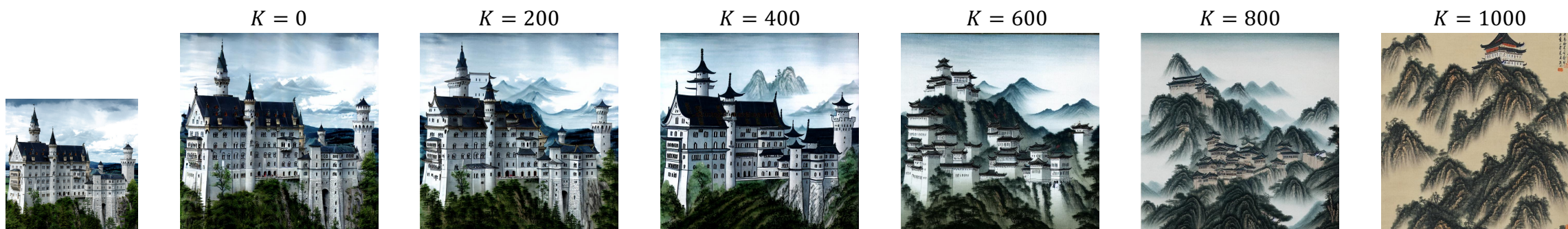
- Image fidelity: LPIPS score,  $1 - \log(\mathcal{L}_{LPIPS})$
- Text alignment: CLIP score  $\uparrow$ ,  $100 \times \text{cosine}(\text{text}, \text{image})$





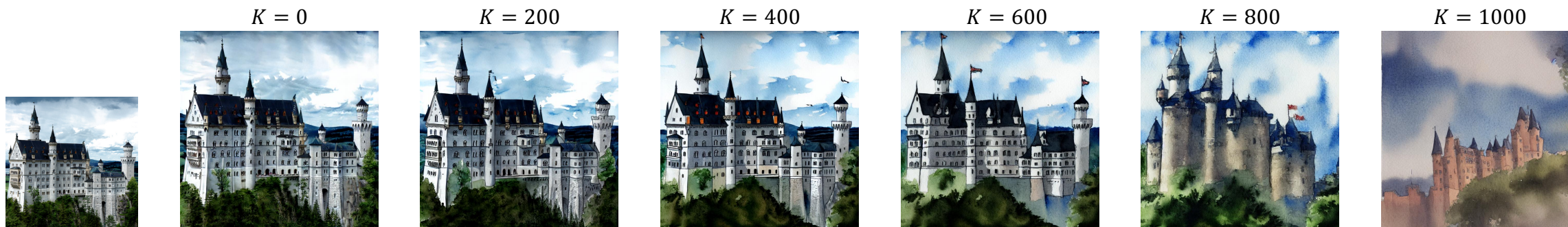
# Analysis: Guidance step $K$

- $v = 0.7$



Source Image

“a Chinese painting painting of a castle” ( $H = 768, W = 768$ )



Source Image

“a watercolor painting of a castle” ( $H = 768, W = 768$ )



# Analysis: Guidance step $K$

- $v = 0.7$

Source Image



$K = 0$



$K = 100$



$K = 200$



$K = 300$



$K = 400$



$K = 500$



$K = 600$



$K = 700$



$K = 800$



$K = 900$

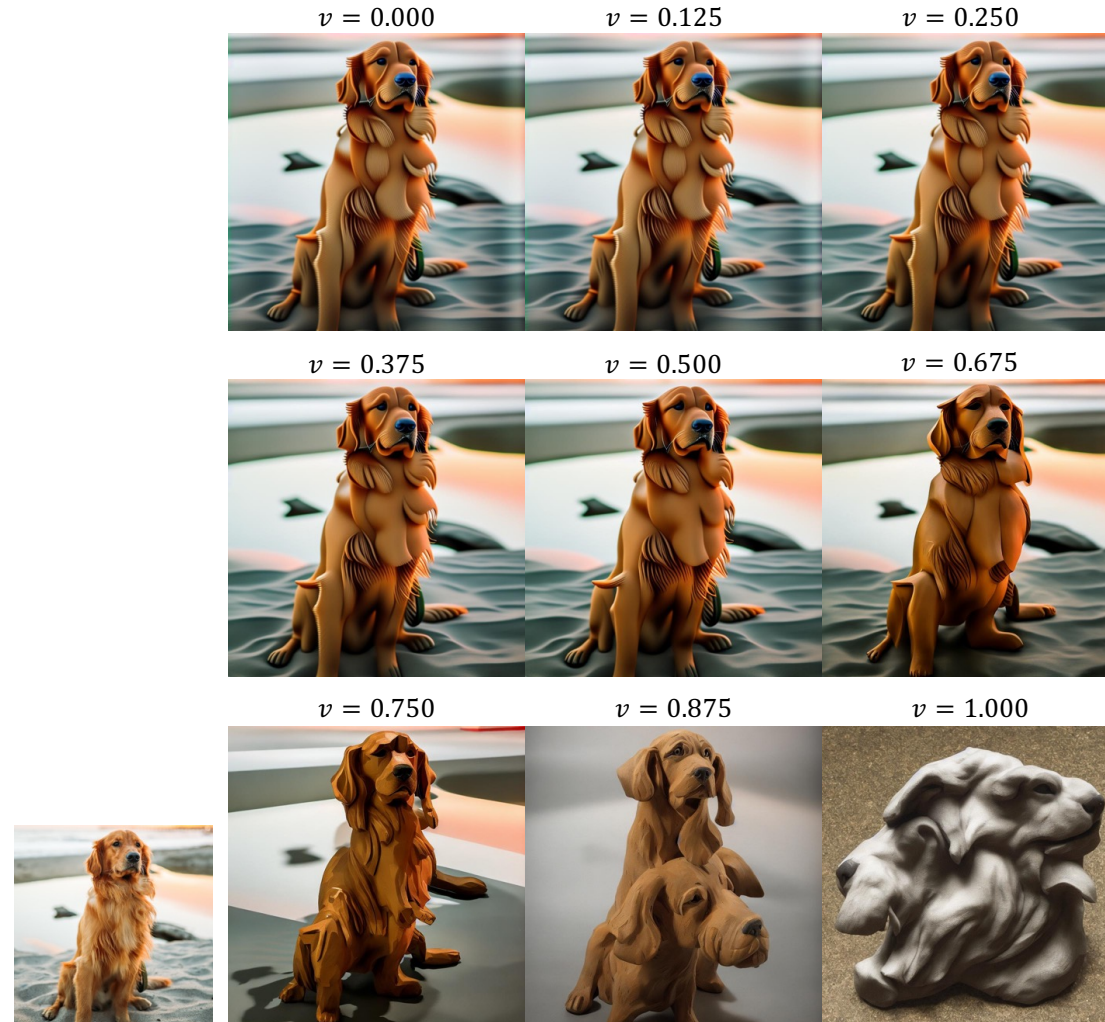


$K = 1,000$



# Analysis: Guidance weight $v$

- $K = 400$



Source Image

*"a sculpture of a dog in the style of Michelangelo"* ( $H = 768$ ,  $W = 768$ )



# Analysis: Guidance weight $v$

- $K = 400$

Source Image



$v = 0.0$



$v = 0.1$



$v = 0.2$



$v = 0.3$



$v = 0.4$



$v = 0.5$



$v = 0.6$



$v = 0.7$



$v = 0.8$



$v = 0.9$



$v = 1.0$



***Thanks!***

