

An Image is Worth One Word: Personalizing Text-to-Image Generation using Textual Inversion

Sangyun Lee



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it



panda mad scientist mixing sparkling chemicals, artstation

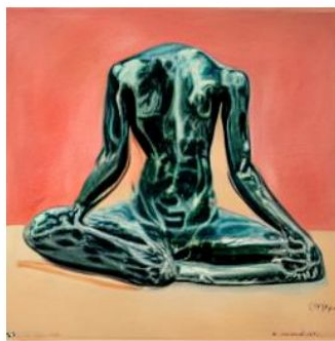


a corgi's head depicted as an explosion of a nebula

It is difficult to fully control the generative process using text prompt only.



→



Input samples $\xrightarrow{\text{invert}}$ “ S_* ”

“An oil painting of S_* ”

“App icon of S_* ”

“Elmo sitting in the same pose as S_* ”

“Crochet S_* ”



→



Input samples $\xrightarrow{\text{invert}}$ “ S_* ”

“Painting of two S_* fishing on a boat”

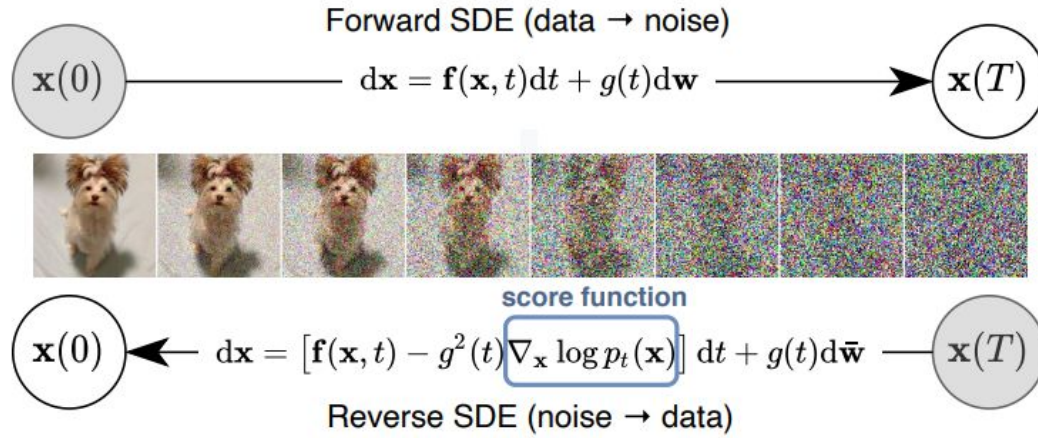
“A S_* backpack”

“Banksy art of S_* ”

“A S_* themed lunchbox”

cat doll wearing a vertical striped clothing, where stripes are black, yellow, red, and blue, and find this via optimization!

Diffusion models



Song et al., 2020

Ho et al., 2020

$$dx = -\frac{1}{2}\beta(t)x dt + \sqrt{\beta(t)} dw.$$

Denoising AEs are score estimators

$$\arg \min_{\theta} \mathbb{E}_i \{ \lambda(i) \mathbb{E}_{q_0(\mathbf{x})} \mathbb{E}_{q(\mathbf{x}_i|\mathbf{x}_0)} [\| \mathbf{s}_{\theta}(\mathbf{x}_i, i) - \nabla_{\mathbf{x}} \log q_i(\mathbf{x}_i|\mathbf{x}_0) \|_2^2] \},$$

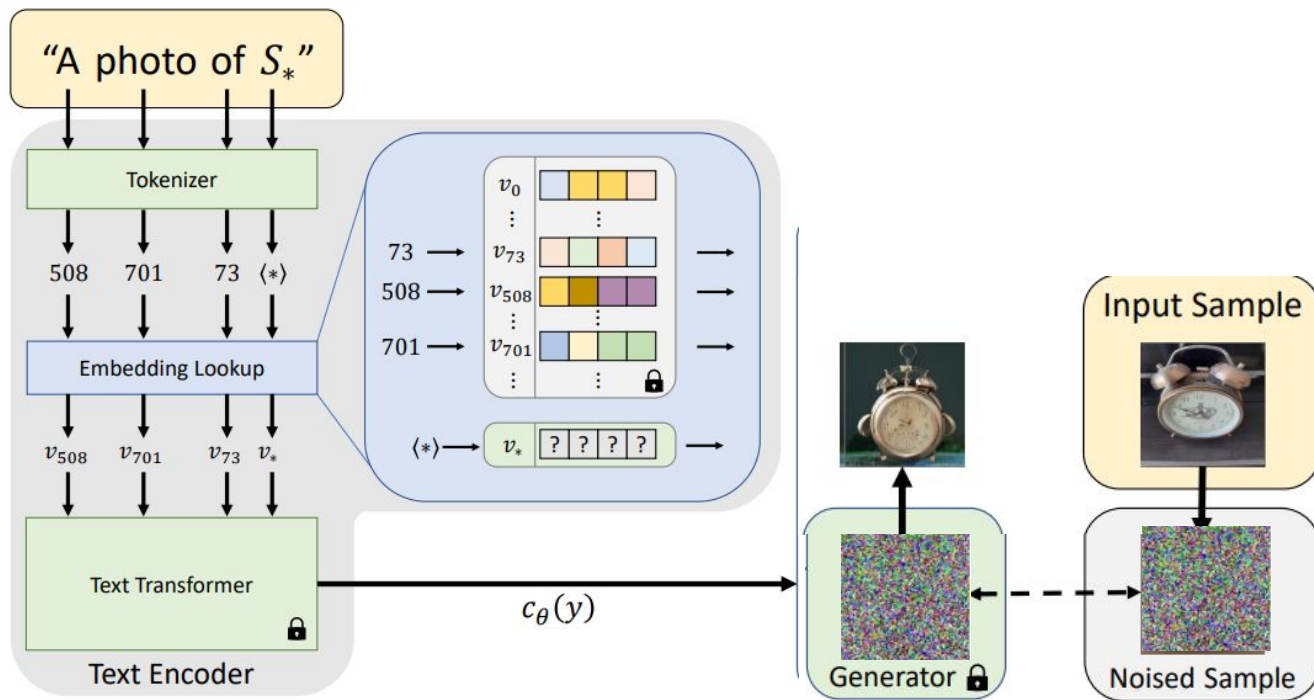
Denoising score matching
[Vincent, 2011]

Ho et al., 2020

$$x_i = \sqrt{\bar{\alpha}_i} x_0 + \sqrt{1 - \bar{\alpha}_i} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I)$$

$$q(x_i|x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_i} x_0, (1 - \bar{\alpha}_i) I)$$

$$\nabla_x \log q(x_i|x_0) = -\frac{\epsilon}{\sqrt{1 - \bar{\alpha}_i}}$$



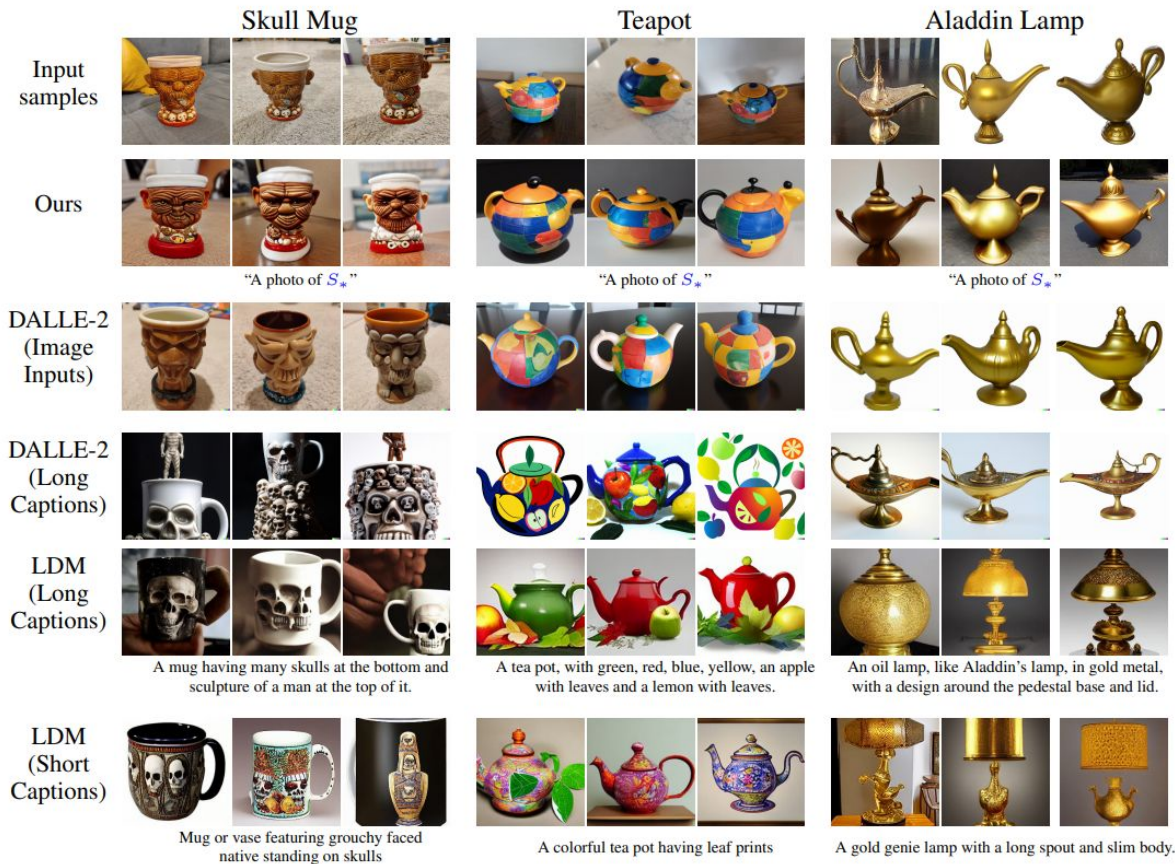


Figure 3: Object variations generated using our method, the CLIP-based reconstruction of DALLE-2 (Ramesh et al., 2022), and human captions of varying lengths. Our method generates variations which are typically more faithful to the original subject.



Input samples



"S* sports car"



"S* made of lego"



"S* onesie"



"da Vinci sketch of S*"



Input samples



"Watercolor painting of S* on a branch"



"A house in the style of S*"



"Grainy photo of S* in angry birds"



"S* made of chocolate"



Input samples



"A mosaic depicting S*"



"Death metal album cover featuring S*"



"Masterful oil painting of S* hanging on the wall"



"An artist drawing a S*"



Figure 6: The textual-embedding space can represent more abstract concepts, including styles. This allows us to discover words which can be used for style-guided generation. Image credits: @QinniArt (top), @David Revoy (bottom). Image reproduction authorized for non-commercial use only.



Sstyle



Sclock



Scat



Scraft



“Photo of *Sclock*
in the style of *Sstyle*”



“Photo of *Scat*
in the style of *Sstyle*”



“Photo of *Scraft*
in the style of *Sstyle*”



“Photo of *Sclock*
in the style of *Scat*”



“Photo of *Sclock*
in the style of *Scraft*”



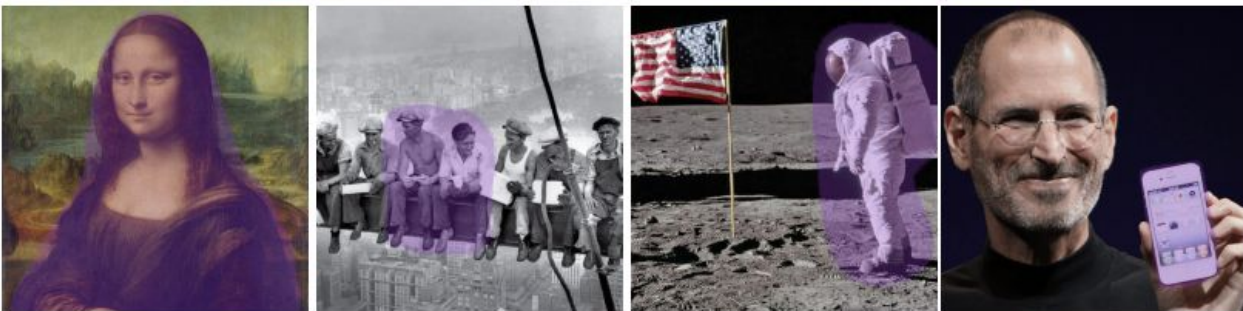
“Photo of *Scat*
in the style of *Scraft*”

Figure 7: Compositional generation using two learned pseudo-words. The model is able to combine the semantics of two concepts when using a prompt that combines them both. It is limited in its ability to reason over more complex relational prompts, such as placing two concepts side-by-side. Image credits: @QinniArt (left), @Leslie Manlapig (right). Reproductions authorized for non-commercial / non-print use respectively.

Input
Samples



Target Image
With Mask



Output
Image



“An oil painting
of S_* ”

“A black and white
photo of S_* ”

“A S_* ”

“A S_* ”

Figure 9: Our words can be used with downstream models that build on LDM. Here, we perform localized image editing using Blended Latent Diffusion (Avrahami et al., 2022a)