



# Continuous control of latent diffusion model through token embedding interpolation

Henry Lin, Pracha Promthaw, Nikhil Gautam Mentor: Dmitry Petrov

UMassAmherst

Manning College of Information & Computer Sciences

## Motivation

In this work, we are trying to investigate the encoded embedding space of certain attributes. To the best of our knowledge, no one has investigated the continuous changes of certain attributes before. Our goal is to change certain attributes continuously.

How does the realistic concept transform into Picasso?

Can we gradually change the level of realistic and Picasso in the following figures?

Realistic



Picasso



## Background

The previous work by (Brack et. al, 2022) introduced Stable Artist which can steer semantics in diffusion latent space. The idea is that by finding the direction of certain attributes, they can steer the semantics of the image into the attributes that they picked. However, the limitation of their approach is that it **doesn't have continuous control over the attributes**.

- Start with an unconditioned image
- find the latent vectors that correspond to different concepts
- manipulate the unconditioned latent space
- steering the unconditioned to conditioned latent space based on the concepts



Brack, Manuel & Schramowski, Patrick & Friedrich, Felix & Hintersdorf, Dominik & Kersting, Kristian. (2022). The Stable Artist: Steering Semantics in Diffusion Latent Space. 10.48550/arXiv.2212.06013.

## Methology

- Input: two texts that describes two different concepts.

We utilized the Stable Diffusion algorithm (Rombach et.al, 2021) and used the Clip model to encode the text, denoising loop, using a decoder from an auto-encoder to decode an embedding vector to an image. We pick two attributes that we want to interpolate and use the interpolation methods to shift between two concepts. The results show that the encoded embedding has the continuous attribute, so it is able to understand the continuous change in the vector space. This could be further adapted into continuous editing in latent vector space for more controllable image generation.

$\underline{E}$ : Clip encoder  
 $\underline{D}$ : Autoencoder's decoder  
Embedding: Convert the input sequence of tokens into a continuous representation.

$\vec{E}_1, \vec{E}_2$  are the encoded text embeddings  
 $\vec{t}_1, \vec{t}_2$  are the text embeddings  
 $\vec{E}_1, \vec{E}_2 \in \mathbb{R}^{768}$   
 $\vec{v} \in \mathbb{R}^{1 \times 4 \times 64 \times 64}$   
 $\text{image} \in \mathbb{R}^{512 \times 512}$

text1, text2  $\rightarrow$   $\vec{t}_1 = \text{embedding}(\text{tokenizer}(\text{text1}))$   
 $\vec{t}_2 = \text{embedding}(\text{tokenizer}(\text{text2}))$

Approach 1:

$\vec{t} = \text{Linear Interpolation}(\vec{t}_1) \rightarrow \vec{E} = E(\vec{t}) \rightarrow \vec{v} = \text{Denoising Loop}(\vec{E}) \rightarrow \text{image} = D(\vec{v})$

Approach 2:

$\vec{t}_1, \vec{t}_2 \rightarrow E_1 = E(\vec{t}_1) \ E_2 = E(\vec{t}_2) \rightarrow T = \text{Gradient Descent}(\vec{t}_1) \rightarrow \text{Non-Linear Interpolation } \vec{t} \subset T \rightarrow \vec{E} = E(\vec{t})$

Linear Interpolation:

$\vec{t} = \alpha t_1 + (1 - \alpha)t_2$

$\alpha \in [0, 1]$

Gradient Descent:

$\vec{E}_{\text{diff}} = \vec{E}_2 - \vec{E}_1$

$L = \sqrt{\vec{E}_{\text{diff}}(\vec{E}_{\text{diff}})^T}$

Algorithm 1 Gradient Descent

```

 $\vec{v}_1 = \text{Embedding}(\text{Tokenizer}(\text{text1}))$ 
 $\vec{\theta}_0 = \text{Encoder}(\vec{v}_1)$ 
 $\vec{v}_2 = \text{Embedding}(\text{Tokenizer}(\text{text2}))$ 
 $\vec{\gamma}_2 = \text{Encoder}(\vec{v}_2)$ 
curves = [ $\vec{\theta}_0$ ]
while not converged do
     $j = 1$ 
     $\vec{\gamma}_1 = \text{Encoder}(\vec{\theta}_{j-1})$ 
     $L = \sqrt{(\vec{\gamma}_1 - \vec{\gamma}_2)(\vec{\gamma}_1 - \vec{\gamma}_2)^T}$ 
     $\vec{\theta}_j = \vec{\theta}_{j-1} - \alpha \frac{\partial L}{\partial \vec{\theta}_{j-1}}$ 
    curves.append( $\vec{\theta}_j$ )
     $j = j + 1$ 
end while
return curves

```

image =  $D(\vec{v}) \leftarrow \vec{v} = \text{Denoising Loop}(\vec{E})$

Rombach, Robin and Blattmann, Andreas and Lorenz, Dominik and Esser, Patrick and Ommer, Björn. (2021). High-Resolution Image Synthesis with Latent Diffusion Models. <https://arxiv.org/abs/2112.10752>

## Conclusion/Discussion

From the results, we can continuously control attributes through two approaches: linear interpolation and non-linear interpolation. However, some change occurs suddenly and didn't depict a sign of continuous editing. For example, in figure 3, the middle figure didn't depict a continuous change. Also, we tried to form a Voronoi diagram on the text embedding space, but the dimension of it is too large. The future work of this project could be continuing to investigate the encoded latent space. For example, we can try to interpolate in the space that is formed by Voronoi Diagram or a space formed by the convex hull.

## Results

- Linear Interpolation: The alpha value controls the depiction of an attribute in the figure.

"A realistic drawing of a cat"  
 "A picasso drawing of a cat"

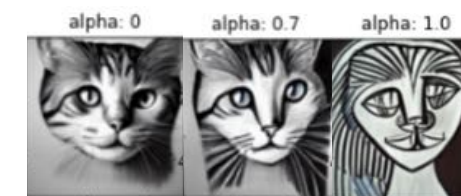


Figure 1: Alpha value of 0 means no Picasso and means realistic.

"A photo of a blue panda"  
 "A photo of a red panda"



Figure 2: Alpha value of 0 means no red and means blue.

"A photo of a white cat" => "A photo of a black cat"



Figure 3: Alpha value of 0 means no black and means black.

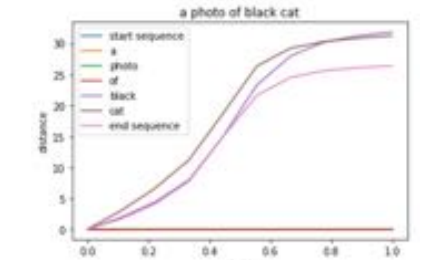


Figure 4: l2 embedding distance of each token

- Non-Linear Interpolation: Figure with a high loss value will have a depiction of attribute 2, and figure with a low loss value will have a depiction of attribute 1.

"A photo of winter landscape" => "A photo of spring landscape"

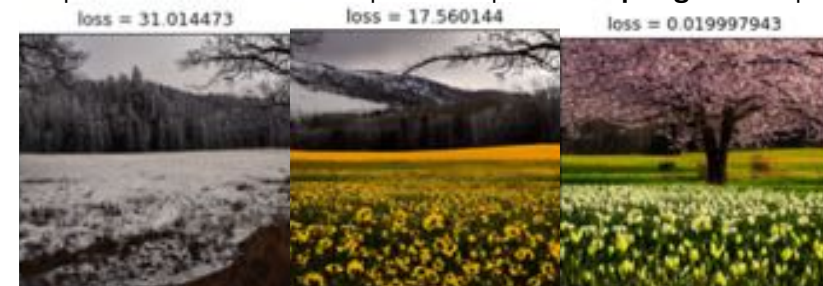


Figure 5: We can see as the loss decreases, the figure depicts more of the feature of spring

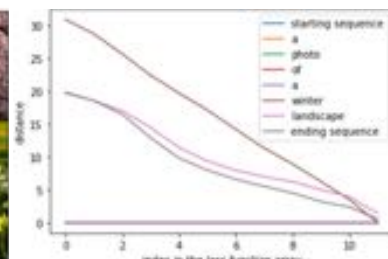


Figure 6: l2 embedding distance of each token

"A photo of Spooky castle" => "A photo of beautiful castle"



Figure 7: We can see as the loss decreases, the figure depicts more of the feature of beautiful.

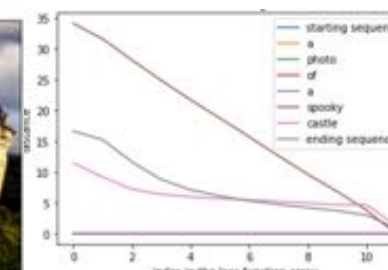


Figure 8: l2 embedding distance of each token