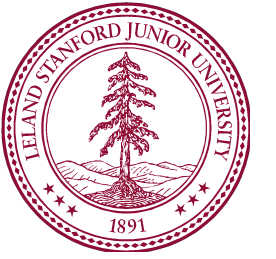


# Text-to-Image Generation



Aleksandr Timashov

June 4, 2021, CS231n Project

## Introduction

The **score function** of a distribution  $p(\mathbf{x})$  is defined as  $\nabla_{\mathbf{x}} \log p(\mathbf{x})$  and a model to approximate score function is called **score-based model**.

We should optimize **Fisher divergence**, using **score matching**:

$$\mathbb{E}_{p(\mathbf{x})} [\|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x})\|_2^2] = \int p(\mathbf{x}) \|\nabla_{\mathbf{x}} \log p(\mathbf{x}) - \mathbf{s}_{\theta}(\mathbf{x})\|_2^2 d\mathbf{x}.$$

**Key challenge:** Score function is not accurate in a low density regions.

**Solution:** Add noise to the data points and train score - based model on a noise perturbed data points.

Relation to **diffusion models**: Transition from finite number of noise scales to continuous scale. In this case we have **continuous-time stochastic process**.

Sample using **annealed Langevin dynamics**:  $\bar{\mathbf{x}}_{i+1} \leftarrow \bar{\mathbf{x}}_i + \frac{\alpha_i}{2} \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \sqrt{\alpha_i} \mathbf{z}_i$

## Problem Statement

The challenge of image generation **usually** accomplished with **adversarial learning**. However adversarial training often suffers from unstable training. I **propose** to use **score-based** generative models to solve this challenge.

**Input:** Text description.

**Output:** Image, based on the text description.

**Main idea** is based on Bayes' rule:  $\nabla_{\mathbf{x}} \log p(\mathbf{x} | \mathbf{y}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x})$ , where  $\mathbf{x}$  is an image, conditioned on text description  $\mathbf{y}$ .

## Dataset

I used **COCO Captions Dataset**, containing over 330,000 images and more than 1.5 million captions describing these images. The split for **train** / **validate** / **test** is default as follows:

- train: 120,000 images;
- validate: 5,000 images;
- test: 205,000 images;



## Methods

I will approach this problem in **three steps**:

- Unconditional **score-based model** trained on COCO dataset;
- **Image captioning model**, based on transformer architecture;
- Inverse problem solving for **controllable generation**;

**Score-based model details:**

- Resizing images to **32 x 32** size and keeping 3 channels;
- Geometric progression of added noises;
- **U-Net architecture** with the same dimension of input and output;
- **Loss function**:  $\frac{1}{L} \sum_{i=1}^L [\|\sigma_i \mathbf{s}_{\theta}(\mathbf{x}_i + \sigma_i \mathbf{z}_i, \sigma_i) + \mathbf{z}_i\|_2^2]$

**Image captioning model details:**

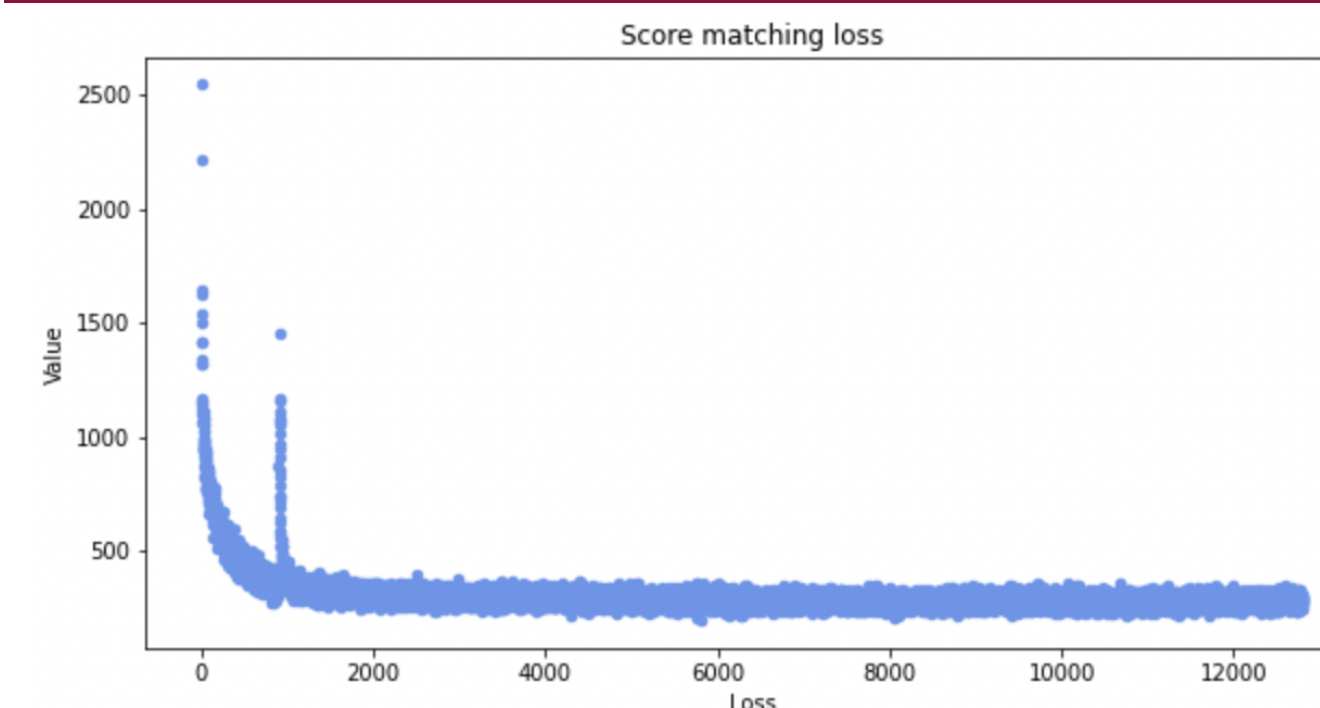
- Resizing images to 32 x 32 size and keep 3 channels;
- **GloVe** pretrained embeddings is used for text embeddings;
- I create a pair of time - dependent training data  $(\mathbf{x}(\mathbf{t}), \mathbf{y})$  by adding noise to  $(\mathbf{x}(\mathbf{0}), \mathbf{y})$ , sampled from COCO dataset;
- I use **2 headed self-attention** in a transformer-based encoder;
- To find score function, I use an idea from Bayes' rule:

$$\nabla_{\mathbf{x}} \log p(\mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_k | \mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{y}_1 | \mathbf{x}) + \sum_{i=2}^k \nabla_{\mathbf{x}} \log p(\mathbf{y}_i | \mathbf{y}_1, \mathbf{y}_2, \dots, \mathbf{y}_{i-1}, \mathbf{x})$$

**Sampling:**

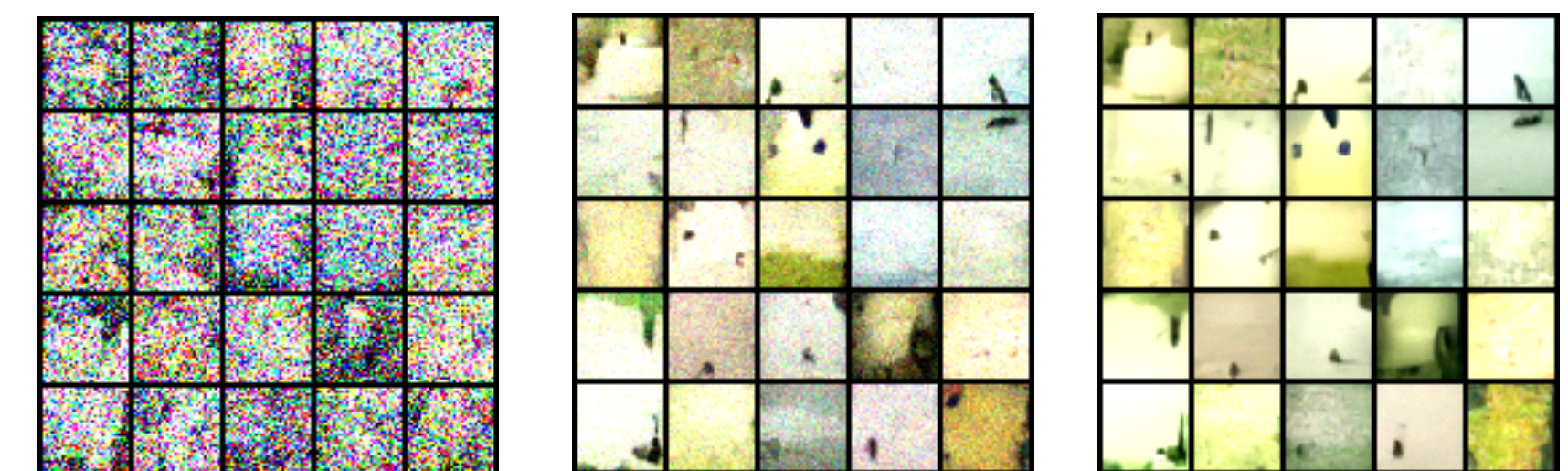
$$\bar{\mathbf{x}}_{i+1} \leftarrow \bar{\mathbf{x}}_i + \frac{\alpha_i}{2} [\nabla_{\mathbf{x}} \log p(\mathbf{x}) + \nabla_{\mathbf{x}} \log p(\mathbf{y} | \mathbf{x})] + \sqrt{\alpha_i} \mathbf{z}_i$$

## Experiments & Analysis



- As we can see, **loss** is converging. Around step number 1,000 Adam optimizer lost local minimum, but after that started converging again.

Below are **some results** starting from noise and converging to clear pictures.



## Conclusions & Future Work

My work shows the great opportunity to experiment **with image generation** using **score-based models**.

I experimented using different model architectures, including **DDPM** and **NCSN**.

Sampling was done using **Annealed Langevin dynamics**.

**Future Work:**

- Experiments with Conceptual Captions dataset;
- Experiments with different models architectures and different sampling techniques;
- Experiments with increased input size;

## References

Yang Song et al. Score-based generative modeling through stochastic differential equations, 2021.

Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. In Advances in Neural Information Processing Systems, pages 11895– 11907, 2019.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020.