# An Implementation of Denoising Diffusion Probabilistic Models

Yug Ajmera
Athrva Pandhare

Penn
Engineering

# Motivation

Recent breakthrough in Generative AI: Open AI's DALL-E 2 and GLIDE, Google's Imagen, and Stability AI's Stable Diffusion.

# Literature Survey



Sohl-Dickstein (2015)



Ho et al. (2020)

# Full Derivation

https://yainnoware.blogspot.com/2022/11/decoding-math-behind-diffusion-models.html

# UNet Architecture

1. This implementation is called the Conditional UNet, primarily because of the time embedding.
2. The architecture contains ResNet blocks near the bottleneck interleaved with attention blocks.
3. Group normalization is used before the Attention blocks.
4. Time embedding (also called the position embedding due to similarities with the position embedding in the Transformer) allows the neural network to produced outputs *conditioned* on the current timestep.

Penn Engineering

# Architecture

# Ablation study

| Schedule | Loss | Inception Score (IS) |
|----------|------|----------------------|
| Linear | L1 | 3.3707 |
| Cosine | L1 | 1.614 |
| Quadratic | L1 | 3.6521 |
| Sigmoid | L1 | 3.304 |
| Linear | L2 | 3.3808 |
| Cosine | L2 | 2.6078 |
| Quadratic | L2 | 3.4592 |
| Sigmoid | L2 | 3.3833 |
| Linear | Smooth-L1 | 3.4503 |
| Cosine | Smooth-L1 | 2.1523 |
| Quadratic | Smooth-L1 | 3.878 |
| Sigmoid | Smooth-L1 | 3.6121 |

Table 1. Results of the Ablation Study



Diffusion Loss

Penn Engineering

# Results

| Model | Inception Score (IS) |
|-------|----------------------|
| Real Test Set | 4.0885 |
| Huber + Quadratic (20 epochs) | 4.1842 |

Table 2. Results with test set

Penn Engineering

# Visualizations on MNIST Fashion dataset

# Extending the DDPM to Complex Dataset

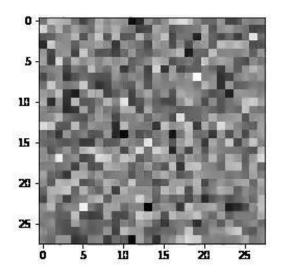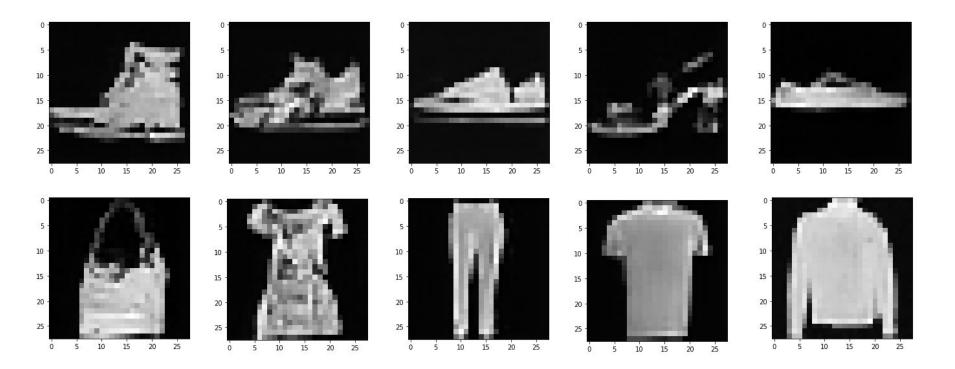We extend the implemented DDPM to a more complex dataset (Stanford Cars dataset).

**Model Changes :**

1.  Increased the depth of the UNet architecture
2.  Increased the number of channels in the UNet architecture.

**Implementation Details :**

1.  Used a linear schedule for beta.
2.  The loss between the predicted noise and the true noise was L2-loss.

# Results of Training on the Cars Dataset



We think that using a Larger model / using Latent Diffusion can produce better results.

Penn Engineering