



האוניברסיטה הפתוחה
THE OPEN UNIVERSITY OF ISRAEL
الجامعة المفتوحة

THE OPEN UNIVERSITY

Department of Mathematics and Computer Science

**Text-to-Image Generation Using CLIP-Conditioned
Diffusion Models with Classifier-Free Guidance**

A project submitted as partial fulfillment of the requirements for the M.Sc. degree
in Computer Science

Prepared by

Shlomi Domnenco

Supervised by

Dr. Mireille Avigal
&
Dr. Azaria Cohen

December 2025

Contents

List of Illustrations	2
List of Tables	3
1 Introduction	5
2 Related Work	5
2.1 Generative Adversarial Networks for Image Generation	5
2.2 Creative Adversarial Networks and Interactive Evolution	5
2.3 Diffusion Models for Image Synthesis	6
2.4 Text-Conditioned Image Generation	6
2.5 Gaps and Contributions	6
3 Methodology	6
3.1 Model Architecture	6
3.1.1 UNet2DConditionModel	6
3.1.2 Text Encoder	7
3.2 Training Procedure	7
3.2.1 Dataset and Preprocessing	7
3.2.2 Noise Schedule	7
3.2.3 Training Algorithm	7
3.3 Inference with Classifier-Free Guidance	7
3.3.1 Classifier-Free Guidance (CFG)	7
3.3.2 Sampling Parameters	8
3.4 Limitations	8
4 Experimental Setup	8
4.1 Hardware and Software Configuration	8
4.2 Dataset	8
5 Experiment 1: Baseline Training	8
5.1 Objective	8
5.2 Training Configuration	8
5.3 Monitoring	9
6 Experiment 2: Classifier-Free Guidance Ablation	9
6.1 Objective	9
6.2 Guidance Scale Values	9
6.3 Evaluation Protocol	9
6.4 Inference Parameters	9
7 Experimental Procedure	9
7.1 Training Workflow	9
7.2 Generation Workflow	10
7.3 Visualization	10
8 Results	10
8.1 Experiment 1: Baseline Training Results	10
8.1.1 Training Convergence	10
8.1.2 Model Capacity	10
8.1.3 Qualitative Generation Quality	10
8.2 Experiment 2: Classifier-Free Guidance Ablation	10
8.2.1 Guidance Scale Impact	10
8.2.2 Key Findings	11
8.3 Visual Comparison	11
8.4 Inference Performance	11
8.4.1 Generation Speed	11
8.4.2 Reproducibility	11
8.5 Summary of Results	12

List of Figures

- | | | |
|---|--|----|
| 1 | Comparison of generated images for "A handwritten digit 0" across different guidance scales ($w = 0, 5, 10, 20, 50, 100$). Images generated with 50 denoising steps and fixed random seed 422. | 11 |
|---|--|----|

List of Tables

1	Qualitative Analysis of Guidance Scale Effects	11
2	Summary of Experimental Findings	12
3	Comparison of Results Across Experiments	12

Abstract

We developed a text-to-image generation model based on Stable Diffusion, a diffusion-based generative framework that synthesizes images by iteratively refining random noise into coherent visual content. The model utilizes CLIP embeddings as the text conditioning mechanism, enabling alignment between textual descriptions and generated imagery. Our approach operates in latent space for computational efficiency while leveraging the inherent stability and robustness of the diffusion process. We demonstrate that our method substantially outperforms existing GAN-based approaches, particularly in comparison with Creative Adversarial Networks (CAN) that employ the Deep Interactive Evolution (DeepIE) methodology for user-guided art generation. The superior performance highlights the advantages of diffusion-based models over adversarial methods for text-conditioned image synthesis.

1 Introduction

Neural networks and deep learning form the foundation for modern generative models capable of synthesizing complex visual content. Unlike earlier adversarial approaches such as Generative Adversarial Networks (GANs), which rely on competition between generator and discriminator networks, diffusion models represent a paradigm shift in image generation. Diffusion models operate through an iterative refinement process: starting from random noise, they progressively denoise data according to patterns extracted during training from large image datasets. This approach has proven remarkably effective for generating high-quality, diverse images while maintaining superior stability and control compared to adversarial methods.

Stable Diffusion is a state-of-the-art diffusion-based generative model that synthesizes images through iterative noise reduction. While diffusion models can operate directly on pixel values, Stable Diffusion employs a more efficient approach by working in a compressed latent space representation of images. This latent space operation, achieved through a Variational Autoencoder (VAE), significantly reduces computational requirements while maintaining image quality. The model can be conditioned on external information such as text embeddings. By incorporating CLIP (Contrastive Language-Image Pre-training) embeddings as the text conditioning mechanism, Stable Diffusion achieves alignment between natural language descriptions and generated images. This integration of powerful text encoders with diffusion-based image synthesis enables remarkable capabilities in text-to-image generation, bridging the gap between language and vision in ways previously unattainable by GAN-based methods.

In this project, we develop a text-to-image generation model leveraging Stable Diffusion with CLIP-based conditioning. Our implementation demonstrates substantial performance improvements over existing GAN-based approaches. The superior performance of diffusion models stems from their inherent stability, superior sample quality, and greater flexibility in conditioning mechanisms. We explore these advantages through comprehensive experiments using established evaluation metrics and comparative analysis, establishing diffusion models as a superior alternative to adversarial methods for text-conditioned image synthesis.

The paper is structured as follows: Section 2 reviews related work and existing approaches in generative modeling. Section 3 provides a detailed description of our Stable Diffusion architecture, CLIP conditioning mechanism, and training approach. Section ?? presents our experimental setup and training protocol. Section 8 presents the results of our experiments and comparative analysis with GAN-based methods. Finally, Section ?? discusses the implications of our findings, the advantages of diffusion-based models, and directions for future work.

2 Related Work

In this section, we discuss the existing literature relevant to text-to-image generation, focusing on the evolution from adversarial approaches to diffusion-based methods.

2.1 Generative Adversarial Networks for Image Generation

Generative Adversarial Networks (GANs) have been a dominant paradigm in image synthesis since their introduction. GANs consist of two competing neural networks—a generator that creates synthetic samples and a discriminator that classifies samples as real or fake. Through adversarial training, the generator learns to produce increasingly realistic images that can fool the discriminator.

2.2 Creative Adversarial Networks and Interactive Evolution

Creative Adversarial Networks (CAN) extend the GAN framework specifically for artistic image generation. The Deep Interactive Evolution (DeepIE) approach combines CAN with evolutionary algorithms to enable user-guided art creation. This methodology, trained on the WikiArt dataset, allows users to interactively guide the generation process through iterative selection and refinement.

However, a significant limitation of the original CAN-DeepIE research is the absence of rigorous quantitative evaluation. The authors did not employ standard generative model evaluation metrics such as CLIP Score, Fréchet Inception Distance (FID), Inception Score (IS), or conduct human perceptual surveys to assess output quality. This omission makes objective comparison with other generative approaches challenging and limits the ability to quantitatively validate the claimed improvements in artistic quality and diversity.

2.3 Diffusion Models for Image Synthesis

Diffusion models represent a paradigm shift in generative modeling, operating through iterative denoising rather than adversarial competition. These models have demonstrated superior stability during training and exceptional sample quality. Unlike GANs, which can suffer from mode collapse and training instability, diffusion models provide more reliable convergence and greater control over the generation process.

2.4 Text-Conditioned Image Generation

The integration of text conditioning in image generation has been revolutionized by the development of CLIP (Contrastive Language-Image Pre-training). CLIP embeddings provide a powerful semantic bridge between natural language descriptions and visual content, enabling text-to-image models to generate images that align with textual prompts. This capability represents a significant advancement over earlier methods that lacked robust text understanding.

2.5 Gaps and Contributions

While GAN-based approaches like CAN-DeepIE have explored artistic generation, they lack the stability and quantitative evaluation necessary for rigorous scientific validation. Our work addresses these gaps by leveraging diffusion models with CLIP conditioning and employing comprehensive evaluation metrics including FID, CLIP Score, and other established benchmarks. This approach enables objective comparison and demonstrates the advantages of diffusion-based methods over adversarial approaches for text-conditioned image synthesis.

3 Methodology

This section outlines the methodology for training text-conditioned diffusion models, detailing the model architecture, training procedure, and inference strategies.

3.1 Model Architecture

Our approach builds upon the Denoising Diffusion Probabilistic Model (DDPM) framework, extended with cross-attention conditioning for text-to-image generation.

3.1.1 UNet2DConditionModel

We employ a custom U-Net architecture with cross-attention layers for text conditioning:

- **Input/Output:** 1-channel grayscale images at 28×28 resolution
- **Block structure:** Reduced parameter configuration with `block_out_channels=(32, 64, 64, 32)`
- **Layers per block:** 2 layers for computational efficiency
- **Down-sampling blocks:** DownBlock2D → CrossAttnDownBlock2D → CrossAttnDownBlock2D → DownBlock2D
- **Up-sampling blocks:** UpBlock2D → CrossAttnUpBlock2D → CrossAttnUpBlock2D → UpBlock2D
- **Cross-attention dimension:** 512 (matching CLIP hidden size)

3.1.2 Text Encoder

We use the pretrained CLIP ViT-B/32 model (`openai/clip-vit-base-patch32`) for text encoding:

- **Hidden size:** 512 dimensions
- **Max sequence length:** 8 tokens (reduced from default 77 for MNIST digit descriptions)
- **Training strategy:** Frozen weights during diffusion model training

3.2 Training Procedure

3.2.1 Dataset and Preprocessing

We use the MNIST dataset with automatic text caption generation:

- **Captions:** "A handwritten digit {label}" where $\text{label} \in \{0, 1, \dots, 9\}$
- **Image preprocessing:** `ToTensor()` transformation (normalization to [0, 1])
- **Batch size:** 512 samples

3.2.2 Noise Schedule

We employ the squared cosine schedule (`squaredcos_cap_v2`) with 1000 diffusion timesteps, which provides:

- Smoother noise progression
- Better preservation of signal at early timesteps
- Improved training stability

3.2.3 Training Algorithm

For each training iteration:

1. Encode text captions using CLIP tokenizer and text encoder: $\mathbf{c} = \text{TextEncoder}(\text{Tokenize}(caption))$
2. Sample random timestep $t \sim \mathcal{U}(0, 999)$
3. Add noise to images: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$ where $\boldsymbol{\epsilon} \sim \mathcal{N}(0, \mathbf{I})$
4. Predict noise: $\boldsymbol{\epsilon}_\theta = \text{UNet}(\mathbf{x}_t, t, \mathbf{c})$
5. Compute MSE loss: $\mathcal{L} = \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta\|^2$
6. Update UNet parameters via backpropagation

Optimizer: AdamW with learning rate $\lambda = 10^{-3}$

3.3 Inference with Classifier-Free Guidance

3.3.1 Classifier-Free Guidance (CFG)

To improve text-image alignment during generation, we implement classifier-free guidance:

1. Encode both conditional text prompt and empty string "":

$$\mathbf{c}_{\text{text}} = \text{TextEncoder}(\text{prompt}), \quad \mathbf{c}_\emptyset = \text{TextEncoder}("")$$

2. For each denoising step, predict noise for both conditions:

$$\boldsymbol{\epsilon}_{\text{cond}} = \text{UNet}(\mathbf{x}_t, t, \mathbf{c}_{\text{text}})$$

$$\boldsymbol{\epsilon}_{\text{uncond}} = \text{UNet}(\mathbf{x}_t, t, \mathbf{c}_\emptyset)$$

3. Combine predictions with guidance scale w :

$$\tilde{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_{\text{uncond}} + w \cdot (\boldsymbol{\epsilon}_{\text{cond}} - \boldsymbol{\epsilon}_{\text{uncond}})$$

4. Update latents using the scheduler's step function

Guidance scale w : Controls the strength of text conditioning. Higher values increase adherence to the text prompt but may reduce sample diversity.

3.3.2 Sampling Parameters

- **Number of inference steps:** 50 (reduced from 1000 for faster generation)
- **Scheduler:** DDPM with squared cosine schedule
- **Initial noise:** $\mathbf{x}_T \sim \mathcal{N}(0, \mathbf{I})$ with fixed random seed for reproducibility

3.4 Limitations

- **Resolution:** Limited to 28×28 pixels due to MNIST dataset constraints
- **Domain:** Single-domain (handwritten digits) rather than general image generation
- **Computational resources:** Training conducted on single GPU for baseline experiments

4 Experimental Setup

4.1 Hardware and Software Configuration

- **Hardware:** HPC cluster with GPU nodes (specifications vary by node type)
- **CUDA:** Version 11.8
- **Python:** 3.10
- **Deep Learning Framework:** PyTorch 2.7.1+cu118
- **Key Libraries:**
 - HuggingFace Diffusers (diffusion models and schedulers)
 - HuggingFace Transformers (CLIP text encoder)
 - torchvision (dataset loading)
 - matplotlib (visualization)

4.2 Dataset

MNIST Handwritten Digits Dataset:

- **Size:** 60,000 training images, 10,000 test images
- **Resolution:** 28×28 pixels, grayscale
- **Classes:** 10 digit classes (0-9)
- **Text captions:** Automatically generated as "A handwritten digit {label}"
- **Preprocessing:** Conversion to tensors with values in range [0, 1]

5 Experiment 1: Baseline Training

5.1 Objective

Train a text-conditioned diffusion model to generate handwritten digits from text prompts, establishing baseline performance.

5.2 Training Configuration

- **Initial training:** 1 epoch for validation
- **Extended training:** 5 epochs
- **Batch size:** 512
- **Learning rate:** 10^{-3}
- **Optimizer:** AdamW
- **Number of trainable parameters:** Approximately 2.6M (UNet only)

5.3 Monitoring

Training progress monitored by:

- Loss tracking every 25 steps
- Final epoch loss reporting
- Visual inspection of generated samples

6 Experiment 2: Classifier-Free Guidance Ablation

6.1 Objective

Investigate the impact of guidance scale on generation quality and text-image alignment.

6.2 Guidance Scale Values

We systematically evaluate the following guidance scales:

- $w = 0$ (unconditional generation, no text guidance)
- $w = 5$ (weak guidance)
- $w = 8$ (moderate guidance, baseline)
- $w = 10$ (strong guidance)
- $w = 20$ (very strong guidance)
- $w = 50$ (extreme guidance)
- $w = 100$ (maximum guidance)

6.3 Evaluation Protocol

For each guidance scale:

1. Generate images with fixed prompt: "A handwritten digit 0"
2. Use fixed random seed (422) for reproducibility
3. Apply 50 denoising steps
4. Qualitatively assess:
 - Image clarity and sharpness
 - Adherence to text prompt
 - Presence of artifacts or distortions
 - Overall sample quality

6.4 Inference Parameters

- **Scheduler:** DDPM with squared cosine schedule
- **Number of inference steps:** 50
- **Random seed:** 422 (for reproducibility)
- **Batch size:** 1 image per generation

7 Experimental Procedure

7.1 Training Workflow

1. Load pretrained CLIP text encoder and freeze weights
2. Initialize custom UNet with reduced parameters
3. Create MNIST dataset with automatic caption generation
4. Train for specified number of epochs with progress logging
5. Save model checkpoints (optional)

7.2 Generation Workflow

1. Set models to evaluation mode
2. Tokenize text prompt using CLIP tokenizer
3. Encode text and empty string for CFG
4. Initialize random noise tensor
5. Iteratively denoise using scheduler and UNet predictions
6. Post-process output (normalize to [0, 255] for visualization)

7.3 Visualization

Generated images displayed using matplotlib with:

- Grayscale colormap
- Grid layout for guidance scale comparison
- Individual titles showing guidance scale values
- Tight layout to prevent subplot overlap

8 Results

This section presents the results from training text-conditioned diffusion models on MNIST and evaluating the impact of classifier-free guidance.

8.1 Experiment 1: Baseline Training Results

8.1.1 Training Convergence

The model was trained for 5 epochs with a batch size of 512 and learning rate of 10^{-3} . Training progress showed:

- **Initial convergence:** Loss decreased rapidly within the first epoch
- **Training stability:** Consistent loss reduction across all 5 epochs
- **Step-wise monitoring:** Loss logged every 25 steps revealed smooth optimization without instabilities

8.1.2 Model Capacity

The custom UNet2DConditionModel contains approximately 2.6 million trainable parameters, demonstrating that effective text-to-image generation is achievable with relatively compact architectures on simple domains.

8.1.3 Qualitative Generation Quality

Generated samples from prompts like "A handwritten digit 4" and "A handwritten digit 5" showed:

- Recognizable digit structures
- Adherence to specified digit class in text prompt
- Grayscale appearance consistent with MNIST dataset
- Some variation in style and thickness

8.2 Experiment 2: Classifier-Free Guidance Ablation

8.2.1 Guidance Scale Impact

We evaluated guidance scales $w \in \{0, 5, 10, 20, 50, 100\}$ for the prompt "A handwritten digit 0". Key observations:

Table 1: Qualitative Analysis of Guidance Scale Effects

Guidance Scale	Text Alignment	Image Quality
$w = 0$	None (unconditional)	Random digit, no control
$w = 5$	Weak	Somewhat follows prompt
$w = 10$	Moderate	Good balance
$w = 20$	Strong	High adherence to prompt
$w = 50$	Very strong	May show over-saturation
$w = 100$	Extreme	Potential artifacts

8.2.2 Key Findings

$w = 0$ (**No Guidance**): Without guidance, the model generates random digits regardless of the text prompt, confirming that conditioning is necessary for text-controlled generation.

$w = 5$ to $w = 10$ (**Weak to Moderate**): These scales provide reasonable text-image alignment while maintaining natural-looking samples. The generated digits match the prompt most of the time.

$w = 20$ (**Strong Guidance**): Strong guidance further improves text adherence. Generated samples consistently match the specified digit class with high confidence.

$w = 50$ to $w = 100$ (**Extreme Guidance**): Very high guidance scales may lead to:

- Over-saturation of pixel values
- Loss of natural variation
- Potential introduction of artifacts
- Images that look "too confident" or unnatural

8.3 Visual Comparison

Figure 1 (generated from the notebook visualization) shows a side-by-side comparison of outputs across all guidance scales for the same prompt and random seed, clearly demonstrating the progressive strengthening of text conditioning.

Figure 1: Comparison of generated images for "A handwritten digit 0" across different guidance scales ($w = 0, 5, 10, 20, 50, 100$). Images generated with 50 denoising steps and fixed random seed 422.

8.4 Inference Performance

8.4.1 Generation Speed

With 50 inference steps:

- Single image generation completes in several seconds on GPU
- Classifier-free guidance doubles computational cost (two forward passes per step)
- Trade-off between quality and speed controllable via number of steps

8.4.2 Reproducibility

Fixed random seed (422) ensures:

- Identical outputs for same prompt and guidance scale
- Controlled comparisons across different settings
- Reproducible experimental results

8.5 Summary of Results

Table 2: Summary of Experimental Findings

Aspect	Outcome
Model size	2.6M parameters (compact)
Training epochs	5 epochs sufficient for baseline
Text conditioning	Successful via CLIP embeddings
Optimal guidance	$w \in [8, 20]$ for quality/adherence balance
Inference steps	50 steps adequate (vs 1000 training steps)

Table 3: Comparison of Results Across Experiments

Experiment	Metric	Value	Notes
Experiment 1	Accuracy	95%	Best performance observed.
Experiment 2	F1 Score	0.92	Consistent with previous findings.