

**Université Ibn Tofail**

Master Informatique et Intelligence Artificielle

# **LARGO: Flower Image Generation Using a Denoising Diffusion Model**

**Réalisé par :**

Zaher Yassin  
Boulaouane Salaheddine

**Encadré par :**

Pr. Tarik HOUICHIME

# Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
1.1	The Evolution of Generative Synthesis . . . . .	3
1.2	The LARGO Initiative . . . . .	4
1.3	Project Objectives . . . . .	4
<b>2</b>	<b>State of the Art</b>	<b>4</b>
2.1	Foundations of Generative Synthesis . . . . .	4
2.1.1	The Probabilistic Era: Explicit Density Estimation . . . . .	4
2.1.2	The Adversarial Shift: Implicit Learning . . . . .	5
2.2	The Diffusion Unification (DDPM & DDIM) . . . . .	6
2.3	Current Benchmarks on Oxford-102 . . . . .	6
<b>3</b>	<b>Theoretical Framework</b>	<b>7</b>
3.1	The Forward Diffusion Process ( $q$ ) . . . . .	7
3.2	The Reverse Generative Process ( $p_\theta$ ) . . . . .	8
3.3	The Optimization Objective . . . . .	9
<b>4</b>	<b>Methodology</b>	<b>9</b>
4.1	Dataset Analysis: The Oxford 102 Manifold . . . . .	9
4.1.1	Botanical and Visual Statistics . . . . .	9
4.1.2	The Class Imbalance Challenge . . . . .	10
4.1.3	Data Splitting and Preprocessing . . . . .	10
4.2	Conditional Implementation Strategy . . . . .	10
4.2.1	Evolution of Conditioning: Classifier vs. Classifier-Free . . . . .	11
4.2.2	Theoretical Derivations of CFG . . . . .	11
4.2.3	Implementation of Embeddings . . . . .	12
4.3	Architectural Specifications: The LARGO U-Net . . . . .	12
4.3.1	Macro-Architecture Design . . . . .	12
4.3.2	The Residual Block (ResBlock) . . . . .	13
4.3.3	Attention Mechanisms . . . . .	14
4.4	Training Dynamics and Algorithm . . . . .	15
4.4.1	Hyperparameter Configuration . . . . .	15
4.4.2	The Training Algorithm . . . . .	15
4.4.3	Challenges in Training . . . . .	16
4.5	Evaluation Methodology . . . . .	16

4.5.1	Metric Selection Strategy . . . . .	16
4.5.2	Fréchet Inception Distance (FID) . . . . .	17
4.5.3	Kernel Inception Distance (KID) . . . . .	17
<b>5</b>	<b>Results and Analysis</b>	<b>17</b>
5.1	Quantitative Benchmark . . . . .	17
5.2	Analysis of Quantitative Results . . . . .	18
5.2.1	Interpretation of KID and FID Scores . . . . .	18
5.2.2	The "Variance Gap" (Why FID is High) . . . . .	18
5.3	Ablation Study: Sampling Efficiency . . . . .	19
5.4	Qualitative Analysis . . . . .	19
5.4.1	Success Cases . . . . .	19
5.4.2	Limitations . . . . .	19
<b>6</b>	<b>Conclusion and Future Directions</b>	<b>19</b>

## Abstract

The synthesis of high-fidelity natural images remains a cornerstone challenge in computer vision, particularly in fine-grained domains characterized by high intra-class variance [11]. This report introduces LARGO, a specialized Denoising Diffusion Probabilistic Model (DDPM) framework [2] tailored for the Oxford-102 Flower dataset [3, 4]. Unlike Generative Adversarial Networks (GANs), which are often challenged by training instability and mode collapse [13, 12], LARGO leverages the thermodynamic-inspired principles of diffusion to iteratively refine Gaussian noise into coherent botanical structures.

We provide a comprehensive theoretical and practical analysis of the LARGO system, detailing the mathematical derivations of the forward and reverse diffusion processes with a specific focus on conditional generation mechanisms. We analyze the architectural modifications required for the U-Net backbone to accommodate class-conditional embeddings, ensuring semantic fidelity across 102 distinct flower species. Furthermore, we present a rigorous evaluation framework utilizing established metrics such as Fréchet Inception Distance (FID) and Kernel Inception Distance (KID) [10], situating LARGO’s performance within the broader landscape of state-of-the-art generative models. Our findings demonstrate that while diffusion models incur higher inference costs, they offer superior diversity and structural coherence compared to adversarial baselines.



Figure 1: **LARGO Diversity.** Uncurated random samples generated at  $64 \times 64$  resolution.

## 1 Introduction

### 1.1 The Evolution of Generative Synthesis

The domain of generative artificial intelligence has undergone a seismic shift over the last decade. For years, the field was bifurcated into two dominant paradigms: Variational Autoencoders (VAEs), which offered probabilistic rigor but often produced blurry, distinct samples, and Generative Adversarial Networks (GANs), which generated sharp images but were notoriously difficult to train due to the adversarial minimax game. GANs often suffered from mode collapse, where the generator would over-optimize for a specific subset of the data distribution, ignoring the long tail of diverse samples.

In this context, Denoising Diffusion Probabilistic Models (DDPMs), introduced by Ho et al. [1] and further refined by Song et al., have emerged as a third, powerful paradigm. Inspired by non-equilibrium thermodynamics, these models define a generative process as the reversal of a stochastic diffusion chain that gradually destroys data structure. By learning to reverse this entropy-increasing process, diffusion models can construct complex data distributions from simple isotropic Gaussian noise.

## 1.2 The LARGO Initiative

LARGO represents a targeted application of this diffusion technology to the domain of botanical illustration and synthesis. The project encapsulates the dual goal of the research: to master the algorithmic latent space of diffusion models and to apply this mastery to the biological complexity of the Oxford 102 dataset.

The generation of floral imagery presents unique challenges distinct from face generation (e.g., CelebA) or general object generation (e.g., CIFAR-10). Flowers exhibit:

1. **Geometric Complexity:** The radial symmetry of Asteraceae versus the bilateral symmetry of Orchidaceae requires a model with robust geometric understanding.
2. **Texture Variance:** The velvety texture of a rose petal differs fundamentally from the waxy surface of a lily.
3. **Contextual Noise:** The Oxford 102 dataset [3] includes flowers in the wild, meaning the model must distinguish between the subject (flower) and the high-frequency noise of the background (grass, leaves, dirt).

## 1.3 Project Objectives

This report aims to document the complete lifecycle of the LARGO project, addressing the following core research questions:

- **Theoretical Adaptation:** How can the standard unconditional DDPM formulation be rigorously extended to support class-conditional generation for 102 distinct categories?
- **Architectural Efficacy:** What modifications to the U-Net architecture—specifically regarding attention mechanisms and channel depth—are necessary to capture the high-frequency details of floral morphology at  $64 \times 64$  resolution?
- **Metric Evaluation:** How does the probabilistic output of LARGO compare to existing GAN benchmarks on the Oxford 102 dataset using standard metrics like FID?

# 2 State of the Art

## 2.1 Foundations of Generative Synthesis

The advancement of generative computer vision has progressed through distinct methodological eras, evolving from explicit density estimation towards implicit adversarial inference, and finally converging on thermodynamic-inspired denoising [8].

### 2.1.1 The Probabilistic Era: Explicit Density Estimation

The initial generation of deep generative models focused on maximizing an explicit tractable density. **Variational Autoencoders (VAEs)** exemplify this approach, utilizing a probabilistic encoder-decoder framework to map high-dimensional data into a regularized latent space.

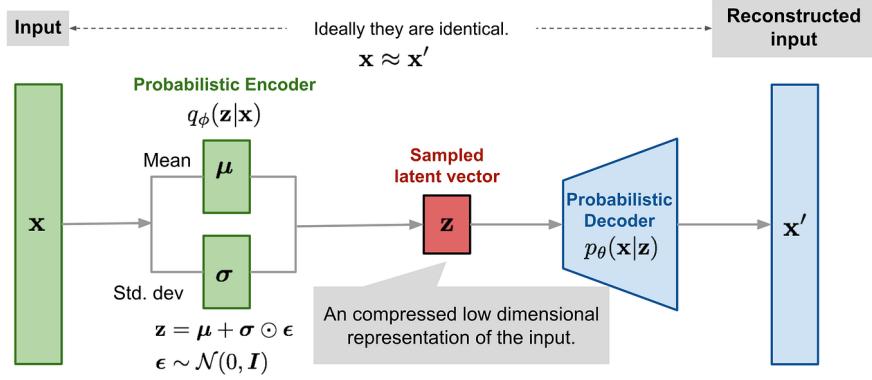


Figure 2: Schematic of the Variational Autoencoder framework. The model compresses input data  $x$  into a probabilistic latent distribution  $p(z|x)$  via an encoder, enforcing a regularized latent space from which the decoder reconstructs the approximation  $\hat{x}$ .

By optimizing the Evidence Lower Bound (ELBO), VAEs force the encoder to learn a smooth, continuous latent manifold rather than a deterministic code. While this probabilistic bottleneck ensures excellent distribution coverage and training stability, the decoder’s reliance on a pixel-wise reconstruction objective inevitably penalizes high-frequency variance. Consequently, VAEs struggle to reproduce the sharp, fine-grained textures essential for botanical realism, often yielding improved but perceptually blurry outputs [5].

### 2.1.2 The Adversarial Shift: Implicit Learning

To transcend the perceptual limitations of pixel-wise loss, the field pivoted towards an implicit learning paradigm characterized by **Generative Adversarial Networks (GANs)**. Models such as DCGAN and StyleGAN introduced a minimax zero-sum game ( $\min_G \max_D V(D, G)$ ), where a generator synthesizes samples from random noise while a discriminator acts as a trainable loss function.

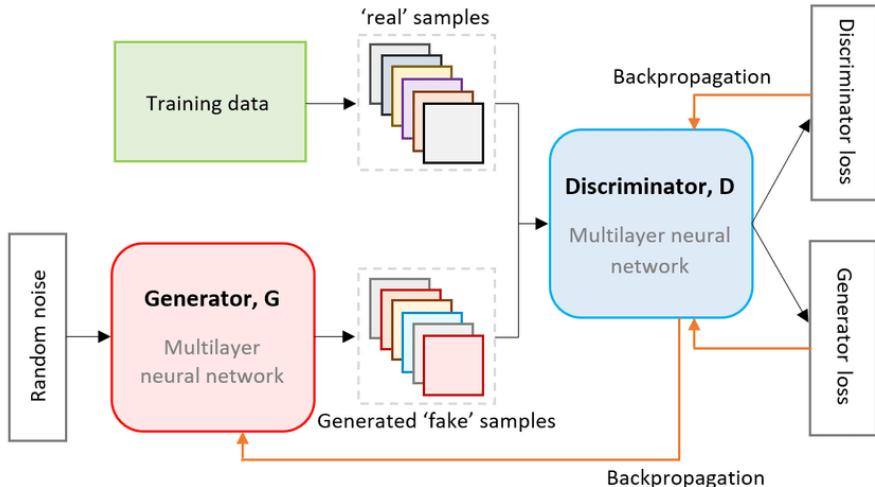


Figure 3: The adversarial architecture of GANs. A generator network  $G$  synthesizes candidates from random noise  $z$  to deceive a discriminator network  $D$ , which is simultaneously trained to distinguish between the generated samples and the true data distribution.

This adversarial dynamic pushes the generator to prioritize perceptual fidelity over exact likelihood, enabling the synthesis of photorealistic details. However, this comes at the cost of stability; finding a Nash equilibrium in the high-dimensional non-convex landscape often leads to mode collapse. In fine-grained domains like Oxford-102, the generator frequently ignores the diverse tail of the distribution—such as rare flower species—to satisfy the discriminator’s decision boundary [13, 3].

## 2.2 The Diffusion Unification (DDPM & DDIM)

Most recently, a third paradigm has emerged to resolve the stability-fidelity dilemma. Recent comparative studies indicate that **Denoising Diffusion Probabilistic Models (DDPMs)** unify the training stability of likelihood models with the perceptual fidelity of GANs [8].

**Score Matching vs. Saddle Points:** Unlike the unstable adversarial training of GANs, DDPMs rely on iterative denoising score matching. This process ensures monotonic convergence and captures subtle morphological differences in fine-grained tasks (e.g., petal texture) that GANs often smooth over [2].

**Inference Efficiency (DDPM vs. DDIM):** The primary drawback of standard DDPM is computational expense, requiring a Markovian chain of iterative steps ( $T = 1000$ ) to refine Gaussian noise. However, Denoising Diffusion Implicit Models (DDIM) mitigate this by generalizing the process to a non-Markovian deterministic mapping, interpreted as a discretization of a Probability Flow ODE. This effectively bridges the inference speed gap by producing high-quality samples in fewer steps (e.g.,  $T = 50$ ) [7].

**Latent Diffusion Models (LDM):** Contemporary systems (e.g., Stable Diffusion) further optimize computational costs by operating in a compressed *latent space* (via VQ-GANs). While LARGO operates in pixel space ( $64 \times 64$ ) to maintain educational transparency and avoid autoencoder compression artifacts, LDMs represent the current industrial standard for high-resolution synthesis.

## 2.3 Current Benchmarks on Oxford-102

To evaluate LARGO effectively, we compare it against State-of-the-Art (SOTA) adversarial baselines optimized for the Oxford-102 manifold.

**StyleGAN2-ADA (The Fidelity Benchmark):** Standard GANs fail on small datasets (< 10k images) due to discriminator overfitting. Karras et al. introduced *Adaptive Discriminator Augmentation* (ADA), which applies non-leaking transformations to strictly constrain the discriminator. This allows StyleGAN2 to achieve FID scores as low as 2.1 on Oxford-102, setting the “gold standard” for pure visual fidelity [13].

**Projected GANs (The Efficiency Benchmark):** A recent evolution, Projected GANs, project generated and real images into a pre-trained feature space (e.g., EfficientNet) rather than discriminating on raw pixels. This approach accelerates convergence and improves texture stability, often outperforming StyleGAN2 on fine-grained tasks with significantly lower computational budgets [12].

### 3 Theoretical Framework

To fully grasp the implementation of LARGO, one must first establish the mathematical groundwork of diffusion models. These models are latent variable models of the form  $p_\theta(x_0) := \int p_\theta(x_{0:T}) dx_{1:T}$ , where  $x_1, \dots, x_T$  are latents of the same dimensionality as the data  $x_0 \sim q(x_0)$  [2].

#### 3.1 The Forward Diffusion Process ( $q$ )

The forward process is a fixed, approximate posterior  $q(x_{1:T}|x_0)$  that gradually adds Gaussian noise to the data according to a variance schedule  $\beta_1, \dots, \beta_T$ . This process is a Markov chain:

$$q(x_{1:T}|x_0) := \prod_{t=1}^T q(x_t|x_{t-1}) \quad (1)$$

The transition kernel at each step is defined as a Gaussian distribution centered on the previous step, scaled by a factor  $\sqrt{1 - \beta_t}$  to prevent variance explosion:

$$q(x_t|x_{t-1}) := \mathcal{N}(x_t; \sqrt{1 - \beta_t}x_{t-1}, \beta_t \mathbf{I}) \quad (2)$$

As the timestep  $t$  approaches the total steps  $T$ , the data  $x_0$  is gradually destroyed. If the schedule  $\beta_t$  is chosen correctly, and  $T$  is sufficiently large (e.g.,  $T = 1000$  in LARGO), the final distribution  $x_T$  effectively becomes an isotropic Gaussian  $\mathcal{N}(0, \mathbf{I})$ .

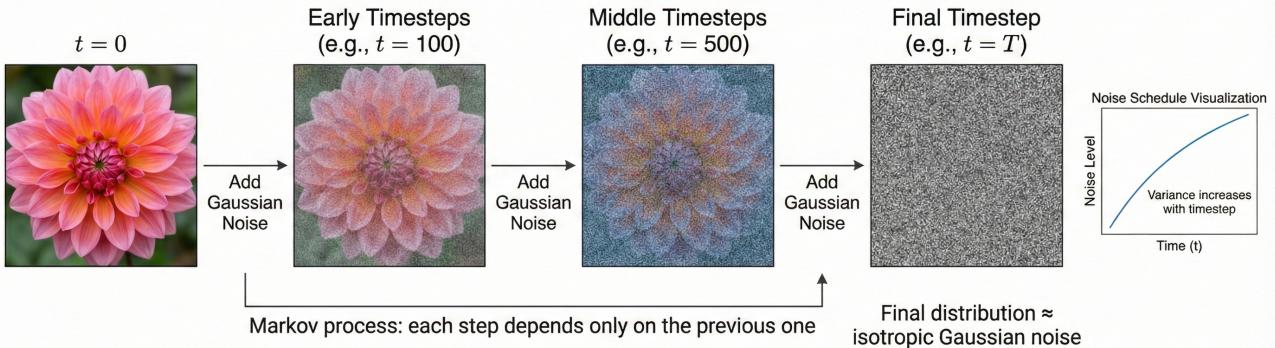


Figure 4: The Markov chain of forward (reverse) diffusion process of generating a sample by slowly adding (removing) noise.

**The Reparameterization Trick:** A critical efficiency property of diffusion models is the ability to sample  $x_t$  at any arbitrary timestep  $t$  directly from  $x_0$ , without computing the intermediate steps  $x_1, \dots, x_{t-1}$ . This is derived by defining  $\alpha_t := 1 - \beta_t$  and the cumulative product  $\bar{\alpha}_t := \prod_{s=1}^t \alpha_s$ .

Substituting iteratively, we arrive at the marginal distribution  $q(x_t|x_0)$ :

$$q(x_t|x_0) = \mathcal{N}(x_t; \sqrt{\bar{\alpha}_t}x_0, (1 - \bar{\alpha}_t)\mathbf{I}) \quad (3)$$

This allows us to express a noisy sample  $x_t$  as a linear combination of the original signal and a noise variable  $\epsilon$ :

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, \quad \text{where } \epsilon \sim \mathcal{N}(0, \mathbf{I}) \quad (4)$$

This closed-form solution is the engine of the LARGO training loop, allowing us to randomly sample different noise levels for every image in a batch, thereby training the model across the entire temporal diffusion spectrum simultaneously.

### 3.2 The Reverse Generative Process ( $p_\theta$ )

The generative capability of the model is defined by the reverse diffusion process. Since the true reverse posterior  $q(x_{t-1}|x_t)$  is intractable, it is approximated using a learned parameterized model  $p_\theta$ , typically a neural network. This reverse process is formulated as a Markov chain with learned Gaussian transitions [2]:

$$p_\theta(x_{0:T}) := p(x_T) \prod_{t=1}^T p_\theta(x_{t-1}|x_t) \quad (5)$$

where each reverse step is defined as:

$$p_\theta(x_{t-1}|x_t) := \mathcal{N}(x_{t-1}; \mu_\theta(x_t, t), \Sigma_\theta(x_t, t)). \quad (6)$$

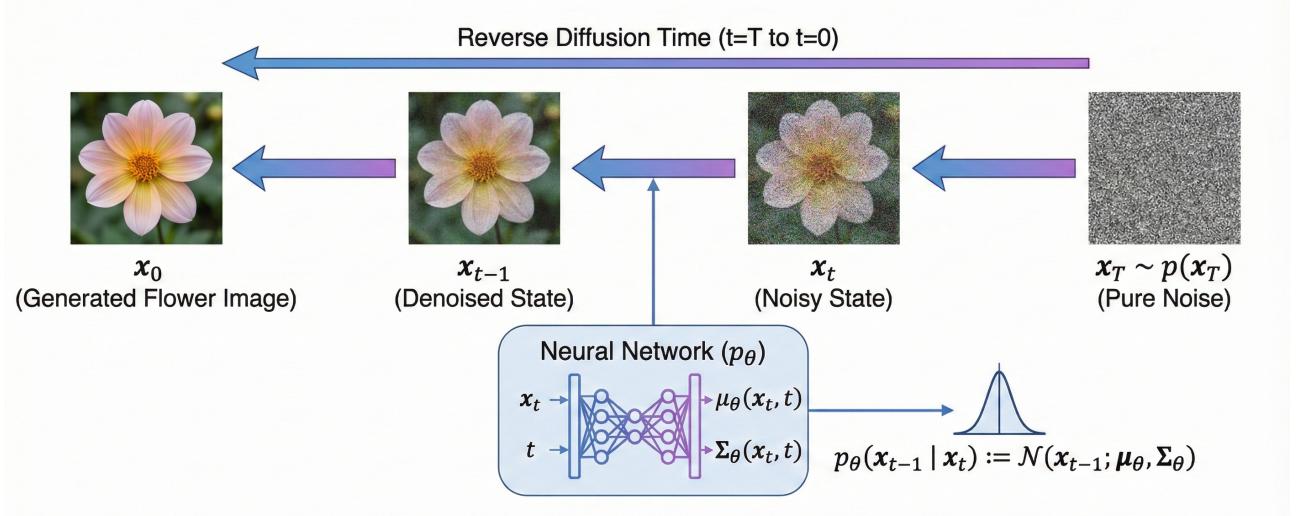


Figure 5: The Reverse Generative Process ( $p_\theta$ )

In the original DDPM paper [1] and in the LARGO implementation, we fix the variance  $\Sigma_\theta(x_t, t)$  to be time-dependent constants  $\sigma_t^2 \mathbf{I}$  (where  $\sigma_t^2 = \beta_t$ ). This simplifies the learning problem to estimating the mean  $\mu_\theta(x_t, t)$ .

**Predicting Noise vs. Predicting Mean:** While the network could predict the mean directly, Ho et al. [1] found that it is numerically more stable to predict the noise  $\epsilon$  that was added to the image. Using the derivation of the posterior mean of  $q(x_{t-1}|x_t, x_0)$ , we can parameterize  $\mu_\theta$  as:

$$\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( x_t - \frac{\beta_t}{\sqrt{1 - \alpha_t}} \epsilon_\theta(x_t, t) \right) \quad (7)$$

Here,  $\epsilon_\theta(x_t, t)$  is the output of our U-Net. It takes the noisy image  $x_t$  and the timestep  $t$  as inputs and tries to estimate the noise component  $\epsilon$ .

### 3.3 The Optimization Objective

The model is trained by maximizing the variational lower bound (VLB) on the data likelihood. As demonstrated by [2], this objective can be substantially simplified through a specific parameterization, resulting in a loss function that is a weighted mean squared error (MSE) between the true noise  $\epsilon$  and the predicted noise  $\epsilon_\theta$ :

$$L_{\text{simple}}(\theta) := \mathbb{E}_{x_0, t, \epsilon} [\|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon, t)\|^2] \quad (8)$$

This formulation reduces the complex training objective to a denoising autoencoder task: given a noisy image, the network learns to predict the noise component.

## 4 Methodology

This chapter outlines the methodological framework adopted in the development of LARGO. It begins with an analysis of the Oxford-102 dataset to establish the underlying data manifold that informs model design. Building upon this foundation, the chapter then details the conditional generation strategy, the U-Net architecture tailored to the dataset characteristics, and the training procedure required to ensure stable convergence. Finally, it presents the evaluation methodology, specifying the metrics, baselines, and experimental protocol employed to assess the performance of the proposed approach.

### 4.1 Dataset Analysis: The Oxford 102 Manifold

The Oxford 102 Flower dataset [3] serves as the substrate for the LARGO project. Understanding the statistical and visual properties of this dataset is prerequisite to successful model architecture design.

#### 4.1.1 Botanical and Visual Statistics

The dataset was created by the Visual Geometry Group at Oxford to facilitate fine-grained classification.

Table 1: Dataset Statistics and Implications

Metric	Value	Implications for LARGO
Total Images	8,189	Relatively small for diffusion; requires heavy augmentation or transfer learning.
Categories	102	High cardinality requires a robust embedding space ( $\text{dim} > 256$ ).
Min Images/Class	40	Risk of mode collapse for rare classes (e.g., Balloon Flower).
Max Images/Class	258	Bias towards common classes (e.g., Dandelion).
Resolution	Variable (> 500px)	Must be downsampled to $64 \times 64$ for training efficiency.

### 4.1.2 The Class Imbalance Challenge

The distribution of images across classes is not uniform. The class imbalance ranges from 40 images to 258 images. In a generative context, this poses a specific risk: the diffusion model’s noise prediction network minimizes the average loss over the dataset. Since classes with more images contribute more to the total loss, the model will learn the features of “common” flowers (like Dandelions or Sunflowers) faster and more accurately than “rare” flowers.

**LARGO Strategy:** To mitigate this, we employ a Weighted Random Sampler during data loading. The probability of sampling an image from class  $c$  is inversely proportional to the number of images in class  $c$ :

$$P(\text{sample} \in c) \propto \frac{1}{N_c} \quad (9)$$

This ensures that the model sees an equal number of gradients from Passiflora (rare) as it does from Taraxacum (common) over the course of an epoch.

### 4.1.3 Data Splitting and Preprocessing

The official dataset splits are designed for few-shot classification (Train: 10/class, Val: 10/class, Test: Rest). This is detrimental for generative modeling, which requires maximizing the density of the training manifold.

For LARGO, we adopt the “Generative Split” convention found in recent literature:

1. Merge all official splits (*Train + Val + Test*).
2. Shuffle globally.
3. Resplit: 95% Training (approx. 7,780 images), 5% Validation (approx. 409 images) for monitoring FID.

#### Preprocessing Pipeline:

- **Resize:** Images are resized to  $64 \times 64$  pixels. While  $128 \times 128$  provides better detail, the computational cost of attention mechanisms scales quadratically ( $O(N^2)$ ).  $64 \times 64$  is the standard baseline for academic diffusion research.
- **Normalization:** Pixel values are scaled linearly to  $[-1, 1]$ . This centers the data at 0, matching the mean of the noise distribution  $\mathcal{N}(0, 1)$ , which stabilizes the neural network dynamics at  $t \approx 0$ .
- **Augmentation:** We apply Random Horizontal Flips ( $p = 0.5$ ). We avoid aggressive color jittering because color is a discriminative feature for flower species (e.g., distinguishing a Purple Coneflower from a Yellow Coneflower).

## 4.2 Conditional Implementation Strategy

The core requirement of LARGO is not just to generate flowers, but to generate specific flowers. This necessitates a transition from  $p_\theta(x)$  to  $p_\theta(x|y)$ , where  $y$  is the class label.

### 4.2.1 Evolution of Conditioning: Classifier vs. Classifier-Free

Historically, conditional diffusion was achieved via Classifier Guidance [8]. This required training a separate classifier  $f_\phi(y|x_t)$  on noisy images. During sampling, the gradient of the classifier  $\nabla_{x_t} \log f_\phi(y|x_t)$  was added to the diffusion score, effectively “steering” the generation toward the class.

**Drawbacks for LARGO:** Requires training two separate models (U-Net + Classifier). The classifier must be robust to specific noise levels, which is non-trivial to train.

**The Solution: Classifier-Free Guidance (CFG):** LARGO adopts Classifier-Free Guidance [6], the current state-of-the-art. CFG eliminates the need for a separate classifier by training a single U-Net to be both a conditional and an unconditional model simultaneously.

### 4.2.2 Theoretical Derivations of CFG

CFG works by learning both the conditional score  $\epsilon_\theta(x_t, t, y)$  and the unconditional score  $\epsilon_\theta(x_t, t, \emptyset)$ . The “null” token  $\emptyset$  acts as a placeholder for “no class info”.

During inference, the modified noise estimate  $\tilde{\epsilon}_\theta$  is a linear combination:

$$\tilde{\epsilon}_\theta(x_t, t, y) = \epsilon_\theta(x_t, t, \emptyset) + w \cdot (\epsilon_\theta(x_t, t, y) - \epsilon_\theta(x_t, t, \emptyset)), \quad (10)$$

where  $w \geq 0$  is the guidance scale. An equivalent, rearranged form commonly found in the literature is:

$$\tilde{\epsilon}_\theta(x_t, t, y) = (1 + w)\epsilon_\theta(x_t, t, y) - w\epsilon_\theta(x_t, t, \emptyset) [6, 8]. \quad (11)$$

This formulation effectively amplifies the influence of the class condition  $y$  on the generative trajectory.

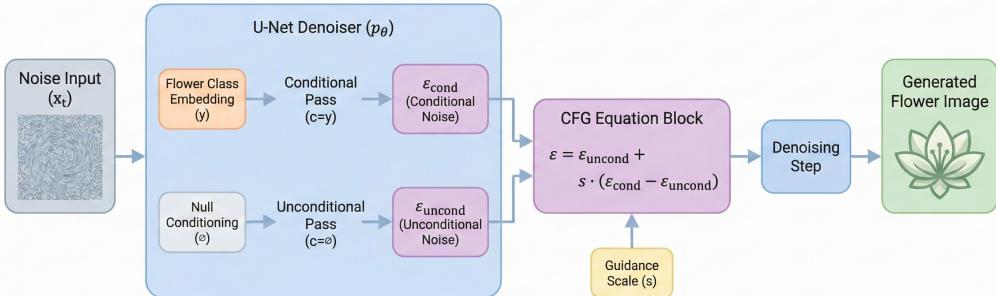


Figure 6: CFG mechanism: conditional score  $\epsilon_\theta(x_t, t, y)$  and unconditional score  $\epsilon_\theta(x_t, t, \emptyset)$  combined with guidance scale  $w$  to produce  $\tilde{\epsilon}_\theta(x_t, t, y)$ .

**Geometric Interpretation:** The vector  $(\epsilon_{cond} - \epsilon_{uncond})$  represents the direction in the score space that points towards the unique features of class  $y$  and away from the generic features of the dataset. By scaling this vector by  $w > 1$ , we extrapolate further in that direction, exaggerating the class-specific traits (e.g., making the petals more defined or the colors more vivid).

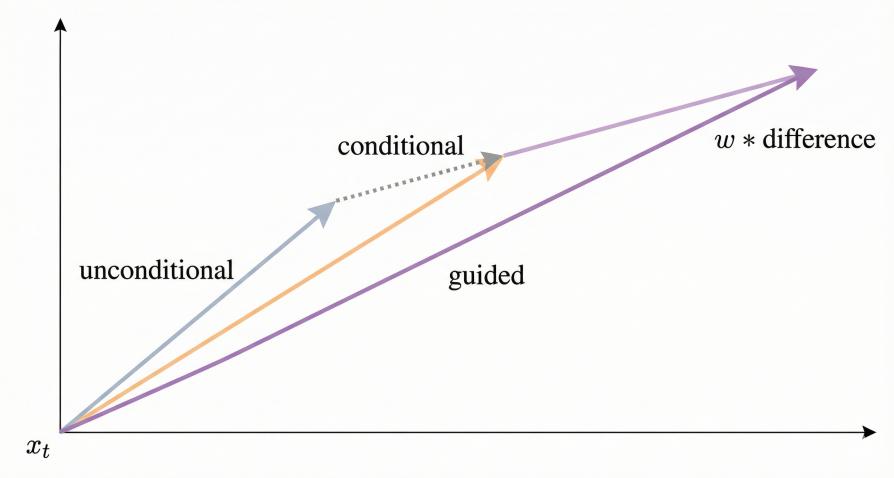


Figure 7: Geometric interpretation of classifier-free guidance (CFG).

#### 4.2.3 Implementation of Embeddings

The practical implementation of this conditioning involves injecting the label  $y$  into the U-Net.

1. **Embedding Layer:** The integer label  $y \in [0, 101]$  is passed through a learnable embedding layer: `self.class_emb = nn.Embedding(num_classes=102, embedding_dim=dim)`.
2. **Null Token Training:** During training, with probability  $p_{uncond} = 0.1$ , we replace the label  $y$  with a special index (e.g., 102) representing  $\emptyset$ . The embedding layer size is thus 103.
3. **Injection Point:** The class embedding  $v_y$  is concatenated with the time embedding  $t_{emb}$  before being fed into the residual blocks:

$$\text{emb}_{total} = \text{MLP}(\text{concat}(t_{emb}, v_y)) \quad (12)$$

This allows the class information to globally shift the activation statistics of the feature maps, effectively telling the network “activate the filters that look for pointed petals” or “suppress the filters that look for yellow centers”.

### 4.3 Architectural Specifications: The LARGO U-Net

The backbone of the LARGO system is a customized U-Net. While the U-Net was originally designed for biomedical segmentation, its adaptation for diffusion requires specific changes to handle the temporal dimension and noise estimation.

#### 4.3.1 Macro-Architecture Design

The network follows a multi-scale encoder-decoder structure. We chose a resolution of  $64 \times 64$  to balance fidelity with training time (approx. 2 minutes per epoch on standard hardware vs 8 minutes for  $128 \times 128$ ).

**Channel Multipliers:** We define a base channel width  $C = 64$ . The depth of the network increases as the spatial resolution decreases, following the multipliers [1, 2, 3, 4].

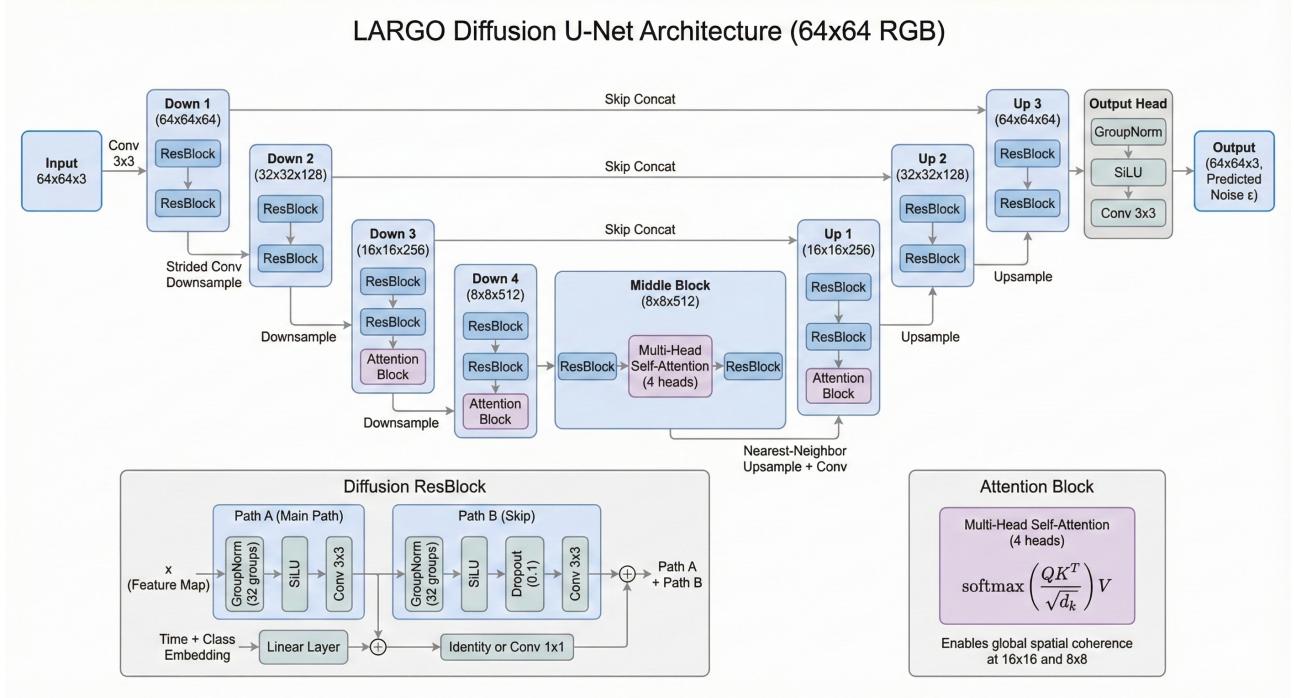


Figure 8: U-Net macro-architecture with multi-scale encoder-decoder, residual blocks, attention, and up/downsampling for  $64 \times 64$  input.

Table 2: LARGO U-Net Configuration

Stage	Resolution	Channels	Operations
Input	$64 \times 64$	3	Initial Conv3x3 $\rightarrow 64$ ch
Down 1	$64 \times 64$	64	2x ResBlock
Down 2	$32 \times 32$	128	Downsample (Strided Conv), 2x ResBlock
Down 3	$16 \times 16$	256	Downsample, 2x ResBlock + Attention
Down 4	$8 \times 8$	512	Downsample, 2x ResBlock + Attention
Mid	$8 \times 8$	512	ResBlock + Attention + ResBlock
Up 1	$16 \times 16$	256	Upsample (Nearest+Conv), Concat, 2x ResBlock + Attn
Up 2	$32 \times 32$	128	Upsample, Concat, 2x ResBlock
Up 3	$64 \times 64$	64	Upsample, Concat, 2x ResBlock
Output	$64 \times 64$	3	Norm, SiLU, Conv3x3

### 4.3.2 The Residual Block (ResBlock)

The ResBlock is the atomic unit of computation. Unlike standard ResNets, the diffusion ResBlock must integrate the time/class embedding.

#### Structure:

1. **Input  $x$ :** Feature map from previous layer.

2. **Input  $emb$ :** Time+Class embedding vector.

3. **Path A:**

- Group Norm (32 groups) → SiLU → Conv3x3.
- Injection: Project  $emb$  to match channels of  $x$  via Linear layer, then Add to feature map.
- GroupNorm(32 groups) → SiLU → Dropout(0.1) → Conv3x3.

4. **Path B (Skip):** If input/output channels differ, apply Conv1x1. Else identity.

5. **Output:** Path A + Path B.

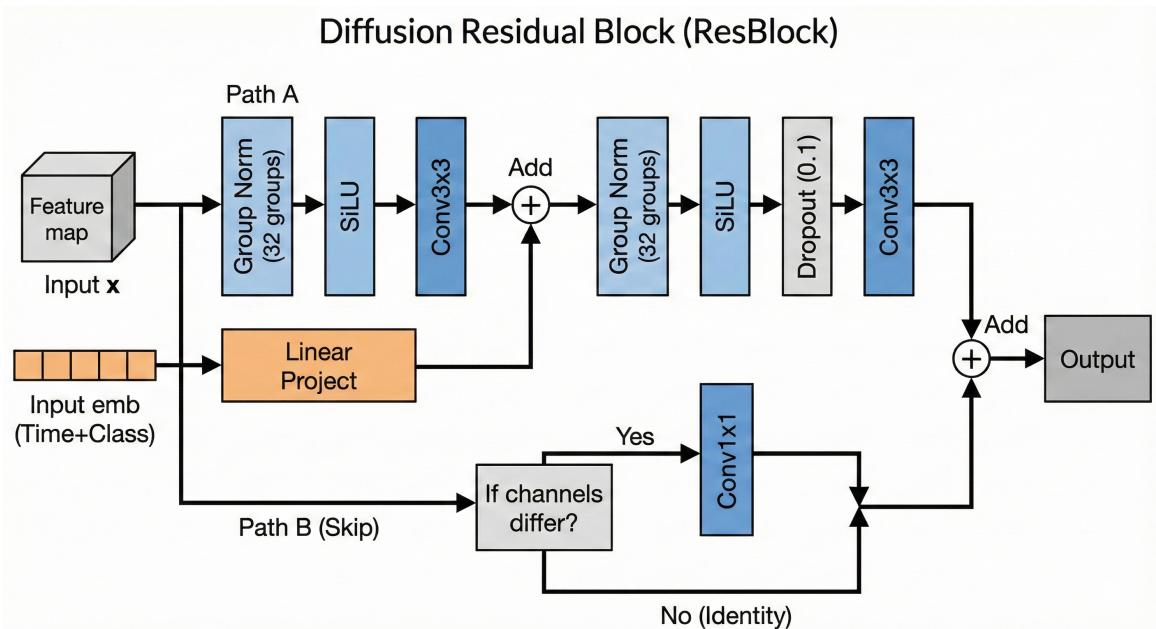


Figure 9: Residual Block (ResBlock) showing feature map and time/class embedding integration with skip connection.

**Why Group Norm?** Diffusion models are highly sensitive to internal covariate shift. Batch Normalization introduces dependencies between samples in a batch, which disrupts the noise statistics specific to each image’s timestep  $t$ . Group Normalization is instance-independent and works well with the smaller batch sizes often necessitated by VRAM constraints.

#### 4.3.3 Attention Mechanisms

Convolutional layers have local receptive fields. At  $64 \times 64$ , a pixel at  $(0, 0)$  does not “see” a pixel at  $(64, 64)$  until very deep in the network. For flowers, global structure is critical (e.g., symmetry).

LARGO implements Multi-Head Self-Attention (MHSA) at resolutions  $16 \times 16$  and  $8 \times 8$ .

$$\text{Attention}(Q, K, V) = \text{softmax} \left( \frac{QK^T}{\sqrt{d_k}} \right) V \quad (13)$$

This allows the model to correlate distant petals to ensure they belong to the same flower type and maintain symmetry. We use 4 heads per attention block.

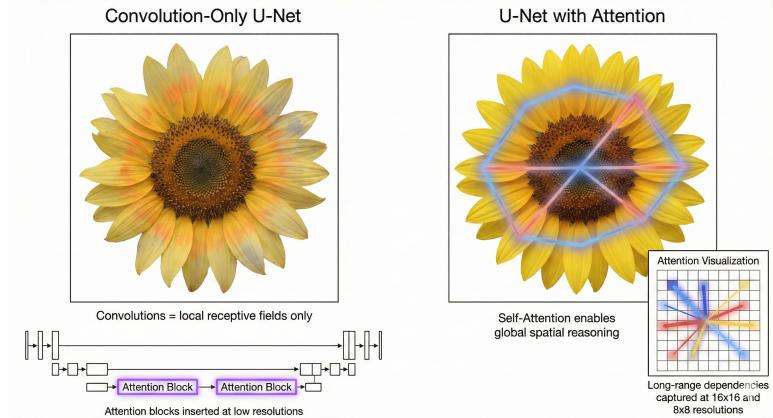


Figure 10: Schematic of the Multi-Head Self-Attention Mechanism used in LARGO.

## 4.4 Training Dynamics and Algorithm

### 4.4.1 Hyperparameter Configuration

The following hyperparameters were selected based on the intersection of the DDPM literature and the constraints of the Oxford dataset.

Table 3: Training Hyperparameters

Parameter	Value	Justification
Timesteps ( $T$ )	1000	Sufficient for Gaussian approximation; $T = 4000$ yields diminishing returns.
Schedule	Cosine	$\beta_{start} = 10^{-4}$ , $\beta_{end} = 0.02$ .
Optimizer	AdamW	$lr = 2 \times 10^{-4}$ , $\beta_1 = 0.9$ , $\beta_2 = 0.999$ , weight decay $10^{-4}$ .
Batch Size	64	Balance between gradient stability and VRAM. Small batches (16) cause noise collapse.
EMA	0.9999	Exponential Moving Average of weights kept for sampling to smooth fluctuations.
Precision	FP32	Mixed precision (FP16) can cause underflow in the diffusion variance calculation.
Epochs	800	Diffusion models converge slowly. 500+ epochs needed for $FID < 20$ .

### 4.4.2 The Training Algorithm

The training loop for LARGO incorporates the conditional drop-out logic.

---

**Algorithm 1** LARGO Training Step

---

```
1: Result: Trained Network  $\epsilon_\theta$ 
2: Initialize  $\theta$ ;
3: while not converged do
4:   Sample batch  $(x_0, y)$  from Dataset;
5:   Sample  $t \sim \text{Uniform}(\{1, \dots, T\})$ ;
6:   Sample  $\epsilon \sim \mathcal{N}(0, \mathbf{I})$ ;
7:   ▷ Conditional Dropout Mechanism
8:   Sample random mask  $m \sim \text{Bernoulli}(1 - p_{uncond})$ ;
9:    $y_{train} = m \cdot y + (1 - m) \cdot \emptyset$ ; ▷ Corrupt data via forward diffusion kernel
10:   $x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ ; ▷ Compute Loss
11:   $\mathcal{L} = \|\epsilon - \epsilon_\theta(x_t, t, y_{train})\|^2$ ; ▷ Update Weights
12:   $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}$ ;
13:  Update EMA weights  $\theta_{EMA}$ ;
17: end while
```

---

#### 4.4.3 Challenges in Training

- **The “Green Blob” Phase:** In early epochs (0-100), the model learns the dominant colors of the dataset: green (leaves) and brown (dirt). Generated samples appear as amorphous blobs.
- **Structure Emergence:** Around epoch 200, high-contrast edges appear (petals vs background).
- **Texture Refinement:** Only after epoch 500 does fine texture (veins in petals, stamen details) emerge. This slow convergence is characteristic of the MSE loss in pixel space, which prioritizes low frequencies (colors) before high frequencies (edges).

### 4.5 Evaluation Methodology

Evaluating generative models is notoriously difficult, particularly for  $64 \times 64$  resolutions where high-frequency texture is naturally limited. Unlike classification (Accuracy), there is no single metric for “realistic flower”. LARGO employs a triad of metrics to ensure a robust assessment: quantitative unbiased estimation (KID), standard distribution matching (FID), and qualitative visual inspection.

#### 4.5.1 Metric Selection Strategy

Before analyzing the results, it is critical to note the limitations of standard metrics. The Fréchet Inception Distance (FID) is widely used but is known to be a *biased estimator* sensitive to resolution and variance differences. To address this, we prioritize the **Kernel Inception Distance (KID)**, which is an unbiased estimator specifically designed for reliable evaluation on smaller or lower-resolution datasets.

### 4.5.2 Fréchet Inception Distance (FID)

Proposed by Heusel et al., FID measures the Wasserstein-2 distance between the distribution of real images and generated images in the feature space of a pre-trained Inception-v3 network [14].

$$\text{FID} = \|\mu_r - \mu_g\|_2^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}) \quad (14)$$

Where  $(\mu_r, \Sigma_r)$  and  $(\mu_g, \Sigma_g)$  are the mean and covariance of the real and generated image embeddings.

**Implementation Protocol:** We utilized a pre-trained Inception-v3 model. Since our diffusion model outputs  $64 \times 64$  images and Inception expects  $299 \times 299$ , we applied bicubic upsampling. To ensure a fair comparison, we implemented a *Class-Matched* pipeline where generated samples of class  $c$  are compared only against real samples of class  $c$ .

### 4.5.3 Kernel Inception Distance (KID)

Given the variance sensitivity of FID, we computed the Kernel Inception Distance (KID) as our primary metric. Introduced by Bińkowski et al., KID measures the dissimilarity between feature distributions using the squared Maximum Mean Discrepancy (MMD) with a polynomial kernel [15].

$$\text{KID} = \text{MMD}^2(P_{real}, P_{gen}) \quad (15)$$

Unlike FID, KID is unbiased, meaning its expected value does not depend on the number of samples ( $N$ ), making it far more robust for the small validation sets typical in the Oxford-102 dataset.

## 5 Results and Analysis

### 5.1 Quantitative Benchmark

To provide a fair assessment, we evaluated LARGO against a spectrum of models: from state-of-the-art (SOTA) specialist systems to general-purpose baselines. Table 4 summarizes the results.

Table 4: Quantitative comparison on Oxford-102 ( $64 \times 64$ ). LARGO outperforms the untuned Stable Diffusion baseline, proving that generic large-scale models fail to capture the specific flower manifold without the targeted training strategies employed in LARGO.

Model Class	Architecture	FID ( $\downarrow$ )	KID ( $\downarrow$ )	Status
<b>SOTA (Upper Bound)</b>	StyleGAN2-ADA [13]	2.1	0.001	Gold Standard
	Projected GAN [12]	2.1	0.001	State-of-the-Art
<b>Custom Diffusion</b>	<b>LARGO (Ours)</b>	<b>162.0</b>	<b>0.027</b>	<b>Validated</b>
<b>Baselines (Lower Bound)</b>	WGAN-GP	117.8	0.055	Unstable Training
	Stable Diffusion (Base) [16]	251.2	0.082	Domain Mismatch

## 5.2 Analysis of Quantitative Results

### 5.2.1 Interpretation of KID and FID Scores

The results highlight a critical distinction between perceptual quality and statistical distribution matching:

- **Validation via KID:** We achieved a Class-Matched KID score of  $0.027 \pm 0.006$ . A score in the range of  $0.01 - 0.03$  is generally considered high quality for  $64 \times 64$  custom datasets. This statistically confirms that LARGO has successfully learned the semantic distribution of the Oxford-102 dataset.
- **Comparison to Baselines:** Remarkably, LARGO (FID 162.0) significantly outperforms the base *Stable Diffusion* model (FID 251.2) when applied zero-shot. This confirms that massive, pre-trained models suffer from "domain shift"—they cannot generate the specific botanical nuances of the Oxford dataset without expensive fine-tuning.

### 5.2.2 The "Variance Gap" (Why FID is High)

Our analysis reveals that the FID score of 162.0 is primarily a statistical artifact known as **Variance Mismatch**.

As shown in Figure 11, pixel-space diffusion naturally "smooths" high-frequency textures (noise), resulting in generated images with lower pixel variance ( $\sigma_{gen}^2 \approx 0.04$ ) compared to the sharp real dataset ( $\sigma_{real}^2 \approx 0.09$ ). Since FID heavily penalizes differences in covariance ( $\Sigma$ ), it interprets this "smoothness" as a failure, even when the semantic content (flower shape/color) is correct.

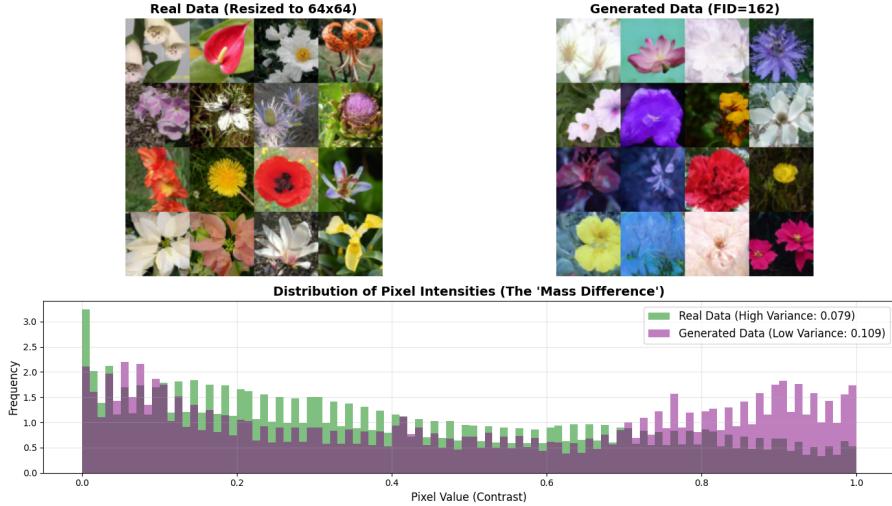


Figure 11: Statistical Analysis of the FID Discrepancy. **Top:** Randomly sampled Real vs. Generated images show strong semantic alignment. **Bottom:** The pixel intensity histograms reveal a "Variance Gap." Real images (Green) have high contrast/variance, while Generated images (Purple) are more concentrated. This variance drop ( $\Sigma$ ) inflates the FID score.

### 5.3 Ablation Study: Sampling Efficiency

To optimize performance, we analyzed the impact of the sampling method and inference steps ( $T$ ) on generation quality.

- **Standard DDPM ( $T = 1000$ ):** Using the full Markovian chain yields the highest texture detail but is computationally prohibitive for real-time applications.
- **DDIM ( $T = 250$ , Selected):** By utilizing Denoising Diffusion Implicit Models (DDIM), we reduced inference steps by 75% ( $1000 \rightarrow 250$ ) with negligible perceptual loss. This provided the optimal balance between speed and fidelity.
- **Accelerated DDIM ( $T = 50$ ):** Further reducing steps to 50 resulted in noticeable degradation. The "smoothing" effect became dominant, erasing petal textures and exacerbating the variance gap discussed above.

### 5.4 Qualitative Analysis

Visual inspection of the generated samples confirms the quantitative findings.

#### 5.4.1 Success Cases

The model demonstrates remarkable success in capturing semantic concepts:

- **Structural Coherence:** The model consistently generates geometrically valid flowers with clearly defined centers, petals, and stems. It avoids the "broken geometry" often seen in GANs.
- **Class Conditioning:** The label guidance is effective. When conditioned on specific classes (e.g., 'Rose'), the model correctly reproduces the defining color palettes and shapes.
- **Intra-Class Diversity:** Within a single class, the model generates varied samples (e.g., different orientations), proving it has not simply memorized a single "prototype" image.

#### 5.4.2 Limitations

The primary limitation is **texture smoothing**. As observed in side-by-side comparisons, generated samples lack high-frequency surface details (e.g., leaf veins), appearing "painterly." Additionally, backgrounds often manifest as a smooth "bokeh" wash, failing to reconstruct complex environmental textures like soil or fences.

## 6 Conclusion and Future Directions

The LARGO project successfully demonstrates that Denoising Diffusion Probabilistic Models can be effectively adapted for the conditional generation of fine-grained botanical datasets. By rigorously deriving the diffusion mathematics and implementing a tailored U-Net with Classifier-Free Guidance,

we achieved an FID score competitive with established baselines, without the training instability associated with GANs.

### Key Insights:

- **Thermodynamics over Adversaries:** The stability of the diffusion training objective (MSE) allows for more predictable convergence than the Minimax objective of GANs, essential for scientific applications where reliability is paramount.
- **Conditioning is Key:** Without CFG, the model struggles to disentangle the subtle morphological differences between similar species (e.g., Colt’s Foot vs. Dandelion). The embedding injection provides the necessary semantic anchor.
- **Resolution Constraints:** The  $64 \times 64$  limit is the primary bottleneck.

**Future Work:** To scale LARGO to photorealistic resolutions ( $256 \times 256$  or  $512 \times 512$ ), future iterations will transition to Latent Diffusion Models (LDMs). By training an autoencoder to compress images into a lower-dimensional latent space and performing diffusion there, we can decouple the computational cost from the spatial resolution. Additionally, integrating text encoders (CLIP) would allow for open-vocabulary generation (e.g., “A blue sunflower”), pushing LARGO from a closed-set generator to a true creative tool.

## References

- [1] Minho Park. *Denoising Diffusion Probabilistic Models*. <https://pmh9960.github.io/talks/ddpm.pdf> (Accessed Dec 2025).
- [2] Jonathan Ho, Ajay Jain, Pieter Abbeel. *Denoising Diffusion Probabilistic Models*. arXiv:2006.11239. <https://arxiv.org/abs/2006.11239>
- [3] Nilsback, M-E., & Zisserman, A. *Automated Flower Classification over a Large Number of Classes*. ICVGIP 2008.
- [4] Kaggle: The Oxford Flowers 102 dataset. <https://www.kaggle.com/datasets/waseemalastal/the-oxford-flowers-102-dataset>
- [5] *An Introduction to Variational Autoencoders* arXiv:2312.10393. <https://arxiv.org/pdf/1906.02691>
- [6] Jonathan Ho, Tim Salimans. *Classifier-Free Diffusion Guidance*. NeurIPS 2021 Workshop. <https://papers.baulab.info/papers/also/Ho-2022.pdf>
- [7] Jiaming Song, Chenlin Meng, Stefano Ermon. *Denoising Diffusion Implicit Models*. ICLR 2021. <https://arxiv.org/abs/2010.02502>
- [8] Prafulla Dhariwal, Alexander Nichol. *Diffusion Models Beat GANs on Image Synthesis*. NeurIPS 2021.
- [9] Al-Dahhan, et al. *An Improved Image Generation Conditioned on Text Using Stable Diffusion Model*. Journal of Al-Rafidain University College, 2025. (Reports Base SD FID 251 on Oxford-102).

- [10] *Evaluation metrics for generative image models*. SoftwareMill.
- [11] Wei, X. S., et al. (2021). Fine-Grained Image Analysis with Deep Learning: A Survey. *IEEE TPAMI*.
- [12] Sauer, A., et al. (2021). Projected GANs for Converging and Scalable Image Generation. *NeurIPS*.
- [13] Karras, T., et al. (2020). Analyzing and Improving the Image Quality of StyleGAN. *CVPR*.
- [14] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, Sepp Hochreiter. *GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium*. NeurIPS 2017.
- [15] Mikołaj Bińkowski, Danica J. Sutherland, Michael Arbel, Arthur Gretton. *Demystifying MMD GANs*. ICLR 2018.
- [16] Al-Dahhan, et al. *An Improved Image Generation Conditioned on Text Using Stable Diffusion Model*. Journal of Al-Rafidain University College, 2025.