# Foundations of Generative Models

Lecture Notes for Summer Course on Diffusion Models at
PES University, India

Shubham Chatterjee

Department of Computer Science

Missouri University of Science and Technology, Rolla, MO, USA

shubham.chatterjee@mst.edu

# Contents

# List of Figures

# List of Tables

# Preface

In recent years, diffusion models have rapidly emerged as one of the most powerful tools in the generative modeling landscape, rivaling and often surpassing traditional methods such as GANs and VAEs in terms of sample quality, flexibility, and theoretical grounding. From producing photorealistic images and editing fine-grained visual details to powering breakthroughs in text-to-image generation with systems like Stable Diffusion and DALL-E, diffusion models now lie at the heart of modern AI-driven creativity and synthesis.

This book aims to demystify diffusion models from the ground up, blending rigorous mathematical foundations with practical implementation insights. It is designed for students, researchers, and engineers who wish to build a principled understanding of how diffusion models work, why they are effective, and how to use and extend them in real-world applications.

The structure of this book reflects a progressive learning journey. We begin with the basics of generative modeling and probabilistic inference, then introduce denoising diffusion probabilistic models (DDPMs) as a stochastic generative framework. From there, we develop the mathematical tools necessary to understand continuous-time formulations, score-based generative modeling, and advanced sampling techniques such as DDIM and ODE solvers.

In later chapters, we explore conditional generation, classifier-free guidance, latent diffusion models, and the practical considerations required to scale diffusion models to high-resolution and multi-modal tasks. Each chapter is accompanied by implementation notes, exercises, and pointers to further reading, aiming to serve both as a conceptual reference and as a hands-on guide.

This book draws on current research, best practices from open-source frameworks, and my own experience teaching these models to students and practitioners. It is written with the belief that mastering generative modeling is not just about understanding equations, but about building intuition, developing critical thinking, and connecting theory to practice.

I hope this book not only teaches you how diffusion models work, but inspires you to explore what else they can become — in research, in industry, and in creative expression.

<div align="right">

Shubham Chatterjee
Assistant Professor
Department of Computer Science
Missouri University of Science and Technology, Rolla, MO, USA

</div>

# Chapter 1

# Preliminaries

## 1.1 The Challenge of Generating Reality

Imagine standing before a vast collection of photographs—portraits, landscapes, street scenes, each capturing a moment of reality. A natural question arises: *How can we generate new, realistic photographs that have never been taken before?* This seemingly simple question touches the heart of one of the most fascinating problems in machine learning and computer vision.

The challenge of generating realistic data extends far beyond photography. Whether we're synthesizing speech, creating artwork, or modeling molecular structures, we face the fundamental question of how to capture and reproduce the patterns that make data "real" or "plausible." This chapter explores the mathematical foundations and computational approaches that make such generation possible.

To understand this challenge, let us begin with a simpler scenario. Consider a collection of points scattered in two-dimensional space, perhaps representing the locations of coffee shops in a city or the positions of stars in a constellation. Given these observed points, how might we sample a new point that fits naturally within this collection? This simplified version of our problem reveals the core statistical principles that underlie all generative modeling.

## 1.2 The Statistical Perspective

### 1.2.1 From Observations to Distributions

The key insight in generative modeling is to adopt a *statistical perspective*. Rather than viewing our collection of data points as isolated observations, we imagine that each point is a sample drawn from some underlying probability distribution. This distribution, which we call the *data distribution*, captures the likelihood that any particular point (or image, or sound) represents genuine, realistic data.

Formally, if we have observed data points $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_n$, we assume these are independent samples from some unknown probability distribution $p(\mathbf{x})$. The function $p(\mathbf{x})$ represents the probability density at point $\mathbf{x}$—regions of high probability correspond to "typical" or "realistic" data, while regions of low probability represent unlikely or unrealistic configurations.

This statistical viewpoint immediately suggests a solution to our generation problem: *if we knew the probability density function $p(\mathbf{x})$, we could sample from*

*it directly to generate new, realistic data points.* The challenge, of course, is that we don't know this distribution—we only have samples from it.

### 1.2.2   The Sampling Problem: From Probability to Data

Before tackling the challenging problem of unknown distributions, let us first understand how to sample from known distributions. This foundational knowledge will prove essential as we build toward more complex generative models.

## 1.3   Sampling from Known Distributions

### 1.3.1   Discrete Distributions and the Cumulative Approach

Consider a discrete probability distribution over $n$ possible values $x_1, x_2, \ldots, x_n$, where each value $x_i$ occurs with probability $p_i$. For this to be a valid probability distribution, we require:

$$p_i \geq 0 \quad \text{for all } i \tag{1.1}$$

$$\sum_{i=1}^{n} p_i = 1 \tag{1.2}$$

The question is: how do we generate samples according to these probabilities? The answer lies in the *cumulative distribution function* (CDF), which transforms probabilities into a form suitable for sampling.

The CDF at point $i$ is defined as:

$$P_i = \sum_{j=1}^{i} p_j \tag{1.3}$$

This represents the probability that a sample falls in the range $\{x_1, x_2, \ldots, x_i\}$. Notice that $P_i$ is monotonically increasing, with $P_0 = 0$ and $P_n = 1$.

**Inverse Transform Sampling**

The elegant solution to sampling from this distribution is *inverse transform sampling*:

1. Draw a uniform random number $u \sim \mathcal{U}(0,1)$

2. Find the smallest index $i$ such that $P_{i-1} \leq u \leq P_i$

3. Return $x_i$

Why does this work? Intuitively, we're "throwing a dart" at the interval $[0,1]$ and seeing which probability "bin" it lands in. Since the bins have sizes proportional to their probabilities, we naturally sample more frequently from high-probability regions.

### 1.3.2 Continuous Distributions

For continuous distributions with probability density function $p(x)$, the concepts extend naturally. The CDF becomes:

$$F_X(x) = \Pr(X \le x) = \int_{-\infty}^{x} p(t)\, dt \tag{1.4}$$

with the property that $F_X(\infty) = 1$.

The inverse transform sampling method for continuous distributions follows the same logic:

1. Compute the CDF $F_X(x)$

2. Draw $u \sim \mathcal{U}(0, 1)$

3. Return $x = F_X^{-1}(u)$

The key requirement is that we must be able to compute the inverse CDF $F_X^{-1}(u)$.

**Example: Sampling from a Quadratic Distribution**

Let's work through a concrete example. Suppose we want to sample from the distribution:

$$p(x) = \frac{3}{8}x^2, \quad x \in [0, 2] \tag{1.5}$$

First, we compute the CDF:

$$F_X(x) = \int_0^x \frac{3}{8}t^2\, dt = \frac{3}{8} \cdot \frac{x^3}{3} = \frac{x^3}{8} \tag{1.6}$$

To find the inverse, we solve $u = \frac{x^3}{8}$ for $x$:

$$x = F_X^{-1}(u) = (8u)^{1/3} = 2u^{1/3} \tag{1.7}$$

Therefore, to sample from this distribution, we:

1. Draw $u \sim \mathcal{U}(0, 1)$

2. Return $x = 2u^{1/3}$

### 1.3.3 When Inverse Transform Fails: Rejection Sampling

Unfortunately, many distributions of interest don't have closed-form inverse CDFs. The normal distribution, gamma distribution, beta distribution, and most complex multimodal distributions fall into this category—their inverse CDFs either don't exist in closed form or are extremely difficult to compute. In such cases, we can use *rejection sampling*, which provides a more general approach at the cost of computational efficiency.

**The Geometric Intuition**

Rejection sampling is built on a beautiful geometric insight: **sampling from a distribution is equivalent to sampling uniformly from the area under its probability density curve**. Instead of trying to invert a complex CDF, we can think of the problem as throwing darts randomly at the area under the probability density function and keeping the $x$-coordinates where they land.

**The Algorithm**

The rejection sampling method works as follows:

1. **Find a proposal distribution** $q(x)$ such that $q(x) \geq p(x)$ for all $x$, where $q(x)$ is easy to sample from (such as uniform, exponential, or normal distributions)

2. **Sample a candidate point**: Draw $x \sim q(x)$

3. **Sample a height**: Draw $h \sim \mathcal{U}(0, q(x))$

4. **Accept or reject**: If $h \leq p(x)$, accept the sample $x$; otherwise, reject it and return to step 2

**Why This Works**

Geometrically, we're sampling uniformly from the area under $q(x)$ and keeping only those samples that fall under $p(x)$. The proposal distribution $q(x)$ acts like a "tent" that completely covers the shape of our target distribution $p(x)$ from above.

The acceptance condition $h \leq p(x)$ ensures that we only keep samples that fall within the area under our target distribution. The rejected samples correspond to the "excess area" between $q(x)$ and $p(x)$. Mathematically, the probability of accepting a sample at point $x$ is:

$$P(\text{accept at } x) = \frac{p(x)}{q(x)} \tag{1.8}$$

This ratio automatically ensures that our final samples follow the correct distribution $p(x)$.

**Efficiency and the Choice of Proposal Distribution**

The efficiency of rejection sampling depends critically on how tightly $q(x)$ bounds $p(x)$. The **acceptance rate** (fraction of samples we keep) equals:

$$\text{Acceptance Rate} = \frac{\int p(x)\, dx}{\int q(x)\, dx} = \frac{1}{\int q(x)\, dx} \tag{1.9}$$

since $p(x)$ integrates to 1.

**Key principle**: The tighter the bound, the higher the acceptance rate and the more efficient the algorithm becomes.

- If $q(x)$ closely approximates $p(x)$, most samples will be accepted → high efficiency

- If $q(x)$ is much larger than $p(x)$, many samples will be rejected → low efficiency, wasted computation

The art of rejection sampling lies in choosing a proposal distribution $q(x)$ that balances two competing requirements:

1. **Ease of sampling**: $q(x)$ should be simple to generate samples from

2. **Tight bounding**: $q(x)$ should be as close as possible to $p(x)$ to minimize rejections

**Advantages and Trade-offs**

**Advantages:**

- **Universal applicability**: Works for any distribution, regardless of complexity

- **Exact sampling**: Produces samples that follow the target distribution exactly (no approximation error)

- **Conceptual simplicity**: The geometric interpretation makes the method intuitive to understand and implement

**Disadvantages:**

- **Variable computational cost**: The number of iterations required is random and potentially unbounded

- **Efficiency dependence**: Performance heavily depends on the quality of the proposal distribution

- **Function evaluation overhead**: Requires evaluating $p(x)$ at each candidate point, which can be expensive for complex distributions

**Practical Considerations**

When implementing rejection sampling, consider these guidelines:

1. **Choose proposal distributions wisely**: Common choices include uniform distributions (for bounded support), exponential distributions (for one-sided unbounded support), and normal distributions (for two-sided unbounded support).

2. **Scale and shift appropriately**: You can multiply $q(x)$ by a constant $M \geq 1$ to ensure $Mq(x) \geq p(x)$. This doesn't change the acceptance probability but may make finding a suitable $q(x)$ easier.

3. **Monitor acceptance rates**: In practice, acceptance rates below 10–20% often indicate that a different approach might be more efficient.

Rejection sampling represents a fundamental trade-off in computational statistics: we sacrifice computational efficiency for complete generality, making it possible to sample from virtually any distribution at the cost of potentially many rejected attempts.

### 1.3.4 The Reparameterization Trick

A particularly important technique for neural networks is the *reparameterization trick*, which allows us to sample from complex distributions by transforming samples from simple ones. This elegant technique has revolutionized deep generative modeling and variational inference by enabling gradient-based optimization through stochastic sampling operations.

**The Problem**

In models like **Variational Autoencoders (VAEs)**, we often need to *sample* from a distribution during training. For example, we might sample a latent variable $z$ from a Gaussian distribution:

$$z \sim \mathcal{N}(\mu, \sigma^2)$$

However, sampling is a **non-differentiable operation**, which prevents us from computing gradients through it. This is a major issue because we rely on gradients to update our model parameters using backpropagation.

**The Insight**

The **reparameterization trick** offers a workaround:
    Instead of sampling $z$ directly from $\mathcal{N}(\mu, \sigma^2)$, we:

1. Sample an auxiliary noise variable $\varepsilon \sim \mathcal{N}(0, 1)$

2. Then define $z$ as:
$$z = \mu + \sigma \cdot \varepsilon$$

This transforms the sampling process into a **deterministic function** of $\mu$, $\sigma$, and $\varepsilon$.

**Why This Helps**

- Now, $z$ is a differentiable function of the model parameters $(\mu, \sigma)$.

- We can apply standard gradient-based optimization, since backpropagation can flow through this deterministic transformation.

**Intuition**

Before: *"Sample z randomly — can't backpropagate."*
After: *"Sample $\varepsilon$ once (fixed randomness), then compute z as a function of parameters — can backpropagate."*

**The Multivariate Normal Case**

The power of the reparameterization trick becomes even more apparent when we extend it to multivariate distributions. Imagine you're building a variational autoencoder that needs to sample from a multivariate normal distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu}$ is a mean vector and $\boldsymbol{\Sigma}$ is a covariance matrix. This scenario arises constantly in modern deep learning—from generating high-dimensional latent representations to modeling complex dependencies between variables.

The challenge remains the same: direct sampling from $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ is non-differentiable, blocking gradient flow. The reparameterization trick extends naturally to solve this problem.

**The Multivariate Reparameterization**  Instead of sampling $\mathbf{z}$ directly from the complex multivariate distribution, we decompose the process:

1. **Sample standard noise**: Draw $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$—a standard multivariate Gaussian with zero mean and identity covariance

2. **Transform deterministically**: Convert this noise using:

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\boldsymbol{\epsilon} \tag{1.10}$$

   where $\boldsymbol{\Sigma}^{1/2}$ is the matrix square root of the covariance matrix

This transformation preserves the desired statistical properties: the resulting $\mathbf{z}$ follows $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ exactly, but now the sampling process is differentiable with respect to the parameters $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

**Why the Matrix Square Root? The Mathematical Intuition**  The choice of $\boldsymbol{\Sigma}^{1/2}$ isn't arbitrary—it comes from a fundamental requirement about how covariances transform. Let's understand why.

When we transform our standard noise $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ using $\mathbf{z} = \boldsymbol{\mu} + \mathbf{A}\boldsymbol{\epsilon}$, we need to determine what matrix $\mathbf{A}$ will give us the correct covariance structure. But why do we need the "correct" covariance structure at all? Because **we want our samples to actually follow the target distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$**. If we get the covariance wrong, we're sampling from an entirely different distribution—one that might have the wrong correlations, wrong variances, or wrong shape entirely. In applications like VAEs or Bayesian neural networks, this would mean learning completely incorrect relationships between variables.

The covariance of the transformed variable is:

$$\begin{align}
\mathrm{Cov}[\mathbf{z}] &= \mathrm{Cov}[\boldsymbol{\mu} + \mathbf{A}\boldsymbol{\epsilon}] \tag{1.11}\\
&= \mathrm{Cov}[\mathbf{A}\boldsymbol{\epsilon}] \quad (\text{since } \boldsymbol{\mu} \text{ is constant}) \tag{1.12}\\
&= \mathbf{A} \cdot \mathrm{Cov}[\boldsymbol{\epsilon}] \cdot \mathbf{A}^T \tag{1.13}\\
&= \mathbf{A} \cdot \mathbf{I} \cdot \mathbf{A}^T \quad (\text{since } \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})) \tag{1.14}\\
&= \mathbf{A}\mathbf{A}^T \tag{1.15}
\end{align}$$

For our transformed variable $\mathbf{z}$ to have the desired covariance $\boldsymbol{\Sigma}$, we need:

$$\mathbf{A}\mathbf{A}^T = \boldsymbol{\Sigma} \tag{1.16}$$

This equation defines $\mathbf{A}$ as a **matrix square root** of $\boldsymbol{\Sigma}$. Just as in the scalar case where $\sigma = \sqrt{\sigma^2}$, we need $\mathbf{A} = \boldsymbol{\Sigma}^{1/2}$ such that when we "square" it (multiply by its transpose), we recover the original covariance matrix.

**The Geometric Intuition**  Think of the standard normal $\mathcal{N}(\mathbf{0}, \mathbf{I})$ as a perfectly round cloud of points. The matrix $\boldsymbol{\Sigma}^{1/2}$ acts like a transformation that:

- **Stretches** the cloud along different directions according to the eigenvalues of $\boldsymbol{\Sigma}$

- **Rotates** the cloud to align with the eigenvectors of $\boldsymbol{\Sigma}$

The reason we need the square root specifically is that covariance measures how much variables vary *together*, which involves products of deviations. When we apply a linear transformation $\mathbf{A}$ to our data, the covariance gets transformed as $\mathbf{A}\mathbf{A}^T$—so to achieve covariance $\boldsymbol{\Sigma}$, we need $\mathbf{A}$ such that $\mathbf{A}\mathbf{A}^T = \boldsymbol{\Sigma}$.

This is analogous to how in the univariate case, if we want variance $\sigma^2$, we multiply by $\sigma$ (the square root), because $(\sigma \cdot \epsilon)$ has variance $\sigma^2 \cdot \mathrm{Var}[\epsilon] = \sigma^2$.

**Computing the Matrix Square Root**  The matrix square root $\boldsymbol{\Sigma}^{1/2}$ is any matrix $\mathbf{A}$ such that $\mathbf{A}\mathbf{A}^T = \boldsymbol{\Sigma}$. This generalizes the scalar case where $\sqrt{\sigma^2} = \sigma$. There are two common approaches:

1. **Cholesky Decomposition (Preferred):** Since covariance matrices are positive definite, we can write $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$, where $\mathbf{L}$ is lower triangular. Then $\boldsymbol{\Sigma}^{1/2} = \mathbf{L}$.

2. **Eigendecomposition:** If $\boldsymbol{\Sigma} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T$, where $\mathbf{Q}$ contains eigenvectors and $\boldsymbol{\Lambda}$ is diagonal, then:
$$\boldsymbol{\Sigma}^{1/2} = \mathbf{Q}\boldsymbol{\Lambda}^{1/2}\mathbf{Q}^T \tag{1.17}$$

In practice, Cholesky decomposition is preferred for its computational efficiency and numerical stability.

**A Concrete Example**  Let's sample from:
$$\boldsymbol{\mu} = \begin{bmatrix} 1 \\ 2 \end{bmatrix}, \quad \boldsymbol{\Sigma} = \begin{bmatrix} 1.0 & 0.8 \\ 0.8 & 1.0 \end{bmatrix} \tag{1.18}$$

**Step 1:** Compute Cholesky decomposition $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$:
$$\mathbf{L} = \begin{bmatrix} 1.0 & 0 \\ 0.8 & 0.6 \end{bmatrix} \tag{1.19}$$

**Step 2:** Sample noise $\boldsymbol{\epsilon} = \begin{bmatrix} 0.5 \\ -1.0 \end{bmatrix} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

**Step 3:** Apply transformation:
$$\mathbf{z} = \boldsymbol{\mu} + \mathbf{L}\boldsymbol{\epsilon} \tag{1.20}$$
$$= \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 1.0 & 0 \\ 0.8 & 0.6 \end{bmatrix} \begin{bmatrix} 0.5 \\ -1.0 \end{bmatrix} \tag{1.21}$$
$$= \begin{bmatrix} 1 \\ 2 \end{bmatrix} + \begin{bmatrix} 0.5 \\ -0.2 \end{bmatrix} = \begin{bmatrix} 1.5 \\ 1.8 \end{bmatrix} \tag{1.22}$$

The resulting sample $\mathbf{z} = [1.5, 1.8]^T$ follows $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, and the entire process is differentiable with respect to $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$.

**Why This Enables Gradient Flow**

To truly appreciate the power of the reparameterization trick, we need to understand the fundamental challenge it solves in machine learning: **how do we optimize neural networks that involve random sampling?** This is not just a theoretical curiosity—it's the core problem that blocked progress in deep generative models until the reparameterization trick was developed.

**The Central Challenge: Gradients Through Expectations**  In many machine learning scenarios, we need to optimize an objective function that involves an expectation over a parameterized distribution. For instance, in variational autoencoders, we might want to maximize:

$$\mathcal{L}(\boldsymbol{\theta}) = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}[f(\mathbf{z})] \tag{1.23}$$

Here's what each component means:

- $\boldsymbol{\theta}$ are the neural network parameters we want to optimize

- $p(\mathbf{z}|\boldsymbol{\theta})$ is a probability distribution whose parameters depend on $\boldsymbol{\theta}$ (e.g., a Gaussian whose mean and variance are outputs of a neural network)

- $f(\mathbf{z})$ is some function we want to maximize in expectation (e.g., the likelihood of reconstructing input data)

- The expectation $\mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}[f(\mathbf{z})]$ represents the average value of $f(\mathbf{z})$ when $\mathbf{z}$ is sampled from $p(\mathbf{z}|\boldsymbol{\theta})$

To optimize this using gradient descent, we need to compute:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}[f(\mathbf{z})] \tag{1.24}$$

The challenge is that **the distribution itself depends on the parameters we're optimizing**. As we update $\boldsymbol{\theta}$, not only does the function $f(\mathbf{z})$ change, but so does the distribution $p(\mathbf{z}|\boldsymbol{\theta})$ we're sampling from.

**The Traditional Approach: Score Function Estimators (And Why They're Problematic)**  The traditional solution uses the **score function estimator** (also called REINFORCE). Using the identity $\nabla_{\boldsymbol{\theta}} p(\mathbf{z}|\boldsymbol{\theta}) = p(\mathbf{z}|\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \log p(\mathbf{z}|\boldsymbol{\theta})$, we can derive:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}[f(\mathbf{z})] = \mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}[f(\mathbf{z}) \nabla_{\boldsymbol{\theta}} \log p(\mathbf{z}|\boldsymbol{\theta})] \tag{1.25}$$

Let's unpack what this means:

- $\nabla_{\boldsymbol{\theta}} \log p(\mathbf{z}|\boldsymbol{\theta})$ is called the **score function**—it measures how much the log-probability changes with respect to the parameters

- We multiply $f(\mathbf{z})$ by this score function and take the expectation

- Intuitively, this weights each sample $\mathbf{z}$ by both its function value $f(\mathbf{z})$ and how much changing $\boldsymbol{\theta}$ would affect the probability of sampling that particular $\mathbf{z}$

**Why this approach is problematic:**

- **High variance**: The score function $\nabla_{\boldsymbol{\theta}} \log p(\mathbf{z}|\boldsymbol{\theta})$ can vary dramatically across samples, leading to very noisy gradient estimates

- **Complex derivatives**: Computing the score function requires taking derivatives of the log-probability, which can be complex for sophisticated distributions

- **Slow convergence**: The high variance means we need many samples to get reliable gradient estimates, slowing down training

**The Reparameterization Approach: Moving the Gradient Inside**   The reparameterization trick offers an elegant alternative. Instead of sampling directly from $p(\mathbf{z}|\boldsymbol{\theta})$, we express samples as a deterministic function of the parameters and some parameter-free noise:

$$\mathbf{z} = g(\boldsymbol{\theta}, \boldsymbol{\epsilon}) \tag{1.26}$$

where $\boldsymbol{\epsilon}$ is sampled from a simple distribution $p(\boldsymbol{\epsilon})$ that doesn't depend on $\boldsymbol{\theta}$.

Now our expectation becomes:

$$\mathbb{E}_{p(\mathbf{z}|\boldsymbol{\theta})}[f(\mathbf{z})] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[f(g(\boldsymbol{\theta}, \boldsymbol{\epsilon}))] \tag{1.27}$$

The key insight is that we can now compute the gradient as:

$$\nabla_{\boldsymbol{\theta}} \mathbb{E}_{p(\boldsymbol{\epsilon})}[f(g(\boldsymbol{\theta}, \boldsymbol{\epsilon}))] = \mathbb{E}_{p(\boldsymbol{\epsilon})}[\nabla_{\boldsymbol{\theta}} f(g(\boldsymbol{\theta}, \boldsymbol{\epsilon}))] \tag{1.28}$$

**Why can we move the gradient inside the expectation?**  Because $p(\boldsymbol{\epsilon})$ doesn't depend on $\boldsymbol{\theta}$! The distribution we're integrating over is now fixed, so we can interchange the gradient and expectation operators.

**Why This Is Revolutionary**   The reparameterized gradient has several crucial advantages that fundamentally changed how we approach stochastic optimization in deep learning:

**1. Direct Backpropagation**: We can compute $\nabla_{\boldsymbol{\theta}} f(g(\boldsymbol{\theta}, \boldsymbol{\epsilon}))$ using standard automatic differentiation. The chain rule flows naturally through the deterministic function $g(\boldsymbol{\theta}, \boldsymbol{\epsilon})$.

This is revolutionary because it treats stochastic sampling as just another differentiable operation. Consider the computation graph: we have

$$\boldsymbol{\epsilon} \rightarrow g(\boldsymbol{\theta}, \boldsymbol{\epsilon}) \rightarrow f(\cdot).$$

Since $g$ is deterministic, automatic differentiation systems like PyTorch or TensorFlow can compute $\frac{\partial f}{\partial \boldsymbol{\theta}}$ by simply applying the chain rule:

$$\frac{\partial f}{\partial \boldsymbol{\theta}} = \frac{\partial f}{\partial \mathbf{z}} \frac{\partial g}{\partial \boldsymbol{\theta}}.$$

No special-purpose gradient estimators are needed—the same backpropagation algorithm that works for standard neural networks works here. This means existing optimized automatic differentiation frameworks can handle stochastic models without modification.

**2. Lower Variance**: Since we're not multiplying by a score function, the gradient estimates typically have much lower variance. The randomness in $\boldsymbol{\epsilon}$ doesn't interact with the parameter derivatives in the problematic way that score function methods do.

The variance reduction is dramatic because we've eliminated the multiplicative interaction between random quantities. In score function methods, we compute

$$f(\mathbf{z}) \cdot \nabla_{\boldsymbol{\theta}} \log p(\mathbf{z}|\boldsymbol{\theta}),$$

where both terms are random and can vary wildly. The reparameterization approach computes

$$\nabla_{\boldsymbol{\theta}} f(g(\boldsymbol{\theta}, \boldsymbol{\epsilon})),$$

where the randomness enters *additively* through $\boldsymbol{\epsilon}$, not multiplicatively. The gradient computation $\frac{\partial f}{\partial \boldsymbol{\theta}}$ varies smoothly with $\boldsymbol{\theta}$ and $\boldsymbol{\epsilon}$, without the explosive behavior of score functions near distribution boundaries or in low-probability regions. Empirically, this often reduces gradient variance by orders of magnitude.

**3. Computational Efficiency**: We avoid computing complex score functions and can leverage efficient automatic differentiation systems.

Score function methods require computing $\nabla_{\boldsymbol{\theta}} \log p(\mathbf{z}|\boldsymbol{\theta})$, which involves:

- Computing the probability density $p(\mathbf{z}|\boldsymbol{\theta})$ (expensive for complex distributions)

- Taking its logarithm (numerical issues for small probabilities)

- Computing derivatives of the log-probability (algebraically complex)

- Implementing these derivatives correctly for each distribution type

The reparameterization trick sidesteps all of this. We only need to implement the transformation $g(\boldsymbol{\theta}, \boldsymbol{\epsilon})$ and let automatic differentiation handle the rest. For a multivariate Gaussian, instead of deriving and implementing the score function $\nabla_{\boldsymbol{\theta}} \log \mathcal{N}(\mathbf{z}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we simply implement

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2} \boldsymbol{\epsilon}$$

and get gradients "for free." This also leverages highly optimized linear algebra routines and GPU parallelization.

**4. Stable Training**: The more stable gradients lead to more reliable optimization and faster convergence.

The stability improvements manifest in several ways:

- **Consistent gradient directions**: Instead of gradient estimates that swing wildly between positive and negative values, reparameterized gradients provide consistent directional information about how to improve the objective

- **Reduced need for gradient clipping**: Score function methods often require aggressive gradient clipping to prevent training instability; reparameterized gradients rarely need such intervention

- **Better interaction with optimizers**: Adaptive optimizers like Adam work much better with the smoother gradient landscapes provided by reparameterization

- **Faster convergence**: Lower variance means each gradient step provides more reliable information, allowing the optimizer to take larger steps and reach good solutions faster

- **More predictable training curves**: Training loss decreases more smoothly and predictably, making it easier to diagnose problems and tune hyperparameters

These advantages collectively explain why the reparameterization trick enabled the explosion of successful deep generative models starting in the mid-2010s. What had been a computational bottleneck became a routine operation, allowing researchers to focus on model architecture and applications rather than struggling with gradient estimation.

**A Concrete Example: Why We Need This**   Consider training a variational autoencoder where:

- An encoder network outputs $\boldsymbol{\mu}(\mathbf{x})$ and $\boldsymbol{\sigma}(\mathbf{x})$ for input $\mathbf{x}$

- We sample $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}(\mathbf{x}), \mathrm{diag}(\boldsymbol{\sigma}^2(\mathbf{x})))$

- A decoder network reconstructs $\hat{\mathbf{x}}$ from $\mathbf{z}$

- We want to maximize $\mathbb{E}_{p(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]$ (reconstruction likelihood)

Without the reparameterization trick, we couldn't train this model end-to-end because gradients couldn't flow through the sampling step. The reparameterization $\mathbf{z} = \boldsymbol{\mu}(\mathbf{x}) + \boldsymbol{\sigma}(\mathbf{x}) \odot \boldsymbol{\epsilon}$ makes the entire pipeline differentiable, enabling the breakthrough that led to modern generative models.

The reparameterization trick didn't just solve a technical problem—it opened the door to an entire generation of sophisticated probabilistic models that would have been impossible to train otherwise.

**Extensions Beyond Gaussians**

While the multivariate normal case is most common, the reparameterization trick applies to many other distributions:

**Exponential Distribution:**

$$z \sim \mathrm{Exp}(\lambda) \iff z = -\frac{\log(\epsilon)}{\lambda}, \quad \epsilon \sim \mathcal{U}(0, 1) \tag{1.29}$$

**Gamma Distribution (using shape-scale parameterization):** For integer shape parameters, we can use:

$$z \sim \text{Gamma}(\alpha, \beta) \iff z = -\frac{1}{\beta}\sum_{i=1}^{\alpha}\log(\epsilon_i), \quad \epsilon_i \sim \mathcal{U}(0,1) \qquad (1.30)$$

**Beta Distribution:** Using the relationship with Gamma distributions:

$$z \sim \text{Beta}(\alpha, \beta) \iff z = \frac{x}{x+y} \qquad (1.31)$$

where $x \sim \text{Gamma}(\alpha, 1)$ and $y \sim \text{Gamma}(\beta, 1)$.

### Applications in Deep Learning

The reparameterization trick is fundamental to several important deep learning architectures:

**Variational Autoencoders (VAEs):** VAEs use the trick to sample from the latent space while maintaining gradient flow through the encoder network. The encoder outputs parameters $\boldsymbol{\mu}$ and $\boldsymbol{\sigma}^2$ of a diagonal Gaussian, and samples are generated as:

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \qquad (1.32)$$

**Normalizing Flows:** These models apply a sequence of invertible transformations to a simple base distribution, essentially extending the reparameterization trick to create very flexible distributions.

**Stochastic Neural Networks:** Any neural network with stochastic layers can benefit from the reparameterization trick to enable end-to-end gradient-based training.

### Computational Advantages

The reparameterization trick offers several key advantages:

- **Low variance gradients**: Unlike score function estimators, reparameterized gradients typically have much lower variance

- **Efficient computation**: Leverages automatic differentiation systems without requiring specialized gradient estimators

- **Stable training**: The deterministic relationship between parameters and samples leads to more stable optimization

- **Scalability**: Works efficiently with modern GPU architectures and vectorized operations

### Limitations and Considerations

While powerful, the reparameterization trick has some limitations:

- **Distribution constraints**: Not all distributions admit simple reparameterizations

- **Computational overhead**: Some reparameterizations (like matrix square roots) can be expensive

- **Numerical stability**: Operations like Cholesky decomposition can be numerically unstable for ill-conditioned covariance matrices

For these cases, researchers often fall back to score function estimators or develop specialized reparameterization schemes.

The reparameterization trick represents a cornerstone technique in modern deep learning, elegantly bridging the gap between stochastic sampling and gradient-based optimization. Its influence extends far beyond its original application in VAEs, enabling an entire generation of sophisticated generative models and probabilistic neural networks.

# 1.4 Why Images Are Hard

## 1.4.1 The Astronomical Complexity of Images

Now let's return to our original problem: generating realistic images. Consider a modest RGB image of size $256 \times 256$ pixels. Each pixel can take any of $256^3$ possible color values, so the total number of possible images is:

$$(256^3)^{256 \times 256 \times 3} = 2^{24 \times 196608} \approx 2^{4.7 \times 10^6} \tag{1.33}$$

This number is so large that it defies comprehension—far exceeding the number of atoms in the observable universe.

Yet, among this vast space of possibilities, only a tiny fraction corresponds to realistic photographs. The challenge of generative modeling for images is thus to identify and sample from this minuscule subset of the full space.

## 1.4.2 The Manifold Hypothesis

The key insight that makes image generation tractable is the *manifold hypothesis*: realistic images don't fill the entire high-dimensional space uniformly. Instead, they lie on or near a much lower-dimensional manifold embedded within the full space. For instance, while our $256 \times 256$ RGB image lives in a space of dimension $256 \times 256 \times 3 = 196{,}608$, the actual "degrees of freedom" in natural images—the fundamental variations in lighting, pose, identity, etc.—might be captured by just hundreds or thousands of parameters.

This observation suggests a powerful approach: instead of trying to model the full high-dimensional distribution $p(\mathbf{x})$ directly, we can:

1. Work in a lower-dimensional *latent space* with distribution $p(\mathbf{z})$

2. Learn a mapping from this latent space to the data space

## 1.5 The Neural Network Approach

### 1.5.1 From Latent Variables to Data

The core idea of neural generative models is to use a neural network as a flexible function approximator to map from a simple latent distribution to the complex data distribution.

Specifically, we:

1. Choose a simple latent distribution $p(\mathbf{z})$, typically $\mathcal{N}(\mathbf{0}, \mathbf{I})$

2. Design a neural network $D_\theta(\mathbf{z})$ (the *decoder* or *generator*) that maps latent vectors to data points

3. Train the network so that $D_\theta(\mathbf{z})$ produces realistic data when $\mathbf{z} \sim p(\mathbf{z})$

Once trained, generation is simple: sample $\mathbf{z} \sim p(\mathbf{z})$ and compute $\mathbf{x} = D_\theta(\mathbf{z})$.

The power of this approach lies in the flexibility of neural networks. Given sufficient capacity and training data, a neural network can approximate virtually any continuous function, allowing us to model complex transformations from simple latent distributions to realistic data.

### 1.5.2 The Autoencoder Architecture

A natural starting point is the *autoencoder* architecture, which consists of two neural networks:

- An *encoder* $E_\phi(\mathbf{x})$ that maps data points to latent representations

- A *decoder* $D_\theta(\mathbf{z})$ that maps latent representations back to data points

The autoencoder is trained to minimize the reconstruction error:

$$\mathcal{L}(\theta, \phi) = \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}}[\|\mathbf{x} - D_\theta(E_\phi(\mathbf{x}))\|^2] \tag{1.34}$$

For generation, we would ideally sample from the latent distribution and decode: $\mathbf{z} \sim p(\mathbf{z})$, $\mathbf{x} = D_\theta(\mathbf{z})$.

### 1.5.3 The Fundamental Challenge

However, there's a critical problem with this naive approach: *how do we ensure that randomly sampled latent vectors $\mathbf{z}$ correspond to realistic data points?*

The encoder learns to map training data to latent representations, but there's no guarantee that the latent space is organized in a meaningful way. Random samples from the latent space might map to unrealistic outputs, since the decoder was only trained on the specific latent codes produced by the encoder on training data.

This challenge leads us to two major approaches in modern generative modeling:

1. **Generative Adversarial Networks (GANs)**: Use an adversarial training process to ensure the generator produces realistic outputs

2. **Variational Autoencoders (VAEs)**: Impose a prior distribution on the latent space and use variational inference for training

## 1.6   Looking Forward

The statistical foundations and neural network approaches we've explored in this chapter form the bedrock of modern generative modeling. The key insights we've established are:

1. **Statistical Perspective**: View data as samples from an underlying distribution

2. **Sampling Techniques**: Understand how to generate samples from known distributions

3. **Dimensionality Challenges**: Recognize why direct modeling of high-dimensional data is intractable

4. **Latent Variable Approach**: Use neural networks to map from simple latent distributions to complex data distributions

5. **Training Challenges**: Identify the need for sophisticated training procedures to ensure meaningful latent representations

These foundations pave the way for the sophisticated generative models that have revolutionized fields from computer vision to natural language processing. GANs tackle the training challenge through adversarial learning, while VAEs address it through probabilistic inference. More recent approaches, such as diffusion models, offer yet another perspective on the generation problem.

Each of these approaches builds upon the fundamental principles we've established: the statistical view of data, the power of neural networks as function approximators, and the insight that generation can be framed as a mapping from simple to complex distributions. As we delve deeper into specific generative models in subsequent chapters, these foundations will provide the conceptual framework for understanding their design, training, and application.

The journey from statistical sampling to state-of-the-art generative models illustrates a beautiful progression in machine learning: from mathematical theory to practical algorithms, from simple probability distributions to complex neural architectures, and from toy problems to systems that can generate remarkably realistic images, text, and other forms of data. This progression continues to unfold, with new models and techniques constantly pushing the boundaries of what's possible in artificial generation of realistic content.

## 1.7   Exercises

1. Implement inverse transform sampling for the exponential distribution with rate parameter $\lambda$: $p(x) = \lambda e^{-\lambda x}$ for $x \geq 0$.

2. Consider a 2D dataset consisting of points arranged in a circle. Describe the challenges of modeling this distribution with a simple autoencoder, and explain why the latent space organization matters for generation.

3. Derive the acceptance probability for rejection sampling when using $q(x) = c \cdot p(x)$ where $c > 1$ is a constant. What is the optimal value of $c$?

4. Explain why the reparameterization trick is essential for training neural networks with stochastic components. What would happen if we tried to backpropagate through a sampling operation directly?

5. Calculate the effective dimensionality of natural images by considering the number of degrees of freedom in typical image transformations (translation, rotation, scaling, lighting changes, etc.). How does this compare to the nominal dimensionality?

# Chapter 2

# Generative Adversarial Networks: Learning Through Competition

## 2.1   From Cooperation to Competition

In the previous chapter, we established the fundamental challenge of generative modeling: how to learn a mapping from a simple latent distribution to the complex data distribution without knowing the true probability density function. We explored autoencoders as a natural starting point, but identified a critical limitation—there's no guarantee that randomly sampled latent vectors will produce realistic outputs.

This chapter introduces a revolutionary approach to this problem: *Generative Adversarial Networks* (GANs), proposed by Ian Goodfellow and colleagues in 2014. Rather than trying to explicitly model the data distribution, GANs learn to generate realistic data through an adversarial training process—a game-theoretic framework where two neural networks compete against each other.

The key insight behind GANs draws inspiration from an unexpected source: the vulnerability of neural networks to adversarial attacks. Before diving into GANs themselves, let's explore this connection, as it provides crucial intuition for understanding why adversarial training works so well for generation.

## 2.2   The Vulnerability of Neural Networks

### 2.2.1   The Fragility of Deep Learning

Despite their remarkable success across diverse applications, neural networks exhibit a surprising and concerning vulnerability. Consider a state-of-the-art image classifier that correctly identifies a panda in a photograph with high confidence. Now, suppose we add a carefully crafted noise pattern to this image—noise so subtle that it's virtually imperceptible to human eyes. Remarkably, this tiny perturbation can cause the same classifier to misclassify the image as an ostrich with equally high confidence.

These *adversarial examples*, first systematically studied by Szegedy et al. in 2013, reveal that the decision boundaries learned by neural networks are far more fragile than we might expect. A small step in input space can lead to dramatically different outputs, suggesting that neural networks learn to exploit statistical

regularities that don't necessarily align with human perception or robust feature representations.

## 2.2.2   The Mathematics of Adversarial Attacks

Formally, given a neural network classifier $f_\theta(\mathbf{x})$ and an input image $\mathbf{x}$ with true label $y$, an adversarial example $\mathbf{x}_{adv}$ satisfies:

$$\|\mathbf{x}_{adv} - \mathbf{x}\|_p < \epsilon \tag{2.1}$$

$$f_\theta(\mathbf{x}_{adv}) \neq y \tag{2.2}$$

**Breaking down these constraints:**

- $\|\mathbf{x}_{adv} - \mathbf{x}\|_p < \epsilon$: This measures how "close" the adversarial example is to the original image. The $L_p$ norm computes distance—for example, $L_2$ norm is Euclidean distance, while $L_\infty$ norm measures the maximum change in any single pixel. We need $\epsilon$ to be small so the change is imperceptible to humans.

- $f_\theta(\mathbf{x}_{adv}) \neq y$: The classifier must give a different prediction on the perturbed image. This is what makes it "adversarial"—a tiny change completely fools the network.

- Why both constraints? Without the first, we could just show the network a completely different image. Without the second, we haven't actually attacked anything. Together, they define the core challenge: find the smallest change that breaks the classifier.

One of the simplest and most effective methods for generating adversarial examples is the *Fast Gradient Sign Method* (FGSM):

$$\mathbf{x}_{adv} = \mathbf{x} + \epsilon \cdot \text{sign}(\nabla_\mathbf{x}\mathcal{L}(f_\theta(\mathbf{x}), y)) \tag{2.3}$$

**Understanding each component:**

- $\mathbf{x}$: The original, correctly classified image

- $\epsilon$: A small step size (controls how much we perturb the image)

- $\nabla_\mathbf{x}\mathcal{L}(f_\theta(\mathbf{x}), y)$: The gradient of the loss with respect to the input image. This tells us how to change each pixel to maximally increase the loss (make the classifier more wrong)

- $\text{sign}(\cdot)$: Takes only the direction of the gradient (positive or negative), ignoring magnitude. This ensures we take equal-sized steps in each pixel dimension

- Why this works: We're asking "which direction should I nudge each pixel to make the classifier most confused?" The gradient points in the direction of steepest increase in loss. By taking a small step in this direction, we can often fool the classifier with minimal visual change.

### 2.2.3  The Arms Race: Attack and Defense

The discovery of adversarial examples naturally leads to an arms race between attackers and defenders:

- **Attackers** develop increasingly sophisticated methods to fool neural networks

- **Defenders** attempt to make networks more robust through adversarial training, regularization, or architectural modifications

This adversarial dynamic suggests a powerful training paradigm: what if we could harness this competition to improve our models? Instead of viewing adversarial examples as a security threat, could we use the adversarial training process constructively?

## 2.3  The GAN Framework

### 2.3.1  From Attack to Generation

The breakthrough insight of GANs is to reframe the adversarial paradigm for generative modeling. Instead of attacking a fixed classifier, we train two networks simultaneously:

1. A **Generator** $G_\theta(\mathbf{z})$ that tries to create realistic data from random noise

2. A **Discriminator** $D_\phi(\mathbf{x})$ that tries to distinguish between real data and generated samples

The generator plays the role of an "adversarial attack network"—but instead of slightly perturbing real images to fool a classifier, it starts from pure noise and tries to synthesize completely realistic images. The discriminator acts as the classifier being attacked, attempting to detect these synthetic images.

### 2.3.2  The Architecture

Let's examine the components in detail:

**Generator Network**

The generator $G_\theta : \mathbb{R}^d \to \mathbb{R}^n$ is a neural network that maps from a $d$-dimensional latent space to the $n$-dimensional data space. Typically:

- Input: $\mathbf{z} \sim p(\mathbf{z})$, where $p(\mathbf{z})$ is a simple distribution (e.g., $\mathcal{N}(\mathbf{0}, \mathbf{I})$)

- Output: $G_\theta(\mathbf{z}) \in \mathbb{R}^n$, a synthetic data sample

- Architecture: Often uses transposed convolutions (for images) to progressively upsample from low-resolution latent codes to high-resolution outputs

**Discriminator Network**

The discriminator $D_\phi : \mathbb{R}^n \to [0, 1]$ is a binary classifier that estimates the probability that its input is real data:

- Input: Either real data $\mathbf{x} \sim p_{data}(\mathbf{x})$ or generated data $G_\theta(\mathbf{z})$

- Output: $D_\phi(\mathbf{x}) \in [0, 1]$, the probability that the input is real

- Architecture: Often uses standard convolutional layers for downsampling and classification

### 2.3.3   The Adversarial Objective

The training objective for GANs is formulated as a minimax game:

$$\min_G \max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))] \qquad (2.4)$$

**Understanding the structure:**

- $\min_G \max_D$: This is a two-player game. The discriminator $D$ wants to maximize the objective (be good at classification), while the generator $G$ wants to minimize it (fool the discriminator). It's like a counterfeiter (generator) trying to make fake money while a detective (discriminator) tries to catch the fakes.

- $V(D, G)$: The "value function" that measures how well each player is doing

- The equation has two terms—let's examine each one separately

**First term: $\mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})]$**

- $\mathbf{x} \sim p_{data}(\mathbf{x})$: We sample real data from our training dataset

- $D(\mathbf{x})$: The discriminator outputs a probability that $\mathbf{x}$ is real (should be close to 1 for real data)

- $\log D(\mathbf{x})$: Taking the log converts probabilities to "log-odds." When $D(\mathbf{x}) = 1$ (perfect classification), $\log D(\mathbf{x}) = 0$. When $D(\mathbf{x}) = 0.5$ (random guessing), $\log D(\mathbf{x}) = \log(0.5) < 0$

- $\mathbb{E}[\cdot]$: We average over all real data samples

- **Intuition**: The discriminator wants this term to be as large as possible (close to 0), meaning it correctly identifies real data as real

**Second term: $\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$**

- $\mathbf{z} \sim p(\mathbf{z})$: We sample random noise from our latent distribution (usually standard normal)

- $G(\mathbf{z})$: The generator creates a fake image from the noise

- $D(G(\mathbf{z}))$: The discriminator outputs the probability that this fake image is real (should be close to 0 for fake data)

- $1 - D(G(\mathbf{z}))$: This flips the probability—now it represents the probability that the fake image is fake (should be close to 1)

- $\log(1 - D(G(\mathbf{z})))$: Again, log of this probability

- **Intuition**: The discriminator wants this term large (close to 0), meaning it correctly identifies fake data as fake

**Why this objective makes sense:**

- When the discriminator is winning: $D(\mathbf{x}) \approx 1$ and $D(G(\mathbf{z})) \approx 0$, so $V(D, G) \approx 0 + 0 = 0$ (maximum possible)

- When the generator is winning: $D(G(\mathbf{z})) \approx 1$ (fake data looks real), so the second term becomes $\log(1 - 1) = \log(0) = -\infty$ (minimum possible)

- At equilibrium: $D(\mathbf{x}) = D(G(\mathbf{z})) = 0.5$ (discriminator can't tell real from fake), and both players have no incentive to change

Let's unpack this objective function piece by piece:

## The Discriminator's Perspective

From the discriminator's viewpoint, this is a standard binary classification problem:

$$\max_D V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))] \qquad (2.5)$$

**The discriminator's goal broken down:**

- **Maximize $D(\mathbf{x})$ for real data**: We want $D(\mathbf{x}) \to 1$, which makes $\log D(\mathbf{x}) \to 0$ (the maximum possible value since $\log(1) = 0$ and $\log(p) < 0$ for $p < 1$)

- **Minimize $D(G(\mathbf{z}))$ for generated data**: We want $D(G(\mathbf{z})) \to 0$, which makes $(1 - D(G(\mathbf{z}))) \to 1$ and $\log(1 - D(G(\mathbf{z}))) \to 0$

- **Why logarithms?** This is exactly the negative log-likelihood for binary classification:

    - Real data has label $y = 1$: contribute $\log D(\mathbf{x}) = \log P(y = 1|\mathbf{x})$
    - Fake data has label $y = 0$: contribute $\log(1 - D(G(\mathbf{z}))) = \log P(y = 0|G(\mathbf{z}))$

- **Perfect discriminator**: Would achieve $V(D^*, G) = 0 + 0 = 0$, meaning it never makes mistakes

This is equivalent to minimizing the binary cross-entropy loss for a classifier with labels $y = 1$ for real data and $y = 0$ for fake data.

**The Generator's Perspective**

The generator aims to minimize the same objective:

$$\min_G V(D, G) = \mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))] \qquad (2.6)$$

**The generator's strategy:**

- **First term is irrelevant**: $\mathbb{E}_{\mathbf{x} \sim p_{data}(\mathbf{x})}[\log D(\mathbf{x})]$ doesn't depend on the generator parameters $\theta_G$, so the generator can't control it

- **Focus on the second term**: The generator effectively tries to minimize $\mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))]$

- **What this means**:

  - To minimize $\log(1 - D(G(\mathbf{z})))$, we need to maximize $D(G(\mathbf{z}))$
  - This means the generator wants the discriminator to think its fake images are real
  - When $D(G(\mathbf{z})) \to 1$, we get $\log(1 - 1) = \log(0) = -\infty$ (the minimum possible)

- **Intuitive interpretation**: The generator is trying to "fool" the discriminator by making fake data so realistic that $D(G(\mathbf{z})) \approx 1$

**Alternative generator loss:** In practice, early in training when generated samples are poor, $D(G(\mathbf{z})) \approx 0$, making $\log(1 - D(G(\mathbf{z}))) \approx \log(1) = 0$. The gradient of this is nearly zero, providing little learning signal. Instead, many implementations use:

$$\max_G \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log D(G(\mathbf{z}))] \qquad (2.7)$$

This is equivalent to the original objective at convergence but provides stronger gradients early in training.

## 2.3.4   The Training Dynamics

The adversarial training process alternates between two phases:

1. **Discriminator Update**: Fix the generator and train the discriminator to better distinguish real from fake data

2. **Generator Update**: Fix the discriminator and train the generator to better fool the discriminator

In the ideal case, this process converges to a Nash equilibrium where:

- The generator produces samples indistinguishable from real data

- The discriminator cannot do better than random guessing ($D(\mathbf{x}) = 0.5$ for all $\mathbf{x}$)

### 2.3.5 Theoretical Analysis

Under ideal conditions, Goodfellow et al. proved that the global optimum of the GAN objective corresponds to the generator learning the true data distribution. Specifically, when the discriminator is optimal for a fixed generator:

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_g(\mathbf{x})} \tag{2.8}$$

**Understanding the optimal discriminator formula:**

- $p_{data}(\mathbf{x})$: The true probability density of real data at point $\mathbf{x}$

- $p_g(\mathbf{x})$: The probability density of generated data at point $\mathbf{x}$ (induced by the generator)

- **Intuition**: This formula tells us the optimal probability that $\mathbf{x}$ is real, given that it came from either the real data or the generator

- **Bayes' rule interpretation**: If we have two sources of data (real and generated) with equal prior probability, this is exactly Bayes' rule for classification

- **Special cases**:

  - If $p_g(\mathbf{x}) = 0$ (generator never produces $\mathbf{x}$): $D_G^*(\mathbf{x}) = 1$ (definitely real)
  - If $p_{data}(\mathbf{x}) = 0$ (real data never has $\mathbf{x}$): $D_G^*(\mathbf{x}) = 0$ (definitely fake)
  - If $p_{data}(\mathbf{x}) = p_g(\mathbf{x})$ (perfect generator): $D_G^*(\mathbf{x}) = 0.5$ (can't tell real from fake)

where $p_g(\mathbf{x})$ is the distribution induced by the generator. Substituting this optimal discriminator back into the objective shows that minimizing the generator loss is equivalent to minimizing the Jensen-Shannon divergence between $p_{data}$ and $p_g$:

$$\min_G V(D_G^*, G) = -\log(4) + 2 \cdot \mathrm{JS}(p_{data}\|p_g) \tag{2.9}$$

**Breaking down the Jensen-Shannon result:**

- $\mathrm{JS}(p_{data}\|p_g)$: Jensen-Shannon divergence between real and generated distributions

- **What is JS divergence?** It measures how different two probability distributions are:

  - $\mathrm{JS}(p\|q) = 0$ when $p = q$ (distributions are identical)
  - $\mathrm{JS}(p\|q) > 0$ when $p \neq q$ (distributions differ)
  - It's symmetric: $\mathrm{JS}(p\|q) = \mathrm{JS}(q\|p)$
  - It's bounded: $0 \leq \mathrm{JS}(p\|q) \leq \log(2)$

- $-\log(4)$: A constant offset that doesn't affect optimization

- $2 \cdot \mathrm{JS}(p_{data}\|p_g)$: The factor of 2 comes from the mathematical derivation

- **Key insight**: Minimizing the GAN objective is equivalent to making the generated distribution as close as possible to the real distribution in terms of JS divergence

- **Optimality**: The global minimum is achieved when $p_g = p_{data}$, meaning the generator perfectly learns the real data distribution

This theoretical result provides strong justification for the GAN approach: at convergence, the generator should recover the true data distribution.

## 2.4 Challenges in GAN Training

While theoretically elegant, GANs are notoriously difficult to train in practice. The minimax optimization problem presents several fundamental challenges that have driven much of the subsequent research in this area.

### 2.4.1 Non-Convergence and Training Instability

**The Problem of Oscillation**

Unlike standard supervised learning where we minimize a single loss function, GAN training involves a two-player game with potentially conflicting objectives. This can lead to oscillatory behavior where the networks never reach a stable equilibrium.

Consider a simplified scenario where the generator and discriminator have equal learning rates and capabilities. As the discriminator improves, it becomes better at identifying fake samples, providing stronger gradients to the generator. However, as the generator improves, it becomes harder for the discriminator to distinguish real from fake data. This back-and-forth dynamic can continue indefinitely without convergence.

**Gradient Vanishing and Exploding**

A particularly problematic scenario occurs when the discriminator becomes too powerful too quickly. If the discriminator perfectly separates real and fake data early in training, it will output probabilities very close to 1 for real data and very close to 0 for fake data. This leads to:

**Mathematical analysis of the vanishing gradient problem:**

The generator loss is $\mathcal{L}_G = \mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]$. The gradient with respect to generator parameters is:

$$\frac{\partial \mathcal{L}_G}{\partial \theta_G} = \mathbb{E}_{\mathbf{z}}\left[\frac{\partial \log(1 - D(G(\mathbf{z})))}{\partial \theta_G}\right] \tag{2.10}$$

Using the chain rule:

$$\frac{\partial \log(1 - D(G(\mathbf{z})))}{\partial \theta_G} = \frac{-1}{1 - D(G(\mathbf{z}))} \cdot \frac{\partial D(G(\mathbf{z}))}{\partial G(\mathbf{z})} \cdot \frac{\partial G(\mathbf{z})}{\partial \theta_G} \tag{2.11}$$

**Why gradients vanish:**

- When $D(G(\mathbf{z})) \to 0$ (discriminator correctly identifies fake data):

  - $1 - D(G(\mathbf{z})) \to 1$, so $\log(1 - D(G(\mathbf{z}))) \to 0$
  - BUT: $\frac{\partial D(G(\mathbf{z}))}{\partial G(\mathbf{z})} \to 0$ because the discriminator is "saturated"

- The discriminator outputs are "flat" in regions where it's confident, providing no gradient signal

- **Analogy**: It's like trying to improve a clearly fake painting when the art expert (discriminator) immediately says "obviously fake" without explaining what's wrong

**Alternative loss to fix vanishing gradients:** Instead of minimizing $\mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]$, maximize:

$$\mathcal{L}_G^{alt} = \mathbb{E}_{\mathbf{z}}[\log D(G(\mathbf{z}))] \tag{2.12}$$

**Why this helps:**

- When $D(G(\mathbf{z})) \to 0$: $\log D(G(\mathbf{z})) \to -\infty$ (strong negative signal)

- The gradient is $\frac{1}{D(G(\mathbf{z}))} \cdot \frac{\partial D(G(\mathbf{z}))}{\partial G(\mathbf{z})} \cdot \frac{\partial G(\mathbf{z})}{\partial \theta_G}$

- As $D(G(\mathbf{z})) \to 0$, the factor $\frac{1}{D(G(\mathbf{z}))}$ grows large, amplifying the gradient

- This provides strong learning signal even when the discriminator is confident the data is fake

- **Vanishing gradients**: $\log(1 - D(G(\mathbf{z}))) \to -\infty$ as $D(G(\mathbf{z})) \to 0$, but the gradient $\frac{\partial}{\partial G} \log(1 - D(G(\mathbf{z})))$ approaches zero

- **Loss of signal**: The generator receives little useful feedback about how to improve

Conversely, if the generator becomes too powerful, it might exploit weaknesses in the discriminator without learning to generate realistic data.

## Hyperparameter Sensitivity

GAN training is extremely sensitive to hyperparameters:

- **Learning rates**: Different learning rates for generator and discriminator can destabilize training

- **Architecture choices**: The relative capacity of the two networks critically affects training dynamics

- **Normalization**: Batch normalization, layer normalization, and other techniques can have dramatic effects on stability

## 2.4.2 Mode Collapse

**The Phenomenon**

Mode collapse is perhaps the most frustrating challenge in GAN training. Instead of learning to generate the full diversity of the training data, the generator "collapses" to producing samples from only a few modes of the data distribution.

For example, when training on a dataset of handwritten digits (0-9), a mode-collapsed GAN might only generate high-quality images of 1s and 7s, completely ignoring the other digits. The generator discovers that it can fool the discriminator by producing very realistic samples from a limited subset of the data distribution.

**Why Mode Collapse Occurs**

Mode collapse arises from the fundamental asymmetry in the GAN objective. The generator only needs to produce samples that fool the current discriminator—it's not explicitly encouraged to cover the entire data distribution. If the generator finds a small region of data space where it can consistently fool the discriminator, it has no incentive to explore other regions.

Mathematically, the generator minimizes:

$$\mathcal{L}_G = \mathbb{E}_{\mathbf{z} \sim p(\mathbf{z})}[\log(1 - D(G(\mathbf{z})))] \tag{2.13}$$

**Understanding why this leads to mode collapse:**

- **What the objective measures**: This is the expected log-probability that generated samples are classified as fake

- **The averaging problem**: The expectation $\mathbb{E}_{\mathbf{z}}[\cdot]$ averages over all latent codes $\mathbf{z}$, but doesn't care about diversity

- **Mathematical example**: Consider two strategies:

  1. Strategy A: Generate perfect samples from 10% of modes: $\mathbb{E}[\log(1 - 0.01)] = \log(0.99) \approx -0.01$

  2. Strategy B: Generate good samples from all modes: $\mathbb{E}[\log(1 - 0.1)] = \log(0.9) \approx -0.105$

- Strategy A achieves a better (lower) loss value despite ignoring 90% of the data distribution!

- **No diversity term**: Unlike maximum likelihood estimation, which would include a term encouraging coverage of all data points, the GAN objective has no explicit diversity reward

- **Local optima**: Once the generator finds a successful mode, gradient descent has no incentive to explore other regions

**Why this is different from other generative models:**

- **Maximum Likelihood**: Tries to maximize $\mathbb{E}_{\mathbf{x} \sim p_{data}}[\log p_g(\mathbf{x})]$, which requires assigning probability mass to all observed data points

- **GANs**: Only require that generated samples fool the discriminator, with no direct penalty for ignoring parts of the data distribution

This objective cares about the expected probability over generated samples, but places no penalty on lack of diversity. A generator that produces perfect samples from 10% of the modes will achieve a better objective value than one that produces slightly imperfect samples from all modes.

**Detecting Mode Collapse**

Mode collapse can be subtle and difficult to detect:

- **Visual inspection**: Generated samples may look high-quality but lack diversity

- **Inception Score (IS)**: Measures both quality and diversity of generated samples

- **Fréchet Inception Distance (FID)**: Compares the distribution of generated samples to real data

### 2.4.3 Evaluation Challenges

Unlike supervised learning where we have clear metrics like accuracy, evaluating GANs is inherently difficult:

- **No ground truth**: We don't have access to the true data distribution

- **Subjective quality**: Visual quality is often evaluated by human inspection

- **Mode coverage**: Ensuring the generator covers all modes of the data distribution

- **Sample quality vs. diversity trade-off**: High-quality samples from few modes vs. diverse samples with lower quality

## 2.5 Practical Considerations and Improvements

Despite these challenges, GANs have achieved remarkable success through various improvements and best practices developed by the research community.

### 2.5.1 Training Techniques

**Alternative Loss Functions**

The original GAN loss can be replaced with more stable alternatives:

- **Wasserstein GAN (WGAN)**: Uses Earth Mover's distance instead of JS divergence

- **Least Squares GAN (LSGAN)**: Replaces sigmoid loss with least squares

- **Hinge Loss**: Adapts the hinge loss from SVMs to the GAN setting

**Architectural Innovations**

- **Deep Convolutional GANs (DCGANs)**: Established best practices for convolutional architectures

- **Progressive GANs**: Gradually increase resolution during training

- **Self-Attention GANs**: Incorporate attention mechanisms for better long-range dependencies

**Regularization Methods**

- **Gradient Penalty**: Regularizes the discriminator's gradients to improve stability

- **Spectral Normalization**: Constrains the Lipschitz constant of the discriminator

- **Batch/Layer Normalization**: Stabilizes training dynamics

## 2.6 The Impact and Legacy of GANs

GANs have fundamentally changed the landscape of generative modeling by demonstrating that adversarial training can produce remarkably realistic synthetic data. Their impact extends far beyond computer vision:

- **Creative Applications**: Art generation, style transfer, and creative tools

- **Data Augmentation**: Generating additional training data for other tasks

- **Scientific Applications**: Drug discovery, climate modeling, and simulation

- **Privacy**: Generating synthetic datasets that preserve statistical properties while protecting individual privacy

More importantly, GANs introduced the powerful concept of adversarial training, which has influenced numerous other areas of machine learning, from domain adaptation to robustness to fairness.

## 2.7 Looking Forward

While this chapter has focused on the foundational GAN framework, the field has continued to evolve rapidly. Subsequent developments include:

- **Conditional GANs**: Incorporating additional information to control generation

- **CycleGANs**: Learning mappings between different domains without paired data

- **StyleGANs**: Achieving unprecedented control over generated image style and content

41

- **BigGANs**: Scaling to large datasets and high resolutions

The next chapter will explore an alternative approach to the challenges we've discussed: Variational Autoencoders (VAEs), which tackle the generation problem through probabilistic inference rather than adversarial training. While GANs learn through competition, VAEs learn through approximation, offering different trade-offs between training stability, theoretical grounding, and generation quality.

The comparison between these two paradigms—adversarial training vs. variational inference—illustrates the rich diversity of approaches in modern generative modeling and sets the stage for understanding even more recent developments like diffusion models.

## 2.8 Exercises

1. Derive the optimal discriminator $D_G^*(\mathbf{x})$ for a fixed generator $G$, and show that the resulting generator objective is equivalent to minimizing the Jensen-Shannon divergence.

2. Explain why the generator loss $\mathbb{E}_{\mathbf{z}}[\log(1 - D(G(\mathbf{z})))]$ can lead to vanishing gradients early in training. Propose an alternative loss function that mitigates this issue.

3. Consider a 1D example where real data comes from a mixture of two Gaussians. Describe how mode collapse might manifest in this setting and why it occurs.

4. Implement a simple GAN for generating 2D points from a ring distribution. Experiment with different learning rates for the generator and discriminator and observe the training dynamics.

5. Compare and contrast the advantages and disadvantages of GANs versus traditional density estimation methods (e.g., kernel density estimation, mixture models) for generative modeling.

# Chapter 3

# Variational Autoencoders

## 3.1 Introduction

In the previous chapter, we explored Generative Adversarial Networks (GANs), which learn to generate realistic data through an adversarial training process. While GANs can produce impressive results, they suffer from several challenges: training instability, mode collapse, and the difficulty of solving minimax optimization problems. These issues naturally lead us to ask: *Is there a way to create generative models without the complications of adversarial training?*

This chapter introduces Variational Autoencoders (VAEs), proposed by Kingma and Welling in 2014, which offer an elegant alternative approach. Instead of pitting two networks against each other, VAEs frame generative modeling as a probabilistic inference problem. They combine the intuitive encoder-decoder architecture we encountered in Chapter 1 with principled probabilistic foundations, creating a theoretically grounded and practically stable training procedure.

The key insight behind VAEs is to model the relationship between data and latent variables explicitly as a conditional probability distribution, then use variational inference to make the intractable problem of learning this distribution tractable. This approach not only avoids the minimax optimization challenges of GANs but also provides a principled framework for understanding what the model has learned.

## 3.2 Mathematical Foundations

Before diving into VAEs, we need to establish several key concepts from probability theory that form the mathematical backbone of variational inference. These concepts might seem abstract at first, but they are the essential tools that make VAEs work.

### 3.2.1 Marginal Distributions

**Definition 3.1** (Marginal Distribution)**.** The marginal distribution of a subset of random variables is the probability distribution obtained by integrating (or summing) over all other variables in the joint distribution.

For continuous variables, if we have a joint distribution $p(\mathbf{x}, \mathbf{z})$, the marginal distribution of $\mathbf{x}$ is:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z}) \, d\mathbf{z} \tag{3.1}$$

**Intuitive understanding:**

- Think of $\mathbf{z}$ as a "hidden" or latent variable that influences $\mathbf{x}$

- The marginal $p(\mathbf{x})$ tells us the probability of observing $\mathbf{x}$ regardless of what value $\mathbf{z}$ takes

- We "marginalize out" $\mathbf{z}$ by considering all possible values it could take, weighted by their probabilities

- **Example**: If $\mathbf{x}$ represents the brightness of a photo and $\mathbf{z}$ represents the time of day it was taken, $p(\mathbf{x})$ tells us how likely different brightness levels are across all times of day

### 3.2.2 Expected Values

**Definition 3.2** (Expected Value). The expected value of a random variable is the average value it takes, weighted by the probability of each outcome.

For a continuous random variable $\mathbf{x}$ with distribution $p(\mathbf{x})$:

$$\mathbb{E}_{p(\mathbf{x})}[\mathbf{x}] = \int \mathbf{x} \cdot p(\mathbf{x}) \, d\mathbf{x} \tag{3.2}$$

**Why expected values matter in VAEs:**

- VAEs involve computing expectations over distributions we can't solve analytically

- We'll approximate these expectations using Monte Carlo sampling: draw samples and average

- **Key insight**: Instead of computing $\mathbb{E}_{p(\mathbf{z})}[f(\mathbf{z})]$ exactly, we approximate it as $\frac{1}{N} \sum_{i=1}^{N} f(\mathbf{z}_i)$ where $\mathbf{z}_i \sim p(\mathbf{z})$

### 3.2.3 Bayes' Rule: The Foundation of Inference

Bayes' rule provides the mathematical foundation for inference—learning about hidden variables from observed data:

$$p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z}) \cdot p(\mathbf{z})}{p(\mathbf{x})} \tag{3.3}$$

**Breaking down each component:**

- $p(\mathbf{z}|\mathbf{x})$ (Posterior): What we want to know—the probability of latent variable $\mathbf{z}$ given observed data $\mathbf{x}$

- $p(\mathbf{x}|\mathbf{z})$ (Likelihood): How likely the data is, given a particular latent variable value

- $p(\mathbf{z})$ (Prior): Our initial belief about the latent variable before seeing data

- $p(\mathbf{x})$ (Evidence/Marginal): The total probability of observing the data, summed over all possible latent values

**The challenge for generative models:**

- We want to learn $p(\mathbf{z}|\mathbf{x})$ to understand what latent factors generated our data

- But computing $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})\,d\mathbf{z}$ is usually intractable

- VAEs solve this by approximating $p(\mathbf{z}|\mathbf{x})$ with a simpler distribution $q_\phi(\mathbf{z}|\mathbf{x})$

### 3.2.4 Kullback-Leibler (KL) Divergence

**Definition 3.3** (KL Divergence). KL divergence measures how different one probability distribution is from another:

$$D_{KL}(p\|q) = \int p(\mathbf{x}) \log \frac{p(\mathbf{x})}{q(\mathbf{x})}\,d\mathbf{x} = \mathbb{E}_{p(\mathbf{x})}\left[\log \frac{p(\mathbf{x})}{q(\mathbf{x})}\right] \tag{3.4}$$

**Understanding KL divergence intuitively:**

- $D_{KL}(p\|q) = 0$ if and only if $p = q$ (distributions are identical)

- $D_{KL}(p\|q) > 0$ when $p \neq q$ (always non-negative)

- **Asymmetric**: $D_{KL}(p\|q) \neq D_{KL}(q\|p)$ in general

- **Interpretation**: How much extra information (in nats) we need when using $q$ to approximate $p$

**Special case for Gaussians:** When $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \sigma^2 \mathbf{I})$ and $q(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \mathbf{0}, \mathbf{I})$:

$$D_{KL}(p\|q) = \frac{1}{2}\left(\|\boldsymbol{\mu}\|^2 + d\sigma^2 - d - d\log\sigma^2\right) \tag{3.5}$$

where $d$ is the dimensionality. This closed form will be crucial for VAE training.

### 3.2.5 Jensen's Inequality: The Key to Tractable Bounds

**Definition 3.4** (Convex Function). A function $f$ is convex if for all $\mathbf{x}_1, \mathbf{x}_2$ and $t \in [0, 1]$:

$$f(t\mathbf{x}_1 + (1-t)\mathbf{x}_2) \leq tf(\mathbf{x}_1) + (1-t)f(\mathbf{x}_2) \tag{3.6}$$

**Theorem 3.1** (Jensen's Inequality). *If $f$ is a convex function and $\mathbf{x}$ is a random variable, then:*

$$f(\mathbb{E}[\mathbf{x}]) \leq \mathbb{E}[f(\mathbf{x})] \tag{3.7}$$

*For concave functions, the inequality flips: $f(\mathbb{E}[\mathbf{x}]) \geq \mathbb{E}[f(\mathbf{x})]$.*

**Why Jensen's inequality is crucial for VAEs:**

- The logarithm function is concave: $\log(\mathbb{E}[X]) \geq \mathbb{E}[\log(X)]$

- This allows us to create a lower bound on $\log p(\mathbf{x})$, which we can optimize

- Without this inequality, we couldn't derive the tractable VAE objective

## 3.3 The VAE Framework

### 3.3.1 The Generative Model

To understand Variational Autoencoders, we must first grasp the elegant probabilistic story they tell about how data comes into existence. Imagine you're trying to explain how faces are generated in nature, or how handwritten digits come to be. VAEs propose a beautifully simple two-step process that captures the essence of data generation.

**The Data Generation Story**

VAEs start with a compelling narrative: all observable data (images, text, audio) emerges from some underlying, unobservable factors that we call **latent variables** or **latent codes**. These latent variables capture the essential, abstract properties that determine what the final data will look like.

For example, when generating a face:

- The latent variables might encode abstract concepts like "how happy the person looks," "the lighting angle," "age," or "hair color"

- These abstract factors then get transformed into the concrete pixels we observe

- Different combinations of these latent factors produce different faces

The VAE formalizes this intuition with a precise probabilistic model:

1. **Sample latent factors**: $\mathbf{z} \sim p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$

2. **Generate observable data**: $\mathbf{x} \sim p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D_\theta(\mathbf{z}), \sigma^2 \mathbf{I})$

**Dissecting the Model Components**

**The Prior Distribution:** $p(\mathbf{z}) = \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{I})$ The prior distribution represents our "default belief" about the latent factors before we see any data. Think of it as the natural distribution of abstract concepts in the world.

**Why a standard normal distribution?**

- **Simplicity**: The standard normal $\mathcal{N}(\mathbf{0}, \mathbf{I})$ has no parameters to learn, making the model simpler

- **Symmetry**: It treats all latent dimensions equally and assumes they're independent

- **Easy sampling**: We can easily generate new latent codes by sampling from this well-known distribution

- **Regularization**: It encourages the model to use the latent space efficiently by preferring codes near the origin

**Intuitive interpretation**: Each dimension of $\mathbf{z}$ is like a "dial" that controls some aspect of the generated data. The standard normal prior says that, by default, all these dials are set to around zero (neutral position), with some natural variation following a bell curve.

**The Likelihood Function:** $p(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; D_\theta(\mathbf{z}), \sigma^2\mathbf{I})$   The likelihood function describes how latent codes get transformed into observable data. This is where the magic happens—where abstract concepts become concrete observations.

**The decoder network $D_\theta(\mathbf{z})$:**

- This is a neural network that acts as a "universal translator" from abstract latent space to data space

- Input: A latent code $\mathbf{z}$ (e.g., a 128-dimensional vector encoding abstract face properties)

- Output: The mean of the data distribution (e.g., pixel intensities for a generated face image)

- Parameters $\theta$: The weights and biases that determine exactly how this translation works

**Why Gaussian noise around the decoder output?**   The Gaussian assumption $\mathcal{N}(\mathbf{x}; D_\theta(\mathbf{z}), \sigma^2\mathbf{I})$ captures the idea that:

- The decoder gives us the "most likely" data for a given latent code

- But there's still some randomness—the same latent code might produce slightly different outputs

- This randomness accounts for factors too complex or subtle to capture in the latent code

- The fixed variance $\sigma^2\mathbf{I}$ assumes this noise is the same across all data dimensions (can be relaxed in practice)

**The Complete Generative Process**

Putting it all together, here's how a VAE believes data is born:

**Step 1 - Sample Abstract Concepts**: Nature first "decides" on the abstract properties by sampling $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For a face, this might determine values like:

$$z_1 = 0.3 \quad \text{(slightly positive emotion)} \tag{3.8}$$
$$z_2 = -1.2 \quad \text{(younger age)} \tag{3.9}$$
$$z_3 = 0.8 \quad \text{(good lighting)} \tag{3.10}$$
$$\vdots \tag{3.11}$$

**Step 2 - Manifest Concrete Data**: These abstract factors then get transformed by the decoder network $D_\theta(\mathbf{z})$ into concrete pixel values, with some Gaussian noise added to account for fine details the model can't capture perfectly.

**Generation in Practice**:

- **To generate new data**: Simply sample $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and compute $D_\theta(\mathbf{z})$

- **To generate specific types of data**: Learn what latent codes correspond to desired properties, then sample around those regions

- **To interpolate between data points**: Take two latent codes and smoothly interpolate between them in latent space

**The Key Insight: Latent Space Organization**

The fundamental assumption underlying VAEs is that **similar latent codes should produce similar data**. This means the decoder network $D_\theta$ must learn to organize the latent space in a meaningful way:

- Nearby points in $\mathbf{z}$-space should correspond to similar-looking images

- Moving smoothly through latent space should produce smooth transformations in data space

- Different directions in latent space should correspond to interpretable changes (e.g., changing age, pose, or expression)

This organization doesn't happen automatically—it emerges from the training process, where the model learns to compress real data into meaningful latent representations and then reconstruct it faithfully. The beauty of the VAE framework is that it provides a principled way to encourage this kind of structured latent space through probabilistic inference, which we'll explore in the next section.

The generative model we've described here is only half the story. To actually train such a model, we need to solve the inverse problem: given observed data $\mathbf{x}$, how do we infer the latent code $\mathbf{z}$ that most likely generated it? This is where the encoder network and variational inference come into play.

## 3.3.2 The Learning Problem: Maximum Likelihood with Latent Variables

Now we face the central challenge of training our generative model. We have a dataset $\{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_N\}$ of real observations (say, a million face images), and we want to learn the decoder parameters $\theta$ so that our model can generate realistic new faces. The natural approach is maximum likelihood estimation: find the parameters that make our observed data as probable as possible under our model.

**The Marginal Likelihood: What We Want to Maximize**

For each data point $\mathbf{x}$, the likelihood under our generative model is:

$$p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z}) \, d\mathbf{z} \tag{3.12}$$

Let's unpack what this integral means intuitively. For a given face image $\mathbf{x}$, we're asking: "What's the probability that our generative model would produce this exact face?" To answer this, we need to consider **all possible latent codes $\mathbf{z}$** that could have generated this face:

- Maybe $\mathbf{z}_1$ encodes "young person, smiling, good lighting" and produces something very close to $\mathbf{x}$

- Maybe $\mathbf{z}_2$ encodes "older person, neutral expression, dim lighting" and produces something far from $\mathbf{x}$

- We need to weight each possibility by how likely that latent code is under the prior $p(\mathbf{z})$

- Then we sum (integrate) over all these possibilities to get the total probability of observing $\mathbf{x}$

## Why This Integral is Computationally Nightmarish

The integral in Equation 3.12 represents one of the most fundamental computational bottlenecks in machine learning:

**The Curse of Dimensionality**:

- Our latent space $\mathbf{z}$ is typically high-dimensional (64, 128, or even 512 dimensions)

- We need to integrate over *every possible value* in this high-dimensional space

- For a 128-dimensional latent space, this means integrating over $\mathbb{R}^{128}$—an unimaginably vast space

**No Closed-Form Solution**:

- The decoder $D_\theta(\mathbf{z})$ is a complex neural network with millions of parameters

- There's no analytical formula for this integral—we can't solve it with pencil and paper

- Even sophisticated symbolic mathematics software can't help us here

**Monte Carlo Estimation Fails**:

- We might try to approximate the integral by sampling many $\mathbf{z}$ values and averaging: $p(\mathbf{x}) \approx \frac{1}{K} \sum_{k=1}^{K} p(\mathbf{x}|\mathbf{z}_k)$ where $\mathbf{z}_k \sim p(\mathbf{z})$

- But most random $\mathbf{z}$ values produce data nothing like $\mathbf{x}$, so $p(\mathbf{x}|\mathbf{z}_k) \approx 0$ for almost all samples

- We'd need an astronomically large number of samples to get even one that contributes meaningfully to the sum

- This is like trying to hit a tiny target in a vast space by throwing darts randomly

**The Core Dilemma**: We can't directly optimize the quantity we care about most—the likelihood of our observed data.

## The Posterior Inference Alternative (Also Problematic)

Perhaps we can approach this differently using Bayes' rule. If we could compute the posterior distribution $p(\mathbf{z}|\mathbf{x})$, we could rearrange to find:

$$p(\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{z}|\mathbf{x})} \tag{3.13}$$

The posterior $p(\mathbf{z}|\mathbf{x})$ would tell us: "Given that we observed this face $\mathbf{x}$, what latent codes $\mathbf{z}$ most likely generated it?" This is exactly the kind of inference we want our model to perform!

**But the posterior is equally intractable**:

- By Bayes' rule: $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$

- The denominator $p(\mathbf{x})$ is the same intractable integral we started with!

- We're stuck in a circular dependency: to compute the posterior, we need the marginal likelihood, but to compute the marginal likelihood, we need to integrate over all possible posteriors

This creates a frustrating deadlock: the two things we most want to compute (the likelihood and the posterior) both depend on each other in intractable ways.

### 3.3.3 The Variational Solution: Breaking the Deadlock

The breakthrough insight of Variational Autoencoders is to sidestep this computational deadlock entirely by introducing an ingenious approximation strategy. Instead of trying to compute the true posterior $p(\mathbf{z}|\mathbf{x})$ directly, we'll **approximate it with a simpler distribution that we can actually work with**.

**Introducing the Variational Distribution**

The key idea is to introduce a *variational distribution* $q_\phi(\mathbf{z}|\mathbf{x})$ that serves as our best guess for what the true posterior $p(\mathbf{z}|\mathbf{x})$ looks like. Think of this as our "approximate inference engine":

- **Input**: An observed data point $\mathbf{x}$ (e.g., a face image)

- **Output**: A probability distribution over latent codes $\mathbf{z}$ that could have generated $\mathbf{x}$

- **Goal**: Make this distribution as close as possible to the true posterior $p(\mathbf{z}|\mathbf{x})$

**The Encoder Network: Learning to Infer**

We implement this variational distribution using a neural network called the **encoder**, parameterized by $\phi$:

$$q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_\phi(\mathbf{x}), \mathrm{diag}(\boldsymbol{\sigma}_\phi^2(\mathbf{x}))) \tag{3.14}$$

Here's how this works:

- **Encoder input**: The observed data $\mathbf{x}$ (e.g., pixel values of a face image)

- **Encoder outputs**: Two vectors:
  - $\boldsymbol{\mu}_\phi(\mathbf{x})$: The mean of our "best guess" latent code
  - $\boldsymbol{\sigma}_\phi(\mathbf{x})$: The uncertainty (standard deviation) in each dimension

- **Interpretation**: "Given this face, I think the latent code is around $\boldsymbol{\mu}_\phi(\mathbf{x})$, but I'm uncertain by amounts $\boldsymbol{\sigma}_\phi(\mathbf{x})$ in each dimension"

**Why a Gaussian approximation?**

- **Simplicity**: Gaussians are fully characterized by mean and variance, making them easy to parameterize with neural networks

- **Tractability**: We can easily sample from Gaussians and compute their densities

- **Differentiability**: Using the reparameterization trick, we can backpropagate through Gaussian sampling

- **Flexibility**: Despite their simplicity, Gaussians can approximate many distributions reasonably well in practice

**The Intuitive Picture: Encoder-Decoder Symmetry**

Now we have a beautiful symmetry in our model:

**Decoder (Generation)**: "Given an abstract concept $\mathbf{z}$, what concrete data $\mathbf{x}$ should I generate?"

$$\mathbf{z} \xrightarrow{D_\theta} \mathbf{x}$$

**Encoder (Inference)**: "Given concrete data $\mathbf{x}$, what abstract concept $\mathbf{z}$ most likely generated it?"

$$\mathbf{x} \xrightarrow{E_\phi} q(\mathbf{z}|\mathbf{x})$$

The encoder and decoder are like two neural networks having a conversation in an abstract language. The encoder translates from the concrete world (pixels, words, sounds) into abstract concepts (latent codes), while the decoder translates back from abstract concepts to the concrete world.

**The Approximation Trade-off**

By introducing the variational distribution $q_\phi(\mathbf{z}|\mathbf{x})$, we've made a crucial trade-off:

**What we gained**:

- A tractable way to perform approximate inference

- The ability to sample latent codes given observed data

- A path toward optimizing our model parameters

**What we sacrificed**:

- Exact inference—our approximation $q_\phi(\mathbf{z}|\mathbf{x})$ will never perfectly match $p(\mathbf{z}|\mathbf{x})$

- We now have more parameters to learn (both $\theta$ for the decoder and $\phi$ for the encoder)

- The quality of our model depends on how well our encoder can approximate the true posterior

The genius of the VAE framework is that it provides a principled way to balance this trade-off, ensuring that our approximation is as good as possible while still remaining computationally tractable. In the next section, we'll see how the Evidence Lower Bound (ELBO) gives us exactly the objective function we need to train both the encoder and decoder jointly.

## 3.4 The Evidence Lower Bound (ELBO)

Now that we have our variational distribution $q_\phi(\mathbf{z}|\mathbf{x})$ as an approximation to the intractable posterior $p(\mathbf{z}|\mathbf{x})$, we need to figure out how to actually train our model. We still can't compute the marginal likelihood $p(\mathbf{x})$ directly, but the variational framework gives us something almost as good: a **lower bound** that we can compute and optimize. This lower bound, known as the Evidence Lower Bound (ELBO), is the mathematical heart of VAEs.

### 3.4.1 Deriving the Lower Bound: A Step-by-Step Journey

The derivation of the ELBO is like a mathematical magic trick—we start with something we can't compute and, through a series of clever manipulations, end up with something we can. Let's walk through this transformation step by step.

**Starting Point: The Intractable Marginal Likelihood**

We begin with what we ultimately want to maximize—the log-likelihood of our observed data:

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}, \mathbf{z}) \, d\mathbf{z} \tag{3.15}$$

Remember, this integral is our nemesis—it requires summing over all possible latent codes in a high-dimensional space, which is computationally impossible.

**Step 1: The Variational Trick - Introducing Our Approximation**

Here comes the first clever move. We're going to multiply and divide by our variational distribution $q_\phi(\mathbf{z}|\mathbf{x})$:

$$\log p(\mathbf{x}) = \log \int \frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} p(\mathbf{x}, \mathbf{z}) \, d\mathbf{z} \tag{3.16}$$

**Why does this seemingly pointless step matter?** This is a profound question that gets to the heart of variational inference strategy. On the surface, we're just multiplying by 1 (since $\frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})} = 1$), so mathematically nothing has changed. But this "trivial" algebraic manipulation sets up one of the most important conceptual shifts in machine learning.

**The Strategic Intuition**: We're deliberately introducing our approximation $q_\phi(\mathbf{z}|\mathbf{x})$ into an equation where it doesn't belong yet, because we have a plan for how to use it. Think of this like a chess move that seems useless now but sets up a powerful combination later.

**What we're really doing**: We're preparing to change the *measure* of integration. The original integral $\int p(\mathbf{x}, \mathbf{z}) \, d\mathbf{z}$ integrates over *all possible* values of $\mathbf{z}$

with equal weight. But we suspect that most values of $\mathbf{z}$ are irrelevant for a given $\mathbf{x}$—only certain regions of latent space matter.

**The key insight**: By introducing $q_\phi(\mathbf{z}|\mathbf{x})$, we're preparing to *reweight* the integration. Instead of considering all $\mathbf{z}$ equally, we'll focus our computational effort on the $\mathbf{z}$ values that our encoder thinks are plausible. This transforms an impossible global search (integrate over all latent space) into a focused local search (integrate over encoder-suggested regions).

**The analogy**: Imagine you're looking for your keys in a dark house. The original integral is like searching every room with equal effort. The variational trick is like first asking "where did I last remember having them?" and then focusing your search on those likely areas. The $\frac{q_\phi(\mathbf{z}|\mathbf{x})}{q_\phi(\mathbf{z}|\mathbf{x})}$ multiplication is the mathematical setup that lets us implement this "focus your search" strategy.

### Step 2: Transforming to an Expectation

Now we rearrange the terms inside the integral:

$$\log p(\mathbf{x}) = \log \int q_\phi(\mathbf{z}|\mathbf{x}) \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \, d\mathbf{z} \tag{3.17}$$

The key insight is recognizing that this integral has the form $\int q(\mathbf{z}) \cdot f(\mathbf{z}) \, d\mathbf{z}$, which is exactly the definition of an expectation $\mathbb{E}_{q(\mathbf{z})}[f(\mathbf{z})]$.

**How do we identify this pattern?** Let's break down our integral to match the expectation definition. By definition, the expectation of a function $f(\mathbf{z})$ with respect to a probability distribution $q(\mathbf{z})$ is: $\mathbb{E}_{q(\mathbf{z})}[f(\mathbf{z})] = \int f(\mathbf{z}) \cdot q(\mathbf{z}) \, d\mathbf{z}$

In our integral $\int q_\phi(\mathbf{z}|\mathbf{x}) \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \, d\mathbf{z}$, we can identify:

- **The probability distribution**: $q(\mathbf{z}) = q_\phi(\mathbf{z}|\mathbf{x})$

- **The function being averaged**: $f(\mathbf{z}) = \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})}$

So our integral has exactly the right structure: $\int \underbrace{q_\phi(\mathbf{z}|\mathbf{x})}_{\text{probability distribution}} \cdot \underbrace{\frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})}}_{\text{function of } \mathbf{z}} \, d\mathbf{z}$

Therefore:

$$\log p(\mathbf{x}) = \log \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \tag{3.18}$$

**Why is this transformation crucial?** We've converted an intractable integral over all of latent space into an expectation over our variational distribution. The profound change is that instead of considering every possible $\mathbf{z}$ equally (as the original integral would), we're now **weighting each z by how likely our encoder thinks it is**, given the observed data $\mathbf{x}$. This dramatically reduces the effective space we need to explore—we can now approximate this expectation by sampling only from regions of latent space that our encoder believes are relevant for the given data point.

### Step 3: Jensen's Inequality - Creating the Lower Bound

Here's where the mathematical magic happens, but first, we need to understand **why** we're about to apply Jensen's inequality. The strategic thinking behind this step reveals the core genius of variational inference.

**The Problem We're Still Facing** After Step 2, we have:

$$\log p(\mathbf{x}) = \log \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

Unfortunately, this is **still intractable**! We've converted the integral to an expectation, but we haven't solved the fundamental optimization problem. Here's why this remains problematic:

**The expectation is still hard to work with**: To compute this, we'd need to sample many $\mathbf{z}$ values, compute $\frac{p(\mathbf{x},\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})}$ for each, average these ratios, and then take the logarithm of the average.

**The killer issue**: The expression log(average of stuff) is not easily differentiable in the way we need for backpropagation. We can't straightforwardly compute gradients with respect to $\theta$ and $\phi$ when we have the logarithm wrapped around an expectation.

**The Strategic Insight: If You Can't Optimize It, Bound It** This is where variational inference reveals its brilliance. The key realization is:

*"We can't compute or optimize* $\log p(\mathbf{x})$ *directly, but what if we could find something that's always* $\leq \log p(\mathbf{x})$, *and when we make this lower bound as large as possible, we're also pushing* $\log p(\mathbf{x})$ *up?"*

Our ultimate goal is to find parameters $\theta$ and $\phi$ that maximize $\log p(\mathbf{x})$. For gradient-based optimization, we need an objective that is:

1. **Computable**: We can evaluate it numerically

2. **Differentiable**: We can compute gradients with respect to our parameters

3. **Relevant**: Optimizing it actually helps with our real goal

The strategy is: instead of trying to maximize the intractable $\log p(\mathbf{x})$, we'll find a **lower bound** that satisfies all three criteria.

**Enter Jensen's Inequality** Jensen's inequality is the perfect tool for this strategy. For any concave function $f$ (like the logarithm):

$$f(\mathbb{E}[X]) \geq \mathbb{E}[f(X)] \tag{3.19}$$

Applied to our problem:

$$\log \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \tag{3.20}$$

**Why Jensen's Inequality Solves Our Problem** Jensen's inequality transforms our intractable expression into something we can actually optimize:

**Left side (still problematic)**:

$$\log \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

This involves taking the logarithm of an expectation, which creates the differentiability issues we discussed.

**Right side (tractable)**:

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

This is an expectation of logarithms, which we can:

- **Compute via sampling**: $\mathbb{E}[g(\mathbf{z})] \approx \frac{1}{K} \sum_{k=1}^{K} g(\mathbf{z}_k)$ where $\mathbf{z}_k \sim q_\phi(\mathbf{z}|\mathbf{x})$

- **Differentiate easily**: We can move gradients inside expectations and use automatic differentiation

- **Optimize with standard methods**: Standard gradient descent algorithms work perfectly

**Understanding Concavity and Why It Creates a Lower Bound**  **What does "concave" mean?** A function is concave if it curves downward—like an upside-down bowl. For the logarithm, this means $\frac{d^2}{dx^2} \log(x) = -\frac{1}{x^2} < 0$.

**Geometric intuition**: Imagine you have several points on a concave curve. If you:

1. Average the $x$-coordinates first, then find the $y$-value on the curve (left side of Jensen's inequality)

2. Find the $y$-value for each point first, then average those $y$-values (right side)

The first approach always gives a higher result because the concave curve bends downward, making the "middle" point higher than the average of the individual points.

**Concrete example**: With $x_1 = 1$, $x_2 = 4$, and $f(x) = \log(x)$:

- $\log \left( \frac{1+4}{2} \right) = \log(2.5) \approx 0.916$

- $\frac{\log(1) + \log(4)}{2} = \frac{0 + 1.386}{2} \approx 0.693$

Indeed, $0.916 > 0.693$, confirming the inequality.

**The Beautiful Irony**  The profound insight is this: **by "giving up" on computing the exact likelihood and settling for a lower bound, we actually gain the ability to optimize the likelihood effectively!**

Think of it like climbing a mountain (maximizing $\log p(\mathbf{x})$) that's covered in fog (intractable). Jensen's inequality gives us a clear path that's always below the mountain. As we walk uphill on this clear path (maximize the ELBO), we're guaranteed to be pushing ourselves up the actual mountain too.

This is the core genius of variational inference: sometimes the best way to solve a hard problem is to solve an easier related problem that guides you toward the solution of the hard one.

**The Final Result: Our Tractable Lower Bound**

Putting it all together, we have:

$$\log p(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \triangleq \mathcal{L}(\theta, \phi; \mathbf{x}) \tag{3.21}$$

This lower bound $\mathcal{L}(\theta, \phi; \mathbf{x})$ is the famous **Evidence Lower Bound (ELBO)**. **Why this bound is revolutionary**:

- **Computable**: Unlike the original marginal likelihood, we can actually evaluate this expression

- **Differentiable**: We can compute gradients with respect to both $\theta$ and $\phi$

- **Approximable**: The expectation can be estimated using Monte Carlo sampling

- **Optimizable**: Maximizing this bound pushes up the true likelihood we care about

## 3.4.2 Understanding the Jensen Gap

A natural question arises: how much are we losing by using this lower bound instead of the true likelihood? It turns out there's a beautiful answer.

The difference between the true log-likelihood and our lower bound is:

$$\log p(\mathbf{x}) - \mathcal{L}(\theta, \phi; \mathbf{x}) = D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}|\mathbf{x})) \tag{3.22}$$

This gap is exactly the KL divergence between our approximate posterior and the true posterior!

**The profound implications**:

- **Perfect approximation $\Rightarrow$ tight bound**: When $q_\phi(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$, the KL divergence is zero and our bound becomes exact

- **Better approximation $\Rightarrow$ tighter bound**: The closer our variational distribution is to the true posterior, the smaller the gap

- **Dual optimization**: By maximizing the ELBO, we simultaneously:

  - Increase the likelihood of our data (what we ultimately want)
  - Improve our posterior approximation (making our inference better)

This dual nature is what makes VAEs so elegant—we're not just learning to generate data, we're also learning to perform inference.

## 3.4.3 Decomposing the ELBO: The Two Fundamental Forces

The ELBO might look like abstract mathematical machinery, but it decomposes into two terms with beautiful, intuitive interpretations. Let's break it down step by step, carefully explaining each transformation:

**Starting with the ELBO**

We begin with our Evidence Lower Bound:

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \tag{3.23}$$

Our goal is to transform this expression into a form that reveals its intuitive meaning and computational structure.

**Step 1: Applying the Chain Rule of Probability**

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \tag{3.24}$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \tag{3.25}$$

**Why this step?** We use the fundamental chain rule of probability: $p(\mathbf{x}, \mathbf{z}) = p(\mathbf{x}|\mathbf{z})p(\mathbf{z})$. This factorization separates the joint probability into:

- $p(\mathbf{x}|\mathbf{z})$: The likelihood of data given latent code (our decoder)

- $p(\mathbf{z})$: The prior probability of the latent code

This separation is crucial because it aligns with our generative model structure and will lead to interpretable terms.

**Step 2: Using Properties of Logarithms**

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right] \tag{3.26}$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \tag{3.27}$$

**How do we get here?** We apply the logarithm properties:

- $\log \frac{A \cdot B}{C} = \log A + \log B - \log C$

- $\log \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} = \log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})$

**Why is this helpful?** Breaking the single logarithm into a sum allows us to separate the expectation in the next step, leading to terms we can interpret individually.

**Step 3: Linearity of Expectation**

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x})] \tag{3.28}$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log q_\phi(\mathbf{z}|\mathbf{x})] \tag{3.29}$$

**The principle used**: Expectation is linear, so $\mathbb{E}[A + B - C] = \mathbb{E}[A] + \mathbb{E}[B] - \mathbb{E}[C]$.

**What we now have**: Three separate expectation terms that we can analyze individually:

1. $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]$: Expected log-likelihood of data under decoder

2. $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{z})]$: Expected log-probability of latent codes under prior

3. $\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log q_\phi(\mathbf{z}|\mathbf{x})]$: Expected log-probability of latent codes under encoder

## Step 4: Recognizing the KL Divergence Pattern

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log q_\phi(\mathbf{z}|\mathbf{x})] \quad (3.30)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}\left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})}\right] \quad (3.31)$$

**How do we combine the last two terms?** We rearrange:

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log q_\phi(\mathbf{z}|\mathbf{x})] \quad (3.32)$$

$$= -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log q_\phi(\mathbf{z}|\mathbf{x})] + \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{z})] \quad (3.33)$$

$$= -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log q_\phi(\mathbf{z}|\mathbf{x}) - \log p(\mathbf{z})] \quad (3.34)$$

$$= -\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}\left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})}\right] \quad (3.35)$$

**Why this rearrangement?** We're setting up to recognize the KL divergence formula.

## Step 5: Identifying the KL Divergence

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}\left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})}\right] \quad (3.36)$$

$$= \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \quad (3.37)$$

**The key recognition**: The second term is exactly the definition of KL divergence:

$$D_{KL}(q\|p) = \mathbb{E}_q\left[\log \frac{q(\mathbf{z})}{p(\mathbf{z})}\right]$$

In our case: $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}\left[\log \frac{q_\phi(\mathbf{z}|\mathbf{x})}{p(\mathbf{z})}\right]$

## Final Interpretable Form

This gives us the final, interpretable decomposition:

$$\boxed{\mathcal{L}(\theta, \phi; \mathbf{x}) = \underbrace{\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]}_{\text{Reconstruction Term}} - \underbrace{D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}))}_{\text{Regularization Term}}} \quad (3.38)$$

**Why Each Step Matters**

**Step 1 (Chain rule)**: Separates our generative model into decoder likelihood and prior, aligning with our model architecture.

**Step 2 (Log properties)**: Converts a single complex logarithm into a sum of simpler terms.

**Step 3 (Linearity)**: Allows us to analyze each component separately, leading to interpretable terms.

**Step 4-5 (KL recognition)**: Reveals that part of our objective is measuring the difference between encoder outputs and our prior assumptions.

**The beautiful result**: What started as abstract mathematical machinery becomes two intuitive, competing forces:

- **Reconstruction**: "Preserve information so we can reconstruct the data"

- **Regularization**: "Keep latent codes well-behaved and close to our prior"

This decomposition shows that VAEs naturally balance compression (regularization) with reconstruction fidelity, leading to meaningful latent representations.

**Term 1: The Reconstruction Term**

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]$$

**What it measures**: How well we can reconstruct the original data from latent representations.

**The process it captures**:

1. The encoder looks at $\mathbf{x}$ and produces a distribution $q_\phi(\mathbf{z}|\mathbf{x})$ over latent codes

2. We sample $\mathbf{z}$ from this distribution

3. We pass $\mathbf{z}$ through the decoder to get a distribution $p(\mathbf{x}|\mathbf{z})$ over reconstructions

4. We measure how likely the original $\mathbf{x}$ is under this reconstruction distribution

**Intuitive interpretation**: "If I encode this data point into latent space and then decode it back, how likely am I to get something close to the original?"

**What maximizing this term encourages**:

- The encoder should produce latent codes that contain all information necessary to reconstruct $\mathbf{x}$

- The decoder should be able to generate high-quality reconstructions from these latent codes

- The autoencoder pipeline (encode then decode) should preserve information

**Term 2: The Regularization Term**

$$D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}))$$

**What it measures**: How different the encoder's output distribution is from our prior assumption about latent codes.
**The tension it captures**:

- $q_\phi(\mathbf{z}|\mathbf{x})$: What the encoder thinks the latent code should be given the data

- $p(\mathbf{z}) = \mathcal{N}(\mathbf{0}, \mathbf{I})$: Our simple prior assumption (standard normal)

- KL divergence: Measures how far apart these two distributions are

**Intuitive interpretation**: "How much does the encoder's guess about latent codes deviate from our simple assumptions?"
**What minimizing this term encourages**:

- The encoder outputs should stay close to the prior distribution

- Latent codes should be "well-behaved"—close to standard normal

- Prevents the encoder from using arbitrary or extreme latent representations

- Ensures we can generate new data by sampling from the prior $p(\mathbf{z})$

**The Beautiful Tension: Compression vs. Reconstruction** These two terms create a fundamental tension that drives the entire learning process:
**The Reconstruction Term says**: "Encode everything! Don't lose any information that might be needed to reconstruct the data perfectly."
**The Regularization Term says**: "Don't be too exotic! Keep your latent codes simple and close to the standard normal distribution."
This tension forces the model to learn a meaningful compromise:

- **Efficient compression**: The model must find a compact representation that captures the most important features

- **Structured latent space**: The representation must be organized according to our prior assumptions

- **Generative capability**: By keeping latent codes close to the prior, we ensure we can generate new data by sampling $\mathbf{z} \sim p(\mathbf{z})$

**The emergence of disentanglement**: In the best case, this tension encourages the model to discover meaningful, disentangled representations where different dimensions of $\mathbf{z}$ correspond to different semantic factors (like pose, lighting, identity in faces).

The ELBO doesn't just give us a tractable objective function—it encodes a principled approach to learning representations that balance reconstruction fidelity with structural constraints. This is why VAEs don't just compress data; they learn to organize it in meaningful ways.

## 3.5  Training VAEs: The Reparameterization Trick

### 3.5.1  The Gradient Problem

To train a VAE, we need to compute gradients of the ELBO with respect to both encoder parameters $\phi$ and decoder parameters $\theta$. The decoder gradients are straightforward, but the encoder gradients present a challenge.

The reconstruction term involves an expectation over the encoder distribution:

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] \tag{3.39}$$

**The gradient challenge:**

- We need $\frac{\partial}{\partial \phi}\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})]$

- But the expectation is over a distribution that depends on $\phi$

- We can't simply move the gradient inside the expectation

- Naive Monte Carlo would give biased gradient estimates

### 3.5.2  The Reparameterization Trick

The solution is to reparameterize the sampling process to separate the randomness from the parameters:

Instead of sampling $\mathbf{z} \sim q_\phi(\mathbf{z}|\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\phi(\mathbf{x}), \boldsymbol{\sigma}_\phi^2(\mathbf{x})\mathbf{I})$, we:

1. Sample $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (independent of $\phi$)

2. Compute $\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\sigma}_\phi(\mathbf{x}) \odot \boldsymbol{\epsilon}$

**Why this works:**

$$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] = \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})}[\log p(\mathbf{x}|\mathbf{z})] \tag{3.40}$$

where $\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\sigma}_\phi(\mathbf{x}) \odot \boldsymbol{\epsilon}$.

**The key advantages:**

- The expectation is now over $\boldsymbol{\epsilon}$, which doesn't depend on $\phi$

- We can move the gradient inside: $\frac{\partial}{\partial \phi}\mathbb{E}_{\boldsymbol{\epsilon}}[f(\mathbf{z})] = \mathbb{E}_{\boldsymbol{\epsilon}}\left[\frac{\partial f(\mathbf{z})}{\partial \phi}\right]$

- Monte Carlo estimation is now unbiased

- Gradients flow through the deterministic computation $\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\sigma}_\phi(\mathbf{x}) \odot \boldsymbol{\epsilon}$

### 3.5.3 The Complete Training Algorithm

For each training example $\mathbf{x}$:

1. **Encoder forward pass**: Compute $\boldsymbol{\mu}_\phi(\mathbf{x})$ and $\boldsymbol{\sigma}_\phi(\mathbf{x})$

2. **Sample**: Draw $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and compute $\mathbf{z} = \boldsymbol{\mu}_\phi(\mathbf{x}) + \boldsymbol{\sigma}_\phi(\mathbf{x}) \odot \boldsymbol{\epsilon}$

3. **Decoder forward pass**: Compute $D_\theta(\mathbf{z})$

4. **Compute ELBO**:

$$\mathcal{L} = \log p(\mathbf{x}|\mathbf{z}) - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \tag{3.41}$$
$$= \log \mathcal{N}(\mathbf{x}; D_\theta(\mathbf{z}), \sigma^2 \mathbf{I}) - D_{KL}(\mathcal{N}(\boldsymbol{\mu}_\phi(\mathbf{x}), \boldsymbol{\sigma}_\phi^2(\mathbf{x})\mathbf{I})\|\mathcal{N}(\mathbf{0}, \mathbf{I})) \tag{3.42}$$

5. **Backpropagation**: Compute gradients and update parameters

**Practical considerations:**

- The KL term has a closed form for Gaussian distributions

- The reconstruction term is typically implemented as squared error: $\|\mathbf{x} - D_\theta(\mathbf{z})\|^2$

- Multiple samples of $\boldsymbol{\epsilon}$ can be used to reduce variance, but often one sample suffices

## 3.6 VAEs vs GANs: Complementary Approaches

### 3.6.1 Fundamental Differences

| Aspect | VAEs | GANs |
|---|---|---|
| Training | Single objective (ELBO) | Minimax game |
| Theoretical foundation | Variational inference | Game theory |
| Stability | Generally stable | Often unstable |
| Mode collapse | Rare | Common problem |
| Sample quality | Often blurry | Can be very sharp |
| Latent space | Structured, interpretable | Less structured |
| Likelihood estimation | Approximate via ELBO | Not directly available |
| Convergence guarantees | Strong theoretical basis | Limited guarantees |

Table 3.1: Comparison of VAEs and GANs

**Why VAEs produce blurrier samples:**

- The Gaussian assumption in $p(\mathbf{x}|\mathbf{z})$ encourages "average" reconstructions

- The ELBO objective penalizes reconstruction error, leading to conservative estimates

- Missing details are "averaged out" rather than hallucinated

- **Trade-off**: Stability and theoretical grounding vs. sample sharpness

**Why GANs can produce sharper samples:**

- No explicit reconstruction requirement—only need to fool the discriminator

- Adversarial training can push toward realistic details

- Generator learns to exploit discriminator weaknesses

- **Trade-off**: Sample quality vs. training stability and mode coverage

## 3.6.2 Complementary Strengths

Rather than viewing VAEs and GANs as competing approaches, it's more productive to see them as complementary tools with different strengths:

**Use VAEs when:**

- You need stable, reliable training

- Interpretable latent representations are important

- You want to estimate likelihoods or perform inference

- You're working with limited computational resources

- You need to avoid mode collapse

**Use GANs when:**

- Sample quality is the primary concern

- You have extensive computational resources for hyperparameter tuning

- You're willing to invest in careful training procedures

- You don't need explicit likelihood estimation

**Hybrid approaches:** Recent research has explored combining the best of both worlds:

- **VAE-GANs**: Use VAE structure with adversarial training for the decoder

- **BiGANs/ALI**: Add an encoder to GANs for inference capabilities

- $\beta$**-VAEs**: Modify the VAE objective to improve disentanglement

## 3.7 Extensions and Variations

### 3.7.1 $\beta$-VAE: Controlling Disentanglement

The standard VAE objective can be modified to encourage more disentangled representations:

$$\mathcal{L}_\beta = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p(\mathbf{x}|\mathbf{z})] - \beta \cdot D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \tag{3.43}$$

**Understanding the $\beta$ parameter:**

- $\beta = 1$: Standard VAE

- $\beta > 1$: Stronger prior matching, encouraging independence between latent dimensions

- $\beta < 1$: Weaker regularization, potentially better reconstruction but less structured latent space

- **Trade-off**: Disentanglement vs. reconstruction quality

### 3.7.2 Conditional VAEs

We can extend VAEs to conditional generation by providing additional information:

$$p(\mathbf{x}|\mathbf{z}, \mathbf{c}) = \mathcal{N}(\mathbf{x}; D_\theta(\mathbf{z}, \mathbf{c}), \sigma^2 \mathbf{I}) \tag{3.44}$$
$$q_\phi(\mathbf{z}|\mathbf{x}, \mathbf{c}) = \mathcal{N}(\mathbf{z}; \boldsymbol{\mu}_\phi(\mathbf{x}, \mathbf{c}), \boldsymbol{\sigma}_\phi^2(\mathbf{x}, \mathbf{c})\mathbf{I}) \tag{3.45}$$

where $\mathbf{c}$ represents conditioning information (e.g., class labels, text descriptions).

### 3.7.3 Hierarchical VAEs

For more complex data, we can use multiple levels of latent variables:

$$\mathbf{z}_1 \sim p(\mathbf{z}_1) \tag{3.46}$$
$$\mathbf{z}_2 \sim p(\mathbf{z}_2|\mathbf{z}_1) \tag{3.47}$$
$$\mathbf{x} \sim p(\mathbf{x}|\mathbf{z}_1, \mathbf{z}_2) \tag{3.48}$$

This allows the model to capture hierarchical structure in the data.

## 3.8 Practical Considerations and Implementation

### 3.8.1 Architecture Design

**Encoder Architecture:**

- Typically uses convolutional layers (for images) or dense layers

- Outputs both mean $\boldsymbol{\mu}_\phi(\mathbf{x})$ and log-variance $\log \boldsymbol{\sigma}_\phi^2(\mathbf{x})$

- Log-variance output ensures positivity: $\boldsymbol{\sigma}_\phi^2(\mathbf{x}) = \exp(\log \boldsymbol{\sigma}_\phi^2(\mathbf{x}))$

**Decoder Architecture:**

- Mirror of encoder (for autoencoders) or task-specific design

- For images: transposed convolutions to upsample from latent space

- Output layer depends on data type (sigmoid for binary, linear for continuous)

## 3.8.2  Training Tips

**KL Annealing:**

- Start with small weight on KL term, gradually increase

- Helps avoid "posterior collapse" where encoder ignores input

- Schedule: $\beta_t = \min(1, t/T)$ where $T$ is annealing steps

**Architecture Balance:**

- Encoder and decoder should have similar capacity

- Too powerful encoder: posterior collapse

- Too powerful decoder: ignores latent variables

**Latent Dimensionality:**

- Higher dimensions: more expressive but harder to interpret

- Lower dimensions: better interpretability but potential information bottleneck

- Typical range: 10-1000 dimensions depending on data complexity

# 3.9  Applications and Impact

## 3.9.1  Scientific Applications

**Drug Discovery:**

- Encode molecular structures into latent space

- Generate new molecules by sampling and decoding

- Optimize properties by moving through latent space

**Astronomy:**

- Compress high-dimensional telescope data

- Generate synthetic observations for rare events

- Discover new types of astronomical objects

**Biology:**

- Model gene expression patterns

- Generate synthetic biological sequences

- Understand cellular state transitions

### 3.9.2 Creative Applications

**Art and Design:**

- Generate artistic styles and variations

- Interactive tools for creative exploration

- Style transfer and image manipulation

**Music:**

- Generate musical compositions

- Model musical style and structure

- Interactive music creation tools

### 3.9.3 Data Science Applications

**Dimensionality Reduction:**

- Nonlinear alternative to PCA

- Preserves local and global structure

- Enables visualization of high-dimensional data

**Anomaly Detection:**

- High reconstruction error indicates anomalies

- Density estimation via ELBO

- Applications in fraud detection, quality control

**Data Augmentation:**

- Generate additional training samples

- Particularly useful for small datasets

- More principled than simple transformations

## 3.10 Theoretical Insights and Connections

### 3.10.1 Information Theory Perspective

The VAE objective can be understood through an information-theoretic lens:

$$\mathcal{L} = \mathbb{E}[\log p(\mathbf{x}|\mathbf{z})] - D_{KL}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \tag{3.49}$$

**Information-theoretic interpretation:**

- **First term**: Negative conditional entropy $-H(\mathbf{x}|\mathbf{z})$—how much information $\mathbf{z}$ provides about $\mathbf{x}$

- **Second term**: Information cost of encoding—how much information we use to represent $\mathbf{z}$

- **Trade-off**: Minimize description length while preserving information about $\mathbf{x}$

### 3.10.2 Connection to Compression

VAEs implement a form of lossy compression:

- **Encoder**: Compresses data $\mathbf{x}$ into latent code $\mathbf{z}$

- **Decoder**: Decompresses latent code back to data space

- **Rate-distortion trade-off**: KL term controls "rate" (compression level), reconstruction term controls "distortion"

### 3.10.3 Bayesian Deep Learning

VAEs bridge deterministic neural networks and Bayesian methods:

- Explicit modeling of uncertainty through distributions

- Principled handling of latent variables

- Connection to variational inference in Bayesian statistics

## 3.11 Limitations and Challenges

### 3.11.1 Posterior Collapse

One of the most significant challenges in VAE training is posterior collapse, where the encoder learns to ignore the input and the decoder learns to generate data independently of the latent code.

**Why it happens:**

- If the decoder is very powerful, it might learn to generate good data without using $\mathbf{z}$

- The KL term then pushes $q(\mathbf{z}|\mathbf{x})$ toward the prior $p(\mathbf{z})$

- Result: $q(\mathbf{z}|\mathbf{x}) \approx p(\mathbf{z})$ for all $\mathbf{x}$, and latent variables become uninformative

**Solutions:**

- KL annealing: gradually increase the weight of the KL term

- Architecture constraints: limit decoder capacity

- Alternative objectives: $\beta$-VAE, InfoVAE, etc.

### 3.11.2  Blurry Samples

The Gaussian assumption in $p(\mathbf{x}|\mathbf{z})$ leads to blurry reconstructions:

- Gaussian distributions have "soft" probability mass

- Multiple possible reconstructions are averaged together

- Sharp details are smoothed out

**Potential solutions:**

- More flexible likelihood models (e.g., autoregressive, flow-based)

- Adversarial training components

- Different loss functions beyond MSE

### 3.11.3  Limited Expressiveness

The variational family $q_\phi(\mathbf{z}|\mathbf{x})$ may be too restrictive:

- Gaussian distributions can't capture complex multimodal posteriors

- Mean-field approximation assumes independence between latent dimensions

- True posterior might be much more complex

**Advances addressing this:**

- Normalizing flows for more flexible posteriors

- Importance weighted autoencoders (IWAE)

- Hierarchical variational models

## 3.12  Looking Forward: The Legacy of VAEs

Variational Autoencoders have had a profound impact on generative modeling and machine learning more broadly. Their key contributions include:

### 3.12.1 Methodological Contributions

- **Reparameterization trick**: Now used throughout probabilistic deep learning

- **Variational inference**: Brought Bayesian methods into the deep learning mainstream

- **Principled latent variable models**: Theoretical foundation for representation learning

### 3.12.2 Conceptual Impact

- **Probabilistic thinking**: Emphasized uncertainty and distributions over point estimates

- **Interpretable representations**: Showed how to learn meaningful latent spaces

- **Theory-practice bridge**: Demonstrated how theoretical insights can lead to practical algorithms

### 3.12.3 Influence on Subsequent Work

VAEs have influenced numerous subsequent developments:

- **Normalizing flows**: More flexible variational families

- **Variational inference**: Modern scalable Bayesian methods

- **Representation learning**: Self-supervised and disentangled representations

- **Diffusion models**: Share the idea of gradually transforming noise into data

The next chapter will explore diffusion models, which represent another major paradigm in generative modeling. While VAEs learn to map from simple to complex distributions in a single step, diffusion models learn to gradually transform noise into data through a series of small steps. This seemingly different approach actually shares deep connections with the variational inference principles we've explored in VAEs.

## 3.13 Exercises

1. **Mathematical Derivation**: Derive the closed-form expression for

$$D_{KL}(\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2 \mathbf{I}) \| \mathcal{N}(\mathbf{0}, \mathbf{I}))$$

and explain each term's contribution to the VAE objective.

2. **Jensen's Inequality**: Show that the "Jensen gap" between $\log p(\mathbf{x})$ and the ELBO is exactly $D_{KL}(q_\phi(\mathbf{z}|\mathbf{x}) \| p(\mathbf{z}|\mathbf{x}))$. What does this tell us about the quality of our variational approximation?

3. **Reparameterization**: Explain why naive Monte Carlo estimation of $\frac{\partial}{\partial \phi}\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[f(\mathbf{z})]$ gives biased gradients, and show how the reparameterization trick solves this problem.

4. **Posterior Collapse**: Design a simple experiment to demonstrate posterior collapse in VAEs. What architectural or training modifications could prevent it?

5. **Implementation**: Implement a simple VAE for MNIST digits. Experiment with different $\beta$ values in $\beta$-VAE and analyze the trade-off between reconstruction quality and latent space structure.

6. **Conditional VAE**: Extend your MNIST VAE to be conditional on digit class. How does this change the latent space organization? Can you perform class-conditional generation?

7. **Comparison**: Train both a VAE and a simple GAN on the same dataset. Compare sample quality, training stability, and latent space interpretability. Under what circumstances would you choose each approach?

# Chapter 4

# Diffusion Models: Gradual Generation Through Denoising

## 4.1 Introduction

In our journey through generative modeling, we have explored two fundamental approaches: GANs, which learn through adversarial competition, and VAEs, which use variational inference to approximate intractable posteriors. Both methods attempt to map directly from a simple latent distribution to the complex data distribution in a single step. This chapter introduces a radically different paradigm: *diffusion models*, which generate data through a gradual process of iterative refinement.

The key insight behind diffusion models is elegantly simple: instead of trying to generate complex data in one leap, we can learn to gradually transform noise into data through a sequence of small, manageable steps. This approach draws inspiration from non-equilibrium thermodynamics, where complex systems evolve through a series of incremental changes rather than sudden transformations.

Diffusion models, particularly Denoising Diffusion Probabilistic Models (DDPMs) introduced by Ho et al. in 2020, have achieved remarkable success in generating high-quality images, often surpassing GANs in both sample quality and training stability. Their success stems from a clever mathematical framework that combines the theoretical rigor of VAEs with a training procedure that avoids many of the pathologies that plague other generative models.

## 4.2 From VAEs to Hierarchical VAEs: Setting the Stage

As we established in Chapter 3, VAEs face several limitations: blurry samples due to Gaussian assumptions, potential posterior collapse, and the restrictive nature of simple variational families. These observations naturally lead us to ask: *what if we could make the generation process more flexible and gradual?*

## 4.2.1 Hierarchical VAEs: The Foundation

Recall that a standard VAE uses a single latent variable $\mathbf{z}$ to capture all the complexity of the data distribution. A natural extension is to use multiple levels of latent variables, creating a hierarchy:

**Notation clarification:** In hierarchical VAEs, we often switch from the standard $\mathbf{z}$ notation to using $\mathbf{x}$ with subscripts to emphasize the sequential, hierarchical nature of the variables. Here, $\mathbf{x}_0$ represents the observed data, while $\mathbf{x}_{1:T}$ are latent variables at different levels of abstraction.

$$\mathbf{x}_0 \ (\text{data}) \tag{4.1}$$

$$\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T \ (\text{latent variables}) \tag{4.2}$$

**The generative process becomes:**

$$p_\theta(\mathbf{x}_0, \mathbf{x}_{1:T}) = p_\theta(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \tag{4.3}$$

**Understanding this decomposition:**

- $p_\theta(\mathbf{x}_T)$: Prior distribution for the "deepest" latent variable

- $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$: Transition distributions from higher to lower levels

- Each step in the hierarchy refines the representation, gradually moving from abstract to concrete

## 4.2.2 Markovian Assumption

To make this hierarchy tractable, we impose a *Markovian* assumption: each latent variable depends only on the previous one in the hierarchy.

**Definition 4.1** (Markov Property). A stochastic process satisfies the Markov property if the future state depends only on the current state, not on the sequence of events that led to it: "What happens next depends only on the state of affairs now."

For our hierarchical VAE:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_{t+1}, \ldots, \mathbf{x}_T) = p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \tag{4.4}$$

$$q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_{t-2}, \ldots, \mathbf{x}_0) = q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1}) \tag{4.5}$$

The variational posterior also factorizes:

$$q_\phi(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q_\phi(\mathbf{x}_t|\mathbf{x}_{t-1}) \tag{4.6}$$

**Why the Markovian assumption helps:**

- Dramatically reduces the number of dependencies we need to model

- Makes both training and inference computationally tractable

- Allows us to sample efficiently by following the chain step by step

- Each transition can be relatively simple, even if the overall transformation is complex

## 4.3 The Diffusion Model Framework

Imagine you're an artist tasked with creating a masterpiece, but instead of starting with a blank canvas, you begin with pure static noise and gradually sculpt it into a beautiful image. This is the revolutionary idea behind diffusion models—they learn to reverse a process of gradual corruption, turning chaos into order.

But here's the clever twist: instead of trying to learn how to corrupt data (which would be complex and unpredictable), diffusion models fix the corruption process to be simple and mathematical, then focus all their learning power on the much harder problem of restoration. This design choice transforms an intractable optimization problem into an elegant and surprisingly effective approach to generative modeling.

### 4.3.1 The Key Innovation: A Tale of Two Processes

To understand why diffusion models work so well, we need to appreciate the fundamental design decision that sets them apart from their predecessors. Traditional generative models like VAEs attempt to learn two complex mappings simultaneously: how to encode data into a latent space (the forward process) and how to decode latent representations back into data (the reverse process). This creates a challenging optimization landscape where both networks must coordinate and adapt to each other.

Diffusion models take a radically different approach, inspired by a simple yet profound insight: **what if we only had to learn one of these processes?**

**The Forward Process: A Predetermined Journey into Chaos**

Instead of learning how to encode data, diffusion models fix the forward process to be a simple, predefined noise injection procedure. This isn't arbitrary chaos—it's a carefully designed mathematical process that gradually and systematically destroys structure while maintaining important properties that make the reverse process learnable.

The forward process is defined as a Markov chain where each step adds a small amount of Gaussian noise:

**The complete forward trajectory**:

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1}) \tag{4.7}$$

**Each individual transition**:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1-\beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \tag{4.8}$$

Think of this as a controlled demolition of your data. Just as demolition experts don't randomly blow up a building but carefully plan each explosion to ensure the structure collapses predictably, the forward process doesn't randomly destroy information—it follows a precise schedule that ensures the destruction is both gradual and reversible.

**Understanding the Demolition Schedule**

Let's examine the anatomy of each forward step to understand why this particular form of corruption is so effective:

**The Noise Schedule** $\beta_t$: These values control the rate of destruction at each timestep.

- $\beta_t \in (0, 1)$ with $\beta_1 \leq \beta_2 \leq \cdots \leq \beta_T$

- Start gentle: $\beta_1 \approx 0.0001$ preserves most structure initially

- End stronger: $\beta_T \approx 0.02$ completes the transformation to noise

- The gradual increase allows the reverse process to learn easier steps first

**Signal Preservation** $\sqrt{1 - \beta_t}$: This coefficient determines how much of the original signal survives each step.

- When $\beta_t$ is small, $\sqrt{1 - \beta_t} \approx 1$, so most signal is preserved

- As $\beta_t$ increases, less signal survives each timestep

- The square root scaling ensures proper variance control throughout the process

**Noise Injection** $\beta_t \mathbf{I}$: This adds fresh Gaussian noise at each step.

- The amount of noise increases as $\beta_t$ grows

- Independent noise in each dimension prevents correlated artifacts

- The specific variance $\beta_t$ (not $\sqrt{\beta_t}$) maintains mathematical elegance

**Why This Controlled Chaos Works**

The beauty of this fixed forward process lies in its mathematical properties:

**Predictable Convergence**: For sufficiently large $T$, any complex data eventually becomes pure noise:

$$q(\mathbf{x}_T|\mathbf{x}_0) \approx \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{4.9}$$

**Universal Coverage**: Every data point, regardless of complexity, ends up in the same noise distribution, ensuring the reverse process can generate any type of data.

**Gradual Transformation**: Each step involves only a small perturbation, making each individual reverse step learnable by a neural network.

**Mathematical Tractability**: The Gaussian nature of each step leads to closed-form expressions we'll derive shortly.

### 4.3.2 The Reparameterization Revelation: Making Randomness Differentiable

Now that we understand the forward process conceptually, we need to make it computationally tractable. Here, the reparameterization trick from VAEs comes to our rescue, but with a new twist that unlocks remarkable efficiency gains.

**From Sampling to Computation: The Reparameterization Deep Dive**

Our forward process involves sampling from a conditional Gaussian distribution:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I}) \tag{4.10}$$

This notation tells us that $\mathbf{x}_t$ follows a multivariate normal distribution with:

- **Mean**: $\boldsymbol{\mu} = \sqrt{1 - \beta_t}\mathbf{x}_{t-1}$

- **Covariance**: $\boldsymbol{\Sigma} = \beta_t\mathbf{I}$

**The Reparameterization Transformation** The reparameterization trick transforms this sampling operation into a deterministic computation. Recall from our VAE chapter that for any multivariate Gaussian $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, we can equivalently write:

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\boldsymbol{\epsilon} \tag{4.11}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\boldsymbol{\Sigma}^{1/2}$ is the matrix square root of the covariance.

**Applying This to Our Diffusion Step** In our case:

- $\boldsymbol{\mu} = \sqrt{1 - \beta_t}\mathbf{x}_{t-1}$

- $\boldsymbol{\Sigma} = \beta_t\mathbf{I}$

The matrix square root of $\boldsymbol{\Sigma} = \beta_t\mathbf{I}$ is particularly simple:

$$\boldsymbol{\Sigma}^{1/2} = (\beta_t\mathbf{I})^{1/2} = \sqrt{\beta_t}\mathbf{I} \tag{4.12}$$

**Why is this true?** Because:

- $(\sqrt{\beta_t}\mathbf{I})(\sqrt{\beta_t}\mathbf{I}) = \beta_t\mathbf{I}^2 = \beta_t\mathbf{I} = \boldsymbol{\Sigma}$

- The square root of a diagonal matrix is the square root of each diagonal element

**The Final Reparameterization** Substituting into our reparameterization formula:

$$\mathbf{x}_t = \boldsymbol{\mu} + \boldsymbol{\Sigma}^{1/2}\boldsymbol{\epsilon}_{t-1} \tag{4.13}$$

$$= \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\mathbf{I}\boldsymbol{\epsilon}_{t-1} \tag{4.14}$$

$$= \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\boldsymbol{\epsilon}_{t-1} \tag{4.15}$$

where $\boldsymbol{\epsilon}_{t-1} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

**Why This Transformation is Crucial** This reparameterization achieves several critical goals:

**1. Preserves Statistical Properties**:

- The mean is preserved: $\mathbb{E}[\mathbf{x}_t] = \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\mathbb{E}[\boldsymbol{\epsilon}_{t-1}] = \sqrt{1 - \beta_t}\mathbf{x}_{t-1}$

- The variance is preserved: $\text{Var}[\mathbf{x}_t] = \text{Var}[\sqrt{\beta_t}\boldsymbol{\epsilon}_{t-1}] = \beta_t\mathbf{I}$

- The distribution remains exactly $\mathcal{N}(\sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$

**2. Enables Gradient Flow**:

- The randomness is now isolated in $\boldsymbol{\epsilon}_{t-1}$, which doesn't depend on any model parameters

- The deterministic transformation $\mathbf{x}_{t-1} \mapsto \sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\boldsymbol{\epsilon}_{t-1}$ is differentiable

- Gradients can flow through this operation during backpropagation

**3. Computational Efficiency**:

- Instead of using specialized sampling algorithms, we use simple arithmetic operations

- The operation vectorizes efficiently on modern hardware (GPUs)

- Memory access patterns are predictable and cache-friendly

**The Variance Preservation Insight** A particularly elegant aspect of this reparameterization is how it maintains controlled variance. If we assume $\mathbf{x}_{t-1}$ has some variance $\sigma_{t-1}^2$, then:

$$\text{Var}[\mathbf{x}_t] = \text{Var}[\sqrt{1 - \beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\boldsymbol{\epsilon}_{t-1}] \tag{4.16}$$

$$= (1 - \beta_t)\text{Var}[\mathbf{x}_{t-1}] + \beta_t\text{Var}[\boldsymbol{\epsilon}_{t-1}] \tag{4.17}$$

$$= (1 - \beta_t)\sigma_{t-1}^2 + \beta_t \tag{4.18}$$

This shows exactly how the variance evolves: we keep a fraction $(1 - \beta_t)$ of the previous variance and add $\beta_t$ units of fresh noise variance.

This transformation gives us the same statistical properties as direct sampling but makes the entire forward process differentiable with respect to any parameters we might want to optimize—a crucial property for the learning algorithms that follow.

### The Architecture of Noise Addition

Let's trace through what happens as we apply this reparameterization step by step, watching how a crisp image gradually dissolves into static:

**The Original Image ($\mathbf{x}_0$)**: A perfect, structured data point—perhaps a clear photograph of a cat.

**First Corruption ($\mathbf{x}_1$)**:

$$\mathbf{x}_1 = \sqrt{1 - \beta_1}\mathbf{x}_0 + \sqrt{\beta_1}\boldsymbol{\epsilon}_0 \tag{4.19}$$

The cat photo now has a barely perceptible amount of grain, but every detail is still clearly visible.

**Second Corruption ($\mathbf{x}_2$):**

$$\mathbf{x}_2 = \sqrt{1 - \beta_2}\mathbf{x}_1 + \sqrt{\beta_2}\boldsymbol{\epsilon}_1 \tag{4.20}$$

$$= \sqrt{1 - \beta_2}(\sqrt{1 - \beta_1}\mathbf{x}_0 + \sqrt{\beta_1}\boldsymbol{\epsilon}_0) + \sqrt{\beta_2}\boldsymbol{\epsilon}_1 \tag{4.21}$$

$$= \sqrt{(1 - \beta_2)(1 - \beta_1)}\mathbf{x}_0 + \sqrt{(1 - \beta_2)\beta_1}\boldsymbol{\epsilon}_0 + \sqrt{\beta_2}\boldsymbol{\epsilon}_1 \tag{4.22}$$

The cat is still recognizable, but the image quality has noticeably degraded.

**The Pattern Emerges**: As we continue this process, we see that $\mathbf{x}_t$ becomes a weighted combination of:

- The original image $\mathbf{x}_0$ (with a coefficient that shrinks over time)

- Multiple independent noise terms that accumulate

### 4.3.3 Mathematical Magic: Forward Jumps Through Time

Here's where diffusion models reveal their most elegant property. Instead of laboriously applying hundreds or thousands of noise injection steps, we can mathematically "jump" directly to any point in the future. This isn't just a computational trick—it's a profound insight that makes efficient training possible.

**Setting Up the Mathematical Framework**

To make our derivations cleaner, let's introduce some notation:

- $\alpha_t = 1 - \beta_t$: The signal preservation factor at each step

- $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$: The cumulative signal preservation from start to timestep $t$

With this notation, our reparameterization becomes:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1} \tag{4.23}$$

**The Key Insight: Gaussian Arithmetic and Why It Matters**

The mathematical breakthrough comes from recognizing that when we have multiple independent Gaussian noise terms, they combine in a beautifully simple way:

$$\sqrt{a}\boldsymbol{\epsilon}_1 + \sqrt{b}\boldsymbol{\epsilon}_2 \sim \mathcal{N}(\mathbf{0}, (a + b)\mathbf{I})$$

where $\boldsymbol{\epsilon}_1, \boldsymbol{\epsilon}_2 \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ are independent.

But why is this property so revolutionary for diffusion models? Let's understand the profound implications.

**The Problem Without Noise Combination**   If we couldn't combine noise terms, our expression for $\mathbf{x}_t$ would look like:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{\alpha_t(1 - \bar{\alpha}_{t-1})}\boldsymbol{\epsilon}' + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1} + \sqrt{\alpha_t\alpha_{t-1}(1 - \bar{\alpha}_{t-2})}\boldsymbol{\epsilon}'' + \dots$$

This would create a computational nightmare:
**Storage Requirements**:

- We'd need to track and store ALL individual noise terms $\boldsymbol{\epsilon}_0, \boldsymbol{\epsilon}_1, \dots, \boldsymbol{\epsilon}_{t-1}$

- Memory usage would grow as $O(t)$ for each forward step

- For $T = 1000$ timesteps, we'd need 1000 times more memory

**Computational Complexity**:

- Computing $\mathbf{x}_t$ would require $O(t)$ operations and $O(t)$ memory

- Training would be computationally intractable: $O(T^2)$ operations for a full epoch

- Each timestep would depend on all previous timesteps, preventing parallelization

**Mathematical Intractability**:

- No closed-form expressions for forward distributions

- Impossible to compute exact KL divergences in loss functions

- No analytical understanding of convergence properties

**The Magic of Gaussian Closure**   The Gaussian closure property—that linear combinations of independent Gaussians are still Gaussian—transforms this chaos into elegance. ALL those accumulated noise terms collapse into a single equivalent term: All accumulated noise $= \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$

This is not an approximation, it is mathematically exact. The combined noise has precisely the same statistical properties as the sum of all individual noise terms.

**Why This Enables Everything**   The ability to collapse all accumulated noise into a single term creates a cascade of computational advantages that transform diffusion models from a theoretical curiosity into a practical powerhouse. Let's examine how this simple mathematical property revolutionizes every aspect of training and inference:
**1. Constant-Time Sampling Revolution**:

- **Without combination**: $O(T)$ sequential steps to reach timestep $T$

- **With combination**: $O(1)$ direct jump to any timestep

- Training transforms from $O(T^2)$ to $O(T)$ complexity

**2. Memory Efficiency Breakthrough**:

- **Without combination**: $O(T)$ memory per sample (store all noise terms)

- **With combination**: $O(1)$ memory per sample (just current noise)

- Enables training on much larger datasets and longer sequences

**3. Training Parallelization**: Instead of sequential dependence:

$$\text{Sequential: } \mathbf{x}_1 = f(\mathbf{x}_0, \boldsymbol{\epsilon}_0) \tag{4.24}$$
$$\mathbf{x}_2 = f(\mathbf{x}_1, \boldsymbol{\epsilon}_1) \quad \text{(depends on } \mathbf{x}_1\text{)} \tag{4.25}$$
$$\mathbf{x}_3 = f(\mathbf{x}_2, \boldsymbol{\epsilon}_2) \quad \text{(depends on } \mathbf{x}_2\text{)} \tag{4.26}$$

We get complete independence:

$$\text{Parallel: } \mathbf{x}_{50} = \sqrt{\bar{\alpha}_{50}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{50}}\boldsymbol{\epsilon} \tag{4.27}$$
$$\mathbf{x}_{200} = \sqrt{\bar{\alpha}_{200}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{200}}\boldsymbol{\epsilon}' \tag{4.28}$$
$$\mathbf{x}_{750} = \sqrt{\bar{\alpha}_{750}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{750}}\boldsymbol{\epsilon}'' \tag{4.29}$$

Each timestep can be computed completely independently!

**4. Mathematical Tractability**:

- Analytical forms for all forward process distributions

- Exact computation of KL divergences in training objectives

- Provable convergence to $\mathcal{N}(\mathbf{0}, \mathbf{I})$ as $t \to \infty$

- Clean theoretical analysis of model properties

**The Broader Significance**   This noise combination property isn't just a computational convenience—it's the fundamental insight that makes diffusion models practical:

**Training Efficiency**: Random timestep sampling becomes possible, leading to much more efficient and stable training.

**Scalability**: Models can handle thousands of timesteps without computational explosion.

**Theoretical Foundation**: Provides the mathematical rigor needed for principled model design and analysis.

**Implementation Simplicity**: Complex sequential processes become simple parallel operations.

Without this property, diffusion models would remain an interesting theoretical curiosity. With it, they become one of the most powerful and practical generative modeling approaches ever developed.

This means that all the noise we've accumulated over multiple timesteps can be collapsed into a single equivalent noise term—a mathematical miracle that enables efficient training!

**The Forward Jump Derivation**

Let's prove this remarkable result using mathematical induction:

**Base Case** ($t = 1$):

$$\mathbf{x}_1 = \sqrt{\alpha_1}\mathbf{x}_0 + \sqrt{1 - \alpha_1}\boldsymbol{\epsilon}_0 = \sqrt{\bar{\alpha}_1}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_1}\boldsymbol{\epsilon}_0 \tag{4.30}$$

This works because $\bar{\alpha}_1 = \alpha_1$.

**Inductive Hypothesis**: Assume that for timestep $t - 1$:

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\boldsymbol{\epsilon}' \tag{4.31}$$

**Inductive Step**: Now apply one more forward step:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1} \tag{4.32}$$

$$= \sqrt{\alpha_t}(\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\boldsymbol{\epsilon}') + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1} \tag{4.33}$$

$$= \sqrt{\alpha_t\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{\alpha_t(1 - \bar{\alpha}_{t-1})}\boldsymbol{\epsilon}' + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1} \tag{4.34}$$

**The Magic Happens**: The two noise terms combine:

$$\sqrt{\alpha_t(1 - \bar{\alpha}_{t-1})}\boldsymbol{\epsilon}' + \sqrt{1 - \alpha_t}\boldsymbol{\epsilon}_{t-1} \sim \mathcal{N}(\mathbf{0}, [\alpha_t(1 - \bar{\alpha}_{t-1}) + (1 - \alpha_t)]\mathbf{I}) \tag{4.35}$$

**Simplifying the Combined Variance**:

$$\alpha_t(1 - \bar{\alpha}_{t-1}) + (1 - \alpha_t) = \alpha_t - \alpha_t\bar{\alpha}_{t-1} + 1 - \alpha_t \tag{4.36}$$

$$= 1 - \alpha_t\bar{\alpha}_{t-1} \tag{4.37}$$

$$= 1 - \bar{\alpha}_t \tag{4.38}$$

**The Final Result**: We arrive at the elegant forward jump formula:

$$\boxed{\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}} \tag{4.39}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\bar{\alpha}_t = \prod_{s=1}^{t} \alpha_s$.

## 4.3.4   The Profound Implications

This forward jump formula is more than just a mathematical convenience—it fundamentally changes how we can train and understand diffusion models.

**Computational Revolution**

**Training Efficiency**: Instead of running a chain of $T$ forward steps (which might be 1000 steps), we can jump directly to any timestep in constant time. During training, we can:

- Sample a random timestep $t$ from $\{1, 2, \ldots, T\}$

- Use the forward jump to instantly create $\mathbf{x}_t$ from $\mathbf{x}_0$

- Train our denoising network on this $(t, \mathbf{x}_t)$ pair

- Repeat with a different random timestep

**Memory Efficiency**: We never need to store intermediate states $\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_{t-1}$—we can generate $\mathbf{x}_t$ directly from $\mathbf{x}_0$.

**Mathematical Insight**

The forward jump formula reveals the signal-to-noise ratio at any timestep:

$$\text{SNR}_t = \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t}$$

This gives us a precise understanding of how much original signal remains versus how much noise has been added. We can design our noise schedules by choosing the $\bar{\alpha}_t$ values to achieve desired signal-to-noise ratios throughout the process.

**The Complete Picture**

Looking at our forward jump formula, we can now see the complete story:

- At $t = 0$: $\bar{\alpha}_0 = 1$, so $\mathbf{x}_0$ is pure signal

- As $t$ increases: $\bar{\alpha}_t$ decreases, signal fades and noise grows

- At $t = T$: $\bar{\alpha}_T \approx 0$, so $\mathbf{x}_T \approx \boldsymbol{\epsilon}$ is pure noise

This transformation from perfect structure to pure chaos is completely deterministic and mathematically tractable. The reverse process—which we'll explore next—must learn to walk backwards along this path, gradually recovering structure from noise.

**The Artist's Perspective**

Returning to our artist metaphor: we now understand exactly how the masterpiece was destroyed. We know that at timestep $t$, exactly $\sqrt{\bar{\alpha}_t}$ of the original painting remains, mixed with $\sqrt{1 - \bar{\alpha}_t}$ amounts of random static. The artist's job is to learn how to remove this static step by step, revealing the hidden masterpiece beneath. This is the setup for one of the most remarkable learning problems in machine learning: given noisy images at various corruption levels, learn to predict and remove the noise. The forward jump formula ensures we can create unlimited training data for this task, sampling any corruption level instantly from any starting image.

The stage is now set for the reverse process—the neural network that must learn to be this digital restoration artist, turning chaos back into order one step at a time.

## 4.4 The ELBO for Diffusion Models: Learning to Reverse Chaos

Now that we understand how diffusion models systematically destroy structure through their fixed forward process, we turn to the harder question: how do we learn to reverse this destruction? Just as with VAEs, we'll use the Evidence Lower Bound (ELBO) framework, but adapted to the sequential nature of diffusion models.

## 4.4.1 Setting Up the Variational Problem

Our goal is to learn a reverse process that can undo the forward corruption, transforming pure noise back into structured data. We want to maximize the likelihood of our training data under our generative model: $\log p_\theta(\mathbf{x}_0)$

However, this involves marginalizing over all possible noise trajectories that could have led to our observed data—a computationally intractable integral over the high-dimensional space of all possible sequences $\mathbf{x}_{1:T}$.

Following the variational inference approach we learned from VAEs, we'll derive a lower bound that we can actually optimize. The fundamental insight remains the same: if we can't compute the exact likelihood, we'll find a tractable lower bound and maximize that instead.

## 4.4.2 Deriving the Diffusion ELBO: Step by Step

Let's carefully work through the derivation of the diffusion ELBO, seeing how the sequential structure leads to interpretable terms.

### Starting with the Variational Lower Bound

Using the same framework from VAEs, the log-likelihood lower bound becomes:

$$\log p_\theta(\mathbf{x}_0) \geq \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\right] \triangleq \mathcal{L} \tag{4.40}$$

Here's what each component represents:

- $p_\theta(\mathbf{x}_{0:T})$: Our learned generative model for the entire sequence

- $q(\mathbf{x}_{1:T}|\mathbf{x}_0)$: The fixed forward process that corrupts data

- The expectation is taken over all possible forward trajectories from $\mathbf{x}_0$

### Exploiting the Markovian Structure

The magic happens when we exploit the Markovian structure of both processes. Our generative model factorizes as:

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T)\prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \tag{4.41}$$

And our forward process factorizes as:

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1}) \tag{4.42}$$

### Telescoping the Logarithm

Substituting these factorizations into our ELBO:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log \frac{p(\mathbf{x}_T)\prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1})}\right] \tag{4.43}$$

$$= \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p(\mathbf{x}_T) + \sum_{t=1}^{T} \log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \sum_{t=1}^{T} \log q(\mathbf{x}_t|\mathbf{x}_{t-1})\right] \tag{4.44}$$

**Rearranging Terms Strategically**

Here's where we need to be clever about regrouping terms. We'll use a technique called "telescoping" to create a form that reveals the structure. Let's work through this step by step, understanding the strategic thinking behind each move.

**The Strategic Goal**   Before diving into algebra, let's understand what we're trying to achieve. We want to rearrange our terms so that we get:

- A **reconstruction term**: How well we recover $\mathbf{x}_0$ from $\mathbf{x}_1$

- A **prior matching term**: How well our endpoint matches the noise distribution

- **Denoising terms**: How well our learned reverse steps match the optimal reverse steps

The key insight is that we need to pair forward and reverse steps so we can compare what we *learn* to what we *should* learn.

**Starting from Our Factorized Form**   We have:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p(\mathbf{x}_T) + \sum_{t=1}^{T}\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \sum_{t=1}^{T}\log q(\mathbf{x}_t|\mathbf{x}_{t-1})\right] \qquad (4.45)$$

**Step 1: Separate the Prior Term**   This separation is strategic—we're identifying terms that will play special roles in our final objective and treating them differently from the "bulk" denoising terms.

**The Prior Term:** $\log p(\mathbf{x}_T)$   The term $\log p(\mathbf{x}_T)$ doesn't involve any learned parameters—it's just our assumption that the endpoint should be $\mathcal{N}(\mathbf{0}, \mathbf{I})$. This term will eventually be paired with something from the forward process to create a "prior matching" constraint.

**What makes it special**: Unlike the other terms in our sum, $\log p(\mathbf{x}_T)$ stands alone—it's not part of a conditional chain. It represents our "boundary condition" at the end of the diffusion process.

**The Reconstruction Term:** $\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$   This term is fundamentally different from all other reverse steps because:

- **It's the only term involving observed data**: $\mathbf{x}_0$ is the actual data we observe, while $\mathbf{x}_2, \mathbf{x}_3, \ldots, \mathbf{x}_T$ are latent variables

- **It represents the "output" of generation**: When generating new data, this is the step that produces the final result

- **It's closest to traditional reconstruction losses**: Like in VAEs, this measures how well we recover clean data from a noisy version

- **It will be treated differently in the loss**: This term won't become a KL divergence like the others—it stays as a likelihood term

**How the Separation Works Algebraically** Starting from:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p(\mathbf{x}_T) + \sum_{t=1}^{T}\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \sum_{t=1}^{T}\log q(\mathbf{x}_t|\mathbf{x}_{t-1})\right] \quad (4.46)$$

**Step 1a: Extract the prior term** The term $\log p(\mathbf{x}_T)$ is already separate, so we just pull it out:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p(\mathbf{x}_T)\right] + \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\sum_{t=1}^{T}\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \sum_{t=1}^{T}\log q(\mathbf{x}_t|\mathbf{x}_{t-1})\right]$$
$$(4.47)$$

**Step 1b: Extract the reconstruction term** From the reverse sum $\sum_{t=1}^{T}\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, let's identify what each term represents:

- $t = 1$: $\log p_\theta(\mathbf{x}_0|\mathbf{x}_1) \leftarrow$ This is our reconstruction term

- $t = 2$: $\log p_\theta(\mathbf{x}_1|\mathbf{x}_2) \leftarrow$ Denoising step

- $t = 3$: $\log p_\theta(\mathbf{x}_2|\mathbf{x}_3) \leftarrow$ Denoising step

- $\vdots$

- $t = T$: $\log p_\theta(\mathbf{x}_{T-1}|\mathbf{x}_T) \leftarrow$ Denoising step

So we can write:

$$\sum_{t=1}^{T}\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \log p_\theta(\mathbf{x}_0|\mathbf{x}_1) + \sum_{t=2}^{T}\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \quad (4.48)$$

Substituting this back:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p(\mathbf{x}_T)\right] \quad (4.49)$$
$$+ \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)\right] \quad (4.50)$$
$$+ \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\sum_{t=2}^{T}\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \sum_{t=1}^{T}\log q(\mathbf{x}_t|\mathbf{x}_{t-1})\right] \quad (4.51)$$

**The Strategic Purpose** By separating these terms, we're setting up a structure where:

1. **Prior term**: Will ensure our process ends at the right distribution

2. **Reconstruction term**: Will ensure we can generate good data from the first denoising step

3. **Remaining terms**: Will ensure all intermediate denoising steps are learned correctly

**The key insight**: Not all terms in our ELBO should be treated the same way. The prior and reconstruction terms have special roles that are different from the intermediate denoising steps, so we separate them early to handle them appropriately in the final loss function.

This separation also makes the telescoping in the next steps possible—by isolating the boundary terms (prior and reconstruction), we can focus on creating proper pairings for the intermediate terms that form the bulk of the learning objective.

**Step 2: Reindex the Forward Process Sum** We want to pair each reverse step $\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ with a corresponding forward step. But there's a mismatch:

- Remaining reverse steps: $t = 2, 3, \ldots, T$ (going from $\mathbf{x}_t$ to $\mathbf{x}_{t-1}$)

- Forward steps: $t = 1, 2, \ldots, T$ (going from $\mathbf{x}_{t-1}$ to $\mathbf{x}_t$)

To create proper pairs for comparison, we need to align the indices.

**The splitting strategy**: Let's rewrite the forward process sum by separating the last term, which will be paired with the prior:

Let's rewrite the forward process sum by changing the index. We'll split off the last term:

$$\sum_{t=1}^{T} \log q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \log q(\mathbf{x}_T|\mathbf{x}_{T-1}) + \sum_{t=1}^{T-1} \log q(\mathbf{x}_t|\mathbf{x}_{t-1}) \tag{4.52}$$

**Why separate** $\log q(\mathbf{x}_T|\mathbf{x}_{T-1})$**?** This term represents the final forward step—the transition to the endpoint. We want to pair this with $\log p(\mathbf{x}_T)$ to create a term about how well our endpoint matches the prior.

**The reindexing step**: Now reindex the remaining sum by substituting $s = t + 1$. When $t$ goes from 1 to $T - 1$, $s$ goes from 2 to $T$:

$$\sum_{t=1}^{T-1} \log q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \sum_{s=2}^{T} \log q(\mathbf{x}_{s-1}|\mathbf{x}_{s-2}) \tag{4.53}$$

Renaming the dummy variable $s$ back to $t$:

$$\sum_{t=1}^{T} \log q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \log q(\mathbf{x}_T|\mathbf{x}_{T-1}) + \sum_{t=2}^{T} \log q(\mathbf{x}_{t-1}|\mathbf{x}_{t-2}) \tag{4.54}$$

**What we've achieved**: Now both our reverse sum and the bulk of our forward sum run from $t = 2$ to $T$, allowing for proper pairing.

**Step 3: Substitute Back and Rearrange** **The substitution**: Substituting our reindexed forward sum back into the expression:

Substituting this back into our expression:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p(\mathbf{x}_T)\right] \tag{4.55}$$

$$+ \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)\right] \tag{4.56}$$

$$+ \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\sum_{t=2}^{T} \log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \log q(\mathbf{x}_T|\mathbf{x}_{T-1}) - \sum_{t=2}^{T} \log q(\mathbf{x}_{t-1}|\mathbf{x}_{t-2})\right]$$
$$\tag{4.57}$$

**What's happening here**: We now have three distinct groups forming:

1. The prior term $\log p(\mathbf{x}_T)$ standing alone

2. The reconstruction term $\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$ standing alone

3. A collection of terms involving intermediate timesteps and the final forward step

**Step 4: Group Terms to Reveal Structure** The strategic pairing: Now we strategically move terms around to create meaningful pairs:

- Pair $\log p(\mathbf{x}_T)$ with $-\log q(\mathbf{x}_T|\mathbf{x}_{T-1})$ to compare endpoint distribution with prior

- Keep $\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$ as the reconstruction term

- Pair each $\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ with $-\log q(\mathbf{x}_{t-1}|\mathbf{x}_{t-2})$ to compare learned and forward steps

Now we strategically pair terms to create recognizable patterns:

$$\mathcal{L} = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p(\mathbf{x}_T) - \log q(\mathbf{x}_T|\mathbf{x}_{T-1})\right] \tag{4.58}$$

$$+ \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)\right] \tag{4.59}$$

$$+ \sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \log q(\mathbf{x}_{t-1}|\mathbf{x}_{t-2})\right] \tag{4.60}$$

**Why This Specific Rearrangement?** The magic of this rearrangement is that we've created three groups of terms that will each become meaningful quantities:

**First group: Prior Matching** $\log p(\mathbf{x}_T) - \log q(\mathbf{x}_T|\mathbf{x}_{T-1})$ compares:

- $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$: Our assumption about where the process should end

- $q(\mathbf{x}_T|\mathbf{x}_{T-1})$: Where the forward process actually goes from the second-to-last step

This will become a KL divergence ensuring our forward process reaches the right endpoint.

**Second group: Reconstruction** $\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)$ is already in perfect form—it measures how well we can reconstruct clean data from the first level of corruption.

**Third group: Denoising Consistency** $\log p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) - \log q(\mathbf{x}_{t-1}|\mathbf{x}_{t-2})$ compares:

- $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$: Our learned reverse step from $\mathbf{x}_t$ to $\mathbf{x}_{t-1}$

- $q(\mathbf{x}_{t-1}|\mathbf{x}_{t-2})$: The forward step from $\mathbf{x}_{t-2}$ to $\mathbf{x}_{t-1}$

**Note**: This pairing isn't quite right yet—we're comparing a reverse step with a forward step that doesn't directly correspond. The next mathematical step will use Bayes' rule to transform this into the proper comparison between our learned reverse step and the true optimal reverse step.

**The telescoping magic**: The telescoping technique ensures that when we sum over all timesteps, terms cancel in exactly the right way to leave us with these three interpretable boundary and comparison terms. We've transformed a complex sum over all timesteps into three distinct learning objectives, each with clear meaning and purpose.

**Converting to KL Divergences**

Using the relationship between expectations and KL divergences, and exploiting conditional independence properties of our Markov chains, we can rewrite this as:

$$\mathcal{L} = \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)}[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\mathcal{L}_0:\text{Reconstruction term}} \tag{4.61}$$

$$- \underbrace{D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0)\|p(\mathbf{x}_T))}_{\mathcal{L}_T:\text{Prior matching term}} \tag{4.62}$$

$$- \underbrace{\sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}\left[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{x}_0)\|p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))\right]}_{\mathcal{L}_{1:T-1}:\text{Denoising matching terms}} \tag{4.63}$$

This decomposition reveals three distinct types of terms, each with a clear role in the learning process.

## 4.4.3 Understanding the Three Forces

The diffusion ELBO reveals three fundamental forces that shape the learning process:

**Term 1: The Reconstruction Force ($\mathcal{L}_0$)**

$\mathcal{L}_0 = \mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)}[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]$

**What it measures**: How well we can reconstruct the original data $\mathbf{x}_0$ from slightly noisy version $\mathbf{x}_1$.

**The intuition**: This is like asking an art restorer: "Given this photo with just a tiny bit of grain, can you recover the original perfectly?" This term ensures our model can handle the final, most delicate step of the restoration process.

**Implementation**: Often simplified to a mean squared error: $\|\mathbf{x}_0 - \boldsymbol{\mu}_\theta(\mathbf{x}_1, 1)\|^2$, where $\boldsymbol{\mu}_\theta$ is our learned denoising function.

**Term 2: The Prior Matching Force ($\mathcal{L}_T$)**

$\mathcal{L}_T = D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0)\|p(\mathbf{x}_T))$

**What it measures**: How well our forward process endpoint matches our assumed prior $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$.

**The beautiful insight**: Since we designed the forward process to converge to $\mathcal{N}(\mathbf{0}, \mathbf{I})$ by construction, this term is approximately zero! We get this "for free" from our careful design of the noise schedule.

**Practical implication**: No parameters to optimize here—the forward process design handles this automatically. This is one of the elegant aspects of diffusion models.

**Term 3: The Denoising Matching Force ($\mathcal{L}_{1:T-1}$)**

$\mathcal{L}_{t-1} = \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}\left[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t,\mathbf{x}_0)\|p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))\right]$

**What it measures**: How well our learned reverse step $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ matches the true optimal reverse step $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$.

**The key insight**: This is the heart of diffusion learning. We're training our neural network to predict what the optimal denoising step would be if it had access to the original clean image $\mathbf{x}_0$.

**Why this formulation is brilliant**: The true reverse distribution $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ conditions on both the noisy observation AND the clean target. This gives us a well-defined learning target for each timestep.

### 4.4.4 The Tractable Reverse Distribution

The remarkable property that makes diffusion models trainable is that when we condition on both $\mathbf{x}_t$ and $\mathbf{x}_0$, the reverse step becomes analytically tractable.

**Deriving the True Reverse Step**

Using Bayes' rule, we can compute:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \tag{4.64}$$

**Why this works**: All three distributions on the right are Gaussian (from our forward process design), so their ratio is also Gaussian. We can compute this exactly using Gaussian arithmetic.

**The magic of conditioning**: While $q(\mathbf{x}_{t-1}|\mathbf{x}_t)$ alone would be intractable, adding the condition on $\mathbf{x}_0$ makes everything computable because we know the forward process exactly.

**The Resulting Gaussian Form**

After working through the Gaussian algebra, we get:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\sigma}_t^2 \mathbf{I}) \tag{4.65}$$

where:

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{x}_t \tag{4.66}$$

$$\tilde{\sigma}_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t \tag{4.67}$$

**Understanding the Optimal Mean**

The mean $\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0)$ has a beautiful interpretation:

**Weighted interpolation**: It's a weighted combination of the clean data $\mathbf{x}_0$ and the noisy observation $\mathbf{x}_t$.

**Adaptive weighting**: The weights depend on the noise level. As $t$ increases (more corruption), the formula relies more heavily on the clean target $\mathbf{x}_0$ and less on the noisy observation $\mathbf{x}_t$.

**Fixed variance**: The variance $\tilde{\sigma}_t^2$ is completely determined by the noise schedule—no learning required.

**Perfect denoising:** This formula tells us exactly what the optimal single denoising step would be if we had access to both the noisy image and the original.

### 4.4.5 The Reparameterization Breakthrough

There's one more crucial insight that makes training practical. Instead of directly predicting the denoised image, we can reparameterize the problem to predict the noise itself.

**From the Forward Jump to Noise Prediction**

Recall our forward jump formula:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_t \tag{4.68}$$

Solving for $\mathbf{x}_0$:

$$\mathbf{x}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}\left(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_t\right) \tag{4.69}$$

**Substituting into the Optimal Mean**

When we substitute this expression for $\mathbf{x}_0$ into our optimal mean formula and simplify, we get:

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \boldsymbol{\epsilon}_t) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\boldsymbol{\epsilon}_t\right) \tag{4.70}$$

**The Noise Prediction Insight**

This reparameterization reveals a profound insight: instead of predicting the denoised image $\mathbf{x}_{t-1}$ directly, we can predict the noise $\boldsymbol{\epsilon}_t$ that was added to create $\mathbf{x}_t$ from $\mathbf{x}_0$.

**Why this is brilliant**:

- **Scale invariance**: Noise has the same scale at all timesteps

- **Simpler learning target**: Predicting noise is often easier than predicting images

- **Better optimization**: Avoids scaling issues that can hurt training

- **Unified objective**: The same network architecture can handle all timesteps

Our learning problem now becomes: **train a neural network $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ to predict the noise that was added at timestep** $t$.

This completes the mathematical foundation of diffusion models. We've transformed the intractable problem of learning to reverse a complex stochastic process into the much more manageable problem of learning to predict noise—a task that neural networks excel at.

## 4.5 Training Objectives: Three Parameterizations

The denoising matching term can be parameterized in three equivalent ways, each offering different insights and practical advantages.

### 4.5.1   Mean Prediction

The most direct approach is to predict the mean $\mu_\theta(\mathbf{x}_t, t)$ of the reverse distribution:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \mu_\theta(\mathbf{x}_t, t), \tilde{\sigma}_t^2 \mathbf{I}) \tag{4.71}$$

The KL divergence becomes:

$$D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) = \frac{1}{2\tilde{\sigma}_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \tag{4.72}$$

**Training objective:**

$$\mathcal{L}_{t-1}^{\text{mean}} = \mathbb{E}_{\mathbf{x}_0, \epsilon_t} \left[ \frac{1}{2\tilde{\sigma}_t^2} \|\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0) - \mu_\theta(\mathbf{x}_t, t)\|^2 \right] \tag{4.73}$$

### 4.5.2   Data Prediction ($\mathbf{x}_0$ Prediction)

Instead of predicting the mean, we can directly predict the original data $\mathbf{x}_0$:

$$\hat{\mathbf{x}}_\theta(\mathbf{x}_t, t) \approx \mathbf{x}_0 \tag{4.74}$$

Using the relationship between $\tilde{\mu}_t$ and $\mathbf{x}_0$, the loss becomes:

$$\mathcal{L}_{t-1}^{\mathbf{x}_0} = \omega_t \mathbb{E}_{\mathbf{x}_0, \epsilon_t} \left[ \|\mathbf{x}_0 - \hat{\mathbf{x}}_\theta(\mathbf{x}_t, t)\|^2 \right] \tag{4.75}$$

where $\omega_t$ is a timestep-dependent weighting factor.

**Intuition:** At each timestep, we observe a noisy version $\mathbf{x}_t$ and try to predict what the original clean image $\mathbf{x}_0$ was.

**Advantages:**

- Direct interpretation as a denoising autoencoder

- Natural connection to image restoration tasks

- Can be easier to evaluate and debug

### 4.5.3   Noise Prediction ($\epsilon$ Prediction)

The most popular parameterization predicts the noise that was added:

$$\epsilon_\theta(\mathbf{x}_t, t) \approx \epsilon_t \tag{4.76}$$

The loss becomes:

$$\mathcal{L}_{t-1}^{\epsilon} = \omega_t^{\epsilon} \mathbb{E}_{\mathbf{x}_0, \epsilon_t} \left[ \|\epsilon_t - \epsilon_\theta(\mathbf{x}_t, t)\|^2 \right] \tag{4.77}$$

**Why noise prediction is preferred:**

- **Scale invariance**: Noise has consistent magnitude across timesteps

- **Training stability**: More stable gradients than data prediction

- **Connection to score matching**: Related to estimating the score function $\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$

- **Empirical performance**: Often achieves the best results in practice

**The simplified objective:** Ho et al. showed that the optimal weighting $\omega_t^{\boldsymbol{\epsilon}} = 1$ often works best:

$$\boxed{\mathcal{L}_{\text{simple}} = \mathbb{E}_{\mathbf{x}_0, \boldsymbol{\epsilon}_t, t}\left[\|\boldsymbol{\epsilon}_t - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2\right]} \tag{4.78}$$

This is remarkably similar to training a denoising autoencoder, but with a specific noise schedule and conditioning on timestep $t$.

# 4.6 Training and Sampling Procedures

## 4.6.1 Training Algorithm

The training procedure for diffusion models is elegantly simple:

---
**Algorithm 1** Diffusion Model Training
---
**repeat**
    Sample $\mathbf{x}_0 \sim q(\mathbf{x}_0)$      ▷ Get training data
    Sample $t \sim \text{Uniform}\{1, \dots, T\}$      ▷ Random timestep
    Sample $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$      ▷ Sample noise
    Compute $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$      ▷ Forward jump
    Compute loss: $\mathcal{L} = \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2$
    Update $\theta$ using gradient descent
**until** convergence

---

**Key observations:**

- No encoder network needed—forward process is fixed

- Each training step requires only one forward pass through the denoising network

- Training is stable—no adversarial dynamics or posterior collapse issues

- Can train on arbitrary timesteps independently

## 4.6.2 Sampling Algorithm

Generation requires running the reverse process sequentially:

---
**Algorithm 2** Diffusion Model Sampling
---
Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ Start from noise $t = T, T-1, \dots, 1$ Predict noise: $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ Compute: $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\right) + \sigma_t\mathbf{z}$ where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and $\sigma_t = \sqrt{\tilde{\sigma}_t^2}$ **return** $\mathbf{x}_0$

---

**Understanding the sampling step:**

- Start with pure noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

- At each step, predict and remove the noise to get a slightly less noisy version

- Add a small amount of noise $\sigma_t \mathbf{z}$ to maintain proper sampling distribution

- The final step ($t = 1 \to 0$) is often deterministic (no added noise)

# 4.7 Connections and Intuitions

## 4.7.1 Relationship to Other Models

**Connection to VAEs:**

- Both use variational inference and ELBO optimization

- Diffusion models are hierarchical VAEs with fixed encoder and specific structure

- The fixed forward process eliminates posterior collapse issues

**Connection to Score-Based Models:**

- Noise prediction is related to score function estimation: $\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t)$

- Both can be viewed as learning to denoise data at different noise levels

- Unified framework shows they're different parameterizations of the same idea

**Connection to Energy-Based Models:**

- The reverse process can be viewed as gradient descent in an energy landscape

- Each denoising step moves toward lower energy (higher probability) regions

## 4.7.2 Why Diffusion Models Work So Well

**1. Divide and Conquer:**

- Instead of learning one complex mapping, learn many simple denoising steps

- Each step only needs to remove a small amount of noise

- Much easier optimization landscape than direct generation

**2. Strong Inductive Biases:**

- Forward process provides natural curriculum learning

- Early timesteps learn global structure, later timesteps learn fine details

- Noise schedule can be tuned for optimal learning dynamics

**3. Stable Training:**

- No adversarial dynamics or mode collapse

- No posterior collapse since forward process is fixed

- Training objective is a simple regression problem

**4. High Sample Quality:**

- Iterative refinement allows for very detailed generation

- Can spend more computation at inference time for better samples

- Natural way to trade off speed vs quality

# 4.8 Practical Considerations and Extensions

## 4.8.1 Noise Schedules

The choice of $\beta_t$ schedule significantly impacts performance:
**Linear Schedule:** $\beta_t = \beta_1 + \frac{t-1}{T-1}(\beta_T - \beta_1)$

- Simple and widely used

- $\beta_1 = 10^{-4}$, $\beta_T = 0.02$ are common choices

**Cosine Schedule:** $\bar{\alpha}_t = \frac{f(t)}{f(0)}$ where $f(t) = \cos\left(\frac{t/T+s}{1+s} \cdot \frac{\pi}{2}\right)^2$

- Slower noise injection early, faster later

- Often improves sample quality

- Better balances easy and hard denoising steps

## 4.8.2 Architecture Considerations

**U-Net Architecture:**

- Most common choice for image generation

- Skip connections help preserve spatial information

- Multi-scale processing matches hierarchical generation

**Timestep Conditioning:**

- Timestep $t$ typically embedded using sinusoidal encoding

- Added to intermediate activations throughout the network

- Essential for the model to understand what noise level to expect

**Attention Mechanisms:**

- Self-attention helps with global coherence

- Cross-attention enables conditional generation

- Particularly important for high-resolution images

### 4.8.3 Acceleration Techniques

**DDIM (Denoising Diffusion Implicit Models):**

- Deterministic sampling procedure

- Can skip timesteps for faster generation

- Trades off between speed and sample quality

**Progressive Distillation:**

- Train smaller models to mimic multi-step diffusion

- Can reduce sampling steps dramatically

- Maintains sample quality while improving speed

### 4.8.4 Conditional Generation

Diffusion models excel at conditional generation by modifying the denoising network:

**Class-Conditional Generation:**

$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, c) \tag{4.79}$$

where $c$ is the class label.

**Text-to-Image Generation:**

$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{c}) \tag{4.80}$$

where $\mathbf{c}$ is a text embedding from models like CLIP.

**Classifier-Free Guidance:** A powerful technique for improving conditional generation quality:

$$\tilde{\boldsymbol{\epsilon}}_\theta(\mathbf{x}_t, t, \mathbf{c}) = (1 + w)\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \mathbf{c}) - w\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \emptyset) \tag{4.81}$$

**Understanding classifier-free guidance:**

- $w > 0$: Guidance weight that controls conditioning strength

- $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t, \emptyset)$: Unconditional noise prediction

- Higher $w$ leads to better conditioning but potentially lower diversity

- Extrapolates in the direction of the conditioning signal

## 4.9 Mathematical Insights

### 4.9.1 The Score Function Perspective

The noise prediction in diffusion models is intimately connected to score-based generative modeling. The score function of a distribution $q(\mathbf{x}_t)$ is defined as:

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t) \tag{4.82}$$

**Connection to noise prediction:** For the forward process, we can show that:

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) = -\frac{\boldsymbol{\epsilon}_t}{\sqrt{1 - \bar{\alpha}_t}} \tag{4.83}$$

Therefore, predicting the noise $\boldsymbol{\epsilon}_t$ is equivalent to estimating the score function up to a scaling factor. This connection provides theoretical justification for why noise prediction works so well and connects diffusion models to the broader framework of score-based generative modeling.

### 4.9.2 Variance Preservation vs. Variance Explosion

The forward process we've described maintains a specific relationship between signal and noise that's worth examining more carefully.

**Variance Preserving (VP) Process:** Our standard formulation:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t \mathbf{I}) \tag{4.84}$$

The variance of $\mathbf{x}_t$ remains approximately constant throughout the process.

**Variance Exploding (VE) Process:** An alternative formulation:

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \mathbf{x}_{t-1}, (\sigma_t^2 - \sigma_{t-1}^2)\mathbf{I}) \tag{4.85}$$

Here, the signal is preserved exactly, but noise variance grows over time.

**Continuous-Time Perspective:** Both can be viewed as discretizations of stochastic differential equations (SDEs), providing a unified framework for understanding different diffusion processes.

### 4.9.3 Information-Theoretic Analysis

From an information theory perspective, the diffusion process can be understood as a gradual compression and decompression of information:

**Forward Process:** Gradually destroys information by adding noise

- Mutual information $I(\mathbf{x}_0; \mathbf{x}_t)$ decreases with $t$

- At $t = T$: $I(\mathbf{x}_0; \mathbf{x}_T) \approx 0$ (pure noise)

- Information is lost in a controlled, reversible manner

**Reverse Process:** Reconstructs information by learning the denoising transitions

- Each step $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ adds back a small amount of information

- The network must learn to "hallucinate" plausible details consistent with the data distribution

- Success depends on the network's ability to capture statistical regularities in the training data

# 4.10 Comparing Generative Paradigms

## 4.10.1 Comprehensive Comparison Table

| Aspect | GANs | VAEs | Diffusion | Autoregressive |
|---|---|---|---|---|
| Sample Quality | Very High | Moderate | Very High | High |
| Training Stability | Low | High | Very High | High |
| Mode Coverage | Poor | Good | Excellent | Excellent |
| Likelihood Estimation | No | Approximate | Approximate | Exact |
| Sampling Speed | Fast | Fast | Slow | Slow |
| Controllability | Moderate | Good | Excellent | Good |
| Theoretical Foundation | Game Theory | Bayesian | Bayesian | Information Theory |
| Memory Requirements | Low | Low | Moderate | High |

Table 4.1: Comparison of major generative modeling paradigms

## 4.10.2 When to Use Each Approach

**Choose Diffusion Models when:**

- High sample quality is paramount

- Training stability is important

- You need excellent mode coverage

- Controllable generation is required

- You can afford longer sampling times

- You want theoretical guarantees about the learning process

**Choose GANs when:**

- You need very fast sampling

- You have expertise in GAN training tricks

- You're working on specific domains where GANs excel (faces, artwork)

- Real-time generation is required

**Choose VAEs when:**

- You need explicit likelihood estimation

- Interpretable latent representations are important

- You want guaranteed training stability

- You're working with limited computational resources

- You need to perform inference (encoding new data)

# 4.11 Recent Advances and Future Directions

## 4.11.1 Architectural Innovations

**Diffusion Transformers (DiTs):**

- Replace U-Net with transformer architecture

- Better scalability to high resolutions and large datasets

- Improved performance on diverse image generation tasks

**Latent Diffusion Models:**

- Apply diffusion process in latent space rather than pixel space

- Dramatically reduces computational requirements

- Enables high-resolution generation with reasonable resources

- Foundation of Stable Diffusion and similar models

## 4.11.2 Applications Beyond Images

**Audio Generation:**

- Speech synthesis with natural prosody

- Music generation with long-term structure

- Audio super-resolution and enhancement

**Video Generation:**

- Temporal consistency through 3D convolutions

- Frame interpolation and prediction

- Video editing and style transfer

**3D and Molecular Generation:**

- 3D shape generation using point clouds or voxels

- Drug discovery through molecular diffusion

- Protein structure prediction and design

**Scientific Computing:**

- Climate modeling and weather prediction

- Fluid dynamics simulation

- Materials science and crystal structure generation

### 4.11.3 Theoretical Developments

**Continuous-Time Formulations:**

- Stochastic differential equation (SDE) frameworks

- Ordinary differential equation (ODE) solvers for deterministic sampling

- Connection to neural ODEs and continuous normalizing flows

**Optimal Transport Perspectives:**

- Understanding diffusion as optimal transport between distributions

- Rectified flows for straight-line generation paths

- Flow matching for more efficient training

**Information-Theoretic Analysis:**

- Rate-distortion theory for understanding compression-generation trade-offs

- Mutual information bounds on generation quality

- Connection to minimum description length principles

## 4.12 Implementation Considerations

### 4.12.1 Training Best Practices

**Numerical Stability:**

- Use mixed precision training to handle small noise values

- Careful handling of $\bar{\alpha}_t$ when it approaches zero

- Gradient clipping to prevent exploding gradients

**Data Preprocessing:**

- Normalize data to $[-1, 1]$ range for optimal training

- Consider data augmentation strategies

- Handle varying image resolutions appropriately

**Hyperparameter Tuning:**

- Number of timesteps $T$: typically 1000 for training, fewer for sampling

- Learning rate scheduling: often benefits from warmup and cosine decay

- Batch size: larger batches generally improve stability

### 4.12.2   Evaluation Metrics

**Quantitative Metrics:**

- **FID (Fréchet Inception Distance):** Measures distribution similarity

- **IS (Inception Score):** Evaluates both quality and diversity

- **Precision and Recall:** Separate assessment of quality vs. coverage

- **CLIP Score:** For text-to-image alignment evaluation

**Qualitative Assessment:**

- Visual inspection of generated samples

- Interpolation quality in latent or noise space

- Conditional generation controllability

- Mode coverage through diverse prompt evaluation

## 4.13   Limitations and Open Challenges

### 4.13.1   Current Limitations

**Computational Cost:**

- Sampling requires many neural network evaluations

- Training can be expensive for high-resolution models

- Memory requirements scale with number of timesteps

**Limited Theoretical Understanding:**

- Optimal noise schedules are still largely empirical

- Relationship between network capacity and generation quality unclear

- Limited theoretical guidance for architecture design

**Mode Representation:**

- May struggle with very rare modes in the data distribution

- Difficulty generating precisely specified content

- Limited control over fine-grained attributes

### 4.13.2  Open Research Questions

**Efficiency and Speed:**

- How can we achieve single-step generation without quality loss?

- What are the fundamental trade-offs between speed and quality?

- Can we learn adaptive noise schedules?

**Theoretical Foundations:**

- What is the optimal number of diffusion steps?

- How does network expressivity relate to generation capability?

- Can we provide convergence guarantees for the training process?

**Applications and Generalizations:**

- How can diffusion models handle discrete data effectively?

- What is the best way to incorporate strong conditional information?

- How can we ensure generated content follows ethical guidelines?

# 4.14  Conclusion: The Diffusion Revolution

Diffusion models represent a paradigm shift in generative modeling, moving from direct mapping approaches (GANs, VAEs) to gradual refinement processes. Their success stems from several key insights:

## 4.14.1  Fundamental Contributions

**1. Decomposition of Complexity:**

- Break down the challenging generation task into many simple denoising steps

- Each step is easier to learn than the full mapping

- Natural curriculum from coarse to fine generation

**2. Stable Training Dynamics:**

- Eliminate adversarial training instabilities

- Avoid posterior collapse through fixed forward process

- Simple regression-based training objective

**3. Theoretical Rigor:**

- Principled derivation from variational inference

- Clear connection to score-based models and optimal transport

- Well-understood approximation guarantees

**4. Practical Excellence:**

- State-of-the-art sample quality across many domains

- Excellent mode coverage and diversity

- Natural framework for conditional generation

### 4.14.2 Impact on the Field

Diffusion models have not only achieved impressive practical results but have also influenced our understanding of generative modeling more broadly:

- **Iterative Refinement:** Shown the power of gradual improvement over direct generation

- **Fixed Stochastic Processes:** Demonstrated that fixing part of the model can improve rather than limit performance

- **Score-Based Methods:** Revitalized interest in score function estimation and energy-based models

- **Continuous Perspectives:** Bridged discrete and continuous views of generative modeling

### 4.14.3 Looking Forward

As we look to the future, diffusion models are likely to continue evolving in several directions:

- **Efficiency:** Faster sampling methods and more efficient architectures

- **Scale:** Larger models trained on more diverse datasets

- **Multimodality:** Integration with other modalities beyond images

- **Applications:** Expansion into scientific computing, drug discovery, and other domains

- **Theory:** Deeper understanding of fundamental limitations and capabilities

The journey from GANs through VAEs to diffusion models illustrates the rapid evolution of generative modeling. Each approach has contributed crucial insights: GANs showed us the power of implicit generation, VAEs provided principled probabilistic foundations, and diffusion models demonstrated that gradual processes can achieve remarkable results.

As the field continues to advance, we can expect even more sophisticated approaches that combine the best aspects of each paradigm while addressing their respective limitations. The future of generative modeling promises to be as exciting and unpredictable as the journey that brought us here.

## 4.15 Exercises

1. **Forward Process Analysis:**

   (a) Derive the forward jump formula $q(\mathbf{x}_t|\mathbf{x}_0)$ step by step for $t = 3$.

   (b) Show that if $\beta_t \in (0, 1)$ for all $t$, then $\bar{\alpha}_T \to 0$ as $T \to \infty$.

   (c) Implement the forward process and visualize how an image gradually becomes noise.

2. **Reverse Process Derivation:**

   (a) Derive the mean $\tilde{\mu}_t(\mathbf{x}_t, \mathbf{x}_0)$ of $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ using Bayes' rule.

   (b) Show the equivalence between the three parameterizations (mean, data, noise prediction).

   (c) Explain intuitively why noise prediction often works better than data prediction.

3. **ELBO Decomposition:**

   (a) Starting from the standard ELBO for hierarchical VAEs, derive the diffusion model ELBO with three terms.

   (b) Explain why the prior matching term $\mathcal{L}_T$ is approximately zero.

   (c) Implement the simplified loss function and train a small diffusion model on MNIST.

4. **Noise Schedule Analysis:**

   (a) Compare linear vs. cosine noise schedules by plotting $\beta_t$ and $\bar{\alpha}_t$.

   (b) Implement both schedules and compare their effect on training dynamics.

   (c) Design your own noise schedule and justify your design choices.

5. **Sampling Procedures:**

   (a) Implement both DDPM and DDIM sampling procedures.

   (b) Compare the trade-off between number of sampling steps and quality.

   (c) Experiment with different amounts of noise injection during sampling.

6. **Conditional Generation:**

   (a) Implement class-conditional diffusion for CIFAR-10.

   (b) Add classifier-free guidance and study its effect on sample quality and diversity.

   (c) Compare conditional diffusion with conditional GANs and VAEs.

7. **Theoretical Connections:**

   (a) Show the connection between noise prediction and score function estimation.

(b) Implement a simple score-based model and compare it with diffusion.

(c) Discuss the relationship between diffusion models and energy-based models.

# Chapter 5

# Continuous-Time Diffusion Models: From Discrete Steps to Continuous Flows

## 5.1   Introduction: Beyond Discrete Timesteps

In our exploration of diffusion models through Chapters 1-4, we have worked exclusively with discrete timesteps $t \in \{1, 2, \ldots, T\}$. While this discrete formulation provides an intuitive framework and enables practical implementation, it raises several fundamental questions: What happens as we increase the number of timesteps $T$? Is there a natural continuous-time limit? And if so, what mathematical tools do we need to understand and work with continuous diffusion processes?

This chapter addresses these questions by developing the continuous-time theory of diffusion models. We will see that the discrete DDPM framework naturally generalizes to *stochastic differential equations* (SDEs), providing a unified mathematical foundation that encompasses not only the models we've studied but also opens doors to new algorithmic possibilities and theoretical insights.

The continuous-time perspective offers several advantages:

- **Mathematical elegance**: SDEs provide a principled framework with well-established theory

- **Flexible sampling**: Continuous-time formulations enable adaptive timestep selection and novel sampling algorithms

- **Theoretical insights**: Continuous analysis reveals deep connections to optimal transport, score-based models, and physics

- **Algorithmic innovations**: The continuous perspective has led to breakthrough methods like probability flow ODEs and neural ODEs

## 5.2   From Discrete to Continuous

### 5.2.1   Limitations of Discrete Timesteps

The discrete DDPM formulation, while successful, has several inherent limitations:

### 1. Fixed Resolution:

- The choice of $T$ is somewhat arbitrary and affects both training and sampling

- Different problems might benefit from different temporal resolutions

- Fixed timesteps may be inefficient—some regions may need finer resolution while others could use coarser steps

### 2. Discrete Artifacts:

- Discrete jumps can introduce artifacts, especially with aggressive noise schedules

- The discrete nature doesn't reflect the underlying continuous processes we're trying to model

- Error accumulation over many discrete steps can degrade sample quality

### 3. Limited Theoretical Tools:

- Discrete analysis is often more complex than continuous analysis

- Many powerful mathematical tools (differential equations, calculus of variations) are naturally continuous

- Connections to physics and other sciences are clearer in continuous time

## 5.2.2 The Continuous Limit: Intuitive Development

To develop intuition for the continuous limit, consider what happens as we increase $T$ while keeping the total diffusion time fixed. Let's normalize time to the interval $[0, 1]$ and define:

$$\Delta t = \frac{1}{T}, \quad t_i = i \cdot \Delta t \tag{5.1}$$

The discrete forward process becomes:

$$\mathbf{x}_{t_{i+1}} = \sqrt{1 - \beta_{t_i}\Delta t}\mathbf{x}_{t_i} + \sqrt{\beta_{t_i}\Delta t}\boldsymbol{\epsilon}_i \tag{5.2}$$

**Understanding the scaling:**

- As $\Delta t \to 0$, we need $\beta_{t_i}\Delta t \to 0$ to maintain stability

- The noise term $\sqrt{\beta_{t_i}\Delta t}$ scales as $\sqrt{\Delta t}$

- This scaling is characteristic of Brownian motion and diffusion processes

**The limiting behavior:** As $T \to \infty$ and $\Delta t \to 0$, the discrete updates become:

$$\mathbf{x}_{t+\Delta t} - \mathbf{x}_t = (\sqrt{1 - \beta_t\Delta t} - 1)\mathbf{x}_t + \sqrt{\beta_t\Delta t}\boldsymbol{\epsilon} \tag{5.3}$$

$$\approx -\frac{1}{2}\beta_t\Delta t\mathbf{x}_t + \sqrt{\beta_t\Delta t}\boldsymbol{\epsilon} \tag{5.4}$$

Dividing by $\Delta t$ and taking the limit:

$$\frac{d\mathbf{x}}{dt} = -\frac{1}{2}\beta_t\mathbf{x} + \sqrt{\beta_t}\frac{d\mathbf{w}}{dt} \tag{5.5}$$

This is our first glimpse of the SDE formulation!

### 5.2.3 Stochastic Differential Equations

**Definition 5.1** (Stochastic Differential Equation)**.** A stochastic differential equation (SDE) has the general form:

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + \mathbf{g}(\mathbf{x}, t)d\mathbf{w} \tag{5.6}$$

where:

- $\mathbf{f}(\mathbf{x}, t)$ is the *drift coefficient* (deterministic part)

- $\mathbf{g}(\mathbf{x}, t)$ is the *diffusion coefficient* (stochastic part)

- $d\mathbf{w}$ represents infinitesimal Brownian motion

**Understanding the components:**
**Drift Term $\mathbf{f}(\mathbf{x}, t)dt$:**

- Represents the deterministic "flow" of the process

- In our context, often represents the systematic shrinking toward zero

- Analogous to a vector field that pushes particles in specific directions

**Diffusion Term $\mathbf{g}(\mathbf{x}, t)d\mathbf{w}$:**

- Represents random fluctuations

- $d\mathbf{w}$ is a Wiener process increment with $d\mathbf{w} \sim \mathcal{N}(\mathbf{0}, dt\mathbf{I})$

- The coefficient $\mathbf{g}(\mathbf{x}, t)$ can depend on both state and time

**Physical Interpretation:** Think of a particle suspended in a fluid:

- The drift term represents systematic forces (like gravity or electric fields)

- The diffusion term represents random collisions with fluid molecules

- The combination produces complex, realistic trajectories

## 5.3 Forward Diffusion SDE

### 5.3.1 SDE Formulation of the Forward Process

The continuous-time forward diffusion process is described by the SDE:

$$\boxed{d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w}} \tag{5.7}$$

For diffusion models, we typically use:

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\mathbf{w} \tag{5.8}$$

**Breaking down this formulation:**
**Drift Coefficient: $\mathbf{f}(\mathbf{x}, t) = -\frac{1}{2}\beta(t)\mathbf{x}$**

- **Linear in x**: The drift scales with the current state

- **Negative sign**: Systematically shrinks the signal toward zero

- **Time-dependent**: $\beta(t)$ controls the rate of shrinkage over time

- **Factor of $\frac{1}{2}$**: Ensures correct variance evolution (derived from matching discrete case)

**Diffusion Coefficient:** $g(t) = \sqrt{\beta(t)}$

- **State-independent**: Noise magnitude doesn't depend on current $\mathbf{x}$

- **Time-dependent**: $\beta(t)$ controls noise injection rate

- **Square root relationship**: Ensures variance grows at rate $\beta(t)$

## 5.3.2  The Ornstein-Uhlenbeck Process

The forward diffusion SDE is a special case of the *Ornstein-Uhlenbeck process*, a fundamental stochastic process with many applications in physics, finance, and biology.

**Definition 5.2** (Ornstein-Uhlenbeck Process)**.** The general Ornstein-Uhlenbeck process is:

$$d\mathbf{x} = -\theta(\mathbf{x} - \boldsymbol{\mu})dt + \sigma d\mathbf{w} \tag{5.9}$$

where $\theta > 0$ is the mean reversion rate, $\boldsymbol{\mu}$ is the long-term mean, and $\sigma$ is the diffusion coefficient.

**Our diffusion SDE as Ornstein-Uhlenbeck:** Setting $\theta = \frac{1}{2}\beta(t)$, $\boldsymbol{\mu} = \mathbf{0}$, and $\sigma = \sqrt{\beta(t)}$:

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\mathbf{w} \tag{5.10}$$

**Key properties of Ornstein-Uhlenbeck processes:**

- **Mean reversion**: The process is pulled toward the mean $\boldsymbol{\mu} = \mathbf{0}$

- **Stationary distribution**: For constant coefficients, converges to $\mathcal{N}(\boldsymbol{\mu}, \frac{\sigma^2}{2\theta}\mathbf{I})$

- **Gaussian transitions**: Transitions remain Gaussian (crucial for tractability)

- **Markovian**: Future evolution depends only on current state

### 5.3.3 Connection to the Discrete Case

To verify our continuous formulation, let's show how it connects to the discrete DDPM equations.

**Discretizing the SDE:** Using the Euler-Maruyama scheme with timestep $\Delta t$:

$$\mathbf{x}_{t+\Delta t} = \mathbf{x}_t + \mathbf{f}(\mathbf{x}_t, t)\Delta t + g(t)\Delta \mathbf{w} \tag{5.11}$$

$$= \mathbf{x}_t - \frac{1}{2}\beta(t)\mathbf{x}_t\Delta t + \sqrt{\beta(t)}\sqrt{\Delta t}\boldsymbol{\epsilon} \tag{5.12}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

**Rearranging:**

$$\mathbf{x}_{t+\Delta t} = \left(1 - \frac{1}{2}\beta(t)\Delta t\right)\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\epsilon} \tag{5.13}$$

**Matching with discrete DDPM:** For small $\Delta t$, we have $\sqrt{1 - \beta\Delta t} \approx 1 - \frac{1}{2}\beta\Delta t$, so:

$$\mathbf{x}_{t+\Delta t} = \sqrt{1 - \beta(t)\Delta t}\mathbf{x}_t + \sqrt{\beta(t)\Delta t}\boldsymbol{\epsilon} \tag{5.14}$$

This exactly matches the discrete DDPM transition with $\beta_t = \beta(t)\Delta t$!

### 5.3.4 Variance Preserving vs. Variance Exploding

The SDE formulation reveals a fundamental choice in how we design the diffusion process.

**Variance Preserving (VP) SDE:**

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\mathbf{w} \tag{5.15}$$

**Why "variance preserving"?**

- The drift term shrinks the signal: $\mathbb{E}[\mathbf{x}_t] \to \mathbf{0}$

- The diffusion term adds noise: $\text{Var}[\mathbf{x}_t]$ increases

- The balance is designed so that $\mathbb{E}[\|\mathbf{x}_t\|^2]$ remains approximately constant

- This prevents the process from either exploding or vanishing too quickly

**Variance Exploding (VE) SDE:**

$$d\mathbf{x} = \sqrt{\frac{d[\sigma^2(t)]}{dt}}d\mathbf{w} \tag{5.16}$$

**Understanding VE:**

- No drift term: $\mathbb{E}[\mathbf{x}_t] = \mathbb{E}[\mathbf{x}_0]$ (signal preserved exactly)

- Pure diffusion: $\text{Var}[\mathbf{x}_t] = \text{Var}[\mathbf{x}_0] + \sigma^2(t)$

- Variance grows without bound as $t \to \infty$

- Often easier to analyze but can be numerically unstable

### 5.3.5 Solving the Forward SDE

For the VP-SDE with time-dependent coefficients, we can find the exact solution using integrating factors.

**The solution:**

$$\mathbf{x}_t = e^{-\frac{1}{2}\int_0^t \beta(s)ds}\mathbf{x}_0 + \int_0^t e^{-\frac{1}{2}\int_s^t \beta(r)dr}\sqrt{\beta(s)}d\mathbf{w}_s \tag{5.17}$$

**Defining $\alpha(t)$:** Let $\alpha(t) = e^{-\frac{1}{2}\int_0^t \beta(s)ds}$. Then:

$$\mathbf{x}_t = \alpha(t)\mathbf{x}_0 + \sqrt{1 - \alpha^2(t)}\boldsymbol{\epsilon} \tag{5.18}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ is independent of $\mathbf{x}_0$.

**This gives us the transition distribution:**

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \alpha(t)\mathbf{x}_0, (1 - \alpha^2(t))\mathbf{I}) \tag{5.19}$$

**Connection to discrete case:** This exactly matches our discrete forward jump formula with $\bar{\alpha}_t = \alpha(t)$!

## 5.4 Probability Flow ODE: Deterministic Paths

### 5.4.1 From Stochastic to Deterministic

One of the most remarkable discoveries in the continuous-time theory is that every diffusion SDE has a corresponding *ordinary differential equation* (ODE) that produces the same marginal distributions at each time $t$.

**Theorem 5.1** (Probability Flow ODE). *For the forward diffusion SDE:*

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w} \tag{5.20}$$

*there exists a corresponding* probability flow ODE*:*

$$\frac{d\mathbf{x}}{dt} = \mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g^2(t)\nabla_{\mathbf{x}} \log p_t(\mathbf{x}) \tag{5.21}$$

*such that if $\mathbf{x}_0^{SDE} \sim p_0$ and $\mathbf{x}_0^{ODE} \sim p_0$, then $\mathbf{x}_t^{SDE} \sim p_t$ and $\mathbf{x}_t^{ODE} \sim p_t$ for all $t$.*

**Understanding this remarkable result:**

- **Same marginals**: Both processes have identical marginal distributions $p_t(\mathbf{x})$

- **Different paths**: Individual trajectories are completely different

- **Deterministic vs. stochastic**: ODE paths are deterministic, SDE paths are random

- **Score dependence**: The ODE requires the score function $\nabla_{\mathbf{x}} \log p_t(\mathbf{x})$

## 5.4.2 Deriving the Probability Flow ODE

For our specific VP-SDE:

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\mathbf{w} \tag{5.22}$$

The corresponding probability flow ODE is:

$$\boxed{\frac{d\mathbf{x}}{dt} = -\frac{1}{2}\beta(t)\mathbf{x} - \frac{1}{2}\beta(t)\nabla_{\mathbf{x}}\log p_t(\mathbf{x})} \tag{5.23}$$

**Breaking down the terms:**

- $-\frac{1}{2}\beta(t)\mathbf{x}$: Same drift as the SDE

- $-\frac{1}{2}\beta(t)\nabla_{\mathbf{x}}\log p_t(\mathbf{x})$: Additional "score-based" drift

- The score term compensates for the missing stochastic component

## 5.4.3 Connection to Score-Based Models

The probability flow ODE reveals a deep connection to score-based generative models. Recall from Chapter 4 that diffusion models can be viewed as learning the score function:

$$\nabla_{\mathbf{x}}\log q(\mathbf{x}_t|\mathbf{x}_0) = -\frac{\boldsymbol{\epsilon}_t}{\sqrt{1 - \bar{\alpha}_t}} \tag{5.24}$$

**Neural network approximation:** If we train $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ to predict the noise, we can approximate:

$$\nabla_{\mathbf{x}}\log p_t(\mathbf{x}) \approx -\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}, t)}{\sqrt{1 - \alpha^2(t)}} \tag{5.25}$$

**Practical probability flow ODE:**

$$\frac{d\mathbf{x}}{dt} = -\frac{1}{2}\beta(t)\mathbf{x} + \frac{1}{2}\beta(t)\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}, t)}{\sqrt{1 - \alpha^2(t)}} \tag{5.26}$$

## 5.4.4 Applications of Probability Flow ODEs

**1. Deterministic Sampling:**

- Solve the ODE backwards from $t = T$ to $t = 0$

- Same marginal distributions as stochastic sampling

- Reproducible results (same noise leads to same output)

- Often faster convergence with adaptive ODE solvers

**2. Latent Space Interpolation:**

- Map data to latent space using forward ODE

- Interpolate in latent space

- Map back to data space using reverse ODE

- Produces meaningful interpolations

**3. Likelihood Computation:** Using the change of variables formula:

$$\log p_0(\mathbf{x}_0) = \log p_T(\mathbf{x}_T) + \int_0^T \nabla \cdot \mathbf{v}(\mathbf{x}_t, t) dt \tag{5.27}$$

where $\mathbf{v}(\mathbf{x}, t)$ is the velocity field of the probability flow ODE.

**4. Neural ODEs Integration:**

- Natural connection to Neural ODE frameworks

- Can use sophisticated ODE solvers (Runge-Kutta, adaptive timesteps)

- Enables training with continuous normalizing flows

## 5.5 Continuous Diffusion

### 5.5.1 Invariant Measures and Stationary Distributions

**Definition 5.3** (Invariant Measure). A probability measure $\pi$ is invariant for a stochastic process if starting from $\pi$ at time $t$, the distribution remains $\pi$ at all future times.

**For the VP-SDE with constant $\beta$:** The stationary distribution is $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

**Proof sketch:**

- If $\mathbf{x}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, then $\mathbb{E}[\mathbf{x}_t] = \mathbf{0}$

- The drift $-\frac{1}{2}\beta\mathbf{x}_t$ doesn't change the mean

- The diffusion maintains the covariance structure

- The balance between drift and diffusion preserves $\mathcal{N}(\mathbf{0}, \mathbf{I})$

**Convergence to invariant measure:** For any starting distribution $p_0$:

$$\lim_{t \to \infty} p_t = \mathcal{N}(\mathbf{0}, \mathbf{I}) \tag{5.28}$$

This explains why the forward process eventually produces pure noise regardless of the initial data distribution.

### 5.5.2 Time-Reversibility

**Theorem 5.2** (Time Reversal of SDEs). *If $\mathbf{x}_t$ follows the forward SDE:*

$$d\mathbf{x} = \mathbf{f}(\mathbf{x}, t)dt + g(t)d\mathbf{w} \tag{5.29}$$

*then the time-reversed process $\mathbf{y}_t = \mathbf{x}_{T-t}$ follows:*

$$d\mathbf{y} = [\mathbf{f}(\mathbf{y}, T-t) - g^2(T-t)\nabla_{\mathbf{y}} \log p_{T-t}(\mathbf{y})]dt + g(T-t)d\bar{\mathbf{w}} \tag{5.30}$$

*where $d\bar{\mathbf{w}}$ is a reverse-time Brownian motion.*

**For our VP-SDE:**

$$\text{Forward:} \quad d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\mathbf{w} \tag{5.31}$$

$$\text{Reverse:} \quad d\mathbf{y} = \left[-\frac{1}{2}\beta(t)\mathbf{y} - \beta(t)\nabla_{\mathbf{y}}\log p_t(\mathbf{y})\right]dt + \sqrt{\beta(t)}d\bar{\mathbf{w}} \tag{5.32}$$

**Key insights:**

- The reverse process is also an SDE with the same diffusion coefficient

- The drift includes an additional score-based term

- This explains why we need to learn the score function for generation

- The reverse SDE is exactly what we're approximating in diffusion model training

### 5.5.3   Marginal Distributions and Fokker-Planck Equation

The evolution of the probability density $p_t(\mathbf{x})$ is governed by the *Fokker-Planck equation*:

$$\frac{\partial p_t}{\partial t} = -\nabla \cdot (\mathbf{f}p_t) + \frac{1}{2}g^2\nabla^2 p_t \tag{5.33}$$

**For our VP-SDE:**

$$\frac{\partial p_t}{\partial t} = \frac{1}{2}\beta(t)\nabla \cdot (\mathbf{x}p_t) + \frac{1}{2}\beta(t)\nabla^2 p_t \tag{5.34}$$

**Physical interpretation:**

- $-\nabla \cdot (\mathbf{f}p_t)$: Advection term (systematic transport)

- $\frac{1}{2}g^2\nabla^2 p_t$: Diffusion term (spreading due to noise)

- The equation describes how probability "flows" through space

**Steady-state solution:** Setting $\frac{\partial p_t}{\partial t} = 0$ and solving gives the invariant measure $\mathcal{N}(\mathbf{0}, \mathbf{I})$.

### 5.5.4   Connections to Optimal Transport

The continuous-time formulation reveals deep connections to optimal transport theory.

**Optimal Transport Problem:** Given two probability distributions $p_0$ and $p_1$, find the "cheapest" way to transport mass from $p_0$ to $p_1$.

**Diffusion as Transport:**

- The probability flow ODE defines a transport map from $p_0$ to $p_T$

- This transport is generally not optimal in the classical sense

- However, it has other desirable properties (e.g., minimum entropy production)

**Schrödinger Bridges:** Recent work shows connections between diffusion models and Schrödinger bridges—optimal stochastic transport problems that minimize entropy production.

# 5.6 Practical Implications and Algorithms

## 5.6.1 Adaptive Timestep Selection

The continuous formulation enables adaptive algorithms that adjust timesteps based on local error estimates:

**Error-Controlled Integration:**

- Use ODE solver error estimates to adapt step sizes

- Take smaller steps when the solution changes rapidly

- Take larger steps in smooth regions

- Can dramatically improve efficiency

**Tolerance-Based Stopping:**

- Set target error tolerance

- Solver automatically chooses steps to meet tolerance

- Trade-off between speed and accuracy

- More principled than fixed timestep schedules

## 5.6.2 Higher-Order Integration Schemes

Beyond Euler-Maruyama, we can use sophisticated numerical methods:

**For SDEs:**

- Milstein scheme (second-order)

- Runge-Kutta for SDEs

- Implicit methods for stiff equations

**For ODEs:**

- Runge-Kutta methods (RK4, Dormand-Prince)

- Adaptive step size control

- High-order methods for smooth solutions

## 5.6.3 Connection to Neural ODEs

The probability flow ODE naturally connects to the Neural ODE framework:

**Neural ODE Formulation:**

$$\frac{d\mathbf{h}}{dt} = f_\theta(\mathbf{h}, t) \tag{5.35}$$

**Diffusion Neural ODE:**

$$\frac{d\mathbf{x}}{dt} = -\frac{1}{2}\beta(t)\mathbf{x} + \frac{1}{2}\beta(t)\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}, t)}{\sqrt{1 - \alpha^2(t)}} \tag{5.36}$$

**Advantages:**

- Can use existing Neural ODE software

- Automatic differentiation through the ODE solver

- Memory-efficient backpropagation

- Continuous normalizing flow connections

## 5.7 Conclusion: The Continuous Perspective

The continuous-time formulation of diffusion models provides a unifying mathematical framework that offers both theoretical insights and practical advantages. Key takeaways include:

**Theoretical Contributions:**

- **Unified Framework**: SDEs encompass many different discrete schemes

- **Deep Connections**: Links to optimal transport, score-based models, and physics

- **Mathematical Rigor**: Well-established SDE theory provides solid foundations

- **Limiting Behavior**: Clear understanding of what happens as $T \to \infty$

**Practical Advantages:**

- **Flexible Sampling**: Adaptive timesteps and sophisticated ODE solvers

- **Deterministic Options**: Probability flow ODEs for reproducible generation

- **Exact Likelihoods**: Change of variables formula for likelihood computation

- **Interpolation**: Meaningful latent space operations

**Future Directions:** The continuous perspective continues to drive innovation in generative modeling:

- **Advanced SDE Formulations**: Non-linear SDEs, state-dependent diffusion coefficients

- **Geometric Integration**: Structure-preserving numerical methods for specific manifolds

- **Multi-Scale Methods**: Hierarchical approaches that resolve different temporal scales

- **Hybrid Approaches**: Combining continuous theory with discrete optimizations

- **Physics-Informed Models**: Incorporating physical constraints and conservation laws

The journey from discrete timesteps to continuous flows illustrates the power of mathematical abstraction in machine learning. By embracing the continuous perspective, we gain not only computational flexibility but also theoretical understanding that guides the development of next-generation generative models.

As we look toward the future, the continuous-time framework provides a solid foundation for addressing fundamental questions in generative modeling: How can we guarantee optimal transport between distributions? What are the fundamental limits of diffusion-based generation? How can we incorporate domain-specific knowledge into the diffusion process? The mathematical tools developed in this chapter provide a starting point for answering these questions and many others.

## 5.8   Exercises

1. **SDE to Discrete Conversion:**

    (a) Start with the VP-SDE: $d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x}dt + \sqrt{\beta(t)}d\mathbf{w}$

    (b) Use the Euler-Maruyama scheme to derive the discrete approximation

    (c) Show that this matches the DDPM forward process for appropriate choice of $\beta_t$

    (d) Implement both versions and compare their behavior for different step sizes

2. **Ornstein-Uhlenbeck Analysis:**

    (a) For the constant coefficient case $d\mathbf{x} = -\theta\mathbf{x}dt + \sigma d\mathbf{w}$, derive the exact solution

    (b) Show that the stationary distribution is $\mathcal{N}(\mathbf{0}, \frac{\sigma^2}{2\theta}\mathbf{I})$

    (c) Compute the autocorrelation function $\mathbb{E}[\mathbf{x}_t \cdot \mathbf{x}_s]$

    (d) Verify numerically by simulating the process

3. **Probability Flow ODE Implementation:**

    (a) Implement the probability flow ODE for a trained diffusion model

    (b) Compare deterministic vs. stochastic sampling on the same dataset

    (c) Analyze the trade-offs in terms of sample quality and diversity

    (d) Experiment with different ODE solvers (Euler, RK4, adaptive methods)

4. **Variance Preserving vs. Exploding:**

    (a) Implement both VP and VE SDEs

    (b) Compare their behavior on a simple 2D dataset

    (c) Analyze the evolution of mean and variance over time

    (d) Discuss the pros and cons of each approach

5. **Fokker-Planck Equation:**

(a) For a 1D VP-SDE with constant $\beta$, write out the Fokker-Planck equation

(b) Find the steady-state solution analytically

(c) Verify by numerical simulation that the process converges to this distribution

(d) Investigate the convergence rate as a function of $\beta$

6. **Time Reversal:**

(a) Derive the reverse-time SDE for the VP process

(b) Explain why the score function $\nabla \log p_t(\mathbf{x})$ appears in the reverse drift

(c) Implement a simple score-based model and verify the time-reversal property

(d) Compare with the standard diffusion model reverse process

7. **Adaptive Timestep Integration:**

(a) Implement an adaptive ODE solver for the probability flow ODE

(b) Compare fixed vs. adaptive timesteps in terms of efficiency and accuracy

(c) Analyze where the solver chooses fine vs. coarse timesteps

(d) Develop heuristics for good initial timestep selection

8. **Continuous Normalizing Flows Connection:**

(a) Show how the probability flow ODE relates to continuous normalizing flows

(b) Implement likelihood computation using the change of variables formula

(c) Compare likelihood estimates from the continuous vs. discrete formulations

(d) Investigate the computational trade-offs

**Software Resources:**

- `diffrax`: JAX library for differential equation solving

- `torchdiffeq`: PyTorch Neural ODE implementations

- `sdeint`: Python SDE integration library

- `DiffEqFlux.jl`: Julia ecosystem for scientific machine learning

The continuous-time perspective represents a mature and rapidly evolving area of research, with new theoretical insights and practical algorithms being developed regularly. The mathematical framework presented in this chapter provides the foundation for understanding and contributing to this exciting field.

# Chapter 6

# Score-Based Generative Models: Learning the Geometry of Data

## 6.1 A Different Perspective on Generation

Imagine you're lost in a mountainous landscape at night, trying to find the highest peak where a bonfire awaits. You can't see the entire terrain, but you have a magical compass that always points toward the steepest uphill direction. By following this compass step by step, you would eventually reach the summit. This is precisely the intuition behind *score-based generative models*—a revolutionary approach that has transformed how we think about generating data.

Throughout our exploration of generative models, we have encountered various strategies for tackling the fundamental challenge of learning data distributions. GANs pit two neural networks against each other in an adversarial game, VAEs compress data through a bottleneck and then reconstruct it, and diffusion models learn to gradually remove noise. Each approach has its merits, but they all share a common goal: learning to generate samples that look like they come from our training data.

Now, we introduce yet another elegant perspective that might initially seem completely different from everything we've seen before. However, as we'll discover, score-based models share surprisingly deep mathematical connections with diffusion models—so deep, in fact, that they can be viewed as different parameterizations of the same underlying mathematical framework. This connection isn't just academically interesting; it has led to some of the most powerful generative models we have today.

The central insight behind score-based models is both simple and profound: instead of trying to learn the probability density $p(\mathbf{x})$ directly—which is notoriously difficult—we learn something else entirely. We learn the *score function* $\nabla_{\mathbf{x}} \log p(\mathbf{x})$, which tells us the direction of steepest increase in log-probability at any point. This seemingly minor shift in perspective opens up entirely new theoretical insights and practical algorithms that have revolutionized generative modeling.

But why would we want to learn gradients instead of probabilities? The answer lies in a fundamental limitation that has plagued generative modeling from the beginning: the normalization problem.

## 6.2   The Motivation: Why Score Functions?

Before diving into the mathematical details, let's understand why someone would think to learn score functions in the first place. The story begins with a fundamental challenge in probabilistic modeling.

When we want to learn a probability distribution $p(\mathbf{x})$ from data, we typically face the *intractable normalization problem*. Most flexible probability models can be written as:

$$p(\mathbf{x}) = \frac{1}{Z}\tilde{p}(\mathbf{x}) \tag{6.1}$$

where $\tilde{p}(\mathbf{x})$ is an unnormalized density that's easy to compute, and $Z = \int \tilde{p}(\mathbf{x})d\mathbf{x}$ is the normalization constant (partition function). The problem? Computing $Z$ requires integrating over the entire data space, which is typically impossible for high-dimensional data.

This is where score functions come to the rescue. Here's the key insight: when we take the gradient of the log-probability, something magical happens:

$$\nabla_{\mathbf{x}} \log p(\mathbf{x}) = \nabla_{\mathbf{x}} \log \left( \frac{1}{Z}\tilde{p}(\mathbf{x}) \right) \tag{6.2}$$

$$= \nabla_{\mathbf{x}}(\log \tilde{p}(\mathbf{x}) - \log Z) \tag{6.3}$$

$$= \nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}) - \nabla_{\mathbf{x}} \log Z \tag{6.4}$$

$$= \nabla_{\mathbf{x}} \log \tilde{p}(\mathbf{x}) \tag{6.5}$$

The normalization constant $Z$ disappears! This is because $Z$ doesn't depend on $\mathbf{x}$, so its gradient is zero. Suddenly, we can work with unnormalized models without ever computing the intractable partition function.

This observation leads us naturally to the formal definition of what we're trying to learn.

## 6.3   The Score Function: Your Probability Compass

### 6.3.1   Mathematical Definition and Intuition

**Definition 6.1** (Score Function). For a probability density function $p(\mathbf{x})$, the *score function* (also called the Stein score) is defined as:

$$\mathbf{s}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) \tag{6.6}$$

Let's unpack this definition piece by piece, because each component tells part of the story:

- $\log p(\mathbf{x})$: Taking the logarithm transforms products into sums, making the mathematics more tractable and numerically stable.

- $\nabla_{\mathbf{x}}$: The gradient operator computes the direction and magnitude of steepest increase.

- $\mathbf{s}(\mathbf{x})$: The resulting vector field that serves as our "probability compass"—pointing toward regions of higher likelihood.

Think of the score function as providing local directions in the probability landscape. At any point $\mathbf{x}$, it tells you: "If you want to find more likely data points, move in this direction."

## 6.3.2 Geometric Intuition: Navigating the Probability Landscape

The score function has a beautiful and intuitive geometric interpretation that makes its power immediately clear. Imagine the probability density $p(\mathbf{x})$ as a mountainous landscape where the height at each point represents the probability of that data point:

- **Peaks**: High-probability regions where real data tends to concentrate (the modes of the distribution)

- **Valleys**: Low-probability regions where real data is rare

- **Score vectors**: Arrows scattered throughout the landscape, each pointing in the direction of steepest uphill climb toward higher probability

This geometric picture immediately reveals several important properties of the score function:

- $\|\mathbf{s}(\mathbf{x})\|$ is large where probability changes rapidly—these are the steep slopes of our probability mountain

- $\|\mathbf{s}(\mathbf{x})\|$ is small in flat regions or near modes—the gentle slopes and peaks

- $\mathbf{s}(\mathbf{x}) = \mathbf{0}$ exactly at local maxima of $p(\mathbf{x})$—at the very tops of peaks, there's no uphill direction

- The score field provides a complete "roadmap" for navigating from any point toward high-probability regions

This last point is crucial for generation: if we can learn the score function accurately, we can start from random noise and follow the score field like a hiker following trail markers, eventually reaching regions where real data lives.

## 6.3.3 Connection to Energy-Based Models: A Unifying Perspective

Now that we understand what score functions are, let's explore a powerful alternative way to think about them that will deepen our intuition. This perspective comes from asking a simple question: *What if we think about probability in terms of energy?*

This might seem like an odd question, but it leads to one of the most elegant unifying principles in machine learning. The idea comes from statistical physics, where systems naturally evolve toward lower-energy states—think of a ball rolling downhill to settle in a valley. What if we could think about data distributions the same way?

Here's the key insight: we can reformulate any probability distribution in terms of an *energy function*. Instead of directly specifying probabilities, we specify how much "energy" each data point has, with the understanding that lower energy corresponds to higher probability.

This energy-based perspective allows us to express many probability distributions in a unified form:

$$p(\mathbf{x}) = \frac{1}{Z} e^{-E(\mathbf{x})} \tag{6.7}$$

Let's unpack this formula:

- $E(\mathbf{x})$ is the *energy function*—it assigns an energy value to each possible data point $\mathbf{x}$

- $e^{-E(\mathbf{x})}$ converts energy to probability: high energy means low probability, low energy means high probability

- $Z = \int e^{-E(\mathbf{x})} d\mathbf{x}$ is the partition function that ensures the probabilities sum to 1

Why is this perspective useful? Because it gives us an intuitive way to think about data: **natural, realistic data points have low energy, while unnatural, unrealistic data points have high energy**. For images, a photo of a cat might have low energy, while random pixel noise has high energy.

This formulation comes from statistical physics, where lower energy corresponds to more stable (and hence more probable) states. In our context, think of energy as a measure of how "unusual" or "unlikely" a data point is—natural images have low energy, while random noise has high energy.

Now, let's see what happens when we compute the score function for an energy-based model:

$$\mathbf{s}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) \tag{6.8}$$

$$= \nabla_{\mathbf{x}} \log \left( \frac{1}{Z} e^{-E(\mathbf{x})} \right) \tag{6.9}$$

$$= \nabla_{\mathbf{x}} (-E(\mathbf{x}) - \log Z) \tag{6.10}$$

$$= -\nabla_{\mathbf{x}} E(\mathbf{x}) \tag{6.11}$$

This reveals a beautiful relationship: **the score function is simply the negative gradient of the energy function**. This connection provides powerful intuition:

- High-probability regions correspond to low energy (like balls settling in valleys)

- The score points toward lower energy, which means higher probability

- Following the score function is like following the path of steepest descent in energy space

- We can learn scores without ever computing the intractable partition function $Z$

This energy perspective also explains why score-based models are so effective: they're essentially learning the "force field" that pushes data points toward more natural, lower-energy configurations.

### 6.3.4 Concrete Examples: Seeing Scores in Action

To make these abstract concepts concrete, let's examine the score functions for distributions we understand well.

**Example 1: The Gaussian Distribution**

Consider a multivariate Gaussian distribution $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$. Let's derive the score function step by step.

**Step 1: Write out the full probability density**

The multivariate Gaussian probability density function is:

$$p(\mathbf{x}) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right) \tag{6.12}$$

where $d$ is the dimensionality of $\mathbf{x}$, and $|\boldsymbol{\Sigma}|$ is the determinant of the covariance matrix.

**Step 2: Take the logarithm**

Taking the natural logarithm:

$$\log p(\mathbf{x}) = \log\left[\frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)\right] \tag{6.13}$$

$$= \log\left[\frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}}\right] + \log\left[\exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right)\right] \tag{6.14}$$

$$= -\frac{d}{2}\log(2\pi) - \frac{1}{2}\log|\boldsymbol{\Sigma}| - \frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \tag{6.15}$$

**Step 3: Identify terms that depend on x**

For computing the gradient $\nabla_{\mathbf{x}} \log p(\mathbf{x})$, we only need terms that depend on $\mathbf{x}$:

$$\log p(\mathbf{x}) = -\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) + \text{const} \tag{6.16}$$

where $\text{const} = -\frac{d}{2}\log(2\pi) - \frac{1}{2}\log|\boldsymbol{\Sigma}|$ contains all terms independent of $\mathbf{x}$.

**Step 4: Compute the gradient**

Now we take the gradient with respect to $\mathbf{x}$. Let $\mathbf{z} = \mathbf{x} - \boldsymbol{\mu}$, so we have:

$$\frac{\partial}{\partial \mathbf{x}}\left[-\frac{1}{2}\mathbf{z}^T \boldsymbol{\Sigma}^{-1} \mathbf{z}\right] \tag{6.17}$$

Using the matrix calculus identity $\frac{\partial}{\partial \mathbf{x}}(\mathbf{x}^T \mathbf{A}\mathbf{x}) = (\mathbf{A} + \mathbf{A}^T)\mathbf{x}$, and noting that $\boldsymbol{\Sigma}^{-1}$ is symmetric (so $\boldsymbol{\Sigma}^{-1} = (\boldsymbol{\Sigma}^{-1})^T$):

$$\nabla_{\mathbf{x}}\left[-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu})^T \boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu})\right] = -\frac{1}{2} \cdot 2\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \tag{6.18}$$

$$= -\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \tag{6.19}$$

**Step 5: Final result**

Therefore, the score function for a multivariate Gaussian is:

$$\mathbf{s}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) = -\boldsymbol{\Sigma}^{-1}(\mathbf{x} - \boldsymbol{\mu}) \tag{6.20}$$

This simple formula reveals beautiful intuition: the score always points toward the mean $\boldsymbol{\mu}$, with strength that depends inversely on the variance. Points far from the mean experience a strong "pull" back toward the center, while points near the mean experience gentle guidance. The covariance matrix $\boldsymbol{\Sigma}$ shapes this force field—high variance directions have weaker pull, while low variance directions have stronger pull.

**Example 2: Mixture of Gaussians**

For a more complex distribution—a mixture of Gaussians—let's derive the score function step by step to see how multiple modes interact.

**Step 1: Define the mixture distribution**

A mixture of $K$ Gaussians is defined as:

$$p(\mathbf{x}) = \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{6.21}$$

where $\pi_k \geq 0$ are the mixture weights satisfying $\sum_{k=1}^{K} \pi_k = 1$, and each component is:

$$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}_k|^{1/2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \boldsymbol{\mu}_k)^T \boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)\right) \tag{6.22}$$

**Step 2: Take the logarithm**

$$\log p(\mathbf{x}) = \log\left(\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right) \tag{6.23}$$

Unlike the single Gaussian case, we cannot simplify this logarithm of a sum.

**Step 3: Apply the chain rule for differentiation**

To find the score function, we compute:

$$\mathbf{s}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x}) = \nabla_{\mathbf{x}} \log\left(\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right) \tag{6.24}$$

$$= \frac{1}{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \cdot \nabla_{\mathbf{x}}\left(\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right) \tag{6.25}$$

**Step 4: Compute the gradient of the sum**

The gradient of the sum is the sum of gradients:

$$\nabla_{\mathbf{x}}\left(\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)\right) = \sum_{k=1}^{K} \pi_k \nabla_{\mathbf{x}} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{6.26}$$

**Step 5: Gradient of each Gaussian component**

For each Gaussian component, we can use our result from Example 1. Since:

$$\nabla_{\mathbf{x}} \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = -\boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k) \tag{6.27}$$

And using the identity $\nabla_{\mathbf{x}} f(\mathbf{x}) = f(\mathbf{x}) \cdot \nabla_{\mathbf{x}} \log f(\mathbf{x})$:

$$\nabla_{\mathbf{x}} \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \cdot \nabla_{\mathbf{x}} \log \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \tag{6.28}$$

$$= \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \cdot (-\boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)) \tag{6.29}$$

**Step 6: Combine everything**

Substituting back into our expression:

$$\mathbf{s}(\mathbf{x}) = \frac{1}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \cdot \sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \cdot (-\boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)) \qquad (6.30)$$

This gives us the final result:

$$\boxed{\mathbf{s}(\mathbf{x}) = \frac{\sum_{k=1}^{K} \pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k) \cdot (-\boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k))}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}} \qquad (6.31)$$

**Step 7: Interpretation using posterior probabilities**

We can rewrite this in a more intuitive form. Define the posterior probability of component $k$ given $\mathbf{x}$:

$$\gamma_k(\mathbf{x}) = \frac{\pi_k \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)}{\sum_{j=1}^{K} \pi_j \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \qquad (6.32)$$

Then the score function becomes:

$$\mathbf{s}(\mathbf{x}) = \sum_{k=1}^{K} \gamma_k(\mathbf{x}) \cdot (-\boldsymbol{\Sigma}_k^{-1}(\mathbf{x} - \boldsymbol{\mu}_k)) \qquad (6.33)$$

**Beautiful interpretation:** The score function is a weighted average of the individual Gaussian score functions, where the weights $\gamma_k(\mathbf{x})$ represent how much each component "believes" it is responsible for the point $\mathbf{x}$. This automatically balances the competing influences of multiple modes:

- Near any particular mode $k$, we have $\gamma_k(\mathbf{x}) \approx 1$ and $\gamma_j(\mathbf{x}) \approx 0$ for $j \neq k$, so that mode dominates and pulls points toward its center $\boldsymbol{\mu}_k$

- In regions between modes, the score represents a weighted combination of pulls from different directions, naturally guiding points toward the nearest or most influential mode

- The strength of each pull is modulated by both the mixture weight $\pi_k$ and how well component $k$ explains the current point

This example illustrates why learning score functions is so powerful: even for complex, multi-modal distributions, the score function provides a unified way to navigate the probability landscape, automatically handling the interactions between different modes without requiring us to explicitly model the mixture structure.

Now comes the exciting part: how do we actually use these score functions to generate new data?

### 6.3.5 Why Score Functions Enable Generation

Here's where the magic happens. One of the most remarkable properties of score functions is that they enable us to sample from a distribution even without knowing the normalization constant—that intractable partition function that has plagued probabilistic modeling for decades.

The key insight is both simple and profound: **if we can estimate $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ accurately, we can sample from $p(\mathbf{x})$ without ever computing $p(\mathbf{x})$ itself**. This is accomplished through a beautiful algorithm called *Langevin dynamics*, which follows the score field while adding controlled noise to explore the distribution.

Think back to our hiker analogy: imagine the hiker not only has a compass pointing uphill, but also occasionally takes random steps in unexpected directions. This combination of following the gradient (systematic exploration) and random steps (stochastic exploration) ensures the hiker doesn't get stuck in local valleys and eventually explores all the peaks in the mountain range.

### 6.3.6 Langevin Dynamics: Following the Score Field

Langevin dynamics provides us with a systematic way to convert our score function into a sampling algorithm. The algorithm is elegantly simple, yet theoretically guaranteed to work.

**procedure** LANGEVINSAMPLING($\mathbf{s}(\mathbf{x})$, $\eta$, $T$)
    **Input:** Score function $\mathbf{s}(\mathbf{x}) = \nabla_{\mathbf{x}} \log p(\mathbf{x})$
           Step size $\eta$, number of steps $T$
    Initialize $\mathbf{x}_0$ from some distribution (e.g., $\mathcal{N}(\mathbf{0}, \mathbf{I})$)
    **for** $t = 0, 1, \ldots, T-1$ **do**
        Sample $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
        $\mathbf{x}_{t+1} = \mathbf{x}_t + \eta \mathbf{s}(\mathbf{x}_t) + \sqrt{2\eta} \boldsymbol{\epsilon}_t$
    **end for**
    **Return:** $\mathbf{x}_T$
**end procedure**

The beauty of this algorithm lies in its update rule, which balances two competing forces:

- $\eta \mathbf{s}(\mathbf{x}_t)$: **Drift term**—systematically moves toward higher probability regions, following our probability compass

- $\sqrt{2\eta} \boldsymbol{\epsilon}_t$: **Diffusion term**—adds controlled randomness to prevent getting trapped in local modes and ensure thorough exploration

The delicate balance between these terms is crucial. Too much drift and we get stuck in the first high-probability region we find. Too much diffusion and we wander randomly without making progress toward realistic data. The algorithm automatically finds the right balance to efficiently explore the entire probability landscape.

**Theoretical guarantee:** As the step size $\eta \to 0$ and the number of steps $T \to \infty$, Langevin dynamics is guaranteed to converge to the target distribution $p(\mathbf{x})$, regardless of where we start. This is a remarkable result—it means that

even if we initialize with pure noise, we'll eventually generate samples that look like real data.

### 6.3.7 Connection to Physics: Nature's Sampling Algorithm

Langevin dynamics isn't just a clever mathematical trick—it's actually how nature works. The algorithm has deep roots in statistical physics, where it describes the motion of particles suspended in a fluid.

Consider a tiny pollen grain floating in water (this is Brownian motion, first observed by botanist Robert Brown in 1827). The grain experiences two forces:

- **Systematic forces**: Perhaps an electric field or gravity pulling it in a particular direction

- **Random collisions**: Water molecules constantly bumping into it from all directions

Over time, the pollen grain settles into an equilibrium distribution that balances these forces. This is exactly what Langevin dynamics simulates in probability space:

- **Particle**: A sample $\mathbf{x}_t$ in data space

- **Systematic force**: The score function $\mathbf{s}(\mathbf{x})$ pulling toward high-probability regions

- **Random collisions**: The noise term $\sqrt{2\eta}\boldsymbol{\epsilon}_t$ providing thermal fluctuations

- **Equilibrium distribution**: Our target data distribution $p(\mathbf{x})$

This physical interpretation explains why the algorithm works so well: we're not fighting against nature, we're using the same principles that govern how physical systems naturally reach equilibrium.

But there's a catch: to use Langevin dynamics, we need to know the score function. And that's where the real challenge begins.

## 6.4 Score Matching: Learning the Probability Compass

### 6.4.1 The Fundamental Chicken-and-Egg Problem

Now we face a classic chicken-and-egg dilemma that initially seems insurmountable:

- We want to learn $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ so we can use Langevin dynamics to sample from $p(\mathbf{x})$

- But to compute $\nabla_{\mathbf{x}} \log p(\mathbf{x})$, we need to know $p(\mathbf{x})$ in the first place

- We don't know $p(\mathbf{x})$—learning it is the whole point of our generative model!

This seems hopeless. How can we learn the gradient of something we don't know? The answer lies in a brilliant technique called *score matching*, which finds a way to learn score functions directly from data without ever explicitly modeling the probability density.

### 6.4.2 The Naive Approach and Why It Fails

Let's start by understanding why the obvious approach doesn't work. We might try to train a neural network $\mathbf{s}_\theta(\mathbf{x})$ to approximate the true score function by minimizing:

$$\mathcal{L}_{\text{naive}} = \mathbb{E}_{p(\mathbf{x})} \left[ \|\mathbf{s}_\theta(\mathbf{x}) - \nabla_{\mathbf{x}} \log p(\mathbf{x})\|^2 \right] \tag{6.34}$$

This loss makes perfect sense: we want our learned score function $\mathbf{s}_\theta(\mathbf{x})$ to be as close as possible to the true score function $\nabla_{\mathbf{x}} \log p(\mathbf{x})$ for every point in the data distribution.

**The fatal flaw:** We don't know $\nabla_{\mathbf{x}} \log p(\mathbf{x})$—that's exactly what we're trying to learn! So we can't compute this loss function at all.

### 6.4.3 Score Matching: The Elegant Mathematical Solution

The breakthrough insight of score matching is that we can mathematically transform this impossible loss function into one that we can actually compute. Through a clever application of integration by parts, we can eliminate the unknown true score from the objective entirely.

**Theorem 6.1** (Score Matching Objective). *The intractable naive loss is mathematically equivalent to:*

$$\mathcal{L}_{SM} = \mathbb{E}_{p(\mathbf{x})} \left[ tr(\nabla_{\mathbf{x}} \mathbf{s}_\theta(\mathbf{x})) + \frac{1}{2} \|\mathbf{s}_\theta(\mathbf{x})\|^2 \right] + const \tag{6.35}$$

This transformation is mathematically beautiful and practically revolutionary:

- The trace term $tr(\nabla_{\mathbf{x}} \mathbf{s}_\theta(\mathbf{x})) = \sum_i \frac{\partial s_{\theta,i}(\mathbf{x})}{\partial x_i}$ only depends on our learned score function and its derivatives

- The squared norm term $\|\mathbf{s}_\theta(\mathbf{x})\|^2$ also only depends on our learned score function

- The unknown true score has vanished completely through the integration by parts

- We can compute this loss using only samples from our training data

In practice, we simply compute:

$$\mathcal{L}_{\text{SM}} = \frac{1}{N} \sum_{i=1}^{N} \left[ tr(\nabla_{\mathbf{x}} \mathbf{s}_\theta(\mathbf{x}_i)) + \frac{1}{2} \|\mathbf{s}_\theta(\mathbf{x}_i)\|^2 \right] \tag{6.36}$$

where $\{\mathbf{x}_i\}_{i=1}^{N}$ are samples from our training dataset.

This is remarkable: we can learn the gradient of a probability distribution without ever learning the distribution itself!

### 6.4.4 Denoising Score Matching: A More Practical Path

While standard score matching is mathematically elegant, computing the trace term $\text{tr}(\nabla_{\mathbf{x}}\mathbf{s}_\theta(\mathbf{x}))$ can be computationally expensive, especially for high-dimensional data. This has led to the development of *denoising score matching*, which provides a more practical alternative with a beautiful intuitive interpretation.

The key insight is deceptively simple: instead of learning the score of the clean data distribution directly, we learn the score of a "noisy" version of the data distribution. This small change makes the problem much more tractable while maintaining the essential properties we need for generation.

**The denoising score matching objective:**

$$\mathcal{L}_{\text{DSM}} = \mathbb{E}_{\mathbf{x}_0 \sim p_{data}(\mathbf{x})}\mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mathbf{x};\mathbf{x}_0,\sigma^2\mathbf{I})}\left[\|\mathbf{s}_\theta(\mathbf{x}) - \nabla_{\mathbf{x}}\log q(\mathbf{x}|\mathbf{x}_0)\|^2\right] \qquad (6.37)$$

Here's what's happening: we take a clean data point $\mathbf{x}_0$, add Gaussian noise to get $\mathbf{x}$, and then train our score network to predict the score of this noisy distribution.

For Gaussian noise $q(\mathbf{x}|\mathbf{x}_0) = \mathcal{N}(\mathbf{x};\mathbf{x}_0,\sigma^2\mathbf{I})$, we can compute the true score analytically:

$$\nabla_{\mathbf{x}}\log q(\mathbf{x}|\mathbf{x}_0) = -\frac{\mathbf{x} - \mathbf{x}_0}{\sigma^2} \qquad (6.38)$$

This gives us a wonderfully simple training procedure:

**procedure** DENOSINGSCOREMATCHING(dataset $\{\mathbf{x}_0^{(i)}\}$, noise level $\sigma$)
    **for** each training iteration **do**
        Sample $\mathbf{x}_0$ from the training data
        Sample $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0},\mathbf{I})$
        Set $\mathbf{x} = \mathbf{x}_0 + \sigma\boldsymbol{\epsilon}$ (add noise)
        Compute loss: $\mathcal{L} = \left\|\mathbf{s}_\theta(\mathbf{x}) + \frac{\boldsymbol{\epsilon}}{\sigma}\right\|^2$
        Update $\theta$ using gradient descent
    **end for**
**end procedure**

**The beautiful connection:** This objective is equivalent to training a neural network to predict the noise that was added to clean data. If you've heard of denoising autoencoders, this should sound familiar—we're essentially training the network to denoise data, but framing it in terms of score functions!

The network learns to look at noisy data and figure out which direction to move to make it more realistic. This is exactly what we want for generation: starting from noise, the score function tells us how to iteratively make it look more and more like real data.

But as we dive deeper into practical implementation, we encounter a new challenge that reveals an important limitation of our approach so far.

## 6.5 The Multi-Scale Challenge: Why One Noise Level Isn't Enough

### 6.5.1 The Goldilocks Problem of Noise

As we begin implementing denoising score matching in practice, we quickly discover that choosing the right noise level $\sigma$ is like the classic Goldilocks story—it's

surprisingly difficult to get it "just right." This isn't just a minor tuning issue; it reveals a fundamental limitation that we must address.

Let's explore what happens when our noise level choice goes wrong:

**When the noise is too small ($\sigma$ is small):**

Imagine trying to learn a score function with very little noise added to the training data. The noisy samples stay very close to the original clean data points, which sounds good at first. Our score network becomes very accurate at predicting scores in high-density regions where real data lives.

However, this creates a critical problem: what happens in the vast empty spaces between data points? In these low-density regions, our network has never seen training examples, so its score predictions become unreliable or even completely wrong. During generation with Langevin dynamics, if our sampling trajectory ever wanders into these poorly-learned regions, it might get completely stuck. The score function might point in random directions, causing the sampler to waste time wandering aimlessly instead of progressing toward realistic data.

**When the noise is too large ($\sigma$ is large):**

Now imagine the opposite extreme: we add so much noise that our training samples are spread out everywhere, covering even the low-density regions between data modes. This solves the coverage problem—our network learns reasonable score predictions throughout the entire space.

But this creates a different problem: we've fundamentally changed the distribution we're learning. The heavily-noised version of our data distribution looks nothing like the original. It's like trying to learn the shape of a mountain by studying it during a thick fog—the smooth, over-blurred version might be easier to navigate, but it doesn't capture the sharp peaks and valleys of the true landscape. Generated samples might be somewhat realistic but lack the fine details and sharp features that make data truly convincing.

This is our Goldilocks dilemma: small noise gives us accuracy but poor coverage, while large noise gives us coverage but poor accuracy. We need a solution that gives us both.

## 6.5.2   The Multi-Scale Solution: Learning All Noise Levels

The breakthrough insight is elegantly simple: instead of trying to find the perfect single noise level, why not learn them all? We can train a single network that knows how to handle every noise level from very small to very large.

This leads us to *noise-conditional score networks* $\mathbf{s}_\theta(\mathbf{x}, \sigma)$—networks that take both the noisy data point and the noise level as inputs. It's like training a universal translator that doesn't just understand one language, but can switch between languages on command.

**The multi-scale training objective:**

$$\mathcal{L}_{\text{multi}} = \mathbb{E}_{\sigma \sim p(\sigma)} \mathbb{E}_{\mathbf{x}_0 \sim p_{data}} \mathbb{E}_{\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \left[ \lambda(\sigma) \left\| \mathbf{s}_\theta(\mathbf{x}_0 + \sigma\boldsymbol{\epsilon}, \sigma) + \frac{\boldsymbol{\epsilon}}{\sigma} \right\|^2 \right] \qquad (6.39)$$

Let's break down what's happening here:

- $p(\sigma)$: We sample noise levels from some distribution (often uniform over a range $[\sigma_{\min}, \sigma_{\max}]$)

- $\lambda(\sigma)$: A weighting function that ensures we balance learning across different noise levels appropriately

- $\mathbf{s}_\theta(\mathbf{x}, \sigma)$: Our network learns to predict scores conditional on the noise level

**Architecture considerations:**

To make this work in practice, we need to tell our neural network what noise level it's dealing with. This is typically done through *noise level embedding*:

- The noise level $\sigma$ is encoded using sinusoidal embeddings (similar to positional encoding in transformers)

- This embedding is injected at multiple layers throughout the network

- The network learns to modulate its behavior based on the noise level, like having different "modes" of operation

This approach is remarkably similar to how diffusion models handle time steps, and as we'll see later, this isn't a coincidence—it's a deep mathematical connection.

### 6.5.3 Annealed Langevin Dynamics: A Coarse-to-Fine Journey

With our noise-conditional score network in hand, we can now design a much more sophisticated sampling procedure called *annealed Langevin dynamics*. The word "annealed" comes from metallurgy, where metal is slowly cooled to remove imperfections and reach a more stable state.

**procedure** ANNEALEDLANGEVIN(noise schedule $\{\sigma_1 > \sigma_2 > \cdots > \sigma_L\}$)

 **Input:** Decreasing sequence of noise levels

 Initialize $\mathbf{x}_0 \sim \mathcal{N}(\mathbf{0}, \sigma_1^2\mathbf{I})$ (start from pure noise)

 **for** $\ell = 1, 2, \ldots, L$ **do**         ▷ For each noise level

  **for** $t = 1, 2, \ldots, T_\ell$ **do**      ▷ Multiple Langevin steps

   Sample $\boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

   $\mathbf{x}_t = \mathbf{x}_{t-1} + \eta_\ell \mathbf{s}_\theta(\mathbf{x}_{t-1}, \sigma_\ell) + \sqrt{2\eta_\ell}\boldsymbol{\epsilon}_t$

  **end for**

  $\mathbf{x}_0 = \mathbf{x}_{T_\ell}$          ▷ Prepare for next noise level

 **end for**

 **Return:** $\mathbf{x}_0$

**end procedure**

This algorithm implements a beautiful coarse-to-fine generation strategy:

- **Large $\sigma$ (coarse)**: We start with high noise levels where the score network has learned to navigate the global structure. The high noise smooths out local details, making it easy to find the general vicinity of data modes

- **Decreasing $\sigma$ (refinement)**: As we reduce the noise level, the score function becomes more precise, gradually revealing finer details of the data distribution

- **Small $\sigma$ (fine)**: At the lowest noise levels, we perform final polishing steps that add realistic details and sharp features

This is like creating a sculpture: you start with rough cuts to establish the basic shape, then gradually use finer tools to add details, and finally polish to achieve the perfect finish.

**Why annealing works so well:**

The annealing strategy embodies a principle from optimization and learning theory: *curriculum learning.* Just as humans learn better when presented with concepts from easy to hard, our sampling algorithm works better when it tackles the generation problem from coarse to fine:

- **Global exploration first**: High noise levels help us explore the entire probability landscape without getting trapped in local modes

- **Gradual refinement**: Each step down in noise level builds on the previous step's progress

- **Avoiding local optima**: The coarse-to-fine approach helps avoid getting stuck in poor local solutions

- **Computational efficiency**: We spend more steps refining promising regions rather than wandering randomly

**Connection to diffusion models:** If this sounds familiar, it should! This is exactly analogous to the reverse diffusion process, where we gradually remove noise from a sample. This parallel hints at a deeper mathematical connection that we're about to explore.

# 6.6 The Great Unification: Score-Based Models Meet Diffusion Models

## 6.6.1 Two Paths to the Same Mountain Peak

One of the most profound discoveries in recent generative modeling is that what initially appeared to be two completely different approaches—score-based models and diffusion models—are actually two different perspectives on the same underlying mathematical framework.

Imagine two hikers taking different paths up the same mountain. One follows the ridgeline, constantly looking at the steepest upward direction (the score-based approach). The other follows a marked trail that gradually ascends by removing obstacles one by one (the diffusion approach). They're using different strategies and different maps, but they're climbing the same mountain and will reach the same peak.

This realization has been transformative for the field, revealing that the noise prediction networks in diffusion models are secretly learning score functions, and that score-based models are implicitly performing a diffusion-like process.

## 6.6.2 The Mathematical Rosetta Stone

Let's uncover this connection step by step. The key insight is that the noise prediction network $\epsilon_\theta(\mathbf{x}_t, t)$ in diffusion models is directly related to the score function through a simple mathematical transformation.

Recall from our study of diffusion models that the forward process is defined by:

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \tag{6.40}$$

Now, let's compute the score function of this distribution:

$$\log q(\mathbf{x}_t|\mathbf{x}_0) = -\frac{1}{2(1 - \bar{\alpha}_t)}\|\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0\|^2 + \text{const} \tag{6.41}$$

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) = -\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{1 - \bar{\alpha}_t} \tag{6.42}$$

Here comes the crucial step. Using the reparameterization from diffusion models:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} \tag{6.43}$$

We can substitute this into our score expression:

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) = -\frac{\sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}}{1 - \bar{\alpha}_t} = -\frac{\boldsymbol{\epsilon}}{\sqrt{1 - \bar{\alpha}_t}} \tag{6.44}$$

This gives us the fundamental equivalence relationship:

$$\boxed{\mathbf{s}_\theta(\mathbf{x}_t, t) = -\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}}} \tag{6.45}$$

This simple equation is the Rosetta Stone that translates between the two frameworks:

- **Learning to predict noise** $\boldsymbol{\epsilon}_\theta$ is mathematically equivalent to learning the score function

- **Diffusion models implicitly learn score functions** without ever calling them that

- **Score-based models** can be viewed as diffusion models with a different parameterization

- **Both approaches** are learning the same underlying probability landscape, just using different coordinates

### 6.6.3 Training Objectives: Different Equations, Same Optimization

To make this connection even more concrete, let's compare the training objectives and show they're essentially identical.

**Score-based training objective:**

$$\mathcal{L}_{\text{score}} = \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}}\left[\|\mathbf{s}_\theta(\mathbf{x}_t, t) - \nabla_{\mathbf{x}_t}\log q(\mathbf{x}_t|\mathbf{x}_0)\|^2\right] \tag{6.46}$$

**Diffusion training objective:**

$$\mathcal{L}_{\text{diffusion}} = \mathbb{E}_{t,\mathbf{x}_0,\boldsymbol{\epsilon}}\left[\|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\|^2\right] \tag{6.47}$$

Now let's show these are equivalent. Substituting our score expression:

$$\mathcal{L}_{\text{score}} = \mathbb{E}\left[\left\|\mathbf{s}_\theta(\mathbf{x}_t, t) + \frac{\boldsymbol{\epsilon}}{\sqrt{1-\bar{\alpha}_t}}\right\|^2\right] \tag{6.48}$$

Using the relationship $\mathbf{s}_\theta(\mathbf{x}_t, t) = -\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{1-\bar{\alpha}_t}}$:

$$\mathcal{L}_{\text{score}} = \mathbb{E}\left[\left\|-\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{1-\bar{\alpha}_t}} + \frac{\boldsymbol{\epsilon}}{\sqrt{1-\bar{\alpha}_t}}\right\|^2\right] \tag{6.49}$$

$$= \frac{1}{1-\bar{\alpha}_t}\mathbb{E}\left[\left\|\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \boldsymbol{\epsilon}\right\|^2\right] \tag{6.50}$$

$$= \frac{1}{1-\bar{\alpha}_t}\mathcal{L}_{\text{diffusion}} \tag{6.51}$$

The two objectives are identical up to a time-dependent weighting factor! This means that training a diffusion model to predict noise is mathematically equivalent to training a score-based model to predict score functions.

### 6.6.4 Tweedie's Formula: The Statistical Foundation

The deep connection between score functions and denoising has roots that go back decades in statistics, embodied in a beautiful result called *Tweedie's formula*. This classical result provides the theoretical foundation for why learning to denoise is equivalent to learning score functions.

**Theorem 6.2** (Tweedie's Formula). *For a random variable with conditional distribution* $\mathbf{x} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, *the posterior mean given the observation* $\mathbf{x}$ *is:*

$$\mathbb{E}[\boldsymbol{\mu}|\mathbf{x}] = \mathbf{x} + \boldsymbol{\Sigma}\nabla_\mathbf{x}\log p(\mathbf{x}) \tag{6.52}$$

This theorem tells us something profound: **the best estimate of the clean signal given a noisy observation involves the score function**. In other words, optimal denoising requires knowledge of the probability landscape.

Applied to our diffusion setting where $q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I})$:

$$\mathbb{E}[\mathbf{x}_0|\mathbf{x}_t] = \frac{1}{\sqrt{\bar{\alpha}_t}}\left(\mathbf{x}_t + (1-\bar{\alpha}_t)\nabla_{\mathbf{x}_t}\log q(\mathbf{x}_t)\right) \tag{6.53}$$

This shows that **estimating the score function is equivalent to estimating the posterior mean**—exactly what denoising networks learn to do! Every time we train a network to remove noise from data, we're implicitly teaching it about the geometry of the probability distribution.

### 6.6.5 A Unified Framework: Three Perspectives, One Truth

The equivalence between score-based models and diffusion models reveals that we've been studying the same mathematical object from different angles. Each perspective offers unique insights and computational advantages:

**Score-based perspective:**

- **Focus**: Learn the score function $\nabla_{\mathbf{x}} \log p(\mathbf{x})$

- **Sampling**: Use Langevin dynamics to follow the probability landscape

- **Intuition**: Emphasizes the geometric structure—thinking of probability as a landscape with peaks and valleys

- **Advantage**: Direct connection to the underlying probability geometry

**Diffusion perspective:**

- **Focus**: Learn to reverse a gradual noising process

- **Sampling**: Use ancestral sampling to step-by-step remove noise

- **Intuition**: Emphasizes the temporal dynamics—thinking of generation as a process that unfolds over time

- **Advantage**: Natural way to think about multi-step generation and controllable sampling

**Continuous perspective (which we'll explore in the next chapter):**

- **Focus**: Describe the process using stochastic differential equations

- **Sampling**: Use ODE/SDE solvers for flexible sampling

- **Intuition**: Emphasizes the mathematical foundation—thinking in terms of continuous-time processes

- **Advantage**: Unified theoretical framework and flexible computational methods

This unification is more than just mathematical elegance—it has practical implications. Techniques developed for one framework often translate directly to the other. The score-based insight of annealed sampling led to better noise schedules for diffusion models. The diffusion model insight of parameterizing with noise prediction led to more stable training for score-based models.

As we continue our journey through generative modeling, this unified understanding will serve as our foundation for exploring even more sophisticated approaches to learning and sampling from complex probability distributions.

## 6.7 From Theory to Practice: Building Score-Based Models

### 6.7.1 Engineering the Score Function Network

Now that we understand the theoretical foundation, let's roll up our sleeves and see how to actually build these models in practice. Converting our mathematical insights into working code requires careful consideration of network architecture, training procedures, and implementation details.

**Designing the network architecture:**

At its core, our score network needs to accomplish a seemingly simple task: given a noisy data point $\mathbf{x}$ and a noise level $\sigma$, predict the score vector $\mathbf{s}_\theta(\mathbf{x}, \sigma) \in \mathbb{R}^d$. However, this simplicity is deceptive—the network must learn incredibly complex patterns to accurately estimate score functions across different noise levels.

The key architectural considerations include:

- **Input processing**: Handle both the data point $\mathbf{x}$ and the conditioning noise level $\sigma$ (or time $t$)

- **Output structure**: Produce a score vector with the same dimensionality as the input data

- **Architectural backbone**: For images, U-Net architectures (borrowed from diffusion models) work exceptionally well due to their multi-scale processing capabilities

**Noise level conditioning—telling the network what it's dealing with:**

One of the most crucial implementation details is how we inform the network about the current noise level. This conditioning must be strong enough that the network can dramatically change its behavior based on the noise level, essentially learning different "modes" of operation.

The standard approach follows these steps:

1. **Embed the noise level**: Use sinusoidal positional encoding to transform $\sigma$ into a rich embedding vector

2. **Inject broadly**: Add this embedding at multiple layers throughout the network, not just at the input

3. **Use adaptive conditioning**: Employ techniques like Adaptive Group Normalization (AdaGN) or Feature-wise Linear Modulation (FiLM) to let the noise embedding modulate the network's processing

**Output parameterization choices:**

Interestingly, we have several mathematically equivalent ways to parameterize what our network outputs:

- **Direct score prediction**: $\mathbf{s}_\theta(\mathbf{x}, \sigma)$ directly outputs the score vector

- **Noise prediction**: $\boldsymbol{\epsilon}_\theta(\mathbf{x}, \sigma)$ predicts the noise, which we convert to score using our equivalence relationship

- **Data prediction**: $\hat{\mathbf{x}}_\theta(\mathbf{x}, \sigma)$ predicts the clean data, from which we can compute the implicit score

In practice, noise prediction often works best because it's numerically more stable and the targets are better behaved across different noise levels.

### 6.7.2 The Training Recipe

Training a score-based model follows a surprisingly simple recipe, but each ingredient matters for the final result.

**procedure** SCOREMODELTRAINING(dataset $\{\mathbf{x}_i\}$, noise schedule $\{\sigma_j\}$)
    **repeat**
        Sample batch $\{\mathbf{x}_i\}$ from training dataset
        Sample noise levels $\{\sigma_j\}$ uniformly from schedule
        Sample noise vectors $\{\boldsymbol{\epsilon}_i\} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
        Create noisy data: $\mathbf{x}_{noisy} = \mathbf{x}_i + \sigma_j \boldsymbol{\epsilon}_i$
        Compute loss: $\mathcal{L} = \lambda(\sigma_j) \|\mathbf{s}_\theta(\mathbf{x}_{noisy}, \sigma_j) + \boldsymbol{\epsilon}_i / \sigma_j\|^2$
        Update parameters $\theta$ using gradient descent
    **until** convergence
**end procedure**

This procedure looks deceptively simple, but several details are crucial for success:

**Loss weighting strategy:** The function $\lambda(\sigma)$ determines how much we care about accurate score prediction at different noise levels. This choice significantly affects the final model quality:

- $\lambda(\sigma) = \sigma^2$: Emphasizes learning at high noise levels (good for global structure)

- $\lambda(\sigma) = 1$: Uniform weighting across all noise levels

- $\lambda(\sigma) = 1/\sigma^2$: Emphasizes learning at low noise levels (good for fine details)

The choice depends on what aspects of generation quality matter most for your application. Many successful models use a balanced approach that emphasizes both ends of the noise spectrum.

### 6.7.3 Sampling Strategies: From Score Functions to Samples

Once we have a trained score network, we need strategies to convert it into actual samples. Different sampling approaches offer different trade-offs between quality, speed, and computational requirements.

**1. Annealed Langevin Dynamics—the gold standard:**

This is our theoretically grounded approach that we've already explored. It provides excellent sample quality but requires many network evaluations:

- Use multiple noise levels with a carefully designed decreasing schedule

- Perform a fixed number of Langevin steps at each noise level

- Results in high-quality samples but slower generation

- Best for applications where quality matters more than speed

**2. Predictor-Corrector Sampling—balancing speed and quality:**
This hybrid approach tries to get the best of both worlds:

- **Predictor step**: Take a large step using the score function (fast progress)

- **Corrector step**: Refine the result with a few Langevin steps (quality improvement)

- Repeat this predictor-corrector cycle at each noise level

- Provides a good balance between sampling speed and quality

**3. Probability Flow ODE—deterministic and fast:**
This approach converts our stochastic sampling process into a deterministic one:

- Convert the score function into a velocity field for an ordinary differential equation

- Use sophisticated ODE solvers (adaptive step sizes, higher-order methods)

- Produces deterministic, reproducible samples

- Often much faster than stochastic sampling methods

# 6.8 The Strengths and Challenges of Score-Based Modeling

## 6.8.1 Why Score-Based Models Excel

After working with score-based models, several key advantages become apparent:
**1. Theoretical elegance meets practical power:**
Score-based models rest on a beautiful mathematical foundation that actually translates into practical benefits:

- Clean theoretical foundation based on well-understood score matching principles

- Direct connections to differential geometry and manifold learning provide geometric insights

- Principled approach that gives us confidence in the method's validity

- Deep connections to physics and statistical mechanics offer intuitive understanding

**2. Unprecedented sampling flexibility:**
Unlike many generative models that lock you into a single sampling procedure, score-based models offer remarkable flexibility:

- Multiple sampling algorithms (Langevin, predictor-corrector, ODE) for different needs

- Easy to adjust speed-quality trade-offs by changing the sampling procedure

- Natural incorporation of domain-specific constraints during sampling

- Elegant extensions to constrained or manifold-valued data

**3. Training that just works:**
One of the most practical advantages is the stability of the training process:

- No adversarial training dynamics—no need to balance competing networks

- Simple regression-based objective that's easy to optimize

- Much less sensitive to hyperparameters than GANs

- Graceful degradation with limited data—doesn't catastrophically fail

**4. Interpretability and debuggability:**
The geometric perspective makes these models much easier to understand and debug:

- Score function provides direct geometric intuition about what the model has learned

- Easy to visualize the probability landscape (at least in low dimensions)

- Clear connection between learning and sampling phases

- Facilitates analysis and debugging when things go wrong

## 6.8.2   The Current Limitations and Ongoing Challenges

Despite their many strengths, score-based models face several important limitations:

**1. Computational demands:**
The multi-step nature of sampling creates significant computational costs:

- Sampling requires many network evaluations (often hundreds or thousands)

- Much slower than single-step generators like GANs for generation

- Memory intensive for high-dimensional problems

- Difficult to parallelize the sampling process effectively

**2. Score estimation challenges:**
Learning accurate score functions is inherently difficult:

- Requires large, powerful networks to achieve good score estimation

- Performance depends heavily on careful noise schedule design

- Becomes increasingly challenging in very high dimensions

- May struggle with complex multimodal distributions

**3. Theory-practice gaps:**
While the theory is elegant, several practical challenges remain:

- Convergence analysis for finite samples and finite steps is complex

- Approximation errors can accumulate during the multi-step sampling process

- Difficult to provide tight sample complexity bounds

- Significant gaps between theoretical guarantees and practical performance

# 6.9 The Cutting Edge: Modern Developments and Extensions

## 6.9.1 Making Training More Efficient

Researchers have developed several techniques to improve training efficiency and sample quality:

**Variance reduction techniques:**
Training can be made more efficient by reducing the variance in gradient estimates:

- **Importance sampling**: Weight different noise levels based on their contribution to final sample quality

- **Adaptive noise scheduling**: Dynamically adjust the distribution of noise levels during training

- **Control variates**: Use known score functions (like Gaussian baselines) to reduce estimation variance

## 6.9.2 Conditional Generation and Control

Score-based models excel at conditional generation, where we want to generate samples that satisfy certain conditions:

**Class-conditional generation:**

$$\mathbf{s}_\theta(\mathbf{x}, \sigma, c) = \nabla_{\mathbf{x}} \log q(\mathbf{x}|\mathbf{x}_0, c) \tag{6.54}$$

This extends our score function to depend on class labels $c$, enabling controlled generation.

**Text-to-image models:**
The combination of score-based models with modern language understanding has produced remarkable results:

- Condition on rich text embeddings from models like CLIP or T5

- Use cross-attention mechanisms to align text and image features

- Employ classifier-free guidance to improve conditioning strength

**Inpainting and editing:**
Score-based models naturally support constrained generation:

- Condition on masked regions for inpainting tasks

- Enable iterative refinement of specific image areas

- Support controllable generation with precise spatial conditioning

### 6.9.3 Beyond Euclidean Spaces

One of the most exciting developments is extending score-based models beyond standard Euclidean data:
**Riemannian score-based models:**
For data that lives on curved manifolds (like spheres or hyperbolic spaces):

- Extend score functions using Riemannian geometry

- Applications to molecular conformations, 3D shape modeling

- Protein structure generation on configuration manifolds

**Graph score-based models:**
For discrete graph structures:

- Define score functions on graph adjacency matrices

- Use permutation-invariant architectures to handle graph symmetries

- Applications to molecular generation, social network modeling

## 6.10 Connections Across the Generative Modeling Landscape

### 6.10.1 The Flow Connection

Score-based models have deep connections to normalizing flows:
**Continuous normalizing flows (CNFs):**

- Use Neural ODEs to define invertible transformations

- Score-based models can be viewed as a special case through the probability flow ODE

- The connection provides new insights into both approaches

**Flow matching:**
A recent development that generalizes both approaches:

- Direct learning of flow fields rather than score functions

- Potentially more efficient training than traditional score matching

- Represents a promising direction for future research

### 6.10.2  Energy-Based Model Connections

The relationship to energy-based models provides another lens for understanding:
**Shared mathematical framework:**

- Both approaches work with energy functions $E(\mathbf{x})$

- Score function is the negative gradient of energy

- Different sampling mechanisms offer complementary advantages

**Hybrid approaches:**

- Combine score-based training with energy-based inference

- Leverage the strengths of both approaches

- Particularly powerful for physics-informed modeling

### 6.10.3  Architectural Innovations

Modern implementations explore various architectural improvements:
**Transformer-based score models:**

- Replace U-Net backbones with transformer architectures

- Better scalability to large-scale problems

- Applications extending beyond images to sequences and other modalities

**Latent score-based models:**

- Apply score-based modeling in learned latent spaces

- Combine with VAE or GAN encoders for efficiency

- Achieve improved computational efficiency without sacrificing quality

## 6.11  Real-World Impact and Applications

### 6.11.1  Transforming Computer Vision

Score-based models have achieved remarkable results across computer vision tasks:
**Image generation and manipulation:**

- High-resolution image synthesis that rivals or exceeds GAN quality

- Sophisticated style transfer and image manipulation capabilities

- Super-resolution and restoration that recover fine details

**Medical imaging breakthroughs:**

- MRI reconstruction from severely undersampled data

- CT denoising and artifact removal for improved diagnostics

- Synthetic medical data generation for training and research

### 6.11.2 Advancing Scientific Computing

The geometric perspective of score-based models makes them particularly well-suited for scientific applications:

**Physics simulations:**

- Accelerating molecular dynamics simulations

- Modeling complex fluid flows and turbulence

- Quantum state preparation and optimization

**Climate and weather modeling:**

- Generating realistic weather patterns for climate studies

- Downscaling global climate models to local predictions

- Uncertainty quantification in climate projections

### 6.11.3 Revolutionizing Drug Discovery

The ability to generate novel molecular structures has profound implications for drug discovery:

**Molecular design:**

- Generate novel drug candidates with desired properties

- Optimize molecular structures for specific targets

- Efficient sampling of molecular conformational spaces

**Protein engineering:**

- Predict protein structures from sequence information

- Analyze folding pathways and kinetics

- Design entirely novel proteins with specified functions

## 6.12 Looking Forward: The Future of Score-Based Modeling

As we reach the end of our exploration of score-based generative models, it's worth reflecting on both what we've learned and where the field is heading.

### 6.12.1 Key Contributions to Generative Modeling

Score-based models have contributed to generative modeling in several fundamental ways:

**1. A new theoretical lens:**

- Provided a rigorous mathematical foundation based on score matching and differential geometry

- Revealed deep connections between seemingly different approaches (diffusion, flows, energy models)

- Offered geometric intuition that makes complex probability distributions more interpretable

- Established clear theoretical understanding of when and why these methods work

**2. Algorithmic innovations:**

- Introduced Langevin dynamics as a principled sampling method for learned distributions

- Developed annealed sampling strategies that dramatically improve sample quality

- Created flexible frameworks that allow easy trade-offs between generation speed and quality

- Connected machine learning to advanced numerical methods from scientific computing

**3. Practical breakthroughs:**

- Achieved state-of-the-art results across numerous domains and applications

- Provided stable, reliable training procedures that avoid many pitfalls of adversarial training

- Created an extensible framework that adapts naturally to diverse applications

- Built bridges between machine learning and scientific computing communities

### 6.12.2 The Road Ahead

Current research in score-based modeling is pushing forward on multiple fronts:
**Efficiency and scale:**

- **Faster sampling**: New deterministic samplers, better ODE solvers, adaptive step sizing

- **Better architectures**: More efficient networks, transformer-based designs, specialized layers

- **Optimized training**: Improved noise schedules, variance reduction, curriculum learning

**Theoretical advances:**

- **Convergence analysis**: Tighter bounds on sampling convergence rates

- **Sample complexity**: Understanding how much data is needed for good score estimation

- **Approximation theory**: Characterizing when and how well neural networks can approximate score functions

**New frontiers:**

- **Scientific modeling**: Physics-informed score models, multi-scale phenomena

- **Engineering design**: Generative models for engineering optimization and design

- **Creative applications**: Music generation, 3D content creation, interactive design tools

**Integration and hybridization:**

- **Hybrid methods**: Combining score-based approaches with other generative paradigms

- **Foundation models**: Large-scale pretraining for diverse downstream tasks

- **Multi-modal models**: Unified frameworks for text, images, audio, and other modalities

### 6.12.3   The Lasting Impact

Perhaps the most important contribution of score-based models isn't any single technical innovation, but rather the *perspective* they've introduced to generative modeling. By focusing on the *geometry* of probability distributions rather than just their values, score-based models have shown us new ways to think about learning and generation.

This geometric view has implications that extend far beyond generative modeling:

- **Understanding representation learning**: How neural networks organize information geometrically

- **Optimization landscapes**: How to navigate complex loss surfaces more effectively

- **Transfer learning**: How knowledge transfers across different but geometrically related domains

- **Interpretability**: How to understand what models learn in terms of geometric structures

As we continue to push the boundaries of what's possible in machine learning and AI, the score-based perspective provides both solid mathematical foundations and a source of continuing inspiration. The journey from simple score functions to sophisticated generative models illustrates a fundamental principle: often the most powerful approaches come not from trying to solve problems directly, but from finding the right way to think about them.

The geometric view of probability, learning, and generation offers a rich landscape for future exploration and discovery. As we've seen throughout this exploration, understanding the *shape* of the probability landscape—through score functions, energy landscapes, and flow fields—provides both practical tools and conceptual frameworks that illuminate the deep structure of learning and generation.

The story of score-based generative models is still being written, and the next chapters promise to be even more exciting than what we've seen so far.

## 6.13 Exercises

1. **Score Function Computation:**

   (a) Compute the score function for a 2D Gaussian distribution $\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$

   (b) Visualize the score field and verify it points toward the mode

   (c) Implement Langevin dynamics and verify convergence to the target distribution

   (d) Experiment with different step sizes and analyze convergence

2. **Score Matching Implementation:**

   (a) Implement both explicit score matching and denoising score matching

   (b) Compare computational costs and training stability

   (c) Apply to a simple 2D mixture of Gaussians dataset

   (d) Analyze the quality of learned score functions

3. **Annealed Langevin Dynamics:**

   (a) Implement annealed Langevin dynamics with multiple noise levels

   (b) Compare with single-noise-level Langevin dynamics

   (c) Experiment with different noise schedules (linear, geometric, cosine)

   (d) Analyze the effect of the number of noise levels on sample quality

4. **Equivalence to Diffusion Models:**

(a) Train both a score-based model and a diffusion model on the same dataset

(b) Convert between score and noise predictions using the derived formulas

(c) Compare samples from both models

(d) Verify that the learned functions are equivalent up to scaling

5. **Tweedie's Formula Verification:**

   (a) Implement Tweedie's formula for Gaussian denoising

   (b) Show that score-based denoising gives the posterior mean

   (c) Compare with other denoising methods (Wiener filtering, neural networks)

   (d) Analyze performance as a function of noise level

6. **Multi-Scale Score Modeling:**

   (a) Implement a noise-conditional score network

   (b) Train on multiple noise levels simultaneously

   (c) Compare with single-noise training

   (d) Analyze how the learned scores vary with noise level

7. **Manifold Score-Based Models:**

   (a) Extend score-based modeling to data on a circle (1-sphere)

   (b) Implement the appropriate Riemannian score matching objective

   (c) Compare with Euclidean score matching applied to embedded data

   (d) Discuss the advantages of the geometric approach

8. **Conditional Score Models:**

   (a) Implement class-conditional score-based generation

   (b) Add classifier-free guidance to improve conditioning

   (c) Compare conditional and unconditional sample quality

   (d) Analyze the effect of guidance strength on diversity vs. fidelity

# Chapter 7

# Training Diffusion Models: ELBO, Losses, and Optimization

## 7.1 Introduction: From Theory to Practice

Having explored the theoretical foundations of diffusion models through continuous-time formulations (Chapter 5) and score-based perspectives (Chapter 6), we now turn to the practical challenge of training these models. This chapter bridges the gap between elegant mathematical theory and efficient computational implementation, focusing on how the Evidence Lower BOund (ELBO) leads to practical training objectives.

The journey from the ELBO to a simple mean squared error loss is one of the most elegant examples of how theoretical insights can dramatically simplify practical implementation. We'll see how the complex variational inference problem reduces to training a neural network to predict noise—a transformation that makes diffusion models both theoretically principled and computationally tractable.

**Chapter objectives:**

- **Connect theory to practice**: Show how ELBO derivation leads to practical loss functions

- **Understand parameterization choices**: Explore different ways to parameterize the reverse process

- **Master the training pipeline**: Provide a complete guide to implementing diffusion model training

- **Bridge mathematical frameworks**: Connect variational inference, score matching, and denoising perspectives

## 7.2 Revisiting Variational Inference and the ELBO

### 7.2.1 Variational Inference Refresher

The fundamental challenge in generative modeling is learning a complex probability distribution $p(\mathbf{x})$ from data. When we introduce latent variables $\mathbf{z}$, the

marginal likelihood becomes:

$$p(\mathbf{x}) = \int p(\mathbf{x}, \mathbf{z})d\mathbf{z} = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z} \tag{7.1}$$

**The intractability problem:**

- Direct computation of this integral is usually impossible

- The posterior $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{p(\mathbf{x})}$ is also intractable

- We need approximation methods to make progress

**Variational inference solution:** Instead of computing the exact posterior, we:

1. Choose a tractable family of distributions $\mathcal{Q}$

2. Find the best approximation $q(\mathbf{z}|\mathbf{x}) \in \mathcal{Q}$ to $p(\mathbf{z}|\mathbf{x})$

3. Use this approximation for inference and learning

## 7.2.2 Evidence Lower Bound (ELBO) Derivation

The ELBO provides a tractable objective that lower-bounds the intractable log-likelihood.

**Starting from the log-likelihood:**

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}, \mathbf{z})d\mathbf{z} \tag{7.2}$$

$$= \log \int q(\mathbf{z}|\mathbf{x})\frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})}d\mathbf{z} \tag{7.3}$$

$$= \log \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}\left[\frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})}\right] \tag{7.4}$$

**Applying Jensen's inequality:** Since log is concave, Jensen's inequality gives us:

$$\log \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}\left[\frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})}\right] \geq \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}\left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})}\right] \tag{7.5}$$

**The ELBO:**

$$\boxed{\mathcal{L}(\mathbf{x}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{x})}\left[\log \frac{p(\mathbf{x}, \mathbf{z})}{q(\mathbf{z}|\mathbf{x})}\right] \leq \log p(\mathbf{x})} \tag{7.6}$$

**The KL gap:** The difference between the true log-likelihood and the ELBO is exactly:

$$\log p(\mathbf{x}) - \mathcal{L}(\mathbf{x}) = D_{KL}(q(\mathbf{z}|\mathbf{x})\|p(\mathbf{z}|\mathbf{x})) \tag{7.7}$$

**Key insights:**

- The ELBO is always a lower bound on the log-likelihood

- Maximizing ELBO simultaneously improves likelihood and posterior approximation

- The bound is tight when $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$

- We can optimize the ELBO without knowing the true posterior

### 7.2.3  ELBO in the DDPM Setup

For diffusion models, the latent variables are the noisy versions $\mathbf{x}_{1:T} = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_T\}$.
   **Forward process (encoder):**

$$q(\mathbf{x}_{1:T}|\mathbf{x}_0) = \prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1}) \tag{7.8}$$

where $q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$.
   **Reverse process (decoder):**

$$p_\theta(\mathbf{x}_{0:T}) = p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) \tag{7.9}$$

where $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{x}_T; \mathbf{0}, \mathbf{I})$ and $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$ are learned.
   **ELBO for diffusion models:**

$$\mathcal{L}(\mathbf{x}_0) = \mathbb{E}_{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\left[\log \frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)}\right] \tag{7.10}$$

**Expanding the ratio:**

$$\frac{p_\theta(\mathbf{x}_{0:T})}{q(\mathbf{x}_{1:T}|\mathbf{x}_0)} = \frac{p(\mathbf{x}_T) \prod_{t=1}^{T} p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{\prod_{t=1}^{T} q(\mathbf{x}_t|\mathbf{x}_{t-1})} \tag{7.11}$$

$$= \frac{p(\mathbf{x}_T)}{q(\mathbf{x}_T|\mathbf{x}_{T-1})} \prod_{t=1}^{T} \frac{p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)}{q(\mathbf{x}_t|\mathbf{x}_{t-1})} \cdot q(\mathbf{x}_{T-1}|\mathbf{x}_{T-2}) \cdots q(\mathbf{x}_1|\mathbf{x}_0) \tag{7.12}$$

This leads to a telescoping product that can be reorganized into interpretable terms.

## 7.3  Deriving the DDPM Training Objective

### 7.3.1  ELBO Decomposition into KL Terms

Through careful algebraic manipulation (using the Markovian structure and properties of Gaussian distributions), the ELBO can be decomposed as:

$$\mathcal{L}(\mathbf{x}_0) = \underbrace{\mathbb{E}_{q(\mathbf{x}_1|\mathbf{x}_0)}[\log p_\theta(\mathbf{x}_0|\mathbf{x}_1)]}_{\text{Reconstruction term } L_0} \tag{7.13}$$

$$- \underbrace{D_{KL}(q(\mathbf{x}_T|\mathbf{x}_0)\|p(\mathbf{x}_T))}_{\text{Prior matching term } L_T} \tag{7.14}$$

$$- \underbrace{\sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}[D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)\|p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t))]}_{\text{Denoising terms } L_{t-1}} \tag{7.15}$$

**Understanding each term:**
**1. Reconstruction term $L_0$:**

- Measures how well we can reconstruct $\mathbf{x}_0$ from the slightly noisy $\mathbf{x}_1$

- Similar to the reconstruction term in VAEs

- In practice, often ignored as its contribution is small

**2. Prior matching term $L_T$:**

- Ensures $q(\mathbf{x}_T|\mathbf{x}_0)$ matches the prior $p(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I})$

- Fixed by design of the forward process (no learnable parameters)

- Becomes negligible as $T \to \infty$

**3. Denoising terms $L_{t-1}$:**

- The heart of diffusion model training

- Each term matches the learned reverse step to the true reverse step

- These are the terms we actually optimize

## 7.3.2 The True Reverse Posterior

The key insight is that the true reverse step $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ is tractable when conditioned on both $\mathbf{x}_t$ and $\mathbf{x}_0$.

**Using Bayes' rule:**

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \tag{7.16}$$

Since all three distributions are Gaussian, their product is also Gaussian:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0), \tilde{\sigma}_t^2 \mathbf{I}) \tag{7.17}$$

**The mean formula:**

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{\sqrt{\bar{\alpha}_{t-1}}\beta_t}{1 - \bar{\alpha}_t}\mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t}\mathbf{x}_t \tag{7.18}$$

**The variance formula:**

$$\tilde{\sigma}_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t}\beta_t \tag{7.19}$$

**Key insight:** The variance is completely determined by the noise schedule and doesn't depend on $\mathbf{x}_t$ or $\mathbf{x}_0$.

## 7.3.3 Parameterizing the Learned Reverse Process

Since the variance is fixed, we only need to learn the mean. We parameterize:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \tilde{\sigma}_t^2 \mathbf{I}) \tag{7.20}$$

**The KL divergence term becomes:**

$$D_{KL}(q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \| p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)) = \frac{1}{2\tilde{\sigma}_t^2}\|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t, t)\|^2 \tag{7.21}$$

**This gives us the training objective:**

$$L_{t-1} = \mathbb{E}_{q(\mathbf{x}_t|\mathbf{x}_0)}\left[\frac{1}{2\tilde{\sigma}_t^2}\|\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) - \boldsymbol{\mu}_\theta(\mathbf{x}_t, t)\|^2\right] \tag{7.22}$$

### 7.3.4 Simplifying to $L_{\text{simple}}$

The breakthrough insight of Ho et al. was to reparameterize this objective in terms of noise prediction.

**Using the forward process reparameterization:**

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} \tag{7.23}$$

where $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$.

**Reparameterizing the mean:** We can show that:

$$\tilde{\boldsymbol{\mu}}_t(\mathbf{x}_t, \mathbf{x}_0) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\boldsymbol{\epsilon}\right) \tag{7.24}$$

**Noise prediction parameterization:** Instead of predicting the mean directly, we predict the noise:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\right) \tag{7.25}$$

**The simplified objective:** Substituting this parameterization and removing constant factors:

$$\boxed{L_{\text{simple}} = \mathbb{E}_{t, \mathbf{x}_0, \boldsymbol{\epsilon}}\left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2\right]} \tag{7.26}$$

where $t \sim \text{Uniform}(1, T)$, $\mathbf{x}_0 \sim q(\mathbf{x}_0)$, $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$.

**Remarkable simplicity:** The complex variational inference problem reduces to training a neural network to predict noise!

### 7.3.5 Connection to Score Matching

The noise prediction objective has a deep connection to denoising score matching.

**Score function relationship:** From Chapter 6, we know that:

$$\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) = -\frac{\boldsymbol{\epsilon}}{\sqrt{1 - \bar{\alpha}_t}} \tag{7.27}$$

**Equivalence:** Learning to predict $\boldsymbol{\epsilon}$ is equivalent to learning the score function:

$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) = -\sqrt{1 - \bar{\alpha}_t}\nabla_{\mathbf{x}_t} \log q(\mathbf{x}_t|\mathbf{x}_0) \tag{7.28}$$

**Unified perspective:** The DDPM objective can be viewed as:

- **Variational inference**: Maximizing ELBO

- **Score matching**: Learning score functions via denoising

- **Denoising**: Predicting noise added to data

All three perspectives describe the same mathematical objective!

## 7.4  Parameterization Choices

### 7.4.1  Common Parameterizations

While noise prediction ($\epsilon$-parameterization) is most common, there are several alternatives, each with different properties.

    **1. Noise prediction ($\epsilon$-parameterization):**

$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \quad \text{predicts } \boldsymbol{\epsilon} \tag{7.29}$$

    **2. Data prediction ($x_\theta$-parameterization):**

$$\hat{\mathbf{x}}_{0,\theta}(\mathbf{x}_t, t) \quad \text{predicts } \mathbf{x}_0 \tag{7.30}$$

    **3. Velocity prediction (v-parameterization):**

$$\mathbf{v}_\theta(\mathbf{x}_t, t) \quad \text{predicts a scaled combination} \tag{7.31}$$

### 7.4.2  Mathematical Forms and Relationships

**Forward process constraint:** All parameterizations must be consistent with:

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} \tag{7.32}$$

**Conversion formulas:**
**From $\epsilon$ to $x_\theta$:**

$$\hat{\mathbf{x}}_{0,\theta}(\mathbf{x}_t, t) = \frac{1}{\sqrt{\bar{\alpha}_t}}\left(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\right) \tag{7.33}$$

**From $x_\theta$ to $\epsilon$:**

$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{1 - \bar{\alpha}_t}}\left(\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\hat{\mathbf{x}}_{0,\theta}(\mathbf{x}_t, t)\right) \tag{7.34}$$

**Velocity parameterization:**

$$\mathbf{v}_\theta(\mathbf{x}_t, t) = \sqrt{\bar{\alpha}_t}\boldsymbol{\epsilon} - \sqrt{1 - \bar{\alpha}_t}\mathbf{x}_0 \tag{7.35}$$

**Reverse process mean:** For any parameterization, the reverse mean is:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\right) \tag{7.36}$$

### 7.4.3  Empirical Comparisons and Trade-offs

$\epsilon$-parameterization (most common):

- **Pros**: Stable training, well-understood, consistent across noise levels

- **Cons**: May struggle with very low noise levels

- **Best for**: General-purpose diffusion models

$x_\theta$-**parameterization:**

- **Pros**: Direct prediction of final output, intuitive

- **Cons**: Can be unstable, especially at high noise levels

- **Best for**: Low-noise scenarios, image-to-image tasks

**v-parameterization:**

- **Pros**: Better numerical properties, balanced across noise levels

- **Cons**: Less intuitive, more complex implementation

- **Best for**: High-resolution images, challenging datasets

**Recent findings:**

- v-parameterization often performs better for high-resolution generation

- $\epsilon$-parameterization remains the standard for most applications

- Choice can significantly impact training dynamics and final quality

## 7.5   Practical Training Pipeline

### 7.5.1   Complete Training Algorithm

---

**Algorithm 3** DDPM Training

---

**Input:** Dataset $\mathcal{D}$, network $\boldsymbol{\epsilon}_\theta$, noise schedule $\{\beta_t\}_{t=1}^T$ **Precompute:** $\alpha_t = 1 - \beta_t$, $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ Sample batch $\{\mathbf{x}_0^{(i)}\}_{i=1}^B$ from $\mathcal{D}$ Sample timesteps $\{t^{(i)}\}_{i=1}^B$ from $\mathrm{Uniform}(1, T)$ Sample noise $\{\boldsymbol{\epsilon}^{(i)}\}_{i=1}^B$ from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ Compute noisy samples: $\mathbf{x}_t^{(i)} = \sqrt{\bar{\alpha}_{t^{(i)}}}\mathbf{x}_0^{(i)} + \sqrt{1 - \bar{\alpha}_{t^{(i)}}}\boldsymbol{\epsilon}^{(i)}$ Predict noise: $\hat{\boldsymbol{\epsilon}}^{(i)} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t^{(i)}, t^{(i)})$ Compute loss: $L = \frac{1}{B}\sum_{i=1}^B \|\boldsymbol{\epsilon}^{(i)} - \hat{\boldsymbol{\epsilon}}^{(i)}\|^2$ Update parameters: $\theta \leftarrow \theta - \eta \nabla_\theta L$ convergence

---

### 7.5.2   Sampling Time Steps

**Uniform Sampling**

The standard approach samples $t \sim \mathrm{Uniform}(1, T)$ for each training example.
**Why Uniform Sampling Works:**

- Provides equal training signal to all noise levels

- Encourages the network to handle the full range of corruption

- Simple and effective in practice

**Alternative Sampling Strategies**

- **Importance sampling**: Weight timesteps by their contribution to the loss

- **Curriculum learning**: Start with easier timesteps, gradually include harder ones

- **Adaptive sampling**: Adjust timestep distribution based on training progress

### 7.5.3   Generating Noisy Input

**Efficient vectorized implementation:**

$$\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} \tag{7.37}$$

**Implementation tips:**

- Precompute $\sqrt{\bar{\alpha}_t}$ and $\sqrt{1 - \bar{\alpha}_t}$ for all timesteps

- Use broadcasting for efficient batch computation

- Store noise schedule parameters on GPU to avoid CPU-GPU transfers

**Numerical considerations:**

- Ensure $\bar{\alpha}_t > 0$ for all timesteps to avoid numerical issues

- Use appropriate precision (float32 usually sufficient)

- Consider mixed precision training for efficiency

### 7.5.4   Computing the Loss

**Standard MSE loss:**
$$L = \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \tag{7.38}$$

**Loss weighting (optional):** Some implementations use weighted losses:

$$L = w_t\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)\|^2 \tag{7.39}$$

**Common weighting schemes:**

- $w_t = 1$ (unweighted, most common)

- $w_t = \frac{1}{1-\bar{\alpha}_t}$ (SNR weighting)

- $w_t = \sigma_t^2$ (variance weighting)

### 7.5.5 Forward and Backward Schedules

**Linear schedule (original DDPM):**

$$\beta_t = \beta_1 + \frac{t-1}{T-1}(\beta_T - \beta_1) \tag{7.40}$$

Typical values: $\beta_1 = 0.0001$, $\beta_T = 0.02$, $T = 1000$.

**Cosine schedule (improved DDPM):**

$$\bar{\alpha}_t = \frac{f(t)}{f(0)}, \quad f(t) = \cos\left(\frac{t/T + s}{1 + s} \cdot \frac{\pi}{2}\right)^2 \tag{7.41}$$

where $s$ is a small offset (typically $s = 0.008$).

**Schedule comparison:**

- **Linear**: Simple, well-understood, can be aggressive at high timesteps

- **Cosine**: More gradual noise injection, often better sample quality

- **Custom**: Can be learned or hand-tuned for specific datasets

**Impact on training:**

- Schedule affects the relative difficulty of different timesteps

- Cosine schedule often leads to better sample quality

- Choice interacts with parameterization and network architecture

## 7.6 Advanced Training Considerations

### 7.6.1 Network Architecture for Diffusion

**U-Net architecture (most common):**

- Encoder-decoder structure with skip connections

- Time embedding injected at multiple layers

- Attention mechanisms at appropriate resolutions

- Multi-scale feature processing

**Time conditioning:**

- Sinusoidal position encoding for timestep $t$

- Learned embedding layers

- Injection via AdaGN (Adaptive Group Normalization)

- Cross-attention for more complex conditioning

**Recent architectural advances:**

- Vision Transformers for diffusion (DiT)

- Efficient architectures (MobileDiffusion)

- Scale-specific optimizations

### 7.6.2 Training Stability and Optimization

**Gradient clipping:**

- Prevents exploding gradients during training

- Typical values: clip norm of 1.0-5.0

- Particularly important for high-resolution models

**Learning rate scheduling:**

- Warmup period to stabilize early training

- Cosine annealing or linear decay schedules

- Typical peak learning rates: $1e - 4$ to $2e - 4$

- AdamW optimizer most commonly used

**Batch size considerations:**

- Larger batches generally improve stability

- Typical range: 64-512 depending on model size and hardware

- Gradient accumulation for effective larger batches

- Batch normalization vs. Group normalization trade-offs

**Mixed precision training:**

- Use float16 for forward pass, float32 for gradients

- Significant speedup with minimal quality loss

- Requires careful handling of numerical stability

- Gradient scaling to prevent underflow

### 7.6.3 Monitoring Training Progress

**Loss curves:**

- Overall training loss should decrease steadily

- Can decompose loss by timestep to identify problematic regions

- Validation loss should track training loss

**Sample quality during training:**

- Generate samples periodically during training

- Use fixed seeds for reproducible evaluation

- Monitor both unconditional and conditional samples

- Track metrics like FID, IS if available

**Timestep analysis:**

- Plot loss as a function of timestep

- Identify if certain noise levels are problematic

- Adjust noise schedule if needed

**Network weights and gradients:**

- Monitor gradient norms to detect training issues

- Check for dead neurons or exploding activations

- Use tools like TensorBoard or Weights and Biases

# 7.7 Derivation of Posterior Mean

For completeness, we provide the detailed derivation of the posterior mean $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$.

## 7.7.1 Detailed Mathematical Derivation

**Starting from Bayes' rule:**

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q(\mathbf{x}_t|\mathbf{x}_{t-1})q(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q(\mathbf{x}_t|\mathbf{x}_0)} \tag{7.42}$$

**Known distributions:**

$$q(\mathbf{x}_t|\mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1}, (1-\alpha_t)\mathbf{I}) \tag{7.43}$$

$$q(\mathbf{x}_{t-1}|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1-\bar{\alpha}_{t-1})\mathbf{I}) \tag{7.44}$$

$$q(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1-\bar{\alpha}_t)\mathbf{I}) \tag{7.45}$$

**Working in log space:**

$$\log q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \tag{7.46}$$

$$\propto \log q(\mathbf{x}_t|\mathbf{x}_{t-1}) + \log q(\mathbf{x}_{t-1}|\mathbf{x}_0) - \log q(\mathbf{x}_t|\mathbf{x}_0) \tag{7.47}$$

$$\propto -\frac{1}{2(1-\alpha_t)}\|\mathbf{x}_t - \sqrt{\alpha_t}\mathbf{x}_{t-1}\|^2 \tag{7.48}$$

$$-\frac{1}{2(1-\bar{\alpha}_{t-1})}\|\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0\|^2 \tag{7.49}$$

$$+\frac{1}{2(1-\bar{\alpha}_t)}\|\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0\|^2 \tag{7.50}$$

**Expanding the quadratics:** After expanding and collecting terms quadratic in $\mathbf{x}_{t-1}$, we get:

$$\log q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \propto -\frac{1}{2\tilde{\sigma}_t^2}\|\mathbf{x}_{t-1} - \tilde{\boldsymbol{\mu}}_t\|^2 \tag{7.51}$$

**Coefficient of quadratic term:**

$$\frac{1}{\tilde{\sigma}_t^2} = \frac{\alpha_t}{1 - \alpha_t} + \frac{1}{1 - \bar{\alpha}_{t-1}} \tag{7.52}$$

Solving for $\tilde{\sigma}_t^2$:

$$\tilde{\sigma}_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \tag{7.53}$$

**Mean computation:** The mean is found by setting the derivative of the log-probability to zero:

$$\tilde{\boldsymbol{\mu}}_t = \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \mathbf{x}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \tag{7.54}$$

### 7.7.2 Alternative Forms of the Mean

**In terms of noise:** Using $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}$:

$$\tilde{\boldsymbol{\mu}}_t = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon} \right) \tag{7.55}$$

**In terms of predicted** $x_\theta$**:** If we predict $\hat{\mathbf{x}}_0 = \frac{1}{\sqrt{\bar{\alpha}_t}}(\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\theta)$:

$$\tilde{\boldsymbol{\mu}}_t = \frac{\sqrt{\bar{\alpha}_{t-1}} \beta_t}{1 - \bar{\alpha}_t} \hat{\mathbf{x}}_0 + \frac{\sqrt{\alpha_t}(1 - \bar{\alpha}_{t-1})}{1 - \bar{\alpha}_t} \mathbf{x}_t \tag{7.56}$$

# 7.8 Practical Implementation Tips

## 7.8.1 Code Organization

**Modular design:**

- Separate noise scheduling from model architecture

- Create reusable sampling and training loops

- Abstract parameterization choices for easy experimentation

**Efficient data handling:**

- Precompute noise schedule coefficients

- Use efficient data loaders with proper augmentation

- Consider mixed precision for memory efficiency

**Hyperparameter management:**

- Use configuration files for reproducible experiments

- Log all hyperparameters with training runs

- Implement learning rate warmup and scheduling

### 7.8.2 Debugging Common Issues

**Training instability:**

- Check gradient norms and clip if necessary

- Verify noise schedule implementation

- Ensure proper time embedding injection

- Monitor loss curves for anomalies

**Poor sample quality:**

- Verify correct parameterization implementation

- Check sampling algorithm correctness

- Ensure sufficient training iterations

- Validate network architecture choices

**Slow convergence:**

- Adjust learning rate and schedule

- Consider architectural improvements

- Check data preprocessing and normalization

- Verify batch size and optimization settings

## 7.9 Connection to Other Training Objectives

### 7.9.1 Weighted ELBO Variants

**Standard ELBO weighting:** Different weighting schemes can be derived from the full ELBO:

$$L_{\text{weighted}} = \mathbb{E}_t \left[ \lambda_t \| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \|^2 \right] \tag{7.57}$$

**SNR weighting:**

$$\lambda_t = \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t} \tag{7.58}$$

**Min-SNR weighting:**

$$\lambda_t = \min \left( \gamma, \frac{\bar{\alpha}_t}{1 - \bar{\alpha}_t} \right) \tag{7.59}$$

### 7.9.2   Progressive Training

**Curriculum learning approach:**

- Start training with easier timesteps (lower noise)

- Gradually include more challenging timesteps

- Can improve training stability and final quality

**Progressive distillation:**

- Train models to match multi-step trajectories in fewer steps

- Enables faster sampling without retraining from scratch

- Recent advance in diffusion model efficiency

## 7.10   Summary and Key Takeaways

This chapter has traced the elegant path from the complex variational inference problem to the simple and practical training objective used in modern diffusion models.

**Theoretical insights:**

- **ELBO foundation**: The ELBO provides the theoretical justification for the training objective

- **KL decomposition**: Breaking the ELBO into interpretable terms reveals the structure of the problem

- **Parameterization choice**: Different ways of parameterizing the reverse process lead to different training objectives

- **Unified framework**: Multiple perspectives (variational inference, score matching, denoising) describe the same objective

**Practical outcomes:**

- **Simple objective**: Complex theory reduces to training a noise prediction network

- **Stable training**: The objective is well-behaved and amenable to standard optimization

- **Flexible implementation**: Multiple parameterizations allow for different trade-offs

- **Scalable approach**: The method scales to high-dimensional problems and large datasets

**Key implementation points:**

- Use the simplified $L_{\text{simple}}$ objective for most applications

- Precompute noise schedule coefficients for efficiency

- Monitor training with both loss curves and sample quality

- Consider different parameterizations for different types of data

**Looking ahead:** The training framework established in this chapter provides the foundation for the advanced topics we'll explore next:

- **Conditional generation**: Extending the framework to conditional models

- **Accelerated sampling**: Using the trained models for fast sampling

- **Advanced architectures**: Scaling to larger models and datasets

- **Applications**: Applying diffusion models to diverse domains

The journey from ELBO to $L_{\text{simple}}$ represents one of the most successful examples of how theoretical insights can lead to practical breakthroughs. The mathematical rigor of variational inference provides confidence in the approach, while the simplicity of the final objective makes it practical and widely applicable.

## 7.11  Exercises

1. **ELBO Derivation Practice:**

   (a) Derive the ELBO for a simple 2-step diffusion process ($T = 2$)
   (b) Show how the KL terms arise from the telescoping product
   (c) Verify that the bound is tight when $q(\mathbf{z}|\mathbf{x}) = p(\mathbf{z}|\mathbf{x})$

2. **Posterior Mean Calculation:**

   (a) Derive the mean and variance of $q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$ from first principles
   (b) Verify the result using properties of Gaussian distributions
   (c) Show the equivalence between the $\mathbf{x}_0$ and $\boldsymbol{\epsilon}$ parameterizations

3. **Implementation Comparison:**

   (a) Implement the $\epsilon$-parameterization training loop
   (b) Implement the $\mathbf{x}_0$-parameterization and compare results
   (c) Test both on a simple 2D dataset and visualize the learned denoising

4. **Noise Schedule Analysis:**

   (a) Implement both linear and cosine noise schedules
   (b) Plot $\beta_t$, $\alpha_t$, and $\bar{\alpha}_t$ for both schedules
   (c) Train models with both schedules and compare sample quality

5. **Training Dynamics:**

   (a) Monitor loss as a function of timestep during training

(b) Implement weighted loss variants and compare convergence

(c) Analyze how different timesteps contribute to the overall training signal

6. **Score Matching Connection:**

(a) Show that the noise prediction network learns the score function

(b) Implement both the DDPM loss and the score matching loss

(c) Verify they give equivalent results on a simple dataset

7. **Parameterization Conversion:**

(a) Implement conversion functions between $\epsilon$, $\mathbf{x}_0$, and $v$ parameterizations

(b) Train a model with one parameterization and convert to others at test time

(c) Verify that all parameterizations produce identical samples

8. **Advanced Training Techniques:**

(a) Implement gradient clipping and analyze its effect on training stability

(b) Experiment with different learning rate schedules

(c) Compare the effect of different batch sizes on training dynamics

**Implementation Resources:**

- Hugging Face Diffusers Library: `github.com/huggingface/diffusers`

- OpenAI Guided Diffusion: `github.com/openai/guided-diffusion`

- PyTorch Lightning Implementation: `github.com/pytorch-lightning/pytorch-lightning`

The training of diffusion models represents a remarkable convergence of theoretical rigor and practical simplicity. The framework developed in this chapter provides the foundation for understanding not only how to train these models effectively, but also why the training procedure works so well across such a diverse range of applications.

# Chapter 8

# Sampling from Trained Diffusion Models: From Noise to Data

## 8.1 Introduction: The Art of Generation

Having mastered the training of diffusion models in Chapter 8, we now turn to the equally important challenge of *sampling*—the process of generating new data from our trained models. This is where the theoretical elegance and practical training procedures we've developed finally bear fruit in the form of synthetic samples.

Sampling from diffusion models is fundamentally different from other generative models. Unlike GANs, which generate samples in a single forward pass, or VAEs, which decode from a simple latent code, diffusion models require an iterative process that gradually transforms noise into data. This iterative nature presents both opportunities and challenges: it allows for fine-grained control over the generation process but can be computationally expensive.

**What sampling means in generative modeling:** Sampling is the process of drawing new data points from the probability distribution learned by our model. For diffusion models, this means starting with pure noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and applying the learned reverse process to gradually "denoise" our way to a realistic sample $\mathbf{x}_0$.

**The fundamental challenge:** We've trained a neural network $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ that can predict the noise added at any timestep. But how do we use this to actually generate samples? The answer lies in carefully implementing the reverse diffusion process, which comes in two main flavors:

- **Stochastic sampling (DDPM)**: Follows the original probabilistic reverse process with injected noise

- **Deterministic sampling (DDIM)**: Traces deterministic paths through the probability flow ODE

**Chapter roadmap:** We'll start with the standard DDPM sampling procedure and understand its stochastic nature. Then we'll explore DDIM's deterministic alternative, analyzing the trade-offs between quality, speed, and diversity. We'll conclude with practical implementation guidance and a preview of guided sampling techniques that enable controllable generation.

## 8.2 DDPM Sampling: Stochastic Reverse Denoising

### 8.2.1 Reverse Process Recap

From Chapter 4, we know that the reverse process is defined by learned transitions:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \tilde{\sigma}_t^2 \mathbf{I}) \tag{8.1}$$

**Mean parameterization:** The mean is parameterized using our trained noise prediction network:

$$\boldsymbol{\mu}_\theta(\mathbf{x}_t, t) = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \right) \tag{8.2}$$

**Variance schedule:** The variance follows the schedule derived from the forward process:

$$\tilde{\sigma}_t^2 = \frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t \tag{8.3}$$

**Key insight:** The mean depends on our learned network $\boldsymbol{\epsilon}_\theta$, while the variance is fixed by the noise schedule. This separation allows us to focus our learning on predicting the mean while maintaining the correct amount of stochasticity.

### 8.2.2 The DDPM Sampling Algorithm

The DDPM sampling procedure is remarkably straightforward given our trained model:

---
**Algorithm 4** DDPM Sampling

---
**Input:** Trained model $\boldsymbol{\epsilon}_\theta$, number of timesteps $T$ Sample $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ Start from pure noise $t = T, T-1, \ldots, 1$ Predict noise: $\hat{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ Compute mean: $\boldsymbol{\mu} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}} \hat{\boldsymbol{\epsilon}} \right)$ $t > 1$ Sample noise: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ $\mathbf{x}_{t-1} = \boldsymbol{\mu} + \sqrt{\frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t} \beta_t} \cdot \mathbf{z}$ $\mathbf{x}_0 = \boldsymbol{\mu}$ No noise in final step **Return:** $\mathbf{x}_0$

---

**Understanding each step:**

**1. Initialization:** We start with pure Gaussian noise $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. This represents the "maximally corrupted" state of our data.

**2. Noise prediction:** At each timestep, we use our trained network to predict what noise was originally added to create $\mathbf{x}_t$ from some $\mathbf{x}_0$.

**3. Mean computation:** Using the predicted noise, we compute the mean of the reverse distribution. This tells us the "expected" previous state $\mathbf{x}_{t-1}$.

**4. Stochastic sampling:** We don't just use the mean; we add calibrated noise to maintain the correct sampling distribution. This stochasticity is crucial for sample diversity.

**5. Final step:** At $t = 1$, we typically don't add noise, making the final step deterministic.

### 8.2.3 The Role of Stochasticity

The noise injection in DDPM sampling serves several important purposes:
**Theoretical correctness:**

- Ensures we're sampling from the correct probability distribution

- Maintains the balance between signal and noise at each timestep

- Guarantees that the sampling distribution matches the training distribution

**Sample diversity:**

- Different noise realizations lead to different sample trajectories

- Prevents mode collapse by exploring the full distribution

- Enables generation of diverse samples from the same starting point

**Error correction:**

- Small errors in noise prediction can compound over many steps

- Injected noise helps "reset" the trajectory and prevent error accumulation

- Acts as a form of regularization during sampling

### 8.2.4 Quality vs. Speed Considerations

DDPM sampling has both strengths and weaknesses:
**Strengths:**

- **High quality**: Produces excellent samples when using many timesteps

- **Diversity**: Stochastic sampling ensures good mode coverage

- **Theoretical guarantees**: Provably samples from the correct distribution

- **Robustness**: Works well across different datasets and architectures

**Weaknesses:**

- **Slow**: Requires many function evaluations (typically 1000)

- **Computational cost**: Each step requires a full neural network forward pass

- **Fixed schedule**: Difficult to adapt the sampling process for different quality/speed trade-offs

## 8.3 DDIM: Deterministic Sampling

### 8.3.1 The Problem with DDPM Sampling

DDPM sampling works beautifully in theory, but it has a critical practical limitation: it requires many timesteps to generate high-quality samples. Each step removes only a small amount of noise, meaning we need hundreds or thousands of function evaluations to go from pure noise to a clean image. This makes DDPM sampling prohibitively slow for many applications.

The natural question arises: *Is DDPM's particular reverse process the only way to denoise? Or can we find alternative paths that are faster while maintaining quality?*

### 8.3.2 DDIM's Key Insight: Breaking the Markovian Chain

The breakthrough insight of DDIM (Denoising Diffusion Implicit Models) comes from questioning a fundamental assumption in DDPM. DDPM assumes that the reverse process is **Markovian**—each denoising step depends only on the current noisy state:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t) \tag{8.4}$$

But what if we allow the reverse process to also look at the original clean image $\mathbf{x}_0$? This gives us the more general form:

$$q(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \tag{8.5}$$

This seemingly small change is actually profound—it **breaks the Markovian assumption** and opens up a much larger family of possible reverse processes.

**Why is this more powerful?** In DDPM's Markovian process, each step must be conservative because it only knows the current noisy state. But if we allow the process to "peek" at what the final clean image should look like, we can take much more direct paths through the denoising trajectory. It's like the difference between navigating in fog (only seeing your current position) versus having a map with your destination marked.

### 8.3.3 The DDIM Construction

Now comes the clever part: DDIM constructs a new forward process that maintains exactly the same marginal distributions as DDPM, but allows for this more flexible reverse process.

**The constraint we must satisfy:** To ensure our new process generates the same types of samples as DDPM, we require that the amount of noise at each timestep remains identical:

$$q_\sigma(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \tag{8.6}$$

This is the same as in DDPM—we're just changing *how* we get between timesteps, not the overall noise schedule.

**The new reverse transition:** With this constraint, we can derive (through marginalization and Bayes' rule) that our new reverse process has the form:

$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\sigma(\mathbf{x}_t, \mathbf{x}_0), \sigma_t^2\mathbf{I}) \tag{8.7}$$

where the mean is:

$$\boldsymbol{\mu}_\sigma(\mathbf{x}_t, \mathbf{x}_0) = \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}} \tag{8.8}$$

**The magic parameter $\sigma_t$:** Here's where DDIM becomes incredibly flexible. We can choose $\sigma_t$ freely (within constraints), allowing us to interpolate between completely different behaviors:

- $\sigma_t^2 = \frac{1-\bar{\alpha}_{t-1}}{1-\bar{\alpha}_t}\beta_t$: We recover exactly the original DDPM process

- $\sigma_t = 0$: We get a completely deterministic process

- Anything in between: We get a hybrid approach

### 8.3.4   Deriving the DDIM Reverse Transition

The derivation involves constructing a non-Markovian forward process and then finding its reverse. Here are the key steps:

**Step 1: Define the non-Markovian forward process**
We want to construct a forward process where:

$$q_\sigma(\mathbf{x}_{1:T}|\mathbf{x}_0) = q_\sigma(\mathbf{x}_T|\mathbf{x}_0) \prod_{t=2}^{T} q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) \tag{8.9}$$

Subject to the constraint that marginals remain the same as DDPM:

$$q_\sigma(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \tag{8.10}$$

**Step 2: Use Bayes' rule to find the reverse transition**
We want to find $q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0)$. Using Bayes' rule:

$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \frac{q_\sigma(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) \cdot q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_0)}{q_\sigma(\mathbf{x}_t|\mathbf{x}_0)} \tag{8.11}$$

**Step 3: We need to define $q_\sigma(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$**
This is where DDIM makes a key choice. We assume:

$$q_\sigma(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\alpha_t}\mathbf{x}_{t-1} + (1 - \alpha_t - \sigma_t^2)\frac{\mathbf{x}_0 - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}}{\sqrt{1 - \bar{\alpha}_{t-1}}}, \sigma_t^2\mathbf{I}) \tag{8.12}$$

**The problem we're trying to solve:** We want to construct a forward process $q_\sigma(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0)$ such that:

1. When we marginalize out $\mathbf{x}_{t-1}$, we get the same marginals as DDPM: $q_\sigma(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I})$

2. We can control the amount of randomness via parameter $\sigma_t$

3. When $\sigma_t^2 = \beta_t$, we recover the original DDPM process

**Building intuition: What should this transition look like?**

Think about what's happening: we're at $\mathbf{x}_{t-1}$ and we want to add noise to get to $\mathbf{x}_t$. But unlike DDPM, we're allowed to "peek" at the clean image $\mathbf{x}_0$.

**In DDPM (Markovian):**

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\boldsymbol{\epsilon} \tag{8.13}$$

This just adds noise to the previous state, with no knowledge of $\mathbf{x}_0$.

**In DDIM (Non-Markovian):** We want something more sophisticated that can use information about $\mathbf{x}_0$ to be "smarter" about how we add noise.

**The DDIM ansatz (educated guess):**

Let's break down the formula piece by piece:

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} + \sqrt{1 - \alpha_t - \sigma_t^2} \cdot \frac{\mathbf{x}_0 - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}}{\sqrt{1 - \bar{\alpha}_{t-1}}} + \sigma_t\boldsymbol{\epsilon} \tag{8.14}$$

**Term 1:** $\sqrt{\alpha_t}\mathbf{x}_{t-1}$ This is the "signal preservation" term, just like in DDPM. We keep most of the previous state.

**Term 2:** $\sqrt{1 - \alpha_t - \sigma_t^2} \cdot \frac{\mathbf{x}_0 - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}}{\sqrt{1-\bar{\alpha}_{t-1}}}$ This is the clever part! Let's understand what $\frac{\mathbf{x}_0 - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}}{\sqrt{1-\bar{\alpha}_{t-1}}}$ represents.

From the DDPM forward process, we know:

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1}}\boldsymbol{\epsilon}_{t-1} \tag{8.15}$$

Rearranging:

$$\boldsymbol{\epsilon}_{t-1} = \frac{\mathbf{x}_{t-1} - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_{t-1}}} \tag{8.16}$$

So $\frac{\mathbf{x}_0 - \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_{t-1}}{\sqrt{1-\bar{\alpha}_{t-1}}} = -\boldsymbol{\epsilon}_{t-1}$

This means Term 2 is:

$$-\sqrt{1 - \alpha_t - \sigma_t^2} \cdot \boldsymbol{\epsilon}_{t-1} \tag{8.17}$$

**Term 3:** $\sigma_t\boldsymbol{\epsilon}$ This adds fresh randomness controlled by $\sigma_t$.

**Putting it together:**

$$\mathbf{x}_t = \sqrt{\alpha_t}\mathbf{x}_{t-1} - \sqrt{1 - \alpha_t - \sigma_t^2} \cdot \boldsymbol{\epsilon}_{t-1} + \sigma_t\boldsymbol{\epsilon} \tag{8.18}$$

**The intuition:**

1. **Keep the signal**: $\sqrt{\alpha_t}\mathbf{x}_{t-1}$ preserves the previous state

2. **Adjust the existing noise**: $-\sqrt{1 - \alpha_t - \sigma_t^2} \cdot \boldsymbol{\epsilon}_{t-1}$ modifies the noise that was already in $\mathbf{x}_{t-1}$, using knowledge of $\mathbf{x}_0$ to figure out what that noise was

3. **Add controlled randomness**: $\sigma_t\boldsymbol{\epsilon}$ adds fresh noise we can control

**Why this works:**

- When $\sigma_t^2 = \beta_t$ and we do the math, this recovers exactly the DDPM transition

- When $\sigma_t = 0$, the process becomes deterministic

- The marginal $q_\sigma(\mathbf{x}_t|\mathbf{x}_0)$ works out to be exactly the same as DDPM (this requires checking, but it does!)

**The key insight:** Instead of just adding random noise like DDPM, DDIM says: "Let me figure out what noise is already in $\mathbf{x}_{t-1}$ (using $\mathbf{x}_0$), adjust it appropriately, and then add just the right amount of fresh randomness."

This is much more sophisticated than DDPM's blind noise addition, which is why DDIM can take much larger steps without losing quality.

**Step 4: Apply Bayes' rule with Gaussian distributions**

Since all three distributions in Bayes' rule are Gaussian, we can use the formula for the product of Gaussians. For Gaussians:

$$\mathcal{N}(x; \mu_1, \Sigma_1) \cdot \mathcal{N}(x; \mu_2, \Sigma_2) \propto \mathcal{N}(x; \mu_3, \Sigma_3) \tag{8.19}$$

where:

$$\Sigma_3^{-1} = \Sigma_1^{-1} + \Sigma_2^{-1} \tag{8.20}$$
$$\mu_3 = \Sigma_3(\Sigma_1^{-1}\mu_1 + \Sigma_2^{-1}\mu_2) \tag{8.21}$$

**Step 5: Compute the Gaussian product**

We have:

$$q_\sigma(\mathbf{x}_t|\mathbf{x}_{t-1}, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \mu_{t|t-1}, \sigma_t^2\mathbf{I}) \tag{8.22}$$
$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0, (1 - \bar{\alpha}_{t-1})\mathbf{I}) \tag{8.23}$$
$$q_\sigma(\mathbf{x}_t|\mathbf{x}_0) = \mathcal{N}(\mathbf{x}_t; \sqrt{\bar{\alpha}_t}\mathbf{x}_0, (1 - \bar{\alpha}_t)\mathbf{I}) \tag{8.24}$$

The numerator becomes a joint Gaussian in $(\mathbf{x}_{t-1}, \mathbf{x}_t)$, and we condition on $\mathbf{x}_t$ to get the desired posterior.

**Step 6: The final result**

After working through the Gaussian algebra (which involves several pages of calculations), we get:

$$q_\sigma(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\sigma(\mathbf{x}_t, \mathbf{x}_0), \sigma_t^2\mathbf{I}) \tag{8.25}$$

where:

$$\boldsymbol{\mu}_\sigma(\mathbf{x}_t, \mathbf{x}_0) = \sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2} \cdot \frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1 - \bar{\alpha}_t}} \tag{8.26}$$

**Intuitive breakdown of the mean formula:**

- $\sqrt{\bar{\alpha}_{t-1}}\mathbf{x}_0$: The "signal" component at timestep $t - 1$

- $\frac{\mathbf{x}_t - \sqrt{\bar{\alpha}_t}\mathbf{x}_0}{\sqrt{1-\bar{\alpha}_t}}$: The "noise direction" (residual after removing expected signal from $\mathbf{x}_t$)

- $\sqrt{1 - \bar{\alpha}_{t-1} - \sigma_t^2}$: The amount of this noise direction to add back

The formula essentially says: "Take the clean image $\mathbf{x}_0$, scale it appropriately for timestep $t-1$, then add back some of the noise in the direction suggested by the current state $\mathbf{x}_t$."

**Why this derivation matters:**

This construction is non-trivial because we're working backwards from desired properties (same marginals, flexible variance) to construct a valid stochastic process. The DDIM authors showed that such a process exists and gave us the explicit formula.

## 8.3.5 The Deterministic Revolution: $\sigma_t = 0$

The most important case is when $\sigma_t = 0$, making the process completely deterministic. In this case, the update rule becomes:

$$\mathbf{x}_{t-1} = \sqrt{\bar{\alpha}_{t-1}} \underbrace{\frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{\bar{\alpha}_t}}}_{\text{predicted } \mathbf{x}_0} + \sqrt{1 - \bar{\alpha}_{t-1}}\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \qquad (8.27)$$

**The intuition behind this formula:**

1. **Predict the clean image**: Use the neural network to predict what the original clean image $\mathbf{x}_0$ looks like

2. **Re-noise to the right level**: Take this predicted clean image and add exactly the right amount of noise for timestep $t-1$

3. **Use consistent noise**: The same predicted noise vector is used throughout, ensuring consistency

This is fundamentally different from DDPM's approach of gradually removing noise. Instead, DDIM makes a prediction about the final answer and then adjusts the noise level accordingly.

---

**Algorithm 5** DDIM Sampling ($\sigma_t = 0$)

---

**Input:** Trained model $\boldsymbol{\epsilon}_\theta$, timestep subset $\{\tau_1, \tau_2, \ldots, \tau_S\}$ Sample $\mathbf{x}_{\tau_S} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ $i = S, S-1, \ldots, 1$ $t = \tau_i$, $s = \tau_{i-1}$ (where $\tau_0 = 0$) Predict noise: $\hat{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$ Predict $\mathbf{x}_0$: $\hat{\mathbf{x}}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\hat{\boldsymbol{\epsilon}}}{\sqrt{\bar{\alpha}_t}}$ Update: $\mathbf{x}_s = \sqrt{\bar{\alpha}_s}\hat{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_s}\hat{\boldsymbol{\epsilon}}$ **Return:** $\mathbf{x}_0$

---

## 8.3.6 The Speed Revolution: Why DDIM Enables Acceleration

Now we can see why DDIM solves the speed problem. Because the process is deterministic and takes direct paths toward the final image, we can skip many intermediate timesteps without losing quality.

**Timestep subsampling:** Instead of using all timesteps $\{1, 2, \ldots, T\}$ like DDPM, we can use just a small subset $\{\tau_1, \tau_2, \ldots, \tau_S\}$ where $S \ll T$.

**Why this works:**

- **No accumulated randomness**: Deterministic steps don't accumulate stochastic errors

- **Consistent predictions**: The same noise prediction is used throughout the trajectory

- **Direct paths**: Instead of many small random steps, we take fewer large deterministic steps

- **ODE connection**: The deterministic process approximates a smooth ODE, which can be solved with fewer steps

**Typical acceleration results:**

- DDPM: 1000 steps for high quality

- DDIM: 50 steps for comparable quality

- Speedup: 20x faster sampling

- Quality: Maintained or sometimes improved

### 8.3.7   Fine-Tuning Stochasticity: The $\eta$ Parameter

In practice, we often parameterize the variance using an $\eta$ parameter:

$$\sigma_t = \eta \sqrt{\frac{1 - \bar{\alpha}_{t-1}}{1 - \bar{\alpha}_t} \beta_t} \tag{8.28}$$

This gives us a simple knob to control the trade-off:

- $\eta = 0$: Pure DDIM (deterministic, fastest)

- $\eta = 1$: Pure DDPM (stochastic, slower but more diverse)

- $0 < \eta < 1$: Hybrid approach (balanced speed and diversity)

**When to use each setting:**

- Use $\eta = 0$ when speed is critical and you want reproducible results

- Use $\eta = 1$ when you need maximum sample diversity

- Use intermediate $\eta$ values when you want to balance speed and diversity

The beauty of DDIM is that it gives us this flexibility—we can choose exactly the right point on the speed-diversity spectrum for our application.

| Aspect | DDPM | DDIM |
|---|---|---|
| Sample Quality (1000 steps) | Excellent | Excellent |
| Sample Quality (50 steps) | Poor | Good |
| Sample Diversity | High | Medium |
| Sampling Speed | Slow | Fast |
| Determinism | No | Yes (with $\eta = 0$) |
| Memory Usage | Standard | Standard |
| Implementation Complexity | Simple | Moderate |

Table 8.1: Comparison of DDPM and DDIM sampling characteristics

## 8.4 Comparing DDPM and DDIM: Trade-offs and Considerations

### 8.4.1 Quality vs. Speed Analysis

### 8.4.2 When to Use Each Method

**Use DDPM when:**

- Sample diversity is crucial

- Computational time is not a constraint

- You need the highest possible quality

- Working with smaller models or simpler datasets

- Theoretical guarantees are important

**Use DDIM when:**

- Speed is important (e.g., real-time applications)

- You need deterministic generation (e.g., for reproducibility)

- Working with high-resolution images or large models

- Interpolation or editing applications

- Quality with fewer steps is sufficient

### 8.4.3 Choosing Sampling Parameters

**Step 1: Choose sampling method**

- Start with DDIM for most applications

- Use DDPM only when maximum quality is needed

**Step 2: Select number of steps**

- DDIM: Start with 50 steps, reduce if speed is critical

- DDPM: Use full 1000 steps or as many as computationally feasible

- Monitor quality degradation as you reduce steps

**Step 3: Tune stochasticity**

- Start with $\eta = 0$ for maximum speed

- Increase $\eta$ if diversity is more important than speed

- Use $\eta = 0.1 - 0.3$ for a good balance

**Step 4: Optimize for your use case**

- Interactive applications: Prioritize speed (DDIM, fewer steps)

- Batch generation: Prioritize quality (more steps, possibly DDPM)

- Real-time applications: Consider distilled models

# 8.5 Advanced Sampling Techniques and Optimizations

## 8.5.1 Higher-Order Methods

While DDPM and DDIM use first-order updates, more sophisticated numerical methods can improve quality or speed:

**Heun's method:**

1. Predict noise: $\boldsymbol{\epsilon}_1 = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$

2. Take provisional step: $\mathbf{x}_{t-1}^* = f(\mathbf{x}_t, \boldsymbol{\epsilon}_1)$

3. Predict noise at new point: $\boldsymbol{\epsilon}_2 = \boldsymbol{\epsilon}_\theta(\mathbf{x}_{t-1}^*, t-1)$

4. Final step: $\mathbf{x}_{t-1} = f(\mathbf{x}_t, \frac{\boldsymbol{\epsilon}_1 + \boldsymbol{\epsilon}_2}{2})$

**Adaptive methods:**

- Estimate local error and adjust step size

- Use smaller steps in regions of high curvature

- Automatically balance quality and speed

### 8.5.2 Progressive Distillation

A recent advance for extreme acceleration is *progressive distillation*:

**Basic idea:** Train student models to match the output of multi-step teacher models in fewer steps.

**Process:**

1. Start with a trained diffusion model (teacher)

2. Train a student to match 2-step teacher outputs in 1 step

3. Use the student as a new teacher for 4→2 step distillation

4. Continue until reaching the desired speed

**Benefits:**

- Can achieve high quality with very few steps (4-8)

- Maintains the benefits of diffusion training

- Provides a clear path to real-time generation

# 8.6 Practical Implementation Guide

## 8.6.1 Common Implementation Pitfalls

| Issue | Solution |
|---|---|
| Incorrect noise scheduling | Ensure $\alpha$ and $\beta$ values are computed correctly; verify $\bar{\alpha}_t$ is monotonically decreasing |
| Numerical stability | Handle edge cases where $\bar{\alpha}_t \to 0$ or $\bar{\alpha}_t \to 1$; use appropriate epsilon values |
| Timestep indexing | Be consistent with 0-indexed vs 1-indexed conventions between training and sampling |
| DDIM step selection | Start with uniform subsampling; experiment with other strategies if needed |
| Memory issues | Use gradient checkpointing, mixed precision, or model quantization |

Table 8.2: Common implementation issues and solutions

## 8.6.2 Debugging and Validation

**Sanity checks:**

- Compare DDIM with $\eta = 1$ to DDPM results

- Check that deterministic sampling produces identical results with same seed

- Verify intermediate denoising steps look reasonable

**Quality assessment:**

- Visual inspection of generated samples

- Quantitative metrics (FID, IS) if available

- Diversity analysis across multiple samples

- Comparison with baseline implementations

**Performance profiling:**

- Measure sampling time per step and total

- Identify computational bottlenecks

- Monitor memory usage during generation

- Compare with theoretical expectations

## 8.6.3 Memory and Computational Optimizations

**Mixed precision sampling:**

- Use float16 for forward passes, float32 for sensitive operations

- Reduces memory usage and increases speed

- Requires careful handling of numerical stability

**Model quantization:**

- Reduce model precision (8-bit, 4-bit) for inference

- Dramatic speedups with minimal quality loss

- Particularly effective for deployment scenarios

**Caching and precomputation:**

- Precompute noise schedule coefficients

- Cache intermediate activations when possible

- Use efficient tensor operations and broadcasting

## 8.7 Preview: Guided Sampling

While this chapter has focused on unconditional generation—producing samples from the learned data distribution without specific control—many applications require *conditional* generation. We might want to generate images from text descriptions, create specific classes of objects, or control artistic style. This need for controllable generation leads us to guided sampling methods.

The most successful approach is *Classifier-Free Guidance* (CFG), which modifies the sampling process to incorporate conditioning information. During sampling, CFG uses a linear combination of conditional and unconditional noise predictions:

$$\boxed{\tilde{\boldsymbol{\epsilon}}(\mathbf{x}_t, y, t) = (1 + \omega)\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, y, t) - \omega\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t)} \tag{8.29}$$

where $y$ is the conditioning information, $\emptyset$ represents unconditional generation, and $\omega \geq 0$ controls the guidance strength.

Key insights about CFG:

- Eliminates the need for separate classifier training by jointly learning conditional and unconditional models

- The guidance scale $\omega$ controls the trade-off between prompt adherence and sample diversity

- Works by extrapolating in the direction from unconditional to conditional predictions

The full treatment of conditional generation, including CFG training procedures, guidance scale scheduling, and advanced conditioning techniques, will be covered in Chapter X on Conditional Generation.

## 8.8 Summary and Key Takeaways

This chapter has explored the crucial transition from trained diffusion models to actual sample generation, revealing how theoretical frameworks translate into practical algorithms.

**Core concepts mastered:**

- **DDPM sampling**: Stochastic reverse process with theoretical guarantees

- **DDIM sampling**: Deterministic alternative enabling dramatic acceleration

- **Trade-offs**: Understanding quality, speed, and diversity relationships

- **Implementation**: Practical considerations for robust sampling

**Practical insights:**

- **Method selection**: DDIM for speed, DDPM for maximum quality

- **Step reduction**: Fewer steps possible with deterministic sampling

- **Parameter tuning**: $\eta$ parameter controls stochasticity-speed trade-off

- **Optimization**: Numerical stability and performance considerations

**Key algorithmic contributions:**

- DDIM's insight that multiple reverse processes can share the same marginals

- The $\eta$ parameter for interpolating between stochastic and deterministic sampling

- The connection between deterministic sampling and ODE integration

- Progressive distillation for extreme acceleration

**Looking ahead:** The sampling techniques covered in this chapter provide the foundation for conditional generation, image editing, and advanced applications. The ability to efficiently and reliably generate high-quality samples transforms diffusion models from theoretical constructs into powerful practical tools.

The evolution from DDPM's careful stochastic sampling to DDIM's deterministic acceleration illustrates how the field rapidly developed practical solutions to theoretical insights. As we move forward to explore conditional generation in the next chapter, these sampling foundations will prove essential for understanding how to control and guide the generation process.

# 8.9 Exercises

1. **DDPM vs DDIM Implementation:**

   (a) Implement both DDPM and DDIM sampling for a trained model
   (b) Compare sample quality using the same random seeds
   (c) Measure wall-clock time for different numbers of steps
   (d) Analyze the trade-off between speed and quality

2. **Stochasticity Analysis:**

   (a) Implement DDIM with controllable $\eta$ parameter
   (b) Generate samples with $\eta \in \{0, 0.25, 0.5, 0.75, 1.0\}$
   (c) Analyze how stochasticity affects sample diversity
   (d) Find the optimal $\eta$ for your dataset and quality metric

3. **Step Reduction Study:**

   (a) Test DDIM sampling with $\{10, 20, 50, 100, 250, 1000\}$ steps
   (b) Compare quality using both visual inspection and quantitative metrics
   (c) Determine the minimum number of steps for acceptable quality
   (d) Experiment with different timestep scheduling strategies

4. **Numerical Stability Analysis:**

   (a) Identify potential numerical issues in the sampling algorithms

(b) Implement safeguards for edge cases (e.g., $\bar{\alpha}_t \to 0$)

(c) Test sampling with extreme noise schedules

(d) Develop robust implementations that handle numerical edge cases

5. **Advanced Sampling Methods:**

   (a) Implement a higher-order sampling method (e.g., Heun's method)

   (b) Compare convergence properties with first-order methods

   (c) Analyze computational cost vs. quality trade-offs

   (d) Experiment with adaptive step size selection

6. **Performance Optimization:**

   (a) Profile your sampling implementation to identify bottlenecks

   (b) Implement optimizations (mixed precision, caching, vectorization)

   (c) Measure speedup and quality impact of optimizations

   (d) Compare your implementation with reference implementations

7. **Timestep Scheduling Comparison:**

   (a) Implement uniform, quadratic, and cosine timestep schedules

   (b) Compare their performance on your dataset

   (c) Analyze how different schedules affect different noise levels

   (d) Develop intuition for when each schedule works best

8. **Deterministic vs. Stochastic Analysis:**

   (a) Implement both deterministic and stochastic sampling with various $\eta$ values

   (b) Analyze the reproducibility and consistency of deterministic sampling

   (c) Study how stochasticity affects interpolation and editing applications

   (d) Compare the diversity-quality trade-offs quantitatively

9. **Progressive Distillation Implementation:**

   (a) Implement a simple progressive distillation setup

   (b) Train a 2-step student to match a 4-step teacher model

   (c) Evaluate the quality degradation vs. speed improvement

   (d) Experiment with different distillation objectives and schedules

**Practical Resources:**

- Hugging Face Diffusers Library: Production-ready implementations of DDPM and DDIM sampling with extensive optimization

- CompVis/stable-diffusion: Reference implementations with practical optimizations and deployment considerations

- OpenAI/guided-diffusion: Research code for classifier guidance and improved sampling techniques

- NVIDIA/DeepLearningExamples: Highly optimized implementations for GPU deployment and scaling

- Google Research/ddpm: Original implementation with detailed documentation and experimental setups

**Implementation Notes:**

- Most modern implementations default to DDIM with 50 steps for a good speed-quality balance

- Consider starting with established libraries before implementing from scratch

- Pay careful attention to numerical precision, especially for long sampling chains

- Profile your implementation to identify bottlenecks specific to your hardware setup

The sampling phase of diffusion models represents where theory meets practice in the most direct way. The algorithms and techniques covered in this chapter form the bridge between the mathematical elegance of diffusion processes and their practical application in generating the remarkable synthetic content that has captured public imagination and transformed creative workflows across industries.

Understanding these sampling fundamentals is essential for anyone working with diffusion models, whether for research, development, or deployment in real-world applications. The iterative nature of diffusion sampling, while computationally demanding, provides unprecedented control over the generation process and enables the fine-grained manipulation that makes these models so powerful for creative and practical applications.

As the field continues to evolve, with new acceleration techniques, improved numerical methods, and novel conditioning approaches, the principles established in this chapter will remain the foundation upon which future innovations are built. The careful balance between theoretical rigor and practical efficiency exemplified in the evolution from DDPM to DDIM serves as a model for how research advances translate into real-world impact.

# Chapter 9

# Conditional Generation in Diffusion Models

## 9.1 Introduction: From Unconditional to Controllable Generation

In the previous chapter, we mastered the art of sampling from trained diffusion models, learning how to transform noise into realistic data through careful implementation of reverse diffusion processes. However, unconditional generation—while theoretically elegant—has limited practical utility. Real-world applications demand *controllable* generation: the ability to specify what we want to create.

Consider the difference between asking a model to "generate an image" versus "generate a photo of a golden retriever playing in a park on a sunny day." The latter represents conditional generation, where we provide additional information to guide the sampling process toward desired outcomes.

**Motivation for conditional generation:**

- **Specificity**: Generate particular categories, styles, or attributes

- **User control**: Enable interactive and directed creative workflows

- **Practical utility**: Make generative models useful for real applications

- **Sample relevance**: Reduce the need to generate many samples to find desired outputs

**Broad conditioning strategies:** We can incorporate conditioning information at two different stages:

- **Training time**: Learn conditional distributions $p(\mathbf{x}|\mathbf{y})$ directly

- **Sampling time**: Guide unconditional models using external information

**Chapter roadmap:** We'll begin with class-conditional diffusion models that learn to generate specific categories during training. Then we'll explore classifier guidance, which uses separate classifiers to steer generation at sampling time. We'll dive deep into classifier-free guidance (CFG), the most successful approach

that elegantly combines conditional and unconditional learning. Finally, we'll examine advanced conditioning mechanisms for text, images, and multimodal control.

## 9.2 Conditional Diffusion Models

### 9.2.1 Class-Conditional Generation

The most straightforward approach to conditional generation is to directly model the conditional distribution $p(\mathbf{x}|\mathbf{y})$, where $\mathbf{y}$ represents conditioning information such as class labels, attributes, or other structured data.

**Mathematical formulation:** Instead of learning the unconditional reverse process $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$, we learn:

$$p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t, \mathbf{y}) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, \mathbf{y}, t), \sigma_t^2 \mathbf{I}) \tag{9.1}$$

The forward process remains unchanged, but our neural network now predicts noise conditioned on both the current state and the conditioning information:

$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t) \tag{9.2}$$

**Training objective:** The denoising objective becomes:

$$\mathcal{L} = \mathbb{E}_{t, \mathbf{x}_0, \mathbf{y}, \boldsymbol{\epsilon}} \left[ \|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)\|^2 \right] \tag{9.3}$$

where $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$ and $(\mathbf{x}_0, \mathbf{y})$ are sampled from the training data.

### 9.2.2 Conditioning Mechanisms

The key challenge in conditional diffusion models is how to effectively incorporate the conditioning information $\mathbf{y}$ into the neural network architecture. Several strategies have proven effective:

**1. Concatenation approaches:**

- **Input concatenation**: Add extra channels to the input, concatenating encoded $\mathbf{y}$ with $\mathbf{x}_t$

- **Embedding concatenation**: Concatenate condition embeddings with timestep embeddings

- **Feature concatenation**: Combine conditioning features at intermediate network layers

**2. Attention-based approaches:**

- **Cross-attention**: Use conditioning information as keys and values in attention layers

- **Self-attention with conditioning**: Include conditioning tokens in self-attention computations

**3. Normalization-based approaches:**

- **Conditional BatchNorm**: Scale and shift normalization parameters based on $\mathbf{y}$

- **Adaptive GroupNorm**: Modulate group normalization using condition-dependent parameters

- **FiLM**: Feature-wise Linear Modulation using conditioning information

## 9.2.3 Architecture Modifications

For class-conditional generation, the most common approach is to modify the U-Net architecture used in unconditional models:

**Conditional U-Net design:**

1. **Embedding layer**: Map class labels to dense embeddings

2. **Timestep fusion**: Combine class embeddings with time embeddings

3. **Conditioning injection**: Add conditioning information at multiple scales throughout the U-Net

---

**Algorithm 6** Conditional U-Net Forward Pass

**Input:** Noisy image $\mathbf{x}_t$, condition $\mathbf{y}$, timestep $t$ Embed timestep: $\mathbf{e}_t = \text{TimeEmbedding}(t)$ Embed condition: $\mathbf{e}_y = \text{ClassEmbedding}(\mathbf{y})$ Combine embeddings: $\mathbf{e} = \mathbf{e}_t + \mathbf{e}_y$ Pass through conditional U-Net: $\boldsymbol{\epsilon} = \text{UNet}(\mathbf{x}_t, \mathbf{e})$ **Return:** Predicted noise $\boldsymbol{\epsilon}$

---

**Implementation details:**

- Class embeddings are typically learned during training

- Embedding dimensions should match timestep embedding dimensions for easy addition

- Conditioning can be injected at multiple resolution levels in the U-Net

- Dropout can be applied to conditioning to improve robustness

## 9.2.4 Training with Labels

Training conditional diffusion models requires paired data $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}_{i=1}^N$ where each data point has an associated label or conditioning information.

---

**Algorithm 7** Conditional Diffusion Training

**Input:** Dataset $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}$, model $\boldsymbol{\epsilon}_\theta$ not converged Sample batch $\{(\mathbf{x}, \mathbf{y})\}$ from dataset Sample timesteps $t \sim \text{Uniform}(1, T)$ Sample noise $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ Compute noisy samples: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x} + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$ Predict noise: $\hat{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)$ Compute loss: $\mathcal{L} = \|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}\|^2$ Update parameters: $\theta \leftarrow \theta - \eta\nabla_\theta\mathcal{L}$

---

**Key considerations:**

- Ensure balanced representation of all classes in training batches

- Consider using label smoothing or mixup techniques for robustness

- Monitor conditional generation quality across different classes

- Validate that the model can ignore conditioning when needed

# 9.3   Classifier Guidance: The Art of Steering Generation

## 9.3.1   The Big Picture: Why Do We Need Guidance?

Imagine you've just learned to drive, but your car only goes in random directions. You know how to operate the vehicle (our trained diffusion model), but you have no steering wheel! That's exactly the situation we face with unconditional diffusion models—they can generate beautiful, realistic samples, but we have no control over *what* they generate.

This is where **classifier guidance** comes to the rescue. It's like adding a steering wheel to our generative car, allowing us to navigate toward specific destinations (desired classes, attributes, or conditions) rather than just wandering randomly through the space of possible images.

**The core insight:** Instead of training an entirely new conditional diffusion model from scratch (which would be like building a completely new car), we can take our existing unconditional model and add a guidance system that steers it during sampling.

## 9.3.2   The Mathematical Foundation: Bayes to the Rescue

Let's build up the mathematics step by step, starting from a fundamental question: *How do we sample from a conditional distribution using tools designed for unconditional sampling?*

**Step 1: The Conditional Score Decomposition**

Suppose we want to sample from $p(\mathbf{x}_t|\mathbf{y})$—the distribution of noisy images at timestep $t$ given that they belong to class $\mathbf{y}$. The key insight comes from Bayes' rule applied to probability densities.

Starting with Bayes' rule for the conditional probability:

$$p(\mathbf{x}_t|\mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{x}_t)p(\mathbf{x}_t)}{p(\mathbf{y})} \tag{9.4}$$

Now, here's the crucial step: let's take the logarithm of both sides and then compute the gradient with respect to $\mathbf{x}_t$:

$$\log p(\mathbf{x}_t|\mathbf{y}) = \log p(\mathbf{y}|\mathbf{x}_t) + \log p(\mathbf{x}_t) - \log p(\mathbf{y}) \tag{9.5}$$

Taking the gradient with respect to $\mathbf{x}_t$:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) - \nabla_{\mathbf{x}_t} \log p(\mathbf{y}) \tag{9.6}$$

Notice that $\log p(\mathbf{y})$ doesn't depend on $\mathbf{x}_t$, so its gradient is zero:

$$\boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)} \tag{9.7}$$

**What does this equation tell us?** It says that the conditional score (the direction we should move to increase the likelihood of generating class $\mathbf{y}$) can be decomposed into two parts:

- $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$: The *unconditional score*—where to go to generate realistic data in general

- $\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)$: The *classification score*—where to go to make the current sample more likely to be classified as class $\mathbf{y}$

This is beautiful! We already know how to approximate the first term (that's what our trained diffusion model does), and the second term is just the gradient of a classifier's log-probability.

### Step 2: Connecting Scores to Noise Predictions

Now we need to translate this insight into something we can use with our DDPM/DDIM sampling procedures. Recall that our diffusion models don't directly predict scores—they predict noise. But there's a beautiful relationship between these two quantities.

From the training of diffusion models, we know that:

$$\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) \approx \boldsymbol{\epsilon} \tag{9.8}$$

where $\boldsymbol{\epsilon}$ is the noise that was added to create $\mathbf{x}_t$.

The connection to scores comes from the reparameterization of the forward process. Since $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$, we can show that:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) = -\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}} \tag{9.9}$$

**Intuitive explanation:** The score points in the direction of increasing probability density. Since $\boldsymbol{\epsilon}_\theta$ predicts the noise that was added, $-\boldsymbol{\epsilon}_\theta$ points toward the clean data. Dividing by $\sqrt{1 - \bar{\alpha}_t}$ adjusts for the noise level at timestep $t$.

### Step 3: The Guided Noise Prediction

Now we can substitute our score decomposition into the noise prediction framework. If we want to sample from the conditional distribution $p(\mathbf{x}_t|\mathbf{y})$, we need to use the conditional score:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = -\frac{\tilde{\boldsymbol{\epsilon}}(\mathbf{x}_t, \mathbf{y}, t)}{\sqrt{1 - \bar{\alpha}_t}} \tag{9.10}$$

Substituting our decomposition:

$$-\frac{\tilde{\boldsymbol{\epsilon}}(\mathbf{x}_t, \mathbf{y}, t)}{\sqrt{1 - \bar{\alpha}_t}} = -\frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)}{\sqrt{1 - \bar{\alpha}_t}} + \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y}|\mathbf{x}_t) \tag{9.11}$$

Solving for $\tilde{\boldsymbol{\epsilon}}$:

$$\boxed{\tilde{\boldsymbol{\epsilon}}(\mathbf{x}_t, \mathbf{y}, t) = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \sqrt{1 - \bar{\alpha}_t}\nabla_{\mathbf{x}_t}\log p_\phi(\mathbf{y}|\mathbf{x}_t)} \tag{9.12}$$

**This is the heart of classifier guidance!** Let's unpack what this equation means:

- We start with the unconditional noise prediction $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t)$

- We subtract a term proportional to the classifier gradient $\nabla_{\mathbf{x}_t}\log p_\phi(\mathbf{y}|\mathbf{x}_t)$

- The classifier gradient points toward regions where the classifier is more confident about class $\mathbf{y}$

- By subtracting this (remember, we want to subtract noise, not add it), we steer the sampling toward the desired class

**Step 4: Adding the Guidance Scale**

In practice, we often want to control how strongly we apply the guidance. This leads us to introduce a guidance scale $\omega$:

$$\boxed{\tilde{\boldsymbol{\epsilon}}(\mathbf{x}_t, \mathbf{y}, t) = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \omega\sqrt{1 - \bar{\alpha}_t}\nabla_{\mathbf{x}_t}\log p_\phi(\mathbf{y}|\mathbf{x}_t)} \tag{9.13}$$

**Why do we need this knob?** Think of $\omega$ as the gain on your steering wheel:

- $\omega = 0$: No steering—pure unconditional sampling

- $\omega = 1$: Theoretically correct steering according to Bayes' rule

- $\omega > 1$: Over-steering—stronger preference for the condition, potentially at the cost of realism

## 9.3.3 The Complete Sampling Algorithm: Putting It All Together

Now that we understand the mathematics, let's see how this works in practice. The beauty of classifier guidance is that it seamlessly integrates with existing sampling algorithms like DDIM.

---

**Algorithm 8** DDIM Sampling with Classifier Guidance

---

1: **Input:** Unconditional model $\epsilon_\theta$, classifier $p_\phi$, target class $\mathbf{y}$, guidance scale $\omega$
2: **Input:** Timestep schedule $\{\tau_1, \tau_2, \ldots, \tau_S\}$ where $\tau_S = T$ and $\tau_0 = 0$
3: Sample $\mathbf{x}_{\tau_S} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$         ▷ Start from pure noise
4: **for** $i = S, S-1, \ldots, 1$ **do**
5:      $t = \tau_i$, $s = \tau_{i-1}$         ▷ Current and next timesteps
6:      // Step 1: Get unconditional noise prediction
7:      $\epsilon_u = \epsilon_\theta(\mathbf{x}_t, t)$
8:      // Step 2: Compute classifier guidance
9:      $\mathbf{g} = \nabla_{\mathbf{x}_t} \log p_\phi(\mathbf{y}|\mathbf{x}_t)$         ▷ Requires backprop through classifier
10:     // Step 3: Apply guidance to noise prediction
11:     $\tilde{\epsilon} = \epsilon_u - \omega\sqrt{1 - \bar{\alpha}_t}\mathbf{g}$
12:     // Step 4: Standard DDIM update with guided noise
13:     $\hat{\mathbf{x}}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\tilde{\epsilon}}{\sqrt{\bar{\alpha}_t}}$         ▷ Predict clean image
14:     $\mathbf{x}_s = \sqrt{\bar{\alpha}_s}\hat{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_s}\tilde{\epsilon}$         ▷ Re-noise for next step
15: **end for**
16: **Return:** $\mathbf{x}_0$

---

**What's happening at each step?**

1. **Unconditional prediction:** We use our trained diffusion model to predict what noise should be removed, ignoring any conditioning

2. **Classifier consultation:** We ask our classifier "if I move in this direction, will the image become more likely to be class $\mathbf{y}$?"

3. **Guidance application:** We adjust our noise prediction to incorporate the classifier's advice

4. **DDIM step:** We take a standard DDIM step using our guided noise prediction

## 9.3.4   Training Classifiers for Guidance: The Noisy Challenge

Here's where things get interesting. We can't just use any classifier—we need one that works well on noisy images at all timesteps of the diffusion process. This is like training a radiologist to diagnose diseases from X-rays with varying amounts of static interference!

**The Training Recipe**

**Why is this challenging?** Standard classifiers are trained on clean images. But during diffusion sampling, we need to classify images that are heavily corrupted with noise. A classifier trained only on clean data will be useless when faced with the noisy intermediate states of the diffusion process.

     **The solution: Noise-aware training**

---

**Algorithm 9** Training a Classifier for Guidance

---

1: **Input:** Clean training dataset $\{(\mathbf{x}_0^{(i)}, \mathbf{y}^{(i)})\}$, noise schedule $\{\bar{\alpha}_t\}$
2: Initialize classifier $p_\phi(\mathbf{y}|\mathbf{x}_t, t)$ $\qquad$ ▷ Note: takes both image and timestep
3: **while** not converged **do**
4: $\qquad$ Sample batch $\{(\mathbf{x}_0^{(i)}, \mathbf{y}^{(i)})\}$ from training set
5: $\qquad$ **for** each sample $(\mathbf{x}_0, \mathbf{y})$ in batch **do**
6: $\qquad\qquad$ $t \sim \text{Uniform}(0, T)$ $\qquad\qquad\qquad$ ▷ Random timestep
7: $\qquad\qquad$ $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ $\qquad\qquad\qquad\qquad$ ▷ Random noise
8: $\qquad\qquad$ $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$ $\qquad\qquad$ ▷ Add noise
9: $\qquad\qquad$ Compute loss: $\mathcal{L} = -\log p_\phi(\mathbf{y}|\mathbf{x}_t, t)$ $\qquad$ ▷ Cross-entropy loss
10: $\qquad$ **end for**
11: $\qquad$ Update $\phi$ to minimize $\mathcal{L}$
12: **end while**

---

**Key insights about this training process:**

- **Timestep conditioning:** The classifier takes both the noisy image *and* the timestep as input, so it knows how much noise to expect

- **Uniform noise sampling:** We train on all noise levels equally, so the classifier is robust across the entire diffusion trajectory

- **Same noise schedule:** We use exactly the same noise schedule that was used to train the diffusion model

**Architecture Considerations**

**What makes a good guidance classifier?**

- **Timestep embedding:** Similar to diffusion models, we need to inject timestep information throughout the network

- **Robustness to noise:** Deeper networks and attention mechanisms help maintain performance under heavy noise

- **Multi-scale features:** Different noise levels might require different receptive fields

- **Regularization:** Heavy dropout and weight decay prevent overfitting to noise patterns

## 9.3.5 The Guidance Scale: A Delicate Balance

The guidance scale $\omega$ is perhaps the most important hyperparameter in classifier guidance. It controls the fundamental trade-off between fidelity to the conditioning and sample diversity.

**Understanding the Trade-off**

**Low guidance ($\omega \to 0$):**

- **Pro:** High sample diversity—images look natural and varied

- **Pro:** Stable sampling—less risk of adversarial artifacts

- **Con:** Poor conditioning—might not generate the desired class

- **Con:** Weak control—similar to unconditional sampling

**High guidance ($\omega \gg 1$):**

- **Pro:** Strong conditioning—clearly generates the desired class

- **Pro:** Predictable behavior—reliable conditioning

- **Con:** Low diversity—samples might look very similar

- **Con:** Artifacts—adversarial-like effects from over-optimization

**The sweet spot ($\omega \in [1, 10]$):** Most applications find optimal performance in this range, but the exact value requires empirical tuning for each domain and classifier quality.

**Dynamic Guidance Schedules**

An advanced technique is to vary the guidance scale throughout the sampling process:

$$\omega(t) = \omega_{\max} \cdot f(t) \tag{9.14}$$

Common schedules include:

- **Linear decay:** $f(t) = t/T$ (stronger guidance early, weaker later)

- **Exponential decay:** $f(t) = e^{-\lambda(T-t)}$ (smooth transition)

- **Step function:** High guidance for certain timestep ranges, low for others

**Intuition:** Early in sampling (high noise), we need strong guidance to establish the basic structure. Later (low noise), we want to preserve fine details and avoid artifacts.

## 9.3.6 Practical Considerations and Limitations

While classifier guidance is powerful, it comes with several practical challenges that practitioners must navigate.

**Computational Overhead: The Price of Control**

**Additional costs during sampling:**

- **Forward pass:** Classifier evaluation at each step

- **Backward pass:** Gradient computation through classifier

- **Memory:** Storing gradients and intermediate activations

- **Time:** Typically 20-50% slower than unconditional sampling

  **Optimization strategies:**

- Use smaller, efficient classifier architectures

- Implement gradient checkpointing to reduce memory

- Cache classifier features when possible

- Use mixed precision arithmetic

**The Classifier Quality Bottleneck**

**Why classifier quality matters so much:**

- **Garbage in, garbage out:** Poor classifiers provide bad guidance signals

- **Noise robustness:** Classifiers must work well across all noise levels

- **Domain alignment:** Training data should match generation targets

- **Calibration:** Overconfident classifiers can cause instability

  **Common failure modes:**

- **Adversarial artifacts:** Overoptimization toward classifier weaknesses

- **Mode collapse:** All samples look similar due to classifier bias

- **Label bleeding:** Classifier confusion causes mixed or incorrect classes

- **Unstable sampling:** Poor gradients cause sampling to diverge

**Conditioning Limitations**

**What works well:**

- **Classification tasks:** Natural fit for categorical conditions

- **Continuous attributes:** Age, pose, brightness (with regression)

- **Spatial conditions:** With appropriate spatial classifiers

  **What's challenging:**

- **Complex text:** Difficult to build robust text classifiers

- **Compositional conditions:** "Red car next to blue house"

- **Abstract concepts:** Style, mood, artistic properties

- **Long-range dependencies:** Global scene understanding

### 9.3.7 When to Use Classifier Guidance

**Classifier guidance shines when:**

- You have a pre-trained unconditional model you want to control

- The conditioning is well-suited to classification (categorical, continuous attributes)

- You can train a robust classifier on noisy data

- Computational overhead is acceptable for your application

- You need modular control (different classifiers for different conditions)

**Consider alternatives when:**

- Conditioning is complex or compositional (text, scenes)

- Speed is critical and overhead is prohibitive

- Training robust classifiers is difficult for your domain

- You're building a system from scratch (conditional training might be better)

In our next section, we'll explore classifier-free guidance, which addresses many of these limitations by eliminating the need for separate classifiers entirely!

## 9.4 Classifier-Free Guidance: The Elegant Solution

### 9.4.1 The Breakthrough: Why CFG Changed Everything

Imagine you're a chef who has just discovered that instead of needing two separate ovens—one for baking bread and another for roasting meat—you can train a single, versatile oven to do both tasks perfectly. That's essentially what Classifier-Free Guidance (CFG) accomplished in the world of diffusion models.

Before CFG, if we wanted controllable generation, we faced a fundamental dilemma:

- **Option 1:** Train separate conditional models for each type of control (expensive and inflexible)

- **Option 2:** Use classifier guidance (requires separate classifier training and introduces computational overhead)

CFG introduced a third, elegant option: *train one model to be both conditional and unconditional, then use the difference between these capabilities to achieve control.*

**The core insight that started a revolution:**

Remember from our discussion of classifier guidance that we wanted to compute:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) \tag{9.15}$$

The brilliant realization was this: *What if we could compute the classifier gradient term without actually having a classifier?*

It turns out we can! Through a beautiful mathematical relationship, the difference between conditional and unconditional scores contains exactly the information we need:

$$\boxed{\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)} \tag{9.16}$$

**Let's understand why this is true:** Starting from Bayes' rule in log form:

$$\log p(\mathbf{x}_t|\mathbf{y}) = \log p(\mathbf{y}|\mathbf{x}_t) + \log p(\mathbf{x}_t) - \log p(\mathbf{y}) \tag{9.17}$$

Taking gradients with respect to $\mathbf{x}_t$:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) - \underbrace{\nabla_{\mathbf{x}_t} \log p(\mathbf{y})}_{=0} \tag{9.18}$$

Rearranging:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \tag{9.19}$$

**What does this mean in plain English?** If we have a model that can compute both conditional and unconditional scores, we can extract the "classifier-like" guidance signal simply by taking their difference! No separate classifier needed.

**Why is this revolutionary?**

- **Unified architecture:** One model, infinite conditioning possibilities

- **No separate training:** No need to train classifiers on noisy data

- **Computational efficiency:** No gradient computation through classifiers

- **Stability:** Avoids adversarial-like optimization issues

- **Flexibility:** Easily extends to text, images, or any conditioning modality

## 9.4.2   Training a Jack-of-All-Trades Model

The key to CFG is training a single model that can seamlessly switch between conditional and unconditional generation. This is like teaching a translator to work both with and without context clues.

### The Conditioning Dropout Technique

The training procedure is elegantly simple: during training, we randomly "forget" the conditioning information for some percentage of samples. This forces the model to learn both conditional and unconditional generation simultaneously.

**Algorithm 10** Classifier-Free Guidance Training

---

1: **Input:** Dataset $\{(\mathbf{x}^{(i)}, \mathbf{y}^{(i)})\}$, unconditional dropout probability $p_{\text{uncond}}$
2: Initialize model $\epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)$        ▷ Takes image, condition, and timestep
3: **while** not converged **do**
4:      Sample batch $\{(\mathbf{x}, \mathbf{y})\}$ from training dataset
5:      Sample timesteps $t \sim \text{Uniform}(1, T)$ for each sample
6:      Sample noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ for each sample
7:      Compute noisy samples: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x} + \sqrt{1 - \bar{\alpha}_t}\epsilon$
8:      **for** each sample $(\mathbf{x}_t, \mathbf{y}, t)$ in batch **do**
9:          **if** random() $< p_{\text{uncond}}$ **then**      ▷ Randomly drop conditioning
10:             Set $\mathbf{y} = \emptyset$      ▷ Use null token for unconditional training
11:          **end if**
12:      **end for**
13:      Predict noise: $\hat{\epsilon} = \epsilon_\theta(\mathbf{x}_t, \mathbf{y}, t)$
14:      Compute MSE loss: $\mathcal{L} = \|\epsilon - \hat{\epsilon}\|^2$
15:      Update parameters: $\theta \leftarrow \theta - \eta\nabla_\theta\mathcal{L}$
16: **end while**

---

### What's happening step by step?

1. **Standard setup:** We start with the usual diffusion training—add noise to images and train the model to predict it

2. **Conditioning dropout:** For a fraction of samples (typically 10-20%), we replace the real condition $\mathbf{y}$ with a special "null" token $\emptyset$

3. **Dual learning:** The model learns to predict noise both when it has conditioning information and when it doesn't

4. **Implicit knowledge:** By learning both tasks, the model implicitly learns the relationship between conditional and unconditional distributions

**Critical Training Considerations**

**The dropout rate ($p_{\text{uncond}}$): Finding the sweet spot**

- **Too low ($< 5\%$):** Model doesn't learn unconditional generation well

- **Too high ($> 30\%$):** Model's conditional generation suffers

- **Just right (10-20%):** Good balance between conditional and unconditional capabilities

**The null token ($\emptyset$): Representing "nothing"**

- **For text:** Use a special [NULL] token or empty string

- **For classes:** Use a special "no class" embedding

- **For embeddings:** Use zero vector or learned null embedding

- **Key insight:** The model must be able to recognize and respond to this null condition

**Architecture requirements:**

- **Graceful degradation:** Model must handle missing conditioning without breaking

- **Condition injection:** Clear pathways for conditioning information throughout the network

- **Null handling:** Proper masking or zero-ing mechanisms for null conditions

### 9.4.3 The CFG Magic Formula: Where Mathematics Meets Elegance

Now comes the beautiful part: how do we actually use our dual-capable model to achieve guidance during sampling?

**Deriving the CFG Equation**

Let's build up to the famous CFG formula step by step. We want to sample from the conditional distribution $p(\mathbf{x}_t|\mathbf{y})$, but with amplified conditioning strength.

From our earlier insight, we know:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \tag{9.20}$$

Now, we want to create a *modified* conditional score that amplifies the conditioning:

$$\tilde{\nabla}_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + (1+\omega)\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) \tag{9.21}$$

where $\omega \geq 0$ is our guidance scale that amplifies the classifier gradient term. Substituting our expression for the classifier gradient:

$$\tilde{\nabla}_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + (1+\omega)[\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)] \tag{9.22}$$

$$= \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) + (1+\omega)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - (1+\omega)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \tag{9.23}$$

$$= (1+\omega)\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t|\mathbf{y}) - \omega\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) \tag{9.24}$$

Now, converting back to noise predictions using the relationship $\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t) = -\frac{\epsilon_\theta(\mathbf{x}_t,t)}{\sqrt{1-\bar{\alpha}_t}}$:

$$-\frac{\tilde{\epsilon}}{\sqrt{1-\bar{\alpha}_t}} = (1+\omega)\left(-\frac{\epsilon_\theta(\mathbf{x}_t,\mathbf{y},t)}{\sqrt{1-\bar{\alpha}_t}}\right) - \omega\left(-\frac{\epsilon_\theta(\mathbf{x}_t,\emptyset,t)}{\sqrt{1-\bar{\alpha}_t}}\right) \tag{9.25}$$

Simplifying:

$$\boxed{\tilde{\epsilon}(\mathbf{x}_t,\mathbf{y},t) = (1+\omega)\epsilon_\theta(\mathbf{x}_t,\mathbf{y},t) - \omega\epsilon_\theta(\mathbf{x}_t,\emptyset,t)} \tag{9.26}$$

**This is the heart of classifier-free guidance!**

**Understanding the Formula Intuitively**

Let's rewrite the CFG formula in a more intuitive form:

$$\tilde{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t) + (1 + \omega)[\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t) - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t)] \tag{9.27}$$

This reveals the geometric interpretation:

1. **Start with unconditional:** $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t)$ gives us the "default" noise prediction

2. **Compute conditioning direction:** $\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t) - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t)$ shows how conditioning changes the prediction

3. **Amplify and apply:** We amplify this direction by $(1 + \omega)$ and apply it to the unconditional baseline

**Geometric analogy:** Imagine you're walking from point A to point B. The unconditional prediction shows you the "default" direction. The conditioning direction shows you how to adjust your path to reach your specific destination. CFG amplifies this adjustment—like taking bigger steps in the right direction!

**Special cases that build intuition:**

- $\omega = 0$: $\tilde{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)$ (standard conditional sampling)

- $\omega = -1$: $\tilde{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t)$ (pure unconditional sampling)

- $\omega > 0$: Extrapolation beyond conditional prediction (amplified conditioning)

## 9.4.4 Implementing CFG: Making Theory Practical

Now let's see how this beautiful theory translates into practical sampling code.

---

**Algorithm 11** DDIM Sampling with Classifier-Free Guidance

---

1: **Input:** Model $\boldsymbol{\epsilon}_\theta$, condition $\mathbf{y}$, guidance scale $\omega$, timestep schedule $\{\tau_1, \ldots, \tau_S\}$
2: Sample $\mathbf{x}_{\tau_S} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ ▷ Start from pure noise
3: **for** $i = S, S - 1, \ldots, 1$ **do**
4: $\quad t = \tau_i$, $s = \tau_{i-1}$ ▷ Current and next timesteps
5: $\quad$ // Step 1: Get both predictions
6: $\quad \boldsymbol{\epsilon}_c = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)$ ▷ Conditional prediction
7: $\quad \boldsymbol{\epsilon}_u = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t)$ ▷ Unconditional prediction
8: $\quad$ // Step 2: Apply classifier-free guidance
9: $\quad \hat{\boldsymbol{\epsilon}} = (1 + \omega)\boldsymbol{\epsilon}_c - \omega\boldsymbol{\epsilon}_u$ ▷ CFG magic formula
10: $\quad$ // Step 3: Standard DDIM update
11: $\quad$ Predict clean image: $\hat{\mathbf{x}}_0 = \frac{\mathbf{x}_t - \sqrt{1 - \bar{\alpha}_t}\hat{\boldsymbol{\epsilon}}}{\sqrt{\bar{\alpha}_t}}$
12: $\quad$ Re-noise for next step: $\mathbf{x}_s = \sqrt{\bar{\alpha}_s}\hat{\mathbf{x}}_0 + \sqrt{1 - \bar{\alpha}_s}\hat{\boldsymbol{\epsilon}}$
13: **end for**
14: **Return:** $\mathbf{x}_0$

---

**What's happening at each step?**

1. **Dual prediction:** We run our model twice—once with the real condition, once with the null condition

2. **Guidance computation:** We apply the CFG formula to combine these predictions

3. **Standard update:** We use the guided prediction in a normal DDIM step

**Implementation Optimizations**

**Computational considerations:**

- **Double forward passes:** CFG requires two model evaluations per step (doubled computational cost)

- **Batching trick:** Stack conditional and unconditional inputs in a single batch for efficiency

- **Memory usage:** Effective batch size doubles during inference

- **Caching opportunities:** Some computation might be shared between conditional/unconditional paths

**Batching optimization example:**

1: // Instead of two separate forward passes:
2: $\boldsymbol{\epsilon}_c = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)$
3: $\boldsymbol{\epsilon}_u = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t)$
4: // Use batched computation:
5: $\mathbf{x}_{\text{batch}} = [\mathbf{x}_t, \mathbf{x}_t]$                                  ▷ Duplicate input
6: $\mathbf{y}_{\text{batch}} = [\mathbf{y}, \emptyset]$                                  ▷ Conditional and null
7: $\boldsymbol{\epsilon}_{\text{batch}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_{\text{batch}}, \mathbf{y}_{\text{batch}}, t)$
8: $\boldsymbol{\epsilon}_c, \boldsymbol{\epsilon}_u = \text{split}(\boldsymbol{\epsilon}_{\text{batch}})$                                  ▷ Extract results

## 9.4.5 CFG vs. Classifier Guidance: The Beautiful Connection

One of the most elegant aspects of CFG is how it relates to classifier guidance. They're not just similar—they're mathematically equivalent under certain conditions!

**The Mathematical Connection**

Recall that classifier guidance uses:

$$\tilde{\boldsymbol{\epsilon}}_{\text{classifier}} = \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, t) - \omega\sqrt{1 - \bar{\alpha}_t}\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) \tag{9.28}$$

While CFG uses:

$$\tilde{\boldsymbol{\epsilon}}_{\text{CFG}} = (1 + \omega)\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t) - \omega\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t) \tag{9.29}$$

These become equivalent when:

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) \approx \frac{\boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \emptyset, t) - \boldsymbol{\epsilon}_\theta(\mathbf{x}_t, \mathbf{y}, t)}{\sqrt{1 - \bar{\alpha}_t}} \tag{9.30}$$

**What does this tell us?** The difference between conditional and unconditional noise predictions contains the same information as a classifier gradient! CFG essentially trains an "implicit classifier" as part of the diffusion model.

**Why CFG Often Works Better**

In practice, CFG frequently outperforms classifier guidance:

**Advantages of the implicit approach:**

- **Joint training:** The "classifier" is trained together with the diffusion model, ensuring perfect alignment

- **No domain gap:** No mismatch between classifier training data and generation targets

- **Stable optimization:** Avoids adversarial-like dynamics of gradient-based guidance

- **Richer representations:** The diffusion model's internal representations are often richer than external classifiers

**Flexibility benefits:**

- **Complex conditioning:** Easily handles text, images, or multi-modal conditions

- **Compositional control:** Better at combining multiple conditions

- **Fine-grained control:** Can condition on subtle attributes that are hard to classify

## 9.4.6   The Guidance Scale: Your Creative Control Knob

The guidance scale $\omega$ is your primary tool for controlling the behavior of CFG. Understanding how to use it effectively is crucial for getting the results you want.

**The Quality-Diversity Trade-off**

**Low guidance ($\omega \in [0, 1]$): The diversity regime**

- **High diversity:** Samples are varied and creative

- **Natural appearance:** Images look realistic and organic

- **Weak conditioning:** May not follow prompts/conditions precisely

- **Inconsistent results:** Output can be unpredictable

**Medium guidance ($\omega \in [2, 7]$): The sweet spot**

- **Balanced trade-off:** Good mix of quality and diversity

- **Reliable conditioning:** Follows conditions while maintaining creativity

- **Production ready:** Most commercial applications use this range

**High guidance ($\omega \in [8, 15]$): The fidelity regime**

- **Strong conditioning:** Precisely follows prompts and conditions

- **Consistent results:** Predictable, reliable outputs

- **Low diversity:** Samples may look similar or repetitive

- **Potential artifacts:** Over-optimization can create unrealistic features

**Very high guidance ($\omega > 15$): The danger zone**

- **Oversaturation:** Colors and features become exaggerated

- **Artifacts:** Unrealistic or uncanny valley effects

- **Mode collapse:** All outputs start looking the same

**Advanced Guidance Strategies**

**Dynamic guidance scheduling:** Instead of using a fixed $\omega$, we can vary it during sampling:

$$\omega(t) = \omega_{\max} \cdot f(t) \tag{9.31}$$

Common scheduling functions:

- **Constant:** $f(t) = 1$ (standard CFG)

- **Linear decay:** $f(t) = \frac{t}{T}$ (stronger guidance early, weaker later)

- **Cosine schedule:** $f(t) = \frac{1 + \cos(\pi t/T)}{2}$ (smooth transition)

- **Step function:** High guidance for structure, low for details

**Intuition behind dynamic guidance:**

- **Early timesteps (high noise):** Need strong guidance to establish correct structure and content

- **Late timesteps (low noise):** Want weaker guidance to preserve fine details and avoid artifacts

**Condition-dependent guidance:** Adjust $\omega$ based on the complexity or specificity of the condition:

- **Simple conditions:** Lower guidance (e.g., "a dog" $\to \omega = 3$)

- **Complex conditions:** Higher guidance (e.g., "a golden retriever wearing a red hat sitting on a blue chair" $\to \omega = 7$)

- **Style conditions:** Medium guidance to balance style transfer with natural appearance

### 9.4.7 Practical Considerations and Best Practices

**Training Best Practices**

**Choosing the dropout rate:**

- **Start with 10-15%** for most applications

- **Increase to 20%** if unconditional performance is poor

- **Decrease to 5-10%** if conditional performance suffers

- **Monitor both** conditional and unconditional sample quality during training

**Null token design:**

- **Text conditioning:** Use a special [NULL] or empty string, ensure tokenizer handles it properly

- **Class conditioning:** Add a special "background" or "no class" category

- **Embedding conditioning:** Use learned null embedding or zero vector

- **Consistency:** Use the same null representation during training and inference

**Architecture considerations:**

- **Condition injection:** Ensure conditioning pathways can be cleanly disabled

- **Normalization:** Be careful with batch norm when mixing conditional/unconditional samples

- **Attention masks:** Properly mask attention when using null tokens

**Common Pitfalls and Solutions**

**Problem: Poor unconditional generation**

- **Symptom:** Unconditional samples are low quality or nonsensical

- **Cause:** Dropout rate too low, or null token not properly implemented

- **Solution:** Increase dropout rate, verify null token handling

**Problem: Weak conditioning at low guidance**

- **Symptom:** Even with guidance, samples don't follow conditions well

- **Cause:** Poor separation between conditional and unconditional distributions

- **Solution:** Check conditioning architecture, ensure sufficient model capacity

**Problem: Artifacts at high guidance**

- **Symptom:** Oversaturated colors, unrealistic features, or repeated patterns

- **Cause:** Over-extrapolation beyond training distribution

- **Solution:** Use lower guidance, implement guidance scheduling, or improve training data

**Problem: Mode collapse with CFG**

- **Symptom:** All samples look very similar despite different seeds

- **Cause:** Guidance scale too high, or conditioning too specific

- **Solution:** Reduce guidance scale, use dynamic guidance, or add diversity regularization

## 9.4.8   When to Use CFG: The Decision Matrix

**CFG is your best choice when:**

- Building a new conditional generation system from scratch

- Working with complex conditioning (text, multi-modal, compositional)

- Need flexible control over the conditioning strength

- Computational efficiency during training is important

- Want state-of-the-art results with minimal complexity

**Consider alternatives when:**

- You already have a well-trained unconditional model and good classifiers

- Working with simple categorical conditioning where classifiers work well

- Need to retrofit conditioning onto existing models without retraining

- Memory constraints make double forward passes prohibitive

**The modern consensus:** CFG has become the dominant approach for conditional diffusion models, powering systems like Stable Diffusion, DALL-E 2, and Midjourney. Its combination of simplicity, effectiveness, and flexibility makes it the go-to choice for most applications.

In our next section, we'll explore how CFG enables sophisticated text-to-image generation and discuss advanced conditioning techniques that push the boundaries of what's possible with diffusion models!

## 9.5 Conditioning Approaches: The Complete Picture

### 9.5.1 A Comprehensive Comparison: Understanding Your Options

Now that we've explored the three major conditioning paradigms, let's step back and see the full landscape. Think of this as choosing the right tool for the job—each approach has its own strengths, weaknesses, and ideal use cases.

| Aspect | Class-Conditional | Classifier Guidance | Classifier-Free |
| --- | --- | --- | --- |
| **Training Complexity** | Low | High | Medium |
| **Inference Speed** | Fast | Slow | Medium |
| **Conditioning Flexibility** | Limited | Medium | High |
| **Implementation Difficulty** | Easy | Hard | Medium |
| **Sample Quality** | Good | Excellent | Excellent |
| **Memory Usage** | Low | High | Medium |
| **Guidance Strength Control** | Fixed | Variable | Variable |
| **Retraining Required** | Yes | No | Yes |
| **External Dependencies** | None | Classifier | None |
| **Modern Adoption** | Legacy | Niche | Standard |

Table 9.1: Detailed comparison of conditional generation approaches across key dimensions

**Decoding the Trade-offs**

Let's dive deeper into what these comparisons really mean in practice:
**Training Complexity Analysis:**

- **Class-conditional (Low):** Just add conditioning to standard training—straightforward and predictable

- **Classifier guidance (High):** Requires training both a diffusion model AND robust noise-aware classifiers

- **CFG (Medium):** Single model training but with conditioning dropout complexity

**Inference Speed Breakdown:**

- **Class-conditional:** Single forward pass per step—as fast as unconditional generation

- **Classifier guidance:** Forward + backward pass through classifier at each step—significant overhead

- **CFG:** Two forward passes (conditional + unconditional)—2x model evaluation cost

**Flexibility Spectrum:**

- **Class-conditional:** Limited to categories seen during training—no creative combinations

- **Classifier guidance:** Bounded by classifier capabilities—good for well-defined attributes

- **CFG:** Handles complex, compositional, and even out-of-distribution conditioning

## 9.5.2   The Decision Framework: When to Use What

Choosing the right conditioning approach is like selecting the right vehicle for a journey—it depends on your destination, constraints, and priorities.

**Class-Conditional Models: The Reliable Workhorse**

**Choose class-conditional when:**

- **Simple, fixed conditioning:** You have well-defined categories (ImageNet classes, MNIST digits)

- **Speed is paramount:** Real-time applications where every millisecond counts

- **Implementation simplicity:** You need a straightforward, proven approach

- **Limited computational resources:** Memory and compute are constrained

- **Legacy compatibility:** Working with existing class-conditional datasets

**Real-world examples:**

- Medical image generation for specific diagnostic categories

- Game asset generation with predefined object types

- Scientific visualization with fixed simulation parameters

- Educational tools with structured content categories

**Classifier Guidance: The Precision Instrument**

**Choose classifier guidance when:**

- **Retrofitting existing models:** You have great unconditional models and want to add control

- **Multiple conditioning objectives:** Need to combine different types of guidance signals

- **Maximum quality is critical:** Willing to pay computational costs for best results

- **Robust classifiers available:** You have or can train excellent noise-aware classifiers

- **Research applications:** Exploring novel conditioning mechanisms

  **Real-world examples:**

- Scientific image generation where multiple physical constraints must be satisfied

- Art generation combining style, color, and composition classifiers

- Medical imaging where multiple diagnostic criteria need simultaneous control

- Research prototypes exploring novel conditioning mechanisms

**Classifier-Free Guidance: The Modern Standard**

**Choose CFG when:**

- **Complex conditioning:** Text prompts, compositional descriptions, or multi-modal inputs

- **Production applications:** Building commercial or user-facing systems

- **Flexible control needed:** Users want variable guidance strength

- **Modern best practices:** Want to follow current state-of-the-art approaches

- **Training from scratch:** Starting a new project without legacy constraints

  **Real-world examples:**

- Text-to-image platforms (Stable Diffusion, Midjourney)

- Content creation tools for artists and designers

- E-commerce product visualization

- Social media content generation

# 9.6 Advanced Conditioning: Beyond the Basics

Now that we understand the fundamental approaches, let's explore the sophisticated conditioning mechanisms that power modern AI applications.

## 9.6.1 Text Conditioning: Teaching Models to Understand Language

Text-to-image generation represents one of the most remarkable achievements in AI—the ability to create visual content from natural language descriptions. Let's unpack how this magic works.

**The Text-to-Image Pipeline: From Words to Pixels**

**The complete journey:**

1. **Text encoding:** Transform natural language into mathematical representations

2. **Cross-modal alignment:** Bridge the gap between text and visual features

3. **Spatial injection:** Integrate text understanding throughout the generation process

4. **Guided sampling:** Use text information to steer the diffusion process

**Step 1: Text Encoding Strategies**

**CLIP: The Vision-Language Pioneer** CLIP (Contrastive Language-Image Pre-training) revolutionized text conditioning by learning aligned representations:

- **Joint training:** Learned on 400M image-text pairs from the internet

- **Contrastive learning:** Positive pairs (matching image-text) pulled together, negative pairs pushed apart

- **Semantic alignment:** Similar concepts in text and vision occupy similar embedding spaces

- **Zero-shot capability:** Can understand concepts not seen during diffusion training

**T5: The Language Understanding Powerhouse** T5 (Text-to-Text Transfer Transformer) brings deep linguistic understanding:

- **Rich semantics:** Better understanding of grammar, context, and nuance

- **Longer sequences:** Can handle detailed, complex prompts

- **Compositional understanding:** Better at parsing complex scene descriptions

- **Language flexibility:** Strong performance across different languages and domains

**Choosing your text encoder:**

- **CLIP for visual alignment:** When image-text consistency is paramount

- **T5 for complex language:** When dealing with detailed, nuanced prompts

- **Domain-specific encoders:** For specialized vocabularies (medical, scientific, artistic)

- **Multilingual models:** For international applications

### Step 2: Cross-Attention—Where Text Meets Vision

The magic happens in cross-attention layers, where spatial features from the U-Net can "look at" and "understand" different parts of the text prompt.

**Mathematical formulation:**

$$\text{Query (from image): } \mathbf{Q} = \mathbf{F}_{\text{spatial}}\mathbf{W}_Q \tag{9.32}$$

$$\text{Key (from text): } \mathbf{K} = \mathbf{F}_{\text{text}}\mathbf{W}_K \tag{9.33}$$

$$\text{Value (from text): } \mathbf{V} = \mathbf{F}_{\text{text}}\mathbf{W}_V \tag{9.34}$$

$$\text{Attention weights: } \mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right) \tag{9.35}$$

$$\text{Attended features: } \mathbf{F}_{\text{out}} = \mathbf{A}\mathbf{V} \tag{9.36}$$

**Intuitive understanding:**

- **Queries:** Each spatial location asks "what should I generate here?"

- **Keys:** Text tokens say "I represent this concept"

- **Values:** Text tokens provide "this is how I should influence generation"

- **Attention:** Spatial locations focus on relevant text concepts

**Why this works so well:**

- **Spatial specificity:** Different image regions can attend to different text concepts

- **Dynamic focus:** Attention patterns change throughout the denoising process

- **Compositional understanding:** Can handle complex prompts with multiple objects and relationships

- **Flexible conditioning:** Same architecture works for simple and complex prompts

**Prompt Engineering: The Art of Communication**

Getting the best results from text-to-image models requires understanding how to communicate effectively with them.

**Anatomy of an effective prompt:**

1. **Main subject:** What is the primary focus? ("A majestic lion")

2. **Setting/environment:** Where is this happening? ("in an African savanna")

3. **Style/medium:** How should it look? ("oil painting style")

4. **Lighting/mood:** What's the atmosphere? ("golden hour lighting")

5. **Technical details:** Camera/artistic specifications ("shallow depth of field")

**Advanced prompting techniques:**

- **Weighted terms:** Use (concept:1.5) or [concept:0.7] to adjust importance

- **Negative prompts:** Specify what to avoid for better results

- **Progressive prompts:** Start simple, then add detail in subsequent generations

- **Style transfer prompts:** "in the style of [artist]" or "as a [medium]"

## 9.6.2   Image Conditioning: Visual Control Mechanisms

Sometimes we want to control generation not through words, but through visual cues—sketches, layouts, depth maps, or other images.

**ControlNet: The Breakthrough in Visual Control**

ControlNet revolutionized image conditioning by showing how to add precise visual controls to pre-trained diffusion models without breaking them.

**The ControlNet insight:** Instead of modifying the original model, add parallel "control pathways" that inject conditioning information at multiple scales.

**Algorithm 12** ControlNet Training Strategy

---

1: **Input:** Pre-trained diffusion model, control dataset $\{(\mathbf{x}, \mathbf{c})\}$
2: <span style="color:red">// Step 1: Preserve original model</span>
3: Freeze all parameters of the original U-Net $\theta_{\text{orig}}$
4: <span style="color:blue">// Step 2: Create control pathway</span>
5: Initialize control modules as copies of U-Net encoder blocks
6: Initialize control pathway parameters $\theta_{\text{control}}$
7: <span style="color:green">// Step 3: Training loop</span>
8: **while** not converged **do**
9:     Sample batch $\{(\mathbf{x}, \mathbf{c})\}$ from control dataset
10:     Add noise: $\mathbf{x}_t = \sqrt{\bar{\alpha}_t}\mathbf{x} + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$
11:     <span style="color:blue">// Process control information</span>
12:     $\mathbf{f}_{\text{control}} = \text{ControlEncoder}(\mathbf{c})$
13:     <span style="color:blue">// Forward through both pathways</span>
14:     $\mathbf{f}_{\text{orig}} = \text{UNet}_{\text{orig}}(\mathbf{x}_t, t)$                          ▷ Frozen pathway
15:     $\mathbf{f}_{\text{combined}} = \mathbf{f}_{\text{orig}} + \mathbf{f}_{\text{control}}$                 ▷ Additive combination
16:     Predict noise: $\hat{\boldsymbol{\epsilon}} = \text{UNet}_{\text{combined}}(\mathbf{f}_{\text{combined}})$
17:     Compute loss: $\mathcal{L} = \|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}\|^2$
18:     Update only $\theta_{\text{control}}$ (keep $\theta_{\text{orig}}$ frozen)
19: **end while**

---

**Why ControlNet works so brilliantly:**

- **Zero-harm principle:** Original capabilities are perfectly preserved

- **Additive control:** Control features are added, not replacing existing knowledge

- **Multi-scale injection:** Control information flows at multiple resolution levels

- **Fast training:** Only control parameters need updating—much faster than full retraining

**Popular Control Modalities**

**Canny Edge Control:** Structure without content

- **Use case:** Maintain composition while changing style or content

- **Preprocessing:** Apply Canny edge detection to control image

- **Strength:** Preserves fine structural details

- **Example:** Convert a photo's edges into a painting with the same composition

**Depth Map Control:** 3D structure guidance

- **Use case:** Control spatial layout and depth relationships

- **Preprocessing:** Extract depth using MiDaS or similar models

- **Strength:** Maintains spatial relationships and perspective

- **Example:** Generate new scenes with the same 3D layout

**Human Pose Control:** Character positioning

- **Use case:** Control human figures and character poses

- **Preprocessing:** Extract pose keypoints using OpenPose

- **Strength:** Precise control over human anatomy and positioning

- **Example:** Generate different people in the same pose

**Segmentation Mask Control:** Object-level layout

- **Use case:** Control where different objects appear

- **Preprocessing:** Create or extract segmentation masks

- **Strength:** Precise object placement and scene composition

- **Example:** Generate scenes with specific object layouts

### 9.6.3 Multi-Modal Conditioning: The Ultimate Control

The most powerful applications combine multiple conditioning modalities—text descriptions with visual controls, audio with video, etc.

**Combination Strategies**

**Feature Concatenation:** Simple but effective

$$\mathbf{f}_{\text{combined}} = \text{Concat}([\mathbf{f}_{\text{text}}, \mathbf{f}_{\text{image}}, \mathbf{f}_{\text{other}}]) \tag{9.37}$$

**Multi-Head Cross-Attention:** Specialized attention for each modality

$$\mathbf{f}_{\text{out}} = \text{Attention}_{\text{text}}(\mathbf{Q}, \mathbf{K}_{\text{text}}, \mathbf{V}_{\text{text}}) + \text{Attention}_{\text{image}}(\mathbf{Q}, \mathbf{K}_{\text{image}}, \mathbf{V}_{\text{image}}) \tag{9.38}$$

**Hierarchical Conditioning:** Different modalities at different scales

- **Global scale:** Text describes overall scene and style

- **Medium scale:** Image control provides spatial layout

- **Fine scale:** Detail conditioning for texture and finish

**Adaptive Fusion:** Learn to weight different conditioning signals

$$\mathbf{w} = \text{MLP}([\mathbf{f}_{\text{text}}, \mathbf{f}_{\text{image}}]) \tag{9.39}$$

$$\mathbf{f}_{\text{combined}} = w_1 \mathbf{f}_{\text{text}} + w_2 \mathbf{f}_{\text{image}} + w_3 \mathbf{f}_{\text{other}} \tag{9.40}$$

**Real-World Multi-Modal Applications**

**Architectural visualization:**

- **Text:** "Modern minimalist house with large windows"

- **Floor plan:** 2D layout specifying room arrangements

- **Style image:** Reference photos for architectural style

- **Result:** Photorealistic renderings matching all constraints

**Fashion design:**

- **Text:** "Elegant evening dress with flowing fabric"

- **Sketch:** Rough design outline and silhouette

- **Color palette:** Specific color scheme constraints

- **Result:** Detailed fashion renderings ready for production

**Product visualization:**

- **Text:** "Luxury smartphone with premium materials"

- **3D model:** Basic geometric constraints

- **Material samples:** Texture and finish references

- **Result:** Marketing-ready product renders

## 9.7 Implementation Mastery: From Theory to Production

Understanding the theory is just the beginning. Let's explore the practical considerations that separate successful implementations from research prototypes.

### 9.7.1 Training Best Practices: The Foundation of Success

**Data Preparation: Quality In, Quality Out**

**Conditioning data quality checklist:**

- **Accuracy:** Are text descriptions or labels accurate and comprehensive?

- **Consistency:** Do similar inputs have similar conditioning information?

- **Coverage:** Does the dataset cover the full range of intended use cases?

- **Balance:** Are different conditioning categories well-represented?

- **Clean pairing:** Do conditioning inputs match their target outputs correctly?

**Data augmentation strategies:**

- **Synonym replacement:** Vary text descriptions while preserving meaning

- **Paraphrasing:** Generate alternative descriptions for the same content

- **Partial conditioning:** Train with incomplete conditioning information

- **Noise injection:** Add controlled noise to conditioning inputs for robustness

## Architecture Decisions: Building for Success

**Conditioning injection points:** Where to add conditioning information

- **Input level:** Concatenate with input features (simple but limited)

- **Embedding level:** Add to timestep embeddings (good for global conditioning)

- **Feature level:** Inject at multiple U-Net scales (most flexible)

- **Attention level:** Use cross-attention mechanisms (best for complex conditioning)

**Capacity considerations:**

- **Model size:** Larger models handle complex conditioning better but cost more

- **Conditioning pathway capacity:** Ensure sufficient parameters for conditioning processing

- **Bottleneck analysis:** Identify and address capacity limitations

- **Computational budget:** Balance capability with inference constraints

## Training Strategies: Curriculum and Progressive Learning

**CFG dropout scheduling:**

- **Warm-up phase:** Start with lower dropout rates to establish basic conditioning

- **Progressive increase:** Gradually increase dropout rate as training progresses

- **Final tuning:** Fine-tune with target dropout rate for best CFG performance

- **Validation monitoring:** Track both conditional and unconditional quality

**Curriculum learning for complex conditioning:**

1. **Simple conditioning first:** Start with basic, unambiguous conditions

2. **Complexity progression:** Gradually introduce more complex, compositional conditions

3. **Multi-modal introduction:** Add additional modalities once basic conditioning is stable

4. **Fine-tuning phase:** Final training on full complexity distribution

### 9.7.2 Inference Optimization: Speed Meets Quality

**Computational Efficiency Strategies**

**Batching optimizations:**

```
1:  // Naive approach (inefficient):
2:  for each sample in batch do
3:      ε_c = model(x, y, t)                              ▷ Conditional
4:      ε_u = model(x, ∅, t)                              ▷ Unconditional
5:      Apply CFG formula
6:  end for
7:  // Optimized approach (efficient):
8:  x_batch = concat([x, x])                              ▷ Duplicate inputs
9:  y_batch = concat([y, ∅])                              ▷ Conditional + null
10: ε_batch = model(x_batch, y_batch, t)                 ▷ Single forward pass
11: ε_c, ε_u = split(ε_batch)                            ▷ Extract results
```

**Memory optimization techniques:**

- **Gradient checkpointing:** Trade computation for memory during training

- **Mixed precision:** Use float16 for most operations, float32 for critical paths

- **Model sharding:** Distribute large models across multiple GPUs

- **Dynamic batching:** Adjust batch sizes based on available memory

**Model optimization for deployment:**

- **Quantization:** Reduce precision to int8 or int4 for inference

- **Pruning:** Remove less important parameters to reduce model size

- **Distillation:** Train smaller student models to match larger teachers

- **Compilation:** Use TensorRT, ONNX, or similar for optimized inference

### 9.7.3 Quality Control and Debugging

**Common Issues and Diagnostic Approaches**

**Quality Assessment Strategies**

**Automated evaluation metrics:**

- **CLIP Score:** Measures text-image alignment for text conditioning

| Problem | Likely Causes | Solutions |
|---------|---------------|-----------|
| Poor conditioning adherence | Low guidance scale; poor conditioning data; insufficient model capacity | Increase guidance scale; audit data quality; check model architecture |
| Oversaturated or artificial results | Guidance scale too high; conditioning conflicts; poor training data diversity | Reduce guidance scale; check for contradictory conditions; improve data coverage |
| Inconsistent generation quality | Conditioning preprocessing errors; distribution shifts; inadequate sampling | Validate preprocessing pipeline; check train/test distribution alignment; tune sampling parameters |
| Slow inference speed | Inefficient batching; redundant computations; suboptimal model size | Optimize batching strategies; implement caching; consider model compression |
| Training instability | Learning rate issues; batch size problems; conditioning dropout imbalance | Adjust learning rate schedule; tune batch size; balance conditional/unconditional training |

Table 9.2: Troubleshooting guide for conditional diffusion model training and deployment

- **Classification accuracy:** Use pre-trained classifiers to verify conditioning adherence

- **Feature similarity:** Compare generated and target features in embedding space

- **Diversity metrics:** Measure variety within conditioned generations

**Human evaluation protocols:**

- **Conditioning adherence:** How well do results match the intended condition?

- **Image quality:** Overall visual quality and realism assessment

- **Preference ranking:** Compare different conditioning approaches

- **Use case validation:** Test with actual intended users and workflows

## 9.8 The Evolution of Conditioning: Past, Present, and Future

### 9.8.1 Historical Perspective: How We Got Here

**The journey from simple to sophisticated:**

1. **Early days (2020-2021):** Simple class conditioning, limited flexibility

2. **Classifier guidance era (2021-2022):** External control, computational overhead

3. **CFG revolution (2022-2023):** Unified approach, practical deployment

4. **Multi-modal age (2023-present):** Complex conditioning, real-world applications

**Key breakthrough moments:**

- **DDPM with class conditioning:** Proof that conditioning was possible

- **Classifier guidance paper:** Showed external control could work

- **CFG introduction:** Eliminated classifier dependence

- **CLIP integration:** Enabled natural language conditioning

- **ControlNet:** Made visual control practical and robust

### 9.8.2 Current State: The Modern Landscape

**Production systems powered by advanced conditioning:**

- **Stable Diffusion:** CFG with CLIP text conditioning and ControlNet extensions

- **Midjourney:** Sophisticated prompt understanding and style control

- **DALL-E 3:** Advanced text comprehension and safety conditioning

- **Adobe Firefly:** Commercial-grade conditioning with copyright awareness

  **Emerging conditioning paradigms:**

- **Compositional conditioning:** Combine multiple objects and relationships

- **Temporal conditioning:** Video generation with motion and timing control

- **3D conditioning:** Spatial and geometric constraints for 3D generation

- **Interactive conditioning:** Real-time user feedback and iterative refinement

### 9.8.3 Future Directions: What's Coming Next

**Research frontiers:**

- **Few-shot conditioning:** Learn new concepts from minimal examples

- **Personalization:** Adapt models to individual users and preferences

- **Hierarchical control:** Multi-level conditioning from global to fine-grained

- **Safety conditioning:** Built-in mechanisms for responsible generation

**Technical challenges being addressed:**

- **Conditioning conflicts:** How to handle contradictory or impossible conditions

- **Efficiency improvements:** Reducing computational overhead of conditioning

- **Training stability:** Making conditioning training more robust and predictable

- **Evaluation methods:** Better metrics for assessing conditioning quality

## 9.9 Summary: Mastering Conditional Generation

### 9.9.1 Key Takeaways

**Fundamental principles mastered:**

- **Conditioning spectrum:** From simple class labels to complex multi-modal control

- **Trade-off management:** Balancing quality, speed, flexibility, and implementation complexity

- **Architectural integration:** How to effectively inject conditioning information

- **Training strategies:** Best practices for robust conditional model training

**Practical skills developed:**

- **Method selection:** Choosing the right conditioning approach for specific use cases

- **Implementation optimization:** Efficient training and inference strategies

- **Quality control:** Debugging and improving conditional generation systems

- **Multi-modal integration:** Combining different types of conditioning effectively

**Impact on the field:**

- **Democratization:** Made controllable generation accessible to non-experts

- **Commercial viability:** Enabled practical applications and business models

- **Creative empowerment:** Gave artists and designers powerful new tools

- **Research acceleration:** Provided foundation for next-generation techniques

### 9.9.2 The Road Ahead: Continuing Your Journey

**Next steps in your learning:**

1. **Hands-on implementation:** Build your own conditional diffusion models

2. **Advanced techniques:** Explore personalization, few-shot learning, and safety

3. **Domain application:** Apply conditioning to your specific use cases

4. **Research contributions:** Identify and solve remaining challenges

**Key principles to remember:**

- **Start simple:** Begin with basic conditioning before adding complexity

- **Validate early:** Test conditioning quality throughout development

- **Think holistically:** Consider the entire pipeline from data to deployment

- **Stay current:** The field evolves rapidly—keep learning and adapting

**The bigger picture:** Conditional generation represents more than just a technical achievement—it's a fundamental shift toward AI systems that can understand and respond to human intent. The techniques we've explored here form the foundation for the next generation of creative AI tools, from personalized content generation to interactive design assistants.

As we continue to push the boundaries of what's possible with diffusion models, the principles of effective conditioning—clear communication between human intent and machine capability—will remain at the heart of creating AI systems that truly serve human creativity and productivity.

**Your role in this evolution:** Whether you're building the next breakthrough application, solving domain-specific challenges, or pushing the research frontier, you now have the foundational knowledge to contribute meaningfully to this rapidly evolving field. The future of conditional generation is limited only by our imagination and our ability to effectively bridge the gap between human intent and machine capability.

**Final thought:** The journey from unconditional random generation to precisely controlled, intent-driven creation represents one of the most significant advances in making AI truly useful for human creativity. You're now equipped to be part of writing the next chapter of this story.

## 9.10 Exercises

1. **Class-Conditional Implementation:**

   (a) Implement a class-conditional diffusion model for CIFAR-10

   (b) Compare different conditioning injection strategies (concatenation vs. embedding addition)

   (c) Analyze generation quality across different classes

   (d) Experiment with conditioning dropout rates during training

2. **Classifier Guidance Study:**

   (a) Train a classifier on noisy CIFAR-10 images across different timesteps

   (b) Implement classifier guidance sampling

   (c) Compare guidance with different classifier architectures and training procedures

   (d) Analyze the effect of guidance scale on sample quality and diversity

3. **CFG Implementation and Analysis:**

(a) Implement classifier-free guidance training and sampling

(b) Experiment with different unconditional dropout rates (5

(c) Compare CFG with classifier guidance on the same dataset

(d) Analyze computational costs and memory usage of both approaches

4. **Text Conditioning Project:**

   (a) Implement a simple text-to-image model using cross-attention

   (b) Experiment with different text encoders (CLIP vs. sentence transformers)

   (c) Develop and test prompt engineering strategies

   (d) Evaluate text-image alignment using CLIP scores

5. **Guidance Scale Optimization:**

   (a) Implement dynamic guidance scheduling (different $\omega$ per timestep)

   (b) Compare constant vs. scheduled guidance approaches

   (c) Develop heuristics for automatic guidance scale selection

   (d) Analyze the relationship between prompt complexity and optimal guidance

6. **Multi-Modal Conditioning:**

   (a) Implement a system that combines text and image conditioning

   (b) Experiment with different fusion strategies

   (c) Handle conflicting conditioning signals gracefully

   (d) Develop evaluation metrics for multi-modal adherence

7. **Controllable Generation Application:**

   (a) Implement an image inpainting system using conditional diffusion

   (b) Add support for different types of masks and conditioning

   (c) Develop an interactive interface for iterative refinement

   (d) Compare different conditioning strategies for editing applications

8. **Evaluation Framework:**

   (a) Implement comprehensive evaluation metrics for conditional generation

   (b) Design human evaluation protocols for assessing conditioning adherence

   (c) Compare quantitative metrics with human judgments

   (d) Develop domain-specific evaluation criteria for your application

**Practical Resources:**

- Hugging Face Diffusers: Production implementations of CFG and ControlNet

- Stability AI Stable Diffusion: Open-source text-to-image with CFG

- CompVis/Latent Diffusion: Research codebase for latent diffusion models

- ControlNet Official Implementation: Advanced image conditioning techniques

This chapter has demonstrated how conditioning transforms diffusion models from interesting research artifacts into powerful tools for creative and practical applications. The techniques covered here represent the foundation of modern controllable generation and continue to drive innovation in AI-assisted content creation across industries.

# Chapter 10

# Latent Diffusion Models: Scaling to High Resolution

## 10.1 Introduction: Why Latent Diffusion?

In previous chapters, we've explored diffusion models that operate directly in the pixel space of images. While this approach has proven successful for generating high-quality samples, it faces significant challenges when scaling to high-resolution images that are practical for real-world applications.

**Core challenges of pixel-space diffusion:**

**1. Computational expense:** High-resolution images (e.g., 512×512 or 1024×1024) contain hundreds of thousands to millions of pixels. Applying diffusion directly in this space requires enormous computational resources:

- Each denoising step operates on the full image resolution

- Memory requirements scale quadratically with image dimensions

- Training and inference become prohibitively expensive for consumer hardware

**2. Low-level noise modeling:** Pixel-space diffusion models often spend considerable capacity modeling low-level details that may not contribute to semantic quality:

- High-frequency noise patterns in individual pixels

- Imperceptible variations that don't affect human perception

- Fine-grained textures that could be handled more efficiently

**3. Perceptual inefficiency:** Human perception is not equally sensitive to all aspects of an image:

- We prioritize semantic content over pixel-perfect accuracy

- Perceptual quality often matters more than L2 reconstruction error

- Compression artifacts in semantically unimportant regions are often acceptable

**The latent diffusion solution:** Latent Diffusion Models (LDMs) address these challenges by introducing a two-stage approach:

1. **Compression stage**: Learn a powerful autoencoder that maps high-resolution images to a compact latent representation

2. **Generation stage**: Perform diffusion in the learned latent space, which is much more computationally efficient

**Key benefits of this approach:**

- **Computational efficiency**: Operating in a lower-dimensional latent space dramatically reduces computational costs

- **Semantic focus**: The latent space naturally emphasizes semantically meaningful features

- **Perceptual quality**: By using perceptual losses during autoencoder training, the model optimizes for human-relevant image quality

- **Scalability**: Enables high-resolution generation on consumer hardware

**Trade-offs and considerations:** While latent diffusion offers significant advantages, it also introduces new considerations:

- **Two-stage training**: Requires careful design and training of both the autoencoder and diffusion components

- **Compression artifacts**: The autoencoder may introduce artifacts that the diffusion model must work with

- **Limited fine detail**: Some pixel-level details may be lost in the compression process

**Chapter roadmap:** We'll begin by exploring the design of latent spaces and autoencoder architectures. Then we'll adapt diffusion processes to work in latent space, including training procedures and inference pipelines. We'll examine the applications and results of latent diffusion, particularly focusing on Stable Diffusion as a prominent example. Finally, we'll discuss the broader impact and future directions of this approach.

## 10.2   Latent Space Design

The success of latent diffusion models critically depends on designing an effective latent space that balances compression efficiency with semantic preservation. This section explores the key components and design decisions for creating robust latent representations.

## 10.2.1  Autoencoder Architecture

The autoencoder forms the foundation of latent diffusion, consisting of an encoder that compresses images and a decoder that reconstructs them.

**Encoder design:** The encoder $E : \mathbb{R}^{H \times W \times C} \to \mathbb{R}^{h \times w \times c}$ maps high-resolution images to compact latent representations:

$$\mathbf{z} = E(\mathbf{x}) \tag{10.1}$$

where the compression factor is typically $f = \frac{H \times W}{h \times w}$, commonly ranging from 4 to 16.

**Decoder design:** The decoder $D : \mathbb{R}^{h \times w \times c} \to \mathbb{R}^{H \times W \times C}$ reconstructs images from latent codes:

$$\hat{\mathbf{x}} = D(\mathbf{z}) \tag{10.2}$$

The goal is to minimize the reconstruction error: $\|\mathbf{x} - D(E(\mathbf{x}))\|$ while maintaining semantic fidelity.

**Common architectural choices:**

**1. Convolutional encoder-decoder:**

- Uses strided convolutions for downsampling in the encoder

- Employs transposed convolutions or upsampling for reconstruction in the decoder

- Includes skip connections to preserve fine details

- Often incorporates residual blocks for better gradient flow

**2. Variational autoencoders (VAE):**

- Introduces probabilistic latent representations: $\mathbf{z} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\sigma}^2)$

- Includes KL regularization to ensure smooth latent space: $\mathcal{L}_{KL} = D_{KL}(q(\mathbf{z}|\mathbf{x})\|\mathcal{N}(0, \mathbf{I}))$

- Provides better interpolation properties and latent space structure

- May sacrifice some reconstruction quality for regularization

**3. Patch-based encoders:**

- Process images as sequences of patches (similar to Vision Transformers)

- Can achieve higher compression ratios

- Particularly effective for capturing long-range dependencies

- May require more sophisticated positional encoding schemes

---

**Algorithm 13** Autoencoder Training

---

**Input:** Dataset $\{(\mathbf{x}^{(i)})\}$, encoder $E$, decoder $D$ not converged Sample batch $\{\mathbf{x}\}$ from dataset Encode: $\mathbf{z} = E(\mathbf{x})$ Decode: $\hat{\mathbf{x}} = D(\mathbf{z})$ Compute reconstruction loss: $\mathcal{L}_{rec} = \|\mathbf{x} - \hat{\mathbf{x}}\|^2$ Compute perceptual loss: $\mathcal{L}_{perc} = \|\phi(\mathbf{x}) - \phi(\hat{\mathbf{x}})\|^2$ Total loss: $\mathcal{L} = \lambda_{rec}\mathcal{L}_{rec} + \lambda_{perc}\mathcal{L}_{perc}$ Update parameters: $\theta \leftarrow \theta - \eta\nabla_\theta\mathcal{L}$

---

### 10.2.2 Perceptual Losses

Traditional reconstruction losses like L2 (mean squared error) optimize for pixel-wise accuracy but may not align with human perceptual quality. Perceptual losses address this by comparing features in learned representation spaces.

**Feature-based perceptual loss:** Using a pretrained network $\phi$ (commonly VGG or similar), the perceptual loss compares intermediate features:

$$\mathcal{L}_{perc} = \sum_l \lambda_l \|\phi_l(\mathbf{x}) - \phi_l(\hat{\mathbf{x}})\|^2 \tag{10.3}$$

where $\phi_l$ denotes features from layer $l$ of the pretrained network.

**Benefits of perceptual losses:**

- **Semantic preservation**: Prioritizes semantically meaningful features over pixel-level details

- **Texture quality**: Better preservation of textures and patterns that matter for visual quality

- **Robustness**: Less sensitive to small spatial misalignments

- **Human alignment**: Better correlation with human perceptual judgments

**Common perceptual loss variants:**

- **VGG loss**: Uses features from pretrained VGG networks, typically from multiple layers

- **LPIPS**: Learned Perceptual Image Patch Similarity, trained specifically for perceptual distance

- **Feature matching**: Compares statistics of feature distributions rather than direct features

**Adversarial losses:** Many successful latent diffusion models also incorporate adversarial training:

$$\mathcal{L}_{adv} = \mathbb{E}_{\mathbf{x}}[\log(D_{disc}(\mathbf{x}))] + \mathbb{E}_{\mathbf{z}}[\log(1 - D_{disc}(D(\mathbf{z})))] \tag{10.4}$$

where $D_{disc}$ is a discriminator network that distinguishes real images from reconstructions.

### 10.2.3 Latent Properties

The design of the latent space significantly impacts the performance of the overall system. Understanding the properties and trade-offs is crucial for effective implementation.

**Benefits of well-designed latent spaces:**

**1. Dimensionality reduction:**

- Typical compression factors of 4-16× reduce computational costs dramatically

- A 512×512×3 image becomes a 64×64×4 latent (16× compression)

- Memory usage scales proportionally to the compression factor

**2. Semantic abstraction:**

- Latent representations focus on semantically meaningful features

- Reduces model capacity spent on imperceptible details

- Enables more efficient learning of high-level structure

**3. Smooth manifold structure:**

- Well-trained latent spaces exhibit smooth interpolation properties

- Similar latent codes correspond to semantically similar images

- Enables meaningful latent space arithmetic and editing

**Challenges and considerations:**
**1. Compression artifacts:**

- Lossy compression may introduce artifacts that the diffusion model inherits

- Balance between compression ratio and reconstruction quality

- Artifacts may manifest as recurring patterns or unrealistic textures

**2. Training instability:**

- Complex objective functions with multiple loss terms require careful tuning

- Adversarial training can be unstable and require specialized techniques

- Stage-wise training helps but adds complexity to the pipeline

**3. Decoder bias:**

- The decoder's inductive biases influence the final generated images

- May limit the diversity or style of possible outputs

- Requires careful architecture design to maintain flexibility

**Latent space evaluation metrics:**

- **Reconstruction quality**: PSNR, SSIM, and perceptual metrics on test images

- **Compression efficiency**: Rate-distortion analysis

- **Semantic preservation**: Classification accuracy on encoded-decoded images

- **Interpolation quality**: Smoothness and meaningfulness of latent interpolations

## 10.3 Diffusion in Latent Space

Once we have established a robust latent representation, we can adapt the diffusion process to operate in this compressed space. This adaptation requires careful consideration of the distributional properties and scaling factors of the latent space.

### 10.3.1 Adapting DDPM to Latents

The core diffusion process remains conceptually the same, but now operates on latent codes rather than raw images.

**Forward process in latent space:** Instead of adding noise directly to images, we add noise to latent representations:

$$\mathbf{z}_t = \sqrt{\bar{\alpha}_t}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon} \tag{10.5}$$

where $\mathbf{z}_0 = E(\mathbf{x}_0)$ is the latent encoding of the original image.

**Reverse process in latent space:** The learned denoising network $\boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t)$ predicts noise in latent space:

$$p_\theta(\mathbf{z}_{t-1}|\mathbf{z}_t) = \mathcal{N}(\mathbf{z}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{z}_t, t), \sigma_t^2\mathbf{I}) \tag{10.6}$$

**Training objective:** The standard DDPM objective translates directly to latent space:

$$\mathcal{L} = \mathbb{E}_{\mathbf{z}_0, t, \boldsymbol{\epsilon}}\left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t)\|^2\right] \tag{10.7}$$

where $\mathbf{z}_t = \sqrt{\bar{\alpha}_t}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$.

**Key architectural considerations:**

- The U-Net architecture is adapted to the latent space dimensions

- Input/output channels match the latent space dimensionality

- Spatial resolution corresponds to the compressed latent size

- Attention mechanisms operate on latent spatial locations

### 10.3.2 Rescaling and Normalization

Latent spaces often have different statistical properties than the standard Gaussian distributions assumed by DDPM. Proper scaling and normalization are crucial for stable training and sampling.

**Latent statistics analysis:** Before training the diffusion model, analyze the distribution of latent codes:

- Compute mean and standard deviation across the dataset: $\boldsymbol{\mu}_{latent}, \boldsymbol{\sigma}_{latent}$

- Check for channel-wise variations in scale

- Identify any systematic biases or outliers

**Normalization strategies:**
**1. Z-score normalization:**

$$\mathbf{z}_{norm} = \frac{\mathbf{z} - \boldsymbol{\mu}_{latent}}{\boldsymbol{\sigma}_{latent}} \tag{10.8}$$

**2. Scaling to unit variance:**

$$\mathbf{z}_{scaled} = \frac{\mathbf{z}}{\boldsymbol{\sigma}_{latent}} \tag{10.9}$$

**3. Component-wise scaling:** Apply different scaling factors to different latent channels if they have significantly different ranges.

**Inference considerations:** During sampling, ensure that the scaling is properly reversed:

1. Sample from the diffusion model in normalized latent space

2. Rescale back to the original latent distribution

3. Decode using the pretrained decoder

### 10.3.3   Training Procedure

Latent diffusion models require a carefully orchestrated two-stage training process to achieve optimal results.

**Stage 1: Autoencoder training**

---

**Algorithm 14** Stage 1: Autoencoder Training

---

**Input:** Image dataset $\{\mathbf{x}^{(i)}\}$ Initialize encoder $E$ and decoder $D$ autoencoder not converged Sample batch $\{\mathbf{x}\}$ from dataset Encode: $\mathbf{z} = E(\mathbf{x})$ Decode: $\hat{\mathbf{x}} = D(\mathbf{z})$ Compute reconstruction loss: $\mathcal{L}_{rec} = \|\mathbf{x} - \hat{\mathbf{x}}\|^2$ Compute perceptual loss: $\mathcal{L}_{perc} = \sum_l \|\phi_l(\mathbf{x}) - \phi_l(\hat{\mathbf{x}})\|^2$ Optionally compute adversarial loss: $\mathcal{L}_{adv}$ Total loss: $\mathcal{L} = \lambda_{rec}\mathcal{L}_{rec} + \lambda_{perc}\mathcal{L}_{perc} + \lambda_{adv}\mathcal{L}_{adv}$ Update $E$ and $D$: $\theta_{E,D} \leftarrow \theta_{E,D} - \eta\nabla\mathcal{L}$ Freeze encoder and decoder parameters

---

**Stage 2: Diffusion model training**

---

**Algorithm 15** Stage 2: Latent Diffusion Training

---

**Input:** Image dataset $\{\mathbf{x}^{(i)}\}$, frozen encoder $E$, diffusion model $\boldsymbol{\epsilon}_\theta$ Compute latent dataset: $\{\mathbf{z}^{(i)} = E(\mathbf{x}^{(i)})\}$ Analyze latent statistics and determine scaling factors diffusion model not converged Sample batch $\{\mathbf{z}\}$ from latent dataset Normalize: $\mathbf{z} \leftarrow \text{normalize}(\mathbf{z})$ Sample timesteps $t \sim \text{Uniform}(1, T)$ Sample noise $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ Compute noisy latents: $\mathbf{z}_t = \sqrt{\bar{\alpha}_t}\mathbf{z} + \sqrt{1 - \bar{\alpha}_t}\boldsymbol{\epsilon}$ Predict noise: $\hat{\boldsymbol{\epsilon}} = \boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t)$ Compute loss: $\mathcal{L} = \|\boldsymbol{\epsilon} - \hat{\boldsymbol{\epsilon}}\|^2$ Update diffusion model: $\theta \leftarrow \theta - \eta\nabla\mathcal{L}$

---

**Training considerations:**

- **Computational efficiency**: Stage 1 requires full-resolution processing, while Stage 2 operates in compressed space

- **Dataset preparation**: Pre-compute latent encodings to speed up diffusion training

- **Memory management**: Latent diffusion enables training larger models due to reduced memory requirements

- **Validation**: Monitor both reconstruction quality and generation quality throughout training

### 10.3.4 Inference Pipeline

The inference process combines sampling in latent space with decoding to produce final images.

---
**Algorithm 16** Latent Diffusion Inference

---
**Input:** Trained models $\epsilon_\theta$, $D$, number of steps $S$ Sample initial noise: $\mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ $t = T, T-1, \ldots, 1$ Predict noise: $\hat{\epsilon} = \epsilon_\theta(\mathbf{z}_t, t)$ Compute mean: $\boldsymbol{\mu} = \frac{1}{\sqrt{\alpha_t}}\left(\mathbf{z}_t - \frac{1-\alpha_t}{\sqrt{1-\bar{\alpha}_t}}\hat{\epsilon}\right)$ $t > 1$ Sample noise: $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ $\mathbf{z}_{t-1} = \boldsymbol{\mu} + \sigma_t \mathbf{z}$ $\mathbf{z}_0 = \boldsymbol{\mu}$ Rescale latent: $\mathbf{z}_0 \leftarrow \text{rescale}(\mathbf{z}_0)$ Decode to image: $\mathbf{x}_0 = D(\mathbf{z}_0)$ **Return:** $\mathbf{x}_0$

---

**Inference optimizations:**

- **DDIM sampling**: Use deterministic sampling for faster generation

- **Fewer steps**: Latent diffusion often works well with 20-50 steps instead of 1000

- **Batched generation**: Generate multiple samples simultaneously

- **Mixed precision**: Use float16 for memory efficiency during inference

## 10.4 Applications and Results

Latent diffusion models have enabled a new class of applications that were previously computationally prohibitive. This section explores the key applications and analyzes their performance characteristics.

### 10.4.1 High-Resolution Image Generation

Latent diffusion models excel at generating high-resolution images that would be impractical with pixel-space approaches.

**Resolution scaling:**

- **Standard diffusion**: Limited to 256×256 on consumer hardware

- **Latent diffusion**: Enables 512×512, 1024×1024, and beyond

- **Compression strategy**: 8× compression allows 4× larger images with same computational cost

**Quality characteristics:**

- **Sharp details**: Perceptual training leads to crisp, detailed outputs

- **Semantic coherence**: Latent space operation maintains global structure

- **Texture quality**: Effective at generating realistic textures and patterns

**Architecture adaptations for high resolution:**

- **Progressive generation**: Some models generate at multiple scales

- **Attention efficiency**: Sparse attention patterns for large latent maps

- **Multi-scale features**: U-Net designs that capture different scales effectively

## 10.4.2 Computational Savings

The computational advantages of latent diffusion are substantial and enable practical deployment.

**Memory requirements:**

| Method | Image Size | Memory per Sample |
|---|---|---|
| Pixel Diffusion | 512×512×3 | 3.1 MB |
| Latent Diffusion (8× compression) | 64×64×4 | 0.065 MB |
| **Reduction Factor** | | **48×** |

Table 10.1: Memory comparison between pixel and latent diffusion

**Training time comparison:**

- **Pixel-space DDPM**: 1000 GPU-hours for high-quality 256×256 model

- **Latent diffusion**: 200 GPU-hours for comparable quality 512×512 model

- **Amortized cost**: Autoencoder training cost is amortized across multiple diffusion models

**Inference speed:**

- **Sampling speed**: 10-50× faster due to smaller intermediate representations

- **Memory efficiency**: Enables larger batch sizes and longer sequences

- **Hardware accessibility**: Consumer GPUs can handle high-resolution generation

### 10.4.3 Sample Quality

Latent diffusion models achieve impressive sample quality, often surpassing pixel-space approaches on perceptual metrics.

**Quantitative evaluation:**

- **FID scores**: Often achieve lower (better) FID than pixel-space models

- **IS scores**: Comparable or superior Inception Scores

- **Perceptual metrics**: Excel on LPIPS and other perceptual distance measures

**Qualitative characteristics:**

- **Sharpness**: Generated images appear sharp and detailed

- **Realism**: High perceptual realism due to perceptual training

- **Diversity**: Maintain good sample diversity despite compression

**Limitations and artifacts:**

- **Decoder artifacts**: May exhibit recurring patterns from the decoder

- **Compression losses**: Some fine details may be lost in compression

- **Hallucinated details**: May generate plausible but incorrect fine details

- **Training data bias**: Inherits biases from both autoencoder and diffusion training data

### 10.4.4 Stable Diffusion Overview

Stable Diffusion represents one of the most successful implementations of latent diffusion, demonstrating the practical impact of this approach.

**Modular architecture:** Stable Diffusion consists of three main components, as referenced in the provided material:

**1. VAE encoder-decoder:**

- Pretrained variational autoencoder for image compression

- $8\times$ spatial compression ($512\times512 \rightarrow 64\times64$)

- 4-channel latent space representation

- Trained with perceptual and adversarial losses

**2. U-Net denoising model:**

- Operates in the $64\times64\times4$ latent space

- Incorporates cross-attention for text conditioning

- Uses classifier-free guidance for controllable generation

- Efficient architecture enabling fast sampling

**3. CLIP text encoder:**

- Encodes text prompts into embeddings

- Provides rich semantic conditioning information

- Enables fine-grained control over generated content

- Pre-trained on large-scale vision-language data

**Cross-attention conditioning:** As described in the lecture material, Stable Diffusion uses cross-attention to incorporate text conditioning:

$$\text{Query } \mathbf{Q} = \mathbf{F}_{\text{spatial}}\mathbf{W}_Q \tag{10.10}$$
$$\text{Key } \mathbf{K} = \mathbf{F}_{\text{text}}\mathbf{W}_K \tag{10.11}$$
$$\text{Value } \mathbf{V} = \mathbf{F}_{\text{text}}\mathbf{W}_V \tag{10.12}$$
$$\text{Attention } \mathbf{A} = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V} \tag{10.13}$$

where $\mathbf{F}_{\text{spatial}}$ represents spatial features from the U-Net and $\mathbf{F}_{\text{text}}$ represents encoded text features.

**Open-source impact:**

- **Accessibility**: Runs on consumer hardware (8GB+ GPU memory)

- **Ecosystem**: Spawned numerous extensions and applications

- **Research catalyst**: Enabled rapid experimentation and development

- **Community contributions**: Active development of improvements and variants

**Extension ecosystem:** The modular design of Stable Diffusion has enabled various extensions:

- **ControlNet**: Adds spatial conditioning capabilities (edges, depth, pose)

- **LoRA**: Low-rank adaptation for efficient fine-tuning

- **DreamBooth**: Personalization techniques for custom subjects

- **Inpainting**: Specialized models for image completion tasks

## 10.5   Zero-Shot Applications

One of the remarkable properties of well-trained latent diffusion models is their ability to perform various image editing and manipulation tasks without additional training. This section explores key zero-shot applications that demonstrate the flexibility of the latent diffusion approach.

227

## 10.5.1 SDEdit: Stochastic Differential Editing

SDEdit, as presented in the lecture material, provides a simple yet powerful approach to image editing through controlled noise injection and denoising.

**Core concept:** SDEdit performs image editing by partially corrupting an input image and then applying the reverse diffusion process:

---
**Algorithm 17** SDEdit Process
---
    **Input:** Input image $\mathbf{x}_0$, edit strength $t_{edit}$, pretrained diffusion model Encode to latent: $\mathbf{z}_0 = E(\mathbf{x}_0)$ Add noise: $\mathbf{z}_{t_{edit}} = \sqrt{\bar{\alpha}_{t_{edit}}}\mathbf{z}_0 + \sqrt{1 - \bar{\alpha}_{t_{edit}}}\boldsymbol{\epsilon}$ $t = t_{edit}, t_{edit} - 1, \ldots, 1$ Apply reverse diffusion step: $\mathbf{z}_{t-1} = \text{DDPM\_step}(\mathbf{z}_t, t)$ Decode result: $\mathbf{x}_{edited} = D(\mathbf{z}_0)$ **Return:** $\mathbf{x}_{edited}$

---

**Realism vs. faithfulness trade-off:** As noted in the lecture material, SDEdit exhibits a fundamental trade-off:

- **Lower** $t_{edit}$: More faithful to input but less realistic/natural

- **Higher** $t_{edit}$: More realistic but deviates more from the original input

- **Optimal range**: Typically $t_{edit} \in [0.2T, 0.8T]$ for good balance

**Applications:**

- **Sketch-to-image**: Convert rough sketches into realistic images

- **Style transfer**: Modify artistic style while preserving content

- **Enhancement**: Improve image quality and details

- **Variation generation**: Create variations of existing images

## 10.5.2 RePaint: Inpainting with Diffusion

RePaint demonstrates how latent diffusion models can perform image inpainting without specialized training, as detailed in the lecture material.

**Core algorithm:** RePaint combines known background regions with generated foreground content at each denoising step:

---
**Algorithm 18** RePaint Inpainting
---
    **Input:** Image $\mathbf{x}$, mask $\mathbf{M}$, pretrained diffusion model Encode: $\mathbf{z}_{bg} = E(\mathbf{x})$ (background image) Initialize: $\mathbf{z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ (foreground noise) $t = T, T - 1, \ldots, 1$ **Step 1:** Denoise foreground: $\mathbf{z}_{fg}^{t-1} = \text{DDPM\_step}(\mathbf{z}_t, t)$ **Step 2:** Add noise to background: $\mathbf{z}_{bg}^{t-1} = \sqrt{\bar{\alpha}_{t-1}}\mathbf{z}_{bg} + \sqrt{1 - \bar{\alpha}_{t-1}}\boldsymbol{\epsilon}$ **Step 3:** Combine using mask: $\mathbf{z}_{t-1} = \mathbf{M} \odot \mathbf{z}_{bg}^{t-1} + (1 - \mathbf{M}) \odot \mathbf{z}_{fg}^{t-1}$ Decode: $\mathbf{x}_{inpainted} = D(\mathbf{z}_0)$ **Return:** $\mathbf{x}_{inpainted}$

---

**Key insights:**

- **Consistent background**: The background is re-noised at each step to maintain consistency with the noise level

- **Seamless blending**: The masking operation ensures smooth transitions between known and generated regions

- **Iterative refinement**: Multiple forward-backward steps may be used for better quality

- **No additional training**: Uses the pretrained unconditional model without modification

**Advantages over specialized inpainting models:**

- **No retraining required**: Works with any pretrained diffusion model

- **Flexible mask shapes**: Handles arbitrary mask geometries

- **High quality**: Produces coherent and realistic inpainting results

- **Controllable**: Edit strength can be adjusted through the number of replay steps

### 10.5.3   Other Zero-Shot Applications

The flexibility of latent diffusion enables numerous other zero-shot applications:
**Image interpolation:**

- Encode two images to latent space: $\mathbf{z}_A = E(\mathbf{x}_A)$, $\mathbf{z}_B = E(\mathbf{x}_B)$

- Interpolate in latent space: $\mathbf{z}_{interp} = (1 - \lambda)\mathbf{z}_A + \lambda\mathbf{z}_B$

- Decode interpolated result: $\mathbf{x}_{interp} = D(\mathbf{z}_{interp})$

- Optionally refine through partial denoising

**Image outpainting:**

- Extend images beyond their original boundaries

- Use masking similar to inpainting but for extended canvas

- Maintain consistency with the original image content

**Super-resolution:**

- Encode low-resolution image to latent space

- Apply denoising with modified noise schedules

- Decode to higher resolution using trained decoder

## 10.6   Advanced Techniques and Extensions

The success of latent diffusion has sparked numerous extensions and improvements that further enhance capabilities and efficiency.

## 10.6.1 ControlNet: Spatial Conditioning

ControlNet, as mentioned in the lecture material, provides a method to add spatial conditioning to pretrained diffusion models without full retraining.

**Architecture design:** ControlNet creates a trainable copy of the U-Net encoder that processes control inputs:

- **Frozen backbone**: Original U-Net parameters remain unchanged

- **Trainable copy**: Duplicate encoder processes control information

- **Zero convolutions**: Special initialization ensuring gradual learning

- **Additive combination**: Control features are added to backbone features

**Zero convolution technique:** As noted in the lecture material, zero convolutions ensure stable training:

$$\mathbf{y}_c = F(\mathbf{x}; \Theta) + Z(F(\mathbf{x} + Z(\mathbf{c}; \theta_1), \theta_2); \theta_3), \theta_4) \tag{10.14}$$

where $Z(\cdot; \theta)$ represents zero convolution layers initialized with zero weights.

**Control modalities:**

- **Canny edges**: Structural control through edge detection

- **Depth maps**: 3D spatial structure conditioning

- **Human pose**: Pose-based character generation

- **Segmentation**: Semantic layout control

- **Scribbles**: Hand-drawn sketch conditioning

## 10.6.2 Low-Rank Adaptation (LoRA)

LoRA provides an efficient method for fine-tuning large diffusion models with minimal parameter overhead.

**Core principle:** LoRA introduces trainable low-rank matrices that modify the behavior of pretrained layers:

$$\mathbf{W}_{modified} = \mathbf{W}_{frozen} + \mathbf{BA} \tag{10.15}$$

where $\mathbf{A} \in \mathbb{R}^{r \times d}$ and $\mathbf{B} \in \mathbb{R}^{d \times r}$ with $r \ll d$.

**Benefits:**

- **Parameter efficiency**: Only train 1

- **Modular**: LoRA modules can be combined and swapped

- **Fast training**: Reduces training time and memory requirements

- **Preserves quality**: Maintains base model capabilities while adding new skills

**Applications:**

- **Style adaptation**: Learn new artistic styles

- **Character consistency**: Generate consistent characters across images

- **Concept learning**: Incorporate new visual concepts

- **Domain adaptation**: Adapt to specific image domains

### 10.6.3 DreamBooth: Subject-Driven Generation

DreamBooth enables personalization of text-to-image models for specific subjects or objects.

**Method overview:**

- **Few-shot learning**: Requires only 3-5 images of target subject

- **Unique identifier**: Binds subject to rare token in vocabulary

- **Class-specific prior**: Maintains general knowledge of object class

- **Fine-tuning**: Updates model weights while preserving capabilities

**Training objective:**

$$\mathcal{L} = \mathbb{E}_{\mathbf{z},\boldsymbol{\epsilon},t}\left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{z}_t, t, P)\|^2\right] + \lambda \mathbb{E}_{\mathbf{z}',\boldsymbol{\epsilon}',t'}\left[\|\boldsymbol{\epsilon}' - \boldsymbol{\epsilon}_\theta(\mathbf{z}'_{t'}, t', P_{class})\|^2\right] \quad (10.16)$$

where $P$ is the personalized prompt and $P_{class}$ is the class-specific prompt.

## 10.7 Implementation Considerations and Best Practices

Successfully implementing latent diffusion models requires careful attention to numerous technical details and design choices.

### 10.7.1 Training Strategies

**Stage coordination:**

- **Autoencoder convergence**: Ensure reconstruction quality before diffusion training

- **Learning rate scheduling**: Use different schedules for different components

- **Batch size considerations**: Balance memory usage with training stability

- **Regularization**: Apply appropriate regularization to prevent overfitting

**Loss function tuning:**

- **Perceptual weight**: Typically $\lambda_{perc} = 0.1$ to $1.0$

- **Adversarial weight**: Start low ($\lambda_{adv} = 0.01$) and increase gradually

- **Reconstruction weight**: Balance pixel accuracy with perceptual quality

- **KL regularization**: Use $\beta = 10^{-6}$ to $10^{-4}$ for VAE training

### 10.7.2 Architecture Choices

**Compression factors:**

| Compression | Latent Size | Memory Savings | Quality Trade-off |
|---|---|---|---|
| 4× | 128×128 | 4× | Minimal loss |
| 8× | 64×64 | 16× | Good balance |
| 16× | 32×32 | 64× | Some detail loss |
| 32× | 16×16 | 256× | Significant loss |

Table 10.2: Compression factor trade-offs for 512×512 images

**Latent dimensions:**

- **Channel count**: Typically 3-8 channels for good reconstruction

- **Spatial resolution**: Balance between detail preservation and efficiency

- **Bottleneck design**: Consider spatial vs. channel compression trade-offs

### 10.7.3 Evaluation Metrics

**Autoencoder evaluation:**

- **PSNR/SSIM**: Traditional reconstruction metrics

- **LPIPS**: Perceptual similarity measures

- **FID**: Distributional similarity of reconstructions

- **User studies**: Human evaluation of reconstruction quality

**Generation evaluation:**

- **FID scores**: Standard metric for generative model quality

- **Inception Score**: Diversity and quality assessment

- **CLIP Score**: Text-image alignment for conditional models

- **Human evaluation**: Subjective quality and preference studies

## 10.8 Summary and Key Takeaways

Latent diffusion models represent a pivotal advancement in generative modeling that has made high-quality, high-resolution image generation accessible and practical.

**Core innovations:**

- **Two-stage approach**: Separate compression and generation for efficiency

- **Perceptual optimization**: Focus on human-relevant image quality

- **Computational efficiency**: Dramatic reduction in training and inference costs

- **Scalable architecture**: Enables generation at previously impractical resolutions

**Technical achievements:**

- **High-resolution generation**: Routine generation of 512×512+ images

- **Quality improvements**: Often superior perceptual quality compared to pixel-space models

- **Training efficiency**: Significantly reduced computational requirements

- **Hardware accessibility**: Consumer GPU deployment capabilities

**Practical impact:**

- **Stable Diffusion**: Democratized access to high-quality text-to-image generation

- **Extension ecosystem**: Enabled ControlNet, LoRA, DreamBooth, and other innovations

- **Zero-shot applications**: SDEdit, RePaint, and other editing capabilities without retraining

- **Research acceleration**: Faster experimentation and development cycles

**Key design principles:**

- **Modular architecture**: Clean separation between compression and generation

- **Perceptual alignment**: Optimize for human perception rather than pixel accuracy

- **Efficient representations**: Latent spaces that preserve semantic information

- **Flexible conditioning**: Cross-attention and other mechanisms for control

**Current limitations and future directions:**

- **Compression artifacts**: Ongoing work to improve autoencoder quality

- **Fine detail generation**: Balancing efficiency with detail preservation

- **Training complexity**: Simplifying the two-stage training process

- **Evaluation metrics**: Better metrics for perceptual quality assessment

**Broader significance:** Latent diffusion models have fundamentally changed the landscape of generative AI by making high-quality image generation accessible to researchers, developers, and end-users. The principles established in this work—efficient representations, perceptual optimization, and modular architectures—continue to influence developments across generative modeling.

The success of models like Stable Diffusion demonstrates how academic research can rapidly translate into practical tools that impact millions of users. The open-source nature of many latent diffusion implementations has fostered an ecosystem of innovation that continues to push the boundaries of what's possible in generative AI.

As we look toward the future, latent diffusion models serve as a foundation for even more advanced applications including video generation, 3D synthesis, and multimodal content creation. The core insights—that efficient representations and perceptual optimization can dramatically improve both quality and accessibility—remain relevant across these emerging domains.

## 10.9   Exercises

1. **Autoencoder Design and Implementation:**

   (a) Design and implement a convolutional autoencoder for image compression

   (b) Experiment with different compression factors ($4\times$, $8\times$, $16\times$)

   (c) Compare reconstruction quality using pixel loss vs. perceptual loss

   (d) Analyze the trade-offs between compression ratio and reconstruction fidelity

2. **Perceptual Loss Investigation:**

   (a) Implement VGG-based perceptual loss using features from multiple layers

   (b) Compare reconstruction results with L2 loss, perceptual loss, and combined losses

   (c) Analyze which image features are preserved or lost with different loss functions

   (d) Conduct user studies to evaluate perceptual quality vs. pixel accuracy

3. **Latent Space Analysis:**

   (a) Train an autoencoder and analyze the statistical properties of the latent space

   (b) Compute mean, variance, and distribution shapes for different latent channels

   (c) Implement and compare different normalization strategies

   (d) Evaluate the impact of normalization on diffusion model training

4. **Two-Stage Training Implementation:**

(a) Implement the complete two-stage training pipeline for latent diffusion

(b) Monitor reconstruction quality throughout autoencoder training

(c) Train a diffusion model in the learned latent space

(d) Compare results with end-to-end training approaches

5. **Computational Efficiency Analysis:**

(a) Measure memory usage and training time for pixel vs. latent diffusion

(b) Analyze the computational breakdown across different components

(c) Implement optimizations for inference speed and memory usage

(d) Evaluate the practical trade-offs for different hardware configurations

6. **Zero-Shot Application Development:**

(a) Implement SDEdit for image editing with controllable noise injection

(b) Develop RePaint for inpainting using pretrained models

(c) Create an image interpolation system using latent space operations

(d) Compare results with specialized models trained for each task

7. **ControlNet Implementation:**

(a) Implement a basic ControlNet architecture with zero convolutions

(b) Train control modules for edge-based or depth-based conditioning

(c) Evaluate the quality and controllability of the resulting system

(d) Experiment with different control modalities and training strategies

8. **Stable Diffusion Analysis:**

(a) Analyze the architecture and training approach of Stable Diffusion

(b) Implement cross-attention conditioning for text-to-image generation

(c) Experiment with different guidance scales and sampling strategies

(d) Evaluate the model's capabilities and limitations across different prompts

**Further Reading:**

- Rombach et al. (2022): "High-Resolution Image Synthesis with Latent Diffusion Models" - The foundational LDM paper

- Esser et al. (2021): "Taming Transformers for High-Resolution Image Synthesis" - VQGAN and perceptual training

- Zhang et al. (2023): "Adding Conditional Control to Text-to-Image Diffusion Models" - ControlNet methodology

- Hu et al. (2021): "LoRA: Low-Rank Adaptation of Large Language Models" - LoRA technique

- Ruiz et al. (2023): "DreamBooth: Fine Tuning Text-to-Image Diffusion Models for Subject-Driven Generation"

- Meng et al. (2022): "SDEdit: Guided Image Synthesis and Editing with Stochastic Differential Equations"

- Lugmayr et al. (2022): "RePaint: Inpainting using Denoising Diffusion Probabilistic Models"

**Practical Resources:**

- Hugging Face Diffusers: Production implementation of Stable Diffusion and variants

- CompVis/stable-diffusion: Original Stable Diffusion implementation and weights

- CompVis/latent-diffusion: Research implementation of the original LDM paper

- lllyasviel/ControlNet: Official ControlNet implementation and models

- cloneofsimo/lora: LoRA implementation for diffusion models

- XavierXiao/Dreambooth-Stable-Diffusion: DreamBooth training scripts

Latent diffusion models have fundamentally transformed the field of generative AI by demonstrating that efficient representations and perceptual optimization can dramatically improve both the quality and accessibility of generative models. The techniques and principles covered in this chapter continue to influence developments across computer vision, natural language processing, and multimodal AI systems.