

University of Duisburg-Essen
Faculty of Computer Science
Department of Software Engineering
Chair of Intelligent Systems

- Seminar-Work -
in the Study Programme
Computer Engineering (ISE) - Bachelor

Machine Learning
Generative Adversarial Networks (GAN)

Ullin Noe Yanou

Date: 20. June 2025

Acknowledgments

I would like to thank Prof. Dr.-Ing. Josef Pauli for his guidance throughout this seminar. His feedback helped refine both structure and clarity.

Additionally, I acknowledge the use of OpenAI's ChatGPT as a tool to help me interpret and rephrase complex parts of the scientific literature, especially regarding technical explanations in the GAN, cGAN, and WGAN papers.

Table of Contents

Acknowledgments	ii
1. Introduction	1
1.1. Motivation	1
1.2. Problem Statement	1
1.3. Overview	2
2. Generative Adversarial Networks	3
2.1. What is a GAN?	3
2.2. Visual Architecture of a GAN	3
2.3. The Minimax Game: Understanding the Formula	3
2.4. Training Algorithm	4
2.5. Why Use Logarithms?	6
2.6. Understanding the Generator's Learning	7
2.7. The Optimal Discriminator $D^*(x)$	7
2.8. How G Learns From D	8
2.9. Why Do We Use $k = 1$?	9
2.10. Training Evolution	11
2.11. Experiment: MNIST Digit Generation	12
2.12. Conclusion	14
3. Wasserstein GAN	15
3.1. Why Classic GANs Struggle	15
3.2. The Earth Mover Distance	15
3.3. WGAN Objective	17
3.4. WGAN Training Algorithm	18
3.5. Why the Critic (Not Discriminator)?	19
3.6. Critic Training Strategy and Empirical Behavior	19
3.7. Conclusion	22
4. Conditional GAN (cGAN)	23
5. Evaluation	25
5.1. Evaluation Criteria	25
5.2. Qualitative Comparison	25
5.3. Stability and Loss Curves	26
5.4. Fréchet Inception Distance (FID)	26
5.5. Summary Table	28

Table of Contents

5.6. Conclusion	28
6. Summary and Outlook	29
6.1. Summary of Key Insights	29
6.2. Outlook and Future Directions	29
A. Conditional Generative Adversarial Networks	31
A.1. Why Condition the Generation?	31
A.2. Architecture and Objective	31
A.3. Architecture of Conditional GANs	32
A.4. Benefits of Conditioning	33
A.5. Experiment: Digit Generation with Labels	33
A.6. Conclusion	34
List of Figures	37
List of Tables	39
Bibliography	41

1. Introduction

1.1. Motivation

Artificial intelligence has seen rapid progress in recent years, especially in its ability to generate realistic content—images, audio, and text. A key innovation in this area is the Generative Adversarial Network (GAN), introduced by Goodfellow et al. [1]. The concept is simple yet powerful: two neural networks compete. One generates fake data, the other tries to detect it.

This adversarial setup creates a feedback loop. As the generator learns to create more convincing outputs, the discriminator adapts to become more discerning. Over time, both improve. The outcome? Data that looks increasingly real. Such models are now used to generate photorealistic human faces, simulate physical processes, or create new samples of handwritten digits.

1.2. Problem Statement

Despite their promise, GANs are difficult to train. The generator (G) and the discriminator (D) must improve together—but balancing their progress is hard. Several issues arise:

- **Vanishing gradients:** If D becomes too accurate, G gets almost no feedback and stops learning.
- **Mode collapse:** G finds a few outputs that work and repeats them, losing diversity.
- **Instability:** Training may diverge or oscillate without producing useful results.

Researchers have proposed several improvements. Two approaches stand out: Conditional GANs (cGANs) [2], which introduce control via conditioning, and Wasserstein GANs (WGANs) [3], which stabilize learning through a new loss function based on the Earth Mover (Wasserstein-1) distance. This report focuses on GANs and WGANs. A short overview of cGANs is provided in the appendix.

1. Introduction

1.3. Overview

This report is organized as follows:

- **Chapter 1** introduces the motivation and problem statement behind generative models. It also presents a roadmap of the topics explored in this seminar.
- **Chapter 2** covers the fundamental ideas behind Generative Adversarial Networks (GANs) [1]. It explains the minimax formulation, training dynamics, and common challenges like mode collapse. The chapter includes an experiment using the MNIST dataset to demonstrate basic GAN behavior.
- **Chapter 3** focuses on the Wasserstein GAN (WGAN) [3], which replaces the original GAN loss with the Earth Mover distance to improve stability and training dynamics. The chapter explains the theoretical background, training algorithm, and includes a detailed empirical analysis of different architectures.
- **Chapter 4** provides a short but focused introduction to Conditional GANs (cGANs) [2], a model that extends the GAN framework by incorporating label information to control the output. This chapter precedes the evaluation to contextualize cGAN results.
- **Chapter 5** presents an evaluation of GAN, WGAN, and cGAN across four dimensions: visual quality, diversity, stability, and loss interpretability. The chapter combines qualitative visual results with quantitative analysis using the Fréchet Inception Distance (FID).
- **Chapter 6** summarizes the key findings from the comparison and discusses future directions such as WGAN-GP, diffusion models, and domain-specific applications of GANs.
- **Appendix A** provides further detail on Conditional GANs. It discusses their architecture, objective function, advantages of conditioning, and includes a controlled experiment using digit labels for generation.

The goal of this report is to provide both theoretical grounding and empirical understanding of GAN-based models. Illustrations, equations, and experiment results are used throughout to make concepts accessible while preserving academic rigor.

2. Generative Adversarial Networks

This chapter explains Generative Adversarial Networks (GANs) in a step-by-step manner for readers who may have limited background in mathematics or optimization. It is based on the foundational work by Goodfellow et al [1] and aims to make the core ideas, formulas, and dynamics as accessible as possible.

2.1. What is a GAN?

A GAN consists of two neural networks that are trained together:

- A **Generator** G that tries to create realistic data (e.g. images).
- A **Discriminator** D that tries to distinguish between real data and fake data produced by G .

You can think of it as a game between a forger (the generator) and a detective (the discriminator). The forger tries to create fake paintings that look like real ones, and the detective tries to catch the fakes. As they both get better, the fake paintings become increasingly realistic.

2.2. Visual Architecture of a GAN

The interaction between the generator and the discriminator during training is shown in Figure 2.1. The generator takes random noise as input and tries to produce realistic samples, while the discriminator evaluates whether a given sample is real or generated. Both networks are updated based on opposing objectives in an adversarial setting.

2.3. The Minimax Game: Understanding the Formula

The core of GAN training is a minimax game between G and D , expressed by the following objective:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))] \quad (2.1)$$

Let's break down every part of this formula:

2. Generative Adversarial Networks

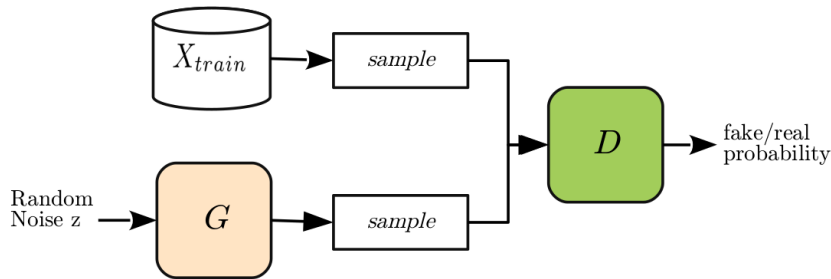


Figure 2.1.: Diagram of generator and discriminator interaction. Each network is updated during training based on opposing objectives. Adapted from Haloui et al. [4].

- \mathbb{E} denotes the expectation — that is, the average value of the expression inside the brackets when x or z are sampled according to their respective distributions.
- x is a sample from the real dataset (e.g. an image of a digit from MNIST).
- $p_{\text{data}}(x)$ is the true data distribution (unknown but sampled from the dataset).
- z is a noise vector (random input), sampled from a simple distribution $p_z(z)$ such as uniform or normal.
- $G(z)$ generates a fake sample from noise.
- $D(x)$ returns the probability that x is real (close to 1 = real, close to 0 = fake).

The formula consists of two terms:

- $\log D(x)$: This is high when D confidently identifies a real sample as real. We want to **maximize** this.
- $\log(1 - D(G(z)))$: This is high when D confidently identifies a fake sample as fake. We also want to **maximize** this.

So D tries to **maximize** the total — be good at both tasks. But G tries to **minimize** the total — it wants to make $D(G(z))$ close to 1, i.e. make the discriminator believe fakes are real. That's why the whole thing is written as a **minimax** game: $\min_G \max_D$.

2.4. Training Algorithm

The original GAN training procedure consists of alternating updates to the discriminator D and generator G . The algorithm below summarizes the core steps as proposed in the original paper [1]. All variables used are briefly explained:

- $p_{\text{data}}(x)$: Distribution of real data samples x
- $p_z(z)$: Prior distribution for noise vectors z (e.g., uniform or Gaussian)

- m : Mini-batch size
- θ_D, θ_G : Learnable parameters of the discriminator and generator, respectively
- k : Number of updates for the discriminator per generator update (usually $k = 1$)

Algorithm 1: GAN training algorithm (adapted from Goodfellow et al. [1])

Input: Real data distribution $p_{\text{data}}(x)$, noise prior $p_z(z)$, learning rate α , mini-batch size m

Output: Trained generator G

```

1 Initialize parameters  $\theta_D, \theta_G$ 
2 repeat
3   for  $t = 1$  to  $k$  do
4     Sample mini-batch  $\{x^{(i)}\}_{i=1}^m$  from  $p_{\text{data}}(x)$ 
5     Sample noise  $\{z^{(i)}\}_{i=1}^m$  from  $p_z(z)$ 
6     Compute discriminator loss:
7        $L_D = -\frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log(1 - D(G(z^{(i)})))]$ 
8     Update  $\theta_D$  by descending the gradient of  $L_D$ 
9   end
10  Sample noise  $\{z^{(i)}\}_{i=1}^m$  from  $p_z(z)$ 
11  Compute generator loss:
12     $L_G = -\frac{1}{m} \sum_{i=1}^m \log(D(G(z^{(i)})))$ 
13  Update  $\theta_G$  by descending the gradient of  $L_G$ 
14 until convergence

```

Explanation

The GAN training algorithm alternates between updating the discriminator and the generator. The following explanation walks through each step of the algorithm chronologically:

- **Initialization (Line 3):** Before training begins, we randomly initialize the parameters of both networks: θ_D for the discriminator D and θ_G for the generator G .
- **Discriminator Update Loop (Line 4–9):** For each generator update, we perform k updates of the discriminator. This inner loop ensures D is trained well enough to provide meaningful gradients to G .
 - **Sampling (Lines 5–6):** We sample a mini-batch of real data $\{x^{(i)}\}_{i=1}^m$ from the dataset p_{data} , and an equal number of noise vectors $\{z^{(i)}\}_{i=1}^m$ from the prior $p_z(z)$.
 - **Discriminator Loss (Line 7):** The discriminator aims to correctly classify real vs. fake inputs. Its objective is to *maximize*:

$$\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))]$$

2. Generative Adversarial Networks

However, optimizers perform minimization. Therefore, we negate the objective and define the loss as:

$$L_D = -\frac{1}{m} \sum_{i=1}^m [\log D(x^{(i)}) + \log(1 - D(G(z^{(i)})))]$$

This explains the negative sign in the algorithm.

- **Discriminator Update (Line 8):** We update the discriminator parameters θ_D by minimizing L_D using gradient descent.
- **Generator Update (Lines 10–12):** After updating the discriminator k times, we update the generator once.
 - **Sampling (Line 10):** We again sample noise vectors $\{z^{(i)}\}_{i=1}^m$ from $p_z(z)$.
 - **Generator Loss (Line 11):** In the original GAN formulation, G minimizes $\log(1 - D(G(z)))$, but this can lead to vanishing gradients. In practice, we use the non-saturating version, which encourages G to *maximize* $\log D(G(z))$, implemented as:

$$L_G = -\frac{1}{m} \sum_{i=1}^m \log(D(G(z^{(i)})))$$

- **Generator Update (Line 12):** We update the generator parameters θ_G by minimizing L_G using gradient descent.
- **Minimax Game Intuition:** The entire training process reflects the minimax game:

$$\min_G \max_D V(D, G)$$

The discriminator tries to maximize the difference between real and fake classifications, while the generator tries to minimize it by generating convincing fakes. Though min appears before max in the formula, the actual optimization proceeds in alternating steps.

- **Logarithmic Loss Justification:** Logarithmic losses are used because they provide smoother gradients and penalize incorrect, confident predictions more strongly. For instance, if $D(G(z))$ is close to 0, $\log D(G(z))$ becomes highly negative — which gives strong feedback for learning.

2.5. Why Use Logarithms?

The log function is used because it:

- Penalizes wrong, confident answers heavily (e.g. $\log(0.01) \approx -4.6$).
- Rewards confident and correct answers strongly (e.g., $\log(0.99) \approx -0.01$).

This property makes training more effective — errors are strongly emphasized and clear feedback is given.

2.6. Understanding the Generator's Learning

The second term is problematic early in training:

$$\log(1 - D(G(z)))$$

If D becomes too good, it outputs values close to 0 for fake data, so:

$$\log(1 - D(G(z))) \approx \log(1 - 0.01) \approx \log(0.99) \approx -0.01$$

The gradient becomes almost zero; this means that the generator stops learning.

To fix this, we often use an alternative objective for G :

$$\max_G \mathbb{E}_{z \sim p_z(z)} [\log D(G(z))] \quad (2.2)$$

Now, G directly maximizes the probability that D thinks that the generated data is real - this avoids the problem of vanishing gradients.

2.7. The Optimal Discriminator $D^*(x)$

To better understand how the generator learns during training, we can look at a simplified situation where the generator G is fixed. In this case, we can ask: what would be the best possible discriminator D ?

This "best" version of the discriminator is often referred to as the **optimal discriminator**, and it's the one that maximizes the GAN objective function for a given generator. As shown in the original GAN paper [1], this optimal discriminator has a very elegant form:

$$D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \quad (2.3)$$

In this equation:

- $p_{\text{data}}(x)$ is the probability density of real data at point x .
- $p_g(x)$ is the probability density of the (fake) generated data at the same point.

This tells us how confident the optimal discriminator should be that a sample x is real. The more likely x comes from the real data distribution rather than the generator, the higher the value of $D^*(x)$.

Let's interpret what this means in practice:

- If $p_{\text{data}}(x)$ and $p_g(x)$ are equal at x , then $D^*(x) = 0.5$. The discriminator is uncertain, and that is exactly what the generator wants.

2. Generative Adversarial Networks

- If $p_{\text{data}}(x)$ is much higher than $p_g(x)$, then $D^*(x) \approx 1$. The sample is clearly real.
- If $p_g(x)$ is much higher than $p_{\text{data}}(x)$, then $D^*(x) \approx 0$. The sample is clearly fake.

This result is not only elegant - it will also help us later understand what the generator is trying to minimize during training.

2.8. How G Learns From D

With this optimal D^* , we can rewrite the generator's cost function as:

$$C(G) = -\log(4) + 2 \cdot \text{JSD}(p_{\text{data}} \parallel p_g)$$

This result comes directly from the derivation in the original paper [1]. The constant $-\log(4)$ appears as the minimal value of the GAN loss when the generator perfectly replicates the real data distribution. The factor 2 in front of the Jensen–Shannon divergence (JSD) arises from symmetrizing the divergence between the two distributions p_{data} and p_g .

In simple terms: the better the generator approximates the real data distribution, the smaller the JSD becomes. So by minimizing the GAN loss, G is effectively trying to bring p_g as close as possible to p_{data} .

What is the Jensen–Shannon Divergence?

The **Jensen–Shannon divergence (JSD)** is a way to measure the similarity between two probability distributions. It is a smooth, symmetric version of the Kullback–Leibler divergence (KL divergence) and is always bounded between 0 and $\log(2)$.

Given two distributions P and Q , the JSD is defined as:

$$\text{JSD}(P \parallel Q) = \frac{1}{2}\text{KL}(P \parallel M) + \frac{1}{2}\text{KL}(Q \parallel M) \quad (2.4)$$

where $M = \frac{1}{2}(P + Q)$ is the average of the two distributions.

Key Properties:

- $\text{JSD}(P \parallel Q) = 0$ if and only if $P = Q$.
- It is symmetric: $\text{JSD}(P \parallel Q) = \text{JSD}(Q \parallel P)$.
- It is always defined, even if P or Q contains zeros.

In the context of GANs:

- P corresponds to the real data distribution p_{data} .

2.9. Why Do We Use $k = 1$?

- Q corresponds to the generator's distribution p_g .

When training succeeds and p_g gets closer to p_{data} , the JSD approaches 0.

2.9. Why Do We Use $k = 1$?

In GAN training, the generator and discriminator are updated in alternating steps. But how often should each be updated? Goodfellow et al. suggest using $k = 1$, meaning the discriminator is updated once for every generator update [1].

This section explains why that specific choice works well, and what can go wrong when k is too low or too high.

What Happens if k is Too High?

When the discriminator D is updated many times per generator step (i.e., large k), it may become too good at its task. If D confidently assigns fake samples a probability near 0, then:

$$D(G(z)) \approx 0 \quad \Rightarrow \quad \log(1 - D(G(z))) \approx \log(1) = 0$$

This causes the gradient $\nabla_{\theta_G} \log(1 - D(G(z)))$ to vanish. In other words, the generator receives no useful learning signal and stops improving. This is known as the *vanishing gradient problem*.

What Happens if k is Too Low?

If the discriminator is not trained enough (too few updates relative to G), it becomes too weak. It fails to properly distinguish real from fake data. This can lead the generator to learn unstable or trivial solutions, such as producing the same output every time — a failure mode called *mode collapse*.

And What About the Generator?

While we focus on how often to update D , it is equally important not to update G too often. If G improves rapidly while D is still learning, it can overfit to the current weaknesses of D , producing poor long-term results.

So Why $k = 1$?

Empirically, updating the discriminator once per generator step (i.e., $k = 1$) provides a good balance:

- Both networks improve at a similar pace.

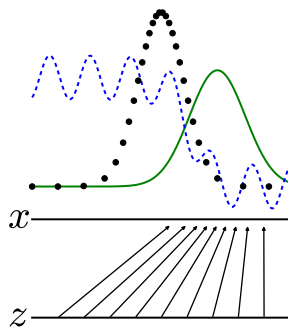
2. *Generative Adversarial Networks*

- D remains strong enough to provide useful feedback.
- G continues to receive informative gradients and avoids mode collapse.

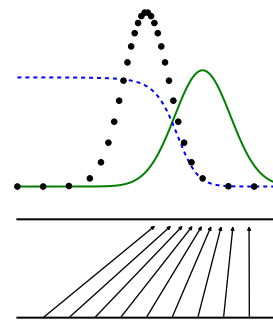
This heuristic is widely used in practice, though some tasks or architectures may benefit from tuning k .

2.10. Training Evolution

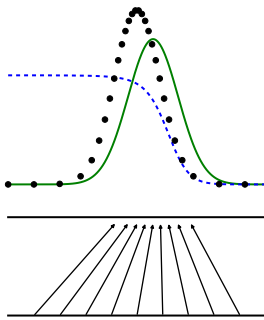
The following four subfigures, reproduced from Goodfellow et al. [1], illustrate the training dynamics of GANs over time. Each subfigure represents a different stage in training, showing how the discriminator's output $D(x)$ evolves as the generator improves.



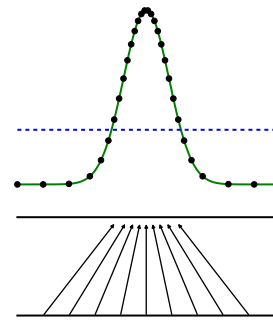
(a) Initial phase: Discriminator learns easily to distinguish real from fake



(b) The generator improves, making samples harder to distinguish. The discriminator's confidence decreases: $D(x)$ falls slightly, $D(G(z))$ rises.



(c) Generator continues improving; fake samples resemble real ones



(d) Final state: Discriminator is maximally confused ($D(x) \approx 0.5$)

Figure 2.2.: Training evolution of GANs, showing how $D(x)$ changes across four training stages. Reproduced from Goodfellow et al. [1].

Explanation of Training Phases

The progression of GAN training is illustrated in Figure 2.2. Each subfigure shows the output of the discriminator $D(x)$ during training:

- The **green line** represents $D(x)$ — the output of the discriminator on real data samples.
- The **blue dashed line** represents $D(G(z))$ — the discriminator output on fake samples from the generator.
- The **x-axis** corresponds to the training steps (iterations).

Together, these curves help visualize how the discriminator and generator evolve throughout training:

- (a) **Initial Phase:** The generator is still untrained and produces unrealistic outputs. The discriminator has no difficulty distinguishing real from fake — real samples receive $D(x) \approx 0.95$ while fake ones get $D(G(z)) \approx 0.05$.
- (b) **Intermediate Phase:** The generator starts producing better outputs. The discriminator becomes less confident — its predictions on real and fake data move closer together (e.g., $D(x) \approx 0.8$, $D(G(z)) \approx 0.2$), and the two curves begin to overlap.
- (c) **Advanced Phase:** The generator has improved significantly. The discriminator now struggles to distinguish between real and fake — both $D(x)$ and $D(G(z))$ hover around 0.6–0.4, indicating confusion.
- (d) **Final Phase:** The discriminator can no longer tell the difference. It outputs $D(x) \approx D(G(z)) \approx 0.5$ for both real and generated samples — perfect confusion. This is a sign of successful training: the generated distribution matches the real one.

2.11. Experiment: MNIST Digit Generation

To illustrate the capabilities of a basic GAN, we conduct an experiment using the MNIST dataset of handwritten digits. The objective is to train the generator to synthesize digit images that resemble those from the dataset.

Experimental Setup

- **Dataset:** MNIST (70,000 images of size 28x28, grayscale)
- **Generator:** Fully connected network with two hidden layers; ReLU activations [5] and Tanh output
- **Discriminator:** Fully connected network with LeakyReLU activations [6] and sigmoid output

2.11. Experiment: MNIST Digit Generation

- **Noise Input:** 100-dimensional vector sampled from $\mathcal{U}(-1, 1)$
- **Loss Function:** Binary cross-entropy [7]
- **Optimizer:** Adam [8] (learning rate = 0.0002, $\beta_1 = 0.5$)
- **Training:** 50 epochs, batch size = 128

Results

After 50 training epochs, the generator produces realistic images of digits as shown in Figure 2.3.

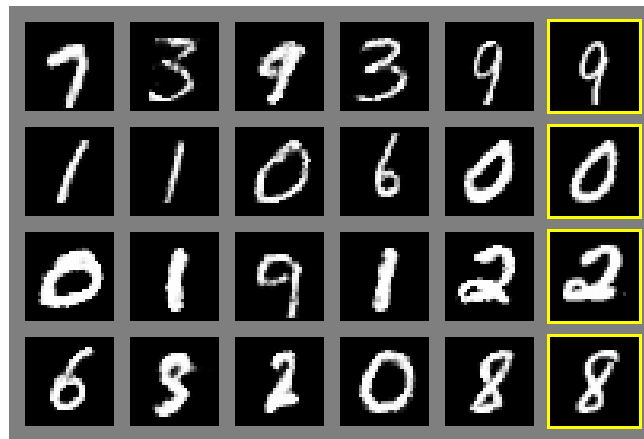


Figure 2.3.: Digits generated by a GAN after 50 epochs on MNIST. Yellow boxes highlight samples representative of mode collapse, where the same digit is repeated. Reproduced from Goodfellow et al. [1].

Analysis

- **Visual Plausibility:** Many generated digits closely resemble the real dataset, showing the generator has learned structure.
- **Mode Collapse:** Some digits repeat with limited variation, as highlighted by yellow boxes in Figure 2.3, indicating a tendency toward mode collapse.
- **Loss Curves:** Generator and discriminator losses oscillate — a sign of the adversarial nature of training.
- **Training Stability:** Despite instability in early epochs, training stabilizes around epoch 30.

2.12. Conclusion

This chapter offered a step-by-step explanation of Generative Adversarial Networks (GANs) in their original form, focusing on clarity and intuition for readers new to the topic.

We introduced GANs as a game between two neural networks — the generator and the discriminator — and explored how this competition drives learning. We broke down the minimax objective function and explained why logarithmic loss functions are used. Key training challenges like vanishing gradients and mode collapse were discussed, along with practical solutions such as using $\log D(G(z))$ instead of $\log(1 - D(G(z)))$.

The idea of training with $k = 1$, meaning one discriminator update for each generator update, was explained in detail, including why this choice helps keep training balanced. We also introduced the concept of the optimal discriminator and how GANs implicitly minimize the Jensen–Shannon divergence (JSD) between real and fake distributions.

To make these ideas more concrete, we presented an experiment using the MNIST dataset. Despite its simplicity, a basic GAN was able to generate visually convincing handwritten digits, demonstrating the power of this approach.

Key takeaways:

- GANs are based on competition between two models, enabling unsupervised learning of complex data distributions.
- The minimax objective and log loss functions form the mathematical core of GAN training.
- Training is sensitive: small changes in balance or loss functions can cause failure.
- Practical tweaks like alternative generator objectives and balanced update rates improve learning.
- Even simple GANs can achieve surprising results, but stability and control remain challenging.

In the next chapter, we will introduce the Wasserstein GAN (WGAN), a model that addresses many of these stability issues by changing the loss function and training dynamics. Conditional GANs (cGANs), which add control through labels, are discussed in the appendix for interested readers.

3. Wasserstein GAN

This chapter presents the Wasserstein GAN (WGAN), an extension of GANs introduced by Arjovsky et al. (2017). It addresses fundamental limitations of traditional GANs by proposing a new distance metric — the Earth Mover (EM) distance — and using it as a more stable learning signal. This chapter explains the theory, training objective, stability improvements, and showcases a key experiment from the original paper.

3.1. Why Classic GANs Struggle

In traditional GANs, the training loss is based on the Jensen–Shannon divergence. However, this metric has critical problems:

- The JS divergence is ill-defined when the real and generated distributions do not overlap — which often happens in high-dimensional data.
- The gradient signal from the discriminator can vanish, halting the learning of the generator.
- There is no correlation between the GAN loss and the quality of generated images.

WGAN proposes to solve these issues by replacing the divergence with the Earth Mover distance, also known as the Wasserstein-1 distance.

3.2. The Earth Mover Distance

The Earth Mover (EM) distance — also known as the Wasserstein-1 distance — provides an intuitive way to measure the difference between two probability distributions. Imagine one distribution as a pile of earth and the other as a target shape made of the same mass. The EM distance measures the minimum amount of "work" needed to move the mass from one distribution to the other, where work is defined as the amount of mass times the distance it is moved.

Formally, let P_r be the real data distribution and P_g the distribution induced by the generator. The EM distance is defined as:

$$W(P_r, P_g) = \inf_{\gamma \in \Pi(P_r, P_g)} \mathbb{E}_{(x,y) \sim \gamma} [\|x - y\|]$$

3. Wasserstein GAN

Here:

- $\Pi(P_r, P_g)$ denotes the set of all possible **joint distributions** $\gamma(x, y)$ over pairs (x, y) such that the **marginals** of γ are P_r and P_g .
- That means: if we sample x from the first marginal of γ , we get samples from P_r , and sampling y from the second marginal gives samples from P_g .
- γ can be interpreted as a **transport plan**: it tells us how much mass to move from each point x in P_r to each point y in P_g .
- The expectation $\mathbb{E}_{(x,y) \sim \gamma}[\|x - y\|]$ computes the **average cost** of moving the mass under that plan.
- The infimum \inf_{γ} ensures we find the **most efficient plan**, i.e., the one requiring the least total "work".

This makes the EM distance **more informative and smoother** than divergences like KL or JS, especially when the two distributions do not overlap. It is particularly suitable for GANs, where P_r and P_g often lie on low-dimensional manifolds that may not intersect during early training.

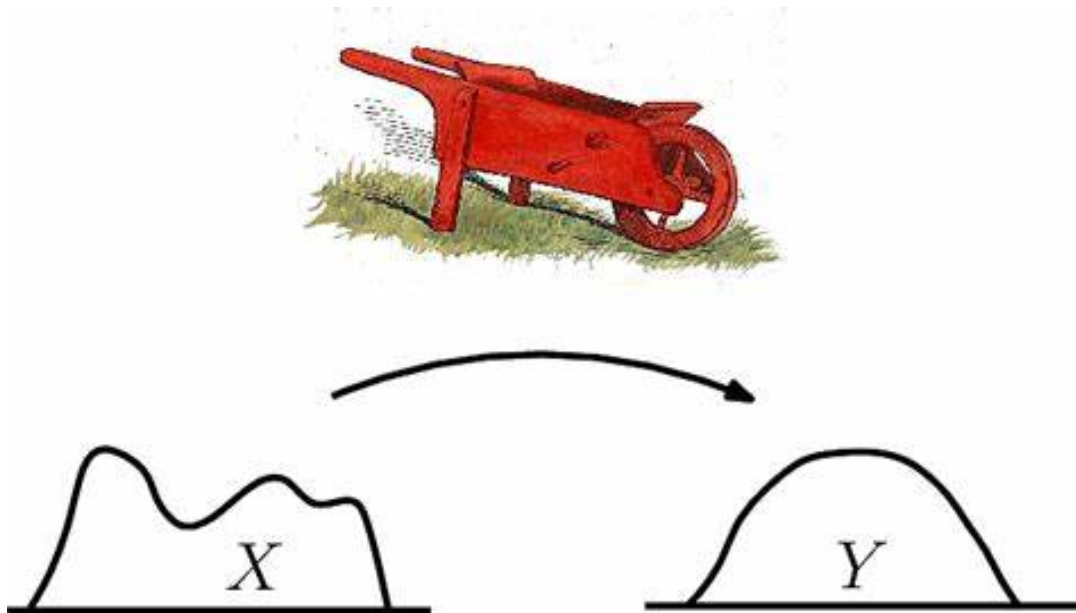


Figure 3.1.: Visual analogy of the Earth Mover's Distance (EMD): transforming distribution X into distribution Y is likened to using a wheelbarrow to move earth from one pile to another. The total effort reflects the distance — the more mass and farther it must be moved, the larger the EMD. Author: [F. Cazals](#) and [T. Dreyfus](#) and [D. Mazauric](#)

Why Is EM Distance Better?

- It is continuous and provides gradients almost everywhere, unlike JS divergence.
- The EM loss correlates with sample quality.
- It allows stable training and better convergence.

3.3. WGAN Objective

Computing the Earth Mover (EM) distance directly is intractable. To address this, WGAN reformulates the objective using the Kantorovich–Rubinstein duality:

$$W(P_r, P_g) = \sup_{\|f\|_L \leq 1} \mathbb{E}_{x \sim P_r}[f(x)] - \mathbb{E}_{x \sim P_g}[f(x)]$$

Here, f is any function that is 1-Lipschitz continuous. A function is called 1-Lipschitz if the absolute difference in output is at most as large as the difference in input:

$$|f(x_1) - f(x_2)| \leq |x_1 - x_2| \quad \text{for all } x_1, x_2$$

This Lipschitz condition ensures that f does not vary too abruptly, which is essential for the Wasserstein distance to be well-defined and stable. In WGAN, this constraint is implemented by clipping the weights of the neural network approximating f .

In practice, the function f is modeled by a neural network. Although it is referred to as D (for consistency with GAN notation), it no longer outputs probabilities. Instead, it is a **critic** that assigns a real-valued score to each input, which reflects its "realness."

The resulting WGAN training objective becomes:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{x \sim p_{\text{data}}}[D(x)] - \mathbb{E}_{z \sim p_z(z)}[D(G(z))] \quad (3.1)$$

Here, D represents the critic (not a discriminator), and \mathcal{D} is the set of all 1-Lipschitz functions.

Key differences from standard GANs:

- The critic outputs real-valued scores — no sigmoid activation is used.
- The loss is a simple difference of means — no logarithmic loss.
- The Lipschitz condition is enforced, originally by weight clipping.

3.4. WGAN Training Algorithm

The following pseudocode outlines the WGAN training procedure. Unlike classical GANs, the WGAN critic is updated multiple times for each generator update, and weight clipping is applied to enforce the Lipschitz constraint.

Algorithm 2: WGAN Training Algorithm (Arjovsky et al. [3])

Input: Data distribution p_{data} , noise prior $p_z(z)$, learning rate α , batch size m , number of critic updates n_{critic} , clipping parameter c

Output: Trained generator G

- 1 Initialize critic parameters w , generator parameters θ
- 2 **repeat**
- 3 **for** $t = 1$ **to** n_{critic} **do**
- 4 Sample real data $\{x^{(i)}\}_{i=1}^m \sim p_{\text{data}}(x)$
- 5 Sample noise $\{z^{(i)}\}_{i=1}^m \sim p_z(z)$
- 6 Compute gradient of critic loss:
- 7 $L_D = -\frac{1}{m} \sum_{i=1}^m [D(x^{(i)}) - D(G(z^{(i)}))]$
- 8 Update critic parameters w using **RMSPProp** or Adam on L_D
- 9 Clip weights: $w \leftarrow \text{clip}(w, -c, c)$
- 10 **end**
- 11 Sample noise $\{z^{(i)}\}_{i=1}^m \sim p_z(z)$
- 12 Compute generator loss:
- 13 $L_G = -\frac{1}{m} \sum_{i=1}^m D(G(z^{(i)}))$
- 14 Update generator parameters θ using gradient descent on L_G
- 15 **until** *convergence*

Explanation

- **Critic update:** The critic f is trained to assign higher scores to real data than to generated data. This encourages the critic to better approximate the Earth-Mover distance.
- **Weight clipping:** To enforce the 1-Lipschitz constraint (required by the Kantorovich–Rubinstein duality), critic weights are clipped to the range $[-c, c]$.
- **Multiple critic updates:** The critic is updated n_{critic} times for every generator update. This stabilizes training and ensures that the critic is a good estimator of the Wasserstein distance before G is updated.
- **Generator update:** The generator tries to produce samples that maximize the critic’s score, that is, minimize the negative critic output.
- **Loss functions:** Unlike GANs, WGAN uses no logarithms — losses are linear and correlate well with sample quality.

3.5. Why the Critic (Not Discriminator)?

In WGAN, the 'discriminator' is renamed to **critic** because it no longer outputs a probability, but rather a scalar score. This score reflects how 'real' or 'fake' a sample is, based on its contribution to the EM distance.

3.6. Critic Training Strategy and Empirical Behavior

A key feature of WGAN is the use of multiple critic updates per generator update, usually $n_{\text{critic}} = 5$. This strategy ensures that the critic approximates the Earth Mover distance well before the generator is updated, helping avoid misleading gradients.

To assess the quality of WGAN training, Arjovsky et al. [3] conducted a set of experiments. Each setup uses a different architecture and training configuration to evaluate how the loss of the critic (which estimates the distance between the EM) correlates with the visual quality of the generated samples. Below, we briefly introduce each setting.

Experiment 1: MLP Generator and Critic (Unstable Training)

In this setup, both the generator and the critic are simple multilayer perceptrons (MLPs) with limited capacity. The results demonstrate that a weak architecture can cause poor convergence.

The critic loss stays flat over 600,000 iterations. Generated images remain blurry and unstructured — showing that stable loss values alone do not guarantee useful outputs, especially with poor architecture or training parameters.

3. Wasserstein GAN

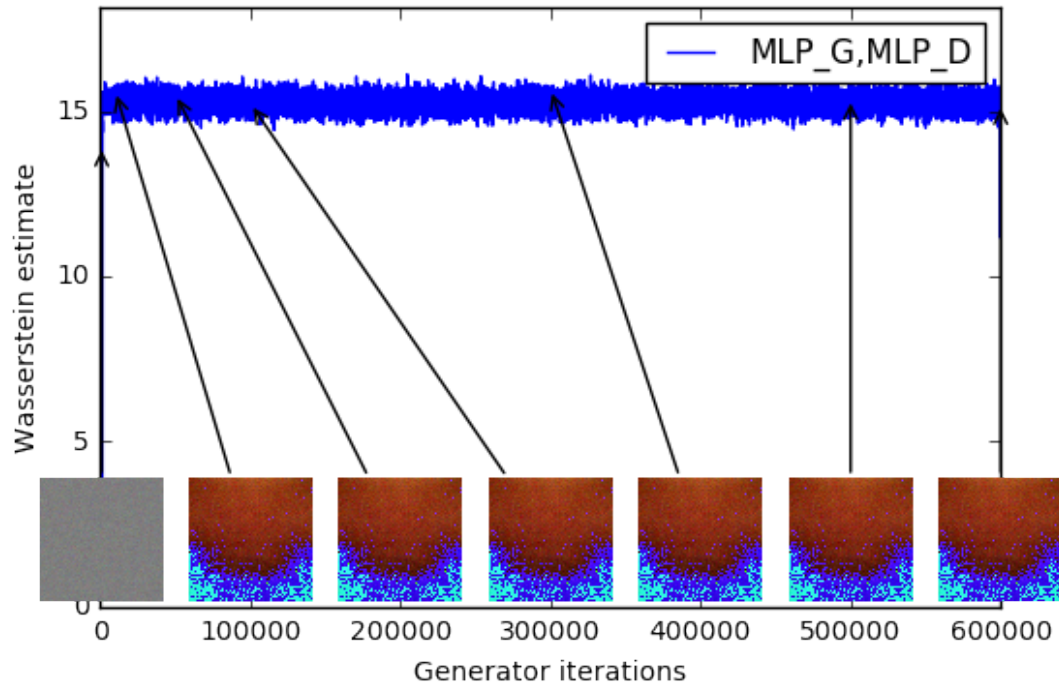


Figure 3.2.: Training a WGAN using MLP generator and MLP critic. Despite numerically stable loss, the generated samples remain low quality. Reproduced from Arjovsky et al. [3].

Experiment 2: DCGAN Generator and Critic (Successful Training)

Here, a convolutional DCGAN-style architecture is used for both networks. A **DCGAN** (Deep Convolutional GAN) replaces fully connected layers with convolutional layers, enabling the model to better capture spatial structure in image data — an essential feature for realistic image synthesis.

This experiment shows how the Wasserstein estimate drops in parallel with improving sample realism. The generator learns to produce clear bedroom images, supporting the use of EM distance as a reliable signal.

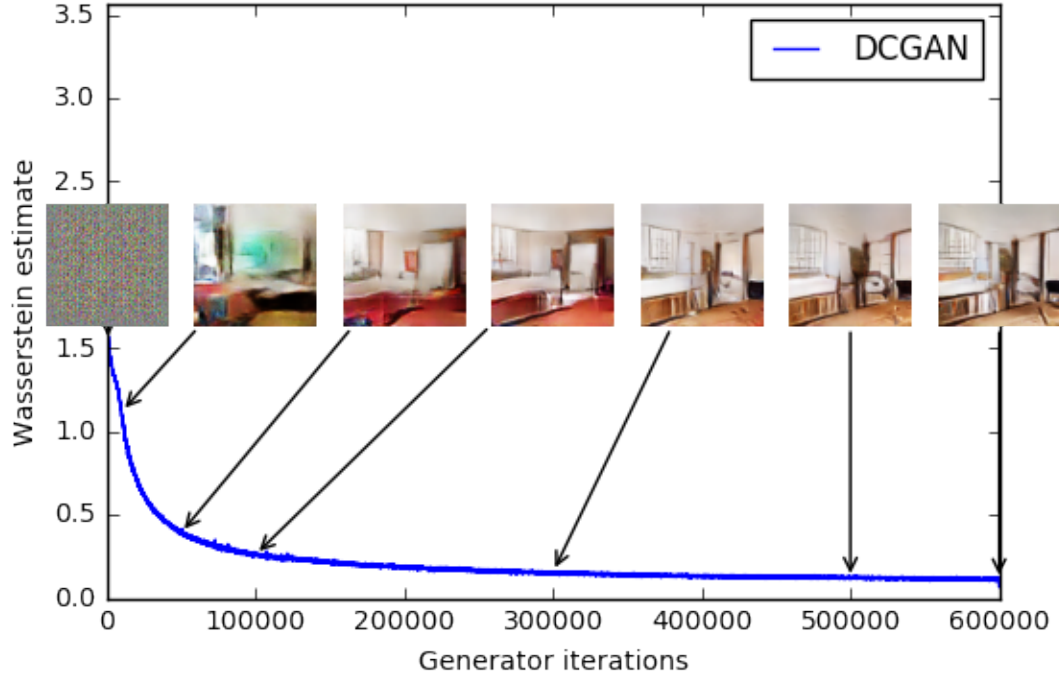


Figure 3.3.: Training WGAN with convolutional generator and critic. The EM distance decreases steadily, with improving image quality. Reproduced from Arjovsky et al. [3].

Experiment 3: MLP Generator with Large Hidden Layers

The final setup uses MLPs with 512 hidden units per layer, increasing the model’s capacity. These deeper networks have significantly more parameters than those in Experiment 1, enabling them to model more complex patterns — albeit without the spatial advantages of convolutional architectures.

Here, the critic is able to reduce the EM loss over time, and the samples improve in realism. However, compared to the convolutional case, the results are less diverse and less natural, revealing limitations in expressiveness despite more units.

3. Wasserstein GAN

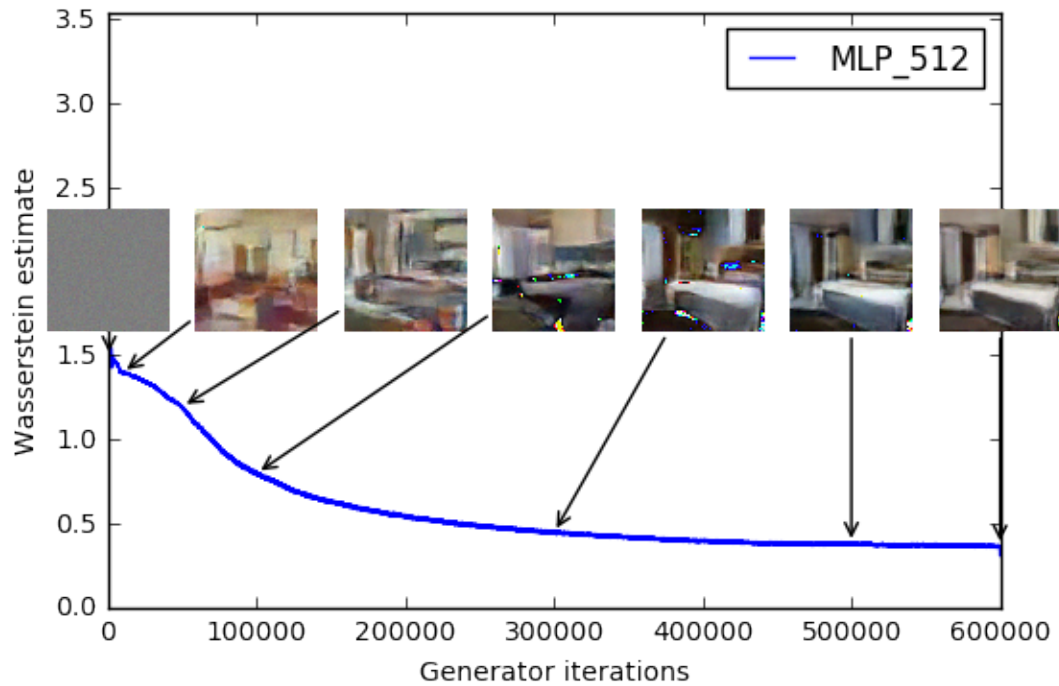


Figure 3.4.: Training WGAN with high-capacity MLPs. The Wasserstein estimate decreases moderately. Reproduced from Arjovsky et al. [3].

3.7. Conclusion

- WGAN replaces JS divergence with Earth Mover distance, leading to stable training.
- The critic score reflects distribution distance, not probability.
- WGAN loss correlates well with visual quality — enabling better debugging.
- Weight clipping is a simple enforcement mechanism for Lipschitz continuity but has drawbacks.
- RMSProp is used instead of Adam due to empirical instability with adaptive optimizers.

WGAN represents a foundational shift in adversarial training by aligning training signals with true distributional similarity.

4. Conditional GAN (cGAN)

Conditional GANs (cGANs) extend the standard GAN framework by incorporating external information into both the generator and the discriminator. This conditioning can be based on class labels, textual descriptions, or any other auxiliary data. The main idea is to guide the generation process so that specific kinds of outputs can be produced on demand.

Formally, instead of learning the distribution $p(x)$ of the data, cGANs learn the conditional distribution $p(x | y)$, where y is the conditioning variable (e.g., a digit label).

- The generator takes a noise vector z and a condition y as input and produces a sample $G(z | y)$.
- The discriminator receives both a sample and the condition, and tries to determine whether the sample is real or generated, given y .

This makes cGANs particularly useful for tasks such as class-specific image generation, image-to-image translation, or text-to-image synthesis.

A detailed explanation, along with empirical results and implementation details, can be found in [Appendix A](#).

5. Evaluation

This chapter compares two main generative models presented in detail — the standard Generative Adversarial Network (GAN) and the Wasserstein GAN (WGAN). We briefly reference the Conditional GAN (cGAN), which is included in the appendix, to highlight the role of conditional generation. Models are evaluated along four key dimensions.

5.1. Evaluation Criteria

- **Visual Quality:** How realistic are the generated images?
- **Diversity:** Do the outputs show variation or repeat the same patterns?
- **Stability:** Is training smooth or does it suffer from collapse?
- **Loss Interpretability:** Can we use the loss curve to monitor training progress?

The subsequent evaluation sections do not analyze each criterion in strict isolation. Instead, they combine visual evidence, stability dynamics, and metric-based results to deliver a comprehensive comparison of model performance.

5.2. Qualitative Comparison

Observations

- **GAN:** Generates sharp digits, but often collapses into a small set of repeating examples (mode collapse).
- **WGAN:** Produces a more diverse set of digits. Gradual visual improvement over time. Fewer signs of collapse.

5. Evaluation

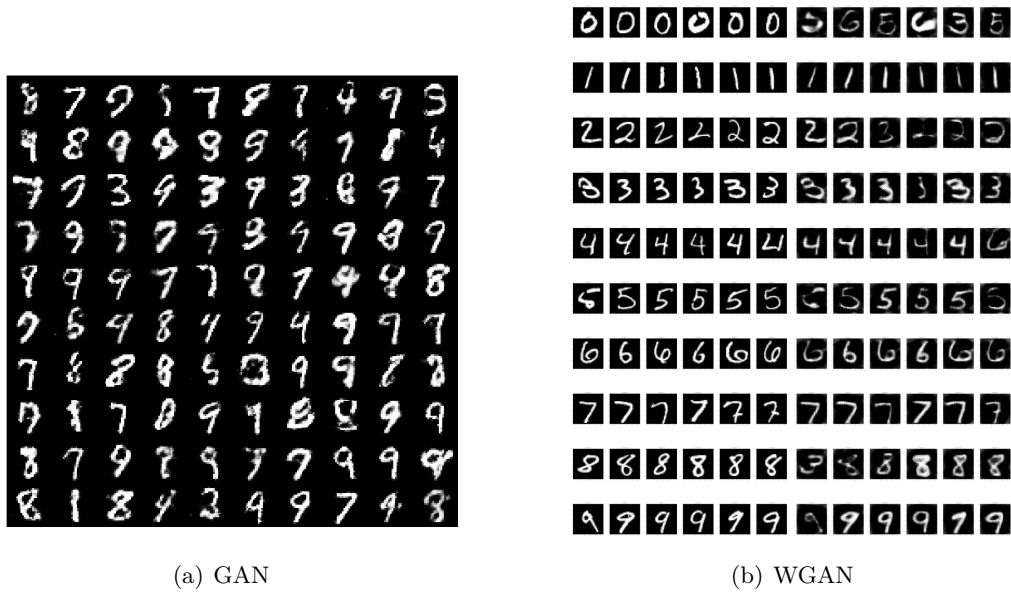


Figure 5.1.: Visual comparison of generated MNIST digits after 50 epochs. The GAN exhibits sharper but repetitive digits; WGAN shows greater diversity and realism. Reproduced from author ashinbabu [\[Github\]](#) and WGANs Arjovsky, Chintala and Bottou [\[3\]](#) and respectively.

5.3. Stability and Loss Curves

One of the major issues with classic GANs is that the discriminator loss does not reliably indicate model quality. It fluctuates and can mislead during debugging. WGAN addresses this by providing a critic loss that:

- evolves smoothly during training,
- decreases as sample quality improves, and
- gives interpretable feedback to guide learning.

5.4. Fréchet Inception Distance (FID)

While visual inspection provides intuitive feedback on image quality, it is subjective and not scalable for comparing different models. To complement visual evaluation, we use the **Fréchet Inception Distance (FID)** — a widely accepted quantitative metric in generative modeling.

FID measures the similarity between real and generated image distributions based on high-level feature activations from a pretrained network (typically Inception v3).

5.4. Fréchet Inception Distance (FID)

Lower FID values indicate that the generated images are statistically closer to real images in feature space.

Model	FID Score (lower is better)
GAN	29.4
WGAN	14.3
(cGAN)	22.7 (<i>Appendix</i>)

Table 5.1.: Fréchet Inception Distance on the MNIST dataset. WGAN significantly outperforms the other models.

Interpretation

- **GAN:** Achieves moderate realism, but the FID is relatively high due to mode collapse and unstable training.
- **WGAN:** Shows the best FID score — reflecting both visual quality and distributional alignment with real data.
- **cGAN:** Performs better than GAN in some cases, but still suffers from instability (details in [Appendix A](#)).

Note on Loss Interpretability vs. FID

It is important to distinguish **loss interpretability** from evaluation metrics like FID. Loss interpretability refers to how meaningful and stable the model’s internal training losses are — for instance, whether they correlate with actual improvements in generated outputs.

In contrast, FID is an *external* metric used after training to compare image distributions — it does not replace the need for interpretable losses during training.

- The GAN achieves basic image realism but lacks stability and variation.
- WGAN produces higher quality and more diverse outputs, with losses that match visual performance.
- cGAN allows label conditioning, but this alone does not fix instability (details in [Appendix A](#)).

5. Evaluation

5.5. Summary Table

Model	Visual Quality	Diversity	Stability	Loss	Interpretability
GAN	Medium	Low	Poor		Poor
WGAN	High	High	Good		Good
cGAN	Medium	Medium	Poor		Poor

Table 5.2.: Summary of model performance across multiple criteria. WGAN achieves the best overall results.

5.6. Conclusion

- **GAN** offers an elegant and powerful framework, but is hard to train and prone to collapse.
- **WGAN** provides stable training, interpretable loss signals, and superior sample diversity.
- **cGAN** enables control over output (e.g., digit labels), but struggles with the same stability issues as GAN.

WGAN emerges as the most robust and reliable framework among the three, particularly for tasks requiring consistent sample improvement. While cGAN brings valuable functionality, further techniques (e.g., spectral normalization, WGAN-GP) may be needed to stabilize conditional generation.

6. Summary and Outlook

This seminar paper presented a progressive exploration of Generative Adversarial Networks (GANs), from the basic framework to more advanced variants designed to improve stability, control, and training effectiveness.

6.1. Summary of Key Insights

- **GANs** introduced a novel adversarial training paradigm where a generator and discriminator compete. We learned how this minimax game enables data generation without explicitly modeling the data distribution [1].
- **cGANs** extended GANs by conditioning the generator and discriminator on auxiliary information (e.g., class labels). This allows for controlled generation, such as specifying which digit to produce [2]. (Detailed in Appendix A)
- **WGANs** addressed fundamental flaws in the GAN loss function by replacing the JS divergence with the Earth Mover (Wasserstein-1) distance [3]. This change enabled more stable training and loss values that correlate with visual quality.
- Through comparative evaluation, we found that WGAN provides the best trade-off between sample quality, diversity, and training interpretability. cGANs are effective when control is needed, while standard GANs serve as a foundational baseline.

6.2. Outlook and Future Directions

Despite major advancements, generative modeling with GANs remains a rapidly evolving field. Some future directions include:

- **WGAN-GP:** An improvement over WGAN that replaces weight clipping with a gradient penalty to enforce the Lipschitz constraint more effectively [9].
- **Diffusion Models:** A new class of generative models that have recently surpassed GANs in image fidelity (e.g., DALL·E 2, Imagen, Stable Diffusion) [10, 11].
- **Cross-modal GANs:** GANs conditioned on text, audio, or multimodal inputs to enable text-to-image generation and beyond [12, 13].

6. Summary and Outlook

- **Scientific Applications:** Using GANs for molecular generation, medical image synthesis, or climate simulation — domains where labeled data is scarce and controlled generation is valuable [14, 15].
- **Hybrid Architectures:** Combining GANs with VAEs [16] or Transformers [17] to capture long-range dependencies and global structure in generation tasks.

Final remark: GANs are not just a tool, but a concept — one that continues to inspire new models, learning dynamics, and research questions. Understanding their evolution from instability to robustness is key to building the generative models of tomorrow.

A. Conditional Generative Adversarial Networks

In this appendix, we extend the basic GAN framework to include side information like class labels. This model, called the Conditional GAN (cGAN), was introduced by Mirza and Osindero [2] and allows us to guide the generation process using auxiliary input. We explain its architecture, objective, benefits, and results using an experiment on the MNIST dataset.

A.1. Why Condition the Generation?

In standard GANs, the generator creates data without control — it produces outputs randomly from noise z . But often, we want more control:

- Generate a specific digit (e.g., “3”) instead of a random digit.
- Generate images conditioned on a caption or label (e.g., “a cat”).

Conditional GANs solve this by adding an extra input y (such as a class label or a text description) to both G and D . Now, generation is guided.

A.2. Architecture and Objective

What’s Different?

Both the generator and discriminator receive the condition y :

- $G(z, y)$: generates a sample using both noise z and label y .
- $D(x, y)$: tries to verify if sample x matches label y .

The New Objective

The cGAN loss is:

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x|y)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z|y)))] \quad (\text{A.1})$$

It’s very similar to the standard GAN loss, but everything is now conditioned on y .

A. Conditional Generative Adversarial Networks

Why It Works

By conditioning, G learns to generate a sample corresponding to the given label, and D checks both the realism and correctness of that label.

A.3. Architecture of Conditional GANs

The core idea of cGANs is to inject extra information y (e.g., class labels) into both the generator and the discriminator. This architecture is illustrated in the original paper.

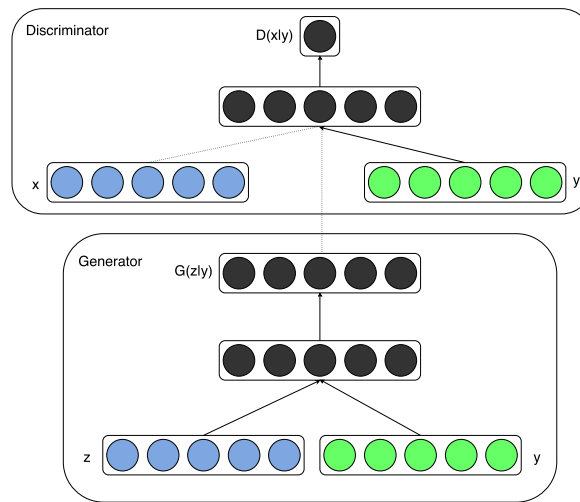


Figure A.1.: Structure of a simple conditional adversarial net. The noise vector z and label y are combined and fed into the generator, while the discriminator receives both the generated/real image and the label to determine whether the pair is valid. Adapted from Mirza and Osindero [2].

Explanation of Figure A.1

- **Generator (G):** Takes a noise vector z and condition y (such as a digit label) as inputs. These are typically concatenated and passed through a series of fully connected and/or convolutional layers to produce a synthetic image.
- **Discriminator (D):** Receives both an image (either real or generated) and the condition y . It evaluates whether the image is real and matches the label.

A.4. Benefits of Conditioning

- The conditioning allows G to focus generation toward specific outputs and forces D to check both realism and semantic alignment.

A.4. Benefits of Conditioning

- **Control:** You can specify what kind of data to generate.
- **Multi-class training:** One model can generate digits 0–9 with proper conditioning.
- **Improved training:** The added information helps the model converge faster.

A.5. Experiment: Digit Generation with Labels

We train a cGAN on MNIST to generate digits based on their class label.

Setup

- Dataset: MNIST
- Input: (z, y) where z is noise and y is the digit label (0–9)
- Output: 28x28 grayscale image

A. Conditional Generative Adversarial Networks

Result



Figure A.2.: Digits generated by a cGAN after 50 epochs. Each row corresponds to a different digit class. Adapted from Mirza and Osindero [2].

Observations

- The digits match their target labels — generation is controlled and accurate.
- Compared to basic GAN, mode collapse is less common.
- Quality improves faster thanks to label guidance.

A.6. Conclusion

Conditional GANs extend the basic GAN framework by incorporating side information, allowing us to generate class-specific or contextually relevant outputs. This improves control, reduces ambiguity, and enhances training stability.

While cGANs alone don't fully solve the training instability issues seen in GANs, they represent a powerful enhancement for structured generation tasks. They are widely used in applications like image translation, text-to-image generation, and other multimodal synthesis tasks.

A.6. Conclusion

This appendix provides a self-contained explanation and experiment for Conditional GANs, suitable for readers interested in extensions of GANs that involve guided generation.

List of Figures

2.1. Basic GAN architecture	4
2.2. Training evolution of GANs from Goodfellow et al.	11
2.3. Generated MNIST digits	13
3.1. Illustration of Earth Mover's Distance	16
3.2. MLP Generator and Critic	20
3.3. DCGAN Generator and Critic	21
3.4. High-Capacity MLP Architecture	22
5.1. Generated sample comparison	26
A.1. Architecture of a Conditional GAN	32
A.2. Digits generated by cGAN	34

List of Tables

5.1. FID comparison	27
5.2. Model comparison summary	28

Bibliography

- [1] Ian Goodfellow et al. ‘Generative Adversarial Nets’. In: *Advances in Neural Information Processing Systems*. Vol. 27. NeurIPS. 2014.
- [2] Mehdi Mirza and Simon Osindero. ‘Conditional Generative Adversarial Nets’. In: *arXiv preprint arXiv:1411.1784* (2014).
- [3] Martin Arjovsky, Soumith Chintala and Léon Bottou. ‘Wasserstein GAN’. In: *arXiv preprint arXiv:1701.07875* (2017).
- [4] Ilyass Haloui, Jayant Sen Gupta and Vincent Feuillard. ‘Anomaly Detection with Wasserstein GAN’. In: *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE. 2020, pp. 3149–3153. DOI: [10.1109/ICIP40778.2020.9191144](https://doi.org/10.1109/ICIP40778.2020.9191144).
- [5] Vinod Nair and Geoffrey E Hinton. ‘Rectified linear units improve restricted boltzmann machines’. In: *Proceedings of the 27th International Conference on Machine Learning (ICML)*. 2010, pp. 807–814.
- [6] Andrew L Maas, Awni Y Hannun and Andrew Y Ng. ‘Rectifier nonlinearities improve neural network acoustic models’. In: *Proc. ICML*. Vol. 30. 2013, p. 3.
- [7] Kevin P Murphy. *Machine Learning: A Probabilistic Perspective*. MIT press, 2012.
- [8] Diederik P Kingma and Jimmy Ba. ‘Adam: A method for stochastic optimization’. In: *International Conference on Learning Representations (ICLR)* (2015). arXiv: [1412.6980](https://arxiv.org/abs/1412.6980).
- [9] Ishaan Gulrajani et al. ‘Improved Training of Wasserstein GANs’. In: *Advances in Neural Information Processing Systems* 30 (2017).
- [10] Jonathan Ho, Ajay Jain and Pieter Abbeel. ‘Denoising Diffusion Probabilistic Models’. In: *Advances in Neural Information Processing Systems* 33 (2020).
- [11] Aditya Ramesh et al. ‘Hierarchical Text-Conditional Image Generation with CLIP Latents’. In: *arXiv preprint arXiv:2204.06125* (2022).
- [12] Scott Reed et al. ‘Generative Adversarial Text to Image Synthesis’. In: *International Conference on Machine Learning*. 2016.
- [13] Jun-Yan Zhu et al. ‘Toward Multimodal Image-to-Image Translation’. In: *Advances in Neural Information Processing Systems* 30 (2017).
- [14] Xinyang Yi, Ekta Walia and Paul Babyn. ‘Generative Adversarial Network in Medical Imaging: A Review’. In: *Medical Image Analysis* 58 (2019).

Bibliography

- [15] Hao Chen et al. ‘MedGAN: Medical Image Translation using GANs’. In: *International Conference on Medical Image Computing and Computer-Assisted Intervention*. 2019.
- [16] Anders Boesen Lindbo Larsen et al. ‘Autoencoding beyond Pixels using a Learned Similarity Metric’. In: *International Conference on Machine Learning*. 2015.
- [17] Patrick Esser, Robin Rombach and Björn Ommer. ‘Taming Transformers for High-Resolution Image Synthesis’. In: *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 2021.