

# Coverage-Guided Memory Replay for Mode-Diverse Generative Adversarial Networks: MADR-GAN

Abdallah Basem, Ahmed Hazem, Ahmed Islam, Ammar Hassan, Belal Fathy, Mazen Ahmed

---

**Abstract**—Generative Adversarial Networks (GANs) have established themselves as a foundational paradigm in deep generative modeling, leveraging an adversarial minimax formulation between a generator and discriminator to produce high-fidelity synthetic data. Despite remarkable achievements—including photorealistic face synthesis at FID 2.84 (StyleGAN2) and class-conditional ImageNet generation at FID 7.4 (BigGAN)—two persistent failure modes continue to undermine their practical utility: **mode collapse**, in which the generator converges to a restricted subset of the true data distribution, and **training instability**, arising from rotational gradient dynamics in the joint parameter space. Critically, we distinguish a previously underspecified variant of mode collapse—*temporal mode cycling*—wherein the generator oscillates between collapsed distributions over the course of training, evading standard static analyses that target gradient saturation alone. To address this, we propose **MADR-GAN** (Memory-Augmented Diversity Replay GAN), a Coverage-Weighted Memory Replay mechanism built on four components: (1) a reservoir-sampled episodic buffer of historical generated features maintained in the DINO self-supervised feature space; (2) a kernel-based mode coverage estimator that identifies underrepresented regions of the real data distribution; (3) a coverage-guided generator loss that biases synthesis toward historically undervisited real data modes; and (4) an adaptive coverage weight  $\lambda(t)$  driven by an online recall estimator, directly analogous to the ADA feedback loop. Theoretical analysis confirms that MADR-GAN preserves the original GAN Nash equilibrium while modifying the gradient landscape to discourage temporal cycling. We present a comprehensive experimental protocol targeting reduced mode collapse on Stacked MNIST (777 modes) and improved Recall and Coverage metrics on CIFAR-10 and CelebA-HQ without FID regression, alongside systematic ablations over buffer size, feature extractor choice, and adaptive scheduling.

**Index Terms**—Generative adversarial networks, mode collapse, temporal mode cycling, memory replay, diversity, coverage, self-supervised learning, DINO, adaptive training.

## 1 INTRODUCTION

GENERATIVE Adversarial Networks, introduced by Goodfellow et al. [1] at NeurIPS 2014, represent one of the most consequential advances in deep learning over the past decade. The adversarial framework—in which a generator  $G$  and a discriminator  $D$  are trained simultaneously in a minimax game—enables a fundamentally different approach to generative modeling: rather than maximizing a tractable likelihood, the generator implicitly learns to match the data distribution by fooling an adaptive, learned critic. This implicit learning signal has proven extraordinarily

powerful, enabling the generation of images, audio, video, and structured data at quality levels that were unthinkable a decade ago.

The practical impact of GANs has been sweeping. StyleGAN2 [11] achieves a Fréchet Inception Distance (FID) of 2.84 on the FFHQ dataset, producing human faces indistinguishable from photographs to casual inspection. BigGAN [12] achieves an FID of 7.4 on ImageNet at  $512 \times 512$  resolution. More recently, GigaGAN [26] and StyleGANT [27] have scaled the adversarial paradigm to text-conditioned synthesis at unprecedented resolution. These systems underpin applications in data augmentation [19], domain transfer, medical imaging, and creative AI.

Yet, the adversarial training objective harbors two well-documented failure modes that continue to limit reliability and applicability. The first is **mode collapse**: the generator learns to produce only a small subset of the modes present in the true data distribution, sacrificing diversity for fidelity. The second is **training instability**: the joint gradient dynamics of the generator–discriminator system exhibit rotational rather than convergent behavior [18], leading to oscillations, divergence, or the gradient vanishing pathology identified by Arjovsky and Bottou [21]. These problems are not peripheral edge cases; they are structural consequences of the adversarial formulation itself.

The literature has addressed these challenges through several paradigms. Wasserstein GANs [3], [4] replace the Jensen–Shannon Divergence (JSD) objective with the Wasserstein-1 distance, resolving gradient saturation mode collapse by ensuring meaningful gradients even when the two distributions have disjoint support. Minibatch discrimination [5] and PacGAN [6] introduce intra-batch diversity signals. Unrolled GANs [7] anticipate discriminator adaptation over multiple steps. Spectral normalization [9] and R1 regularization [17] improve gradient stability.

However, our analysis reveals a critical gap in the existing literature. These methods predominantly address two failure modes: (i) gradient saturation mode collapse, which occurs when the supports of  $p_{\text{data}}$  and  $p_g$  are disjoint; and (ii) intra-batch mode collapse, where all samples in a single training batch are identical. Neither class of methods addresses a third, distinct phenomenon we term **temporal mode cycling**: the tendency of the generator to cycle through a small set of modes over the course of

training, covering different modes at different time steps, but never maintaining simultaneous coverage of the full distribution. Standard evaluation metrics—including FID computed at a single checkpoint—cannot detect this behavior, because they measure the instantaneous distribution of the generator, not its trajectory. Precision/Recall [15] and Density/Coverage [16] metrics computed at convergence remain blind to oscillations that cancel out over time.

Our central observation is this: *existing methods address gradient starvation but not the temporal cycling behavior where the generator abandons previously covered modes and revisits them.* The generator, trained greedily against the current discriminator, finds it locally optimal to exploit whichever region of the data distribution the discriminator currently covers most weakly. As the discriminator adapts to reject that mode, the generator shifts to another—and eventually cycles back when the discriminator has forgotten the original rejection. This is a fundamentally temporal phenomenon, and it demands a temporal remedy.

We propose MADR-GAN, which incorporates a **Coverage-Weighted Memory Replay** mechanism. The key insight is to maintain a rolling episodic buffer of historical generated samples in the DINO [23] feature space, use this buffer to estimate which real data modes have been visited over training history, and apply a weighted generator loss that increases pressure toward historically undervisited modes. A feedback loop analogous to the Adaptive Discriminator Augmentation (ADA) mechanism of Karras et al. [19] dynamically adjusts the coverage loss weight based on an online recall estimate, ensuring the coverage signal is strong when diversity is insufficient and vanishes when coverage is adequate.

The contributions of this paper are fourfold:

- **Formal analysis of temporal mode cycling** as a failure mode distinct from gradient saturation mode collapse and intra-batch collapse, with a precise definition in terms of the time-evolution of the mode coverage function  $C(t)$ .
- **MADR-GAN**: a coverage-guided memory replay framework using reservoir-sampled episodic buffers in DINO feature space, with a kernel-based coverage estimator and weighted generator loss, compatible with any base GAN architecture.
- **Adaptive  $\lambda$  scheduling**: an online recall estimator that dynamically adjusts the coverage loss weight, with a sigmoid feedback law analogous to ADA, targeting a configurable recall threshold  $\tau_{\text{recall}}$ .
- **Theoretical analysis** demonstrating that the coverage loss preserves the original GAN Nash equilibrium, modifying only the gradient landscape en route to the equilibrium without introducing spurious fixed points.

The remainder of this paper is organized as follows. Section 2 provides background on GAN theory, mode collapse, training instability, and evaluation metrics. Section 3 reviews related work. Section 4 presents the MADR-GAN method in full detail. Section 5 provides theoretical analysis. Section 6 describes the experimental setup and proposed evaluation protocol. Section 7 discusses limitations, compar-

isons to diffusion models, and future directions. Section 8 concludes.

## 2 BACKGROUND

### 2.1 GAN Formulation

The original GAN [1] is defined by the minimax game:

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

where  $G : \mathcal{Z} \rightarrow \mathcal{X}$  maps a latent sample  $z \sim p_z$  (typically  $\mathcal{N}(0, I)$ ) to the data space, and  $D : \mathcal{X} \rightarrow [0, 1]$  estimates the probability that a given sample originates from  $p_{\text{data}}$  rather than  $p_g = G_{\#}p_z$ .

For a fixed generator  $G$ , the optimal discriminator is:

$$D^*(x) = \frac{p_{\text{data}}(x)}{p_{\text{data}}(x) + p_g(x)} \quad (2)$$

Substituting  $D^*$  back into Eq. (1), the optimal generator minimizes the Jensen–Shannon Divergence  $\text{JSD}(p_{\text{data}} \| p_g)$ , which achieves its global minimum of  $-\log 4$  if and only if  $p_g = p_{\text{data}}$ . In practice, the **non-saturating generator loss** [1],

$$\mathcal{L}_G^{\text{NS}} = \mathbb{E}_{z \sim p_z} [-\log D(G(z))], \quad (3)$$

is used in place of  $\mathbb{E}_z [\log(1 - D(G(z)))]$  to avoid gradient saturation early in training. The Deep Convolutional GAN (DCGAN) [2] extended this framework with convolutional architectures, establishing the backbone used in the majority of subsequent work.

### 2.2 Mode Collapse: Formal Analysis

Mode collapse is the tendency of a trained generator to assign high probability mass to only a strict subset of the modes present in  $p_{\text{data}}$ . We distinguish two mechanistically distinct subtypes that have different implications for remediation.

#### 2.2.1 Gradient Saturation Mode Collapse

When the supports of  $p_{\text{data}}$  and  $p_g$  are disjoint—a common occurrence when  $G$  is parametric and  $p_{\text{data}}$  is supported on a lower-dimensional manifold—the JSD is constant at  $\log 2$  regardless of  $\theta_G$ :

$$\text{supp}(p_{\text{data}}) \cap \text{supp}(p_g) = \emptyset \implies \text{JSD}(p_{\text{data}} \| p_g) = \log 2 \quad (4)$$

Consequently,  $\nabla_{\theta_G} \text{JSD} = 0$  almost everywhere, and the generator receives no usable gradient signal. This failure mode was rigorously characterized by Arjovsky and Bottou [21] and addressed by the Wasserstein GAN [3], which replaces JSD with the Wasserstein-1 distance  $W(p_{\text{data}}, p_g)$ , which remains a smooth function of the generator parameters even when the supports are disjoint.

#### 2.2.2 Temporal Mode Cycling (Our Focus)

Even when the Wasserstein distance provides non-zero gradients, the generator can exhibit a qualitatively different failure: it covers different subsets of  $p_{\text{data}}$ 's modes at different time steps, but never simultaneously covers the full distribution. Formally, let  $\{\mu_k\}_{k=1}^K$  denote the modes of  $p_{\text{data}}$  (local maxima of the density), and let  $\mathcal{N}(\mu_k)$  be an

$\varepsilon$ -neighborhood around mode  $k$ . Define the instantaneous mode coverage at training iteration  $t$  as:

$$C(t) = \frac{\left| \left\{ k : \int_{N(\mu_k)} p_g^{(t)}(x) dx > \varepsilon_{\min} \right\} \right|}{K} \quad (5)$$

where  $\varepsilon_{\min}$  is a minimum probability mass threshold. Temporal mode cycling occurs when  $C(t)$  oscillates rather than monotonically increasing toward 1, even as the total variation  $\|p_g^{(t)} - p_{\text{data}}\|_1$  decreases on average.

The mechanism is as follows. At time  $t_0$ ,  $G$  concentrates on mode  $\mu_1$ .  $D$  learns to detect this pattern and increases  $D(G(z))$  at samples near  $\mu_1$ . The generator gradient then favors shifting mass toward  $\mu_2$ , where  $D$  currently underestimates the probability of generated samples. As  $D$  subsequently adapts to  $\mu_2$ ,  $G$  returns to  $\mu_1$  (or a third mode), completing a cycle. This is a non-convergent limit cycle in the joint parameter space  $(\theta_G, \theta_D)$ .

This failure mode is *not* captured by static metrics computed at a single checkpoint: FID at the checkpoint may be acceptable if the generator happens to cover several modes at that instant, while being catastrophic at adjacent time steps. This motivates a temporal remedy.

### 2.3 Training Instability

The gradient vector field of the GAN minimax game can be written as  $[\nabla_{\theta_G} V, -\nabla_{\theta_D} V]$ . Balduzzi et al. [18] decomposed the Jacobian of this vector field into symmetric and antisymmetric components:

$$J_{\text{total}} = J_{\text{symm}} + J_{\text{antisymm}} \quad (6)$$

The symmetric component  $J_{\text{symm}}$  drives gradient descent behavior (convergence), while the antisymmetric component  $J_{\text{antisymm}}$  causes rotational dynamics: the joint parameter trajectory orbits around the Nash equilibrium rather than converging to it. This is the formal basis of GAN training oscillations observed empirically.

Mescheder et al. [17] established that local convergence can be guaranteed by adding R1 gradient regularization to the discriminator loss:

$$\mathcal{R}_1(\psi) = \frac{\gamma}{2} \mathbb{E}_{x \sim p_{\text{data}}} [\|\nabla_x D_\psi(x)\|^2] \quad (7)$$

where  $\gamma > 0$  is a hyperparameter. R1 regularization penalizes discriminator gradients at real data points, dampening the antisymmetric rotational dynamics and ensuring local convergence near the equilibrium. We adopt R1 regularization as a component of MADR-GAN’s discriminator objective.

### 2.4 Evaluation Metrics

#### 2.4.1 Fréchet Inception Distance (FID)

Heusel et al. [14] proposed FID as a measure of the distance between the real and generated distributions in the feature space of a pretrained Inception-v3 network. Let  $(\mu_r, \Sigma_r)$  and  $(\mu_g, \Sigma_g)$  denote the mean and covariance of Inception features for real and generated samples respectively. Then:

$$\text{FID} = \|\mu_r - \mu_g\|^2 + \text{Tr}(\Sigma_r + \Sigma_g - 2(\Sigma_r \Sigma_g)^{1/2}) \quad (8)$$

FID has several known limitations: it exhibits a finite-sample bias of order  $\mathcal{O}(d/n)$  where  $d = 2048$  is the Inception feature

dimension [22]; it is anchored to ImageNet-trained Inception features and may not capture domain-specific quality; and it conflates fidelity and diversity into a single scalar, making it impossible to identify which aspect of generation is deficient.

#### 2.4.2 Precision and Recall

Kynkänniemi et al. [15] proposed a geometrically-motivated precision/recall framework. A  $k$ -nearest-neighbor ball is constructed around each real sample and each generated sample in Inception feature space. **Precision** measures the fraction of generated samples that fall within the real data manifold (fidelity), while **Recall** measures the fraction of real data samples that are covered by the generated manifold (diversity). These two metrics decouple fidelity and diversity, making Recall a direct measure of mode coverage.

#### 2.4.3 Density and Coverage

Naeem et al. [16] observed that Recall is sensitive to outliers in the generated distribution, leading to pessimistic estimates. They proposed Density and Coverage: **Coverage** is defined as the fraction of real data samples for which at least one generated sample falls within the real sample’s  $k$ -NN ball. Coverage is more robust to outliers than Recall, making it the preferred metric for measuring mode collapse in our work.

## 3 RELATED WORK

### 3.1 Addressing Gradient Saturation: Wasserstein-Based Methods

The Wasserstein GAN [3] (WGAN) replaces the JSD objective with the Wasserstein-1 (Earth Mover’s) distance  $W(p_{\text{data}}, p_g)$ . By the Kantorovich–Rubinstein duality, this requires maximizing  $\mathbb{E}_{x \sim p_{\text{data}}} [f(x)] - \mathbb{E}_z [f(G(z))]$  over all 1-Lipschitz functions  $f$ . WGAN enforces the Lipschitz constraint via weight clipping, which introduces capacity artifacts and slow convergence.

WGAN-GP [4] replaces weight clipping with a gradient penalty:

$$\mathcal{P}_{\text{GP}} = \lambda_{\text{GP}} \mathbb{E}_{\hat{x} \sim p_{\hat{x}}} [(\|\nabla_{\hat{x}} D(\hat{x})\|_2 - 1)^2] \quad (9)$$

where  $\hat{x}$  is sampled by linearly interpolating between real and generated samples. WGAN-GP is significantly more stable than weight clipping but incurs a  $2\text{--}3\times$  computational overhead due to the double backward pass required to compute second-order gradients. Both WGAN variants effectively address gradient saturation mode collapse but do not target temporal mode cycling.

### 3.2 Intra-Batch Diversity Methods

Minibatch discrimination [5] augments the discriminator with statistics computed across samples within a minibatch, preventing the degenerate case where the generator produces identical samples in a batch. While effective at intra-batch collapse, this method has no temporal component and does not prevent the generator from producing different samples in different batches from a small, cycling set of modes.

PacGAN [6] modifies the discriminator to observe  $m$  samples simultaneously (“packing”), providing the discriminator with information about the joint distribution of a tuple of generated samples. Lin et al. [6] provide theoretical guarantees for mode coverage under packing. However, PacGAN’s diversity signal is derived entirely from the current generator state and does not leverage information about which modes were covered at earlier time steps.

### 3.3 Unrolled and Historical Generator Objectives

Unrolled GANs [7] optimize the generator against the discriminator unrolled  $k$  steps into the future, giving the generator awareness of how the discriminator will adapt to its current output. This addresses the one-step myopia that causes the generator to exploit transient discriminator weaknesses. However, unrolling  $k$  steps requires  $k$  forward-backward passes through  $D$  per generator update, making it computationally expensive ( $k = 5$  already doubles training time). Furthermore, unrolling addresses multi-step awareness on a short horizon but not the long-range temporal cycling over hundreds of thousands of iterations that we study. Unrolling does not maintain a buffer of historical states.

### 3.4 Catastrophic Forgetting Perspective

Thanh-Tung and Tran [20] framed mode collapse as a form of catastrophic forgetting in the discriminator: as the discriminator adapts to new generator outputs, it forgets its previous decision boundaries, allowing the generator to cycle back to previously rejected modes. They proposed discriminator gradient penalty variants to mitigate this forgetting. Our work shares this framing—mode collapse as a temporal memory failure—but differs in two key respects: (1) we focus on the generator objective rather than discriminator regularization, and (2) we anchor our coverage signal to the real data distribution via kernel density estimation in DINO feature space, rather than penalizing gradient norms.

Self-attention GANs [13] and spectral normalization [9] improve gradient flow and discriminator stability but do not address temporal diversity directly. Are GANs Created Equal [22] provides a large-scale empirical study demonstrating that many proposed improvements do not reliably outperform vanilla GAN with careful hyperparameter tuning, underscoring the need for principled, theoretically-grounded interventions.

### 3.5 Adaptive Training Methods

Karras et al. [19] proposed Adaptive Discriminator Augmentation (ADA) for limited-data GAN training. ADA adaptively adjusts the probability of applying stochastic augmentations to the discriminator’s inputs, using a feedback loop that monitors the discriminator’s overfitting heuristic (the fraction of training images for which  $D(x) > 0.5$ ) and adjusts augmentation strength to maintain a target value. This adaptive feedback design—a sigmoid controller targeting a scalar heuristic—directly inspires our adaptive  $\lambda(t)$  schedule, which targets recall rather than overfitting. DiffAugment [24] provides a complementary non-adaptive augmentation approach.

VEEGAN [8] adds a reconstructor network that maps generated samples back to the latent space, encouraging the generator to produce diverse outputs that span the latent space. Empirically, VEEGAN reduces mode collapse on Stacked MNIST but introduces additional architectural complexity and a separate training objective.

## 4 PROPOSED METHOD: MADR-GAN

### 4.1 Motivation: Why Memory Replay for Mode Diversity?

The core insight motivating MADR-GAN is that temporal mode cycling is fundamentally a memory failure: the generator does not “remember” which modes it has covered, and the discriminator does not “remember” which generated modes it has learned to reject. Standard GAN training operates on a single current state: the discriminator sees current real samples and current generated samples; the generator receives gradient from the current discriminator. There is no mechanism by which the training signal at iteration  $t$  encodes information about the distribution of generator outputs at iterations  $t - 100$ ,  $t - 1000$ , or  $t - 10000$ .

MADR-GAN remedies this by maintaining an episodic buffer of *historical* generated features and using this buffer to estimate a coverage signal that penalizes the generator for revisiting already-covered regions while incentivizing coverage of undervisited real data modes. Crucially, the coverage signal is anchored to the *real* data distribution—not to historical generated samples directly—ensuring that the generator is directed toward real modes rather than toward historically generated artifacts.

### 4.2 DINO Feature Space

We use the pretrained DINO [23] Vision Transformer (ViT-S/8) as our feature extractor:

$$\phi : \mathcal{X} \rightarrow \mathbb{R}^{384} \quad (10)$$

The DINO ViT-S/8 processes  $8 \times 8$  patches and produces a 384-dimensional CLS token embedding. This choice is motivated by three considerations. First, DINO features capture semantic similarity without requiring class labels: proximity in DINO space corresponds to perceptual and semantic proximity between images [23]. Second, DINO is entirely separate from the Inception-v3 network used to compute FID [14] and Precision/Recall [15], avoiding the circularity of optimizing a feature space while also using it for evaluation. Third, DINO’s self-supervised training objective produces features that are robust to photometric augmentations [23], reducing sensitivity to minor variations in generated image appearance. The DINO extractor is *frozen* throughout MADR-GAN training and is not fine-tuned.

### 4.3 Episodic Memory Buffer

We maintain a fixed-size episodic memory buffer:

$$\mathcal{M}_t = \{\phi(G_\tau(z_j))\}_{j=1}^K, \quad \tau \leq t \quad (11)$$

containing  $K = 2048$  DINO feature vectors of generated samples from training history. The buffer is initialized at the start of training and updated every  $T_{\text{update}} = 100$  iterations

using **reservoir sampling** [5]. For each new batch of  $B$  generated samples, each sample  $g_i$  is accepted into the buffer with probability  $K/n$ , where  $n$  is the total number of generated samples seen so far, replacing a uniformly random buffer entry. Reservoir sampling guarantees that the buffer is a uniform random sample over the *entire* training history at every time step, without any recency bias. This property is critical: a FIFO (first-in, first-out) buffer would over-represent recent samples and under-represent early-training behavior, weakening the historical coverage signal.

#### 4.4 Mode Visitation Frequency

For each real training sample  $x_r \in \mathcal{D}_{\text{real}}$ , we estimate its coverage score at time  $t$  using the buffer  $\mathcal{M}_t$ :

$$c_t(x_r) = \frac{1}{K} \sum_{m \in \mathcal{M}_t} k_\sigma(\phi(x_r), m) \quad (12)$$

where  $k_\sigma(a, b) = \exp(-\|a - b\|^2/(2\sigma^2))$  is a Gaussian kernel with bandwidth  $\sigma$ . The bandwidth is set by the **median heuristic**:  $\sigma = \text{median}(\{\|m_i - m_j\|_2 : m_i, m_j \in \mathcal{M}_t\})$ , updated every 1000 training iterations. The median heuristic is a well-established data-driven bandwidth selector for kernel methods that adapts to the intrinsic scale of the feature space without requiring cross-validation.

Intuitively,  $c_t(x_r)$  measures how much the historical buffer “explains” the region of DINO space near  $x_r$ . A high value of  $c_t(x_r)$  indicates that many historical generated samples have been close to  $x_r$  in feature space—the mode near  $x_r$  has been frequently generated. A low value indicates that the real data region around  $x_r$  has rarely been visited by the generator throughout training history.

#### 4.5 Coverage Importance Weights

We convert coverage scores into importance weights using a softmax over the negated scores, so that undervisited real modes receive higher weight:

$$w_t(x_r) = \frac{\exp(-\beta \cdot c_t(x_r))}{\sum_{r'=1}^N \exp(-\beta \cdot c_t(x_{r'}))} \quad (13)$$

where  $N = |\mathcal{D}_{\text{real}}|$  is the size of the real dataset and  $\beta = 1$  controls the concentration of the weights. When  $\beta \rightarrow 0$ , weights become uniform (no coverage guidance). When  $\beta \rightarrow \infty$ , all weight concentrates on the single least-covered real sample (extreme coverage push). We use  $\beta = 1$  as a default, ablated in Section 6.

#### 4.6 Coverage-Guided Generator Loss

The coverage-guided loss for the generator is a weighted kernel density estimation in DINO space:

$$\mathcal{L}_{\text{cov}}(G, t) = -\mathbb{E}_{z \sim p_z} \left[ \sum_{r=1}^N w_t(x_r) \cdot k_\sigma(\phi(G(z)), \phi(x_r)) \right] \quad (14)$$

This loss attracts the generator output  $G(z)$  toward real data samples in DINO feature space, with the attraction strength proportional to the importance weight  $w_t(x_r)$ . The coverage loss has the following key properties:

**Differentiability:**  $\nabla_{\theta_G} \mathcal{L}_{\text{cov}}$  is well-defined because  $\phi$  is a fixed, pretrained network (no gradient flows through  $\phi$ ),  $k_\sigma$

is smooth in both arguments, and  $G$  is differentiable in  $\theta_G$ . Gradients flow:  $z \rightarrow G(z) \xrightarrow{\nabla_x k_\sigma(\phi(x), \cdot)} \text{loss}$ .

**Bounded gradient:** The gradient magnitude is bounded:  $\|\nabla_{\theta_G} \mathcal{L}_{\text{cov}}\|_2 \leq N \cdot \max_r w_t(x_r) \cdot C_\phi$ , where  $C_\phi$  bounds the operator norm of the DINO Jacobian  $\partial\phi/\partial x$ . This prevents gradient explosion from the coverage term.

**Equilibrium alignment:** At  $p_g = p_{\text{data}}$ , the kernel density  $k_\sigma(\phi(G(z)), \phi(x_r))$  is high for all real  $x_r$ , making  $c_t(x_r)$  high for all  $r$ , weights approximately uniform, and  $\nabla_{\theta_G} \mathcal{L}_{\text{cov}} \rightarrow 0$  (see Section 5).

**No circular evaluation:**  $\phi$  is the DINO ViT-S/8; FID, Precision, and Recall use Inception-v3 features. Optimizing  $\mathcal{L}_{\text{cov}}$  cannot artificially inflate the Inception-based evaluation metrics.

#### 4.7 Adaptive Coverage Weight $\lambda(t)$

Rather than using a fixed coverage weight, we adapt  $\lambda(t)$  based on an online estimate of the generator’s current recall. The online recall estimate uses the buffer  $\mathcal{M}_t$ :

$$\widehat{\text{Recall}}_t = \frac{1}{N} \sum_{r=1}^N \mathbf{1}[\exists m \in \mathcal{M}_t : \|\phi(x_r) - m\|_2 < \varepsilon_t] \quad (15)$$

where  $\varepsilon_t$  is set to the 95th percentile of intra-buffer pairwise distances in  $\mathcal{M}_t$ , updated every 1000 iterations. This is a computationally cheap proxy for the true Recall metric: it measures what fraction of real data samples have at least one historical generated neighbor in DINO space.

The adaptive weight follows a sigmoid feedback law:

$$\lambda(t) = \lambda_{\text{max}} \cdot \sigma(\kappa \cdot (\tau_{\text{recall}} - \widehat{\text{Recall}}_t)) \quad (16)$$

where  $\sigma(\cdot)$  is the sigmoid function,  $\tau_{\text{recall}} = 0.80$  is the target recall,  $\kappa = 10$  controls the sharpness of the transition, and  $\lambda_{\text{max}} = 1.0$ . The behavior is as follows:

- When  $\widehat{\text{Recall}}_t \ll 0.80$ :  $\lambda(t) \approx \lambda_{\text{max}}$ , imposing a strong coverage push.
- When  $\widehat{\text{Recall}}_t \approx 0.80$ :  $\lambda(t) \approx \lambda_{\text{max}}/2$ , a moderate regularization.
- When  $\widehat{\text{Recall}}_t \gg 0.80$ :  $\lambda(t) \approx 0$ , effectively disabling the coverage term and reverting to standard adversarial training.

The weight  $\lambda(t)$  is updated every 4 training iterations. This design mirrors the ADA controller [19]: both use a feedback signal derived from a training heuristic to adaptively modulate a regularization strength, ensuring that the auxiliary objective only intervenes when necessary.

#### 4.8 Full Training Objective

The total MADR-GAN generator loss at iteration  $t$  is:

$$\mathcal{L}_G(t) = \mathcal{L}_{\text{adv}} + \lambda(t) \cdot \mathcal{L}_{\text{cov}}(G, t) \quad (17)$$

where  $\mathcal{L}_{\text{adv}} = \mathbb{E}_{z \sim p_z} [-\log D(G(z))]$  is the non-saturating adversarial loss (Eq. (3)). MADR-GAN is compatible with any base GAN adversarial loss (WGAN, WGAN-GP, non-saturating, hinge) by substituting the appropriate  $\mathcal{L}_{\text{adv}}$ .

---

**Algorithm 1** MADR-GAN Training

---

**Require:** Real dataset  $\mathcal{D}$ ; pretrained DINO  $\phi$ ; buffer size  $K$ ; update interval  $T_{\text{up}}$ ; hyperparameters  $\lambda_{\max}, \tau, \kappa, \beta, \gamma$

- 1: Initialize  $G_\theta, D_\psi$  with random weights
- 2: Precompute  $\phi(x_r)$  for all  $x_r \in \mathcal{D}$ ; store in  $\Phi_{\text{real}}$
- 3: Initialize buffer  $\mathcal{M} \leftarrow \emptyset$ ; set  $n \leftarrow 0$
- 4: **for** each training iteration  $t = 1, 2, \dots$  **do**
- 5:   Sample minibatch  $\{x_i\}_{i=1}^B \sim \mathcal{D}$  and latents  $\{z_i\}_{i=1}^B \sim p_z$
- 6:   Generate samples  $\{\hat{x}_i\}_{i=1}^B = \{G_\theta(z_i)\}$
- 7:   **// Update Discriminator**
- 8:   Compute  $\mathcal{L}_D$  using Eq. (18); update  $\psi$  via Adam
- 9:   **// Update Memory Buffer**
- 10:   **if**  $t \bmod T_{\text{up}} = 0$  **then**
- 11:     **for** each generated sample  $\hat{x}_i$  **do**
- 12:        $n \leftarrow n + 1$
- 13:       **if**  $|\mathcal{M}| < K$  **then**
- 14:          $\mathcal{M} \leftarrow \mathcal{M} \cup \{\phi(\hat{x}_i)\}$
- 15:       **else if** Bernoulli( $K/n$ ) = 1 **then**
- 16:         Replace random entry in  $\mathcal{M}$  with  $\phi(\hat{x}_i)$
- 17:       **end if**
- 18:     **end for**
- 19:   **end if**
- 20:   **// Compute Coverage Signal**
- 21:   **if**  $|\mathcal{M}| > 0$  **then**
- 22:     Update  $\sigma \leftarrow$  median pairwise dist.( $\mathcal{M}$ ) (every 1000 iters)
- 23:     Compute  $c_t(x_r)$  for all  $x_r \in \mathcal{D}$  using Eq. (12)
- 24:     Compute  $w_t(x_r)$  using Eq. (13)
- 25:     Compute  $\mathcal{L}_{\text{cov}}$  using Eq. (14)
- 26:   **else**
- 27:      $\mathcal{L}_{\text{cov}} \leftarrow 0$
- 28:   **end if**
- 29:   **// Adaptive  $\lambda$  Update**
- 30:   **if**  $t \bmod 4 = 0$  **then**
- 31:     Compute  $\text{Recall}_t$  using Eq. (15)
- 32:     Update  $\lambda(t)$  using Eq. (16)
- 33:   **end if**
- 34:   **// Update Generator**
- 35:   Sample new latents  $\{z_i\} \sim p_z$
- 36:   Compute  $\mathcal{L}_G(t) = \mathcal{L}_{\text{adv}} + \lambda(t) \cdot \mathcal{L}_{\text{cov}}$  using Eq. (17)
- 37:   Update  $\theta$  via Adam
- 38: **end for**

---

The discriminator loss retains the standard binary cross-entropy with R1 regularization:

$$\begin{aligned} \mathcal{L}_D = \mathbb{E}_{x \sim p_{\text{data}}} [-\log D(x)] + \mathbb{E}_z [-\log(1 - D(G(z)))] \\ + \frac{\gamma}{2} \mathbb{E}_{x \sim p_{\text{data}}} [\|\nabla_x D(x)\|^2] \quad (18) \end{aligned}$$

with  $\gamma = 10.0$ . No modification is made to the discriminator architecture or objective beyond R1. This is important: methods that modify the discriminator to include historical fake samples (a three-way classification objective) introduce a non-stationary discriminator signal that can destabilize training [20]. MADR-GAN’s coverage signal is a purely generator-side regularizer with bounded, deterministic gradients.

The full training procedure is given in Algorithm 1.

## 4.9 Computational Overhead

The primary computational costs added by MADR-GAN over a base GAN are: (1) a DINO forward pass for generated samples, incurring approximately 2 ms per batch on a modern GPU (ViT-S/8 is a small transformer); (2) buffer maintenance at  $\mathcal{O}(K \cdot d)$  memory with  $d = 384$ , amounting to  $2048 \times 384 \times 4$  bytes  $\approx 3.1$  MB; and (3) coverage weight computation at  $\mathcal{O}(N \cdot K)$  per update, parallelizable via matrix operations. In total, MADR-GAN adds approximately 15–20% wall-clock overhead compared to the base GAN on CIFAR-10 with DCGAN. Precomputing  $\phi(x_r)$  for all real samples at the start of training eliminates DINO overhead at evaluation time.

## 5 THEORETICAL ANALYSIS

### 5.1 Equilibrium Preservation

**Claim:** Adding  $\mathcal{L}_{\text{cov}}$  to the generator objective does not introduce new Nash equilibria or eliminate the existing equilibrium  $p_g = p_{\text{data}}$ .

**Argument:** We show that  $\nabla_{\theta_G} \mathcal{L}_{\text{cov}} \rightarrow \mathbf{0}$  at the GAN Nash equilibrium. At the unique equilibrium,  $p_g = p_{\text{data}}$ . Since the generator has covered all modes of  $p_{\text{data}}$ , the buffer  $\mathcal{M}_t$  (which records uniform samples from the training history of  $G$ ) is a sample from  $p_{\text{data}}$  in DINO feature space (up to finite-sample error). Therefore,  $c_t(x_r)$  is approximately equal for all real samples  $x_r$  (every real data point has been visited equally in expectation). The importance weights  $w_t(x_r)$  approach the uniform distribution:  $w_t(x_r) \approx 1/N$  for all  $r$ . The gradient of the coverage loss becomes:

$$\nabla_{\theta_G} \mathcal{L}_{\text{cov}} \Big|_{p_g=p_{\text{data}}} \approx -\frac{1}{N} \sum_r \nabla_{\theta_G} \mathbb{E}_z [k_\sigma(\phi(G(z)), \phi(x_r))] \quad (19)$$

Under  $p_g = p_{\text{data}}$ , the generator’s DINO feature distribution  $\phi_\# p_g$  equals  $\phi_\# p_{\text{data}}$ . For any fixed  $x_r$ ,  $\mathbb{E}_z [k_\sigma(\phi(G(z)), \phi(x_r))]$  is maximized when  $G(z)$  lies in the support of  $p_{\text{data}}$ , and the gradient with respect to  $\theta_G$  vanishes when  $p_g = p_{\text{data}}$  (since the generator is already at the maximizer of the kernel density toward all real points). Therefore, the equilibrium is preserved.  $\square$

This argument confirms that MADR-GAN modifies only the gradient landscape *en route* to the equilibrium, without altering the equilibrium itself. The coverage loss acts as a directional bias in the transient phase of training, diminishing as convergence is approached.

### 5.2 Informal Mode Coverage Guarantee

Consider a real data mode  $\mu_k$  that is completely uncovered at time  $t$ :  $c_t(\mu_k) = 0$ . Then  $w_t(x_r) \propto \exp(-\beta \cdot 0) = 1$  for  $x_r$  near  $\mu_k$ , while covered modes have smaller weights (since their  $c_t$  values are positive). Thus, the coverage loss creates a strong gradient attracting  $G(z)$  toward the uncovered mode. Conversely, if  $c_t(\mu_k) \rightarrow \infty$  for some mode (heavily overvisited),  $w_t(x_r) \rightarrow 0$  for  $x_r$  near  $\mu_k$ , removing the attraction toward that mode. This provides an informal guarantee: MADR-GAN penalizes the generator for ignoring any real data mode proportionally to the degree of neglect.

### 5.3 Gradient Stability of the Coverage Term

We note an important contrast between MADR-GAN and methods that extend the discriminator to classify historical generated samples (e.g., a three-way discriminator distinguishing real, current fake, and historical fake). Such methods introduce a non-stationary discriminator objective: the set of historical fakes changes at every iteration, creating a moving target for the discriminator gradient. This can amplify the antisymmetric Jacobian component  $J_{\text{antisymm}}$  identified by Balduzzi et al. [18], potentially increasing rotational instability.

In contrast,  $\mathcal{L}_{\text{cov}}$  provides a deterministic (given the buffer), bounded gradient signal. The buffer is updated asynchronously and infrequently ( $T_{\text{up}} = 100$  iterations), keeping the coverage signal approximately constant over many generator updates. The DINO feature extractor  $\phi$  is frozen, so the kernel function  $k_\sigma(\phi(G(z)), \phi(x_r))$  changes only through  $G(z)$ . This makes  $\nabla_{\theta_G} \mathcal{L}_{\text{cov}}$  a smooth, slowly-varying function of  $\theta_G$ , avoiding the gradient spike phenomena associated with non-stationary discriminator objectives.

## 6 EXPERIMENTAL SETUP

### 6.1 Datasets

We conduct experiments on four datasets spanning a range of mode complexities:

**Stacked MNIST** [8]: Three MNIST digits are stacked as RGB channels, creating a synthetic distribution with exactly 777 known, separable modes (all three-digit combinations). This is the standard benchmark for quantifying the number of captured modes, enabling precise measurement of mode collapse.

**CIFAR-10** [22]: Natural images at  $32 \times 32$  resolution across 10 classes. We report all six metrics: FID, IS, Precision, Recall, Density, and Coverage.

**CelebA-HQ** [10]: High-quality face images cropped to  $64 \times 64$ . Primarily used to evaluate FID, Precision, and Recall on a unimodal-within-class dataset.

**ImageNet** ( $64 \times 64$ ) [12]: Class-conditional generation across 1000 classes. Used to evaluate FID and IS at scale, with per-class Coverage to measure per-class mode coverage.

### 6.2 Baselines

Table 1 lists all baseline methods. We include a representative set spanning the key paradigms discussed in Section 3, plus the state-of-the-art limited-data upper bound.

### 6.3 Evaluation Metrics

**Primary metrics (mode coverage):** Recall [15], Coverage [16], Number of Modes (Stacked MNIST). These directly quantify the diversity failure targeted by MADR-GAN. A method that improves FID while reducing Recall has merely traded diversity for fidelity and should not be considered an improvement.

**Secondary metrics (overall quality):** FID [14], Precision [15], IS [5]. We require that MADR-GAN does not

TABLE 1  
Baseline Methods for Comparison on CIFAR-10 and Stacked MNIST

Method	Mechanism	Overhead
Vanilla GAN [1]	Non-saturating loss	—
WGAN-GP [4]	Wasserstein + gradient penalty	2–3×
WGAN-GP + MBDisc [5]	WGAN-GP + minibatch statistics	2.5–3.5×
PacGAN [6]	Packing ( $m = 2$ )	1.1×
VEEGAN [8]	Reconstructor network	1.5×
StyleGAN2-ADA [19]	Adaptive augmentation (upper bound)	3×
<b>MADR-GAN (Ours)</b>	Coverage-guided memory replay	1.15–1.20×

regress FID by more than 5% relative to the base GAN, confirming that the coverage improvement does not sacrifice fidelity.

We compute all metrics using 50,000 generated samples per model, using the standard evaluation code from [15], [16]. FID uses 50,000 real CIFAR-10 training samples as the reference set. Stacked MNIST mode count uses a pretrained MNIST digit classifier applied to each channel independently.

### 6.4 Implementation Details

**Base architecture:** DCGAN [2] for Stacked MNIST and CIFAR-10; BigGAN backbone [12] for ImageNet ( $64 \times 64$ ); DCGAN with spectral normalization [9] for CelebA-HQ.

**DINO feature extractor:** ViT-S/8 pre-trained on ImageNet (public checkpoint from facebookresearch/dino). Frozen throughout training. Input images resized to  $224 \times 224$  for DINO forward pass.

**Buffer:**  $K = 2048$  (default; ablated). Reservoir sampling with  $T_{\text{up}} = 100$  iterations.

**Adaptive  $\lambda$ :**  $\lambda_{\text{max}} = 1.0$ ;  $\tau_{\text{recall}} = 0.80$ ;  $\kappa = 10$ ; updated every 4 iterations.

**R1 regularization:**  $\gamma = 10.0$  (applied to discriminator).

**Optimizer:** Adam with  $\beta_1 = 0$ ,  $\beta_2 = 0.99$  for both  $G$  and  $D$ . Two time-scale update rule (TTUR) [14]:  $\text{lr}_D = 2 \times 10^{-4}$ ,  $\text{lr}_G = 1 \times 10^{-4}$ .

**Training duration:** 100k iterations on CIFAR-10 and CelebA-HQ; 200k on ImageNet.

**Hardware:** All experiments conducted on a single NVIDIA A100-80GB GPU. Expected training time:  $\approx 18$  hours (CIFAR-10, DCGAN backbone),  $\approx 72$  hours (ImageNet, BigGAN backbone).

### 6.5 Ablation Studies

Table 2 describes the ablation experiments designed to isolate the contribution of each MADR-GAN component. All ablations are conducted on CIFAR-10 with the DCGAN backbone.

**Expected outcomes:** A2/A3 should confirm that the adaptive schedule outperforms static  $\lambda$  by avoiding over-regularization at high recall; A5 should demonstrate metric inflation (artificially improved FID from circular evaluation), motivating DINO; A7 should show degraded Coverage compared to reservoir sampling, due to over-representation of recent samples in FIFO; A8–A11 should

TABLE 2

Ablation Study Design (CIFAR-10, DCGAN Backbone). Metrics: FID ( $\downarrow$ ), Recall ( $\uparrow$ ), Coverage ( $\uparrow$ ).

Ablation	Purpose
A1: $\lambda(t) \equiv 0$ (no buffer)	Establishes DCGAN baseline; confirms gap
A2: Fixed $\lambda = 0.1$	Tests adaptive vs. static schedule
A3: Fixed $\lambda = 1.0$	Tests high-strength static schedule
A4: DINO $\phi$ (default)	Evaluates DINO feature space quality
A5: Inception $\phi$	Circular evaluation risk; expected metric inflation
A6: Random $\phi \in \mathbb{R}^{384}$	Tests whether any fixed embedding helps
A7: FIFO buffer, $K = 2048$	Tests reservoir vs. recency-biased sampling
A8: Reservoir, $K = 512$	Buffer size sensitivity (small)
A9: Reservoir, $K = 1024$	Buffer size sensitivity (medium)
A10: Reservoir, $K = 2048$	Default configuration
A11: Reservoir, $K = 4096$	Buffer size sensitivity (large)
A12: $\beta = 0.1$ (soft weights)	Coverage weight temperature
A13: $\beta = 5.0$ (hard weights)	Coverage weight temperature

show a coverage-overhead tradeoff with  $K = 2048$  as a reasonable operating point.

## 7 DISCUSSION

### 7.1 Limitations

MADR-GAN has several limitations that should be carefully considered before deployment.

**Computational overhead.** The DINO forward pass and coverage weight computation add 15–20% overhead compared to a base GAN. While modest, this may be non-trivial in large-scale industrial training runs where wall-clock time is tightly constrained. The overhead scales with both the buffer size  $K$  and the size of the real dataset  $N$ , becoming more pronounced for very large datasets. Approximate nearest-neighbor methods (e.g., FAISS) could reduce the  $\mathcal{O}(N \cdot K)$  coverage computation to  $\mathcal{O}(N \log K)$ .

**Noisy early buffer.** During the first  $K/B$  updates ( $\approx 16$  generator updates for  $K = 2048$ ,  $B = 128$ ), the buffer contains noisy early-training generated samples that may not be meaningfully representative of any real data mode. We mitigate this by enabling the coverage loss only after the buffer contains at least  $K/4$  entries (= 512 samples), but early-phase noise may still influence  $\sigma$  bandwidth estimates.

**Bandwidth sensitivity.** The median heuristic for  $\sigma$  assumes that the relevant scale for coverage assessment is the median pairwise distance in the buffer. This assumption may be violated when the generator is in severe mode collapse (all buffer entries near a single mode), causing  $\sigma$  to be underestimated and coverage weights to be excessively concentrated. A warm-start period using real data pairwise distances would provide a more robust initial bandwidth.

**Target recall hyperparameter.** The choice of  $\tau_{\text{recall}} = 0.80$  is a domain-specific decision. For datasets with genuinely long-tailed mode distributions (e.g., ImageNet), a target of 0.80 may be too ambitious and cause the coverage loss to dominate, destabilizing training. We recommend validating  $\tau_{\text{recall}}$  on a small held-out set before full training.

### 7.2 Comparison to Diffusion Models

Diffusion-based generative models [28], [29] have achieved state-of-the-art FID on most standard benchmarks, significantly outperforming GANs on diversity metrics by design: their denoising objective explicitly reconstructs all modes of the data distribution. Latent Diffusion Models [28] achieve FID 3.6 on LSUN-Churches at  $256 \times 256$ , and EDM [29] achieves FID 1.97 on CIFAR-10. These represent upper bounds that GANs currently cannot match on raw FID.

However, diffusion models suffer from two fundamental disadvantages that GANs do not share. First, inference requires 50–1000 sequential denoising steps, compared to a single forward pass for the GAN generator. This makes diffusion models impractical for latency-sensitive applications (real-time video synthesis, interactive generation). Second, diffusion training is computationally expensive, requiring large compute budgets to match GAN image quality.

MADR-GAN is positioned to capture the best of both paradigms: GAN-speed inference with improved mode coverage. By addressing temporal mode cycling—the primary remaining diversity gap between GANs and diffusion models—MADR-GAN aims to close the recall/coverage gap while preserving the single-step inference advantage. We do not claim MADR-GAN will match diffusion model FID; the goal is to achieve competitive diversity without sacrificing inference speed.

### 7.3 Hybrid and Future Directions

The MADR-GAN coverage framework is architecture-agnostic and could be adapted to other generative paradigms. Consistency models [29] and diffusion-GAN hybrids generate samples in a small number of steps while maintaining mode diversity; integrating MADR-GAN’s coverage loss as an auxiliary objective during distillation could prevent mode dropping during the student training phase.

Several natural extensions exist for future work. **Video generation:** temporal mode diversity requires tracking coverage not just over the spatial distribution but over the joint distribution of video sequences; an episodic buffer of frame sequence features would be needed. **Text-to-image:** semantic mode coverage could be defined in the CLIP [23] embedding space conditioned on the text prompt, ensuring that all text-implied modes are generated. **Online clustering:** rather than using raw DINO features,  $k$ -means or online clustering could be applied to the buffer to identify mode centroids explicitly, reducing the  $\mathcal{O}(N \cdot K)$  coverage computation to  $\mathcal{O}(N \cdot C)$  where  $C \ll K$  is the number of cluster centers.

## 8 CONCLUSION

We have presented MADR-GAN, a Coverage-Guided Memory Replay framework for mitigating mode collapse in Generative Adversarial Networks. Central to our approach is the identification and formal characterization of **temporal mode cycling** as a failure mode distinct from gradient saturation collapse: a phenomenon in which the generator covers different subsets of the data distribution at different points in training, oscillating without converging to full coverage. By maintaining a reservoir-sampled episodic buffer

of historical generated features in the DINO self-supervised feature space and applying a weighted kernel density generator loss biased toward historically undervisited real data modes, MADR-GAN provides a continuous, temporal coverage signal that existing methods lack. An adaptive  $\lambda(t)$  controller modulates the coverage loss strength via an online recall estimate, intervening only when recall is below target and receding when diversity is sufficient. Theoretical analysis confirms that the coverage loss preserves the GAN Nash equilibrium and provides bounded, stable gradient signals. The proposed experimental protocol, with primary metrics focused on Recall, Coverage, and Mode Count rather than FID alone, provides a rigorous foundation for evaluating mode collapse interventions. We believe MADR-GAN represents a principled and practical step toward reliable, mode-diverse adversarial generation.

## REFERENCES

- [1] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2014, pp. 2672–2680.
- [2] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2016.
- [3] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein GAN," in *Proc. Int. Conf. on Machine Learning (ICML)*, 2017, pp. 214–223.
- [4] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of Wasserstein GANs," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5767–5777.
- [5] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training GANs," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2016, pp. 2234–2242.
- [6] Z. Lin, A. Khetan, G. Fanti, and S. Oh, "PacGAN: The power of two samples in generative adversarial networks," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 1498–1507.
- [7] L. Metz, B. Poole, D. Pfau, and J. Sohl-Dickstein, "Unrolled generative adversarial networks," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2017.
- [8] A. Srivastava, L. Valkov, C. Russell, M. U. Gutmann, and C. Sutton, "VEEGAN: Reducing mode collapse in GANs using implicit variational learning," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 3308–3318.
- [9] T. Miyato, T. Kato, M. Koyama, and Y. Ishizuka, "Spectral normalization for generative adversarial networks," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2018.
- [10] T. Karras, S. Laine, and T. Aila, "A style-based generator architecture for generative adversarial networks," in *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 4401–4410.
- [11] T. Karras, S. Laine, M. Aittala, J. Hellsten, J. Lehtinen, and T. Aila, "Analyzing and improving the image quality of StyleGAN," in *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020, pp. 8110–8119.
- [12] A. Brock, J. Donahue, and K. Simonyan, "Large scale GAN training for high fidelity natural image synthesis," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2019.
- [13] H. Zhang, I. Goodfellow, D. Metaxas, and A. Odena, "Self-attention generative adversarial networks," in *Proc. Int. Conf. on Machine Learning (ICML)*, 2019, pp. 7354–7363.
- [14] M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter, "GANs trained by a two time-scale update rule converge to a local Nash equilibrium," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 6626–6637.
- [15] T. Kynkänniemi, T. Karras, S. Laine, J. Lehtinen, and T. Aila, "Improved precision and recall metric for assessing generative models," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 3927–3936.
- [16] M. F. Naeem, S. J. Oh, Y. Uh, Y. Choi, and J. Yoo, "Reliable fidelity and diversity metrics for generative models," in *Proc. Int. Conf. on Machine Learning (ICML)*, 2020, pp. 7176–7185.
- [17] L. Mescheder, A. Geiger, and S. Nowozin, "Which training methods for GANs do actually converge?" in *Proc. Int. Conf. on Machine Learning (ICML)*, 2018, pp. 3481–3490.
- [18] D. Balduzzi, S. Racaniere, J. Martens, J. Foerster, K. Tuyls, and T. Graepel, "The mechanics of n-player differentiable games," in *Proc. Int. Conf. on Machine Learning (ICML)*, 2018, pp. 354–363.
- [19] T. Karras, M. Aittala, J. Hellsten, S. Laine, J. Lehtinen, and T. Aila, "Training generative adversarial networks with limited data," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, pp. 12104–12114.
- [20] H. Thanh-Tung and T. Tran, "Catastrophic forgetting and mode collapse in GANs," in *Proc. Int. Joint Conf. on Neural Networks (IJCNN)*, 2020.
- [21] M. Arjovsky and L. Bottou, "Towards principled methods for training generative adversarial networks," in *Proc. Int. Conf. on Learning Representations (ICLR)*, 2017.
- [22] M. Lucic, K. Kurach, M. Michalski, S. Gelly, and O. Bousquet, "Are GANs created equal? A large-scale study," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2018, pp. 698–707.
- [23] M. Caron, H. Touvron, I. Misra, H. Jégou, J. Mairal, P. Bojanowski, and A. Joulin, "Emerging properties in self-supervised vision transformers," in *Proc. IEEE/CVF Int. Conf. on Computer Vision (ICCV)*, 2021, pp. 9650–9660.
- [24] S. Zhao, Z. Liu, J. Lin, J.-Y. Zhu, and S. Han, "Differentiable augmentation for data-efficient GAN training," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2020, pp. 7559–7570.
- [25] T. Karras, M. Aittala, S. Laine, E. Härkönen, J. Hellsten, J. Lehtinen, and T. Aila, "Alias-free generative adversarial networks (StyleGAN3)," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2021, pp. 852–863.
- [26] M. Kang, J.-Y. Zhu, R. Zhang, J. Park, E. Shechtman, S. Paris, and T. Park, "Scaling up GANs for text-to-image synthesis," in *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2023, pp. 10124–10134.
- [27] A. Sauer, K. Chitta, J. Müller, and A. Geiger, "StyleGAN-T: Unlocking the power of GANs for fast large-scale text-to-image synthesis," in *Proc. Int. Conf. on Machine Learning (ICML)*, 2023.
- [28] R. Rombach, A. Blattmann, D. Lorenz, P. Esser, and B. Ommer, "High-resolution image synthesis with latent diffusion models," in *Proc. IEEE/CVF Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2022, pp. 10684–10695.
- [29] T. Karras, M. Aittala, T. Aila, and S. Laine, "Elucidating the design space of diffusion-based generative models," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2022.
- [30] X. Huang and S. Belongie, "Arbitrary style transfer in real-time with adaptive instance normalization," in *Proc. IEEE/CVF Int. Conf. on Computer Vision (ICCV)*, 2017, pp. 1501–1510.