

# Paper Critique: Diffusion Models Beat GANs on Image Synthesis

Dan Peng

Department of Computer Science

April 11, 2024

danpeng@unc.edu

## 1. Research Problem

### 1.1. What research problem does the paper address?

This paper tackles the challenge of high-quality image synthesis by demonstrating how diffusion models can dethrone GANs as the reigning champions of image generation. The researchers show that with the right architecture and a clever "classifier guidance" technique, diffusion models can generate images that are not just comparable to GANs but actually superior.

### 1.2. What is the motivation of the research work?

The motivation stems from a frustrating paradox in generative modeling: GANs produce stunning images but are notoriously finicky to train and often miss parts of the data distribution. Meanwhile, likelihood-based models like diffusion models are more stable and complete but historically couldn't match GANs' image quality. The authors sought to bridge this gap, aiming to create the best of both worlds—a model with GAN-level quality and diffusion-level reliability.

## 2. Technical Novelty

### 2.1. What are the key technical challenges identified by the authors?

The authors faced a triple threat of challenges: (1) finding the sweet spot in architecture design that would maximize diffusion model performance, (2) inventing a way to let diffusion models trade off diversity for fidelity like GANs can, and (3) tackling the computational elephant in the room—diffusion models require multiple sampling steps while GANs generate images in one go.

### 2.2. How significant is the technical contribution of the paper? If you think that the paper is incremental, please provide references to the most similar work

This work represents a breakthrough moment for image synthesis, not just an incremental improvement. For years, the narrative has been that GANs are unbeatable for photo-realism, while other approaches make different trade-offs. By demonstrating diffusion models achieving state-of-the-art FID scores (the gold standard metric for image quality), the authors flip this narrative on its head. Their classifier guidance technique is particularly innovative, offering a new way to control the quality-diversity trade-off.

### 2.3. Identify 1-5 main strengths of the proposed approach.

- Their architectural refinements serve up a masterclass in model design, showing how attention mechanisms and residual connections can be optimized for diffusion models
- The classifier guidance technique is brilliantly simple yet effective - like adding a GPS to a wandering explorer, it steers the diffusion process toward more realistic outputs
- Unlike the temperamental GANs, these models don't suffer from mode collapse or training instability - they're the steady tortoises that ultimately outpace the hare

### 2.4. Identify 1-5 main weaknesses of the proposed approach.

- The sampling speed remains the Achilles' heel - even with optimizations, generating an image still requires multiple network passes, making real-time applications challenging
- The classifier guidance technique, while powerful, introduces a dependency on labeled data, limiting its ap-

plication to domains where such supervision is available

### 3. Empirical Results

#### 3.1. Identify 1-5 key experimental results, and explain what they signify.

- The guided diffusion model achieves an FID of 4.59 on ImageNet 256×256, handily outperforming BigGAN-deep's 6.95 - signifying that the torch has officially been passed from GANs to diffusion models for state-of-the-art image synthesis
- The models maintain higher "recall" scores even when tuned for quality, indicating they're not just cherry-picking easy examples but actually representing more of the true data distribution than GANs

#### 3.2. Are there any weaknesses in the experimental section (i.e., unfair comparisons, missing ablations, etc)?

The experimental section leaves a few stones unturned. First, while the authors rightfully celebrate beating GAN metrics, they don't fully explore the computational cost difference at inference time—how many more GPU-hours are needed to generate 10,000 images with diffusion versus GANs? Additionally, the paper doesn't investigate the latent space properties of their models compared to GANs, which is crucial for applications like image editing and interpolation. For a complete picture, these practical considerations should have been addressed.

### 4. Summary

The authors have pulled off what many thought impossible - dethroning GANs from their long-held position at the summit of image synthesis. Their classifier guidance technique is particularly elegant, offering a simple yet effective knob to tune the fidelity-diversity trade-off. The visual results speak for themselves - crisp, diverse images that capture details GANs might miss.

### 5. QA Prompt for a Paper Discussion

#### 5.1. Discussion Question

Why are diffusion models fundamentally better suited for capturing diverse data distributions compared to GANs, and how does this theoretical advantage translate to practical improvements?

#### 5.2. Your Answer

Think of GANs as artists trying to learn by having a critic constantly judge their work - they improve drastically where criticism is strongest but might neglect areas

rarely critiqued. Diffusion models, however, work more like restoration experts learning to remove noise from every part of the image equally.

This fundamental difference means diffusion models naturally cover the entire data distribution rather than just focusing on the most visually striking modes. In practice, this translates to generating not just the typical golden retriever in perfect lighting, but also the uncommon angles, unusual breeds, and edge cases that GANs might ignore. The paper's higher "recall" metrics confirm this intuition, showing diffusion models generate more diverse outputs while their classifier guidance technique cleverly bridges the quality gap, giving us the best of both worlds.