

A DEEP DIVE INTO SPECTRAL NORMALIZED DCGAN

ABSTRACT. Image generation using Deep Neural Networks (DNN) represents a compelling yet underexplored frontier within the realm of computer science. The rapid advancements in artificial intelligence, particularly in the field of deep learning[1], have paved the way for novel approaches to image synthesis, leveraging models such as Generative Adversarial Networks (GAN)[2], Variational Auto-Encoders (VAE)[3, 4], and diffusion models[5, 6]. In this study, we undertook the task of implementing a Spectrally Normalized Deep Convolutional GAN (SN DCGAN)[7] from the ground up to train and generate images based on CIFAR-10 datasets[8, 9]. The subsequent section of this paper provides a detailed exposition of the generated images and accompanying analysis. Our findings, presented and compared in this work, aim to uncover more efficient methods for producing authentic and high-quality images.

1. INTRODUCTION

Generating realistic images is a fundamental challenge in computer vision and machine learning, and our study addresses this challenge through the implementation of a specialized model known as the Spectral Normalization Deep Convolutional Generative Adversarial Network (SN DCGAN)[7]. Tailored for proficiency in training on the CIFAR-10 dataset[8, 9], a widely recognized benchmark in computer vision, this model incorporates spectral normalization, a critical enhancement designed to tackle common issues encountered in generative model training, including mode collapse and instability. Spectral normalization plays a pivotal role in stabilizing the discriminator by normalizing the weights of convolutional layers, thereby ensuring Lipschitz continuity. This normalization process prevents the discriminator from excessively reacting to minor variations in the input data, contributing to more robust and reliable training. The primary goal of the SN-DCGAN is to train a generator network capable of producing synthetic images of high quality that closely mirror those present in the CIFAR-10 dataset. CIFAR-10 is a well-established dataset consisting of 60,000 32x32 color images across ten different classes, making it a challenging yet essential dataset for evaluating image generation models. The significance of the SN DCGAN's task extends to a variety of applications, ranging from data augmentation to style transfer and the generation of novel visual content. With our study we aim to make valuable contribution to the broader understanding of deep learning techniques and their implications for artificial intelligence and computer graphics.

Taking evaluation a step further, our study incorporates the Frechet Inception Distance (FID)[10, 11] metric to quantitatively assess the performance of the generator. FID serves as a robust measure, providing insights into both the quality and diversity of the generated samples. Specifically, FID calculates the Frechet distance between the distribution of features extracted from real and generated images, offering a comprehensive measure of how well the generated images match the characteristics of the CIFAR-10 dataset. A lower FID score indicates better performance, suggesting that the generator produces images with characteristics closely aligned to the real dataset. Additionally, the study's flexibility shines through in experiments involving hyperparameter tuning and the comparison of different learning rates. These experiments underscore the versatility of our study in optimizing the training process for generative models, showcasing its adaptability to various scenarios and datasets, including the specific challenges posed by CIFAR-10.

2. BACKGROUND

2.1. Generative Adversarial Networks and Image Synthesis: Generating realistic images has been a fundamental challenge in the fields of computer vision and machine learning. One approach to tackle this challenge is through Generative Adversarial Networks (GANs)[2], a class of models that consist of a generator and a discriminator in a competitive setup. The generator creates synthetic data, and the discriminator evaluates its authenticity. The interplay between these components drives the improvement of the generator over time. Traditional GANs, like DCGAN (Deep Convolutional GAN)[12], have shown promise but often face issues such as mode collapse and instability during training. Spectral Normalization (SN) is a technique introduced to stabilize training by normalizing the weights of convolutional layers, preventing the discriminator from reacting excessively to small changes in input data. Our study presents a Spectral Normalization Deep Convolutional Generative Adversarial Network (SN DCGAN) specifically tailored for training on the CIFAR-10 dataset.

2.2. CIFAR-10 Dataset: The CIFAR-10 dataset[8, 9] is a widely recognized benchmark in computer vision. It consists of 60,000 32x32 color images across ten different classes, with each class containing 6,000 images. The dataset is commonly used for image classification tasks and presents unique challenges due to its small image size. CIFAR-10 serves as the backdrop for evaluating image generation models, and the SN DCGAN aims to train a generator capable of producing synthetic images that closely resemble those found in this dataset. The challenges posed by CIFAR-10, including the diversity of classes and the limited image resolution, make it an ideal test-bed for assessing the performance and generalization capabilities of generative models.

3. RELATED WORK

3.1. Evolution of Generative Adversarial Networks (GANs): The foundation of the Spectral Normalization Deep Convolutional Generative Adversarial Network (SN DCGAN) presented in our study is rooted in the progression of Generative Adversarial Networks (GANs). Initially proposed by Ian Goodfellow et al in 2014[2], GANs introduced a paradigm-shifting approach to generative modeling. However, the evolution of GANs led to the identification of challenges, such as mode collapse and training instability. To address these issues, Spectral Normalization (SN) was introduced by Takeru Miyato et al in 2018[7]. SN DCGAN leverages this advancement to stabilize the training process by normalizing the weights of convolutional layers, thereby enhancing the overall robustness of the model. In the context of the broader literature, our study aligns with recent advancements like Progressive GANs, introduced by Karras et al. in 2018[13]. Progressive GANs aim to increase the resolution of generated images gradually during training, enhancing stability and mitigating challenges associated with mode collapse. The choice of spectral normalization in SN DCGAN complements these efforts by providing an effective means of stabilizing training dynamics, addressing issues encountered in earlier GAN iterations.

3.2. Quantitative Evaluation with Frechet Inception Distance (FID): The incorporation of the Frechet Inception Distance (FID) as the evaluation metric in our study is inspired by the need for a robust quantitative measure of image generation quality. Proposed by Martin Heusel and his team in 2017[11], FID measures the distance between distributions of real and generated images in the feature space of a pre-trained neural network. By adopting FID, the code introduces a sophisticated and widely accepted metric to assess the fidelity and diversity of the synthetic images produced by the SN DCGAN. This choice aligns with contemporary practices in the evaluation of generative models.

In the realm of quantitative evaluation metrics for generative models, our study distinguishes itself by focusing on FID. While alternatives like Inception Score (IS), introduced by Salimans et al in 2016[14], and Kernel Inception Distance (KID), introduced by Binkowski et al in 2018[15], exist, FID has demonstrated effectiveness in capturing both image quality and diversity. This specific choice reflects our study’s commitment to leveraging state-of-the-art evaluation techniques, contributing to the broader understanding of generative model performance.

4. APPROACH AND MODEL ARCHITECTURE

Similar to traditional GANs, our SN-DCGAN comprises two pivotal components: the Discriminator (D) and the Generator (G), both structured as convolutional neural networks. Training for both D and G is conducted within a unified training loop, and their weights are updated after each iteration using the Adam (Adaptive Moment Estimation) optimizer. In alignment with experimental practices from relevant literature, we initialized the random weights for both G and D with mean 0 and standard deviation 0.02 [12]. The first and second-order momentum parameters for Adam are set to (0.5, 0.999) [7]. The chosen loss function for our network is the Binary Cross Entropy (BCE) loss.

In our implementation, D is initially trained to capture the distributions of the training dataset before G commences its training. Subsequently, G is trained to generate synthetic samples, which are then fed back to the D network for binary classification. It is important to note a heuristic aspect: owing to the inherent instability observed in DCGANs, both D and G are trained simultaneously in our implementation. This concurrent training strategy is employed to address the challenges associated with the dynamic interplay between the discriminator and generator in achieving stable convergence.

4.1. Discriminator. D in our model estimates the probability $p(y|x)$ of a sample x being classified with label $y = 1$ (real) or $y = 0$ (fake). Consequently, both real and fake data are incorporated into D in mini-batches during training.

Our D model is composed of 5 convolutional layers, with Spectral Normalization and Batch Normalization applied to each layer. Drawing inspiration from suggestions in the paper by Radford et al. [12], we eschewed pooling layers in favor of strided convolutions (stride = 2). Additionally, we employed LeakyReLU activation for all layers, with the leak set to 0.2. The use of LeakyReLU, characterized by a relatively steep slope, helps mitigate the vanishing gradient issue often encountered in GAN networks. At each convolutional layer, the input size of a CIFAR image (64 x 64) is halved, resulting in an output size of 1 x 1—a scalar value. A Sigmoid function is subsequently applied to this output to rescale it to the probability range of [0, 1]. Conversely, the number of channels doubles through each convolutional layer. This architectural choice is designed to effectively capture and process the features essential for discriminating between real and fake samples.

4.2. Generator. G in our architecture is responsible for generating a fake sample $G(z)$ given a latent space z with a normal distribution of dimensions (100, 1, 1). The G model is constructed with 4 fractional-strided convolutional layers, each equipped with Batch Normalization and the ReLU activation function. The ultimate output is a 3-channel RGB sample ($3 \times 64 \times 64$), mirroring the structure of the training data. To ensure the generated sample falls within the appropriate pixel value range, a Tanh activation function is applied to the output. This design enables G to effectively synthesize realistic samples that resemble the characteristics of the training data.

4.3. Spectral Normalization. Spectral normalization is a regularization technique used to avoid vanishing gradients and stabilize the training of D model [7]. To control the Lipschitz constant of D , spectral normalization is applied to each layer of D . For a linear layer $g(h) = W_h$, the Lipschitz norm $\|g\|_{\text{Lip}} = \sup_h \sigma(\nabla g(h)) = \sup_h \sigma(W) = \sigma(W)$ is the largest singular value of the weight matrix W . Given that the Lipschitz norm of LeakyReLU is 1, we use the inequality $\|g_1 \circ g_2\|_{\text{Lip}} \leq \|g_1\|_{\text{Lip}} \cdot \|g_2\|_{\text{Lip}}$ to obtain the bound on the Lipschitz norm of the D function $\|f\|_{\text{Lip}}$:

$$\|f\|_{\text{Lip}} \leq \prod_{l=1}^{L+1} \|W^l\|_{\text{Lip}} \quad (1)$$

We constrain each layer of D to have the Lipschitz norm equals to 1 with the following spectral normalization [7]:

$$W_{\text{SN}}(W) := \frac{W}{\sigma(W)}. \quad (2)$$

We now have a D function with the Lipschitz norm $\|f\|_{\text{Lip}}$ bounded by 1.

4.4. Loss Function. The following value function $V(G, D)$ captures the competitive nature between G and D [2].

$$\min_G \max_D V(D, G) = \mathbb{E}_{x \sim p_{\text{data}}(x)} [\log D(x)] + \mathbb{E}_{z \sim p_z(z)} [\log(1 - D(G(z)))]. \quad (3)$$

Given sample $x \sim p_{\text{data}}(x)$, we denote $D(x) \in [0, 1]$ the probability of x coming from $p_{\text{data}}(x)$. The probability of distribution for sample $G(z)$ passing through D is denoted as $D(G(z)) \in [0, 1]$. To be consistent with the paper [2], we negate the $-$ sign of the conventional BCE loss and maximize the likelihood instead. This is equivalent to minimizing the BCE loss function. Sample size N is also omitted for simpler explanation.

$$L(y, \hat{y}) = \sum_{i=1}^N [y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i)] \quad (4)$$

For sample $x \sim p_{\text{data}}(x)$, the true label is $y = 1$ and the predicted label is $\hat{y} = D(x)$. The loss simplifies to $L(1, D(x)) = \log(D(x))$, note that the function reaches maximum when $D(x) = 1$. Similarly, for sample $z \sim p_z(z)$ with label $y = 0$ and predicted label $D(G(z))$, the loss $L(0, D(G(z))) = \log(1 - D(G(z)))$ which peaks at $D(G(z)) = 0$. Our objective is for D to maximize the correct labeling for both $x \sim p_{\text{data}}(x)$ and $G(z)$, that is, $\max_D [\log D(x) + \log(1 - D(G(z)))]$. As a result, this will push $D(x)$ being close to 1 (real) and $D(G(z))$ to 0 (fake). This is the exact behavior we expect D to possess.

The objective of G is to fool D such that $D(G(z)) = 1$. This is equivalent to $\min_G [\log(1 - D(G(z)))]$ which drops to $-\infty$ when $D(G(z))$ approaches 1. In practice, this could be problematic, since at the start of the training, G has not learnt the dataset distribution, therefore, $D(G(z))$ is more likely to be close to 0. As a result, when back-propagating, the gradient is also close to 0. A simple modification for the G loss function is $\max_G [\log(D(G(z)))]$ which is applied in our implementation.

5. EXPERIMENTAL RESULTS

Our SN-DCGAN model was implemented and trained using PyTorch. We apply the Adam solver with momentum parameters $\beta_1 = 0.5$ and $\beta_2 = 0.999$. The initial learning rate of $1e-4$ was used in the training for 100 epochs. The learning rate was reduced to $1e-5$ and $1e-6$, with 100 epochs respectively. A batch size of 128 was used. Table 1 shows the performance of the SN-DCGAN on three different learning rates during the training measured by 1) generator loss and discriminator loss, 2) FID score, and 3) discriminator accuracy score.

Learning Rate	$1e-4$	$1e-5$	$1e-6$
Generator Loss	3.0062	2.8577	4.2558
Discriminator Loss	0.3962	0.2060	0.0286
FID Scores	41.0076	117.2166	335.2663
Discriminator Accuracy Scores	0.5971	0.7393	0.5000

TABLE 1. Training results with different learning rates

5.1. Generator Loss and Discriminator Loss. As we concurrently trained the generator and discriminator, we observed an inherent competition between the two entities in the initial phases. A reduction in loss for one corresponds to an increase in loss for the other, reflecting their adversarial nature. Over numerous epochs, both the generator and discriminator progressively improve, leading to the convergence of their losses to relatively stable values. Achieving a state where both the generator and discriminator stabilize and consistently produce desirable results poses a challenge. This challenge arises from the dynamic nature of their relationship: as the generator enhances its performance, the discriminator’s task becomes more arduous, making it increasingly difficult to discern between real and generated images. Figure 1 illustrates the intricate interplay between generator loss and discriminator loss, showcasing their evolving relationship over 100 epochs with a learning rate of $1e-4$.

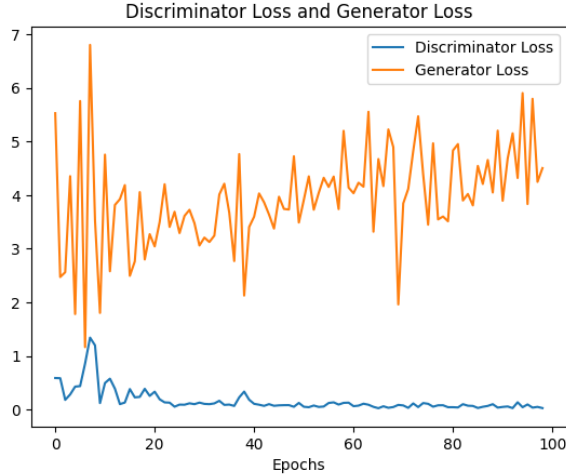


FIGURE 1. A figure showing generator loss and discriminator loss for SN-DCGAN.

5.2. FID Scores. Table 1 displays the FID scores categorized by learning rate. Notably, a higher learning rate of $1e-4$ corresponds to a lower FID score, indicating superior performance. Conversely, a lower learning rate of $1e-6$ exhibits a relatively worse FID score compared to the learning rate of $1e-5$. Considering the inherent nature of FID, where lower scores signify enhanced quality and diversity, we deduce that our SN-DCGAN model achieves optimal performance at a learning rate of $1e-4$ in this experiment.

5.3. Discriminator Accuracy Scores. The training dynamics of the SN-DCGAN model unfold as a min-max game between the generator and discriminator. The objective is for the generator to minimize the min-max GAN loss function, while the discriminator aims to maximize it. When the generator consistently succeeds, the discriminator attains an accuracy of 50% indicative of a scenario where it becomes equally challenging to differentiate between real and generated images.

As illustrated in Table 1, following 100 epochs of training across the three learning rates, the discriminator accuracy scores converge toward the target value of 0.5. This convergence signifies that the images generated by the generator pose a formidable challenge for the discriminator, meeting the training objective of the SN-DCGAN model.

6. DISCUSSION

6.1. Interpretation of Results In the context of the broader literature on generative models, our implementation of the Spectral Normalization Deep Convolutional Generative Adversarial Network (SN DCGAN) has demonstrated promising results in generating realistic images, particularly on the challenging CIFAR-10 dataset. The incorporation of spectral normalization has proven effective in stabilizing the training process and mitigating common issues such as mode collapse and training instability. The competitive dynamics between the generator and discriminator, as reflected in the loss functions and accuracy scores, contribute to the overall success of the model.

6.2. Comparison with Existing Approaches Our approach aligns with the broader evolution of Generative Adversarial Networks (GANs), addressing challenges identified in earlier models. Spectral normalization, introduced to tackle mode collapse and instability, complements recent advancements such as Progressive GANs. The choice of using Frechet Inception Distance (FID) as a quantitative evaluation metric positions our work within the contemporary trend of assessing image generation models. Comparative analysis against alternative metrics, such as Inception Score (IS) and Kernel Inception Distance (KID), could further enrich the evaluation landscape.

6.3. Limitations and Challenges Despite the success observed, our SN DCGAN model exhibits some limitations. One notable challenge is the generation of images with unrealistic color tones and loss of fine details in certain cases. Addressing these issues may require exploring novel algorithms or refining the existing architecture. Additionally, the sensitivity of GANs to hyperparameters necessitates careful consideration, as demonstrated by the varied performance under different learning rates.

6.4. Areas for Improvement If given more time, data, or resources, several avenues for improvement present themselves. Firstly, a more extensive hyperparameter search could be conducted to fine-tune the model’s performance. Exploring alternative architectures or training strategies, such as the use of attention mechanisms or progressive growing techniques, may contribute to enhanced image quality. Furthermore, increasing the dataset size or experimenting with more diverse datasets could improve the model’s generalization capabilities.

6.5. Future Work For a more long-term investment in solving the specific challenges encountered, future work could involve an in-depth exploration of data augmentation techniques tailored to address color tone inconsistencies and detail loss. Additionally, investigating methods to make the model more interpretable and controllable, perhaps through conditional generation or style-based approaches, could open up new avenues for practical applications.

In conclusion, while our SN DCGAN model marks a significant step forward, there remains ample room for refinement and exploration, highlighting the dynamic and evolving nature of research in generative image synthesis.

REFERENCES

- [1] Y. LeCun, Y. Bengio, and G. Hinton. Deep learning. *Intelligent Control and Automation*, 7(4):436–444, 2016.
- [2] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Machine Learning (stat.ML); Machine Learning (cs.LG)*, 26(61):1406, 2014.
- [3] Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. *Machine Learning (stat.ML); Machine Learning (cs.LG)*, 13(12):6114, 2013.
- [4] Diederik P. Kingma and Max Welling. An introduction to variational autoencoders. *Machine Learning (stat.ML); Machine Learning (cs.LG)*, 12(4):307–392, 2019.
- [5] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Machine Learning (cs.LG); Machine Learning (stat.ML)*, 11(23):9, 2006.
- [6] Ling Yang, Zhilong Zhang, Yang Song, Shenda Hong, Runsheng Xu, Yue Zhao, Wentao Zhang, Bin Cui, and Ming-Hsuan Yang. Diffusion models: A comprehensive survey of methods and applications. *Machine Learning (cs.LG); Artificial Intelligence (cs.AI); Computer Vision and Pattern Recognition (cs.CV)*, 22(9):796, 2022.
- [7] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *Machine Learning (cs.LG); Computer Vision and Pattern Recognition (cs.CV)*, 14(41):39, 2018.
- [8] Alex Krizhevsky. Learning multiple layers of features from tiny images. *IEEE Signal Processing Magazine*, 5(8):11, 2012.
- [9] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar-10 (canadian institute for advanced research). *IEEE Signal Processing Magazine*, 10(11):99, 2012.
- [10] Wikipedia contributors. Frechet inception distance. Accessed: 2004-07-22.
- [11] Martin Heusel, Hubert Ramsauer, Thomas Unterthiner, Bernhard Nessler, and Sepp Hochreiter. Gans trained by a two time-scale update rule converge to a local nash equilibrium. *Machine Learning (cs.LG); Computer Vision and Pattern Recognition (cs.CV)*, 17(6):850, 2017.
- [12] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *Machine Learning (cs.LG); Computer Vision and Pattern Recognition (cs.CV)*, 22(50):32, 2015.
- [13] Tero Karras, Timo Aila, Samuli Laine, and Jaakko Lehtinen. Progressive growing of gans for improved quality, stability, and variation. *Neural and Evolutionary Computing (cs.NE); Machine Learning (cs.LG); Machine Learning (stat.ML)*, 17(10):196, 2017.
- [14] Tim Salimans, Ian Goodfellow, Wojciech Zaremba, Vicki Cheung, Alec Radford, and Xi Chen. Improved techniques for training gans. *Machine Learning (cs.LG); Computer Vision and Pattern Recognition (cs.CV); Neural and Evolutionary Computing (cs.NE)*, 16(6):3498, 2016.
- [15] Mikolaj Binkowski, Danica J. Sutherland, Michael Arbel, and Arthur Gretton. Demystifying mmd gans. *Machine Learning (stat.ML); Machine Learning (cs.LG)*, 18(1):1401, 2018.

7. APPENDIX

7.1. Code Repository Reference

During the development of this study, we have referenced a few research papers which are listed in the reference part at the end of this report. Below shows a few code repositories we have referenced for code implementation:

- DCGAN Implementation 1
- DCGAN Implementation 2
- SN-DCGAN Implementation

7.2. CIFAR-10 Image Results and Discussion

Figure 2 illustrates the training outcomes for our SN-DCGAN model on the CIFAR-10 dataset. In Figure (a), we present original CIFAR-10 images spanning various categories, while Figure (b) showcases images generated/reconstructed through the training of our DC-GAN network. As evident in Figure 2(b), the generated CIFAR-10 images exhibit a high degree of authenticity, a noteworthy achievement considering the lower resolution of the original CIFAR-10 images.

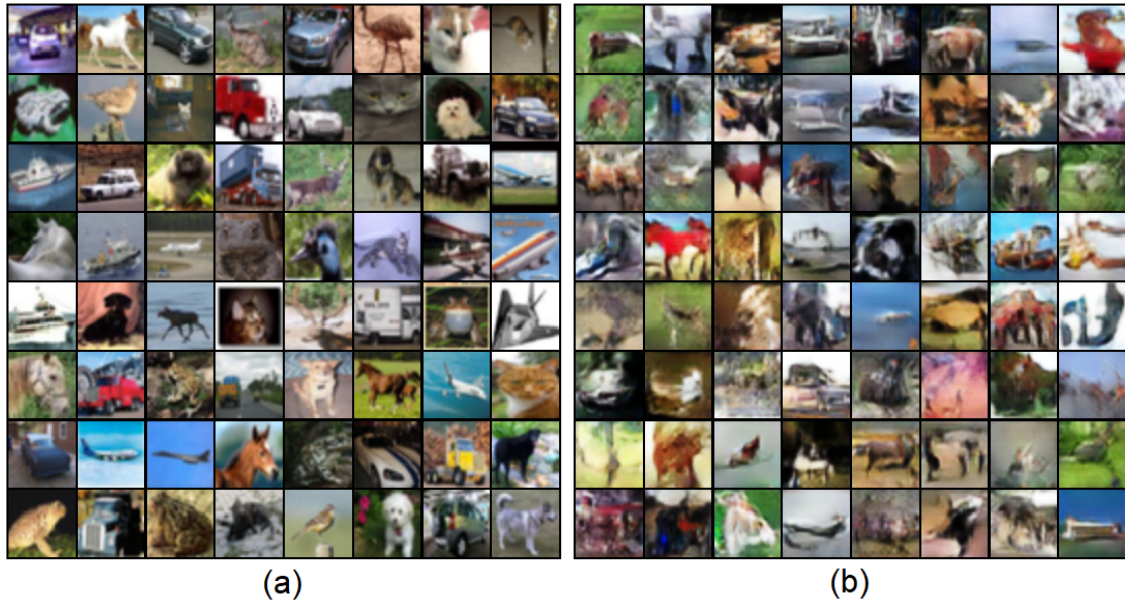


FIGURE 2. A figure showing (a) image inputs from CIFAR-10 dataset; (b) generated/reconstructed image outputs from generative adversarial network trained on CIFAR-10 dataset.