# COMPLEXITY CONTROLLED GENERATIVE ADVERSARIAL NETWORKS

**Himanshu Pant, Jayadeva and Sumit Soman**
Department of Electrical Engineering
Indian Institute of Technology
Delhi, India
jayadeva@ee.iitd.ac.in

## ABSTRACT

One of the issues faced in training Generative Adversarial Nets (GANs) and their variants is the problem of mode collapse, wherein the training stability in terms of the generative loss increases as more training data is used. In this paper, we propose an alternative architecture via the Low-Complexity Neural Network (LCNN), which attempts to learn models with low complexity. The motivation is that controlling model complexity leads to models that do not overfit the training data. We incorporate the LCNN loss function for GANs, Deep Convolutional GANs (DCGANs) and Spectral Normalized GANs (SNGANs), in order to develop hybrid architectures called the LCNN-GAN, LCNN-DCGAN and LCNN-SNGAN respectively. On various large benchmark image datasets, we show that the use of our proposed models results in stable training while avoiding the problem of mode collapse, resulting in better training stability. We also show how the learning behavior can be controlled by a hyperparameter in the LCNN functional, which also provides an improved inception score.

*Keywords* Generative Adversarial Nets · Deep Convolutional Generative Adversarial Nets · Low-Complexity Neural Network · Mode Collapse · Model Complexity

## 1 Introduction

Generative Adversarial Networks (GANs) have recently emerged as an interesting area of research in the deep learning domain. GANs conventionally comprise of a generator and a discriminator. The generator attempts to continually generate fake images, while the discriminator network aims to differentiate between real and fake images. As the network is trained, it is expected that the ability of the generator and discriminator in their respective tasks would both improve, based on the concept of a *"zero sum game."*

In terms of learning theory, Generative Adversarial Networks (GANs) [1] are a class of learning models that involve the training of a generative model $G$ and a discriminative model $D$. $G$ tries to capture the generator distribution and $D$ tries to estimate whether a training sample belongs to the data, rather than the distribution modeled by $G$. Specifically, given a prior $p_z(z)$ on input noise variable and data mapping $P(z; \phi_G)$ and $P(x; \phi_D)$ for the generator and discriminator models (which are conventionally Multi-Layer Perceptrons), GAN attempts to solve the optimization problem given by

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

GANs have tried to overcome the limitations faced by techniques such as Maximum Likelihood Estimation (MLE) in the computation of probabilities (which tend to become intractable). However, they face issues in training stability, and restrictions in visualization when compared to techniques such as Convolutional Neural Networks (CNNs).

Other variants of GANs include Least Squared GAN [2], which uses a least squared loss function for the discriminator in order to address the vanishing gradient problem faced with conventional GANs. LS-GANs show improved image quality

and better training stability compared to the baseline. Variants have also been developed using Integral Probability Metrics (IPM) such as Mean and Covariance Feature Matching GAN [3] and Fisher GAN [4]. Other approaches include Maximum Mean Discrepancy (MMD) [5, 6], margin adaptation [7], Boundary Equilibrium GAN [8], Energy-Based GAN [9], probabilistic GAN [10], Bayesian GAN [11], conditional GAN [12], Wasserstein GAN [13, 14], among others.

Deep Convolutional GANs (DCGANs) were proposed [15] as an approach to implement GANs with Convolutional Neural Networks (CNNs) for unsupervised image classification. The convolutional layers implemented in the generator network has been shown to give better results than conventional GANs. The network proposed in [15] applies convolutional layers to reshape the network via a 'project and reshape' approach and the same has been shown to be beneficial for large-scale scene generation. The use of the CNN-based approach has also been shown to have improved training stability resulting in optimally trained discriminator networks while also facilitating visualization via filters.

One of the primary challenges elucidated in the concluding section of [15] is the problem of *"collapse of a subset of filters to a single oscillating node"* over prolonged training, which is termed as *mode collapse*. In this paper, we discuss a mechanism for addressing this problem in GANs and DCGANs, motivated by an approach for reducing model complexity in neural networks.

Spectral Normalized GANs (SN-GANs) [16] introduce a normalization scheme that uses a single hyperparameter in order to prevents the discriminator from being unstable during training, particularly in high dimensional spaces. The optimal discriminator is given by

$$D_G^*(x) = \text{sigmoid}(f^*(x)) \tag{2}$$

where

$$f^*(x) = log(q_{data}(x) - log(p_G(x))) \tag{3}$$

The spectral normalization scheme adopted by the SN-GAN normalizes the weight matrix $W$ so that its Lipschitz constant $\sigma(W) = 1$, which is realized as

$$\hat{W}_{SN}(W) = \frac{W}{\sigma(W)} \tag{4}$$

The computation of the spectral norm is made computationally efficient in the case of SN-GAN by using the power iteration method. SN-GANs have been shown to give improved results compared to other stabilization and weight normalization techniques in literature.

## 2 Low Complexity Neural Network

The LCNN [17] aims to introduce the notion of controlling model complexity in a neural network by modifying the cost function. Specifically, we minimize the error with explicit $L_2$ regularizers on weights along with the model complexity term.

Consider a family of large margin neural network classifiers in which each has $k$ layers and one output neuron, where all neurons have real valued weights and a constant bias input each. All neurons use a continuous, monotonically non-decreasing activation function. The class of any input sample is determined by thresholding the output of the neuron in the final layer keeping the weights in previous layers fixed. Let the number of neurons in the penultimate layer be $n$. Let the margin of the classifier obtained by thresholding the output neuron be at least $d_{min} > 0$. Then, the VC dimension $\gamma$ of this neural network satisfies

$$\gamma \leq 1 + \text{Min}\left(\frac{4R^2}{d_{min}^2}, n\right) \tag{5}$$

where $R$ denotes the radius of the smallest sphere enclosing $\left\{V_{(k-1),1}^i, V_{(k-1),2}^i, ..., V_{(k-1),n}^i\right\}$, $i = 1, 2, ..., M$.

The set of input samples $x^i$, $i = 1, 2, ..., M$, has a finite radius $R_{input} < \infty$. The class of any input sample is determined by thresholding the output of the neuron in the final layer. The outputs of the penultimate layer neurons when the $i - th$ sample is presented at the input, are denoted by $\{V_{(k-1),1}^i, V_{(k-1),2}^i, ..., V_{(k-1),n}^i\}$. The weight vector of the output neuron is denoted by $w$, and its bias input is denoted by $b$. In other words, when the $i - th$ training sample

is presented at the input, the net input to the final layer neuron is given by

$$net^i = \sum_{j=1}^{n} w_j \, V_{(k-1),j}^i \; + \; b = \mathbf{w}^T \mathbf{V}_{(k-1)}^i + b, \tag{6}$$

and the class predicted by the classifier is determined from the sign of $net^i$. Further, assume that any member of the family of classifiers so defined has a margin of at least $d_{min}$, where $d_{min}$ is strictly positive. Then, the VC dimension $\gamma$ of this classifier is bounded from above by

$$\gamma \le 1 + \text{Min}\left( C \sum_{i=1}^{M} (net^i)^2 \; , \; n \right), \tag{7}$$

where $C$ is a real positive constant.

For binary classification with labels $y_i \in \{-1, 1\}$ or $y_i \in \{0, 1\}$, the empirical error $E_{emp}$ is given as

$$E_{emp} = \left[ \frac{1}{M} \sum_{i=1}^{M} \left( -log(\text{softmax}(net^i)) \right) \right] \tag{8}$$

where $M$ denotes the number of training samples and $f(net^i)$ denotes the activation function $f(\cdot)$ on net input to neuron $i$ given by $net^i$.

In regularized neural networks, this is modified by adding a term proportional to $\|w\|^2$ at individual weights (weight decay). In order to minimize model complexity, the Low Complexity Neural Network (LCNN) classifier uses the modified error functional $E$ of the form

$$E = E_{emp} + C \sum_{i=1}^{M} (net^i)^2 \tag{9}$$

## 3 GANs with LCNN

As mentioned in the Introduction, the generator and discriminator in a GAN compete for improved performance. The generator loss determines the performance of the generator, while the discriminator loss is indicative of the error being made in distinguishing the real and adversarial examples. One of the primary reasons for mode collapse and unstable training is the divergence of the generator and discriminator losses over the course of their training. Our approach, based on the LCNN, involves introducing an additional term in the objective function of the GAN which addresses the issue of mode collapse. The motivation is that controlling model complexity leads to models that do not overfit the training data.

We add the LCNN loss functional in the final layer of the generator and discriminator which allows for minimization of the activations of the neurons in the final layer, driving them closer to the origin. During the course of training, when the difference of the error between the discriminator and generator diverges, the LCNN term $(\min \|w^T x + b\|)$ tries to maintain the balance between discriminator and generator errors by pushing activations generated by both real and generated samples towards minimum, hence making it difficult for the discriminator to distinguish between real and generated samples. This improves stability and the quality of generated images as more and more data are presented to the GAN. Also, since we conventionally use the tan-sigmoid or log-sigmoid activation functions, and these have maximum gradient near the origin, the vanishing gradient problem is mitigated during backpropagation, which also results in improved training stability.

### 3.1 LCNN-GAN

The modified objective function for the LCNN-GAN is obtained by adding an additional term to the conventional objective function of GANs which can be represented as

$$\min_G \max_D \mathbb{E}_{x \sim p_{data}(x)} log[D(x)] + \mathbb{E}_{z \sim p_z(z)} log[1 - D(G(z))] +$$
$$\min_D (C_1 L_C(D(x)) + C_2 L_C(D(G(x)))) \tag{10}$$

3

where $L_C$ is the LCNN term applied on the last layer of $D$. Specifically, $L_C = \sum_{i=1}^{M} (net^i)^2$ for the discriminator. Here $C_1$ and $C_2$ are hyper-parameters that can be tuned by the user. We call the model using the objective function above as the LCNN-GAN.

Similar analogy can also be extended for obtaining the hybrid LCNN-incorporated objective functions for the DC-GANs and for SN-GANs in order to obtain the LCNN-DCGAN and LCNN-SNGAN respectively. The final layers in these hybrid models use the softmax activation function with cross entropy loss.

## 4   Experiments

We now present results on the LCNN incorporated approaches for GANs, DC-GANs and SN-GANs in the following subsections. The conventional architectures proposed in literature for these GAN models have been used in our experiments. The results are presented in terms of generator and discriminator loss, as well as in terms of representative images. Subsequently, we also present results in terms of inception score for the various LCNN-incorporated GAN models discussed in this paper. For these experiments, the values of hyper-parameters have been set as the same for both $C_1$ and $C_2$ and have been indicated as $C$ in the plots comparing the generator and discriminator loss.

### 4.1   Results on LCNN-GAN

In Fig. 1 we compare mode collapse for GAN in Fig. 1a and the case for GAN with LCNN in Fig. 1b for the MNIST dataset. It can be seen that the generative loss is contained in case of the LCNN-GAN, as opposed to the conventional GAN, thereby avoiding mode collapse.



(a) Mode Collapse in GAN                                    (b) LCNN-GAN
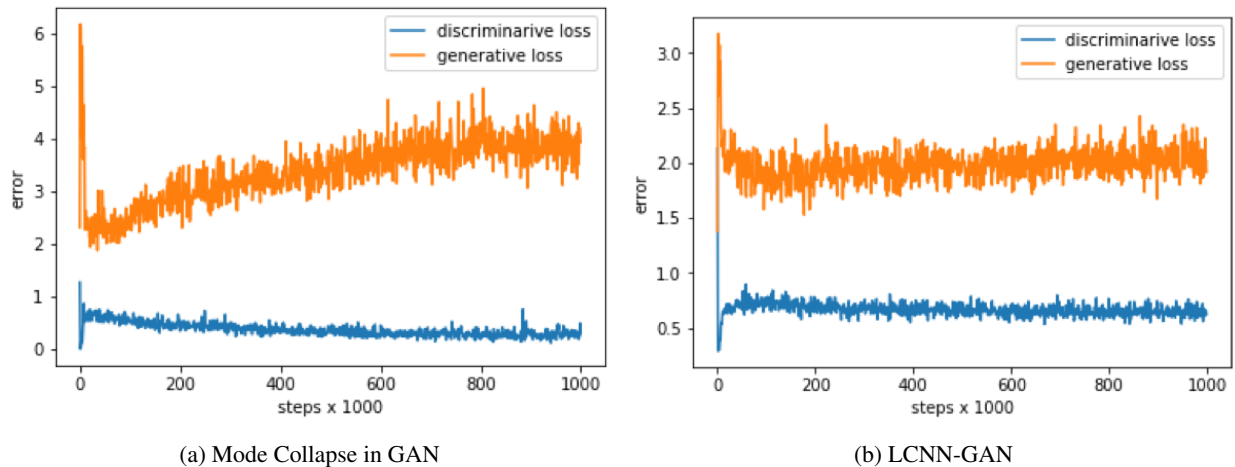
Figure 1: Mode Collapse Analysis for MNIST Data

We also show the fake images generated by the GAN and LCNN-GAN to illustrate how mode collapse is avoided in case of LCNN-GAN. As can be seen from Fig. 3 for the LCNN-GAN, the images generated as the number of mini batches are increased (Figs. 3a - 3f) retains the characteristics of the individual digits. However, in case of the conventional GAN, with an increase in the number of mini batches processed, (Figs. 2a - 2d), mode collapse results in all images being similar to the digit 1. This illustrates that mode collapse is avoided in case of the LCNN-GAN.
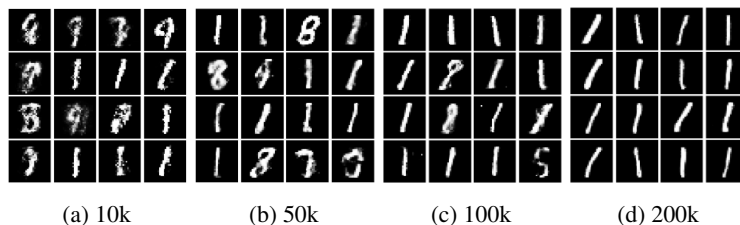


(a) 10k              (b) 50k              (c) 100k              (d) 200k

Figure 2: GAN Generated Fake MNIST Samples at different training stages

(a) 10k     (b) 50k     (c) 100k     (d) 200k     (e) 500k     (f) 1000k
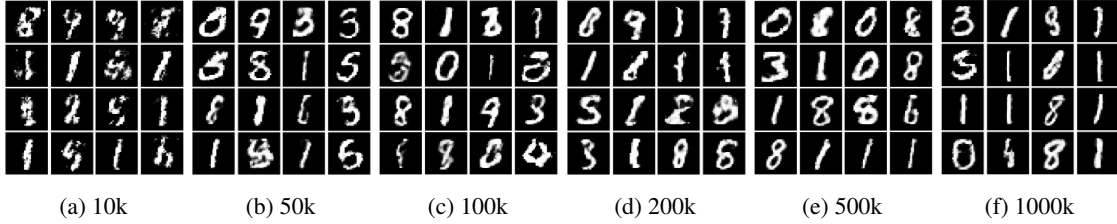
Figure 3: GAN-LCNN Generated Fake MNIST Samples at different training stages (numbers indicate mini batches processed)

We discuss how the hyperparameter C in the LCNN-GAN affects balance between the generative and discriminative loss. In Fig. 4 we compare the GAN and LCNN-GAN for C=0.01. It can be seen that the loss is contained in case of LCNN-GAN.
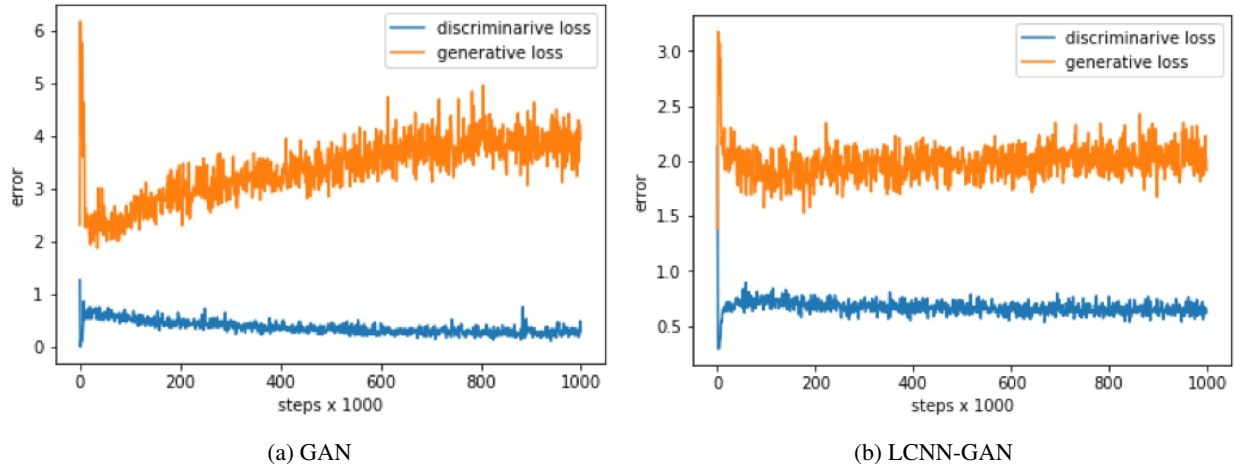


(a) GAN                          (b) LCNN-GAN

Figure 4: Balance between generator and discriminator for C=0.01

The effect of hyperparameter C can also be seen in Fig. 5, where we compare the loss for C=0.1 and C=0.2 in Figs. 5a and 5c respectively. Representative images from the MNIST dataset are also shown alongside in Figs. 5b and 5d respectively. It can be seen that the hyperparameter $C$ plays an important role in balancing the error between generator and discriminator. Also, by varying $C$ during training we can control generator and discriminator behavior.

Since MNIST has 10 classes, most of the energy (variance) content of MNIST will be encompassed by top 10 eigenvalues. On doing PCA on MNIST train data set it was found that top 10 eigenvalues correspond to 83% energy. For analyzing variance and modes generated in GAN and LCNN-GAN in each iteration, we calculated how many eigenvectors correspond to 83% energy content for GAN and LCNN-GAN. In a way, we are deducing that how many modes GAN and LCNN-GAN are capable of generating in fake images.

Results are shown in Fig.6, for GAN in Fig. 6a and LCNN-GAN in Fig. 6b. It can be seen that the LCNN-GAN retains many more modes when compared to GAN, providing further evidence in favor of addressing the mode collapse problem faced in GANs.

## 4.2 Results on LCNN-DCGAN

### 4.2.1 Dataset Used

CelebFaces Attributes Dataset (CelebA) is a large-scale face attributes dataset with more than 200K celebrity images, each with 40 attribute annotations. The images in this dataset cover large pose variations and background clutter. CelebA has large diversities, large quantities, and rich annotations, including

1. 10,177 identities
2. 202,599 face images
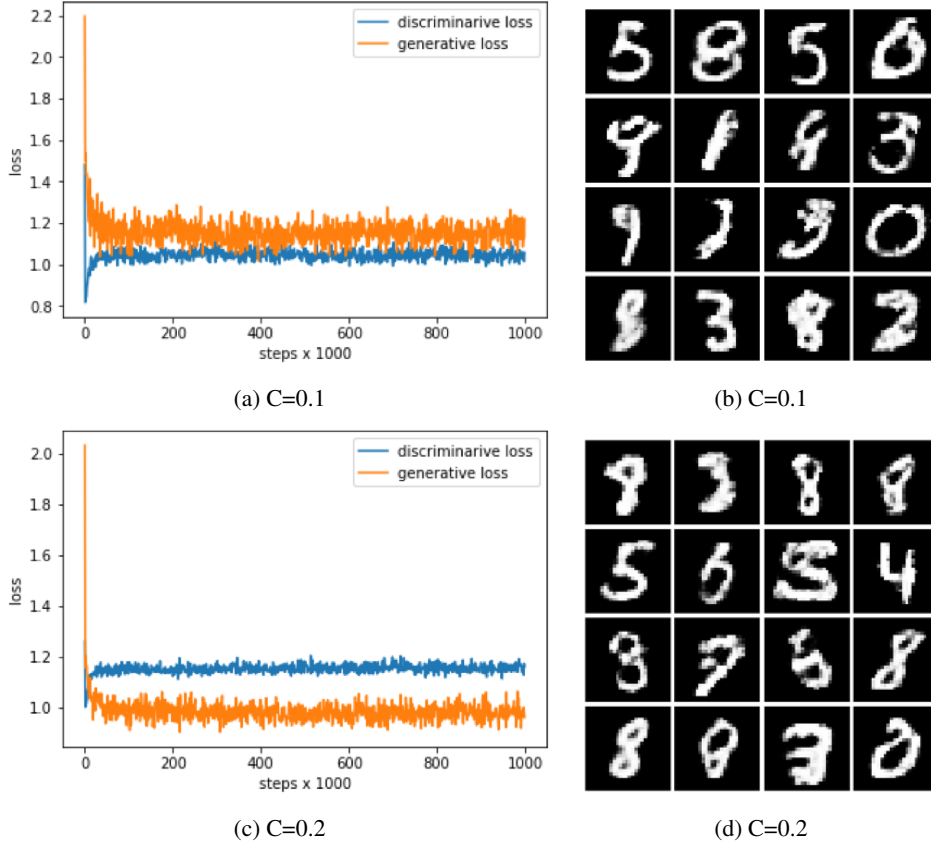
(a) C=0.1



(b) C=0.1



(c) C=0.2



(d) C=0.2

Figure 5: Balance between generator and discriminator for different values of C for LCNN-GAN

3. 5 landmark locations, 40 binary attributes annotations per image.

In Fig. 7 we show the training stability for the DCGAN and LCNN-DCGAN in Figs. 7a and 7b respectively. It can be seen that in case of the DCGAN, the generative and discriminative losses tend to rapidly vary, while it is contained in case of the LCNN-DCGAN. It may be noted here that the values of hyperparameter C have varied linearly from 0.01 to 0.1 for the LCNN-DCGAN plot.

In Fig. 8, we show representative images from the CelebA dataset for the case of DCGAN before (Figs. 8a-8c) and after (Figs. 8d-8f) mode collapse, and compare them with those obtained from the LCNN-DCGAN (Figs. 8g-8i). It can be seen that the mode collapse occurring in case of DCGANs is controlled in case of using the LCNN-DCGAN.

## 4.3 CIFAR-10

As a baseline, we show the generative and discriminative loss for the DCGAN on the CIFAR-10 dataset in Fig. 9.

Further, fake images generated using the DCGAN before and after mode collapse are also shown in Fig. 10. It can be seen that after mode collapse (10c-10e), the images tend to represent only few of the classes.

We now compare the generative and discriminative loss in case of LCNN-DCGAN in Fig. 11. The training stability is compared for different values of C and it can be seen that compared to the baseline training stability for DCGAN, there is relatively lower generative loss in case of using the LCNN-DCGAN.

In terms of images from the CIFAR-10 dataset, few representative images are shown in Fig. 12 for the LCNN-DCGAN. It can be seen that mode collapse is avoided when compared to conventional DCGAN.
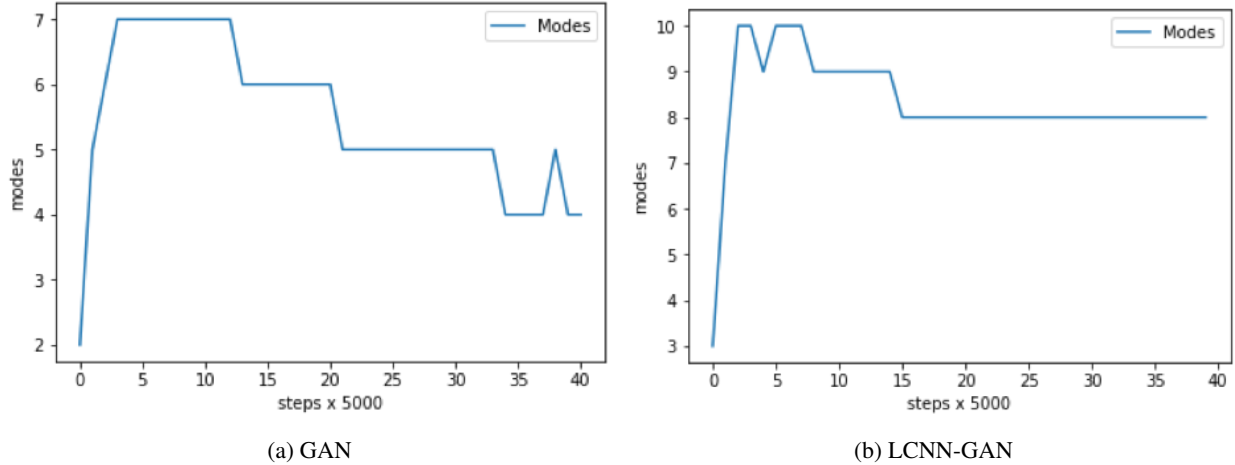
6

(a) GAN          (b) LCNN-GAN

Figure 6: Energy Content and Modes study for GAN and LCNN-GAN
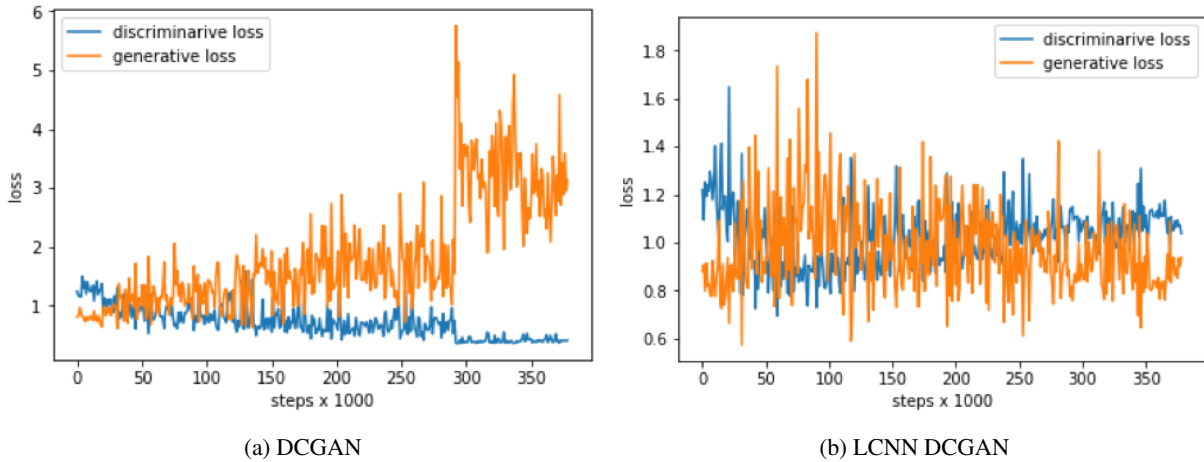


(a) DCGAN          (b) LCNN DCGAN

Figure 7: Training Stability for DCGAN and LCNN-DCGAN

## 4.4 SVHN

We finally present results on the SVHN dataset. The baseline generative and discriminative loss are shown in Fig. 13 for the DCGAN. Further, representative images for the SVHN dataset are shown in Figs. 14a-14b and 14c-14d respectively.

In case of the LCNN-DCGAN, the training stability for C=0.01 is shown in Fig. 15. It can be seen that generative loss is significantly contained in case of the LCNN-DCGAN.

Finally, representative images from the SVHN dataset for the LCNN-DCGAN are shown in Fig. 16. It can be seen that mode collapse is avoided when compared to the baseline images obtained using DCGAN.

## 4.5 Quantitative Evaluation of LCNN-GAN models

The inception score is a quantitative metric for determining the performance of GANs. It can be used to determine whether the images have variety and how close these images represent actual images. The score is high when both these factors are true, and low when the converse holds. The Inception Score is based on Google's inception classifier and the lowest possible score is zero. In essence, the inception score compares the label distribution of each image with the entire set of images to determine the difference between these distributions in terms of the Kullback-Leibler (KL) divergence.

Table 1 compares the mean and standard deviation of inception scores for various models for the CIFAR-10 dataset. The inception score has been computed on 5000 generated images which have been split into ten parts and mean and

7

(a) DCGAN 1 (Before)     (b) DCGAN 2 (Before)     (c) DCGAN 3 (Before)

(d) DCGAN 1 (After)     (e) DCGAN 2 (After)     (f) DCGAN 3 (After)

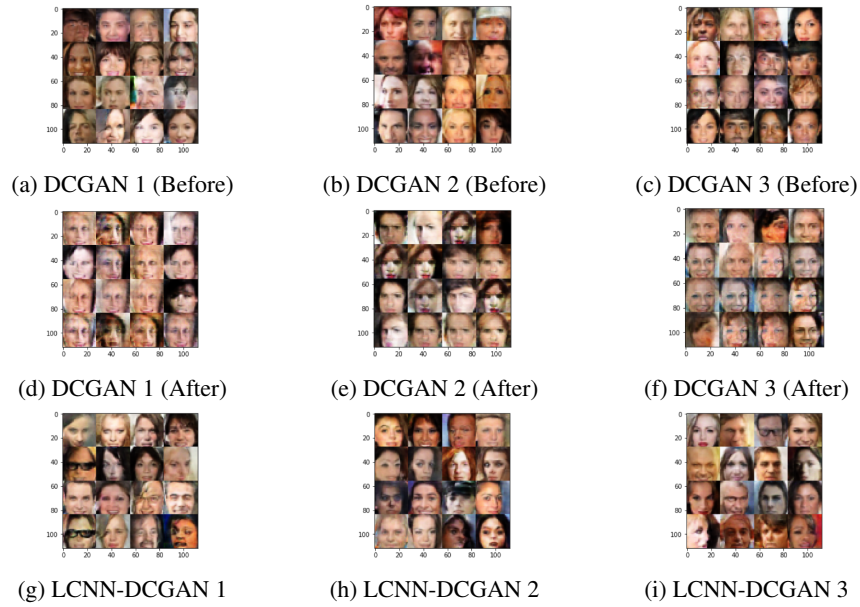(g) LCNN-DCGAN 1     (h) LCNN-DCGAN 2     (i) LCNN-DCGAN 3

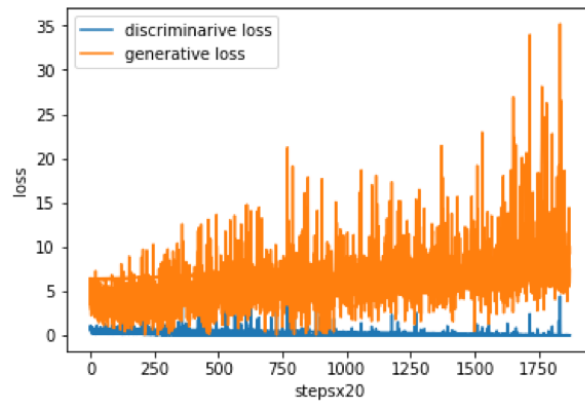Figure 8: Fake Images Generated from DCGAN before and after Unstable/Mode Collapse compared with LCNN-DCGAN



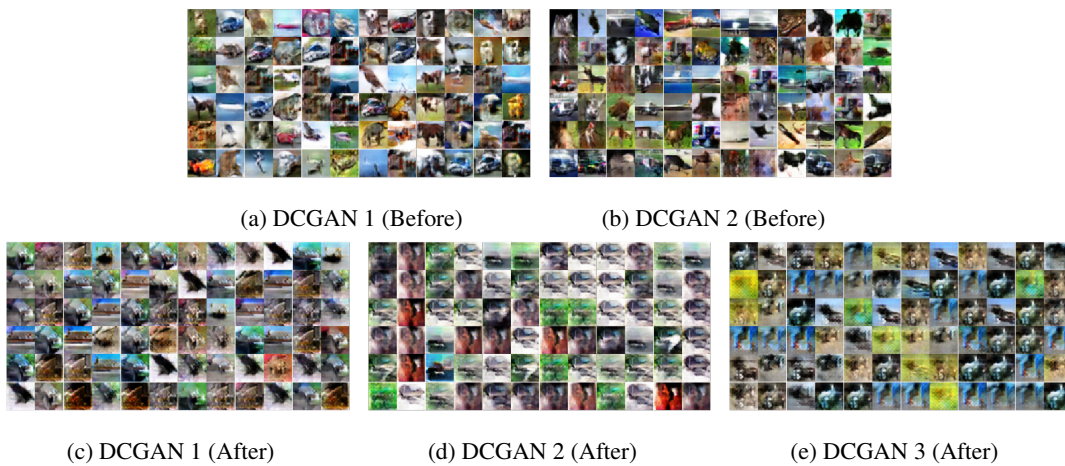Figure 9: Generative and Discriminative Loss for DCGAN on CIFAR-10



(a) DCGAN 1 (Before)     (b) DCGAN 2 (Before)

(c) DCGAN 1 (After)     (d) DCGAN 2 (After)     (e) DCGAN 3 (After)

Figure 10: Fake Images Generated from DCGAN Before and After Unstable/Mode Collapse for CIFAR-10 dataset

(a) C=0.1
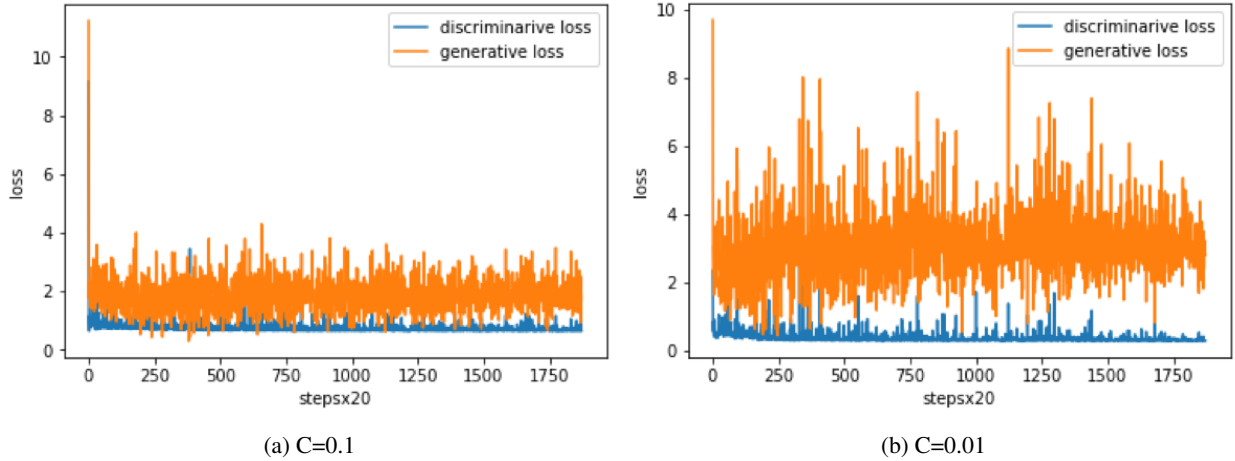
(b) C=0.01

Figure 11: Training Stability for LCNN-DCGAN for different values of C.



(a) LCNN-DCGAN 1      (b) LCNN-DCGAN 2      (c) LCNN-DCGAN 3



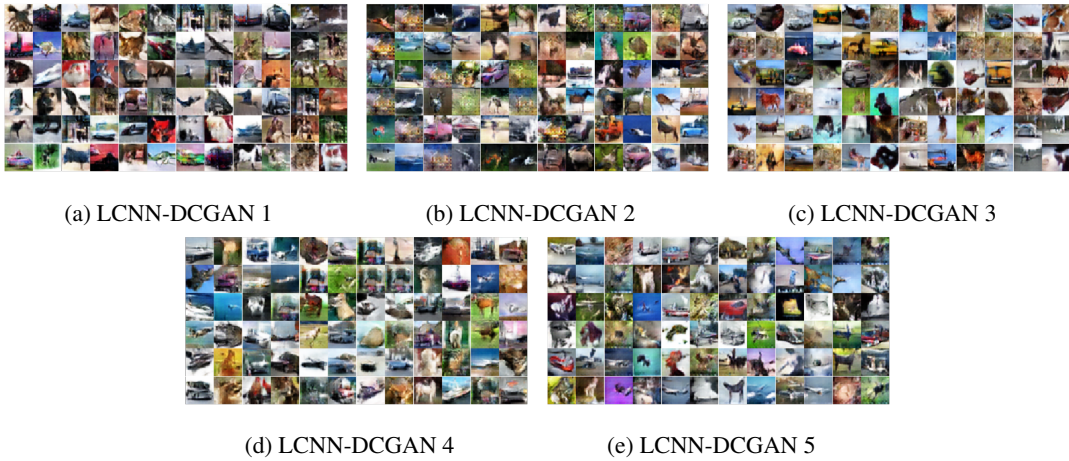(d) LCNN-DCGAN 4      (e) LCNN-DCGAN 5

Figure 12: Fake Images Generated from LCNN-DCGAN for CIFAR-10 dataset

standard deviation for inception scores have been reported. The corresponding learning rate has also been indicated for reference.

Table 1: Results on inception score for CIFAR-10

| S.No. | Network Architecture | Mean | Std. Dev. | Learning Rate |
|---|---|---|---|---|
| 1 | DCGAN | 5.8 | 0.25 | 0.0005 |
| 2 | DCGAN-LCNN (C=0.01) | 6.04 | 0.16 | 0.0005 |
| 3 | DCGAN-LCNN (C=0.05) | 6.33 | 0.2 | 0.0005 |
| 4 | DCGAN-LCNN (C=0.5) | 6.35 | 0.2 | 0.0005 |
| 5 | DC-SNGAN | 7.4 | 0.26 | 0.0002 |
| 6 | DC-SNGAN-LCNN (C=0.01) | 7.68 | 0.24 | 0.0002 |
| 7 | DC-SNGAN-LCNN (C=0.05) | 7.63 | 0.21 | 0.0002 |

## 5 Conclusion

This paper shows that controlling model complexity in GANs significantly alleviates the problem of mode collapse. We use the LCNN model complexity control approach to achieve this objective. Hybrid models obtained by incorporating the LCNN loss function into various GAN models, including the conventional GAN, DC-GAN and SN-GAN yield demonstrable and consistent improvements and do not show mode collapse despite prolonged training. Our approach has
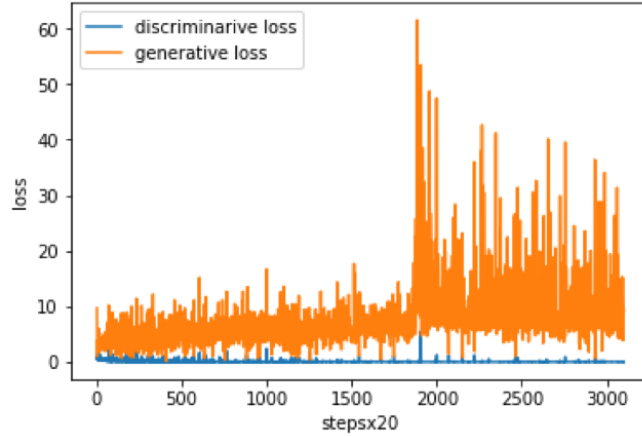
Figure 13: DCGAN: Mode Collapse/Unstable after 1900x20 mini Batches.



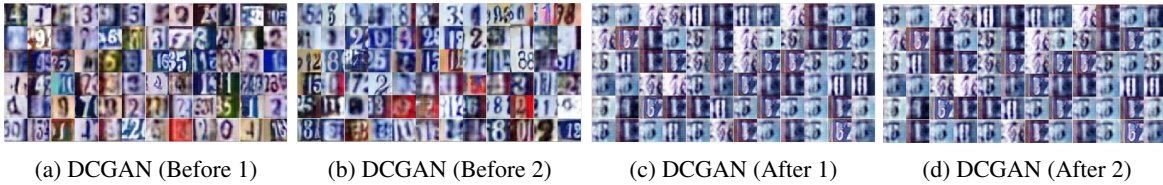(a) DCGAN (Before 1)     (b) DCGAN (Before 2)     (c) DCGAN (After 1)     (d) DCGAN (After 2)

Figure 14: Fake Images Generated from DCGAN Before Unstable/Mode Collapse for SVHN Dataset

shown to provide better control over generator and discriminator losses during training by means of the hyperparameter used in the LCNN functional. The claims have been established by results in terms of modes of generated images, generator/discriminator losses; the inception scores for benchmark datasets show that the use of the LCNN-based loss function avoids mode collapse and provides improved performance across the board for the proposed hybrid models.

## References

[1] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.

[2] Xudong Mao, Qing Li, Haoran Xie, Raymond YK Lau, Zhen Wang, and Stephen Paul Smolley. Least squares generative adversarial networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2794–2802, 2017.

[3] Youssef Mroueh, Tom Sercu, and Vaibhava Goel. Mcgan: Mean and covariance feature matching gan. *arXiv preprint arXiv:1702.08398*, 2017.

[4] Youssef Mroueh and Tom Sercu. Fisher gan. In *Advances in Neural Information Processing Systems*, pages 2513–2523, 2017.

[5] Chun-Liang Li, Wei-Cheng Chang, Yu Cheng, Yiming Yang, and Barnabás Póczos. Mmd gan: Towards deeper understanding of moment matching network. In *Advances in Neural Information Processing Systems*, pages 2203–2213, 2017.

[6] Yujia Li, Kevin Swersky, and Rich Zemel. Generative moment matching networks. In *International Conference on Machine Learning*, pages 1718–1727, 2015.

[7] Ruohan Wang, Antoine Cully, Hyung Jin Chang, and Yiannis Demiris. Magan: Margin adaptation for generative adversarial networks. *arXiv preprint arXiv:1704.03817*, 2017.

[8] David Berthelot, Thomas Schumm, and Luke Metz. Began: boundary equilibrium generative adversarial networks. *arXiv preprint arXiv:1703.10717*, 2017.

[9] Junbo Zhao, Michael Mathieu, and Yann LeCun. Energy-based generative adversarial network. *arXiv preprint arXiv:1609.03126*, 2016.
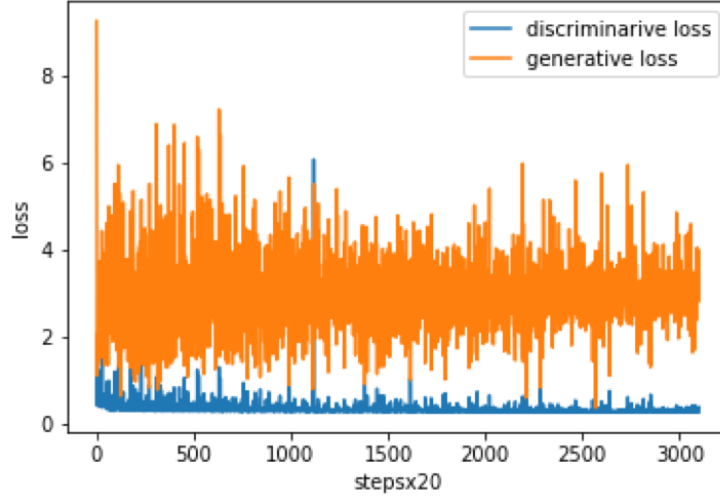
Figure 15: LCNN DCGAN: Stable Training Value of C = 0.01



(a) LCNN-DCGAN 1 (b) LCNN-DCGAN 2 (c) LCNN-DCGAN 3

Figure 16: Fake Images Generated from LCNN-DCGAN

[10] Hamid Eghbal-zadeh and Gerhard Widmer. Probabilistic generative adversarial networks. *arXiv preprint arXiv:1708.01886*, 2017.

[11] Yunus Saatci and Andrew G Wilson. Bayesian gan. In *Advances in neural information processing systems*, pages 3622–3631, 2017.

[12] Mehdi Mirza and Simon Osindero. Conditional generative adversarial nets. *arXiv preprint arXiv:1411.1784*, 2014.

[13] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein gan. *arXiv preprint arXiv:1701.07875*, 2017.

[14] Marc G Bellemare, Ivo Danihelka, Will Dabney, Shakir Mohamed, Balaji Lakshminarayanan, Stephan Hoyer, and Rémi Munos. The cramer distance as a solution to biased wasserstein gradients. *arXiv preprint arXiv:1705.10743*, 2017.

[15] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015.

[16] Takeru Miyato, Toshiki Kataoka, Masanori Koyama, and Yuichi Yoshida. Spectral normalization for generative adversarial networks. *arXiv preprint arXiv:1802.05957*, 2018.

[17] Himanshu Pant, Mayank Sharma, Abhimanyu Dubey, Sumit Soman, Suraj Tripathi, Sai Guruju, Nihal Goalla, et al. Learning neural network classifiers with low model complexity. *arXiv preprint arXiv:1707.09933*, 2017.