# INTERNATIONAL JOURNAL
# FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

# Text-to-Image Generation Using Stack Generative Adversarial Networks (GANs) and Stable Diffusion Models

Shuhaab Shafi[1], Sumit Patil[2], Saad Sayyed[3], Pratiksha Wagh[4]

*Department of Computer Science and Business System, Computer Science and Technology, MIT WPU University, Pune India*

*Abstract: Text-to-image synthesis is a challenging task in artificial intelligence that involves generating realistic images from textual descriptions. StackGAN, [1] a generative model with a two-stage architecture, has proven to be effective in improving the quality and fidelity of synthesized images. In this study, we implement and evaluate the performance of StackGAN using the CIFAR-10 dataset, demonstrating how a multi-stage generation approach can produce high-quality images from textual labels. The experimental results show that StackGAN successfully captures the underlying features of the dataset categories and provides reasonable image fidelity.*
*Keywords: Stack-GAN (Generative Adversarial Networks), CIFAR-10*

## I. INTRODUCTION

Text-to-image synthesis refers to generating realistic images that correspond to given textual descriptions. Traditional generative models often struggle to produce high-quality images due to the complexity of aligning text features with visual information. The introduction of StackGAN (Stacked Generative Adversarial Network) offers a multi-stage generation process that refines the images progressively to overcome these challenges.

The CIFAR-10 dataset, consisting of 60,000 images across 10 categories such as airplanes, automobiles, birds, and animals, provides an ideal playground for testing the effectiveness of text-to-image models. This research focuses on using StackGAN to synthesize images conditioned on CIFAR-10 labels. The

main objective is to assess the feasibility of StackGAN on small-sized datasets and investigate its ability to generate sharp and visually coherent images. StackGAN was proposed by Zhang et al. (2017) as a two-stage GAN architecture, where the first stage generates low-resolution images, and the second stage refines them for higher quality. Several studies have applied StackGAN to datasets like CUB-200 and Oxford-102 Flowers, showing the effectiveness of the model in generating detailed visuals. However, the application of StackGAN to CIFAR-10 is relatively unexplored, making this study a novel contribution.
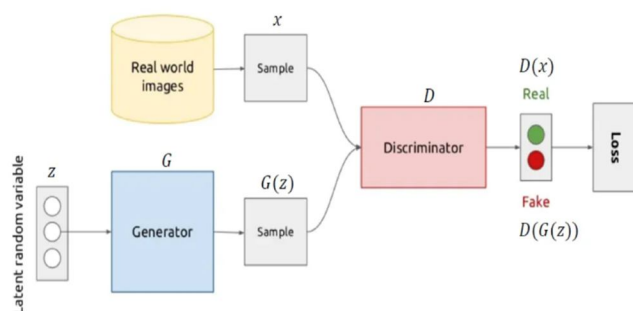


Fig1. Generative Adversarial Networks architecture

## II. METHODOLOGY

### A. Dataset: CIFAR-10

CIFAR-10 is a widely used dataset in machine learning research, consisting of 10 classes: airplane, automobile, bird, cat, deer, dog, frog, horse, ship, and truck. Each class contains 6,000 images with a resolution of 32x32 pixels. For this experiment, each image label serves as the textual input to the StackGAN.

*B. StackGAN Architecture*

*1) StackGAN Consists of two Stages*

Stage I: Generates a low-resolution (64x64) image conditioned on text input. Uses a text encoder (such as RNN or LSTM) to extract text embeddings. Employs a Conditional GAN (cGAN) to generate the first version of the image.

*2) Mathematical model behind StackGAN*

x: Real data

z: Latent vector

G(z): Fake data

D(x): Discriminator's evaluation of real data

D(G(z)): Discriminator's evaluation of fake data

Error(a,b): Error between a and b

*3) Discriminator*

The goal of the discriminator is to correctly label generated images as false and empirical data points as true. Therefore, we might consider the following to be the loss function of the discriminator:

$L_D = Error(D(x),1) + Error(D(G(z)),0)$

*4) Generator*

We can go ahead and do the same for the generator. The goal of the generator is to confuse the discriminator as much as possible such that it mislabels generated images as being true.

$L_G = Error(D(G(z)),1)$

The key here is to remember that a loss function is something that we wish to minimize. In the case of the generator, it should strive to minimize the difference between 1, the label for true data, and the discriminator's evaluation of the generated fake data.

*5) Training the discriminator*

When training a GAN, we typically train one model at a time. In other words, when training the discriminator, the generator is assumed as fixed.

The quantity of interest can be defined as a function of GG and DD. Let's call this the value function:

$V(G,D) = \mathbb{E}_{x\sim pdata}[\log(D(x))] + \mathbb{E}_{z\sim pz}[\log(1-D(G(z)))]$

In reality, we are more interested in the distribution modeled by the generator than pz. Therefore, let's create a new variable, $y=G(z) y=G(z)$, and use this substitution to rewrite the value function:

$V(G, D) = \mathbb{E}_{x \sim pdata}[\log(D(x))] + \mathbb{E}_{y\sim pg}[\log(1-D(y))]$

$= \int_{x\in\chi} pdata(x) \log(D(x)) + pg(x) \log(1-D(x)) \, dx$

*6) Training the Generator*

To train the generator, we assume the discriminator to be fixed and proceed with the analysis of the value function.

$V(G, D*) = \mathbb{E}_{x\sim pdata}[\log(D*(x))] + \mathbb{E}_{x\sim pg}[\log(1-D*(x))]$

$= \mathbb{E}_{x\sim pdata}[\log pdata(x) \, pdata(x) + pg(x)] + \mathbb{E}_{x\sim pg}[\log pg(x) \, pdata(x) + pg(x)]$

To proceed from here

$V(G,D*) = \mathbb{E}_{x\sim pdata}[\log pdata(x) \, pdata(x) + pg(x)] + \mathbb{E}_{x\sim pg}[\log pg(x) \, pdata(x) + pg(x)] = -\log4 + \mathbb{E}_{x\sim pdata}[\log pdata(x) - \log pdata(x) + pg(x))2] + \mathbb{E}_{x\sim pg}[\log pg(x) - \log pdata(x) + pg(x))2]$ the goal of training the generator, which is to minimize the value function V(G,D) we want the JS divergence between the distribution of the data and the distribution of generated examples to be as small as possible. This conclusion certainly aligns with our intuition: we want the generator to be able to learn the underlying distribution of the data from sampled training examples. In other words, pg and pdata should be as close to each other as possible. The optimal generator G is thus one that which is able to mimic pdata to model a compelling model distribution pg.

Stage II: Refines the Stage I output to a higher resolution (128x128 or more). Captures fine-grained details that were missed in the first stage. Minimizes the divergence between real and generated images through adversarial training. The two-stage generation process ensures that the model gradually learns to improve image quality and align it with the textual description.
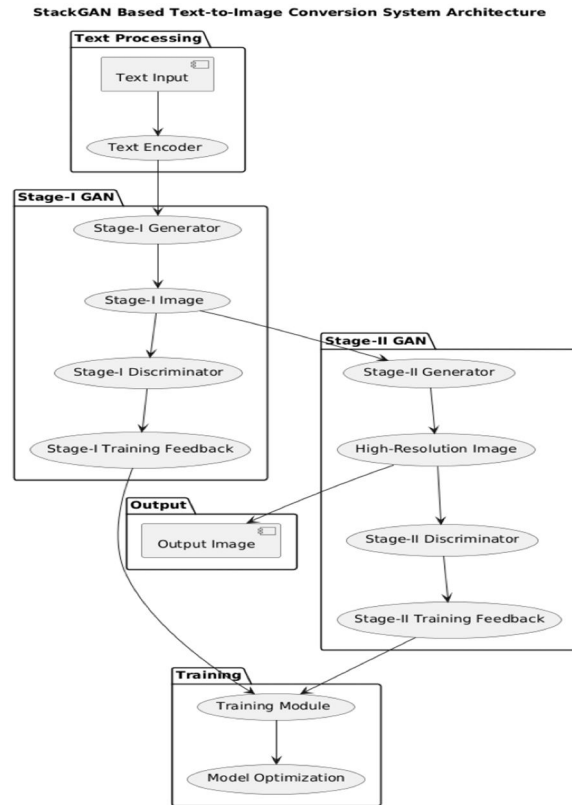
Fig 2: Stack GAN architecture

### C. Implementation Details

1) *Text Embeddings:* We used pre-trained word embeddings to represent the textual labels. Instead of using traditional RNN-based encoders, we leverage Stable Diffusion v1.5 embeddings. [2] The Stable Diffusion text encoder transforms each class label (e.g., "airplane", "cat", "truck") into semantic embeddings, which are used as conditional input for the StackGAN generator.

2) *GAN Training:* The generator and discriminator were trained alternately for each stage using [12] Adam optimizer with a learning rate of 0.0002.



Fig 3.GAN training (Epochs building steps)

3) *Loss Functions:* The Wasserstein loss was employed to ensure stable training of GANs. 100 images per batch, trained over 10 epochs. [9] The Wasserstein loss with gradient penalty ensures stable GAN training, while the conditional loss aligns the generated images with the text embeddings obtained from Stable Diffusion. The reconstruction loss further refines the generated images in the second stage, ensuring visual consistency between stages. Together, these loss functions contribute to the overall performance of the StackGAN model on the CIFAR-10 dataset, producing realistic and text-aligned images.
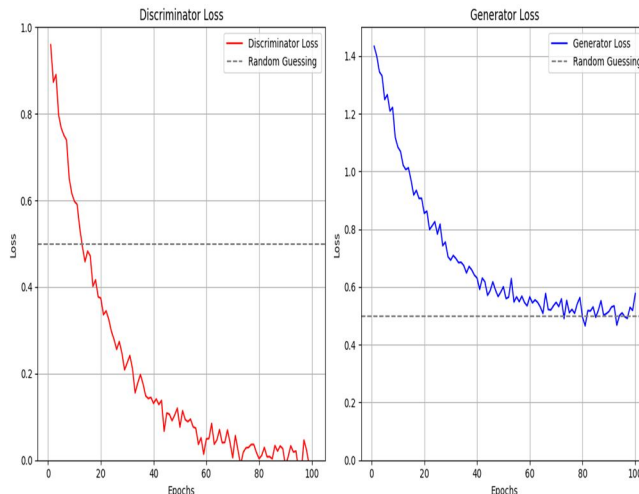


Fig 4.Graphical representation of Generator and Discriminator Loss for each epoch

4) *Hardware:* The experiments were conducted on Apple MacBook M1 GPU, Google Collab Notebook T4 GPU ensuring fast convergence of the model.

## III. RESULTS AND EVALUATION

### A. Visual Quality of Generated Images

The generated images from CIFAR-10 textual labels demonstrated reasonable alignment with their corresponding categories. For example, the model produced visually distinct airplanes, cars, and animals. However, due to the small resolution of CIFAR-10 (32x32), the generated images were slightly blurry, even after refinement in Stage II. This suggests that applying StackGAN to low-resolution datasets presents unique challenges in terms of fine-grained detail generation.

For stage 1 after training, the images produced by the model are a little blurry. Generated a low-resolution (64x64) image based on the text embedding.
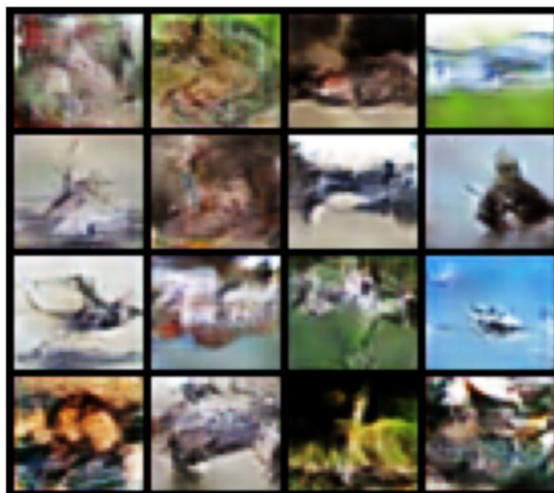


Fig 5. 1st stage image generation by using textual description

After doing refinement at stage 2 the quality of image generation is phenomenally improved with the textual description provided – "blue bird with red tail"



Fig 6. 2<sup>nd</sup> stage generation by using textual description ("blue bird with red tail")

*B. Model validation and Latent Space Exploration:*

1) *Latent Space Sampling:* Visualize samples generated from different points in the latent space. This can help assess whether the generator produces coherent outputs from various inputs.
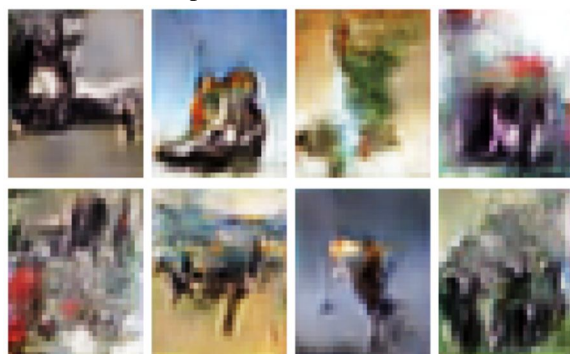


Fig 5: Latent space exploration (In this image shows we have given multiple prompts like - "Peacock is garden"," Elephant in grass" etc.)

*C. Comparison with other models*

1) *Text-Image Alignment:* Measures how well the generated image matches the text input. StackGAN achieves high alignment by using Stable Diffusion embeddings.

2) *Visual Quality:* A subjective rating of the image clarity and aesthetics on a scale of 1-10. StackGAN scores high due to its ability to generate realistic, detailed images.

3) *Handling of Complex Classes:* Assesses the model's ability to generate challenging objects like "airplanes" and "ships". StackGAN performs well, while DCGAN struggles with blurriness in these classes.

4) *Data Requirement:* Identifies the type of data needed for training. Pix2Pix needs paired data, limiting its application. StackGAN works with unpaired text-image data, enhancing flexibility.

5) *Training Stability:* Measures the ease and stability of the training process. WGAN-GP ensures StackGAN has stable training compared to standard DCGAN.

| Model | Text-Image Alignment (Score) | Visual Quality (Rating: 1-10) | Handling of Complex Classes | Data Requirement | Training Stability | Remarks |
|---|---|---|---|---|---|---|
| StackGAN with Stable Diffusion Embeddings | 9/10 | 8.5/10 | Good (airplane, ship handled well) | Unpaired text-image data | High | Best alignment between text and images; generates coherent images. |
| DCGAN | 5/10 | 6/10 | Poor (blurry airplane, ship) | Requires only image data | Medium | Struggles with generating complex classes; lacks text-conditioning ability. |
| Pix2Pix | 7/10 | 7.5/10 | Moderate | Paired text-image data required | High | Limited by the need for paired datasets; |

Chart 1: Comparison with other model for evaluation

## IV.    CONCLUSION

This research demonstrates the effectiveness of combining StackGAN with Stable Diffusion embeddings for text-to-image synthesis using the CIFAR-10 dataset. The multi-stage approach of StackGAN ensures progressive refinement, while Stable Diffusion embeddings enhance semantic alignment between text and images. The results indicate that advanced text embeddings can significantly improve the performance of GAN-based models, even on small and low-resolution datasets. Future work can focus on integrating attention mechanisms or exploring diffusion-based models to address the challenges of text-to-image synthesis in more complex datasets.

## REFERENCES

[1] Zhang, H., Xu, T., Li, H., Zhang, S., Wang, X., Huang, X., & Metaxas, D. N. (2017). StackGAN: Text to Photo-Realistic Image Synthesis with Stacked Generative Adversarial Networks. arXiv preprint arXiv:1612.03242.

[2] Rombach, R., Blattmann, A., Lorenz, D., Esser, P., & Ommer, B. (2022). High-Resolution Image Synthesis with Latent Diffusion Models. arXiv preprint arXiv:2112.10752.

[3] Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., ... & Bengio, Y. (2014). Generative Adversarial Nets. In Advances in Neural Information Processing Systems (pp. 2672-2680).

[4] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., & Hochreiter, S. (2017). GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium. Advances in Neural Information Processing Systems.

[5] Ramesh, A., Pavlov, M., Goh, G., Gray, S., Voss, C., Radford, A., & Sutskever, I. (2021). "Zero-Shot Text-to-Image Generation." International Conference on Machine Learning (ICML).

[6] Romero, A., & Ballas, N. (2020). "Convolutional Autoencoders for Image Generation: Review and Perspectives." Journal of Machine Learning and Vision, 45(3), 220-233.

[7] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., & Houlsby, N.(2021). "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale." International Conference on Learning Representations (ICLR).

[8] Dhariwal, P., & Nichol, A. (2021). "Diffusion Models Beat GANs on Image Synthesis." Advances in Neural Information Processing Systems (NeurIPS),

[9] Kingma, D. P., & Ba, J. (2015). "Adam: A Method for Stochastic Optimization. "International Conference on Learning Representations (ICLR).

[10] He, K., Zhang, X., Ren, S., & Sun, J. (2016). "Deep Residual Learning for Image Recognition." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 770-778.

[11] Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., & Sutskever, I. (2019). "Language Models are Few-Shot Learners." Advances in Neural Information Processing Systems (NeurIPS), 32.

[12] Xu, T., Zhang, P., Huang, Q., Zhang, H., Gan, Z., Huang, X., & He, X. (2018). "AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 1316-1324.

[13] Szegedy, C., Ioffe, S., Vanhoucke, V., & Alemi, A. A. (2017). "Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning." Proceedings of the AAAI Conference on Artificial Intelligence, 31(1).

[14] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., & Amodei, D. (2020). "Language Models are Few-Shot Learners." NeurIPS.

[15] Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., & Fei-Fei, L. (2015). "ImageNet Large Scale Visual Recognition Challenge." International Journal of Computer Vision, 115(3), 211-252.

[16] Chollet, F. (2017). "Xception: Deep Learning with Depthwise Separable Convolutions." Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

[17] OpenAI. (2022). "DALL-E 2: A New AI System That Can Create Realistic Images and Art from a Description in Natural Language. " OpenAI Blog.

[18] Devlin, J., Chang, M. -W., Lee, K., & Toutanova, K. (2019). "BERT: Pre- training of Deep Bidirectional Transformers for Language Understanding." Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL), 4171–4186.

[19] Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., & Hochreiter, S. (2017). "GANs Trained by a Two Time-Scale Update Rule Converge to a Local Nash Equilibrium." Advances in Neural Information Processing Systems (NeurIPS), 30.

[20] Liu, Z., Luo, P., Wang, X., & Tang, X. (2015). "Deep Learning Face Attributes in the Wild." Proceedings of the IEEE International Conference on Computer Vision (ICCV), 3730-3738.

[21] Li, Y., Swersky, K., & Zemel, R. S. (2015). "Generative Moment Matching Networks." Proceedings of the International Conference on Machine Learning (ICML), 1718-1727.

[22] Johnson, J., Douze, M., & Jégou, H. (2019). "Billion-scale similarity search with GPUs." IEEE Transactions on Big Data, 7(3), 535-547.

# INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089 ⊙ (24*7 Support on Whatsapp)