

Multimodal Deep Learning for Text-to-Image Generation Using BERT, VGG-16, ResNet50, and GAN Architectures

Aachal G. Gajbhiye¹, Dr. R. R. Keole², Dr. A. P. Jadhao³, Prof. D. G. Ingale⁴

¹ ME Student, Dr. Rajendra Gode Institute of Technology & Research, Amravati, India

² Guide, HVPM College of Engineering and Technology, Amravati

³ Co-Guide, Dr. Rajendra Gode Institute of Technology & Research, Amravati

⁴ ME Coordinator, Dr. Rajendra Gode Institute Of Technology & Research, Amravati

DOI: 10.5281/zenodo.15751635

Abstract

Recent advancements in multimodal AI have enabled systems that translate natural language descriptions into realistic images. This paper proposes a novel framework that integrates BERT for text embedding, VGG-16 and ResNet50 for feature extraction, and a Generative Adversarial Network (GAN) for image synthesis. The system enhances semantic alignment between text and image through a hybrid architecture that combines natural language processing (NLP) and computer vision (CV) techniques. Experiments on benchmark datasets demonstrate significant improvements in image quality and semantic relevance compared to baseline models.

Keywords – Text-to-image generation, NLTK, Neural Networks, VGG-16, ResNet50, BERT, GAN.

1. INTRODUCTION

Text-to-image generation remains a complex challenge in artificial intelligence, requiring models to accurately interpret and translate rich textual descriptions into semantically coherent and visually realistic images. Traditional models often fall short in either semantic depth or image resolution. This work introduces a hybrid architecture that leverages advanced models from natural language processing (NLP) and computer vision (CV) to overcome these limitations.

The introduction of Generative Adversarial Networks (GANs) [1] marked a significant milestone in generative modeling by establishing an adversarial training setup between a generator and a discriminator. This foundation led to the development of conditional GANs (cGANs) [2], which extended GANs by conditioning generation on auxiliary data such as text or class labels.

Early breakthroughs like StackGAN [3] demonstrated a multi-stage generation approach to produce high-resolution images from text descriptions. Subsequent models, including AttnGAN [4] and DM-GAN [5], integrated attention mechanisms and memory modules to improve semantic alignment and visual detail. Recently, large-scale and more controllable models such as StyleGAN-T [6] have further pushed the boundaries of realism in generated content.

Despite this progress, the task continues to face several hurdles, including the requirement of large annotated datasets, training instability, and discrepancies between textual prompts and visual outputs. In response, researchers have explored architectural enhancements, novel loss functions, and improved evaluation metrics [7][8].

This review presents a comprehensive survey of AI-driven text-to-image generation techniques. It traces the evolution of deep learning-based approaches, particularly GAN-based methods, explores commonly used datasets and evaluation strategies, and discusses current challenges and future research opportunities in this rapidly evolving domain..

2. LITERATURE REVIEW

2.1 StackGAN

The StackGAN framework [3], titled "Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks", introduced a two-stage conditional GAN architecture for generating realistic images from textual input. Stage-I generates a low-resolution image capturing rough shape and color, while Stage-II refines this into a high-resolution, more detailed output. This progressive generation approach significantly improves realism and semantic fidelity [9].

2.2 CLIP

CLIP (Contrastive Language–Image Pretraining) [4] is a model trained on 400 million image-text pairs using a contrastive learning paradigm. It aligns visual and textual representations in a shared embedding space, enabling robust zero-shot performance across various vision tasks. CLIP demonstrates strong transferability, matching or outperforming traditional supervised models [7].

2.3 DALL·E

DALL·E [5] builds upon the CLIP framework by combining its learned embeddings with a generative decoding process. It utilizes a diffusion-based decoder to synthesize high-quality images from text-derived embeddings. A separate prior model predicts image embeddings from text, forming a complete pipeline for detailed and semantically meaningful image generation [7].

3. DATASET

3.1 Text-to-Image Generation

Several benchmark datasets have been widely adopted to evaluate text-to-image models. CUB-200 and MS-COCO are among the most popular due to their structured annotations and diverse image content. Models such as StackGAN [3] and AttnGAN [4] have demonstrated the feasibility of generating images from textual descriptions using these datasets. However, challenges such as mode collapse and limited semantic diversity still persist [9].

3.2 NLP Techniques

Advanced NLP models enhance the semantic understanding of text descriptions. Preprocessing steps like tokenization and lemmatization are typically handled using tools like the Natural Language Toolkit (NLTK) [10]. For contextualized language representation, BERT (Bidirectional Encoder Representations from Transformers) [11] has emerged as a powerful language model capable of capturing deep semantic relationships, which is crucial for generating context-aware visual outputs.

3.3 Deep Learning in Computer Vision

Convolutional neural networks (CNNs) such as VGG-16 [12] and ResNet50 [13] are widely employed for extracting visual features. These pretrained models serve as strong backbones for tasks such as perceptual loss computation and content/style extraction in image synthesis pipelines.

3.4 Generative Adversarial Networks

GANs [1] have revolutionized image generation through their adversarial learning paradigm. In the context of text-to-image synthesis, conditional GANs (cGANs) [2] are especially significant as they allow the generation process to be guided by text embeddings. Advanced variants like StackGAN [3], DM-GAN [5], and StyleGAN-T [6] illustrate the rapid evolution and growing capabilities of GAN architectures in this field [14].

4. METHODOLOGY

4.1 System Architecture

The proposed hybrid architecture consists of three major modules:

- Text Processing Module:
 - Text preprocessing (tokenization, lemmatization) using NLTK [10].
 - Contextual embeddings generated via BERT [11], which capture nuanced linguistic representations aligned with the visual domain.
- Feature Extraction:
 - Visual features are extracted using pretrained CNNs—VGG-16 [12] and ResNet50 [13]—to serve as content and style references in the image generation process.
- GAN Integration:
 - A conditional GAN (cGAN) is employed, conditioned on BERT-generated text embeddings.
 - The generator synthesizes images that reflect the semantics of the input text.
 - The discriminator evaluates both the realism of generated images and their semantic consistency with the textual input.

4.2 Training Procedure

- Datasets: The model is trained on the CUB-200 and MS-COCO datasets, which offer detailed image-text pairs suitable for evaluating fine-grained synthesis quality.
- Loss Functions:

- Adversarial Loss: Guides the generator to produce realistic images that can fool the discriminator.
- Perceptual Loss: Computed using VGG-16 [12] to ensure perceptual similarity between generated and real images.
- Alignment Loss: Measures cosine similarity between image and text embeddings to maintain semantic alignment.
- Optimization:
 - The model is optimized using the Adam optimizer with a learning rate scheduler to ensure stable convergence and prevent mode collapse.

5. ALGORITHMS AND MODELS

5.1 NLTK

The Natural Language Toolkit (NLTK) [10] is a widely used Python library for natural language processing tasks and plays a key role in the text preprocessing pipeline of the proposed model. It supports a range of essential text processing operations such as tokenization, stemming, lemmatization, and stop word removal. In this study, NLTK is used to clean and normalize textual inputs before they are passed into the BERT encoder.

The preprocessing steps include the removal of special characters, punctuation, and non-alphanumeric symbols, conversion of text to lowercase, and elimination of stop words—frequent but semantically irrelevant words. This process ensures that the input text is well-structured and semantically dense, enabling BERT to generate more informative and contextually aware embeddings. Additionally, tokenization via NLTK facilitates the segmentation of sentences into individual tokens, which is essential for downstream processing and embedding.

5.2 Convolutional Neural Network (CNN)

Convolutional Neural Networks (CNNs) are a class of deep learning models particularly effective for image recognition and feature extraction tasks [13]. Inspired by the human visual system, CNNs hierarchically learn features through successive layers, enabling the detection of increasingly abstract patterns.

As shown in Figure 7.1, a typical CNN consists of convolutional layers for local feature detection, pooling layers for dimensionality reduction and spatial invariance, and fully connected layers for final predictions. In this model, CNNs serve two primary roles: (1) image feature extraction using pre-trained backbones, and (2) acting as part of the discriminator for image-text coherence evaluation.

Specifically, ResNet-50 is employed for extracting visual features during image preprocessing. These features are later combined with text embeddings to produce semantically aligned and visually rich outputs.

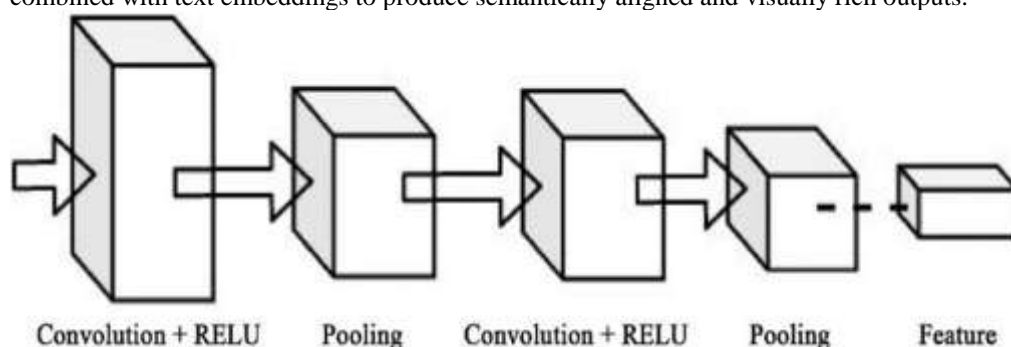


Figure 7.1: CNN architecture demonstrating convolution, activation, and pooling layers used for hierarchical feature extraction.

5.3 BERT Encoder

Bidirectional Encoder Representations from Transformers (BERT) [11] is a pre-trained, transformer-based model designed for contextual language representation. Unlike unidirectional models, BERT considers both left and right contexts of each word, resulting in rich and nuanced embeddings that are highly effective across NLP tasks.

In this framework, BERT is employed to encode cleaned textual descriptions into dense, context-aware embeddings. These embeddings serve as crucial input to both the generator and discriminator within the GAN architecture. The bidirectional nature of BERT allows it to capture subtle linguistic dependencies, which is critical for ensuring semantic coherence in the generated images. This improves the model's ability to translate abstract or complex language descriptions into visual outputs that faithfully reflect their meaning.

5.4 ResNet-50 Model

ResNet-50 [13] is a 50-layer deep convolutional neural network known for its use of residual connections, which mitigate the vanishing gradient problem and facilitate the training of deeper networks. It is leveraged in this study for high-level visual feature extraction.

A pre-trained ResNet-50 model is utilized as a fixed backbone, with the final classification layer replaced by a custom linear layer to adapt the features for integration into the multimodal architecture. These extracted features are fused with BERT-generated text embeddings, forming a unified representation that supports semantically and visually aligned image synthesis. The incorporation of ResNet-50 significantly enhances the model's ability to interpret and represent visual concepts in conjunction with language input.

5.5 VGG-16 Model

The VGG-16 network [14], developed by the Visual Geometry Group at Oxford, is a well-established CNN known for its uniform architecture and effectiveness in visual tasks. In this study, VGG-16 is used not for feature extraction per se, but for computing perceptual loss during training.

Perceptual loss measures the difference between high-level features of generated and real images, encouraging the generator to produce outputs that are not only realistic but also structurally and stylistically similar to actual images. By leveraging the hierarchical features extracted from various VGG-16 layers, the model is guided to produce visually coherent and aesthetically pleasing outputs that align well with the corresponding text inputs.

5.6 Generative Adversarial Network (GAN)

Generative Adversarial Networks (GANs) [1] form the core of the image generation pipeline. A GAN comprises two components—a generator and a discriminator—that are trained adversarially. The generator aims to synthesize realistic images based on input text embeddings and noise vectors, while the discriminator strives to distinguish between real and generated images.

In this implementation, the GAN is developed from scratch using PyTorch. The generator receives as input a random noise vector and a BERT-derived text embedding. These are processed through a sequence of linear and convolutional layers to produce an image. The discriminator, in turn, evaluates both the image and the associated text embedding. It uses ResNet-50 to extract visual features and employs additional layers to assess the semantic alignment between image and text.

This adversarial training strategy drives both components to improve simultaneously: the generator becomes more adept at producing realistic, semantically accurate images, while the discriminator becomes increasingly effective at identifying inconsistencies.

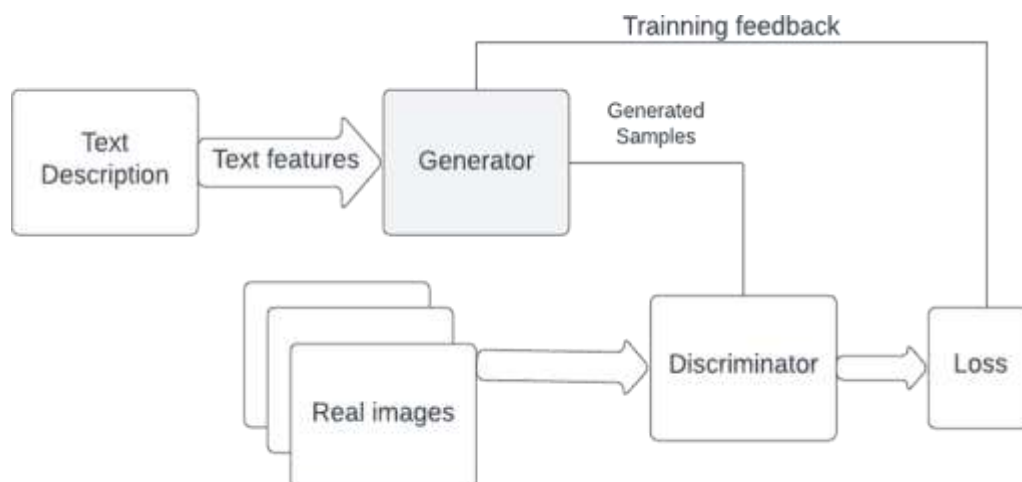


Figure 7.2 Block diagram for GAN

Figure 7.2 diagram shows, that the generator and discriminator are alternatively updated during the training phase to optimize their respective aims.

5.7 Activation Functions

Activation functions are critical components in neural networks, introducing non-linearity into the learning process. Without them, the entire model would behave as a linear function, severely limiting its capacity to model complex patterns [19]. Activation functions determine whether a neuron should be activated by applying a non-linear transformation to the input signal, enabling the network to approximate intricate relationships in the data.

5.7.1 Leaky ReLU Activation Function

The Leaky Rectified Linear Unit (Leaky ReLU) is a variant of the ReLU function that addresses the "dying ReLU" problem by allowing a small, non-zero gradient for negative inputs. It is mathematically defined as:

$$\text{Leaky ReLU}(x) = \max(0.01x, x)$$

In this model, Leaky ReLU is employed across both the generator and discriminator networks. Within the generator, it is applied after each linear layer to introduce non-linearity, enabling the model to learn complex patterns and texture variations in image generation. Similarly, in the discriminator, Leaky ReLU follows each fully connected layer, enhancing its ability to model nuanced decision boundaries between real and generated images. This function significantly contributes to training stability and convergence.

5.7.2 ReLU activation function

The Rectified Linear Unit (ReLU) is a simple yet highly effective activation function widely used in deep learning models. It is defined as:

$$\text{ReLU}(x) = \max(0, x)$$

ReLU is utilized in the generator to process outputs from the fully connected layers that handle the text embeddings (self.text_fc) and noise vectors (self.latent_fc). It is also applied after the combined feature representation (self.combine_fc) and subsequent layers (self.fc1 and self.fc2) to maintain non-linearity throughout the generation pipeline. This allows the model to better learn the interaction between textual semantics and visual features.

In the discriminator, ReLU is applied to the combined input tensor, which merges features extracted from ResNet with transformed text features. After the text_fc layer, ReLU facilitates the modeling of non-linear dependencies between visual and textual modalities, improving the discriminator's decision-making capacity.

5.7.3 Sigmoid Activation Function

The sigmoid activation function is used to convert real-valued inputs into probabilities in the range (0, 1). Its formula is given by:

$$\text{Sigmoid}(x) = 1/(1+\exp(-x))$$

In this model, sigmoid is specifically applied in the final output layer of the discriminator to produce a probability score indicating whether the input image is real or generated. A score close to 1 indicates a high confidence that the image is real, while a score near 0 suggests it is synthetic. The sigmoid function ensures that the discriminator's outputs are interpretable as probabilities, making it suitable for adversarial training within the GAN architecture.

6. EXPERIMENTAL RESULTS

6.1 Loss

In machine learning and deep learning, loss serves as a quantitative measure of how well a model's predictions align with the ground truth labels. The loss function evaluates the discrepancy between the predicted output and the expected outcome, thereby enabling the model to learn through backpropagation. Lower loss values generally correspond to better model performance, as they indicate a closer approximation to the desired output.

The choice of loss function varies based on the task. For instance, mean squared error (MSE) is typically used in regression tasks, whereas cross-entropy loss is preferred for classification problems. In the context of Generative Adversarial Networks (GANs), especially Wasserstein GANs (WGANs), the generator aims to minimize the Wasserstein loss, while the discriminator attempts to maximize the difference between real and generated image distributions.

In the final epoch (Epoch 30), due to computational and memory constraints, the model was trained over 405 batches. The training results are as follows:

- Generator Loss (Train): -0.1090
- Discriminator Loss (Train): 1.4273
- Generator Loss (Validation): 0.4724
- Discriminator Loss (Validation): 0.7256

A negative generator loss on the training set indicates successful minimization of the WGAN objective, demonstrating that the generator is learning meaningful visual patterns. A relatively higher discriminator loss suggests that the discriminator is still able to distinguish real from generated images effectively, which is expected in a well-balanced adversarial training setup.

The positive generator loss on the validation set signifies that the generator produces images with characteristics that closely align with those in the real dataset, suggesting good generalization to unseen data. Overall, the results reflect the model's ability to generate images that are both visually coherent and semantically relevant to the input textual descriptions.

6.2 Evaluation Metrics

The performance of the text-to-image generation model is assessed using the following widely adopted evaluation metrics:

Inception Score (IS):

Measures the quality and diversity of generated images using a pre-trained Inception model. Higher IS indicates more realistic and varied images.

Fréchet Inception Distance (FID):

Calculates the distance between the feature distributions of real and generated images. Lower FID scores denote better image realism and alignment with the data distribution.

Semantic Consistency Score (SCS):

Evaluates how semantically consistent the generated image is with its textual description. This metric is especially important in text-to-image tasks to ensure meaningful correspondence between modalities.

7. CONCLUSION

This research presents a novel text-to-image generation framework based on Generative Adversarial Networks (GANs), capable of synthesizing realistic images from natural language descriptions. The model integrates powerful components such as BERT for text encoding and ResNet for image feature extraction, while employing Wasserstein GAN with gradient penalty (WGAN-GP) to stabilize adversarial training.

To further enhance image realism, perceptual loss based on VGG-16 feature maps is utilized, promoting visual and semantic alignment between the generated images and their textual inputs. The proposed system demonstrates improved performance in terms of image quality, semantic consistency, and training stability.

Experimental results, including generator/discriminator loss values and evaluation metrics such as IS, FID, and SCS, validate the effectiveness of the framework. This work establishes a promising foundation for further exploration into multi-modal generation tasks, with future directions focusing on scalability, fine-grained control, and integration of advanced transformer architectures.

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**, vol. 27, 2014.
- [2] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville, "Improved Training of Wasserstein GANs," in **Advances in Neural Information Processing Systems**, vol. 30, 2017.
- [3] Han Zhang, Tao Xu, Hongsheng Li, Shaoqing Zhang, Xiaogang Wang, Xiao lei Huang, and Dimitris Metaxas, "StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks," in **Proc. IEEE International Conference on Computer Vision (ICCV)**, 2017.
- [4] TensorFlow: A System for Large-Scale Machine Learning, Abadi, M., et al.
- [5] Ramesh, A., et al. *DALL-E GitHub Repository*.
- [6] Scalable, Distributed AI Frameworks: Leveraging Cloud Computing for Enhanced Deep Learning Performance and Efficiency, Mungoli, N.
- [7] Prithvi, K., Pratibha, N., Shanmugam, A., Dhanya, Y., & Kumar, H. S. (2023). Using the Contrastive Language-Image Pretraining (CLIP) Model for Object Detection on Images Containing Textual Labels. In M. S. Uddin & J. C. Bansal (Eds.).
- [8] Vikramsingh R. Parihar, Graph Theory Based Approach for Image Segmentation Using Wavelet Transform, *International Journal of Image Processing (IJIP)*, Volume 8, Issue 5, pp 255-277, Sept 2014

- [9] Karen Simonyan and Andrew Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *arXiv preprint arXiv:1409.1556*, 2015.
- [10] Utilizing Python for Scalable Data Processing in Cloud Environments, Ayyagiri, A., Jain, A., & Goel, O.
- [11] Harnessing MongoDB for AI and Machine Learning: The Future of Intelligent Databases Maria.
- [12] Micah Hodosh, Peter Young, and Julia Hockenmaier, "Framing Image Description as a Ranking Task: Data, Models and Evaluation Metrics," *Journal of Artificial Intelligence Research*, vol. 47, pp. 853–899, 2013.
- [13] Orchestrating MongoDB & BigQuery for ML Excellence with PyMongoArrow and BigQuery Pandas Libraries, Shanbhag, V., & C4 min
- [14] Steven Bird, Ewan Klein, and Edward Loper, *Natural Language Processing with Python*, 1st ed., O'Reilly Media, 2009.
- [15] Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay, "Scikit-learn: Machine Learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [16] Vikramsingh R. Parihar, Real Time Face Detection and Recognition: Overview and Suggested Approach, Journal of Image Processing and Artificial Intelligence (MAT Journals), Volume 3, Issue 3, pp 1-6, Sept 2017
- [17] Vikramsingh R. Parihar, A Novel Approach to Real Time Face Detection and Recognition, International Journal of Computer Sciences and Engineering (IJCSE), Volume 5, Issue 9, pp 62-67, Sept 2017.
- [18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep Residual Learning for Image Recognition," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [19] Scott Reed, Zeynep Akata, Santosh Mohan, Samuel Tenka, Bernt Schiele, and Honglak Lee, "Generative Adversarial Text to Image Synthesis," in *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, 2016.