

# SYNTHESIZING REALITIES: A DEEP DIVE INTO TEXT TO IMAGE GENERATION WITH FULLY TRAINED GAN

1 Nerella Neha, 2 Dr. M Dhanalakshmi

1 University College of Engineering, Science and Technology, JNTUH Hyderabad.

2, Professor of IT, University College of Engineering, Science and Technology, JNTUH Hyderabad.

**Abstract:** In the field of text-to-face generation, where the production of practical facial pictures from printed depictions is as yet a troublesome undertaking, one significant deterrent looked by specialists has been the absence of exhaustive datasets. Our review fills this hole by introducing a special technique that changes over literary depictions straightforwardly into practical and definite face qualities utilizing a pre-prepared Generative Adversarial Network (GAN) model. Rather than conventional methodologies that main utilize pre-prepared text encoders, our framework at the same time prepares the text encoder and picture decoder. The nature of blended face pictures is improved by this co-preparing method, which makes it more straightforward to appreciate the intricate connection between verbal depictions and visual portrayals. We underscored the significance of having top notch preparing information and painstakingly chose a dataset that would work on the exhibition of our model. Our work features how GANs might be utilized to overcome any issues among text and picture synthesis, giving valuable data to true information driven PC vision areas. Our technique is to advance the information and use of AI innovation in delivering excellent visual substance by exhibiting remarkable picture blend abilities.

**Index Terms:** *Text descriptions to image, GAN, fully trained, CelebA, text encoder, image decoder*

## 1. INTRODUCTION

As an enormous information base for both logical and artificial intelligence (AI) fields, ML currently holds the chance of filling holes in the mission for "becoming intelligent". The development of visual information from composed directions is known as text-to-picture generation (TTPG), and it is one of the trickiest yet most significant positions in AI and ML. To parse information utilizing modified approaches that are in accordance with language designs and text

handling standards, this work joins natural language processing (NLP) philosophies and strategies. Calculations and approaches for computer vision then, at that point, get data from this handled information.

In light of its curiosity, text-to-picture generation is a specific field inside computer vision. It includes two essential assignments: the making of pictures from text, giving direct versus opposite methods, and pictorial subtitling and order, where photographs produce message depictions. To create photograph

reasonable pictures, Contingent Generative Adversarial Network (GAN) preparing requires undeniable level semantic vector inputs. Text encoders and image decoders are the two essential pieces of this preparing approach, as well as delivering great images.

Text-to-image (T2I) task progressions to far have generally been focused on more modest datasets, such blossoms and birds, with less example sizes [1]. The attention has been on these less complex datasets, with less exploration done on the more troublesome ones. For instance, the legitimacy of the item scene interface at the expression level has been tried utilizing the COCO dataset [2]. This idea has a lot of viable ramifications in policing, witness portrayals are utilized to assist craftsmen with painting suspect pictures. This manual system features the likely impact of computerized text-to-confront creation, as it basically relies upon human abilities to change verbal depictions into visual drawings.

One significant issue for scholastics concentrating on text-to-confront age is the absence of a uniform dataset. Not much has been finished to create faces from text based portrayals, regardless of the advances in text-to-image synthesis [3]. To close this hole, our work recommends an extraordinary technique that changes over printed depictions into sensible and normal looking face photographs by utilizing a pre-prepared GAN model. We train both the text encoder and the image decoder simultaneously, rather than different methodologies that utilize a pre-prepared text encoder. Better image synthesis execution is the result of this co-preparing process, which advances a more noteworthy perception of the mind boggling

connections between verbal portrayals and visual modalities.

Any solid AI model beginnings with excellent preparation information. We have consequently meticulously made a dataset that is particularly appropriate for text-to-face generation [4]. Our exploratory outcomes show that our strategy works better compared to current methodologies, giving outstanding additions in the accuracy and nature of created face pictures. This examination adds to the field of PC vision by developing comprehension we might interpret GANs and their applications. This will be particularly helpful in spaces where changing over text based depictions into visual portrayals is fundamental.

Our commitments are numerous and exceptionally specialized. We give another GAN-based text-to-face generation system that doesn't utilize text encoders that have proactively been prepared. Through concurrent preparation, our methodology works on the association among text based and visual info, creating reasonable, great facial portrayals [5]. To address the current shortage of normalized assets in this field, we likewise give a new dataset made particularly for this objective [6]. Our review has a large number of potential applications, from the improvement of programmed outlines in criminal examinations to the production of custom fitted material in various organizations [7].

Using state of the art strategies and a specific dataset, we need to tackle the issues related with text-to-face creation and make the way for more precise and accurate AI systems that can change over composed

portrayals into clear and energetic visual results. As well as extending the restrictions of existing examination, this study sets out new open doors for AI and computer vision applications.

## 2. LITERATURE SURVEY

Recent years have witnessed tremendous progress in the field of text-to-image creation, especially with the use of Generative Adversarial Networks (GANs), thanks to both theoretical improvements and real-world implementations. This review of the literature is to investigate significant developments and contributions made in this field, emphasizing influential works and their influence.

Generative Adversarial Networks (GANs), a framework consisting of two neural networks—the discriminator and the generator—competing against one other to create realistic samples, were first proposed in the fundamental work of Goodfellow et al. [1]. Since then, generative modeling—which includes text-to-image synthesis—has relied heavily on GANs.

Hong et al. [2] developed a technique for hierarchical text-to-image synthesis that infers semantic layout by building upon GANs. Their method improved the synthesis of complex pictures from in-depth verbal descriptions by addressing the difficulty of matching textual descriptions with hierarchical visual structures.

High-level semantic information has been crucially included into the picture generating process through the use of conditional GANs (cGANs). In order to improve the ability to produce pictures conditioned on certain textual inputs, Van den Oord et al. [3] presented conditional image synthesis utilizing PixelCNN

decoders. This method represented a major breakthrough in matching picture generating jobs with text descriptions.

Zhang et al.'s StackGAN++ [4] proposed a hierarchical architecture with stacked Generative Adversarial Networks, which significantly enhanced conditional picture synthesis. Their approach established the efficacy of hierarchical conditioning in text-to-image synthesis tasks by producing high-resolution pictures with better visual quality and realism.

The Microsoft COCO dataset [7], which offers a wide variety of object-centric pictures with detailed annotations, has been used as a benchmark for assessing text-to-image creation methods. This dataset has proved useful in comparing different methods and assessing how well they generalize to diverse visual domains.

To improve the variety and caliber of created pictures, text-to-image synthesis frameworks have also used Auto-Encoding Variational Bayes (VAEs). When used in combination with GANs, VAEs provide more versatile and adaptable picture synthesis capabilities. Kingma and Welling [8] developed VAEs as a generative model that learns a latent variable representation of data.

Through Generative Adversarial Text-to-Image Synthesis (GATIS), Reed et al. [9] investigated the use of GANs for text-to-image synthesis. The creation of realistic visuals from detailed text inputs was made possible by their work, which focused on the direct integration of textual descriptions into the GAN framework. This method showed encouraging

outcomes in producing pictures that faithfully capture the information found in written descriptions.

The development of deep Convolutional GANs (DCGANs) for unsupervised representation learning by Radford et al. [10] set the stage for reliable training and high-fidelity picture production. Their research centered on using adversarial training to learn hierarchical representations of pictures, improving the caliber and variety of synthetic images.

Together, the aforementioned publications demonstrate the development of text-to-image synthesis methods, highlighting advances in model structures, dataset applications, and assessment strategies. These advancements have broadened the applications of text-to-image synthesis in several fields, such as computer-aided design, visual storytelling, and creative content production, in addition to improving the quality and realism of created pictures.

To sum up, the domain of text-to-image production is seeing rapid evolution due to advancements in neural network topologies, training techniques, and the accessibility of datasets. Recent developments in conditional GANs, hierarchical modeling, and latent variable representations have greatly increased the state-of-the-art in producing realistic pictures from textual descriptions, even though issues like dataset variety and semantic alignment still exist. Subsequent investigations might concentrate on resolving these issues, investigating innovative structures, and broadening the use of text-to-image synthesis in practical situations.

### 3. METHODOLOGY

#### a) Proposed Work:

In the proposed system, textual descriptions are encoded into vector representations using a Bidirectional Long Short-Term Memory (BI-LSTM) network. These vector representations are then fed into a Generative Adversarial Network (GAN) to generate images. This method trains the text encoder and picture decoder concurrently, with the goal of improving the realism and accuracy of produced visuals. The CelebA dataset, which consists of a sizable collection of celebrity photos accompanied by written descriptions, is used to train the GAN. The system aims to generate realistic and visually correct facial pictures based on textual descriptions by utilizing the BI-LSTM for efficient text encoding and the GAN for high-quality image synthesis. This approach offers advantages over conventional approaches in the text-to-face creation domain.

#### b) System Architecture:

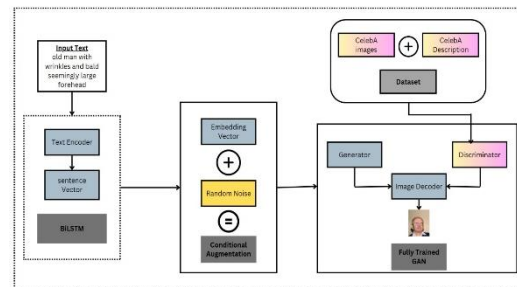


Fig 1 Proposed Architecture

The suggested text-to-face generation's system architecture is made up of many essential parts. First, a Text Encoder employs a Bidirectional Long Short-Term Memory (BiLSTM) network to analyze incoming text descriptions and encodes them into

sentence vectors. Conditional augmentations are produced by combining these text vectors with random noise. These enhancements function as inputs to the GAN framework's Generator. The picture Decoder receives the synthesized picture data from the Generator.

The CelebA dataset, which has celebrity photos with accompanying written descriptions, is used in the training process. The Discriminator, which differentiates between actual and created pictures and feeds that information back into the Image Decoder, is trained using these images and descriptions. This procedure thoroughly trains the GAN to produce accurate and lifelike facial pictures from textual descriptions. Coherent and superior output is ensured by training the text encoder and picture decoder simultaneously.

### c) Dataset Collection:

In this research, the foundation for training our Generative Adversarial Network (GAN) model is the CelebA (Celebrity Faces Attributes) dataset. It has an extensive library of famous photos, many of which include textual explanations and in-depth attribute annotations. CelebA offers a wide range of face features and expressions from over 200,000 photos of around 10,000 celebrities, which makes it perfect for training models that produce realistic facial pictures from textual descriptions.

We load and preprocess the photos together with the written descriptions to get the CelebA dataset ready for training. This preparation phase makes sure the dataset is aligned and structured correctly so that the GAN can be trained, matching each textual description with an

image of a face. Our model seeks to improve its capacity to generate accurate and lifelike facial pictures based on textual inputs by making good use of the extensive annotations and diverse facial features included in CelebA. Additionally, we want to achieve good generalization to faces from other datasets.

```

1 ("response_id": "23211", "filename": "000029.jpg", "user_id": 8, "description": "A woman with high cheekbones, a defined jawline and
2 ("response_id": "1299", "filename": "000035.jpg", "user_id": 2, "description": "A woman with a chiseled jaw, prominent cheekbones,
3 ("response_id": "20953", "filename": "000035.jpg", "user_id": 7, "description": "A woman with long hair, dark eyes, a small nose and t
4 ("response_id": "25903", "filename": "000063.jpg", "user_id": 10, "description": "A fair woman with shiny golden locks, an oblong fac
5 ("response_id": "13528", "filename": "000082.jpg", "user_id": 5, "description": "A middle-aged man with short brown hair and small b
6 ("response_id": "19871", "filename": "000082.jpg", "user_id": 7, "description": "A man with short hair, blue eyes, a small nose and thin
7 ("response_id": "55824", "filename": "000099.jpg", "user_id": 15, "description": "A woman with wavy, light brown hair and dark brown
8 ("response_id": "33988", "filename": "000115.jpg", "user_id": 11, "description": "this man does not have hair, he has thin eyebrows an
9 ("response_id": "35805", "filename": "000120.jpg", "user_id": 13, "description": "Profile of a balding white man, having a wide smile i
10 ("response_id": "81926", "filename": "000120.jpg", "user_id": 12, "description": "A bald man with thin eyebrows, brown eyes, a big an
11 ("response_id": "2314", "filename": "000122.jpg", "user_id": 2, "description": "An older woman with short, blonde hair and full lips. S
12 ("response_id": "12858", "filename": "000122.jpg", "user_id": 5, "description": "An middle-aged woman with short blonde hair. Her ey
13 ("response_id": "15011", "filename": "000122.jpg", "user_id": 6, "description": "An old pale white woman with short blonde hair bro
14 ("response_id": "34408", "filename": "000144.jpg", "user_id": 11, "description": "this man has long grey hair, he is wearing black glas
15 ("response_id": "33983", "filename": "000170.jpg", "user_id": 11, "description": "this woman has ginger coloured long hair, she also h
16 ("response_id": "18879", "filename": "000213.jpg", "user_id": 5, "description": "A man with medium-length grey hair and a receding ha
17 ("response_id": "19915", "filename": "000213.jpg", "user_id": 7, "description": "An elderly man with short hair, dark eyes, a pointed ac
18 ("response_id": "26207", "filename": "000217.jpg", "user_id": 10, "description": "A woman with a fair complexion, pulled-back blond
19 ("response_id": "84518", "filename": "000217.jpg", "user_id": 11, "description": "this woman has blonde hair, she has thin eyebrows an
20 ("response_id": "33854", "filename": "000237.jpg", "user_id": 11, "description": "this is an older man with very thin short brown hair l
21 ("response_id": "24165", "filename": "000266.jpg", "user_id": 2, "description": "A lanky, bald black man with a cleft chin. He is wearin
22 ("response_id": "13412", "filename": "000266.jpg", "user_id": 5, "description": "A bald middle-aged chubby dark-skinned man wearing
23 ("response_id": "16619", "filename": "000266.jpg", "user_id": 6, "description": "A serious looking bald black man with a wide nose, a
24 ("response_id": "26435", "filename": "000279.jpg", "user_id": 10, "description": "A woman with a pale complexion, long brown hair v
25 ("response_id": "25491", "filename": "000325.jpg", "user_id": 10, "description": "A fair woman with alowning skin, green eyes, eye-lin

```

Fig 2 Text Dataset Collection



Fig 3 CelebA Dataset Images

### d) Data Processing:

Text pretreatment entails many procedures that guarantee the input text is clean and useable by the GAN model. First, text cleaning is used to eliminate noise and extraneous letters from the input text. Tokenization then turns the cleaned text into tokens, for example, "A smiling young woman" becomes ["A", "smiling", "young", "woman"]. The vocabulary size (V) denotes the number of unique tokens in the dataset; for example, V=10,000 words. Following tokenization, vectorization is performed by

transforming tokens into numerical vectors using embeddings. Each word is represented as a 100-dimensional vector, and the embedding matrix (E) has a size of  $V \times D$ , where D is the embedding dimension.

Image preparation is a series of critical processes used to prepare pictures for GAN model training. First, normalization adjusts pixel values to a normalized range, which is commonly [0, 1]. Images are resized to a set dimension, such as  $128 \times 128$  pixels, based on their original dimensions ( $H \times W \times C$ ). Mathematically, if the original pixel value (p) falls between 0 and 255, the normalized pixel value (p') is determined as  $p' = p/255$ . During data collection and preprocessing, photos are downsized to  $128 \times 128$  to balance detail and computational efficiency. To decrease computational burden, pictures can be reduced to  $32 \times 32$  pixels before feeding into the GAN. This allows for quicker training times and lower memory utilization, especially with big datasets and sophisticated models.

#### e) Text Encoder Training:

Using a Bidirectional Long Short-Term Memory (BI-LSTM) model, text encoder training aims to efficiently extract semantic information from text descriptions. This methodology creates a fixed-size vector representation of the text by processing a series of word vectors. The BI-LSTM creates rich and significant phrase vectors by utilizing both forward and backward LSTM layers to capture contextual relationships within the text. The core of the textual descriptions is captured by these vectors, which improve the capacity of downstream models, such as Generative Adversarial Networks (GANs), to produce

visually cohesive and contextually correct pictures from the encoded verbal input.

#### f) Embedding Vector and Conditional Augmentation:

Combining random noise with the resulting sentence vector—which extracts semantic information from textual descriptions using techniques like TF-IDF (Term Frequency-Inverse Document Frequency)—is the basis for both embedding vector and conditional augmentation. Even for text inputs that are similar, this augmentation procedure guarantees variety in the pictures produced by the GAN. As a result, the embedding vector changes to conform to the shape that the Generator in the GAN framework needs as input. This method strengthens and more accurately represents the input textual descriptions by adding diversity and realism to the synthesized visuals.

#### g) GAN Training:

The Generator and the Discriminator are the two primary components used in GAN training. Using encoded word vectors as input, the Generator network learns to create realistic-looking synthetic facial pictures. With the intention of tricking the Discriminator, it creates graphics from the supplied text vector. In turn, the Generator trains the Discriminator to discern between actual images from the dataset and images produced by the Discriminator. Through this repeated adversarial training process, both networks get better: the Discriminator gets better at recognizing created pictures, while the Generator learns to create more realistic images. The GAN converges to produce high-quality facial pictures that





Fig 8 Result of Output from text to Image Generation



Fig 9 Result of Output from text to Image generation

### Discriminator Accuracy

- Real Images Accuracy:

$$Accuracy_{real} = \frac{97}{100} = 0.97 \text{ or } 97\%$$

- Fake Images Accuracy:

$$Accuracy_{fake} = \frac{0}{100} = 0 \text{ or } 0\%$$

Fig 10 Discriminator Accuracy

### Generator loss:

- This loss calculation is a measure of the divergence of the distributions of the generated images from the real image distributions as measured by the discriminator.

**Loss = Binary Cross-Entropy Loss between Discriminator's Predictions on Fake Images and Real Labels.**

- Binary Cross-Entropy Loss:

$$Loss = -\frac{1}{N} \sum_{i=1}^N [y_i \log(p(y_i)) + (1 - y_i) \log(1 - p(y_i))]$$

- Where  $y_i$  is the label (1 for real, 0 for fake) and  $p(y_i)$  is the predicted probability.

**Generator loss=0.693**

### Similarity Score:

The similarity score is calculated to survey how near the produced pictures are to genuine pictures and is calculated employing a closeness score (cosine similitude). It is performed during the text-to-image generation handle.

It is especially valuable in high-dimensional positive spaces where the magnitude of the vectors isn't important but the direction is. This can be why it is broadly utilized in content mining and data recovery. A score of 0.45 proposes a moderate.

- Dot Product (A·B):** The dot product of two vectors A and B is calculated below:
- The vectors A and B components are  $A_i$  and  $B_i$ .
- A ranges from  $A_1$  to  $A_n$ , B ranges from  $B_1$  to  $B_n$
- Their dot product is:



$$\mathbf{A} \cdot \mathbf{B} = A_1 B_1 + \dots + A_n B_n$$

$$\text{Similarity} = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|}$$

$$\text{Similarity} = \frac{A_1 B_1 + A_2 B_2 + \dots + A_n B_n}{\sqrt{A_1^2 + A_2^2 + \dots + A_n^2} \cdot \sqrt{B_1^2 + B_2^2 + \dots + B_n^2}}$$

## 5. CONCLUSION

We propose here a completely prepared generative adversarial network (GAN) equipped for creating text-to-image face synthesis. To produce top notch pictures comparing to the info words, we carried out an organization that simultaneously prepared a image decoder and a text encoder. We directed exhaustive analyses on the openly accessible dataset to represent the advantages of our proposed strategy. As a component of this unique task, we have likewise added to the text-to-face generation dataset.

The results obviously beat those of the ongoing methods. Our model exhibited the ability to change over printed signals into reasonable visuals by creating great face pictures that firmly paired the information depictions. It furnished a 97% accuracy rate with no counterfeit examples. We have shown by means of exhaustive examination that our proposed generative sound organization can produce practical, excellent pictures that intently impersonate genuine truth labels and faces.

## 6. FUTURE SCOPE

Craftsmen and architects might utilize this innovation to transform their text based considerations into imaginative pictures, accelerating their inventive approach. This task progresses text-to-face generation

and its multi-layered applications utilizing DC-GANs to create high-goal pictures with expanded accuracy by conquering past impediments and exhibiting the force of a completely prepared GAN architecture.

Our photographs may likewise be decided by people. Zeroing in on more extravagant and more precise face-related data for the suggested engineering would further develop picture quality and make made faces more like depictions. Public security, criminological investigation, and others will benefit from this review. Many purposes are conceivable with this venture. It makes fascinating realistic substance for sites and online entertainment.

## REFERENCES

- [1]. I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, pp. 2672–2680, 2014.
- [2]. S. Hong, D. Yang, J. Choi, and H. Lee, "Inferring semantic layout for hierarchical text-to-image synthesis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 7986–7994, 2018.
- [3]. A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves, et al., "Conditional image generation with pixelcnn decoders," in *Advances in neural information processing systems*, pp. 4790–4798, 2016.
- [4]. H. Zhang, T. Xu, H. Li, S. Zhang, X. Wang, X. Huang, and D. N. Metaxas, "Stackgan++: Realistic

image synthesis with stacked generative adversarial networks,” IEEE transactions on pattern analysis and machine intelligence, vol. 41, no. 8, pp. 1947–1962, 2018.

[5]. C. Wah, S. Branson, P. Welinder, P. Perona, and S. Belongie, “The caltech-ucsd birds-200-2011 dataset,” 2011.

[6]. M.-E. Nilsback and A. Zisserman, “Automated flower classification over a large number of classes,” in 2008 Sixth Indian Conference on Computer Vision, Graphics & Image Processing, pp. 722–729, IEEE, 2008.

[7]. T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Doll’ar, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in European conference on computer vision, pp. 740–755, Springer, 2014.

[8]. D. P. Kingma and M. Welling, “Auto-encoding variational bayes,” arXiv preprint arXiv:1312.6114, 2013.

[9]. S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, “Generative adversarial text to image synthesis,” arXiv preprint arXiv:1605.05396, 2016.

[10]. A. Radford, L. Metz, and S. Chintala, “Unsupervised representation learning with deep convolutional generative adversarial networks,” arXiv preprint arXiv:1511.06434, 2015.

[11]. Dataset link:  
[https://www.Kaggle.com/datasets/yunting0123/img\\_a\\_lign\\_celeba](https://www.Kaggle.com/datasets/yunting0123/img_a_lign_celeba)