

# ENHANCING TEXTUAL DESCRIPTION USING DEEP LEARNING WITH RC-GAN GENERATION

<sup>1</sup>R. Jagadeeswari, <sup>2</sup>A. Srujan Reddy, <sup>3</sup>N. Venkata Mahesh Kumar Reddy, <sup>4</sup>N. Bharath, <sup>5</sup>S. Raaga Sindhu

<sup>1</sup>Assistant Professor, <sup>2,3,4,5</sup> Students

Department Of Computer Science And Engineering  
Bharath Institute Of Higher Education And Research  
Chennai, India 600073

**Abstract-** One method employed to create pixels representations corresponding to verbal descriptions is known as "text-to-image generation," impacting a broad spectrum of applications and research domains, including exploration, design, computer-aided drafting, image restoration, labelling, and portraiture. The primary challenge lies in consistently generating realistic photographs under specified conditions. Existing algorithms for script-to-image conversion often fail to accurately align images with the accompanying text descriptions. In our investigation, we addressed this challenge by devising a deep learning architecture named the recurrent convolutional generative adversarial network (RC-GAN) specifically fine-tuned for producing semantically concise snapshots. RC-GAN effectively translates conceptual visual elements from linguistic cues to pixel-level representations, bridging the gap between advancements in text and image modelling. To train our proposed model, we utilised various datasets, employing metrics such as inception score and peak signal-to-noise ratio (PSNR) to analyse its performance. Empirical findings showcase an inception required score and a PSNR required measurement, indicating that our model can produce more realistic images from verbal descriptions. Our future endeavour involves training the proposed system on various datasets to further improve its capabilities.

**Key Words:** Deep Learning (DL), Convolutional Neural Network (CNN), Recurrent Neural Network (RNN), Generation of Image, Generative Adversarial Networks (GAN).

## I. INTRODUCTION:

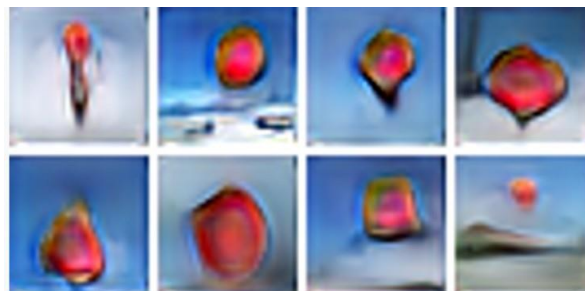
Individuals rapidly evoke images in their brains to imagine the parts of tales they examine or absorb. Usually addressed as "seeing with the mind's eye, inner picture formation is important for a range of intellectual operations, such as recalling, deliberation, and idea processes. A noteworthy progress in user ingenuity is the creation of technology that can translate written characterisation into graphics and understand the relation between perception and language. The milestone of deep learning and artificial intelligence has enabled expansion growth in recent years in picture execution techniques and digital vision implementations. The fabrication of script from picture is one of these growing sectors. The work of building visually realistic pictures from script inputs is called as text-to-image. Textual depiction extraction from visual content, also known as T2I synthesis, stands as the antithesis to image captioning. I2T synthesis entails the generation of a linguistic narrative derived from an initial visual input. When crafting Text-to-Image systems, the model ingests a textual description authored by a human and produces an RGB representation aligning with the provided text. T2I synthesis stands as a pivotal realm of inquiry, owing to its vast potential across numerous domains. Perusing images, altering visual content, crafting visual compositions, adding textual descriptions, sketching lifelike portraits, designing for industrial purposes, and manipulating visual content stand out as prevalent uses for translating (I) textual descriptions into realistic imagery.

The progression of generative adversarial networks (GANs) has illustrated impressive efficacy in photo generation, photo enhancement, data amplification, and photo transformation. GANs leverage convolutional neural networks (CNNs) rooted in deep learning principles. They comprise a duo of neural networks: one dedicated to data generation and another to discerning true from false data. Employing principles from game theory, GANs are geared towards teaching the apparatus to generate specimens and the differentiator to distinguish veracity. To enhance the realism of generated images, textual information undergoes encoding using Recurrent Neural Networks (RNNs), while convolutional layers facilitate photo deciphering. Our innovation, the Recurrent Convolution GAN (RC-GAN), presents a streamlined yet potent structure for translating textual descriptions into synthesised images. Instructed on the various Dataset, the model validates the authenticity of the synthesised images. The primary advancements elucidated by this study encompass:

- Constructing a sophisticated deep learning architecture, RC-GAN, aimed at fabricating images of heightened realism.
- Constructing lifelike visuals based on provided textual depictions.
- Enhancing the benchmarks and Peak Signal-to-Noise Ratio (PSNR) of images synthesised from textual input.

## II. BACKGROUND

Goodfellow initially introduced Generative Adversarial Networks (GANs) in 2014, but Reed and colleagues in 2016 were the pioneers in utilising them for the generation of images from textual input. Salimans and team suggested refining training methodologies for models that had not undergone training yet, leading to enhanced results across various datasets which has taken as references from different publications devised an vigilance-driven recurrent neural network architecture. The approach involved learning linguistic-dot relationships through an attentional feature extractor and Pixel-boundary relationships via an autoregressive decoder. Liu et al. introduced a flexible context-driven photo creation framework and performed extensive experiments across various conditional generation tasks. Gao et al. proposed an efficient strategy dubbed Lightweight Dynamic Conditional GAN (LD-CGAN), which extracted textual attributes and furnished image features via multi parameter characteristic extraction. Dong et al. Trained a model for text-based image generation without manual inspection. Berrahal et al. concentrated on advancing text-to-image acquisition utilisations by employing a Deep Fusion GAN (DF-GAN) for generating person images from textual descriptions. Zhang et al. put forward Cross-Domain GAN Fusion (CF-GAN), which aimed to translate words to photo with richer context sensitive. Generally, contemporary text-to-illustrations synthesis approaches heavily rely on extensive parameterisation and computationally intensive operations to produce high-resolution images, resulting in unstable and costly training procedures. In 2015, scholars from the University of Toronto unveiled the inaugural contemporary text-to-image framework known as align DRAW. This advancement expanded upon the pre-existing DRAW architecture, which employed an iterative variational auto encoder complemented by an attention mechanism, to incorporate textual sequences. Despite yielding outputs characterised by haziness and a lack of photorealism, align DRAW demonstrated proficiency in accommodating supplementary tasks, such as envisioning scenarios like a "Stop sign soaring into a cerulean sky," and could generalise objects not explicitly present in its training data, such as a crimson school bus. This underscores its capacity to transcend mere memorisation of training set information.



Eight visuals emerged following the textual cue 'A stop sign soars against the azure heavens' using Align DRAW (2015), magnified to unveil intricate details.

In 2016, Reed, Akata, Yan, and colleagues pioneered the utilisation of adversarial generative networks for the conversion of text into images. Leveraging models trained on specialised, locality-focused datasets, they utilised descriptive cues such as "a dark-feathered avian with a prominent, curved beak" to generate images depicting birds and flowers that were deemed "convincingly realistic." However, when employing a framework instructed on the most heterogeneous COCO data pool, the resultant photos exhibited inconsistency in their level of detail and were described as lacking in satisfactory quality. Subsequent advancements in this field encompass the utilisation of methodologies such as VQGAN+CLIP, XMC-GAN, and GauGAN2.



Illustrations depicting a "Stop sign soaring through the azure sky" grace both DALL·E 2 (upper) from April 2022, and DALL·E 3 (lower, from September 2023). Open AI introduced DALL-E, a transformative system for converting text to images, in January 2021, garnering significant public attention. Its successor, DALL·E 2, emerged in April 2022, boasting enhanced capabilities in crafting intricate and lifelike visuals, subsequently achieving a stable deployment by August 2022. This breakthrough came amid the proliferation of text-to-image technologies. Additionally, the advent of text-to-video platforms such as Runway, Make-A-Video, Imaginary Video, Midjourney, and Phenaki, in August 2022, ushered in a new era of multimedia content creation from textual or visual prompts. Notably, advancements in August 2022 underscored the adaptability of large-scale text-to-image frameworks, showcasing the feasibility of tailored customisation. Personalising text-to-image conversions enables the incorporation of novel concepts by training the model with a limited set of images depicting previously unseen objects. This entails translating textual descriptions into corresponding image representations, thereby expanding the model's repertoire.

### III. DATA ARRAY

#### 3.1. Data Array

The dataset utilised was the Oxford-102 Flowers dataset, comprising 7125 diverse flower photos categorised into 98 classes. Each class encompasses between 30 to 50 images, accompanied by 10 corresponding textual descriptions per image. Within this investigation, we analysed 7000 images for the training phase. The architecture underwent training for 250 epochs utilising particular data pool.

#### 3.2. Data grooming

Through the duration of data gathering and retrieving, a pool of 7180 images featuring diverse settlements alongside consistent texts were amassed. To standardise the textual data, the NLTK token extractor was employed, facilitating the conversion of text sentences into individual words. These words were then organised into tagged lists, subsequently transformed into subtitle tags. In order to achieve uniformity, all images were resized to dimensions of  $128 \times 128$  upon upload. Both the training and test images underwent this resizing process. As part of the preprocessing stage, with both the word catalog and photographs being integrated into the techniques. A crucial component of training the text-to-image model involves utilising datasets that combine images with accompanying textual descriptions. One such prominent dataset is COCO (Common Object in Context), which was introduced in 2014, comprising around 98,000 images showcasing various objects, each accompanied by five human-annotated captions. Additionally, there exist smaller datasets like Oxford-120 Flowers and CUB-200 Birds, every single one containing roughly 9000 images, focusing exclusively on blooms and fowls respectively. Leveraging these datasets for training purposes is often deemed advantageous due to their more specialised subject matter, facilitating the development of high-quality text-to-image models.



Instances of visual content coupled with textual descriptions extracted from three widely utilised public datasets for educational purposes. Text-based blueprints for generating images.

#### IV. EXPERIMENTS AND RESULTS

This section entails the empirical examination of both the floral specimens and the synthesised visuals. The proposed framework underwent training on an Nvidia 1070 Ti GPU, equipped with 32 GB of memory, operating under the Windows 10 environment. The optimisation of generator and separation weights ensued through the utilisation of the Adam improver, employing a segment size of 64 and a learning coefficient at 0.0002. Figure 2 illustrates the originated pictures alongside textual descriptors from the authentic data pool. To gauge the productiveness the invented model, we computed baseline metrics such as IS and PSNR values. These metrics serve as benchmarks, delineating the diversity and fidelity of the generated imagery. PSNR, serving as an indicator of image fidelity, quantifies the signal-to-noise ratio between two images, typically denoted in decibels. The PSNR value augments proportionally with the enhancement in image fidelity, as it derives from the equation:  $PSNR = 10\log_{10} (R^2 / MSE)$ , where R2 represents the maximum pixel value and MSE denotes the mean squared error.

To authenticate the suggested method, outcomes are juxtaposed against prevailing frameworks such as GAN-INT-CLS, Stack GAN, Stack GAN++, HDGAN, and Dual AttnGAN using various data pools.

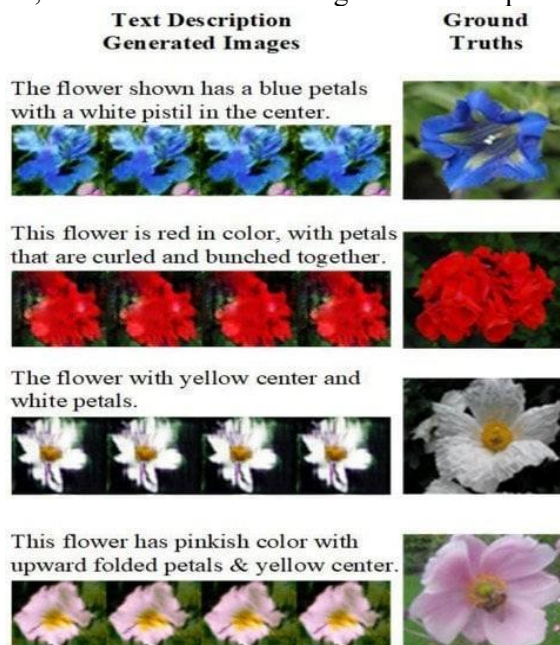


Figure 2. Embed textual explanations and accompanying visuals featuring essential details.



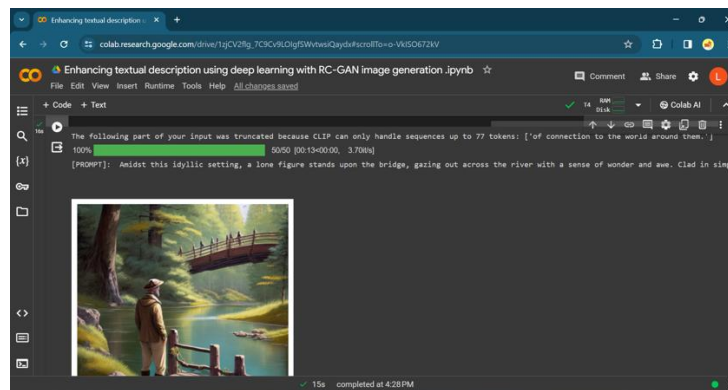


Figure 3. An Image generated with the instructions given in prompt

In this evocative image generated using RCGAN technology, we are transported to a tranquil scene, where nature's beauty intertwines with human contemplation. The picturesque setting captures the essence of serenity, with lush foliage reflected in the gentle ripples of the river below.

At the heart of this scene stands a solitary figure, positioned upon a bridge that spans the tranquil waters. The figure's stance exudes a sense of quiet introspection, as they gaze outwards towards the horizon with a mixture of wonder and awe. Cloaked in simplicity, the individual becomes a focal point amidst the harmonious blend of earth and water.

The use of deep learning techniques in generating this image imbues it with a remarkable level of detail and realism. From the delicate interplay of light and shadow to the subtle nuances of expression, every element is meticulously crafted to evoke a sense of immersion and emotion.

As viewers, we are invited to linger in this moment of solitude, allowing ourselves to be enveloped by the beauty and tranquility of the scene. It serves as a poignant reminder of the profound connection between humanity and the natural world, and the profound sense of peace that can be found in moments of quiet contemplation.

## V. CONCLUSION

Text-to-picture creation stands as a trending theme in contemporary discussions within the realms of computerised vision and linguistic understanding. In order to fabricate visually authentic and contextually coherent pictures, we unveiled a deep learning sophisticated architecture, termed RC-GAN, and delineated its functionality within the frameworks of both computerised vision and linguistic understanding. This architectural construct underwent a training regimen employing text encoding and subsequent image decoding. A comprehensive array of experiments conducted on floral dataset manifested that the developed GAN architecture yields superior image quality in contrast to pre-existing models, showcasing the highest recorded Inception Score (IS). The efficacy of our proposed methodology underwent scrutiny through confrontation leading-edge approaches utilising IS metrics. The assessment and juxtaposition of text-to-picture model efficacy present a multifaceted challenge, involving the evaluation of numerous desirable attributes. Analogous to all generative picture frameworks, it is appealing that the synthesised pictures exhibit realism (within the framework of they appear to plausibly originate from the training set) while also encompassing diverse stylistic variations. Text-to-picture models are emblematic of the aspiration for semantic alignment between the generated images and their textual descriptors. Various evaluation frameworks have emerged to scrutinise these aspects, encompassing both automated algorithms and human-driven assessments. A prevalent algorithmic measure employed to gauge image fidelity and diversity is the initial score (IS), derived from the tag distribution anticipated by a pre-installed Inception v3 image classifier when enacted to snapshots spawned by a text to picture prototype.

## REFERENCES:

- [1]. Kosslyn, S.M.; Ganis, G.; Thompson, W.L. The neurological foundation of visual representations. Published in the journal "Nature Reviews Neuroscience" in the year 2001, volume 2, pages 635–642. [Accessed via Google Scholar] [CrossRef] [PubMed]
- [2]. Karpatio, A.; Fei-Fei, L. Deep visual-semantic alignments for image description generation. Proceedings of IEEE Conference on Computer Vision and Pattern Recognition, Boston, MA, USA, 7-12 June 2015; page 3128–3137. [GoogleScholar]
- [3]. Erhan, D.; Bengio, S.; Vinyals, O.; Toshev, A. Generate a neural picture and demonstrate it. IEEE Conference on Computer Vision and Pattern Recognition, June 7–12, 2015, Boston, MA, USA, Proceedings, 3156–3164.[Scholar on Google]
- [4]. Salakhudinov, R.; Zemel, R.; Bengio, Y.; Xu, K.; Ba, J.; Kiros, R.; Cho, K.; Courville, A. Presenting, focusing, and enlightening: Neural generation of images and titles with visual attention. Published in the proceedings of the

- International Conference on Machine Learning, held in Lille, France, from July 6–11, 2015. Volume 20, Pages 2048–2057. [Accessed via Google Scholar]
- [5]. Albawi, S.; Muhammad, T.A.; Al-Zawi, S. Focal convolution. Presented at the 2017 International Conference on Engineering and Technology (ICET) in Antalya, Turkey, from October 21–23. Published in August 2017; pp. 1–6. [Accessed via Google Scholar]
- [6]. Kim, P. Spatial Transformational Neural Network. Utilizing MATLAB for Deep Learning; Published by Springer in Berlin/Heidelberg, Germany, in the year 2017; pages 121–147. [Accessed via Google Scholar]
- [7]. Generative adversarial networks by Goodfellow, I. J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y. Published on arXiv with reference numbers arXiv:1406.2661 and arXiv:2014.061. Accessed through Scholar Google and Cross Reference platforms.
- [8]. Reed, S.; Akata, Z.; Yan, X.; Logeswaran, L.; Schiele, B.; Lee, H. Creating contrasting images from textual descriptions. Published on arXiv in 2016 with the reference arXiv:1605.05396. Accessed via Google Scholar.
- [9]. Salimans, T.; Bonulo, I.; Zaremba, W.; Cheung, V.; Radford, A.; Chen, X. Employing advanced methodologies for achieving improvements. Presented at the Advances in Neural Information Processing Systems 29 (NIPS 2016) conference held in Barcelona, Spain, from 5–10 December 2016. Accessed via Google Scholar.
- [10]. Zia, T.; Arif, S.; Murtaza, S.; Ullah, M.A. Generating images from text descriptions utilizing recurrent neural networks with attention mechanisms. Published on arXiv in 2020 with the reference arXiv:2001.06658. Accessed via Google Scholar
- [11]. Gao, L.; Chen, D.; Zhao, Z.; Shao, J.; Shen, H.T. A Lightweight Dynamic Conditional GAN for Text-to-Image Synthesis with Pyramid Attention. Pattern recognition. 2021, 110, 107384. [Google Scholar] [CrossRef]
- [12]. Dong, Y.; Zhang, Y.; Ma, L.; Wang, Z.; Luo, J. Unsupervised Generation of Images from Textual Inputs. Published in the journal "Pattern Recognition" in 2021, volume 110, page 107573. Accessed via Google Scholar.
- [13]. Berrahal, M.; Azizi, M. Crafting an optimal model for generating images from textual descriptions using techniques from generative adversarial networks. Published in the "Indonesia Journal of Electrical Engineering and Computer Science" in 2022, volume 25, pages 972–979. Accessed via Google Scholar.
- [14]. Zhang, Y.; Han, S.; Zhang, Z.; Wang, J.; Bi, H. CF-GAN: Introducing a cross-domain adversarial network for synthesizing images from text descriptions by merging features. Published in the journal "Visual Computing" in 2022, pages 1–11. Accessed via Google Scholar.