



ISSN: 2454-9940



**INTERNATIONAL JOURNAL OF APPLIED
SCIENCE ENGINEERING AND MANAGEMENT**

E-Mail :
editor.ijasem@gmail.com
editor@ijasem.org

www.ijasem.org

TEXT TO IMAGE GENERATOR USING DEEP LEARNING

¹ N. Anjamma, ² K. Nikhitha, ³ K. Sai Preethi, ⁴ K. Rakesh

¹ Assistant Professor, anjammanomula123@gmail.com

^{1,2,3,4} Teegala Krishna Reddy Engineering College

² Nikkirdykothur28@gmail.com, ³ Kottesaipreethi@gmail.com, ⁴ rakeshvarma9010@gmail.com

Abstract: Text to image synthesis refers to the method of generating images from the input text automatically. Deciphering data between picture and text is a major issue in artificial intelligence. Automatic image synthesis is highly beneficial in many ways. Generation of the image is one of the applications of conditional generative models. For generating images, GAN (Generative Adversarial Models) are used. Recent progress has been made using Generative Adversarial Networks (GAN). The conversion of the text to image is an extremely appropriate example of deep learning. It is particularly interesting to us because it programmatically synthesizes one data type into another, generating a photorealistic image based off a phrase. One field of application for our model is to aid language learning since children understand the meaning of a phrase better if they see images that correspond to text phrases.

Index terms - *Text-to-Image Synthesis, Generative Adversarial Networks (GAN), Conditional Generative Models, Deep Learning, Image Generation.*

1. INTRODUCTION

Text-to-image synthesis is a cutting-edge technology that generates images based on textual descriptions. This process, which bridges the gap between language and visual content, holds significant potential in various applications, ranging from aiding language learning to enhancing creative processes in industries like advertising and entertainment. At the core of this innovation are Generative Adversarial Networks (GANs), a class of generative models capable of creating new, previously unseen content. GANs consist of two competing neural networks: a generator that creates images and a discriminator that evaluates their authenticity, thereby improving the model's accuracy through continuous iteration and competition [1]. The result is the generation of realistic images that align with the input text, such

as generating a "flower with pink petals" from a descriptive phrase.

The power of GANs in text-to-image synthesis lies in their ability to transform abstract concepts or descriptions into visual representations, mimicking how humans associate words with images. However, one of the challenges in this domain is the complexity of mapping a single text description to a variety of potential visual representations. For example, the phrase "a dog in a park" can be interpreted and visualized in many different ways, with variations in color, size, and background. Deep learning models, especially GANs, are trained to handle such ambiguities by learning patterns from large datasets, thereby improving their ability to generate more accurate and diverse images for a wide range of textual descriptions [2].

Despite the promise of text-to-image synthesis, there are still challenges that need to be addressed. One significant issue is the issue of diversity in image generation. While GANs have made significant strides in creating photorealistic images, the generated content may sometimes lack variety, especially when dealing with highly abstract or complex descriptions. To overcome this, researchers continue to focus on refining the model architectures and training techniques, such as using larger, more diverse datasets and incorporating techniques like attention mechanisms to enhance the model's ability to focus on relevant features in the text [3].

In the future, text-to-image synthesis is expected to become a powerful tool for various applications, such as in virtual reality, educational tools, and creative design. By providing an intuitive way to generate images from simple text, this technology could revolutionize fields like content creation and language learning, enabling users to gain a better understanding of words and phrases by visualizing them in real-time. As deep learning and GAN technologies continue to advance, the potential for text-to-image synthesis to impact multiple industries grows, paving the way for more personalized and immersive experiences [4].

2. LITERATURE REVIEW

Text-to-image synthesis has seen significant advancements in recent years, particularly with the rise of Generative Adversarial Networks (GANs), which have been pivotal in generating realistic images from textual descriptions. The foundational work in this area stems from the exploration of generative models that have the capacity to learn data distributions and generate novel instances. A key component of these models, GANs consist of two neural networks: a generator that creates

synthetic data and a discriminator that evaluates it against real data, thereby driving the generator to produce more realistic outputs over time [1]. These models have been widely adopted for text-to-image synthesis due to their ability to map textual descriptions into photorealistic images.

In early attempts at text-to-image generation, models like the Generative Adversarial Text-to-Image Synthesis proposed by Reed et al. [4] were developed to directly translate natural language descriptions into images. The authors introduced a conditional GAN where the generator is conditioned on the textual input, allowing it to generate images that reflect the content described. This model laid the groundwork for later works by demonstrating that textual information could guide image generation, which is especially useful in domains such as virtual reality, content creation, and education. Reed's model used word-level embeddings and image features to improve the alignment between the generated image and textual description, marking an important step in improving synthesis accuracy.

Following this, StackGAN, introduced by Zhang et al. [3], took a novel approach by stacking GANs to improve the quality and resolution of the generated images. In their architecture, the first GAN generates a low-resolution image from the text, and the second GAN refines it to a high-resolution image. This stacking mechanism allowed StackGAN to create more detailed and photorealistic images, significantly enhancing the visual fidelity of the generated outputs. The method was particularly effective in handling complex descriptions, as it combined the benefits of a coarse-to-fine image generation process. StackGAN's success in creating high-quality images

demonstrated the importance of multi-stage approaches for text-to-image synthesis.

Another important contribution was made by Zhang et al. in their development of AttnGAN, a model that incorporated attention mechanisms into the GAN framework to address fine-grained image details [6]. AttnGAN uses an attention mechanism to focus on different parts of the text description during the image generation process. By attending to specific words and phrases in the input text, the model is able to generate images that more accurately reflect the nuances of the description. This attention mechanism is especially helpful in cases where the textual description involves multiple objects or intricate details, as it allows the generator to pay closer attention to relevant aspects of the input. The introduction of attention mechanisms in text-to-image synthesis has significantly improved the model's ability to capture fine-grained details and produce more coherent and contextually accurate images.

While these models have made substantial progress in generating high-quality images from text, they still face challenges related to ambiguity in natural language. Text descriptions can often be vague or contain multiple interpretations, which makes it difficult for a single model to consistently generate the most appropriate image. Yadav and Vishwakarma [1] highlight the difficulty of addressing such ambiguities, as deep learning models like GANs must be trained on large datasets to understand the various configurations and contexts that can arise from a single description. For example, the description "A bird on a tree" could lead to several valid image outputs, depending on the tree type, bird species, and the surrounding environment. Overcoming this challenge requires sophisticated training methods that can account for

such variations and still generate contextually appropriate images. One potential solution involves enhancing models with larger and more diverse training datasets that can encompass a broader range of visual and textual scenarios.

To mitigate these issues, researchers have explored alternative architectures and training methods. Gregor et al. [2] introduced the Draw model, a recurrent neural network designed to generate images by drawing them step by step. This approach can improve the generation of complex images by progressively refining the output over multiple iterations, similar to how a human might draw an image. Although Draw is not strictly a GAN, its use of recurrent processes has influenced subsequent text-to-image models by showing the benefits of iterative image creation. Such models offer a different perspective on the process, focusing on refining the image as it is generated, which could be valuable for handling ambiguous or complex textual inputs.

The challenge of generating diverse images from the same text description has also been addressed by techniques like conditional generation and Variational Autoencoders (VAEs). These models introduce latent variables that allow for the generation of multiple plausible images from a single input description. By conditioning on the text and introducing a probabilistic model, these approaches aim to explore a broader space of possible images. Yadav and Vishwakarma [1] note that these models can produce more varied outputs, which is essential for applications like content generation, where multiple interpretations of a description are often required.

In addition to these methods, recent advances in multi-modal learning and reinforcement learning

have further contributed to the refinement of text-to-image synthesis models. These approaches integrate feedback mechanisms into the generation process, allowing models to learn from both textual and visual data concurrently. This has proven to be effective in improving the coherence and realism of the generated images. Furthermore, new techniques in domain adaptation and transfer learning have helped models become more generalized, enabling them to generate high-quality images across various domains with minimal retraining.

Overall, the development of text-to-image synthesis using GANs and related models has made significant strides in creating realistic and contextually relevant images from textual descriptions. While challenges such as ambiguity in text and the need for large-scale datasets remain, ongoing research continues to improve the capabilities of these models. With advancements in attention mechanisms, multi-stage generation processes, and probabilistic models, text-to-image synthesis holds great promise for applications in creative industries, education, and beyond. As the technology matures, its integration into real-world applications will likely transform the way we interact with both text and images, opening up new opportunities for visual content creation and learning.

3. METHODOLOGY

The proposed model Stable Diffusion is a powerful, open-source text-to-image generation model. While there exist multiple open-source implementations that allow you to easily create images from textual prompts, Keras CV's offers a few distinct advantages. The "Stable Diffusion 2" model is a text-to-image model developed by Stability AI. It is a latent diffusion model that uses a fixed, pretrained

text encoder (OpenCLIP-ViT/H) to generate and modify images based on text prompts. The model takes a latent seed and a text prompt as input, and the latent seed is used to generate random latent image representations of the text prompt. The model can generate high-resolution, photorealistic images from text prompts and is trained on a subset of the LAION-5B dataset. It uses a UNet backbone and a v-objective loss function during training.

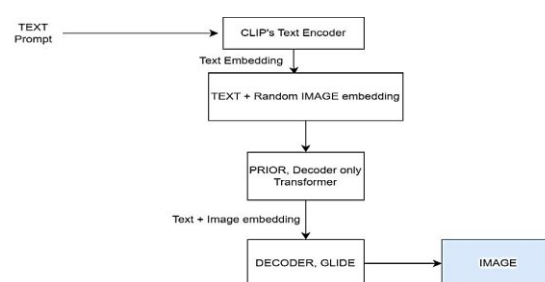


Fig.1 Proposed Architecture

The architecture of a deep learning text-to-image model regularly involves two main parts, which include a text encoder and an image generator. The text encoder converts the inputted text prompt into a compact and meaningful representation, which is then delivered to the image generator to create a final image. Text-to-image models have been built using a variety of architectures. The text encoding step may be performed with a recurrent neural network such as a long short-term memory (LSTM) network, though transformer models have since become a more popular option. For the image generation step, conditional generative adversarial networks have been commonly used, with diffusion models also becoming a popular option in recent years. Rather than directly training a model to output a high-resolution image conditioned on a text embedding, a popular technique is to train a model to generate low-resolution images, and use one or more auxiliary deep learning models to upscale it, filling in finer details. Text-to-image models are trained on

large datasets of (text, image) pairs, often scraped from the web. With their 2022 Imagen model, Google Brain reported positive results from using a large language model trained separately on a text-only corpus (with its weights subsequently frozen), a departure.

i) Dataset:

A widely used dataset for text-to-image generation is the Microsoft COCO (Common Objects in Context) dataset. It contains over 330,000 images, each paired with multiple descriptive captions. These captions provide detailed descriptions of the scenes, objects, and actions within the images, making it a valuable resource for training text-to-image models. The dataset is commonly used to evaluate the performance of image generation models, as it includes diverse and complex visual content across various contexts and environments.

ii) Text to Image Generator:

A Text-to-Image Generator using deep learning typically utilizes a Conditional Generative Adversarial Network (CGAN) to generate images from textual descriptions. The system comprises three main components: a Text Encoder, a Generator, and a Discriminator. The Text Encoder converts the input text into a numerical representation, often using embeddings or recurrent networks such as LSTMs or GRUs. This encoded vector captures the semantic meaning of the text and is passed to the Generator.

The Generator, usually built with deep convolutional or deconvolutional layers, takes the text vector as input to generate an image. Initially, the generated image is of low resolution, which is progressively refined using successive layers of deconvolution to upscale the image. The

Discriminator evaluates both real and generated images, conditioned on the input text, to determine whether an image is real or fake.

The training process involves an adversarial setup where the Generator aims to produce images that deceive the Discriminator, while the Discriminator improves its ability to differentiate between real and fake images. This iterative process enhances the Generator's ability to produce realistic images from new text inputs. Once trained, the system can generate high-quality images based on any given textual description.

4. EXPERIMENTAL RESULTS



Fig.2 Astronaut in space

A text-to-image generator using deep learning works by training a model on large datasets of images and their textual descriptions. The model learns to associate words with visual patterns, enabling it to generate realistic images based on input text prompts

5. CONCLUSION

In conclusion, a text-to-image generator is an innovative application of generative models, particularly within the framework of Generative Adversarial Networks (GANs). This technology addresses the challenge of translating textual

descriptions into visually coherent and realistic images. The generator, at the heart of the system, synthesizes visual representations from text, refining its output through adversarial training. The discriminator provides feedback, creating an iterative improvement loop that helps the generator produce images that are indistinguishable from real-world examples. Implementing a text-to-image generator requires careful data preprocessing, thoughtful neural network design, and the use of effective training methodologies. The selection of a suitable dataset, efficient network structures, and relevant evaluation metrics are crucial to the system's success. Optional techniques like stable diffusion further enhance image quality, providing additional control over the generation process. The deployment of a text-to-image generator offers immense potential for applications in creative content generation, as well as in situations where visual information is scarce or unavailable, such as aiding language learning or enhancing accessibility in various fields.

6. FUTURE SCOPE

The future of text-to-image generators is promising, with advancements in generative models, deep learning, and natural language processing driving progress. Key areas include enhancing image realism, offering fine-grained control, and integrating multimodal capabilities like audio and video generation. Cross-domain translation, zero-shot learning, and commonsense reasoning will improve adaptability and contextual understanding. Ethical considerations, such as addressing biases, will ensure fairness. These advancements have the potential to revolutionize industries like gaming, education, and virtual reality, fostering new possibilities in interactive content creation.

REFERENCES

- [1] Ankit Yadav¹, Dinesh Kumar Vishwakarma², Recent Developments in Generative Adversarial Networks: A Review (Workshop Paper),2020.
- [2] Gregor, K., Danihelka, I., Graves, A., Rezende, D., and Wierstra, D. Draw: A recurrent neural network for image generation. In ICML, 2021.
- [3] Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, "StackGAN: Text to Photo-realistic Image Synthesis with Stacked Generative Adversarial Networks" in Rutgers University and Lehigh University August 2020.
- [4] Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, Honglak Lee, "Generative Adversarial Text to Image Synthesis" in University of Michigan and Max Planck Institute for Informatics June 2021.
- [5] Stian Bodnar, Jon Shapiro, "Text to Image Synthesis Using Generative Adversarial Networks" in The University of Manchester May 2021.
- [6] Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. Attngan: Fine grained text to image generation with attentional generative adversarial networks. CoRR, abs/1711.10485, 2020. Mehdi Mirza, Simon Osindero, Conditional Generative AdversarialNets, 2021.
- [7] Goodfellow, I.J.; Pouget-Abadie, J.; Mirza, M.; Xu, B.; Warde-Farley, D.; Ozair, S.; Courville, A.; Bengio, Y. Generative adversarial networks. arXiv 2014, arXiv:1406.2661.

[8] Reed, S.; Akata, Z.; Yan, X.; Logeswaran, L.; Schiele, B.; Lee, H. Generative adversarial text to image synthesis. arXiv 2016, arXiv:1605.05396.

[9] Salimans, T.; Goodfellow, I.; Zaremba, W.; Cheung, V.; Radford, A.; Chen, X. Improved techniques for training gans. In Proceedings of the Advances in Neural Information Processing Systems 29 (NIPS 2016), Barcelona, Spain, 5–10 December 2016.

[10] Zia, T.; Arif, S.; Murtaza, S.; Ullah, M.A. Text-to-Image Generation with Attention Based Recurrent Neural Networks. arXiv 2020, arXiv:2001.06658.