AI BASED TEXT TO IMAGE GENERATION USING CNN

P. B. Vikhe¹, Sakshi Bagade², Komal Kadam³, Sakshi Vikhe⁴, Kote Pratiksha⁵

¹ Professor, Computer Engineering, Pravara Rural Engineering College, Maharashtra, India

² Student, Computer Engineering, Pravara Rural Engineering College, Maharashtra, India

³ Student, Computer Engineering, Pravara Rural Engineering College, Maharashtra, India

⁴ Student, Computer Engineering, Pravara Rural Engineering College, Maharashtra, India

⁵Student, Computer Engineering, Pravara Rural Engineering College, Maharashtra, India

ABSTRACT

This paper presents a novel approach to generate images from textual descriptions using Convolutional Neural Networks (CNN). The system interprets user input in natural language and synthesizes a visually coherent image that represents the input text. By leveraging a deep learning-based encoder-decoder architecture, the model converts linguistic features into spatial image representations. The CNN acts as a generative decoder that progressively constructs images from learned text embeddings. This technique finds practical applications in art generation, accessibility tools, content creation, and more. The implementation demonstrates promising results with moderate computational needs and strong generalization for simple to moderately complex image concepts.

This technology has broad applications across various domains, including photo-searching, digital art, computer-aided design, and image reconstruction. A key component in this process is the Convolutional Neural Network (CNN), a deep learning architecture renowned for its effectiveness in image recognition and feature extraction. CNNs excel in processing visual data, making them ideal for tasks requiring high accuracy in image classification and generation.

Keyword: - Text-to-Image, Deep Learning, CNN, Generative Model, Image Synthesis, Natural Language Processing, Artificial Intelligence.

1. INTRODUCTION

The ability to generate images from text represents a critical advancement in the field of AI, combining natural language understanding with computer vision. Recent breakthroughs in deep learning, especially Convolutional Neural Networks (CNNs), have enabled systems to understand and map complex relationships between words and visual elements. This paper explores how CNNs can be applied to generate images based on textual descriptions, bypassing traditional graphics tools. The goal is to enable machines to "draw" what is described in language, thereby bridging the semantic gap between text and visual representation. Text-to-image generation holds immense potential across a wide range of applications including entertainment, accessibility, education, and virtual reality. It allows non-artistic users to bring their imagination to life simply by describing what they envision in words. With the growing availability of large-scale datasets and pre-trained embeddings, training effective models is now more feasible than ever. However, generating high-resolution, semantically accurate images remains a complex challenge due to the ambiguity and variability inherent in natural language. This research addresses these challenges by utilizing CNN-based architectures that effectively translate encoded text features into structured visual data. Human cognition naturally converts narratives into mental imagery-a process fundamental to memory, reasoning, and creativity. This innate ability to "see with the mind's eye" inspires our technology, which bridges language and visualization by transforming text into precise visual representations. By capturing this connection between words and images, we're advancing how people express and develop ideas visually.

2. LITERATURE SURVEY:

2.1 Paper Name: Zero-Shot Text-to-Image Generation

Authors: Aditya Ramesh et al. (OpenAI – DALL \cdot E)

Year: 2021

Summary: This paper introduced DALL·E, a transformer-based model capable of generating diverse and realistic images from text prompts in a zero-shot setting. Unlike traditional CNN or GAN-based methods, DALL·E uses a discrete variational autoencoder (dVAE) to tokenize images and combines it with a GPT-like model to learn text-to-image mappings. It demonstrated the ability to generate novel concepts and plausible scenes, even for prompts never seen during training. This work significantly broadened the scope of text-to-image generation, showing potential for creative design, education, and visualization tasks.

2.2 Paper Name: T2F: Text to Face Generation Using Deep Neural Networks

Authors: Anish Mittal, R. Soundararajan

Year: 2019

Summary: This research focused on generating realistic human facial images from descriptive textual inputs using deep neural networks. The model converts text into embeddings and feeds them into a CNN decoder to create face images that match the description. Applications include law enforcement, facial composite sketches, and personalized avatar creation. The authors emphasize handling diverse human features such as skin tone, facial hair, age, and emotion from relatively vague textual cues, addressing a domain-specific challenge in text-to image synthesis.

2.3Paper Name: Object-driven Attentive Generative Adversarial Network for Fine-Grained Image Generation from Text

Authors: Hao Dong et al. Year: 2020

Summary: This work presents an object-driven approach that uses a parsing model to identify objects mentioned in text, then guides the GAN to generate images with coherent object structures and layouts. The model integrates object-level attention with bounding box prediction, enabling it to generate complex scenes containing multiple entities. Compared to traditional GANs, this model exhibits improved structural coherence and spatial awareness, making it particularly useful for scene-based generation tasks.

2.4Paper Name: AttnGAN: Fine-Grained Text to Image Generation with Attentional Generative Adversarial Networks

Authors: Tao Xu et al. Year: 2018

Summary: AttnGAN introduced an attention mechanism that enables the model to focus on specific words when generating different parts of the image, allowing for fine-grained text-to-image synthesis. The model utilizes a multi-stage refinement process, where images are progressively improved and sharpened at each stage using both global sentence-level and word-level features. Experiments showed that AttnGAN outperformed earlier GAN-based models in both inception score and human evaluation, making it a major advancement in generating complex and detailed images from natural language.

3. METHODOLOGY

The proposed system for AI-based text-to-image generation adopts an encoder-decoder architecture, where the input is a natural language description and the output is a synthetic image generated using a Convolutional Neural Network (CNN). This model aims to capture the semantic features of the text and translate them into visual representations through a structured generation process.

3.1 EXISTING SYSTEM:

3.1.1 Text Processing and Embedding Module:

• Text Preprocessing:

• Input descriptions are tokenized and cleaned to remove stop words, punctuation, and special characters.

• Words are lowercased and stemmed/lemmatized to unify variations.

• Word Embedding:

• Each tokenized word is converted into a vector using pre-trained embedding models like GloVe, Word2Vec, or BERT.

• The word vectors are aggregated (e.g., averaged or passed through an RNN/LSTM) to form a single feature vector that captures the overall semantic meaning of the sentence.

• Dimensionality Reduction:

• The embedded vector is reduced to a manageable size using a fully connected (dense) layer before being passed to the generator.

3.1.2 CNN-Based Image Generator (Decoder):

The CNN model is designed to accept the text feature vector as input and produce an image in pixel form. It typically involves:

• Dense Layer Expansion: The text vector is expanded to a higher dimensional space to seed the image generation process.

• **Reshape Layer:** The vector is reshaped into a small feature map (e.g., 8×8×128), serving as the starting point for up sampling.

• **Transposed Convolutional Layers:** Also known as deconvolution layers, these up sample the feature maps to increase resolution (e.g., from 8×8 to 256×256).

• Batch Normalization and Activation (ReLU/Tanh): Used after each convolution to stabilize training and enhance non-linearity.

• **Final Output Layer**: A final convolutional layer with tanh activation is used to produce a 3-channel RGB image scaled between [-1, 1].

3.1.3 Training and Optimization Pipeline:

• Dataset:

• Uses paired datasets like MS-COCO, Oxford-102 Flowers, or CUB-200, which include both images and corresponding captions.

• Loss Functions:

• **Pixel-Level Loss (MSE or L1):** Measures pixel-wise difference between the generated and ground truth images.

• **Perceptual Loss:** Compares high-level features from intermediate layers of a pre-trained network (e.g., VGG 16).

• **Text-Image Matching Loss:** Encourages semantic consistency between the input text and the generated image using a discriminator or cosine similarity.

• Optimization:

• The model is trained using Adam optimizer with learning rate scheduling and gradient clipping to prevent instability.

• Epochs and Batching:

• Training is performed over multiple epochs with batch sizes ranging from 16 to 64 depending on GPU capacity.

3.1.4 Image Refinement and Output:

• Post-processing:

• The raw output image may be passed through a super-resolution module or denoising CNN to improve quality.

• Result Evaluation:

• Generated images are evaluated using Inception Score (IS), FID score, and human assessment for semantic alignment.

• User Output:

• The generated image is displayed or saved, optionally allowing users to adjust resolution or style parameters.

3.1.5 User Interface:



Fig 3.1: Architecture of AI based Text to Image Generation

3.2 PROPOSED SYSTEM:

We propose a CNN-based text-to-image generation model that follows an encoder-decoder structure.

• **Text Encoder:** Converts text input into fixed-length feature vectors using word embeddings (e.g., Word2Vec or BERT).

• CNN Generator (Decoder): Transforms these text embeddings into images using stacked convolutional layers with up sampling.

• Modules:

- Text Preprocessing: Tokenize and normalize text input.
- Embedding Layer: Map words to dense vectors representing semantic meaning.
- CNN Generator: Accepts concatenated or transformed text features and produces image matrices.

• Loss Functions: Includes reconstruction loss (MSE) and perceptual loss to align image realism and text relevance.

4. APPLICATIONS OF AI BASED TEXT TO IMAGE GENERATION

• Creative Art and Design: Enables artists and designers to instantly visualize ideas from textual descriptions, aiding rapid prototyping and concept development.

• Accessibility Tools for the Visually Impaired: Converts text into images to help visually impaired users interpret descriptions through assistive technologies or tactile graphics.

• E-Commerce Visualization: Generates product visuals from descriptions, allowing dynamic previews and aiding in product design and user engagement.

• Educational and Training Content: Transforms academic or instructional text into visual aids, supporting enhanced learning through image-based understanding.

• Entertainment and Gaming: Creates characters, scenes, or assets from narrative input, accelerating content generation in games and storytelling.

• **Personalized Avatar and Media Content Creation:** Produces avatars or icons from brief descriptions, enabling personalized content in social media and virtual platforms.

5. CONCLUSION

In conclusion, the AI-Based Text to Image Generation system using CNN demonstrates how deep learning can effectively bridge the gap between language and vision. By translating semantic information from text into coherent visual representations, the system supports innovative applications across art, education, e-commerce, and accessibility. The proposed architecture leverages CNNs for efficient image generation, achieving a balance between quality, interpretability, and computational performance. Although challenges remain in handling abstract text, context ambiguity, and generating high-resolution details, the results show that CNN-based models can reliably generate contextually accurate images for a wide range of descriptive inputs. Future improvements may include integrating attention mechanisms, GANs, or diffusion models for enhanced realism and versatility.

6. ACKNOWLEDGEMENTS

We would like to express our sincere gratitude to Prof. P. B. Vikhe For their continuous support and guidance throughout the project. We also extend our thanks to Pravara Rural Engineering College, Loni for providing us with the necessary resources and infrastructure to carry out this research. Finally, we appreciate the valuable feedback from our project coordinator Prof. P. B. Vikhe which helped improve the quality of this work.

7. REFERENCES

[1]. Ramesh, A., Pavlov, M., Goh, G., et al. (2021). "Zero-Shot Text-to-Image Generation." Proceedings of the 38th International Conference on Machine Learning, 139, 8821-8831. DOI: 10.5555/3495724.3495860.

[2]. Tandjoura, M., & Ghods, M. (2022). "Generative Models for Text-to-Image Synthesis: A Review." Journal of Computer Vision and Image Understanding, 214, 103309. DOI: 10.1016/j.jcviu.2022.103309.

[3]. Ding, Y., & Yu, Z. (2022). "Text to Image Generation with Deep Learning: A Compre- hensive Survey." ACM Computing Surveys, 54(3), Article 52. DOI: 10.1145/3482247.

[4]. Li et al. (2024) – "Efficient Text-to-Image CNNs for Edge Devices" (ACM MM 2024).