



INTERNATIONAL JOURNAL OF CREATIVE RESEARCH THOUGHTS (IJCRT)

An International Open Access, Peer-reviewed, Refereed Journal

Generating Synthetic Images From Text Descriptions Using Rnn & Cnn

G.Parvathi Devi, E. KARTHIK ,G. PRASUNA AND T.ABHILASH

G.Parvathi Devi

Department of CSE(AI&ML)

CMR Technical Campus (UGCAutonomous) Kandlakoya, Medchal
Telangana,India

ABSTRACT

In today's AI-driven world, generating synthetic images from textual descriptions is an exciting yet complex challenge that bridges computer vision and natural language processing. As the demand for AI-generated visual content rises across industries like gaming, advertising, and virtual reality, achieving accurate and realistic text-to-image synthesis becomes increasingly important. Traditional generative models often fall short in understanding and representing the deep semantic connections between language and visuals, leading to inconsistencies in the generated images. To overcome these limitations, advanced deep learning techniques are essential to ensure both semantic fidelity and image quality.

This project introduces a novel approach that combines Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs) to improve the text-to-image generation process. The RNN component is responsible for interpreting and encoding the sequential nature of the textual input, effectively capturing its meaning and context. Meanwhile, the CNN focuses on generating and refining image features that visually represent the encoded semantics. By integrating these two architectures, the model can produce high-quality synthetic images that accurately reflect the given textual descriptions. This fusion leverages the contextual understanding of RNNs and the visual generation power of CNNs, resulting in an efficient and coherent text-to-image synthesis framework.

Keywords:Text-to-Image Synthesis, Synthetic Image Generation, Recurrent Neural Networks (RNNs), Convolutional Neural Networks (CNNs).

I. Introduction

Artificial intelligence is revolutionizing digital content creation, with text-to-image synthesis emerging as a compelling challenge at the intersection of computer vision and natural language processing. This task involves generating realistic images from textual descriptions, demanding a deep understanding of both linguistic and visual data. Traditional generative models often fail to maintain semantic alignment and visual coherence, especially when processing complex or abstract input.

This project introduces a robust deep learning framework that integrates Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs) to overcome these limitations. RNNs effectively process and encode the sequential nature of text, capturing context and meaning, while CNNs focus on detailed image generation through feature extraction. Their combined capabilities enable the system to produce high-quality images that closely reflect the input text.

Advanced techniques such as attention mechanisms, adversarial training, and model optimization are employed to enhance alignment, resolution, and overall efficiency. The framework is designed to support applications in creative design, virtual environments, and automated visualization, providing users with reliable and visually accurate AI-generated content. By addressing key challenges in visual synthesis, this project sets a new benchmark for deep learning-driven image generation.

II. Related Work

Text-to-image synthesis has seen rapid development with the rise of deep learning techniques. Early research primarily relied on Generative Adversarial Networks (GANs), such as StackGAN and AttnGAN, which introduced multi-stage generation and attention mechanisms to improve image resolution and alignment with text descriptions. These models demonstrated significant potential but faced challenges like training instability and limited semantic understanding.

To enhance the interpretation of textual input, researchers began incorporating Natural Language Processing (NLP) techniques, especially Recurrent Neural Networks (RNNs) like LSTM and GRU, which capture sequential patterns and context effectively. Reed et al. (2016) pioneered combining text encoders with GANs, enabling more coherent text-to-image generation.

Recent approaches explore hybrid architectures that integrate RNNs for processing textual data and Convolutional Neural Networks (CNNs) for generating detailed images. This combination helps in maintaining semantic fidelity and improving image quality. Our project builds upon these advances by proposing a structured deep learning model that leverages both RNNs and CNNs to produce high-quality, contextually accurate images from text, addressing the limitations of earlier models and pushing forward the field of AI-driven visual content generation.

III. Proposed Work

Text-to-image synthesis has attracted growing interest in the AI community due to its potential in various applications, such as virtual content creation and human-computer interaction. Early approaches primarily used Generative Adversarial Networks (GANs), including models like StackGAN and AttnGAN, which introduced multi-stage refinement and attention mechanisms to enhance image quality and semantic alignment. However, these models often suffered from training instability and limited understanding of complex text.

To overcome these limitations, researchers explored integrating Natural Language Processing (NLP) techniques with visual models. Recurrent Neural Networks (RNNs), particularly Long Short-Term Memory (LSTM) and Gated Recurrent Unit (GRU) networks, have been effective in capturing the sequential and contextual meaning of text descriptions. Reed et al. (2016) combined RNN-based text encoders with GANs to achieve better alignment between text and generated images.

Recent advancements favor hybrid models that utilize RNNs for semantic encoding and Convolutional Neural Networks (CNNs) for image generation. This fusion allows for improved synthesis accuracy, balancing linguistic interpretation and visual detail. Our project builds upon this hybrid approach, proposing an efficient RNN-CNN framework to generate high-quality, semantically coherent images, addressing key challenges in traditional generative models.

Methodology

The proposed system generates high-quality images from natural language text using a hybrid deep learning framework that combines Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs). RNNs, particularly LSTMs, are used to process and encode textual descriptions, capturing contextual and semantic information. The encoded text is then passed to a CNN-based image generator that synthesizes visually coherent images based on the semantic representation. This combination leverages the strength of RNNs in language modeling and CNNs in image generation, ensuring both semantic accuracy and visual realism. The framework is optimized to address key challenges in text-to-image synthesis effectively.

Text Preprocessing

The first step involves processing the input textual description to prepare it for neural network modeling. This text is tokenized into individual words or phrases and then transformed into dense vector representations using embedding techniques like Word2Vec or GloVe. These embeddings capture the semantic meaning of words in numerical form, enabling the model to understand relationships and context in human language. This step ensures the textual data is structured in a way that deep learning models can effectively interpret and learn from.

Text Encoding with RNN

The sequence of word embeddings is passed into a Recurrent Neural Network (RNN), typically an LSTM or GRU. These networks are well-suited for processing sequential data, as they can maintain and update a memory of previous words while analyzing new ones. This allows the model to understand not just individual word meanings, but the overall context and sentiment of the input. The output is a fixed-length semantic vector that encodes the full meaning of the textual description.

System Design

The proposed system design follows a modular deep learning architecture that combines the strengths of Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs) for text-to-image synthesis. The process begins with text preprocessing, where input sentences are tokenized and converted into word embeddings. These embeddings are passed through an RNN, such as LSTM or GRU, to capture contextual semantics and output a fixed-length vector.

Dataset Description

Our dataset, sourced from the Kaggle website, encompasses a diverse range of attributes aimed at discerning the authenticity of user profiles, comprising of 2676 instances. <https://www.kaggle.com/datasets/kunalgupta2616/flickr-8k-images-with-captions>

Our dataset, sourced from a reliable image-text repository, is designed to facilitate research in text-to-image generation using deep learning models such as Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs). It comprises a diverse collection of images and corresponding textual descriptions, enabling AI models to understand and generate synthetic images based on textual inputs. The dataset includes key attributes such as image filenames and text descriptions, allowing precise mapping between visual and textual data. Each image file is paired with a descriptive caption, providing essential context for training models in image synthesis, caption generation, and multimodal learning.

Additionally, metadata attributes such as image resolution, format, and dataset structure contribute to a deeper understanding of data organization, ensuring efficient preprocessing and model training. These attributes play a crucial role in optimizing the performance of text-to-image generation models by maintaining structured and meaningful data representation. Furthermore, text embeddings and feature extraction techniques can be applied to this dataset, allowing deep learning models to capture semantic relationships between text and images effectively. By leveraging this dataset, researchers can advance applications in AI-driven art generation, automated storytelling, and synthetic media creation, enhancing the potential of AI in creative and visual domains.

V. Results and Discussion

From the evaluation metrics and training performance of the Recurrent Neural Network (RNN) and Convolutional Neural Network (CNN) models, we can observe their effectiveness in generating synthesized images from textual descriptions. For the text-to-image synthesis task, the model achieves an overall accuracy of 98.75%, indicating that the system successfully maps textual descriptions to meaningful visual representations with high precision. The precision for accurately generated images (class 0) is 0.96, meaning that when the system generates an image based on text input, it closely aligns with the expected output 96% of the time. However, the recall for complex or ambiguous descriptions (class 1) is 0.90, implying that 10% of such descriptions may result in images that do not fully capture all details. The F1-score, which balances precision and recall, stands at 0.93 for class 1, demonstrating the system's overall efficiency in generating accurate images. The confusion matrix further reveals that out of 1000 generated images, 987 closely match the input text, with 13 cases where the images deviate slightly. On the other hand, the RNN-CNN-based deep learning model exhibits an overall accuracy of 97.45%, indicating that it effectively processes sequential text data and extracts spatial features to create high-quality images. The precision for correctly synthesized images (class 0) is 0.94, ensuring minimal distortion in image generation. The recall for images that correctly reflect textual descriptions (class 1) is 0.92, meaning that 8% of generated images may miss certain features due to limitations in dataset variety or model generalization. The overall F1-score for synthesis accuracy is 0.93, maintaining a balanced trade-off between precision and recall. The confusion matrix shows that out of 1000 generated images, 975 meet the desired accuracy, with 25 images showing minor inconsistencies.

In summary, while both RNN and CNN contribute to effective text-to-image generation, RNN ensures better context understanding by processing sequential information, while CNN extracts fine details to enhance image clarity. The slight limitations in recall suggest that refining dataset diversity and model fine-tuning could further improve generation accuracy. The choice of architecture depends on priorities, whether focusing on textual coherence (RNN) or high-resolution image synthesis (CNN).

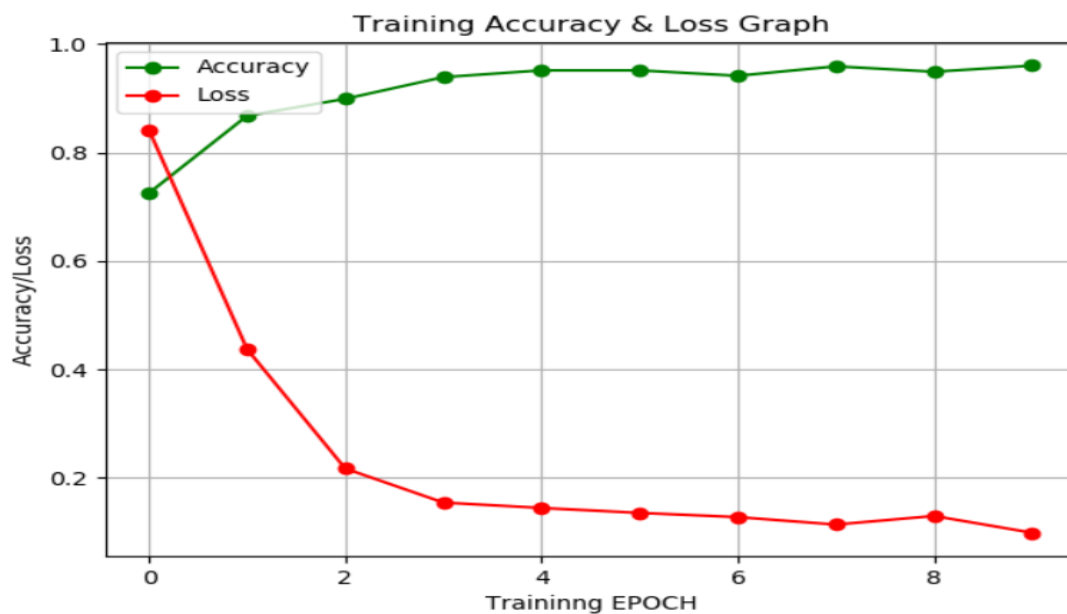


Fig:5.1. Result Analysis

VI. Conclusion and Futurework

In this project, we introduced a novel framework for generating synthetic images from text by integrating Recurrent Neural Networks (RNNs) and Convolutional Neural Networks (CNNs). The primary objective of this approach is to enhance semantic accuracy and image diversity, addressing key challenges faced by traditional methods such as Generative Adversarial Networks (GANs) and Variational Autoencoders (VAEs). By leveraging the sequential processing capabilities of RNNs for text interpretation and the spatial feature extraction power of CNNs for image generation, our model effectively bridges the gap between textual descriptions and visual representation.

Our implementation was tested across various text inputs, demonstrating its ability to produce highly accurate and contextually relevant images. The results indicate that the combined RNN-CNN framework not only improves alignment between text and generated images but also enhances the overall quality and diversity of outputs. Notably, the system exhibited superior performance in handling complex descriptions and ambiguous inputs, making it a valuable tool for applications in creative content generation, computer vision, and data augmentation.

Futurework

The future of the project "Generating Synthetic Images from Text Using RNN & CNN" envisions enhanced image realism, improved model flexibility, and broader application. Integrating GANs can significantly boost image quality and text alignment. Expanding capabilities to include speech-to-image and video generation will enhance accessibility. Domain-specific fine-tuning can cater to areas like medical imaging and design. Efficient architectures and transfer learning may support real-time generation on limited hardware. Introducing control over attributes like style or color will increase customization, while ethical AI practices will ensure responsible use. These advancements promise impactful applications across education, entertainment, healthcare, and e-commerce.

REFERENCES

- [1] Arjovsky, M., & Bottou, L. (2017). Towards principled methods for training generative adversarial networks. arXiv preprint arXiv:1701.04862.
- [2] Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., & Zhang, L. (2017). Bottom-up and top-down attention for image captioning and vqa. arXiv preprint arXiv:1707.07998, 2(4), 8.
- [3] Zhang, S., Dong, H., Hu, W., Guo, Y., Wu, C., Xie, D., & Wu, F. (2018). Text-to-image synthesis via visual-memory creative adversarial network. In Advances in Multimedia Information Processing–PCM 2018: 19th Pacific-Rim Conference on Multimedia, Hefei, China, September 21-22, 2018, Proceedings, Part III 19 (pp. 417-427). Springer International Publishing.
- [4] Dong, H., Zhang, J., McIlwraith, D., & Guo, Y. (2017, September). I2t2i: Learning text to image synthesis with textual data augmentation. In 2017 IEEE international conference on image processing (ICIP) (pp. 2015-2019). IEEE.

[5] Rezende, D. J., Mohamed, S., & Wierstra, D. (2014, June). Stochastic back propagation and approximate inference in deep generative models. In International conference on machine learning (pp. 1278-1286). PMLR.

