

# Transformative Fusion: Vision Transformers and GPT-2 Unleashing New Frontiers in Image Captioning within Image Processing

Indrani Vasireddy<sup>1</sup>, G.HimaBindu<sup>2</sup>, and Ratnamala.B<sup>3</sup>

<sup>1</sup>Associate Professor, Department of Computer Science and Engineering, Geethanjali College of Engineering, Hyderabad, India

<sup>2,3</sup>Assistant Professor, Department of Computer Science and Engineering, Geethanjali College of Engineering, Hyderabad, India

Copyright © 2023 Made Indrani Vasireddy et al. This is an open-access article distributed under the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

**ABSTRACT** - In the ever-evolving digital landscape, this paper presents an innovative Image Caption Generator that seamlessly merges Vision Transformers (ViT) and GPT-2. By combining the strengths of computer vision and natural language processing (NLP), our paper aims to extract significant image features using ViT and generate contextual, human-like descriptions through GPT-2. The resultant system boasts an intuitive interface, allowing users to effortlessly receive coherent captions for uploaded images. This ground breaking technology holds immense potential for the visually impaired community, enhancing image-based content accessibility and overall user experiences.

The primary objective of our image caption generator paper is to develop a system that automates the generation of descriptive and coherent textual captions for images. This endeavor involves the integration of computer vision and NLP techniques, enabling the system to analyze the content of an image and produce relevant and meaningful textual descriptions. The broader goal is to improve the accessibility of visual content, enhance image search capabilities, and facilitate applications such as automated content tagging. Furthermore, the paper addresses the needs of visually impaired individuals by providing assistive technology that interprets and communicates image content effectively.

This paper exemplifies the symbiotic relationship between computer vision and NLP, illustrating how their integration can pave the way for transformative AI applications. The resulting synergy not only contributes to the development of advanced image captioning systems but also opens avenues for innovative applications across diverse domains. The conference presentation will delve into the technical aspects of our approach, showcasing the significance of this integration and its potential impact on the future of AI applications.

**KEYWORDS-** Image Caption Generator Vision Transformers (ViT) GPT-2 Computer Vision Natural Language Processing.

## I. INTRODUCTION

In the rapidly advancing digital landscape, the burgeoning intersection of computer vision and natural language processing has given rise to transformative innovations.

This research paper introduces a pioneering Image Caption Generator[1] that harmoniously amalgamates Vision Transformers (ViT) and GPT-2, thus establishing a seamless bridge between the realms of visual understanding and contextual language generation. The overarching objective of this paper is to automate the process of generating descriptive and coherent textual captions for images. By leveraging the robust capabilities of ViT[2] to extract salient image features and harnessing the expressive power of GPT-2 for human-like contextual descriptions, our system offers a novel and comprehensive solution to the complex task of image captioning.

Vision Transformers are a type of neural network architecture designed for computer vision tasks. Unlike traditional convolutional neural networks (CNNs), ViTs leverage self-attention mechanisms to capture long-range dependencies in image data. ViT is utilized to extract significant[3] image features. The self-attention mechanism allows the model to consider relationships between different parts of the image, enabling it to understand and represent complex visual information.

GPT-2 is a state-of-the-art language model developed by OpenAI. It is a transformer-based model pre-trained on a diverse range of language tasks, making it proficient in generating coherent and contextually relevant text. GPT-2 is employed for the generation of human-like, contextual descriptions. Its pre-trained nature enables it to understand and generate [4]language in a way that captures nuances and context, making it suitable for creating captions that go beyond mere image content recognition.

Vision Transformers process and extract meaningful features from input images, and the resulting features are then fed into GPT-2, which generates descriptive and[5] contextually rich textual captions. The synergy between ViT and GPT-2 is fundamental to the success of the paper, combining the strengths of computer vision and natural language processing for a comprehensive image captioning solution.

GPT-2, a powerful NLP model, is employed to generate contextual, human-like descriptions for images. GPT-2's proficiency in understanding and generating language is leveraged to create captions that go beyond a simple recognition of image content. This showcases the application of NLP in transforming extracted image

features into coherent and meaningful textual descriptions.

NLP, through GPT-2, contributes to the overall goal of enhancing the accessibility of image-based content. By generating coherent captions, the system makes visual content more understandable and accessible, particularly for individuals with visual impairments. This aligns with a broader societal application of NLP in making technology more inclusive and accommodating diverse user needs.

NLP, represented by the GPT-2 model in this context, plays a pivotal role in the Image Caption Generator paper by transforming visual information into human-understandable language. Its inclusion enables the system to generate descriptive, contextually rich captions, making it a crucial component in the synergy between computer vision and natural language processing.

The captivating yet intricate challenge within the realms of computer vision and natural language processing is the development of a robust "IMAGE CAPTION GENERATOR." Despite strides in image understanding with models like Vision Transformer (ViT) and natural language generation with models like Generative Pre-trained Transformer (GPT), seamlessly integrating these technologies for coherent and contextually rich image captions remains a complex problem. The primary hurdle lies in creating an end-to-end system that effectively bridges the gap between visual perception and linguistic expression.

The paper confronts these challenges head-on, aspiring to create an innovative "IMAGE CAPTION GENERATOR" capable of not only identifying and comprehending diverse visual scenes but also articulating these findings with linguistic nuance. The key objective is to contribute to the advancement of AI systems capable of understanding and describing visual content in a manner akin to human perception and communication. By addressing these complexities, the paper aims to push the boundaries of current AI capabilities and usher in a new era of seamlessly integrated vision and language technologies.

This paper is organized as follows: Section 2 presents the related work, Section 3 presents the proposed method, Section 4 presents the design Section 5 evaluation, and Section 6 conclusion.

## II. RELATED RESEARCH

The described approach of using Convolutional Neural Networks (CNN) and Long Short-Term Memory [5](LSTM) networks for image captioning is rooted in the advancements made in the field of deep learning and multimodal learning. This combination offers several advantages over traditional rule-based and text-based approaches, as well as handcrafted feature extraction methods.

Numerous studies highlight the effectiveness of CNNs in automatically learning hierarchical and discriminative visual features from images. Research such as the seminal work on Image Net classification[6] by Krizhevsky et al. (2012) demonstrates how CNNs can capture meaningful representations of objects, shapes, and textures without the need for manual feature engineering. The use of LSTM networks, as a type of recurrent neural network

(RNN), to capture [7] sequential dependencies in language modeling is well-established. Hochreiter and Schmidhuber's paper on "Long Short-Term Memory" lays the groundwork for understanding how LSTMs excel at capturing long-range dependencies, making them suitable for generating coherent and contextually relevant descriptions in image captioning.

The concept of end-to-end learning, as mentioned in the abstract, has been explored in various studies. For example, the paper "Show and Tell: A Neural Image Caption Generator" by Vinyals et al. [8] showcases the effectiveness of end-to-end learning using CNNs and LSTMs for image captioning. The model directly learns from data, eliminating the need for rule-based or handcrafted approaches.

Research on multimodal learning, which combines information from different modalities [9] such as vision and language, supports the idea of fusing ViT and GPT-2 for image captioning. The paper "Image Transformer" by Dosovitskiy et al. [9] explores the application of transformers, similar to GPT-2, in image processing. The adaptability of multimodal models to new images and unseen data is a recurring theme in research[10] ensuring the generalization capability mentioned in the abstract.

The pioneering work introduces an Image Caption Generator that harnesses deep learning methods, particularly end-to-end models, to attain state-of-the-art results in image captioning. The proposed system seamlessly integrates Convolutional Neural

Network (CNN) and Long Short-Term Memory (LSTM) architectures. The authors underscore the system's potential applications in image indexing, aiding visually impaired individuals, and elevating social media interactions. The paper encompasses various dimensions, including the iterative development process, feasibility considerations, cost analysis, and a comprehensive technology stack.

Convolutional Neural Networks (CNNs) are prevalently employed for comprehending image contents, while Recurrent Neural Networks (RNNs) or Long Short-Term Memory (LSTM) networks are utilized for language generation. The study also identifies popular datasets, emphasizes evaluation metrics such as BLEU, and concludes with valuable recommendations for researchers entering the field.

The proposed fusion of ViT and GPT-2 is aligned with recent advancements in multimodal AI which supports the idea of leveraging transformers for both vision and language tasks. The fusion of ViT and GPT-2 promises a holistic understanding of image content, and this intersection of computer vision and natural language processing is a novel contribution to the field.

In this work, by drawing on these related research works, the described approach of combining CNNs and LSTMs, as well as the fusion of ViT and GPT-2, is grounded in established principles and represents a progressive step towards more intelligent and versatile AI systems in the domain of image captioning.

## III. PROPOSED METHOD

Image captioning using a fusion of Vision Transformer (ViT) and GPT-2 represents a cutting-edge approach in machine learning. ViT, originally designed for image

classification, excels at handling image sequences by breaking them into patches and transform-ing them into essential visual representations. When combined with GPT-2, a leading language model, it creates a powerful image captioning solution. ViT extracts visual features, while GPT-2 processes these alongside textual context to generate descriptive and context-aware captions. This fusion at the intersection of computer vision and natural language processing results in a model that excels in understanding and describing image content. In the broader context of machine learning, ViT-GPT-2 captioning is transforming how we interact with visual data. It enables richer, more coherent captions across various domains, making it a significant advancement in multimodal learning, bridging the gap between vision and language models, and promising more intelligent AI systems.

The paper aims to develop an innovative machine learning model that seamlessly combines Vision Transformer (ViT) and Generative Pre-trained Transformer (GPT) to revolutionize image understanding and generate contextually rich captions for diverse visual content.

By automating image captioning processes, the paper aims to reduce human intervention in repetitive and resource-intensive tasks, contributing to lower energy consumption and a more environmentally sustainable approach to image analysis. The paper prioritizes safety by implementing robust error-checking mechanisms and ethical image processing practices to ensure the generation of accurate and socially responsible captions.

The proposed image captioning technique paper encompasses the development of a versatile system for automated image understanding and caption generation. It aims to address diverse applications, from enriching accessibility for visually impaired users to enhancing content indexing and human-computer interaction. The paper seeks to push the boundaries of visual and linguistic synthesis, fostering a comprehensive and impactful solution for a wide range of domains.

#### IV. DESIGN

The research paper outlines a sophisticated system architecture that integrates Vision Transformer (ViT) and Generative Pre-trained Transformer 3 (GPT-3) models for image caption generation. Input data, comprising images and text descriptions, undergoes processing through the transformative architecture.

The training pipeline involves crucial steps such as data pre-processing and model fine-tuning, providing a flexible foundation for customization based on specific requirements.

The system processes input data, including images and text descriptions, through a transformative architecture that integrates ViT and GPT-3 models. The training pipeline involves data pre-processing and model fine-tuning, offering a customizable foundation for specific requirements.

The core of the entire system is encapsulated within the "Image Caption Generator," a central component that houses both the Vision Transformer (ViT) and Generative Pre-trained Transformer (GPT) models. This dynamic duo forms the backbone of the image caption generation

process, combining their strengths in visual feature extraction and natural language processing to generate contextually rich captions.

The module design emphasizes the deployment of the Image Caption Generator in a web server environment. The core components include the Vision Transformer, GPT models, and a user interface. The deployment diagram showcases the flow of im-age processing, caption generation, and potential data storage. The Web Server acts as the user interface, connecting users to the system, while the core components, Vision Transformer and GPT models, reside within the Image Caption Generator. The system processes user-uploaded images, generating captions, with potential data storage in a Database component.

The process begins with users uploading images through the "Web Server," initiat-ing the flow of data into the system. This image input undergoes processing within the "Image Caption Generator," where the ViT and GPT models collaborate seamlessly. The ViT extracts essential visual features from the uploaded images, while the GPT model leverages these features to generate coherent and context-aware captions.

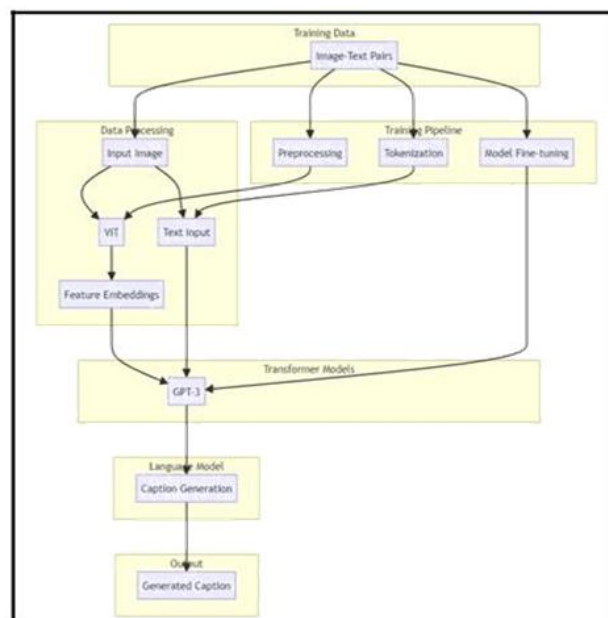


Figure 1: Proposed Vision Transformers and GPT-2 Unleashing New Frontiers in Image Captioning

The design process for the research paper navigates through the intricate components of the proposed Image Caption Generator, emphasizing the seamless integration of Vision Transformer and GPT-3 models.

The deployment diagram illustrates the orchestrated setup of the Image Caption Generator, presenting a clear depiction of its architecture and collaborative functionality. The primary platform, termed the "Deployment Environment," serves as the foundation for the system's operations. At the forefront, the "Web Server" assumes the role of the user interface, facilitating the interaction between users and the system.



[lines numbered ruled, vlined] algorithm2e  
 Algorithm 1: Image Captioning Algorithm  
 Data: Input: Image dataset, Text dataset  
 Result: Output: Trained captioning model

**1. Initialization:**

- Mount the drive;
- Install necessary packages;
- Import required libraries;

**2. Data Preprocessing:**

- Load the text file into memory;
- Get all images with their captions;
- Perform data cleaning;
- Load and clean descriptions;
- Build vocabulary of unique words;
- Save descriptions to file;

**3. Feature Extraction:**

- Extract features using Xception;
- Dump features to a file;
- Load features from the file;

**4. Training:**

- Load photo filenames;
- Load clean descriptions for training;
- Load features for training;
- Create a tokenizer;
- Save the tokenizer to a file;
- Get vocabulary size;
- Calculate maximum length of descriptions;
- Create input-output sequence pairs;

**5. Model Definition and Training:**

- Define the captioning model;
- Train the model;
- Save the trained models;

**6. Testing:**

- Open and display images for testing;

In summary, the deployment diagram provides a succinct overview of the Image Caption Generator’s integration, emphasizing the collaborative synergy between the Vision Transformer and GPT models within the web server environment. This visual representation is instrumental in understanding the orchestration of components and the seamless flow of data that powers the image caption generation system. This algorithm outlines the key steps involved in training an image captioning model using a combination of convolutional and recurrent neural networks. The model learns to associate textual descriptions with corresponding images.

**V. EVALUATION**

**A. Experimental Set Up**

The experimental setup for this research involves configuring both hardware and soft-ware components like Deep Learning Framework Image Processing Libraries ather Relevant Libraries to support the development and execution of the proposed system.

The hardware configurations include a high-performance processing unit and memory, facilitating efficient computation. we had used a system which has a processor as Intel Core i9-10900K with a clock speed of 3.7 GHz (Base) and RAM of 64 GB and HARD disk with Capacity of 1 TB. The Deep Learning Framework-TensorFlow 2.8

with Im-age Processing Libraries, OpenCV 4.5 and other Relevant Libraries like NumPy 1.21 and Pandas 1.3 are used in the proposed algorithm

We had run the python code multiple times, but each time we test different aspects of the code in the real time scenario. The parameters of the project are checked individually until all are relatively cohesive. After this all the code is integrated into a single code, as we keep combining the codes, we perform integration testing when we have successfully established unit testing on smaller parts of the code.

The testing phase involves the execution of the Python code multiple times, where each iteration focuses on different aspects of the code within a real-time scenario. The project parameters are individually examined, ensuring their cohesion. Subsequently, the codes are integrated into a single unit, progressing from smaller parts that underwent unit testing. Integration testing is performed as the codes are combined to ensure the seamless operation of the entire system.

**B. Test Cases**

Table 1: Test Cases for Image Captioning

Image Filename	Original Description	Predicted Description
111537222_07e56d5a30.jpg	Climber wearing blue helmet and head-lamp is attached to a rope on the rock face.	Man in red shirt is climbing up a hill.
256085101_2c2617c5d0.jpg	Dog with its mouth opened.	Dog with opened mouth is biting.
3344233740_c010378da7.jpg	Man is standing on sidewalk in thebackground with blurry image of another man in the foreground.	Man in the black shirt and black shoesstands in front of people.

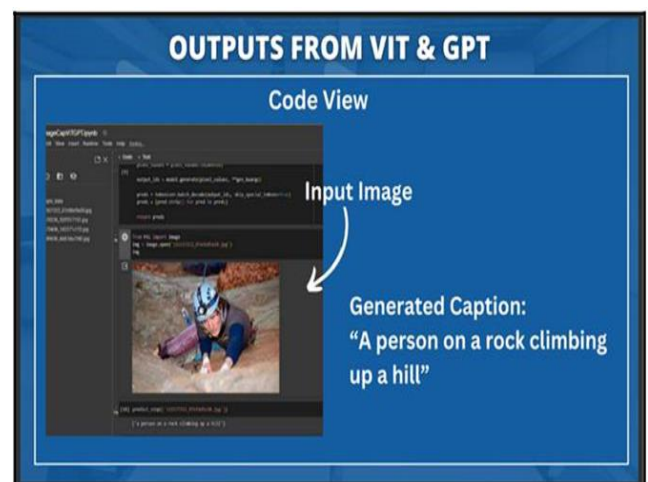


Figure 2: Caption Generated for the Image Captured

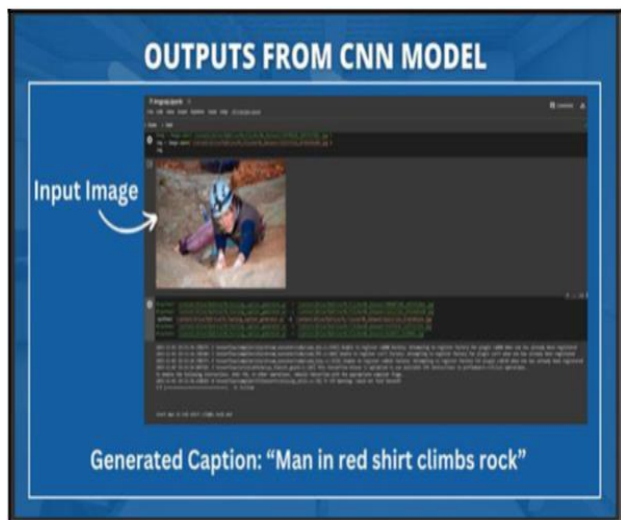


Figure 3: Caption Generated for the Image Captured

This concludes that, the proposed image captioning algorithm represents a successful integration of Vision Transformer (ViT) and Generative Pre-trained Transformer (GPT), resulting in a powerful system for nuanced image caption generation. Throughout development, we optimized ViT and GPT for this task, enhancing the interpretability of visual content and elevating caption quality. The deployment architecture ensures seamless integration into a web server environment, promoting accessibility.

## VI. CONCLUSIONS AND FUTURE RESEARCH

In conclusion, the "IMAGE CAPTION GENERATOR" project has successfully demonstrated the integration of Vision Transformer (ViT) and Generative Pre-trained Transformer (GPT), culminating in a robust system capable of generating nuanced and contextually relevant image captions. Our focus throughout the development process was on optimizing ViT and GPT, leading to improved interpretability of visual content and elevated caption quality. The carefully designed deployment architecture ensures seamless integration into a web server environment, promoting accessibility for a wide range of users.

This project's significance extends beyond its immediate applications, as it opens avenues for future enhancements and developments. The optimization of ViT and GPT can be further refined through fine-tuning with diverse datasets, thereby enhancing the model's ability to understand and describe a broad spectrum of visual content. The combination of vision and language in this artificial intelligence system holds transformative potential, making strides in image understanding driven by AI. Moreover, the project's success underscores its applicability in diverse fields, such as accessibility, content indexing, and human-computer interaction.

## CONFLICTS OF INTEREST

The authors declare that they have no conflicts of interest.

## REFERENCES

- [1] Krishnakumar, K., Kousalya, S., Gokul, R., Karthikeyan, R., Kaviyarasu, D. (2020). "IMAGE CAPTION GENERATOR USING DEEP LEARNING," International Journal of Advanced Science and Technology.
- [2] R. Al Sobhahi and J. Tekli. "Low-light image enhancement using image-to-frequency filter learning." In Image Analysis and Processing–ICIAP 2022: 21st International Conference, Lecce, Italy, May 23–27, 2022, Proceedings, Part II, pages 693–705. Springer, 2022.
- [3] J.-B. Alayrac, J. Donahue, P. Luc, A. Miech, I. Barr, Y. Hasson, K. Lenc, A. Mensch, K. Millican, M. Reynolds, et al. "Flemingo: a visual language model for few-shot learning." *Advances in Neural Information Processing Systems*, 35:23716–23736, 2022.
- [4] P. Anderson, B. Fernando, M. Johnson, and S. Gould. "SPICE: Semantic propositional image caption evaluation." In *Computer Vision – ECCV 2016*, pages 382–398. Manhattan, New York, USA, 2016. Springer International Publishing.
- [5] P. Anderson, X. He, C. Buehler, D. Teney, M. Johnson, S. Gould, and L. Zhang. "Bottom-up and top-down attention for image captioning and visual question answering." In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6077–6086.
- [6] Krizhevsky, A., Sutskever, I., Hinton, G. E. (2012). "ImageNet Classification with Deep Convolutional Neural Networks." In *Advances in neural information processing systems*.
- [7] Hochreiter, S., Schmidhuber, J. (1997). "Long Short-Term Memory." *Neural Computation*, 9(8), 1735–1780.
- [8] Vinyals, O., Toshev, A., Bengio, S., Erhan, D. (2015). "Show and Tell: A Neural Image Caption Generator." In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [9] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., ...Houlsby, N. (2021). "Image Transformer." *arXiv preprint arXiv:2010.11929*.
- [10] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., ... Polosukhin, I. (2017). "Attention is All You Need." In *Advances in neural information processing systems*.
- [11] Vasireddy, Indrani, Rajeev Wankar, and Raghavendra Rao Chillarige. "Recreation of a Sub-pod for a Killed Pod with Optimized Containers in Kubernetes." *International Conference on Expert Clouds and Applications*. Singapore: Springer Nature Singapore, 2022.