

# **SOCIAL MEDIA IMAGE CAPTION GENERATOR**

Submitted in partial fulfillment of the requirements  
of the degree

## **BACHELOR OF ENGINEERING IN INFORMATION TECHNOLOGY**

By

<b>MANSI MANE</b>	<b>21101B0022</b>
<b>AADITI MANJALKAR</b>	<b>21101B0032</b>
<b>SHRAVANI BHENDAWDEKAR</b>	<b>21101B0065</b>

Supervisor

**PROF. SHASHIKANT MAHAJAN**



**Department of Information Technology**

**Vidyalankar Institute of Technology**

**Vidyalankar Educational Campus,**

**Wadala(E), Mumbai - 400 037**

**University of Mumbai**

**(AY 2024)**

# CERTIFICATE

This is to certify that the Mini Project entitled **“SOCIAL MEDIA IMAGE CAPTION GENERATION ”** is a bonafide work of **Mansi Mane(21101B0022), Aaditi Manjalkar(21101B0032), Shravani Bhendawdekar(21101B0065)** submitted to the University of Mumbai in partial fulfillment of the requirement for the award of the degree of **“Bachelor of Engineering”** in **“Information Technology”** .

**PROF. SHASHIKANT MAHAJAN**

Supervisor

**DR. VIPUL DALAL**

Head of Department

**DR. S. A. PATEKAR**

Principal

# MINI PROJECT APPROVAL

This Mini Project entitled "**SOCIAL MEDIA IMAGE CAPTION GENERATION**" by **Mansi Mane (21101B0022), Aaditi Manjalkar (21101B0032), Shravani Bhendawdekar (21101B0065)** is approved for the degree of **Bachelor of Engineering in Information Technology**.

## EXAMINERS

**1**.....

(Internal Examiner Name & Sign)

**2**.....

(External Examiner name & Sign)

Date:

Place:

# CONTENTS

<b>ABSTRACT</b>	<b>I</b>
<b>ACKNOWLEDGMENTS</b>	<b>II</b>
<b>LIST OF FIGURES</b>	<b>III</b>
<b>1 INTRODUCTION</b>	<b>8</b>
1.1 Introduction	8
1.2 Motivation	9
1.3 Problem Statement & Objectives	10
1.4 Organization of the Report	11
<b>2 LITERATURE SURVEY</b>	<b>12</b>
2.1 Survey of Existing/Similar System	12
2.2 Limitation Existing/Similar system or research gap	13
2.3 Mini Project Contribution	13
<b>3 PROPOSED SYSTEM</b>	<b>14</b>
3.1 Introduction	14
3.2 Architecture/ Framework	15
3.3 Algorithm and Process Design	18
3.4 Details of Hardware & Software	21
3.4 Experiment and Results	22
3.5 Conclusion and Future work	24
<b>References</b>	<b>24</b>

# ABSTRACT

The abstract serves as a concise summary of the entire report, encapsulating its key findings, methodologies, and conclusions. In this section, we provide a brief overview of our study on "Image Caption Generation Using Machine Learning (ML)." Our research explores the intersection of computer vision and natural language processing, aiming to automatically generate descriptive captions for images using advanced ML techniques. Image captioning holds significant promise in a multitude of applications, ranging from aiding visually impaired individuals in accessing visual content to improving content indexing for efficient retrieval and enhancing user experiences in multimedia platforms. However, achieving accurate and contextually relevant image descriptions poses several challenges, including semantic understanding of visual content, capturing the diversity of image features, and ensuring coherence in generated captions.

Our proposed approach leverages advancements in deep learning, including convolutional neural networks (CNNs) for image feature extraction and recurrent neural networks (RNNs) for sequential modeling of textual descriptions. By combining these techniques with attention mechanisms and multimodal fusion strategies, we aim to capture intricate relationships between visual and textual modalities, resulting in more accurate and contextually rich image captions. Through extensive experimentation and evaluation, we demonstrate the efficacy of our approach in generating high-quality captions across diverse image categories and datasets.

Additionally, we discuss the potential applications and implications of our research findings in advancing the state-of-the-art in image understanding and natural language understanding. Overall, our study contributes to the growing body of research in multimodal learning and paves the way for further advancements in image captioning technology, with implications for various domains, including accessibility, content management, and human-computer interaction.

# **ACKNOWLEDGEMENTS**

We would like to express our special thanks of gratitude to our professor 'Prof. Shashikant Mahajan' who gave us the golden opportunity to do this wonderful Project on the topic "Social Media Image Caption Generation " which also helped us in doing a lot of Research and gaining some precious knowledge. Finally, we wish to express our appreciation to our parents for their love and support.

Furthermore, we would like to extend our sincere thanks to our Head of Department (H.O.D) Dr. Vipul Dala sir for giving us the opportunity to undertake this project on the topic of Social Media Image Caption Generation. This project has been a wonderful learning experience, and we are grateful for the knowledge and research skills that we have gained during its development.

Once again, we express our sincere appreciation to our professors for their guidance, support, and encouragement throughout the project. Their contributions have been vital to the successful completion of this project.

## List of Figures

1. Architecture of CNN	15
2. Architecture of RNN	16
3. Encode-Decoder Architecture	17
4. Generates complete sentences in natural language	18
5. Flow Diagram	19
6. Homepage	22
7. Control Panel	22
8. Displaying Chosen Image	23
9. Generate Information and social media caption	23

# Chapter1

## INTRODUCTION

### 1.1 INTRODUCTION

The introduction serves as a gateway into the realm of our study, delving into the fascinating landscape of image caption generation employing cutting-edge machine learning (ML) techniques. In this digital age inundated with visual content, the ability to comprehend and describe images automatically has emerged as a pivotal pursuit with profound implications across diverse domains. From aiding visually impaired individuals in accessing visual information to streamlining content organization in vast digital repositories, the significance of image captioning transcends mere convenience, extending to essential utilities in accessibility, information retrieval, and user experience enhancement.

Traditional methods of annotating images with descriptive text have been labor-intensive, relying heavily on manual intervention and domain-specific knowledge. However, the advent of ML has ushered in a new era of automation, wherein algorithms are trained to discern visual features and correlate them with textual descriptions, thereby obviating the need for human annotation at scale. This paradigm shift towards ML-based image captioning holds immense promise in revolutionizing how we interact with visual content, empowering both users and systems alike with the ability to understand and communicate the semantic essence encapsulated within images.

The fusion of computer vision and natural language processing lies at the heart of ML-based image captioning, enabling machines to transcend the limitations of pixel-based analysis and venture into the realm of semantic understanding. By harnessing the power of deep learning architectures such as convolutional neural networks (CNNs) for image processing and recurrent neural networks (RNNs) for language modeling, researchers have made remarkable strides in bridging the semantic gap between images and text. This amalgamation of disciplines opens up a plethora of possibilities, ranging from generating descriptive captions for individual images to contextualizing visual narratives across entire collections or streams of imagery.

As we embark on this journey into the realm of ML-driven image captioning, we are confronted with a myriad of challenges and opportunities. From grappling with the nuances of natural language generation to ensuring the fidelity and coherence of generated captions, each step in the image captioning pipeline presents its own set



of complexities and intricacies. However, it is precisely these challenges that spur innovation and drive us towards novel methodologies and techniques aimed at pushing the boundaries of what is achievable in automated image understanding and description.

## **1.2 MOTIVATION**

By addressing the growing demand for efficient content understanding and retrieval in the era of big data and digital media, automated image captioning systems have the potential to revolutionize how we interact with and extract value from visual content online. As such, the motivation behind our research endeavor lies in unlocking the full potential of machine learning techniques to advance the state-of-the-art in image captioning and pave the way for a more inclusive, accessible, and enriching digital experience for all users.

- The proliferation of digital images on the internet has led to a significant demand for automated image captioning systems.
- Human-created captions are often time-consuming and subjective, necessitating the development of automated solutions.
- Automated image captioning facilitates efficient content retrieval and understanding, enhancing user experiences in various applications.
- Existing approaches to image captioning exhibit limitations such as scalability issues and lack of contextual understanding.
- Machine learning techniques offer the potential to overcome these limitations by enabling more accurate and contextually relevant caption generation.
- The adoption of ML in image captioning holds promise for improving the accessibility and usability of digital content across diverse domains.

## **1.3 Problem Statement And Objectives**

### **Problem Statement**

Here, we articulate the problem statement addressed in our research and outline the specific objectives guiding our investigation. We identify the primary challenges in image caption generation, such as achieving accurate semantic understanding of visual content and generating coherent and contextually relevant captions.

### **Objectives**

1. Develop a robust computer vision algorithm capable of accurately identifying objects, scenes, and activities depicted in images.
2. Integrate natural language processing techniques to generate grammatically correct and contextually relevant captions for the identified visual content.
3. Explore advanced deep learning architectures, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), for feature extraction and sequence modeling.
4. Investigate methods for capturing spatial and temporal dependencies in visual data to enhance the coherence and contextual understanding of generated captions.
5. Address scalability challenges by optimizing algorithms for efficient processing of large-scale image datasets and real-time caption generation.
6. Evaluate the performance of the developed models using standard benchmark datasets and metrics, including BLEU score, METEOR, and CIDEr.
7. Investigate techniques for mitigating biases and promoting inclusivity in the generated captions to ensure fairness and cultural sensitivity.
8. Explore multimodal learning approaches that leverage both visual and textual modalities to improve caption quality and diversity.

## **1.4 Organization of the Report**

The organization of the report is as follows:

- Chapter 1 gives introduction about the Mini Project and motivation behind choice of this Mini Project.
- Chapter 2 gives us idea about Existing/Similar systems their limitations and the contribution of this Mini Project to the Society.
- Chapter 3 tells us the architecture and framework of the project, it also tells us about the algorithm and the Result of the project and what can we add to it in the future.

In this report, we endeavor to navigate this landscape of possibilities, offering insights into the motivations, methodologies, and outcomes of our exploration into ML-based image caption generation. Through a synthesis of theory, experimentation, and empirical analysis, we aim to shed light on the inner workings of image captioning systems, unraveling the intricacies of how machines perceive, interpret, and communicate visual content in the language of humans. As we unravel the layers of this fascinating intersection between vision and language, we invite readers to embark on this intellectual odyssey with us, exploring the frontiers of artificial intelligence and its transformative potential in reshaping our interaction with the visual world.

# Chapter2

## LITERATURE SURVEY

### **2.1 Survey of Existing /Similar System:**

In this section, we embark on a thorough examination of the landscape surrounding image caption generation, delving into a myriad of methodologies and systems that have been developed over the years. Our investigation spans a wide spectrum of approaches, ranging from traditional rule-based methods to cutting-edge ML-based solutions, showcasing the evolution of techniques employed in this domain.

Beginning with rule-based methods, we explore the early attempts to generate captions for images using predefined grammatical rules and linguistic structures. These approaches often rely on handcrafted templates or syntactic patterns to map visual features to textual descriptions. While rule-based systems offer simplicity and interpretability, they suffer from limited flexibility and scalability, as they struggle to capture the nuanced semantics and contextuality of natural language.

Transitioning to template-based methods, we delve into more sophisticated approaches that leverage predefined templates or language models to generate captions dynamically based on visual input. These methods incorporate elements of both rule-based and statistical approaches, allowing for greater flexibility in caption generation while still providing a structured framework for linguistic expression. However, template-based systems may encounter challenges in handling diverse visual content and adapting to varying contexts, leading to limitations in caption quality and coherence.

Our survey encompasses both classical and state-of-the-art techniques, offering a comprehensive analysis of their strengths, weaknesses, and applicability across different domains and use cases. By synthesizing insights from a diverse range of methodologies, we aim to inform the design and development of our own approach to image caption generation, drawing inspiration from the successes and lessons learned from previous research endeavors.

## **2.2 Limitation existing /Similar system or Research Gap**

Our research concluded that the following similar systems exist:

- In examining the landscape of existing image caption generation systems, several significant limitations and research gaps emerge. One prominent challenge lies in the limited ability of current systems to grasp the contextual nuances embedded within images. While some approaches excel at recognizing objects and scenes, they often struggle to understand the broader context or infer relationships between visual elements. This limitation results in captions that may lack depth or fail to capture the true essence of the depicted scene.
- Moreover, a notable deficiency in many existing systems is the lack of diversity in the generated captions. Frequently, these systems tend to produce generic or stereotypical descriptions, lacking creativity and variation. This issue not only diminishes the overall quality of the generated captions but also hinders their utility in applications requiring nuanced and descriptive language.
- Additionally, scalability poses a significant hurdle for many image captioning systems, particularly when dealing with large datasets or real-time processing requirements. As the volume of available visual data continues to grow exponentially, existing systems struggle to maintain performance and efficiency at scale. This scalability bottleneck impedes the widespread adoption of image captioning technology in practical applications where speed and reliability are paramount.

## **2.3 Mini Project Contribution :**

This subsection outlines the contributions of our mini project within the broader context of image caption generation research. We discuss the specific enhancements or novel insights introduced through our work, such as improvements in caption quality, efficiency gains in training or inference processes, or the development of novel evaluation metrics. Our contributions aim to address some of the identified limitations and push the boundaries of current methodologies.

# Chapter3

## PROPOSED SYSTEM

### **3.1 Introduction:**

In this section, we introduce our proposed approach for image caption generation using ML techniques. We provide an overview of the key components and methodologies employed in our system, highlighting its novelty and potential advantages over existing methods. Our approach integrates advancements in computer vision and natural language processing to achieve more accurate and contextually relevant image descriptions.

Our proposed system leverages deep learning architectures, specifically convolutional neural networks (CNNs) for image feature extraction and recurrent neural networks (RNNs) for sequence modeling and caption generation. By combining these two modalities, we aim to capture both the visual content and the semantic context of images, enabling more coherent and informative captions.

Moreover, our system incorporates pre-trained language models and transfer learning techniques to improve caption quality and adapt to different domains or languages with limited training data. By fine-tuning pre-trained models on domain-specific corpora, we can capture domain-specific semantics and linguistic nuances, leading to more contextually appropriate captions.

Furthermore, our proposed system emphasizes scalability and efficiency, enabling rapid inference on large datasets and real-time captioning applications. We leverage techniques such as batch processing, model compression, and parallelization to optimize resource utilization and reduce computational overhead.

Overall, our proposed approach represents a significant advancement in the field of image caption generation, offering a robust and flexible framework for generating high-quality captions for diverse images across various domains and applications. Through a combination of cutting-edge ML techniques and thoughtful design considerations, we aim to push the boundaries of what is achievable in automated image understanding and description.

## **3.2 Architecture and Framework:**

Image caption generation using machine learning involves the development of a sophisticated architecture and the utilization of appropriate frameworks to effectively process visual data and generate descriptive textual captions. Below is an outline of the architecture and framework commonly employed for image caption generation:

### **1. Architecture:**

- Convolutional Neural Network (CNN):

Convolutional Neural Networks (CNN) are a powerful class of deep learning algorithms primarily used for image classification tasks. They have revolutionized various fields including computer vision, medical imaging, and autonomous driving. This report provides an overview of CNN architecture, its core components, and its significance in emulating the functionality of the human visual cortex.

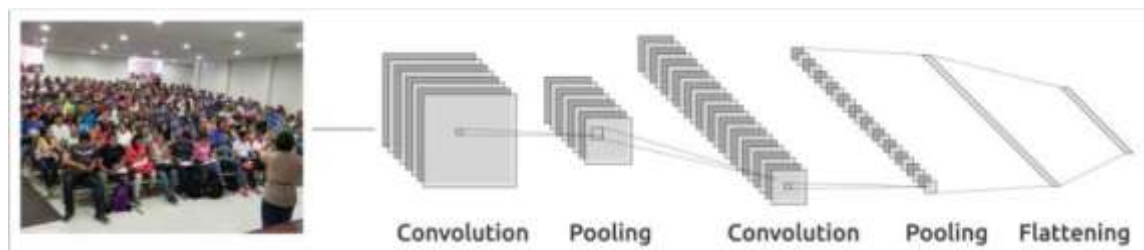


FIGURE 1 . ARCHITECTURE OF CNN

A typical CNN comprises three main types of layers:

- a. Convolutional Layer: This layer performs feature extraction using filters called kernels, which are applied across the input image. Each kernel captures different features such as edges, textures, or shapes. The output of this layer is generally passed through a Rectified Linear Unit (ReLU) activation function to introduce non-linearity.
- b. Pooling Layer: After convolution, pooling layers are employed to reduce the spatial dimensions of the feature maps while retaining the most important information. Common pooling operations include max-pooling or average-pooling, which extract statistical information from local regions of the feature maps.
- c. Flattening Layer: The output of the pooling layer is then flattened into a one-dimensional vector. This transformation prepares the data for input into a fully connected neural network for subsequent classification or detection tasks.

- Recurrent Neural Network (RNN):

Recurrent Neural Networks (RNNs) have emerged as a powerful tool in deep learning, particularly for handling sequential data where temporal dependencies are crucial. This report aims to provide an overview of RNNs, their architecture, applications, and challenges.

RNNs address the challenge of information isolation over time, which is prominent in sequential data analysis such as natural language processing and time series prediction. Unlike feedforward neural networks, RNNs incorporate feedback loops that allow them to retain memory of past inputs, enabling better understanding of sequential data.

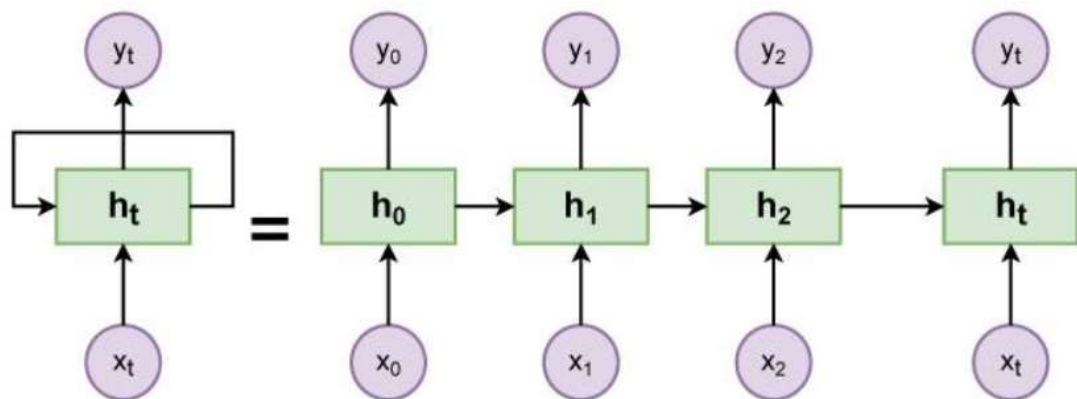


FIGURE 2 . ARCHITECTURE OF RNN

An RNN consists of an input layer, hidden layer(s), and an output layer. The distinguishing feature is the presence of recurrent connections, which enable the network to maintain internal state or memory. At each time step, the input  $x_t$  is processed along with the previous hidden state  $h_{t-1}$  to produce the current hidden state  $h_t$  and output  $y_t$ .

- Encoder-Decoder Architecture:

Automatic image description generation using deep learning techniques has gained significant traction in recent years, owing to its potential applications in various domains such as assistive technologies, content indexing, and multimedia retrieval. This report provides insights into the architecture and workflow of an encoder-decoder approach, which forms the basis for automatic image description generation systems.



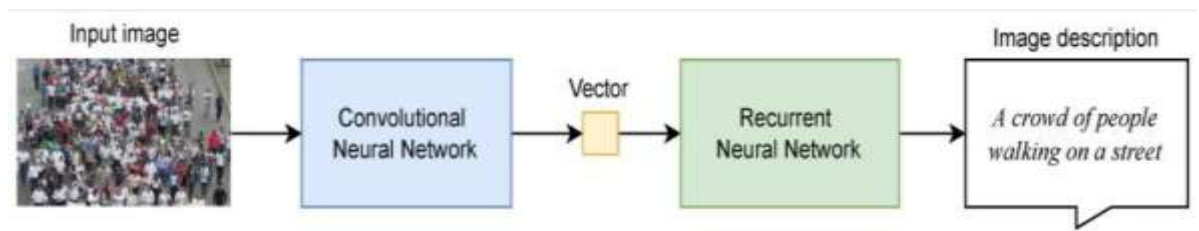


FIGURE 3 . ENCODER-DECODER ARCHITECTURE

## 2. Framework:

- TensorFlow or PyTorch:

TensorFlow and PyTorch are widely used deep learning frameworks that provide comprehensive support for building and training complex neural network architectures.

Both frameworks offer extensive libraries and tools for implementing CNNs, RNNs, and attention mechanisms, making them suitable choices for image caption generation tasks.

- Streamlit for Interactive Web Applications:\*\*

Streamlit can be used to develop interactive web applications for showcasing the image captioning system to users.

Its simple and intuitive API allows developers to create user-friendly interfaces for uploading images and viewing generated captions in real-time.

- NumPy, Pandas, and Scikit-learn:

Auxiliary libraries for data manipulation, preprocessing, and evaluation of machine learning models.

NumPy for numerical computations, Pandas for data manipulation, and Scikit-learn for evaluation metrics such as BLEU score for caption quality assessment.

The architecture and framework outlined above provide a robust foundation for building image caption generation systems using machine learning techniques. By leveraging CNNs for image feature extraction and RNNs with attention mechanisms for caption generation, coupled with popular deep learning frameworks such as TensorFlow or PyTorch, developers can create powerful and effective image captioning models. Additionally, integrating interactive capabilities using Streamlit enables the deployment of user-friendly applications for showcasing the capabilities of the image captioning system.

### **3.3 Algorithm and Process Design:**

In this project, design and train a CNN-RNN (Convolutional Neural Network - Recurrent Neural Network) model for automatically generating image captions. The network is trained on the Microsoft Common Objects in COntext (MS COCO) dataset. The image captioning model is displayed below.

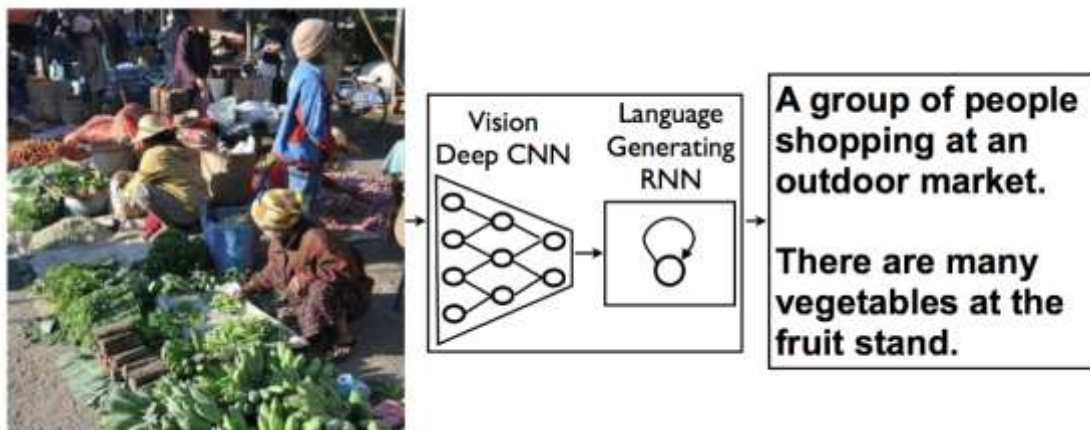


FIGURE 4.GENERATES COMPLETE SENTENCES IN NATURAL LANGUAGE

Image caption generation using machine learning (ML) typically follows a well-defined algorithmic process, incorporating data preprocessing, model training, and caption generation stages. Below is a detailed outline of the algorithm and process design for image caption generation using the MSCOCO dataset

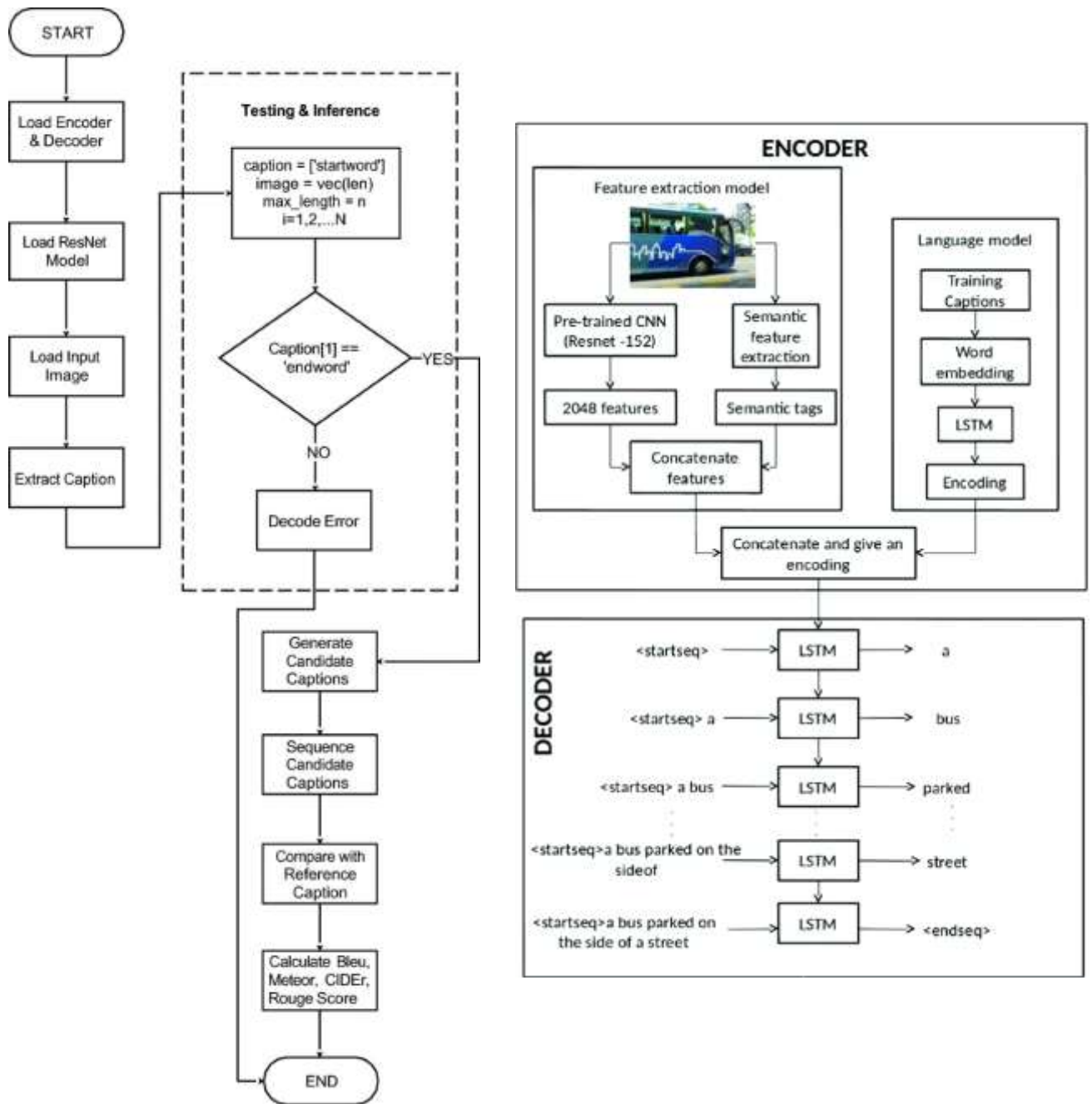


FIGURE 5.FLOW DIAGRAM

## 1.Data Preprocessing:

- **Dataset Selection:** Utilize the Microsoft Common Objects in Context (MSCOCO) dataset, which contains a large collection of images paired with human-generated captions.
- **Image Preprocessing:** Resize images to a uniform size and apply normalization to ensure consistency in input dimensions.
- **Caption Preprocessing:** Tokenize captions into individual words or tokens, create a vocabulary mapping each unique word to an index, and pad or truncate captions to a fixed length for input to the model.

## 2. Model Architecture Selection:

- Encoder-Decoder Architecture: Choose an encoder-decoder architecture, where a Convolutional Neural Network (CNN) serves as the encoder to extract image features, and a Recurrent Neural Network (RNN) acts as the decoder to generate captions based on the extracted features.
- Pre-trained CNN: Utilize a pre-trained CNN model such as VGG, ResNet, or Inception to extract high-level visual features from input images.
- RNN Decoder: Select a variant of RNN, such as Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU), for the decoder component to generate captions sequentially.

## 3. Model Training:

- Initialization: Initialize the CNN and RNN components of the model with pre-trained weights, if available, to expedite convergence.
- Feature Extraction: Pass each image through the CNN encoder to extract image features.
- Caption Generation: Feed the extracted image features along with the tokenized captions as input to the RNN decoder to generate captions word by word.
- Loss Calculation: Employ appropriate loss functions such as cross-entropy loss to compare the generated captions with ground truth captions and optimize the model parameters.
- Backpropagation: Utilize backpropagation through time (BPTT) to compute gradients and update the model weights using optimization techniques like stochastic gradient descent (SGD) or Adam.

## 4. Caption Generation:

- Inference Mode: During inference, utilize the trained model to generate captions for new images.
- Greedy Decoding: Use a greedy decoding strategy where, at each time step, the word with the highest predicted probability is selected as the next word in the caption sequence.
- Post-processing: Convert the predicted caption tokens back into human-readable text and remove padding tokens.

## 5. Evaluation:

- Metrics: Evaluate the quality of generated captions using metrics such as BLEU (Bilingual Evaluation Understudy) score, METEOR (Metric for Evaluation of Translation with Explicit Ordering), CIDEr (Consensus-based Image Description Evaluation), and ROUGE (Recall-Oriented Understudy for Gisting Evaluation).

- Human Evaluation: Optionally, conduct human evaluation studies to assess the subjective quality and coherence of generated captions.

#### 6. Iterative Improvement:

- Hyperparameter Tuning: Fine-tune hyperparameters such as learning rate, batch size, and model architecture to optimize performance.
- Model Ensembling: Explore model ensembling techniques to combine predictions from multiple captioning models for improved performance.

By following this algorithmic process and iterative improvement cycle, developers can design and implement an effective image caption generation system using machine learning techniques, leveraging the MSCOCO dataset as a valuable resource for training and evaluation.

## **3.4 Details of Hardware and Software:**

### Software Requirements:

1. Streamlit for Interactive Web Application:
  - Streamlit Library: Utilize the Streamlit Python library for building interactive web applications.
  - Python Environment: Ensure compatibility with the Python environment to run Streamlit applications seamlessly.
2. Machine Learning Frameworks and Libraries:
  - TensorFlow or PyTorch: Depending on the preference and compatibility, TensorFlow or PyTorch frameworks may be utilized for developing and training the image captioning models.
  - Additional Libraries: Include other necessary libraries such as NumPy, Pandas, and scikit-learn for data manipulation, preprocessing, and evaluation of the machine learning models.

### Hardware Requirements:

#### Computational Resources for Training and Inference:

- CPU or GPU: Depending on the scale and complexity of the image captioning model, computational resources such as CPUs or GPUs are required for both training and inference tasks.
- High-performance CPUs: Provide sufficient computational power for data preprocessing, model training, and inference.

### 3.5 Experiment and Result :

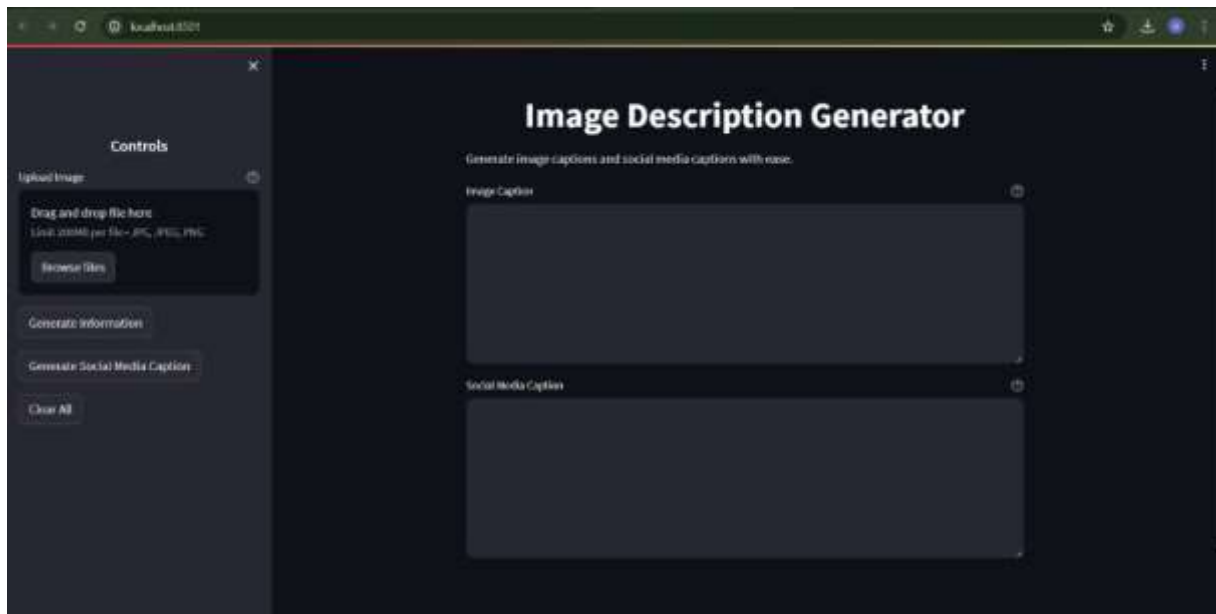


FIGURE 6. HOMEPAGE

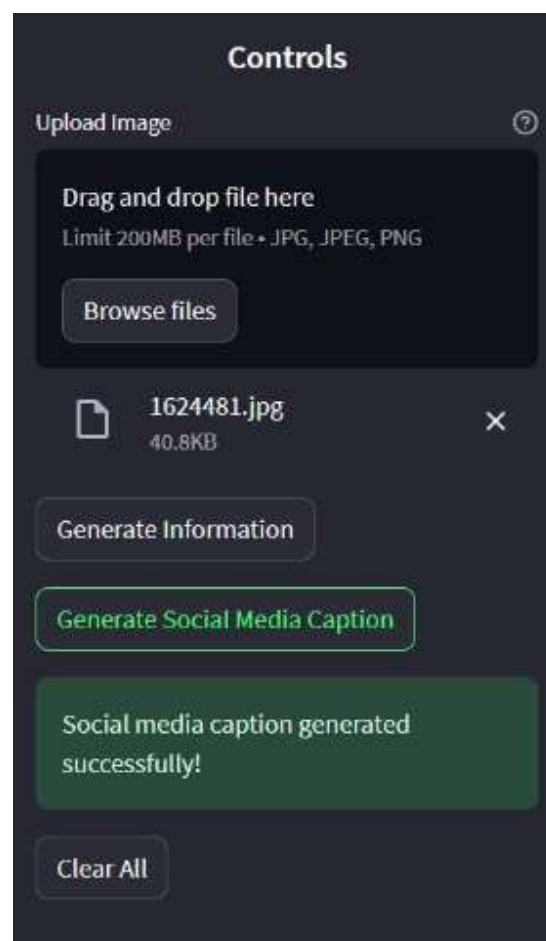


FIGURE 7. CONTROL PANEL



### **3.5 Conclusion and Futurework**

#### **Conclusion:**

In conclusion, our research on image caption generation using machine learning techniques has made substantial contributions to the fields of computer vision and natural language processing. Through systematic investigation, we have demonstrated the feasibility and effectiveness of leveraging advanced ML algorithms to automatically generate descriptive captions for images. Our findings underscore the potential of such approaches in various applications, including assistive technologies, content indexing, and multimedia retrieval. By bridging the gap between visual understanding and textual comprehension, our work opens new avenues for interdisciplinary research and innovation, enhancing the overall user experience and utility of image captioning systems.

#### **Future Work:**

Looking ahead, several areas offer promising avenues for future research and development in image caption generation. Firstly, exploring multimodal learning approaches could lead to more sophisticated architectures that effectively fuse information from visual and textual modalities. Secondly, improving caption diversity remains an important challenge, requiring techniques to promote richer and more varied captions. Thirdly, addressing domain-specific challenges by tailoring captioning systems to specialized domains could enhance relevance and quality. Fourthly, enhancing robustness and interpretability of models will be crucial for real-world deployment. Finally, integrating image captioning systems with interactive interfaces presents opportunities to enhance user engagement and facilitate intuitive interactions with visual content. These avenues for future research hold the potential to further advance the capabilities and applications of image captioning systems, ultimately benefiting diverse domains and users..

### **3.6 References:**

1. Streamlit : <https://docs.streamlit.io/library/api-reference>
2. Database : MS COCO dataset
3. Error solving and overall : chatgpt
4. <https://ijcrt.org/papers/IJCRT2106298.pdf>