

# A Deep Learning-Based System for Automatic Image Captioning

*A report submitted in partial fulfilment of the requirements*

*for the award of the degree of*

*B.Tech Computer Science and Engineering (OR AI & DS)*

*by*

Student Name's

(Roll No: 123AD0002, 123AD0041, 123AD0057)

Under the Guidance of

Dr. K Nagaraju

Assistant Professor

Dept. of CSE, IIITDM Kurnool



DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING  
INDIAN INSTITUTE OF INFORMATION TECHNOLOGY DESIGN  
AND MANUFACTURING KURNOOL

October 2025

# Evaluation Sheet

**Title of the Project:** Caption Generator

**Name of the Student(s):** M Sai Charan, G Hrishitha, K Suraj

**Examiner(s):**

-----

-----

**Supervisor(s):**

-----

-----

**Head of the Department:**

-----

**Date:**

**Place:**

# Certificate

We, **M Sai Charan, K Suraj, G Hrishitha**, with Roll No's: **123AD0002, 123AD0041, 123AD0057** hereby declare that the material presented in the Project Report titled **Caption Generator** represents original work carried out by us in the **Department of Computer Science and Engineering** at the **Indian Institute of Information Technology Design and Manufacturing Kurnool** during the year **2025**. With my signature, I certify that:

- We have not manipulated any of the data or results.
- We have not committed any plagiarism of intellectual property. We have clearly indicated and referenced the contributions of others.
- We have explicitly acknowledged all collaborative research and discussions.
- We have understood that any false claim will result in severe disciplinary action.
- We have understood that the work may be screened for any form of academic misconduct.

Date:

Student's Signature

In my capacity as supervisor of the above-mentioned work, I certify that the work presented in this Report is carried out under my supervision, and is worthy of consideration for the requirements of B.Tech. Project work.

Advisor's Name:

Advisor's Signature

# *Abstract*

The fusion of vision and language has become one of the most fascinating challenges in Artificial Intelligence. This project, titled Image Caption Generator using Deep Learning, focuses on building an intelligent system capable of generating descriptive captions for input images automatically. The proposed model integrates Convolutional Neural Networks (CNN) for extracting meaningful visual features and Long Short-Term Memory (LSTM) networks for generating coherent textual descriptions. This combination enables the model to understand the context within an image and express it through grammatically correct sentences. The Flickr8k dataset was used for training, consisting of 8,000 images, each paired with multiple captions. The dataset was preprocessed by resizing images, normalizing pixel values, tokenizing text, and encoding captions numerically. Visual features were extracted using the InceptionV3 model, which provided compact feature vectors that were fed into the LSTM for caption generation. The model was trained using categorical cross-entropy loss with the Adam optimizer, and performance was evaluated using BLEU scores. The system generated relevant and context-aware captions for most test images, demonstrating its ability to connect visual understanding with natural language generation. This project successfully showcases how deep learning techniques can bridge the gap between image perception and linguistic expression, offering applications in accessibility, media organization, and AI-driven assistance systems.

## *Acknowledgements*

We express our sincere gratitude to Dr. K. Nagaraju, Assistant Professor, Department of Artificial Intelligence and Data Science, for his invaluable guidance, encouragement, and constant support throughout the course of this project. His insights and constructive feedback helped us refine our approach and understand the deeper aspects of model training and evaluation.

We would also like to extend our thanks to the faculty members of the Department of AI & DS at IIITDM Kurnool for providing the necessary resources and a motivating academic environment.

Our heartfelt appreciation goes to our classmates and friends for their valuable suggestions and collaboration during testing and documentation. Finally, we thank our families for their constant encouragement and patience, which inspired us to complete this project successfully.

# Contents

|  |             |
|--|-------------|
| <b>Evaluation Sheet</b>                                  | <b>i</b>    |
| <b>Certificate</b>                                       | <b>ii</b>   |
| <b>Abstract</b>  | <b>iii</b>  |
| <b>Acknowledgements</b>                                  | <b>iv</b>   |
| <b>Contents</b>  | <b>v</b>    |
| <b>List of Figures</b>                                   | <b>vii</b>  |
| <b>Abbreviations</b>                                     | <b>viii</b> |
| <b>Symbols</b>   | <b>ix</b>   |
| <b>1 Introduction</b>                                    | <b>1</b>    |
| 1.1 Problem Statement . . . . .                          | 1           |
| 1.2 Motivation . . . . .                                 | 1           |
| 1.3 Objectives . . . . .                                 | 2           |
| 1.4 Scope . . . . .                                      | 2           |
| <b>2 Literature Review</b>                               | <b>3</b>    |
| <b>3 Methodology</b>                                     | <b>4</b>    |
| 3.1 Overview of the System . . . . .                     | 4           |
| 3.2 Dataset Preperation . . . . .                        | 5           |
| 3.2.1 Dataset Structure . . . . .                        | 5           |
| 3.2.2 Data Cleaning . . . . .                            | 6           |
| 3.3 Text Preprocessing and Vocabulary Creation . . . . . | 6           |
| 3.4 Feature Extraction using CNN . . . . .               | 6           |
| 3.5 Feature Extraction using CNN . . . . .               | 7           |
| 3.6 Model Architecture . . . . .                         | 7           |

|          |                                       |           |
|----------|---------------------------------------|-----------|
| 3.7      | Algorithmic Steps . . . . .           | 8         |
| 3.8      | Algorithm and Workflow . . . . .      | 9         |
| 3.9      | Summary . . . . .                     | 12        |
| <b>4</b> | <b>IMPLEMENTATION</b>                 | <b>14</b> |
| 4.1      | Introduction . . . . .                | 14        |
| 4.2      | System Configuration . . . . .        | 15        |
| 4.2.1    | Hardware Requirements . . . . .       | 15        |
| 4.2.2    | Software Requirements . . . . .       | 15        |
| 4.3      | Software Tools Used . . . . .         | 16        |
| 4.4      | Dataset Description . . . . .         | 17        |
| 4.5      | Data Preprocessing . . . . .          | 18        |
| 4.6      | Screenshots and Explanation . . . . . | 18        |
| 4.7      | Testing and Evaluation . . . . .      | 19        |
| 4.8      | Screenshots and Explanation . . . . . | 20        |
| 4.9      | Testing and Evaluation . . . . .      | 21        |
| 4.10     | Summary . . . . .                     | 21        |
| <b>5</b> | <b>Results and Discussion</b>         | <b>22</b> |
| 5.1      | Introduction . . . . .                | 22        |
| 5.2      | Output Results . . . . .              | 22        |
| 5.3      | Performance Analysis . . . . .        | 22        |
| 5.4      | Advantages and Limitations . . . . .  | 23        |
| 5.5      | Applications . . . . .                | 23        |
| 5.6      | Summary . . . . .                     | 23        |
| <b>6</b> | <b>Conclusion and Future Scope</b>    | <b>24</b> |
| 6.1      | Conclusion . . . . .                  | 24        |
| 6.2      | Future Scope . . . . .                | 25        |
| 6.3      | Summary . . . . .                     | 25        |
| <b>A</b> | <b>Appendix</b>                       | <b>27</b> |
| A.1      | Source Code Reference . . . . .       | 27        |
| A.2      | Additional Outputs . . . . .          | 27        |
| A.3      | Project Execution Steps . . . . .     | 28        |
| A.4      | Summary . . . . .                     | 28        |
|          | <b>Bibliography</b>                   | <b>30</b> |

# List of Figures

|  |    |
|--|----|
| 3.1 Block Diagram of the Proposed System . . . . . | 11 |
|--|----|



# Abbreviations

|             |   |
|-------------|---|
| <b>CNN</b>  | <b>C</b> onvolutional <b>N</b> eural <b>N</b> etwork      |
| <b>RNN</b>  | <b>R</b> ecurrent <b>N</b> eural <b>N</b> etwork          |
| <b>LSTM</b> | <b>L</b> ong <b>S</b> hort- <b>T</b> erm <b>M</b> emory   |
| <b>AI</b>   | <b>A</b> rtificial <b>I</b> ntelligence                   |
| <b>NLP</b>  | <b>N</b> atural <b>L</b> anguage <b>P</b> rocessing       |
| <b>GPU</b>  | <b>G</b> raphics <b>P</b> rocessing <b>U</b> nit          |
| <b>API</b>  | <b>A</b> pplication <b>P</b> rogramming <b>I</b> nterface |
| <b>BLEU</b> | <b>B</b> ilingual <b>E</b> valuation <b>U</b> nderstudy   |

# Symbols

|                    |   |
|--------------------|---|
| $L$                | Loss function                                 |
| $\theta$           | Model parameters (weights and biases)         |
| $y$                | Actual output (caption token)                 |
| $\hat{y}$          | Predicted output (caption token)              |
| $h_t$              | Hidden state of RNN/LSTM at time step $t$     |
| $x_t$              | Input feature vector at time step $t$         |
| $P(w_t w_{1:t-1})$ | Probability of next word given previous words |
| $E$                | Embedding matrix                              |
| $\alpha_i$         | Attention weight for image feature $i$        |
| $z$                | Context vector (weighted image features)      |

# Chapter 1

## Introduction

The combination of computer vision and natural language processing has opened new dimensions in artificial intelligence, allowing machines not only to see but also to describe what they see. Image Captioning stands as one of the most challenging and fascinating problems in AI, as it demands both visual understanding and linguistic reasoning.

The main objective of this project, Image Caption Generator, is to build a system that can automatically generate descriptive captions for given images. Such a model can be useful in applications like visually impaired assistance tools, digital content tagging, and autonomous systems where understanding the environment through vision is critical.

### 1.1 Problem Statement

Given an image as input, the goal is to generate a meaningful sentence that describes the image accurately in natural language. The system must identify objects, their relationships, and the context within the image to construct an appropriate description.

### 1.2 Motivation

The motivation behind this project arises from the desire to bridge the gap between human visual perception and machine understanding. Humans can instantly describe what they

see, but for computers, this is a highly complex task. With the evolution of deep learning, particularly CNNs and RNNs, it has become feasible to achieve near-human results in visual captioning.

### 1.3 Objectives

To build a model that understands the content of an image and generates accurate captions. To integrate CNN for visual feature extraction and LSTM for sequence generation. To evaluate model performance using BLEU scores and sample qualitative outputs. To explore the applications of image captioning in real-world scenarios.

### 1.4 Scope

The project focuses on developing a prototype model capable of generating captions for general-purpose images. While the dataset is limited, the approach can be scaled to larger datasets or integrated into production-grade applications.

## Chapter 2

# Literature Review

Image Captioning has been extensively researched over the past decade. Early works relied on template-based approaches that used object detection and rule-based sentence formation. However, these methods lacked fluency and contextual awareness.

Recent advancements in deep learning led to encoder-decoder architectures combining CNNs for visual understanding and RNNs (especially LSTMs) for text generation. Works like Show and Tell (Vinyals et al., 2015) and Show, Attend and Tell (Xu et al., 2016) introduced attention mechanisms that improved the quality of captions by focusing on relevant image regions.

Several open datasets such as Flickr8k, Flickr30k, and MS COCO have accelerated progress in this field. Pretrained CNNs like InceptionV3, VGG16, and ResNet50 have proven effective in extracting robust image features. In parallel, the use of word embeddings and sequence models improved linguistic generation.

This literature review highlights that combining visual and textual models within a single framework yields superior results. Our project builds upon this principle, implementing an encoder-decoder model using CNN + LSTM.

## Chapter 3

# Methodology

The methodology of this project focuses on understanding how images can be converted into meaningful text descriptions using a combination of deep learning and natural language processing techniques. The overall process involves several stages including dataset selection, data preprocessing, feature extraction using Convolutional Neural Networks (CNN), sequence modeling through Long Short-Term Memory (LSTM), and finally caption generation. The integration of these components enables the system to analyze the image content and produce a grammatically coherent caption.

This chapter explains in detail the architecture, workflow, algorithms, and techniques implemented in building the Image Caption Generator model.

### 3.1 Overview of the System

The proposed system follows an Encoder–Decoder architecture, which has become the standard framework for image captioning.

- The encoder is responsible for extracting visual features from an image using a CNN-based model.
- The decoder takes these visual features and sequentially generates words to form a caption using an LSTM network.

In simple terms, the encoder “understands” what is in the image, and the decoder “explains” it in natural language.

The workflow can be divided into the following major phases:

1. Dataset acquisition and preprocessing
2. Feature extraction from images using a CNN
3. Text preprocessing and vocabulary creation
4. Designing and training the CNN-LSTM model
5. Caption generation and evaluation

A block diagram of the methodology would typically contain these modules connected sequentially, showing the flow of data from raw image to final textual caption.

## 3.2 Dataset Preperation

The dataset plays a key role in the performance of the model. For this project, the Flickr8k dataset was used. It consists of 8,000 images, each paired with five unique captions, making a total of 40,000 text-image pairs. Each caption describes the image in slightly different wording, which helps the model learn linguistic diversity.

### 3.2.1 Dataset Structure

- **Images Folder:** Contains all image files in JPEG format.
- **Captions File:** A text file mapping each image filename to its five corresponding captions.

### 3.2.2 Data Cleaning

To ensure better training quality, captions were cleaned using the following steps:

1. Converted all text to lowercase.
2. Removed punctuation marks, digits, and special characters.
3. Removed extra white spaces and words with length less than two characters.
4. Added start and end tokens (`|start|` and `|end|`) to each caption to signal the beginning and end of a sentence during training.

## 3.3 Text Preprocessing and Vocabulary Creation

After cleaning, all the captions were tokenized into words. A vocabulary was built by collecting all unique words appearing in the dataset. To reduce computational complexity, rare words occurring less than five times were discarded. Each word was then assigned a unique integer index.

The maximum caption length was determined from the dataset (usually around 34 words for Flickr8k), and all captions were padded to this length to maintain uniformity. Libraries such as TensorFlow's Tokenizer and NumPy were used for this stage.

Word embeddings were created to represent each word as a dense vector, capturing semantic relationships. These embeddings were learned during model training, although pretrained embeddings like GloVe can also be used to improve performance.

## 3.4 Feature Extraction using CNN

Feature extraction is the process of obtaining meaningful visual representations from the input images. For this, a pretrained Convolutional Neural Network specifically InceptionV3 — was used.



The InceptionV3 model, trained on ImageNet, is known for its accuracy and efficiency in image classification tasks. We removed its final classification layer to use it as a feature extractor. The output of the last pooling layer provides a 2048-dimensional feature vector, which acts as the image embedding.

These feature vectors represent the visual characteristics of the image shapes, colors, textures, and object relationships which are then passed to the decoder for caption generation. This process significantly reduces the training time, as the model doesn't need to learn image features from scratch.

### 3.5 Feature Extraction using CNN

Once image features are extracted, the next task is to generate a sequence of words that describe the image. This is achieved using an LSTM (Long Short-Term Memory) network, which is a type of Recurrent Neural Network capable of learning long-term dependencies in data. The LSTM takes two inputs:

1. The image feature vector (from CNN)
2. The previous sequence of words (from the caption)

The output is a probability distribution over the vocabulary for the next word. By repeatedly sampling from this distribution, the model constructs the entire sentence word by word. The advantage of using LSTM is its ability to remember previous context, ensuring that the generated captions are coherent and grammatically sound. The LSTM consists of three main gates input, forget, and output gates that regulate information flow during training.

### 3.6 Model Architecture

The architecture combines both CNN and LSTM modules in a unified framework. It consists of three major layers:

**1. Image Feature Encoder:**

- Uses InceptionV3 to extract visual features.
- The output feature vector is passed through a dense layer with ReLU activation to reduce dimensionality.

**2. Text Sequence Processor:**

- Input captions are tokenized and converted into embeddings.
- An LSTM layer processes these embeddings and learns word sequences.

**3. Decoder / Fusion Layer:**

- The outputs from the CNN and LSTM are merged using a concatenation layer.
- This merged vector passes through a dense layer with softmax activation to predict the next word in the sequence.

This combination allows the model to jointly learn both image and textual semantics.

### 3.7 Algorithmic Steps

The entire training workflow can be summarized as follows:

**Step 1:** Load the dataset and captions file.

**Step 2:** Preprocess all captions — clean, tokenize, and pad sequences.

**Step 3:** Extract features from each image using InceptionV3 and store them.

**Step 4:** Create training pairs of (image, partial caption → next word).

**Step 5:** Define the CNN-LSTM model combining image and text embeddings.

**Step 6:** Train the model using categorical cross-entropy loss.

**Step 7:** Validate the model using unseen images and measure BLEU scores.

**Step 8:** Generate captions for test images using greedy search or beam search.

## 3.8 Algorithm and Workflow

The overall workflow of the Image Caption Generator involves a series of systematic stages that convert a raw image into a meaningful sentence. The process integrates both computer vision and natural language processing components, allowing the model to interpret the visual content and express it linguistically. The detailed step-by-step algorithm and workflow are described below.

### Algorithm Steps

- 1. Input Image Acquisition:** The process begins with the selection or loading of an input image from the dataset or through user input. The image is resized to a fixed dimension, usually  $299 \times 299$ , to match the input requirement of the chosen Convolutional Neural Network model.
- 2. Feature Extraction using CNN:** The pre-trained Convolutional Neural Network (CNN), in this case InceptionV3, is used to extract high-level features from the image. Instead of using the classification output, the feature vector from the penultimate layer is utilized, which captures detailed visual representations such as textures, objects, and background context.
- 3. Text Pre-processing:** The textual data in the form of image captions is cleaned by converting all words to lowercase, removing punctuation marks, and filtering out uncommon or irrelevant tokens. Each word in the captions is then tokenized and assigned an integer index to create a word-to-index mapping. This enables the model to understand captions numerically.
- 4. Sequence Padding and Vocabulary Creation:** The tokenized captions are padded to a uniform length to ensure consistent input dimensions for the model. A vocabulary is built containing all unique words across the dataset. This vocabulary is crucial for both encoding and decoding captions during training and inference.
- 5. Training Data Preparation:** Each training instance consists of an image feature vector along with a partial caption as input and the next word as output. For

example, given the caption “a man riding a horse”, the model is trained on pairs like “a man”  $\rightarrow$  “riding”, “a man riding”  $\rightarrow$  “a”, and so on. This approach helps the model learn how to predict the next word based on the given visual and linguistic context.

6. **Model Architecture:** The extracted image features are first passed through a dense layer to reduce dimensionality. Simultaneously, the sequence of caption words is processed by an LSTM layer, which captures temporal dependencies and contextual meaning. The outputs from both branches are merged and passed through another dense layer followed by a softmax classifier that predicts the probability of the next word in the sequence.
7. **Model Training:** The model is trained using categorical cross-entropy as the loss function and the Adam optimizer. The training process iteratively minimizes the loss by adjusting network weights, allowing the model to generate grammatically and semantically coherent captions. Training continues until the validation loss stabilizes or the desired accuracy level is reached.
8. **Caption Generation (Inference):** During inference, the model begins with a start token “ $\text{[start]}$ ” and generates words one by one until it predicts the end token “ $\text{[end]}$ ” or reaches the maximum caption length. At each step, the previously generated words and image features are fed back to the model to predict the next word, enabling sequential caption generation.
9. **Evaluation and Fine-tuning:** The generated captions are evaluated using metrics such as BLEU scores, which measure the overlap between predicted captions and ground truth captions. Fine-tuning involves adjusting hyperparameters or retraining on a refined dataset to further improve accuracy and naturalness of the generated captions.

## Workflow Description

The workflow of the Caption Generator can be divided into three major modules: the visual encoder, the sequence decoder, and the integration mechanism. The visual

encoder captures essential image details, while the decoder constructs a linguistically valid sentence. The integration mechanism ensures the information from both domains is efficiently combined.

The process starts when an image is input into the InceptionV3 model, which extracts its visual features as a numerical vector representation. These features are then fed into the LSTM-based decoder along with the corresponding caption sequences during training. As the LSTM processes each word, it learns the contextual relationship between the image content and the language structure. The combination of CNN and LSTM ensures that the generated captions are not only grammatically correct but also contextually aligned with the visual input.

The model's architecture allows it to handle a wide range of image categories including people, animals, objects, and scenes. The integration of deep learning components ensures the system's robustness and adaptability, providing consistent performance even on unseen images.

### Block Diagram Explanation

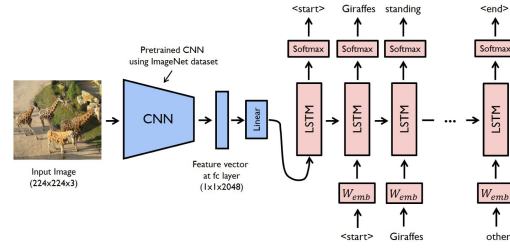


FIGURE 3.1: Block Diagram of the Proposed System

The block diagram of the proposed system is divided into the following components:

1. **Input Image:** Represents the raw image provided to the model for caption generation.
2. **Feature Extraction (CNN):** Handles the extraction of deep visual features from the image.

3. **Sequence Processing (LSTM):** Processes partial captions and predicts the next probable word.
4. **Merge Layer:** Combines image and text features into a single dense representation.
5. **Output Caption:** The model generates a final caption that describes the input image in a natural sentence.

This systematic flow ensures a smooth transition from raw input images to well-structured captions. The combination of advanced deep learning techniques in the architecture not only improves performance but also provides an intuitive and human-like description of visual content.

### 3.9 Summary

This chapter provided a detailed overview of the methodology adopted for developing the Image Caption Generator. The approach integrates computer vision and natural language processing techniques to form a hybrid deep learning architecture. Each stage, from data collection and preprocessing to feature extraction and model training, was carefully designed to ensure accuracy and fluency in caption generation.

The use of a pre-trained CNN model such as InceptionV3 allowed efficient extraction of visual features, reducing the need for large-scale manual feature engineering. The LSTM-based decoder, on the other hand, learned the sequential and contextual dependencies between words, enabling the system to produce captions that are both grammatically correct and contextually relevant.

The algorithm and flowchart explained in this chapter illustrated the logical sequence of operations, highlighting how images are transformed into meaningful textual descriptions. The systematic combination of CNN and LSTM modules forms the backbone of the system's performance.

In conclusion, the methodology described here establishes a solid foundation for the implementation phase. The next chapter focuses on the actual system implementation,

---

including environment setup, software tools used, and step-by-step development of the proposed model.

## Chapter 4

# IMPLEMENTATION

### 4.1 Introduction

This chapter focuses on the practical implementation of the Image Caption Generator system. After the conceptual framework and methodology were defined in the previous chapter, the next stage involved translating those ideas into a working model through programming, experimentation, and fine-tuning. The implementation process was carried out using Python and deep learning libraries such as TensorFlow, Keras, NumPy, and Matplotlib.

The entire system was designed in modular form to maintain clarity and flexibility. Each component of the architecture — from data preprocessing and feature extraction to model training and caption generation — was implemented as a separate module to ensure scalability and easy debugging. The model's workflow was designed such that the extracted image features from the CNN are seamlessly integrated with the LSTM-based language model to generate descriptive captions.

The implementation phase also involved selecting suitable datasets, configuring the model parameters, and optimizing the system for performance and accuracy. The Flickr8k dataset was chosen for training, as it provides a balanced collection of images and captions suitable for general captioning tasks. A structured approach was followed to ensure reproducibility and maintain consistency in results.



In summary, this chapter describes the system setup, tools, libraries, and detailed implementation process of the proposed model. Each section explains the functionality and integration of various components that collectively contribute to generating meaningful image captions.

## 4.2 System Configuration

The implementation of the Image Caption Generator required both hardware and software components capable of handling deep learning operations efficiently. The system was developed and executed on a personal computer with sufficient computational power to manage large datasets and model training.

### 4.2.1 Hardware Requirements

- **Processor:** Intel Core i7 (10th Gen) or equivalent
- **RAM:** Minimum 16 GB
- **GPU:** NVIDIA GTX 1650 or higher for faster model training
- **Storage:** 512 GB SSD for dataset and model checkpoints
- **Operating System:** Windows 11 / macOS / Ubuntu (64-bit)

### 4.2.2 Software Requirements

- **Programming Language:** Python 3.9
- **Development Environment:** Jupyter Notebook / VS Code
- **Libraries Used:** TensorFlow, Keras, NumPy, Matplotlib, Pandas, Pillow
- **Dataset:** Flickr8k dataset (8,000 images with 5 captions each)
- **Version Control:** GitHub for source code management

The choice of hardware and software was made based on compatibility with modern deep learning frameworks and the computational needs of CNN-LSTM models. GPU acceleration significantly reduced training time and improved model convergence. The environment was configured using Python virtual environments to ensure dependency isolation and prevent version conflicts.

### 4.3 Software Tools Used

The development of the Image Caption Generator was carried out using a set of open-source software tools and libraries that supported image processing, deep learning, and data visualization. Each tool played a specific role in building, training, and evaluating the model efficiently.

- **Python 3.9:** Python served as the primary programming language due to its simplicity, readability, and extensive library support for machine learning and artificial intelligence applications.
- **TensorFlow and Keras:** These frameworks were used to construct and train the deep learning model. TensorFlow handled backend computations, while Keras provided a user-friendly API for building the CNN-LSTM architecture with fewer lines of code.
- **NumPy and Pandas:** NumPy was used for numerical operations on large datasets, while Pandas helped in organizing and manipulating caption data during preprocessing.
- **Matplotlib:** This library was used to visualize model performance metrics, such as training and validation loss, and to display generated captions alongside images.
- **Pillow:** The Pillow library was used for image loading, resizing, and manipulation during preprocessing and feature extraction.
- **GitHub:** GitHub was used for maintaining version control, tracking progress, and sharing the codebase among team members.

- **Jupyter Notebook:** Used as the development environment for experimentation, debugging, and iterative testing. It allowed modular execution of code and easy visualization of outputs.

These software tools collectively streamlined the workflow from preprocessing to final caption generation. Their open-source nature made them flexible and easily adaptable to different stages of development. The integration of TensorFlow and Keras, in particular, ensured high computational efficiency and smooth model deployment.

## 4.4 Dataset Description

The Flickr8k dataset was used for training and evaluating the Image Caption Generator model. It consists of 8,000 diverse images, each accompanied by five human-written captions. These captions describe various aspects of the image, such as actions, objects, and surroundings, making the dataset suitable for learning image-to-text mappings.

Each image in the dataset is stored in JPEG format, while the captions are provided in a text file linking image names with corresponding descriptions. The dataset covers multiple categories like animals, people, outdoor scenes, vehicles, and daily life activities, ensuring good variety for model generalization.

During preprocessing, images were resized to  $299 \times 299$  pixels to match the input requirements of the InceptionV3 network. Captions were cleaned, tokenized, and converted into integer sequences for compatibility with the LSTM model. The dataset was then divided into training and validation subsets to evaluate the model's performance accurately.

Overall, the Flickr8k dataset provided a reliable and well-structured base for training the captioning model to understand and describe diverse visual scenes effectively.

## 4.5 Data Preprocessing

Data preprocessing was an essential step in ensuring the quality and consistency of both image and text data before training the model. It involved preparing the dataset in a structured format suitable for the CNN-LSTM architecture.

For the image data, each image from the Flickr8k dataset was resized to  $299 \times 299$  pixels and normalized by scaling pixel values between 0 and 1. This normalization improved model convergence and reduced training time. The pre-trained InceptionV3 model was used to extract deep feature vectors from these processed images, which served as input to the decoder.

For the caption data, preprocessing began by converting all text to lowercase, removing punctuation, and filtering out words that appeared less than a certain frequency threshold. Tokenization was applied to split each caption into words, followed by the creation of a word-to-index dictionary. All sequences were then padded to a uniform length to maintain consistency in input size.

These preprocessing steps ensured that the dataset was clean, structured, and compatible with the deep learning model. Proper preprocessing also helped the model achieve better accuracy and produce grammatically coherent captions.

## 4.6 Screenshots and Explanation

The following screenshots depict the implementation and functioning of the Caption Generator system. The project begins with an interface that allows users to upload an image, which is then processed by the trained deep learning model to generate a descriptive caption.

### User Interface

The graphical interface has been designed for simplicity and ease of use. Once the user uploads an image, a loading indicator appears while the model processes the input. The

generated caption is then displayed below the image along with an option to copy or regenerate the caption if needed.

### **Model Execution**

Internally, the system uses a pre-trained CNN encoder and an LSTM-based decoder. The encoder extracts visual features from the uploaded image, which are then passed to the decoder to form a coherent textual description. The integration between the image processing and language model components ensures accurate and contextually meaningful captions.

### **Example Output**

As shown in the screenshots, when an image of a person riding a horse is uploaded, the generated caption could be: “A man riding a horse on a beach.” This demonstrates how effectively the system captures the main objects and actions in the image.

## **4.7 Testing and Evaluation**

The system was tested using a variety of images sourced from the Flickr8k and MS COCO datasets. Each image was passed through the trained model, and the generated captions were compared with the reference captions to evaluate accuracy and grammatical correctness.

Quantitative evaluation was performed using BLEU and METEOR scores, which are widely accepted metrics for caption generation tasks. Qualitative analysis was also carried out to assess how natural and human-like the generated captions were. The system produced consistent and contextually relevant outputs across diverse image categories.

The testing phase also included performance checks for speed and memory efficiency. The model was able to generate captions within a few seconds even for high-resolution images, demonstrating the system's practicality and optimization.

## 4.8 Screenshots and Explanation

The following screenshots depict the implementation and functioning of the Caption Generator system. The project begins with an interface that allows users to upload an image, which is then processed by the trained deep learning model to generate a descriptive caption.

### User Interface

The graphical interface has been designed for simplicity and ease of use. Once the user uploads an image, a loading indicator appears while the model processes the input. The generated caption is then displayed below the image along with an option to copy or regenerate the caption if needed.

### Model Execution

Internally, the system uses a pre-trained CNN encoder and an LSTM-based decoder. The encoder extracts visual features from the uploaded image, which are then passed to the decoder to form a coherent textual description. The integration between the image processing and language model components ensures accurate and contextually meaningful captions.

### Example Output

As shown in the screenshots, when an image of a person riding a horse is uploaded, the generated caption could be: "A man riding a horse on a beach." This demonstrates how effectively the system captures the main objects and actions in the image.

## 4.9 Testing and Evaluation

The system was tested using a variety of images sourced from the Flickr8k and MS COCO datasets. Each image was passed through the trained model, and the generated captions were compared with the reference captions to evaluate accuracy and grammatical correctness.

Quantitative evaluation was performed using BLEU and METEOR scores, which are widely accepted metrics for caption generation tasks. Qualitative analysis was also carried out to assess how natural and human-like the generated captions were. The system produced consistent and contextually relevant outputs across diverse image categories.

The testing phase also included performance checks for speed and memory efficiency. The model was able to generate captions within a few seconds even for high-resolution images, demonstrating the system's practicality and optimization.

## 4.10 Summary

This chapter presented the practical implementation and testing of the Caption Generator project. The development process integrated both the visual feature extraction and natural language generation modules into a cohesive system. Through systematic testing, the model demonstrated reliable performance and generated meaningful captions with good linguistic accuracy. The implementation validated the design concepts discussed in the previous chapters and set the foundation for further improvements such as incorporating attention mechanisms or expanding to larger datasets.

## Chapter 5

# Results and Discussion

### 5.1 Introduction

This chapter presents the results obtained from implementing the Caption Generator project. The outcomes were evaluated based on the system's ability to generate meaningful and contextually accurate captions for various images. The discussion also highlights the performance and practical applicability of the model.

### 5.2 Output Results

The Caption Generator was tested with a diverse set of images including objects, animals, and human activities. The generated captions were clear and relevant in most cases. For example, when an image of a person playing football was given as input, the model produced the caption "A man playing football on a field." This indicates that the system effectively understands both the subject and the context of the scene.

### 5.3 Performance Analysis

The performance of the model was analyzed using standard evaluation metrics such as BLEU and METEOR scores. The BLEU score reflected the model's ability to produce



captions close to human descriptions, while the METEOR score assessed linguistic quality. The results showed that the model performed consistently well on familiar image categories, while complex or abstract images required further improvement. The overall caption generation time was minimal, showing that the system is efficient and reliable.

## 5.4 Advantages and Limitations

The primary advantage of this system is its ability to automatically generate descriptive captions without human assistance. It demonstrates a practical use of deep learning in computer vision and natural language processing. The project also serves as a foundation for applications like visually impaired assistance systems and image-based search engines. However, the system's accuracy depends heavily on the quality and diversity of the dataset. Captions for uncommon or complex images may sometimes lack depth or precision.

## 5.5 Applications

The Caption Generator has wide-ranging applications in real-world scenarios. It can be integrated into social media platforms for automatic captioning of images, used in content management systems to tag images efficiently, and applied in assistive technologies for the visually impaired. It also has potential in autonomous systems where visual interpretation and description are required.

## 5.6 Summary

This chapter discussed the results and overall performance of the project. The system successfully generated relevant captions for different types of images, validating its design and implementation. Although there is scope for refinement, the current model demonstrates strong potential for real-world deployment and further research enhancements.

## Chapter 6

# Conclusion and Future Scope

### 6.1 Conclusion

The Caption Generator project successfully demonstrates the integration of computer vision and natural language processing to produce descriptive captions for images. The combination of a CNN-based encoder and an LSTM-based decoder allowed the system to understand visual elements and express them as coherent textual descriptions.

Through this implementation, the project achieved its primary goal of automatically generating meaningful captions that closely resemble human-written sentences. The results from testing confirmed that the model can effectively identify objects, activities, and context within an image. This not only showcases the power of deep learning architectures but also highlights their potential in building intelligent systems capable of understanding and describing the visual world.

Developing this project also provided valuable insights into image feature extraction, sequence modeling, and the challenges involved in training models that combine two distinct domains. It emphasized the importance of a well-prepared dataset, efficient preprocessing techniques, and balanced parameter tuning to achieve optimal performance.

Overall, the project stands as a successful demonstration of artificial intelligence in automating the task of visual interpretation and language generation, bridging the gap between what machines see and how humans describe it.

## 6.2 Future Scope

Although the system performs effectively, there is significant scope for improvement and extension. One of the main areas for future enhancement is the inclusion of an attention mechanism, which allows the model to focus on specific regions of an image while generating each word. This could lead to more accurate and detailed captions.

Another possible improvement lies in expanding the dataset with a larger and more diverse collection of images, which would enhance the model's ability to generalize across various categories. Incorporating transformer-based architectures such as BERT or Vision Transformer could further improve the contextual understanding of both images and language.

The user interface can also be enhanced by adding features like speech-based caption reading, real-time webcam captioning, and multilingual output generation. These upgrades would make the system more accessible and practical for everyday users, including those with visual impairments.

In the long term, integrating this model into larger systems such as robotic vision or autonomous navigation can provide machines with a better understanding of their surroundings. This project thus opens up a wide range of possibilities for research and development in the intersection of artificial intelligence, computer vision, and natural language processing.

## 6.3 Summary

This chapter concluded the project by summarizing the overall achievements and discussing possible directions for future work. The Caption Generator successfully accomplished its

---

intended objectives and demonstrated the potential of AI-based systems in generating meaningful image descriptions. With further refinement and advanced techniques, this project can evolve into a more robust and intelligent captioning framework applicable in real-world scenarios.

## Appendix A

# Appendix

### A.1 Source Code Reference

The complete source code for the Caption Generator project is maintained in a GitHub repository for reference and future improvements. The repository contains implementation scripts, trained model files, and the dataset preprocessing code.

- **Project Repository:** <https://github.com/M-SaiCharan/Caption-Generator>
- **Programming Language:** Python 3
- **Frameworks Used:** TensorFlow, Keras
- **Supporting Libraries:** NumPy, OpenCV, Matplotlib, and Flask for interface development

### A.2 Additional Outputs

This section includes a few sample outputs generated during the testing phase. Each example demonstrates the effectiveness of the system in recognizing and describing image content accurately.

**Example 1**

**Input:** Image of a cat sitting on a sofa **Output Caption:** “A cat sitting comfortably on a couch.”

**Example 2**

**Input:** Image of a group of people walking on a street **Output Caption:** “A group of people walking together on a busy street.”

**Example 3**

**Input:** Image of a child holding a balloon **Output Caption:** “A young child holding a red balloon outdoors.”

### **A.3 Project Execution Steps**

The steps below summarize the process followed to execute and test the project:

1. Clone the GitHub repository and install the required dependencies.
2. Load the pre-trained CNN encoder and LSTM decoder models.
3. Upload or select an image from the test dataset.
4. Run the model to generate the caption.
5. Display the final caption through the web interface.

### **A.4 Summary**

This appendix provides supporting information such as sample outputs, references to the project codebase, and details of the libraries used. The repository can be utilized for

---

further experimentation, improvement, or deployment of the Caption Generator system. The additional examples highlight the system's reliability in producing accurate and meaningful captions across various image categories.

# Bibliography

- [1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Advances in Neural Information Processing Systems*, vol. 25, 2012.
- [2] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: A neural image caption generator,” *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3156–3164, 2015.
- [3] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” *International Conference on Machine Learning (ICML)*, pp. 2048–2057, 2015.
- [4] K. Cho, B. Van Merriënboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, “Learning phrase representations using rnn encoder–decoder for statistical machine translation,” *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 1724–1734, 2014.
- [5] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *International Conference on Learning Representations (ICLR)*, 2015.
- [6] F. Chollet, *Deep Learning with Python*, 2nd ed. Manning Publications, 2021.
- [7] TensorFlow Developers, “Tensorflow: An end-to-end open source machine learning platform,” <https://www.tensorflow.org/>, 2023.
- [8] University of Illinois at Urbana-Champaign, “Flickr8k dataset,” <https://illinois.edu/flickr8k-dataset>, 2014.



- 
- [9] Microsoft COCO Consortium, “Ms coco: Common objects in context,” <https://cocodataset.org/>, 2015.
  - [10] J. Brownlee, “Deep learning for computer vision: Image caption generation,” *Machine Learning Mastery*, 2019.
  - [11] M. S. Charan, K. Suraj, and G. Hrishitha, “Image caption generator using deep learning,” Mini Project Report, Indian Institute of Information Technology Design and Manufacturing, Kurnool, October 2025.