

RAJIV GANDHI INSTITUTE OF TECHNOLOGY
GOVERNMENT ENGINEERING COLLEGE
KOTTAYAM-686 501



DEPARTMENT OF COMPUTER APPLICATIONS

20MCA246 - MAIN PROJECT REPORT

**DEEP LEARNING-BASED IMAGE CAPTION GENERATION BASED
ON CONTEXT**

Submitted By

ASHNA C P

(KTE22MCA-2015)

Under the Guidance of

Prof. Shilpa M Thomas



APJ ABDUL KALAM TECHNOLOGICAL UNIVERSITY
THIRUVANANTHAPURAM

APRIL 2024

DECLARATION

I, undersigned hereby declare that the project report entitled “**DEEP LEARNING-BASED IMAGE CAPTION GENERATION BASED ON CONTEXT**”, submitted for partial fulfillment of the requirements for the award of the degree of Master of Computer Applications of the APJ Abdul Kalam Technological University, Kerala is a bonafide work done by me under the supervision of **Prof. Shilpa M Thomas**. This submission represents my ideas in my own words and where ideas or words of others have been included, I have adequately and accurately cited and referenced the original sources. I also declare that I have adhered to the ethics of academic honesty and integrity and have not misrepresented or fabricated any data or idea or fact or source in my submission. I understand that any violation of the above will be a cause for disciplinary action by the institute and/or the University and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been obtained. This report has not been previously formed as the basis for the award of any degree, diploma, or similar title of any other University.

Place: Kottayam

Date: 22/04/2024

ASHNA C P

ACKNOWLEDGEMENT

I want to express my gratitude to everyone who has supported me throughout the endeavour. First and foremost, I give thanks to God Almighty for His mercy and blessings, for without His unexpected direction, this would still be only a dream.

I sincerely thank **Dr. Prince A**, Principal, Rajiv Gandhi Institute of Technology, Kottayam, for providing the environment in which this project could be completed.

I owe a huge debt of gratitude to **Dr. Reena Murali**, Head of the Department of Computer Applications, for granting permission and making available all of the facilities needed to complete the project properly.

I am grateful to my project guide **Prof. Shilpa M Thomas**, for her helpful criticism of my thesis.

I also express my sincere thanks to the Project Co-ordinators **Dr. John C John** and **Prof. Sreejith V P**, for the constructive suggestions and inspiration throughout the project.

Finally, I would like to take this chance to express my gratitude to the faculty and technical staff of the Department of Computer Applications.

ASHNA C P

ABSTRACT

This project focuses on developing an Image Caption Generator utilizing deep learning methodologies, specifically employing a Convolutional Neural Network (CNN) for image feature extraction and a Long Short-Term Memory (LSTM) network for language modeling. The primary objective is to automate the generation of descriptive captions for visual content, addressing the challenge of effectively associating visual information with textual descriptions. It uses the Flickr8K dataset for efficient model development without compromising performance. The model is trained on pairs of images and corresponding captions, enabling it to learn the complex relationship between visual features and textual context. The CNN extracts high-level features from the images, which are then utilized by the LSTM to generate captions word by word. The training process allows the model to understand and generate coherent descriptions based on the visual signals provided by the images. The project includes the successful generation of human-readable captions for diverse visual content, validated through evaluation metrics such as BLEU, which assess the quality and fluency of the generated captions compared to human-written references. These findings demonstrate the effectiveness of the model in understanding visual content and translating it into natural language descriptions. In conclusion, the project not only addresses the challenge of automatic image captioning but also explores potential future advancements. The integration possibilities with real-time image processing systems and exploration of advanced language models represent promising avenues for further research and development. Ultimately, the project contributes to enhancing human-computer interaction and multimedia understanding across various domains by bridging the gap between visual perception and linguistic expression, thus paving the way for improved content-based image retrieval, assistive technologies, and enriched social media content.

Keywords: Image Caption Generator, CNN, LSTM, Flickr8k dataset, BLEU

Contents

DECLARATION	i
ACKNOWLEDGEMENT	ii
ABSTRACT	iii
LIST OF FIGURES	vi
LIST OF TABLES	vii
LIST OF ABBREVIATIONS	viii
1 INTRODUCTION	1
1.1 Need for the project	1
1.2 Objective	2
1.3 Scope of the project	2
1.4 Organization of Thesis	2
2 LITERATURE REVIEW	3
2.1 Existing System	3
2.2 Study on existing system	3
2.3 Gap Identification	4
3 PROPOSED SYSTEM	5
3.1 Features of proposed system	5

3.1.1	System architecture	5
3.2	MATERIALS AND METHODS	8
3.2.1	Tools	8
3.2.2	Flickr8k dataset	8
3.2.3	Architectures employed for image caption generation	8
4	RESULTS AND ANALYSIS	12
4.1	Result	12
4.2	Analysis	12
4.2.1	Data Cleaning	13
4.2.2	Feature Extraction	13
4.2.3	Tokenization	14
4.2.4	Model Architecture	14
4.3	Model Validation	16
5	CONCLUSION	17
6	FUTURE SCOPE	18
	BIBLIOGRAPHY	19
	ANNEXURE	21

LIST OF FIGURES

3.1	Block diagram ^[5]	6
3.2	CNN Architecture ^[10]	9
3.3	VGG16 Architecture ^[11]	9
3.4	Gates in LSTM ^[12]	10
3.5	CNN-LSTM model ^[13]	11
4.1	Model architecture	15
6.1	Actual and Predicted captions	23
6.2	Caption for an image	23
6.3	GitHub History	24

LIST OF TABLES

2.1	Literature Survey	4
4.1	Data cleaning of captions	13
4.2	BLEU Scores	16

LIST OF ABBREVIATIONS

Abbreviations	Definition
CNN	Convolutional Neural Network
VGG16	Visual Geometry Group 16
LSTM	Long Short-Term Memory
BLEU	Bilingual Evaluation Understudy

Chapter 1

INTRODUCTION

This chapter introduces the project's core motivations, highlighting its aim to develop an Image Caption Generator using deep learning techniques to bridge the gap between visual perception and linguistic expression.

1.1 Need for the project

The need for the Image Caption Generation project arises from the growing importance of bridging visual understanding with linguistic expression in artificial intelligence and computer vision domains. In today's digital landscape, where visual content proliferates across online platforms and digital repositories, the significance of bridging visual understanding with linguistic expression through Image Caption Generation becomes increasingly pronounced. This project addresses the imperative need for machines to interpret and articulate visual content in human-readable language, thereby enhancing accessibility, content organization, and user engagement in computer vision domains.

In content management systems and image databases, automated caption generation streamlines content retrieval and organization, amplifying usability and searchability for users navigating vast repositories of visual data. Furthermore, Image Caption Generation enriches multimedia understanding by empowering machines to interpret and express visual content in natural language. This capability extends its utility across diverse fields such as content-based image retrieval, social media content enrichment, and automated image annotation, augmenting the efficiency and effectiveness of various applications.

By addressing the pressing need for automated image captioning, this project contributes significantly to advancing human-computer interaction paradigms. It renders visual content more accessible and comprehensible to a broader audience, fostering innovation and exploration in artificial intelligence and computer vision research, and paving the way for enhanced human-machine collaboration in navigating and understanding the ever-expanding visual landscape of the digital era.

1.2 Objective

The primary objective of this project is to develop an Image Caption Generator using deep learning techniques. By leveraging convolutional neural networks (CNNs) for image feature extraction and recurrent neural networks (RNNs) like Long Short-Term Memory (LSTM) for language modeling, the system aims to automatically generate coherent and contextually relevant captions for a wide array of images.

1.3 Scope of the project

The scope of this project encompasses the development and implementation of an Image Caption Generator using deep learning techniques, with a focus on robustness and contextual relevance in caption generation. The project involves leveraging large-scale image-caption datasets like Flickr8k for model training and validation, encompassing a diverse range of visual content.

1.4 Organization of Thesis

In the subsequent chapters, the document delves into a comprehensive exploration of the existing systems and literature, highlighting gaps identified in current research. Following this, Chapter 3 presents the proposed system, detailing its design, methodology, and tools utilized. Chapter 4 then presents the results obtained and initiates discussions around them, analyzing findings in the context of the proposed system's objectives. Finally, Chapter 5 concludes the document, summarizing key insights, implications, and avenues for future research.

As I conclude this chapter, the groundwork has been laid for understanding the pivotal need and objectives of the image caption generator using CNN and LSTM.

Chapter 2

LITERATURE REVIEW

The literature survey chapter offers a comprehensive overview of existing research and studies relevant to the image caption generator stages, providing insights into current methodologies and challenges in the field.

2.1 Existing System

In the current landscape of image caption generation, the reliance on manual annotation methods remains common. Despite advancements in machine learning and neural networks, automated image captioning using CNN and LSTM remains limited. Research has progressed in understanding image features and natural language processing, but deploying scalable systems for accurate captions is lacking. Manual annotation introduces errors and hampers efficiency, necessitating the development of automated solutions to bridge this gap.

2.2 Study on existing system

M. Israk Ahmed et al.^[1] proposed Context-based Image Caption using Deep Learning. utilizes Resnet101 for feature extraction and context coding for image processing. The model incorporates SCST and experimental LSTM for captioning. They have achieved an accuracy of 78.3% with 113,287 images as a dataset.

Xu et al ^[2] proposed Show, Attend and Tell: Neural Image Caption Generation with Visual Attention. This paper introduces an attention-based model for generating image descriptions, focusing on relevant image regions while generating captions. use GoogLeNet model. They have achieved an accuracy of 96.88% with 600 images as a dataset.

Anderson et al ^[3] proposed Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering paper. they use CNN and LSTM model. This work explores a combined bottom-up and top-down approach for attention in image captioning, enhancing the model's ability to focus on relevant visual features. They have achieved an accuracy of 70.3% with 600 images as a dataset.

You et al ^[4] proposed Image Captioning with Semantic Attention. They have used GoogleNet model. This study delves into semantic attention mechanisms, allowing the model to attend to semantically meaningful regions when generating image captions. They have achieved an accuracy of 75.06% accuracy with 30000 images as a dataset.

Table 2.1: Literature Survey

Sl No.	Title	Author	Model	Accuracy	No. of images in dataset	Features
1	Context-based Image Caption using Deep Learning	M. Israk Ahmed et al.	Resnet101	78.3%	113,287	Utilizes contextual information for captioning
2	Show, Attend and Tell: Neural Image Caption Generation with Visual Attention	Fu yuesheng et al	GoogLeNet	96.88%	600	Incorporates visual attention mechanisms
3	Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering	Anderson et al	CNN	70.3%	5100	Utilized both top-down and bottom-up mechanisms.
4	Image Captioning with Semantic Attention	You et al	GoogleNet	75.06%	30000	Employs semantic attention for improved captioning

2.3 Gap Identification

In the subsequent chapters, the document delves into a comprehensive exploration of the existing systems and literature, highlighting gaps identified in current research. Following this, Chapter 3 presents the proposed system, detailing its design, methodology, and tools utilized. Chapter 4 then presents the results obtained and initiates discussions around them, analyzing findings in the context of the proposed system's objectives. Chapter 5 concludes the document, summarizing key insights and implications. Finally, in Chapter 6, the future scope of the project will be explored, outlining potential areas for further research and advancement within the field.

As I conclude this chapter, the groundwork has been laid for understanding the pivotal need and objectives of the image caption generator using CNN and LSTM.

Chapter 3

PROPOSED SYSTEM

The proposed system is a comprehensive solution designed to address the challenges of image caption generation. Leveraging advanced deep learning techniques, including state-of-the-art Convolutional Neural Network (CNN) architectures, such as VGG16 and LSTM. The system aims to develop accurate captions on images. This chapter outlines the key features of the proposed system and the tools and technologies employed in its development.

3.1 Features of proposed system

The proposed system for image caption generator uses two different neural networks to generate the captions. The first neural network is Convolutional Neural Network(CNN), which is used to train the images as well as to detect the objects in the image with the help of various pre-trained models like VGG. The second neural network used is Recurrent Neural Network(RNN) based Long Short Term Memory(LSTM), which is used to generate captions from the generated object keywords. As, there is lot of data involved to train and validate the model, generalized machine learning algorithms will not work. Deep Learning has been evolved from the recent times to solve the data constraints on Machine Learning algorithms. GPU based computing is required to perform the Deep Learning tasks more effectively.

By providing appropriate, expressive, and fluid subtitles, Deep Neural Networks can tackle the problems and accelerate the creation of subtitles. Users of social media will no longer have to waste hours searching for subtitles on Google with the system we offer. This technology will provides an easy-to-use platform for social network users to upload selected photographs. Photos of any size can be uploaded and also can read the caption out in English. Tensor flows and algorithms can be used by neural networks to solve any problem and provide appropriate, expressive, and fluent subtitles. It is feasible to calculate automatic metrics efficiently. The time spent searching for captions can be minimized as they will be automatically generated.

3.1.1 System architecture

The system architecture for image caption generation using CNN and LSTM extract high-level features from input images via a pre-trained CNN, which captures spatial information

about objects and textures. These features are then fed into an LSTM network that sequentially processes them to generate captions. Word embeddings represent the vocabulary and contextual information, enriching the caption generation process. During training, the model learns to minimize the discrepancy between predicted and ground truth captions using optimization techniques like Adam. In inference, captions are generated by sampling words from the LSTM output distribution. Evaluation metrics such as BLEU assess the quality and relevance of generated captions. This combined CNN-LSTM architecture achieves contextually relevant image captions, finding extensive use and demonstrating state-of-the-art performance in image captioning tasks.

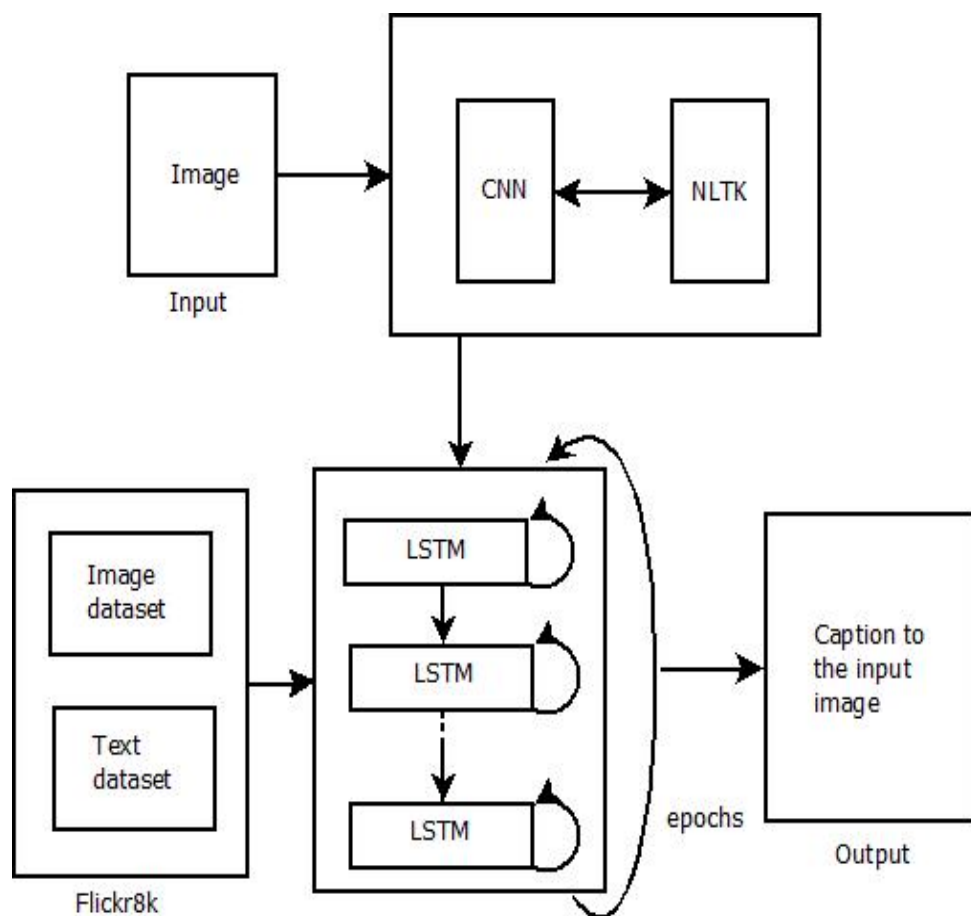


Figure 3.1: Block diagram^[5]

The project involves a sequential process of data collection, training, testing, and model evaluation.

Data Collection This project begins with the acquisition of a suitable dataset, focusing on high-quality images paired with descriptive captions. The Flickr 8K dataset is chosen for its comprehensive collection of images and corresponding human-written captions, providing a

rich source for training and evaluation. The dataset is carefully curated to ensure diversity in visual content and linguistic expressions, essential for robust model development. Each image-caption pair serves as a training example, facilitating the learning process of the Image Caption Generator model.

Training During the training phase, the collected dataset is utilized to train the Image Caption Generator model. The Convolutional Neural Network (CNN) is employed to extract meaningful features from images, while the Long Short-Term Memory (LSTM) network learns to generate descriptive captions based on these extracted features. The model is trained iteratively using optimization techniques such as backpropagation and gradient descent, adjusting its parameters to minimize the discrepancy between generated and ground truth captions. This process allows the model to learn the intricate associations between visual content and textual descriptions, enabling it to produce accurate and coherent captions for unseen images.

Testing Following training, the performance of the Image Caption Generator model is assessed through rigorous testing procedures. A separate subset of the dataset, distinct from the training data, is reserved for testing purposes to ensure unbiased evaluation. The trained model is deployed to generate captions for unseen images in the test set, and the quality of the generated captions is evaluated against human-written references using metrics such as BLEU. This testing phase serves to validate the generalization capability of the model and assess its performance in real-world scenarios, providing insights into its effectiveness and areas for improvement.

Model Evaluation In the final stage of the project, the performance of the Image Caption Generator model is comprehensively evaluated. Evaluation metrics such as BLEU are employed to quantify the quality, fluency, and relevance of the generated captions compared to human-authored references. The model's ability to capture semantic meaning, linguistic diversity, and contextual relevance is assessed, providing valuable insights into its strengths and limitations. Additionally, qualitative analysis and user feedback may be solicited to gain further understanding of the model's performance and refine its capabilities. This evaluation process ensures that the Image Caption Generator meets quality standards and effectively bridges the gap between visual perception and linguistic expression.

3.2 MATERIALS AND METHODS

3.2.1 Tools

1. **Google Colab :** This project was done using a cloud-based platform provided by Google that allows to write and execute Python code in a Jupyter Notebook environment directly on Google's servers. It offers free access to computing resources, including CPU, GPU, and TPU, as well as integration with Google Drive for storing and sharing notebooks.
2. **TensorFlow and Keras:** This is a deep learning framework that provide efficient implementations of CNNs and LSTMs, along with other neural network components. It offers high-level APIs for building and training complex models, making them popular choices for implementing image captioning systems.
3. **CPU: Intel(R) Core(TM) i3-10110U CPU @ 2.10GHz:** The computational power for this study was provided by an Intel(R) Core(TM) i3-10110U CPU running at 2.10GHz. The CPU played a vital role in model training, feature extraction, and the overall execution of the deep learning algorithms. The efficiency of the CPU contributes to the timely processing of image data and the training of complex neural network architectures.

3.2.2 Flickr8k dataset

The Flickr8k dataset consists of 8091 images, with five captions provided for each image. Each caption offers a clear description of the entities and events present in the image. The dataset depicts a variety of events and scenarios and doesn't include images containing well-known people and places, which makes the dataset more generic. This dataset is available on Kaggle and has a size of 1GB. For image caption generation the dataset is split into 7282 images for training, 809 for validation, and 809 for testing.

3.2.3 Architectures employed for image caption generation

3.2.3.1 CNN

A Convolutional Neural Network (CNN) is a powerful deep learning architecture used for analyzing grid-structured data like images. It comprises three key layers: convolutional, pool-

ing, and fully-connected layers. Convolutional layers extract features from input data by applying filters that detect edges, corners, and textures, generating feature maps. Pooling layers reduce feature map size, conserving computational resources by combining neuron cluster outputs into single neurons. Fully-connected layers link all neurons between layers, performing final classification based on extracted features. Additional components include dropout layers to prevent overfitting by randomly deactivating neurons during training, and activation functions to introduce non-linearity for learning complex relationships. An example for CNN is LeNet-5, consists of seven layers including convolutional, pooling, and fully-connected layers, used for basic image classification and often serves as a foundation for more sophisticated models.

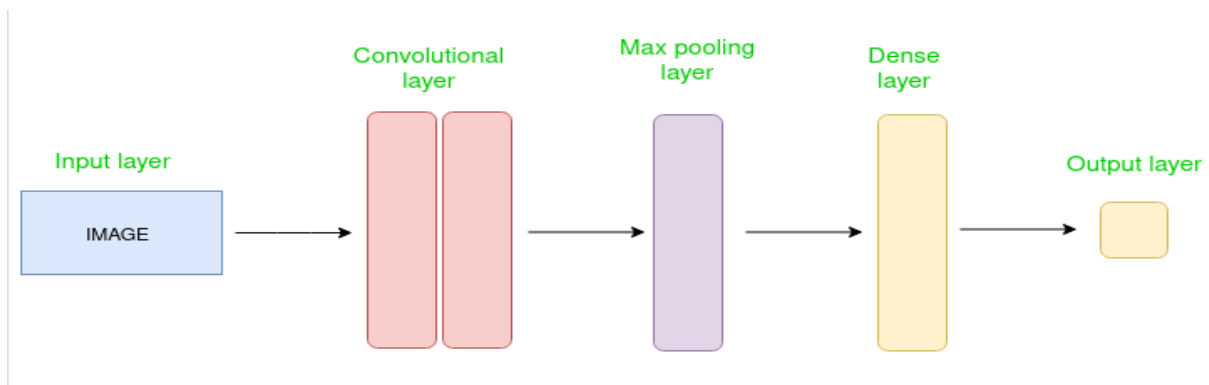


Figure 3.2: CNN Architecture^[10]

3.2.3.2 VGG16

Vgg16 is a widely used CNN architecture trained on the ImageNet database. It was developed by Karen simonyan and Andrew zissarman. The 16 in vgg16 implies that the architecture has 16 layers. It has a 92.7% top 5 accuracy on the ImageNet database which is a dataset containing over 14 million images in 1000 categories. VGG16 is used in many deep learning image classification techniques and is popular due to its ease of implementation.

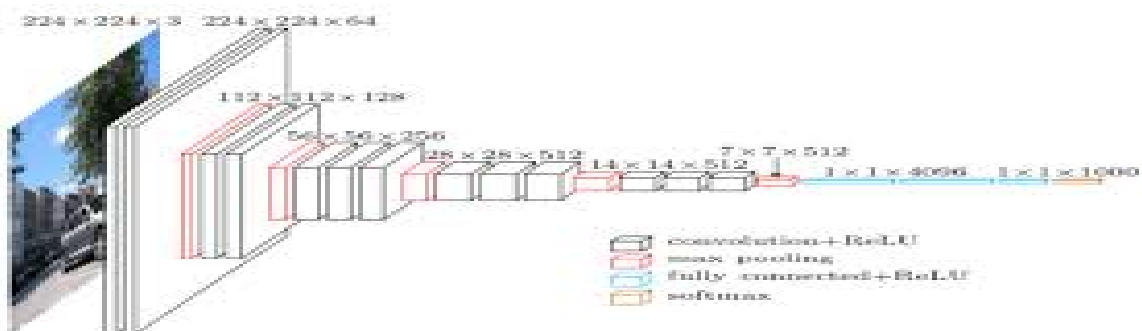


Figure 3.3: VGG16 Architecture^[11]

3.2.3.3 LSTM

The Long Short-Term Memory (LSTM) architecture is comprised of three major gates: the Forget gate, the Input gate, and the Output gate.

- **Forget gate:** This gate is responsible for filtering data by deciding what information to discard from the cell state. It evaluates the relevance of the current input and the previous cell state to determine which information is important for the current prediction task and which can be safely forgotten. By selectively discarding unnecessary information, the forget gate helps optimize the performance of the LSTM model.
- **Input gate:** The input gate marks the beginning of the LSTM process. It controls the flow of new input data into the cell state. The gate evaluates the incoming data and determines which parts are relevant to the current task. It then selectively updates the cell state with this new information, allowing the LSTM to adapt to the input data over time.
- **Output gate:** Once the LSTM has processed the input data and updated its internal state accordingly, the output gate determines how to present the final result. It regulates the flow of information from the cell state to the output of the LSTM network. By controlling the output in a proper manner, the output gate ensures that the desired result is effectively conveyed to downstream layers or output nodes.

These three gates work together to enable LSTMs to store information for longer periods and overcome the limitations of traditional Recurrent Neural Networks (RNNs).

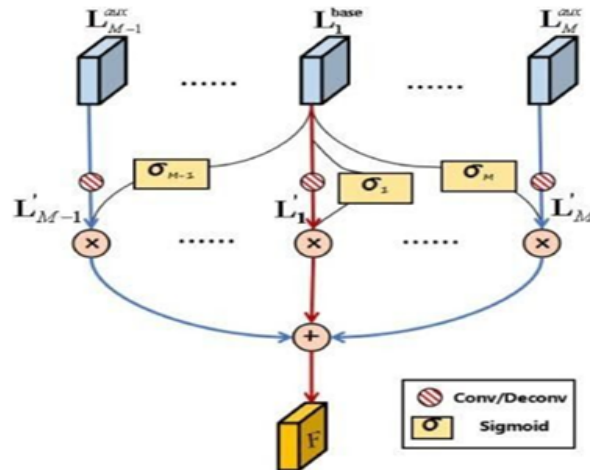


Figure 3.4: Gates in LSTM ^[12]

3.2.3.4 CNN-LSTM Model

In order to prepare an image caption generation model, we will be summing up the two different architectures. It is further called as CNN-LSTM model. So, in this we will be using these two architectures to get the caption for the input pictures.

CNN - It has been used to extract the important features from the input picture. To do this, we have taken a pre-trained model for our consideration named Xception.

LSTM - It has been used to store the data or the features from the CNN model and further process it and to support in the generation of a good caption for the picture.

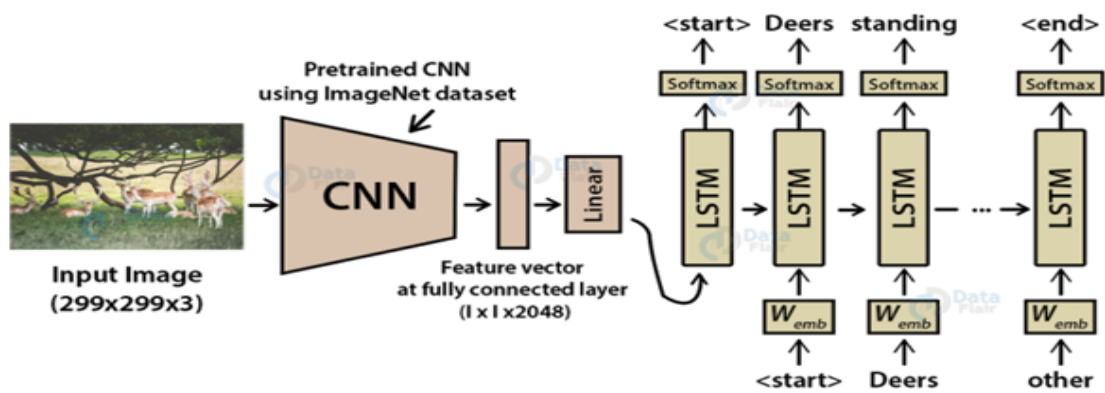


Figure 3.5: CNN-LSTM model ^[13]

Chapter 4

RESULTS AND ANALYSIS

This chapter presents the results of image captioning model. It showcases various outcomes obtained through tasks like data cleaning, feature extraction, Tokenization, model creation and model validation.

4.1 Result

The result of image caption generation using a combination of CNN and LSTM architectures shows that the model can create captions that closely match real captions and understand how words in a sentence relate to each other. The captions generated for images in the dataset closely resembled the actual captions provided in the Flickr8k dataset also the model's performance extends beyond dataset images, as it generated accurate captions for external images. The performance of image captioning model was evaluated using BLEU scores, a standard metric for assessing the quality of generated captions. It serve as metrics for evaluating the similarity between the generated captions and the reference captions. A BLEU-1 score of 0.54 indicates relatively good performance in capturing individual words within the generated captions. However, the lower BLEU-2 score suggests there's room for improvement in accurately generating longer phrases that effectively reflect the content depicted in the images. The project widely-used Flickr8K dataset as a benchmark for both training and evaluating the image caption generation model. Notably, evaluation methodology involved comparing the model's generated caption for a given image with corresponding five captions in the dataset. These results demonstrate the effectiveness of the model in generating accurate and contextually relevant captions for a diverse range of images. Despite the overall success, further research is needed to address remaining challenges and improve the robustness of the model in handling complex visual scenes. The findings highlight the potential of the approach to advance the field of image captioning and facilitate multimodal learning applications.

4.2 Analysis

This section thoroughly assesses the implemented image captioning system, covering essential stages from data cleaning to model validation. It evaluates the effectiveness of each phase, including text standardization, feature extraction, tokenization, and architectural composition. Model validation, through BLEU scores, offers insights into caption accuracy and coherence.

Overall, this analysis provides a clear understanding of the system’s performance and its potential practical use.

4.2.1 Data Cleaning

Flickr8k dataset contains multiple descriptions described for a single image. In the data preparation phase, each image id is taken as key and its corresponding captions are stored as values in a dictionary. In order to make the text dataset work in machine learning or deep learning models, raw text should be converted to a usable format. The following text cleaning steps are done before using it for the project: Removal of punctuations, Removal of numbers, Removal of single length words, Conversion of uppercase to lowercase characters. Table 4.1 shows samples of captions after data cleaning.

Table 4.1: Data cleaning of captions

Original Captions	Captions after Data Cleaning
A child in a pink dress is climbing up a set of stairs in an entry way .	startseq child in pink dress is climbing up set of stairs in an entry way endseq
A girl going into a wooden building.	startseq girl going into wooden building endseq
A little girl climbing into a wooden playhouse .	startseq little girl climbing into wooden playhouse endseq
A little girl climbing the stairs to her playhouse .	startseq little girl climbing the stairs to her playhouse endseq
A little girl in a pink dress going into a wooden cabin .	startseq little girl in pink dress going into wooden cabin endseq

4.2.2 Feature Extraction

The VGG16 model is utilized as a feature extractor to obtain high-level visual representations from input images. By removing the model’s last classification layer, feature vectors of length 4096 was obtained, capturing rich visual information crucial for generating descriptive captions. The model, originally trained on the ImageNet dataset, requires input images to be resized to (224, 224) dimensions. Preprocessing steps, including converting images to arrays and normalizing them using the ‘preprocess_input’ function, ensured compatibility with the VGG16 model. Extracted features were stored in a dictionary structure, mapping image IDs to their corresponding feature vectors, facilitating the association with captions during model training. While feature extraction demanded significant computational resources and time, the

effectiveness of the VGG16-extracted features in generating accurate and contextually relevant captions validated their utility for image captioning task.

4.2.3 Tokenization

In order to facilitate the processing of textual data by image captioning model, the tokenizer function provided by the Keras library to tokenize vocabulary is used. Given that computers understand numerical data rather than English words, tokenization involves mapping each word in the vocabulary to a unique index value. This process allows the model to interpret textual input as numerical data, enabling it to learn associations between image features and textual descriptions. The resulting tokenized vocabulary was saved to a "tokenizer.pkl" pickle file for future use. The vocabulary consisted of 8485 words, each mapped to a specific index value. Additionally, calculated the maximum length of the descriptions, which plays a crucial role in determining the model structure parameters. Here the maximum length of descriptions was determined to be 35 tokens, providing guidance for the design and configuration of our image captioning model.

4.2.4 Model Architecture

The model architecture employed for text processing comprises several key layers, each serving a specific purpose in the data processing pipeline. The system takes two inputs, likely numerical representations of text data, and proceeds through a series of layers:

- **InputLayer:** Responsible for preprocessing the raw input data and preparing it for subsequent layers.
- **Embedding:** Converts the numerical representation of text data into a denser vector representation using pre-trained word embedding models such as Word2Vec or GloVe. This allows the model to capture semantic relationships between words.
- **Dropout:** Randomly drops a certain percentage of elements from the activation tensor to prevent overfitting during training, enhancing the model's generalization capability.

- **LSTM (Long Short-Term Memory):** A type of recurrent neural network (RNN) specifically designed to handle sequential data, such as text. The stacked LSTM layers enable the model to learn long-term dependencies in the data, facilitating more accurate predictions.
- **Dense:** A fully-connected layer that maps the data to a lower-dimensional space, enabling further processing and feature extraction.
- **Add:** Combines the outputs of the two LSTM layers, possibly enhancing the model's representation capabilities by capturing complementary information from different layers.
- **Dense_1 and Dense_2:** Additional fully-connected layers used for further processing and refinement of the data.

The final output of the system is a dense vector of size (None, 8485), likely representing a probability distribution over a vocabulary of 8485 words. This output vector can be utilized for various natural language processing (NLP) tasks, including machine translation and text summarization. This architecture demonstrates the system's capability to process large volumes of text data, extract meaningful insights, and perform tasks such as machine translation and text summarization. The fig 4.1 shows the model architecture.

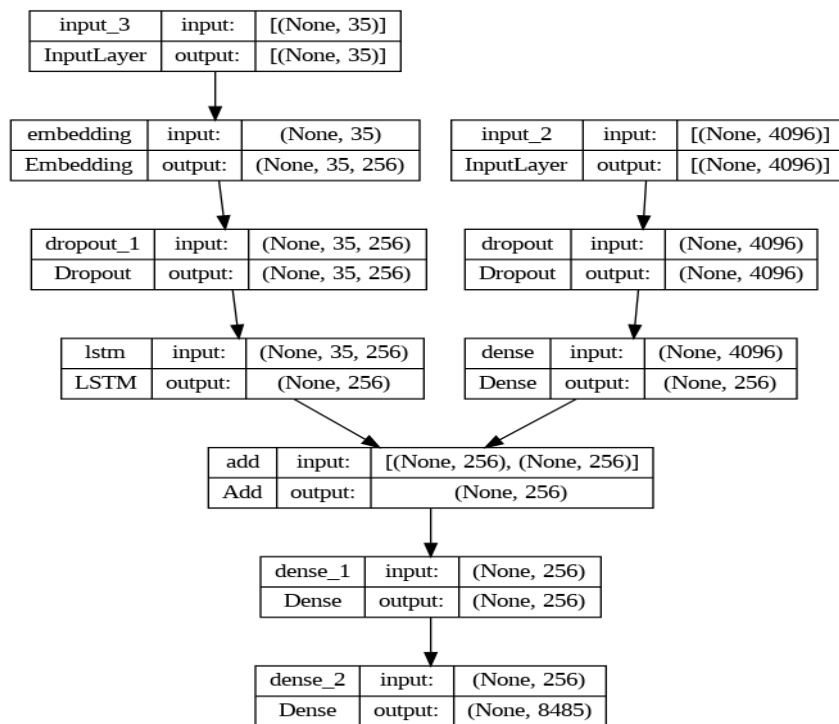


Figure 4.1: Model architecture

4.3 Model Validation

During the model validation process, the image captioning model was trained over multiple epochs, tracking its performance by calculating BLEU scores at the end of each epoch. The final BLEU-1 and BLEU-2 scores are 0.537540 and 0.312828 by training on 40 epochs. It provides insights into the performance of the image caption generation model. BLEU scores serve as metrics for evaluating the similarity between the generated captions and the reference captions. A BLEU-1 score of 0.54 indicates relatively good performance in capturing individual words within the generated captions. However, the lower BLEU-2 score suggests there's room for improvement in accurately generating longer phrases that effectively reflect the content depicted in the images. Through this validation process, it aimed to ensure that the model produces accurate and coherent captions for a diverse range of images, enhancing its utility in real-world applications such as image captioning systems and assistive technologies. The table 4.2 shows different BLEU scores obtained over different epochs.

Table 4.2: BLEU Scores

Epoch	BLEU-1	BLEU-2
10	0.318014	0.149546
20	0.392435	0.203196
30	0.520644	0.310496
40	0.537540	0.312828

Chapter 5

CONCLUSION

In conclusion, the project aimed to address the challenge of image caption generation using a combination of Convolutional Neural Networks (CNN) for feature extraction and Long Short-Term Memory (LSTM) networks for sequence generation. By training the model on the Flickr8k dataset and validating its performance through BLEU scores, the project effectively showcased the model's capability to generate captions that closely match ground truth captions. The project contributes to the advancement of computer vision and natural language processing by providing a robust solution for automatically generating captions for images. This technology has broad applications in various domains, including assistive technologies for visually impaired individuals, content indexing for image databases, and enhancing user experiences in photo-sharing platforms. However, it's important to acknowledge some limitations of the project. One limitation is the reliance on pre-existing datasets, which may not fully represent the diversity of real-world images. Additionally, the model's performance may be affected by factors such as image quality, complexity, and context, which could lead to inaccuracies in caption generation. Further research and improvements in model architecture and training methodologies are needed to address these limitations and enhance the accuracy and robustness of image caption generation systems.

Chapter 6

FUTURE SCOPE

In considering the future scope of image caption generator, several avenues for research and development emerge, presenting opportunities to enhance its capabilities and applicability. One promising direction involves the integration of additional modalities, such as audio or text descriptions, into the existing model architecture. This multi-modal approach could enrich the generated captions with more comprehensive and contextually rich information, thereby improving their utility across diverse applications. Furthermore, there is scope to refine the model's captioning precision by focusing on fine-grained details and nuances present in images. This could be achieved through the exploration of datasets like Flickr30k, with more diverse and detailed annotations, enabling the model to produce captions with greater accuracy and specificity. Additionally, incorporating attention mechanisms within the LSTM network holds promise for dynamically focusing on different regions of the image during caption generation, thereby enhancing the alignment between visual features and words in the captions. Moreover, advancements in long-term dependency modeling techniques could further improve the model's ability to capture complex relationships within input sequences, contributing to the generation of more coherent and contextually relevant captions. Expanding the model's applicability to cross-domain captioning tasks, such as medical imaging or satellite imagery, represents another fruitful avenue for exploration, offering opportunities to address domain-specific challenges and requirements. Furthermore, the development of interactive captioning systems, where users can provide feedback on generated captions for iterative refinement, holds potential for enhancing user engagement and satisfaction. Finally, exploring alternative evaluation metrics beyond traditional BLEU scores, such as CIDEr and SPICE, could provide a more comprehensive assessment of caption quality, leading to insights that further refine the model's performance. By pursuing these future research directions, the image caption generator can continue to evolve and contribute to advancements in computer vision and natural language processing, with broad implications for diverse real-world applications.

BIBLIOGRAPHY

- [1] M. Israk Ahmed et al., “Context-based Image Caption using Deep Learning,” in *Proceedings of the International Conference on Intelligent Computing and Signal Processing (ICSP)*, vol. 2021, pp. 1-34, 2021, doi: <https://doi.org/10.1109/ICSP51882.2021.9408871>.
- [2] Kelvin Xu et al., “Show and tell: A neural image caption generator,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, vol. 2015, pp. 677-686, 2015, doi: <https://doi.org/10.1109/CVPR.2015.729>.
- [3] Peter Anderson et al., “Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering,” in *Computer Vision and Pattern Recognition*, vol. 2018, pp. 1-15, 2018, doi: <https://doi.org/10.48550/arXiv.1707.07998>.
- [4] Jamin Young et al., “Image Captioning with Semantic Attention,” in *Proceedings of a Computer Vision Conference*, vol. 2022, pp. 1-9, 2022.
- [5] chen2014image Jianhui Chen, Wenqiang Dong, and Minchen Li, “Image caption generator based on deep neural networks,” 2014. Available online: <chrome-extension://efaidnbmnnnibpca/jpcglclefindmkaj/https://ijcrt.org/papers/IJCRT2112031.pdf>.
- [6] Kelvin Xu, Jimmy Ba, Rasmus Kjeldgaard, and Arun Teller, “Show and tell: A neural image caption generator,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, vol. 2015, pp. 677-686, 2015, doi: <https://doi.org/10.1109/CVPR.2015.729>.
- [7] Oriol Vinyals, Alexander Toshev, and Samy Bengio, “Show and attend: Neural image caption generation with global context attention,” in *Proceedings of the IEEE international conference on computer vision*, vol. 2015, pp. 3156-3164, 2015, doi: <https://doi.org/10.1109/ICCV.2015.391>.
- [8] Andrej Karpathy and Li Fei-Fei, “Deep visual-semantic alignments for generating image descriptions,” *IEEE transactions on pattern analysis and machine intelligence*, vol. 39, no. 8, pp. 1838-1853, 2017, doi: <https://doi.org/10.1109/TPAMI.2016.2598817>.

- [9] Jeff Donahue, Lisa Anne Hendricks, Ross Girshick, and Trevor Darrell, “Long short-term memory networks for machine reading,” in *Proceedings of the 31st International Conference on Machine Learning*, vol. 1, pp. 1150-1158, 2015.
- [10] ”CNN architecture” <https://www.geeksforgeeks.org/introduction-convolution-neural-network/>
- [11] Vgg16,”<https://towardsdatascience.com/step-by-step-vgg16-implementation-in-keras-for-beginners-a833c686ae6c>
- [12] Gates in LSTM,”<https://www.sciencedirect.com/science/article/abs/pii/S0031320319304327>
- [13] CNN-LSTM model <https://www.ijraset.com/research-paper/image-captioning-generator-using-cnn-and-lstm>
- [14] Changwan Zhou, Yifan Sun, Yuxuan Liu, Zhuang Liu, and Jun Tang, “Attention-based lstm for collaborative filtering,” in *Proceedings of the 26th ACM Conference on Information and Knowledge Management*, pp. 1417-1426, 2017.

APPENDIX

Sample Code

Model Creation

```
1 # encoder model
2 # image feature layers
3 inputs1 = Input(shape=(4096,))
4 fe1 = Dropout(0.4)(inputs1)
5 fe2 = Dense(256, activation='relu')(fe1)
6 # sequence feature layers
7 inputs2 = Input(shape=(max_length,))
8 se1 = Embedding(vocab_size, 256, mask_zero=True)(inputs2)
9 se2 = Dropout(0.4)(se1)
10 se3 = LSTM(256)(se2)
11
12 # decoder model
13 decoder1 = add([fe2, se3])
14 decoder2 = Dense(256, activation='relu')(decoder1)
15 outputs = Dense(vocab_size, activation='softmax')(decoder2)
16
17 model = Model(inputs=[inputs1, inputs2], outputs=outputs)
18 model.compile(loss='categorical_crossentropy', optimizer='adam')
19
20 # plot the model
21 plot_model(model, show_shapes=True)
```

Generate captions for image

```
1  def idx_to_word(integer, tokenizer):
2      for word, index in tokenizer.word_index.items():
3          if index == integer:
4              return word
5      return None
6  def predict_caption(model, image, tokenizer, max_length):
7      in_text = 'startseq'
8      for i in range(max_length):
9          sequence = tokenizer.texts_to_sequences([in_text])[0]
10         sequence = pad_sequences([sequence], max_length)
11         yhat = model.predict([image, sequence], verbose=0)
12         yhat = np.argmax(yhat)
13         word = idx_to_word(yhat, tokenizer)
14         if word is None:
15             break
16         in_text += " " + word
17         if word == 'endseq':
18             break
19     return in_text
```

Project Screenshots

Caption generated for an image in dataset

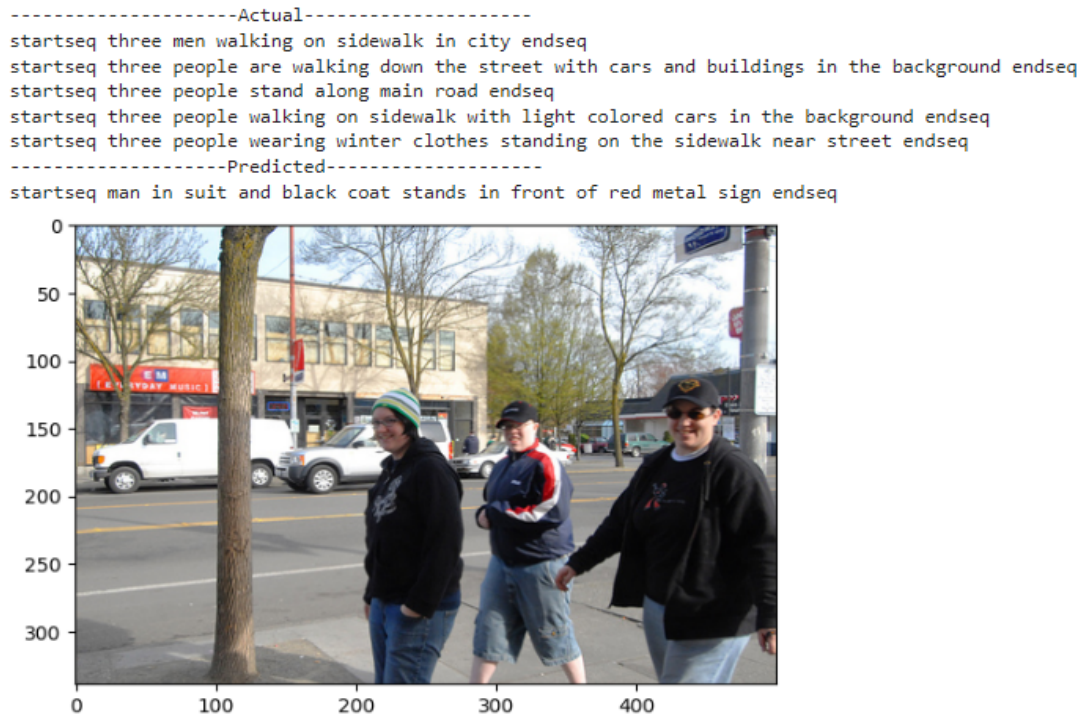


Figure 6.1: Actual and Predicted captions

Caption generated for external image



Generated Caption: two dogs are playing together outside

Figure 6.2: Caption for an image

GitHub History

Bitbucket

Your work

Pull requests

Repositories

Projects

People

More

Create

Q Search

?

⚙

AP

</> imagecaption

</> Source

Commits

Branches

Pull requests

Pipelines

Deployments

Jira issues

Security

Downloads

Repository settings

student student / Imagecaption / imagecaption

Commits

Compare

Clone

Search commits

All branches

Author	Commit	Message	Date
Ashna C P	642ae2c	new created online with Bitbucket	2024-03-31
Ashna C P	e9fdb98	successfully committed	2024-03-17
Ashna C P	e8a6e7d	successfully committed	2024-03-17
Ashna C P	39b8618	added system architecture	2024-03-17
Ashna C P	abaa00	added enco_deco architecture	2024-03-17
Ashna C P	f8b66f6	added dfd	2024-03-17
Ashna C P	86596a4	successfully committed	2024-03-17
Ashna C P	44c73dc	successfully committed	2024-02-25
Ashna C P	0ec6b26	successfully committed	2024-02-25
Ashna C P	606d500	Initial commit	2024-02-17

Q Search

ENG IN

3:32 PM 11-Apr-24

Figure 6.3: GitHub History