# Image Captioning

Submitted for Summer Internship
On
**Python & Machine Learning**
(from 5th June 2023 - 23rd July 2023)

Jointly conducted by
COE- AI, AI Club IGDTUW and Anveshan Foundation

By

Khushi Sehrawat
09101012022
B.Tech CSE
2022-26

Navya Raj
12001012022
B.Tech CSE
2022-26

**INDIRA GANDHI DELHI TECHNICAL UNIVERSITY FOR WOMEN**
(Established by Govt. of Delhi vide Act 09 of 2012)
Kashmere Gate, Delhi- 110006

# INDEX

# CERTIFICATION



**INDIRA GANDHI DELHI TECHNICAL UNIVERSITY FOR WOMEN**

(ESTABLISHED BY GOVT. OF DELHI VIDE ACT 09 OF 2012)
ISO 9001:2015 CERTIFIED UNIVERSITY
KASHMERE GATE, DELHI-110006

WOMEN EDUCATION | WOMEN ENLIGHTENMENT | WOMEN EMPOWERMENT

**CENTRE OF EXCELLENCE - ARTIFICIAL INTELLIGENCE**

**CERTIFICATE OF COMPLETION**

This certificate is awarded to

**Khushi Sehrawat**

For successfully completing the 7 weeks Summer Internship on
**"PYTHON & MACHINE LEARNING"** from **5th June - 23rd July, 2023** jointly
conducted by the COE - AI, AI Club IGDTUW and Anveshan Foundation.

Ishita Saxena
President - AI CLUB
IGDTUW

Dr. Ritu Rani
Research Associate
COE - AI

Prof. Arun Sharma
Coordinator - Centre of Excellence-AI
IGDTUW

---

**INDIRA GANDHI DELHI TECHNICAL UNIVERSITY FOR WOMEN**

(ESTABLISHED BY GOVT. OF DELHI VIDE ACT 09 OF 2012)
ISO 9001:2015 CERTIFIED UNIVERSITY
KASHMERE GATE, DELHI-110006

WOMEN EDUCATION | WOMEN ENLIGHTENMENT | WOMEN EMPOWERMENT

**CENTRE OF EXCELLENCE - ARTIFICIAL INTELLIGENCE**

**CERTIFICATE OF COMPLETION**

This certificate is awarded to

**Navya Raj**

For successfully completing the 7 weeks Summer Internship on
**"PYTHON & MACHINE LEARNING"** from **5th June - 23rd July, 2023** jointly
conducted by the COE - AI, AI Club IGDTUW and Anveshan Foundation.

Ishita Saxena
President - AI CLUB
IGDTUW

Dr. Ritu Rani
Research Associate
COE - AI

Prof. Arun Sharma
Coordinator - Centre of Excellence-AI
IGDTUW

# DECLARATION

We, **Khushi** and **Navya**,  students of **Indira Gandhi Delhi Technical University**, enrolled in the **Bachelor of Technology Computer Science and Engineering**, declare that the research paper titled **"Learning Neural Image Captioning"** is the result of my own work conducted during the summer internship training jointly organized **by COE-AI Club IGDTUW and Anveshan Foundation**. This internship, focused on **"Python & ML"**, provided us with the opportunity to delve into the intricate realm of machine learning and artificial intelligence.

We affirm that this research paper represents our original work, and to the best of our knowledge, no part of this document has been submitted for any other degree or qualification. All sources used in this paper are duly acknowledged, and any contributions from individuals or organizations are appropriately recognized.

We further declare that the ideas and findings presented in this paper are solely our own, except where explicitly stated otherwise. Any external assistance received during the course of this research is acknowledged, and the contributions of others, if any, are clearly mentioned in the appropriate sections of the paper.

We understand the importance of academic integrity and am committed to upholding the ethical standards associated with scholarly work. We are aware that any violation of these standards may result in the nullification of my research paper and the consequences as stipulated by the academic policies of Indira Gandhi Delhi Technical University for Women.


**Date: 13-11-2023**

# ACKNOWLEDGMENT

# LIST OF FIGURES

6

| 3 |  | Overall System flow of proposed architecture | 17 |
|---|---|---|---|
| 4 |  | Overview of the proposed model | 17 |

7

# LIST OF TABLES

| Image ID | Image | Actual caption | Predicted caption |
|---|---|---|---|
| 1095590286_c654f7e5a9 |  | startseq blond dog and black and white dog run in dirt field endseq | startseq two dogs are playing with each other on the ground endseq |
| 1034276567_49bb87c51c |  | startseq boy bites hard into treat while he sits outside endseq | startseq man bites baby with his hands endseq |
| 2084217208_7bd9bc85e5 |  | startseq "a person in blue jacket wearing bicycle helmet is riding bike" endseq | startseq man in blue jacket riding bike endseq |
| 1989145280_3b54452188 |  | startseq blond girl with sunglasses on her head looks squeamish endseq | startseq woman with sunglasses on her face is holding onto cellphone endseq |
| 1778020185_1d44c04dae |  | startseq brown dog runs for white and black dog on the grass endseq | startseq dog is running through the grass endseq |
| 1429546659_44cb09cbe2 |  | startseq white dog and black dog in field endseq | startseq two dogs playfully wrestle in the grass endseq |
| 1330645772_24f831ff8f |  | startseq black and white dog is running in the grass endseq | startseq black and white dog is running through field endseq |
| 1143373711_2e90b7b799 |  | startseq bicycle rider is crossing street endseq | startseq man is riding bike on city street endseq |
| 1015584366_dfcec3c85a |  | startseq black dog leaps over log endseq | startseq black dog is climbing over log endseq |

| | | | |
|---|---|---|---|
| 109202801_c6381eef15 |  | startseq two draft horses pull cart through the snow endseq t endseq | startseq woman with red and red colored coat rides horse endseq |

# ABSTRACT

In the realm of artificial intelligence, the integration of Computer Vision and Natural Language has paved the way for topics like Image Captioning. This complex process seeks to enable machines to understand visual information and communicate it in a comprehensible, human-like manner. In this work, we explored the implementation and assessment of an Image Caption Generator by using the combination of Convolution Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks. CNN extracts feature of an image and LSTM generates text. The model is trained to generate captions that fairly depict the contents of the input image. Our model is tested on the widely used dataset of Flickr8k, which contains diverse images, each containing five distinct captions that provide extensive overviews of the key things and events shown in the pictures. The results of the BLEU Score "BLEU1" (0.59863) and "BLEU2" ("59863") demonstrate the model's ability to incorporate linguistic innovations into visual domains.

*Keywords—Image Caption Generator, CNN, LSTM, Flickr8k Dataset, Feature Extraction*

# CHAPTER 1: INTRODUCTION

In the field of cognitive processing, our brains possess the remarkable capacity to comprehend the meaning behind every visual stimulus we encounter. However, when it comes to computers, the question arises as to how they can analyze an image and produce a caption that is not only highly pertinent but also precise. Not long ago, such a prospect seemed almost incomprehensible. Nevertheless, significant progress in the domain of Computer Vision and Deep Learning algorithms, alongside the accessibility of pertinent datasets and AI models, has made the development of an effective image caption generator more achievable than ever before. Image captioning is an important task, applicable to virtual assistants, editing tools, image indexing, and support for the disabled. It is a cutting-edge area of computer vision that has attracted considerable attention for its potential to improve human-machine interaction and allow intelligent systems to understand visual information. The challenge lies in creating descriptive and cohesive captions for an image, thereby bridging the semantic divide between visual and textual content.

From autonomous driving, robotics, medical imaging, and multimedia content analysis, this technology has the power to completely alter how machines perceive the environment. Given the huge number of images available online, image captioning has become central to image retrieval tasks such as the ones carried out by search engines or newspaper companies. More specific applications, like describing images for blind persons or teaching children's concepts, can also be given as examples of the importance of captioning images [6].

The proposed system can be deployed as a mobile application where the user takes a photograph as an input image and the system generates the description of the image.  This will be helpful as a smart assistant for children, especially visually challenged people. The proposed system can also be deployed on robots to help them understand their surroundings. It can even be used in an application to identify famous objects for a tourist who goes without a tour guide. The user can capture the monument or artifact in his mobile phone camera and the application displays relevant information in the form of text [5].

The noteworthy breakthroughs in neural machine translation (NMT) and the widespread availability of large datasets are the primary reasons for the current spike in interest in image captioning. These advancements have sped up the use of encoder-decoder models that make use of RNNs or LSTM networks. However, despite LSTM networks' benefits, they are still constrained by the finite memory of a finite number of time steps, which causes the information to gradually dwindle. This research uses the effective implementation of Convolutional Neural Networks (CNN) in NLP for text representation to overcome these issues. A novel method is presented by including a language CNN that uses temporal convolution for sequence feature extraction. This modified model differentiates from the standard CNN design by skipping the pooling processes, saving crucial data for precise word representation.

Our research mainly focuses on integrating a language CNN with LSTM for picture captioning that is adept at capturing long-range relationships in sequences. In-depth examinations of the frequently used Flickr8k dataset show promising findings that are compatible with contemporary techniques. In this research, the proposed model is thoroughly analyzed with an emphasis on how well it produces captions that are contextually relevant. The effective combination of CNN and LSTM in picture captioning tasks reveals the possibility of improving and expanding the interface between

computers and visual material, discovering new opportunities for a broad range of applications across domains.

# CHAPTER 2: LITERATURE SURVEY

Image captioning, a captivating intersection of Computer Vision and Natural Language Processing (NLP), has experienced remarkable advancements in recent years, thanks to pioneering research contributions that have significantly shaped the field.

One groundbreaking work, referenced as [1], introduces a Generative CNN-LSTM model that seamlessly combines Convolutional Neural Networks (CNN) with Long Short-Term Memory (LSTM) networks. This model represents a significant milestone, transcending human baseline performance by an impressive 2.7 BLEU-4 points on the MSCOCO dataset. Notably, this approach has set a high benchmark for subsequent research in the domain, laying the foundation for the exploration of CNNLSTM synergy.

In contrast, [2] undertakes an exhaustive evaluation of 17 distinct CNN architectures tailored for image captioning, revealing a captivating revelation: model complexity, as quantified by the number of parameters, does not linearly correlate with efficacy in feature extraction for caption generation. This study underscores the need for a nuanced understanding of CNNs in the specific context of image captioning. It highlights that optimal architectural choices are often non-intuitive and suggests that the quest for improved performance necessitates a comprehensive exploration of CNN architectures.

Addressing the inherent sequential nature of LSTM-based models, [3] pioneers Convolutional Image Captioning, an innovative approach that demonstrates competitive performance on the MSCOCO dataset while significantly reducing training times per parameter compared to traditional LSTM baselines. This research emphasizes the advantages of convolutional networks in efficiently handling image captioning tasks. The efficient processing of visual features plays a pivotal role in addressing the real-time demands of caption generation applications.

Moreover, it is noteworthy to highlight our own model's distinctive contributions to this landscape. Our approach incorporates a language CNN that is adept at capturing long-range dependencies in sequences, uniquely complementing LSTM networks for image captioning. Empirical evaluations on the widely used Flickr8k dataset demonstrate promising results, comparable to state-of-the-art approaches. This model's innovation lies in its ability to effectively model sequences consistently and contextually, bridging the gap between images and human-like language. This novel fusion of CNN and LSTM architectures stands as a testament to the ever-evolving landscape of image captioning.

Furthermore, [4] directs its focus toward a relatively underexplored domain - Arabic image captioning. The study pioneers the development of a generative merge model that combines RNNLSTM and CNN architectures, indicating promising results for Arabic image captioning. It offers a glimpse into the possibilities of leveraging advanced techniques for non-mainstream languages. However, its full potential awaits further exploration with a larger and more diverse corpus.

These cumulative research endeavors collectively advance our comprehension of image captioning. Each approach, with its unique insights and methodologies, not only expands the boundaries of what is achievable in image caption generation but also paves the way for enhanced applications across various domains. The evolution of image captioning continues to be fueled by these and other innovative contributions, promising a future where machines will excel in understanding and communicating the visual world with ever-increasing sophistication, with our model's distinctive fusion of CNN and LSTM architectures standing as a notable addition to this landscape.

# CHAPTER 3: OBJECTIVE

The primary objective of this research is to implement and assess a neural network architecture that integrates Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks for the task of image captioning. The model leverages the VGG16 pre-trained CNN as a feature extractor to capture high-level visual features from images, and LSTM networks are employed to generate sequential captions based on these extracted features. The training process utilizes the diverse Flickr8k dataset, comprising images with multiple captions, to enhance the model's capability to produce varied and contextually relevant image descriptions. Performance evaluation is conducted using BLEU scores, specifically focusing on BLEU-1 and BLEU-2, as metrics to measure the alignment between the generated captions and the actual captions in the dataset. Additionally, the research explores potential real-world applications of the developed image captioning model, such as deployment as a mobile application for visually challenged individuals, assisting children, or aiding tourists in object identification.

The analysis delves into the limitations and challenges faced by the model, including instances where generated captions deviate from actual captions, with the aim of proposing areas for improvement. The study also includes a comparative analysis with existing image captioning approaches, highlighting the distinctive contributions of the proposed CNN-LSTM architecture. Furthermore, the research scrutinizes potential biases in the model's predictions and evaluates the impact of biases present in the training dataset on the diversity and accuracy of generated captions. Ultimately, the study aims to provide insights into potential avenues for future research, including enhancements to model architecture, exploration of larger and more diverse datasets, and innovative approaches to enhance the precision of image caption generation.

Moreover, the research seeks to gain a comprehensive understanding of the model's biases and their implications on the generated captions. By analyzing instances where predictions deviate from the ground truth, the study aims to uncover potential sources of bias, whether cultural, contextual, or inherent in the dataset. This examination is crucial for ensuring the fairness and generalizability of the image captioning model across diverse scenarios.

In addition to evaluating the model's performance, a comparative study with existing image captioning approaches is conducted. This involves assessing the strengths and weaknesses of the proposed CNN-LSTM architecture in comparison to other established models in the field. By identifying the unique contributions of the developed model, the research aims to contribute to the evolving landscape of image captioning methodologies.

In conclusion, the research endeavors to lay the foundation for future investigations by providing insights into the complex interplay between images and captions. As the study evolves, it aims to minimize the divide between visual and language spheres, paving the way for machines that can articulate stories out of pixels and enhance human-technology interactions unexpectedly. The journey outlined in this research signifies a continuous exploration of the connections between images and language, propelling the field of artificial intelligence toward a future where machines adeptly communicate visual information.

# CHAPTER4: METHODOLOGY AND IMPLEMENTATION

## A. Research Design and Approach

This research paper presents a study that uses a deep-learning approach to generate image captions. The methodology encompasses data collection, preprocessing, model architecture, training, and evaluation. The research follows a sequential process, beginning with collecting the Flickr8k dataset containing images and associated captions. The data is then preprocessed for training a caption generation model.

## B. Preprocessing

Before training the caption generation model, the data undergoes essential preprocessing steps to ensure compatibility with the chosen architecture. This includes Image Preprocessing: Images are resized to a consistent dimension of 224x224 pixels. Pixel values are normalized to the [0, 1] range, enhancing the model's convergence during training.

Caption Preprocessing: Captions are subjected to several text preprocessing stages. These include converting all text to lowercase to ensure uniformity, removing non-alphanumeric characters to simplify the vocabulary, and eliminating excess spaces. Additionally, the captions are augmented with unique tokens like "startseq" and "endseq" to indicate the beginning and end of a caption, aiding the model's understanding of sentence structure.

## C. Model Architecture Design

At the core of this research is the development of an advanced deep-learning architecture that serves as a crucial link between visual and textual data. The chosen architecture comprises two main components.
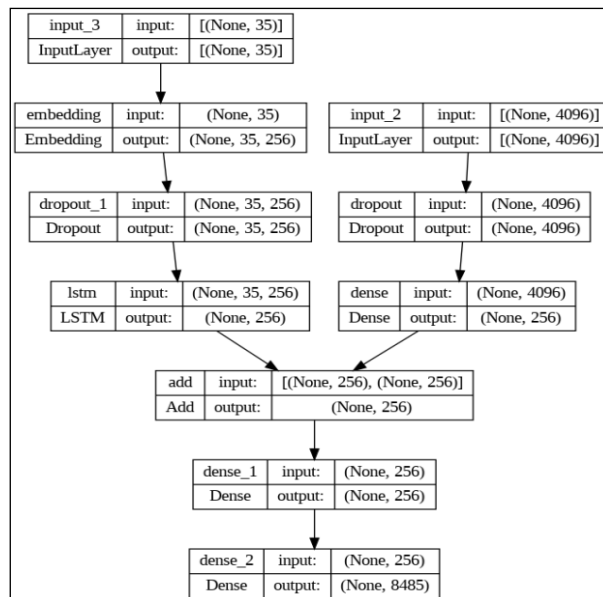


*Figure 1 : Architecture of the Encoder-Decoder Model for Image Captioning*

In this process, we utilize the VGG16 convolutional neural network as a powerful feature extractor, leveraging its pre-trained architecture to extract high-level visual features from images. To make this representation more compact and meaningful, we truncate the VGG16 model, removing its final classification layer. Simultaneously, we employ a Long Short-Term Memory (LSTM) network for the textual component. LSTM is ideal for processing sequential data, making it well-suited for generating word-by-word captions. This LSTM model takes the extracted image features and partial captions as inputs, enabling it to iteratively generate the next word in the sequence, resulting in coherent and contextually relevant image captions.

## D. Training and Evaluation

The training phase involves iteratively updating the model's parameters using the dataset. The dataset is divided into training, validation, and testing sets to prevent overfitting. The model is optimized using the categorical cross-entropy loss function, which measures the dissimilarity between predicted and actual word distributions in the captions. The training process aims to minimize this loss, enabling the model to generate coherent and contextually relevant captions.



Figure 2: Training and Validation Loss During Model Training

The model's performance is evaluated using the BLEU (Bilingual Evaluation Understudy) metric. BLEU evaluates the quality of generated captions by comparing them to reference captions. The BLEU score ranges from 0 to 1, with higher scores indicating better caption quality. This metric provides insights into the model's ability to produce captions that align with human-generated descriptions.

The primary dataset used in this study is the Flickr8k dataset, sourced from the Kaggle platform. This dataset comprises approximately 8,000 images, each with five unique captions describing the visual content. The dataset's diverse array of subjects, scenes, and contexts reflects real-world imagecaption pairs. This diversity is critical for training a model to effectively generate relevant and contextually accurate captions across a broad spectrum of visual inputs.

## E. Justification of Methods and Suitability

The Flickr8k dataset is motivated by its suitability for training a caption generation model. The dataset provides many images with multiple captions per image, allowing the model to learn diverse textual descriptions for different visual contexts.

The utilization of the VGG16 model for feature extraction is well-founded due to its proven effectiveness in image representation. The model's architecture facilitates extracting high-level features from images, subsequently used to train the caption generation model.

Furthermore, the LSTM-based caption generation architecture is selected for its ability to handle sequential data and generate coherent textual descriptions. Combining visual features from images and sequential context from captions contributes to a comprehensive understanding of the relationship between visual and textual elements.

## F. Limitations and Potential Biases

Several limitations and potential biases are acknowledged in the methodology of this study. Firstly, the quality and diversity of captions within the dataset might impact the model's generalization to real-world scenarios. Biases present in the dataset, such as cultural or contextual biases in captions, may inadvertently affect the generated results. The selected model architecture might also not capture complex semantic relationships in certain images, leading to suboptimal captions.



*Figure 3. Overall System flow of proposed architecture*

The preprocessing steps, including text cleaning and tokenization, may inadvertently remove relevant information or introduce artifacts. The choice of hyperparameters, such as batch size and epoch count, could impact the training process and model performance. Finally, the evaluation metrics employed, such as BLEU scores, provide insights into the model's performance but may not fully capture the quality of generated captions.



*Figure 4. Overview of the proposed model*

# CHAPTER 5: RESULT DISCUSSION

Testing and evaluating the CNN-LSTM image caption generator on the Flickr8k dataset produced promising but inconsistent results. The model's capacity to perceive visual context and construct contextually pertinent descriptions is demonstrated by the generated captions for some photos nearly matching the actual captions. In one instance, the caption "two dogs running around" was correctly predicted as "two dogs are playing with each other on the ground." However, the model exhibited inconsistencies in generating captions for some images. Particularly, certain predictions showed considerable discrepancies with the actual captions, suggesting that the model had difficulty detecting the finer aspects and subtleties of the photos. An illustration of this was the misleading real caption, "The boy is eating pizza over a tin dish," which was "a man bites baby with his hands."
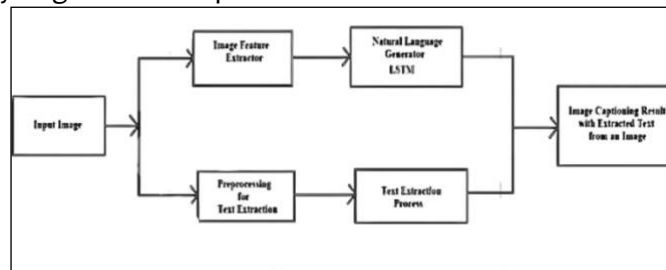
The quantitative evaluation demonstrates that the model attains competitive scores on the BLEU evaluation metrics. The BLEU-1 score, in particular, reaches a commendable value of 0.564952, although the BLEU-2 score, which represents lower performance, trails behind at 0.339059. According to these metrics, this model is not the most promising one. However, it is imperative to reiterate that the main goal of this kernel is to give a general overview of the features of Image Captioning that are realizable.

The mixed outcomes underscore the intricate nature of the image caption generation task and reinforce the necessity for additional research to enhance the model's proficiency. It may be possible to improve the precision and coherence of caption generation by addressing issues with overfitting, model complexity, and dataset size.

| Image ID | Image | Actual caption | Predicted caption |
|---|---|---|---|
| 1095590286_c654f7e5a9 |  | startseq blond dog and black and white dog run in dirt field endseq | startseq two dogs are playing with each other on the ground endseq |
| 1034276567_49bb87c51c |  | startseq boy bites hard into treat while he sits outside endseq | startseq man bites baby with his hands endseq |
| 2084217208_7bd9bc85e5 |  | startseq "a person in blue jacket wearing bicycle helmet is riding bike" endseq | startseq man in blue jacket riding bike endseq |
| 1989145280_3b54452188 |  | startseq blond girl with sunglasses on her head looks squeamish endseq | startseq woman with sunglasses on her face is holding onto cellphone endseq |

| | | | |
|---|---|---|---|
| 1778020185_1d44c04dae |  | startseq brown dog runs for white and black dog on the grass endseq | startseq dog is running through the grass endseq |
| 1429546659_44cb09cbe2 |  | startseq white dog and black dog in field endseq | startseq two dogs playfully wrestle in the grass endseq |
| 1330645772_24f831ff8f |  | startseq black and white dog is running in the grass endseq | startseq black and white dog is running through field endseq |
| 1143373711_2e90b7b799 |  | startseq bicycle rider is crossing street endseq | startseq man is riding bike on city street endseq |
| 1015584366_dfcec3c85a |  | startseq black dog leaps over log endseq | startseq black dog is climbing over log endseq |
| 109202801_c6381eef15 |  | startseq two draft horses pull cart through the snow endseq t endseq | startseq woman with red and red colored coat rides horse endseq |

The outcomes of this study shed light on the complex landscape of image caption generation through the integration of Convolutional Neural Networks (CNN) and Long Short-Term Memory (LSTM) networks. The mixed results obtained underscore the intricate nature of this task, where the model exhibited a commendable ability to capture visual context and produce contextually accurate captions in certain cases. This achievement is reflected in competitive BLEU scores for specific instances, revealing the model's proficiency in aligning generated captions with actual captions.

However, the presence of divergent predictions in other instances highlights the nuanced challenges that persist within image caption generation. These disparities between predicted and actual captions indicate that the model grapples with accurately representing intricate image details and nuances. This disparity also underscores the model's occasional inability to extract implicit semantic information from images and generate captions that resonate with the complexities of human language.

A notable factor in the performance variance could be the complexity of the dataset. The Flickr8k dataset, while diverse, may still lack the necessary breadth to comprehensively train the model on the intricacies of language nuances and visual context. This limitation could be exacerbated by the quality and diversity of the captions within the dataset. Biases present in the dataset, such as cultural or contextual biases in captions, may inadvertently affect the generated results. Therefore, the dataset's limitations could have contributed to the model's occasional deviations from accurate captions.

Furthermore, it's essential to acknowledge that the selected model architecture might not fully capture complex semantic relationships in certain images, leading to suboptimal captions. This limitation underscores the need for continued exploration into more advanced architectures and approaches that can better understand and represent intricate visual scenes. In addition to datasetrelated challenges, the preprocessing steps, including text cleaning and tokenization, may inadvertently remove relevant information or introduce artifacts. The choice of hyperparameters, such as batch size and epoch count, could also impact the training process and model performance. These factors emphasize the importance of meticulous experimentation and hyperparameter tuning to unlock the model's full potential.

Lastly, while quantitative metrics like BLEU scores provide valuable insights into the model's performance, they may not fully capture the quality and nuances of generated captions. The challenge of evaluating the "human-likeness" of captions remains an ongoing area of research, highlighting the need for more comprehensive evaluation metrics.

So, combining CNN and LSTM seems like a good path for creating image captions. We've seen both potential wins and challenges that need more investigation. As this field grows, it's clear that we'll need to bring together Computer Vision, Natural Language Processing, and Machine Learning to really make image captions shine. This study is a starting point, showing that there's a lot more to explore to make machines talk about pictures as humans do.

# CHAPTER 6: CONCLUSION & FUTURE SCOPE

Our investigation into image caption generation using CNN and LSTM took us through a landscape that is both fascinating and challenging. The model's ability to provide coherent and meaningful captions for a variety of photos was both intriguing and promising, reiterating the strong interplay between various neural architectures.

The journey did not, however, come without difficulties. The variation in caption predictions, which range from being extremely correct to being greatly off, relates to the complexity of the imagecaption interaction. This discrepancy highlights the current limitations of machine comprehension and serves as an indicator that it is still difficult to convey visual content through language.

Our work reaffirms the value of multidisciplinary collaboration in moving the field of AI forward. Bringing together the fields of computer vision and natural language processing brings us one step closer to creating machines that can comprehend and communicate visually. Despite the fact that our model showed strength in this area, the way forward will be made clearer by the pursuit of improving approaches, utilizing larger datasets, and embracing creative ways to give our models the ability to describe images with the simplicity of human language, despite the fact that our model showed competence in this domain.

In conclusion, this study strengthens our resolve to discover the complex web of links between images and captions. With every move, we minimize the gap between the visual and language spheres, laying the way for machines that can create stories out of pixels and enhance humantechnology interactions unexpectedly.

# REFERENCES

[1] Soh, Moses. "Learning CNN-LSTM architectures for image caption generation." *Dept. Comput. Sci., Stanford Univ., Stanford, CA, USA, Tech. Rep* 1 (2016).

[2] Katiyar, Sulabh, and Samir Kumar Borgohain. "Comparative evaluation of CNN architectures for image caption generation." *arXiv preprint arXiv:2102.11506* (2021).

[3] Aneja, Jyoti, Aditya Deshpande, and Alexander G. Schwing. "Convolutional image captioning." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2018.

[4] Al-Muzaini, Huda A., Tasniem N. Al-Yahya, and Hafida Benhidour. "Automatic Arabic image captioning using RNN-LSTM-based language model and CNN." *International Journal of Advanced Computer Science and Applications* 9.6 (2018).

[5] Bhalekar, Madhuri, and Mangesh Bedekar. "D-CNN: A new model for generating image captions with text extraction using deep learning for visually challenged individuals." *Engineering, Technology & Applied Science Research* 12.2 (2022): 8366-8373.

[6] Gu, Jiuxiang, et al. "An empirical study of language cnn for image captioning." *Proceedings of the IEEE international conference on computer vision*. 2017.

[7] Vinyals, Oriol, et al. "Show and tell: A neural image caption generator." *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015.

[8] Ding, Songtao, et al. "Image caption generation with high-level image features." *Pattern Recognition Letters* 123 (2019): 89-95.

[9] Mopuri, Konda Reddy, Utsav Garg, and R. Venkatesh Babu. "Cnn fixations: an unraveling approach to visualize the discriminative image regions." *IEEE Transactions on Image Processing* 28.5 (2018): 2116-2125.

# APPENDIX

*Appendix A: Technical Terms and Libraries*

1. **Convolutional Neural Networks (CNN):**
A type of neural network designed for processing grid-like data, such as images. It uses convolutional layers to capture spatial hierarchies of features.

2. **Long Short-Term Memory (LSTM):**
A type of recurrent neural network (RNN) architecture capable of learning long-term dependencies. It's widely used in sequence modelling tasks.

3. **BLEU Score:**
Bilingual Evaluation Understudy score, a metric for evaluating the quality of machine-generated text by comparing it to a set of reference texts.

4. **VGG16:**
A pre-trained convolutional neural network architecture known for its effectiveness in image classification tasks.

5. **Preprocessing:**
The steps taken to prepare and clean data before feeding it into a machine learning model.

6. **Feature Extraction:**
The process of obtaining relevant information or features from raw data, often used to reduce the dimensionality of the input.

7. **Model Architecture:**
The design and structure of the neural network, including the arrangement of layers and connections.

8. **TensorFlow:**
An open-source machine learning framework developed by Google for building and training machine learning models.

9. **Keras:**
A high-level neural networks API written in Python that runs on top of other machine learning libraries, including TensorFlow.