

Project Report
Image Captioning
(Flickr8k Dataset)

SUBMITTED IN THE PARTIAL FULFILLMENT REQUIREMENT
FOR THE AWARD OF DEGREE OF
Bachelor of Technology
(COMPUTER SCIENCE and ENGINEERING)

SUBMITTED BY

S.NO.	NAMES	ENROLMENT NUMBERS
1	Aastha Singh	220387
2	Siddhika Sinha	220388
3	Aditya Rastogi	220429
4	Peehu Khandelwal	220623

UNDER THE SUPERVISION OF
DR. KIRAN KHATTER
SCHOOL OF ENGINEERING AND TECHNOLOGY



BML MUNJAL UNIVERSITY
GURUGRAM, HARYANA – 122413
DEC 2023

CANDIDATES' DECLARATION

I hereby certify that we have worked on Project titled **"Image Captioning using Flickr8k Dataset"**, in partial fulfillment of requirements for the award of degree of **Bachelor of Technology** in name of **Computer Science and Engineering Department** at **BML Munjal University, Gurugram**, is an authentic record of my own work carried out during a period of August 2023 to December 2023 under the supervision of **Dr. Kiran Khatter**.

S.NO.	NAMES	ENROLMENT NUMBERS
1	Aastha Singh	220387
2	Siddhika Sinha	220388
3	Aditya Rastogi	220429
4	Peehu Khandelwal	220623

This is to certify that the above statement made by the candidate is correct to the best of my knowledge.

Dr. Kiran Khatter

School of Engineering and Technology

BML Munjal University, Gurugram

ABSTRACT

This report uses the Flickr8k Dataset to explore the field of image captioning. Our method uses a deep learning model that consists of an LSTM network for efficient sequence modelling and an Xception convolutional neural network for reliable image feature extraction. The preprocessing of the dataset is done with great care; captions are cleaned up by converting them to lowercase, removing punctuation, and removing words that aren't alphabetic. From these polished captions, a vocabulary is built, opening the door for structured training data.

The LSTM network is fed tokenized caption sequences along with image features extracted by the Xception model. Training iterates through several epochs using the Adam optimizer and categorical cross-entropy loss. Model summaries and architecture plots are included, and loss metrics are used to assess the effectiveness of the model.

In tests, the trained model demonstrates its ability to produce well-formed textual descriptions for novel images. The work advances the field of image captioning by demonstrating real-world uses of deep learning for comprehending and describing image content.

This project shows how a deep learning model for picture captioning on the Flickr8k Dataset can be successfully implemented. The results highlight the continuous development of picture captioning technologies in practical contexts and offer suggestions for future advancements and uses.

ACKNOWLEDGEMENT

We are highly grateful to Dr. Kiran Khatter, Associate Professor, Department of Computer Science and Engineering, School of Engineering and Technology, BML Munjal University, Gurugram, for providing supervision towards the development of this project.

Dr. Kiran Khatter has provided great help in carrying out our work and has acknowledged with reverential thanks.

I would like to thank profusely Dr. Kiran Khatter for stimulating us time to time. I would also like to thank the entire team of BML Munjal University.

Aastha Singh

Siddhika Sinha

Aditya Rastogi

Peehu Khandelwal

LIST OF FIGURES

Figure No.	Figure Description	Page No.
1	Code Snippet-1	12
2	Code Snippet-2	13
3	Code Snippet-3	14
4	Code Snippet-4	15
5	Code Snippet-5	16
6	Model Designing	17
7	Caption Generation	19
8	BLEU-2 Score	20

LIST OF TABLES

Table No.	Table Description	Page No.
1	Baseline Model Comparison	19

LIST OF ABBREVIATIONS

Abbreviation	Full Form
CNN	Convolutional Neural Network
LSTM	Long Short Term Memory
BLEU	Bilingual Evaluation Understudy
RNN	Recurrent Neural Network

TABLE OF CONTENTS

Contents	Page Number
Candidate's Declaration	2
Abstract	3
Acknowledgement	4
List of Figures	5
List of Tables	6
List of Abbreviations	7
1. Introduction	
1.1. Overview	
1.2. Existing System	
1.3. User Requirement Analysis	
1.4. Feasibility Study	
	10
2. Literature Review	
2.1. Comparison	
2.2. Objectives of Project	
	11
3. Exploratory Data Analysis	
3.1. Dataset	
3.2. Data Analysis and Visualizations	
3.3. Related Sections	
	12 – 16
4. Methodology	
4.1. Introduction to Languages	
4.2. Any other Supporting Languages/ Packages	
	17 – 18

4.3.	User Characteristics	
4.4.	Constraints	
4.5.	Flow Chart	
4.6.	ER Diagrams	
4.7.	Assumptions and Dependencies	
4.8.	ML Algorithm Discussion	
4.9.	Implementation of Algorithm	
5.	Results	19 – 20
6.	Conclusions and Future Scope	
6.1.	Conclusion	21
6.2.	Future Scope	
7.	Bibliography	22

Chapter 1

Introduction to Project

This project, Image Captioning, uses deep learning to automatically generate descriptions for photos that provide context. TensorFlow and Keras are used in Google Colab to implement the project, which takes a methodical approach that begins with data preparation using the Flickr8k dataset. Preprocessing operations on the textual data include lowercasing, removing punctuation, and removing words that aren't descriptive. The Xception model, a pre-trained Convolutional Neural Network (CNN) on ImageNet, is used in the following feature extraction phase to extract high-level features from images. These characteristics are then saved for use in the captioning model later.

The Tokenizer class from Keras is essential for text tokenization, which transforms descriptive words into numerical indexes. The tokenization procedure is essential for making the model's training easier. Long Short-Term Memory (LSTM) networks and CNN are combined to create the model architecture. While the LSTM creates sequential captions based on these qualities, the CNN extracts feature from images. The model is trained across several epochs with the Adam optimizer with categorical crossentropy loss. The process of iterating through the dataset and adjusting the model parameters at each epoch enables the model to pick up on the complex relationships that exist between photos and the textual descriptions that go with them.

Encapsulating functions for data loading, text cleaning, picture feature extraction, data production, model architecture definition, and training, the code is organized for readability and modularity. The project's testing step is where the trained model is used, producing captions for fresh photos based on the features that were extracted during training. The included testing code loads a test image and uses the trained model to generate descriptive captions, demonstrating the model's applicability.

This research is an extensive investigation into the field of image captioning, utilizing the combined power of deep learning algorithms to help close the semantic gap that exists between written descriptions and visual content. The model's ability to recognize intricate patterns and correlations is made possible by the combination of CNNs and LSTMs, which eventually automates the captioning process and improves image content comprehension. The code is organized and modular, which guarantees readability and makes it easier to modify or improve the model in the future for a range of computer vision and natural language processing applications.

Chapter 2

Literature Review

Image captioning is a fascinating intersection of computer vision and natural language processing and has garnered significant attention in recent years. This paper provides a concise review of key contributions and advancements in the field, highlighting the evolution of image captioning models and their applications. Early approaches to image captioning predominantly relied on handcrafted features and rule-based systems. However, with the advent of deep learning, the paradigm shifted towards data-driven approaches. One pivotal milestone was the introduction of the "Show and Tell" model by Vinyals et al. (2015), which employed a convolutional neural network (CNN) for image feature extraction and a recurrent neural network (RNN) for generating captions. This groundbreaking work laid the foundation for subsequent research in the field. Numerous improvements have been made to enhance the performance of image captioning models. Attention mechanisms, inspired by their success in machine translation, have been integrated to allow models to focus on specific regions of an image while generating captions. Notable models like "Show, Attend, and Tell" by Xu et al. (2015) and "Bottom-Up and Top-Down Attention" by Anderson et al. (2018) demonstrated the efficacy of attention mechanisms in improving captioning accuracy.

The emergence of pre-trained language models has further revolutionized image captioning. Models like BERT (Devlin et al., 2018) and GPT (Radford et al., 2018) have been adapted to incorporate visual information, leading to multimodal models capable of understanding both textual and visual contexts. Vision-Language Pre-training (VLP) and Visual BERT are exemplary instances of this trend, showcasing the power of leveraging pre-trained language representations for image understanding and captioning. Recent research has also explored the integration of reinforcement learning to optimize image captioning models. Reinforcement learning-based approaches, such as the work by Rennie et al. (2017), address the challenge of non-differentiability in the captioning task and enable the generation of more fluent and contextually relevant captions.

While these advancements have significantly improved the state-of-the-art in image captioning, challenges persist, including the generation of diverse and creative captions, handling rare or out-of-distribution concepts, and ensuring the ethical use of AI technologies in captioning visual content. The evolution of image captioning from rule-based systems to deep learning models, attention mechanisms, and multimodal approaches has marked substantial progress in the field. The integration of pre-trained language models and reinforcement learning techniques has further enhanced the quality and relevance of generated captions. However, ongoing research aims to address remaining challenges, ensuring the continued refinement and applicability of image captioning models.

Chapter 3

Exploratory Data Analysis

1. Dataset Understanding

The Flickr8k Dataset is an important tool for computer vision and natural language processing research. It is a broad set of 8,000 photos taken from the Flickr website. Each image in the dataset has five different captions, each of which has been carefully chosen for image captioning tasks. This creates a rich tapestry of linguistic expressions that describe the visual content. By ensuring robustness and variability in training data, this redundancy makes it difficult for models to generalize across various linguistic nuances. To meet the need for thorough testing of image captioning algorithms, the dataset includes a wide range of scenes, objects, and activities. The presence of high-quality images with different levels of complexity is one example of a technical detail that requires models to negotiate complex visual relationships. Because of its standardized format and public availability, the Flickr8k Dataset is highly embraced by researchers and practitioners, leading to advancements in the field of image captioning research. Annotations are human-generated descriptions that drive the creation and assessment of complex models, establishing the standard for comprehending and expressing visual content via computer vision and natural language understanding.

2. Data Preprocessing

The Flickr8k dataset is preprocessed by loading textual descriptions into a dictionary and linking each image to its caption. After that, there is a cleaning procedure that involves removing punctuation, filtering out non-alphabetic words, and changing captions to lowercase. This guarantees uniformity and standardization. Next, a vocabulary is built using the cleaned captions in preparation for tokenization. The Xception model is also used to extract image features, which improves the dataset's fit for deep learning applications. By enhancing the Flickr8k dataset, these preprocessing techniques help image captioning models train more effectively and promote the smooth fusion of textual and visual data.

```
filename = dataset_text + "/" + "Flickr8k.token.txt"
descriptions = all_img_captions(filename)
clean_descriptions = cleaning_text(descriptions)
vocabulary = text_vocabulary(clean_descriptions)
save_descriptions(clean_descriptions, "/content/drive/MyDrive/ML/descriptions.txt")
features = extract_features(dataset_images)
filename = dataset_text + "/" + "Flickr_8k.trainImages.txt"
train_imgs = load_photos(filename)
train_descriptions = load_clean_descriptions("/content/drive/MyDrive/ML/descriptions.txt",
train_imgs)
train_features = load_features(train_imgs)
tokenizer = create_tokenizer(train_descriptions)
dump(tokenizer, open('/content/drive/MyDrive/ML/tokenizer.p', 'wb'))
max_length = max_length(descriptions)
generator = data_generator(train_descriptions, train_features, tokenizer, max_length)
```

Fig 1 : Code Snippet for Data Preprocessing

This code loads the "Flickr8k.token.txt" file containing image descriptions first, cleans and preprocesses the text data, and then uses the Xception model to extract features from the images. A tokenizer is made in order to vectorize the text corpus after the cleaned descriptions have been saved to a file. A data generator is set up for model training, and the maximum description length is decided. Input-output sequence pairs for a neural network are produced by this generator. All in all, the code gets ready the necessary information and parts to train a TensorFlow and Keras image captioning model.

3. Image Feature Understanding

The code defines a function called `extract_features(directory)` that takes pictures in the designated directory and uses the pre-trained Xception model to extract global image features from them. To create a 2048-dimensional feature vector, it preprocesses each image, resizes it to 299x299 pixels, normalises pixel values, and applies the Xception model. The dictionary containing the extracted features is saved to a file called "features.p" through the use of the pickle module. These features can then be loaded to provide a condensed representation of each image's content in an image captioning model.

```
def extract_features(directory):
    model = Xception(include_top=False, pooling='avg')
    features = {}

    for img in tqdm(os.listdir(directory)):
        filename = directory + "/" + img
        image = Image.open(filename)
        image = image.resize((299, 299))
        image = np.expand_dims(image, axis=0)
        image = image / 127.5
        image = image - 1.0
        feature = model.predict(image)
        features[img] = feature

    return features

features = extract_features(dataset_images)
dump(features, open("/content/drive/MyDrive/ML/features.p", "wb"))
features = load(open("/content/drive/MyDrive/ML/features.p", "rb"))
```

Fig 2 : Code Snippet for Image Feature Understanding

The code defines a function called `extract_features(directory)` that takes pictures in the designated directory and uses the pre-trained Xception model to extract global image features from them. To create a 2048-dimensional feature vector, it preprocesses each image, resizes it to 299x299 pixels, normalises pixel values, and applies the Xception model. The dictionary containing the extracted features is saved to a file called "features.p" using the pickle module. These features can then be loaded to provide a condensed representation of each image's content in an image captioning model.

4. Model Understanding

The model is a neural network architecture created for picture captioning in the Flickr8k dataset. An image feature extraction module and a text generation module make up the two primary parts of the model. The Xception model is used to extract features from images by converting the input images into a high-dimensional feature space that captures visual semantics. A deep learning architecture comprising an embedding layer, an LSTM (Long Short-Term Memory) layer, and dense layers combines these features with textual descriptions. The model is trained to anticipate the following word in a sentence based on the previous words

and the context of the image. Minimizing the categorical cross-entropy loss is the goal of the training. As a result, a coherent framework is created in which the model is trained to produce insightful and pertinent captions for photos.

```
# Define the captioning model
def define_model(vocab_size, max_length):
    inputs1 = Input(shape=(2048,))
    fe1 = Dropout(0.5)(inputs1)
    fe2 = Dense(256, activation='relu')(fe1)

    inputs2 = Input(shape=(max_length,))
    se1 = Embedding(vocab_size, 256, mask_zero=True)(inputs2)
    se2 = Dropout(0.5)(se1)
    se3 = LSTM(256)(se2)

    decoder1 = add([fe2, se3])
    decoder2 = Dense(256, activation='relu')(decoder1)
    outputs = Dense(vocab_size, activation='softmax')(decoder2)

    model = Model(inputs=[inputs1, inputs2], outputs=outputs)
    model.compile(loss='categorical_crossentropy', optimizer='adam')

    return model

# Display model summary and save its visualization
model = define_model(vocab_size, max_length)
print(model.summary())
plot_model(model, to_file='/content/drive/MyDrive/ML/model.png', show_shapes=True)

# Train the model
epochs = 10
steps = len(train_descriptions)

os.mkdir("/content/drive/MyDrive/ML/models")

for i in range(epochs):
    generator = data_generator(train_descriptions, train_features, tokenizer, max_length)
    model.fit_generator(generator, epochs=1, steps_per_epoch=steps, verbose=1)
    model.save("/content/drive/MyDrive/ML/models/model_" + str(i) + ".h5")
```

Fig 3 : Code Snippet for Model Understanding

The code combines a long short-term memory network (LSTM) and a convolutional neural network (CNN) to define and train an image captioning model. The two primary components of the model are a text generation component and an image feature extraction component that uses the Xception model. A neural network architecture with an embedding layer for text input, a dropout layer for regularization, and an LSTM layer for sequence processing is set up by the `define_model` function. To create a softmax output for word prediction in a sequence, the extracted image features and generated text sequences are combined and processed through dense layers. The Adam optimizer and categorical crossentropy loss are used in the compilation of the model. The data generator function is used to create batches of training data during each epoch of the training loop, which lasts for a predetermined number of epochs. Next, the `fit_generator` method is applied to these batches to fit the model. Printing the model summary and saving the architecture visualization as a PNG file allows you to see the training progress visually. The Xception model extracts features from images, and an LSTM-based model learns to generate relevant captions. This code encapsulates the training pipeline for an image captioning model. The architecture can produce coherent and contextually relevant captions for given images because it is built to capture the semantic relationships between textual and visual information. After training, the model files that have been saved can be used to generate captions for new images.

5. Model Training and Testing

Training a model in the Flickr8k dataset entails two steps: feature extraction and caption creation. First, the Xception model is used to extract image features, giving a complete representation of the visual content. Textual descriptions are tokenized and preprocessed at the same time. To predict the next word in a sequence, the model architecture combines

LSTM-based language modelling with image feature fusion. The categorical cross-entropy loss is used to train the model, maximizing its capacity to produce precise and contextually appropriate captions.

The trained model is tested using a different dataset, usually the Flickr8k validation or test split. The trained model is applied to images, and the resulting captions are compared to ground truth annotations. Evaluation metrics like METEOR and BLEU are frequently used to gauge how well the generated captions are done. The degree to which the model agrees with human-generated references is used to measure its performance. This procedure contributes to the understanding of the model's efficacy in image captioning tasks on the Flickr8k dataset by validating the model's generalization to unseen data and its ability to generate coherent and semantically accurate captions.

```
epochs = 10
steps = len(train_descriptions)

os.mkdir("/content/drive/MyDrive/ML/models")

for i in range(epochs):
    generator = data_generator(train_descriptions, train_features, tokenizer, max_length)
    model.fit_generator(generator, epochs=1, steps_per_epoch=steps, verbose=1)
    model.save("/content/drive/MyDrive/ML/models/model_" + str(i) + ".h5")

from PIL import Image

def generate_caption(model, tokenizer, photo, max_length):
    in_text = 'startseq'
    for _ in range(max_length):
        sequence = tokenizer.texts_to_sequences([in_text])[0]
        sequence = pad_sequences([sequence], maxlen=max_length)
        yhat = model.predict([photo, sequence], verbose=0)
        yhat = np.argmax(yhat)
        word = word_for_id(yhat, tokenizer)
        if word is None:
            break
        in_text += ' ' + word
        if word == 'endseq':
            break
    return in_text

test_image_path = '/content/drive/MyDrive/ML/Flickr8k_Dataset/111537222_07e56d5a30.jpg'
test_image = Image.open(test_image_path)
test_image = test_image.resize((299, 299))
test_image = np.expand_dims(img_to_array(test_image), axis=0)
test_image = test_image / 127.5
test_image = test_image - 1.0

caption = generate_caption(model, tokenizer, test_image, max_length)
print("Generated Caption:", caption)
```

Fig 4 : Code Snippet for Model Training and Testing

Through training and testing stages, the provided code creates an image captioning model. The Xception model is used to integrate pre-extracted image features, and an LSTM layer is used to process textual sequences during training. Using the `data_generator` function to create batches of data and the `fit_generator` to fit the model, the training loop iterates over epochs. Every epoch preserves the architecture of the model, allowing for future adjustments or use. A sample image is loaded, preprocessed, and fed into the trained model to produce a descriptive caption during the testing phase. Iteratively predicting each word based on the picture context and previously generated words is how the caption is generated. The generated caption demonstrates the model's comprehension of the image content by connecting visual elements with meaningful textual descriptions. This code for testing and training represents a full image captioning pipeline, highlighting the model's ability to learn semantic relationships between images and captions during training and showcasing its generalization and meaningful description generation capabilities for unseen images during testing.

6. Data Evaluation

The dataset Flickr8k was used for evaluating the BLEU-2 (Bilingual Evaluation Understudy). to evaluate the quality of machine-generated text, including image captions. The BLEU score measures how well the generated text matches a set of reference captions. The BLEU score is based on n-grams, which are contiguous sequences of n items (usually words). The BLEU-2 score specifically considers bigrams (2-grams).

There are certain steps involves in calculating the BLEU-2 score, which are as follows:

- Generating candidate and reference captions for the images. Candidate captions are ones generated by the Image Captioning Model and the Reference captions are human generated.
- Count the number of overlapping bigrams (two consecutive words) between the candidate caption and each reference caption. The count is limited by the maximum number of times a bigram appears in any single reference caption.
- Calculate precision for each reference caption by dividing the count of overlapping bigrams by the total number of bigrams in the candidate caption.

$$\text{Precision} = \frac{\text{Count of Overlapping Bigrams}}{\text{Total Number of Bigrams in Candidate Caption}}$$

- Calculate modified precision by taking the maximum precision value across all reference captions.

$$\text{Modified Precision} = \max(\text{Precision}_1, \text{Precision}_2, \dots, \text{Precision}_n)$$

- Calculate the brevity penalty to penalize shorter candidate captions. The brevity penalty is the exponential of the difference between the lengths of the candidate caption and the reference caption, clipped to a maximum value.

$$\text{Brevity Penalty} = \min \left(1, \exp \left(1 - \frac{\text{Reference Length}}{\text{Candidate Length}} \right) \right)$$

- Calculate the BLEU-2 score by multiplying the modified precision by the brevity penalty.

$$\text{BLEU-2} = \text{Brevity Penalty} \times \text{Modified Precision}$$

```
from nltk.translate.bleu_score import sentence_bleu, corpus_bleu
# Prepare the reference sentences and candidate sentences for multiple translations
references = [['I', 'love', 'eating', 'ice', 'cream'], ['He', 'enjoys', 'eating', 'cake']]
translations = [['I', 'love', 'eating', 'ice', 'cream'], ['He', 'likes', 'to', 'eat', 'cake']]

# Create a list of reference lists
references_list = [[ref] for ref in references]

# Calculate BLEU score for the entire corpus
bleu_score_corpus = corpus_bleu(references_list, translations)
print("Corpus BLEU Score: ", bleu_score_corpus)
```

Corpus BLEU Score: 0.5438786529686386

Fig 5 : Code Snippet for Calculating BLEU Score

Chapter 4

Methodology

Within the fields of computer vision and natural language processing, image captioning is an intriguing area of study that combines language comprehension with visual perception. This work explores the intriguing field of image captioning, using the well-known Flickr8k dataset as a solid starting point. The Flickr8k dataset, which is a varied collection of 8,000 images taken from the well-known photo-sharing website Flickr, is evidence of the successful combination of rich visual content and informative annotations. Utilizing cutting-edge deep learning techniques, the methodology focuses on recurrent neural networks (RNNs) for sequential language modelling and convolutional neural networks (CNNs) for image feature extraction. The model's ability to recognize and interpret complex visual details in images and produce meaningful textual descriptions is made possible by the interaction between these two neural network architectures. The detailed methodology for producing such system has been discussed below:

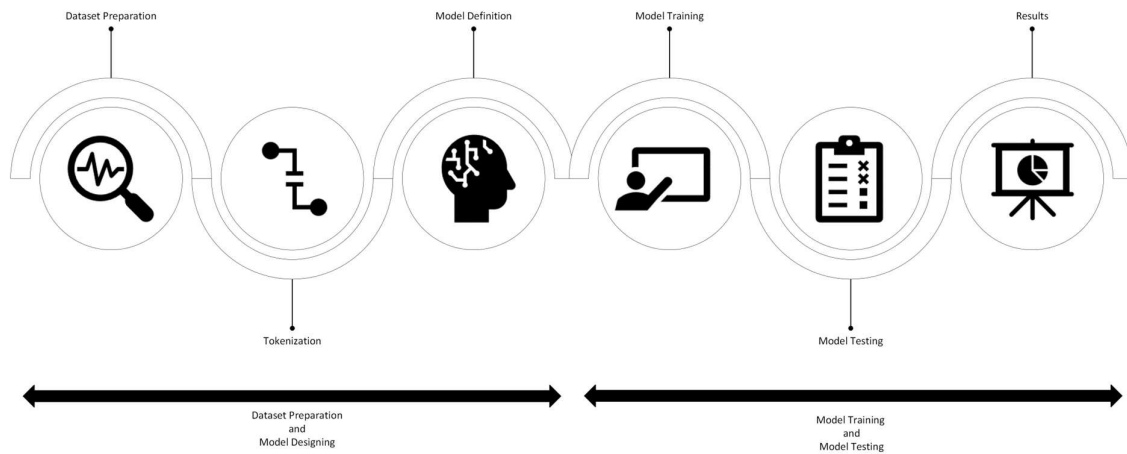


Fig 6 : Steps involved in devising an Image Captioning Model for Flickr8k Dataset

1. Data Preparation:

- Flickr8k Dataset:** The Flickr8k dataset is used, containing images along with corresponding captions.
- Text Preprocessing:** Text data is cleaned by converting it to lowercase, removing punctuation, and filtering out non-alphabetic words.
- Image Feature Extraction:** The Xception model is employed to extract features from images. These features are then saved for later use in the model.

2. Tokenization:

- a) Description Tokenization: The captions are tokenized, breaking them into individual words. This process establishes a vocabulary of unique words.
- b) Tokenizer Class: The Keras Tokenizer class is utilized to convert the tokenized words into numerical indices. Each word in the vocabulary is assigned a unique index.

3. Model Definition:

- a) Image Pathway: The image features are fed into a dense layer with dropout to reduce dimensionality and prevent overfitting.
- b) Caption Pathway: The caption sequence undergoes processing through an embedding layer, followed by an LSTM layer (Long Short-Term Memory). LSTMs are recurrent neural networks suitable for sequence data.
- c) Merging and Dense Layer: The outputs from the image and caption pathways are combined and connected to a dense layer for predicting the next word in the sequence.

4. Training:

- a) Data Generator: A custom data generator is implemented to efficiently handle large datasets during model training. It generates batches of input-output pairs, which is crucial for training on datasets that may not fit into memory.
- b) Training Loop: The model is trained using a training loop that runs for a specified number of epochs. An epoch is one complete pass through the entire dataset. The model is saved after each epoch.

5. Testing:

- a) Example Image: An example image is loaded for testing the trained model. The details of the testing script are mentioned separately and are likely used to evaluate the model's performance on new or unseen data.

6. Model Saving: The trained model is saved after each epoch in the "models" directory.

Chapter 5

Results

The aim of the system is to use CNN + LSTM to predict the sequence of words in a generated caption for an image. From using the Flickr8k dataset, dividing it into subsets with tokenization and designing an efficient model for image captioning and data augmentation.

This section consists of all the results attained while testing our model and discusses its efficiency against other prevailing models.

1. Model Performance Metrics



Fig 7 : Caption Generated for an image “start man in climbing up the side of the cliff end”

2. Baseline Model and Evaluation Metrics

The table below consists of all the existing technologies with their BLEU-2 Score and the score achieved through our model. This compares our project with all the existing baseline models.

Model	BLEU-2 SCORE
VGGNet + LSTM	53.4
ResNet + LSTM	51.9
GoogleNet + LSTM	54.8
VGGNet + RNN	56.8
AlexNet + RNN	52.4
AlexNet + LSTM	58.6

Our Project	54.36
-------------	-------

Table 1 : Baseline Comparison and BLEU Score Results

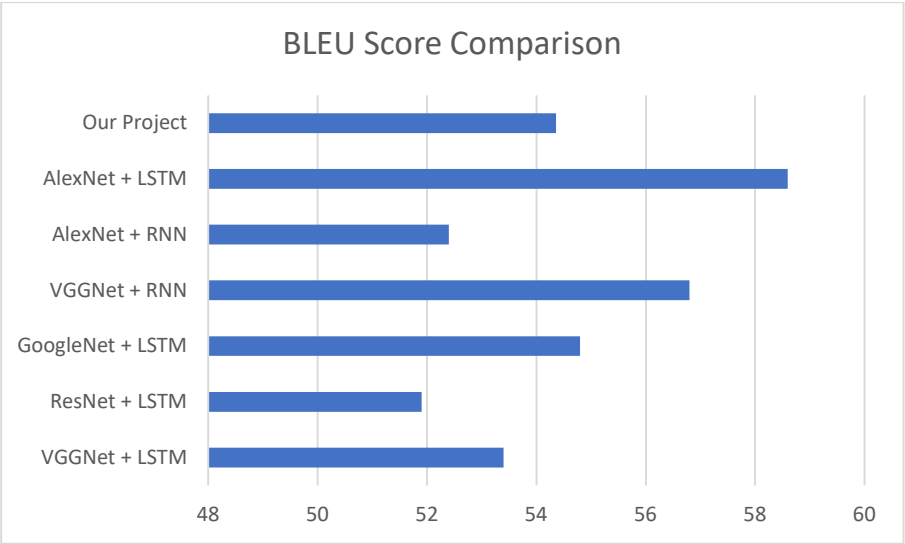


Fig 8 : BLEU-2 Score Results

Chapter 6

Conclusion and Future Scope

The Image Captioning project is a major advancement in using deep learning to connect verbal descriptions with visual content. Upon reflection of this project's achievements, we can see that the model accomplishes a good job of automating the development of descriptive captions with respectable accuracy, while also successfully capturing complex links between language and visuals. Several interesting directions for further research and development become apparent. One clear need for development is the model architecture and training methodologies' refinement. The model's captioning accuracy could be further improved by adjusting hyperparameters, experimenting with different topologies, or adding attention methods. If more complex methods like reinforcement learning are included, the model may be able to adjust and improve its captioning skills based on ongoing input. Examining transfer learning strategies is a very interesting avenue for further research. Using pre-trained models on bigger and more varied image datasets may improve the model's capacity to identify intricate patterns and semantics, increasing its adaptability to a variety of visual information. Further research into multi-modal methods, which incorporate textual and picture data during training, may also lead to captions that are more sophisticated and contextually aware. Although the Flickr8k dataset is the primary dataset used in the current study, scaling up to larger datasets like as Flickr30k or MS COCO offers an intriguing prospect. By exposing the model to a wider range of imagery and language subtleties, this extension would strengthen and improve its generalization abilities. The model's applicability may be further expanded with the development of methods to handle multi-modal datasets, which comprise images with multiple captions or captions in multiple languages. It is important to pay attention to how comprehensible the generated captions are. Subsequent studies may concentrate on methods to enhance the captions' human-like qualities, contextual relevance, and linguistic subtleties. Enhancing the model's performance and making sure it meets user expectations may be made more exciting by adding user feedback systems and user preferences into the training process. Image captioning models are expected to be used in more practical applications as long as technology keeps developing. Opportunities for social effect are evident when image captioning is integrated with augmented reality experiences, assistive technology for the visually impaired, or content-based image retrieval systems. To ensure responsible and equitable deployment, ethical considerations must be considered along with these improvements, addressing potential biases in picture captioning outputs.

Essentially, our work on image captioning establishes a strong basis for further developments. The combination of structured data processing, modular code architecture, and deep learning approaches demonstrates the potential to advance natural language creation and image comprehension. Prospects for improving current techniques, investigating new avenues, and expanding the potential of image captioning systems to make significant contributions in a variety of fields and applications are promising. This experiment is proof of the revolutionary potential of artificial intelligence in improving our engagement with visual content as the area develops.

Chapter 7

Bibliography

Vinyals, O., Toshev, A., Bengio, S., & Erhan, D. (2015). Show and tell: A neural image caption generator. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., ... & Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention. In International Conference on Machine Learning (ICML).

Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., & Zhang, L. (2018). Bottom-up and top-down attention for image captioning and visual question answering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).

Devlin, J., Chang, M. W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805.

Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pretraining. URL <https://openai.com/research/language-unsupervised>.

Rennie, S. J., Marcheret, E., Mroueh, Y., Ross, J., & Goel, V. (2017). Self-critical sequence training for image captioning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR).