

# **ImageSpeak: Generating Captions from Pixels**

**A Project Work Synopsis**

*Submitted in the partial fulfillment for the award of the degree of*

**BACHELOR OF ENGINEERING  
IN  
COMPUTER SCIENCE WITH SPECIALIZATION IN  
ARTIFICIAL INTELLIGENCE AND MACHINE LEARNING**

**Submitted by:**

21BCS6285 Shivani

21BCS6284 Ashmandeep Kaur

21BCS6280 Tarushi Sandeep Gupta

**Under the Supervision of:**

**Mr. Nirmalya Basu**



**CHANDIGARH  
UNIVERSITY**

Discover. Learn. Empower.

**CHANDIGARH UNIVERSITY, GHARUAN, MOHALI - 140413,**

**PUNJAB**

**August, 2023**

# Abstract

The Image Caption Generator project presents a novel application of artificial intelligence, merging computer vision and natural language processing to enable machines to understand visual content and generate contextually relevant captions. This project addresses the challenge of translating images into descriptive text, offering a versatile solution with far-reaching implications across multiple domains.

By harnessing the power of convolutional neural networks (CNNs) for image analysis and recurrent neural networks (RNNs) for language generation, the project creates a cohesive system capable of perceiving and describing intricate visual scenes. The project focuses on training the model using large and diverse datasets comprising images paired with human-generated captions. Through extensive training and optimization, the model learns to recognize visual features and associate them with appropriate linguistic expressions.

The generated captions go beyond mere factual descriptions, aiming to capture the essence of the image, contextual relationships, and emotions conveyed. This contextual understanding opens doors to applications like aiding visually impaired individuals, enriching content accessibility, improving search engine capabilities, and enhancing social media engagement.

In conclusion, the Image Caption Generator project exemplifies the synergy between computer vision and natural language processing, demonstrating the transformative impact of AI-driven technologies on enhancing the human-computer interaction landscape. The project's outcomes lay the groundwork for future advancements, encouraging the development of multimodal AI systems capable of comprehending and generating diverse forms of digital content.

Keywords: Deep Learning, CNN, RNN, LSTM

# Table of Contents

|   |    |
|---|----|
| Title Page                                      | i  |
| Abstract  | ii |
| 1. Introduction                                 |    |
| 1.1 Problem Definition                          |    |
| 1.2 Project Overview                            |    |
| 1.3 Hardware Specification                      |    |
| 1.4 Software Specification                      |    |
| 2. Literature Survey                            |    |
| 2.1 Existing System                             |    |
| 2.2 Proposed System                             |    |
| 2.3 Literature Review Summary                   |    |
| 3. Problem Formulation                          |    |
| 4. Research Objective                           |    |
| 5. Methodologies                                |    |
| 6. Experimental Setup                           |    |
| 7. Conclusion                                   |    |
| 8. Tentative Chapter Plan for the proposed work |    |
| 9. Reference                                    |    |

# 1. INTRODUCTION

## 1.1 Problem Definition

Develop an image caption generator that automatically generates descriptive and contextually relevant captions for a given input image. The goal of this task is to create a model that can effectively understand the content of an image and generate coherent and accurate textual descriptions that capture the key elements, objects, and activities depicted in the image. The generated captions should not only be grammatically correct but also semantically meaningful and coherent, reflecting the overall scene and context of the image. The challenge lies in training a model that can comprehend various visual concepts and their relationships within the images, and then translate this understanding into natural language captions that enhance the viewer's comprehension of the image content. The evaluation of the image caption generator will be based on the quality of the generated captions in terms of relevance, accuracy, fluency, and overall coherence with the corresponding images.

## 1.2 Problem Overview

The Image Caption Generator project addresses the pressing issue of enabling machines to comprehend and articulate visual content. While humans effortlessly interpret images, translating this ability into AI systems remains a challenge. Traditional methods fail to capture the complexity of visual scenes, resulting in shallow and inaccurate descriptions.

The fundamental problem lies in the gap between image understanding and linguistic expression. Computers struggle to glean contextual information, emotions, and relationships embedded in images, limiting their ability to generate meaningful captions. This limitation impedes progress in domains like accessibility, e-commerce, and content enrichment.

Moreover, biases in training data can seep into generated captions, perpetuating stereotypes and misinformation. Ethical considerations are thus pivotal in building responsible AI systems.

The project's scope encompasses designing an image caption generator that leverages neural networks, merging computer vision for image analysis and natural

language processing for generating captions. By training on diverse datasets, the system learns to recognize visual features and translate them into coherent textual descriptions. The goal is to create an AI tool that not only produces accurate captions but also does so ethically, unbiased, and contextually.

In conclusion, the Image Caption Generator project seeks to close the gap between visuals and language, revolutionizing AI's interaction with visual content. Through this, it aims to contribute to a more comprehensive and inclusive digital environment.

### **1.3 Hardware Specification**

High performance CPUs GPUs : CPUs like Intel Core i9 or AMD Ryzen 9, and GPUs like NVIDIA GeForce RTX 30 series or AMD Radeon RX 6000 series .

Memory(RAM): A minimum of 16 GB to 32 GB of RAM is recommended to handle large data and complex images.

Reliable Internet Connection A stable, high-speed Internet connection is required.

### **1.4 Software Specification**

Programming Language: There are many different programming languages Which can be used for ML. Python is the most commonly used language for ML.It had a rich ecosystem of libraries and data mining tools,Images and ML. So it is a good choice for developing an image caption generator

Integrated Development Environment (IDEs): used are Visual Studio Code, Google Collab and Jupyter Notebook.

## 2. LITERATURE SURVEY

### 2.1 Existing System

Deep learning approaches have been investigated in the context of image caption generators in a number of publications. "Show and Tell," [1] a crucial model for deep learning-based image caption generation, was introduced by Google Research in 2015. It made a significant addition to the field by showing how neural networks can automatically generate poetic descriptions for photographs. The model incorporates two different neural network types: a convolutional neural network (CNN) and a long short-term memory (LSTM) network. CNN examines the input image to derive its visual characteristics. These attributes are transformed into a fixed-length vector that serves to encode the visual information of the image and reflect its content. In a recurrent neural network, the encoded picture vector acts as the LSTM's initial hidden state.

The expansion of the "Show and Tell" paradigm that integrates an attention mechanism is called "Show, Attend, and Tell," [2] and it was developed by academics at the University of Montreal in 2015. This attention method enables the model to concentrate on various aspects of the image as it generates words, enhancing the calibre and coherence of the captions that are produced. The "Show, Attend, and Tell" concept can be explained in the following manner. This model, which is comparable to "Show and Tell," combines a Convolutional Neural Network (CNN) and a Long Short-Term Memory (LSTM) network. The incorporation of an attention mechanism, however, is the main innovation. To extract the input image's visual properties, CNN analyses it. These features serve to represent the image's content and are applied to the creation of a fixed-length vector that encodes the visual data of the image. The LSTM decoder incorporates the attention mechanism. The LSTM generates each word of the caption by dynamically focusing on various areas of the image rather than just the original image vector.

The "Show, Attend, and Tell" [3] paradigm is an extension of the "Show and Tell" paradigm that incorporates an attention mechanism. It was created in 2015 by researchers at the University of Montreal. The model may focus on different facets of the image while producing words using this attention strategy, which improves the quality and coherence of the captions that are generated. The idea of "Show,

Attend, and Tell" can be explained as follows. This model, which is similar to "Show and Tell," combines an LSTM network and a convolutional neural network. The key novelty, though, is the addition of an attention mechanism. CNN examines the input image to derive its visual characteristics. These characteristics represent the content of the image.

By combining a two-step attention process, the "Bottom-Up and Top-Down Attention" [4] approach, created by Google Research in 2018, improves image captioning. Bottom-Up Uses an object detection network to recognise crucial image sections. A feature vector that captures the relevance and visual content of each region serves as its representation. Top-Down Attention: Generates captions word-by-word using an LSTM-based decoder. It pays attention to the indicated image regions at each stage, concentrating on those that are important for producing each word.

## 2.2 Proposed System

The system for an image caption generator using deep learning aims to automatically generate descriptive captions for images by combining convolutional neural networks (CNNs) for image feature extraction and recurrent neural networks (RNNs) for language generation.

## 2.3 Literature Review Summary (Minimum 7 articles should refer)

| Year and Citation | Article/ Author  | Tools/ Software                   | Technique                     | Source        | Evaluation Parameter |
|-------------------|--|-----------------------------------|-------------------------------|---------------|----------------------|
| 2015              | "Show and Tell: A Neural Image Caption Generator" by Oriol Vinyals et al | Tensorflow, Python, Numpy, Pandas | CNN, RNN, Attention mechanism | Research Gate | BLEU, METEOR         |

|      |  |   |   |                |                      |
|------|--|---|---|----------------|----------------------|
| 2015 | "Neural Image Caption Generation with Visual Attention" by Kelvin Xu et al.                                    | Deep Learning Frameworks, Cuda And cuDNN, LaTeX | CNN,RNN ,Softmax activation, BLEU           | Research Gate  | BLEU, METEOR , CIDEr |
| 2017 | "Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning" by Jiasen Lu et al.      | Does not describe Specific tool                 | Reinforcement learning ,attention mechanism | Google Scholar | ROUGE , BLEU         |
| 2017 | "Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering" by Peter Anderson et al. | Python, Pytorch ,Tensorflow                     | Bottom-up attention ,Top-down attention,    | arXiv          | ROUGE, SPICE,BL EU   |
| 2018 | Up-Down: A Convolutional Encoder-Decoder Architecture for Multimodal Semantic Composition" by Qi Wu et al.     | NLP libraries, Teensorflow, Pytorch             | CNN, RNN, Beam Search                       | Google Scholar | Accuracy, SPICE      |



|      |  |                                   |   |                |   |
|------|--|-----------------------------------|---|----------------|---|
| 2019 | Image Captioning with Semantic Attention" by Long Mai et al. | Tensorflow, NLP Libraries,        | RNN, CNN, Beam Search                           | Google Scholar | Diversity metrics, Informativ e metrics |
| 2020 | Image Captioning Transformer" by Jie Lei et al               | Do not describe the specific tool | Transform er architecture , positional encoding | Google Scholar | ROUGE,S PICE,Hum an Evaluation          |

### **3. PROBLEM FORMULATION**

The core challenge of the Image Caption Generator project lies in enabling machines to comprehend visual content and generate meaningful textual descriptions. This complex task involves merging computer vision and natural language processing.

The initial hurdle is designing algorithms that enable AI systems to interpret images like humans, recognizing objects, scenes, emotions, and relationships between elements. The subsequent challenge is crafting coherent and contextually relevant captions that resonate with human communication. This involves training AI models to generate grammatically correct sentences that capture the essence of visual scenes.

Ethical considerations are vital, as the project aims to mitigate biases in training data and generated captions, ensuring responsible and unbiased AI usage. Ultimately, the goal is to create an AI model that bridges the gap between visual understanding and linguistic expression, transforming how machines perceive and communicate visual content.

### **4. OBJECTIVES**

The objectives of an image caption generator encompass creating an AI system that understands visual content, generates contextually relevant captions, and ensures responsible and ethical AI usage. The primary goals include enhancing the AI model's visual comprehension, refining natural language generation for coherent captions, and addressing biases.

Creativity and emotional context are sought to enrich captions, while practical applications span content accessibility, search enhancement, aiding the visually impaired, and more. Continuous improvement ensures the AI model remains cutting-edge and aligned with evolving needs.

## 5. METHODOLOGY

The methodology for developing an effective image caption generator involves a systematic approach that combines computer vision, natural language processing, data preprocessing, model architecture design, training, and evaluation. Here's a concise outline of the methodology:

### Data Collection and Preprocessing:

- Gather a diverse dataset of images paired with corresponding captions.
- Preprocess the images by resizing, normalizing, and augmenting to enhance model robustness.
- Tokenize the captions into words and create a vocabulary for mapping words to numerical indices.

### Feature Extraction and Embedding:

- Utilize Convolutional Neural Networks (CNNs) to extract visual features from the images. These features capture shapes, textures, and objects within the images.
- Embed words in the captions using word embeddings like Word2Vec or GloVe, transforming them into numerical vectors for language processing.

### Model Architecture Design:

- Design the architecture of the model, which typically comprises an image encoder and a text decoder.
- The image encoder uses the CNN to process images and extract visual features.
- The text decoder employs Recurrent Neural Networks (RNNs) or Transformer-based models to generate captions word by word.

### Training:

- Initialize the model's parameters and use the preprocessed data for training.

- During training, feed the images into the image encoder and the partial captions into the text decoder to predict the next word.
- Optimize the model using loss functions like cross-entropy, minimizing the difference between predicted and actual words.

#### Ethical Considerations:

- Address potential biases in the training data and captions to ensure ethical and responsible caption generation.
- Implement techniques to reduce biases and avoid generating harmful or inaccurate content.

#### Evaluation:

- Evaluate the model's performance using metrics like BLEU (Bilingual Evaluation Understudy) scores, METEOR, CIDEr, and ROUGE to measure caption quality, coherence, and relevance.
- Incorporate human evaluations to assess the captions' fluency and accuracy.

#### Fine-Tuning and Optimization:

- Fine-tune the model based on evaluation results to improve caption quality and contextual understanding.
- Experiment with hyperparameters, model architecture variations, and training techniques to achieve optimal performance.

#### Deployment and Real-time Use:

- Deploy the trained model to generate captions for new images in real-time applications.
- Ensure scalability and efficiency for different platforms and systems.

## 6.EXPERIMENTAL SETUP

Setting up experiments for an image caption generator involves configuring the hardware, software, dataset, model architecture, training, and evaluation processes. Here's a concise outline of the experimental setup:

Hardware and Software:

- Hardware: Utilize a computer with sufficient computational power, including GPUs, to accelerate model training.
- Software: Install necessary frameworks like TensorFlow, PyTorch, or Keras for model development and training.

Dataset:

- Data Collection: Gather a diverse dataset of images paired with humangenerated captions. Common datasets include MS COCO, Flickr30k, and Conceptual Captions.
- Data Preprocessing: Resize images to a consistent size and normalize pixel values. Tokenize captions into words and create a vocabulary mapping.

Model Architecture:

- Image Encoder: Implement a pre-trained CNN (e.g., ResNet, VGG) to extract visual features from images.
- Text Decoder: Design an RNN (LSTM or GRU) or transformer-based model to generate captions word by word.

Training:

- Loss Function: Use cross-entropy loss to minimize the difference between predicted and actual words in the captions.
- Optimizer: Apply optimization algorithms like Adam or RMSprop to update model parameters during training.

#### Batching:

- Organize data into batches for efficient training.
- Learning Rate Scheduling: Adjust the learning rate during training to ensure convergence and avoid overshooting.

#### Ethical Considerations:

- Bias Mitigation: Implement techniques to identify and mitigate biases in training data and generated captions.
- Content Filtering: Ensure that the generated captions are respectful, unbiased, and align with ethical standards.

#### Evaluation:

- Metrics: Utilize evaluation metrics like BLEU, METEOR, CIDEr, and ROUGE to assess caption quality, coherence, and relevance.
- Human Evaluation: Incorporate human assessments to gauge the fluency, accuracy, and context of generated captions.

#### Fine-Tuning and Optimization:

- Hyperparameter Tuning: Experiment with hyperparameters such as learning rate, batch size, and model architecture variations.
- Regularization: Apply techniques like dropout or L2 regularization to prevent overfitting.

#### Deployment and Real-time Use:

- Deployment: Deploy the trained model on appropriate platforms, considering scalability and real-time processing requirements.
- API Integration: Create APIs to allow users to input images and receive generated captions in real-time.

## 7. CONCLUSION

In conclusion, the development and evolution of the image caption generator represent a remarkable convergence of computer vision and natural language processing. This innovative technology has reshaped our interaction with visual content by enabling machines to bridge the gap between visual perception and linguistic expression.

Through the integration of neural network architectures, attention mechanisms, and ethical considerations, image caption generators have transcended traditional boundaries. They are now capable of not only identifying objects and scenes within images but also generating contextually relevant, coherent, and often creatively articulated textual descriptions.

The journey of image caption generators has been marked by advancements in data preprocessing, model design, and evaluation metrics. The incorporation of large datasets, attention mechanisms, and transformer-based models has elevated the quality of generated captions, making them more nuanced, coherent, and aligned with human understanding.

However, challenges remain, including addressing biases in training data, ensuring responsible AI usage, and fine-tuning models for optimal performance. The ethical dimensions of image caption generators underscore the need for ongoing vigilance and continuous improvement to avoid perpetuating harmful stereotypes or inaccuracies.

As image caption generators continue to expand their applications across diverse industries such as social media, e-commerce, education, and healthcare, their impact on communication, accessibility, and engagement cannot be overstated. These generators empower us to navigate the visual landscape with enhanced understanding, catering to both human and machine interpretations.

The ongoing efforts to refine this technology hold the promise of not only enhancing our interaction with images but also fostering responsible and equitable AI-powered communication in our digital age.

## **8. TENTATIVE CHAPTER PLAN FOR THE PROPOSED WORK**

### **CHAPTER 1: INTRODUCTION**

An Image Caption Generator is an innovative application of artificial intelligence that merges the realms of computer vision and natural language processing. It tackles the challenge of endowing machines with the ability to understand visual content and subsequently translate it into coherent and contextually relevant textual descriptions. By analyzing the intricate details, objects, scenes, and emotions within images, these generators aim to bridge the gap between visual perception and linguistic expression. The introduction of neural network architectures, like Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), has revolutionized this field, enabling AI systems to learn from vast datasets and generate captions that align with human understanding. As technology advances, image caption generators hold immense potential for enhancing accessibility, communication, and engagement across domains ranging from social media and e-commerce to education and healthcare.

### **CHAPTER 2: LITERATURE REVIEW**

- The paper "Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning" contributes to the field of image captioning by proposing a mechanism that enables models to focus on the most relevant parts of an image while generating captions. This approach enhances the quality of generated captions and showcases the importance of adaptive attention in image understanding tasks.
- The paper "Neural Image Caption Generation with Visual Attention" contributes to the field of image captioning by introducing a visual attention mechanism to neural network-based models. This mechanism enhances the quality of generated captions by allowing the model to focus on relevant image regions while generating each word. The work paves the way for more sophisticated and context-aware image captioning models that align visual and linguistic information more effectively.
- The paper "Show and Tell: A Neural Image Caption Generator" presents a pioneering approach to image captioning using deep learning techniques,



specifically LSTM-based recurrent neural networks. The paper demonstrates the effectiveness of this approach in generating coherent and contextually relevant captions for images. The work contributes to the development of the field of image captioning and highlights the potential of neural networks in bridging the gap between computer vision and natural language processing tasks.

### **CHAPTER 3: OBJECTIVE**

- Enable AI to comprehend visual content and generate relevant captions.
- Ensure contextually coherent and creative language generation.
- Address biases and uphold ethical standards in both training and captions.
- Enhance practical applications such as content accessibility and engagement.
- Foster continuous improvement to remain adaptive and cutting-edge.

### **CHAPTER 4: METHODOLOGIES**

- Selecting an appropriate diverse dataset of images paired with human-generated captions. Also Analyzing the dataset to understand its distribution, image-caption relationships, and potential biases.
- Preprocessing the images (resizing, normalization) and tokenize the captions.
- Selecting a suitable model architecture for image caption generator. This involves combining CNNs for image feature extraction and RNNs or Transformers for caption generation.
- Training the image caption generator using the training set, monitoring loss and evaluation metrics.
- Evaluating the model's performance using metrics to measure the quality of generated captions.
- Fine-tuning the model by adjusting hyperparameters and optimizing the model architecture to achieve better performance.
- Present the results of your model's performance, including visualizations and qualitative assessment of generated captions.

## **CHAPTER 5: EXPERIMENTAL SETUP**

- **Data and Preprocessing:** Curate a diverse image-caption dataset (e.g., MS COCO) and preprocess images (resize, normalize) and captions (tokenize) for training.
- **Model Architecture:** Implement a CNN-based image encoder and an RNN or transformer-based text decoder for generating captions.
- **Training and Evaluation:** Train the model using appropriate loss functions, optimizers, and learning rate schedules. Evaluate performance using metrics like BLEU, METEOR, CIDEr, and human assessments.

## **CHAPTER 6: CONCLUSION AND FUTURE SCOPE**

- The image caption generator merges computer vision and language processing, translating visual content into coherent text.
- Its progress through neural networks and attention mechanisms has revolutionized content interpretation and communication.
- While challenges remain, ongoing advancements ensure responsible AI use and a transformative impact on accessibility and engagement.
- We can enhance the predictions by using more training examples. For example using Flickr32k dataset which has 32000 images
- Implement visual attention techniques, which focuses on interesting parts of the image

# REFERENCES

- [1] 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* (pp. 3156-3164).
- [2] Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., ... & Bengio, Y. (2015, June). Show, attend and tell: Neural image caption generation with visual attention. In *International conference on machine learning* (pp. 2048-2057). PMLR.
- [3] Lu, J., Xiong, C., Parikh, D., & Socher, R. (2017). Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 375-383).
- [4] 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* (pp. 6077-6086).
- [5] Li, Z., Tran, Q., Mai, L., Lin, Z., & Yuille, A. L. (2020). Context-aware group captioning via self-attention and contrastive features. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition* (pp. 3440-3450).
- [6] Kumar, N. K., Vigneswari, D., Mohan, A., Laxman, K., & Yuvaraj, J. (2019, March). Detection and recognition of objects in image caption generator system: A deep learning approach. In 2019 5th International Conference on Advanced Computing & Communication Systems (ICACCS) (pp. 107-109). IEEE.
- [7] Tanti, M., Gatt, A., & Camilleri, K. P. (2018). Where to put the image in an image caption generator. *Natural Language Engineering*, 24(3), 467-489.
- [8] Kinghorn, P., Zhang, L., & Shao, L. (2018). A region-based image caption generator with refined descriptions. *Neurocomputing*, 272, 416-424.