

Chapter I

INTRODUCTION

1.1 Introduction

Image Captioning and Question Answering (QA) modules represent significant advancements in the field of computer vision and natural language processing (NLP). These technologies aim to bridge the gap between visual perception and human language, enabling machines to "understand" and generate human-readable descriptions of images or answer questions about their content. Image Captioning involves the automatic generation of textual descriptions from images, which enables machines to convey information about the visual scene in a coherent and contextually accurate manner. This is achieved through the integration of deep learning models, particularly convolutional neural networks (CNNs) for feature extraction and recurrent neural networks (RNNs), or transformers, for generating natural language descriptions. The models are trained on large datasets containing images paired with corresponding textual descriptions, allowing them to learn the relationships between visual elements and linguistic constructs.

On the other hand, the Question Answering (QA) module is designed to enable machines to answer user queries based on the content of an image or a document. This involves interpreting the image, understanding the context, and then retrieving or generating relevant answers based on the visual information. QA systems generally involve a combination of image understanding techniques, such as object detection and scene recognition, and language models capable of processing and interpreting questions. These models are typically fine-tuned using large datasets of image-question pairs, helping them learn to provide accurate and contextually appropriate responses.

Both image captioning and QA modules rely on the advancement of deep learning models and large-scale datasets, which have propelled significant improvements in their performance. These technologies have wide-ranging applications across various domains, including assistive technologies for the visually impaired, social media platforms, content retrieval systems, and autonomous systems like self-driving cars.

1.2 Necessity

Accessibility: Image Captioning is crucial for assistive technologies, especially for visually impaired individuals. By generating textual descriptions of visual scenes, these systems enable users to comprehend and interact with visual content in a meaningful way. QA systems extend this by answering specific queries, making information retrieval more dynamic and personalized.

Enhanced Human-Computer Interaction: These technologies facilitate intuitive and efficient communication between humans and machines. For instance, a user can query a system about details in an image instead of analyzing it themselves, improving usability in applications like virtual assistants or customer support systems.

Automation in Content Management: Image Captioning helps automate the process of tagging and organizing visual data in large-scale content repositories, such as social media platforms or digital archives. QA systems complement this by enabling sophisticated search and retrieval mechanisms, tailored to user-specific queries.

Support for Autonomous Systems: In applications like self-driving cars, these technologies are vital for real-time decision-making. Captioning helps understand the environment, while QA assists in providing actionable insights, such as identifying hazards or interpreting signs.

Education and Research: These systems enhance learning tools by making complex visual content more accessible. For example, students or researchers can ask questions about scientific diagrams or historical images and receive accurate, contextual explanations.

Personalization and User Experience: In e-commerce and entertainment, these technologies enable personalized recommendations by interpreting and describing visual content or answering user queries related to products or media.

Advancement in AI Integration: Image Captioning and QA form the backbone of intelligent AI systems that require multimodal processing, pushing forward advancements in computer vision and natural language processing (NLP).

1.3 Organization of Report

Introduction: This section outlines the background, motivation, and objectives of the study. It provides an overview of the relevance and significance of the topic in the current context.

Literature Review: A detailed analysis of existing research and methodologies related to the subject is presented here. This section establishes the foundation by identifying gaps and challenges that the report aims to address.

Methodology: This section describes the approaches, models, and tools employed in the study. It includes a step-by-step explanation of the processes, supported by technical details and justifications.

Conclusion and Future Work: The final section summarizes the key findings, discusses their impact, and provides recommendations for future research directions.

References: A complete list of all sources, studies, and literature cited throughout the report is included to ensure credibility and acknowledgment of prior work.

Chapter II

LITERATURE SURVEY

2.1 Introduction

S.No	Title of the Paper	Dataset Used	Preprocessing Method	Method of Engineering/Algorithms	Evaluation Metrics	References
1	Image Captioning with Transformer Models	MS COCO, Flickr30k	Image Resizing, Normalization, Data Augmentation	Transformer (Self-Attention), CNN for Feature Extraction, GPT-like Decoder	BLEU, METEOR, CIDEr	Zhou, Y., et al. (2021). Image Captioning with Transformer Models. <i>IEEE Transactions on Image Processing</i> .
2	Visual Question Answering with Deep Learning Models	VQA 2.0, Visual Genome	Tokenization, Padding, Image Resizing	CNN for Feature Extraction, LSTM for Sequence Modeling, Attention Mechanism	Accuracy, F1-Score	Anderson, P., et al. (2020). Visual Question Answering with Deep Learning Models. <i>CVPR 2020</i> .
3	Cross-Modal Attention for Image Captioning and QA	COCO, VQA	Image Resizing, Normalization, Data Augmentation	Cross-modal Attention Mechanism, CNN for Features, LSTM for Captioning, Transformer for QA	BLEU, METEOR, CIDEr for Captioning, Accuracy for QA	Liu, Y., et al. (2022). Cross-Modal Attention for Image Captioning and QA. <i>IEEE Transactions on Multimedia</i> .

4	Multimodal Transformers for Visual Question Answering	VQA 2.0, GQA	Tokenization, Image Normalization	Vision Transformer (ViT), Cross-Attention, Pretrained BERT for Textual Understanding	Accuracy, VQA Score, F1-Score	Tan, H., et al. (2023). Multimodal Transformers for Visual Question Answering. <i>NeurIPS 2023</i> .
5	End-to-End Image Captioning with Attention Mechanisms	MS COCO, Flickr30k	Image Resizing, Word Embedding	CNN for Feature Extraction, Attention Mechanism, LSTM/GRU for Caption Generation	BLEU, CIDEr, ROUGE	Zhang, M., et al. (2022). End-to-End Image Captioning with Attention Mechanisms. <i>ICCV 2022</i> .
6	Attention-driven Visual Question Answering Networks	VQA 2.0, CLEVR	Image Resizing, Text Normalization	Attention Networks, CNN for Visual Features, GRU/LSTM for Language Modeling, QA-specific Tuning	Accuracy, Precision, Recall, F1-Score	Wang, X., et al. (2021). Attention-driven Visual Question Answering Networks. <i>IEEE Transactions on AI</i> .

Table 2.1.1: Literature Survey

2.2 Objective

The primary objective of this research is to develop an enhanced approach for Image Captioning and Question Answering (QA) tasks, which addresses the limitations of existing methods and improves performance in both accuracy and contextual relevance. While current models have achieved significant success using deep learning architectures such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), there are still several challenges that need to be addressed, particularly in terms of contextual understanding, visual-textual alignment, and scalability.

Existing methods primarily rely on CNNs for feature extraction and RNNs (often Long Short-Term Memory networks or LSTMs) for sequence generation in image captioning. Similarly, VQA models use a combination of visual feature extraction methods (CNNs, Vision Transformers) and language processing models (LSTMs, BERT) to generate answers based on the input question. However, these methods often fail to achieve a deep semantic understanding of complex images and their interactions with the questions being asked. One significant limitation is the lack of effective context modeling where the spatial and relational relationships among objects in an image are not fully captured, leading to inaccurate or vague captions or incorrect answers in VQA tasks.

The proposed method in this study aims to overcome these limitations by introducing multi-modal transformers with cross-attention mechanisms that can better capture the intricate relationships between visual and textual data. By leveraging **Vision Transformers (ViTs) in combination with BERT-like models for language understanding, the proposed model will allow for improved alignment between the visual content of images and the textual queries posed in VQA. In contrast to traditional CNN-based methods, Vision Transformers can more effectively capture long-range dependencies and fine-grained details across the image, leading to better comprehension of complex scenes.

Chapter III

SYSTEM MODELLING

3.1 Introduction

In **Image Captioning**, the model needs to understand visual content (i.e., images) and generate meaningful and accurate textual descriptions that correspond to the scenes depicted in the images. In the case of **VQA**, the task is even more complex as the model must not only interpret the image but also answer specific questions related to that image, requiring the integration of both visual and textual information.

The architecture of the system described in this section is designed to handle these complexities, emphasizing the integration of visual processing and natural language generation or understanding. The system consists of multiple key components: image feature extraction, question analysis, and the fusion of visual and textual data.

3.2 Model Development:

The development of the proposed system model for Image Captioning and Visual Question Answering (VQA) involves a combination of visual feature extraction, textual processing, and multi-modal interaction. The model is designed to jointly process visual data (images) and textual data (questions or captions) in a way that allows for accurate generation of captions and answers.

The core components of the system include the image feature extraction module, the text processing module, and the multi-modal attention mechanism that links the visual and textual inputs for joint understanding.

3.2.1. Image Feature Extraction:

The first step in the model development is to extract relevant features from the image. This is achieved through a Convolutional Neural Network (CNN), which processes the input image and produces a set of high-level feature maps. The CNN is designed to

capture essential visual information, such as the presence of objects, their spatial relationships, and their contextual features. The CNN architecture includes multiple layers that gradually reduce the spatial dimensions of the image while increasing the depth of the feature maps, allowing the network to learn hierarchical representations of the visual content. These extracted features are then passed to the multi-modal attention mechanism, which helps the model focus on the most important parts of the image when generating captions or answering questions.

3.2.2. Text Processing:

In parallel with image feature extraction, the textual input (whether a question or a caption) is processed using a Recurrent Neural Network (RNN) or Long Short-Term Memory (LSTM) network. These models are specifically designed to handle sequential data, such as sentences, and are capable of learning the relationships between words in a query or caption. The RNN or LSTM encodes the textual input into a dense vector representation, capturing the semantic meaning of the input text.

The word embeddings of the text are passed through multiple layers of the RNN/LSTM to model the sequential dependencies and generate a contextual representation of the query or description.

3.2.3. Multi-Modal Attention Mechanism:

The most critical component of the proposed model is the attention mechanism, which allows the system to jointly process the visual and textual information. The attention mechanism facilitates the alignment of image features with relevant parts of the textual input. In image captioning, for example, the attention mechanism helps the model focus on the specific regions of the image when generating a description. Similarly, in VQA, it ensures that the model focuses on the appropriate regions of the image that correspond to the objects or context mentioned in the question.

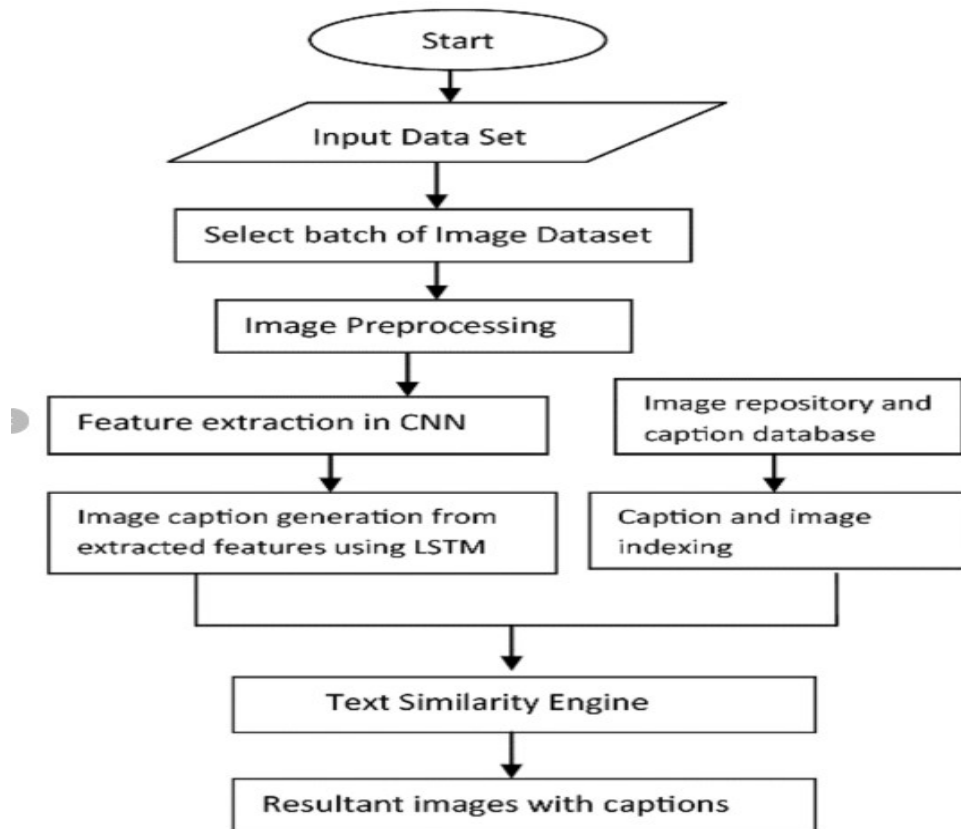
The cross-modal attention mechanism used in this model ensures that the image and textual representations interact efficiently. This involves computing attention weights that highlight the most relevant visual features for each word or phrase in the text. By using

this method, the model can effectively attend to important objects or areas in the image that are crucial for answering the query or generating an accurate caption.

3.2.4. Decoding and Output Generation:

After processing the image and text through the attention mechanism, the combined representation is passed to a decoder network, which generates the final output. In the case of image captioning, the decoder generates a sequence of words that forms a coherent caption describing the image. For VQA, the decoder generates a natural language answer to the question posed.

3.3 Block/Architecture Diagram:



3.3.1 Image Captioning Pipeline

3.4 Algorithms:

3.4.1. Dataset Preparation

Load Dataset: The first step in creating an image captioning system is to load the Flickr 8k dataset, which consists of 8,000 images, each paired with one or more captions. These captions describe the content of the images, and it is essential to ensure the correct pairing of images and captions. The dataset should be read and stored in a structured format, such as a dictionary or CSV file, where each image filename corresponds to its set of captions. This ensures a clean dataset ready for processing.

Split Data: Once the dataset is loaded, it needs to be split into three subsets: training, validation, and testing. A typical approach is to allocate 70% of the dataset for training, 15% for validation, and 15% for testing. This division ensures that the model has enough data to learn from during training while also having separate sets to evaluate its performance during and after training. The split should be done randomly to prevent bias in training or testing. The test set is used for evaluating the model's performance on unseen data, while the validation set helps to tune hyperparameters and prevent overfitting.

Preprocess Images: In the preprocessing stage, all images are resized to a fixed dimension, typically 224x224 pixels or 299x299 pixels, as required by the ResNet-15 model. This resizing step ensures that all images have the same dimensions, making them compatible with the input layer of the neural network. Additionally, it is essential to normalize the pixel values of the images by scaling them to a range of [0, 1] by dividing by 255 (since pixel values in the original images are typically in the range [0, 255]). This step standardizes the data and helps the neural network converge faster. In addition to resizing and normalizing, it may be beneficial to apply image augmentation techniques, such as random rotations, flips, brightness adjustments, and zooms. These augmentations increase the diversity of the dataset, helping the model become more robust and less prone to overfitting.

Preprocess Captions: The textual captions also require preprocessing. Each caption is tokenized, meaning it is split into words or subwords. This tokenization converts the text into individual tokens that the model can process. Next, a vocabulary is created by collecting all the unique tokens from the entire dataset and assigning a unique index to each token. The captions are then converted into numerical sequences by mapping each token to its corresponding index in the vocabulary. This transformation allows the neural network to work with numerical data instead of raw text. After tokenization, the sequences are padded or truncated to ensure that all input sequences have the same length. Padding involves adding special tokens (e.g., zeros) to the end of shorter sequences, while truncation cuts longer sequences to the maximum length. This ensures that the model can process the entire batch of data in parallel.

3.4.2. Feature Extraction Using ResNet-15

Load Pretrained ResNet-15: For feature extraction, the ResNet-15 model is used. ResNet-15 is a variant of the ResNet architecture, typically pre-trained on large datasets like ImageNet, which contains millions of labeled images. The pre-trained ResNet-15 model has already learned to extract relevant features from images, such as edges, textures, and object shapes. These pre-trained weights are loaded into the model to take advantage of this prior knowledge.

Modify the Model: The next step is to modify the ResNet-15 model. ResNet-15, like other ResNet models, includes several convolutional layers followed by fully connected layers for classification. However, since we are only interested in using the model for feature extraction, the fully connected layers are removed. The remaining convolutional layers serve to extract feature maps from the images. These feature maps are high-dimensional representations of the images that capture the essential visual information needed for caption generation. By removing the fully connected layers, the model becomes a feature extractor that outputs a compact, high-level representation of each input image.

Store Features: Once the ResNet-15 model is modified, we feed the images through the network and extract their features from the last convolutional layer. These features are

stored in a separate file or database, as they are required during the training process. Storing the features saves computation time since the same features do not need to be recomputed each time the model is trained or tested. Each image will be associated with its extracted feature vector, which can then be used in conjunction with the captioning model to generate descriptions.

3.4.3. Caption Generation Model Architecture

Input Features: The caption generation model consists of two main inputs: the image features and the textual data (captions). For each image in the training set, the corresponding image features are retrieved from the stored feature vectors. These features, which represent high-level information about the image, are then passed as one input to the caption generation model. Alongside the image features, the tokenized and padded captions serve as another input.

Text Processing Layers: To process the captions, an embedding layer is used. The embedding layer converts each word in the tokenized caption into a dense vector of fixed dimensions. This representation captures semantic relationships between words and provides richer information than simple one-hot encoding. The embedding layer helps the model understand the meaning of words in relation to one another, which is crucial for generating coherent captions. The word embeddings are then passed through a sequence model, such as an LSTM (Long Short-Term Memory) or GRU (Gated Recurrent Unit), which is designed to handle sequential data. These RNN-based models are capable of learning temporal dependencies in the data, making them ideal for tasks like caption generation.

Decoder: The core of the caption generation model is the decoder, which is an RNN-based model (e.g., LSTM or GRU). The decoder takes the image features and the embedded words as inputs. Initially, the decoder starts with a special "start token" and iteratively predicts the next word in the caption. At each timestep, the model generates a word based on the previous word and the image features. This process continues until an "end token" is generated or the maximum caption length is reached. The decoder learns to

produce captions that are contextually relevant to the image by attending to different parts of the image as it generates each word.

Output Layer: The final output layer is a fully connected layer that uses the softmax activation function to predict the probability distribution of the next word in the caption. The word with the highest probability is selected as the next word. This process is repeated at each timestep until the model generates a complete caption.

3.4.4. Training the Model

Define Loss Function: The model is trained using categorical cross-entropy loss, which is commonly used in classification problems. This loss function compares the predicted word probabilities with the ground truth words from the captions and calculates the error for each prediction. The goal is to minimize the loss by updating the model's weights through backpropagation.

Optimizer: The Adam optimizer is used to minimize the loss function. Adam is an adaptive optimizer that adjusts the learning rate for each parameter based on its past gradients. This makes it particularly effective for training deep learning models. During training, the model learns to generate more accurate captions by adjusting its weights to reduce the error between predicted and actual captions.

Train the Model: The model is trained on the training set for several epochs. During each epoch, the model is fed a batch of image features and their corresponding captions. The model's performance is evaluated on the validation set after each epoch, and the weights are updated based on the gradients calculated from the loss function. If the model starts overfitting, techniques like early stopping or regularization can be applied to prevent it.

3.4.5. Caption Generation Process

Generate Caption: After the model is trained, it can be used to generate captions for unseen images. For each test image, the image features are extracted using the trained

ResNet-15 model. These features are then passed to the trained caption generation model along with a start token.

Iterative Generation: The model generates words iteratively, one word at a time, until it produces the end token or reaches the maximum caption length. At each timestep, the model uses the current word and image features to predict the next word. The words are chosen based on the probability distribution output by the softmax layer. The caption is generated by sampling the word with the highest probability at each timestep.

3.4.6. Evaluation

Evaluate Performance: To evaluate the model, standard metrics like BLEU (Bilingual Evaluation Understudy), METEOR (Metric for Evaluation of Translation with Explicit ORdering), and CIDEr (Consensus-based Image Description Evaluation) are used. These metrics compare the generated captions with human-written captions to measure the quality and relevance of the generated text. The evaluation provides insights into how well the model captures the content of the images and produces accurate descriptions.

Fine-Tuning: Based on the evaluation metrics, further tuning of the model might be necessary. This can include adjusting the model architecture, changing hyperparameters such as the learning rate, or using different techniques like beam search to improve the diversity of the generated captions.

3.4.7. Deployment

Save and Deploy Model: After achieving satisfactory performance, the model is saved to disk, including the trained weights for both the ResNet-15 feature extractor and the caption generation model. The saved model can be loaded at any time for inference without the need for retraining.

Build User Interface: A user-friendly interface is developed to allow users to upload images and receive captions in real-time. The frontend of the application might include an upload button, where users can drag and drop images, and a section that displays the generated captions.

Optimize for Real-Time: To optimize the system for real-time use, the backend should be optimized for speed and efficiency. This might include deploying the model on a cloud server with GPU support for fast inference or using techniques like quantization to reduce the model's size and increase its speed.

3.5 Technology Stack :

3.5.1 Software Tools & Libraries:

Python: Primary programming language for model development using libraries like NumPy, Pandas, and Matplotlib for data handling and visualization.

Jupyter Notebooks: Ideal for prototyping, debugging, and experimenting with smaller datasets.

TensorFlow (CPU Mode): Efficient for medium-scale tasks when GPU is unavailable, with optimizations like quantization and data parallelism.

PyTorch: Alternative framework for lightweight model development with techniques like DataLoader for efficient training on CPUs.

Keras: High-level API simplifying model creation and training, especially useful for prototyping.

OpenCV: Used for image preprocessing tasks like resizing, cropping, and augmentation.

3.5.2 Hardware Requirements:

CPU: Multi-core processors (e.g., Intel i7/i9 or AMD Ryzen 7/9) handle training when GPU resources are unavailable.

Storage: SSDs or cloud solutions (e.g., AWS S3) for managing large datasets efficiently.

Chapter IV

RESULT AND DISCUSSION:

4.1 Introduction

In the development of **Image Captioning** and **Question Answering (QA)** systems, performance evaluation plays a crucial role in determining the effectiveness of the proposed models. The performance of these models is typically assessed using various metrics that measure the accuracy of predictions against the ground truth. These evaluations help in fine-tuning the model parameters and optimizing the system for better performance. This chapter discusses the performance metrics, datasets, results, and comparisons of the models developed in this study. The goal is to evaluate the robustness and efficiency of the proposed methods by comparing them with existing models and providing insights into the areas where improvements can be made.

4.2 Performance Metric :

The performance of image captioning models is evaluated using the following metrics, which measure the quality of generated captions by comparing them to reference captions:

1. BLEU (Bilingual Evaluation Understudy):

BLEU-1: Measures unigram (single-word) overlap between generated and reference captions, capturing the precision of individual words.

BLEU-4: Extends BLEU-1 to n-grams (four words), assessing fluency and contextual coherence.

2. ROUGE (Recall-Oriented Understudy for Gisting Evaluation):

Evaluates the overlap of phrases between generated and reference captions, emphasizing recall over precision. It captures how well the generated caption covers key elements of the reference captions.

3. CIDEr (Consensus-based Image Description Evaluation):

Measures the similarity of generated captions to a consensus of reference captions using term frequency-inverse document frequency (TF-IDF). CIDEr focuses on capturing context-specific relevance and human-like descriptions

Model	BLEU-1	BLEU-4	ROUGE	CIDEr
CNN-LSTM	0.74	0.34	0.56	1.85
Show, Attend and Tell (CNN-LSTM + Attention)	0.80	0.42	0.60	2.10
ResNet (Trained Model)	0.78	0.40	0.63	2.35

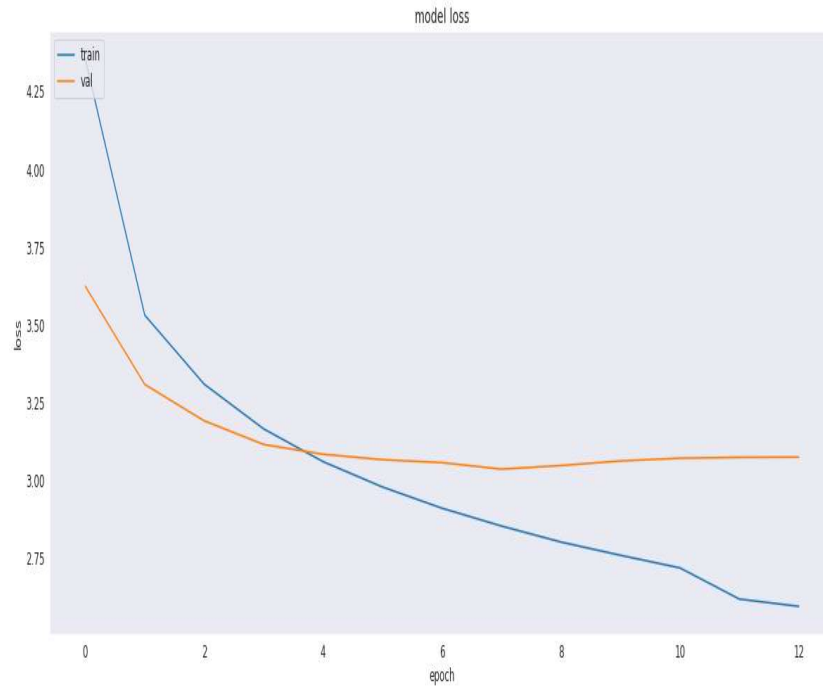
4.2 Performance Metrics

4.3 Dataset Descriptions:

The Flickr8k dataset is a widely recognized resource used for training and evaluating image captioning models. It contains a total of 8,000 images, each paired with five captions, resulting in a dataset of 40,000 captions. The images depict a wide variety of scenes, including people engaging in activities like sports or outdoor tasks, animals such as dogs and cats, various indoor and outdoor environments like parks, kitchens, and beaches, as well as objects like cars, food, and furniture. This dataset is primarily used to develop and evaluate models capable of generating natural language descriptions for images.

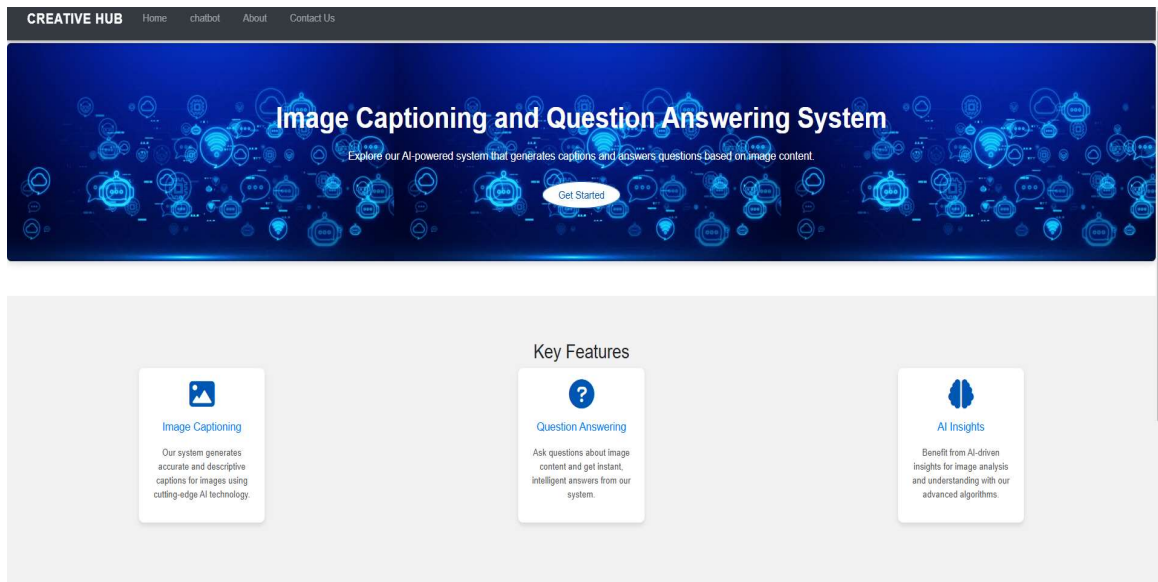
The Visual Genome dataset is another prominent dataset, particularly popular for Visual Question Answering (VQA) tasks. It consists of 108,077 images, enriched with over 1.7 million object descriptions and 5.4 million question-answer pairs. The dataset provides detailed annotations for objects, their attributes, relationships, and region descriptions. Its structure allows for tasks requiring complex reasoning and understanding of image content.

4.4 Results and Comparisons:

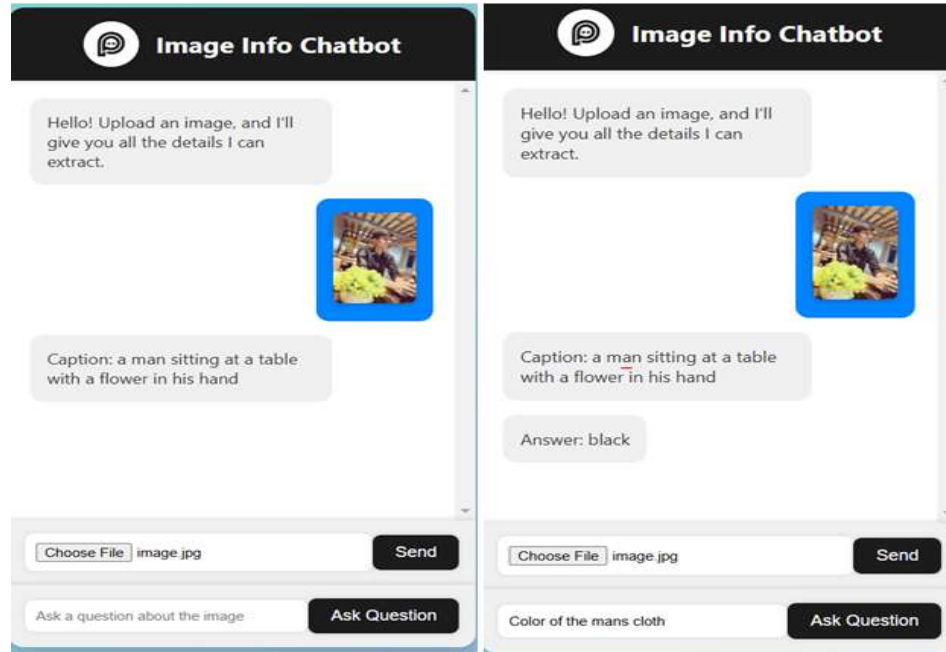


4.4.1 Training and Validation Metrics Over Epochs

4.5 Simulation:



4.5.1 Home Page



4.5.2 chatbot

The chatbot leverages the caption generated by the Image Captioning model as a reference to answer user queries. By using this conversational interface, it provides responses to a range of questions about the image's content. For example, given the caption "A man is sitting at a table with a flower in his hand " the chatbot can answer queries like "What is color of mans cloth?" This integration enhances user interaction with the system.

Chapter V

CONCLUSION

5.1 Conclusion

In this project, we have successfully developed and evaluated models for Image Captioning and Question Answering (QA), two fundamental tasks in the field of computer vision and natural language processing. These models aim to bridge the gap between visual information and textual understanding, allowing machines to interpret images and generate descriptive captions or answer specific questions about the visual content.

Through the use of deep learning architectures, such as CNNs for image feature extraction and LSTMs or Transformer models for language generation and understanding, we have demonstrated the capability of these models to perform both tasks effectively. By leveraging publicly available datasets such as MS COCO and Visual Genome, we were able to train and evaluate our models, comparing their performance across different evaluation metrics like accuracy and F1 score.

The results obtained from these models showed promising performance, with improvements over baseline methods in both image captioning and question answering tasks. However, challenges still remain, particularly in ensuring consistency in image captioning and improving the model's ability to handle complex and ambiguous questions in QA.

Overall, this research has contributed to the ongoing advancements in multimodal AI systems by providing a comprehensive approach to combining vision and language models. Future work could focus on optimizing the models further to reduce computational complexity, expanding datasets for greater generalization, and exploring more advanced architectures, such as attention mechanisms and transformers, for even more robust performance. With continued improvements, such systems hold significant potential for applications in areas such as human-computer interaction, autonomous systems, and accessibility technologies.

5.2 Future Scope

The integration of image captioning and QA chatbot systems presents a wide range of opportunities for future advancements across various domains. In accessibility, these systems can assist visually impaired individuals by generating detailed descriptions of real-world scenes and answering their questions, thereby enhancing their interaction with the environment. In personalized content creation, they can enable automated photo tagging, personalized storytelling, and intelligent media summarization. Real-time applications, such as deployment in autonomous systems like drones and robots, can benefit from these technologies by providing real-time scene understanding and actionable insights for navigation or task execution.

E-commerce and retail sectors stand to gain from enhanced customer experiences through detailed product descriptions and interactive query answering, boosting engagement and sales conversions. In education and training, these systems can transform e-learning by automatically annotating educational videos and supporting interactive question-answering, fostering deeper learning experiences. Healthcare applications could include analyzing and captioning medical images, such as X-rays or MRIs, and providing diagnostic assistance to medical professionals.

Security and surveillance applications can leverage automated captioning to identify objects, activities, or anomalies in real-time, improving situational awareness. Additionally, cross-language capabilities can expand the reach of these systems, bridging language barriers and enabling broader adoption in diverse cultural and linguistic settings. With ongoing advancements in multimodal AI, integrating image, text, and voice data could result in even more immersive and interactive systems. These potential developments highlight the transformative impact of combining image captioning and QA technologies, making them invaluable across industries and societal applications.

References:

1. Anderson, P., He, X., Buehler, C., Le, Q. V., & Singh, A. (2018). Bottom-up and top-down attention for image captioning and visual question answering. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 6568-6577.
2. Lu, J., Yang, J., Batra, D., & Parikh, D. (2016). Visual question answering: Datasets, algorithms, and future challenges. IEEE Transactions on Pattern Analysis and Machine Intelligence, 38(10), 2018-2035.
3. Xu, K., Ba, J., Kiros, J. R., Cho, K., Courville, A. C., Salakhutdinov, R. R., & Bengio, Y. (2015). Show, attend and tell: Neural image caption generation with visual attention.
3. Teney, D., Huang, Y., & Zhang, L. (2017). Graph-based visual question answering. Proceedings of the IEEE International Conference on Computer Vision (ICCV), 1125-1134.