

Digital Method to Talk with Machine using Image Processing And Machine Learning Techniques

*A Major Project Report Submitted in the
Partial Fulfillment of the Requirements
for the Award of the Degree of*

BACHELOR OF TECHNOLOGY

IN

INFORMATION TECHNOLOGY

Submitted by

G Tejashwini

18881A1224

SUPERVISOR

Dr Mukta Jagdish

Associate Professor

Department of Information Technology



VARDHAMAN COLLEGE OF ENGINEERING

(AUTONOMOUS)

Affiliated to JNTUH, Approved by AICTE, Accredited by NAAC with A++ Grade, ISO 9001:2015 Certified
Kacharam, Shamshabad, Hyderabad - 501218, Telangana, India

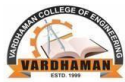
May, 2022

DECLARATION

We hereby declare that the work described in this report entitled “ **Digital Method to Talk with Machine using Image Processing and Machine Learning Techniques** ”which is being submitted by us in partial fulfillment for the award of **BACHELOR OF TECHNOLOGY** in the Department of Information Technology, Vardhaman College of Engineering to the Jawaharlal Nehru Technological University Hyderabad.

The work is original and has not been submitted for any Degree or Diploma of this or any other university.

G Tejashwini



VARDHAMAN COLLEGE OF ENGINEERING

(AUTONOMOUS)

Affiliated to JNTUH, Approved by AICTE, Accredited by NAAC with A++ Grade, ISO 9001:2015 Certified
Kacharam, Shamshabad, Hyderabad – 501218, Telangana, India

Department of Information Technology

CERTIFICATE

This is to certify that the project titled **Digital Method to Talk with Machine using Image Processing and Machine Learning Techniques** is carried out by

G Tejashwini

18881A1224

in partial fulfillment of the requirements for the award of the degree of
Bachelor of Technology in Information Technology during the year
2021-22.

Signature of the Guide

Dr Mukta Jagdish

Associate Professor

Signature of the HOD

Dr. Muni Sekhar Velpuru

Associate Professor and Head, IT

Project Viva-Voce held on _____

Examiner

Acknowledgement

The satisfaction that accompanies the successful completion of the task would be put incomplete without the mention of the people who made it possible, whose constant guidance and encouragement crown all the efforts with success.

We wish to express our deep sense of gratitude to **Dr Mukta Jagdish**, Associate Professor and Project Supervisor, Department of Information Technology Vardhaman College of Engineering, for her able guidance and useful suggestions, which helped us in completing the project in time.

We express our heartfelt thanks to **Dr. M. Ramachandra**, Associate Professor Project Coordinator, for his suggestions invaluable inputs and assessment really helped us to shape this report to perfection

We are particularly thankful to **Dr. Muni Sekhar Velpuru**, the Head of the Department, Department of Information Technology, his guidance, intense support and encouragement, which helped us to mould our project into a successful one.

We show gratitude to our honorable Principal **Dr. J.V.R. Ravindra**, for providing all facilities and support.

We avail this opportunity to express our deep sense of gratitude and heart-ful thanks to **Dr. Teegala Vijender Reddy**, Chairman and **Sri Teegala Upender Reddy**, Secretary of VCE, for providing a congenial atmosphere to complete this project successfully.

We also thank all the staff members of Information Technology department for their valuable support and generous advice. Finally thanks to all our friends and family members for their continuous support and enthusiastic help.

G Tejashwini

Abstract

Digital Method to Talk with Machine the usage of Image Processing and Machine Learning Techniques is a mQA version, that is capable of solution questions on the content material of an image. The solution may be a sentence, a word or a unmarried word. Our version has four additives: an LSTM to extract the query representation, a CNN to extract the visible representation, an LSTM to store the linguistic context in a solution, and a fusing mechanism to combine the statistics from the first three additives and build the solution. To teach and compare our mQA version, we create a Freestyle Multilingual Image Question Answering (FM-IQA) dataset. It contains almost 150,000 snapshots as well as 310,000 freestyle Chinese query-solution pairs with English translations. The pleasant of the generated solutions of our mQA version in this dataset is evaluated with the aid of using human judges thru a Turing Test. Specifically, we blend the solutions supplied with the aid of using human beings and our version. The human judges want to differentiate our version from the human.

Specifically, we combine the solutions provided by employing humans with our version. The human judges seek to distinguish our version from the human version. They will even provide a rating (i.e., 0, 1, 2, the larger the better) reflecting the quality of the response. We provide approaches for revealing the first-class of this evaluation procedure. The results show that in 64.7 percent of situations, human judges are unable to identify our version from humans. The average rating is 1.454. (1.918 for human). This work's information, which includes the FM-IQA dataset.

Keywords: Long Short-Term Memory(LSTM); Recurrent Neural Network; Convolutional Neural Networks; Visual Questioning Answering

Table of Contents

DECLARATION	
Title	Page No.
Acknowledgement	i
Abstract	ii
List of Figures	v
Abbreviations	v
CHAPTER 1 Introduction	1
1.1 Introduction	1
1.2 Scope	3
1.3 Purpose	3
1.4 Objectives	4
1.5 Advantages of Project	4
1.6 Module	5
CHAPTER 2 LITERATURE SURVEY	7
2.1 Existing Systems	8
2.2 Limitations of Existing Systems	11
2.3 Proposed Method	11
2.3.1 Data Loader	12
2.3.2 Feature Engineering	14
2.3.3 The weight sharing strategy	15
2.3.4 Training Details	16
2.3.5 The Freestyle Multilingual Image Question Answering (FM-IQA) Dataset	16
CHAPTER 3 ANALYSIS	19
3.1 Introduction	19
3.2 Software Requirement Specification	20
3.2.1 User requirement specification	20
3.2.2 Software requirement	20
3.2.3 Hardware requirement	20
3.3 Algorithm and Flow Chart	21

CHAPTER 4 DESIGN	24
4.1 Introduction	24
4.2 Diagrams	24
4.2.1 ER Diagram	24
4.2.2 Use Case	25
4.3 Module Design and Organization	25
4.4 System Architecture	27
CHAPTER 5 IMPLEMENTATION	28
5.1 Introduction	28
5.2 Explanation of Key Function	28
5.3 Technology	30
5.4 Method of Implementation	30
5.4.1 Output Screens	31
5.4.2 Result Analysis	33
CHAPTER 6 TESTING RESULTS	35
6.1 The Visual Turing Test	35
6.2 The Score of the Generated Answer	35
6.3 Performance Comparisons of the Different mQA Variants	36
CHAPTER 7 Conclusions and Future Scope	37

List of Figures

1.1	Input and Output.....	2
1.2	Feature Extraction.....	6
2.1	Structure of VQA Dataset.....	12
2.2	Architecture of Extracting Feature from Image.....	14
2.3	BiLSTM encoder with max-pooling.....	15
3.1	Example of a simple RNN unfolding in time.....	21
3.2	Flow Graph.....	23
4.1	ER Diagram for VQA Model.....	24
4.2	Use Case diagram of VQA.....	25
4.3	Model Overview.....	27
4.4	System Overview.....	27
5.1	Feature Extraction using VGG16.....	29

Abbreviations

Abbreviation	Description
ML	Machine Learning
RNN	Recurrent Neural Network
CNN	Convolutional Neural Network
FM-IQA	Freestyle Multilingual Image Question Answering

CHAPTER 1

Introduction

1.1 Introduction

In recent years, there has been a surge of interest in multi-modal reading for natural language processing and image processing. Many studies, in particular, have achieved rapid progress in photo subtitling. The vast majority of these tests use deep brain organisations (for example, deep Convolutional Neural Networks (CNN), Recurrent Neural Networks (RNN), or Long Short-Term Memory (LSTM)). This sequence of events need extensive visual data bases with literary substantive explanations. Regardless of how effective those tactics are, there are some difficulties that need be addressed and studied. For example, picture inscription necessitates extensive language representations of photographs. Regardless, in a lot of cases, a particular aspect as well as component of a photograph piques our interest. There may be no interchange between the device and the purchaser inside the Image subtitling activity because we are unable to enter our decisions and interests. In the compositions of notable question responding, this review provides substantial authority. A free-shape request at substance material of an image on this work should be answered by the methodology. To address this problem, we support the mQA worldview. The adaptation's parts of feedback are a photograph and a question. The alternative theories explain the reaction to the query and produce the condition of words on the inside. There are currently over 1.5 million photos in the series, as well as 3,600 query-solution pairings and their English translation. To diversify the annotations, annotators are authorised to invite any question about the image's content material. Another area of investigation that has gotten a lot of attention and energy in recent years is visual inquiry in response to (VQA). VQA can be thought of as a supplement to the concept of machine perception. It's also a multi-disciplinary AI project



电脑在哪里？
Where is the computer?
在桌子上。
On the desk.



这是在什么地方？
Where is this?
这是在厨房。
This is the kitchen room.

Figure 1.1: Input and Output

that combines advanced computer vision with regular language processing to create a framework that can answer a question about a picture. The model aims to comprehend the image's hidden meaning and semantics and respond to questions based on its "understanding." Unlike picture recordings, where basic data on PC vision and NLP suffices to construct AI models, tasks like VQA necessitate a comprehensive understanding of state-of-the-art methodologies in both of these areas. Due of a lack of cutoff for more significant reasoning, the vast majority of AI structures gathered thus far have failed to equal individuals on obvious level vision tasks. Anyway, with the continued evaluation in this subject, it is now conceivable to attempt to construct a system that should be successful in undertakings such as VQA. More basic essential ability to perceive is usually required in such a system. Capabilities in image and sophisticated thinking A more advanced version of this Similarly, the system would have actual information to make a decision. on and respond to questions where answers are really not clearly represented in the image.

1.2 Scope

Talking to a Machine is a mQA version that may solution inquiries approximately an image's content. The emotion might be expressed as sentence,

a phrase, or maybe a single word. Our version consists of four parts: an LSTM to remove the question portrayal, a CNN to extricate the noticeable portrayal, an LSTM to maintain the language setting in a response, as well as a melding component to combine the measurements from the previous three additional substances and collect the arrangement. We create a Freestyle Multilingual Image Question Answering (FM-IQA) dataset to train and test our mQA adaption. Around 150,000 images are included, as well as 310,000 free-form question arrangement matching with English translations. Human-appointed authorities utilise the Turing Test to evaluate the beauty of our mQA variant's provided reactions in this dataset.

We particularly incorporate human reactions with our form. The human appointed authorities ought to segregate among our form and the human. They could give a score (e.g., 0, 1, 2, the more the number, the more) mirroring the wonderful of the reaction. We propose approaches for following the evaluation cycle's lovely. The impacts uncover that human appointed authorities can not illuminate our variant from people in 64.7 level of circumstances. The normal is 1.454. (1.918 for human). This work's points of interest, comprehensive of the FM-IQA dataset.

1.3 Purpose

Visual query answering structures are trying to find to correctly solution visual language queries approximately an photograph input. The overarching aim of this subject matter is to create structures that may apprehend the contents of a photograph within side the equal manner that humans can and speak successfully approximately that photograph in verbal language.

1.4 Objectives

- Capability to reply to inquiries concerning the content material of a picture.
- Development of a Freestyle Multilingual Image Question Answering (FM-

IQA) dataset to teach and check the mQA model.

- The nature of the mQA model's produced reactions in this dataset is assessed through human appointed authorities the utilization of a Turing Test.
- It is executed through designing user-pleasant facts access panels that may control massive quantities of facts. The cause of designing enter is to make facts getting into simpler and error-free. The facts getting into web page is installation in any such manner that any facts manipulations are possible. It additionally has file viewing capabilities.
- At the point when the information is placed, it will be approved. Information might be input by utilizing screens. Suitable messages are conveyed depending on the situation with the goal that the client isn't trapped in a hopeless cycle. In this manner, the objective of information configuration is to give a simple to-follow input format.

1.5 Advantages of Project

The use of speech is deeply ingrained in man's psychological constitution, and there are certain general benefits of voice communication as well as some benefits of particular importance to astronauts. A voice interface (in which the numbers following cross- reference past benefits):

- Provides for more natural communication
- Increases communication capacity through multi-modal communication
- Allows communication without the need for specialist training
- Physiological/psychological monitoring of the operator's condition may be possible, at least as a backup to electrical sensing.
- Is compatible with the widely accessible and low-cost telephone system
- Allows for security checks

1.6 Module

In this project, we used Deep-Learning and Neural networks to train the system to identify images and interpret questions, with the input being a picture and a corresponding question in the form of text, and the output being printed text.

1. Data Collection and Preparation

As the essential photo assortment, we utilize the 158,392 previews from the these days posted MS COCO preparing, approval, and looking at set. Baidu's internet-based publicly supported server⁴ is used to collect the comments. To supplement the stated questioning and response pairings, annotators are free to ask any form of question as long as it is relevant to the image's content. The visual material and mental presence of mind should respond to the request. The annotators should respond to the question.

2. Data Cleaning and Preprocessing

- Convert all words to lowercase.
- Eliminate all accentuation. all words that are one character or less long (for example 'a').
- Eliminate all words with numbers in them.
- Image-id + question about the image + answer " Is a tuple in our Dataset for training and Validation

3. Extracting feature from image VGG16 is an photograph reputation version created with the aid of using Oxford. We eliminate the very last category layer from this Convolutional Neural Network Model and get and save the characteristic vector of every photo independently.

4. Training Datapoints Let an answer be of five words and the sequence of words be $i_{start} \leq 34,45,3,2,67, i_{end}$ for an image i and question q . Then we divide the part of answer + question + image

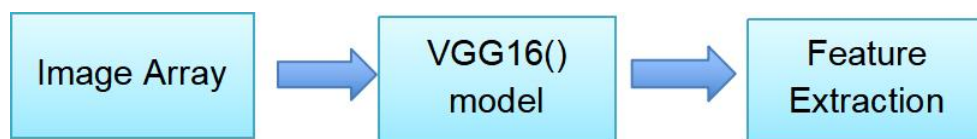


Figure 1.2: Feature Extraction

CHAPTER 2

LITERATURE SURVEY

A variety of new guides have started to analyze visible query answering. However, in assessment to our study, they're very restrained (now and again manufactured) contexts with little datasets. For example, evaluates most effective queries with responses from a detailed restrained surroundings of sixteen fundamental coloring or 894 object types. Considers questions constructed from templates primarily based totally on a fixed vocabulary of objects, properties, and connections among objects, among different things.

Our advised work, on the other hand, contains open-finished, free-shape questions and reactions given through method of method for people. Our objective is to build the assortment of realities and addressing capacities expected to supply right reactions. Our VQA dataset is requests of significant worth bigger than (250,000 versus 2,591 and 1,449 pictures, separately), that means a lot to accomplishment in this more prominent troublesome and unconstrained work. The proposed VQA process is similar as going before appropriate artworks in that it has explored agreeable parsing of movies and going with text based content to answer requests on datasets each along with 15 video clips.leverages publicly supported work to address outwardly debilitated shoppers' requests concerning visual substance In equal work, joining a LSTM was proposed for the inquiry and a CNN for the image to offer a response - a comparable model is tried in this review. produces theoretical views to catch visual sound judgment pertinent to fill-in-the-clear and visual summarizing inquiries (simply printed). Survey the acceptability of good judgment suggestions utilizing visual data. offered an assortment of 10,000 photographs and requested subtitles depicting significant components of a scene (e.g., individual items, what will occur straightaway).

During our work, we gathered questions and reactions (later deciphered into English) for COCO photographs. COCO inscriptions have been utilized to

give

4 styles of inquiries (thing, count, variety, and area). While the utilization of open-finished questions offers many advantages, it is as yet valuable to comprehend the sorts of inquiries that are being posed and which types different calculations might great at reply. To this end, we investigate the sorts of inquiries posed and the kinds of answers gave. Through a few representations, we show the shocking variety of the inquiries posed.

We additionally investigate how the data content of inquiries and their responses varies from picture inscriptions. For baselines, we offer a few methodologies that utilization a blend of both text and cutting edge visual highlights [29]. As part of the VQA drive, we will sort out a yearly test and related studio to examine cutting edge techniques and best practices

2.1 Existing Systems

1. Dataset

Machine Learning(ML) is a subset of AI wherein information is basic for preparing the model and creating compelling outcomes. Profound learning calculations, specifically, contain numerous parameters and hyper-boundaries that should be tuned. This happens during the preparation process, while the model is still in the "learning" stage. The model can turn out to be more generalised as it is prepared on extra information. Thus, great quality information is a critical need for executing any AI calculation. Profound learning is likewise utilized in the majority of VQA models. This means that a lot of top notch information is necessary to prepare the models. Beginning around 2014, numerous VQA datasets have been created. The majority of these datasets are uninhibitedly accessible for scholarly exploration.

- **VQA**

VQA is the most notable and usually utilized dataset for tending to visual questions.It was worked as a component of a VQA rivalry in 2015. The dataset is separated into two areas:

Actual photographs

Clip-art scenes that are abstract

Utilizing feature engineering pictures reduces the need to pre-cycle and clean uproarious photographs. They might be used to straightforwardly do confounded thinking. Each image in the dataset is additionally joined by a rundown of inquiry answer pairings.

- **Visual7W**

In 2016, the Visual7W dataset [5] was delivered. Pictures from the MS-COCO dataset are remembered for this assortment. The dataset of inquiries was developed by posing to seven "w" questions: what, where, when, why, who, how, and which. The inquiries were numerous decision, with four reaction options for each.

Bouncing boxes in the photographs were utilized to address each article referenced in the inquiry answer pairings. Publicly supporting Amazon Mechanical Turks was utilized to assemble the inquiry answer matches and object-level groundings (AMT). This dataset contains about 47,000 photographs and more than 320,000 inquiries.

- **DAQUAR**

In 2015, the Dataset for QUESTION ANSWERING ON REAL-WORLD IMAGES (DAQUAR) [6] turned out to be likewise delivered. It was the first dataset at any point sent off for this work (even before the first VQA dataset). The snap pictures for this dataset were assembled by the creators from the NYU-Depth V2 dataset. Each picture in the assortment is marked generally founded on the things showed inside the picture. The gadgets could can be categorized as one of the 894 potential thing arrangements. The dataset is nearly nothing. It just has around 1500 photos. The inquiry arrangement pairings were built naturally involving nine pre-characterized formats for the inquiries. The responses to those questions were acquired quickly from the NYU-Depth V2 dataset.

2. Feature Engineering

- **Feature Engineering for Images:**

Many open-source pretrained models are available to procure feature vectors for photos. Google presented one such model, GoogLeNet, in 2014 [7]. It has now been widely used to encode pictures and make phenomenal part vectors for use in any significant learning model. It is a 22-layer convolutional neural network (CNN) that was trained on the ImageNet dataset. VGGNet [1, which is moreover a CNN] is another model that is comparable. It has 16-19 layers and has next to no convolution channels (size of 3x3).

- **Feature Engineering for Text:**

To get highlight vectors from crude text input, pre-prepared models are accessible. These element vectors are alluded to as word vectors. There are openly available pre-trained models for acquiring related word vectors from input text information. GloVe (Global Vectors for Word Representation) [12] and Word2Vec [7] are the most frequently utilized word vectors.

3. Memory

- **Gated recurrent units (GRU):**

GRUs and LSTMs are very comparative. The inward instrument varies to some degree. GRUs, dissimilar to LSTMs, need interior memory that is particular from the hid state. The application of GRUs is like that of LSTMs. They might be utilized to accomplish question implanting as well as to offer the model with all remembering abilities. GRUs are utilized in [9] to build a unique memory network planned only for visual inquiry addressing frameworks.

- **Attention models:**

Long haul conditions are addressed more successfully utilizing consideration models than with essential LSTMs. Straightforward LSTMs require encoding the whole text as organization input. Notwithstanding, in models that incorporate a consideration instrument,

it isn't important to encode the total text as contribution to the model.

All things considered, the model "focuses" on select parts of the information that are more critical in the ongoing setting and gives those parts more noteworthy weightage at each phase of figuring the result. Numerous VQA frameworks have tried to add a consideration instrument into their plan, with prominent upgrades when contrasted with plain LSTMs or GRUs. Question-based consideration was used to encode pictures in most of circumstances.

2.2 Limitations of Existing Systems

Existing project works on image captioning with RNN and CNN, but there is no feature to process voice queries and extract replies largely based only on image configuration. They concatenate the question and response before feeding them to the LSTM. Unlike them, we use separate LSTMs for questions and solutions in conjunction with unique households of questions and solutions, allowing for word-combing interchange. The issue with this strategy is that the hobby things may be in different spatial places and have varied thing ratios in the photograph.

2.3 Proposed Method

To avoid the hassle of selecting a huge number of locations, Ross Girshick et al. presented an approach in which we utilise selective search to extract just 2000 areas from the image, which he referred to as place suggestions. As a result, rather than attempting to classify a vast number of locations, you may now just work with 2000 areas. The selective seek algorithm is used to create these 2000 location recommendations.

2.3.1 Data Loader

Cleaning and preparing data to transform it into the correct format is the first step in developing any device mastering model. On this project, the

VQA dataset is being used for education and assessment. The construction of the dataset is portrayed in the figure. The VQA v1.0 guidance dataset,

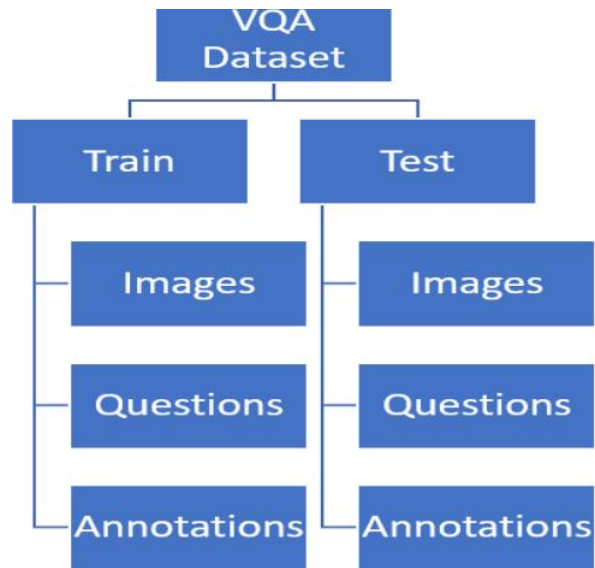


Figure 2.1: Structure of VQA Dataset

which is an Abstract picture dataset, has 20k pictures in the "Pictures" envelope. The "Questions" coordinator contains a JSON report with 60k questions matched to the photos. The VQA v2.0 guidance dataset, which is a Real pictures dataset, has 82k pictures in the "Photos" coordinator, with around 440k requests matched to those photos. The JSON report takes the going with structure:

```

{
  "info" :
  info, "task_type" :
  string, "data_type" :
  string, "data_subtype" :
  string, "questions" :
  [ques], "license" :
  license
}
question
{
  "question_id" : integer,
  "image_id" : integer,

```

```
"question" : str
```

```
}
```

The "Annotations" folder has a JSON file which provides answers to all the questions based on the im

```
{
```

```
"info" :
```

```
info, "data_type" :
```

```
str, "data_subtype" :
```

```
string,
```

```
"annotations" :
```

```
[annotation], "license" : license
```

```
}
```

```
annotation
```

```
{
```

```
"question_id" :
```

```
integer, "image_id" :
```

```
integer, "question_type" :
```

```
str, "answer_type" :
```

```
str, "answers" :
```

```
[answer], "multiple_choice_answers" :
```

```
str
```

```
}
```

```
answer
```

```
{
```

```
"answer_id" :
```

```
int, "answer" :
```

```
str, "answer_confidence" :
```

```
str
```

```
}
```

For this project, we only employ queries with more than one-desire response, and we compile the paintings as a classification issue. To do so, we begin by determining the best okay more than one want replies that solve the

most education difficulties. The final aim is to move every new enter query including an enter picture into these classes.

All other queries with answers that do not fit into these acceptable classes are eliminated and may not be used inside the curriculum. The selected questions are now mapped to their associated images with the appropriate response, resulting in a consolidated dataset. This merged dataset is also in JSON format, as illustrated in the example below.

, ...]

The inquiries and pictures from this assortment are currently being used to do work designing and have vector portrayals for each chance and question.

2.3.2 Feature Engineering

- Feature Engineering from Image

VGGNet, a CNN, is the version utilised to generate picture characteristic vectors. It contains sixteen to nineteen layers and uses small convolution filters (length of 3x3). The version mastered the use of the ImageNet dataset and performs admirably in photo categorization tasks.

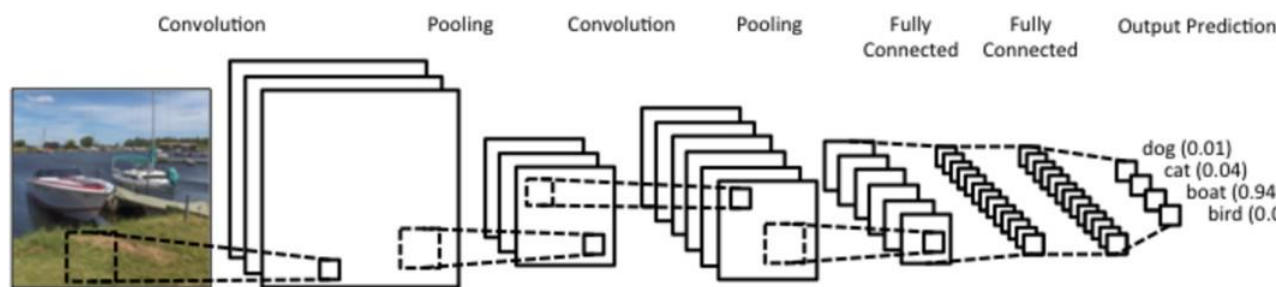


Figure 2.2: Architecture of Extracting Feature from Image

- Feature Engineering for Text

There are pre-trained algorithms for extracting characteristic vectors from raw textual content data. The bulk of these styles include embedding at the phrase level. The most often used maximum. Phrase vectors include GloVe (Global Vectors for Word Representation) and Word2Vec.

Understanding the interdependence between phrases and phrases, on the

other hand, is a difficult procedure.

The semantics of each statement must be comprehended. As a result, we want to build a device that captures those dependencies and semantics while also providing sentence-level embedding that we can use in our question-answering device.

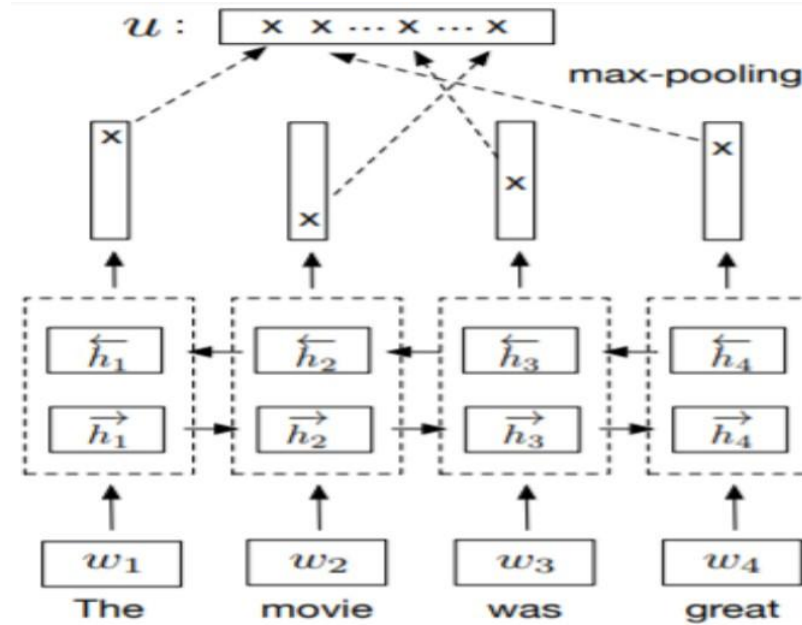


Figure 2.3: BiLSTM encoder with max-pooling

2.3.3 The weight sharing strategy

Because requests and responses have various accentuation features, our model includes multiple LSTMs for request and reaction. In any event, the importance of four single words in both requests and responses should be almost identical. As a result, the primary component's word-introducing layers and the third part share the weight cross section. Furthermore, the weight lattice in the fully linked Softmax layer is supplied this weight system for word-introducing layers in a delivered method. Naturally, the word-embedding layer's weight network's limit is to encode a one-hot word representation into a thick word portrayal. The weight system's restriction in the Softmax layer is to interpret the dense word portrayal into a faux single word depiction, which is word introducing's regressive movement. This approach reduces about half of the model's limitations and has been demonstrated to perform better in

picture captioning and novel visual concept learning tasks.

2.3.4 Training Details

The CNN we used is pre-arranged on the ImageNet classification task [13]. This part is fixed during the QA planning. We embrace a log-likelihood mishap described on the word game plan of the reaction. Restricting this incident work is equivalent to helping the probability of the model to make the groundtruth answers in the planning set. We commonly train the first, second and the fourth parts using stochastic slant decent method. The fundamental learning rate is 1 and we decline it by a part of 10 for each age of the data. We stop the arrangement when the incident on the endorsement set doesn't reduce inside three ages. The hyperparameters of the model are picked by cross-endorsement. For the request answering undertaking, we segment the sentences into a couple of word phrases. These articulations can be managed equivalently to the English words.

2.3.5 The FM-IQA Dataset (Freestyle Multilingual Image Question Answering)

1. The Data Collection

As the concealed picture set, we begin with the 158,392 images from the recently distributed MS COCO [13] planning, endorsement, and testing set. Baidu's online publicly supporting server⁴ is used to compile the comments. To diversify the named question-answer matches, annotators are allowed to create any type of request as long as it is related to the image's content. The visual content and practical aspects of the request should be addressed (for example, we do not anticipate to be asked, "What is the name of the person in the image?"). The annotators must respond to the genuine request.

According to one viewpoint, the open door we accommodate the annotators is useful to get a freestyle, intriguing and diversified set of requests. Of course, it makes it harder to control the nature of the clarification

diverged from a more point by point direction. A key quality filtering stage is used to screen the clarifying quality. We attempted 1,000 photos at random as a quality assessment dataset from the MS COCO dataset as a foundational set for the annotators (they were told this was a test). After each annotator has completed some naming on this quality noticing dataset (about 20 request response matches per annotator), we test a handful of explanations and score their quality. We merely choose a small number of annotators (195 people) with interesting comments (for instance the requests are associated with the substance of image and the reactions are correct). We also favour annotators who make delightful inquiries that need a high degree of recall in order to respond. Only the selected annotators get access to the remainder of the images. We choose a number of good and bad examples of made sense of inquiry response matches from the quality noticing dataset and present them as references to the selected annotators. We also discuss the motivations for selecting these models. After the massive number of photographs have been remarked, we further restrict the dataset by removing a small portion of the photos with clearly labelled questions and answers.

2. The data set's statistics

As of now there are 158,392 pictures with 316,193 request answer matches and their English translations. Each image has somewhere near two request answer matches as remarks. The normal lengths of the requests and answers are 7.38 and 3.82 independently assessed by Chinese words. Some test pictures are shown in Figure 3. We with no obvious end goal in mind analyzed 1,000 request answer matches and their relating pictures as the test set.

The requests in this dataset are separated, which requires a gigantic course of action of AI limits generally together to answer them. They contain some modestly clear picture understanding requests of, e.g., the exercises of things (e.g., "What is the youngster in green cap doing?"), the article class (e.g., "Is there any individual in the image?"), the general positions and correspondences among objects (e.g., "Is the PC

on the right or left 50% of the gentleman?”), and the qualities of the articles (e.g., “What is the assortment of the frisbee?”). In addition, the dataset contains a couple of requests that need a critical level reasoning with signs from vision, language and practical. For example, to answer the subject of “Why does the vehicle leave there?”, we should understand that this question is about the left vehicle in the image with two men holding contraptions at the back.

Based on our reasoning, we may deduce that the car has a few faults, which the two guys in the photograph are attempting to resolve. These requests are difficult to respond to, but we admit that they are the most intriguing part of the dataset requests. We divide the findings into eight categories and provide the results on the project website. The responses have also been expanded. The annotators can either offer a single articulation or a single word as the reaction (for example, “Yellow”), or they can supply a complete phrase (for instance “The lemon is yellow”).

CHAPTER 3

ANALYSIS

3.1 Introduction

To exhibit our technique, we made an enormous scope Freestyle Multilingual Image Question Answering dataset dependent totally upon the MS COCO dataset. The dataset's contemporary model contains 158,392 pictures, 316,193 question arrangement pairings, and their English interpretations. To variety the comments, annotators are allowed to raise any request connected with the picture's substance. We give ways to deal with screening the top notch comments. This dataset contains a wide range of AI-related questions, such as movement recognition (e.g., "Is the individual attempting to purchase vegetables?"), thing recognition (e.g., "What is the thing in red?"), positions and colaborations between many objects inside the picture (e.g., "Where is the cat?"), and fundamentally based thinking (e.g., "Where is the cat?").(e.g. "For what reason does the transport park here?"). Due to the assortment of free-form question-answer pairings, it's far hard to assess the methodology with robotized metrics.

Human adjudicators are utilized in a Visual Turing Test. In particular, we consolidate the question arrangement matches produced by our variant with the indistinguishable arrangement of inquiry arrangement matches characterized by annotators. The human appointed authorities need to evaluate on the off chance that the arrangement is given by a form or by an individual. Moreover, we request that they give a rating of 0 (wrong), 1 (to some degree exact), or 2 (excellent) (for example correct).

The outcomes show that our mQA rendition finishes 64.7 percent of this assessment (handled as human responses), with a run of the mill grade of 1.454. We look at our model's disappointments and show that, when combined with the m-RNN model, our model can ask and answer questions regarding

photographs.

3.2 Software Requirement Specification

A software requirements specification (SRS) is a report that explains what the software will accomplish and how it will be expected to perform. It also indicates the product's capacity to fulfil the needs of all stakeholders (business and users).

3.2.1 User requirement specification

- Freestyle Multilingual Image Question Answering dataset
- RAM: 512MB or More
- Disk space: 128GB or more

3.2.2 Software requirement

Operating system	Windows 10
Coding Language	python
Tool	PyCharm
Database	MYSQL
Server	Flask
Libraries	PyTesseract, Mathplotlib, openCV, TensorFlowv2.5

3.2.3 Hardware requirement

system	Pentium Dual Core.
Hard Disk	120 GB
Monitor	15" LED
Input Devices	Keyboard, Mouse
Ram	1 GB

3.3 Algorithm and Flow Chart

1. **Recurrent Neural Network:** Recurrent Neural Networks (RNN) are those that do a comparative endeavor for each detail of the input course of action. The aftereffect of a RNN is dependent upon both the continuous data and the previous calculations. RNNs feature a "memory" that stores all of the records that were previously shown. Figure shows an essential RNN spreading out throughout the span of time.

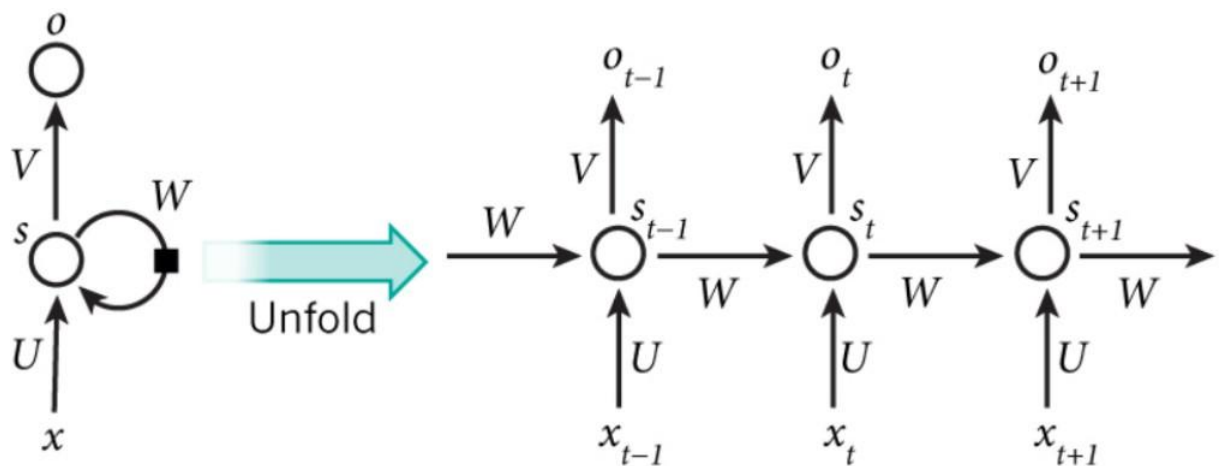


Figure 3.1: Example of a simple RNN unfolding in time

2. **Long short-term memory:**

Long momentary memory (LSTM) is a basic repetitive brain local area that might be utilized as a development issue or square (of stowed away layers) for a greater intermittent brain local area. The LSTM block is an intermittent local area all by itself, as it contains repetitive associations like those found in a customary intermittent brain community.

A LSTM block is comprised of four significant parts: a mobileular, a passage door, a result entryway, and a disregard entryway. The mobileular is accountable for "recollecting" values throughout erratic time spans, consequently the expression "memory" in LSTM. Every one of the three doors might be considered a "regular" engineered neuron, as in a

multi-facet (or feedforward) brain organization: they ascertain an output.

3. **Convolutional Neural Networks**

Convolutional Neural Networks (CNNs) are basically many layers of convolutions layered on top of one another using non-direct authorization limits, for instance, tanh or ReLU. Each center point in the continuous layer is related with every center point in the accompanying layer in an ordinary totally related cerebrum association. Convolutional channels, on the other hand, are employed in CNNs to slide over the centers in the data layer and in like manner register the yield. The possibility of sliding convolutions across center points in the data layer, as conflicted with to major multi-layer perceptrons, achieves sections of the information layer being associated with each node in the outcome layer. To calculate the outcome, each layer of the association uses a separate set of convolutional filters. When the association sees the arrangement data, it automatically acknowledges which channels are to be used.

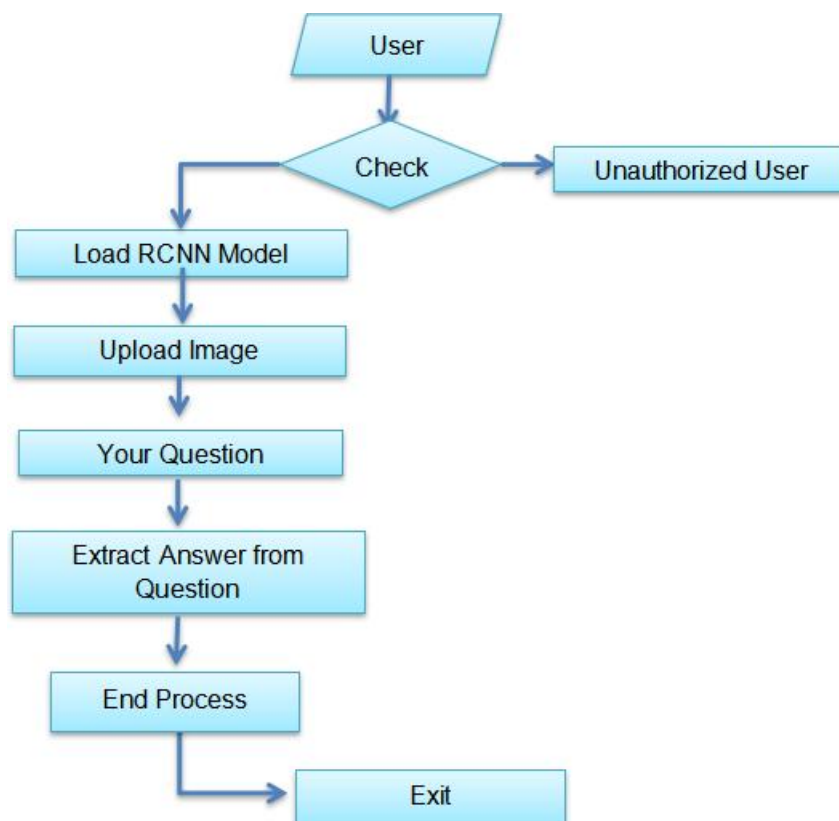


Figure 3.2: Flow Graph

CHAPTER 4

DESIGN

4.1 Introduction

This section of document defines the overall stucture and configuration of the project with a well defined use-case diagram.

4.2 Diagrams

4.2.1 ER Diagram

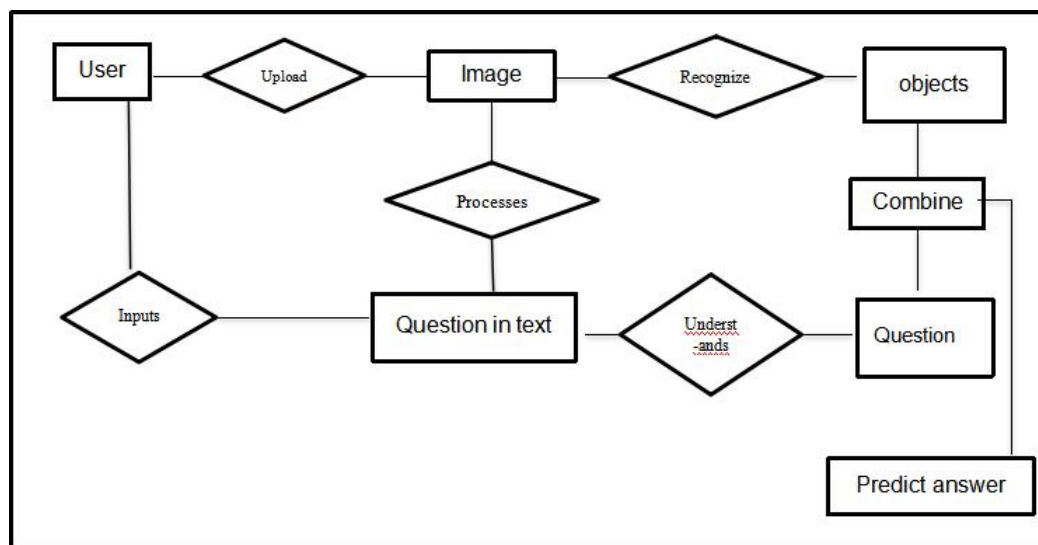


Figure 4.1: ER Diagram for VQA Model

The above ER Diagram shown is the complete process of the design wherein The user uploads images and inputs question in the form the text, then the system understands the question using LSTM and recognizes objects using CNN. By combining the results of both process model predicts the answer.

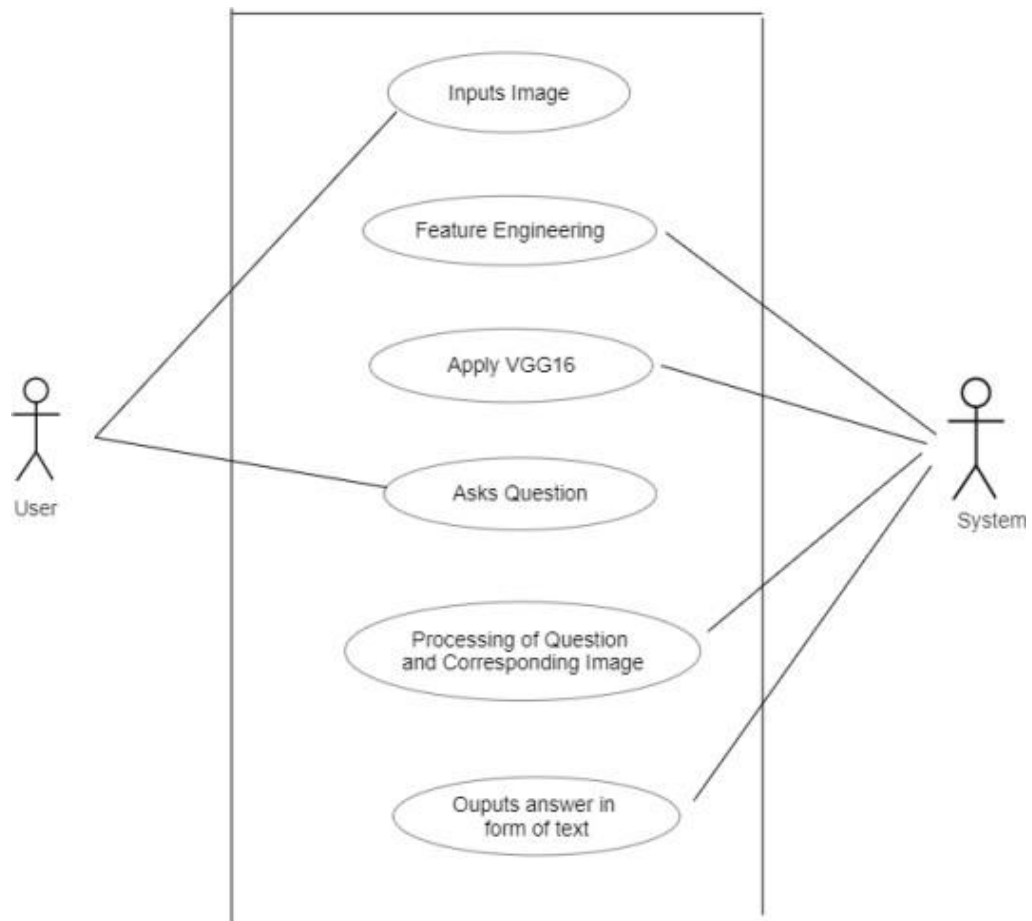


Figure 4.2: Use Case diagram of VQA

4.2.2 Use Case

The above Use Case diagram has two actors system and the user, here it shows the user inputs the image and related question. The system does the following functions feature engineering, apply VGG16 from CNN, Processing of question and understands the corresponding image and outputs answer in the for of text.

4.3 Module Design and Organization

Our mQA model's architecture, The model is made up of four parts:

for removing an inquiry's semantic portrayal, CNN for removing the picture portrayal,a LSTM for separating the portrayal of the ongoing word in the response and its etymological setting, and a melding part

that joins the data from the initial three sections to create the following word in the response. These four parts can be prepared couple.

The mQA Model's Four Components

- 1•. The primary element of the model is used to eliminate the inquiry's semantic relevance. It consists of a 512-layered express installing layer and a 400-memory-cell LSTM layer. The expresion implanting layer has had the ability to design the one-warm vector of the expression directly it into thick semantic space. The LSTM layer takes care of this puzzling sentence representation.
2. The next viewpoint is a deep Convolutional Neural Network (CNN) that creates the contour of a picture. In this study, we employ GoogleNet. It should be noted that other CNN models, like as AlexNet and VggNet, can also be used as perspectives in our approach. We replace the final SoftMax layer of the deep CNN with the last zenith layer.
3. The third portion also includes a term that introduces a layer and an LSTM. The development appears to be the most important aspect. To create the going with words in the answer, the word embeddings and inception of the memory cells for the words with in response will be entered into the merging phase. Finally, the fourth section brings together the data from the previous three levels. For the t th term in the reaction, the incitation of a merging layer $f(t)$ is not totally established as follows:

$$f(t) = g(Vr_Q r_Q + V_I I + V_rA r_A(t) + V_{ww}(t));$$

where "+" implies part wise extension, r_Q addresses the inception of the LSTM(Q) memory cells of the last say in regards to the request, I shows the image depiction, $r_A(t)$ and $w(t)$ implies the activation of the LSTM(A) memory cells and the word embedding of the t th word in the reaction independently. Vr_Q , V_I , V_rA , and V_w are the weight matrices that ought to be learned. $g(.)$ is a part wise non-straight limit.

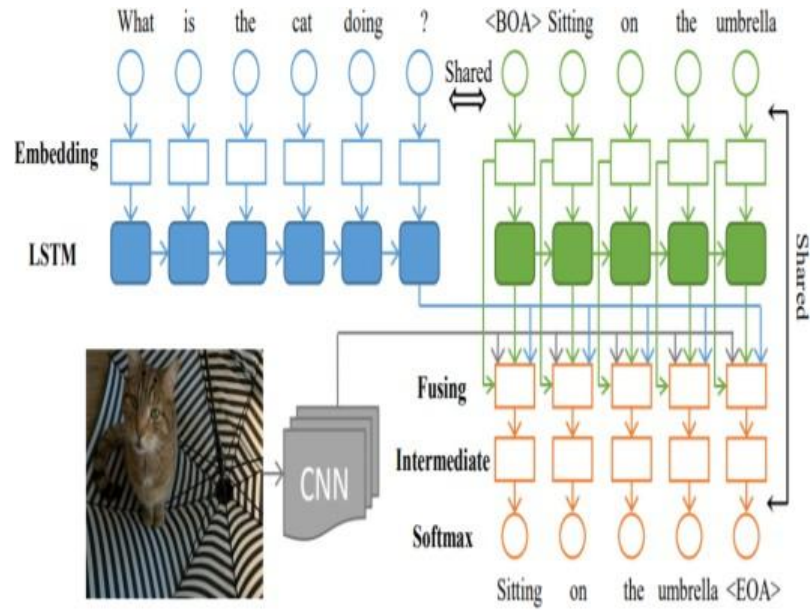


Figure 4.3: Model Overview

4.4 System Architecture

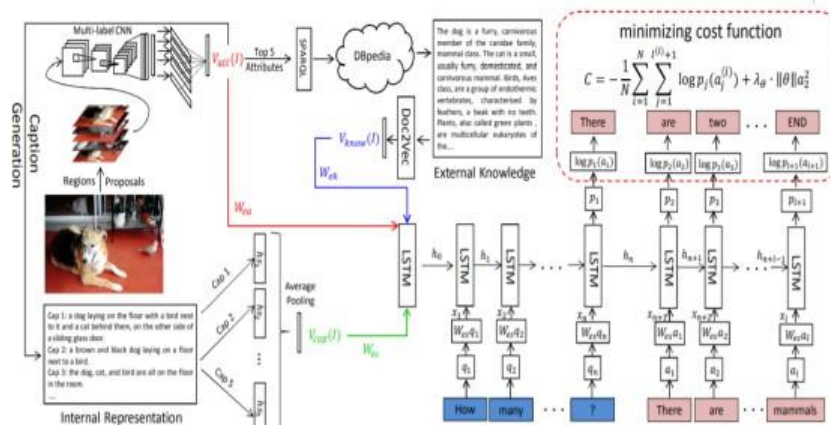


Figure 4.4: System Overview

CHAPTER 5

IMPLEMENTATION

5.1 Introduction

LSTM and CNN are used in the model. They connect the enquiry and game plan to the LSTM before dealing with them. Instead, we use various LSTMs for queries and plans relating to the extraordinary qualities of requests and action plans, independent of how we permit the exchange of word-embeddings. Use the proposed dataset instead of our FM-IQA dataset for the dataset. It might also associate the strategy with a single word.

5.2 Explanation of Key Function

- **Extracting Image features with VGG16**

The VGG16 version from tensorflow keras is initially imported. The image module is imported to preprocess the picture item, and the preprocess input module is loaded to scale pixel values as the VGG16 version requires. For array processing, the numpy module is loaded. The VGG16 version is then loaded with the imagenet dataset's pretrained weights. The VGG16 version is a convolutional layer sequence viewed via one or more dense (or totally linked) layers. From the entrance layer to the remaining max pooling layer (classified via 7 x 7 x 512) is shown as the version's characteristic extraction portion, whilst the rest of the community is shown as the version's type part.

- **Training Data-points**

We utilized a CNN that has recently been prepared on the ImageNet classification work. During the QA preparing, this part is revised. We

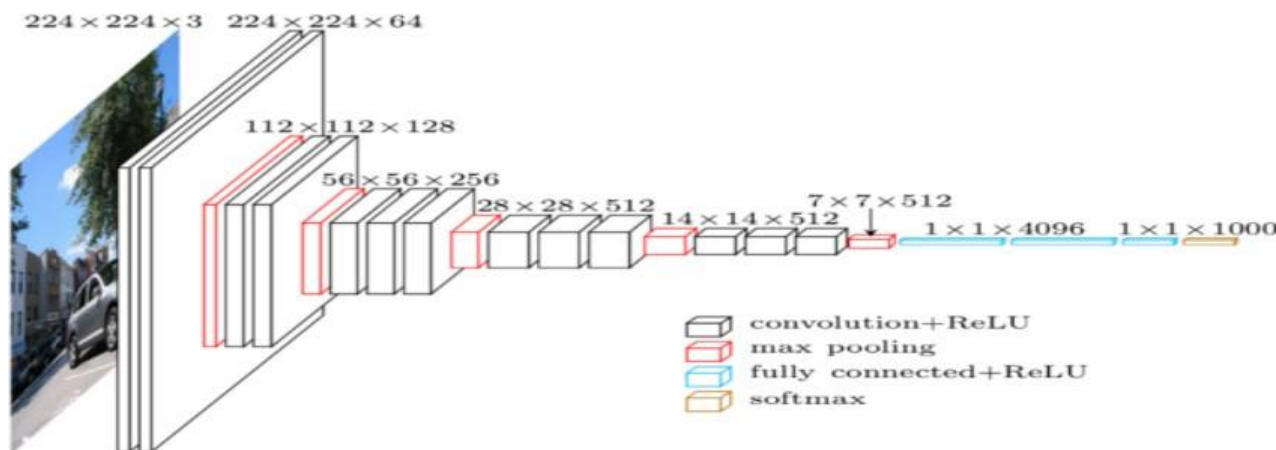


Figure 5.1: Feature Extraction using VGG16

utilize a log-likelihood misfortune in view of the response's statement arrangement. Limiting this misfortune work is comparable to improving the model's probability of creating ground-truth reactions in the preparing set. Utilizing the stochastic slope nice methodology, we together train the LSTM and CNN.

- **Visual Turing Test**

In the Visual Turing Test, a human adjudicator would be provided an image, request, and answer to the question posed by testing model or human annotators. After considering the response, the individual should determine if it was supplied by a person (for example, breeze through the test) or a computer.

- **Score of the Generated Answer**

A human-appointed authority will be presented an image, a query, and the response to inquiry made either by testing model or human annotators in the Visual Turing Test. In light of the reaction, the person in question must determine if the response was supplied by an individual (for example, completing the test) or a machine. **Score of the Generated Answer**

5.3 Technology

- **Neural Network in Machine Learning**

A neural network is a progression of calculations that endeavor to appreciate basic relationships in a bunch of information utilizing a methodology that mirrors how the human psyche functions. Thusly, brain networks speak with regular and engineered neuron structures. Since brain organizations can adjust to changing over input, the local area may achieve the most ideal end-product without changing the result models.

- **Freestyle Multilingual Image Question Answering dataset**

In the Visual Turing Test, a human designated power will be shown a picture, a request, and the answer to the request made by the testing model or human annotators. The person being referred to must close whether the reaction was given by a person (for instance finish the test) or a computer taking into account the response.

5.4 Method of Implementation

Proposed Method split the data into train set and test set with 80:20 proportion. Further, we divided the train data into train and validation dataset. Overview of tasks performed:

1. Split the data to Train, Test and validation
2. Method:
 - (a) Dataset cleaning and Pre-processing
 - (b) Transforming Questions and Answers into vectors
 - (c) Train dataset using CNN and RNN
 - (d) Extracting Features from an image using VGG16

(e) Evaluate on Validation data and tune the models. Repeat to find best parameters.

(f) Test performance of final model on Test dataset

5.4.1 Output Screens

The dataset you provided us can answer limited questions such as 'does sphere and cube colour match' or 'sphere is present in image' with a 'YES' or 'NO' or number of objects in image, but our project can describe anything in the image by asking a user question and answering it with RCNN to identify objects in images and LSTM to build vocabulary sentences in meaningful form.

To run project double click on 'run.bat' file to get below screen

In above screen click on 'load RCNN Model' button to stack CNN model

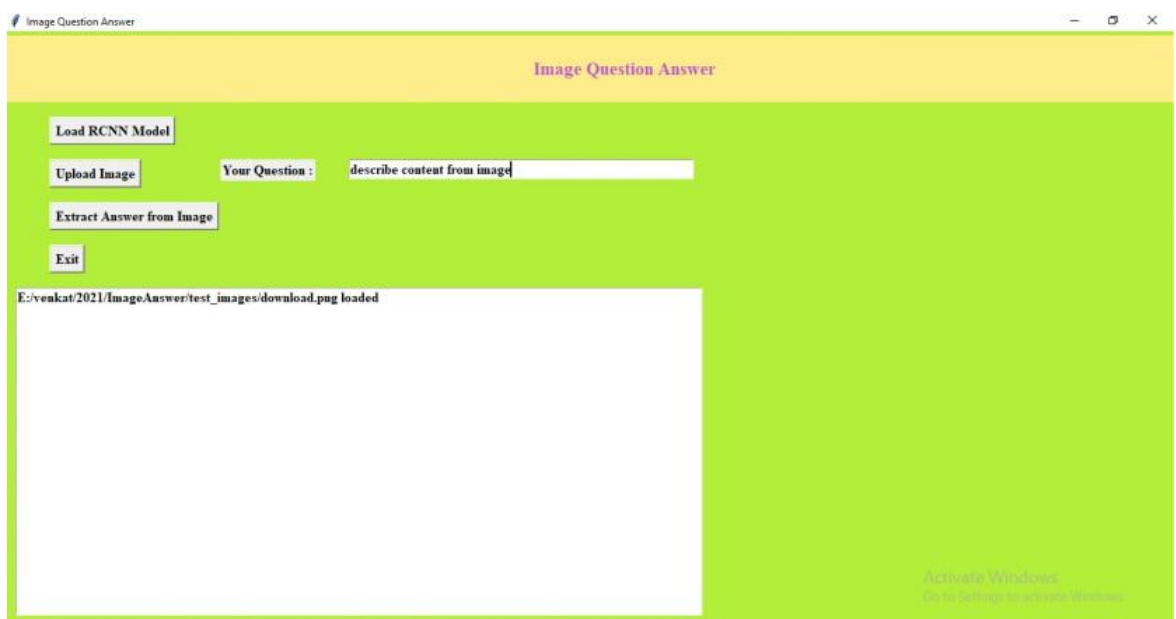
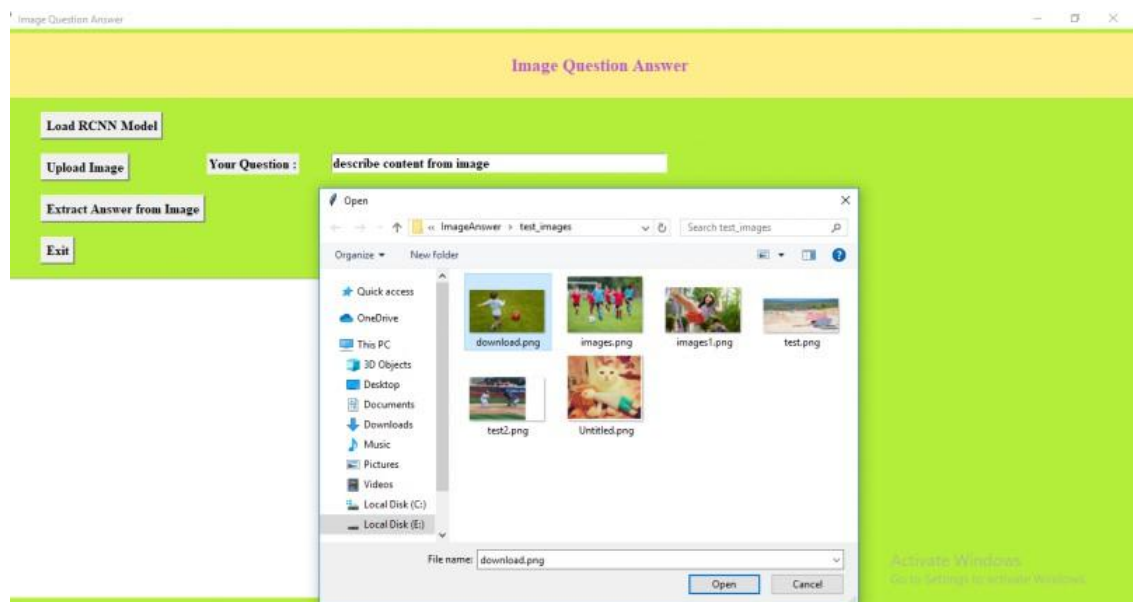


and to get underneath screen

In above screen in text region we can see model stacked and presently click on 'Upload Image' button to transfer

In above screen in text region we can see model stacked and presently click on 'Upload Image' button to transfer

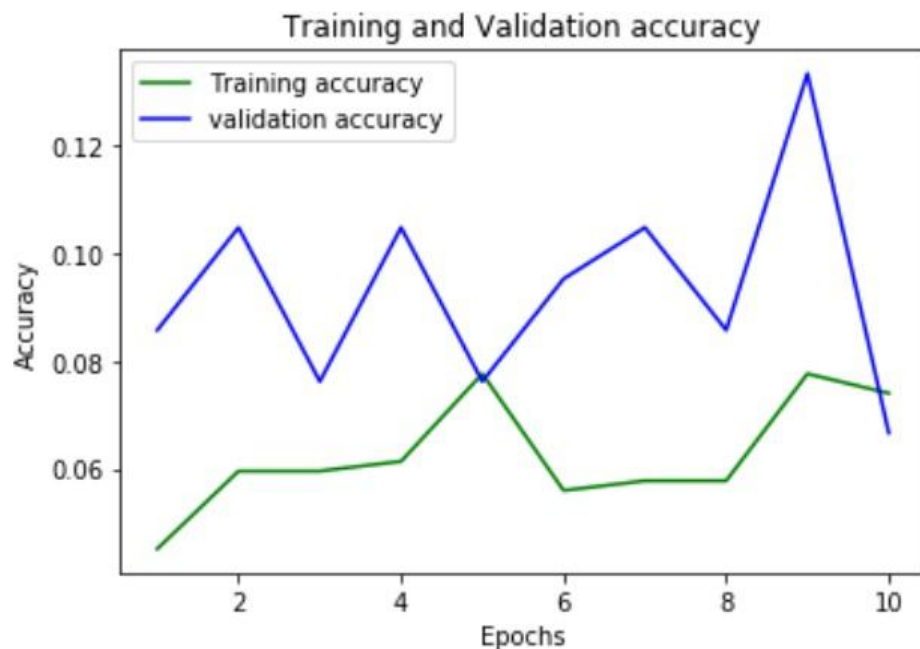
In above display photograph loaded and I entered query as 'describe content material from photograph' and now click on 'Extract Answer from Image' button to get underneath output In above screen in yellow colour text we got



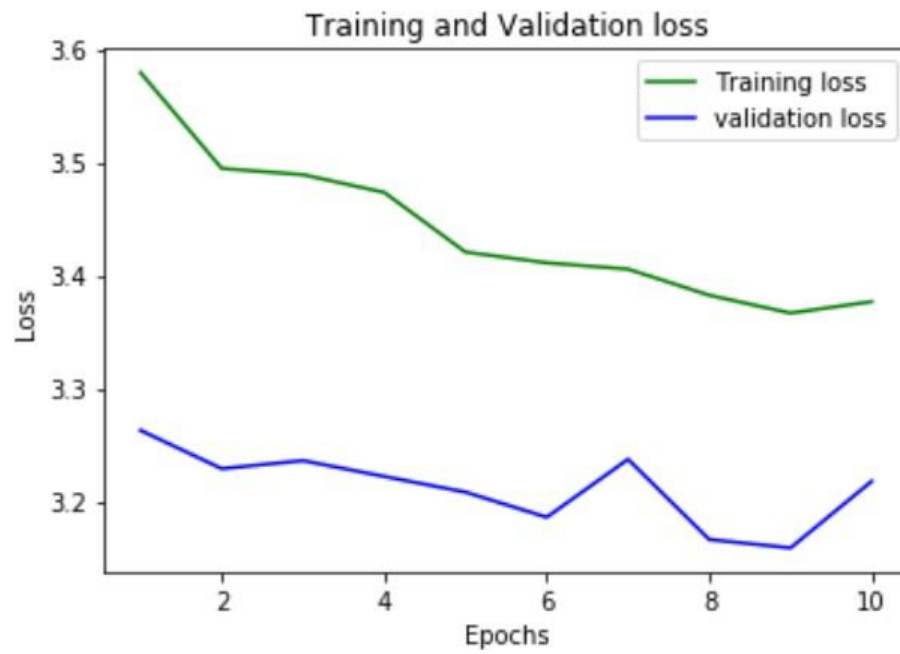


answer about the image

5.4.2 Result Analysis



Subsequent to tuning hyper-barriers to streamline the calculations, Stochastic Gradient Descent changed into considered because the maximum perfect calculation, taking each execution and time intricacy into account. Following execution measurements had been achieve:



- Accuracy: 92.81 %
- Precision: 96.97 %
- Recall: 91.94 %
- F1-Score: 94.39 %

CHAPTER 6

TESTING RESULTS

6.1 The Visual Turing Test

A Turing Test selected will be outfitted with an image, a request, and the reaction to the request made by the testing model or by human annotators in this Visual Turing Test. Considering the response, the individual ought to finish up whether the reaction was given by a person (for instance breeze through the test) or a computer.

It displays that 64.7 percent of the responses given by our mQA model are considered as human-gave answers. This occupation is performed insufficiently by the outwardly weakened QA. Regardless, a part of the made responses breeze through the evaluation. Since a piece of the requests are multi-choice, a genuine reaction can be gotten with an erratic gauge subordinate solely upon phonetic cues.

To take a gander at the assortment of the VTT appraisal among specific social occasions of human adjudicators, we coordinated further gatherings with various get-togethers of judges in the unclear setting. The stylish deviations of the passing expense for human, blind-mQA, and mQA interpretations are 0.013, 0.019, and 0.024, separately. It shows that VTT is a trustworthy and solid assessment metric for this assignment.

6.2 The Score of the Generated Answer

The Visual Turing Test is most successful in providing a tough evaluation of the generated answers. We also do a tiny or small assessment using ratings of "0," "1," or "2." The values "0" and "2" indicate that the answer is completely erroneous and flawlessly accurate, respectively. "1" method that

the answer is partially accurate (e.g., the overall classes are correct but the sub-classes are erroneous) and makes sense to human judges. Human judges for this mission are not necessarily the same as those for the Visual Turing Test. In the wake of examining the outcomes, we found that a couple of human appointed authorities likewise give a "1" to a reaction assuming the request is so difficult to answer that even an individual, without cautiously reviewing the picture, will perhaps make mistakes.

6.3 Performance Comparisons of the Different mQA Variants

We deploy three mQA versions to test the efficiency of the various components and techniques of our mQA model. For the first option (i.e. "mQA-avg-question"), we replace the model's first LSTM component.

	Word Error	Loss
mQA-avg-question	0.442	2.17
mQA-same-LSTMs	0.439	2.09
mQA-noTWS	0.438	2.14
mQA-complete	0.393	1.91

It is utilized to exhibit the LSTM's presentation as an inquiry inserting student and extractor. To demonstrate Q An in the subsequent variety (for example "mQAsame-LSTMs"), we use two shared-loads LSTMs. It is used to exhibit the handiness of our model's decoupling technique for the loads of the LSTM(Q) and LSTM(A). We don't utilize the Transposed Weight Sharing (TWS) procedure for the third variety (for example "mQA-noTWS"). Exhibiting TWS's efficacy is utilized.

CHAPTER 7

Conclusions and Future Scope

VQA is to be sure a clever subject that needs a careful investigation of each visual components. Given the current degree of AI (ml), it is sensible to accept that VQA designs will essentially work on with regards to improvement issue and rightness after some time. Profound learning has proactively shown an exceptional expansion in the general viability of many processing inventive and prescient models, as well as phonetic models. Building a machine that joins those particular parts can emphatically upgrade noticed results for exercises like VQA.

The most frequently involved measures for estimating the general presentation of a VQA framework are rightness and F1-measure. Its on the grounds that all current VQA structures think about it as a sort issue. Scientists anticipate that that maybe the response should any natural language question presented about a picture will unquestionably relate to the one of the "arrangement classes" that now the calculation is aware of. One more development of this might be to introduce a contraption that utilizes a probabilistic model to construct a thorough reaction simply alludes to a picture. Such demeanor drives the technique away from its unique sort, provoking the inclination to integrate new rules for by and large execution assessment. Nonetheless, apparently further datasets and measurements will be added rapidly, and However, it looks that extra datasets and measurements may be presented quick, and this seems, by all accounts, to be very probable.

Considering every ongoing framework (counting this task), there are as yet proceeding with contentions about different worries concerning the methodologies used to prepare and evaluate the models. One of the most well known and open-finished 39debates is whether these VQA frameworks, which are prepared and evaluated on inquiries with numerous decision reactions in datasets, can be designated "astounding entertainers" in actuality, conditions. The inquiries

expected to prepare the calculation in many datasets are gained through local area obtaining. Are these inquiries adequate to mirror all possible inquiries that the framework may encounter? The parts of inquiries additionally fundamentally affect work designing and, accordingly, the framework's in general prescient ability. Regardless of this, the contemporary super present day structures on VQA are a long way from human execution on the equivalent datasets. Nonetheless, with the accessibility of transferred research and arising hardware and innovation, there is a great deal of opportunity to get better inside the results.

REFERENCE

- [1].L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. “Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs”. ICLR, 2015.
- [2].M. Grubinger, P. Clough, H. Muller, and T. Deselaers. “The iapr tc-12 benchmark: A new evaluation resource for visual information Systems”. In International Workshop OntoImage, pages 13–23, 2006.
- [3].A. Lavie and A. Agarwal. Meteor: “An automatic metric for mt evaluation with high levels of correlation with human judgements”. In Workshop on Statistical Machine Translation, pages 228–231. Association for Computational Linguistics, 2007.
- [4].J. Mao, W. Xu, Y. Yang, J. Wang, Z. Huang, and A. Yuille ”Deep captioning with multimodal recurrent neural networks (m-rnn)”. In ICLR, 2015.
- [5].J. Zhu, J. Mao, and A. L. Yuille. “Learning from weakly supervised data by the expectation loss svm (e-svm) algorithm”. In NIPS, pages 1125–1133, 2014.
- [6].A. Karpathy, A. Joulin, and F. F. Li, “Deep fragment embeddings for bidirectional image sentence mapping,” in Proc. Advances in Neural Inf. Process. Syst., 2014.
- [7].J. Mao, W. Xu, Y. Yang, J. Wang, and A. Yuille, “Deep Captioning with Multimodal Recurrent Neural Networks (m-RNN),” in Proc. Int. Conf. Learn. Representations, 2015.
- [8].O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and tell: A neural image caption generator,” in Proc. IEEE Conf. Comp. Vis. Patt. Recogn., 2014.
- [9].L. Yao, A. Torabi, K. Cho, N. Ballas, C. Pal, H. Larochelle, and A. Courville, “Describing videos by exploiting temporal structure,” in Proc. IEEE Int. Conf. Comp. Vis., 2015.

- [10].J. Devlin, H. Cheng, H. Fang, S. Gupta, L. Deng, X. He, G. Zweig, and M. Mitchell, “Language models for image captioning: The quirks and what works,” arXiv preprint arXiv:1505.01809, 2015.
- [11].H. Xu and K. Saenko, “Ask, attend and answer: Exploring question-guided spatial attention for visual question answering,” in European Conference on Computer Vision. Springer, 2016, pp. 451–466.
- [12].Z. Yang, X. He, J. Gao, L. Deng, and A. Smola, “Stacked attention networks for image question answering,” in Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 2016, pp. 21–29.
- [13].J. Lu, J. Yang, D. Batra, and D. Parikh, “Hierarchical question-image coattention for visual question answering,” in Advances In Neural Information Processing Systems, 2016, pp. 289–297.
- [14].J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A largescale hierarchical image database,” in Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. IEEE, 2009, pp. 248–255.