

An
Industry Oriented Mini Project Report On

VISUAL DESCRIPTION GENERATOR

Submitted in Partial Fulfillment of the Requirements for the Award of Degree

BACHELOR OF TECHNOLOGY

IN

**COMPUTER SCIENCE AND ENGINEERING
(DATA SCIENCE)**

Submitted By

SANDA RAKSHITHA	217Z1A6752
ARKALA VISHAL	227Z5A6701
DUDEKALA JALEEL	217Z1A6718

Under the Guidance of
Mr. S. Srikanth Reddy
Assistant Professor



SCHOOL OF ENGINEERING
Department of Computer Science and Engineering

**NALLA NARASIMHA REDDY
EDUCATION SOCIETY'S GROUP OF INSTITUTION
(AN AUTONOMOUS INSTITUTION)**

Approved by AICTE, New Delhi, Chowdariguda (V) Korremula 'x' Roads,
Via Narapally, Ghatkesar (Mandal) Medchal (Dist), Telangana-500088
2024 - 2025



NALLA NARASIMHA REDDY
Education Society's Group of Institutions - Integrated Campus
(UGC AUTONOMOUS INSTITUTION)



SCHOOL OF ENGINEERING
DEPARTMENT OF COMPUTER SCIENCE AND ENGINEERING

CERTIFICATE

This is to certify that the project report titled “**Visual Description Generator**” is being submitted by **S. Rakshitha(2171A6752)**, **A. Vishal(227Z5A6701)**, and **D. Jaleel (217Z1A6718)** in Partial fulfilment for the award of **Bachelor of technology in Computer Science & Engineering (Data Science)** is a record bonafide work carried out by them. The results embodied in this report have not been submitted to any other University for the award of any degree.

Internal Guide

(Mr. S. Srikanth Reddy)

Head of the Department

(Dr. K. Rameshwaraiah)

Submitted for Viva voce Examination held on

External Examiner

DECLARATION

We S. Rakshitha, A. Vishal, and D. Jaleel are students of **Bachelor of Technology in Computer Science and Engineering (Data Science), Nalla Narasimha Reddy Education Society's Group Of Institutions**, Hyderabad, Telangana, hereby declare that the work presented in this project work entitled “**Visual Description Generator**” is the outcome of our own bonafide work and is correct to the best of our knowledge and this work has been undertaken taking care of engineering ethics. It contains no material previously published or written by another person nor material which has been accepted for the award of any other degree or diploma of the university or other institute of higher learning.

By

S. Rakshitha 217Z1A6752

A. Vishal 227Z1A6701

D. Jaleel 217Z1A6718

Date:

Signature:

ACKNOWLEDGEMENT

We express our sincere gratitude to our guide **Mr. S. Srikanth Reddy**, Assistant Professor in Computer Science and Engineering Department, NNRESGI, who motivated throughout the period of the project and also for his valuable and intellectual suggestions apart from his adequate guidance, constant encouragement right throughout our work.

We would like to express our profound gratitude to our mini project In-charge **Mrs. B Sriveni**, Assistant Professor in Computer Science and Engineering Department, NNRESGI, for her support and guidance in completing our project and for giving us this opportunity to present the project work.

We profoundly express thanks to **Dr. K. Rameshwaraiah**, Head of Computer Science and Engineering Department, NNRESGI, for his cooperation and encouragement in completing the project successfully.

We wish to express our sincere thanks to **Dr. G. Janardhana Raju**, Dean School of Engineering, NNRESGI, for providing the facilities for completion of the project.

We wish to express our sincere thanks to **Dr. C. V. Krishna Reddy**, Director NNRESGI, for providing the facilities for completion of the project.

Finally, we would like to thank overall mini-project coordinator, **Members of Project Review Committee (PRC)**, all the faculty members and supporting staff of the Department of Computer Science and Engineering, NNRESGI, for extending their help in all circumstances.

By

S. Rakshitha **217Z1A6752**

A. Vishal **227Z5A6701**

D. Jaleel **217Z1A6718**

ABSTRACT

Visual description generation is the process of creating detailed and contextually accurate natural language descriptions from visual content, such as images. Traditional models have struggled with capturing complex visual details and maintaining consistency in generated descriptions. In this project, we introduce a visual description generator that leverages Vision Transformers (ViT) for image feature extraction, addressing these limitations and pushing the boundaries of descriptive accuracy and coherence. Vision Transformers (ViT) offer a powerful alternative to traditional convolutional neural networks by treating image patches as sequences and applying transformer architecture for feature extraction. This approach enables the model to capture global and local relationships within the visual content, leading to richer and more detailed feature representations. These features are then used in conjunction with a language model to generate descriptive text that is both contextually relevant and coherent. The model is trained on diverse benchmark datasets, ensuring its ability to generate high-quality descriptions across various image domains. The evaluation of generated descriptions is conducted using metrics such as BLEU which measure grammatical correctness, semantic relevance, and overall quality. The resulting visual description generator demonstrates significant improvements in both descriptive depth and contextual understanding, marking a step forward in the field of automated visual content interpretation.

Keywords: Visual description generation, Vision Transformers (ViT), Deep learning techniques, Image feature extraction, Natural language generation, BLEU.

TABLE OF CONTENTS

CONTEXT	Page No
LIST OF FIGURES	i
LIST OF TABLES	ii
LIST OF ABBREVIATIONS	iii
1. INTRODUCTION	1
1.1 MOTIVATION	1
1.2 PROBLEM STATEMENT	1
1.3 PURPOSE	1
1.4 SCOPE	2
1.5 PROJECT OBJECTIVE	2
1.6 LIMITATIONS	2
2. LITERATURE SURVEY	3
2.1 INTRODUCTION	3
2.2 EXISTING SYSTEM	3
2.3 DISADVANTAGES	4
2.4 PROPOSED SYSTEM	4
2.4 ADVANTAGES	5
3. SYSTEM ANALYSIS	6
3.1 FUNCTIONAL REQUIREMENTS	6
3.2 NON-FUNCTIONAL REQUIREMENTS	6
3.3 INTERFACE REQUIREMENTS	7

4. SYSTEM DESIGN	9
4.1 UML DIAGRAMS	9
4.2 MODULES	15
5. IMPLEMENTATION AND RESULT	16
5.1 METHOD OF IMPLEMENTATION	16
5.2 MODULES	17
5.3 EXTENSION OF KEY FUNCTION	20
5.4 OUTPUT SCREENS	24
6. SYSTEM TEST	26
6.1 TYPES OF TESTS	26
6.2 VARIOUS TESTCASE SCENARIOS	28
7. CONCLUSION AND FUTURE ENHANCEMENT	30
7.1 PROJECT CONCLUSION	30
7.2 FUTURE ENHANCEMENT	30
8. REFERENCES	32
8.1 PAPER REFERENCES	32
8.2 WEBSITES	33

LIST OF FIGURES

Figure No	Figure Name	Page No
4.1.1	DFD Level 0	10
4.1.1.1	DFD Level 1	10
4.1.2	Use Case Diagram	11
4.1.3	Class Diagram	12
4.1.4	Sequence Diagram	13
4.1.5	Activity Diagram	14
5.3.1	Setup	20
5.3.1.1	Importing Libraries	20
5.3.2	Model Loading	21
5.3.2.1	Set the Device to GPU	21
5.3.3	Image Captioning Parameters	21
5.3.4	Predict and Show Function	22
5.3.4.1	Generate the Captions	22
5.3.5	Example Usage	23
5.4.1	Model training and building	24
5.4.2	Output when a file is uploaded	24
5.4.3	Output when a folder is uploaded	25

LIST OF TABLES

Table No	Name Of the Table	Page No.
6.2	Test Cases	28 - 29

LIST OF ABBREVIATIONS

S. No	Abbreviation	Definition
1	AI	Artificial Intelligence
2	BLEU	Bilingual Evaluation understudy
3	CNN	Convolutional Neural Network
4	DFD	Data Flow Diagram
5	GPT	Generative Pre-trained Transformer
6	GPU	Graphics Processing Unit
7	GUI	Graphical User Interface
8	IDE	Integrated Development Environment
9	LSTM	Long Short-Term Memory
10	ML	Machine Learning
11	PIL	Python Imaging Library
12	UML	Unified Modeling Language
13	ViT	Vision Transformers

CHAPTER 1

INTRODUCTION

1.1 MOTIVATION

In recent years, as the digital world continues to increase, the need for accessing the content has become important. For individuals with visual impairments, understanding the visual content on websites and social media can be a significant challenge. The main goal of this project is to develop a Visual Description Generator that can automatically generate captions for the images. This project aims to bridge the gap between visual and textual information, making digital content accessible to all the users.

1.2 PROBLEM STATEMENT

Currently, individuals with visual impairments face significant challenges in accessing and understanding visual content on websites, social media, and other digital platforms. The process of manually creating image descriptions is inconsistent and often limited in scope. The lack of automatic and reliable tools for generating these descriptions has created a gap in accessibility. This project resolves this issue by developing a system that can accurately and contextually describe images, providing a more comprehensive solution for users who depend on text-based descriptions.

1.3 PURPOSE

The goal of this project is to create a tool that uses the power of artificial intelligence to make a description of the image provided. Using advanced algorithms in computer vision and natural language processing, the system will understand images and generate an accurate description of the provided image. This will allow the users with visual impairments to gain insights from visual content, make accurate decisions, and participate in the digital world seamlessly. This program also aims to demonstrate how AI can be used for good deeds and increase availability for a wider audience.

1.4 SCOPE

With the increase in the advancement of technology, there is a significant scope to improve accessibility in the digital content through automated systems. The scope of this project is to develop a Visual Description Generator that can be implanted in real-world applications, such as websites and social media platforms. The project aims to make visual content more understandable to the users with visual impairments, ultimately playing an important role in contributing to a more inclusive digital framework.

1.5 PROJECT OBJECTIVE

The main objective of this project is to create a system that uses AI technologies to generate accurate and contextually meaningful descriptions for images. This system will be designed to help users with visual impairments understand visual content more effectively, thus improving their overall digital experience. The project aims to demonstrate the power of smart technology to solve accessibility challenges and promote thoroughness in the digital environment.

1.6 LIMITATIONS

The development of a Visual Description Generator requires access to large datasets of images and their respective descriptions for training the AI models. Moreover, there may be varying challenges related to the accuracy of descriptions in highly complex or abstract images. Ensuring the system's ability to handle diverse and distinct visual content will be an ongoing challenge that requires further research and development.

CHAPTER 2

LITERATURE SURVEY

2.1 INTRODUCTION

As digital content becomes increasingly visual, the need for accessibility tools that fill the gap between visual and textual information is more important than ever. This is particularly important for individuals with visual impairments who rely on text descriptions to recognize images on websites, social media platforms, and other digital environments. The rapid advancements in technology have provided new opportunities to address this need. Visual Description Generators, powered by artificial intelligence, are emerging as a capable solution to enhance digital convenience by automatically generating descriptive text for images.

The current challenge lies in the inconsistency and limited scope of manually created descriptions, which often fail to provide complete and accurate information. This gap can lead to important barriers for users with visual impairments. The intention of this project is to develop a Visual Description Generator that can accurately and contextually describe images, making digital content more inclusive and reachable. By using machine learning techniques, this system seeks to analyze visual data and produce meaningful and accurate descriptions.

2.2 EXISTING SYSTEM

CNN-LSTM Architectures: Already existing image captioning systems use CNNs to extract image features and feed incorporate them into an LSTM to generate captions. Even though it is effective, these models struggle with complex or abstract images.

Transformer-based Models: More recent methods have shifted to using transformers, which have proven superior in tasks that require modeling long-range dependencies, such as image captioning. Transformers are used in both image feature extraction and text generation due to their ability to handle large-scale data and parallelize computations efficiently.

2.3 DISADVANTAGES

Existing image captioning systems have several disadvantages that limit their effectiveness. Many of these models lack transparency, making it difficult to understand how they arrive at specific captions. Additionally, they often generate generic and repetitive descriptions, failing to capture the unique features of the images. This results in a lack of diversity and relevance in the captions produced. Furthermore, traditional models struggle with contextual understanding, leading to captions that may be grammatically correct but semantically disconnected from the visual content. These limitations underscore the need for more advanced techniques to enhance the quality and relevance of image captions.

2.4 PROPOSED SYSTEM

The proposed system aims to overcome the limitations of existing systems by using advanced machine learning algorithms. The proposed system uses Vision Transformer (ViT) for extracting features from images and GPT-2 for generating captions. The system will be trained on significant datasets containing various images paired with detailed descriptions. This training will allow the model to understand and describe a variety of visual content accurately. The goal of the project is not only to provide accurate and detailed descriptions but also to adapt to different circumstances and user needs, making digital content more reachable to a wider audience.

2.4 ADVANTAGES

High-quality feature extraction: ViT can extract more useful features from image compared to CNN's.

Better language generation: GPT-2 generates more easy and contextually accurate descriptions as it is pre-trained on massive text datasets.

Scalability: The model can handle large datasets and complex image types, making it suitable for real-world applications.

CHAPTER 3

SYSTEM ANALYSIS

3.1 FUNCTIONAL REQUIREMENTS

1. Image Input:

The system allows users to upload a single image or a directory containing multiple images for generating captions.

2. Image Preprocessing:

The system uses a pre-trained ViT (Vision Transformer) model to preprocess the images and convert them into pixel values. Images are converted to RGB mode if necessary to ensure compatibility with the model.

3. Caption Generation:

The system leverages the VisionEncoderDecoderModel (ViT-GPT2) to generate captions for the input images and uses beam search to improve the quality of generated captions.

4. Caption Decoding:

The system decodes the tokenized captions using GPT-2 tokenizer into human understandable text.

5. Display and Output:

The system displays each image alongside the generated caption. Outputs captions in a clean format, without any unnecessary tokens or formatting issues.

3.2 NON-FUNCTIONAL REQUIREMENTS

1. Performance:

The system should efficiently process the images, even when batch processing multiple images from a directory. The system should generate captions for an image within a few seconds.

2. Scalability:

The system must be able to handle single images as well as multiple images in batch without a significant drop in performance.

3. Accuracy:

Captions generated should be meaningful, coherent, and accurate to the content of the image. Beam search should help improve caption generation quality.

4. Usability:

Users should be able to easily input images (either single or in batches) without technical knowledge. The captions should be intuitive and meaningful to non-technical users.

5. Security:

Ensure safe handling of images uploaded by users, especially if deployed on a server.

3.3 INTERFACE REQUIREMENTS

1. User Requirements:

- Basic knowledge of uploading files (images) from a local device.
- Ability to interpret and validate the generated captions.

2. System Requirements:

Software Requirements:

- Operating System: Windows, Linux, macOS.
- Programming Language: Python 3.x.
- IDE: Google Colab, VS Code, Jupyter Notebook.
- Required Libraries: transformers, PIL, torch, matplotlib, huggingface_hub, etc.

Hardware Requirements:

- CPU: Minimum 2-core processor (e.g., Intel Core i3).
- RAM: 8GB or more (for handling image data).
- GPU: (Optional) CUDA-supported GPU for faster caption generation.
- Storage: At least 5GB of available space for model downloads.

3. Performance Requirements:

The system should load models and process images efficiently. Caption generation should not take more than a few seconds per image on modern hardware.

CHAPTER 4

SYSTEM DESIGN

4.1 UML DIAGRAMS

UML stands for Unified Modeling Language. UML is a standardized general- purpose modeling language in the field of object-oriented software engineering. The standard is managed, and was created by, the Object Management Group. The goal is for UML to become a common language for creating models of object oriented computersoftware. In its current form UML is comprised of two major components: a Meta-modeland a notation. In the future, some form of method or process may also be added to; or associated with, UML.

The Unified Modeling Language is a standard language for specifying, Visualization, Constructing and documenting the artifacts of software system, as well asfor business modeling and other non-software systems. The UML represents a collectionof best engineering practices that have proven successful in the modeling of large and complex systems.

The UML is a very important part of developing objects oriented software and the software development process. The UML uses mostly graphical notations to express the design of software projects.

Goals:

The Primary goals in the design of the UML are as follows:

- Provide users a ready-to-use, expressive visual modeling Language so that they can develop and exchange meaningful models.
- Provide extendibility and specialization mechanisms to extend the core concepts.
- Be independent of particular programming languages and development process.
- Provide a formal basis for understanding the modeling language.
- Support higher level development concepts such as collaborations, frameworks, patterns and components.

4.1.1 Data Flow Diagram

A data flow diagram (DFD) maps out the flow of information for any process or system. It uses defined symbols like rectangles, circles and arrows, plus short text labels, to show data inputs, outputs, storage points and the routes between each destination. There are 2 levels in the below diagram.

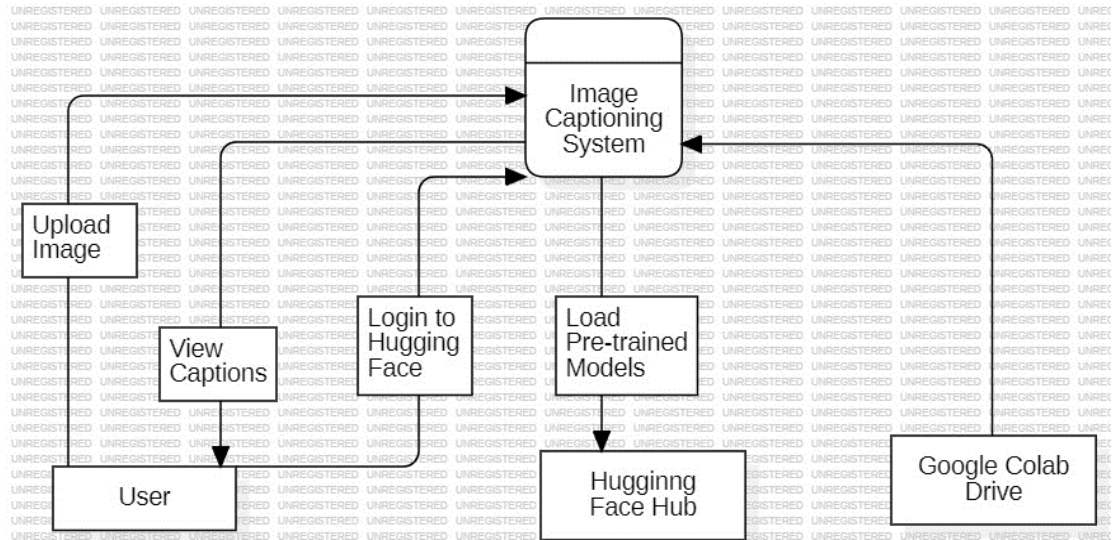


Fig:4.1.1 Data Flow Diagram: Level 0

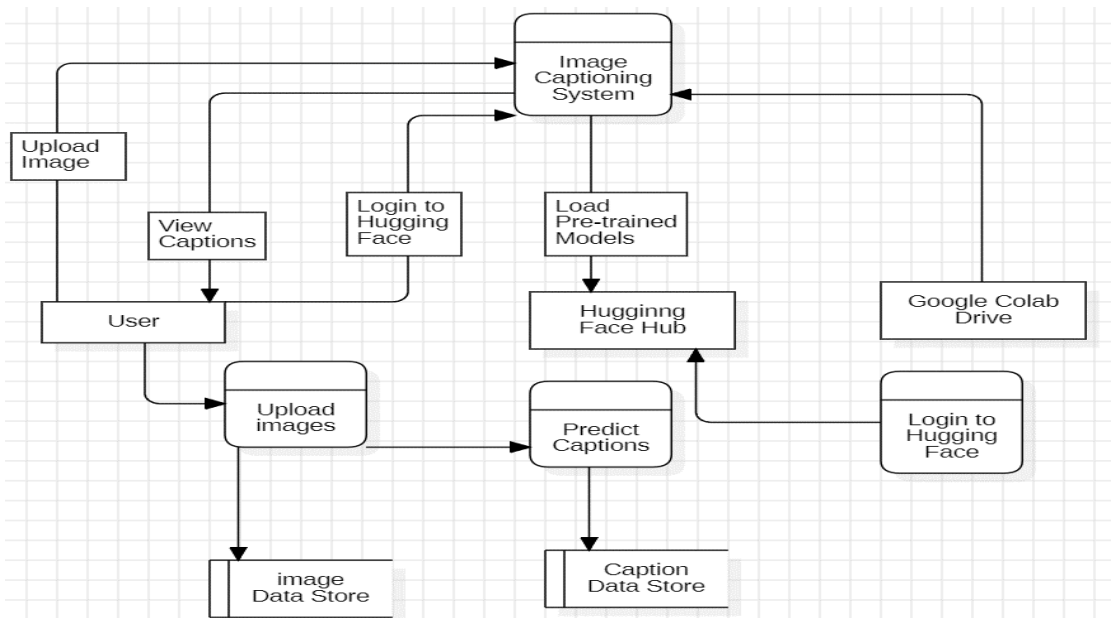


Fig :4.1.1.1 Data Flow Diagram: Level 1

4.1.2 Use Case Diagram:

A use case diagram in the Unified Modeling Language (UML) is a type of behavioral diagram defined by and created from a Use-case analysis. Its purpose is to present a graphical overview of the functionality provided by a system in terms of actors, their goals (represented as use cases), and any dependencies between those use cases. The main purpose of a use case diagram is to show what system functions are performed for which actor. Roles of the actors in the system can be depicted.

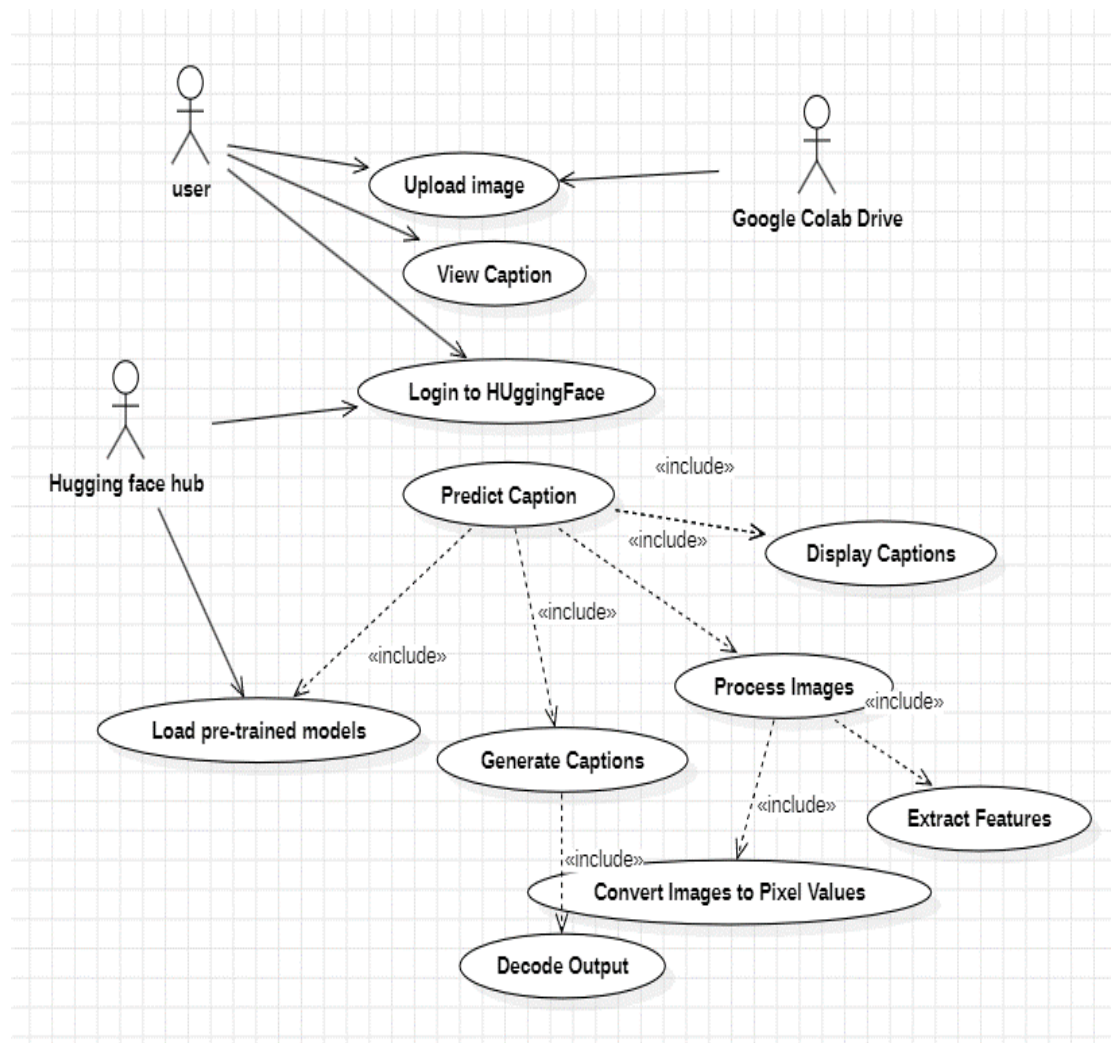


Fig : 4.1.2 Use Case Diagram

4.1.3 Class Diagram:

In software engineering, a class diagram in the Unified Modeling Language (UML) is a type of static structure diagram that describes the structure of a system by showing the system's classes, their attributes, operations (or methods), and the relationships among the classes. It explains which class contains information.

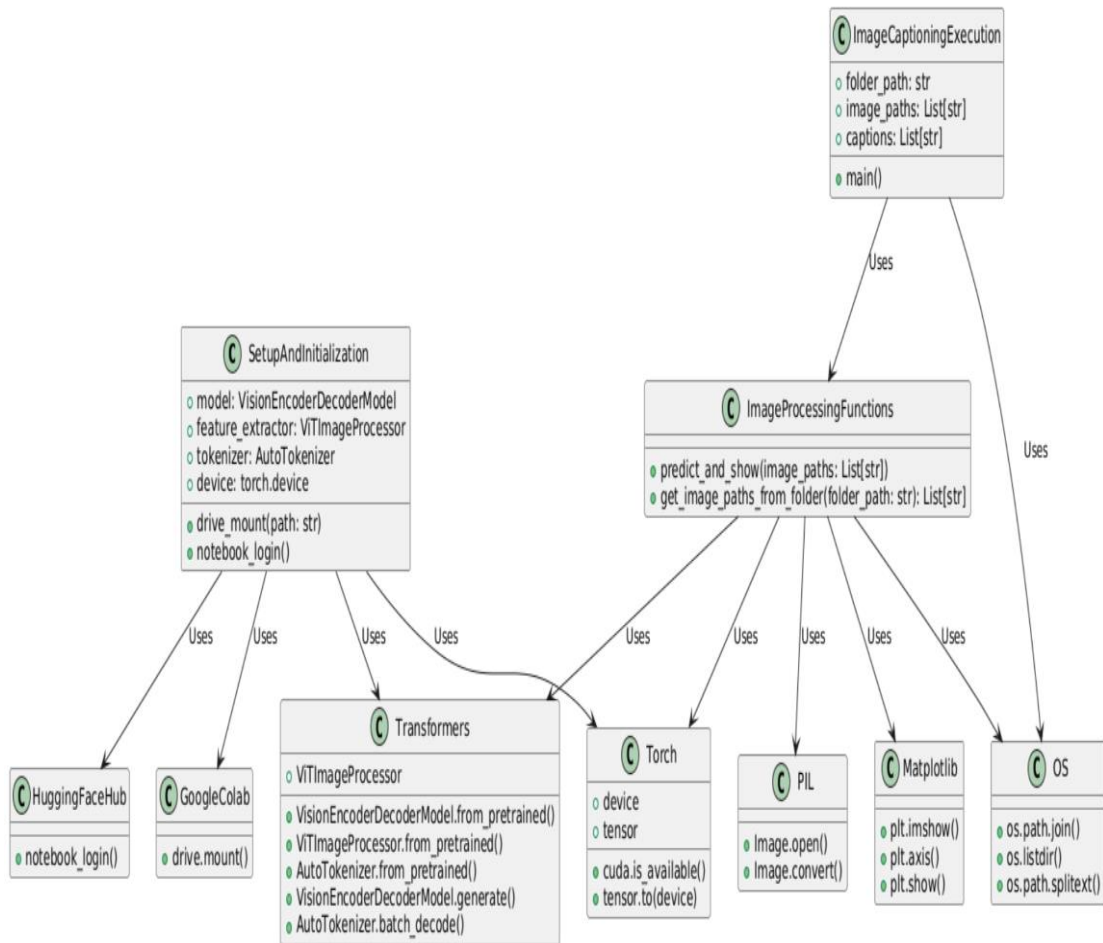


Fig :4.1.3 Class Diagram

4.1.4 Sequence Diagram:

A sequence diagram in Unified Modelling Language (UML) is a kind of interaction diagram that shows how processes operate with one another and in what order. It is a construct of a Message Sequence Chart. Sequence diagrams are sometimes called event diagrams, event scenarios, and timing diagrams.

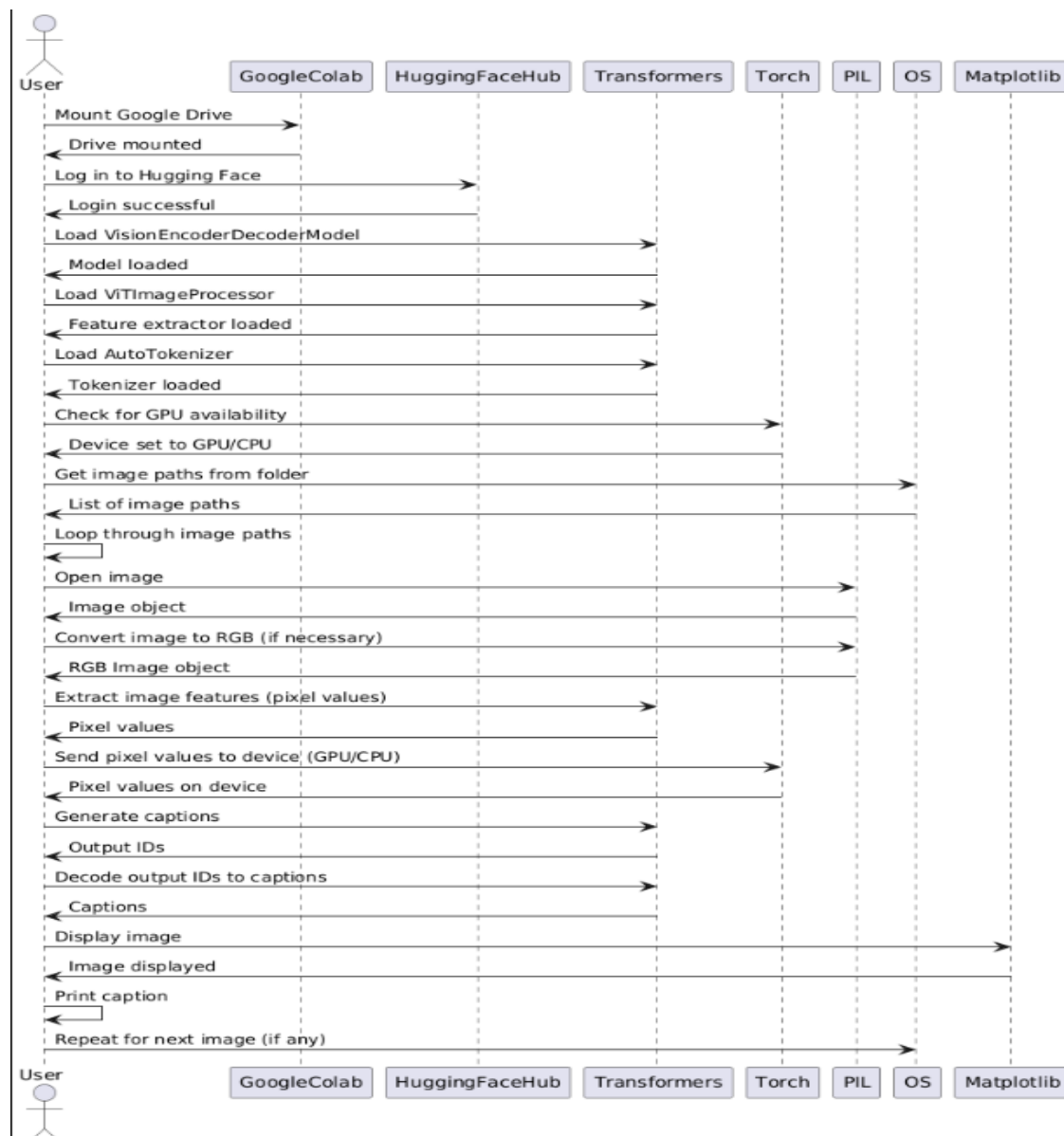


Fig:4.1.4 Sequence Diagram

4.1.5 Activity Diagram:

Activity diagrams are graphical representations of workflows of stepwise activities and actions with support for choice, iteration and concurrency. In the Unified Modeling Language, activity diagrams can be used to describe the business and operational step-by-step workflows of components in a system. An activity diagram shows the overall flow of control.

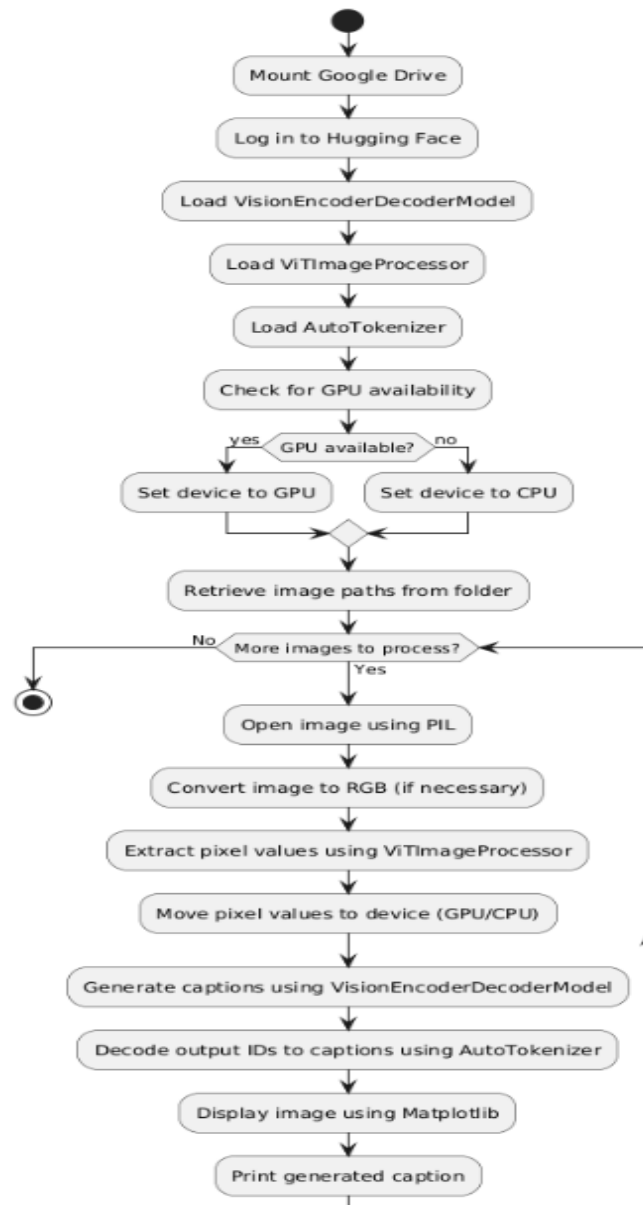


Fig:4.1.5 Activity Diagram

4.2 MODULES

1. Image Input:

Users can input either a single image file or a directory of images.

2. Image Preprocessing:

Converts images to RGB mode (if necessary) and uses the ViTImageProcessor to extract pixel values.

3. Caption Generation:

The system uses a pre-trained ViT-GPT2 model to generate captions. Implements beam search to ensure higher-quality captions.

4. Caption Decoding:

The Auto Tokenizer decodes the generated caption into human-readable text, removing special tokens.

5. Display Results:

Each image is displayed alongside its generated caption using matplotlib. Results are presented in a user-friendly format.

6. Error Handling:

If an invalid file type is inputted, the system throws an error message. Handles directory inputs gracefully, processing each image sequentially.

CHAPTER 5

IMPLEMENTATION AND RESULT

5.1 METHOD OF IMPLEMENTATION

The implementation phase focuses on transforming the theoretical design into a working system. This section explains the implementation of key modules, the methodology used, and presents the results generated by the system.

5.1.1 What is Python

Python is a high-level, general-purpose programming language known for its readability, simplicity, and versatility. It is widely used across various domains, including web development, data science, artificial intelligence, machine learning, automation, and more. Python's design emphasizes code readability, allowing developers to write clean, logical code for both small and large-scale projects.

Here's an overview of Python:

Key Features of Python:

Easy to Read and Write: Python has a simple syntax that mimics natural language, making it easy for beginners to learn and use.

Interpreted Language: Python code is executed line-by-line, which makes debugging easier. You don't need to compile the code before running it.

Dynamically Typed: In Python, you don't need to declare the type of a variable. The type is inferred at runtime.

Open Source: Python is free to use, and its source code is open, allowing anyone to contribute.

The biggest strength of Python is huge collection of standard libraries which can be used for the following –

- Machine Learning
- GUI Applications (like Kivy, Tkinter, PyQt etc.)
- Web frameworks like Django (used by YouTube, Instagram, Dropbox)Image processing (like Opencv, Pillow)

- Web scraping (like Scrapy, BeautifulSoup, Selenium) Test frameworks
- Multimedia

5.1.2 What is Machine Learning

Machine learning (ML) is a branch of artificial intelligence (AI) that enables computers to learn from and make decisions based on data, without being explicitly programmed for every task. In machine learning, algorithms are trained on data, allowing them to recognize patterns, make predictions, and improve their performance over time. The key idea is that the system can learn from experience (data) and generalize from it to new situations.

Key Concepts in Machine Learning:

1.Data Input: Machine learning models require data to learn from. This data can be labeled (supervised learning) or unlabeled (unsupervised learning).

2.Algorithms: These are mathematical models that adjust themselves to improve performance based on the data. Common algorithms include linear regression, decision trees, neural networks, and support vector machines.

3.Training: The process where the model learns from the data by adjusting parameters to minimize errors.

4.Evaluation: Assessing the model's performance on unseen data to ensure it generalizes well.

5.2 MODULES

5.2.1 Google Colab Integration:

This module allows you to access your Google Drive directly from a Colab notebook, making it easier to load files (in this case, images) stored in your drive.

drive.mount('/content/drive'):

This command mounts your Google Drive to the /content/drive directory in your Colab environment, allowing you to read and write files stored in your Google Drive.

5.2.2 Transformers:

The transformers module in Python, provided by Hugging Face, is a popular library that simplifies working with state-of-the-art natural language processing (NLP) and computer vision models. It provides a unified interface for using various pre-trained models such as BERT, GPT, T5, Vision Transformers (ViT), and more. The library is designed to make it easy to download, fine-tune, and apply transformer-based models for tasks like text classification, translation, image captioning, and more.

Key Features of the transformers Module

- Pre-trained Models
- TokenizersCustom
- Training Compatibility with PyTorch and TensorFlow

5.2.3 Torch (PyTorch)

The torch module is a core part of the PyTorch library, widely used for machine learning and deep learning tasks. PyTorch revolves around tensors, which are multidimensional arrays similar to NumPy arrays but optimized for GPU computation, making them ideal for high-performance tasks. One of the key strengths of PyTorch is its dynamic computational graph, which allows for on-the-fly construction and modification of neural networks, offering flexibility during development.

PyTorch includes a powerful automatic differentiation tool called Autograd, which enables automatic calculation of gradients, making it easier to optimize neural networks. It also provides a rich set of APIs for creating, training, and evaluating machine learning models, supporting a wide variety of architectures, from simple linear models to complex convolutional or recurrent neural networks.

5.2.4 PIL (Python Imaging Library)

The PIL module refers to the Python Imaging Library, which was a popular library for image processing tasks in Python. However, it's important to note that PIL has been succeeded by Pillow, which is a maintained and more feature-rich fork of PIL.

Pillow offers a wide range of image processing capabilities, including opening, manipulating, and saving various image formats. It supports operations such as resizing, cropping, rotating, filtering, and color adjustments. Pillow also provides tools for drawing shapes and text on images, making it useful for both basic image manipulation and more complex image generation tasks.

5.2.5 Matplotlib

Matplotlib is a widely-used Python library for creating static, animated, and interactive visualizations. It provides a flexible and comprehensive suite of plotting functions that enable users to generate a variety of plots and charts, such as line graphs, bar charts, scatter plots, histograms, and heatmaps.

The library is particularly known for its pyplot module, which offers a MATLAB-like interface for creating and customizing plots with ease. With matplotlib, users can control every aspect of their visualizations, including axes, labels, legends, and colors. It also supports integration with various backends for rendering plots in different environments, such as Jupyter notebooks, web applications, or standalone windows.

5.2.6 Hugging Face Hub

The Hugging Face Hub is a platform provided by Hugging Face for hosting and sharing machine learning models, datasets, and other resources. It serves as a central repository where users can:

Access Pretrained Models: Users can find and utilize a wide range of pretrained models for tasks such as natural language processing, computer vision, and speech recognition. These models are often developed and shared by the community, enabling easier and faster experimentation.

Share Models and Datasets: Researchers and developers can upload and share their own models and datasets with the community. This fosters collaboration and allows others to build upon existing work.

5.2.7 OS Module

The OS module in Python provides a way to interact with the operating system's functionalities and manage file and directory operations. It includes functions for:

File and Directory Management: You can create, delete, and rename files and directories, as well as navigate the file system. For example, you can list the contents of a directory, check if a path exists, or change the current working directory.

Path Manipulation: The module includes utilities for working with file paths, such as joining and splitting paths, and extracting file extensions or base names.

Environment Variables: It allows you to access and modify environment variables, which are used to configure system-wide settings and behaviour.

Overall, the OS module is essential for performing low-level operating system operations and file system management within Python scripts.

5.3 EXTENSION OF KEY FUNCTION

5.3.1 Setup and Imports

```
from google.colab import drive
drive.mount('/content/drive')
```

Fig:5.3.1 Setup

This mounts Google Drive to the Colab environment, allowing you to access files stored in your Google Drive from within Colab.

Imports various libraries

```
# Import necessary libraries
from transformers import VisionEncoderDecoderModel, ViTImageProcessor, AutoTokenizer
import torch
from PIL import Image
import matplotlib.pyplot as plt
from huggingface_hub import notebook_login
import os
# Log in to Hugging Face to access models and datasets
notebook_login()
```

Fig:5.3.1.1 Importing Libraries

5.3.2 Model Loading

```
[ ] # Load the pre-trained model for vision-encoder-decoder tasks
    model = VisionEncoderDecoderModel.from_pretrained("nlpconnect/vit-gpt2-image-captioning")
    # Load the feature extractor (image processor) associated with the model
    feature_extractor = ViTImageProcessor.from_pretrained("nlpconnect/vit-gpt2-image-captioning")
    # Load the tokenizer associated with the model for decoding text output
    tokenizer = AutoTokenizer.from_pretrained("nlpconnect/vit-gpt2-image-captioning")
```

Fig:5.3.2 Model Loading

Loads a pre-trained vision-encoder-decoder model, its associated feature extractor (image processor), and tokenizer from Hugging Face's model hub.

```
# Set the device to GPU if available, otherwise use CPU
device = torch.device("cuda" if torch.cuda.is_available() else "cpu")
model.to(device) # Move the model to the specified device
```

Fig:5.3.2.1 Set the device to GPU

Sets the device to GPU if available; otherwise, uses CPU. Moves the model to the specified device for computations.

5.3.3 Define Image Captioning Parameters

```
[ ] # Define image captioning parameters
    max_length = 16 # Maximum length of the generated caption
    num_beams = 4 # Number of beams for beam search (used to generate captions with higher quality)
    gen_kwargs = {"max_length": max_length, "num_beams": num_beams}
```

Fig:5.3.3 Image captioning Parameters

Sets parameters for caption generation:

max_length: Maximum length of the caption.

num_beams: Number of beams for beam search to enhance caption quality.

5.3.4 Define the predict and show Function

Defines a function to handle image processing and caption generation.

```

def predict_and_show(image_paths):
    images = [] # List to store the images

    # Check if image_paths is a directory or a single file
    if os.path.isdir(image_paths):
        # Process all images in the directory
        for filename in os.listdir(image_paths):
            file_path = os.path.join(image_paths, filename)
            if file_path.lower().endswith(('.png', '.jpg', '.jpeg', '.bmp', '.gif')):
                i_image = Image.open(file_path)
                if i_image.mode != "RGB":
                    i_image = i_image.convert(mode="RGB")
                images.append(i_image)
    else:
        # Process a single image file
        i_image = Image.open(image_paths)
        if i_image.mode != "RGB":
            i_image = i_image.convert(mode="RGB")
        images.append(i_image)

```

Fig:5.3.4 Predict and Show Function

Checks if image_paths is a directory or a single file.

If it's a directory, processes all images in the directory.

If it's a single image file, processes that single image.

```

# Convert the list of images to pixel values using the feature extractor
pixel_values = feature_extractor(images=images, return_tensors="pt").pixel_values
pixel_values = pixel_values.to(device)

# Generate captions for the images using the model
output_ids = model.generate(pixel_values, **gen_kwargs)

# Decode the generated output into human-readable text
preds = tokenizer.batch_decode(output_ids, skip_special_tokens=True)

# Clean up the predictions by stripping any leading/trailing spaces
preds = [pred.strip() for pred in preds]

#return preds
for i_image, pred in zip(images, preds):
    plt.imshow(i_image)
    plt.axis('off')
    plt.show()
    print(f"Generated Caption: {pred}")

```

Fig: 5.3.4.1 Generate the Captions

Converts images to pixel values using the feature extractor and moves the tensors to the specified device.

Decodes the generated output into human-readable text and removes any leading/trailing spaces.

Displays each image and its corresponding generated caption using matplotlib.

5.3.5 Example Usage

Sets the path to a directory containing images.

Calls the `predict_and_show` function with the directory path and prints the generated captions.

```
# For a folder:  
image_paths = '/content/drive/MyDrive/project pictures'  
  
# For a single image:  
# image_paths = '/content/single_image.jpg'  
  
captions = predict_and_show(image_paths)  
print("Generated Captions:", captions)
```

Fig: 5.3.5 Example Usage

5.4 OUTPUT SCREENS

```
/usr/local/lib/python3.10/dist-packages/huggingface_hub/utils/_token.py:89: UserWarning:
The secret 'HF_TOKEN' does not exist in your Colab secrets.
To authenticate with the Hugging Face Hub, create a token in your settings tab (https://huggingface.co/settings/tokens), set it as secret in your Google Colab.
You will be able to reuse this secret in all of your notebooks.
Please note that authentication is recommended but still optional to access public models or datasets.
warnings.warn(
config.json: 100% ██████████ 4.61k/4.61k [00:00<00:00, 195kB/s]
pytorch_model.bin: 100% ██████████ 982M/982M [00:08<00:00, 196MB/s]
preprocessor_config.json: 100% ██████████ 228/228 [00:00<00:00, 14.6kB/s]
tokenizer_config.json: 100% ██████████ 241/241 [00:00<00:00, 12.8kB/s]
vocab.json: 100% ██████████ 798k/798k [00:00<00:00, 19.5MB/s]
merges.txt: 100% ██████████ 456k/456k [00:00<00:00, 9.52MB/s]
tokenizer.json: 100% ██████████ 1.36M/1.36M [00:00<00:00, 14.6MB/s]
special_tokens_map.json: 100% ██████████ 120/120 [00:00<00:00, 4.34kB/s]
/usr/local/lib/python3.10/dist-packages/transformers/tokenization_utils_base.py:1601: FutureWarning: `clean_up_tokenization_spaces` was not set. It will be set to `True` by default in an upcoming version of Transformers. To prevent this warning, you should explicitly set `clean_up_tokenization_spaces` to `False` in your code. This will have no effect on your current training or inference, but we recommend you set it to `False` before the next release of Transformers.
warnings.warn(
```

Fig:5.4.1 Model training and building



Generated Caption: a large group of people posing for a picture

Fig:5.4.2 Output when a file is uploaded



Generated Caption: a large jetliner flying over a large body of water



Generated Caption: a man riding a wave on top of a surfboard

5.4.3 Output when a folder is uploaded

CHAPTER 6

SYSTEM TEST

The purpose of testing is to discover errors. Testing is the process of trying to discover every conceivable fault or weakness in a work product. It provides a way to check the functionality of components, sub-assemblies, assemblies and/or a finished product. It is the process of exercising software with the intent of ensuring that the Software system meets its requirements and user expectations and does not fail in an unacceptable manner. There are various types of test. Each test type addresses a specific testing requirement.

6.1 TYPES OF TESTS

Unit testing

Unit testing involves the design of test cases that validate that the internal program logic is functioning properly, and that program inputs produce valid outputs. All decision branches and internal code flow should be validated. It is the testing of individual software units of the application. It is done after the completion of an individual unit before integration. This is a structural testing, that relies on knowledge of its construction and is invasive. Unit tests perform basic tests at component level and test a specific business process, application, and/or system configuration. Unit tests ensure that each unique path of a business process performs accurately to the documented specifications and contains clearly defined inputs and expected results.

Integration testing

Integration tests are designed to test integrated software components to determine if they actually run as one program. Testing is event driven and is more concerned with the basic outcome of screens or fields. Integration tests demonstrate that although the components were individually satisfactory, as shown by successfully unit testing, the combination of components is correct and consistent. Integration testing is specifically aimed at exposing the problems that arise from the combination of components.

Functional test

Functional tests provide systematic demonstrations that functions tested are available as specified by the business and technical requirements, system documentation, and user manuals. Organization and preparation of functional tests is focused on requirements, key functions, or special test cases. In addition, systematic coverage pertaining to identify Business process flows; data fields, predefined processes, and successive processes must be considered for testing. Before functional testing is complete, additional tests are identified and the effective value of current tests is determined.

System Test

System testing ensures that the entire integrated software system meets requirements. It tests a configuration to ensure known and predictable results. An example of system testing is the configuration-oriented system integration test. System testing is based on process descriptions and flows, emphasizing pre-driven process links and integration points.

White Box Testing

White Box Testing is a testing in which in which the software tester has knowledge of the inner workings, structure and language of the software, or at least its purpose. It is purpose. It is used to test areas that cannot be reached from a black box level.

Black Box Testing

Black Box Testing is testing the software without any knowledge of the inner workings, structure or language of the module being tested. Black box tests, as most other kinds of tests, must be written from a definitive source document, such as specification or requirements document, such as specification or requirements document. It is a testing in which the software under test is treated, as a black box. you cannot “see” into it. The test provides inputs and responds to outputs without considering how the software works.

6.2 VARIOUS TESTCASE SCENARIOS

Test Case ID	Test Case Description	Test Case Steps	Test Data	Expected Result	Actual Result	Status
TC01	Mount Google Drive.	1.Use drive.mount('/content/drive'). 2.Provide access to Google Drive.	Drive mount path /content/drive.	Google Drive should be mounted successfully.	Google Drive mounted successfully.	Success
TC02	Load Pre-trained Vision Transformer Model.	1.Load the pre-trained Model using VisionEncoder-Decoder-Model.from_pre-trained. 2. Load the image processor and tokenizer using the respective libraries.	Model: nlp-connect/vit-gpt2-image-captioning	The pre-trained model, feature extractor, and tokenizer should be loaded successfully.	Model, feature extractor, and tokenizer loaded successfully.	Success

TC03	Predict Caption for Single Image	1. Call the predict_and_show function with a single image path. 2. The model should generate a caption for the image.	Image file in /content/single_image.jpg	The caption should be generated and displayed successfully for the image.	Caption generated and displayed successfully.	Success
TC04	Predict Captions for Multiple Images	1. Call the predict_and_show function with a directory of images. 2. The model should generate captions for all images.	Directory path /content /Drive /MyDrive/project pictures	Captions should be generated and displayed successfully for all images in the directory.	Captions generated and displayed successfully.	Success

6.2 Test Cases

CHAPTER 7

CONCLUSION AND FUTURE ENHANCEMENT

7.1 PROJECT CONCLUSION

The **Visual Description Generator with Vision Transformers** successfully implements an advanced image captioning system using state-of-the-art vision and language models. By integrating a pre-trained Vision Transformer (ViT) model and GPT-2, the system is able to generate accurate and contextually relevant descriptions for both individual images and batches of images. The project demonstrated the effective use of Hugging Face transformers for vision-based tasks, leveraging beam search to improve caption quality. The system's ability to handle various image formats (such as non-RGB images) and its integration with cloud services like Google Drive make it versatile and accessible.

7.2 FUTURE ENHANCEMENT

1. Interactive User Interface (UI):

Develop a graphical user interface (GUI) to allow users to upload images directly from their devices and receive captions in real-time, enhancing usability for non-technical users.

2. Multi-language Support:

Extend the model to generate captions in multiple languages by integrating multilingual transformers. This would make the system more accessible to a global audience.

3. Real-time Captioning for Videos:

Enhance the model to support real-time video description by generating captions for each frame or key frames of a video. This could be useful for accessibility in media content or video summarization.

4. Optimization for Mobile and Edge Devices:

Optimize the model for mobile and edge devices by reducing its size or using lighter versions of the transformer models, making it suitable for real-time applications on lower-powered devices.

CHAPTER 8

REFERENCES

8.1 PAPER REFERENCES

- [1] Xu, Kelvin, et al. "Show, attend and tell: Neural image caption generation with visual attention.", International Conference on Machine Learning (ICML), 2015.
- [2] Radford, Alec, et al. "Learning Transferable Visual Models from Natural Language Supervision.", Proceedings of the International Conference on Machine Learning (ICML), 2021.
- [3] Vaswani, Ashish, et al, "Attention is All You Need." Advances in Neural Information Processing Systems (NeurIPS), 2017.
- [4] Dosovitskiy, Alexey, et al. "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale.", International Conference on Learning Representations (ICLR), 2021.
- [5] Chen, Xinlei, et al. "UNITER: Universal Image-Text Representation Learning.", European Conference on Computer Vision (ECCV), 2020.
- [6] Huang, Lun, et al. "Attention on Attention for Image Captioning.", IEEE/CVF International Conference on Computer Vision (ICCV), 2019.
- [7] Chen, Lin, et al. "SCA-CNN: Spatial and channel-wise attention in convolutional networks for image captioning.", IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.
- [8] Rennie, Steven J., et al. "Self-critical sequence training for image captioning.", IEEE Conference on Computer Vision and Pattern Recognition (CVPR), 2017.
- [9] Zhou, Luowei, et al. "Unified Vision-Language Pre-Training for Image Captioning and VQA.", AAAI Conference on Artificial Intelligence, 2020.
- [10] Li, Yikang, et al. "Visual Semantic Reasoning for Image-Text Matching.", IEEE/CVF International Conference on Computer Vision (ICCV), 2019.

8.2 WEBSITES

- [1] <https://www.kaggle.com/datasets/umongsain/vit-gpt2-image-captioning>
- [2] <https://ieeexplore.ieee.org/document/10551257>
- [3] <https://ieeexplore.ieee.org/document/10099161>
- [4] <https://www.google.com.in>
- [5] <https://claude.ai>
- [6] <https://www.geeksforgeeks.org>