

Image Describer Net:

A DEEP LEARNING APPROACH TO IMAGE CAPTION GENERATION

Introduction:

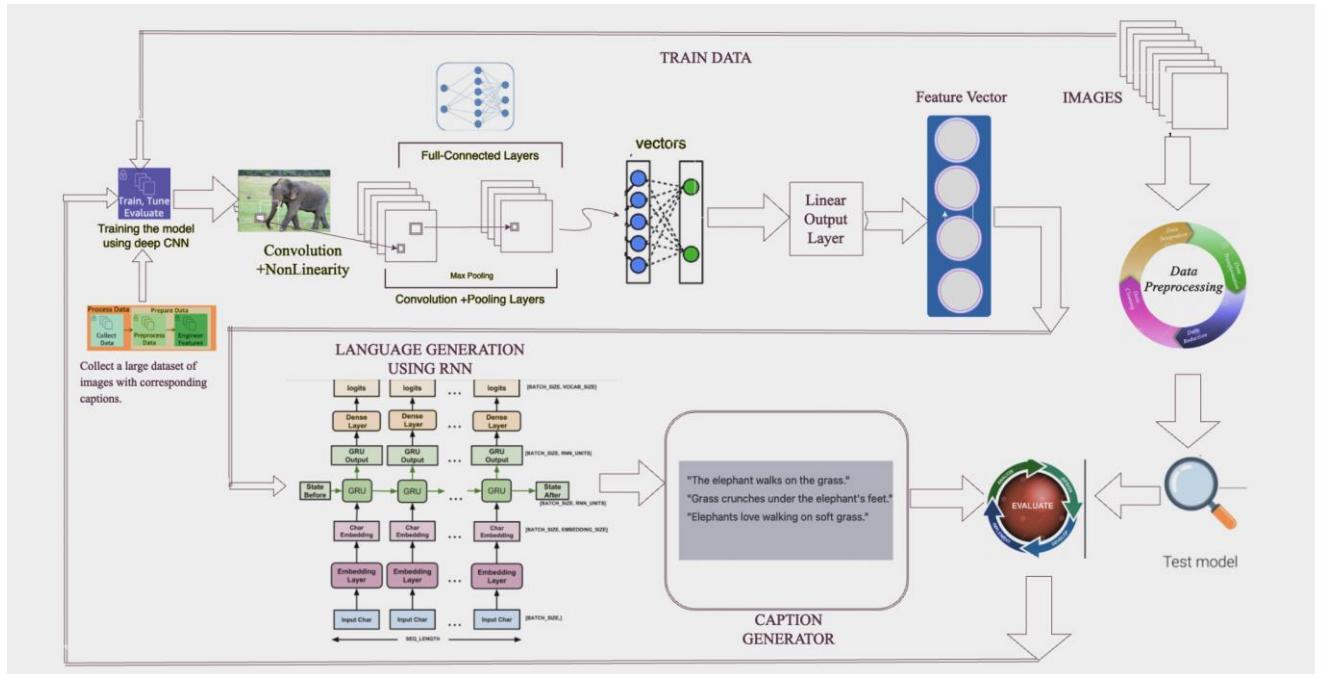
In the realm of artificial intelligence and computer vision, the task of image caption generation involves creating a human-readable textual description based on the content of an image. While this task is intuitive for humans, it poses a significant challenge for computers, as it necessitates the comprehension of visual content and the ability to articulate this understanding in natural language. The advent of deep learning methods, particularly Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs) like Long Short-Term Memory (LSTM), has revolutionized the field, outperformed classical approaches and achieved state-of-the-art results in automatically generating image captions.

This project harnesses the power of deep neural network models, specifically CNNs and LSTMs, to automatically generate descriptive captions for images, exemplifying the advancements in deep learning techniques. The success of these models is contingent upon extensive datasets and computational resources, enabling the development of sophisticated models capable of interpreting images and translating that understanding into coherent and contextually relevant language. By employing CNNs for image feature extraction and LSTM for sequential data processing, the project illustrates the synergy between computer vision and natural language processing, showcasing the potential of deep learning in bridging the gap between visual content and textual descriptions.

As the project unfolds, it exemplifies the growth and significance of deep learning in image caption generation, demonstrating the effectiveness of techniques like CNNs and LSTMs. The utilization of large datasets and powerful computing resources underscores the data-driven nature of these approaches, emphasizing the importance of scale in training models that can seamlessly generate captions for diverse images. The project serves as an insightful exploration into the application of deep learning for image captioning,

shedding light on the evolving landscape of AI and its impactful role in understanding and describing visual content.

Technical Architecture:



DEFINITIONS:

- **Embedding Layer:** In image captioning, the Embedding layer is like a dictionary that converts words into vectors. It helps the model understand the relationships between different words in the captions.
- **LSTM Layer:** The Long Short-Term Memory (LSTM) layer is a type of memory unit in the neural network that allows the model to remember information from earlier words in the captions, helping it generate coherent and contextually relevant sentences.
- **Dense Layer:** The Dense layer produces the final output, predicting the next word in the caption. It takes the information learned by the previous layers and maps it to the vocabulary, providing the most likely word based on the given context.
- **Vocabulary Size:** This is the total number of unique words in your dataset. It helps the model understand the range of words it might encounter in captions.

- **Embedding Dimensions:** The number of dimensions in which each word will be represented as a vector. Higher dimensions allow the model to capture more intricate relationships between words.
- **Adam Optimizer:** Adam is an optimization algorithm that adjusts the model's weights during training, helping it converge to an optimal solution efficiently.
- **Categorical Cross entropy Loss:** This is a measure of how well the model's predictions match the actual captions. It penalizes the model more when it makes larger mistakes, guiding it to improve its caption predictions.
- **VGG16:** VGG16 is a deep convolutional neural network architecture widely used for image classification. It consists of 16 weight layers, including convolutional and fully connected layers, and is known for its simplicity and effectiveness. Pre-trained on large datasets like ImageNet, VGG16 can extract intricate features from images, making it valuable for various computer vision tasks.
- **CNN (Convolutional Neural Network):** A type of deep neural network designed for visual pattern recognition. It uses convolutional layers to automatically learn hierarchical representations of features from input images, making it effective for tasks like image classification.
- **LSTM (Long Short-Term Memory):** A type of recurrent neural network (RNN) architecture that excels in capturing long-term dependencies in sequential data. LSTMs are well-suited for tasks involving time-series data or natural language processing, where context and memory over time are crucial.
- **RNN (Recurrent Neural Network):** A class of neural networks designed to work with sequential data by maintaining hidden states that capture information from previous inputs. While effective for certain tasks, traditional RNNs struggle with long-term dependencies, leading to the development of more advanced architectures like LSTMs.
- **FEATURE EXTRACTION:** Feature extraction refers to the process of transforming raw input data, often high-dimensional and complex, into a reduced set of relevant and informative features. These features capture essential patterns or characteristics, facilitating more efficient and effective analysis by machine learning algorithms.
- **TOKENIZER:** A tokenizer is a linguistic tool that breaks down a text into individual units, such as words or sub words, to facilitate further analysis or processing. It plays a crucial role in natural language processing tasks by converting raw text into a structured format, enabling

efficient computational understanding and manipulation of language. Tokenization is a fundamental step in text preprocessing for various language-related applications, including machine learning and information retrieval.

- **Corpus BLEU:** Corpus BLEU is a metric used for evaluating the quality of machine-generated text by measuring the overlap between the generated output and reference texts, commonly employed in tasks like machine translation.
- **SoftMax and RELU:** SoftMax is an activation function that converts a vector of raw scores into a probability distribution, while RELU (Rectified Linear Unit) is an activation function commonly used in neural networks, introducing non-linearity by outputting the input directly if positive and zero otherwise.

Project Objectives: By the project's completion, participants will achieve the following milestones:

Mastering Convolutional Neural Network (CNN) Concepts and Techniques:

Acquire a comprehensive understanding of the fundamental principles and techniques underlying Convolutional Neural Networks, a key technology in image processing and computer vision.

Deepened Insight into Image Data:

Develop a profound comprehension of image data, exploring its characteristics, formats, and the challenges associated with processing visual information.

Proficiency in Data Pre-processing Techniques:

Demonstrate the ability to pre-process and clean image data using a variety of techniques, gaining hands-on experience in enhancing data quality and preparing it for effective machine learning model training.

Web Application Development with Flask:

Learn the process of building a web application using the Flask framework, a lightweight and efficient web development tool. Understand the principles of web application architecture, routing, and integration of machine learning models into web-based interfaces.

Explore Additional Concepts in Computer Vision:

Expand knowledge beyond CNNs by exploring additional concepts in computer vision, such as feature extraction, image augmentation, and transfer learning, enriching the understanding of advanced techniques used in real-world applications.

Hands-on Experience with Deployment:

Gain practical experience in deploying machine learning models within a web application environment, ensuring a holistic understanding of the end-to-end development lifecycle.

Effective Communication of Technical Concepts:

Develop communication skills to articulate complex technical concepts clearly, enabling effective collaboration and knowledge sharing within the project team and with stakeholders.

Understanding Model Evaluation Metrics:

Familiarize yourself with model evaluation metrics specific to image-related tasks, gaining insights into assessing model performance and making informed decisions during the development process.

This project is designed to equip participants not only with technical expertise but also with the practical skills needed to apply these concepts in real-world scenarios, fostering a well-rounded proficiency in the intersection of deep learning, computer vision, and web development.

Project Flow:

- Data Collection.

Create Train and Test Folders.

- Data Pre-processing.

- Model Building

Import the model building Libraries

Initializing the model

Adding Input Layer

Adding Hidden Layer

Adding Output Layer

Configure the Learning Process

- Training and testing the model

Save the Model

Application Building

Create an HTML file

Build Python Code

Pre-requisites:

We've successfully set up the Kaggle API in Google Colab by installing the Kaggle package, creating the Kaggle directory, copying my API key, and adjusting permissions. With these configurations, I downloaded the "Flickr8k" dataset directly to my Colab environment using the Kaggle CLI. This streamlined process allows for seamless access to Kaggle datasets, facilitating data-driven projects and analyses within my Colab notebook.

```
✓ [13] ! pip install kaggle
Requirement already satisfied: kaggle in /usr/local/lib/python3.10/dist-packages (1.5.16)
Requirement already satisfied: six>=1.10 in /usr/local/lib/python3.10/dist-packages (from kaggle) (1.16.0)
Requirement already satisfied: certifi in /usr/local/lib/python3.10/dist-packages (from kaggle) (2023.7.22)
Requirement already satisfied: python-dateutil in /usr/local/lib/python3.10/dist-packages (from kaggle) (2.8.2)
Requirement already satisfied: requests in /usr/local/lib/python3.10/dist-packages (from kaggle) (2.31.0)
Requirement already satisfied: tqdm in /usr/local/lib/python3.10/dist-packages (from kaggle) (4.66.1)
Requirement already satisfied: python-slugify in /usr/local/lib/python3.10/dist-packages (from kaggle) (8.0.1)
Requirement already satisfied: urllib3 in /usr/local/lib/python3.10/dist-packages (from kaggle) (2.0.7)
Requirement already satisfied: bleach in /usr/local/lib/python3.10/dist-packages (from kaggle) (6.1.0)
Requirement already satisfied: webencodings in /usr/local/lib/python3.10/dist-packages (from bleach->kaggle) (0.5.1)
Requirement already satisfied: text-unidecode>=1.3 in /usr/local/lib/python3.10/dist-packages (from python-slugify->kaggle)
Requirement already satisfied: charset-normalizer<4,>=2 in /usr/local/lib/python3.10/dist-packages (from requests->kaggle)
Requirement already satisfied: idna<4,>=2.5 in /usr/local/lib/python3.10/dist-packages (from requests->kaggle) (3.4)

✓ [14] ! mkdir ~/.kaggle
mkdir: cannot create directory '/root/.kaggle': File exists

✓ [15] ! cp kaggle.json ~/.kaggle/

✓ [16] ! chmod 600 ~/.kaggle/kaggle.json

✓ [17] ! kaggle datasets download -d adityajn105/flickr8k
Downloading flickr8k.zip to /content
99% 1.03G/1.04G [00:09<00:00, 132MB/s]
100% 1.04G/1.04G [00:09<00:00, 116MB/s]
```

IMPORT PACKAGES

STEP-1:

To build machine learning models in Google Colab, essential Python packages are installed using the '!pip install' command. The code imports essential Python

libraries for data manipulation, machine learning, and image processing. NumPy and Pandas are utilized for efficient data handling, Scikit-learn for machine learning tasks, and PIL along with Matplotlib for image processing and visualization. The tqdm library is included for displaying progress bars during iterations. Operating system-related tasks, random number generation, and file operations are facilitated by os, random, and shutil libraries. Additionally, zipfile is imported, suggesting potential involvement in handling compressed files or datasets. This comprehensive set of libraries equips the user with versatile tools for various data science and machine learning tasks.

STEP-2:

The code imports essential modules from TensorFlow and Keras libraries for building and training neural networks. Sequential and Model classes facilitate model creation, while load_model allows loading pre-trained models. Various layers, such as Embedding, LSTM, Dense, and Dropout, are imported for constructing neural network architectures. The VGG16 model, a powerful pre-trained convolutional neural network for computer vision, is included. Functions for preprocessing images, handling sequences, and plotting model architectures are also imported. The corpus_bleu function from NLTK is brought in for evaluating text generation tasks. This set of imports provides a comprehensive toolkit for developing and accessing machine learning models in the context of image and sequence processing tasks.

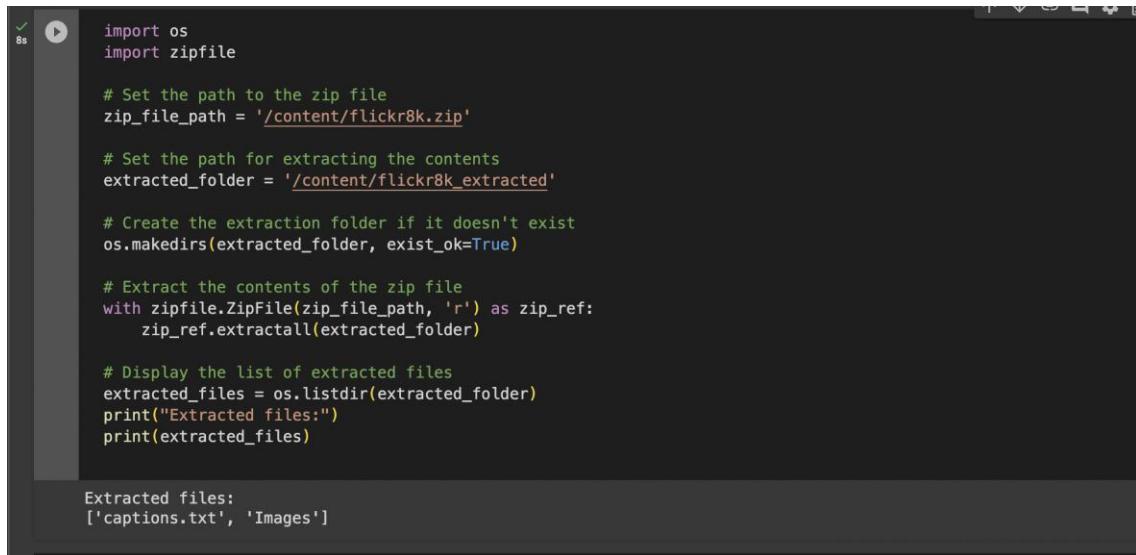
STEP-3

Step 3 encapsulates the entire process, from installing necessary packages to incorporating code snippets for constructing a machine learning model in Google Calab. This comprehensive pipeline enables users to seamlessly work on image caption generation projects in a collaborative and cloud-based environment.

```
[7] Importing necessary libraries for machine learning:  
import numpy as np  
import pandas as pd  
from sklearn.model_selection import train_test_split  
from sklearn.metrics import accuracy_score, classification_report  
from PIL import Image  
import matplotlib.pyplot as plt  
from tqdm.notebook import tqdm  
from tqdm import tqdm  
import os  
import random  
import shutil  
import zipfile  
  
[8] Importing TensorFlow and Keras for deep learning tasks:  
import tensorflow as tf  
from tensorflow import keras  
from tensorflow.keras.utils import plot_model  
from tensorflow.keras.models import Sequential, Model  
from tensorflow.keras.models import load_model  
from tensorflow.keras.preprocessing.sequence import pad_sequences  
from tensorflow.keras.layers import Embedding, LSTM, Dense, Dropout, Input, add  
from tensorflow.keras.applications.vgg16 import VGG16  
from tensorflow.keras.preprocessing.image import load_img, img_to_array, ImageDataGenerator  
from tensorflow.keras.applications.vgg16 import preprocess_input  
from tensorflow.keras.utils import to_categorical  
from nltk.translate.bleu_score import corpus_bleu  
  
[9] Importing Flask for web application development:  
from flask import Flask, render_template, request
```

Data Collection:

I established the path for the dataset zip file, naming it "flickr8k.zip". Following that, I designated a folder named "flickr8k_extracted" to serve as the destination for extracting the contents. If this folder didn't exist, I took the initiative to create it. Utilizing the zipfile module, I executed the extraction process, retrieving all files from "flickr8k.zip" and placing them neatly into "flickr8k_extracted". To offer transparency about the dataset's structure after this extraction, I thoughtfully listed and displayed the names of the extracted files.



```
import os
import zipfile

# Set the path to the zip file
zip_file_path = '/content/flickr8k.zip'

# Set the path for extracting the contents
extracted_folder = '/content/flickr8k_extracted'

# Create the extraction folder if it doesn't exist
os.makedirs(extracted_folder, exist_ok=True)

# Extract the contents of the zip file
with zipfile.ZipFile(zip_file_path, 'r') as zip_ref:
    zip_ref.extractall(extracted_folder)

# Display the list of extracted files
extracted_files = os.listdir(extracted_folder)
print("Extracted files:")
print(extracted_files)
```

Extracted files:
['captions.txt', 'Images']

Data preparation process

Data preparation process, where you unzip and organize the dataset for further analysis or model training in your project.

Mapping creation process:

I established paths, extracted contents from the Flickr8k dataset, and created a mapping dictionary linking image IDs to their respective captions. This groundwork sets the stage for subsequent tasks in building a machine learning model for image captioning.

In this step, I aggregated all captions from the mapping dictionary into a list named **all_captions**. This list contains textual descriptions associated with images, forming a comprehensive dataset for training a machine learning model in subsequent stages of the project.

```

✓ 0s # Set the path to the extracted folder
extracted_folder = '/content/flickr8k_extracted'

# Set the path to the captions file
captions_file_path = os.path.join(extracted_folder, 'captions.txt')

# Create a mapping dictionary
mapping = {}

# Process lines in the captions file
with open(captions_file_path, 'r') as captions_file:
    captions_doc = captions_file.read()

    # Process lines
    for line in tqdm(captions_doc.split('\n'), desc="Creating Mapping"):
        # Split the line by comma,
        tokens = line.split(',')

        if len(tokens) < 2:
            continue

        # Extract image_id and caption
        image_id, caption = tokens[0], tokens[1:]

        # Remove extension from image ID
        image_id = image_id.split('.')[0]

        # Convert caption list to string
        caption = " ".join(caption)

        # Create a list if needed
        if image_id not in mapping:
            mapping[image_id] = []

        # Store the caption
        mapping[image_id].append(caption)

Creating Mapping: 100% 40457/40457 [00:00<00:00, 172140.50it/s]

```

```

✓ 0s [13] len(mapping)
8092

✓ 0s [14] # Initialize a list to store all captions
all_captions = []

# Iterate over keys in the mapping dictionary
for key in mapping:
    # Iterate over captions in the current key
    for caption in mapping[key]:
        # Append the caption to the list
        all_captions.append(caption)

✓ 0s [15] len(all_captions)
40456

```

Initializing the model:

VGG16 model

The provided code loads the VGG16 model, a well-known convolutional neural network (CNN) pre-trained on the ImageNet dataset for image classification. It then restructures the model by removing its last layer, which is typically a dense layer responsible for image classification. This modification transforms the VGG16 model into a feature extractor, preserving its ability to capture high-level features from images.

The significance of this modification lies in the concept of transfer learning. By excluding the final classification layer, the modified model can be

used as a powerful feature extractor for various image-related tasks, such as image retrieval or image captioning. This process leverages the knowledge acquired by VGG16 during its extensive training on ImageNet, allowing for effective reuse of learned features in different applications without the need for extensive retraining. This is particularly useful when dealing with limited labeled data or computational constraints.

```

✓ [24]
14s
# Load the VGG16 model
base_model = VGG16(weights='imagenet')

# Restructure the model by removing the last layer
model = Model(inputs=base_model.inputs, outputs=base_model.layers[-2].output)

# Summarize the modified model
print(model.summary())

Downloading data from https://storage.googleapis.com/tensorflow/keras-applications/vgg16/vgg16_weights_tf_dim_ordering_tf_553467096/553467096 [=====] - 8s 0us/step
Model: "model"

Layer (type)          Output Shape         Param #
=====
input_1 (InputLayer)   [(None, 224, 224, 3)]    0
block1_conv1 (Conv2D)  (None, 224, 224, 64)    1792
block1_conv2 (Conv2D)  (None, 224, 224, 64)    36928
block1_pool (MaxPooling2D) (None, 112, 112, 64)  0
block2_conv1 (Conv2D)  (None, 112, 112, 128)   73856
block2_conv2 (Conv2D)  (None, 112, 112, 128)   147584
block2_pool (MaxPooling2D) (None, 56, 56, 128)  0
block3_conv1 (Conv2D)  (None, 56, 56, 256)    295168
block3_conv2 (Conv2D)  (None, 56, 56, 256)    590080
block3_conv3 (Conv2D)  (None, 56, 56, 256)    590080
block3_pool (MaxPooling2D) (None, 28, 28, 256)  0
block4_conv1 (Conv2D)  (None, 28, 28, 512)   1180160
block4_conv2 (Conv2D)  (None, 28, 28, 512)   2359808
block4_conv3 (Conv2D)  (None, 28, 28, 512)   2359808
block4_pool (MaxPooling2D) (None, 14, 14, 512)  0
block5_conv1 (Conv2D)  (None, 14, 14, 512)   2359808
block5_conv2 (Conv2D)  (None, 14, 14, 512)   2359808
block5_conv3 (Conv2D)  (None, 14, 14, 512)   2359808
block5_pool (MaxPooling2D) (None, 7, 7, 512)  0
flatten (Flatten)     (None, 25088)        0
fc1 (Dense)           (None, 4096)        102764544
fc2 (Dense)           (None, 4096)        16781312
=====
Total params: 134260544 (512.16 MB)
Trainable params: 134260544 (512.16 MB)
Non-trainable params: 0 (0.00 Byte)
None
✓ 0s completed at 19:24

```

Feature Extraction:

The provided code is contributing to the image captioning pipeline by extracting visual features from images using the VGG16 model. The code initializes a

directory path to the folder containing images, loads each image, preprocesses it, and extracts its features using the VGG16 model. These features are then stored in a dictionary with image IDs as keys.

This step is crucial for image captioning models as it bridges the gap between image content and natural language. The extracted features serve as a rich representation of the visual content, allowing subsequent models to understand the visual context of the images. In the broader image captioning workflow, these features can be used alongside textual information to train models that generate descriptive captions for images. The VGG16 model, with its pre-trained weights on ImageNet, provides a robust feature extraction backbone for this task.

```
# Replace '/path/to/base_directory' with the actual path to your base directory
BASE_DIR = "/content/flickr8k_extracted/Images"

# Constructing the path to the "Images" directory
directory = os.path.join(BASE_DIR, "Images")

# Now, 'directory' contains the full path to the "Images" directory
print(directory)

/content/flickr8k_extracted/Images/Images

[26]: from tqdm.notebook import tqdm
      import os
      from tensorflow.keras.preprocessing.image import load_img, img_to_array
      from tensorflow.keras.applications.vgg16 import preprocess_input

# Assuming BASE_DIR is the base directory where "Images" is located
BASE_DIR = "/content/flickr8k_extracted"

# Constructing the path to the "Images" directory
directory = os.path.join(BASE_DIR, "Images")

# Create a dictionary to store features
features = {}

# Loop through images in the directory
for img_name in tqdm(os.listdir(directory), desc="Extracting Features"):
    # Load the image from file
    img_path = os.path.join(directory, img_name)
    image = load_img(img_path, target_size=(224, 224))

    # Convert image pixels to a numpy array
    image = img_to_array(image)

    # Reshape data for the model
    image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))

    # Preprocess image for VGG
    image = preprocess_input(image)

    # Extract features
    feature = model.predict(image, verbose=0)

    # Get image ID
    image_id = img_name.split(".")[0]

    # Store feature
    features[image_id] = feature
```

"Caption Mapping" or "Data Parsing":

This code snippet is part of the data preprocessing phase. It reads the content of the 'captions.txt' file, extracts image IDs and captions, removes file extensions, and creates a mapping between image IDs and corresponding captions. This step is crucial for organizing and structuring the textual information associated with each image, laying the foundation for subsequent steps in the image captioning pipeline.

```
✓ os [27] # Open 'captions.txt' file in read mode
with open(os.path.join(BASE_DIR, 'captions.txt'), 'r') as f:
    # Skip the first line (header)
    next(f)

    # Read the remaining content into captions_doc
    captions_doc = f.read()

✓ os [28] # Create a mapping of image IDs to captions
mapping = {}

# Process lines
for line in tqdm(captions_doc.split("\n")):
    # Split the line by comma(,)
    tokens = line.split(',')

    # Check if the line has enough tokens
    if len(tokens) < 2:
        continue

    # Extract image_id and caption
    image_id, caption = tokens[0], tokens[1:]

    # Remove extension from image ID
    image_id = image_id.split('.')[0]

    # Convert caption list to string
    caption = " ".join(caption)

    # Create a list if needed
    if image_id not in mapping:
        mapping[image_id] = []

    # Store the caption
    mapping[image_id].append(caption)

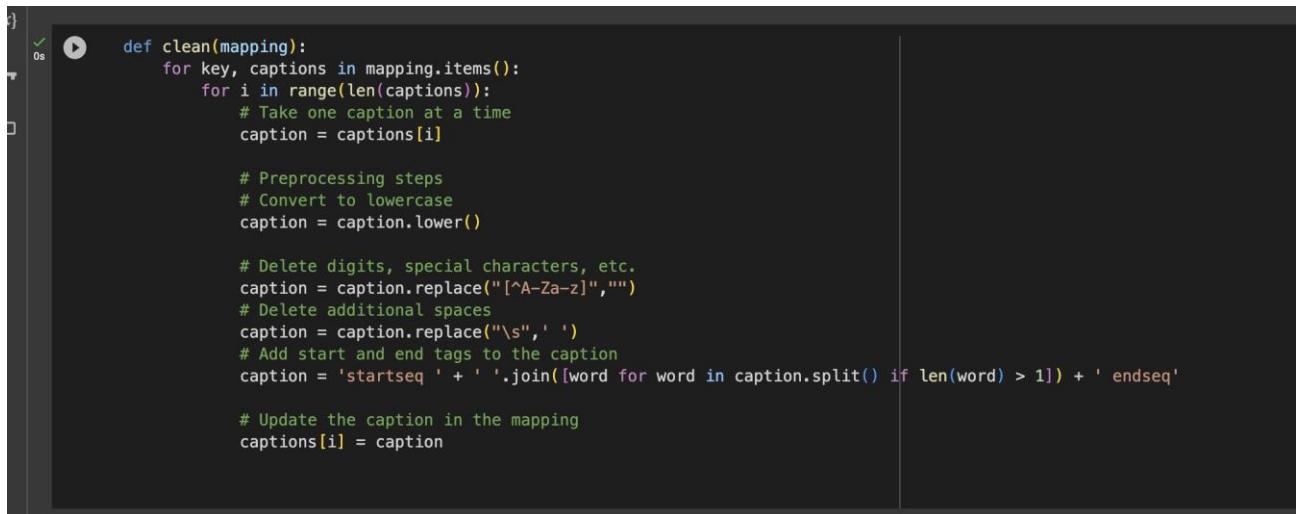
100% [██████████] 40456/40456 [00:00<00:00, 227632.23it/s]
```

```
✓ os [29] len(mapping)
8091
```

The expression `len(mapping)` returns the number of items (key-value pairs) in the mapping dictionary. In this context, it means that there are 8,091 unique image IDs in your dataset, each associated with one or more captions. This count represents the total number of images for which you have captions in your dataset.

Text Cleaning:

This code defines a function called `clean` that takes a mapping of image IDs to captions and performs several preprocessing steps on each caption. The preprocessing steps include converting the text to lowercase, removing digits and special characters, eliminating extra spaces, and adding start and end tags to the caption. These steps are commonly performed to standardize and clean textual data before training a model.



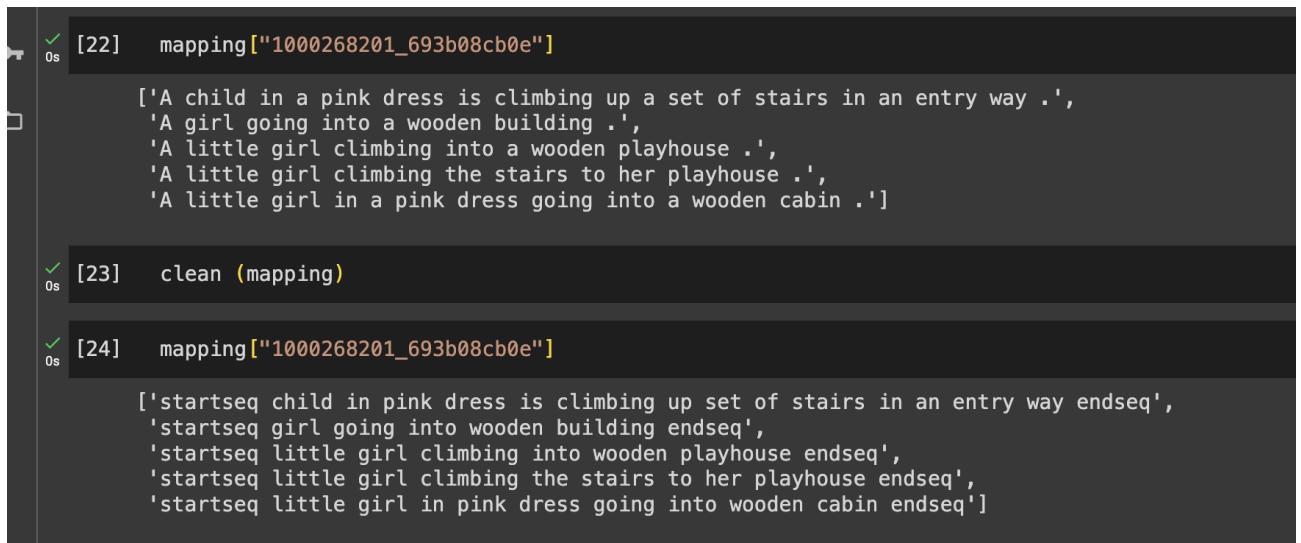
```
def clean(mapping):
    for key, captions in mapping.items():
        for i in range(len(captions)):
            # Take one caption at a time
            caption = captions[i]

            # Preprocessing steps
            # Convert to lowercase
            caption = caption.lower()

            # Delete digits, special characters, etc.
            caption = caption.replace("[^A-Za-z]", "")
            # Delete additional spaces
            caption = caption.replace("\s", ' ')
            # Add start and end tags to the caption
            caption = 'startseq ' + ' '.join([word for word in caption.split() if len(word) > 1]) + ' endseq'

            # Update the caption in the mapping
            captions[i] = caption
```

Before clean and after cleaning sentences look like this:



```
[22] mapping["1000268201_693b08cb0e"]
['A child in a pink dress is climbing up a set of stairs in an entry way .',
 'A girl going into a wooden building .',
 'A little girl climbing into a wooden playhouse .',
 'A little girl climbing the stairs to her playhouse .',
 'A little girl in a pink dress going into a wooden cabin .']

[23] clean (mapping)

[24] mapping["1000268201_693b08cb0e"]
['startseq child in pink dress is climbing up set of stairs in an entry way endseq',
 'startseq girl going into wooden building endseq',
 'startseq little girl climbing into wooden playhouse endseq',
 'startseq little girl climbing the stairs to her playhouse endseq',
 'startseq little girl in pink dress going into wooden cabin endseq']
```

Data Aggregation:

The provided code is collecting all captions from the mapping dictionary and consolidating them into a single list named **all_captions**. Each caption is associated with a specific image ID, and this step helps gather all captions together for further processing. This is often a data preparation step, and the resulting list can be used for tasks such as tokenization, vocabulary creation, or training language models.

```
allCaptions = []

# Iterate through the mapping
for key in mapping:
    # Iterate through captions for each image ID
    for caption in mapping[key]:
        # Append the caption to the list
        allCaptions.append(caption)

# Example: Print the first 5 captions
print(allCaptions[:5])

[33] len(allCaptions)
40455

[34] allCaptions[:10]
['startseq child in pink dress is climbing up set of stairs in an entry way endseq', 'startseq girl going into wooden buil
```

"Text Tokenization" or "Vocabulary Building"

This code snippet utilizes the TensorFlow Keras library to perform tokenization on a collection of captions. Tokenization involves converting textual data into numerical form by assigning a unique integer to each word in the text. In this context, the code is specifically tokenizing captions, thereby building a vocabulary and transforming the textual information into sequences of integers. This process is crucial for preparing text data to be fed into machine learning models, particularly those involved in natural language processing tasks such as image captioning.

```
from tensorflow.keras.preprocessing.text import Tokenizer

# Tokenize the text
tokenizer = Tokenizer()
tokenizer.fit_on_texts(allCaptions)
vocab_size = len(tokenizer.word_index) + 1

# Example: Print vocabulary size
print("Vocabulary Size:", vocab_size)

<>
[35] Vocabulary Size: 8485
```

The output indicates that the vocabulary size is 8485, which means there are 8485 unique words in the tokenized captions. Each word is assigned a unique integer index. This numerical representation is essential for training

machine learning models, as it allows the model to process and understand the textual information in a format that it can work with.

The output indicates that the maximum length of the captions in the dataset is 35 words. This information is crucial for data preprocessing, particularly when dealing with sequence data like captions.

```
✓ 0s [29] # Get the maximum length of the caption available
          max_length = max(len(caption.split()) for caption in allCaptions)
          print("Maximum Caption Length:", max_length)
```

```
Maximum Caption Length: 35
```

Data splitting:

Train-Test Split:

The "train-test split" is a fundamental practice in machine learning aimed at assessing the performance of a model on unseen data. This process involves dividing a dataset into two subsets: a training set, used to train the model by exposing it to labelled examples, and a testing set, held back to evaluate the model's performance on new, unseen instances. The training set enables the model to learn patterns and relationships within the data, while the testing set serves as an independent benchmark to gauge how well the model generalizes to previously unseen samples. The train-test split is essential for ensuring that the model's evaluation reflects its ability to make accurate predictions on new, real-world data, providing valuable insights into its effectiveness and potential for broader application.

```
✓ 0s [48] image_ids = list(mapping.keys())
          split = int(len(image_ids) * 0.90)
          train = image_ids[:split]
          test = image_ids[split:]
```

Sequence Encoding and Padding:

The provided code is a data generator for training a neural network model for image captioning. It takes as input a set of image keys, a mapping of image IDs to captions, image features, a tokenizer, maximum sequence length, vocabulary size, and batch size. The generator produces batches of training data containing

input sequences, partial captions, and corresponding output sequences. It iterates through the image keys, encoding the captions, splitting them into input and output pairs, padding the input sequences, and one-hot encoding the output sequences. The generator yields batches continuously for training the image captioning model. In summary, this code contributes to the data preprocessing phase of the image captioning pipeline, preparing the input data for training the neural network.

```

def data_generator(data_keys, mapping, features, tokenizer, max_length, vocab_size, batch_size):
    ...X1, X2, y = list(), list(), list()
    ...n=0
    while 1:
        for key in data_keys:
            n+=1
            captions = mapping[key]

            for caption in captions:
                # Encode the sequence
                seq = tokenizer.texts_to_sequences([caption])[0]

                # Split the sequence into x, y pairs
                for i in range(1, len(seq)):
                    # Split into input and output pairs
                    in_seq, out_seq = seq[:i], seq[i]

                    # Pad input sequence
                    in_seq = pad_sequences([in_seq], maxlen=max_length)[0]

                    # Encode output sequence
                    out_seq = to_categorical([out_seq], num_classes=vocab_size)[0]

                    # Store the sequences
                    X1.append(features[key][0]) # Assuming it's an image feature vector
                    X2.append(in_seq)
                    y.append(out_seq)
            if n == batch_size:
                X1, X2, y = np.array(X1), np.array(X2), np.array(y)
                yield [X1, X2], y # Yield a tuple of inputs
                X1, X2, y = list(), list(), list()
                n =0 # Reset the counter to 0 after yielding a batch

```

ENCODING AND DECODING:

Encoding (in the Encoder Model):

- **Purpose:** The encoder is responsible for extracting meaningful features from the input data, which, in this case, is an image.
- **Details:** The image features are fed into the encoder (**inputs1**). The layers (**fel** and **fe2**) process these features to create a condensed representation (**fe2**) that captures essential information about the input image. This condensed representation serves as a rich feature vector that summarizes the content of the image.

Decoding (in the Decoder Model):

- **Purpose:** The decoder takes the encoded information and generates a meaningful output, which, in this case, is a textual caption for the image.
- **Details:** The decoder takes two inputs—image features (**fe2**) from the encoder and sequence features (**se3**) from the input caption sequence. These features are combined (**decoder1**) and then processed through

additional layers (**decoder2** and **outputs**) to generate a probability distribution over the vocabulary. This distribution represents the likelihood of each word in the vocabulary being the next word in the caption. During training, the model learns to adjust its parameters to maximize the likelihood of generating the correct caption.

In summary, encoding is the process of extracting relevant features from input data (image), and decoding is the process of using these features to generate meaningful output (caption). The combination of encoding and decoding allows the model to learn a mapping between images and their corresponding captions, enabling it to generate captions for new, unseen images. This is a common approach in image captioning tasks, where the goal is to automatically describe the content of images in natural language

Encoder Model (inputs1):

- **Inputs:** The model takes, as input, a vector of length 4096, representing the features extracted from an image.
- **Dropout Layer:** Introduces dropout to prevent overfitting by randomly setting a fraction of input units to zero during training.
- **Dense Layer (fe2):** A fully connected layer with 256 neurons and ReLU activation. This processes the image features.

Decoder Model (inputs2):

- **Inputs:** Takes a sequence of word indices (captions) with a maximum length of 20.
- **Embedding Layer (sel):** Transforms the input sequences into dense vectors of fixed size (256 dimensions in this case). It also masks zero values.
- **Dropout Layer (se2):** Introduces dropout for regularization in the sequence processing.
- **LSTM Layer (se3):** A Long Short-Term Memory layer with 256 units. It processes the sequential information in the captions.

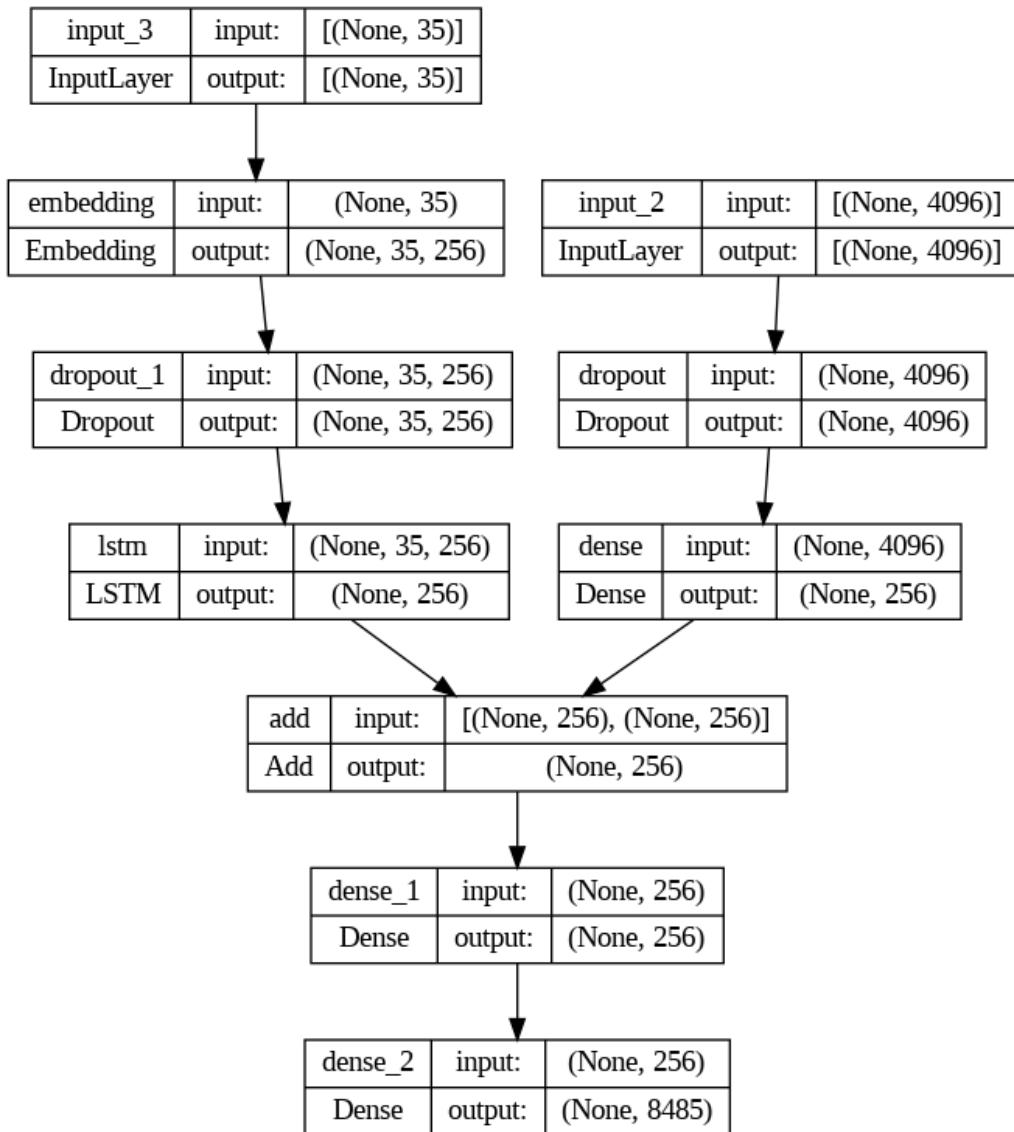
```
Q ✓ 1s
{x} # Encoder model
# Image feature layers
inputs1 = Input(shape=(4096,))
fe1 = Dropout(0.4)(inputs1)
fe2 = Dense(256, activation="relu")(fe1)

# Sequence feature layers
inputs2 = Input(shape=(max_length,))
se1 = Embedding(vocab_size, 256, mask_zero=True)(inputs2)
se2 = Dropout(0.4)(se1)
se3 = LSTM(256)(se2)

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

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

# Plot the model
plot_model(model, show_shapes=True)
```



Training phase:

The provided code snippet involves training a model for a specified number of epochs using a generator-based approach. The key elements include defining parameters such as the number of epochs, batch size, and steps per epoch. The training process iterates through each epoch, generating batches of training data using a custom data generator that combines image features, tokenized captions, and target sequences. The model is then trained for one epoch using the generated batches, and the process is repeated for the specified number of epochs. This training approach is crucial in refining the model's parameters and improving its ability to generate accurate captions for images. In succinct terms, the code represents the training phase of a machine learning model for image captioning, contributing to the model's learning and adaptation to the provided data.

Number of Epochs: The number of epochs refers to the number of times the entire dataset is passed forward and backward through the neural network during training. One epoch is completed when every sample in the training set has been processed once. Increasing the number of epochs allows the model to see the data more times, potentially improving its performance, but it's essential to monitor for overfitting.

Batch Size: Batch size represents the number of training samples utilized in one iteration. During each epoch, the dataset is divided into batches, and the model's parameters are updated based on the average gradient of the loss calculated on each batch. Smaller batch sizes require less memory but may lead to less stable updates. Larger batch sizes may provide more stable updates but require more memory.

Steps per Epoch: Steps per epoch is the number of batches processed in each epoch. It is typically set to the total number of training samples divided by the batch size. This parameter influences how often the model's parameters are updated. A larger value for steps per epoch results in more frequent updates but may extend the training time. It is crucial to choose an appropriate value to balance computational efficiency and model convergence.

```

✓ [32] #train the model
4h epochs = 20
batch_size = 32
steps = len(train) // batch_size

for i in range(epochs):
    #create data generator
    generator = data_generator(train, mapping, features, tokenizer, max_length, vocab_size, batch_size)
    #fit for one epoch
    model.fit(generator, epochs=1, steps_per_epoch = steps, verbose=1)

227/227 [=====] - 774s 3s/step - loss: 5.2214
227/227 [=====] - 770s 3s/step - loss: 4.0046
227/227 [=====] - 771s 3s/step - loss: 3.5804
227/227 [=====] - 769s 3s/step - loss: 3.3168
227/227 [=====] - 786s 3s/step - loss: 3.1177
227/227 [=====] - 774s 3s/step - loss: 2.9685
227/227 [=====] - 775s 3s/step - loss: 2.8572
227/227 [=====] - 776s 3s/step - loss: 2.7648
227/227 [=====] - 776s 3s/step - loss: 2.6827
227/227 [=====] - 781s 3s/step - loss: 2.6141
227/227 [=====] - 779s 3s/step - loss: 2.5526
227/227 [=====] - 783s 3s/step - loss: 2.4997
227/227 [=====] - 777s 3s/step - loss: 2.4502
227/227 [=====] - 777s 3s/step - loss: 2.4093
227/227 [=====] - 783s 3s/step - loss: 2.3660
227/227 [=====] - 785s 3s/step - loss: 2.3274
227/227 [=====] - 792s 3s/step - loss: 2.2883
227/227 [=====] - 782s 3s/step - loss: 2.2635
227/227 [=====] - 783s 3s/step - loss: 2.2260
227/227 [=====] - 779s 3s/step - loss: 2.1990

```

Model Saving:

The code provided is saving the trained model to a file named "IMAGE_CAPTION_best_model.h5" in the specified directory (WORKING DIR). This step is essential in machine learning workflows, especially after training a model. Saving the model allows you to reuse it later for making predictions on new data without having to retrain the model each time. It involves storing the model's architecture, weights, and other configuration details to a file, making it convenient for future use without the need for retraining.

```

Model Saving

✓ [33] # Replace "WORKING DIR" with the actual path to the directory where you want to save the model
      model.save("./content/IMAGE_CAPTION_best_model.h5")

/usr/local/lib/python3.10/dist-packages/keras/src/engine/training.py:3079: UserWarning: You are saving your model as an HDF5 file via `model.save()`. This f:
      saving_api.save_model()

```

Caption Generation.

The provided code is a crucial component of a project focused on image captioning. The function predict caption is designed to generate descriptive captions for images using a pre-trained neural network model. It takes an image, a tokenizer (used for mapping words to numerical indices), and a specified maximum length for the generated sequence as inputs. The process begins with an initial sequence starting with the tag 'startseq,' and the model iteratively predicts the next word in the sequence. The idx_to_word function assists in converting predicted indices back into their corresponding words. This sequence generation process continues until the model predicts the end tag 'endseq' or reaches the specified maximum length. Overall, this code contributes to the generation of meaningful captions for images, enhancing the interpretability and utility of machine learning models in various applications.

The screenshot shows a code editor with two snippets of Python code. The first snippet, line [34], defines a function idx_to_word that takes an integer and a tokenizer, then iterates through the tokenizer's word_index.items() to find the word corresponding to the integer, returning it or None if not found. The second snippet, line [35], defines a function predict_caption that generates a caption for an image. It starts with 'startseq', iterates over the max_length, encodes the input sentence, pads the sequence, predicts the next word, and decodes it back into a word. It stops if the word is 'endseq' or if no word is found. The code uses numpy.argmax to get the index with high probability and idx_to_word to convert the index to a word.

```
[34] def idx_to_word(integer,tokenizer):
    for word,index in tokenizer.word_index.items():
        if index == integer:
            return word
    return None

[35] # generate caption for an image
def predict_caption(model,image ,tokenizer,max_length):
    #add start tag for generation process
    in_text = 'startseq'
    #iterate over the max length of sequence
    for i in range(max_length):
        # encode input sentence
        sequence = tokenizer.texts_to_sequences([in_text])[0]
        # pad the sequence
        sequence = pad_sequences([sequence],max_length)
        # predict next word
        yhat = model.predict([image,sequence],verbose=0)
        # get index with high probability
        yhat = np.argmax(yhat)
        # convert index to word
        word = idx_to_word(yhat, tokenizer)
        #stop if wrd not found
        if word is None:
            break
        # append wrd as input for generating next word
        in_text += " " + word
        # stop if we read end tag
        if word =='endseq':
            break
    return in_text
```

CALCULATES BLEU SCORES:

Code calculates BLEU scores, a widely used metric for assessing the quality of machine-generated text, in the context of an image captioning project. The BLEU-1 score measures the precision of individual words in the generated captions, while the BLEU-2 score extends this evaluation to consider both unigrams and bigrams. The actual and predicted lists represent tokenized reference and machine-generated captions, respectively. These scores provide a quantitative measure of how well the model-generated captions align with human references. In image captioning, where the goal is to generate contextually relevant and linguistically accurate descriptions for images, BLEU scores offer valuable insights into the model's performance, aiding in the

iterative improvement of the captioning system by providing a standardized measure of caption quality.

```
CALCULATES BLEU SCORES

▶   from nltk.translate.bleu_score import corpus_bleu
    # validate with test data
    actual , predicted =list(),list()

    for key in tqdm(test):
        # get actual caption
        captions = mapping[key]
        # predict the caption for image
        y_pred = predict_caption(model, features[key],tokenizer,max_length)
        #split into words
        actual_captions = [caption.split() for caption in captions]
        y_pred = y_pred.split()
        #append to the list
        actual.append(actual_captions)
        predicted.append(y_pred)

    #calculate BLEU Score
    print("BLEU-1: %f" % corpus_bleu(actual,predicted,weights=(1.0,0,0,0)))
    print("BLEU-2: %f" % corpus_bleu(actual,predicted,weights=(0.5,0.5,0,0)))

[?] 100%|██████████| 810/810 [08:23<00:00,  1.61it/s]
BLEU-1: 0.540999
BLEU-2: 0.324626
```

The BLEU score is a metric commonly used for evaluating the quality of machine-generated text, such as translations or image captions. It ranges from 0 to 1, where a higher score indicates better performance. The BLEU score is computed based on the precision of n-grams (contiguous sequences of n items) in the generated text compared to the reference text.

In your specific output:

- "BLEU-1: 0.54" means that, on average, the precision of unigrams (individual words) in your generated captions compared to the reference captions is approximately 54%.
- "BLEU-2: 0.32" means that, on average, the precision of bigrams (two consecutive words) in your generated captions compared to the reference captions is approximately 32%.

"Caption Evaluation":

Load and Display Image:

- Takes an image name as input, constructs the full path to the image, and loads the image using the Python Imaging Library (PIL).
- Displays the loaded image using Matplotlib.

Display Actual Captions:

- Retrieves the actual captions for the given image from a mapping (**mapping**) based on the image's ID.
- Prints the actual captions to provide a reference for comparison.

Predict Caption:

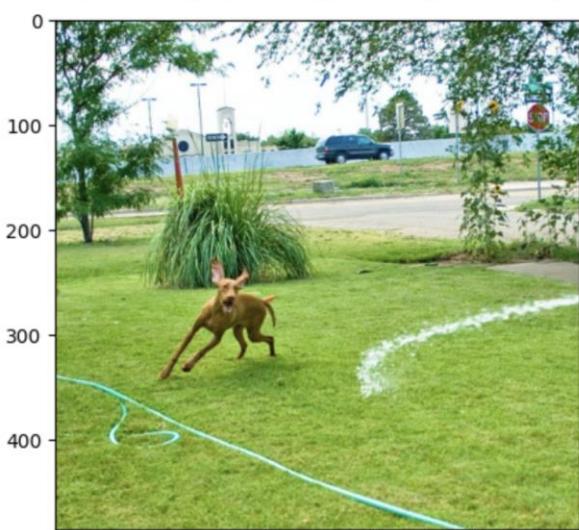
- Calls the **predict_caption** function (assumed to be defined elsewhere) to generate a caption for the given image using a pre-trained captioning model (**model**) and tokenizer (**tokenizer**).
- Prints the predicted caption.

Caption Evaluation

```
[ ]  from PIL import Image
      import matplotlib.pyplot as plt
      def generate_caption(image_name):
          # load the image
#          image_name = "1001773457_577c3a7d70.jpg"
          image_id = image_name.split('.')[0]
          img_path = os.path.join(BASE_DIR, "Images", image_name)
          image = Image.open(img_path)
          captions = mapping[image_id]
          print('-----Actual-----')
          for caption in captions:
              print(caption)
          y_pred=predict_caption(model,features[image_id],tokenizer , max_length)
          print('-----Predicted-----')
          print(y_pred)
          plt.imshow(image)
```

EXECUTING THE MODEL:

```
[ ]  generate_caption("1019077836_6fc9b15408.jpg")
-----Actual-----
startseq brown dog chases the water from sprinkler on lawn endseq
startseq brown dog plays with the hose endseq
startseq brown dog running on lawn near garden hose endseq
startseq dog is playing with hose endseq
startseq large brown dog running away from the sprinkler in the grass endseq
-----Predicted-----
startseq brown dog is running in the grass endseq
```



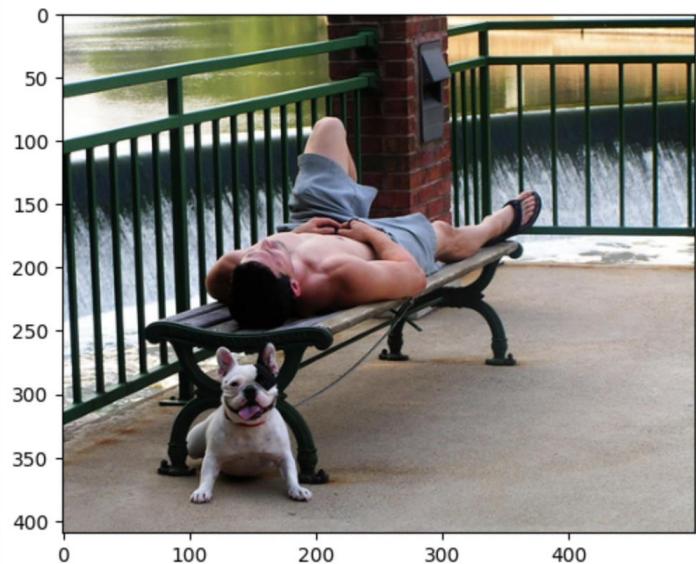
```
generate_caption("1003163366_44323f5815.jpg")
```

-----Actual-----

```
startseq man lays on bench while his dog sits by him endseq  
startseq man lays on the bench to which white dog is also tied endseq  
startseq man sleeping on bench outside with white and black dog sitting next to him endseq  
startseq shirtless man lies on park bench with his dog endseq  
startseq man laying on bench holding leash of dog sitting on ground endseq
```

-----Predicted-----

```
startseq man lays on bench with his dog endseq
```



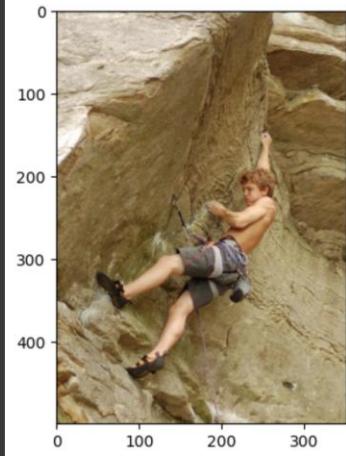
```
generate_caption("236095031_5cb17dc54a.jpg")
```

-----Actual-----

```
startseq boy rock climbs endseq  
startseq man is rock climbing without shirt endseq  
startseq young boy rock climbing endseq  
startseq the boy looked down as he climbed the steep rock face endseq  
startseq the climber is hanging off vertical cliff with gear but no shirt endseq
```

-----Predicted-----

```
startseq shirtless man climbs rock wall endseq
```



✓ 2s completed at 20:40

+ Code + Text

generate_caption("238512430_30dc12b683.jpg")

-----Actual-----

```
startseq boat is illuminated at night endseq
startseq large ship on the water at dusk endseq
startseq man is sitting on the end of hit boat overlooking the water endseq
startseq person sits on the front deck of ship and city and another ships is lit up in the distance endseq
startseq ship in harbor at night with city skyline behind endseq
```

-----Predicted-----

```
startseq man is sitting on dock at sunset endseq
```



0
50
100
150
200
250
300
0 100 200 300 400

generate_caption("96973080_783e375945.jpg")

-----Actual-----

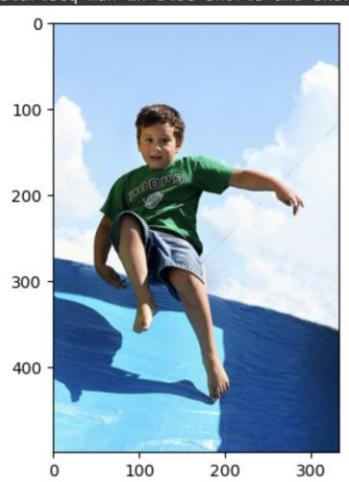
```
startseq brown dog is running over snow near leafless trees endseq
startseq brown dog runs across snowy field endseq
startseq dog is running through the snow endseq
startseq dog runs through the snow endseq
startseq medium sized brown dog is running through an open white wooded area endseq
```

-----Predicted-----

```
startseq dog runs through the snow endseq
```



0
50
100
150
200
250
300
0 100 200 300 400

```
2s  generate_caption("3712742641_641282803e.jpg")  
-----Actual-----  
startseq black and white dog is swimming in lake endseq  
startseq black and white dog is swimming in large green lake endseq  
startseq dog swimming in pond endseq  
startseq dog treads through dark green water endseq  
startseq the black and white dog is swimming in the river endseq  
-----Predicted-----  
startseq dog is in the water endseq  
  
0 50 100 150 200 250 300 350  
0 100 200 300 400  
  
generate_caption("3741827382_71e93298d0.jpg")  
-----Actual-----  
startseq boy in green shirt above something blue endseq  
startseq small boy getting tossed into the air endseq  
startseq young boy jumps around on blue mat with half-smile on his face endseq  
startseq young boy plays in front of blue object endseq  
startseq the boy wearing the green shirt is climbing on blue inflatable endseq  
-----Predicted-----  
startseq man in blue shorts and shorts is jumping into the air endseq  
  
0 100 200 300 400  
0 100 200 300 400
```

✓ 3s completed at 07:13

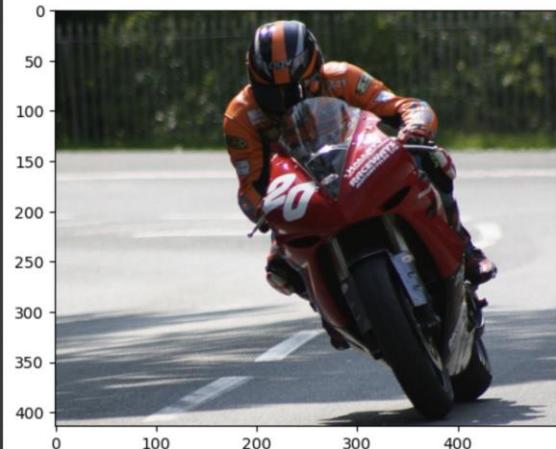
5s generate_caption("166321294_4a5e68535f.jpg")

Actual

```
startseq man is riding on red motorcycle endseq
startseq motorcycle driver dressed in orange gear swerves to the right endseq
startseq motorcyclist on red speed bike leans into sharp turn endseq
startseq motorcyclist crouches low as he rounds turn endseq
startseq this person is on red motorcycle endseq
```

Predicted

```
startseq man in red helmet is on the motorcycle endseq
```



A photograph of a motorcycle rider in an orange racing suit and helmet leaning into a turn on a red motorcycle. The motorcycle has the number '20' on its front fairing. The background shows a road and some greenery.

2s generate_caption("216172386_9ac5356dae.jpg")

Actual

```
startseq bike sits atop rise with mountains in the background endseq
startseq man wearing red uniform and helmet stands on his motorbike endseq
startseq motocross bike is being ridden over rocks endseq
startseq motocross biker about to descend endseq
startseq the motorcyclist has reached the summit endseq
```

Predicted

```
startseq man in red jacket is riding bicycle on dirt hill endseq
```



A photograph of a motocross biker in a red suit and helmet riding a dirt bike over a rocky terrain. The biker is leaning forward, and the background shows a mountainous landscape under a cloudy sky.

```

4s generate_caption("225699652_53f6fb33cd.jpg")
-----
Actual
startseq two dolphins flying headfirst into beautiful tropical blue lake endseq
startseq two dolphins jumped out of the water in this zoo endseq
startseq two dolphins jumping into the water endseq
startseq two dolphins jump out of the blue water with palm trees behind them endseq
startseq two dolphins jump out of the water together endseq
-----
Predicted
startseq two dolphins jumping into the water endseq


```

✓ 4s completed at 20:38

APPLICATION BUILDING:

CAPTION CRAFT:

Now that we trained our model, we tried to build a web application (**Caption Craft**) based on the knowledge we got from this project. This web app utilizes a Flask framework to deploy an image captioning model. Users interact with a user-friendly HTML interface, selecting an image for caption generation. The chosen image undergoes preprocessing before being fed into a pre-trained model. The model predicts a descriptive caption for the image, and the results, along with the original image, are displayed on a HTML page. This seamless flow allows users to effortlessly upload images and receive corresponding captions generated by the trained model.

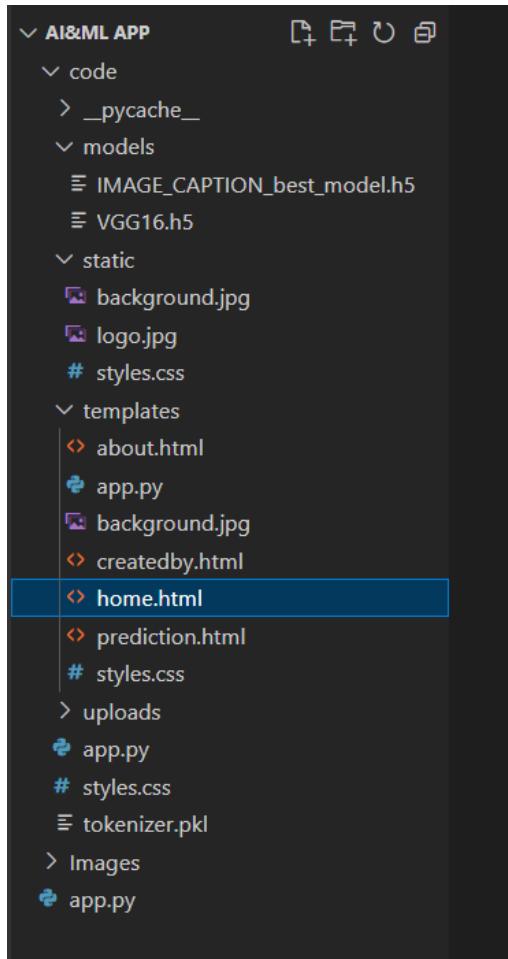
STEP 1:

In our web app we use HTML to create the front-end part of the web page.

- **home.html:** This page serves as the main entry point, providing a welcoming interface for users to navigate through the application.

- **about.html**: Users can access this page to gain insights into the project, offering details about the image captioning system and its functionality.
- **createdby.html**: This page offers information about the creators and contributors behind the development of the image captioning web app, acknowledging the individuals involved in building the application.
- These HTML pages are complemented by JavaScript (main.js) and CSS (main.css) to improve functionality and enhance the visual appeal of the web app. The combination of these components provides a comprehensive and user-friendly interface for uploading images and receiving generated captions.

- **STEP 2:**



- Flask provides a lightweight and flexible framework for building web applications, handling various aspects of web development and allowing

developers to focus on implementing application logic.

ENABLE FLASK IN COMMAND PROMPT:

```
pip install flask
cd "path\to\your\file"
set FLASK_APP=app.py
flask run
```

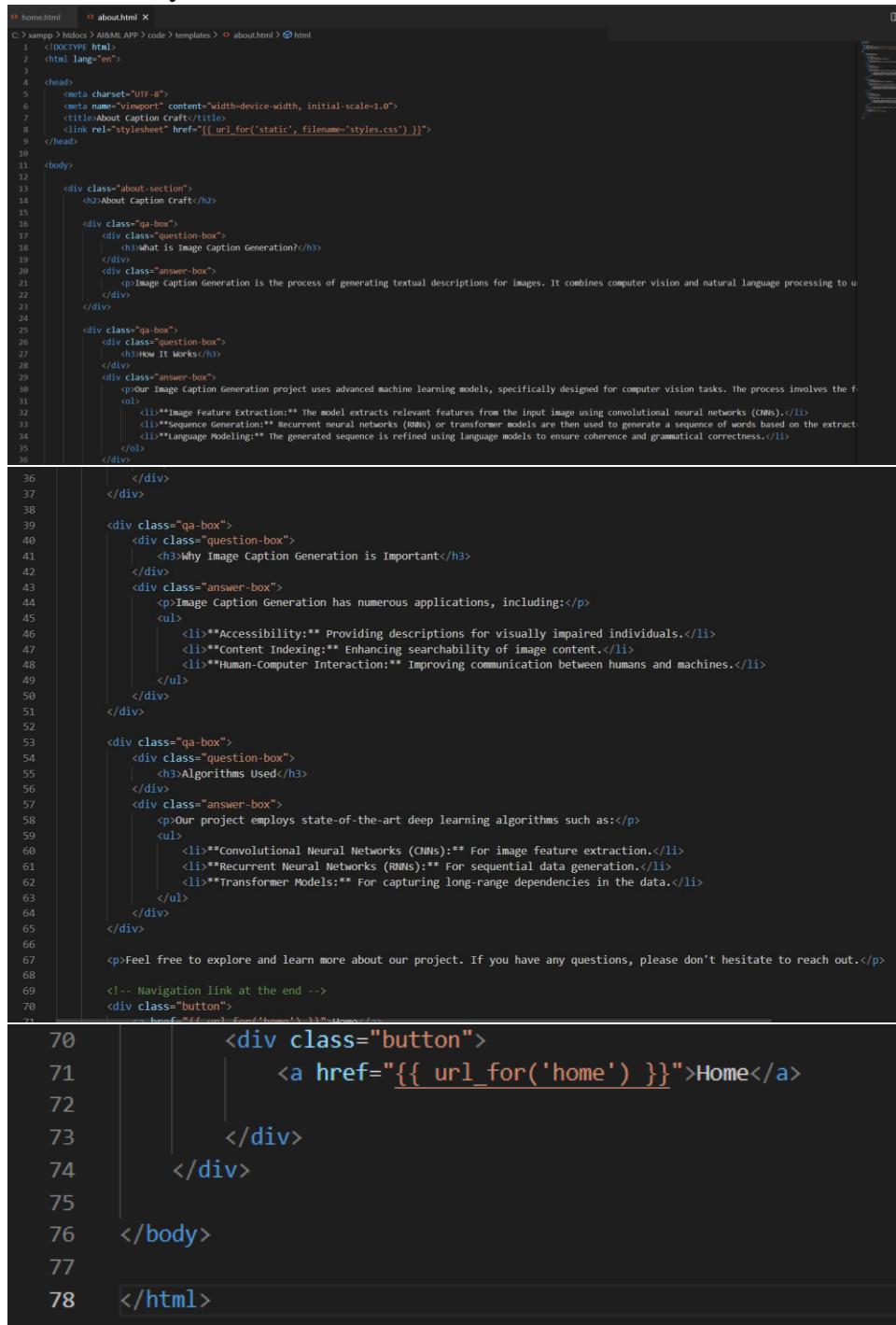
STEP 3:

Home.html: Likely contains the main content for the home page, including any introductory information, navigation links, or featured content.

```
C:\> xampp > htdocs > AI&ML APP > code > templates > home.html > html > head
1  <!DOCTYPE html>
2  <html lang="en">
3
4  <head>
5      <meta charset="UTF-8">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7      <title>Caption Craft</title>
8      <link rel="stylesheet" href="{{ url_for('static', filename='styles.css') }}">
9
10 </head>
11
12 <body>
13
14     <header>
15         
16         <h1>Caption Craft</h1>
17     </header>
18
19     <nav>
20         <a href="{{ url_for('home') }}>Home</a>
21
22         <a href="{{ url_for('about') }}>About</a>
23         <a href="{{ url_for('createdby') }}>Created By</a>
24
25     </nav>
26
27     <section class="slide" id="home">
28         <h2>Welcome to Caption Craft</h2>
29         <p>This is the home page of Caption Craft. Explore our image caption generation project and get started!</p>
30         <p>
31
32         <a href="{{ url_for('predict') }}" class="button">Get Started</a>
33
34     </p>
35     </section>
36
37     <section class="slide" id="about">
38         <h2>About Caption Craft</h2>
39         <p>Explain completely about image caption generation, how it works, algorithms used, etc.</p>
40     </section>
41
42     <section class="slide" id="createdby">
43         <h2>Created By</h2>
44         <p>Meet the developers:</p>
45         <p>BHAVANA</p>
46         <p>BHAVANI</p>
47     </section>
48
49 </body>
50
51 </html>
```

STEP 4:About.html:

presents information about the application, its purpose, or details about its functionality.



```
C:\>xamp>htdocs>AI&ML APP>code>templates>about.html <html lang="en">
1   <!DOCTYPE html>
2   <html lang="en">
3     <head>
4       <meta charset="UTF-8">
5       <meta name="viewport" content="width=device-width, initial-scale=1.0">
6       <title>About Caption Craft</title>
7       <link rel="stylesheet" href="{{ url_for('static', filename='styles.css') }}>
8     </head>
9
10    <body>
11      <div class="about-section">
12        <h2>About Caption Craft</h2>
13
14        <div class="qa-box">
15          <div class="question-box">
16            <h3>What is Image Caption Generation?</h3>
17            <div>
18              <div class="answer-box">
19                <p>Image Caption Generation is the process of generating textual descriptions for images. It combines computer vision and natural language processing to u
20              </div>
21            </div>
22          </div>
23
24        <div class="qa-box">
25          <div class="question-box">
26            <h3>How It Works</h3>
27          </div>
28          <div class="answer-box">
29            <p>Our Image Caption Generation project uses advanced machine learning models, specifically designed for computer vision tasks. The process involves the f
30            <ol>
31              <li>**Image Feature Extraction:** The model extracts relevant features from the input image using convolutional neural networks (CNNs).</li>
32              <li>**Sequence Generation:** Recurrent neural networks (RNNs) or transformer models are then used to generate a sequence of words based on the extract
33              <li>**Language Modeling:** The generated sequence is refined using language models to ensure coherence and grammatical correctness.</li>
34            </ol>
35          </div>
36
37        </div>
38
39        <div class="qa-box">
40          <div class="question-box">
41            <h3>Why Image Caption Generation is Important</h3>
42          </div>
43          <div class="answer-box">
44            <p>Image Caption Generation has numerous applications, including:</p>
45            <ul>
46              <li>**Accessibility:** Providing descriptions for visually impaired individuals.</li>
47              <li>**Content Indexing:** Enhancing searchability of image content.</li>
48              <li>**Human-Computer Interaction:** Improving communication between humans and machines.</li>
49            </ul>
50          </div>
51
52
53        <div class="qa-box">
54          <div class="question-box">
55            <h3>Algorithms Used</h3>
56          </div>
57          <div class="answer-box">
58            <p>Our project employs state-of-the-art deep learning algorithms such as:</p>
59            <ul>
60              <li>**Convolutional Neural Networks (CNNs):** For image feature extraction.</li>
61              <li>**Recurrent Neural Networks (RNNs):** For sequential data generation.</li>
62              <li>**Transformer Models:** For capturing long-range dependencies in the data.</li>
63            </ul>
64          </div>
65
66
67        <p>Feel free to explore and learn more about our project. If you have any questions, please don't hesitate to reach out.</p>
68
69
70        <!-- Navigation link at the end -->
71        <div class="button">
72          <div>
73            <a href="{{ url_for('home') }}>Home</a>
74          </div>
75        </div>
76
77
78      </body>
79
80    </html>
```

STEP 5:

Createdby.html

includes details about the creator of the application, acknowledging contributions or providing contact information.

```

home.html      about.html      createdby.html ×
C: > xampp > htdocs > AI&ML APP > code > templates > createdby.html > html

1  <!DOCTYPE html>
2  <html lang="en">
3
4  <head>
5      <meta charset="UTF-8">
6      <meta name="viewport" content="width=device-width, initial-scale=1.0">
7      <title>Created By</title>
8      <link rel="stylesheet" href="{{ url_for('static', filename='styles.css') }}>
9
10 </head>
11
12 <body>
13
14     <div class="createdby-section">
15         <h2>Developed By</h2>
16
17         <div class="developer-info" id="developer1">
18             <div class="developer-names">
19                 <span>BHAVANA</span>
20             </div>
21
22             <div class="additional-info">
23                 <p>Email: bhavanaperecharla@gmail.com</p>
24                 <p>Date of Birth: DECEMBER 24, 2003</p>
25                 <p>Team ID: 591719</p>
26                 <p>Student ID: 21BCE9184</p>
27                 <p>College Name: VIT AP College</p>
28             </div>
29         </div>
30
31         <div class="developer-info" id="developer2">
32             <div class="developer-names">
33                 <span>BHAVANI</span>
34             </div>
35
36             <div class="additional-info">
37                 <p>Email: bhavaniponnnapali@gmail.com</p>
38
39                 <p>Date of Birth: DECEMBER 23, 2003</p>
40
41                 <p>Team ID: 591719</p>
42                 <p>Student ID: 21BCE9043</p>
43                 <p>College Name: VIT AP College</p>
44             </div>
45         </div>
46
47         <!-- Navigation link at the end -->
48         <div class="button">
49             <a href="{{ url_for('home') }}>Home</a>
50         </div>
51     </div>
52
53 </html>

```

STEP 6:

app.py

app.py is the main script of the Flask application. It defines the core

functionality, including routes for different URLs (home, about, createdby, predict), file upload handling, and integration with a machine learning model. The script orchestrates the application, allowing it to handle incoming requests, process data, and dynamically render HTML templates, creating a functional web-based image captioning application.

```
app.py  ●  home.html  about.html  createdby.html
C: > xampp > htdocs > AI&ML APP > code > templates > app.py
 1  # Import necessary libraries
 2  from flask import Flask, render_template, request, jsonify
 3  from werkzeug.utils import secure_filename
 4  from tensorflow.keras.models import load_model
 5  from tensorflow.keras.preprocessing import image
 6  from tensorflow.keras.applications.inception_v3 import preprocess_input
 7  from tensorflow.keras.preprocessing.sequence import pad_sequences
 8  from flask_cors import CORS
 9  import pickle
10  import numpy as np
11  import os
12  import base64
13
14  # Create Flask app
15  app = Flask(__name__, template_folder='templates', static_folder='static')
16  CORS(app)
17
18  # Load your image captioning model
19  model_path = 'C:\\xampp\\htdocs\\AI&ML APP\\code\\models\\IMAGE_CAPTION_best_model.h5'
20  model = load_model(model_path, compile=False)
21
22  tokenizer_path = 'C:\\xampp\\htdocs\\AI&ML APP\\code\\tokenizer.pkl' # Specify the ac
23
24  with open(tokenizer_path, 'rb') as tokenizer_file:
25      tokenizer = pickle.load(tokenizer_file)
26
27  max_length = 35
28
29  # Function to preprocess image for model input
30  def preprocess_image(img_path):
31      img = image.load_img(img_path, target_size=(229, 229))
32      img_array = image.img_to_array(img)
33      img_array = np.expand_dims(img_array, axis=0)
34      img_array = preprocess_input(img_array)
35
36      return img_array
```

```
35     return img_array
36
37 # Function to generate caption from image features
38 def generate_caption_from_features(img_array):
39     # Assuming your model takes image features as input
40     caption = predict_caption(model, img_array, tokenizer, max_length)
41     return caption
42
43 # Function to predict caption
44 def predict_caption(model, image, tokenizer, max_length):
45     in_text = 'startseq'
46     for _ in range(max_length):
47         sequence = tokenizer.texts_to_sequences([in_text])[0]
48         sequence = pad_sequences([sequence], maxlen=max_length)
49         yhat = model.predict([image, sequence], verbose=0)
50         yhat = np.argmax(yhat)
51         word = idx_to_word(yhat, tokenizer)
52         if word is None:
53             break
54         in_text += " " + word
55         if word == 'endseq':
56             break
57     return in_text
58
59 # Helper function to convert index to word
60 def idx_to_word(integer, tokenizer):
61     for word, index in tokenizer.word_index.items():
62         if index == integer:
63             return word
64     return None
65
66 # Route for home page
67 @app.route('/')
68 def home():
69     return render_template('home.html')
```

```
69 |     return render_template('home.html')
70 |
71 | @app.route('/createdby')
72 | def createdby():
73 |     return render_template('createdby.html')
74 |
75 | @app.route('/about')
76 | def about():
77 |     return render_template('about.html')
78 |
79 | UPLOAD_FOLDER = 'uploads'
80 | ALLOWED_EXTENSIONS = {'png', 'jpg', 'jpeg', 'gif'}
81 |
82 | app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER
83 |
84 | def allowed_file(filename):
85 |     return '.' in filename and filename.rsplit('.', 1)[1].lower() in
86 |
87 |
88 | @app.route('/predict', methods=["GET", "POST"])
89 | def predict():
90 |     if request.method == "POST":
91 |         try:
92 |             # Get the uploaded image file
93 |             uploaded_file = request.files['image']
94 |
95 |             # Check if the file is allowed
96 |             if uploaded_file and allowed_file(uploaded_file.filename):
97 |                 # Save the uploaded image to a temporary location
98 |                 filename = secure_filename(uploaded_file.filename)
99 |                 temp_path = os.path.join(app.config['UPLOAD_FOLDER'],
100 |                                         filename)
101 |                 uploaded_file.save(temp_path)
102 |
103 |                 # Preprocess the image for model input
104 |                 img_array = preprocess_image(temp_path)
105 |                 """ Call the prediction function here """
106 |         except Exception as e:
107 |             print(e)
```

```
69 |     return render_template('home.html')
70 |
71 | @app.route('/createdby')
72 | def createdby():
73 |     return render_template('createdby.html')
74 |
75 | @app.route('/about')
76 | def about():
77 |     return render_template('about.html')
78 |
79 | UPLOAD_FOLDER = 'uploads'
80 | ALLOWED_EXTENSIONS = {'png', 'jpg', 'jpeg', 'gif'}
81 |
82 | app.config['UPLOAD_FOLDER'] = UPLOAD_FOLDER
83 |
84 | def allowed_file(filename):
85 |     return '.' in filename and filename.rsplit('.', 1)[1].lower() in
86 |
87 |
88 | @app.route('/predict', methods=["GET", "POST"])
89 | def predict():
90 |     if request.method == "POST":
91 |         try:
92 |             # Get the uploaded image file
93 |             uploaded_file = request.files['image']
94 |
95 |             # Check if the file is allowed
96 |             if uploaded_file and allowed_file(uploaded_file.filename):
97 |                 # Save the uploaded image to a temporary location
98 |                 filename = secure_filename(uploaded_file.filename)
99 |                 temp_path = os.path.join(app.config['UPLOAD_FOLDER'],
100 |                                         filename)
101 |                 uploaded_file.save(temp_path)
102 |
103 |                 # Preprocess the image for model input
104 |                 img_array = preprocess_image(temp_path)
105 |
106 |                 # Generate caption for the image
```

STEP 7: **prediction.html:**

prediction.html is an HTML template used by the Flask application to dynamically generate the prediction page. It likely includes placeholders for displaying the predicted image, caption, and any additional information. This template is rendered by Flask in response to the **/predict** route, providing a user-friendly interface to view the results of the image captioning process.

```

C: > xampp > htdocs > AI&ML APP > code > templates > prediction.html > ...
1   <!DOCTYPE html>
2   <html lang="en">
3
4   <head>
5       <meta charset="UTF-8">
6       <meta name="viewport" content="width=device-width, initial-scale=1.0">
7       <title>Welcome to the Prediction Page</title>
8       <link rel="stylesheet" href="{{ url_for('static', filename='styles.css') }}">
9   </head>
10
11  <body>
12
13      <div class="prediction-section">
14          <h2>WELCOME TO THE PREDICTION PAGE</h2>
15
16          <p>UPLOAD THE PICTURE TO GET CAPTION</p>
17
18          <form id="predictionForm" enctype="multipart/form-data" action="{{ url_for('predict') }}" method="post">
19              <label for="imageUpload">Upload an Image:</label>
20              <input type="file" id="imageUpload" name="image" accept="image/*" required>
21              <button type="submit">Predict</button>
22          </form>
23
24      </div>
25
26      <h2>Predicted Image</h2>
27
28          {% if data.image %}
29              
30          {% else %}
31              <p>No image to display</p>
32
33          {% endif %}
34
35          <h2>Predicted Caption:</h2>
36              <p>{{ data.caption }}</p>
37
38      </body>
39
40          <!-- Navigation link at the top -->
41          <div class="button">
42              <a href="{{ url_for('home') }}>Home</a>
43          </div>
44
45
46  </body>
47  </html>

```

STEP 8: styles.css

(CSS) file used to define the visual presentation and styling of HTML elements in the web application. It likely contains rules for styling various components such as fonts, colors, layout, and responsiveness

```
C: > xampp > htdocs > AI&ML APP > code > templates > # styles.css > ↵ body
1  body {
2      font-family: 'Raleway', sans-serif;
3      margin: 0;
4      padding: 0;
5      background: url('background.jpg') no-repeat center center fixed;
6      background-size: cover;
7      color: ■#fff;
8      line-height: 1.5; /* Adjust line height as needed */
9  }
10
11 header {
12     text-align: center;
13     padding: 1em 0;
14 }
15
16 .logo-container img {
17     max-width: 100%;
18     height: auto;
19     display: inline-block; /* To prevent extra space below the image */
20 }
21
22
23 nav {
24     text-align: center;
25     margin-top: 20px;
26 }
27
28 nav a {
29     text-decoration: none;
30     color: ■#fff;
31     padding: 10px 20px;
32     margin: 0 10px;
33     background-color: ■#4285f4;
34     border: 1px solid ■#4285f4;
35     border-radius: 5px;
36     font-family: 'Lato', sans-serif;
```

```
36     font-family: 'Lato', sans-serif;
37 }
38
39 .slide {
40     padding: 50px 0;
41     text-align: center;
42     border: 2px solid #4285f4;
43     border-radius: 10px;
44     margin: 20px;
45     box-shadow: 0 0 10px rgba(0, 0, 0, 0.1);
46     color: #fff;
47 }
48
49 .button {
50     display: inline-block;
51     padding: 10px 20px;
52     font-size: 16px;
53     text-align: center;
54     text-decoration: none;
55     background-color: #4285f4;
56     color: #fff;
57     border-radius: 5px;
58     cursor: pointer;
59     font-family: 'Montserrat', sans-serif;
60 }
61
62 .about-section {
63     padding: 50px 20px;
64     text-align: left;
65     border: 2px solid #4285f4;
66     border-radius: 10px;
67     margin: 20px;
```

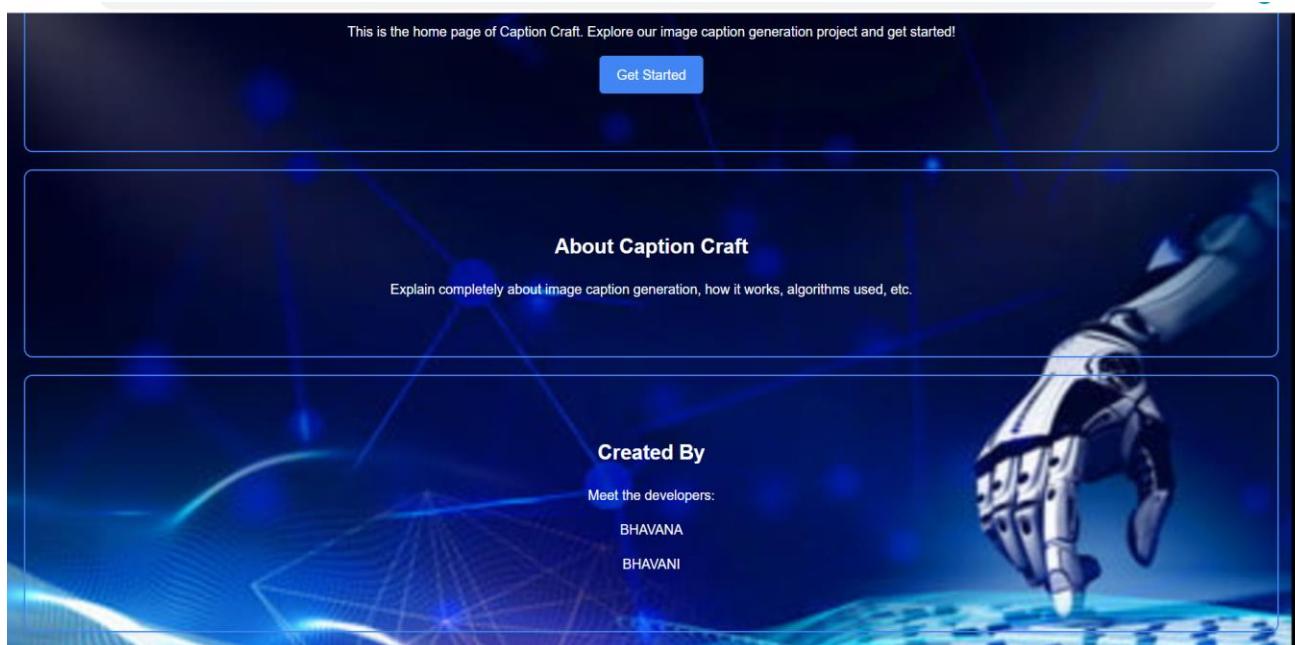
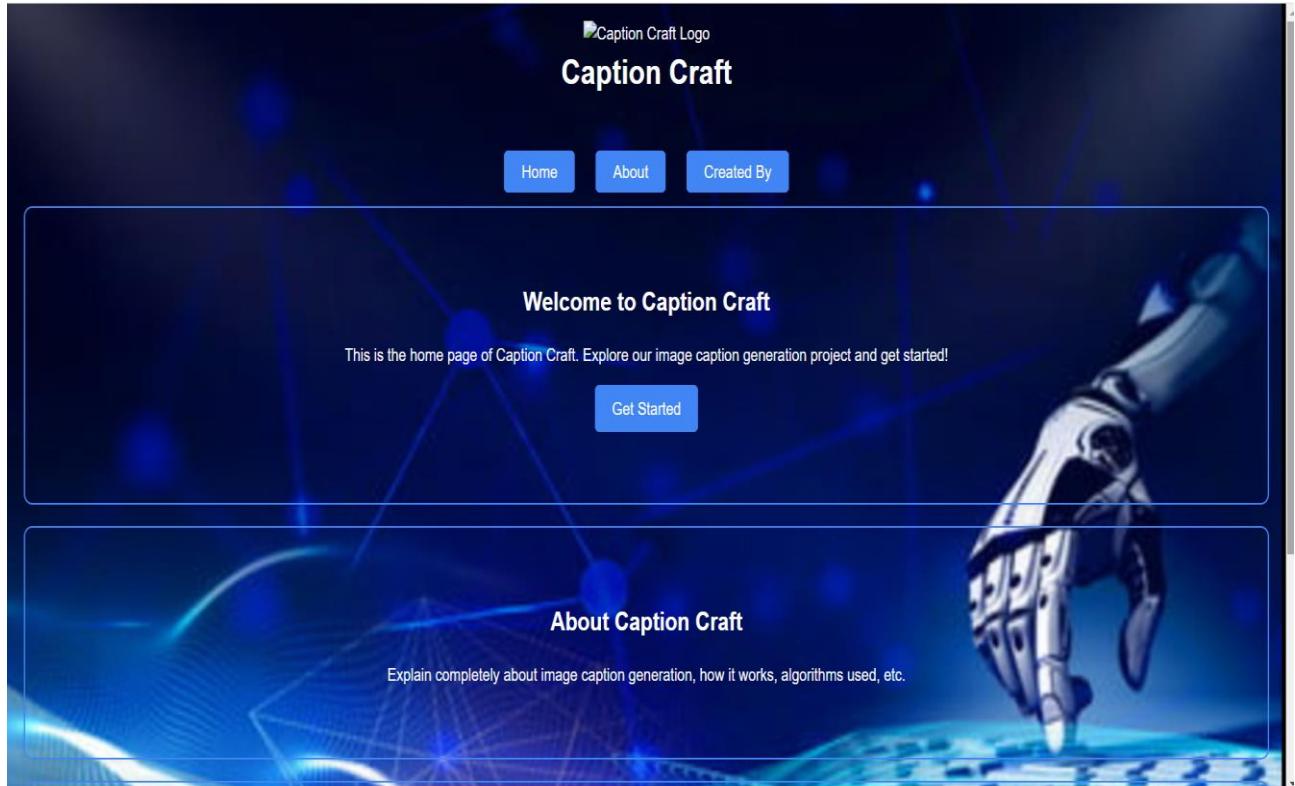
```
67     margin: 20px;
68     box-shadow: 0 0 10px rgba(0, 0, 0, 0.1);
69 }
70
71 .about-section h2 {
72     color: #fff;
73 }
74
75 .qa-box {
76     display: flex;
77     flex-direction: column;
78     margin-bottom: 20px;
79 }
80
81 .question-box,
82 .answer-box {
83     background-color: #fff;
84     padding: 10px;
85     border-radius: 5px;
86     margin-bottom: 10px;
87 }
88
89 .question-box h3 {
90     color: #000;
91 }
92
93 .answer-box p,
94 .answer-box ol,
95 .answer-box ul {
96     color: #000;
97 }
98
99 .createdby-section {
```

```
106  }
107
108 .createdby-section h2 {
109   color: #4285f4;
110   font-size: 48px;
111   font-family: 'Arial', sans-serif;
112   border-bottom: 2px solid #4285f4;
113   padding-bottom: 10px;
114 }
115
116 .developer-info {
117   margin-top: 30px;
118 }
119
120 .developer-names {
121   display: flex;
122   justify-content: center;
123   align-items: center;
124 }
125
126 .developer-names span {
127   color: #fff;
128   font-size: 36px;
129   font-family: 'Verdana', sans-serif;
130   margin: 0 20px;
131   border: 2px solid #4285f4;
132   padding: 10px;
133   border-radius: 10px;
134 }
135
136 .additional-info {
137   margin-top: 20px;
138   color: #fff;
139 }
140
141 .additional-info p {
```

```
140
141 .additional-info p {
142   margin: 10px 0;
143 }
144
145 .prediction-section {
146   padding: 100px 20px;
147   text-align: center;
148   border: 2px solid #4285f4;
149   border-radius: 10px;
150   margin: 20px;
151   box-shadow: 0 0 10px rgba(0, 0, 0, 0.1);
152 }
153
154 .prediction-section h2 {
155   color: #4285f4;
156   font-size: 48px;
157   font-family: 'Arial', sans-serif;
158   border-bottom: 2px solid #4285f4;
159   padding-bottom: 10px;
160 }
161
162 .prediction-section p {
163   margin-top: 20px;
164 }
165
166 #predictionForm {
167   margin-top: 20px;
168   display: flex;
169   flex-direction: column;
170   align-items: center;
171 }
172
173 label {
174   color: #fff;
175   font-size: 18px;
176   margin-bottom: 20px;
```

```
170     |   align-items: center;
171   }
172
173   label {
174     |   color: ■#fff;
175     |   font-size: 18px;
176     |   margin-bottom: 20px;
177   }
178
179   #imageUpload {
180     |   margin-top: 10px;
181     |   margin-bottom: 20px;
182   }
183
184   button {
185     |   background-color: ■#4285f4;
186     |   color: ■#fff;
187     |   border: 1px solid ■#4285f4;
188     |   padding: 10px 20px
189 }
```

STEP 9: HOME PAGE :



STEP :10 ABOUT PAGE (INFORMATION REGARDING IMAGE CAPTION GENERATION:)

About Caption Craft

What is Image Caption Generation?

Image Caption Generation is the process of generating textual descriptions for images. It combines computer vision and natural language processing to understand the content of an image and describe it in human-readable language.

How It Works

Our Image Caption Generation project uses advanced machine learning models, specifically designed for computer vision tasks. The process involves the following steps:

1. **Image Feature Extraction:** The model extracts relevant features from the input image using convolutional neural networks (CNNs).
2. **Sequence Generation:** Recurrent neural networks (RNNs) or transformer models are then used to generate a sequence of words based on the extracted features.
3. **Language Modeling:** The generated sequence is refined using language models to ensure coherence and grammatical correctness.

Why Image Caption Generation is Important

Image Caption Generation has numerous applications, including:

- **Accessibility:** Providing descriptions for visually impaired individuals.
- **Content Indexing:** Enhancing searchability of image content.
- **Human-Computer Interaction:** Improving communication between humans and machines.

Algorithms Used

Our project employs state-of-the-art deep learning algorithms such as:

- **Convolutional Neural Networks (CNNs):** For image feature extraction.
- **Recurrent Neural Networks (RNNs):** For sequential data generation.
- **Transformer Models:** For capturing long-range dependencies in the data.

Feel free to explore and learn more about our project. If you have any questions, please don't hesitate to reach out.

[Home](#)

STEP :11

DEVELOPED BY

Developed By

BHAVANA

Email: bhavanaperecharla@gmail.com

Date of Birth: DECEMBER 24, 2003

Team ID: 591719

Student ID: 21BCE9184

College Name: VIT AP College

BHAVANI

Email: bhavaniponnnapali@gmail.com

Date of Birth: DECEMBER 23, 2003

Team ID: 591719

Student ID: 21BCE9043

BHAVANA

Email: bhavanaperecharla@gmail.com

Date of Birth: DECEMBER 24, 2003

Team ID: 591719

Student ID: 21BCE9184

College Name: VIT AP College

BHAVANI

Email: bhavaniponnnapali@gmail.com

Date of Birth: DECEMBER 23, 2003

Team ID: 591719

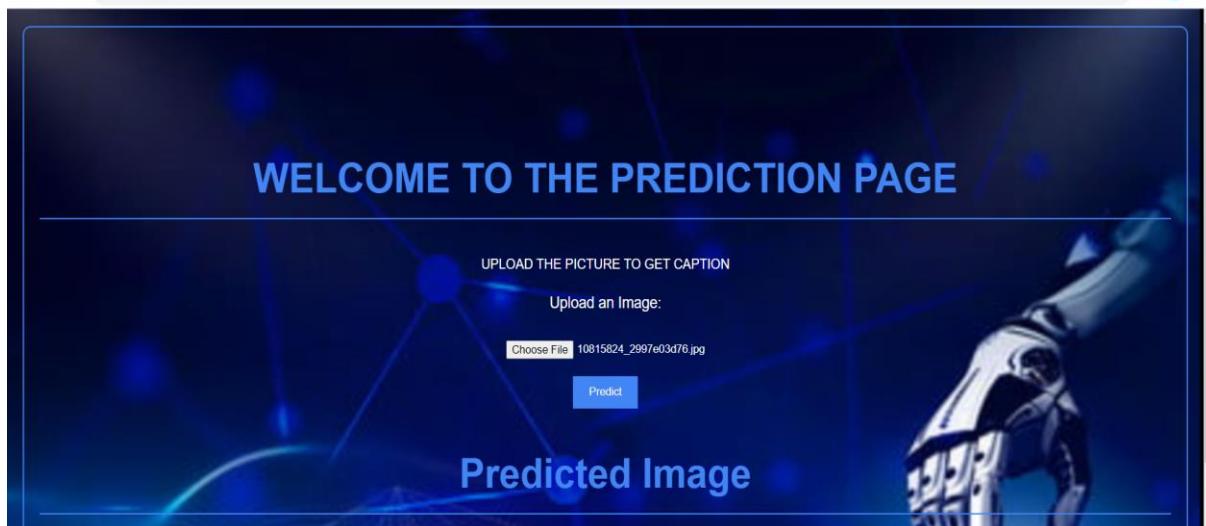
Student ID: 21BCE9043

College Name: VIT AP College

[Home](#)

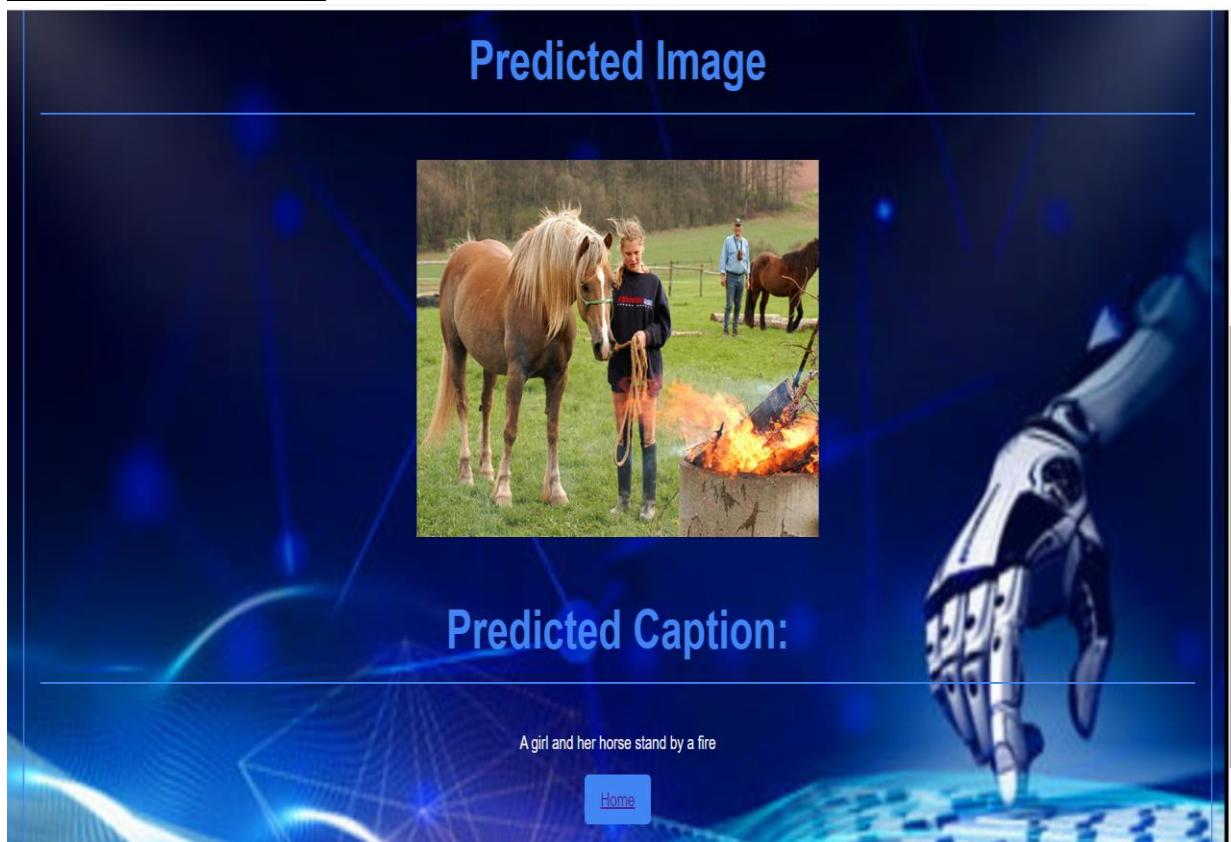
STEP :12

PREDICTION PAGE



Predicted Image

PREDICTION 1:



Predicted Caption:

A girl and her horse stand by a fire

[Home](#)

PREDICTION 2:

UPLOAD THE PICTURE TO GET CAPTION

Upload an image:

Choose File 934375844_...49fed18.jpg

Predict

Predicted Image



A small baby sits on a bed and smiles

[Home](#)

Predicted Caption:

PREDICTION 3:

Predicted Image



Predicted Caption:

A child in a pink dress is climbing up a set of stairs in an entryway.

[Home](#)

PREDICTION 4:



Predicted Caption:

A man is surfing on board

[Home](#)

PREDICTION 5:

Predicted Image



Predicted Caption:

A dog under a sheet

[Home](#)

❖ THANK YOU

