

Image Caption Generation

Project report

Project ID - 591848

Team Members –

Lavanya Gudala

Jyoshil Sree Jayanth Chandana

Chandpasha Alamuri

1. INTRODUCTION

1.1 Project Overview

Image caption generation, a compelling application of deep learning, involves the development of models capable of providing descriptive and meaningful textual descriptions for input images. Image caption generation is a task in computer vision and natural language processing where the goal is to generate a textual description (caption) for an input image. This involves understanding the content of the image and expressing it in natural language. Convolutional neural networks (CNNs) and long short-term memory networks (LSTMs) are effective deep learning paradigms that combine the skills of LSTMs in sequence generation and CNNs in image feature extraction to generate image captions. This model is an example of an advanced method for giving input photos textual descriptions. Because CNNs are good at extracting hierarchical information from images, they may identify complex patterns and spatial relationships. These networks function as a basis for comprehending visual content by extracting pertinent elements from the input image. LSTMs, which are excellent at processing sequential input and maintaining long-range dependencies, are then fed the extracted features. In this project, we will see how deep neural network models can be used to automatically generate descriptions for images, such as photographs. The model can understand the image's content and context thanks to the synergy between CNNs and LSTMs, resulting in insightful descriptions. The CNN-LSTM combination is a state-of-the-art deep learning model that uses neural networks' power to enable automated picture captioning. This greatly expands AI's capacity to comprehend and describe visual content.

1.2 Purpose

The purpose of an image caption generation project, using a combination of Convolutional Neural Networks (CNNs) and Long Short-Term Memory (LSTM) networks, is to create a system that can automatically generate descriptive captions for images. contribute to the broader goal of creating AI systems that can understand and interact with the world in a more human-like manner. They leverage the capabilities of deep learning to extract meaningful information from images and generate coherent and contextually relevant textual descriptions, making visual content more accessible and interpretable by machines. there are many useful cases which these are useful like Content Indexing, Assistance for Autonomous Systems, Research and Development and Social Media Enhancement.

2. LITERATURE SURVEY

2.1 Existing problem

Research in image captioning has evolved from traditional methods to contemporary deep-learning approaches. Early methods often relied on handcrafted features and rule-based systems, while recent approaches leverage the power of deep neural networks. The integration of CNNs in image processing has significantly influenced the field of computer vision. CNNs excel in feature extraction, capturing hierarchical representations of images. The works of [cite key studies] showcase the effectiveness of CNNs in understanding visual content, a crucial aspect for subsequent caption generation. Despite the progress made, challenges persist in image captioning. Ambiguity in scenes, diversity in content, and the need for fine-grained understanding pose significant hurdles. Researchers have addressed these challenges through techniques like attention mechanisms [cite works], yet achieving a holistic solution remains an ongoing pursuit. The evaluation of image captioning models involves metrics like BLEU, METEOR, and CIDEr. However, the limitations of these metrics, as discussed in [cite relevant studies], underscore the complexity of objectively assessing the quality of generated captions.

2.2 References

Yue Zhang,¹ and Xiaosheng Yu have summarized the related methods and focused on the attention mechanism, which plays an important role in computer vision and is recently widely used in image caption generation tasks. Songtao Ding a, Shiru Qu proposed a novel image captioning model based on high-level image features. We combine low-level information, such as image quality, with high-level features, such as motion classification and face recognition to detect attention regions of an image. Xinlei Chen, C. Lawrence Zitnick explored the bi-directional mapping between images and their sentence-based descriptions. Critical to our approach is a recurrent neural network that attempts to dynamically build a visual representation of the scene as a caption is being generated or read. X.Jia, E.Gavves, B.Fernando, and T.Tuytelaars added semantic information extracted from the image as extra input to each unit of the LSTM block, to guide the model towards solutions that are more tightly coupled to the image content. Sulabh Katiyar, and Samir Kumar Borgohain have observed that the model complexity of the Convolutional Neural Network, as measured by several parameters, and the accuracy of the model on the Object Recognition task does not necessarily correlate with its efficacy on feature extraction for Image Caption Generation task.

2.3 Problem Statement Definition

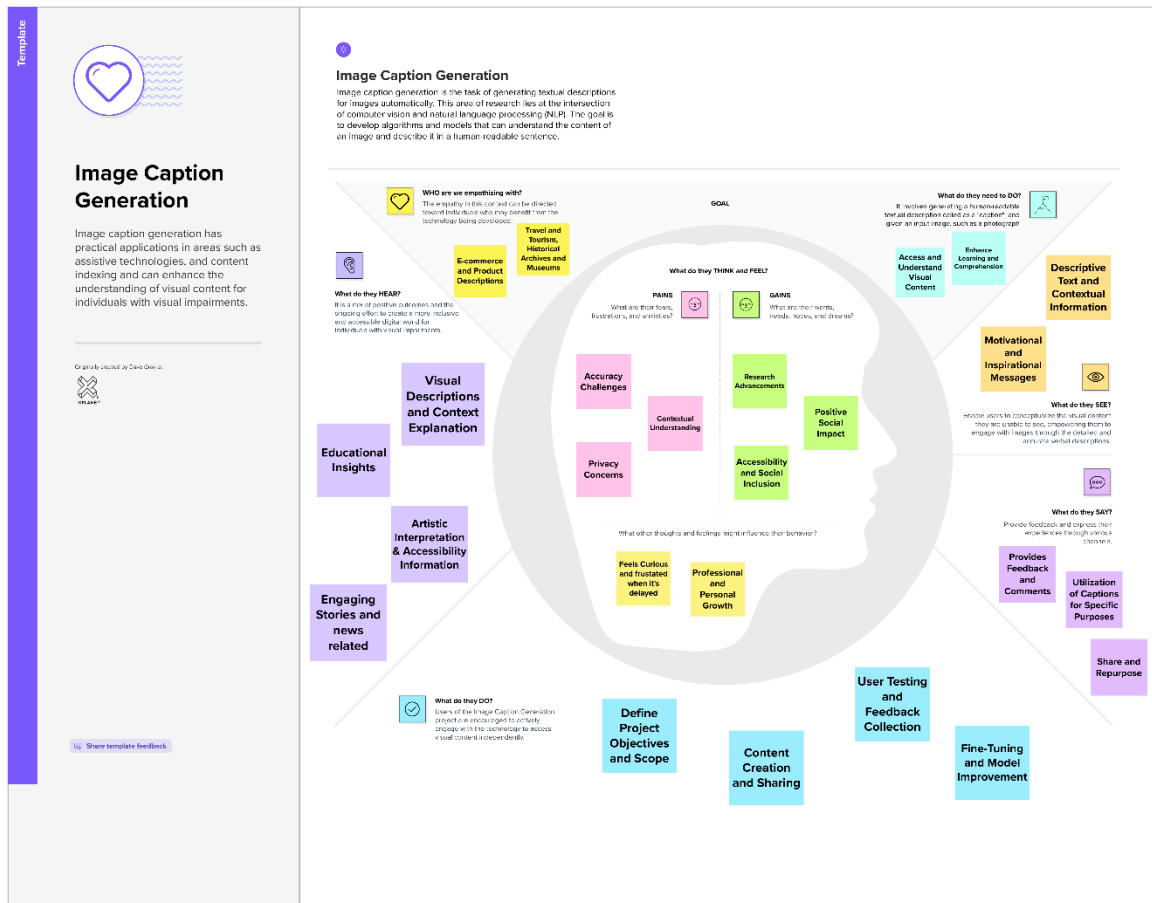
This project endeavours to tackle the persistent challenges in image caption generation, aiming to enhance the quality and contextual relevance of generated descriptions. Current models often fall short in comprehensive scene understanding, leading to captions that may miss intricate details and fail to grasp the overall narrative of complex images. Ambiguity handling remains a critical issue, as existing systems struggle with scenarios involving overlapping objects or unclear spatial arrangements, generating ambiguous or vague descriptions. Additionally, the limited adaptability of image captioning models to diverse domains and the prevalence of redundant captions pose significant hurdles. This project seeks to address these challenges by developing an advanced image captioning system. By leveraging the combined strength of Convolutional Neural Networks (CNNs) for robust feature extraction and Long Short-Term

Memory networks (LSTMs) for sequence generation, the project aims to produce more accurate, diverse, and contextually relevant captions, contributing to the refinement of image captioning technology.

3. IDEATION & PROPOSED SOLUTION

3.1 Empathy Map Canvas

An Empathy Map Canvas is a visual tool designed to deepen the understanding of users' experiences and perspectives in the context of a particular product or service. Divided into sections like "Says," "Thinks," "Feels," "Does," "Pains," and "Gains," the canvas prompts teams to consider various aspects of the user journey



3.2 Ideation & Brainstorming

Ideation and brainstorming are crucial stages in the creative process, where teams generate and explore a multitude of ideas to solve a specific problem or address a particular challenge. Brainstorming sessions often involve rapid idea generation, allowing for a wide range of concepts to be shared. The goal is to encourage collaboration and unlock innovative solutions by leveraging the collective creativity of the team

Brainstorm - Idea Listing and Grouping

2
Brainstorm

Write down any ideas that come to mind that address your problem statement.

10 minutes

TIP
You can select a sticky note and hit the pencil (switch to sketch) icon to start drawing.

JYOSHIL

Develop a bot that automatically generates captions and relevant questions for images to make them more engaging on social media platforms like Instagram and Twitter	Create a system that can extract insights by generating concise and actionable guiding descriptions for the images used in their blog posts, improving SEO and user engagement	Create a coaching app that generates personalized captions for food images, including recipe instructions, nutritional facts, and serving suggestions
---	--	---

For a wide range, argument, a feature that generates persistence and a futuristic text for product images, the ad design looks like high-end, sophisticated, and beautiful	Asks in question for brand identity that can generate more and description options for their mobile phones, showing images, experiences, and top internet for the generation	Showing a text that automatically generated information and spending descriptions for the whole images, providing details about the property, neighborhood, and amenities
--	--	---

Developing a narrative automatically generates descriptive text for educational strategies, techniques and resources needed to engage and motivate students for learning.

Design a pattern that generates references to uniformly distributed random numbers of all integer orders, and enables to search efficiently for the understanding of these patterns.

PASHA

Develop a lexical database that generates descriptive text for educational games, testing students' knowledge of words.	Create an application that generates captions for images of science and math words, exploring the precision, fluidity, and ambiguity in a student-friendly manner.	Design a language learning platform that generates captions and suggestions for images, testing students' ability to build vocabulary and understand new words in context.
---	--	--

Design a problem that generates information and richly relevant evidence for images of an individual's assets, and reflects to students the complexity of these assets.	Build a system that can add explanatory text to visual math problems, making it easier for students to understand and solve problems at concepts.	Create an environment for the generation, display, organization for linked images surrounding students to analyze and interpret individual assets (see figure).
---	---	---

On
auton
let
cop
is
emp
man
state

LAVANYA

<p>Create a mobile app that generates personalized and detailed captions for travel photos, including locations, food containers, and memorable experiences</p>	<p>Develop a healthy meal app that provides ongoing, customized list of recipes, ingredients, vitamins, flavors, and the value aspects each of the eating</p>	<p>Develop an app that provides cultural insights and trivia in image captions, enhancing tourists' understanding of local customs and traditions</p>
---	---	---

- Build a tool for information professionals that generates relevant options for a range of academic, health, or engineering articles, capturing the risk and uncertainty
- Design an application that helps business decision makers choose articles, all the while capturing the risk and uncertainty to those business goals to encourage exploration
- Build a tool that generates eco-friendly travel tips in 10 minutes
#ecotraveling responsible and sustainable travel practices

3 Group ideas

Take turns sharing your ideas while clustering similar or related notes as you go. Once all sticky notes have been grouped, give each cluster a sentence-like label. If a cluster is bigger than six sticky notes, try and see if you can break it up into smaller sub-groups.

Ⓢ 20 minutes

TIP Add customizable tags to sticky notes to make it easier to find, browse, organize, and categorize important ideas as themes within your mural.

Social Media and Content Creation

<ul style="list-style-type: none"> Develop a tool that automatically generates catchy and relevant captions for images to make them more engaging on social media platforms like Instagram and Twitter 	<ul style="list-style-type: none"> Create a mobile app that generates personalized and detailed captions for travel photos, including locations, travel companions, and memorable experiences
---	--

Educational Enhancement:

<p>Develop a test that automatically generates descriptive text for assessment items, keeping test items and educators better engaged and effective materials for students.</p>	<p>Create an application that generates captions for images of science experiments, maintaining the precision, sophistication, and accuracy in a student-friendly manner.</p>
---	---

Food and Culinary Experiences:

Developed a Finder's design app that provides engaging descriptions for images of local dishes, highlighting ingredients, flavors, and the cultural significance of the cuisine

Travel and Exploration:

For online stores, implement a feature that generates persuasive and informative text for product images, including details like features, specifications, and benefits.

Design a portfolio that generates, integrates and effectively presents evidence for progress of all individual students, and attributes to specific students understanding of these subjects.



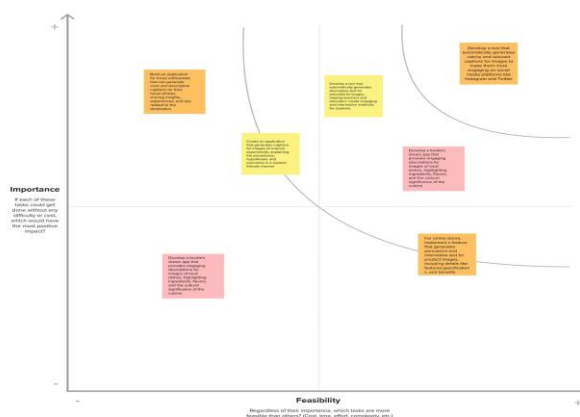
Idea prioritization

Prioritize

Your team should all be on the same page about what's important moving forward. Place your ideas on this grid to determine which ideas are important and which are feasible.

Ⓢ 20 minutes

TIP Participants can use their cursors to point at where sticky notes should go on the grid. The facilitator can confirm the spot by using the laser pointer holding the **H** key on the keyboard.



4. REQUIREMENT ANALYSIS

4.1 Functional requirement

Functional requirements describe the essential capabilities, features, and interactions that the system must exhibit to fulfil its intended purpose. These requirements are specific and measurable, forming the basis for the development and evaluation of the system. The system shall provide Quick Processing which is the secure user authentication to ensure that only authorized users can access the platform. The advanced customization functional requirements pertain to the system's capacity to provide users with tailored and flexible options for personalizing their experience. These features contribute to user satisfaction and engagement by allowing individuals to adapt to the system.

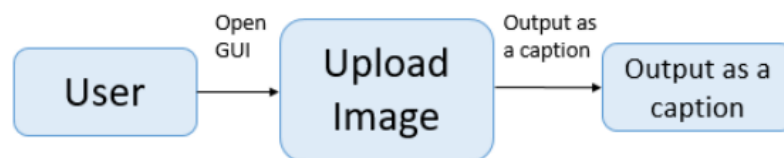
Related to analytics and insights focus on providing users with the capability to gain a deeper understanding of their interactions, activities, and trends within the system. These features aim to empower users with data-driven insights to enhance decision-making and optimize their experience. The functional requirements associated with cost-effective solutions focus on optimizing the system to minimize resource utilization, enhance operational efficiency, and reduce overall expenses while maintaining high performance and user satisfaction. Data Analysis Tools aims to provide users with robust capabilities for extracting meaningful insights from the data generated within the system. These tools empower users to analyze trends, make informed decisions, and derive valuable knowledge from the platform's data. The system shall implement predictive analytics to anticipate user preferences, trends, or potential issues. Predictive analytics features aim to enhance the user experience by proactively suggesting relevant content or highlighting emerging trends.

4.2 Non-Functional requirements

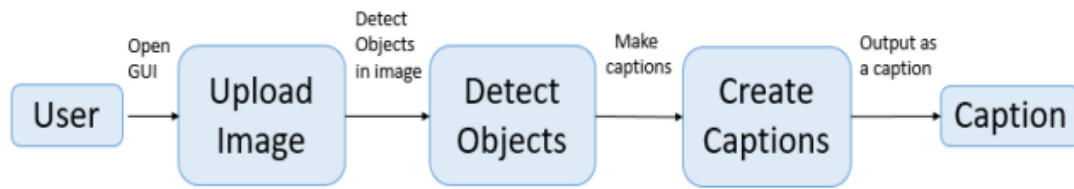
Non-functional requirements are aspects of a system that describe its qualities or characteristics, rather than specific behaviors. The system shall respond to user requests within an acceptable timeframe, with a maximum response time of 2 seconds for image uploads and caption generation. Fault Tolerance, In the event of a server failure, the system should automatically redirect traffic to redundant servers, minimizing service disruption. The system architecture shall be designed with modularity in mind, allowing for easy updates, additions, or modifications to individual components without affecting the entire system. The system shall undergo load testing to verify its performance under expected and peak loads, ensuring it meets performance requirements.

5. PROJECT DESIGN

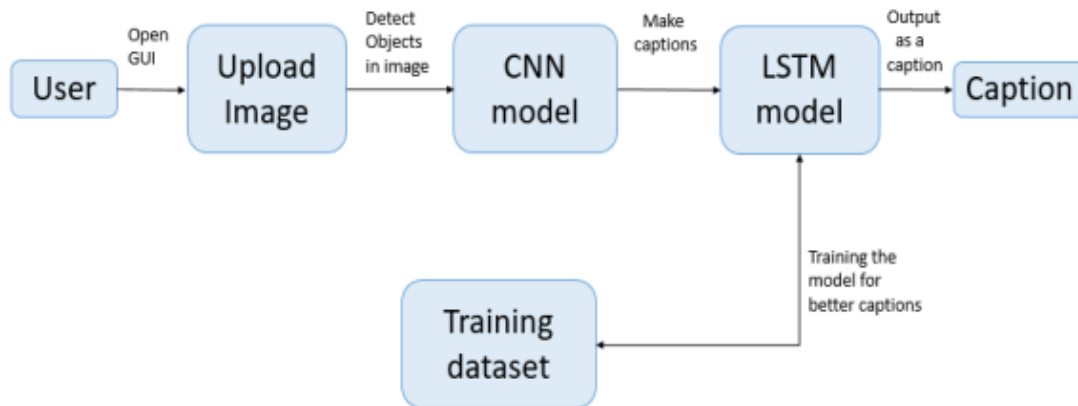
5.1 Data Flow Diagrams & User Stories



Data Flow Diagram Level 0



Data Flow Diagram Level 1



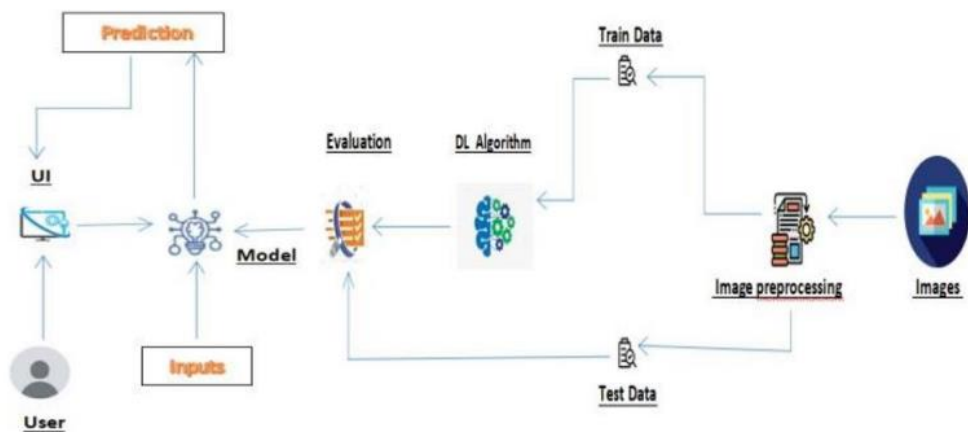
Final Data Flow Diagram Level 2

User stories-

Sprint	Functional Requirement (Epic)	User Story Number	User Story / Task	Story Points	Priority	Team Members
Sprint-1	Include an image upload option for the caption generation	ICG-1	As a user, I can upload an image for caption generation	3	High	Jyoshil
Sprint-2	Include a language selection menu or dropdown for supporting multiple languages for caption generation	ICG-2	As a user, I can select the language for generated captions	2	Medium	Lavanya
Sprint-3	Overlay generated caption on the uploaded image And allow users to adjust the position and style of the caption	ICG-3	As a user, I can preview the generated caption on the image	2	High	Pasha
Sprint-4	Provide a download button	ICG-4	As a user, I can download the	1	Low	Jyoshil

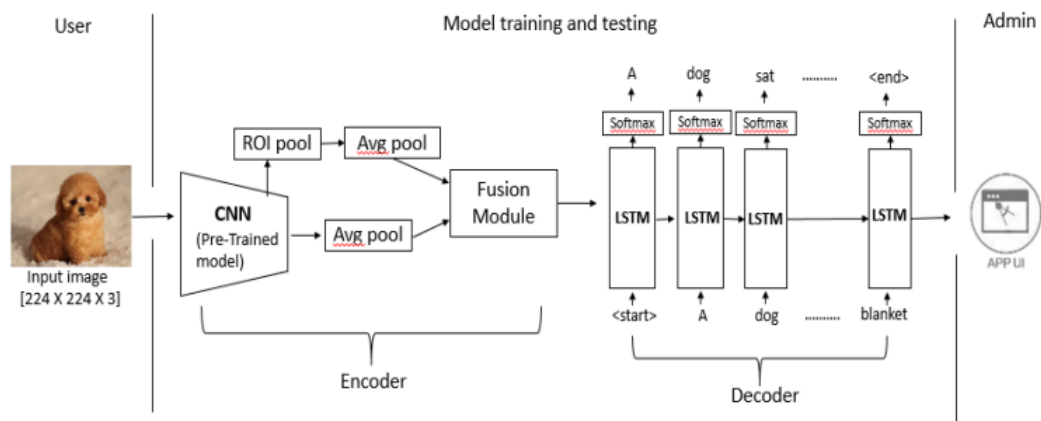
	for the generated caption and Support various download formats (e.g., TXT, PDF)		generated caption			
Sprint-5	Integrate with social media APIs for sharing	ICG-5	As a user, I can share the generated caption on social media	2	Medium	Lavanya
Sprint-6	Select and integrate a machine learning model for caption generation	ICG-6	Implement machine learning model for image caption generation	3	High	Pasha

5.2 Solution Architecture



6. PROJECT PLANNING & SCHEDULING

6.1 Technical Architecture



6.2 Sprint Planning & Estimation

All sprints

⚡ ☆ ↗ Complete sprint ⋮

Q

LG

A

CJ

Epic ▾ Sprint ▾

GROUP BY None ▾ ↺ Insights ⚙

TO DO 2

As a user, I want to have the option to choose from a variety of machine learning models for caption generation.

SELECT AND INTEGRATE A MACHINE LEA...

ICG-17 3 A

As a user, I want the system to integrate with cloud-based machine learning services

SELECT AND INTEGRATE A MACHINE LEA...

ICG-18 3 A

IN PROGRESS 2

As a user, i want to have Support for various download formats (e.g., TXT, PDF)

PROVIDE A DOWNLOAD BUTTON FOR TH...

ICG-12 1 CJ

As a user, I want the integration to provide analytics or insights on the performance of the content

INTEGRATE WITH SOCIAL MEDIA APIS FO...

ICG-15 2 LG

DONE 8 ✓

As a user, I want to preview the generated caption on the image

OVERLAY GENERATED CAPTION ON THE ...

ICG-8 ✓ 2 A

As a user, I want the ability to set a default language for captions, so that I don't have to manually select it every time I access content.

LANGUAGE SELECTION DROPDOWN OR ...

ICG-6 ✓ 2 LG

+

Timeline

🗨 Give feedback ➦ Share 📄 Export ⋮

Q

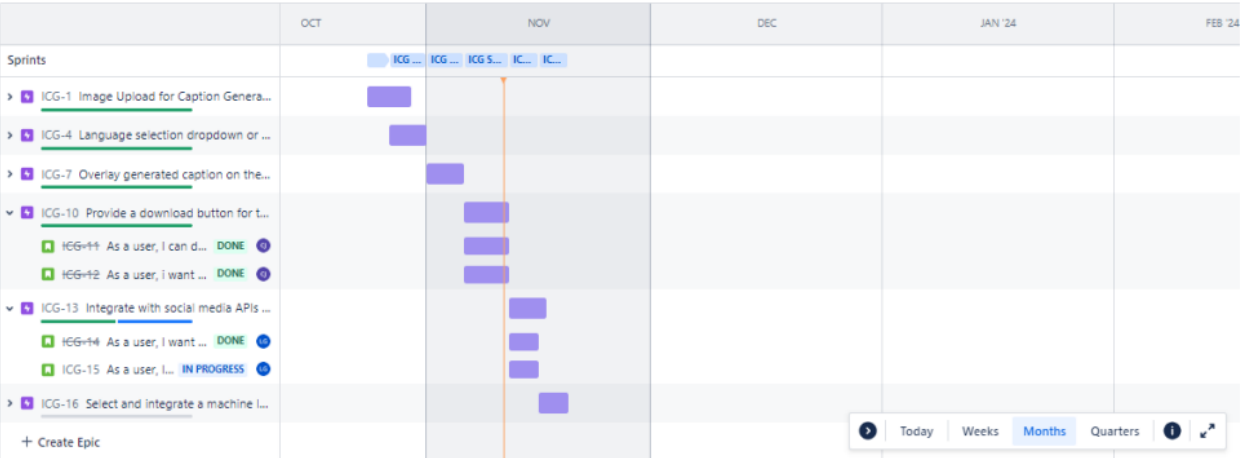
LG

A

CJ

Status category ▾ Epic ▾

⚙ View settings



Backlog

Search: LG CJ A Epic Insights

Epic

Issues without epic

- > Image Upload for Caption Generation
- > Language selection dropdown or menu
- > Overlay generated caption on the uploaded image and adjust the position and size

+ Create epic

ICG Sprint 1 24 Oct – 28 Oct (2 issues) 0 Complete sprint

- ICG-2 As a user, I want image Upload option for Caption G... IMAGE UPLO... DONE 3
- ICG-3 As a user, I want to receive suggestions for alternativ... IMAGE UPLO... DONE 3

+ Create issue

ICG Sprint 2 27 Oct – 31 Oct (2 issues) 0 Complete sprint

Implement a language selection dropdown or menu for supporting multiple languages for caption generation.

- ICG-6 As a user, I want the ability to set a default lang... LANGUAGE S... DONE 2
- ICG-5 As a user, I want to select the language to generate t... LANGUAGE S... DONE 2

+ Create issue

6.3 Sprint Delivery Schedule

Sprint	Total Story Points	Duration	Sprint Start Date	Sprint End Date (Planned)	Story Points Completed (as on Planned End Date)	Sprint Release Date (Actual)
Sprint-1	20	4 Days	24 Oct 2023	28 Oct 2023	17	29 Oct 2023
Sprint-2	20	3 Days	28 Oct 2023	31 Nov 2023	11	01 Nov 2023
Sprint-3	20	4 Days	01 Nov 2023	05 Nov 2023	9	05 Nov 2023
Sprint-4	20	5 Days	06 Nov 2023	11 Nov 2023	13	11 Nov 2023
Sprint-5	20	3 Days	12 Nov 2023	15 Nov 2023	12	16 Nov 2023
Sprint-6	20	3Days	16 Nov 2023	19 Nov 2023	10	19 Nov 2023

7. CODING & SOLUTIONING (Explain the features added in the project along with code)

7.1 Feature of images

```
In [12]: features = {}

directory = os.path.join(BASE_DIR, 'images')

for img_name in tqdm(os.listdir(directory)):
    #Load the image from file
    img_path = directory + '/' + img_name
    image = load_img(img_path, target_size=(224,224))
    #convert image pixels to numpy array
    image = img_to_array(image)
    #reshape data for 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
```

First, it initializes the empty dictionary called features that will be used to store the extracted features from images. Next it constructs a directory path by joining the BASE_DIR variable (which should be defined somewhere in the code) with the 'images' subdirectory. The tqdm function is used to create a progress bar to visualize the progress of the loop. Next within the loop, the code constructs the full path of the image file and loads the image using the load_img function from a library. The loaded image is converted into a NumPy array using the img_to_array function. The reshape function is used to add an extra dimension to the array. The image data is preprocessed to be suitable for the VGG model. Next the image ID is extracted from the filename by splitting the filename at the '.' (dot) and taking the first part. Finally, the extracted feature vector (feature) is stored in the features dictionary using the image ID as the key.

7.2 Features of captions

```
In [16]: # create mapping of image to captions
mapping = {}
# process lines
for line in tqdm(captions_doc.split('\n')):
    # split the line by comma(,)
    tokens = line.split(',')
    if len(tokens) < 2:
        continue
    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 list if needed
    if image_id not in mapping:
        mapping[image_id] = []
    # store the caption
    mapping[image_id].append(caption)

0% | 0/40456 [00:00<?, ?it/s]
```

First, it initializes an empty dictionary called mapping. This dictionary will be used to store image IDs as keys and corresponding captions as values. Next, the code enters a loop that iterates over each line in the captions_doc. Next, Each line is split into a list of tokens using a comma as the delimiter. Next, it extracts the image ID and the caption from the list of tokens. Next, If the image ID is not already present in the mapping dictionary, an empty list is created as the value associated with that image ID. Finally, the caption is appended to the list associated with the corresponding image ID in the mapping dictionary.

8. PERFORMANCE TESTING

8.1 Performance Metrics

```
In [ ]: # train the model
epochs = 40
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 [=====] - 84s 325ms/step - loss: 5.2246 - accuracy: 0.1499
227/227 [=====] - 64s 283ms/step - loss: 3.9249 - accuracy: 0.2573
227/227 [=====] - 64s 280ms/step - loss: 3.4750 - accuracy: 0.2935
227/227 [=====] - 64s 282ms/step - loss: 3.2052 - accuracy: 0.3140
227/227 [=====] - 64s 283ms/step - loss: 3.0115 - accuracy: 0.3300
227/227 [=====] - 63s 277ms/step - loss: 2.8616 - accuracy: 0.3461
227/227 [=====] - 63s 277ms/step - loss: 2.7426 - accuracy: 0.3601
227/227 [=====] - 64s 280ms/step - loss: 2.6475 - accuracy: 0.3724
227/227 [=====] - 64s 279ms/step - loss: 2.5666 - accuracy: 0.3827
227/227 [=====] - 64s 280ms/step - loss: 2.4960 - accuracy: 0.3930
227/227 [=====] - 63s 277ms/step - loss: 2.4401 - accuracy: 0.4004
227/227 [=====] - 65s 285ms/step - loss: 2.3874 - accuracy: 0.4083
227/227 [=====] - 65s 284ms/step - loss: 2.3407 - accuracy: 0.4141
227/227 [=====] - 61s 270ms/step - loss: 2.2996 - accuracy: 0.4212
227/227 [=====] - 64s 283ms/step - loss: 2.2604 - accuracy: 0.4273
227/227 [=====] - 64s 279ms/step - loss: 2.2232 - accuracy: 0.4326
227/227 [=====] - 64s 279ms/step - loss: 2.1894 - accuracy: 0.4381
227/227 [=====] - 64s 280ms/step - loss: 2.1606 - accuracy: 0.4440
227/227 [=====] - 64s 283ms/step - loss: 2.1304 - accuracy: 0.4482
227/227 [=====] - 63s 275ms/step - loss: 2.1028 - accuracy: 0.4530
227/227 [=====] - 63s 277ms/step - loss: 2.0749 - accuracy: 0.4578
227/227 [=====] - 62s 270ms/step - loss: 2.0533 - accuracy: 0.4610
227/227 [=====] - 65s 288ms/step - loss: 2.0313 - accuracy: 0.4651
227/227 [=====] - 65s 285ms/step - loss: 2.0099 - accuracy: 0.4683
227/227 [=====] - 63s 275ms/step - loss: 1.9887 - accuracy: 0.4725
227/227 [=====] - 64s 282ms/step - loss: 1.9715 - accuracy: 0.4754
227/227 [=====] - 67s 293ms/step - loss: 1.9512 - accuracy: 0.4797
227/227 [=====] - 65s 287ms/step - loss: 1.9356 - accuracy: 0.4827
227/227 [=====] - 64s 281ms/step - loss: 1.9167 - accuracy: 0.4859
227/227 [=====] - 63s 277ms/step - loss: 1.9038 - accuracy: 0.4881
227/227 [=====] - 65s 285ms/step - loss: 1.8895 - accuracy: 0.4906
227/227 [=====] - 66s 292ms/step - loss: 1.8752 - accuracy: 0.4935
227/227 [=====] - 64s 282ms/step - loss: 1.8610 - accuracy: 0.4958
227/227 [=====] - 65s 285ms/step - loss: 1.8462 - accuracy: 0.4990
227/227 [=====] - 65s 284ms/step - loss: 1.8327 - accuracy: 0.5010
227/227 [=====] - 66s 291ms/step - loss: 1.8180 - accuracy: 0.5041
227/227 [=====] - 64s 283ms/step - loss: 1.8052 - accuracy: 0.5071
227/227 [=====] - 63s 277ms/step - loss: 1.7930 - accuracy: 0.5097
227/227 [=====] - 63s 275ms/step - loss: 1.7791 - accuracy: 0.5121
227/227 [=====] - 63s 277ms/step - loss: 1.7662 - accuracy: 0.5145

```
In [ ]: 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)))
```

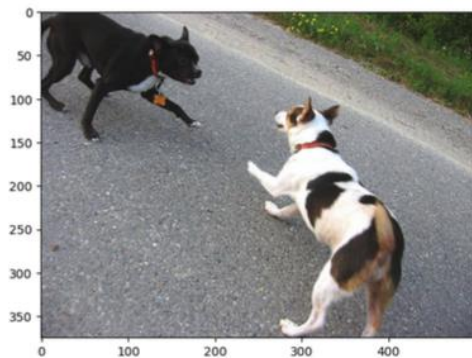
```
BLEU-1: 0.537093
BLEU-2: 0.312298
BLEU-1: 0.537468
BLEU-2: 0.312138
BLEU-1: 0.537433
BLEU-2: 0.311986
BLEU-1: 0.537506
BLEU-2: 0.312160
BLEU-1: 0.537093
BLEU-2: 0.311778
BLEU-1: 0.536933
BLEU-2: 0.311662
BLEU-1: 0.536907
BLEU-2: 0.311554
BLEU-1: 0.537044
BLEU-2: 0.311544
BLEU-1: 0.537404
BLEU-2: 0.311599
BLEU-1: 0.537772
BLEU-2: 0.311948
```

9. RESULTS

9.1 Output Screenshots

```
In [ ]: generate_caption("1001773457_577c3a7d70.jpg")
```

```
-----Actual-----
start black dog and spotted dog are fighting end
start black dog and tri-colored dog playing with each other on the road end
start black dog and white dog with brown spots are staring at each other in the street end
start two dogs of different breeds looking at each other on the road end
start two dogs on pavement moving toward each other end
-----Predicted-----
start black dog and spotted dog playing on pavement end
```



```
In [ ]: generate_caption("1002674143_1b742ab4b8.jpg")
```

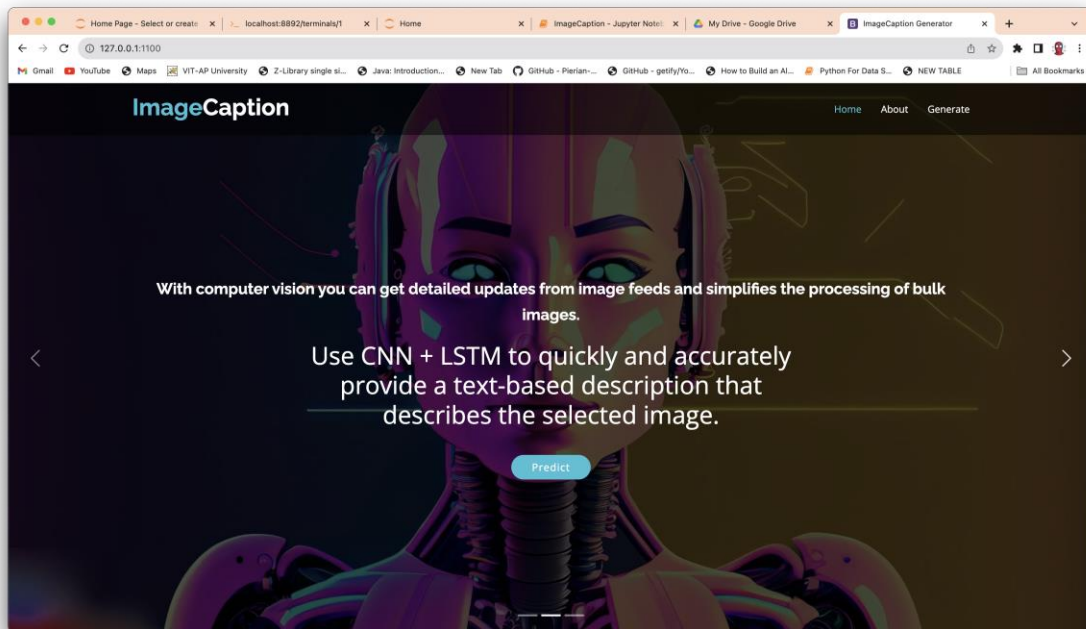
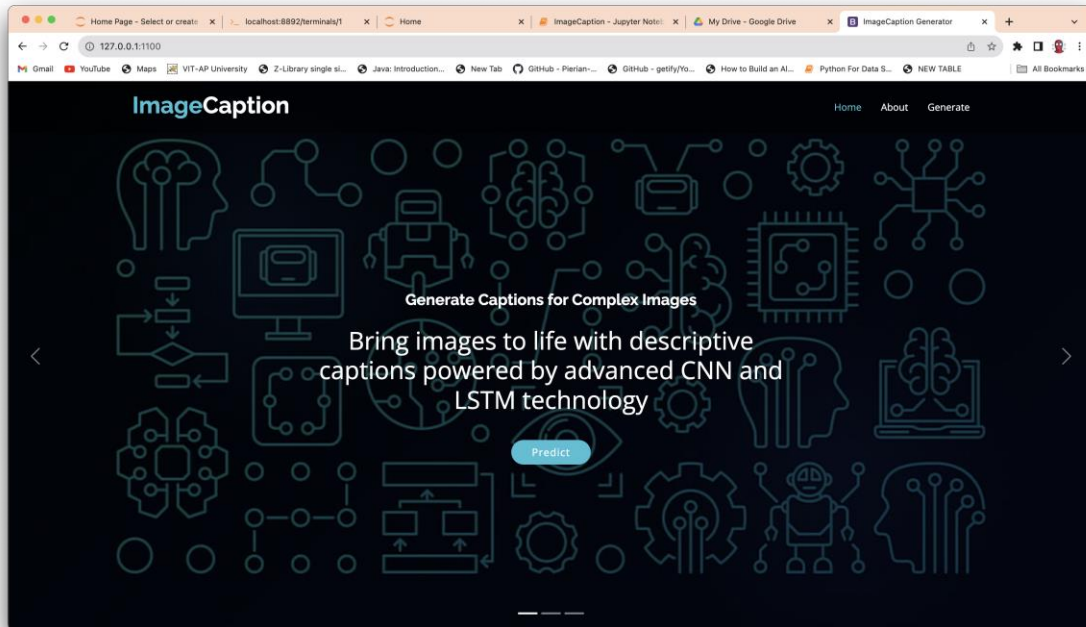
```
-----Actual-----  
start little girl covered in paint sits in front of painted rainbow with her hands in bowl end  
start little girl is sitting in front of large painted rainbow end  
start small girl in the grass plays with fingerpaints in front of white canvas with rainbow on it end  
start there is girl with pigtales sitting in front of rainbow painting end  
start young girl with pigtales painting outside in the grass end  
-----Predicted-----  
start two children playing in plastic colored plastic structure end
```

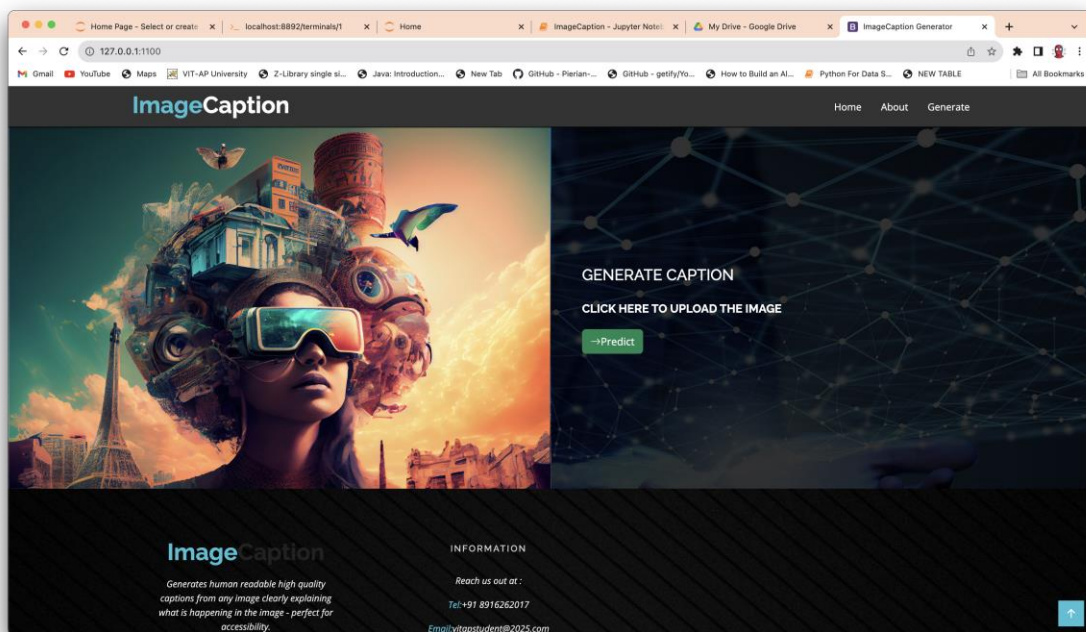
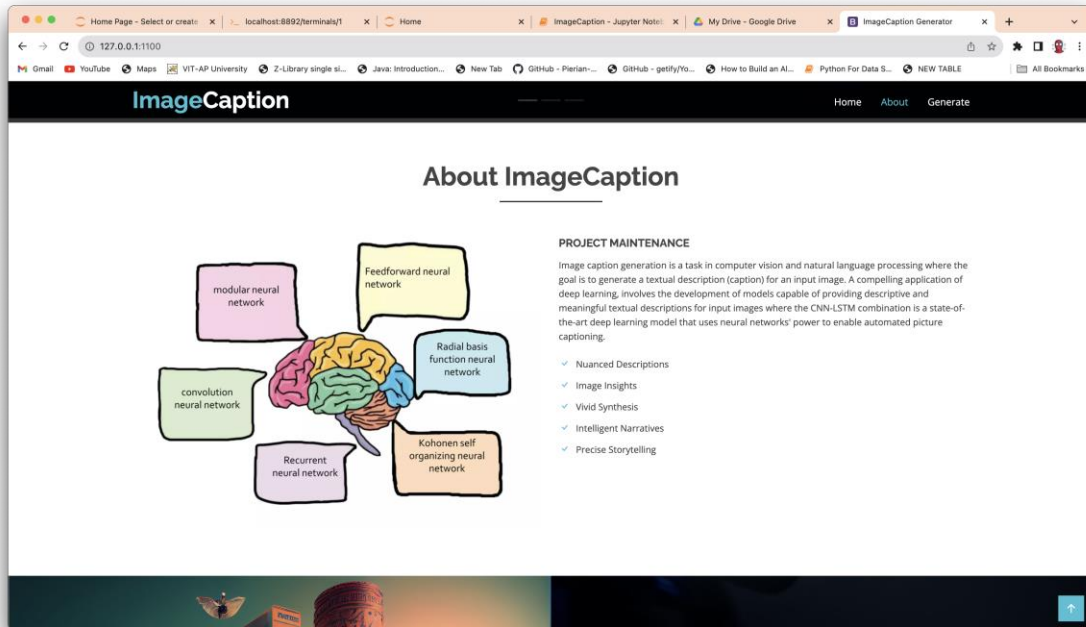


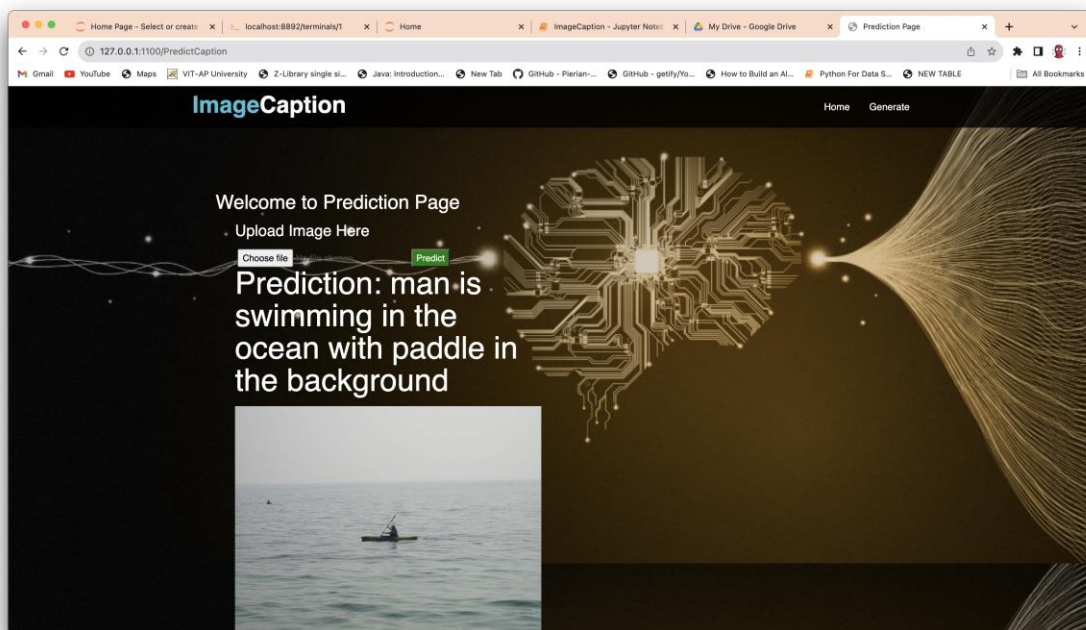
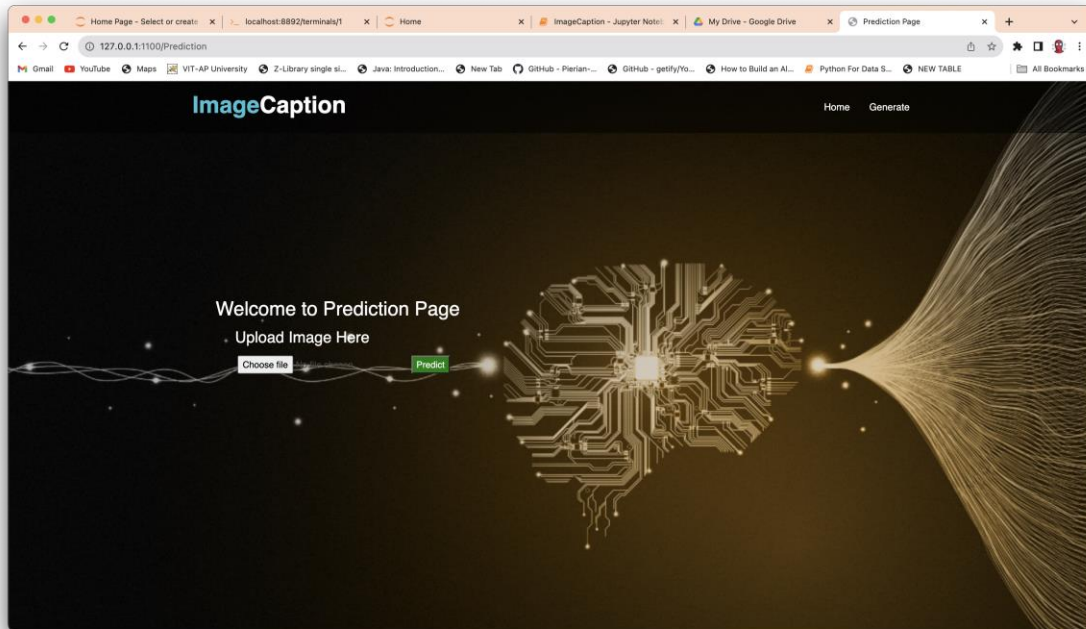
```
In [ ]: generate_caption("101669240_b2d3e7f17b.jpg")
```

```
-----Actual-----  
start man in hat is displaying pictures next to skier in blue hat end  
start man skis past another man displaying paintings in the snow end  
start person wearing skis looking at framed pictures set up in the snow end  
start skier looks at framed pictures in the snow next to trees end  
start man on skis looking at artwork for sale in the snow end  
-----Predicted-----  
start man in skis is displaying paintings in the snow end
```









10. ADVANTAGES & DISADVANTAGES

Advantages-

- Image captions can make visual content accessible to individuals with visual impairments, allowing them to understand and engage with images through text
- Image captions can enhance the searchability of images.

- Enhanced User Experience
- Image captioning can be applied in assistive technologies
- Image captions can serve as concise summaries of visual content, helping users quickly grasp the key information or context conveyed by an image.

Disadvantages-

- Generating accurate and contextually relevant captions for images can be challenging due to the inherent ambiguity
- Image captioning models may overfit to the specific characteristics of the training data, leading to limitations in their ability to generalize
- Captioning models may struggle with generating captions for uncommon or novel scenarios
- Captioning models may not always provide factually accurate information

11. CONCLUSION

Image caption generator recognizes the image context and describes it in English language. It comprises of steps like data cleaning, extracting features, tokenizing the vocabulary, defining, training, validating and testing the model. In this report the generation of caption for the image using CNN with LSTM is depicted. The advantage of training using CNN with LSTM are the ability to handle sequences of arbitrary length, and more importantly, the end-to-end maximization of the joint probability of the source and target sentence, have produced output in machine translation. The efficiency of the model is calculated using BLUE Score. This system is data-driven, therefore it can't anticipate phrases which are not in its vocabulary. The following tasks can be performed in the future. Text-to speech technology, which automatically reads aloud to visually challenged people. Translating images directly into sentences, rather than creating image captions. Static images can only provide information about one particular instant in time to blind people, thus creating video captions may theoretically provide continuous real-time information to blind people. Using CNN and LSTM the developed system contains characteristics that predict an image's caption.

12. FUTURE SCOPE

The future scope of the image caption generator presents a myriad of opportunities for advancement and refinement. Beyond its current capabilities, there is potential for the integration of cutting-edge technologies to elevate the system's performance and broaden its applications. Exploring text-to-speech technology stands as a promising avenue, enhancing accessibility for visually impaired individuals by enabling the automatic vocalization of generated captions. Ongoing research could focus on refining the efficiency and accuracy of the CNN-LSTM model, potentially experimenting with novel architectures and optimizing hyperparameters. Multimodal approaches, integrating additional modalities such as audio, present an avenue for more comprehensive understanding and description of diverse content. Evaluating the system using metrics beyond the BLEU Score could provide a nuanced assessment of its performance. User-centric improvements guided by feedback and experiences can further enhance the system's usability, while ethical considerations and inclusive design principles must be central to its evolution. In essence, the future scope of the image

caption generator encompasses a holistic approach to innovation, user experience, and ethical considerations, paving the way for a more advanced and inclusive technology landscape.

13. **APPENDIX**

Source Code –

<https://colab.research.google.com/drive/1vLAB4Ivf0C1RirRupFtOGY25FV6-OZMw?usp=sharing>

GitHub link – <https://github.com/smartinternz02/SI-GuidedProject-643890-1700545152>

Project Demo Link –

<https://drive.google.com/file/d/1j-QEA2TzXGxOAh50eYfhnaIYXJeUvAAX/view?usp=sharing>

