

Image Caption Generator with Accessibility Features

Prathyusha Reddy Nandikonda
Computer Science
University of Massachusetts
Lowell, MA

prathyushareddy_nandikonda@student.uml.edu

Ranga Rohit Nallamolu
Computer Science
University of Massachusetts
Lowell, MA

rangarohit_nallamolu@student.uml.edu

Naga Sai Tejaswi Gandu
Computer Science
University of Massachusetts
Lowell, MA

nagasaitejaswi_gandu@student.uml.edu

Abstract—Our project aims to revolutionize accessibility for visually impaired individuals by leveraging cutting-edge computer vision and natural language processing techniques to generate accurate and detailed descriptions of visual content. By developing an innovative image caption generator, we seek to bridge the gap between the visual world and those with visual impairments, enabling them to interact more effectively with digital content. This system holds immense potential beyond aiding visually impaired users; it extends to applications in robotic vision, business, image indexing, social media, and various natural language processing domains. Utilizing deep learning methodologies, our approach strives to emulate human-like image description capabilities, ultimately enhancing user experiences by automating the generation of descriptive captions for images and environments.

Keywords—Accessibility, Visually impaired individual, Deep learning methodologies, Image caption generator, Natural language processing.

I. INTRODUCTION

In today's visually driven digital realm, ensuring accessibility for individuals with visual impairments remains a pressing challenge. Our project is dedicated to developing a novel solution—an image caption generator—that aims to address this crucial gap by enabling visually impaired individuals to access and understand visual information effectively.

The foundation of our project is an effective approach that makes use of the Flickr 8k dataset, comprising 8000 images, each accompanied by five descriptive captions. This dataset serves as the foundational resource for our research, providing a diverse range of images and their corresponding textual descriptions.

Our approach unfolds in two fundamental stages. The initial phase involves employing Convolutional Neural Networks (CNN) to extract essential features from the

images, facilitating a deeper comprehension of their visual content. Subsequently, we utilize Recurrent Neural Networks (RNN), specifically Long Short Term Memory (LSTM) architectures, to generate coherent and meaningful textual descriptions based on the extracted image features.

By merging state-of-the-art technologies in image processing and language generation, our goal is to train a sophisticated model capable of accurately describing visual content, thereby enhancing accessibility for individuals with visual impairments by converting it into speech. Our aspiration is to contribute to a more inclusive digital landscape, where everyone, regardless of visual abilities, can engage with and comprehend visual information effortlessly.

II. MOTIVATION

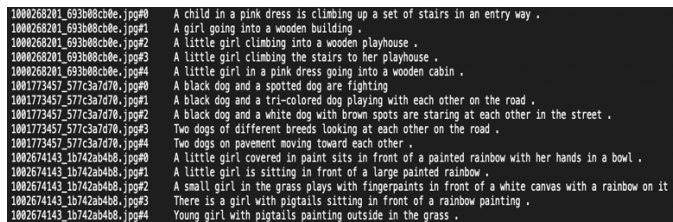
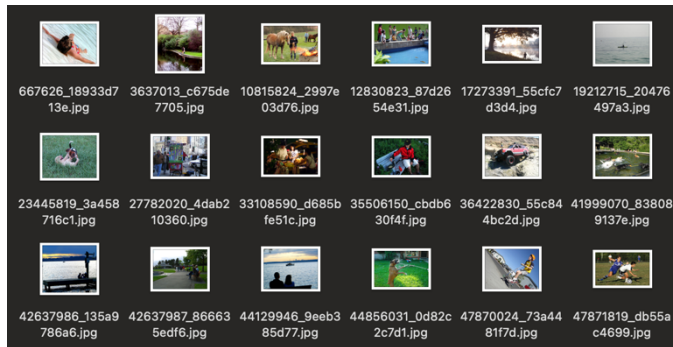
Our project's motivation is to empower visually impaired individuals in the digital landscape. By harnessing computer vision and natural language processing, we aim to provide accurate image captions, fostering a more inclusive experience for those with visual impairments. We strive to break down barriers, allowing these individuals to access and comprehend visual content effortlessly, enhancing their interaction with the digital world. Our goal is to ensure equal access and participation for everyone, enabling a more meaningful and inclusive online experience for individuals with visual impairments.

III. RELATED WORK

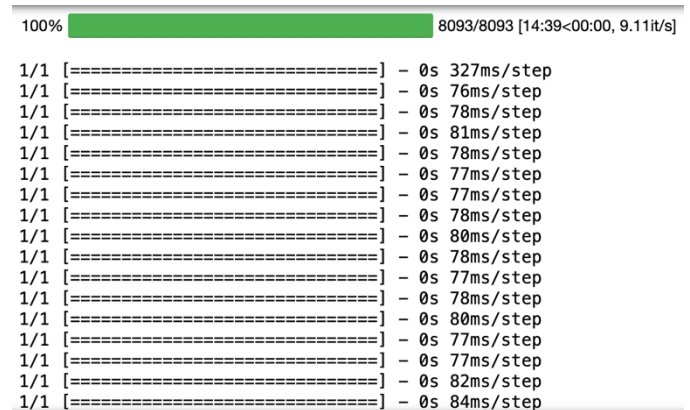
The related work in the field of image caption generation using CNN and RNN techniques has showcased significant advancements in making visual content accessible. ^[1]Panicker et al. (2021) contributed to the International Journal of Innovative Technology and Exploring Engineering with their research on an image caption generator. Their utilization of CNN and RNN methodologies demonstrated promising results in generating descriptive captions for images.

Both research papers underscore the efficacy of employing Convolutional Neural Networks (CNN) for feature extraction from images and Recurrent Neural Networks (RNN) for generating coherent textual descriptions. These findings form a critical foundation for our project, inspiring and guiding our efforts to develop an image caption generator with enhanced accessibility features for individuals with visual impairments.

The Flickr_8K dataset comprises around 8,000 images designated for training, with a set of images—1,000 each—allocated for development and testing purposes, separated into Flickr_8k.trainImages.txt, Flickr_8k.devImages.txt, and Flickr_8k.testImages.txt, respectively. Additionally, the Flickr8k.token.txt file within the Flickr8k_text folder contains five captions for each image in the dataset, providing diverse textual descriptions for the images.



3. Feature Extraction using Transfer Learning: Our project implements transfer learning by utilizing the Xception model pre-trained on the ImageNet dataset. Extracted features from images are obtained by removing the classification layer, generating 2048 feature vectors for each image. These features are stored in a pickle file for future use.



6. Data Generator Creation: Due to memory constraints, we employed a generator method to yield input and output sequences, managing the large dataset efficiently. This generator method helps in training the model with batches of data.

7. Model Architecture Definition: The model architecture comprises three main parts:

1. **Feature Extractor:** Initially, the image's extracted features (with a size of 2048) are processed through a dense layer to reduce dimensions to 256 nodes.
2. **Sequence Processor:** Textual inputs go through an embedding layer followed by an LSTM layer to handle the sequence of words in the captions.
3. **Decoder:** The outputs from the above layers are combined and processed by a dense layer to make the final prediction. The final layer contains nodes equal to the size of our vocabulary, facilitating word predictions.

8. Model Training: Training involves utilizing the 6000 images in batches, where input and output sequences are fed into the model using the `model.fit_generator()` method. The trained model is saved for future use.

```
Dataset: 6000
Descriptions: train= 6000
Photos: train= 6000
Vocabulary Size: 7577
Description Length: 32
Model: "model_1"
```

Layer (type)	Output Shape	Param #	Connected to
input_6 (InputLayer)	[(None, 32)]	0	[]
input_5 (InputLayer)	[(None, 2048)]	0	[]
embedding_1 (Embedding)	(None, 32, 256)	1939712	['input_6[0][0]']
dropout_2 (Dropout)	(None, 2048)	0	['input_5[0][0]']
dropout_3 (Dropout)	(None, 32, 256)	0	['embedding_1[0][0]']
dense_3 (Dense)	(None, 256)	524544	['dropout_2[0][0]']
lstm_1 (LSTM)	(None, 256)	525312	['dropout_3[0][0]']
add_25 (Add)	(None, 256)	0	['dense_3[0][0]', 'lstm_1[0][0]']
dense_4 (Dense)	(None, 256)	65792	['add_25[0][0]']
dense_5 (Dense)	(None, 7577)	1947289	['dense_4[0][0]']

```

Total params: 5002649 (19.08 MB)
Trainable params: 5002649 (19.08 MB)
Non-trainable params: 0 (0.00 Byte)

6000/6000 [=====] - 454s 76ms/step - loss: 4.5037
1/6000 [.....] - ETA: 8:16 - loss: 4.2338

/Users/prathyushareddynandikonda/anaconda3/lib/python3.11/site-packages/keras:
ing: You are saving your model as an HDF5 file via `model.save()`. This fil
end using instead the native Keras format, e.g. `model.save('my_model.keras
saving_api.save_model(

6000/6000 [=====] - 443s 74ms/step - loss: 3.6593
6000/6000 [=====] - 447s 74ms/step - loss: 3.3756
6000/6000 [=====] - 544s 91ms/step - loss: 3.2055
6000/6000 [=====] - 443s 74ms/step - loss: 3.0855
6000/6000 [=====] - 472s 79ms/step - loss: 2.9949
6000/6000 [=====] - 451s 75ms/step - loss: 2.9266
6000/6000 [=====] - 430s 72ms/step - loss: 2.8747
6000/6000 [=====] - 437s 73ms/step - loss: 2.8305
6000/6000 [=====] - 432s 72ms/step - loss: 2.7944
```

Fig 4: Training the Model

9. Testing the Model: We defined a `test_model` which loads the trained model to generate predictions. Predicted index values are transformed back to words using the `tokenizer.p` file, allowing for human-readable captions. It also predicts the BLEU scores showing the accuracy of the model.

10. Integrating accessibility features: This step involves implementing Google Text-to-Speech (TTS) functionality

to convert generated textual captions into audio descriptions for enhanced accessibility.

By incorporating the Google TTS API, our project enables the conversion of text-based image captions into spoken audio files. This process entails sending the generated caption text to Google's TTS service through API requests, which in turn produces an audio rendition of the description.

The resulting audio file serves as an additional accessibility feature, providing visually impaired users with the option to listen to spoken descriptions of the images.

VI. EXPERIMENTAL RESULTS AND DISCUSSIONS

In order to evaluate our model, BLEU scores were taken into consideration which demonstrates reasonable performance, indicating a moderate level of linguistic similarity between the generated captions and the ground truth annotations. The generated caption "A man in a red shirt standing on top of a snowy mountain" showcases the model's ability to capture basic elements in the image. The generated caption is then converted to speech and stored in a .mp3 file which ensures accessibility for individuals with visual impairments.

The BLEU-1 score, indicating the similarity of unigrams between generated and reference captions, stands at 0.3549, showcasing a moderate level of overlap. Meanwhile, the BLEU-2 score, measuring bigram similarity, registers at 0.1884, indicating a lesser degree of match between generated and reference text pairs. Furthermore, the BLEU-3 and BLEU-4 scores, assessing trigram and 4-gram overlap, respectively, demonstrate even lower similarity levels, recording values of 0.1272 and 0.0536.

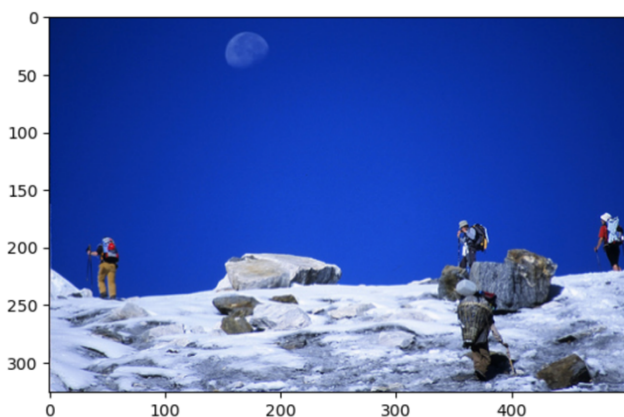
```

100% ██████████ 1000/1000 [03:58:00:00, 4.53it/s]

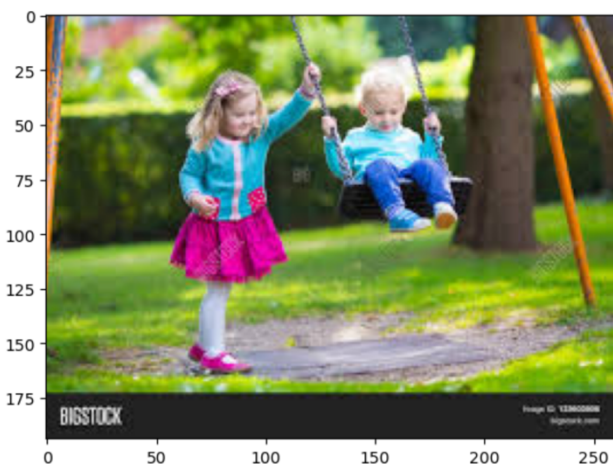
BLEU-1: 0.354917
BLEU-2: 0.188429
BLEU-3: 0.127171
BLEU-4: 0.053565
```

Fig 5: BLEU scores

```
<matplotlib.image.AxesImage at 0x38895b010>
```



```
<matplotlib.image.AxesImage at 0x4031d3010>
```



[https://www.google.com/imgres?imgurl=https%3A%2F%2Fstatic3.bigstockphoto.com%2F3%2F3%2F1%2F1%2F1%2F133603808.jpg&tbid=OAN2Xn8b6CRMVtMY&view=12abUKEwjlZGL7-YyDAX4FFkFH0RfCFeQMygrelUARkA0I&imgref=https%3A%2F%2Fwww.google.com%2Fimgres%2Fimgres%2F133603808%2Fphoto-kind=photo-and-girl-on-a-playground-child-playing-outdoors-in-summer-kids-play-on-school-yard-having-a-happy-in-kindergarten-or-preschool-children-having-fun-at-day-care-play-ground-toddler-on-a-swing&docid=0gRT4dVc3QJkMw=1500&w=1120&q=two%20girls%20playing%20in%20a%20garden%20pictr&ved=2abUKEwjlZGL7-YyDAX4FFkFH0RfCFeQMygrelUARkA0I]

