# ⌄ COMS4995AML Project: Image Captioning

**Group 25**

Building a simple image captioning model using pre-trained VGG16 and LSTM

```
1  import pandas as pd
2  import numpy as np
3
4  import nltk
5  from nltk.tokenize import word_tokenize
6  from nltk.tag import pos_tag
7  nltk.download("punkt")
8  nltk.download("averaged_perceptron_tagger")
9  from nltk.translate.bleu_score import corpus_bleu
10
11 from sklearn.model_selection import train_test_split
12
13 import tensorflow as tf
14 from tensorflow.keras.applications.vgg16 import VGG16, preprocess_input
15 from tensorflow.keras.models import Model
16 from tensorflow.keras.preprocessing.image import load_img, img_to_array
17 from tensorflow.keras.utils import Sequence
18 from tensorflow.keras.preprocessing.text import Tokenizer
19 from tensorflow.keras.preprocessing.sequence import pad_sequences
20 from tensorflow.keras.layers import Input, Embedding, LSTM,\
21  Dense, Concatenate, Dropout, Flatten, Activation, Add
22 from tensorflow.keras.optimizers import Adam
23 from tensorflow.keras.activations import softmax
24 from tensorflow.keras.callbacks import Callback, EarlyStopping, ModelCheckpoin
```

```
[nltk_data] Downloading package punkt to /root/nltk_data...
[nltk_data]   Package punkt is already up-to-date!
[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data]     /root/nltk_data...
[nltk_data]   Package averaged_perceptron_tagger is already up-to-
[nltk_data]       date!
```

## ⌄ 1. Preprocessing Data

```
1 from google.colab import drive
2
3 drive.mount("./drive")
```

    Mounted at ./drive

```
1 %cd drive/MyDrive/COMS4995AML_Project
```

    /content/drive/MyDrive/COMS4995AML_Project

Preprocessing the data

```
1 # Reading the filtered data
2 flickr_df = pd.read_excel("./captions.xlsx")
3
4 # Removing all punctuations
5 regex = r"[^\s\w]"
6 flickr_df["comment"] = flickr_df["comment"].str.replace(regex, "", regex=True)
7
8 # Converting all words to lower case
9 def lower_case(text):
10   return text.lower()
11 flickr_df["comment"] = flickr_df["comment"].apply(lower_case)
12
13 flickr_df.head()
```

|   | image_name | comment |
|---|---|---|
| **0** | 1001773457.jpg | a black dog and a white dog with brown spots ... |
| **1** | 1001773457.jpg | a black dog and a tricolored dog playing with... |
| **2** | 1001773457.jpg | a black dog and a spotted dog are fighting |
| **3** | 1003163366.jpg | a man sleeping on a bench outside with a whit... |
| **4** | 1003163366.jpg | a man lays on the bench to which a white dog ... |

```
1 print("Number of images in the dataset: ",
2       len(set(list(flickr_df["image_name"]))))
3 print("Number of captions in the dataset: ",
4       len(list(flickr_df["comment"])))
```

```
Number of images in the dataset:  2494
Number of captions in the dataset:  9403
```

Adding a category column for each caption. Categories include "dog", "cat", "horse", and "cow".

```
 1 # Categories
 2 categories = [["dog", "puppy"], ["cat", "kitten"], ["horse"], ["cow"]]
 3
 4 # Using NLTK to parse the comment column and detect
 5 # the presence of key words for our categories
 6 def classify_animal(text):
 7   if isinstance(text, str):
 8     words = word_tokenize(text)
 9     tagged_words = pos_tag(words)
10     for word, tag in tagged_words:
11       if tag in ["NN", "NNS"]:
12         for c in categories:
13           if word.lower() in c:
14             return c[0]
15     return "other"
16
17 # flickr_df_copy will have a "category" column added
18 flickr_df_copy = flickr_df.copy(deep = True)
19 flickr_df_copy["category"] = flickr_df["comment"].apply(classify_animal)
20 flickr_df_copy.head()
```

# ∨ 2. Splitting train_df, val_df, test_df

We have an imbalanced dataset where there are a lot of images in the "dog" category and a few images in the other categories. Hence, we used the stratify strategy to preserve the proportions from the dataset. While the category column in the table is never used directly in training the captioning program, it nonetheless is essential for our stratified splitting at this stage.

Around 53.3% of the data are for training, 26.7% of the data are for validation, and 20% of the data are for testing.

```
1 dev_df, test_df = train_test_split(flickr_df_copy, test_size = 0.2,
2                                    stratify = flickr_df_copy["category"],
3                                    shuffle = True)
4 train_df, val_df = train_test_split(dev_df, test_size = 0.333,
5                                     stratify = dev_df["category"],
6                                     shuffle = True)
```

We saved train_df, val_df, and test_df so that we don't have to split the data every time we run the notebook.

```
1 train_df.to_csv("train_df.csv", index = False)
2 val_df.to_csv("val_df.csv", index = False)
3 test_df.to_csv("test_df.csv", index = False)
```

```
1 train_df = pd.read_csv("train_df.csv")
2 val_df = pd.read_csv("val_df.csv")
3 test_df = pd.read_csv("test_df.csv")
```

```
1 train_df.shape, val_df.shape, test_df.shape
```

```
   ((5017, 3), (2505, 3), (1881, 3))
```

# ∨ 3. Computer Vision: use VGG-16 model for image feature extraction

```
1 # Define a VGG-16 model for feature extraction
2 # load pre-trained weights from ImageNet
3 vgg16_model = VGG16(weights="imagenet")
4
5 # We omit the very last fully connected layer in VGG16
6 # since our categories aren't the same as ImageNet ones
7 vgg16_model = Model(inputs=vgg16_model.inputs,
8                     outputs=vgg16_model.layers[-2].output)
9
10 # VGG16 is frozen since its weights don't require any
11 # further training for our program
12 for layer in vgg16_model.layers:
13   layer.trainable = False
```

```
1 def predict_image_feature(image):
2   return vgg16_model.predict(image)
3
4 # This function opens an image file and extracts its feature for VGG16
5 def extract_features(filename):
6   try:
7     image = load_img(filename, target_size=(224, 224))
8     image = img_to_array(image)
9     image = image.reshape((1, image.shape[0], image.shape[1], image.shape[2]))
10    image = preprocess_input(image)
11
12    # Prediction using the VGG-16 model
13    feature = predict_image_feature(image)
14    return feature
15
16  # In case a training image is corrupted and can't be opened
17  except:
18    print("Failed to Open: ", filename)
19    return None
```

```
1 DATA_DIR = "./Images"
2 feature_dict = {}
```

Extracting image embeddings for train_df

```python
1 # Extracting image embeddings from train_df image names. The dictionary
2 # feature_dict stores processed images since some of the images correspond
3 # to multiple comments. This speeds up the processing
4
5 image_embeddings = []
6 for img_path in train_df["image_name"]:
7   if img_path not in feature_dict:
8     img = extract_features(DATA_DIR + "/" + img_path)
9     if img is None:
10       pass
11     feature_dict[img_path] = img
12   else:
13     img = feature_dict[img_path]
14   image_embeddings.append(img)
15
16 train_image_embeddings = np.array(image_embeddings)
```

```python
1 # Saving the NumPy array for train image embeddings
2 np.save("train_image_embeddings", train_image_embeddings)
```

```python
1 train_image_embeddings = np.load("train_image_embeddings.npy", allow_pickle=Tr
```

Extracting image embeddings for val_df

```python
1 # Extracting image embeddings for validation images in val_df
2 image_embeddings = []
3 for img_path in val_df["image_name"]:
4   if img_path not in feature_dict:
5     img = extract_features(DATA_DIR + "/" + img_path)
6     if img is None:
7       pass
8     feature_dict[img_path] = img
9   else:
10     img = feature_dict[img_path]
11   image_embeddings.append(img)
12
13 val_image_embeddings = np.array(image_embeddings)
```

```python
1 np.save("val_image_embeddings", val_image_embeddings)
```

```
1 val_image_embeddings = np.load("val_image_embeddings.npy", allow_pickle=True)
```

Extracting image embeddings for test_df

```
1 # Extracting image embeddings for test_df
2 image_embeddings = []
3 for img_path in test_df["image_name"]:
4   if img_path not in feature_dict:
5     img = extract_features(DATA_DIR + "/" + img_path)
6     if img is None:
7       img = np.zeros((1, 4096))
8     feature_dict[img_path] = img
9   else:
10    img = feature_dict[img_path]
11  image_embeddings.append(img)
12
13 test_image_embeddings = np.array(image_embeddings)
```

```
1 np.save("test_image_embeddings", test_image_embeddings)
```

```
1 test_image_embeddings = np.load("test_image_embeddings.npy", allow_pickle=True
```

```
1 test_image_embeddings.shape, train_image_embeddings.shape, val_image_embedding
```

```
((1881, 1, 4096), (5017, 1, 4096), (2505, 1, 4096))
```

## ⌄ 4. Define the ImageCaptionGenerator Class

Since the input data and labels for our model are multimodal, a customized generator is required for training the model. In the generator defined below, we pass batches of image_embeddings, text_embedding, next_token_embedding. This way, the model can learn to use image embeddings and partial text to generate the next caption.

```python
1  class ImageCaptionGenerator(tf.keras.utils.Sequence):
2    def __init__(self, img_embeddings, captions, tokenizer,
3                  max_length, batch_size=32):
4      self.img_embeddings = img_embeddings
5      self.captions = captions
6      self.tokenizer = tokenizer # the tokenizer instance used for encoding
7      self.max_length = max_length # maximum length of generated captions
8      self.batch_size = batch_size
9
10     def __len__(self): # A required function for the interface
11       return len(self.captions) // self.batch_size
12
13     # Another required function of the interface which
14     # returns one batch of input to the model
15     def __getitem__(self, idx):
16       start = idx * self.batch_size
17       end = (idx + 1) * self.batch_size
18       batch_images = self.img_embeddings[start:end] # batch of images
19       batch_captions = self.captions[start:end] # batch of captions
20
21       X_images, X_seq, y = [], [], []
22       for img_emb, cap in zip(batch_images, batch_captions):
23
24         # captions are encoded by the tokenizer
25         seq = self.tokenizer.texts_to_sequences([cap])[0]
26         for i in range(1, len(seq)):
27
28           # input for text is the encoding of the partial sequence while the
29           # label for text is the one-hot-encoding of the next token
30           in_seq, out_seq = seq[:i], seq[i]
31
32           # padding is required to keep sizes consistent
33           in_seq = pad_sequences([in_seq],
34                                  maxlen=self.max_length, padding='post')[0]
35           out_seq = tf.keras.utils.to_categorical([out_seq],
36                                                   num_classes=vocab_size)[0]
37           X_images.append(img_emb)
38           X_seq.append(in_seq)
39           y.append(out_seq)
40       return [np.array(X_images), np.array(X_seq)], np.array(y)
```

```
1 train_captions = train_df["comment"]
2 val_captions = val_df["comment"]
3 test_captions = test_df["comment"]
```

```
1 # Tokenizer instance for text encoding
2 tokenizer = Tokenizer()
3
4 # Instance is trained on the training captions
5 tokenizer.fit_on_texts(train_captions)
6
7 # word_index provides dictionary for encodings so
8 # so it can be used to define size of our vocabulary
9 vocab_size = len(tokenizer.word_index) + 1
```

Using the tokenizer defined above, we convert our lists of captions into encoded lists of captions. Note that the Tokenizer is trained on the training captions and it is then used below to convert sequences of texts into sequences of encodings for the model

```
1 train_sequences = tokenizer.texts_to_sequences(train_captions)
2 val_sequences = tokenizer.texts_to_sequences(val_captions)
3 test_sequences = tokenizer.texts_to_sequences(test_captions)
```

```
1 batch_size = 64
2 max_length = 20 # chosen from observation
```

Finally, a generator can be constructed from each of training, validation, and testing data. These generators will be used during the model training and evaluation

```python
1 # Generators can now be passed to model.fit()
2 train_generator = ImageCaptionGenerator(
3     np.squeeze(train_image_embeddings, axis=1),
4     train_captions, tokenizer, max_length, batch_size)
5
6 val_generator = ImageCaptionGenerator(
7     np.squeeze(val_image_embeddings, axis=1),
8     val_captions, tokenizer, max_length, batch_size)
9
10 test_generator = ImageCaptionGenerator(
11     np.squeeze(test_image_embeddings, axis=1),
12     test_captions, tokenizer, max_length, batch_size)
```

```python
1 vocab_size = len(tokenizer.word_index) + 1
2 vocab_size #our vocabulary size
```

```
2602
```

## ⌄ 5. The caption generation model

```python
1 embedding_dim = 300 # Embedding dimension for the text data
2 img_feature_dim = 4096  # Feature dimensions from VGG16
3 units = 512 # Number of LSTM units
4
5 inputs1 = Input(shape=(4096,))
6 fe1 = Dropout(0.25)(inputs1) # Dropout to prevent overfitting
7 fe2 = Dense(units, activation='relu')(fe1)
8
9 inputs2 = Input(shape=(max_length,))
10 se1 = Embedding(vocab_size, embedding_dim, mask_zero=True)(inputs2)
11 se2 = Dropout(0.25)(se1) # Dropout to prevent overfitting
12 se3 = LSTM(units, recurrent_dropout=0.35)(se2)
13
14 # Combining text and image input for classification
15 # This allows the fully connected layers to take both image and
16 # text into account when generating the next token in the caption
17 decoder1 = Add()([fe2, se3])
18 decoder2 = Dense(units, activation='relu')(decoder1)
19 outputs = Dense(vocab_size, activation='softmax')(decoder2)
20
21 model = Model(inputs=[inputs1, inputs2], outputs=outputs)
```

```
22 model.compile(loss="categorical_crossentropy", optimizer='adam')
23 model.summary()
```

WARNING:tensorflow:Layer lstm_1 will not use cuDNN kernels since it doesn't me
Model: "model"

| Layer (type) | Output Shape | Param # | Connected |
|---|---|---|---|
| input_4 (InputLayer) | [(None, 20)] | 0 | [] |
| input_3 (InputLayer) | [(None, 4096)] | 0 | [] |
| embedding_1 (Embedding) | (None, 20, 300) | 780600 | ['input_4 |
| dropout_2 (Dropout) | (None, 4096) | 0 | ['input_3 |
| dropout_3 (Dropout) | (None, 20, 300) | 0 | ['embeddi |
| dense_1 (Dense) | (None, 512) | 2097664 | ['dropout_ |
| lstm_1 (LSTM) | (None, 512) | 1665024 | ['dropout_ |
| add (Add) | (None, 512) | 0 | ['dense_1 'lstm_1[( |
| dense_2 (Dense) | (None, 512) | 262656 | ['add[0][( |
| dense_3 (Dense) | (None, 2602) | 1334826 | ['dense_2 |

```
Total params: 6140770 (23.43 MB)
Trainable params: 6140770 (23.43 MB)
Non-trainable params: 0 (0.00 Byte)
```

## 6. Training

Finally, the caption generation model can be trained. The optimal model weights will be stored in "model_weights.h5". Here, the optimal model weights will be defined as the model weights that achieved the lowest loss on the validation data. These weights will be restored once the training is completed. Moreover, the model will monitor validation loss and if there are no improvements in the validation loss for 3 epochs, early stopping is invoked to prevent overfitting to the training data.

```
1 from tensorflow.keras.callbacks import EarlyStopping, ModelCheckpoint
2
3 early_stopping = EarlyStopping(
4     monitor='val_loss',
5     patience=3,
6     verbose=1,
7     restore_best_weights=True)
8
9 checkpoint = ModelCheckpoint(
10     'model_weights.h5',
11     save_best_only=True,
12     save_weights_only=True,
13     monitor='val_loss',
14     verbose=1)
15
16 # Model is trained for a maximum of 50 epochs
17 model_hist = model.fit(train_generator, validation_data = val_generator,
18                        epochs=50, batch_size=64,
19                        callbacks=[early_stopping, checkpoint])
```
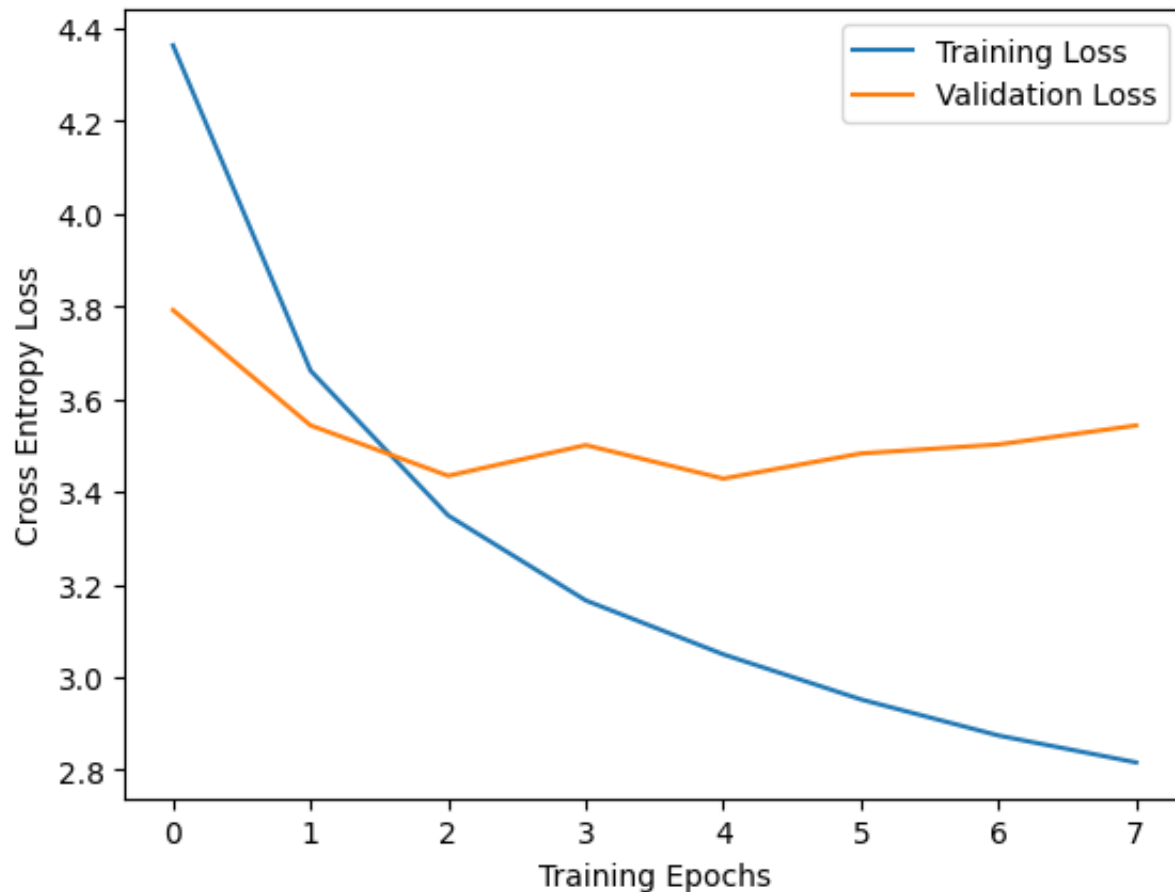
```
Epoch 1/50
5013/5017 [============================>.] – ETA: 0s – loss: 4.3631
Epoch 1: val_loss improved from inf to 3.79195, saving model to model_weights.
5017/5017 [=============================] – 88s 16ms/step – loss: 4.3630 – va
Epoch 2/50
5014/5017 [============================>.] – ETA: 0s – loss: 3.6609
Epoch 2: val_loss improved from 3.79195 to 3.54319, saving model to model_weig
5017/5017 [=============================] – 74s 15ms/step – loss: 3.6612 – va
Epoch 3/50
5017/5017 [=============================] – ETA: 0s – loss: 3.3488
Epoch 3: val_loss improved from 3.54319 to 3.43446, saving model to model_weig
5017/5017 [=============================] – 75s 15ms/step – loss: 3.3488 – va
Epoch 4/50
5017/5017 [=============================] – ETA: 0s – loss: 3.1657
Epoch 4: val_loss did not improve from 3.43446
5017/5017 [=============================] – 72s 14ms/step – loss: 3.1657 – va
Epoch 5/50
5015/5017 [============================>.] – ETA: 0s – loss: 3.0494
Epoch 5: val_loss improved from 3.43446 to 3.42844, saving model to model_weig
5017/5017 [=============================] – 71s 14ms/step – loss: 3.0493 – va
Epoch 6/50
5013/5017 [============================>.] – ETA: 0s – loss: 2.9516
Epoch 6: val_loss did not improve from 3.42844
5017/5017 [=============================] – 73s 15ms/step – loss: 2.9521 – va
Epoch 7/50
5016/5017 [============================>.] – ETA: 0s – loss: 2.8742
Epoch 7: val_loss did not improve from 3.42844
5017/5017 [=============================] – 76s 15ms/step – loss: 2.8743 – va
Epoch 8/50
5014/5017 [============================>.] – ETA: 0s – loss: 2.8164Restoring r

Epoch 8: val_loss did not improve from 3.42844
5017/5017 [=============================] – 73s 15ms/step – loss: 2.8162 – va
Epoch 8: early stopping
```

## ⌄ 7. Visualizing loss

Visualizing the validation and train loss obtained from running our model for the above epochs:

```
1 from matplotlib import pyplot as plt
2
3 plt.plot(model_hist.history['loss'], label='Training Loss')
4 plt.plot(model_hist.history['val_loss'], label='Validation Loss')
5 plt.legend()
6 plt.xlabel("Training Epochs")
7 plt.ylabel("Cross Entropy Loss")
8 plt.title('Training and Validation loss per epoch')
9 plt.show()
```



From the above plot, we observe a clear decrease in the training loss, and a decrease and then slight increase towards the end for the validation loss. As reflected in the early stopping, the model gradually overfits to the training data, and is prevented from doing that by our callbacks.

## ⌄ 8. Results: Generating captions

Since the optimal model weights were saved, they can now be loaded and used immediately for evaluation

```
1 %cd 'Code Submission'
```

    /content/drive/MyDrive/COMS4995AML_Project/Code Submission

```
1 model.load_weights('model_weights.h5')
```

```
1  def generate_caption(model, tokenizer, photo_feature, max_length):
2    # This is a special token to indicate start of the generated caption
3    in_text = 'startseq'
4
5    # Expanding dimension formats input image into a batch of 1
6    # which is required for model since it was trained on batches
7    photo_feature = np.expand_dims(photo_feature, axis=0)
8
9    for i in range(max_length):
10     # Encoding the sequence of tokens
11     sequence = tokenizer.texts_to_sequences([in_text])[0]
12     # Padding until max_length is reached
13     sequence = pad_sequences([sequence], maxlen=max_length)
14
15     # Next token is generated by the model prediction
16     yhat = model.predict([photo_feature, sequence], verbose=0)
17     # The most likely token is chosen
18     yhat = np.argmax(yhat, axis=-1)
19
20     word_index = int(yhat[0])
21     # The predicted token is mapped back to a word based on the tokenizer
22     word = tokenizer.index_word.get(word_index, '')
23
24     # Generation stops when the special token for end of sequence
25     # or the empty string is seen
26     if word == '' or word == 'endseq':
27       break
28
29     # Predicted token is appended to our caption
30     in_text += ' ' + word
31
32   #Removing the startseq and endseq
33   final_caption = in_text.split()[1:-1]
34   final_caption = ' '.join(final_caption)
35   return final_caption
```

Below, a few example captions are provided alongside with the reference captions and images from the dataset. The BLEU metric is used to provide a quantitative measure of the quality of the caption

Example 1:

```
 1 from PIL import Image
 2
 3 # Since image embeddings have already been calculated, they can be used
 4 # directly for prediction without need for more pre-processing
 5 photo_feature = np.squeeze(test_image_embeddings, axis = 1)[300]
 6 caption = generate_caption(model, tokenizer, photo_feature, max_length)
 7 actual_captions = \
 8 list(test_df[test_df["image_name"] == test_df["image_name"][300]]["comment"])
 9 print("Generated Caption:", caption)
10 print("Actual Caption: ", actual_captions, "\n")
11 print("Bleu Score: ", sentence_bleu(actual_captions, caption), "\n")
12 Image.open("../Images/" + test_df["image_name"][300])
```

```
Generated Caption: horse jumping rider in a rodeo rodeo rodeo a rodeo a rodeo
Actual Caption:  ['a bucking brown horse next to a falling cowboy']

Bleu Score:  0.16840394156133062
```



Example 2:

```
1 photo_feature = np.squeeze(test_image_embeddings, axis = 1)[2]
2 caption = generate_caption(model, tokenizer, photo_feature, max_length)
3 actual_captions = \
4 list(test_df[test_df["image_name"] == test_df["image_name"][2]]["comment"])
5 print("Generated Caption:", caption)
6 print("Actual Caption: ", actual_captions, "\n")
7 print("Bleu Score: ", sentence_bleu(actual_captions, caption), "\n")
8 Image.open("../Images/" + test_df["image_name"][2])
```

```
Generated Caption: brown dog running through grass with a orange ball in its m
Actual Caption:  ['the little brown dog runs past another dog on the grass']

Bleu Score:  0.2821075432377299
```



## Example 3:

```
1 photo_feature = np.squeeze(test_image_embeddings, axis = 1)[1876]
2 caption = generate_caption(model, tokenizer, photo_feature, max_length)
3 actual_captions = \
4 list(test_df[test_df["image_name"] == test_df["image_name"][1876]]["comment"])
5 print("Generated Caption:", caption)
6 print("Actual Captions: ")
7 for cap in actual_captions:
8   print("  ", cap)
```

```
 9 bleu_scores = []
10 for i in range(len(actual_captions)):
11   bleu_scores.append(sentence_bleu([actual_captions[i]], caption))
12 print("Average Bleu Score: ",
13      sum(bleu_scores)/len(actual_captions), "\n")
14 print("Highest Bleu Score: ", max(bleu_scores), "\n")
15 Image.open("../Images/" + test_df["image_name"][1876])
```

```
Generated Caption: dog running through the water with a stick in its mouth is
Actual Captions:
    a spotted black and white dog splashes in the water
    a black and white dog is running though water whilst bearing its teeth
    a black and white dog is running and splashing in water
Average Bleu Score:  0.28953074057302325

Highest Bleu Score:  0.38403413539641235
```



1