
Image Caption Generator

Jasmeetsingh Khalsa

Cpts 534 Neural Network Design and Application

Final Report

Spring 2022

11743550

1 Introduction

We know how easy it is for our minds to tell what a given image is about, but can a computer tell what the image is representing? Computer vision researchers worked on this a lot and they considered this impossible till now! With the help of deep learning techniques, availability of huge datasets and computer power, we can build models that can generate captions for an image.

In this project we will see a Image Caption Generator which is a challenging Neural Network problem where a textual description must be generated for a given photograph. It requires both methods from computer vision to understand the content of the image and a language model from the field of natural language processing to turn the understanding of the image into words in the right order. The primary aim would be to build a Network model that would generate captions for an input image which explains the context of the image. To build the model, we will use basic Neural Networks such as Convolutional Neural Networks and LSTMs (a type of Recurrent Neural Network) along with a technique called attention mechanism which recently played an important role in computer vision and is recently widely used in image caption generation tasks.

2 Literature Review

Recently the development of image description system has drawn increasing attention and it has become one of the most important topics in Computer Vision. Early image description generation methods aggregate image information using static object class libraries in the image and modeled using statistical language models. In paper 'Generating image descriptions using dependency relational patterns', Aker and Gaizauskas [1] use a dependency model to summarize multiple web documents containing information related to image locations and propose a method for automatically tagging geotagged images. In paper 'Corpus-Guided Sentence Generation of Natural Images' Yang et al. [2] proposed a language model trained from the English Gigaword corpus to obtain the estimation of motion in the image and the probability of collocated nouns, scenes, and prepositions and used these estimates as parameters of the hidden Markov model. Here, the image description was obtained by predicting the most likely nouns, verbs, scenes, and prepositions that make up the sentence.

Relation to our work: The methods described in the above papers are brainstorming and have their own characteristics, but all have the common disadvantage that none of them make any intuitive feature observations on objects or actions in the image, nor do they give an end-to-end mature general model to solve this problem. This is what we try to achieve through our work in this project. We will propose a model that will make feature observations on objects and will extract useful information for prediction.

33 The initial step of the project is to create a baseline model using a simple combination of CNN along
34 with a RNN [3]. This technique will successfully generate captions for the input image. This is
35 a classic application of Neural networks. However, the problem with a 'classic' image captioning
36 model is that when the model is trying to generate the next word of the caption, this word is usually
37 describing only a part of the image. It is unable to capture the essence of the entire input image. Using
38 the whole representation of the image h to condition the generation of each word cannot efficiently
39 produce different words for different parts of the image. This is exactly where deep learning technique
40 called Attention mechanism is helpful.

41 **3 Dataset**

42 Out of the vast varieties of datasets available for this problem like Flickr 8k , Flickr 30k, MS COCO,
43 etc. The Flickr8K dataset will be used for the model training of image caption generators. Flickr8k
44 image comes from Yahoo's photo album site Flickr, which contains 8,000 photos, 6000 image training,
45 1000 image verification, and 1000 image testing.

46 This dataset is readily available and can be directly downloaded from the internet. The downloading
47 process takes some time due to the large size(1GB) of the dataset. 8091 images are stored inside the
48 Flickr_8k_Dataset folder and the text files with captions of images are stored in the Flickr_8k_text
49 folder. In short the flicker dataset consists of:

- 50 1. Flickr_8k_Dataset: Folder with 8091 images.
- 51 2. Flickr_8k_text: Folder with Text files with captions of images

52 **4 Baseline**

53 A basic image captioning system would encode the image, using a pre-trained Convolutional Neural
54 Network(ENCODER) that would produce a hidden state h . Then, it would decode this hidden state
55 by using a LSTM(DECODER) and generate recursively each word of the caption.

56 To build a baseline image caption generator model we would simply merge CNN with LSTM.
57 Therefore:

58 Image Caption Generator Model(CNN-RNN model) = CNN + LSTM.

- 59 • CNN- To extract features from the image. A pre-trained model called Xception is used for
60 this.
- 61 • LSTM- To generate a description from the extracted information of the image.

62 First, we write functions to load the datasets for model training. Here 3 functions are primarily written,
63 1. load_photos(fname), 2. load_clean_descriptions(fname, image), 3. load_features(photos).

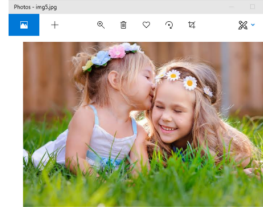
64 Then we tokenize the vocabulary to a numerical representation for machines to understand. This is
65 done by mapping each word of vocabulary with a separate unique index value.

66 Next step is to define a CNN-RNN model. Here, we used Keras model in order to define the structure
67 of the model. This includes:

- 68 1. A feature extractor with a dense layer, which will extract the feature from the images of size
69 2048 and will decrease the dimensions to 256 nodes.
- 70 2. A sequence processor followed by the LSTM layer, the textual input is handled by this
71 embedded layer.
- 72 3. Decoder which will merge the output of the above two layers and process the dense layer to
73 make the final prediction.

74 Now, we train the Image Caption Generation Model. Here, we generate the input and output sequences
 75 to train our model with 6000 training images. We create a function named `model.fit_generator()` to fit
 76 the batches to the model.

77 Finally comes testing, Here the task is to test the model accuracy by inputting test image data. Here,
 78 we see that the proposed baseline successfully generates captions for the test image.



```

_core\python\framework\indexed_slices.py:424: UserWarning: Converting sparse Index
edSlices to a dense Tensor of unknown shape. This may consume a large amount of me
mory.
  "Converting sparse IndexedSlices to a dense Tensor of unknown shape. "
start two girls are playing in the grass end

```

81 5 Main Approach

82 The code for the main approach is available in Github. I tried to keep the Colab notebook as
 83 descriptive as possible so that anyone can understand the purpose of each cell.

84 5.1 Concept of Attention Mechanism

85 In Attention Mechanism, the image is first divided into n parts, and we compute with a CNN
 86 representations of each part h_1, \dots, h_n . When the RNN is generating a new word, the attention
 87 mechanism is focusing on the relevant part of the image, so the decoder only uses specific parts of
 88 the image.

89 There are two types of Attention Mechanism:

- 90 1. Global Attention: Attention is placed on all source positions.
- 91 2. Local Attention: Attention is placed only on a few source positions.

92 **Global attention** takes into consideration all encoder hidden states to derive the context vector ($c(t)$).
 93 In order to calculate $c(t)$, we compute $a(t)$ which is a variable length alignment vector. The alignment
 94 vector is derived by computing a similarity measure between h_t and \bar{h}_s where h_t is the source hidden
 95 state while \bar{h}_s is the target hidden state. Similar states in encoder and decoder are actually referring
 96 to the same meaning.

97 We use a score function which is a content-based function using which we calculate the similarity
 98 between the hidden states of the target and the source.

$$score(h_t, \bar{h}_s) = \begin{cases} h_t^T \bar{h}_s \\ h_t^T W_a \bar{h}_s \\ v_a^T \tanh(W_a [h_t; \bar{h}_s]) \end{cases}$$

99 As Global attention focus on all source side words for all target words, it is computationally very
 100 expensive and is impractical when translating for long sentences. To overcome this deficiency **Local**
 101 **attention** chooses to focus only on a small subset of the hidden states of the encoder per target word.

$$f_{ATT} = V_{attn}^T * \tanh(U_{attn} * h_j + W_{attn} * s_t)$$

102 Where,

$$V_{\text{attn}}^T \in R^d, U_{\text{attn}} \in R^{d \times d}, W_{\text{attn}} \in R^{d \times d}, s_{t-1} \in R^d, h_j \in R^d$$

103 Now, general Score is given as:

$$e_{jt} = f_{\text{ATT}}(s_{t-1}, h_j)$$

104 Where, e_{jt} means at every t^{th} timestep of ENCODER, how important j^{th} is the pixel location in the
 105 input image. s_{t-1} is previous state of DECODER. h_j is state of ENCODER. f_{ATT} is a simple Feed
 106 Forward Neural Network, which is a linear transformation of input ($U_{\text{attn}} * h_j + W_{\text{attn}} * s_t$) and
 107 then a non-linearity(tanh) on top of that, and then again one more transformation(V_{attn}^T). This is a
 108 scalar quantity.

$$f_{\text{ATT}} = V_{\text{attn}}^T * \tanh(U_{\text{attn}} * h_j + W_{\text{attn}} * s_t)$$

109 Where,

$$V_{\text{attn}}^T \in R^d, U_{\text{attn}} \in R^{d \times d}, W_{\text{attn}} \in R^{d \times d}, s_{t-1} \in R^d, h_j \in R^d$$

110 To get the Probability distribution we do Softmax:

$$\alpha_{jt} = \text{Softmax}(e_{jt})$$

111

$$\text{i.e. } \alpha_{jt} = \frac{e^{e_{jt}}}{\sum_{k=1}^{T_x} e^{e_{kt}}}, \text{ such that } \sum_{j=1}^{T_x} \alpha_{jt} = 1 \text{ and } \alpha_{ij} \geq 0$$

112 Now, when we know the input, we need to feed Weighted sum combination of input to Decoder.

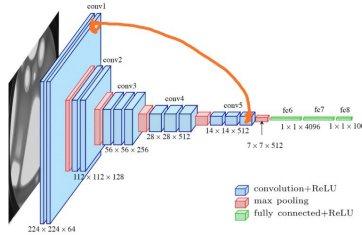
$$C_t = \sum_{j=1}^T \alpha_{jt} h_j \text{ such that } \sum_{j=1}^T \alpha_{jt} = 1 \text{ and } \alpha_{ij} \geq 0$$

113 where, C_t is a context vector, i.e. the Weighted sum of input.

$$S_t = \text{RNN}(S_{t-1}, [e(\hat{y}_{t-1}), C_t])$$

114 where, S_{t-1} is previous state of Decoder. $e(\hat{y}_{t-1})$ is previous predicted word. C_t is the context
 115 vector, i.e. the Weighted sum of input.

116 Usually for images, we use representations from one of the fully connected layers. But consider
 117 a image where a man is throwing a frisbee. When we say the word 'man' that means we need to
 118 focus only on man in the image, and when we say 'throwing' then we have to focus on his hand
 119 in the image. Similarly, when we say 'frisbee' we have to focus only on the frisbee in the image.
 120 This means 'man', 'throwing' and 'frisbee' comes from different pixels in image. Now consider a
 121 VGG-16 (pre-trained model used for main approach) representation of man throwing a frisbee below,



122

123 Looking at the figure, we can see that the convolution layers do not exactly contain any information
 124 in it. However, every location of convolution layers corresponds to some location of image. For
 125 example, Consider the output of 5th convolution layer of VGGNet is a 14*14*512 size feature map.
 126 This 5th layer has 14*14 pixel locations which corresponds to certain portion in image, that means
 127 we have in total 196 such pixel locations. And finally, we can treat these 196 locations (each having
 128 512 dimensional representation). The model will then learn an attention over these locations (Which
 129 in turn corresponds to actual locations in images).

130 **Description of code:**

131 In code, we start with creating the utility functions. The purpose of creating the utility functions
132 is to: 1. Load the files. 2. To clean data i.e removing punctuations, single characters, numerical
133 values from text. Next we create dataframe from the raw data and perform some basic exploratory
134 data analysis. After this we performed preprocessing of images and captions. In preprocessing, we
135 are preprocessing the captions (adding '<start>' and '<end>' tags to every caption), so that the ML
136 model understands the starting and ending of each caption.

137 Next step is to define the pre-trained image model (VGG-16). The VGG16 model was pre-trained
138 on the ImageNet data-set for classifying images. The VGG model contains a convolutional part and
139 a fully connected part which is used for the image classification. Then, we prepared the images
140 and created the image dataset i.e. reshape every image to 224*224*3 shape before feeding it to
141 VGG16 model. For captions, we performed tokenizations where we tokenize the captions and created
142 vocabulary of words present in our data corpus. Then we created vector notations for each word in
143 our vocabulary. For words not appearing in the vocabulary, we gave <unk> notation. After getting
144 sequences to the words in captions, the sequences are of varied length. So, we need pad sequences
145 to pad with the maximum length of the captions. To perform a train-test split, we split the dataset
146 (images and captions) into 80:20 ratio i.e.[train:test]

147 **Implementing Attention Mechanism:** The entire step-by-step process of applying Attention ac-
148 cording to Minh-Thang Luong's paper: "Effective Approaches to Attention-based Neural Machine
149 Translation":

150 **Global Attention:**

- 151 1. Producing the Encoder Hidden States: Encoder produces hidden states of each element in
152 the input sequence.
- 153 2. Decoder RNN: The previous decoder hidden state and decoder output is passed through the
154 Decoder RNN to generate a new hidden state for that time step.
- 155 3. Calculating Alignment Scores: Using the new decoder hidden state and the encoder hidden
156 states, alignment scores are calculated.
- 157 4. Softmaxing the Alignment Scores: the alignment scores for each encoder hidden state are
158 combined and represented in a single vector and subsequently softmaxed.
- 159 5. Calculating the Context Vector: the encoder hidden states and their respective alignment
160 scores are multiplied to form the context vector.
- 161 6. Producing the Final Output: The context vector is concatenated with the decoder hidden
162 state generated in step 2 as passed through a fully connected layer to produce a new output.
- 163 7. The process (Steps 2-6) repeats itself for each time step for the decoder until an token is
164 produced or output is past the specified maximum length.

165 **Local Attention:**

- 166 1. Producing the Encoder Hidden States: Encoder produces hidden states of each element in
167 the input sequence.
- 168 2. Calculating Alignment Scores: Between the previous decoder hidden state and each of the
169 encoder's hidden states are calculated.
- 170 3. Softmaxing the Alignment Scores: The alignment scores for each encoder hidden state are
171 combined and represented in a single vector and subsequently softmaxed
- 172 4. Calculating the Context Vector: The encoder hidden states and their respective alignment
173 scores are multiplied to form the context vector.
- 174 5. Decoding the output: The context vector is concatenated with the previous decoder output
175 and fed into the Decoder RNN for that time step along with the previous decoder hidden
176 state to produce a new output.

177 6. The process (Steps 2-5) repeats itself for each time step for the decoder until an token is
178 produced or output is past the specified maximum length.

179 Finally comes, Selecting Optimizer, defining Loss Function and setting checkpoints. In training,
180 The Encoder output, the hidden state(initialized to 0) and the Decoder input(which is the <start>
181 token) are passed to the Decoder. The Decoder returns predictions and the Decoder hidden state. The
182 Decoder hidden state is then passed back into the model and the predictions are used to calculate the
183 loss. While training, we use the 'Teacher Forcing Technique' to decide the next input of the Decoder.
184 Teacher Forcing is a technique where the target word is passed as the next input to the Decoder.
185 This technique helps to learn the correct sequence or correct statistical properties from the sequence,
186 quickly. Final step is to calculate the gradient and apply it to the optimizer and backpropagate.

187 Testing step is similar to training step, it is just that we do not update the gradients, and provide the
188 predicted output as decoder input to next RNN cell at next time steps. Test step is required to find out
189 whether the model built is overfitting or not.

190 6 Evaluation Metric

191 **The evaluation of Captioning Model:** The evaluation function is similar to the training loop, except
192 we don't use Teacher forcing here. The input to Decoder at each time step is its previous predictions,
193 along with the hidden state and the Encoder output.

194 The following are the two methods to evaluate the captions:

195 1. Greedy Approach

196 2. Beam Search

197 **Greedy Approach:** Is also called as Maximum Likelihood Estimation (MLE) i.e. we select that
198 word which is most likely according to the model for the given input. And sometimes this method is
199 also called as Greedy Search, as we greedily select the word with maximum probability.

200 **Beam Search:** Here we take top k predictions, feed them again in the model and then sort them
201 using the probabilities returned by the model. So, the list will always contain the top k predictions. In
202 the end, we take the one with the highest probability and go through it till we encounter *< end >* or
203 reach the maximum caption length.

204 To evaluate the results obtained by the main approach, we will use BLEU Score. BLEU measure will
205 help to evaluate the result of the test set generated captions. BLEU is simply taking the fraction of
206 n-grams in the predicted sentence that appears in the ground-truth.

207 BLEU is a well-acknowledged metric to measure the similarity of one hypothesis sentence to multiple
208 reference sentences. Given a single hypothesis sentence and multiple reference sentences, it returns
209 value between 0 and 1. The metric close to 1 means that the two are very similar.

210 7 Results and Analysis

211 In conclusion, in this project we have successfully seen that baseline model works perfectly fine
212 and it generates proper captions without applying the attention mechanism. After completing the
213 implementation of main approach, we saw that the model generates captions with improved accuracy
214 and scores.

215 So in all, we can say that the native first-cut model in main approach, without any rigorous hyper-
216 parameter tuning does a decent job in generating captions for images. We must understand that
217 the images used for testing must be semantically related to those used for training the model. For
218 example, if we train our model on the images of cats, dogs, etc. We must not test it on images of air
219 planes, waterfalls, etc. This is an example where the distribution of the train and test sets will be very
220 different and in such cases no Neural Network model in the world will give good performance.

In our main approach, when we compare the results in code for Greedy search and Beam search. We see that Beam Search generated better results than Greedy Search. This might be because beam search makes two improvements over greedy search:

1. With Greedy search, we took just the single best word at each position. In contrast, Beam search expands this and takes the best "N" words.
2. With Greedy search, we consider each position in isolation. Once we had identified the best word for that position, we did not examine what came before it (i.e. in the previous position), or after it. In contrast, Beam search picks the "N" best sequences so far and considers the probabilities of the combination of all of the preceding words along with the word in the current position.

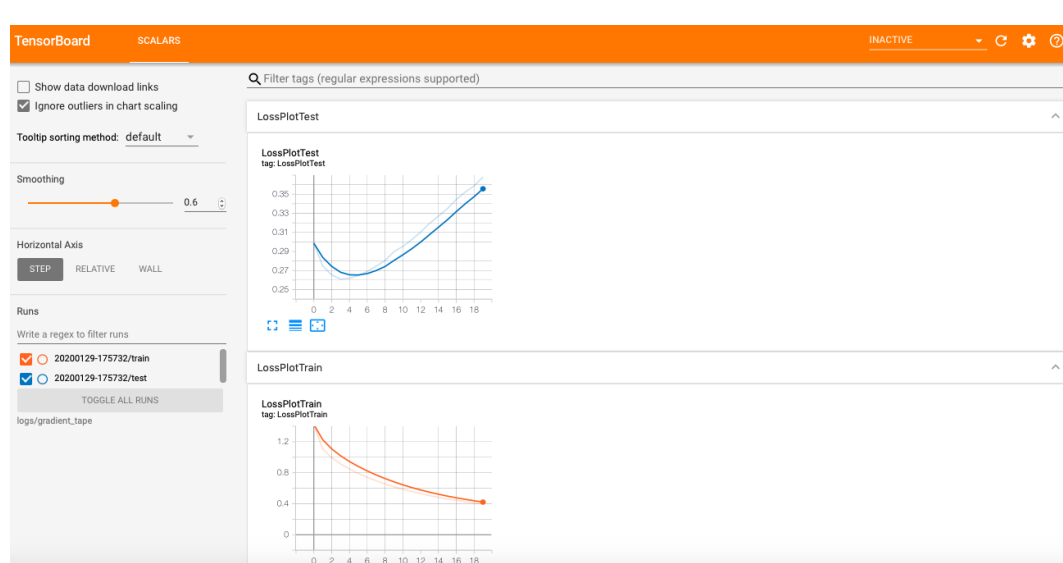
Performance of baseline compared to main approach:

The baseline is short and simple to setup and shows a reasonable chance of providing decent results. The implementation of baseline shows us how hard the problem of image caption generation is. The accuracy scores of BLEU scores shows this difficulty.

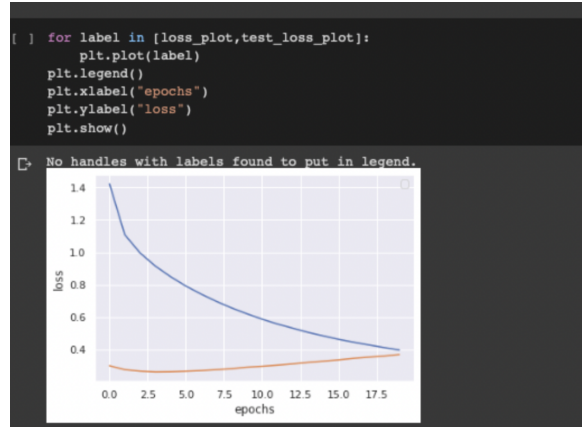
Looking at the BLEU scores of our Baseline and Main Approach, we can say that in average the baseline's BLEU score is 35% comparing this to the average BLEU scores of beam search for all k predictions in main approach i.e. 71.86%. This significant difference in scores of both baseline and main approach shows the power of learning certain features of an image then training a model on those features instead of training a model on a random part of an image. This proves the overall aim of this project.

8 Error Analysis

When training our model in main approach, we used tensorboard logs to visualize changes in our model and training. Tensorboard logs is a tool for providing the measurements and visualizations needed during the machine learning workflow. It enables tracking experiment metrics like loss and accuracy, visualizing the model graph, projecting embeddings to a lower dimensional space.



We plotted the train and test losses to check for overfitting. The following graph shows overfitting. This might be due to less training data. Training the model with larger datasets like MS-COCO or Flickr30 could help us solve this issue.



251

252 9 Future Work

253 There could be many future possible applications for this applications. Image Caption Generators
 254 could be used for recommendations in editing applications, usage in virtual assistants, for image in-
 255 dexing, for visually impaired persons, for social media, and several other natural language processing
 256 applications.

257 This is just a particular first-cut solution to the Image Caption Generator and a lot of future modifica-
 258 tions can be made to improve this solution. Improvements like:

- 259 • Using a larger dataset to improve the accuracy and performance
- 260 • Changing the model architecture. Such as adding BatchNormalization Layer, Dropouts,
 261 etc.)
- 262 • Doing more hyper parameter tuning (learning rate, batch size, number of layers, number of
 263 units, dropout rate, batch normalization etc.)
- 264 • Keep researching on this topic and optimizing the solution even further.
- 265 • API'fy this model using web frameworks like FLASK or DJANGO and deploy it in AWS.

266 10 Link to Code

267 Both Baseline and the Main Approach is available on my Github page.

268 **Link:** https://github.com/Foster1466/CPTS_534-Neural-Network-Design-and-Application

269

270

271 As datasets are too large to download, I have provided links of dataset in both notebooks. Just for
 272 simplicity, I am also providing the links here:

273 https://github.com/jbrownlee/Datasets/releases/download/Flickr8k/Flickr8k_Dataset.zip

274 https://github.com/jbrownlee/Datasets/releases/download/Flickr8k/Flickr8k_text.zip

275

276

277

278

References

1. A. Aker and R. Gaizauskas, “Generating image descriptions using dependency relational patterns,” in Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics, vol. 49, no. 9, pp. 1250–1258, Uppsala, Sweden, July 2010.
2. Y. Yang, C. L. Teo, H. Daume, and Y. Aloimonos, “Corpus-guided sentence generation of natural images,” in Proceeding of the Conference on Empirical Methods in Natural Language Processing, pp. 444–454, Edinburgh, UK, July 2011.
3. Global Attention : <https://arxiv.org/pdf/1508.04025.pdf>
4. Local Attention : <https://arxiv.org/pdf/1502.03044.pdf>
5. BLEU score: <https://machinelearningmastery.com/calculate-bleu-score-for-text-python/>
6. AppliedAiCourse : <https://www.appliedaicourse.com/lecture/11/applied-machine-learning-online-course/4150/attention-models-in-deep-learning/8/module-8-neural-networks-computer-vision-and-deep-learning>
7. Neural Machine Translation(Research Paper) : <https://arxiv.org/pdf/1409.0473.pdf>
8. Image Captioning with Visual attention: <https://www.tensorflow.org/tutorials/text/imagecaptioning>
9. Blogs : <https://towardsdatascience.com/intuitive-understanding-of-attention-mechanism-in-deep-learning-6c9482aecf4f>
10. Deep Learning Lectures by Prof. Mitesh M Khapra(IIT Madras)
11. Intuitive Understanding of attention mechanism in deep learning: <https://towardsdatascience.com/intuitive-understanding-of-attention-mechanism-in-deep-learning-6c9482aecf4f>
12. CS231n: Video by Andrej Karpathy on Image Captioning : <https://www.youtube.com/watch?v=NfnWJUyUJYUlist=PLkt2uSq6rBVctENoVBg1TpCC7OQi31A1C>