

000  
001  
002  
003  
004  
005  
006  
007  
008  
009  
010  
011  
012  
013  
014  
015  
016  
017  
018  
019  
020  
021  
022  
023  
024  
025  
026  
027  
028  
029  
030  
031  
032  
033  
034  
035  
036  
037  
038  
039  
040  
041  
042  
043  
044  
045  
046  
047  
048  
049  
050  
051  
052  
053

---

# Novel Image Captioning

---

**Daniel Thirman**  
CS 224D  
DThirman@Stanford.edu

## Abstract

In this paper I describe a model which is used to generate novel image captions for a previously unseen image by using a combination of a recurrent neural network and a convolutional neural network. This model is trained on the flickr30k and MS-COCO datasets of images and captions and scores a perplexity that is comparable to that of state of the art implementations. The network is evaluated by computing its perplexity as a function of how well the language model scores the sentence and how likely the sentence is given the image. In this paper, I present three different implementation a baseline model that learns its own embeddings for the vocabulary along with a basic RNN, a model that uses the pre-trained GloVe word vectors, as well as model that uses the GloVe word vectors as well as using Gated Recurrent Units to solve the vanishing gradient problem in the network.

## 1 Introduction

There are many applications for being able to automatically obtain a description for a picture. Some of these include image recognition for autonomous robotic systems which need to be able to recognize and process what they see, creating an image database that is search-able by keywords without having to manually tag and describe the different images that are added to the database, as well as many other possible applications. These applications are becoming increasingly relevant as technology processes so being able to determine the sentence level description of an image is a very important and interesting task. In the past several years, there have been many attempts by top researchers in the field to solve the problem and they have implemented models that are able to produce very good results that create natural sentences that describe the content of the images very well. Many well known companies, such as Microsoft and Google, have published models recently that attempt to solve the problem. In this paper I create a model that is inspired by other researchers' work on the problem and compare my results to the state of the art results in the field.

There have been some attempts that rather than generating new captions to describe the image instead choose to try find the best caption in their database that matches to the image. While this method is guaranteed to express captions that our naturally written and clear to understand, it fails to describe images that have unique combinations of objects or objects presented in an unusual way. Novel image captions are captions that are generated by the model from a combination of the image features and a language model instead of matching to an existing captions. Generating novel image captions solves both of the problems of using existing captions and as such is a much more interesting and useful problem.

The model I chose to implement to solve this problem is a multimodal neural network composed of a convolutional neural network and a recurrent neural network. The CNN is used to determine the image features and the RNN is used to generate the language model. These two networks are then combined with a weighted sum and the result is then used to predict the next word in the sentence that is used to describe the image that is inputted. I implemented and ran tests on three different versions of this model. The first is the baseline model which uses a basic RNN and learns the vector embeddings for the vocabulary by itself, the second is a model that uses the pre-trained GloVe word

054 vectors instead of learning them itself, and the final model continues to use the GloVe word vectors  
055 but also uses Gated Recurrent Units in the RNN to increase performance and alleviate the vanishing  
056 gradient problem. Each successive model scores a lower perplexity and produces better results.

057 The data set that is used to train this network is a combination of MS-COCO and the Flickr30k  
058 data sets. Both data sets contain large amount of images and captions that can be used to train and  
059 evaluate the network.  
060

## 061 **2 Related Work**

062 This problem has been one of interest to many researchers in the past few years. There have been  
063 many attempts to solve the problem using a variety of approaches. The models that are implemented  
064 in this paper are inspired by many of the other work performed by other researchers in the field.

065 Recently, Microsoft has released the MS-COCO database for anyone to use to work on this problem  
066 and has set up and challenges and offered prizes to those who do well on the problem so there have  
067 been many people working on the same database obtaining strong results.  
068

## 069 **3 Approach**

### 070 **3.0.1 Baseline model**

071 The model that I chose to implement for this project is a multimodal neural network that is composed  
072 of a convolutional neural network that is used to detect images features as well as a recurrent neural  
073 network, that is trained as a language model and predicts the next word given the context of the  
074 previous words and the image features supplied by the CNN.  
075

076 For this project I chose focus on the RNN instead of on CNN, since the RNN is more relevant to  
077 natural language processing which is the topic of the project, as such I chose to make use of a pre-  
078 trained CNN that would be able to extract image features and I would not have to train one from  
079 scratch. The pre-trained CNN I chose to use was one called AlexNet. AlexNet is a state of the art  
080 CNN implementation that is used to predict the probabilities of whether classes of images are  
081 contained inside a given image. There are 1000 classes of objects and running AlexNet on an image  
082 will return a vector of size 1000 which is the probability of each class.  
083

084 The output of the CNN is then fed into the RNN to predict the next output word of the caption. The  
085 RNN is a single layer network with a hidden state of length 512. Different lengths were tried for the  
086 size of the hidden state however 512 produced better results than other lengths. At each time step  
087 the next hidden state is computed using a combination of the previous state, the word vector for the  
088 input word, and the CNN output for the image that the caption is being generated for. The equation  
089 for the next hidden state is equal to:

$$090 \text{State}_i = \sigma(\text{State}_{i-1} * H + \text{Word}_i * I + \text{Image} * N) + b$$

091 Where H is (Hidden size X Hidden size), State is the hidden state of the RNN, I is (Word  
092 Embedding Size X Hidden Size), Word is the word vector obtained by looking up the current input  
093 word in embedding, N is (Size of CNN output X Hidden Size), and Image is the output from the  
094 CNN, and b is the bias. A drawing of one node in the RNN is found in figure 1.  $\sigma$  is sigmoid  
095 function which is used as the non-linearity of the RNN.

096 The output of each time step is then projected to predict the probability of each word appearing next  
097 in the sentence description of the word. The final projection into the vocab size is computed in a  
098 multimodal layer that computes a weighted the outputs of the projections of the RNN and the CNN.  
099 The equation for the projection at each time step is equal to:

$$100 \text{output} = A * (\text{RNN}_i * U) + (1 - A) * (\text{CNN} * V) + b$$

101 Where A is a variable learned by the network that ensures calculates how much of the RNN  
102 or CNN output should be included in the final projection.  $\text{RNN}_i$  is the output of the RNN at time  
103

108  
109  
110  
111  
112  
113  
114  
115  
116  
117  
118  
119  
120  
121  
122  
123  
124  
125  
126  
127  
128  
129  
130  
131  
132  
133  
134  
135  
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146  
147  
148  
149  
150  
151  
152  
153  
154  
155  
156  
157  
158  
159  
160  
161

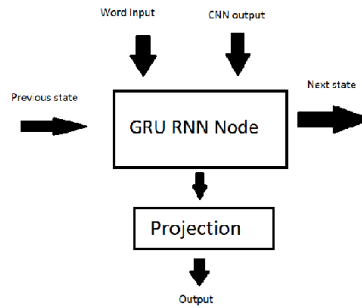


Figure 1: One node in the RNN

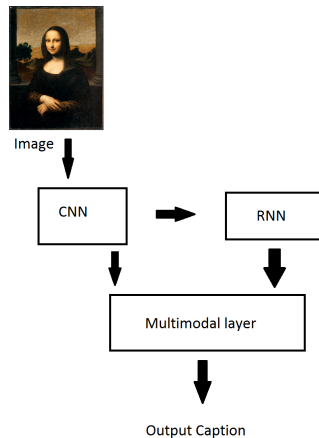


Figure 2: Entire multimodal neural network

step  $i$ ,  $U$  is (Hidden Size X Vocab Size),  $CNN$  is the output of the CNN for the given image,  $V$  is (Size of CNN output X Vocab Size) and  $b$  is the bias.

This weighted sum of RNN and CNN output ensures that both the image features and the language model both factor in the final output probability at each time step and the best word for both inputs is most likely to be selected. The structure of the whole network is shown in figure 2.

### 3.0.2 GloVe Vector based model

The previous model did not generate particularly good results and while its captions might capture some of the semantic meaning of the image the language model was relatively poor. I reasoned that this was because although I had lots of examples in my data set most captions were relatively short so there was not enough text to adequately train the word vectors that were being used in the embedding of the vocab. So, I decided to use the pre-trained GloVe word vector. GloVe word vectors are word vectors that are trained based off the co-occurrences of word pairs in a very large corpus. Since the corpus used to train them was substantially larger than the corpus of sentences used in the image captions I had, using these captions would result in a better representation of the word similarities and would create a better language model.

### 3.0.3 Gated Recurrent Unit Model

While the previous model did do a better job than the baseline model of expressing captions it still had its problems. Since the captions generated were not of insignificant length, one problem that this network suffered from was the vanishing gradient problem. In the vanishing gradient problem the values of the gradients diminish rapidly as they are back-propagated to previous states and they

162 have smaller and smaller changes on the variables in the states until the change is insignificant and  
 163 the variables at states a few steps away do not train at all. One solution to this is to use a gated  
 164 recurrent unit or a GRU. GRU's are a more complicated node structure than a base RNN in which  
 165 the value of the previous state is scored for how important it is in the next state and only as much  
 166 as is optimal is used to compute the next state. This helps alleviate the vanishing gradient problem  
 167 as the values of the gradients are no longer exponentiated at each step in back propagation and the  
 168 gradients no longer diminish so quickly. The new equation for the next hidden state is

$$169 \quad z = \sigma(State_{i-1} * H_z + Word_i * I_z + Image * N_z) + b_z$$

$$170 \quad r = \sigma(State_{i-1} * H_r + Word_i * I_r + Image * N_r) + b_r$$

$$171 \quad \hat{h} = \tanh(r * State_{i-1} * H + Word_i * I + Image * N) + b$$

$$172 \quad State_i = z * State_{i-1} + (1 - z) * \hat{h}$$

173 Here  $z$  is the update gate which is used to calculate how much of the new state should come from  
 174 the previous state and how much of it should come from the new value that it is calculating, and  $r$  is  
 175 the reset gate which calculates how much of the previous state should go into the new value that is  
 176 being calculated,  $\hat{h}$ . All other values and matrices have the same values and dimensions as they did  
 177 in the previous baseline model.

## 178 4 Experiment

### 179 4.1 Data set

180 The data set that I have used to train my model is a combination of the MS-COCO and Flickr30k data  
 181 sets. MS-COCO is a publicly available data set of images and captions provided by Microsoft that  
 182 is used to test and train network for exactly this purpose. The data set contains over 300 thousand  
 183 images and along with each image contains a sentence level description of the image and a list of  
 184 classes of objects that are contained in each image.

185 Flickr30k is a data set distributed by researchers at the University of Illinois. It contains 30 thousand  
 186 images taken from the website Flickr and contains multiple captions for each image for a total of  
 187 nearly 150,000 captions. By combining these two data sets I had nearly 500,000 captions and over  
 188 300,000 images to train the network on.

### 189 4.2 Evaluation

190 The model was scored by using the computing the perplexity based off the how well the sentence fit  
 191 the language model and how well it worked to describe the image that it had been given to caption.  
 192 The score of the language model was calculated based off of how similar successive word vectors  
 193 were to each other. Since, word that occur more closely to each other are more likely to have a  
 194 similar value since the vectors are created based off their co-occurences, so similar word vectors are  
 195 more likely to have a low perplexity. The score for how it fits the image is created by calculating  
 196 how likely it is for the word to be selected based off of the image. To accomplish this I went through  
 197 the data set and calculated likely each class of object is given that a given word is included in its  
 198 caption. Then to compute the perplexity based off how well each word fits the image, the likelihood  
 199 of class given a caption word is compared to the likelihood of the classes given the image. The  
 200 equation for calculating the loss is as follows:

$$201 \quad J(\theta) = - \sum_{i=1}^{\|V\|} y_i^{(t)} * \log(y_i^{(t)}) - \sum_{j=1}^{\|I\|} \sum_{i=1}^{\|V\|} P(y_i^{(t)} | Image_j) * \log(y_i^{(t)})$$

#### 202 4.2.1 Results

203 The perplexity scores for each of the different models is shown in the following table.

216  
217  
218  
219  
220  
221  
222  
223  
224  
225  
226  
227  
228  
229  
230  
231  
232  
233  
234  
235  
236  
237  
238  
239  
240  
241  
242  
243  
244  
245  
246  
247  
248  
249  
250  
251  
252  
253  
254  
255  
256  
257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269



Figure 3: Example images 1,2, and 3

Model	Train Perp.	Valid Perp.	Test Perp.
Baseline	65.23	11.18	8.3925
GloVe Model	152.63	70.23r	32.93
GloVe + GRU 136.16	50.65	26.42	

While the scores for the baseline model are initially relatively high and not particularly good, this is to be expected as the model had many problems such as not being able to train its own word vectors sufficiently well and falling victim to vanishing gradient problem. However, the successive models each do better than the previous versions and the final results are relatively good.

Current state of the art implementations scored and trained on the MS-COCO databases have scored a perplexity of 14.23. My value is comparable to that value and while mine is a little higher that is to be expected as I have had a much smaller amount of time and computing power and there were a decent amount concessions I had to make in the interest of what was possible in the given time given my resources. Examples In this section I will go through a few captions generated by my models and point out the parts that my model does well at and the parts that it does poorly on. Here are some captions created by the GloVe and GRU models, I chose not to include the baseline captions as they are worse than the captions produced by these and rather hard to understand.

GloVe Based Model

270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311  
312  
313  
314  
315  
316  
317  
318  
319  
320  
321  
322  
323

1. round in a young girl in the background a women on a large red pour a boy and woman a restaurant the preparing a girl with a circle in man is relaxing a man hat
  2. they are colored dressed look at another person small group of people bending down a soccer ball wtwo people are gathered on the open in its forefront
  3. two kids playing his tan room guy in pink is one a crd bucket a man wearing a bench with front of a sheath in pat of others and man holding a counter
- GRU Model
1. a woman is wearing a black shirt and woman in a red dress shirt and a woman are working on a restaurant a young woman holds a fruit on a black of a woman at a table with a wooden of woman
  2. boy in white are practicing walking down in street with many children in yellow uniforms four men and a yellow belt belt boys are playing karate on a line a competition player in a field soccer children play
  3. sits in front of a building a man in jeans and sunglasses is standing a girl in a stands in the middle of the street a young woman is standing he is looking at a woman in a city street an older woman is sitting on a sidewalk

Looking at the captions it can be seen that while the sentences can be somewhat awkward in some places they still do a good job representing what is going on in images, while still making some mistakes. Also, it can be seen how the results from the GRU based network are better than the results than the simpler model in many cases. For example, in image 1 both models are able to reflect that there women in the picture the GRU model is able to give more details such as that there is fruit rather than just that there is something round as simpler model suggests, or that it is get that the people in the third image are out in a sidewalk and not inside a room with tan walls. Both models seem to do a good job getting the colors that are in the image though, by recognizing the color clothes that the people are wearing this can be seen by getting the woman in red in image 1, the tan walls in image 2, and the white and yellow uniforms in image 3. However, both networks still make some mistakes like how the GRU network sees a man in sunglasses and jeans even though there is no one dressed like that. So, while there is still some mistakes being made even the state of the art model makes lots of mistakes, so this model can still be seen as relative success as it does do a sufficient job captioning and expressing what is going in images.

## 5 Future Work

One of the main improvements that could be done for this model would be to train the CNN along with RNN. Since I did not have the time or computing power to back propagate the errors to the CNN I used the pre-trained AlexNet network. While that does produce good results the network definitely perform better if it could train the CNN at the same time and the errors would affect the values in the CNN.

Another area for improvement would be to try more complicated ways of scoring the language model to get better results. If the model was switched to scoring the words in caption based on a window it would probably significantly help make the language model more robust and help make the end result better and more natural to read and understand.

After that I would try to use a more complicated RNN structure. A multi-layer RNN may be able to better learn the features of the image and the language model and perform better results.

## References

[1] Mao, Junhua, Xu, Wei, Yang, Yi, Wang, Jiang, and Yuille, Alan. Deep captioning with multimodal recurrent neural networks (m-rnn). arXiv:1412.6632, December 2014.

324 [2] Vinyals, Oriol, Toshev, Alexander, Bengio, Samy, and Erhan, Dumitru. Show and tell: A neural image  
325 caption generator. arXiv:1411.4555, November 2014.  
326 [3] Karpathy, Andrej and Li, Fei-Fei. Deep visual-semantic alignments for generating image descriptions.  
327 arXiv:1412.2306, December 2014.  
328 [4] Chen, Xinlei, and Zitnick, C Lawrence. Learning a Recurrent Visual Representation for Image Caption  
329 Generation. arXiv:1411.5654, November 2014.  
330 [5] Jeffrey Pennington, Richard Socher, and Christopher D. Manning. 2014. GloVe: Global Vectors for Word  
331 Representation  
332 [6] Peter Young, Alice Lai, Micah Hodosh and Julia Hockenmaier. From image descriptions to visual denota-  
333 tions: New similarity metrics for semantic inference over event descriptions, Transactions of the Association  
334 for Computational Linguistics, 2(Feb):67-78, 2014. **15**(7):5249-5262.  
335  
336  
337  
338  
339  
340  
341  
342  
343  
344  
345  
346  
347  
348  
349  
350  
351  
352  
353  
354  
355  
356  
357  
358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377