# Image Captioning COCO Dataset with LSTM Network

**Dylan Haar**
dhaar@ucsd.edu

**Taehyung Kim**
tak088@ucsd.edu

**Tae Kun Kim**
tkkim@ucsd.edu

University of California, San Diego
La Jolla, CA, USA

## Abstract

In this paper, we use pretrained convolutional network, resnet 50, and LSTM model for image captioning with COCO dataset. Since the original dataset is too large, we used about $1/5$th part of original data. In order to compare the performance of models, we made two different models: baseline model and Vanilla RNN. Also, we made two different type of architecture for training and two different method of generating captions for further experiment. We recorded the performance of our models with loss, BLEU-1 score, and BLEU-4 score on validation data and test data. Although we tried to tune the hyperparameters to find the best one, all the models were keep overfitting so we used the default setting for best hyperparameters. The baseline model's loss was 0.026 and BLEU-1 score and BLEU-4 score was 71.1% and 7.2%, respectively. LSTM with second architecture had 1.44 for it loss and BLEU-1 score and BLEU-4 score was 66.2% and 7.6%, respectively. Finally, the Vanilla RNN model's loss was 1.62 and BLEU-1 score and BLEU-4 score was 66.8% and 7.7%, respectively. Based on the result, we were able to discover that performance of baseline model outweighs performance of other type of model and architecture.

## 1 Introduction

Among all different types of neural networks, Recurrent neural network shows its strength on sequential data. Depending on the architecture of the RNN models, different tasks can be performed. For example, image classification can be done with one-to-one structure, sentiment analysis with many-to-one structure, image captioning with one-to-many structure, etc. Since the task for this project is image captioning, we used one-to-many structure while building the RNN model. Before using the RNN model, we used pretrained CNN model, resnet 50, as encoder which encodes the image and use the output of encoder as an input of RNN model which is used as decoder model to generate captions for each image. We used subset of data from "COCO2015Image Captioning Task". Specifically, we used approximately 82k images and 410k captions and generate captions during testing procedure. As we are generating the captions from image, in order to check the performance of model, we used BLEU-1 score, and BLEU-4 score which are commonly used for machine translation and give score of generated caption by comparing it with true caption of image.

## 2 Background/ Related Work

We got a lot of help from the articles below and were able to develop our model and generate captions based on the knowledge on the articles. Specifically, from the pytorch documentation, we were able to know the output shapes of LSTM network and how to use data for teacher forcing. Also, pytorch documentation of embedding layer helped us understand the process of embedding function and what we should use for function's parameter.

1

https://pytorch.org/docs/stable/generated/torch.nn.Embedding.html
https://pytorch.org/docs/stable/generated/torch.nn.LSTM.html        https://medium.com/analytics-vidhya/cnn-lstm-architecture-and-image-captioning-2351fc18e8d7
https://calvinfeng.gitbook.io/machine-learning-notebook/supervised-learning/recurrent-neural-network/recurrent_neural_networks

## 3 Methods

For this image captioning task, we used pair of different model type, CNN-LSTM. Firstly, we utilized pretrained CNN model, ResNet 50, to encode the image by changing the last layer of model from performing softmax to linear layer. By using CNN model, we get 300-dimensional linear feature vector which would be a input for LSTM network. In LSTM network, we used nn.Embedding function which automatically performs one-hot encoding and transforms to linear vector then using the output of embedding layer as an input of LSTM network with 2 layers. Finally, the output of LSTM network was used in fully connected layer and softmax to get the performance of the model or generate captions.

For LSTM network, we made two different architectures. For both of the architectures, we used a 'teacher forcing' method during the training procedure. Instead of using previous output of LSTM cell, we input word from captions to all LSTM cells. In baseline architecture, as I mentioned above, encoder (ResNet 50) outputs a feature vector with given images. We input the encoded images to embedding layer where we used nn.Embedding from pytorch so that we can obtain word embeddings that goes through LSTM network. After obtaining output from LSTM network and its linear layer, we used Cross entropy function which automatically performs softmax and NLLLoss and gives loss for the model.

In second architecture, most of procedures are same. However, we concatenated the feature vector from encoder model and word embeddings from embedding layer and pass the concatenated data to LSTM network. Also, we passed input images at each time step and used '¡pad¿' as input along with the image features.

We also implemented Vanilla RNN as a decoder for image captioning. The overall structure was same with the decoder model with LSTM network. We used nn.Embedding to get word embedding and, again, concatenated feature vector from encoder model and word embeddings to pass into Vanilla RNN model.

For generating captions, we utilized two different approaches: deterministic and stochastic. After the data is passed through all the layers, we have list of probabilities. With deterministic approach, it selects the maximum probability for every time step. However, stochastic approach samples from the probabilities and by using temperature variable which controls the sampling procedure of stochastic approach.

In order to find the best hyperparameters, we performed experiments by changing some values such as batch size, number of epochs, learning rate, number of hidden size, number of embedding size, and value for temperature in stochastic caption generation.

## 4 Results

After trying two different approaches for generating captions and comparing the results with loss and BLEU scores, we were able to find out that stochastic approach works better than deterministic approach. Therefore, we recorded best hyperparameters, loss, and BLEU scores for baseline model with stochastic approach.

### 4.1 Best Hyperparameters

| Model type | Feature | Value |
|---|---|---|
| Baseline | Batch size | 64 |
| Baseline | Epochs | 15 |
| Baseline | Learning Rate | 0.0005 |
| Baseline | Hidden Size | 512 |
| Baseline | Embedding Size | 300 |
| Architecture 2 | Batch size | 64 |
| Architecture 2 | Epochs | 15 |
| Architecture 2 | Learning Rate | 0.0005 |
| Architecture 2 | Hidden Size | 800 |
| Architecture 2 | Embedding Size | 400 |
| Vanilla RNN | Batch size | 64 |
| Vanilla RNN | Epochs | 15 |
| Vanilla RNN | Learning Rate | 0.0005 |
| Vanilla RNN | Hidden Size | 900 |
| Vanilla RNN | Embedding Size | 450 |

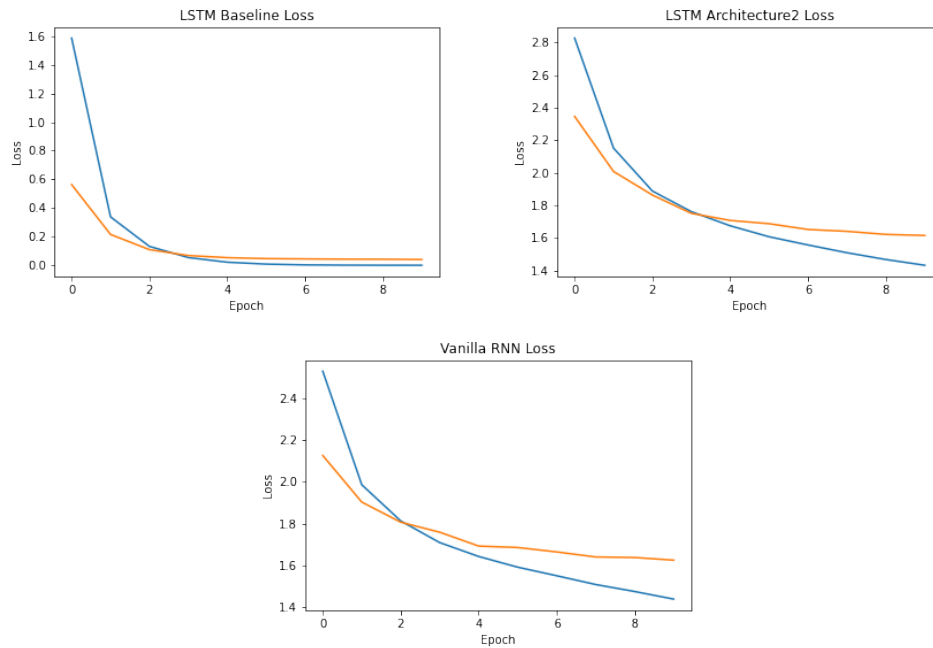Table 1: Best Hyperparameters

## 4.2 Plot



Figure 1: Training and validation loss changes over epochs

| Model Type | Loss |
|---|---|
| Baseline | 0.026 |
| Architecture 2 | 1.44 |
| Vanilla RNN | 1.62 |

Table 2: Cross Entropy Loss

3

## 4.3 BLEU Score

| Model Type | BLEU-1 | BLEU-4 |
|---|---|---|
| Baseline | 71.1% | 7.2% |
| Architecture 2 | 66.2% | 7.6% |
| Vanilla RNN | 66.8% | 7.7% |

Table 3: BLEU Scores

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

## 4.4   Examples

[Good Caption with RNN (temperature = 0.001)]



Actual captions:
a group of people standing on top of a snow covered slope.
three people hugging and posing for a picture while skiing.
the three people are posing for a picture before skiing.
three people are posing on skis in snow.
three friends are posed on the ski slope.

Predicted caption: a group of people standing on a snow covered slope.

[Good Caption with RNN (temperature = 0.001)]



Actual captions:
a vase filled with reddish orange roses and shaving gel.
the large flowers fade from orange to pink.
a bouquet of roses sits in a container next to items.
there are some flowers in a vase on a table
a bouquet of roses in a juice jug on a bedside table.

Predicted caption: a vase filled with flowers on a table.

[Good Caption with RNN (temperature = 0.001)]



Actual captions:
a tour bus is parked along a curb.
a red bus is parked on the side of the road.
a tour bus is parked on a tree lined street.
a bus sits parked next to a stand of trees.
a red white and black bus and some bushes and trees

Predicted caption: a bus is parked on the side of a road.

Actual captions:
two men play wii in a cramped living room.

two men standing in a living room with remotes.
two men standing in a living room playing games.
two men playing video games on the nintendo wii.
two men standing in a living room holding game controllers.

Predicted caption: a man in a white shirt and a tie and a tie

[Bad Caption with RNN (temperature = 0.001)]



Actual captions:
a shopping center with parking lot along a busy street.
safeway has sprung for a roadside clock in front of their plaza.
a bunch of traffic passes by a shopping center
the clock tower actually has a open air bottom.
a small clock tower from across the very busy street.

Predicted caption: a street sign on a pole in front of a building.

486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539

[Bad Caption with RNN (temperature = 0.001)]



Actual captions:
a group of people with bananas tied to their bicycles.
a group of people and motorcycles with bananas on the side.
a crowd of people standing near motorcycles carrying bananas.
a line of parked bikes with bananas attacted to them
bikes parked side by side with lots of green bananas on them.


Predicted caption: a man is riding a horse in a parade.

# 5 Discussion

In this report, we learned that image captioning works best with the combination of ResNet50 and LSTM in terms of the Cross Entropy loss and BLEU 1 score. Although the team had hypothesized that the Architecture 2 model will perform best, it is not surprising that the baseline model performs just as well as thanks to its cell state attributes — the model must have learned to understand that the image features are important when producing captions for images, and therefore saved the image information in its cell states. What was surprising, however, that the LSTM model actually outperformed architecture 2 since architecture 2 explicitly feeds the model with image features. To understand this phenomenon, the team has agreed to think that the model was smart enough to understand what parts of image features need to be saved in the cell state, allowing the model to better perform compared to when it is fed in with complicated image features.

Making sense of how the vanilla RNN model performs worse than the baseline model was not as difficult: the vanilla model does not pass on image features throughout its cells; thus, it is expected that the model has difficulty understanding the context of the image when generating captions. What was challenging to realize was why the vanilla model perform better than architecture 2. Although no experiment has been done in this research paper to prove this speculation, the team postulates that feeding in the same image over and over again to an LSTM model confuses the model to understand why such redundant input data is fed it. That is, the model tries to learn the importance of the repetition of data, when it should not be using its resources on such a matter; when the model could have saved its resources for optimizing the retention of the image data (like in the baseline model), the model had to learn (which it failed to) why the repetition of the input image is important.

It is still dubious to the team why the BLEU 4 score does not follow the pattern of Cross Entropy loss and the BLEU 1 score: in terms of the BLEU 4 score, the baseline model performs worst. Although it is still difficult for the team to put a finger on it, Haar et al think that the small variance in the BLEU 4 score across the three models fails to convince us that the baseline model actually performs worse than the other two models. That is, the team think that the error range of the BLEU 4 score makes it superficially look like the baseline model performs worst, when the true story may be different.

# 6 Team Contribution

**Dylan Haas**
Dylan, along with his two teammates, created the baseline model. He also worked with Tae Kun to ensure that all the code runs smoothly in the main file. Dylan also spent time working on building the vanilla model.

**Taehyung Kim**
By being present in all meetings, Taehyung helped the team create the baseline model and he also worked on creating architecture 2. Taehyung also worked on putting together this report.

**Tae Kun Kim**
Tae Kun was a part of all scheduled meetings. He contributed to building the baseline model and part of the process of putting comments on all teh code. Along with Dylan, Tae Kun made sure that everyone's code work come together and runs smoothly in the main file.

# 7 References

[1] Gary Cottrell, Recurrent nets: Part 1, lecture notes, University of California, San Diego, delivered 2022.

[2] Gary Cottrell, Lecture 4A, Recurrent nets: Part 1, lecture notes, University of California, San Diego, delivered 2022.

[3] Gary Cottrell, Lecture 5, covnets part 1, lecture notes, University of California, San Diego, delivered 2022.

[4] Gary Cottrell, Lecture 6, covnets part 2, lecture notes, University of California, San Diego, delivered 2022.