

---

# Image Captioning Experiments with RNN Models on COCO's dataset

---

000  
001  
002  
003  
004  
005  
006  
007  
008  
009           **Wenbo Hu, Wenxiao Li, Huaning Liu**  
010           University of California, San Diego  
011           w1hu@ucsd.edu, we1032@ucsd.edu, h9liu@ucsd.edu  
012  
013  
014  
015

## Abstract

016           In this project, we are going to design and implement a recurrent neural network  
017           that generates captions for images in the provided dataset. This framework con-  
018           tains an encoder (convolutional) that basically takes the images as input, convert-  
019           ing them into features; and a decoder (LSTM) that actually predicts and generates  
020           the captions. Based on the framework, several related models, including baseline  
021           model, Vanilla RNN and Architecture2 are configured, trained, tuned, and tested  
022           based on the given data. After hyper-tuning for these three models, related test  
023           BLEU-1 and BLEU-4 scores are outputted as the model performance. In this re-  
024           port, related statistics, an in-depth description of methodology and further discus-  
025           sion will be provided. In short, the tuned baseline model with LSTM gives a test  
026           BLEU-1 and BLEU-4 score of 70.916 and 6.544; the tuned Vanilla RNN model  
027           gives a test BLEU-1 score and BLEU-4 score of 64.5 and 5.5; and the tuned ar-  
028           chitecture 2 model gives a test BLEU-1 score and BLEU-4 score of 69.474 and  
029           7.270. Further discussion for this result follows up.  
030  
031

## 1 General Introduction

032           Generating captions based on images is a general practice for deep learning on computer vision. In  
033           the real world, when we are to search text by images or search images by images, texts take less time  
034           for the engine to process, and it also better represents what people are trying to express. Therefore,  
035           converting images to text that describe and represent it will be meaningful given such potential  
036           applications, and it would become the main problem for this project. To build a good network for  
037           this, choosing good trainable datasets and proper network architecture are both important. To tackle  
038           this, the dataset we are going to utilize this time comes from COCO 2015 Image Captioning Task,  
039           where COCO is a “large-scale detection, segmentation, and captioning dataset”. In this dataset, for  
040           train and test data respectively, about 82k and 3k images with each respect to several captions are  
041           provided. Given such a large size of data and the upcoming complex architectures, we only use  $\frac{1}{5}$   
042           of it. Also, to better evaluate such image captioning models, different from the previous project, we  
043           are going to use BLEU (bilingual evaluation understudy) score as the evaluation metrics. This score  
044           actually represents the quality of text which has been machine-translated from the natural language  
045           (images here) to another. It actually takes a very high correlation with human judgments of text  
046           quality. Accordingly, the BLEU-n score refers to the precision of text prediction under n-gram.  
047           In this case, both BLEU-1 and BLEU-4 are presented to better evaluate our models. In this way,  
048           with the dataset and metrics getting well-prepared, we get started on the framework and its related  
049           models.  
050  
051

## 2 Related Works

052           Since this dataset comes from COCO 2015 Image Captioning Task, several reflection articles and  
053           papers emerge to address this topic from different perspectives, including hidden image captioning,

054 better text generations, etc. Based on this project's problems and planning works this time, Oriol  
055 Vinyals and his team's work, **Show and Tell: Lessons learned from the 2015 MSCOCO Image**  
056 **Captioning Challenge** [1], gives us a good reference. In this paper, based on COCO 2015 image  
057 captioning dataset, the team built a deep learn CNN as well as language generating RNN together as  
058 the overall architecture. Furthermore, other related experiments regarding hyper-tuning and different  
059 approaches are provided and analyzed. Besides, extended from our planning, they also discuss the  
060 possibilities of using different evaluation metrics, such as METEOR and CIDEr. This work gives us  
061 a matching initiative, whose ideas and sub-methodologies will definitely be discussed to contribute  
062 to our work.

### 063 **3 Methods**

#### 064 **3.1 LSTM Model (from Baseline)**

065 For the baseline LSTM model, it first uses a pretrained convolutional network (resnet50) as the  
066 encoder the get the feature vectors, then construct a network with consecutive LSTM cells as the  
067 decoder that take the feature vectors and generate a sentence of words signified by starting with  
068 the word “⟨ start ⟩” and ending by “⟨ end ⟩”. The feature vectors of the image will be passed to  
069 the decoder initially, generating the ⟨ start ⟩ symbol, and a hidden state. Then, for the following  
070 LSTMs, the hidden state of the last LSTM, together with the embedding of the last word generated,  
071 will be passed in. nn.Embeddings is here used to address this problem by configuring a one-hot  
072 encoding as well as linear output. This procedure stops when reaching max length, or the ⟨ end ⟩  
073 symbol is generated. The outputs from the decoder are sampled and generated using greedy search,  
074 which also depends on the value of deterministic. If it goes true, the word directly takes the one with  
075 the largest; while for false, a probabilistic series is generated by softmax to decide each word. For  
076 the hyper-tuning process, we focus on the adjustment of the embedding size, and hidden size. The  
077 embedding size is adjusted within the range of 250 and 750, and the hidden size is adjusted within  
078 the range of 256 and 1024. After tuning, the best parameters we find for the baseline LSTM model  
079 is settled as below.

#### 080 **3.2 Vanilla RNN**

081 The overall architecture and logic of the Vanilla RNN model follow the pattern of the baseline  
082 model with LSTM. It first uses a pretrained convolutional network (resnet50) as the encoder the get  
083 the feature vectors, then construct a network with RNN (trained by BPTT) cells as the decoder that  
084 take the feature vectors and generate a sentence of words signified by starting with the word “⟨ start  
085 ⟩” and ending by “⟨ end ⟩”. nn.Embeddings is here used to address this problem by configuring a  
086 one-hot encoding as well as linear output. The outputs from the decoder are sampled and generated  
087 using greedy search, which also depends on the value of deterministic. If it goes true, the word  
088 directly takes the one with the largest; while for false, a probabilistic series is generated by softmax  
089 to decide each word. For the hyper-tuning process, we focus on the adjustment of the optimizer,  
090 embedding size, and hidden size. The optimizer has candidates of SGD and Adam, the embedding  
091 size is adjusted within the range of 250 and 750, and the hidden size is adjusted within the range  
092 of 256 and 1024. As usual, we found that when parameters go extreme, overfitting happens. After  
093 tuning, the best parameters we find for the Vanilla RNN model is settled as below.

#### 094 **3.3 Arch 2**

095 For architecture 2, it first uses a pretrained convolutional network (resnet50) as the encoder the get  
096 the feature vectors, then construct a network with consecutive LSTM cells as the decoder that  
097 take the feature vectors and generate a sentence of words signified by starting with the word “⟨  
098 start ⟩” and ending by “⟨ end ⟩”. The feature vectors of the image, concatenated with the word  
099 embedding of the ⟨ pad ⟩ symbol, will be passed to the decoder initially, generating the ⟨ start ⟩  
100 symbol, and a hidden state. Then, for the following LSTMs, the hidden state of the last LSTM, and  
101 feature vectors of the image, concatenated with the embedding of the last word generated, will be  
102 passed in. This procedure stops when reaching max length, or the ⟨ end ⟩ symbol is generated. The  
103 outputs from the decoder are sampled and generated using greedy search, which also depends on  
104 the value of deterministic. If it goes true, the word directly takes the one with the largest; while  
105

108 for false, a probabilistic series is generated by softmax to decide each word. The hyperparameters  
109 of architecture 2 are set to be the same as the baseline LSTM model, for better observations of the  
110 impact caused by difference in architectures.  
111  
112  
113  
114  
115

## 116 **4 Results: Best Hyperparameters** 117

### 120 **4.1 LSTM Model (from Baseline)** 121

123 The best hyperparameter we achieved for the Vanilla RNN model will be: for experiment-wise, the  
124 number of epochs is settled to be 10, with a learning rate of 5e-4. For model-wise, the hidden size  
125 is adjusted to be 1024, with an embedding size of 500.  
126  
127  
128  
129  
130

### 131 **4.2 Vanilla RNN** 132

134 The best hyperparameter we achieved for the Vanilla RNN model will be: for experiment-wise, the  
135 number of epochs is settled to be 10, with a learning rate of 5e-4. For model-wise, the hidden size  
136 is adjusted to be 512, with an embedding size of 550.  
137  
138  
139  
140

### 141 **4.3 Arch 2** 142

144 To better compare Architecture 2 to the baseline LSTM model, we stick all the hyperparameters to  
145 the best hyperparameters of the baseline LSTM model: the number of epochs is settled to be 10,  
146 with a learning rate of 5e-4. For model-wise, the hidden size is 1024, with an embedding size of  
147 1000(500 for image features, 500 for word embeddings).  
148  
149  
150  
151  
152

## 153 **5 Results: Visualization and Output Report** 154

### 155 **5.1 LSTM Model (from Baseline)** 156

157 The plot 1 shows the training loss and validation loss vs the number of epochs for the baseline LSTM  
158 model. Accordingly, the cross-entropy loss on the test set is reported to be about 2.5210.  
159  
160  
161

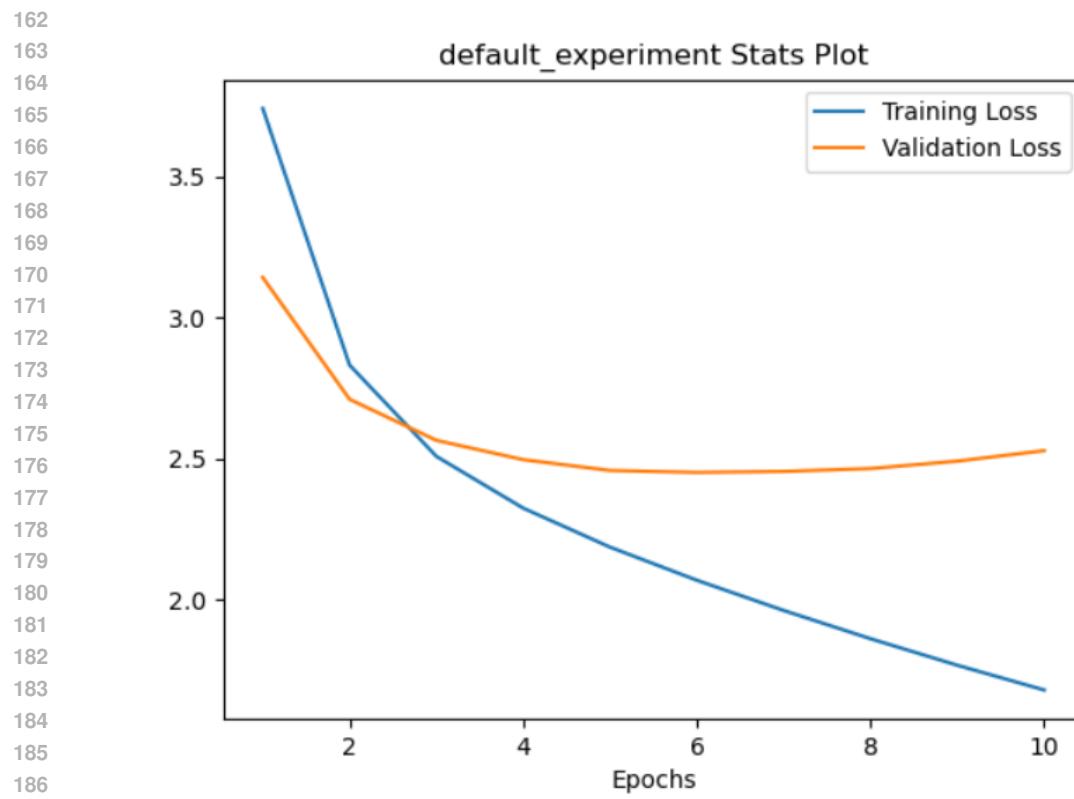


Figure 1: Baseline(LSTM) Model Loss Plot

## 5.2 Vanilla RNN

The plot 2 shows the training loss and validation loss vs the number of epochs for the Vanilla RNN model. Accordingly, the cross-entropy loss on the test set is reported to be about 2.49.

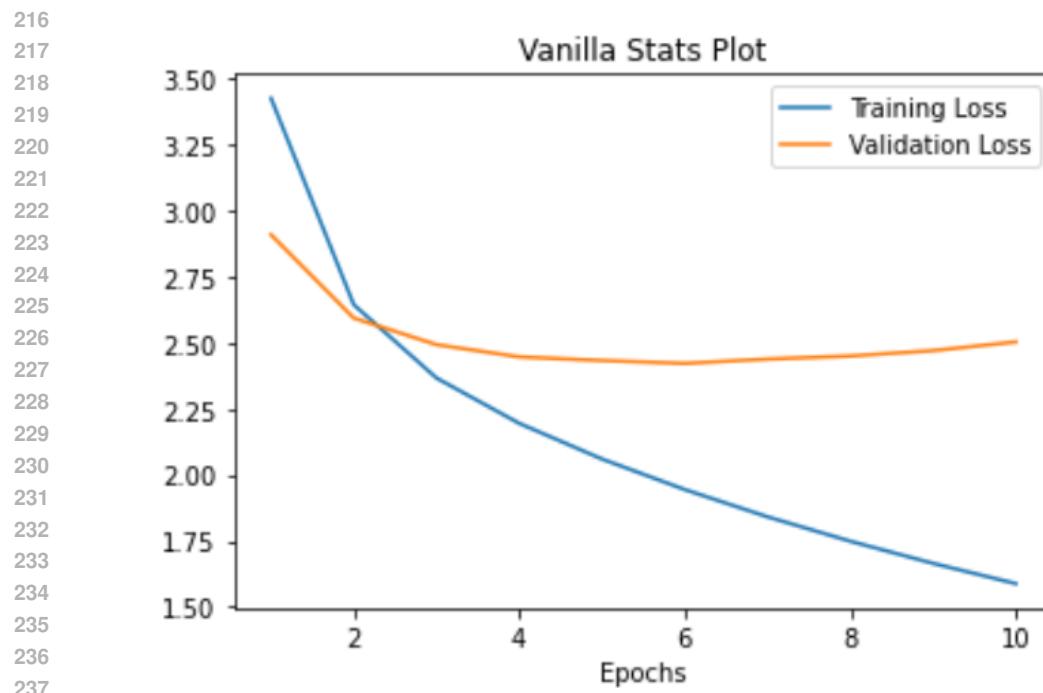


Figure 2: Vanilla RNN Model Loss Plot

### 5.3 Arch 2

260  
261  
262  
263  
264  
265  
266  
267  
268  
269 The plot 3 shows the training loss and validation loss vs the number of epochs for the architecture 2. Accordingly, the cross-entropy loss on the test set is reported to be about 2.4545.

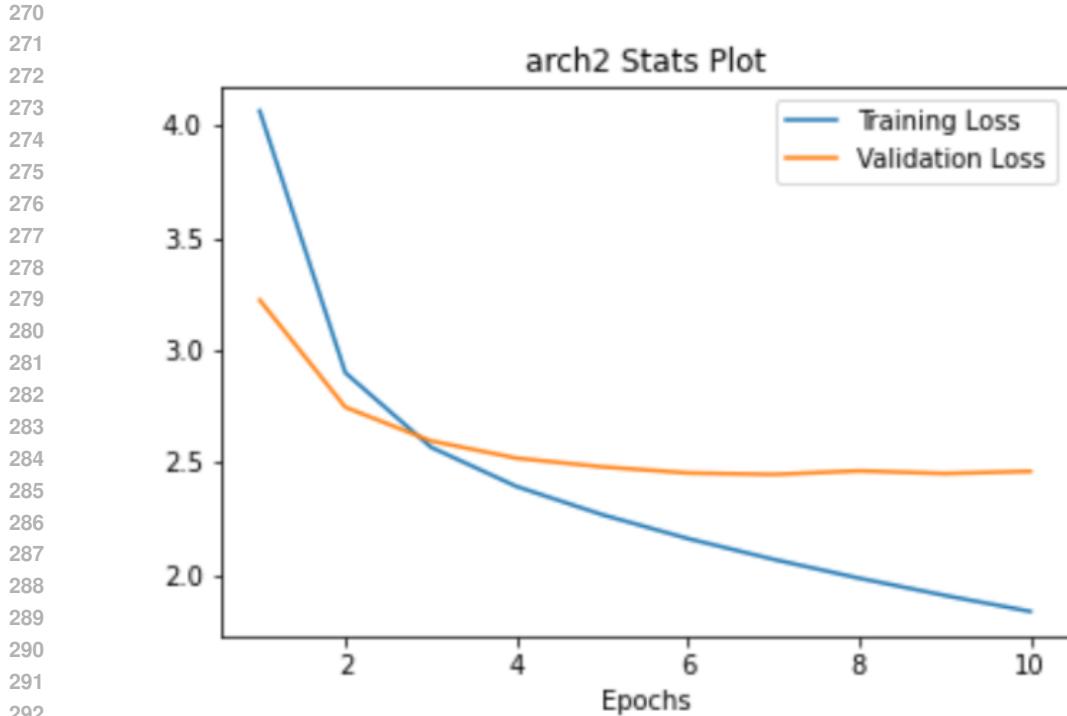


Figure 3: Arch 2 Model Loss Plot

## 6 Results: BLEU Scores Report

### 6.1 LSTM Model (from Baseline)

For the LSTM model with default parameters (hidden size = 300, embedding size = 512), the BLEU-1 score on the test set is reported to be about 66.577, and the BLEU-4 score on the test set is reported to be about 5.636. For the LSTM model with tuned parameters (hidden size = 500, embedding size = 1024), the BLEU-1 score on the test set is reported to be about 70.916, and the BLEU-4 score on the test set is reported to be about 6.544.

### 6.2 Vanilla RNN

For the tuned Vanilla RNN model, the BLEU-1 score on the test set is reported to be about 64.5, and the BLEU-4 score on the test set is reported to be about 5.5.

### 6.3 Arch 2

For the LSTM model with tuned parameters (hidden size = 500, embedding size = 1024), the BLEU-1 score on the test set is reported to be about 69.474, and the BLEU-4 score on the test set is reported to be about 7.270.

## 324 7 Image Caption Examples

325

326

327

328

329

330

331

332

333



With temperature default = 0.4:



With temperature = 0.001

Actual captions: [‘A L shaped couch with a variety of pillows on it, in front of a television.’, ‘A bottle of wine, place settings, and some apples on a coffee table in front of a couch.’, ‘THERE IS A PHOTO OF A LIVING ROOM WITH COUCH AND TV’, ‘a living room filled with couches, tv, coffee table, and a door way’, ‘a white couch tan pillows coffee table and a television set’]

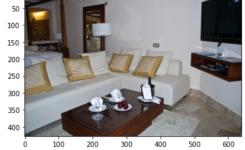
Predicted sentence: a room with a couch , coffee table , and a television .

bleu1 score: 92.3076923076923

340

341

With temperature = 5



Actual captions: [‘A L shaped couch with a variety of pillows on it, in front of a television.’, ‘A bottle of wine, place settings, and some apples on a coffee table in front of a couch.’, ‘THERE IS A PHOTO OF A LIVING ROOM WITH COUCH AND TV’, ‘a living room filled with couches, tv, coffee table, and a door way’, ‘a white couch tan pillows coffee table and a television set’]

Predicted sentence: royal makeshift durham couple aims newlywed quilt also coronation flattened shapes nicely fish fancy square elegant colored patriotic

bleu1 score: 0

With temperature deterministic



bleu1 score: 92.3076923076923

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

Figure 4: Good Caption 1

With temperature default = 0.4:



Actual captions: [‘a bathroom that has a sink and a toilet in it’, ‘A bathroom with marble counter top and red wall treatment.’, ‘A bathroom with a red color scheme on the walls.’, ‘A bathroom that has a tub, a toilet, a sink, and a mirror.’, ‘A bathroom with a tub, toilet, sink, a counter and a cabinet in it.’]

Predicted sentence: a bathroom with a toilet , sink and mirror in it .

bleu1 score: 100.0

With temperature = 0.001



Actual captions: [‘a bathroom that has a sink and a toilet in it’, ‘A bathroom with marble counter top and red wall treatment.’, ‘A bathroom with a red color scheme on the walls.’, ‘A bathroom that has a tub, a toilet, a sink, and a mirror.’, ‘A bathroom with a tub, toilet, sink, a counter and a cabinet in it.’]

Predicted sentence: a bathroom with a toilet , sink and shower .

bleu1 score: 81.43536762323636

With temperature = 5



Actual captions: [‘a bathroom that has a sink and a toilet in it’, ‘A bathroom with marble counter top and red wall treatment.’, ‘A bathroom with a red color scheme on the walls.’, ‘A bathroom that has a tub, a toilet, a sink, and a mirror.’, ‘A bathroom with a tub, toilet, sink, a counter and a cabinet in it.’]

Predicted sentence: curtains man given headset customers instructors about ballplayer snapshots plain walking clothing lloyd ramp officials and stove

bleu1 score: 5.6556556556556556

With temperature deterministic



Actual captions: [‘a bathroom that has a sink and a toilet in it’, ‘A bathroom with marble counter top and red wall treatment.’, ‘A bathroom with a red color scheme on the walls.’, ‘A bathroom that has a tub, a toilet, a sink, and a mirror.’, ‘A bathroom with a tub, toilet, sink, a counter and a cabinet in it.’]

Predicted sentence: a bathroom with a toilet , sink and shower .

bleu1 score: 81.43536762323636

Figure 5: Good Caption 2

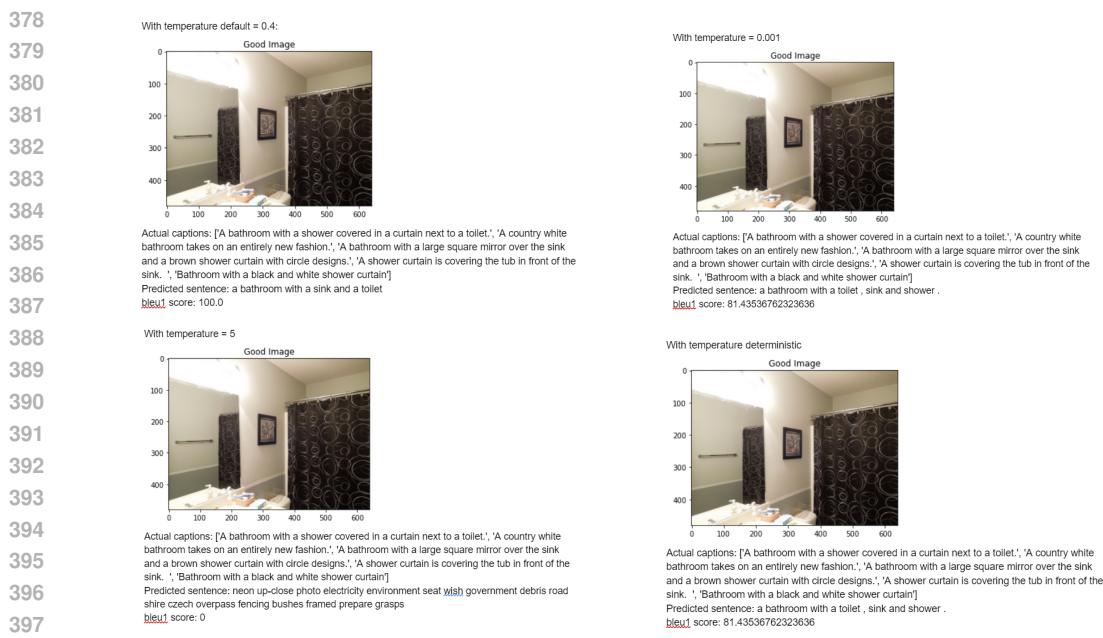


Figure 6: Good Caption 3



Figure 7: Bad Caption 1

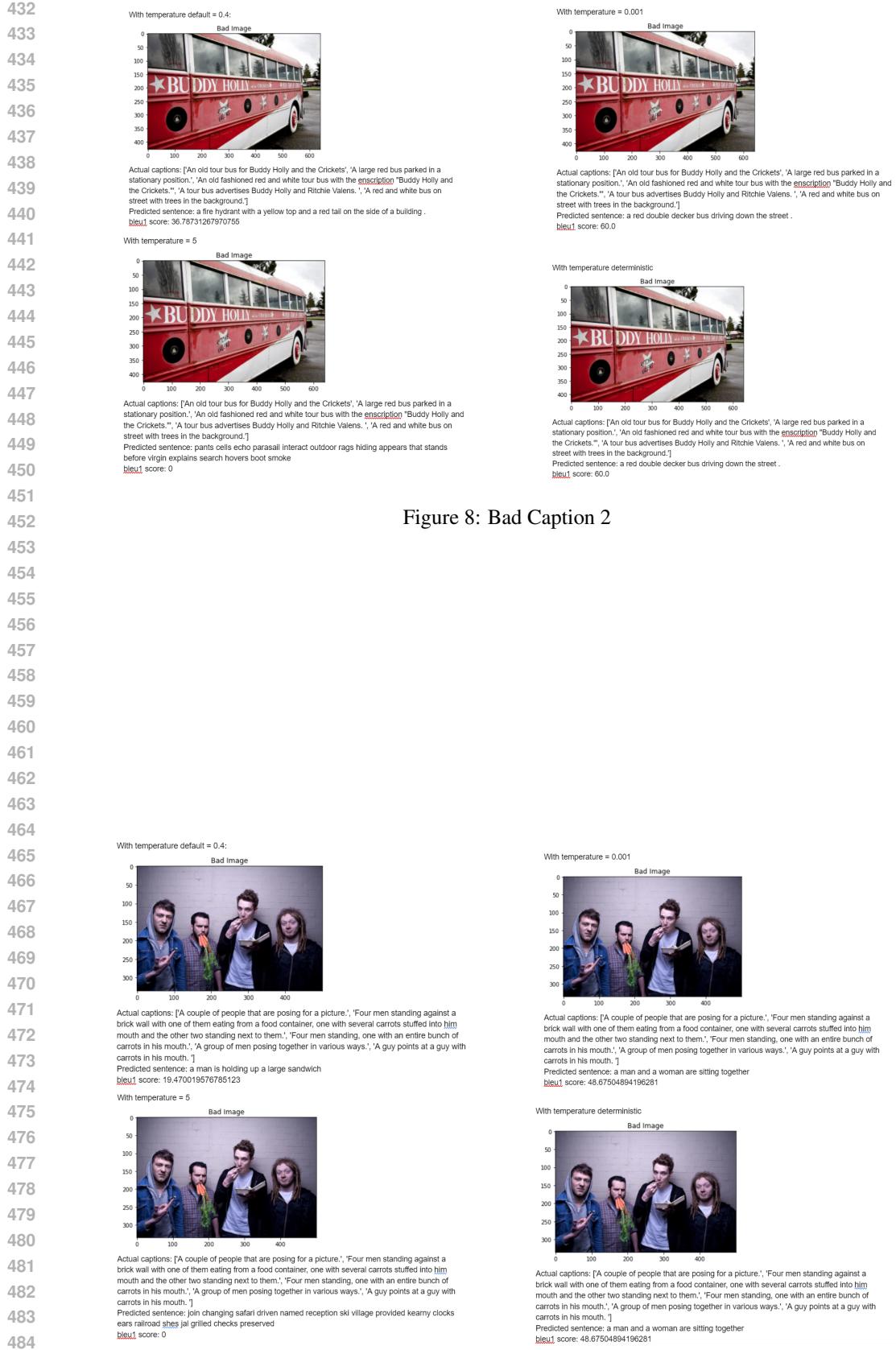


Figure 9: Bad Caption 3

## 7.2 Vanilla RNN

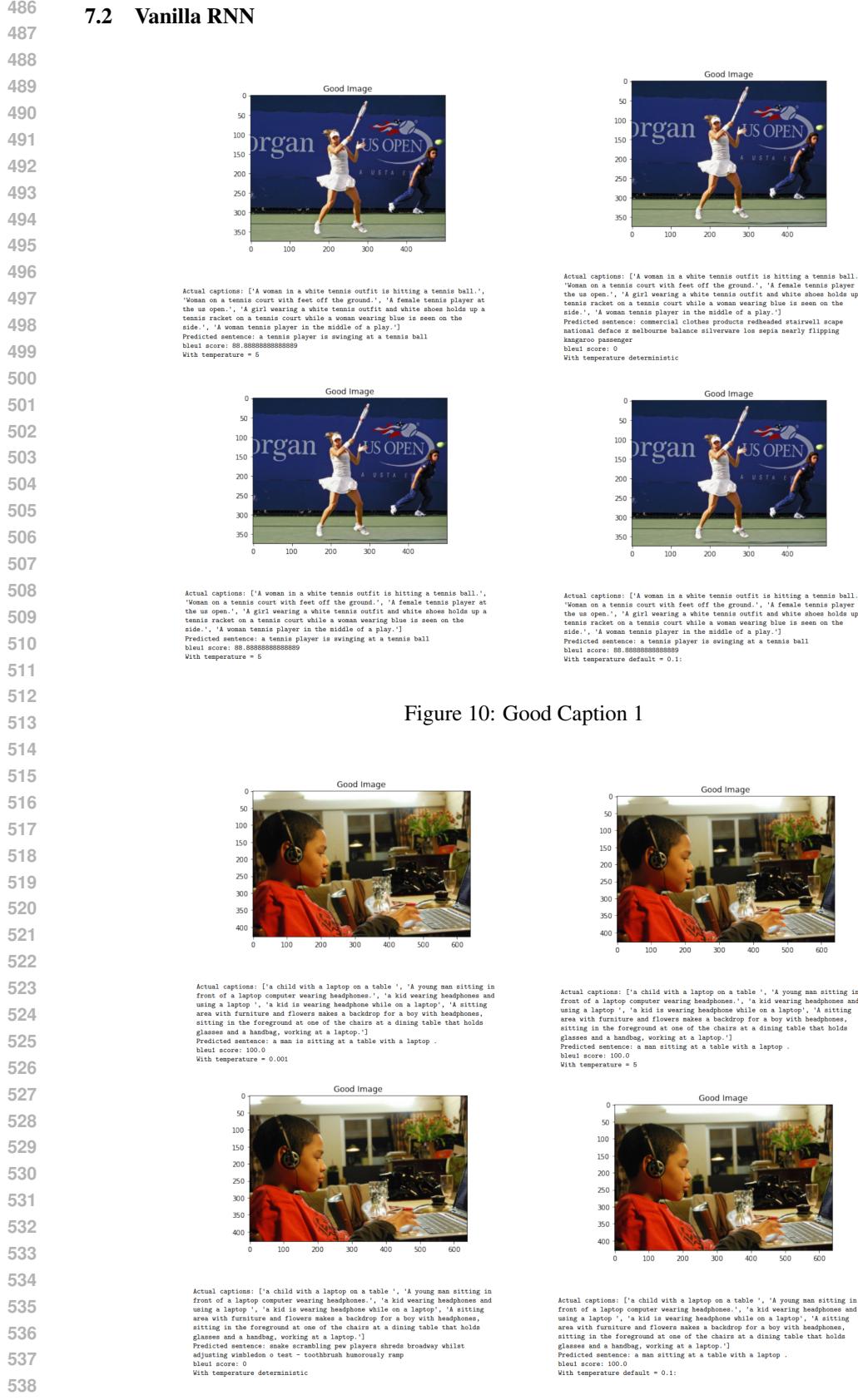


Figure 10: Good Caption 1

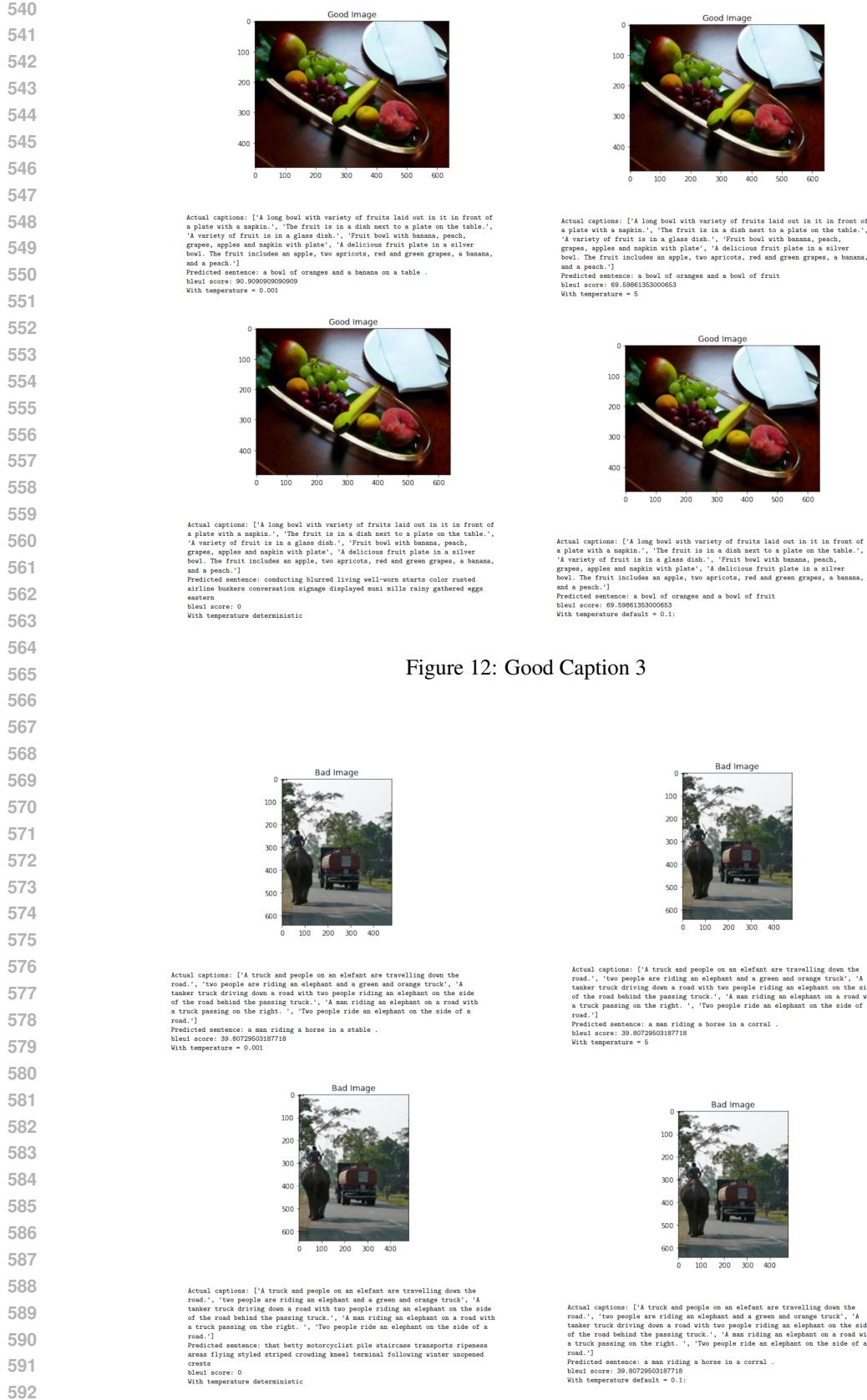


Figure 12: Good Caption 3

Figure 13: Bad Caption 1



Figure 14: Bad Caption 2

Figure 15: Bad Caption 3

### 7.3 Arch 2

648

649

650

651

652

653

With temperature default = 0.4:



Actual captions: [A bathroom sink shaped like a glass bowl.', 'A man that is standing with a mug in front of a mirror.', 'a man taking a picture of himself in a bathroom mirror', 'A man standing in a bathroom taking a picture of his self in the mirror with his cell phone', 'a man taking a selfie in a little bathroom mirror']

Predicted sentence: a man is standing in front of a mirror .

`bleu1` score: 100.0

With temperature = 5



Actual captions: [A bathroom sink shaped like a glass bowl.', 'A man that is standing with a mug in front of a mirror.', 'a man taking a picture of himself in a bathroom mirror', 'A man standing in a bathroom taking a picture of his self in the mirror with his cell phone', 'a man taking a selfie in a little bathroom mirror']

Predicted sentence: dedicated coffee wanting print cleaner one bottled plywood ordering spaghetti key bowls tapestries filling bazaar security similarly squid

`bleu1` score: 0

With temperature = 0.001



Actual captions: [A bathroom sink shaped like a glass bowl.', 'A man that is standing with a mug in front of a mirror.', 'a man taking a picture of himself in a bathroom mirror', 'A man standing in a bathroom taking a picture of his self in the mirror with his cell phone', 'a man taking a selfie in a little bathroom mirror']

Predicted sentence: a man is taking a picture of herself in a mirror .

`bleu1` score: 91.66666666666666

With temperature deterministic



Actual captions: [A bathroom sink shaped like a glass bowl.', 'A man that is standing with a mug in front of a mirror.', 'a man taking a picture of himself in a bathroom mirror', 'A man standing in a bathroom taking a picture of his self in the mirror with his cell phone', 'a man taking a selfie in a little bathroom mirror']

Predicted sentence: a man is taking a picture of herself in a mirror .

`bleu1` score: 91.66666666666666

Figure 16: Good Caption 1

673

674

675

676

677

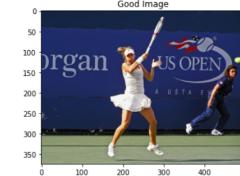
678

679

680

681

With temperature default = 0.4:

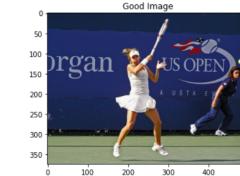


Actual captions: [A woman in a white tennis outfit is hitting a tennis ball.', 'Woman on a tennis court with feet off the ground.', 'A female tennis player at the us open.', 'A girl wearing a white tennis outfit and white shoes holds up a tennis racket on a tennis court while a woman wearing blue is seen on the side.', 'A woman tennis player in the middle of a play.]

Predicted sentence: a woman is hitting a tennis ball with her racket .

`bleu1` score: 90.90909090909090

With temperature = 5



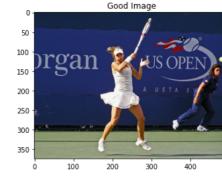
Actual captions: [A woman in a white tennis outfit is hitting a tennis ball.', 'Woman on a tennis court with feet off the ground.', 'A female tennis player at the us open.', 'A girl wearing a white tennis outfit and white shoes holds up a tennis racket on a tennis court while a woman wearing blue is seen on the side.', 'A woman tennis player in the middle of a play.]

Predicted sentence: sash pilot soaps beyond magnetic arms frowning ins tent got fashioned

progress players lounging tailing wildlife fault gallop

`bleu1` score: 0

With temperature = 0.001

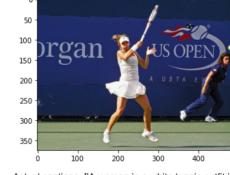


Actual captions: [A woman in a white tennis outfit is hitting a tennis ball.', 'Woman on a tennis court with feet off the ground.', 'A female tennis player at the us open.', 'A girl wearing a white tennis outfit and white shoes holds up a tennis racket on a tennis court while a woman wearing blue is seen on the side.', 'A woman tennis player in the middle of a play.]

Predicted sentence: a tennis player is hitting a ball on the court .

`bleu1` score: 100.0

With temperature deterministic



Actual captions: [A woman in a white tennis outfit is hitting a tennis ball.', 'Woman on a tennis court with feet off the ground.', 'A female tennis player at the us open.', 'A girl wearing a white tennis outfit and white shoes holds up a tennis racket on a tennis court while a woman wearing blue is seen on the side.', 'A woman tennis player in the middle of a play.]

Predicted sentence: a tennis player is hitting a ball on the court .

`bleu1` score: 100.0

Figure 17: Good Caption 2

700

701

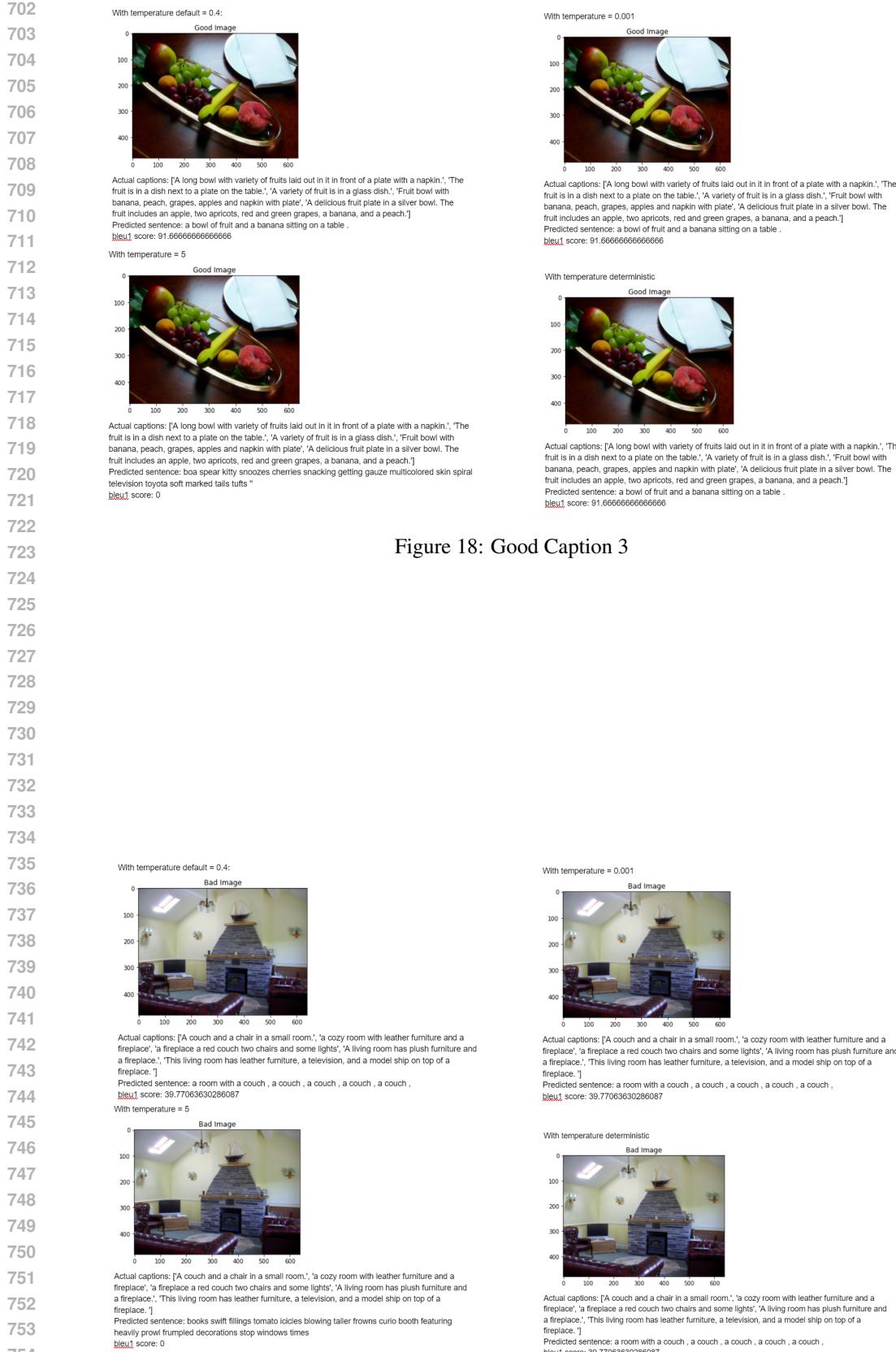


Figure 18: Good Caption 3

Figure 19: Bad Caption 1

756 With temperature default = 0.4:  
 757   
 758 Bad Image  
 759  
 760 Actual captions: ['Two microwave ovens one marked with man symbol and the other for woman.', 'Two microwaves sitting side by side on a countertop are marked with signs printed with the symbols for man and woman.', 'a couple of microwaves that are next to each other', 'Two microwaves side by side with pictures of a man and woman silhouette', 'The microwaves are marked with female and male signs.'];  
 761 Predicted sentence: a kitchen with a stove and a stove  
 bleu1 score: 29.20502936517768  
 762  
 763 With temperature = 0.001  
 764   
 765 Bad Image  
 766  
 767 Actual captions: ['Two microwave ovens one marked with man symbol and the other for woman.', 'Two microwaves sitting side by side on a countertop are marked with signs printed with the symbols for man and woman.', 'a couple of microwaves that are next to each other', 'Two microwaves side by side with pictures of a man and woman silhouette', 'The microwaves are marked with female and male signs.'];  
 768 Predicted sentence: a kitchen with a stove and a stove  
 bleu1 score: 29.20502936517768  
 769  
 770 With temperature = 5  
 771   
 772 Bad Image  
 773  
 774 Actual captions: ['Two microwave ovens one marked with man symbol and the other for woman.', 'Two microwaves sitting side by side on a countertop are marked with signs printed with the symbols for man and woman.', 'a couple of microwaves that are next to each other', 'Two microwaves side by side with pictures of a man and woman silhouette', 'The microwaves are marked with female and male signs.'];  
 775 Predicted sentence: jeep nicely coverit multiple car guards takeoff sign shelf specially managers cleaning fives suited suit nah gone low  
 bleu1 score: 0  
 776  
 777  
 778  
 779  
 780 With temperature default = 0.4:  
 781   
 782 Bad Image  
 783  
 784 Actual captions: ['A swimming pool with lanes marked for racing on the water's edge across from a large city.', 'An Olympic style swimming pool, with many dividers', 'an Olympic pool sitting next to a river', 'A large swimming pool is next to a large harbor', 'A swimming pool that is next to a ocean.'];  
 785 Predicted sentence: a bench and a bench on a sunny day  
 bleu1 score: 33.333333333333  
 786  
 787 With temperature = 5  
 788   
 789 Bad Image  
 790  
 791 Actual captions: ['A swimming pool with lanes marked for racing on the water's edge across from a large city.', 'An Olympic style swimming pool, with many dividers', 'an Olympic pool sitting next to a river', 'A large swimming pool is next to a large harbor', 'A swimming pool that is next to a ocean.'];  
 792 Predicted sentence: sundown occupant insignia rodeo w walkway showing quarters scape  
 793 chateau salls accord upon banner message stand moored  
 794 bleu1 score: 0  
 795  
 796  
 797  
 798  
 799  
 800  
 801  
 802

803 **8 Discussion and Findings**  
 804  
 805 **8.1 Why does the deterministic approach not work well?**  
 806  
 807 Deterministic approach only takes the maximum probability of a word at each step. While for  
 808 stochastic processes, setting a temperature parameter makes the modeler control the stochastic  
 809 attribute of the distribution. If temperature is close to 0, then the distribution is nearly the same as  
 the deterministic approach, if the temperature is approaching infinity, then the distribution is almost

Figure 20: Bad Caption 2

  
 800 With temperature = 0.001  
 801 Bad Image  
 802  
 803 Actual captions: ['A swimming pool with lanes marked for racing on the water's edge across from a large city.', 'An Olympic style swimming pool, with many dividers', 'an Olympic pool sitting next to a river', 'A large swimming pool is next to a large harbor', 'A swimming pool that is next to a ocean.'];  
 804 Predicted sentence: a bench and a bench on a beach  
 bleu1 score: 37.5  
 805  
 806 With temperature deterministic  
 807   
 808 Bad Image  
 809  
 810 Actual captions: ['A swimming pool with lanes marked for racing on the water's edge across from a large city.', 'An Olympic style swimming pool, with many dividers', 'an Olympic pool sitting next to a river', 'A large swimming pool is next to a large harbor', 'A swimming pool that is next to a ocean.'];  
 811 Predicted sentence: a bench and a bench on a beach  
 bleu1 score: 37.5

Figure 21: Bad Caption 3

810 a uniform distribution. So it gives a chance of all the output words according to their probability  
811 distribution. The deterministic approach selects the most confident word. However, sometimes we  
812 don't want the next word only to be very focused on the previous word, we want the next word to  
813 also consider image or the whole correlation in the sentence which could be the second or the third  
814 most confident word. Thus, the stochastic approach is better and more reasonable.

## 817 8.2 LSTM Model (from Baseline)

818  
819 From the results of the baseline models, we are able to observe that choosing a proper temperature  
820 to generate captions is crucial. When the temperature is too big, such as 5, words will be selected  
821 in a roughly uniform distribution, and the predictions do not make any sense. When temperature is  
822 very small, such as 0.001, it is almost certain that the words with the largest output value will be  
823 selected. From the outputs shown, we can see that the output captions for deterministic settings are  
824 the same as the output captions when temperature equals 0.001. There are advantages for adding a  
825 little bit of uncertainty to generate captions by introducing a temperature such as 0.4. For instance,  
826 in the first good picture of the bathroom, it is able to say there is a mirror. However, when it is set  
827 to deterministic, both scenarios will have output "a bathroom with a toilet, sink and shower." This  
828 is probably due to the inherent connection between bathroom, toilet, sink, and shower(often appear  
829 simultaneously). And, the deterministic decoder may output them all together, and make mistakes.  
830 Admittedly, we also observe that deterministic sometimes perform better than temperature = 0.4,  
831 such as the buddy holly bus picture. However, the difference in bleu1 score is only reflected by  
832 predicting the word "bus", which is possibly due to temperature = 0.4 giving a bad choice because  
833 of luck. This scenario hardly occurs. Plus, considering the better performance of bleu1 score when  
834 temperature = 0.4 (71 to 66 bleu1 score), we believe deterministic predictions are worse.

## 835 8.3 Vanilla RNN

836 The result of RNN, 0.645 of BLEU-1 score, leads it to be the worst model. The architecture of  
837 RNN suffers the problem of gradient exploding and gradient vanishing when training with BPTT.  
838 So the weights updating in the cells may not reach the ideal situation. Thus, Vanilla performance is  
839 not good. On the other hand, from the perspective of loss plots, whatever combo of parameters we  
840 attempted, slight or obvious overfitting during the training process will always be happening. This  
841 might be another reason that vanilla RNN becomes the worst model among the three, which could  
842 be caused by the adjustment of layers of the framework compared with the baseline. Furthermore,  
843 considering the intervening of deterministic over the training process, it is another key element  
844 that drags the performance of this model down, so we preferably use the ones without temperature  
845 deterministic. Then, considering the examples we selected, regarding the variation of temperature  
846 between the default of 0.1, 0.001, and 5, not many big gaps of BLEU scores are detected, which is  
847 possibly due to the fair luck we have during the training.

## 848 8.4 Arch2

849 For images in architecture 2, the deterministic sampler still falls into the case of combining words  
850 that often appear simultaneously in the real world, to its actual output. For instance, in the image with  
851 a swimming pool, the deterministic sampler predicts "a bench and a bench on a beach." However,  
852 "beach" does not appear in the picture. "Beach" is predicted due to its close connection to "bench".  
853 And, we also believe that architecture 2 does not perform as good as the baseline LSTM model,  
854 since we detected a "repeating word pattern" in many of our outputs. For instance, for the image  
855 with several couches, architecture 2 outputs "a room with a couch, a couch, a couch, a couch, a couch,".  
856 We believe this output is due to consistently introducing the image features to the LSTM  
857 inputs. When the image contains lots of couches, the image features may strongly incline to the  
858 output "couches" in different time periods. Therefore, the word "couches" is predicted again and  
859 again, which is definitely not a desired output. Overall, we believe that the baseline LSTM model is  
860 the best among these three models.

864  
865  
**9 Team Contributions**

866 Wenbo Hu: Participated in implementation of all networks. Coding of stochastic generating caption,  
867 lstm and vanilla baseline. Tuning hyperparameters and finalizing results.

868 Wenxiao Li: wrote the code for samplers, bleu scores, good/bad image generation, tuned hyperpa-  
869 rameters for the baseline LSTM model. Implemented the learning and testing procedures of archi-  
870 tecture 2. Wrote reports for baseline LSTM and architecture 2.

871 Huaning Liu: Co-implementing for baseline model as well as vanilla RNN model. Hypertuning for  
872 RNN model. Organize the latex report, and write abstract, introduction, related work, model, and  
873 output analysis for Vanilla RNN model.

874  
875 **References**  
876

- 877 [1] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, “Show and Tell: Lessons learned from the  
878 2015 MSCOCO Image Captioning Challenge,” *arXiv e-prints*, p. arXiv:1609.06647, Sep. 2016.