

# Image Caption Generator Bot Based On Deep Neural Networks

Mr. Mukund Upadhyay, Ms. Shallu Bashambu

Submitted: 05-06-2021

Revised: 18-06-2021

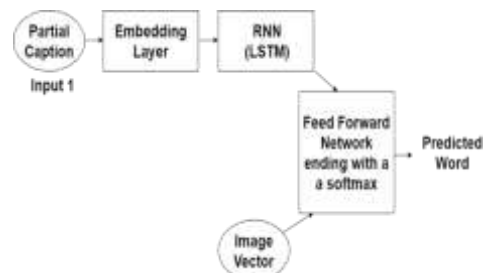
Accepted: 20-06-2021

## ABSTRACT:

In this project, we systematically analyze a deep neural networks based image caption generation method. With an image as the input, the method can output an English sentence describing the content in the image. We analyze three components of the method: convolutional neural network (CNN), recurrent neural network (RNN) and sentence generation. By replacing the CNN part with the state-of-the-art architectures, we find the VG-GNet performs best according to the BLEU score. We also propose a simplified version of the Gated Recurrent Units (GRU) as a new recurrent layer, implementing by both MATLAB and C++ in Caffe. The simplified GRU achieves comparable result when it is compared with the long short-term memory (LSTM) method. But it has few parameters which saves memory and is faster in training. Finally, we generate multiple sentences using Beam Search. The experiments show that the modified method can generate captions comparable to the state-of-the-art methods with less training memory.

## I. INTRODUCTION

Automatically describing the content of images using natural languages is a fundamental and challenging task. It has great potential impact. For example, it could help visually impaired people better understand the content of images on the web. Also, it could provide more accurate and compact information of images/videos in scenarios such as image sharing in social network or video surveillance systems. This project accomplishes this task using deep



**Figure 1:** Image caption generation pipeline. The framework consists of a convolutional neural network (CNN) followed by a recurrent neural network (RNN). It generates an English sentence from an input image.

neural networks. By learning knowledge from image and caption pairs, the method can generate image captions that are usually semantically descriptive and grammatically correct.

Human beings usually describe a scene using natural languages which are concise and compact. However, machine vision systems describe a scene by taking an image which is a two-dimensional array. From this perspective, Vinyalet et al. (Vinyalet et al.,) model the image captioning problem as a language translation problem in their Neural Image Caption (NIC) generator system. The idea is mapping the image and captions to the same space and learning a mapping from the image to the sentences. Donahue et al. (Donahue et al.,) proposed a more general Long-term Recurrent Convolutional Network (LRCN) method. The LRCN method not only models the one-to-many (words) image captioning, but also models many-to-one action generation and many-to-many video description. They also provide publicly available implementation based on Caffe framework (Jia et al., 2014), which further boosts the research on image captioning. This work is based on the LRCN method.

Although all the mappings are learned in an end-to-end framework, we believe the benefits of better understanding of the system by analyzing different components separately. Fig.1 shows the pipeline. The model has three components. The first component is a CNN which is used to understand the content of the image. Image understanding answers the typical questions in computer vision such as “What are the objects?”, “Where are the objects?” and “How are the objects interactive?”. For example, the CNN has to recognize the “teddy bear”, “table” and their relative locations in the image. The second component is a RNN which is used to generate a sentence given the visual feature. For example, the RNN has to generate a sequence of probabilities of words given two words “teddy bear, table”. The third component is used to generate a sentence by exploring the combination of the probabilities. This component is less studied in the reference paper (Donahue et al.).

This project aims at understanding the impact of different components of the LRCN method (Donahue et al.). We have following contributions: understand the LRCN method at the implementation level.

analyze the influence of the CNN component by replacing three CNN architectures (two from author’s and one from our implementation). analyze the influence of the RNN component by replacing two RNN architectures (one from author’s and one from our implementation). analyze the influence of sentence generation method by comparing two methods (one from author’s and one from our implementation).

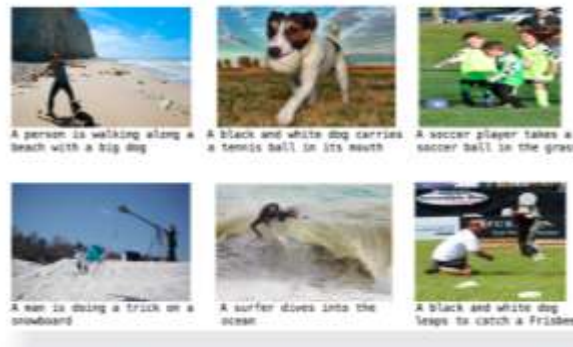
## II. RELATED WORK

Automatically describing the content of an image is a fundamental problem in artificial intelligence that connects computer vision and natural language processing. Earlier methods first generate annotations (i.e., nouns and adjectives) from images (Sermanet et al., 2013; Russakovsky et al., 2015), then generate a sentence from the annotations (Gupta and Mannem, 2015). Donahue et al. (Donahue et al., 2015) developed a recurrent convolutional architecture suitable for large-scale visual learning, and demonstrated the value of the models on three different tasks: video recognition, image description and video description. In these models, long-term dependencies are incorporated into the network state updates and are end-to-end trainable. The limitation is the difficulty of understanding the intermediate result. The LRCN method is further

developed to text generation from videos (Venugopal et al., 2015).

Instead of one architecture for three tasks in LRCN, Vinyals et al. (Vinyals et al., 2015) proposed a neural image caption (NIC) model only for the image caption generation. Combining the GoogLeNet and single layer of LSTM, this model is trained to maximize the likelihood of the target description sentence given the training images. The performance of the model is evaluated qualitatively and quantitatively. This method was ranked first in the MS COCO Captioning Challenge (2015) in which the result was judged by humans. Comparing LRCN with NIC, we find three differences that may indicate the performance differences. First, NIC uses GoogLeNet while LRCN uses VGGNet. Second, NIC inputs visual feature only into the first unit of LSTM while LRCN inputs the visual feature into every LSTM unit. Third, NIC has simpler RNN architecture (single layer LSTM) than LRCN (two factored LSTM layers). We verified that the mathematical models of LRCN and NIC are exactly the same for image captioning. The performance difference lies in the implementation and LRCN has trade-off between simplicity and generality, as it is designed for three different tasks.

Instead of end-to-end learning, Fang et al. (Fang et al., 2015) presented a visual concepts based method. First, they used multiple instance learning to train visual detectors of words that commonly occur in captions such as nouns, verbs, and adjectives. Then, they trained a language model with a set of over 400,000 image descriptions to capture the statistics of word usage. Finally, they re-ranked caption candidates using sentence-level features and a deep multi-modal similarity model. Their captions have equal or better quality 34% of the time than those written by human beings. The limitation of the method is that it has more human controlled parameters which make the system less reproducible. We believe the web application **captionbot** (Microsoft, 2015) is based on this method.



**Figure 2:** This image shows a group of picture with their captions generated

Karpathy et al. (Karpathy and Fei-Fei) proposed a visual-semantic alignment (VSA) method. The method generates descriptions of different regions of an image in the form of words or sentences (see Fig. 2). Technically, the method replaces the CNN with Region-based convolutional Networks (RCNN) so that the extracted visual features are aligned to particular regions of the image. The experiment shows that the generated descriptions significantly outperform retrieval baselines on both full images and on a new dataset of region-level annotations. The experiment is performed from human beings using Amazon's Mechanical Turk (AMT). We manually checked some examples by side-by-side comparing the image and corresponding sentences. We found the captions are very expressive and diverse. The COCO Caption is the largest image caption corpus at the time of writing. There are 413,915 captions for 82,783 images in training, 202,520 captions for 40,504 images in validation and 379,249 captions for 40,775 images in testing. Each image has at least 5 captions. The captions for training and validation are publicly available while the captions for testing is reserved by the authors. In the experiment, we use all the training data in the training process and 1,000 randomly selected validation data in the testing process.

### III. DESCRIPTION OF PROBLEM

**Task** In this project, we want to build a system that can generate an English sentence that describes objects, actions or events in an RGB image:

$$S = f(I) \quad (1)$$

where  $I$  is an RGB image and  $S$  is a sentence,  $f$  is the function that we want to learn.

**Corpus** We use the MS COCO Caption (Chen et al., 2015) as the corpus. The captions are gathered where the  $\theta$  is dropped for convenience,  $S_t$  is

the word at step  $t$ .

The model has two parts. The first part is a CNN which maps the image to a fixed-length visual feature. The visual feature is embedded to the input of the RNN.

$$v = W_v(\text{CNN}(I)) \quad (4)$$

where  $W_v$  is the visual feature embedding. The visual feature is fixed for each step of the RNN.

In the RNN, each word is represented as a one-hot vector  $S_t$  of dimension equal to the size of the dictionary.  $S_0$  and  $S_N$  are for special start and stop words. The word embedding parameter is  $W_s$ :

$$x_t = W_s S_t, t \in \{0 \dots N-1\} \quad (5)$$

In this way, the image and words are mapped to the same space. After the internal processing of the RNN, the features  $v$ ,  $x_t$  and internal hidden parameter  $h_t$  are decoded into a probability to predict the word at current time:

$p = \text{LSTM}(v, x, h), t \in \{0, \dots, N-1\}$  (6) Because the sentence with higher probability does not necessarily mean this sentence is more accurate than other candidate sentences, post-processing methods such as **Beam Search** is used to generate more sentences and pick top-K sentences.

#### IV. METHOD

For image caption generation, LRCN maximizes the probability of the description giving the image: This method generates more diverse and accurate descriptions than the whole image methods such as

LRCN and NIC. The limitation is that the method consists of two separate models. This method is further developed to dense captioning (Johnson et al., 2016) and image based question and answer system (Zhu et al., 2016).

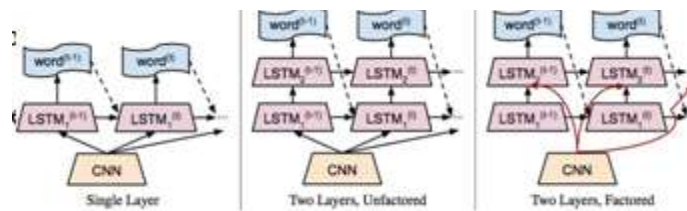
$$\theta^* = \underset{\theta}{\text{argmax}} \log p(S|I; \theta) \quad (2)$$

where  $\theta$  are the parameters of the model,  $I$  is an image, and  $S$  is a sample sentence. Let the length of the sentence be  $N$ , the method applies the chain rule to model the joint probability over  $S_0, \dots, S_N$ :

$$\log p(S|I) = \sum_{t=0}^{N-1} \log p(S_t | I, S_{0:t-1}) \quad (3)$$

#### 4.1 Convolutional neural network

In this project, a convolutional neural network (CNN) maps an RGB image to a visual feature vector. The CNN has three most-used layers: convolution, pooling and fully-connected layers. Also, Rec-



The most right two-layers factored LSTM is used in architecture.

**Figure 3:** Three variations of the LRCN image captioning architecture. Figure from (Donahue et al.,).

#### 4.2 Recurrent neural network

To prevent the gradients vanishing problem, the long short-term memory (LSTM) method is used as the RNN component. A simplified LSTM updates the hidden state using the non-linear active function. The ReLU is faster than the traditional  $f(x) = \tanh(x)$  or  $f(x) = (1 + e^{-x})^{-1}$ . Dropout layer is used to prevent overfitting. The dropout sets the output of

$$o_t = \sigma(W_{x_o}x_t + W_{h_o}h_{t-1} + b_o)$$

$$g_t = \phi(W_{x_c}x_t + W_{h_c}h_{t-1} + b_c) c_t = f_t \odot c_{t-1} + i_t \otimes g_t$$

$$h_t = o_t \odot \phi(c_t) \quad (7)$$

for time step  $t$  given inputs  $x_t, h_{t-1}$ , and  $c_{t-1}$  are:

$$i_t = \sigma(W_{x_i}x_t + W_{h_i}h_{t-1} + b_i)$$

$$f_t = \sigma(W_{x_f}x_t + W_{h_f}h_{t-1} + b_f)$$

each hidden neuron to zero with a probability (i.e., 0.5). The “dropped out” neurons do not contribute to the forward pass and do not participate in back-propagation.

The AlexNet (Krizhevsky et al., 2012), VGGNet (Simonyan and Zisserman, 2014) and GoogLeNet (Szegedy et al., 2015) are three widely used deep convolutional neural network architectures. They share the convolution pooling fully-connection loss function pipeline but with different shapes and connections of layers, especially the convolution layer. AlexNet is the first deep convolutional neural network used in large scale image classification. VGGNet and GoogLeNet achieve the start-of-the-art performance in ImageNet recognition challenge 2014 and 2015.

When the CNN combines the RNN, there are special considerations of convergence since both of them have millions of parameters. For example, Vinyals et al. (Vinyals et al.,) found that it is better to fix the parameters of the convolutional layer as the parameters trained from the ImageNet. As a result, only the non-convolution layer parameters in CNN and the RNN parameters are actually learned from caption examples.

where  $\sigma(x) = (1 + e^{-x})^{-1}$  and  $\phi(x) = 2\sigma(2x)$

In addition to a hidden unit  $h_t \in \mathbb{R}^N$ , the LSTM includes an input gate  $i_t \in \mathbb{R}^N$ , forget gate  $f_t \in \mathbb{R}^N$ , output gate  $o_t \in \mathbb{R}^N$ , input modulation gate  $g_t \in \mathbb{R}^N$ , and memory cell  $c_t \in \mathbb{R}^N$ . These additional cells enable the LSTM to learn extremely complex and long-term temporal dynamics. Additional depth can be added to LSTMs by stacking them on top of each other. Fig. 3 shows three versions of LSTMs. The two-layer factored LSTM achieves the best performance and is used in the method.

In this project, we proposed a simplified version of GRU in section 5.1 which also avoids the vanishing gradient problem and can be easily implemented in Caffe based on the current Caffe LSTM framework. We also provide the MATLAB program in the Appendices verifying our derivation of BPTT on the original GRU model.

**4.3 Sentence generation**  
The output of LSTM is the probability of each word in the vocabulary. **Beam search** is used to generate sentences. Beam search is a heuristic search algorithm that explores a graph by expanding the most promising node in a limited set. In addition to beam search, we also use k-best search to generate sentences. It is very similar to the time

synchronous Viterbi search. The method iteratively selects the k best sentences from all the candidate sentences up to time t, and keeps only the resulting best k of them.

## V. IMPLEMENTATION

are:

$$\begin{aligned} z &= \sigma(U_z x_t + W_z s_{t-1} + b_z) \\ r &= \sigma(U_r x_t + W_r s_{t-1} + b_r) \\ h &= \tanh(U_h x_t + W_h (s_{t-1} \odot r) + b_h) \\ s_t &= (1-z) \odot h + z \odot s_{t-1} \end{aligned} \quad (8)$$

**Preprocessing** Because we want to keep the architecture of the CNN, the input image are randomly cropped to the size of 224x224. As a result, only part of the images are used in training at particular iteration. Because one image will be cropped multiple times in the training, the CNN can probably see the whole image in the training (once for part of the image).

However, the method only sees part of the image in the testing except the dense cropping is also used (our project does not use dense crop). For the sentences, the method first creates a vocabulary only from the training captions and removes lower frequency words (less than 5). Then, words are represented by one-hot vectors.

### 5.1 Caffe architecture

Caffe (Jia et al., 2014) provides a modifiable framework for the state-of-the-art deep learning algorithms. It is implemented using C++ and also provides Python and MATLAB interfaces. Caffe model (network) definitions are written as a configuration where  $z$  is the update gate,  $r$  is the reset gate.  $s$  is used as both hidden states and cell states. With few parameters, GRU can reach a comparable performance to LSTM (Jozefowicz et al.,). To implement GRU, we first wrote a MATLAB program to check our BPTT<sup>2</sup> gradient derivation. This is due to the fact that automatic differentiation in Caffe is not supported at layer units level. Followed by our derivation, the calculated gradients only deviate from the numerical gradients by around  $10^{-5}$  relatively. However, implementing GRU in Caffe is not straight forward since Caffe is based on a complicated software architecture trying to provide convenience for assembling, not further developing. This is the bottleneck of GRU implementation. We have tried a number of implementations based on the original GRU (Equ. 8), with no good results. Finally we simplified GRU model inspired by the simplified LSTM in (Donahue

et al., ). We omit the reset gate and add a transfer gate to make it easily fit into the current Caffe LSTM framework as:

$$z = \sigma(U_z x_t + W_z s_{t-1} + b_z)$$

files using the Protocol Buffer Language<sup>1</sup> so that the net representation and implementation are separated.

$$h = \tanh(U_h x_t + W_h s_{t-1} + b_h)$$

(9)

ated. The separation abstracts from memory underlying location in CPU or GPU so that switching between a CPU and GPU implementation is exactly by one function call. However, these separation makes the implementation less convenient as we will show in the next paragraph.

### 5.2 Simplify and implement GRU in Caffe

In Caffe, a **layer** is the fundamental unit of computation. A **blob** is a wrapper over the actual data providing synchronization capability between CPU and GPU. We tried to implement the Gated Recurrent Units (GRU) (Cho et al., 2014) in Caffe. The GRU updates for time step  $t$  given inputs  $x_t, s_{t-1}$

$$c_t = (1-z) \odot h + z \odot c_{t-1}$$

$$s_t = c_t$$

Note that the omitted reset gate won't bring back the vanishing gradient problem which we see in traditional RNN because we still have the update gate acting as a weight between the previous state and the current processed input. The added transfer gate,  $c$ , seems to be less useful, but it is actually very important for calculating the gradient in the framework. The parameter gradients in an RNN within a single step,  $t$ , depends not only on  $\partial L_t / \partial s_t$ , but also  $\partial L_t / \partial s_{t-i}$  where  $i=1, 2, \dots, t$ . In Caffe,  $\partial L_t / \partial s_t$  is calculated by outer layers automatically, while  $\partial L_t / \partial s_{t-i}$  need to be calculated by inside layer unit. To hold and transfer these two parts of gradient to

<https://developers.google.com/protocol-buffers/docs/proto>  
<sup>2</sup>Backward propagation through time

CNNs	layer	#param	memory	B-4	
	Method	B-1	B-2	B-3	B-4
AlexNet	8	60	0.9	0.253	
VGGNet	16	138	11.6	<b>0.294</b>	
GoogLeNet	22	12	5.8	0.211	

**Table 1: Quantitative comparison of CNNs.** The number of parameter (#param) is in the unit of million, and the training memory is in the unit of Gb. In experiment, we found that the BLEU 4 performance is positively related to the number of parameters.

thenexttimestep,weuseanotherintermediatevari-

AlexNet  
+LSTM  
 AlexNet  
+GRU  
 VGGNet  
+LSTM  
 VGGNet  
+GRU

0.650	0.467	0.324	0.221
0.623	0.433	0.292	0.194
0.588	0.406	0.264	0.168
0.583	0.393	0.256	0.168

**Table2: AlexNet, VGGNet with different RNN models.** Our GRU model achieves comparable result with the LSTM model, but with less parameter and training time. The beam size is 1.

able, which is the added transfer gate c. This is just an engineering issue that might not be avoided while developing new models in Caffe. The theory is always clear and concise (see Appendices for the MATLAB program verifying our BPTT derivation to the original GRU).

### 5.3 Training method

The neural network is trained using the mini-patch stochastic gradient descent (SGD) method. The base learning rate is 0.01. The learning rate drops 50% in every 20,000 iterations. Because the number of training samples is much smaller than the number of parameters of the neural network, overfitting is our big concern. Besides the dropout layer, we fixed the parameters of the convolutional layers as suggested by (Vinyal et al.,). All the networks are trained in a Linux machine with a Tesla K40c graphics card with 12Gb memory.

### 5.4 Quantitative result

**Evaluation metrics** We use BLEU (Papineni et al., 2002) to measure the similarity of the captions generated by our method and human beings.

BLEU is a popular machine translation metric that analyzes the co-occurrences of n-grams between the candidate and reference sentences. The unigram scores (B-1) account for the adequacy of the translation, while longer n-gram scores (B-2, B-3, B-4) account for the fluency.

**Different CNNs** Table 1 compares the performance of three CNN architectures (the RNN part uses LSTM). The VGGNet achieves the best perfor-

mance (BLEU 4) and GoogLeNet has the lowest score. It is out of our expectation at first because GoogLeNet achieves the best performance in the ImageNet classification task. We discussed this phenomenon with our fellow students. One of them pointed out that despite its slightly weaker classification performance, the VGGNet features other forms of those of GoogLeNet in multiple transfer learning tasks (Karpathy, 2015). A downside of the VGGNet is that it is more expensive to evaluate and it uses a lot more memory (11.6 Gb) and parameters (138 million). It takes more time to train VGGNet and GoogleNet than AlexNet (about 8 hours vs 4 hours).

**Different RNNs** Table 2 compares the performance of LSTM and GRU. The GRU model achieves comparable results with less parameters and training time.

**Different sentence generation methods** Table 3 also analyzes the impact of beam size in the Beam Search for different CNN architectures. In general, larger beam size achieves higher BLEU score.

This phenomenon is much more obvious in the VGGNet than other two CNNs. When the beam size is 1, AlexNet outperforms VGGNet. When the beam size is 10, the VGGNet outperforms AlexNet. The most probable reason is that AlexNet is good at detecting a single or few objects in an image while VGGNet is good at detecting multiple objects in the same image. When the beam size becomes larger, the VGGNet-based method can generate more accurate sentences.

#beam	B-1	B-2	B-3	B-4
AlexNet				
1	<b>0.650</b>	0.467	0.324	0.221
5	<b>0.650</b>	0.467	0.343	0.247
10	0.644	<b>0.474</b>	<b>0.347</b>	<b>0.253</b>
VGGNet				
1	0.588	0.406	0.264	0.168
5	0.632	0.450	0.310	0.212
10	<b>0.681</b>	<b>0.513</b>	<b>0.390</b>	<b>0.294</b>
GoogLeNet				
1	0.533	0.353	0.222	0.139
5	0.568	0.385	0.262	0.180
10	<b>0.584</b>	<b>0.410</b>	<b>0.292</b>	<b>0.211</b>

**Table 3: AlexNet, VGGNet and GoogleNet with different beam sizes.** Using AlexNet, the impact of the number of beam size is not significant. Using the VGG net, the impact is significant. Using the GoogLeNet net, the impact is moderate. The best scores are highlighted.

Method	B-1	B-2	B-3	B-4
LRCN	0.669	0.489	0.349	0.249
NIC	N/A	N/A	N/A	0.277
VSA	0.584	0.410	0.292	0.211
This project	0.681	0.513	0.390	0.294

**Table 4: Evaluation of image caption of different methods.** LRCN is tested on the validation set (5,000 images). NIC is tested on the validation set (4,000 images). VSA is tested on the test set (40,775 images). This project is tested on the validation set (1,000 images for B-1, B-2, B-3, and 100 images for B-4).

(VGGNet) Aman and woman sitting at a table with a pizza.  
 (GoogLeNet) A group of people sitting at a dinner table.  
 When beam size is 5, the captions are as follows, (AlexNet) A group of people sitting at a table. (VGGNet) Aman and woman sitting at a table with food.  
 (GoogLeNet) A group of people sitting at a dinner table.  
 When beam size is 10, the captions are as follows, (AlexNet) A group of people sitting at a table. (VGGNet) Aman and woman sitting at a table. (GoogLeNet) A group of people sitting at a dinner table.  
 From the result listed above, we can see that when the beam size is fixed, VGGNet can generate captions with more details. When the beam size increases,

the captions become short and detailed information disappears.  
 Although the sentence generated by our method has the highest probability, we don't know if there are other sentences that can describe the image better. So we use 3-best search to explore the top 3 captions. For Fig. 4, the captions generated by GoogLeNet with beam size 5 using 3-best search are listed as follows,  
 A group of people sitting at a dinner table.  
 A group of people sitting around a dinner table.  
 A group of people sitting at a dinner table with plates of food.  
 The above captions are listed in probability descending order. We can see that the third sentence is actually the best one, although it does not have the



**Comparison with other systems** Table 4 compares BLEU scores of the results from LRCN, NIC, VSA and this project. The BLEU score of the result of this project is comparable or better than those from other systems although our project is tested on less data set (1,000 images).

### 5.5 Qualitative result

Taking Fig. 4 as an example, we analyze the captions generated by AlexNet, VGGNet and GoogLeNet.

When beam size is 1, the captions are as follows, (AlexNet) A group of people sitting at a table with a pizza.

highest probability. This is because when the sen-

tence is long, it is more probable to make mistakes. So, sentences with high probability sometimes

tend to be short, which may miss some detailed information. However, it does not mean that the sentence with the highest probability is bad. In most cases we observed, sentences with the highest probability are good enough to describe an image while long sentences often include redundant information and often make grammatical mistakes.

Fig. 5 shows the good examples of the sentences generated by this project. Most of them successfully describe the main objects and events in images. Fig. 6 shows failed examples of the system. The errors

Task	Wenqiang	Minchen	Jianhui
CNN			100%
GRU		70%	30%
BeamSearch	100%		
Writing	40%	20%	40%

**Table 5:** Division of work. These only measure the implementation and experimenting workload. All the analyses and discussions are conducted by all of us.



**Figure 4:** Sample image for qualitative analysis.

are mainly from object mis-detections such as an airplane is mis-detected as a kite (row 3 column 1), cellphones are detected as laptop (row 4 column 2). The generated sentences are also has minor grammar error. For example, "A motorcycle with a motorcy-cle" (row 4 column 3) is hard to understand.

## VI. LESSONS LEARNED AND FUTURE WORK

This project provides a valuable learning experience. First, the LRCN method has a sophisticated pipeline so that modifying part of the pipeline is complica-

ted than we expected. We learned how to use one of the most popular deep learning frameworks Caffe through the project.

Second, mathematics and the knowledge of particular software architecture are equally important for the success of the project. Although we implemented the MATLAB version of GRU very early before the deadline of the project, we spent a large amount of time on implementing the GRU layer in Caffe. The benefit is that we learned valuable first-hand experience on the developing level of Caffe in-

stead of purely using existing layers in Caffe. Third, working in a team, we could discuss and refine a lot of initial ideas. We could also anticipate problems that could become critical of the cases we were working alone. Table 5 roughly shows the work division among team members.

## VII. EVALUATION

The project is successful. We have finished all the goals before the deadline. The system can generate sentences that are semantically correct according to the image. We also proposed a simplified version of GRU that has less parameters and achieves comparable result with the

LSTM method.

The strength of the method is on its end-to-end learning framework. The weakness is that it requires large number of human labeled data which is very expensive in practice. Also, the current method still has considerable errors in both object detection and sentence generation.

## VIII. CONCLUSION AND FUTURE WORK

We analyzed and modified an image captioning method LRCN. To understand the method deeply, we decomposed the method to CNN, RNN, and sentence generation. For each part, we modified or replaced the component to see the influence on the final result. The modified method is evaluated on the COCO caption corpus. Experiment results show that: first the VGGNet outperforms the AlexNet and GoogleLeNet in BLEU score measurement; second, the simplified GRU model achieves comparable results with more complicated LSTM model; third, increasing the beam size increase the BLEU score in general but does not necessarily increase the quality of the description which is judged by humans.

**Future work** In the future, we would like to explore methods to generate multiple sentences with different content. One possible way is to combine interesting region detection and image captioning.

## REFERENCES

- [1]. [Chen et al.2015] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollar, and C Lawrence Zitnick. 2015. Microsoft coco captions: Data collection and evaluation server. arXiv preprint arXiv:1504.00325.
- [2]. [Cho et al.2014] Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using rnn encoder-decoder for statistical machine translation. arXiv preprint arXiv:1406.1078.
- [3]. [Donahue et al.] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In IEEE CVPR.
- [4]. [Fang et al.] Hao Fang, Saurabh Gupta, Forrest Iandola, Rupesh K Srivastava, Li Deng, Piotr Dollar, Jianfeng Gao, Xiaodong He, Margaret Mitchell, John C Platt, et al. From captions to visual concepts and back. In IEEE CVPR.
- [5]. [Gupta and Mannem] Ankush Gupta and Prashanth Man-nem. From image annotation to image description. In Neural information processing.
- [6]. [Jia et al.2014] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. 2014. Caffe: Convolutional architecture for fast feature embedding. In Proceedings of the ACM International Conference on Multimedia, pages 675–678.
- [7]. [Johnson et al.2016] Justin Johnson, Andrej Karpathy, and Li Fei-Fei. 2016. Densecap: Fully convolutional localization networks for dense captioning. In IEEE CVPR.
- [8]. [Jozefowicz et al.] Rafal Jozefowicz, Wojciech Zaremba, and Ilya Sutskever. An empirical exploration of recurrent network architectures. In ICML.
- [9]. [Karpathy and Fei-Fei] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In IEEE CVPR.
- [10]. [Karpathy2015] Andrej Karpathy. 2015. Cs231n: Convolutional neural networks for visual recognition. <http://cs231n.github.io/>. [Online; accessed 11-April-2015].
- [11]. [Krizhevsky et al.2012] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet classification with deep convolutional neural networks. In Advances in neural information processing systems.

- [12]. [Microsoft] Microsoft.captiobot.<https://www.captiobot.ai/>. [Online;accessed22-April-2015].
- [13]. [Papineni et al.2002] KishorePapineni,SalimRoukos,Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a methodfor automatic evaluation of machine translation.InProceedingsofthe40thannualmeetingonassociationforcomputationallinguistics, pages311–318.
- [14]. [Russakovsky et al.2015] OlgaRussakovsky,JiaDeng,Hao Su,Jonathan Krause,Sanjeev Satheesh,SeanMa,ZhihengHuang,AndrejKarpathy,AdityaKhosla,Michael Bernstein, et al.2015.Imagenet large scalevisualrecognitionchallenge.IJCV,115(3):211–252.
- [15]. [Sermanet et al.2013] Pierre Sermanet, David Eigen, XiangZhang,MichaëlMathieu,RobFergus,andYannLeCun.2013.Overfeat:Integrated recognition, localizationanddetectionusingconvolutionalnetworks.arXivpreprintarXiv:1312.6229.
- [16]. [Simonyan and Zisserman2014] KarenSimonyanandAndrew Zisserman.2014.Very deep convolutionalnetworks for large-scale image recognition.arXivpreprintarXiv:1409.1556.
- [17]. [Szegedy et al.2015] ChristianSzegedy,WeiLiu,Yangqing Jia, Pierre Sermanet, Scott Reed, DragomirAnguelov, Dumitru Erhan, Vincent Vanhoucke, andAndrewRabinovich.2015.Goingdeeperwiththeconvolutions.InIEEECVPR.
- [18]. [Venugopalan et al.] SubhashiniVenugopalan,MarcusRohrbach,JeffreyDonahue,RaymondMooney,TrevorDarr ell,andKateSaenko.Sequencetosequence-videototext.InIEEEICCV.
- [19]. [Vinyals et al.] Oriol Vinyals, Alexander Toshev, SamyBengio, and Dumitru Erhan.Show and tell: A neuralimagecaptiongenerator.InIEEECVPR.
- [20]. [Zhu et al.2016] Yuke Zhu, Oliver Groth, Michael Bern-stein, and Li Fei-Fei.2016.Visual7W: GroundedQuestionAnsweringinImages.InI