

Long Short-Term Memory with Read-only Unit in Neural Image Caption Generator

Aghasi Poghosyan
Institute for Informatics and
Automation Problems of NAS RA
Yerevan, Armenia
agasy18@gmail.com

Hakob Sarukhanyan
Institute for Informatics and
Automation Problems of NAS RA
Yerevan, Armenia
hakop@ipia.sci.am

Abstract—Automated caption generation for digital images is one of the fundamental problems in artificial intelligence. Most of the existing works use Long Short-Term Memory as a recurrent neural network cell to solve this task. After training, their deep neural models can generate an image caption. But there is an issue, the next predicted word of the caption depends mainly on the last predicted word, rather than the image content. In this paper we present model that can automatically generate an image description and is based on a recurrent neural network with modified LSTM cell with an additional gate responsible for image features. This modification results in generation of more accurate captions. We have trained and tested our model on MSCOCO image dataset by using only images and their captions.

Keywords—Deep learning, image caption generation, RNN, LSTM.

I. INTRODUCTION

Automated image caption generation is a difficult task that can help visually impaired people better understand the content of images on the web. Also, it can have a great impact on search engines and in robotics. This task is significantly harder than the well-studied image classification [1] or object recognition.

Generated image caption must contain not only image object names, but their properties, relations, and actions. Moreover, the generated caption must be expressed through natural language like English. This means, that already pre-trained neural language model needs an additional visual information to generate an image caption.

There are number of works approaching this problem. Some of them [2] [3] [4] offer combining the existing image object detection and sentence generation systems. But there is a more efficient solution [5] that offers a joint model. It takes an image and generates the caption, which describes the image adequately.

Last achievements in statistical machine translation were actively used in image caption generation tasks. The reason for this is mainly the proven achievement of greater results when using a powerful sequential model trained by maximizing the probability of the correct translation for the input sentence. These models [6] [7] [8] are based on Recurrent Neural Networks (RNNs). The model encodes variable length input into fixed length vector representation. This representation enables conversion of the input sentence into the target sentence or the input image into the target image caption. The last model was being trained to maximize $p(S|I)$ likelihood to generate the target sequence of words $S = \{S_1, S_2, \dots\}$ for an input image I , where each word S_t comes from a given dictionary, that describes the image adequately.

Modern public image datasets like MSCOCO [9] release 2014 contains 82,783 training, 40,504 validation, and 40,775 testing images. Each image form training and validation sets has 4-5 captions. The most frequently used 100 words make up 40.5% of total words used in captions (see Fig. 1) (top frequently used pronouns, determiners and words that used less than 4 times aren't included in calculations).

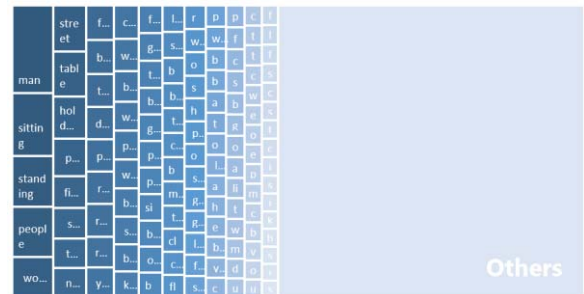


Fig. 1. MSCOCO images caption words distribution. Rectangle area represents each word's usage.

This distribution enables RNN to ignore the image scene or objects when predicting subsequent word of the sentence. This way, the next word prediction generally depends on the previous word. The generated sequence mainly contains words with high frequency of being met in the dataset (see Fig. 2).

This is why the generated caption (see Fig. 3) contains irrelevant words like 'man', 'standing', 'street', 'bench,' 'holding', 'cellphone'.

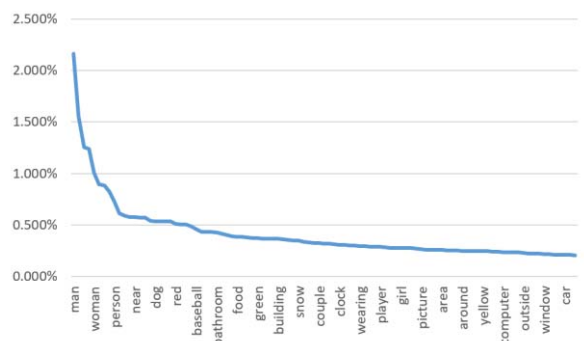


Fig. 2. MSCOCO image's captions top meeting words (chart shows sampled words from top 100, excluding pronouns and determiners)



Fig. 3. “a man and a woman sitting on a bench with their cell phones”.
Caption generated with [5]

II. MODEL

Machine translation based models that can generate image descriptions actively use a recurrent neural network. We will maximize the probability of the correct caption for the given image (1)

$$\theta^* = \arg \max \sum_{(I,S)} \log p(S|I; \theta). \quad (1)$$

In formulation (1) θ represents the parameters of our model and S is its correct caption for the given image I .

If we have a sentence $S = \{S_0, S_1, \dots, S_N\}$ with the length of N (where each $S_i, 0 \leq i < N$ is the index of the word in words dictionary, S_0 and S_N are correspondingly special *start* and *end* word indexes, indicating the beginning and end of the sentence), then we can apply the chain rule to calculate the joint probability (2) over S_0, S_1, \dots, S_N ,

$$\log p(S|I; \theta) = \sum_{t=0}^N \log p(S_t|I, S_0, \dots, S_{t-1}; \theta), \quad (2)$$

where (I, S) is a training example pair. While training, we optimize the sum of the log probabilities for the whole training set using stochastic gradient descent [10].

$p(S_t|I, S_0, \dots, S_{t-1}; \theta)$ probability will correspond to the t step (iteration) of Recurrent Neural Network (RNN) based model. The variable number of words that are conditioned upon, up to $t - 1$ is expressed by a fixed length hidden state or memory h_t . After every iteration for the new input, x_t memory will be updated (3) by using a non-linear function f .

$$h_{t+1} = f(h_t, x_t). \quad (3)$$

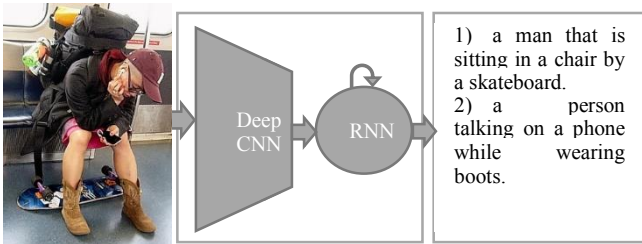


Fig. 4. Model scheme that generates complete captions in natural language from an input image. (MSCOCO sample)

For f we use a Long-Short Term Memory (LSTM), which has shown state-of-the-art performance on sequence generation tasks, such as translation or image caption generation. Model (see Fig. 4) consists of the feed forward deep convolutional neural network (CNN) that feeds RNN.

One of the best Convolutional Neural Networks (CNN) is *Google Inception* [11], that has been widely used in object classification and object detection tasks. Furthermore, there are

works [12] that have done CNN transfer learning for object classification for such tasks as scene classification.

There are high-level features that describe image semantic content like objects, their properties and relations in CNN [13]. In this work, we will select *Mixed_7c* layer from *Google Inception* and append *average pooling* layer which will have 2048-dimensional output for image description. Also, we will append *fully connected* neural layer with N_e neurons, which will convert 2048-dimensional vector into N_e dimensional vector. N_e is an image-words embedding vector's dimensionality [14]. The output vector x_{-1} of fully connected layer (4) will be the first feed vector for RNN,

$$x_{-1} = \text{Mixed}_{7c} * W_i + b_i, \quad (4)$$

where $W_i \in R^{2048 \times N_e}$ and $b_i \in R^{N_e}$ are trainable parameters for image embedding.

We also have lookup embedding matrix $W_l \in R^{D \times N_e}$, where D is the dictionary's words count. Each row of the matrix represents a word embedding in image-word embedding space. Each x_i (where $i \geq 0$) is corresponding row at index (S_i) (5).

$$x_i = W_e^{S_i}. \quad (5)$$

A. LSTM

Long Short-Term Memory (LSTM) is an RNN cell. It helps in solving RNN training time problems like vanishing and exploding gradients [15], which is a significant problem for RNNs. LSTM is commonly used in machine translation, sequence generation and image description generation tasks. Work [5] uses recurrent neural network with an LSTM cell to generate image caption.

From construction perspective LSTM is a memory cell c encoding knowledge at every iteration of what inputs have been seen up to this iteration. Later this knowledge is used for subsequent word generation (10, 11). Behavior of the cell is controlled by three *gates* – *input gate*, *output gate* and *forget gate*. Each gate is a vector of real number elements ranging from 0 to 1. In particular (see Fig. 5), *forget gate* is responsible for controlling whether to forget the cell's old value, *input gate* controls the permission for reading a new input value and finally *output gate* controls the permission to output the new value from the cell. This is done by multiplying the given gate with the corresponding value (9, 10). The definition of the LSTM is as follows:

$$i_t = \sigma(W_{ix}x_t + W_{im}m_{t-1}), \quad (6)$$

$$f_t = \sigma(W_{fx}x_t + W_{fm}m_{t-1}), \quad (7)$$

$$o_t = \sigma(W_{ox}x_t + W_{om}m_{t-1}), \quad (8)$$

$$c_t = f_t \odot c_{t-1} + i_t \odot h(W_{cx}x_t + W_{cm}m_{t-1}), \quad (9)$$

$$m_t = o_t \odot c_t, \quad (10)$$

$$p_{t+1} = \text{softmax}(W_{pm} * m_t). \quad (11)$$

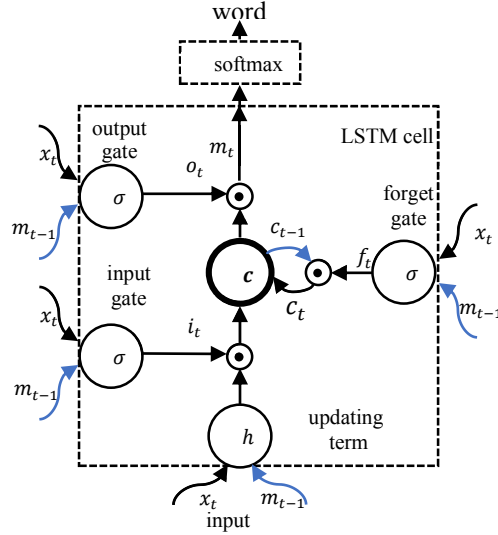


Fig. 5. LSTM: the memory block contains a cell c which is controlled by three gates.

In (7-11) equations i_t, o_t, f_t are *input, output* and *forget gates* correspondingly, c_t is a cell memory in step t and m_t is an output of the LSTM for step (iteration) t . $W_{ix}, W_{im}, W_{fx}, W_{fm}, W_{ox}, W_{om}, W_{cx}, W_{cm}$ are trainable parameters (variables) of the LSTM. \odot represents the product with a gate value. Sigmoid $\sigma(\cdot)$ and hyperbolic tangent $h(\cdot)$ are nonlinearities of the LSTM. The equation (11) will produce a probability distribution p_{t+1} over all words in the dictionary, where W_{pm} is a trainable parameter.

B. Training

The LSTM model is trained to predict the probability for the next word of an image caption after it has observed all the previous words in the captions and image features. For easier training LSTM is represented in unrolled form (see Fig. 6), which is a copy of the LSTM memory for the image and each word of the sentence. Also all LSTMs share the same parameters.

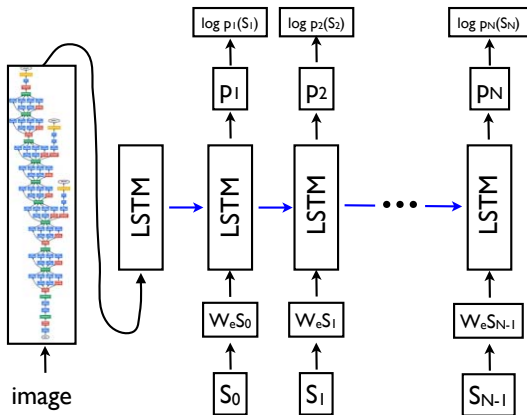


Fig. 6. LSTM model combined with a CNN image embedder (as defined in [5]) and words embedding.

Thus, x_{-1} is the first input for the first LSTM. Initial state of the LSTM is c_{-1} zero-filled memory. For the next LSTMs, inputs correspond to the word embedded vectors. Also, all recurrent connections are converted into feed-forward connections.

For the input image I and the image's true caption $S = \{S_0, S_1, \dots, S_N\}$, the unrolling procedure is:

$$x_{-1} = CNN(I), \quad (12)$$

$$x_t = W_e S_t, \quad (13)$$

$$p_{t+1} = LSTM(x_t), \quad (14)$$

where each word S_t is the row of square identity $N_d \times N_d$ matrix at corresponding index, where N_d is the dictionary size.

Note that both the image and the words are mapped to the same space. Vision CNN's last fully connected layer maps the image content to the embedding space. Also words are embedded by words embedding (5).

Loss [5] is the sum of the negative log likelihood of the correct word at each step:

$$L(I, S) = -\sum_{t=1}^N \log p(S_t). \quad (15)$$

After training the model by minimizing loss with gradient descent, we will have all the parameters of the LSTM, the top layer of the image embedder CNN and word embedding W_e .

C. LSTM with Read-only Unit

After one million training iterations on MSCOCO image dataset by using Tensorflow [16] framework, we got a loss function graphic (see green graphic in Fig. 7).

Experiments have shown that the LSTM's next prediction mainly depends on the previous word, that is why the LSTM generates words that are not associated with an input image.

We will add an additional unit to the LSTM, that will help to create a new LSTM state which depends on the image content. This additional unit will be a new gate the value of which will be included in the state calculation (see Fig. 8).

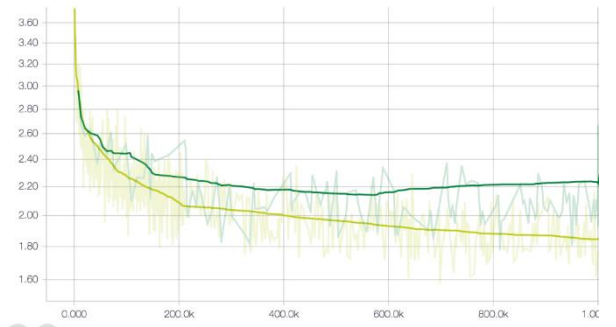


Fig. 7. Train loss function, LSTM – green, LSTM with read-only memory – yellow.

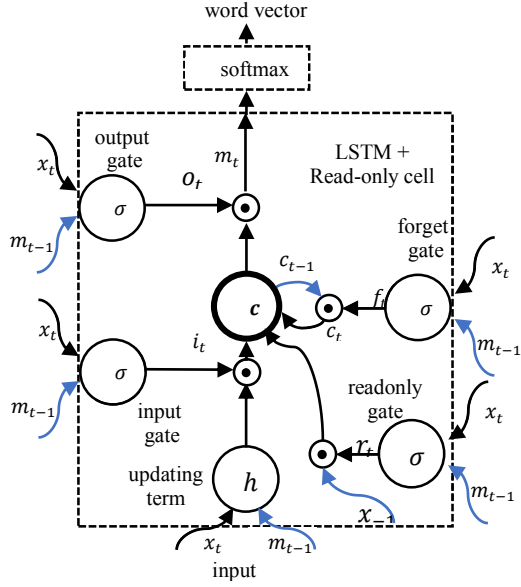


Fig. 8. LSTM: the memory block contains a cell c which is controlled by four gates (additional read-only memory).

$$r_t = \sigma(W_{rx}x_t + W_{rm}m_{t-1}), \quad (16)$$

$$c_t = f_t \odot c_{t-1} + r_t \odot x_{t-1} + i_t \odot h(W_{cx}x_t + W_{cm}m_{t-1}). \quad (17)$$

In (16) r_t is the read-only gate with W_{rx} and W_{rm} additional trainable parameters. New state c_t is calculated as shown in (17). After retraining with the new gate, we can see the new training loss (see yellow graphic in Fig. 7). We have achieved better results for training and for evaluation. Evaluation loss values for LSTM and LSTM with Read-only Unit models are 2.05 and 1.90 accordingly.

D. Model loss optimization

To minimize loss function for our model, we use Tensorflow framework which computes gradients for each trainable parameter (variable) with backpropagation.

1) Fprop: visit nodes in topo-sort order

- Compute value of node given predecessors

2) Bprop:

- initialize output gradient = 1
- visit nodes in reverse order:

Compute gradient wrt each node using gradient wrt successors (18)

$$\frac{\partial z}{\partial x} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x}. \quad (18)$$

We use mini-batch stochastic gradient descent to update trainable parameters (18).

$$\theta^{new} = \theta^{old} - \alpha \Delta_{\theta} I_{t:t+B}(\theta). \quad (19)$$

We clip gradients with norm clipping whenever they explode. The pseudo-code for the gradients norm clipping is:

$$\begin{aligned} \hat{g} &\leftarrow \frac{\partial \varepsilon}{\partial \theta} \\ \text{if } \|\hat{g}\| &\geq \text{threshold then} \\ \hat{g} &\leftarrow \frac{\text{threshold}}{\|\hat{g}\|} \hat{g} \end{aligned}$$

E. Inference

We inference by using Beam Search which gives us variants for the best scored sentence after many predictions. The pseudo-code for Beam Search with N size is:

SCORES := initial scores

SEQUENCES := start words

while (not end word for in all *SEQUENCES*)

LAST_WORDS := *SEQUENCES*[*LAST*]

NEXT_WORDS := *PREDICT*(*LAST_WORDS*)

NEW_SCORES := calculate sequence score for each word from *NEXT_WORDS* with its corresponding last score from *SCORES*

TOP_SCORES := *TOP_N*(*NEW_SCORES*, N)

TOP_WORDS := select words form *NEXT_WORDS* corresponding to *TOP_SCORES*

SCORES := *TOP_SCORES*

SEQUENCES append *TOP_WORDS*

foreach (*WORD* in *TOP_WORDS*)

if (*WORD* is end word)







don't process corresponding sequence from *SEQUENCES* later

Some examples of inference with 3 beam size are presented in TABLE .

CONCLUSION

In this work a Recurrent Neural Network (RNN) with Long Short-Term Memory (LSTM) cell and a Read-only Unit has been developed. An additional unit has been added to work [5] that increases the model accuracy. Two models, one with the LSTM and the other with the LSTM and Read-only Unit have been trained on the same MSCOCO image train dataset. The best (average of minimums) loss values are 2.15 for LSTM and 1.85 for LSTM with Read-only Unit. MSCOCO image test dataset has been used for testing. Loss values for LSTM and LSTM with Read-only Unit model test are 2.05 and 1.90, accordingly. These metrics have shown that the new RNN model can generate the image caption more accurately.

TABLE I. IMAGE CAPTION GENERATED BY LSTM AND LSTM WITH READ-ONLY UNIT MODELS.

	LSTM	LSTM + Read-only cell
	<ol style="list-style-type: none"> 1) a couple of bikes parked next to each other. 2) a couple of bikes that are next to a building. 3) a couple of bikes parked next to a building 	<ol style="list-style-type: none"> 1) a building with a red door and a white door. 2) a building with a red door and a white door 3) a building with a red brick wall and a white building
	<ol style="list-style-type: none"> 1) a woman in a bikini standing on a beach. 2) a woman in a bikini standing on a beach holding an umbrella. 3) a young girl in a bikini standing on a beach. 	<ol style="list-style-type: none"> 1) a woman in a dress and a hat on a beach. 2) a girl in a dress and a hat on a beach 3) a woman in a dress and a hat on a beach
	<ol style="list-style-type: none"> 1) a woman sitting at a table with a laptop. 2) a woman sitting at a table with a laptop 3) a woman sitting at a table in front of a laptop . 	<ol style="list-style-type: none"> 1) a woman sitting at a table in a dress 2) a woman is sitting at a table with a hat. 3) a woman sitting at a table with a hat on .
	<ol style="list-style-type: none"> 1) a group of people flying kites in the sky. 2) a group of people flying kites in the air. 3) a group of kites flying in the sky. 	<ol style="list-style-type: none"> 1) a bunch of animals that are in the grass. 2) a bunch of animals that are standing in the grass. 3) a bunch of animals that are standing in the sand
	<ol style="list-style-type: none"> 1) a woman sitting on a beach with an umbrella. 2) a woman sitting on a beach chair holding an umbrella. 3) a woman sitting on a beach with an umbrella 	<ol style="list-style-type: none"> 1) a woman sitting on the beach under an umbrella. 2) a woman sitting on the beach with an umbrella. 3) a woman sitting on the beach under an umbrella
	<ol style="list-style-type: none"> 1) a man and a woman sitting on a bench with their cell phones. 2) a man and a woman sitting on a bench . 3) a couple of people that are sitting down 	<ol style="list-style-type: none"> 1) a man and a woman sitting on a bus. 2) a man and a woman are sitting on a bus. 3) a man and a woman are sitting on a train.

REFERENCES

- [1] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet large scale visual recognition challenge," *arXiv preprint arXiv:1409.0575*, 2014.
- [2] A. Farhadi, M. Hejrati, M. A. Sadeghi, P. Young, C. Rashtchian, J. Hockenmaier, and D. Forsyth, "Every picture tells a story: Generating sentences from images," in *ECCV*, pp. 12--29, 2010.
- [3] G. Kulkarni, V. Premraj, S. Dhar, S. Li, Y. Choi, A. C. Berg, and T. L. Berg, "Baby talk: Understanding and generating simple image descriptions," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 35, no. 12, pp. 2891--2903, 2011.
- [4] A. Karpathy and L. Fei-Fei, "Deep Visual-Semantic Alignments for Generating Image Descriptions," *IEEE Conference on Computer Vision and Pattern Recognition*, vol. 39, no. 4, pp. 664--676, 2015.
- [5] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan, "Show and Tell: Lessons learned from the 2015 MSCOCO Image Captioning Challenge," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 39, no. 4, pp. 652-663, 2017.
- [6] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *arXiv:1409.0473*, 2014.
- [7] K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," *arXiv preprint arXiv:1406.1078*, 2014.
- [8] I. Sutskever, O. Vinyals, and Q. Le, "Sequencetosequence learning with neural networks," *Advances in neural information processing systems*, pp. 3104-3112, 2014.
- [9] T.Y. Lin, M. Maire, S. Belongie, J. Hays, and P. Perona, "Microsoft coco: Common objects in context," in *European Conference on Computer Vision*, 2014, pp. 740--755.
- [10] L. Bottou, "Large-scale machine learning with stochastic gradient descent," in *Proceedings of COMPSTAT'2010*, pp. 177--186, 2010.
- [11] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens and Z. Wojna, "Rethinking the inception architecture for computer vision," *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2818--2826.
- [12] J. Donahue, Y. Jia, O. Vinyals, J. Hoffman, N. Zhang, E. Tzeng, and T. D. Decaf, "A deep convolutional activation feature for generic visual recognition," in *International conference on machine learning*, 2014, pp. 647--655.
- [13] A. Poghosyan and H. Sarukhanyan, "Image Visual Similarity Based on High Level Features of Convolutional Neural Networks," *Mathematical Problems of Computer Science*, vol. 45, pp. 138--142, 2016.
- [14] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *arXiv preprint arXiv:1301.3781*, 2013.
- [15] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735--1780, 1997.
- [16] M. Abadi, A. Ashish, B. Paul, B. Eugene, C. Zhifeng, C. Craig, S. C. Greg et al. "Tensorflow: Large-scale machine learning on heterogeneous distributed systems." *arXiv preprint arXiv:1603.04467* 2016.