

# SuperCaptioning: Image Captioning Using Two-dimensional Word Embedding

Baohua Sun  
Gyr Falcon Technology Inc.  
Milpitas, CA  
baohua.sun@gyrfalcontech.com

Lin Yang  
Gyr Falcon Technology Inc.  
Milpitas, CA

Michael Lin  
Gyr Falcon Technology Inc.  
Milpitas, CA

Charles Young  
Gyr Falcon Technology Inc.  
Milpitas, CA

Patrick Dong  
Gyr Falcon Technology Inc.  
Milpitas, CA

Wenhan Zhang  
Gyr Falcon Technology Inc.  
Milpitas, CA

Jason Dong  
Gyr Falcon Technology Inc.  
Milpitas, CA

## Abstract

Language and vision are processed as two different modal in current work for image captioning. However, recent work on Super Characters method shows the effectiveness of two-dimensional word embedding, which converts text classification problem into image classification problem. In this paper, we propose the SuperCaptioning method, which borrows the idea of two-dimensional word embedding from Super Characters method, and processes the information of language and vision together in one single CNN model. The experimental results on Flickr30k data shows the proposed method gives high quality image captions. An interactive demo is ready to show at the workshop.

## 1 Introduction

Image captioning outputs a sentence related to the input image. Current methods process the image and text separately [4, 13, 15, 14, 5, 6, 1, 2]. Generally, the image is processed by a CNN model to extract the image feature, and the raw text passes through embedding layer to convert into one-dimensional word-embedding vectors, e.g. a 300x1 dimensional vector. And then the extracted image feature and the word embedding vectors will be fed into another network, such as RNN, LSTM, or GRU model, to predict the next word in the image caption sequentially.

Super Characters method [9] is originally designed for text classification tasks. It has achieved state-of-the-art results on benchmark datasets for multiple languages, including English, Chinese, Japanese, and Korean. It is a two-step method. In the first step, the text characters are printed on a blank image, and the generated image is called Super Characters image. In the second step, the Super Characters image is fed into a CNN model



(a) “Four men in life jackets are riding in a bright orange boat”.



(b) “A woman in a black coat walks down the sidewalk holding a red umbrella”.



(c) “A man in a boat on a lake with mountains in the background”.



(d) “Four performers are performing with their arms outstretched in a ballet”.

Figure 1: Examples of generated image captions using the proposed SuperCaptioning method.

for classification. The CNN model is fine-tuned from pre-trained ImageNet model. The extensions of Super Characters method [8, 12, 11] also prove the effectiveness of two-dimensional embedding on different tasks.

In this paper, we address the image captioning problem by employing the two-dimensional word embedding from the Super Characters method, and the resulting method is named as SuperCaptioning method. In this method, the input image and the raw text are combined together through two-dimensional embedding, and then fed into a CNN model to sequentially predict the words in the image caption. The experimental results on Flickr30k shows that the proposed method gives high quality image captions. Some examples given by SuperCaptioning method are shown in Figure 1.

## 2 The Proposed SuperCaptioning Method

The SuperCaptioning method is motivated by the success of Super Characters method on text classification tasks. Super Characters method converts text into images. So it will be very natural to combine the input image and the image of the text together, and feed it into one single CNN model to predict the next word in the image caption sentence.

Figure 2 illustrates the proposed SuperCaptioning method. The caption is predicted sequentially by predicting the next word in multiple iterations. At the beginning of the caption prediction, the partial caption is initialized as null, and the input image is resized to occupying a designed portion (e.g. top) of a larger blank image as shown in Figure 2. Then the text of the current partial caption is drawn into the the other portion (e.g. bottom) of the larger image as well. The resulting image is called the SuperCaptioning image, which is then fed into a CNN model to classify the next word in the caption. The CNN model is fine-tuned from the ImageNet pre-trained model. The iteration continues until the next word is EOS (End Of Sentence) or the word count reaches the cut-length. Cut-length is defined as the maximum number of words for the caption.

Squared English Word (SEW) method is used to represent the English word in a squared space. For example, the word “child” occupies the same size of space as the word “A”, but each of its alphabet will only occupies  $\{1/\text{ceil}[\text{sqrt}(N)]\}^2$  of the word space, where  $N$  is five for “child” which has five alphabets,  $\text{sqrt}(\cdot)$  stands for square root, and  $\text{ceil}[\cdot]$  is rounding to the top.

The data used for training is from Flickr30k<sup>1</sup>. Each image in Flickr30k has 5 captions by different people, and we only keep the longest caption if it is less than 14 words as the ground truth caption for the training data. After this filtering, 31,333 of the total 31,783 images are left.

After comparing the accuracy of experimental results using different configurations for the font size, cut-length, and resizing of the input image, we finally set the font size to 31, cut-length to 14 words, and resizing the image size to 150x224 in the fixed-size SuperCaptioning image with 224x224 pixels, as shown in Figure 2.

The training data is generated by labeling each SuperCaptioning image as an example of the class indicated by its next caption word. EOS is labeled to the SuperCaptioning image if the response sentence is finished. The model used is SE-net-154 [3] pre-trained on ImageNet<sup>2</sup>. We fine-tuned this model on our generated data set by only modifying the last layer to 11571 classes, which is the vocabulary size of all the selected captions.

Figure 1 shows that the proposed SuperCaptioning method captions the number of objects in the image, as shown in Figure 1a “**Four men** ...”; and it also captions the colors of overlapping objects, as shown in Figure 1b “A woman in a **black coat** ... holding a **red umbrella**”; it captions the big picture of the background, as shown

<sup>1</sup><http://shannon.cs.illinois.edu/DenotationGraph/data/flickr30k.html>

<sup>2</sup><https://github.com/hujie-frank/SENet>

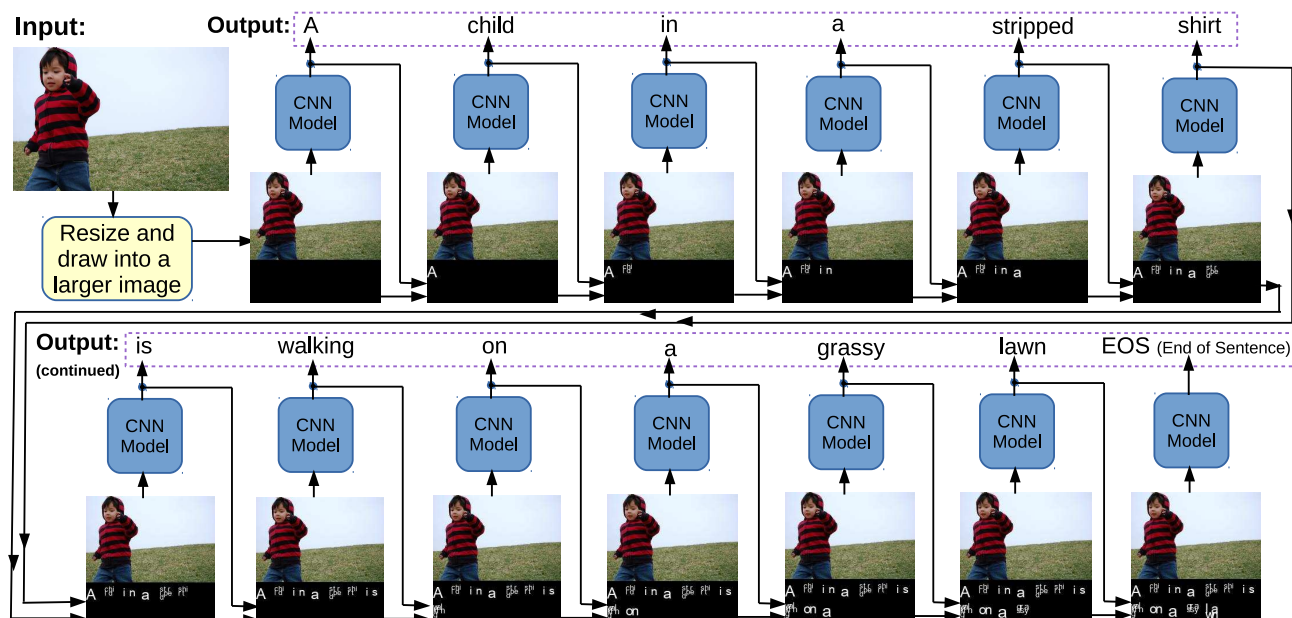


Figure 2: SuperCaptioning method illustration. The output caption is “A child in a striped shirt is walking on a grassy lawn”.

in Figure 1c “... with **mountains in the background**”; and it also captions the detailed activity, as shown in Figure 1d “... with their **arms outstretched** in a ballet”.

### 3 Conclusion

In this paper, we propose the SuperCaptioning method for image captioning using two-dimensional word embedding. The experimental results on Flickr30k shows that the SuperCaptioning method gives high quality image captions. The proposed method could be used for on-device image captioning applications with low-power CNN accelerator becoming more and more available [10, 7]. An interactive demo is ready to show at the workshop.

### References

- [1] Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. Bottom-up and top-down attention for image captioning and visual question answering. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 60776086, 2018.
- [2] MD Hossain, F Sohel, MF Shiratuddin, and H Laga. A comprehensive survey of deep learning for image captioning. ACM Computing Surveys, 51(6):136, 2019.
- [3] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-excitation networks. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 71327141, 2018.
- [4] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 31283137, 2015.
- [5] Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 375383, 2017.
- [6] Steven J Rennie, Etienne Marcheret, Youssef Mroueh, Jerret Ross, and Vaibhava Goel. Self-critical sequence training for image captioning. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, pages 7008 7024, 2017.
- [7] Baohua Sun, Daniel Liu, Leo Yu, Jay Li, Helen Liu, Wenhan Zhang, and Terry Torng. Mram co-designed processing-in-memory cnn accelerator for mobile and iot applications. arXiv preprint arXiv:1811.12179, 2018.

- [8] Baohua Sun, Lin Yang, Catherine Chi, Wenhan Zhang, and Michael Lin. Squared english word: A method of generating glyph to use super characters for sentiment analysis. arXiv preprint arXiv:1902.02160, 2019.
- [9] Baohua Sun, Lin Yang, Patrick Dong, Wenhan Zhang, Jason Dong, and Charles Young. Super characters: A conversion from sentiment classification to image classification. In Proceedings of the 9th Workshop on Computational Approaches to Subjectivity, Sentiment and Social Media Analysis, pages 309315, 2018.
- [10] Baohua Sun, Lin Yang, Patrick Dong, Wenhan Zhang, Jason Dong, and Charles Young. Ultra power-efficient cnn domain specific accelerator with 9.3 tops/watt for mobile and embedded applications. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops, pages 16771685, 2018.
- [11] Baohua Sun, Lin Yang, Michael Lin, Charles Young, Jason Dong, Wenhan Zhang, and Patrick Dong. Superchat: Dialogue generation by transfer learning from vision to language using two-dimensional word embedding and pretrained imagenet cnn models. arXiv preprint arXiv:1905.05698, 2019.
- [12] Baohua Sun, Lin Yang, Wenhan Zhang, Michael Lin, Patrick Dong, Charles Young, and Jason Dong. Supertml: Two-dimensional word embedding and transfer learning using imagenet pretrained cnn models for the classifications on tabular data. arXiv preprint arXiv:1903.06246, 2019.
- [13] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: Lessons learned from the 2015 mscoco image captioning challenge. IEEE transactions on pattern analysis and machine intelligence, 39(4):652663, 2016.
- [14] Ting Yao, Yingwei Pan, Yehao Li, Zhaofan Qiu, and Tao Mei. Boosting image captioning with attributes. In Proceedings of the IEEE International Conference on Computer Vision, pages 48944902, 2017.
- [15] Quanzeng You, Hailin Jin, Zhaowen Wang, Chen Fang, and Jiebo Luo. Image captioning with semantic attention. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 46514659, 2016.