

## AI534 HW5: Machine Learning Paper Review

Name: JuHyun Kim

Paper Title: Show and tell: A neural image caption generator

### Non-technical background -Who

Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan are the authors affiliated with Google, a renowned technology company. Among them, Oriol Vinyals stands out as the Principal Investigator (PI), leading the Google Brain team. This team is notable for its groundbreaking work in the field of artificial intelligence and machine learning

### Where

This paper was published in the 2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), which took place in Boston, Massachusetts. The paper has made a significant impact in the field, as indicated by its citation count of 2749 times according to IEEE. Regarding recognitions or awards, there is no specific mention of this paper receiving any such honors directly. The paper's impact, gauged by the number of citations, suggests it continues to influence research. While exact numbers on the trend of citations over time are not provided, the fact that it is still being cited in recent works implies a sustained or possibly increasing influence in related fields. As for media coverage, while this paper was referenced in a few discussions about deep learning or image captioning, there doesn't appear to be significant coverage by prominent media outlets specifically focused on this paper.

### Reproducibility

In this paper, they use several key datasets for the experiments. Each comprises images and English sentences describing them. The datasets are Pascal VOC 2008, Flickr8k, Flickr30k, MSCOCO, and SBU Captioned Photo Dataset. Pascal VOC 2008 is part of the PASCAL Visual Object Classes Challenge series. Flickr8k and Flickr30k contain thousands of images from Flickr, each annotated with sentences describing the image, and are accessible online for research in image captioning. MSCOCO (Microsoft Common Objects in Context) is a comprehensive dataset for object detection, segmentation, and captioning. The SBU Captioned Photo Dataset consists of Flickr images with descriptive captions. Most of these datasets can be available online.

The paper provides a clear train/dev/test split for these datasets, facilitating reproducibility and methodical evaluation.

Dataset name	size		
	train	valid.	test
Pascal VOC 2008 [6]	-	-	1000
Flickr8k [26]	6000	1000	1000
Flickr30k [33]	28000	1000	1000
MSCOCO [20]	82783	40504	40775
SBU [24]	1M	-	-

While I couldn't find the authors' code or demo, I could see the sample codes based on this study. So, I think this paper inspired others in the field. In addition, these resources, along with the

detailed test result scores presented in the paper, suggest that the study is quite reproducible for researchers and practitioners in the field.

## **Core - What**

The authors are addressing the problem of automatically generating natural language descriptions for images. This task involves not only recognizing the objects within an image but also understanding their relationships, attributes, and the activities they are involved in, and then expressing this information coherently in a natural language like English. This problem is not entirely new; it has been a subject of interest in the field of computer vision for some time. However, the approach presented in the paper represents a significant advancement in terms of methodology and effectiveness.

Previous approaches mainly focused on video, involved complex systems that combined visual primitive recognizers with structured formal languages (like And-Or Graphs or logic systems). These systems were then converted to natural language through rule-based methods. While these systems could describe images "in the wild," they were heavily hand-designed, relatively brittle, and demonstrated only in limited domains like traffic scenes or sports. A more recent approach of work has dealt with ranking descriptions for a given image. These approaches are based on co-embedding images and text in the same vector space, and for an image query, descriptions are retrieved that lie close to the image in this space. However, such methods are limited in their ability to describe previously unseen compositions of objects, even if the individual objects have been recognized in training data. Moreover, they do not address the challenge of evaluating the quality of a generated description.

Therefore, the previous approaches often relied on template-based methods or separate models for image analysis and text generation, which limited the fluency and accuracy of the generated captions. In this paper, the authors aim to overcome these limitations by creating more coherent and contextually relevant descriptions.

## **Why**

The authors chose to address the challenge of automatically generating natural language descriptions for images due to the task's inherent complexity and its significant real-world applications. This problem blends the intricacies of computer vision and natural language processing, making it a technically challenging endeavor. It's not just about identifying objects within an image; it's about understanding their interrelations, attributes, and the activities they're involved in, and then effectively translating this visual information into coherent, natural language. Beyond the intellectual challenge, this technology holds substantial practical and societal value, particularly in enhancing accessibility for the visually impaired. By generating accurate descriptions of images, such systems can greatly improve how visually impaired individuals' access and interpret visual content on the web, thereby enhancing their digital experience and independence. Methodologically, this work marked a departure from previous approaches that pieced together various sub-solutions for individual aspects of the problem, like object recognition or text generation. The authors introduced a unified, end-to-end trainable system, representing a significant innovation in methodology. This approach not only streamlined the process but also enhanced the efficiency and effectiveness of generating image descriptions.

## How

The researchers introduced an innovative approach to the challenge of automatic image description generation, leveraging advancements in both computer vision and natural language processing. Their method centers around the Neural Image Caption (NIC) model, an end-to-end trainable system that seamlessly integrates a deep convolutional neural network (CNN) for image analysis with a recurrent neural network (RNN) for generating text descriptions.

The process begins with image encoding using CNN. This network effectively processes the input image, extracting and condensing its features into a compact representation. Acting as an encoder, the CNN transforms the image into a fixed-length vector, capturing the essential visual information. Following this, the RNN comes into play for sentence generation. It takes CNN's output - the distilled image representation - and uses it to craft a corresponding sentence description. RNNs are particularly adept at sequence generation tasks, such as language modeling, because they can retain information across sequences. This characteristic is key to producing coherent and contextually appropriate sentences. And to make the above RNN more concrete two crucial design choices are to be made: what is the exact form of  $f$  and how are the images and words fed as inputs. For  $f$  we use a Long-Short Term Memory (LSTM) net, which has shown state-of-the-art performance on sequence tasks such as translation.

The LSTM is designed to receive image features extracted by CNN as its initial state. This setup forms the foundation for the model to predict a sequence of words, thereby constructing a sentence that describes the image. The process involves the LSTM learning to produce words sequentially, with each word being conditioned on the image features and the sequence of previously generated words. This approach ensures that the captions generated are not only coherent but also relevant to the context of the image.

Training the LSTM is a critical aspect of this research. It involves using pairs of images and their corresponding sentences. The LSTM is trained to predict the next word in a sentence based on the current word and the context derived from the image. This training process is pivotal in helping the LSTM learn how to generate sentences that are both coherent and contextually relevant, based on the image features. The overarching goal of this training is to maximize the likelihood of producing the correct sentence for a given image, thereby enhancing the model's capability to generate accurate and meaningful image captions.

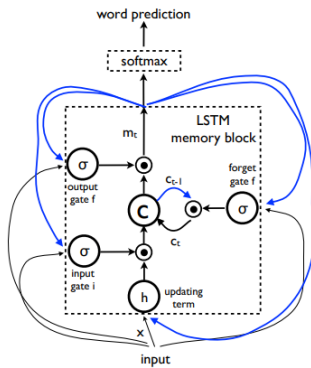


Figure 2. LSTM: the memory block contains a cell  $c$  which is controlled by three gates. In blue we show the recurrent connections – the output  $m$  at time  $t - 1$  is fed back to the memory at time  $t$  via the three gates; the cell value is fed back via the forget gate; the predicted word at time  $t - 1$  is fed back in addition to the memory output  $m$  at time  $t$  into the Softmax for word prediction.

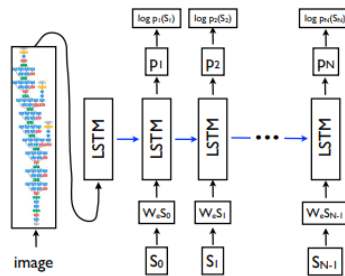


Figure 3. LSTM model combined with a CNN image embedder (as defined in [12]) and word embeddings. The unrolled connections between the LSTM memories are in blue and they correspond to the recurrent connections in Figure 2. All LSTMs share the same parameters.

## Wow

The results illustrate the robust performance of the NIC model in the realm of image captioning. The model was rigorously evaluated across various standard datasets and demonstrated impressive results, as detailed in several tables and figures within the paper. In their quantitative evaluation, the NIC model significantly surpassed both random and nearest neighbor baselines on the MSCOCO dataset. Notably, it achieved a BLEU-4 score of 27.7, a METEOR score of 23.7, and a CIDER score of 85.5. These scores were remarkably close to those obtained by human raters, indicating the model's high level of accuracy and relevance in caption generation. Additionally, on different datasets like the Flickr datasets and the Pascal dataset, the model showed competitive performance with the state-of-the-art results, as evidenced in Table 2.

Metric	BLEU-4	METEOR	CIDER
NIC	<b>27.7</b>	<b>23.7</b>	<b>85.5</b>
Random	4.6	9.0	5.1
Nearest Neighbor	9.9	15.7	36.5
Human	21.7	25.2	85.4

Table 1. Scores on the MSCOCO development set.

Approach	PASCAL (xfer)	Flickr 30k	Flickr 8k	SBU
Im2Text [24]				11
TreeTalk [18]				19
BabyTalk [16]	25			
Tri5Sem [11]			48	
m-RNN [21]		55	58	
MNLM [14] <sup>5</sup>		56	51	
SOTA	25	56	58	19
NIC	<b>59</b>	<b>66</b>	<b>63</b>	<b>28</b>
Human	69	68	70	

Table 2. BLEU-1 scores. We only report previous work results when available. SOTA stands for the current state-of-the-art.

Approach	Image Annotation			Image Search		
	R@1	R@10	Med r	R@1	R@10	Med r
DeFrag [13]	13	44	14	10	43	15
m-RNN [21]	15	49	11	12	42	15
MNLM [14]	18	55	8	13	52	10
NIC	<b>20</b>	<b>61</b>	<b>6</b>	<b>19</b>	<b>64</b>	<b>5</b>

Table 4. Recall@k and median rank on Flickr8k.

Approach	Image Annotation			Image Search		
	R@1	R@10	Med r	R@1	R@10	Med r
DeFrag [13]	16	55	8	10	45	13
m-RNN [21]	18	51	10	13	42	16
MNLM [14]	<b>23</b>	<b>63</b>	<b>5</b>	<b>17</b>	<b>57</b>	<b>8</b>
NIC	17	56	7	<b>17</b>	<b>57</b>	<b>7</b>

Table 5. Recall@k and median rank on Flickr30k.

Furthermore, the NIC model demonstrated its versatility beyond caption generation, performing well in related tasks such as image annotation and image search. This was evidenced by its competitive recall rates compared to other methods, as presented in Tables 4 and 5. The paper also included a qualitative evaluation, where the cumulative distribution of human BLEU scores showed that the NIC model approximated human performance closely, particularly on the Flickr8k dataset.

In discussing their results, the authors acknowledged some limitations and challenges. For instance, they highlighted the issue of overfitting due to the relatively small size of high-quality datasets compared to larger ones like ImageNet. To combat this, they employed strategies like initializing the CNN weights with a pre-trained model and experimenting with dropout. The impact of transfer learning from larger to smaller or differently labeled datasets was also examined. The findings suggested that performance improved with larger datasets and high-quality labels, and the model could generalize across different domains to a certain extent. An important aspect of the model's capability was its diversity in generation. The NIC model was shown to be capable of producing a diverse set of high-quality captions, often matching the diversity seen in human-generated captions. This indicated the model's robustness and adaptability in caption generation.

Overall, the paper's findings demonstrated that the NIC model is a highly effective tool for image captioning. It performed close to human levels on standard metrics and showcased good generalization and diversity in the generated captions. However, the authors also pointed out the

need for improvements, particularly in developing better evaluation metrics and addressing the challenges of overfitting in the context of limited high-quality training data.

### **Further - But**

While the model achieved significant results, it also had limitations. Including potential overfitting due to the relatively small size of high-quality image datasets compared to datasets for other AI tasks like ImageNet. This overfitting issue is partly addressed through techniques like pre-training and dropout, but it remains a challenge. The evaluation metrics, while standard, may not fully capture human-like quality, as suggested by the lower scores from human evaluators. This discrepancy points to a need for better metrics that align more closely with human judgment. Additionally, while the model shows good generalization in some cases, domain mismatches and dataset quality still affect performance, indicating room for improvement in transfer learning capabilities and robustness across varied data.

### **More**

For the follow-up of this research, I think improving the evaluation of the metrics. Current metrics like BLEU often fail to capture the nuances of human language. Future work could involve creating metrics that consider the diversity of expressions in human language, the relevance of the caption to the image, and the fluency or grammatical correctness of the generated text. These metrics could leverage recent advances in language models to provide a more nuanced assessment of the captions' quality, potentially using embeddings to capture semantic similarity more effectively than n-gram matching.

### **All**

The message from this paper is that it presents a significant step in merging deep learning with image processing to create coherent captions, a testament to the powerful synergy between convolutional and recurrent neural networks. Through this paper, I got interested in the potential of neural networks to not just classify images but to also understand and describe them in human-like language.

### **Relevance**

The concepts of neural networks and k-nearest neighbors (KNN) covered in the course were instrumental in understanding the paper. The ability to relate these concepts to practical applications, like sorting out similar sentences, enriched the comprehension of image-based input processing.

To enhance my understanding of the paper, I explored additional materials including detailed information on the datasets used by the authors. And looked over the sample code sample codes based on this study that I mentioned in the Reproducibility part. Furthermore, I reviewed the related work "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention" by Xu et al., which provided valuable insights into the integration of attention mechanisms in neural image captioning.