# Image Captioning With Visual Attention

Aniket Das, Devansh Shringi, Yatin Dandi

*Abstract*—**We study the problem of image captioning, the task of generating natural language descriptions of images. We analyze encoder decoder architectures that use a novel combination of a Convolutional Neural Network to obtain a vectorial representation of the image, and a Recurrent Neural Network to decode the vectorial representation into natural language sentences. We focus particularly on, "Show, Attend and Tell : Neural Image Caption Generation with Visual Attention" [1], a model incorporating visual attention, allowing it, at every timestep to automatically select regions of the image relevant for predicting a target word.**

## I. PROBLEM DESCRIPTION AND MOTIVATION

THE problem of generating descriptions of images is central to the task of contextual understanding. The task requires the model to be capable of determining what features in an image represent objects, understanding the complex relationships between said objects and representing them by a natural language description, thereby compressing large amounts of visual data into descriptive language.

Recent work on image captioning approaches the problem by means of encoder-decoder architectures, which have been widely successful in sequence to sequence training for neural machine translation. [2] performs the machine translation task by using two RNNs, the encoder, which encodes the sequence of words in the sentence into a fixed length vector representation, and the decoder which decodes the vector representation into corresponding words in the target language language. By viewing image captioning as a task analogous to "translating" an image into a sentence, we can see why the encoder decoder paradigm is well suited to this problem. Its effectiveness for image captioning is seen in [3] , which uses a CNN to encode the image into a vector and an LSTM based decoder to generate the descriptions given the encoding vector, and jointly trains them to maximise the likelihood of the target description sentence given the training image. The ability of CNNs to extract rich representations of images suited for a variety of computer vision tasks has been well established in the recent years making it a natural choice for the encoder network. [3] in particular uses features from an intermediate layer of a GoogLeNet network pre-trained on the ImageNet dataset.

## II. ATTENTION IN DEEP NEURAL NETWORKS

Attention Mechanism in Deep Neural Networks is largely inspired by the mechanism of visual attention in the human brain, that, given a task, allows us to select a subset of the available information for further processing. In the context of human vision, this primarily involves focusing on a certain region of the image in high resolution, blurring out the surroundings, then adjusting the focus to different parts over time as they become relevant. As with the encoder decoder framework, the attention mechanism has shown outstanding results in the domain of neural machine translation. [4] performs machine translation by encoding the input sentence into a sequence of vectors, and then using attention, chooses a subset of these vectors to feed into the decoder, learning to "attend" to different parts of the sentence at each step. [1] applies this to images by focusing, at each timestep, on the relevant part of the image when generating a particular word. The encoder and the decoder are jointly trained to maximise the likelihood of the target caption given the training image. The authors present two variants of the attention framework, deterministic soft attention, which is trainable directly by backpropagation, and stochastic hard attention, which is not directly trainable by backpropagation. One can either proceed by maximizing a variational lower bound on the log probability of the caption given the training image by backpropagation, or use REINFORCE. Incorporating visual attention into the model helps it capture log range dependencies and alleviates the vanishing gradient problem in sequence models. Furthermore, attention renders interpretability to the model. By visualising the attention component over the image at every timestep, one can observe what region of the image the network chooses to focus upon when generating a particular word. We observe in our implementation that the alignments learnt by the model strongly correlates to human intuition and in cases where it generates poor quality or incorrect captions, the visualization aids in inferring what elements caused the model to do so

## III. NETWORK ARCHITECTURE AND EXPERIMENTS

The encoder is a CNN which extracts $L$ feature vectors from an intermediate convolutional layer. Apart from the VGG-19 encoder pre-trained on ImageNet that [1] uses, we also experiment with the ResNet-56 and DenseNet as encoders

$$a = \{\mathbf{a_1}, \cdots, \mathbf{a_L}\}, \mathbf{a_i} \in \mathbb{R}^D$$

An LSTM is used for the decoder network to generate a caption $y$ encoded as a sequence of 1-of-K words

$$y = \{\mathbf{y_1}, \cdots, \mathbf{y_C}\}, \mathbf{y_i} \in \mathbb{R}^K$$

$$\begin{pmatrix} \mathbf{i_t} \\ \mathbf{f_t} \\ \mathbf{o_t} \\ \mathbf{g_t} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ tanh \end{pmatrix} T_{D+m+n,n} \begin{pmatrix} \mathbf{Ey_{t-1}} \\ \mathbf{h_{t-1}} \\ \mathbf{\hat{z}_t} \end{pmatrix}$$

$$\mathbf{c_t} = \mathbf{f_t} \odot \mathbf{c_{t-1}} + \mathbf{i_t} \odot \mathbf{g_t}$$

$$\mathbf{h_t} = \mathbf{o_t} \odot \mathbf{tanh}$$

The vector $\hat{z}_t$ is the context vector that captures the visual information associated with a particular input location, representing the relevant part of the image input at time t. The function $\phi$ computes $z_t$ from the annotation vectors $\mathbf{a_1}, \mathbf{a_2} \cdots \mathbf{a_L}$ corresponding to features extracted at different image locations. This is known as the attention mechanism, of which the two variants are soft attention and hard attention. For each location, a positive weight $\alpha_i$ is generated by the attention model $f_{att}$. Depending on whether soft attention or hard attention is used, $\alpha_i$ can either be represented as the probability that location i is the right place to look when generating the caption or the relative importance to be given to location i when mixing the $\mathbf{a_i}$ together. $f_{att}$ is chosen to be an MLP.

$$e_{ti} = f_{att}(\mathbf{a_i}, \mathbf{h_t} - \mathbf{1})$$

$$\alpha_{ti} = \frac{exp(e_{ti})}{\sum_{k=1}^{L} exp(e_{tk})}$$

$$\hat{\mathbf{z}_t} = \phi(\{\mathbf{a_i}\}, \{\alpha_i\},$$

The initial cell state and hidden state vectors are initialized by

$$\mathbf{c_0} = f_{init,c}(\frac{1}{L}\sum_{i}^{L}\mathbf{a_i})$$

$$\mathbf{h_0} = f_{init,h}(\frac{1}{L}\sum_{i}^{L}\mathbf{a_i})$$

Hard: Let the location variable $\mathbf{s_t}$ represent where the model decides to focus when generating the word in timestep t. $\mathbf{s_t}$ is a one hot encoding with the $i_t h$ location set to 1 if it is the one being used to extract visual features. Treating the attention locations as intermediate latent variables we can assign a multinoulli distribution parameterised by $\alpha_i$.

$$p(s_{ti} = 1|s_{j<t}, \mathbf{a}) = \alpha_{t,i}$$

$$\hat{\mathbf{z}_t} = \sum_{i} s_{t,i}\mathbf{a_i}$$

The model is trained by maximising the variational lower bound $L_s$ on the marginal log likelihood log $p(y|a)$ of observing the sequence of words y given image features a

$$L_s = \sum_{s} p(s|\mathbf{a})log(p(\mathbf{y}|s,\mathbf{a}))$$

$$\leq \sum_{s} p(s|\mathbf{a})p(\mathbf{y}|s,\mathbf{a})$$

$$= log(p(\mathbf{y}|\mathbf{a}))$$

In the soft attention variant, instead of a stochastic attention requiring sampling, a weighted mean of the annotation vectors is used. This makes the model deterministic and hence learning is performed using standard back-propagation.

$$\mathbb{E}_{p(s_t|a)}[\hat{z}_t] = \sum_{i=1}^{L} \alpha_{t,i}\mathbf{a_i}$$

The regularisation term encourages $\sum_t \alpha_{ti}$ to be close to 1. This encourages the model to pay equal attention to every part of the image. The final loss is:

$$L_d = -log(P(\mathbf{y}|\mathbf{x})) + \lambda \sum_{i}^{L}(1 - \sum_{t}^{C} \alpha_{ti})^2$$

## IV. EXPERIMENTS

We ran our experiments on the MS COCO Dataset, with a vocabulary size of 25000 and word embeddings of dimension 512. 196 image annotation vectors are used each with dimension 512 and the RNN is unrolled for 16 timesteps. The regularization term $\lambda$ is set to 1.0. The model is trained using the Adam optimizer with a learning rate of 0.001. Here are some of the realistic image captions generated by our model :



Fig. 1. A red bus is near a curb on a brick street near a lane of shops.



Fig. 2. A couple of people that are riding on a motorcycle on a road.

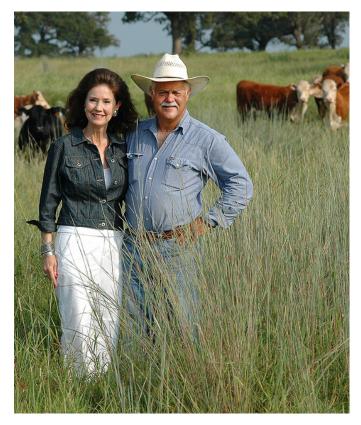Fig. 3.  A group of people in front of a building on a city street.



Fig. 4.  A man and woman stand in tall grass in a field with some cattle.

## REFERENCES

[1] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhutdinov, R. Zemel and Y. Bengio. Show, Attend and Tell: Neural Image Caption Generation With Visual Attention *Proceedings of the 32nd International Conference on Machine Learning*, vol. 37, pp. 2048-2057, Jul. 2015.

[2] K. Cho, B. Merrienboer, C. Gulcehre, D. Bahdanau, F. Bougares, H. Schwenk, Y. Bengio. Learning Phrase Representations using RNN Encoder-Decoder for Statistical Machine Translation *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing*, Istanbul, Turkey, Oct. 2014

[3] O. Vinyals, A. Toshev, S. Bengio, D. Erhan. Show and Tell: A Neural Image Caption Generator *IEEE Conference on Computer Vision and Pattern Recognition*, Oct. 2015

[4] D. Bahdanau, K. Cho, Y. Bengio. Neural Machine Translation by Jointly Learning to Align and Translate *International Conference on Learning Representations*, 2015.