

Human Eye - Task 7

2017202090 Liu Yifan

Problem Definition

Given a picture I , the purpose of this problem is to generate a series of words $X = \{x_1, x_2, x_3, \dots, x_N\}$ which can describe the given picture I reasonably.

Baseline Method

The baseline idea is based on *Show & Tell: A Neural Image Caption Generator*. The method showed in the paper is a supervised learning method by using large dataset which has pairs of images and descriptions to train a generative model. The general idea is to extract features from image using CNN and then put the features and word embeddings of the corresponding description into a RNN.[1]

Dataset

1. MS-COCO 2017

Training Set: 118,287 images (each image with 5 descriptions)

Validation Set: 5000 images

2. Google conceptual captions

Training Set: 3,318,333 images (each image with 1 description)

Data Preprocess

We need to apply preprocess to both image data and text data(the caption of the related image) in both training phase and inference phase.

- For **Image**:
 1. Resize
 2. Apply data augmentation(such as random crop, random flip and so on)
- For **Text**:
 1. Remove infrequent words (frequency < 5) to build vocabulary.
 2. Word embedding to get the representation of each word. (Embedding Size = 300, the default embedding size of *spaCy* language model)
 3. Do sequence padding.

Model

The model has two parts, the first one is a convolution neural network serving as an **Encoder** which is used to extract features from images, and the second part is a recurrent neural network serving as a **Decoder** to generate the reasonable sequence.

Encoder

Use inception_v3 pretrained model with the top layer retrained based on training dataset.

Trainable parameters: weights of the top fully connected layer.

Decoder

Use 1-layer LSTM/GRU and one linear layer as the decoder

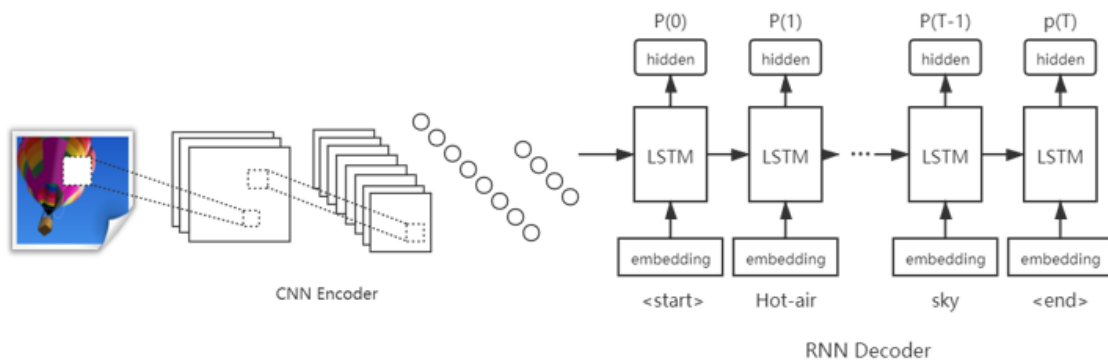
Input sequence:

x_0 is the extracted features from the encoder above.

$X = \{x_1, x_2, x_3 \dots x_N\}$ is the embedding sequence of the given image.

Trainable parameters: LSTM/GRU parameters (the embeddings are derived from spaCy language model and here are set to be not trainable)

The following figure describes the model architecture overview:



Training Phase

The goal of training phase is to maximize the joint probability which describes the occurrence of the given text sequence.

Loss function: use cross entropy loss

Optimizer: Adam optimizer

Inference Phase

The goal of inference phase is to generate sequence of words to describe the given picture.

Input & Output:

x_0 is the extracted features from encoder.

x_1 is the output of decoder when input is x_0 .

x_2 is the output of decoder when input is x_1 .

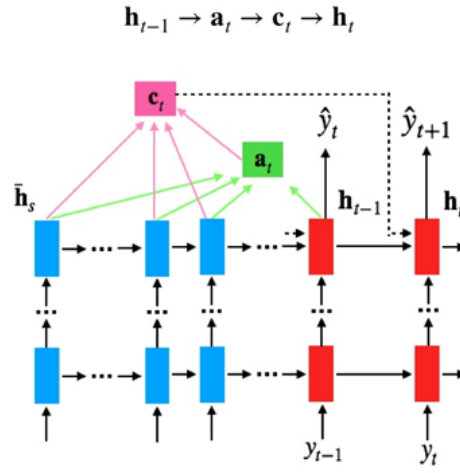
...

Until the inferred sequence meet **<END>** symbol or the length of the inferred sequence meets `max_length`.

As for the criterion for deciding which word to choose for next inference step, now this method uses **Greedy Search** (just choose the word having maximum occurrence probability). The future method will apply **Beam Search** as well.

Attention

Bahdanau Attention



The picture above shows **Bahdanau Attention** applied in Machine Translation scene. The blue blocks refer to sentence in source language and the red blocks refer to target language.

In image captioning scene, we divide each image into several regions, and extract features of each region. These features can be seen as h_s . The y_t in the picture above represents the input word vector of target sentence.

This process can be described as (soft attention)[2]:

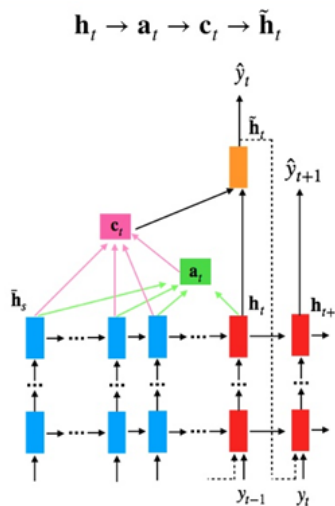
$$e_{ti} = f_{\text{att}}(\mathbf{a}_i, \mathbf{h}_{t-1})$$

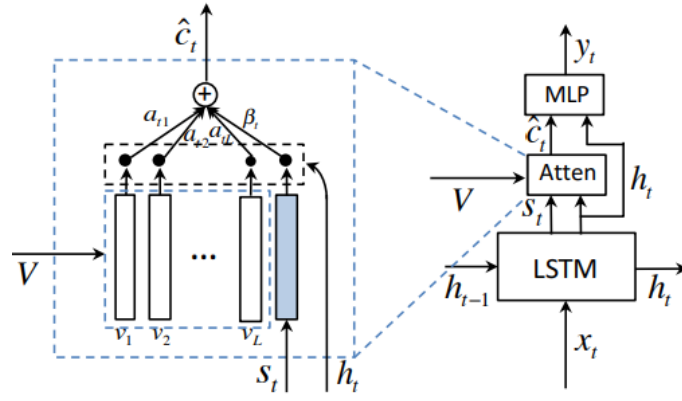
$$\alpha_{ti} = \frac{\exp(e_{ti})}{\sum_{k=1}^L \exp(e_{tk})}$$

$$\hat{\mathbf{z}}_t = \phi(\{\mathbf{a}_i\}, \{\alpha_i\})$$

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

Luong Attention + Visual Sentinel

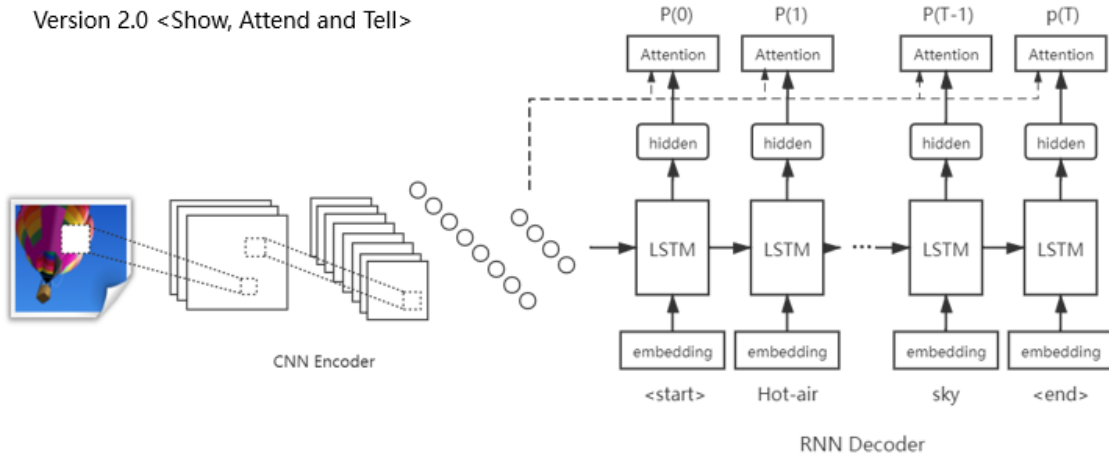




Bahdanau Attention use h_{t-1} to apply attention to calculate context vector, but actually we care more about the which region in the image is the most important part related to the **current** word but not the last word, so the Luong Attention was proposed to improve the attention mechanism.

This model architecture can be described as follows:

Version 2.0 <Show, Attend and Tell>



Results of the method using Bahdanau Attention shows that when the input word is not related to visual conceptions or meaningless such as 'is', 'the', 'a', the attention does not work well, thus results in repeating of meaningless words.

Based on this problem, some work[3] proposed **Visual Sentinel**, a module serving as a trade-of-tool to decide the weight of attended context vector and LSTM/GRU state vector. So if the state in memory cell is more important, the generation of the current word will depend on the sequence before more on the given image, which is more reasonable for the words not related to certain visual conceptions.

The process above can be described as:

$$z_t = w_h^T \tanh(W_v V + (W_g h_t) \mathbb{1}^T)$$

$$\alpha_t = \text{softmax}(z_t)$$

$$c_t = \sum_{i=1}^k \alpha_{ti} v_{ti}$$

$$g_t = \sigma(W_x x_t + W_h h_{t-1})$$

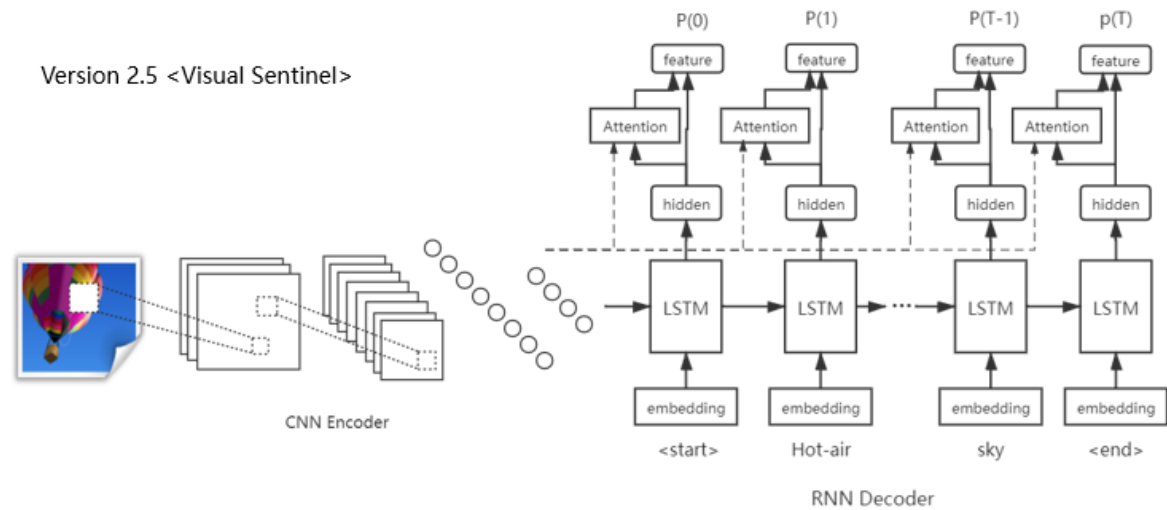
$$s_t = g_t \odot \tanh(m_t)$$

$$\hat{\alpha}_t = \text{softmax}([z_t; w_h^T \tanh(W_s s_t + (W_g h_t))])$$

$$\beta_t = \alpha_t[k + 1].$$

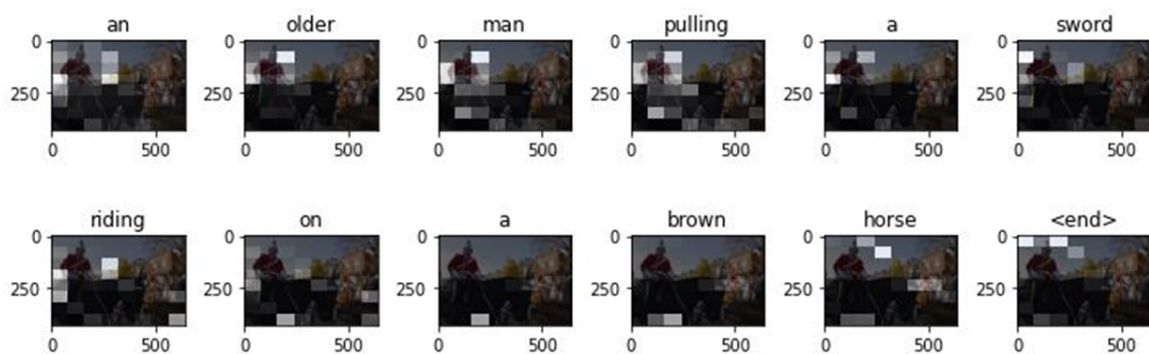
where V is the features of regions in the image, h_t is the hidden output of LSTM/GRU, \hat{c}_t is the adapted context vector.

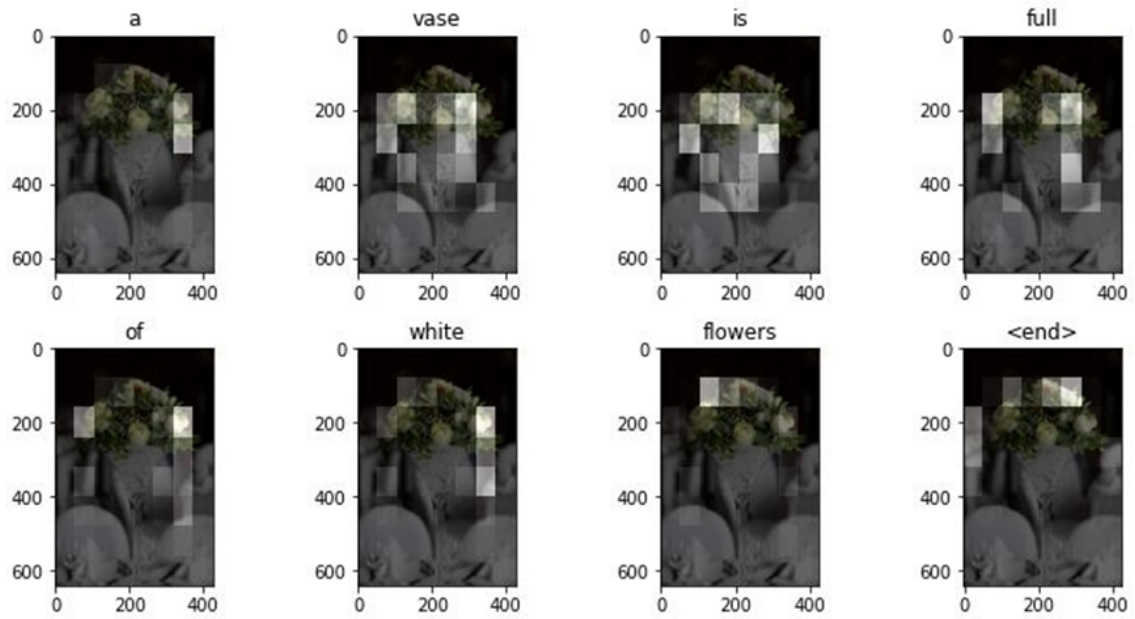
The above process is described as follows:



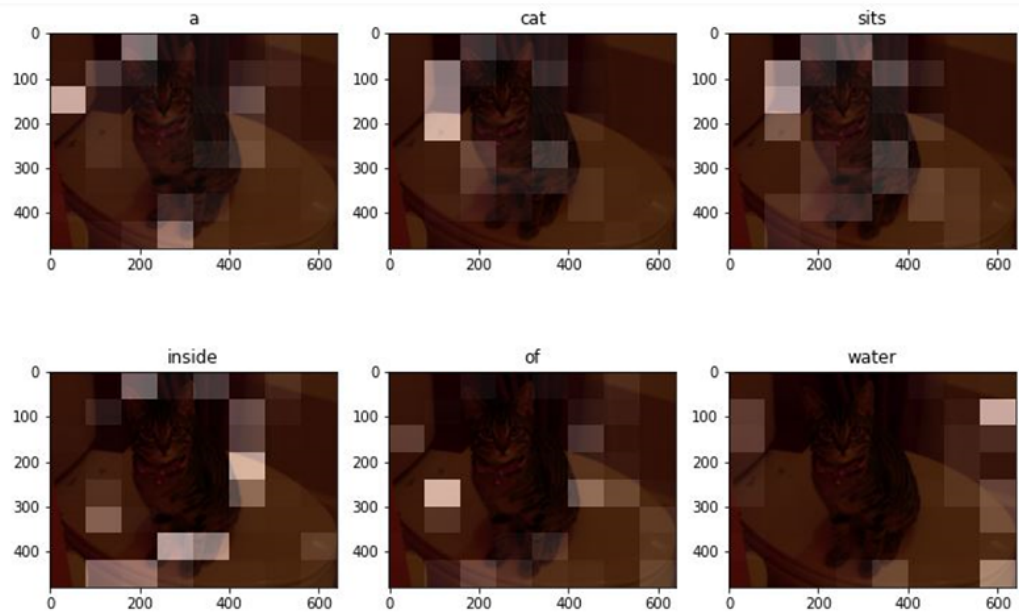
Results

- good example





- fair example



Integrating Knowledge and Other model

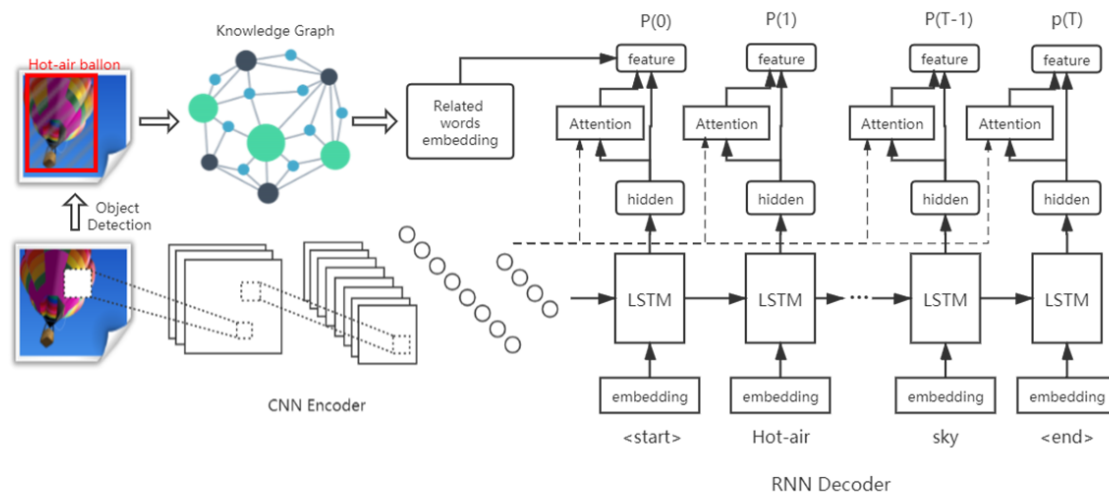
As the results above describes, there are a lot of space to improve the expression of the "human eye".

Firstly, the scene information are virtually not preserved in generated captions when the scene is difficult to detect or there are so many elements in one image.

Secondly, some object in the image is mis-classified.

Based on the statements above, integrating object detection and external knowledge should be a good tool to address these problems. [5]

The models architecture overview is as follow:



Integrating Object Detection

To detect the objects in one image, I use tensorflow hub pretrained object detection model which is built based on Faster R-CNN with Resnet152 as its feature extractor.

For each image, first apply the object detection model on the image, we will obtain a list of object with different confidence score. Then discard object when its confidence score less than the configured threshold (here I used 0.7 as its threshold).

At each time, calculate the cosine similarity between the input word and the filtered objects to find the most related object.

Integrating Knowledge Graph

Use the most related object obtained above as the input of knowledge graph. Then search for its related words (there is a edge between them) with relativeness (the weights of the specific edge).

Especially, we only extract the words which have edge type labeled as `AtLocation` to integrate scene information from external knowledge graph.

After getting the related scene words, use word embedding to get the vector representation of these words. Then add it to the `logit` layer of the model.

Results

A boy in front of the fence in the farm is watering pumpkins.



boy



fence



pumpkins



farm

A kid is playing toy with a woman.



kid



toy



woman

Future Word Ideas

1. Change CNN

The above change are generally applied on the RNN decoder part, but if the dataset used to train can not give more diverse descriptions of images, it will be helpless to just improve RNN part. Getting more specific information from images is probably a good way.

- Use Object Detection Network to get local region separated feature extraction
- Use Scene Classification Network to get global feature extraction
- Concatenate the local features(#proposed regions) and global scene features(#1)

2. Named Entity Recognition

We can see that the second result image showed in last section contains "LEGO" toy but our method failed to describe it. So another idea[4] is to apply NER on images to get celebrities, landmarks, and famous brands. This will make the generated caption more expressive.

3. Use New Dataset

The way of integrating knowledge graph above can not exploit the rich information of knowledge graph. Another better way to integrate external knowledge is possibly just use a sub-knowledge graph to task place of giving a text descriptions.

Stanford Vision Lab has published one manually labeled dataset which has <image, graph(<entity, relation>)> pairs. We can use graph embedding and pretrained language models to generated more reasonable captions. Because the interactions among the objects in one image can be really complex, the ability of sub-knowledge graph on a specific image that it can clearly describe all the relationships and interactions can be a strong tool to make computer more understand the meaning of the content.

Reference

[1] Vinyals, Oriol, et al. "Show and Tell: Lessons Learned from the 2015 MSCOCO Image Captioning Challenge." IEEE Transactions on Pattern Analysis & Machine Intelligence 39.4(2015):652-663.

[2] Xu, Kelvin, et al. "Show, Attend and Tell: Neural Image Caption Generation with Visual Attention." Computer Science (2016):2048-2057.

[3] Lu J, Xiong C, Parikh D, et al. Knowing When to Look: Adaptive Attention via a Visual Sentinel for Image Captioning[C]// IEEE Conference on Computer Vision and Pattern Recognition. IEEE, 2017:3242-3250.

[4] Tran, Kenneth & He, Xiaodong & Zhang, Lei & Sun, Jian. (2016). Rich Image Captioning in the Wild. 434-441. 10.1109/CVPRW.2016.61.

[5] Feicheng Huang, Zhixin Li, Haiyang Wei, Canlong Zhang, Huifang Ma: Boost image captioning with knowledge reasoning. Mach. Learn. 109(12): 2313-2332 (2020)