

# Visual Question Answering with Question-Image Hierarchical Co-Attention

Yibin Hu

March 22, 2024

Tulane University

## Abstract

This report analyses and reproduces the result in the paper Hierarchical Question-Image Co-Attention for Visual Question Answering[1] and compare the co-attention model with a bag of words-image baseline model, and co-attention model is then fine-tuned to do a harder classification task and be compared with pure image model. The VQA problem requires the model to reason about the question based on given image and choose the answer from a list of possible choices. The result is currently missing and to be reported.

## 1 Introduction

Visual question answering[2] is problem that incorporate computer vision and natural language processing. A model for such problem need to understand the relation between image and text question, and this makes attention necessary to this type of question, as attention will identify which and how any part of image or text affect the result and recognize variants of the equivalent question formulated in different way to have the same meaning. The co-attention model has two main components, co-attention which correlate and attends the image and question feature, and question hierarchy which separate the text into three levels as word, phrase, and question level, with each upper level being convoluted from lower level, this tries to capture information contained at different scales.

## 2 Method

VQA problem is formulated as a classification problem, so each distinct answer appeared in the training set will correspond to a class.

In this section, the baseline bag of words-image model[3] and co-attention model is described.

### 2.1 Bag of Words + Images

The bag of words-image model concatenate image features extracted from pretrained visual model GoogleNet and one-hot encoded word feature, and feed this feature vector into a softmax layer to get prediction.

## 2.2 Hierarchical Co-attention

Given a question of  $T$  words, its one-hot encoding is  $Q = \{q_1, \dots, q_T\}$ , and from this the embedding of each word forms word-level feature  $Q^w = \{q_1^w, \dots, q_T^w\}$ . From word-level feature, the phrase-level feature  $Q^p$  is formed by computing the inner product of the word vectors with filters of three window sizes: unigram, bigram and trigram as

$$\hat{q}_{s,t}^p = \tanh(W_c^s q_{t:t+s-1}^w), \quad s \in \{1, 2, 3\}$$

and then max-pooled as

$$q_t^p = \max(\hat{q}_{1,t}^p, \hat{q}_{2,t}^p, \hat{q}_{3,t}^p), \quad t \in \{1, 2, \dots, T\}$$

From the phrase-level feature  $Q^p$ , the question-level feature  $q_t^s \in Q^s$  is encoded as LSTM hidden state at  $t$ . The hierarchical representation is illustrated in figure 1.a.

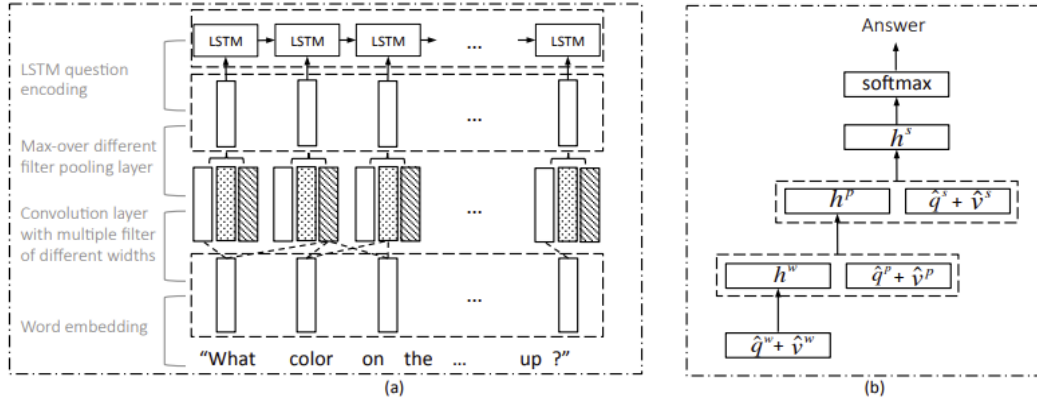


Figure 1: a) Hierarchical question encoding b) Encoding for predicting answers

For each level of word feature  $Q^r \in R^{d_r \times T}$  ( $r \in w, p, s$ ) and image feature  $V \in R^{d_2 \times T}$ , affinity matrix  $C \in R^{N \times T}$  is computed as

$$C = \tanh(Q^T W_c V)$$

then the attention vector for image and feature is trained as

$$\begin{aligned} H^v &= \tanh(W_v V + (W_q Q) C), H^q = \tanh(W_q Q + (W_v V) C^T) \\ a^v &= \text{softmax}(w_{hv}^T H^v), a^q = \text{softmax}(w_{hq}^T H^q) \end{aligned}$$

which are used to form weight-summed image and question feature as

$$\hat{v} = \sum_{n=1}^N a_n^v v_n, \quad \hat{q} = \sum_{t=1}^T a_t^q q_t$$

The parallel co-attention is done at each level in the hierarchy, leading to  $\hat{v}^r$  and  $\hat{q}^r$  where  $r \in \{w, p, s\}$ . Illustrated in figure 1.b, the final prediction is recursively encoded by a MLP from

hierarchical features as

$$\begin{aligned} \mathbf{h}^w &= \tanh(\mathbf{W}_w(\hat{\mathbf{q}}^w + \hat{\mathbf{v}}^w)) \\ \mathbf{h}^p &= \tanh(\mathbf{W}_p[(\hat{\mathbf{q}}^p + \hat{\mathbf{v}}^p), \mathbf{h}^w]) \\ \mathbf{h}^s &= \tanh(\mathbf{W}_s[(\hat{\mathbf{q}}^s + \hat{\mathbf{v}}^s), \mathbf{h}^p]) \\ \mathbf{p} &= \text{softmax}(\mathbf{W}_h \mathbf{h}^s) \end{aligned}$$

with  $\mathbf{p}$  being used to classify the most likely answer.

### 3 Dataset

The model will be trained and evaluated on VQA dataset.

### 4 Experiments

### 5 Discussion

### 6 Conclusion

## References

- [1] J. Lu, J. Yang, D. Batra, and D. Parikh, "Hierarchical question-image co-attention for visual question answering," in *Advances In Neural Information Processing Systems*, pp. 289–297, 2016.
- [2] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. Vqa: Visual question answering. In *ICCV*, 2015.
- [3] Zhou, Bolei, et al. "Simple baseline for visual question answering." *arXiv preprint arXiv:1512.02167* (2015).