



Towards Good Practice for Visual Question Answering

Zhe Wang¹, Xiaoyi Liu², Liangjian Chen¹, Limin Wang⁴, Yu Qiao³, Xiaohui Xie¹, Charless Fowlkes¹

¹ICS UC Irvine, ²EECS UC Irvine, ³SIAT CAS, ⁴EE ETH

Visual Question Answering

Question: Why was the hand of the woman over the left shoulder of the man?

- A:** They were together and engaging in affection
A: The woman was trying to get the man's attention
A: The woman was trying to scare the man
A: The woman was holding on to the man for balance



Visual question answering (VQA) tasks are of significant interest due to their potential as a strong test of image understanding systems and in probing the connection between language and vision. Despite much recent innovation, general VQA is far from a solved problem.

Our Approach

We explore three mechanisms for improving VQA performance

(i) POS Tag Guided Attention:

- Some words (i.e., nouns, verbs and adjectives) should matter more than others (e.g., the conjunctions)
- We use a small set of seven POS categories (numbers, nouns, adjectives, verbs, wh-pronouns, wh-adverbs, other)

(ii) Convolutional N-Gram:

We propose using a convolutional n-gram to combine contextual information over multiple words represented as vectors. Contextual features for different window sizes are pooled to obtain a new word representation: $\tilde{e}_i = \text{maxpool}(F_L, F_{L-1}, \dots, F_1)$.

The final question / answer sentence is represented by an average of word representations $\mathbf{x}_Q = \frac{1}{M} \sum_{i=1}^M \tilde{e}_i$.

(iii) Triplet Attention:

We derive an spatial attention weight from the question and answer representations.

$\text{att}_I = \text{norm}(\lambda \times \text{att}_{Q-I} + \text{att}_{A-I})$ where $\text{norm}(x) = \frac{x}{\sum(x)}$, att_{Q-I} att_{A-I} are attention weights from questions/answers to images.

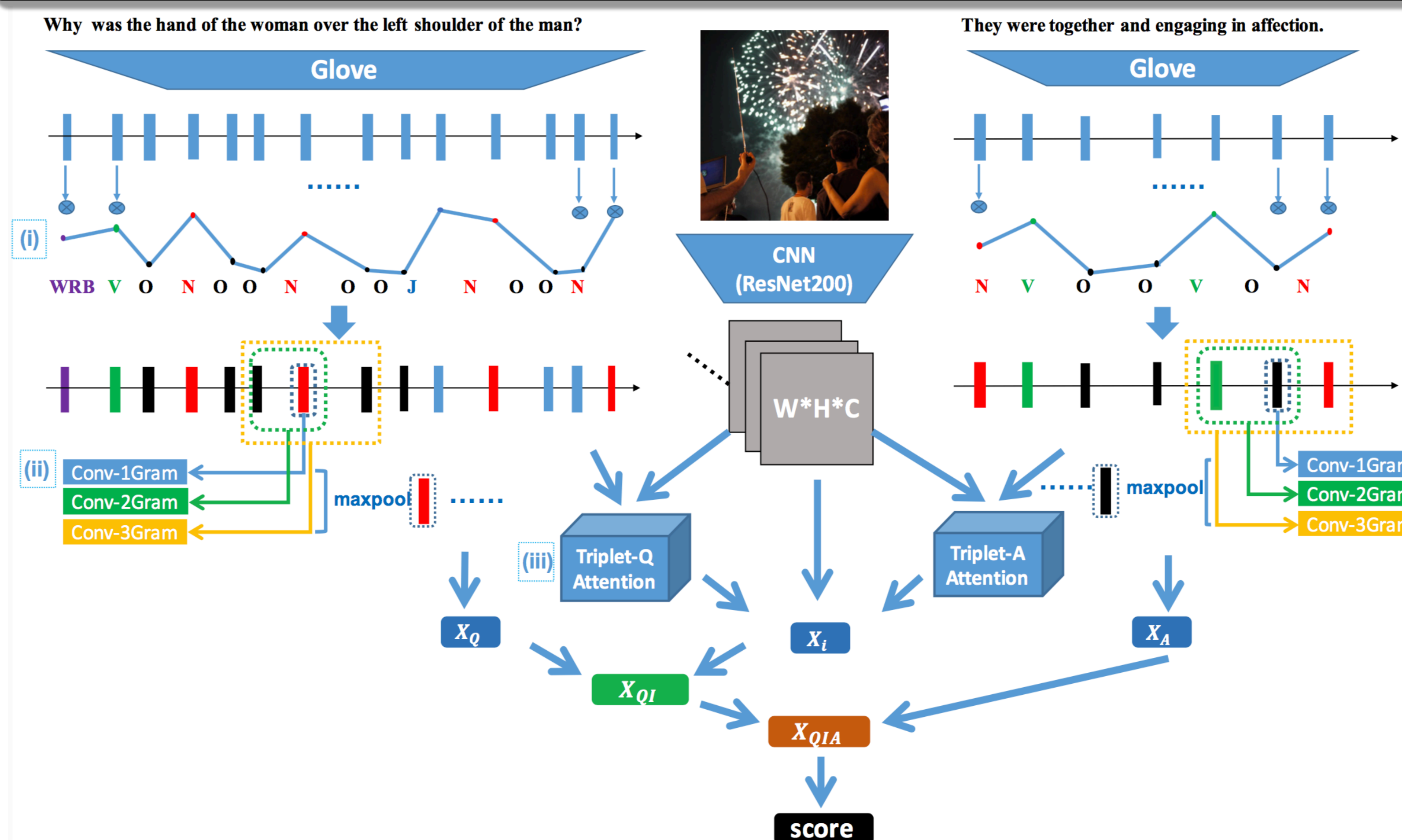
Given image features $\mathbf{X}_I = \text{relu}(\mathbf{W}_I \mathbf{X}_{I,\text{raw}} + \mathbf{b}_I)$
 and question features $\mathbf{X}_Q = [\tilde{e}_1, \dots, \tilde{e}_M]$

We compute an affinity matrix $\mathbf{A} = \text{softmax}(\mathbf{X}_Q^T \times \mathbf{X}_I)$
 and a Question-Image attention vector $\text{att}_{Q-I} = \text{maxpool}(\mathbf{A})$

References

- Akira Fukui, Dong Huk Park, Daylen Yang, Anna Rohrbach, Trevor Darrell, and Marcus Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP*, 2016.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- Allan Jabri, Armand Joulin, and Laurens van der Maaten. Revisiting visual question answering baselines. In *ECCV*, 2016.
- Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. Hierarchical question-image co-attention for visual question answering. In *NIPS*, 2016.
- Damien Teney and Anton van den Hengel. Zero-shot visual question answering. 2016. arXiv:1611.05546.
- Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. Visual7W: Grounded Question Answering in Images. In *CVPR*, 2016.
- Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. Visual genome: Connecting language and vision using crowdsourced dense image annotations. 2016. arXiv:1602.07332.

Our Pipeline for VQA



General framework:

- Extract vector representations of the image, question and candidate answer using deep neural network.
- Score the compatibility using a two layer network:

- \mathbf{x}_Q Question sentence descriptor
- \mathbf{x}_{A_i} Answer descriptor for ith answer
- \mathbf{x}_I Images descriptor
- \odot Hadamard product(element-wise multiplication)

$$\mathbf{x}_{QI} = \mathbf{x}_Q \odot \mathbf{x}_I.$$

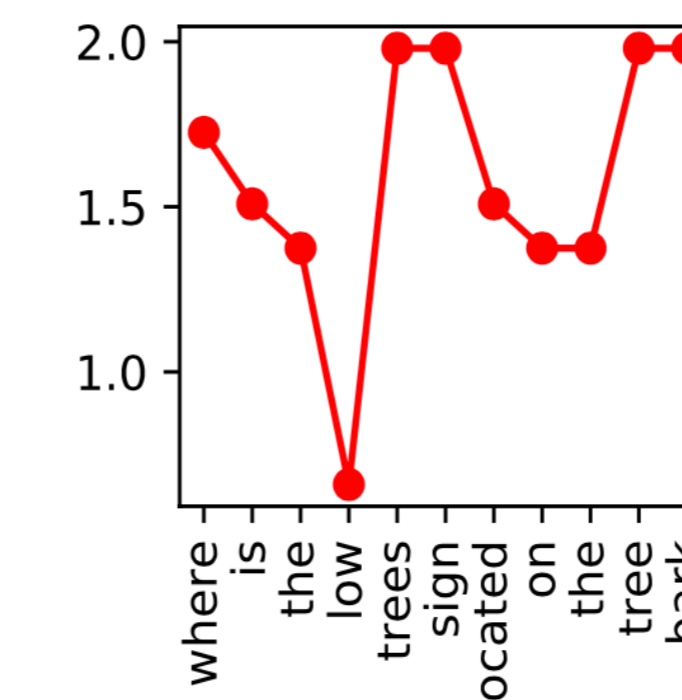
$$\mathbf{x}_{QIA_i} = \tanh(\mathbf{W}_{QI} \mathbf{x}_{QI} + \mathbf{b}_{QI}) \odot \mathbf{x}_{A_i}.$$

$$p_i = \text{sigmoid}(\mathbf{W}_{QIA} \mathbf{x}_{QIA_i} + \mathbf{b}_{QIA}).$$

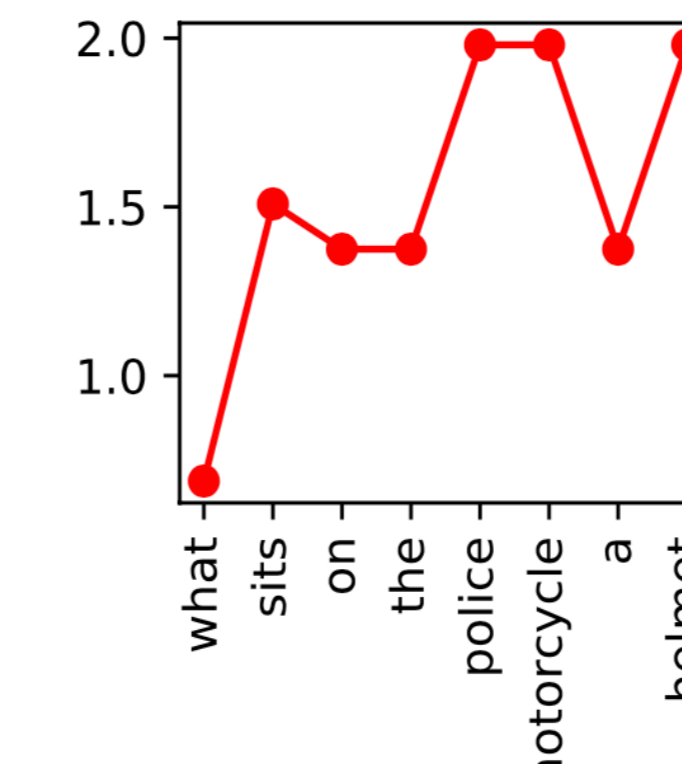
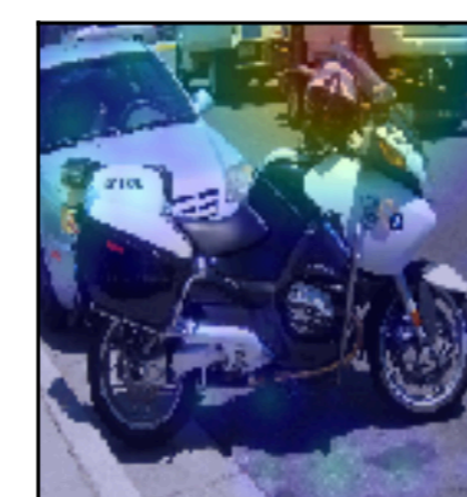
$$i^* = \arg \max_{i=1, \dots, N} p_i.$$

Visualizing Attention

Question: Where is the "low trees" sign located?
A: On the tree bank



Question: What sits on the police motor cycle?
A: A Helmet



Input Image

Triplet attention

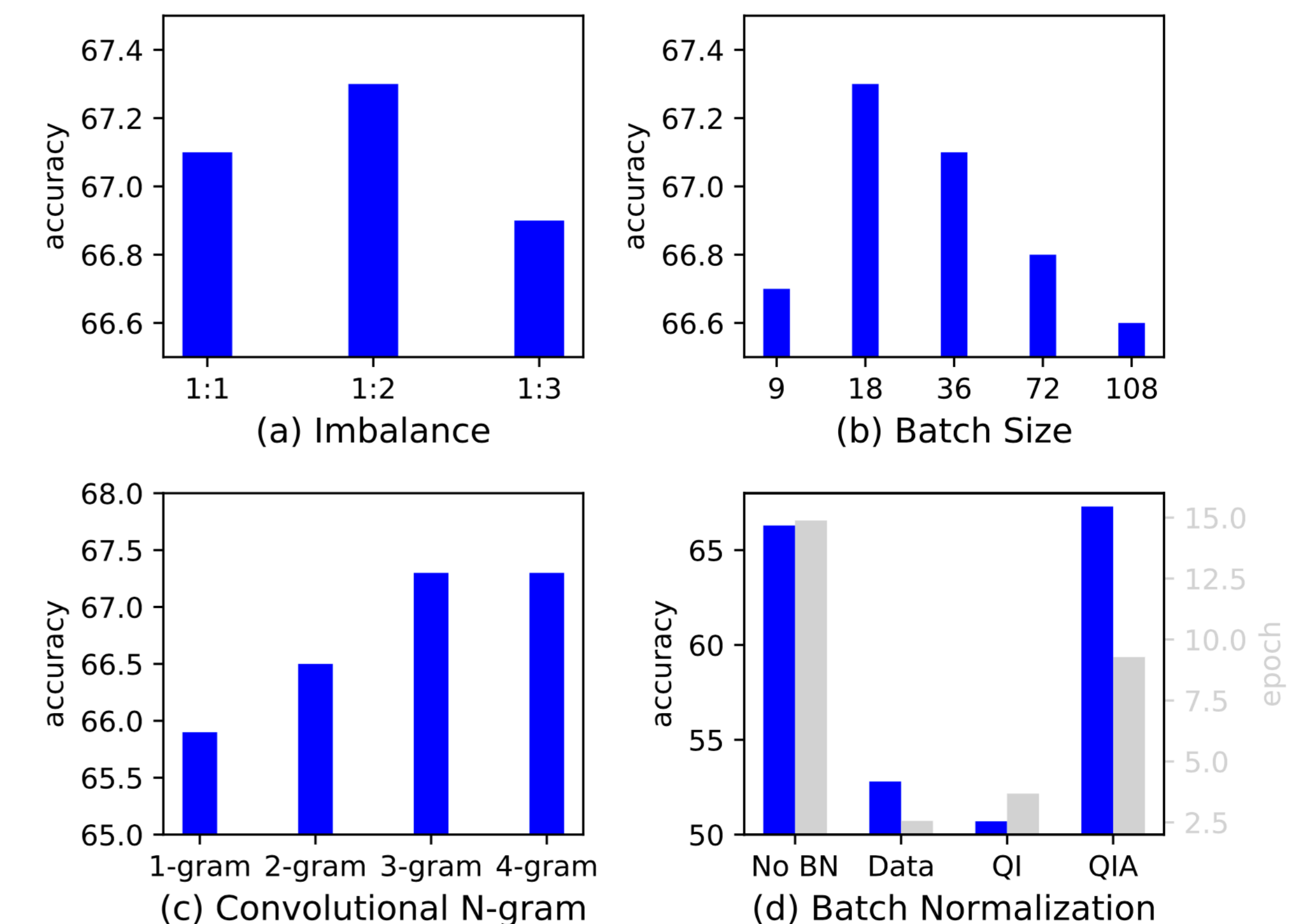
POS-tag guided attention

Triplet attention focuses on image region (e.g. sign on the tree bark) and relevant key words in the question and answer are assigned high weights

Results & Analysis

Method	Visual 7W
Our Baseline	65.6
+POS tag guided attention(POS-Att)	66.3
+Convolutional N-Gram (Conv N-Gram)	66.2
+POS-Att +Conv N-Gram	66.6
+POS-Att +Conv N-Gram +Triplet attention-Q	66.8
+POS-Att +Conv N-Gram +Triplet attention-A	67.0
+POS-Att +Conv N-Gram +Triplet attention-Q+A (Full model)	67.3

Verification of our proposed (1) POS tag guided attention, (2) Conv N-Gram and (3) Triplet Attention step by step. Integrating them all further improves the performance.



Exploration of good practice (1) Handling data imbalance (2) adjusting batch size (3) parameter to adjust convolutional n-gram and (4) where to add batch normalization [2].

Method	Visual 7W Telling	VQA Real Multi Choice
Co-Attention [4]	-	66.1
Attention-LSTM [6]	55.6	-
MCB [1]	-	65.4
MCB + Att [1]	62.2	-
Ensemble of 7 Att models [1]	-	70.1
Zero-shot [5]	65.7	-
MLP [3]	64.8	65.2
Full model (7*7 Resnet feature)	67.3	68.3
Full model (14*14 Resnet feature)	68.1	-

In comparison with the related work:

- We outperform the state-of-the-art performance on visual7w and get competitive performance on VQA.
- Recent state-of-the-art work in [1] used an ensemble of 7 models and trained with additional data (the Visual Genome dataset [7]), performing slightly better than our model on VQA suggesting our simpler model is still quite competitive in terms of computation trade-off.