

# Leveraging Visual Question Answering for Image-Caption Ranking

Xiao Lin<sup>(✉)</sup> and Devi Parikh

Bradley Department of Electrical and Computer Engineering,  
Virginia Tech, Blacksburg, USA  
{linxiao,parikh}@vt.edu

**Abstract.** Visual Question Answering (VQA) is the task of taking as input an image and a free-form natural language question about the image, and producing an accurate answer. In this work we view VQA as a “feature extraction” module to extract image and caption representations. We employ these representations for the task of image-caption ranking. Each feature dimension captures (imagines) whether a fact (question-answer pair) could plausibly be true for the image and caption. This allows the model to interpret images and captions from a wide variety of perspectives. We propose score-level and representation-level fusion models to incorporate VQA knowledge in an existing state-of-the-art VQA-agnostic image-caption ranking model. We find that incorporating and reasoning about consistency between images and captions significantly improves performance. Concretely, our model improves state-of-the-art on caption retrieval by 7.1 % and on image retrieval by 4.4 % on the MSCOCO dataset.

**Keywords:** Visual question answering · Image-caption ranking · Mid-level concepts

## 1 Introduction

Visual Question Answering (VQA) is an “AI-complete” problem that requires knowledge from multiple disciplines such as computer vision, natural language processing and knowledge base reasoning. A VQA system takes as input an image and a free-form open-ended question about the image and outputs the natural language answer to the question. A VQA system needs to not only recognize objects and scenes but also reason beyond low-level recognition about aspects such as intention, future, physics, material and commonsense knowledge. For example (*Q*: Who is the person in charge in this picture? *A*: Chef) reveals the most important person and occupation in the image. Moreover, answers to multiple questions about the same image can be correlated and may reveal

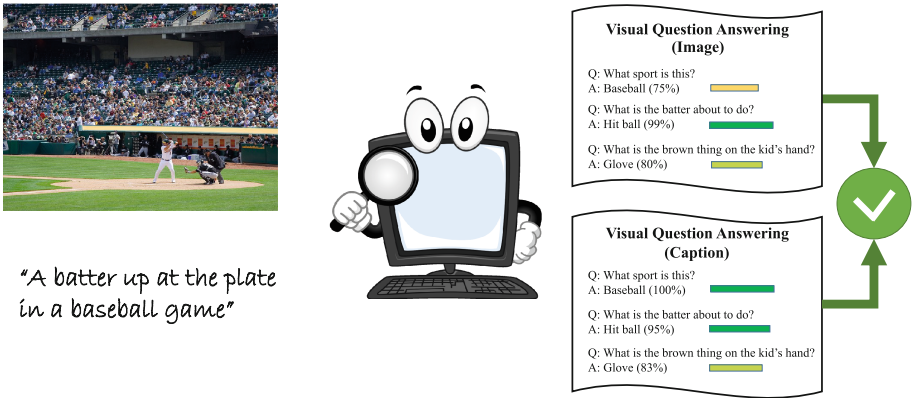
---

**Electronic supplementary material** The online version of this chapter (doi:[10.1007/978-3-319-46475-6\\_17](https://doi.org/10.1007/978-3-319-46475-6_17)) contains supplementary material, which is available to authorized users.

more complex interactions. For example (*Q*: What is this person riding? *A*: Motorcycle) and (*Q*: What is the man wearing on his head? *A*: Helmet) might reveal correlations observable in the visual world due to safety regulations.

Today’s VQA models, while far from perfect, may already be picking up on these semantic correlations of the world. If so, they may serve as an implicit knowledge resource to help other tasks. Just like we do not need to fully understand the theory behind an equation to use it, can we already use VQA knowledge captured by existing VQA models to improve other tasks?

In this work we study the problem of using VQA knowledge to improve image-caption ranking. Consider the image and its caption in Fig. 1. Aligning them not only requires recognizing the batter and that it is a baseball game (mentioned in the caption), but also realizing that a batter up at the plate would imply that a player is holding a bat, posing to hit the baseball and there might be another player nearby waiting to catch the ball (seen in the image). Image captions tend to be generic. As a result, image captioning corpora may not capture sufficient details for models to infer this knowledge.



**Fig. 1.** Aligning images and captions requires high-level reasoning *e.g.* “a batter up at the plate” would imply that a player is holding a bat, posing to hit the baseball and there might be another player nearby waiting to catch the ball. There is rich knowledge in Visual Question Answering (VQA) corpora containing human-provided answers to a variety of questions one could ask about images. We propose to leverage knowledge in VQA by using VQA models learned on images and captions as “feature extraction” modules for image-caption ranking.

Fortunately VQA models try to explicitly learn such knowledge from a corpus of images, each with associated questions and answers. Questions about images tend to be much more specific and detailed than captions. The VQA dataset of [1] in particular has a collection of free-form open-ended questions and answers provided by humans. These images also have associated captions [31].

We propose to leverage VQA knowledge captured by such corpora for image-caption ranking by using VQA models learned on images and captions as “feature

extraction” schemes to represent images and captions. Given an image and a caption, we choose a set of free-form open-ended questions and use VQA models learned on images and captions to assess probabilities of their answers. We use these probabilities as image and caption features respectively. In other words, we embed images and captions into the space of VQA questions and answers using VQA models. Such VQA-grounded representations interpret images and captions from a variety of different perspectives and imagine beyond low-level recognition to better understand images and captions.

We propose two approaches that incorporate these VQA-grounded representations into an existing state-of-the-art<sup>1</sup> VQA-agnostic image-caption ranking model [23]: fusing their predictions and fusing their representations. We show that such VQA-aware models significantly outperform the VQA-agnostic model and set state-of-the-art performance on MSCOCO image-caption ranking. Specifically, we improve caption retrieval by 7.1% and image retrieval by 4.4%.

This paper is organized as follows: Section 2 introduces related works. We first introduce VQA and image-caption ranking tasks as our building blocks in Sect. 3, then detail our VQA-based image-caption ranking models in Sect. 4. Experiments and results are reported in Sect. 5. We conclude in Sect. 6.

## 2 Related Work

**Visual Question Answering.** Visual Question Answering (VQA) [1] is the task of taking an image and a free-form open-ended question about the image and automatically predicting the natural language answer to the question. VQA may require fine-grained recognition, object detection, activity recognition, multi-modal and commonsense knowledge. Large datasets [1, 17, 35, 42, 57] have been made available to cover the diversity of knowledge required for VQA. Most notably the VQA dataset [1] contains 614,163 questions and ground truth answers on 204,721 images of the MSCOCO [31] dataset.

Recent VQA models [1, 17, 33, 36, 42, 61] explore state-of-the-art deep learning techniques combining Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs). [1] also explores a slight variant of VQA that answers a question about the image by reading a caption describing the image instead of looking at the image itself. We call this variant VQA-Caption.

VQA is a challenging task in its early stages. In this work we propose to use both VQA and VQA-Caption models as implicit knowledge resources. We show that current VQA models, while far from perfect, can already be used to improve other multi-modal AI tasks; specifically image-caption ranking.

**Semantic Mid-Level Visual Representations.** Previous works have explored the use of attributes [5, 15, 55], parts [3, 58], poselets [4, 59], objects [30], actions [43] and contextual information [9, 18, 50] as semantic mid-level representations for visual recognition. Benefits of using such semantic mid-level visual

---

<sup>1</sup> To the best of our knowledge on MSCOCO [31], [23] has the state-of-the-art caption retrieval performance. [33] has the state-of-the-art image retrieval performance.

representations include improving fine-grained visual recognition, learning models of visual concepts without example images (zero-shot learning [29, 38]) and improving human-machine communication where a user can explain the target concept during image search [25, 28], or give a classifier an explanation of labels [10, 39]. Recent works also explore using word embeddings [46] and free-form text [12] as representations for zero-shot learning of new object categories. [21] proposes scene graphs for image retrieval. [2] proposes using abstract scenes as an intermediate representation for zero-shot action recognition. Closest to our work is the use of objects, actions, scenes [14], attributes and object interactions [27] for generating and ranking image captions. In this work we propose to use free-form open-ended questions and answers as mid-level representations and we show that they provide rich interpretations of images and captions.

**Commonsense Knowledge for Visual Reasoning.** Recently there has been a surge of interest in visual reasoning tasks that require high-level reasoning such as physical reasoning [19, 60], future prediction [16, 40, 54], object affordance prediction [62] and textual tasks that require visual knowledge [32, 44, 51]. Such tasks can often benefit from reasoning with external commonsense knowledge resources. [63] uses a knowledge base learned on object categories, attributes, actions and object affordances for query-based image retrieval. [53] learns to anticipate future scenes from watching videos for action and object forecasting. [32] learns to imagine abstract scenes from text for textual tasks that need visual understanding. [44, 51] evaluate the plausibility of commonsense assertions by verifying them on collections of abstract scenes and real images, respectively, to leverage the visual common sense in those collections. Our work explores the use of VQA corpora which have both visual (image) and textual (captions) commonsense knowledge for image-caption ranking.

**Images and Captions.** Recent works [6, 22, 23, 34, 37, 56] have made significant progress on automatic image caption generation and ranking by applying deep learning techniques for image recognition [26, 45, 49] and language modeling [7, 48] on large datasets [8, 31]. Algorithms can now often generate accurate, human-like natural-language captions for images. However, evaluating the quality of such automatically generated open-ended image captions is still an open research problem [13, 52].

On the other hand, ranking images given captions and ranking captions given images require a similar level of image and language understanding, but are amenable to automatic evaluation metrics. Recent works on image-caption ranking mainly focus on improving model architectures. [23, 37] study different architectures for projecting CNN image representations and RNN caption representations into a common multi-modal space. [34] uses multi-modal CNNs for image-caption ranking. [22] aligns image and caption fragments using CNNs and RNNs. Our work takes an orthogonal approach to previous works. We propose to leverage knowledge in VQA corpora containing questions about images and associated answers for image-caption ranking. Our proposed VQA-based image and caption representations provide complementary information to those learned using previous approaches on a large image-caption ranking dataset.

### 3 Building Blocks: Image-Caption Ranking and VQA

In this section we present image-caption ranking and VQA modules that we build on top of.

#### 3.1 Image-Caption Ranking

The image-caption ranking task is to retrieve relevant images given a query caption, and relevant captions given a query image. During training we are given image-caption pairs  $(I, C)$  that each corresponds to an image  $I$  and its caption  $C$ . For each pair we sample  $K - 1$  other images in addition to  $I$  so the image retrieval task becomes retrieving  $I$  from  $K$  images  $I_i, i = 1, 2 \dots K$  given caption  $C$ . We also sample  $K - 1$  random captions in addition to  $C$  so the caption retrieval task becomes retrieving  $C$  from  $K$  captions  $C_i, i = 1, 2 \dots K$  given image  $I$ .

Our image-caption ranking models learn a ranking scoring function  $S(I, C)$  such that the corresponding retrieval probabilities:

$$P_{im}(I|C) = \frac{\exp(S(I, C))}{\sum_{i=1}^K \exp(S(I_i, C))} \quad P_{cap}(C|I) = \frac{\exp(S(I, C))}{\sum_{i=1}^K \exp(S(I, C_i))} \quad (1)$$

are maximized. Let  $S(I, C)$  be parameterized by  $\theta$  (to be learnt). We formulate an objective function  $L(\theta)$  for  $S(I, C)$  as the sum of expected negative log-likelihoods of image and caption retrieval over all image-caption pairs  $(I, C)$ :

$$L(\theta) = \mathbb{E}_{(I, C)}[-\log P_{im}(I|C)] + \mathbb{E}_{(I, C)}[-\log P_{cap}(C|I)] \quad (2)$$

Recent works on image-caption ranking often construct  $S(I, C)$  by combining a vectorized image representation which is usually hidden layer activations in a CNN pretrained for image classification, with a vectorized caption representation which is usually a sentence encoding computed using an RNN in a multi-modal space. Such scoring functions rely on large image-caption ranking datasets to learn knowledge necessary for image-caption ranking and do not leverage knowledge in VQA corpora. We call such models VQA-agnostic models.

In this work we use the publicly available state-of-the-art image-caption ranking model of [23] as our baseline VQA-agnostic model. [23] projects a  $D_{x_I}$ -dimensional CNN activation  $x_I$  for image  $I$  and a  $D_{x_C}$ -dimensional RNN latent encoding  $x_C$  for caption  $C$  to the same  $D_{x_C}$ -dimensional common multi-modal embedding space as unit-norm vectors  $t_I$  and  $t_C$ :

$$t_I = \frac{W_I x_I}{\|W_I x_I\|_2} \quad t_C = \frac{x_C}{\|x_C\|_2} \quad (3)$$

The multi-modal scoring function is defined as their dot product  $S_t(I, C) = \langle t_I, t_C \rangle$ .

The VQA-agnostic model of [23] uses the 19-layer VGGNet [45] ( $D_{x_I} = 4096$ ) for image encoding and an RNN with 1024 Gated Recurrent Units [7] ( $D_{x_C} = 1024$ ) for caption encoding. The RNN and parameters  $W_I$  are jointly learned on the image-caption ranking training set using a margin-based objective function.

### 3.2 VQA

VQA is the task of given an image  $I$  and a free-form open-ended question  $Q$  about  $I$ , generating a natural language answer  $A$  to that question. Similarly, VQA-Caption task proposed by [1] takes a caption  $C$  of an image and a question  $Q$  about the image, then generates an answer  $A$ . In [1] the generated answers are evaluated using  $\min(\frac{\# \text{ humans that provided } A}{3}, 1)$ . That is,  $A$  is 100% correct if at least 3 humans (out of 10) provide the answer  $A$ .

We closely follow [1] and formulate VQA as a classification task over top  $M = 1000$  most frequent answers from the training set. The oracle accuracies of picking the best answer for each question within this set of answers are 89.37% on training and 88.83% on validation. During training, given triplets of image  $I$ , question  $Q$  and ground truth answer  $A$ , we optimize the negative log-likelihood (NLL) loss to maximize the probability of the ground truth answer  $P_I(A|Q, I)$  given by the VQA model. Similarly given triplets of caption  $C$ , question  $Q$  and ground truth answer  $A$ , we optimize the NLL loss to maximize the VQA-Caption model probability  $P_C(A|Q, C)$ .

Following [1], for a VQA question  $(I, Q)$  we first encode the input image  $I$  using the 19-layer VGGNet [45] as a 4,096-dimensional image encoding  $x_I$ , and encode the question  $Q$  using a 2-layer RNN with 512 Long Short-Term Memory (LSTM) units [20] per layer as a 2,048-dimensional question encoding  $x_Q$ . We then project  $x_I$  and  $x_Q$  into a common 1,024-dimensional multi-modal space as  $z_I$  and  $z_Q$ :

$$z_I = \text{Tanh}(W_I x_I + b_I) \quad z_Q = \text{Tanh}(W_Q x_Q + b_Q) \quad (4)$$

As in [1] we then compute the representation  $z_{I+Q}$  for the image-question pair  $(I, Q)$  by element-wise multiplying  $z_I$  and  $z_Q$ :  $z_{I+Q} = z_I \odot z_Q$ . The scores  $s_A$  for 1,000 answers are given by:

$$s_A = W_s z_{I+Q} + b_s \quad (5)$$

We jointly learn the question encoding RNN and parameters  $\{W_I, b_I, W_Q, b_Q, W_s, b_s\}$  during training.

For the VQA-Caption task given caption  $C$  and question  $Q$ , we use the same network architecture and learning procedure as above, but using the most frequent 1,000 words in training captions as the dictionary to construct a 1,000 dimensional bag-of-words encoding for caption  $C$  as  $x_C$  to replace the image feature  $x_I$  and compute  $z_C, z_{C+Q}$  respectively.

The VQA and VQA-Caption models are learned on the train split of the VQA dataset [1] using 82,783 images, 413,915 captions and 248,349 questions. These models achieve VQA validation set accuracies of 54.42% (VQA) and 56.28% (VQA-Caption), respectively. Next, they are used as sub-modules in our image-caption ranking approach.

## 4 Approach

To leverage knowledge in VQA for image-caption ranking, we propose to represent the images and the captions in the VQA space using VQA and VQA-Caption models. We call such representations VQA-grounded representations.

### 4.1 VQA-Grounded Representations

Let's say we have a VQA model  $P_I(A|Q, I)$ , a VQA-Caption model  $P_C(A|Q, C)$  and a set of  $N$  questions  $Q_i$  and their plausible answers (one for each question)  $A_i$ ,  $i = 1, 2, \dots, N$ . Then given an image  $I$  and a caption  $C$ , we first extract the  $N$  dimensional VQA-grounded activation vectors  $u_I$  for  $I$  and  $u_C$  for  $C$  such that each dimension  $i$  of  $u_I$  and  $u_C$  is the log probability of the ground truth answer  $A_i$  given a question  $Q_i$ .

$$u_I^{(i)} = \log P_I(A_i|Q_i, I) \quad u_C^{(i)} = \log P_C(A_i|Q_i, C), i = 1, 2, \dots, N \quad (6)$$

For example if the  $(Q_i, A_i)$  pairs are  $(Q_1: \text{What is the person riding?}, A_1: \text{Motorcycle})$  and  $(Q_2: \text{What is the man wearing on his head?}, A_2: \text{Helmet})$ ,  $u_I^{(1)}$  and  $u_C^{(1)}$  verify if the person in image  $I$  and caption  $C$  respectively is riding a motorcycle. At the same time  $u_I^{(2)}$  and  $u_C^{(2)}$  verify whether the man in  $I$  and  $C$  is wearing a helmet. Figure 1 shows another example.

In cases where there is not a man in the image or the caption, *i.e.* the assumption of  $Q_i$  is not met,  $P_I(A_i|Q_i, I)$  and  $P_C(A_i|Q_i, C)$  may still reflect if there *were* a man or if the assumption of  $Q_i$  *were* fulfilled, could he be wearing a helmet. In other words, even if there is no person present in the image or mentioned in the caption, the model may still assess the plausibility of a man wearing a helmet or a motorcycle being present. This imagination beyond what is depicted in the image or caption can be helpful in providing additional information when reasoning about the compatibility between an image and a caption. We show qualitative examples of this imagination or plausibility assessment for selected  $(Q, A)$  pairs in Fig. 2 where we sort images and captions based on  $P_I(A|Q, I)$  and  $P_C(A|Q, C)$ . Indeed, scenes where the corresponding fact  $(Q, A)$  (e.g., man is wearing a helmet) is more likely to be plausible are scored higher.<sup>2</sup>

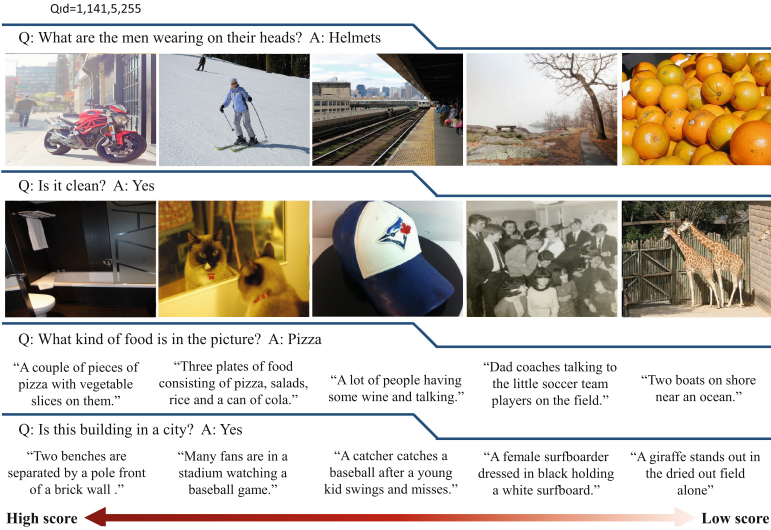
Based on the activation vectors  $u_I$  and  $u_C$ , we then compute the VQA-grounded vector representations  $v_I$  and  $v_C$  for  $I$  and  $C$  by projecting  $u_I$  and  $u_C$  to a  $D_u$ -dimensional vector embedding space:

$$v_I = \sigma(W_{u_I} u_I + b_{v_I}) \quad v_C = \sigma(W_{u_C} u_C + b_{v_C}) \quad (7)$$

Here  $\sigma$  is a non-linear activation function.

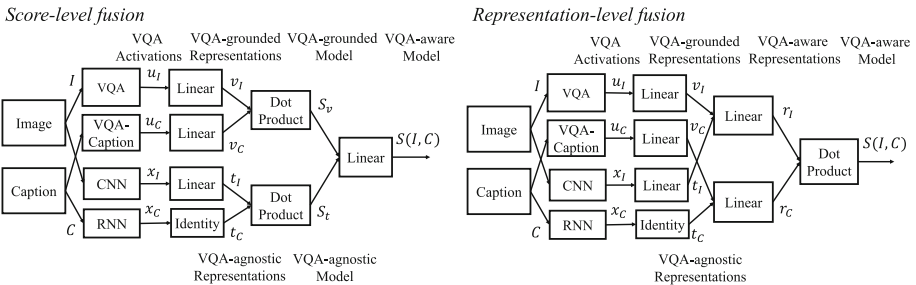
By verifying question-answer pairs on image  $I$  and caption  $C$  and computing vector representations on top of them, the VQA-grounded representations  $v_I$  and  $v_C$  explicitly project image and caption into VQA space to utilize

<sup>2</sup> Nonetheless, checking if a question applies to the target image and caption is also desirable. Contemporary work [41] has looked at modeling  $P(Q|I)$ , and can be incorporated in our approach as an additional feature



**Fig. 2.** Images and captions sorted by  $P_I(A|Q, I)$  and  $P_C(A|Q, C)$  assessed by our VQA (top) and VQA-Caption (bottom) models respectively. Indeed, images and captions that are more plausible for the  $(Q, A)$  pairs are scored higher.

knowledge in the VQA corpora. However, that comes at a cost of losing information such as the sentence structure of the caption and image saliency. These information can also be important for image-caption ranking. As a result, we find VQA-grounded representations are most effective when they are combined with baseline VQA-agnostic models, so we propose two strategies for fusing VQA-grounded representations with baseline VQA-agnostic models: combining their prediction scores or score-level fusion (Fig. 3 left) and combining their representations or representation-level fusion (Fig. 3 right).



**Fig. 3.** We propose score-level fusion (left) and representation-level fusion (right) to utilize VQA for image-caption ranking. They use VQA and VQA-Caption models as “feature extraction” schemes for images and captions and use those features to construct VQA-grounded representations. The score-level fusion approach combines the scoring functions of a VQA-grounded model and a baseline VQA-agnostic model. The representation-level fusion approach combines VQA-grounded representations and VQA-agnostic representations to produce a VQA-aware scoring function.



## 4.2 Score-Level Fusion

A simple strategy to combine our VQA-grounded model with a VQA-agnostic image-ranking model is to combine them at the score level. Given image  $I$  and caption  $C$ , we first compute the VQA-grounded score as the dot product between the VQA-grounded representations of image and caption  $S_v(I, C) = \langle v_I, v_C \rangle$ . We then combine it with the VQA-agnostic scoring function  $S_t(I, C)$  to get the final scoring function  $S(I, C)$ :

$$S(I, C) = \alpha S_t(I, C) + \beta S_v(I, C) \quad (8)$$

We first learn  $\{W_{u_I}, b_{u_I}, W_{u_C}, b_{u_C}\}$  on the image-caption ranking training set, and then learn  $\alpha$  and  $\beta$  on a held out validation set to avoid overfitting.

## 4.3 Representation-Level Fusion

An alternative to combining the VQA-agnostic and VQA-grounded representations at the score level is to inject the VQA-grounding at the representation level. Given the VQA-agnostic  $D_t$ -dimensional image and caption representations  $t_I$  and  $t_C$  used by the baseline model, we first compute the VQA-grounded representations  $v_I$  for image and  $v_C$  for caption introduced in Sect. 4.1. And then they are combined with VQA-agnostic representations to produce VQA-aware representations  $r_I$  for image  $I$  and  $r_C$  for caption  $C$  by projecting them to a  $D_r$ -dimensional multi-modal embedding space as follows:

$$r_I = \sigma(W_{t_I} t_I + W_{v_I} v_I + b_{r_I}) \quad r_C = \sigma(W_{t_C} t_C + W_{v_C} v_C + b_{r_C}) \quad (9)$$

The final image-caption ranking score is then

$$S(I, C) = \langle r_I, r_C \rangle \quad (10)$$

In experiments, we jointly learn  $\{W_{u_I}, b_{u_I}, W_{u_C}, b_{u_C}\}$  (for projecting  $u_I$  and  $u_C$  to the VQA-grounded representations  $v_I, v_C$ ) with  $\{W_{t_I}, W_{v_I}, b_{r_I}, W_{t_C}, W_{v_C}, b_{r_C}\}$  (for computing the combined VQA-aware representations  $r_I$  and  $r_C$ ) on the image-caption ranking training set by optimizing Eq. 2.

Score-level fusion and representation-level fusion models are implemented as multi-layer neural networks. All activation functions  $\sigma$  are  $ReLU(x) = \max(x, 0)$  (for speed) and dropout layers [47] are inserted after all  $ReLU$  layers to avoid overfitting. We set the dimensions of the multi-modal embedding spaces  $D_v$  and  $D_r$  to 4,096 so they are large enough to capture necessary concepts for image-caption ranking. Optimization hyperparameters are selected on the validation set. We optimize both models using RMSProp with batch size 1,000 at learning rate  $1e-5$  for score-level fusion and  $1e-4$  for representation-level fusion. Optimization runs for 100,000 iterations with learning rate decay every 50,000 iterations.

Our main results in Sect. 5.1 use  $N = 3000$  question-answer pairs, sampled 3 questions per image with their ground truth answers with respect to their original images from 1,000 random VQA training images. We discuss using different numbers of question-answer pairs  $N$  and different strategies for selecting the question-answer pairs in Sect. 5.4.

## 5 Experiments and Results

We report results on MSCOCO [31] which is the largest available image-caption ranking dataset. Following the splits of [22, 23] we use all 82,783 MSCOCO train images with 5 captions per image as our train set, 413,915 image-caption pairs in total. Note that this is the same split as the train split in the VQA dataset [1] we used to train our VQA and VQA-Caption models. The validation set consists of 1,000 images sampled from the original MSCOCO validation images. The test set consists of 5,000 images sampled from the original MSCOCO validation images that were not in the image-caption ranking validation set. Same as the train set, there are 5 captions available for each validation and test image.

We follow the evaluation metric of [22] and report caption and image retrieval performances on the first 1,000 test images following [22–24, 33, 37]. Given a test image, the caption retrieval task is to find any 1 out of its 5 captions from all 5,000 test captions. Given a test caption, the image retrieval task is to find its original image from all 1,000 test images. We report recall@(1, 5, 10): the fraction of times a correct item was found among the top (1, 5, 10) predictions.

### 5.1 Image-Caption Ranking Results

Table 1 shows our main results on MSCOCO. Our score-level fusion VQA-aware model using  $N = 3000$  question-answer pairs (“ $N = 3000$  score-level fusion VQA-aware”) achieves 46.9% caption retrieval recall@1 and 35.8% image retrieval recall@1. This model shows an improvement of 3.5% caption and 4.8% image retrieval recall@1 over the state-of-the-art VQA-agnostic model of [23].

Our representation-level fusion approach adds an additional layer on top of the VQA-agnostic representations, resulting in a deeper model, so we experiment with adding an additional layer to the VQA-agnostic model for a fair comparison. That is equivalent to representation-level fusion using  $N = 0$  question-answer pair (“ $N = 0$  representation-level fusion”, *i.e.* deeper VQA-agnostic). Comparing with the VQA-agnostic model of [23], adding this additional layer improves performance by 2.4% caption and 2.6% image retrieval recall@1.

By leveraging VQA knowledge our “ $N = 3000$  representation-level fusion VQA-aware” model achieves 50.5% caption retrieval recall@1 and 37.0% image retrieval recall@1, which further improves 4.7% and 3.4% over the  $N = 0$  VQA-agnostic representation-level fusion model. These improvements are consistent with our score-level fusion approach so this shows that the VQA corpora consistently provide complementary information to image-caption ranking.

To the best of our knowledge, the  $N = 3000$  representation-level fusion VQA-aware result is the best result on MSCOCO image-caption ranking and significantly surpasses previous best results by as much as 7.1% in caption retrieval and 4.4% image retrieval recall@1.

Our VQA-grounded model alone (“ $N = 3000$  score-level fusion VQA-grounded only”) achieves 37.0% caption and 26.2% image retrieval recall@1. This indicates that the VQA activations  $u_I$  and  $u_C$  which evaluate the

**Table 1.** Caption retrieval and image retrieval performances of our models compared to baseline models on MSCOCO image-caption ranking test set. Powered by knowledge in VQA corpora, both our score-level fusion and representation-level fusion VQA-aware approaches outperform state-of-the-art VQA-agnostic models by a large margin

MSCOCO						
Approach	Caption retrieval			Image retrieval		
	R@1	R@5	R@10	R@1	R@5	R@10
Random	0.1	0.5	1.0	0.1	0.5	1.0
DVSA [22]	38.4	69.9	80.5	27.4	60.2	74.8
FV (GMM+HGLMM) [24]	39.4	67.9	80.9	25.1	59.8	76.6
<i>m</i> -RNN-vgg [37]	41.0	73.0	83.5	29.0	42.2	77.0
<i>m</i> -CNN <sub>ENS</sub> [33]	42.8	73.1	84.1	32.6	68.6	82.8
Kiros <i>et al.</i> [23] (VQA-agnostic)	43.4	75.7	85.8	31.0	66.7	79.9
N = 3000 score-level fusion VQA-grounded only	37.0	67.9	79.4	26.2	60.1	74.3
N = 3000 score-level fusion VQA-aware	46.9	78.6	88.9	35.8	70.3	<b>83.6</b>
N = 0 representation-level fusion VQA-agnostic	45.8	76.8	86.1	33.6	67.8	81.0
N = 3000 representation-level fusion VQA-aware	<b>50.5</b>	<b>80.1</b>	<b>89.7</b>	<b>37.0</b>	<b>70.9</b>	82.9

plausibility of facts (question-answer pairs) in images and captions are informative representations.

Figure 4 shows qualitative results on image retrieval comparing our approach ( $N = 3000$  score-level fusion) with the VQA-agnostic model. By looking at several top retrieved images from our model for the failure case (last column), we find that our model seems to have picked up on a correlation between bats and helmets. It seems to be looking for helmets in retrieved images, while the ground truth image does not have one.

We also experiment with using the hidden activations available in the VQA and VQA-Caption models ( $z_I$  and  $z_C$  in Sect. 3.2) as image and caption encodings in place of the VQA activations ( $u_I$  and  $u_C$  in Sect. 4.1). Using these hidden activations of the VQA models is conceptually similar to using the hidden activations of CNNs pretrained on ImageNet as features [11]. These features achieve 46.8% caption retrieval recall@1 and 35.2% image retrieval recall@1 for score-level fusion, and 49.3% caption retrieval recall@1 and 37.9% image retrieval recall@1 for representation-level fusion which are as good as our semantic features  $u_I$  and  $u_C$ . This shows that our semantically meaningful features,  $u_I$  and  $u_C$ , performs as well as their corresponding non-semantic representations  $z_I$  and  $z_C$  using both score-level fusion and representation-level fusion. Note that such hidden activations may not always be available in different VQA models and the semantic features have the added benefit of being interpretable (*e.g.*, Fig. 2).



**Fig. 4.** Qualitative image retrieval results of our score-level fusion VQA-aware model (middle) and the VQA-agnostic model (bottom). The true target image is highlighted (green if VQA-aware found it, red if VQA-agnostic found it but VQA-aware did not). (Color figure online)

## 5.2 Ablation Study

As an ablation study, we compare the following four models: (1) full representation-level fusion: our full  $N = 3000$  representation-level fusion model that includes both image and caption VQA representations; (2) caption-only representation-level fusion: the same representation-level fusion model but using the VQA representation only for the caption,  $v_C$ , and not for the image; (3) image-only representation-level fusion: the same model but using the VQA representation only for the image,  $v_I$ , and not for the caption; (4) deeper VQA-agnostic: The  $N = 0$  representation-level fusion model described earlier that does not use VQA representations for neither the image nor the caption.

Table 2 summarizes the results. We see that incrementally adding more VQA-knowledge improves performance. Both caption-only and image-only models outperform the  $N = 0$  deeper VQA-agnostic baseline. The full representation-level fusion model which combines both representations yields the best performance.

**Table 2.** Ablation study evaluating the gain in performance as more VQA-knowledge is incorporated in the model

MSCOCO						
Approach	Caption retrieval			Image retrieval		
	R@1	R@5	R@10	R@1	R@5	R@10
Deeper VQA-agnostic	45.8	76.8	86.1	33.6	67.8	81.0
Caption-only representation-level fusion	47.3	77.3	86.6	35.5	69.3	81.9
Image-only representation-level fusion	47.0	80.0	89.6	36.4	70.1	82.3
Full representation-level fusion	<b>50.5</b>	<b>80.1</b>	<b>89.7</b>	<b>37.0</b>	<b>70.9</b>	<b>82.9</b>

### 5.3 The Role of VQA and Caption Annotations

In this work we transfer knowledge from one vision-language task (*i.e.* VQA) to another (*i.e.* image-caption ranking). However, VQA annotations and caption annotations serve different purposes.

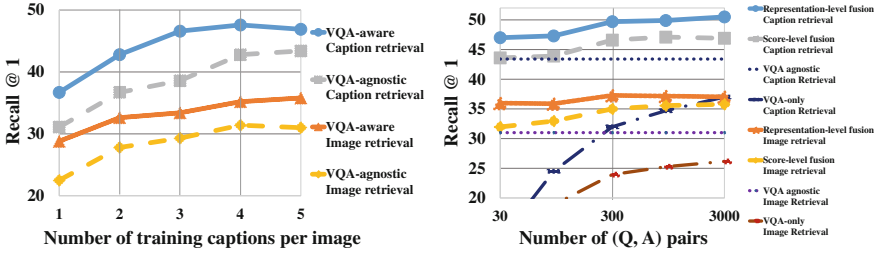
The target language to be retrieved is caption language, and not VQA language. [1] showed qualitatively and quantitatively that the two languages are statistically quite different (in terms of information contained, and in terms of nouns, adjectives, verbs, etc. used). As a result, VQA can not be thought of as providing additional “annotations” for the captioning task. Instead, VQA provides different perspectives/views of the images (and captions). It provides an additional feature representation. To better utilize this representation for an image-caption ranking task, one would still require sufficient ground truth caption annotations for images. In fact, with varying amounts of ground truth (caption) annotations, the VQA-aware representations show improvements in performance across the board. See Fig. 5 (left).

A better analogy of our VQA representation is hidden activations (*e.g.*, fc7) from a CNN trained on ImageNet. Having additional ImageNet annotations would improve the fc7 feature. But to map this fc7 feature to captions, one would still require sufficient caption annotations. Conceptually, caption annotations and category labels in ImageNet play two different roles. The former provides ground truth for the target task at hand (image-caption ranking), and having additional annotations for the target application typically helps. The latter helps learn a better image representation (which may provide improvements in a variety of tasks).

### 5.4 Number of Question-Answer Pairs

Our VQA-grounded representations extract image and caption features based on question-answer pairs. It is important for there to be enough question-answer pairs to cover necessary aspects for image-caption ranking. We experiment with using  $N = 30, 90, 300, 900, 3000$  ( $Q, A$ ) pairs (or facts) for both score-level and representation-level fusion. Figure 5 (right) shows caption and image retrieval performances of our approaches with varying  $N$ . Performance of both score-level and representation-level fusion approaches improve quickly from  $N = 30$  to  $N = 300$ , and then starts to level off after  $N = 300$ .

An alternative to sampling 3 question-answer pairs per image on 1,000 images to get  $N = 3000$  questions is to sample 1 question-answer pair per image from 3,000 images. Sampling multiple ( $Q, A$ ) pairs from the same image provides correlated ( $Q, A$ ) pairs. For example ( $Q$ : What are these animals?  $A$ : Giraffes) and ( $Q$ : Would this animal fit in a house?  $A$ : No). Using such correlated ( $Q, A$ ) pairs, the model could potentially better predict if there is a giraffe in the image by jointly reasoning if the animal looks like a giraffe and the if the animal would fit in a house, if the VQA and VQA-Caption models have not already picked up such correlations. In experiments, sampling 3 question-answer pairs per image



**Fig. 5. Left:** caption retrieval and image retrieval performances of the VQA-agnostic model compared with our  $N = 3000$  score-level fusion VQA-aware model trained using 1 to 5 captions per image. The VQA representations in the VQA-aware model provide consistent performance gains. **Right:** caption retrieval and image retrieval performances of our score-level fusion and representation-level fusion approaches with varying number of  $(Q, A)$  pairs used for feature extraction.

for correlated  $(Q, A)$  pairs does not significantly outperform sampling 1 question-answer pair per image which performs (47.7 %, 35.4 %) (image, caption) recall@1 using  $N = 3000$  score-level fusion, so we hypothesize that our VQA and Caption-QA models have already captured such correlations.

## 6 Conclusion

VQA corpora provide rich multi-modal information that is complementary to knowledge stored in image captioning corpora. In this work we take the novel perspective of viewing VQA as a “feature extraction” module that captures VQA knowledge. We propose two approaches – score-level and representation-level fusion – to integrate this knowledge into an existing image-caption ranking model. We set new state-of-the-art by improving caption retrieval by 7.1 % and image retrieval by 4.4 % on MSCOCO.

Improved individual modules, *i.e.*, VQA models and VQA-agnostic image-caption ranking models, end-to-end training, and an attention mechanism that selects question-answer pairs (facts) in an image-specific manner may further improve the performance of our approach.

**Acknowledgment.** This work was supported in part by the Allen Distinguished Investigator awards by the Paul G. Allen Family Foundation, a Google Faculty Research Award, a Junior Faculty award by the Institute for Critical Technology and Applied Science (ICTAS) at Virginia Tech, a National Science Foundation CAREER award, an Army Research Office YIP award, and Office of Naval Research YIP award to D. P. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the U.S. Government or any sponsor.

## References

1. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: VQA: visual question answering. In: ICCV (2015)
2. Antol, S., Zitnick, C.L., Parikh, D.: Zero-shot learning via visual abstraction. In: ECCV (2014)
3. Berg, T., Belhumeur, P.N.: POOF: part-based one-vs.-one features for fine-grained categorization, face verification, and attribute estimation. In: CVPR (2013)
4. Bourdev, L., Maji, S., Brox, T., Malik, J.: Detecting people using mutually consistent poselet activations. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part VI. LNCS, vol. 6316, pp. 168–181. Springer, Heidelberg (2010)
5. Branson, S., Wah, C., Schroff, F., Babenko, B., Welinder, P., Perona, P., Belongie, S.: Visual recognition with humans in the loop. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part IV. LNCS, vol. 6314, pp. 438–451. Springer, Heidelberg (2010)
6. Chen, X., Lawrence Zitnick, C.: Mind’s eye: a recurrent visual representation for image caption generation. In: CVPR (2015)
7. Cho, K., van Merriënboer, B., Bahdanau, D., Bengio, Y.: On the properties of neural machine translation: encoder-decoder approaches (2014). arXiv preprint [arXiv:1409.1259](https://arxiv.org/abs/1409.1259)
8. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: ImageNet: a large-scale hierarchical image database. In: CVPR (2009)
9. Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: ICCV (2015)
10. Donahue, J., Grauman, K.: Annotator rationales for visual recognition. In: ICCV (2011)
11. Donahue, J., Jia, Y., Vinyals, O., Hoffman, J., Zhang, N., Tzeng, E., Darrell, T.: DeCAF: a deep convolutional activation feature for generic visual recognition (2013). arXiv preprint [arXiv:1310.1531](https://arxiv.org/abs/1310.1531)
12. Elhoseiny, M., Saleh, B., Elgammal, A.: Write a classifier: zero-shot learning using purely textual descriptions. In: ICCV (2013)
13. Elliott, D., Keller, F.: Comparing automatic evaluation measures for image description. In: Proceedings of 52nd Annual Meeting of the Association for Computational Linguistics, pp. 452–457 (2014)
14. Farhadi, A., Hejrati, M., Sadeghi, M.A., Young, P., Rashtchian, C., Hockenmaier, J., Forsyth, D.: Every picture tells a story: generating sentences from images. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part IV. LNCS, vol. 6314, pp. 15–29. Springer, Heidelberg (2010)
15. Farhadi, A., Endres, I., Hoiem, D., Forsyth, D.: Describing objects by their attributes. In: CVPR (2009)
16. Fouhey, D.F., Zitnick, C.L.: Predicting object dynamics in scenes. In: CVPR (2014)
17. Gao, H., Mao, J., Zhou, J., Huang, Z., Wang, L., Xu, W.: Are you talking to a machine? Dataset and methods for multilingual image question answering. In: NIPS (2015)
18. Gupta, A., Davis, L.S.: Beyond nouns: exploiting prepositions and comparative adjectives for learning visual classifiers. In: Forsyth, D., Torr, P., Zisserman, A. (eds.) ECCV 2008, Part I. LNCS, vol. 5302, pp. 16–29. Springer, Heidelberg (2008)
19. Hamrick, J., Battaglia, P., Tenenbaum, J.B.: Internal physics models guide probabilistic judgments about object dynamics. In: Proceedings of 33rd Annual Meeting of the Cognitive Science Society, Boston, MA (2011)

20. Hochreiter, S., Schmidhuber, J.: Long short-term memory. *Neural Comput.* **9**(8), 1735–1780 (1997)
21. Johnson, J., Krishna, R., Stark, M., Li, L.J., Shamma, D., Bernstein, M., Fei-Fei, L.: Image retrieval using scene graphs. In: *CVPR* (2015)
22. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: *CVPR* (2015)
23. Kiros, R., Salakhutdinov, R., Zemel, R.: Unifying visual-semantic embeddings with multimodal neural language models. In: *TACL* (2015)
24. Klein, B., Lev, G., Sadeh, G., Wolf, L.: Associating neural word embeddings with deep image representations using Fisher vectors. In: *CVPR* (2015)
25. Kovashka, A., Parikh, D., Grauman, K.: WhittleSearch: image search with relative attribute feedback. In: *CVPR* (2012)
26. Krizhevsky, A., Sutskever, I., Hinton, G.E.: Imagenet classification with deep convolutional neural networks. In: *NIPS* (2012)
27. Kulkarni, G., Premraj, V., Dhar, S., Li, S., Choi, Y., Berg, A.C., Berg, T.L.: Baby talk: understanding and generating simple image descriptions. In: *CVPR* (2011)
28. Kumar, N., Berg, A.C., Belhumeur, P.N., Nayar, S.K.: Describable visual attributes for face verification and image search. In: *IEEE TPAMI* (2011)
29. Lampert, C.H., Nickisch, H., Harmeling, S.: Learning to detect unseen object classes by between-class attribute transfer. In: *CVPR* (2009)
30. Li, L.J., Su, H., Fei-Fei, L., Xing, E.P.: Object bank: a high-level image representation for scene classification and semantic feature sparsification. In: *NIPS* (2010)
31. Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft COCO: common objects in context. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) *ECCV 2014, Part V. LNCS*, vol. 8693, pp. 740–755. Springer, Heidelberg (2014)
32. Lin, X., Parikh, D.: Don’t just listen, use your imagination: leveraging visual common sense for non-visual tasks. In: *CVPR* (2015)
33. Ma, L., Lu, Z., Li, H.: Learning to answer questions from image using convolutional neural network (2015). arXiv preprint [arXiv:1506.00333](https://arxiv.org/abs/1506.00333)
34. Ma, L., Lu, Z., Shang, L., Li, H.: Multimodal convolutional neural networks for matching image and sentence. In: *ICCV* (2015)
35. Malinowski, M., Fritz, M.: A multi-world approach to question answering about real-world scenes based on uncertain input. In: *NIPS* (2014)
36. Malinowski, M., Rohrbach, M., Fritz, M.: Ask your neurons: a neural-based approach to answering questions about images. In: *ICCV* (2015)
37. Mao, J., Xu, W., Yang, Y., Wang, J., Huang, Z., Yuille, A.: Deep captioning with multimodal recurrent neural networks (M-RNN). In: *ICLR* (2015)
38. Parikh, D., Grauman, K.: Relative attributes. In: *ICCV* (2011)
39. Parkash, A., Parikh, D.: Attributes for Classifier Feedback. In: Fitzgibbon, A., Lazebnik, S., Perona, P., Sato, Y., Schmid, C. (eds.) *ECCV 2012, Part III. LNCS*, vol. 7574, pp. 354–368. Springer, Heidelberg (2012)
40. Pirsivash, H., Vondrick, C., Torralba, A.: Inferring the why in images. *CoRR* abs/1406.5472 (2014). <http://arXiv.org/abs/1406.5472>
41. Ray, A., Christie, G., Bansal, M., Batra, D., Parikh, D.: Question relevance in VQA: identifying non-visual and false-premise questions (2016). arXiv preprint [arXiv:1606.06622](https://arxiv.org/abs/1606.06622)
42. Ren, M., Kiros, R., Zemel, R.: Exploring models and data for image question answering. In: *NIPS* (2015)
43. Sadanand, S., Corso, J.J.: Action bank: a high-level representation of activity in video. In: *CVPR* (2012)



44. Sadeghi, F., Divvala, S.K., Farhadi, A.: VisKE: visual knowledge extraction and question answering by visual verification of relation phrases. In: CVPR (2015)
45. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. CoRR abs/1409.1556 (2014)
46. Socher, R., Ganjoo, M., Manning, C.D., Ng, A.: Zero-shot learning through cross-modal transfer. In: NIPS (2013)
47. Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., Salakhutdinov, R.: Dropout: a simple way to prevent neural networks from overfitting. JMLR **15**, 1929–1958 (2014)
48. Sutskever, I., Vinyals, O., Le, Q.V.: Sequence to sequence learning with neural networks. In: NIPS (2014)
49. Szegedy, C., Liu, W., Jia, Y., Sermanet, P., Reed, S., Anguelov, D., Erhan, D., Vanhoucke, V., Rabinovich, A.: Going deeper with convolutions. In: CVPR (2015)
50. Tang, K., Paluri, M., Fei-fei, L., Fergus, R., Bourdev, L.: Improving image classification with location context. In: ICCV (2015)
51. Vedantam, R., Lin, X., Batra, T., Zitnick, C.L., Parikh, D.: Learning common sense through visual abstraction. In: ICCV (2015)
52. Vedantam, R., Zitnick, C.L., Parikh, D.: Cider: Consensus-based image description evaluation. In: CVPR (2015)
53. Vondrick, C., Pirsivash, H., Torralba, A.: Anticipating the future by watching unlabeled video (2015). arXiv preprint [arXiv:1504.08023](https://arxiv.org/abs/1504.08023)
54. Walker, J., Gupta, A., Hebert, M.: Patch to the future: unsupervised visual prediction. In: CVPR (2014)
55. Wang, Y., Mori, G.: A discriminative latent model of object classes and attributes. In: Daniilidis, K., Maragos, P., Paragios, N. (eds.) ECCV 2010, Part V. LNCS, vol. 6315, pp. 155–168. Springer, Heidelberg (2010)
56. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhutdinov, R., Zemel, R., Bengio, Y.: Show, attend and tell: neural image caption generation with visual attention. In: ICML (2015)
57. Yu, L., Park, E., Berg, A.C., Berg, T.L.: Visual madlibs: fill in the blank image generation and question answering (2015). arXiv preprint [arXiv:1506.00278](https://arxiv.org/abs/1506.00278)
58. Zhang, N., Farrell, R., Iandola, F., Darrell, T.: Deformable part descriptors for fine-grained recognition and attribute prediction. In: ICCV (2013)
59. Zhang, N., Paluri, M., Ranzato, M., Darrell, T., Bourdev, L.: PANDA: pose aligned networks for deep attribute modeling. In: CVPR (2014)
60. Zheng, B., Zhao, Y., Yu, J., Ikeuchi, K., Zhu, S.C.: Beyond point clouds: scene understanding by reasoning geometry and physics. In: CVPR (2013)
61. Zhou, B., Tian, Y., Sukhbaatar, S., Szlam, A., Fergus, R.: Simple baseline for visual question answering (2015). arXiv preprint [arXiv:1512.02167](https://arxiv.org/abs/1512.02167)
62. Zhu, Y., Fathi, A., Fei-Fei, L.: Reasoning about object affordances in a knowledge base representation. In: Fleet, D., Pajdla, T., Schiele, B., Tuytelaars, T. (eds.) ECCV 2014, Part II. LNCS, vol. 8690, pp. 408–424. Springer, Heidelberg (2014)
63. Zhu, Y., Zhang, C., Re, C., Fei-Fei, L.: Building a large-scale multimodal knowledge base for visual question answering (2013). arXiv preprint [arXiv:1310.1531](https://arxiv.org/abs/1310.1531)