

# Mean Box Pooling: A Rich Image Representation and Output Embedding for the Visual Madlibs Task

Ashkan Mokarian  
ashkan@mpi-inf.mpg.de

Mateusz Malinowski  
mmalinow@mpi-inf.mpg.de

Mario Fritz  
mfritz@mpi-inf.mpg.de

Scalable Learning and Perception  
Max Planck Institute for Informatics  
Saarbrücken, Germany

## Abstract

We present Mean Box Pooling, a novel visual representation that pools over CNN representations of a large number, highly overlapping object proposals. We show that such representation together with nCCA, a successful multimodal embedding technique, achieves state-of-the-art performance on the Visual Madlibs task. Moreover, inspired by the nCCA’s objective function, we extend classical CNN+LSTM approach to train the network by directly maximizing the similarity between the internal representation of the deep learning architecture and candidate answers. Again, such approach achieves a significant improvement over the prior work that also uses CNN+LSTM approach on Visual Madlibs.

## 1 Introduction

Question answering about real-world images is a relatively new research thread [2, 5, 14, 15] that requires a chain of machine visual perception, natural language understanding, and deductive capabilities to successfully come up with an answer on a question about visual content. Although similar in nature to image description [3, 8, 27] it requires a more focused attention to details in the visual content, yet it is easier to evaluate different architectures on the task. Moreover, in contrast to many classical Computer Vision problems such as recognition or detection, the task does not evaluate any internal representation of methods, yet it requires a holistic understanding of the image. Arguably, it is also less prone to over-interpretations compared with the classical Turing Test [16, 25].

To foster progress on this task, a few metrics and datasets have been proposed [2, 4, 14, 20]. The recently introduced Visual Madlibs task [32] removes ambiguities in question or scene interpretations by introducing a multiple choice “filling the blank” task, where a

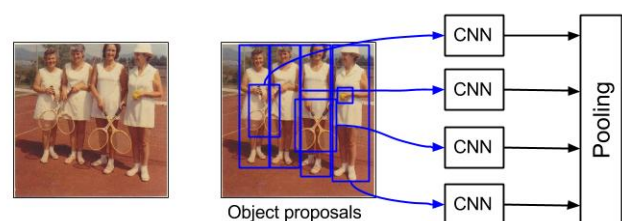


Figure 1: Illustration of proposed Mean Box Pooling representation.

machine has to complete the prompted sentence. Such completed sentence is next matched against four ground truth answers. Thanks to such a problem formulation, a traditional accuracy measure can be used to monitor the progress on this task. Due to its unambiguous evaluation, this work focuses on this task.

**Contributions.** We present two main contributions.

*Mean Box Pooling:* We argue for a rich image representation in the form of pooled representations of the objects. Although related ideas have been explored for visual question answering [22], and even have been used in Visual Madlibs [32], we are first to show a significant improvement of such representation by using object proposals. More precisely, we argue for an approach that pools over a large number, highly overlapping object proposals. This, arguably, increases the recall of extracting bounding boxes that describe an object, but also allows for multi-scale and multi-parts object representation. Our approach in the combination with the Normalized Correlation Analysis embedding technique improves on the state-of-the-art of the Visual Madlibs task.

*Text-Embedding Loss:* Motivated by the popularity of deep architectures for visual question answering, that combine a global CNN image representation with an LSTM [7] question representation [4, 13, 17, 20, 29, 30, 31], as well as the leading performance of nCCA on the multi-choice Visual Madlibs task [32], we propose a novel extension of the CNN+LSTM architecture that chooses a prompt completion out of four candidates (see Figure 4) by measuring similarities directly in the embedding space. This contrasts with the prior approach of [32] that uses a post-hoc comparison between the discrete output of the CNN+LSTM method and all four candidates. To achieve this, we directly train an LSTM with a cosine similarity loss between the output embedding of the network and language representation of the ground truth completion. Such an approach integrates more tightly with the multi-choice filling the blanks task, and significantly outperforms the prior CNN+LSTM method [32].

## 2 Related Work

Question answering about images is a relatively new task that switches focus from recognizing objects in the scene to a holistic “image understanding”. The very first work [14] on this topic has considered real world indoor scenes with associated natural language questions and answers. Since then different variants and larger datasets have been proposed: FM-IQA [4], COCO-QA [20], and VQA [2]. Although answering questions on images is, arguably, more susceptible to automatic evaluation than the image description task [3, 8, 27], ambiguities in the output space still remain. While such ambiguities can be handled using appropriate metrics [14, 15, 17, 26], Visual Madlibs [32] has taken another direction, and handles them directly within the task. It asks machines to fill the blank prompted with a natural language description with a phrase chosen from four candidate completions (Figure 4). In general, the phrase together with the prompted sentence should serve as the accurate description of the image. With such problem formulation the standard accuracy measure is sufficient to automatically evaluate the architectures. The first proposed architecture [14] to deal with the question answering about images task uses image analysis methods and a set of hand-defined schemas to create a database of visual facts. The mapping from questions to executable symbolic representations is done by a semantic parser [12]. Later deep learning approaches for question answering either generate [4, 17] answers or predict answers [13, 20] over a fixed set of choices. Most recently, attention based architectures, which put weights on a fixed grid

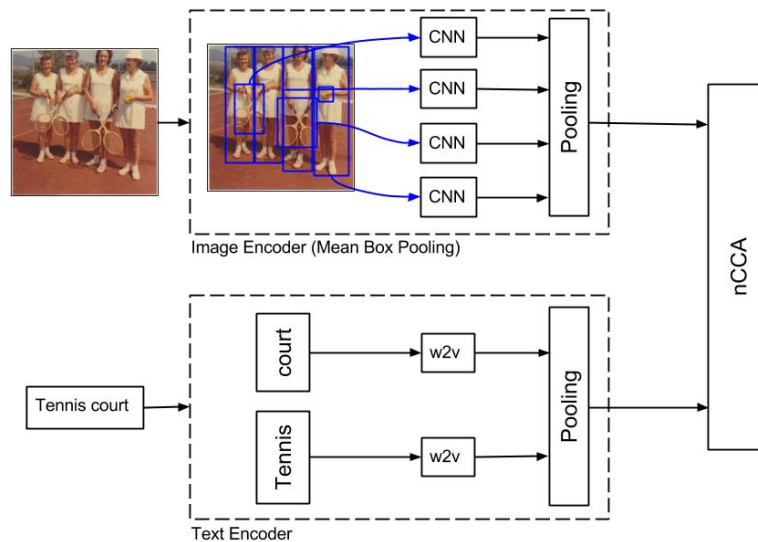


Figure 2: Overview of our full model, i.e. our proposed image representation using Mean Box Pooling, text encoding using average of Word2Vec representations, and normalized CCA for learning the joint space.

over the image, yield state of the art results [29, 30, 31]. Another, more focused “hard” attention, has also been studied in the image-to-text retrieval scenario [9] as well as fine-grained categorization [33], person recognition [19] and zero-shot learning [1]. Here representations are computed on objects, visual fragments or parts, that are further aggregated to form a visual representation. Closer to our work, [22] use Edge Boxes [34] to form memories [28] consisting of different image fragments that are either pooled or “softmax” weighted in order to provide the final score. However, in contrast to [22], our experiments indicate a strong improvement by using object proposals. While a majority of the most recent work on visual question answering combine LSTM [7] with CNN [11, 23, 24] by concatenation or summation or piece-wise multiplication, Canonical Correlation Analysis (CCA and nCCA) [6] have also been shown to be a very effective multimodal embedding technique [32]. Our work further investigates this embedding method as well as brings ideas from CCA over to an CNN+LSTM formulation.

### 3 Method

We use normalized CCA (nCCA) [6] to embed the textual embedding of answers and the visual representation of the image into a joint space, where candidate sentence completions are compared to the image. Furthermore, we also extend popular in the VQA community CNN+LSTM approach by learning to compare in the answer space.

In Section 3.1, we propose a richer representation of the entire image obtained by pooling of CNN representations extracted from object proposals. Figure 1 illustrates the proposed Mean Box Pooling image representation and Figure 2 illustrates our whole method. In Section 3.3, we describe nCCA approach to encode two modalities into a joint space in greater details. In Section 3.4, we also investigate a CNN+LSTM architecture. Instead of generating a prompt completion that is next compared against candidate completions in a post-hoc process, we propose to choose a candidate completion by directly comparing candidates in the embedding space. This puts CNN+LSTM approach closer to nCCA with a tighter integration with the multi-choice Visual Madlibs task. This approach is depicted in Figure 3.

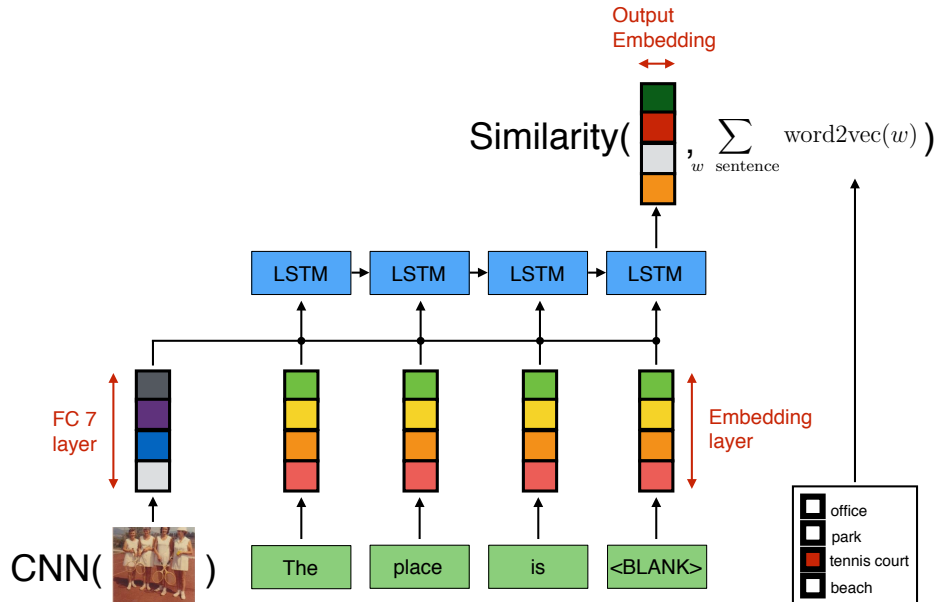


Figure 3: CNN+LSTM architecture that learns to choose the right answer directly in the embedding space. The output embedding is jointly trained with the whole architecture via backpropagation.

### 3.1 Mean Box Pooling Image Representation

Figure 1 illustrates our proposed image representation, which starts from extracting object proposals from the raw image. Next such object proposals are encoded via a CNN, and pooled in order to create a feature vector representation of the image.

**Extracting Region Proposals.** Since questions are mainly about salient parts of the image, it seems reasonable to use object detection in order to extract such parts from the image. At the same time, however, it is important to not miss any object in the image. Moreover, arguably, sampling a context of the objects and capturing multi-scale, multi-parts properties seem to be important as well. Given all these reasons, we choose to use Edge Boxes [34] in order to generate a set of object bounding box proposals for feature extraction.

Edge Boxes extract a number of bounding boxes along with a score for each bounding box that is interpreted as a confidence score that the bounding box contains an object. In our study, two hyper parameters are important: Non-Maxima Suppression and the number of proposals. The latter defines how many object proposals we want to maintain and hence implicitly influence recall of the proposals, while the former defines a threshold  $\beta$  such that all predicted bounding boxes with the intersection over union greater than  $\beta$  are removed. In practice, the lower the  $\beta$  the more spread the object proposals are.

**Feature Extraction.** Once the object proposals are extracted, we use output of the “fc7” layer of the VGG network [23] on the extracted image crops to encode the proposals. VGG is among the best performing recognition architectures on the large scale object recognition task [21].

**Pooling for Image Representation.** Our final image representation is constructed by pooling the encoded object proposals together with the global image representation. Since we do not want to associate any particular order over the extracted object proposals, we investigate popular order-less pooling schemes.

## 3.2 Pooling for Answer Representation.

We encode each word in the answer with a 300 dimensional word embedding [18]. The embedded words are next mean pooled to form a vector representation of the answer. Note that, we do not encode prompts as they follow the same pattern for each Visual Madlibs category.

## 3.3 Multimodal Embedding

We use the Normalized Canonical Correlation Analysis (nCCA) to learn a mapping from two modalities: image and textual answers, into a joint embedding space. This embedding method has shown outstanding performance on the Visual Madlibs task [32]. At the test time, given the encoded image, we choose an answer (encoded by the mean pooling over word2vec words representations) from the set of four candidate answers that is the most similar to the encoded image in the multimodal embedding space. Formally, the Canonical Correlation Analysis (CCA) maximizes the cosine similarity between two modalities (also called views) in the embedding space, that is:

$$W_1^*, W_2^* = \arg \max_{W_1, W_2} \text{tr}(\hat{X}^T \hat{Y})$$

$$\text{subject to } \hat{X}^T \hat{X} = \hat{Y}^T \hat{Y} = I$$

where  $\text{tr}$  is the matrix trace,  $\hat{X} := XW_1$ ,  $\hat{Y} := YW_2$ , and  $X, Y$  are two views (encoded images, and textual answers in our case). Normalized Canonical Correlation Analysis (nCCA) [6] has been reported to work significantly better than the plain CCA. Here, columns of the projection matrices  $W_1$  and  $W_2$  are scaled by the p-th power ( $p=4$ ) of the corresponding eigen values. The improvement is consistent with the findings of [32], where nCCA performs better than CCA by about five percentage points in average on the hard task.

## 3.4 CNN+LSTM with Text-Embedding Loss

We present our novel architecture that extends prior approaches on question answering about images [4, 13, 17, 20, 29, 30, 31] by learning similarity between candidate labels and internal output embedding of the neural network. Figure 3 depicts our architecture. Similarly to prior work, we encode an image with a CNN encoder that is next concatenated with (learnable) word embeddings of the prompt sentence, and fed to a recurrent neural network. We use a special ‘<BLANK>’ token to denote the empty blank space in the image description. On the other side, for each completion candidate  $s$  we compute its representation by averaging over word2vec [18] representations of the words contributing to  $s$ . However, in contrast to the prior work [32], instead of comparing the discrete output of the network with the representation of  $s$ , we directly optimize an objective in the embedding space. During training we maximize the similarity measure between the output embedding and the representation of  $\sigma$  by optimizing the following objective:

$$\Theta^* = \arg \max_{\Theta} \sum_i \frac{\text{embedding}(x_i; \Theta)^T (\sum_{w \in \hat{s}_i} \text{word2vec}(w))}{\|\text{embedding}(x_i; \Theta)\| \|\sum_{w \in \hat{s}_i} \text{word2vec}(w)\|},$$

which is a cosine similarity between the representation of the available during the training correct completion  $\hat{s}_i$ , and an output embedding vector of the i-th image-prompt training



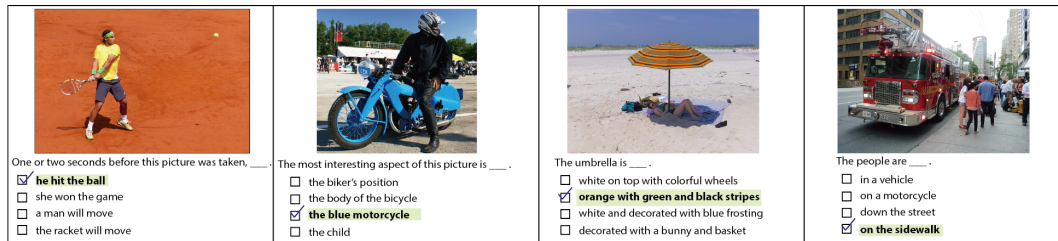


Figure 4: Some examples of the multi-choice filling the blank Visual Madlib task [32].

instance  $x_i$ ;  $\Theta$  denotes all the parameters of the architecture. At test time, we choose a completion  $\hat{s}$  by:

$$\hat{s} = \arg \max_{s \in \mathcal{S}} \frac{\text{embedding}(x; \Theta^*)^T (\sum_{w \in \mathcal{S}} \text{word2vec}(w))}{\|\text{embedding}(x; \Theta^*)\| \|\sum_{w \in \mathcal{S}} \text{word2vec}(w)\|},$$

where  $\mathcal{S}$  denotes a set of available candidate prompt completions,  $x$  is the image-prompt pair fed to the network, and  $\Theta^*$  denotes all the learnt parameters.

## 4 Experimental Results

We evaluate our method on the multiple choice task of the Visual Madlibs dataset. The dataset consists of about 360k descriptions, spanning 12 different categories specified by different types of templates, of about 10k images. The selected images from the MS COCO dataset comes with rich annotations. In the multi-choice scenario a textual prompt is given (every category follows the same, fixed template) with a blank to be filled, together with 4 candidate completions (see Figure 4). Every category represents a different type of question including scenes, affordances, emotions, or activities (the full list is shown in the first column of Table 1). Since each category has fixed prompt, there is no need to include the prompt in the modeling given the training is done per each category. Finally, Visual Madlibs considers an easy and difficult tasks that differ in how the negative 3 candidate completions (distractors) are chosen. In the easy task, the distractors are randomly chosen from three descriptions of the same question type from other images. In the hard task, 3 distractors are chosen only from these images that contain the same objects as the given question image, and hence it requires a more careful and detailed image understanding. We use ADAM gradient descent method [10] with default hyper-parameters.

**Different Non Maxima Suppression Thresholds.** Table 1 shows the accuracy of our CCA model with mean edge-box pooling for two different Non Maxima Suppression (NMS) thresholds  $\beta$ . In the experiments we use up to 100 object proposals, as we have observed saturation for higher numbers. From Table 1, we can see that for both tasks, easy and hard, higher NMS thresholds are preferred. More precisely, the threshold 0.75 outperforms 0.30 in average by 2.4 percentage points for the easy task, and by 1.8 percentage points for the hard task. We have also experimented with max pooling, but mean pooling has performed by about 0.5 percentage points better in average in all our experiments. The experiments, counterintuitively, suggest that many selected bounding boxes with high overlap are still beneficial in achieving better performance. Further experiments use mean pooling and the NMS threshold 0.75.

	Easy Task		Hard Task	
	0.30	0.75	0.30	0.75
NMS thresholds:	0.30	0.75	0.30	0.75
Image’s scenes	84.6	<b>86.2</b>	67.8	<b>69.0</b>
Image’s emotion	52.0	<b>52.5</b>	38.8	<b>39.4</b>
Image’s past	78.9	<b>80.8</b>	53.0	<b>54.6</b>
Image’s future	78.6	<b>81.1</b>	55.2	<b>56.1</b>
Image’s interesting	76.8	<b>78.2</b>	53.5	<b>54.2</b>
Object’s attribute	60.4	<b>62.4</b>	43.1	<b>45.7</b>
Object’s affordance	80.4	<b>83.3</b>	62.5	<b>63.6</b>
Object’s position	73.5	<b>77.4</b>	53.6	<b>56.3</b>
Person’s attribute	53.1	<b>56.0</b>	42.0	<b>44.2</b>
Person’s activity	79.9	<b>83.0</b>	63.2	<b>65.5</b>
Person’s location	82.4	<b>84.3</b>	63.8	<b>65.2</b>
Pair’s relationship	71.0	<b>75.3</b>	51.8	<b>55.7</b>
Average	72.6	<b>75.0</b>	54.0	<b>55.8</b>

Table 1: Accuracies computed for different Non Maxima Suppression thresholds (NMS) on the easy and hard tasks of the Visual Madlibs dataset. Mean pooling and 100 object proposals are used in the experiments. Results in %.

	Easy Task					Hard Task			
	10	25	50	100		10	25	50	100
Scenes	84.5	85.5	86.0	<b>86.2</b>	Scenes	68.0	68.6	68.9	<b>69.0</b>
Emotion	49.9	51.6	52.1	<b>52.5</b>	Emotion	37.9	38.1	38.8	<b>39.4</b>
Past	78.7	80.0	80.6	<b>80.8</b>	Past	52.8	53.9	54.3	<b>54.6</b>
Future	78.7	79.7	80.7	<b>81.1</b>	Future	54.4	55.0	55.8	<b>56.1</b>
Interesting	75.4	77.2	77.9	<b>78.2</b>	Interesting	51.9	53.6	53.7	<b>54.2</b>
Obj. attr.	59.0	60.9	61.7	<b>62.4</b>	Obj. attr.	43.7	44.0	44.9	<b>45.7</b>
Obj. aff.	81.2	82.4	83.0	<b>83.3</b>	Obj. aff.	62.4	63.0	63.4	<b>63.6</b>
Obj. pos.	75.4	76.6	77.4	<b>77.5</b>	Obj. pos.	55.1	55.5	56.3	<b>56.3</b>
Per. attr.	51.4	53.3	55.0	<b>56.0</b>	Per. attr.	41.6	42.2	43.0	<b>44.2</b>
Per. act.	80.7	82.2	82.9	<b>83.0</b>	Per. act.	63.7	64.7	65.3	<b>65.5</b>
Per. loc.	82.9	83.9	84.0	<b>84.3</b>	Per. loc.	64.2	64.8	64.8	<b>65.2</b>
Pair’s rel.	72.4	73.9	74.6	<b>75.3</b>	Pair’s rel.	53.6	54.5	54.9	<b>55.7</b>
Average	72.5	73.9	74.7	<b>75.0</b>	Average	54.1	54.8	55.3	<b>55.8</b>

Table 2: Accuracies computed for different number of Edge Box proposals on the easy and hard tasks of the Visual Madlibs dataset. The NMS threshold 0.75 and mean-pooling is used for all the experiments. Results in %.

**Different number of object proposals.** The maximal number of object proposals is the second factor of Edge Boxes that we study in this work. A larger number of proposals tend to cover a larger fraction of the input image. Moreover, the higher number together with the higher NMS threshold can assign proposals to both an object, and its parts, effectively forming a multi-scale and multi-parts object representation. Table 2 shows the accuracy of our model with different number of Edge Box proposals. The experiments suggest using a larger numbers of proposals, however the gain diminishes with the larger numbers.

**Comparison to the state-of-the-art.** Guided by the results of the previous experiments, we compare nCCA that uses Edge Boxes object proposals (nCCA (ours)) with the state-of-the-arts on Visual Madlibs (nCCA [32]). Both models use the same VGG Convolutional Neural Network [23] to encode images (or theirs crops), and word2vec to encode words. The models are trained per category (a model trained over all the categories performs inferior on the hard task [32]). As Table 3 shows using a large number of object proposals

	Easy Task		Hard Task	
	nCCA (ours)	nCCA [32]	nCCA (ours)	nCCA [32]
Scenes	86.2	<b>86.8</b>	69.0	<b>70.1</b>
Emotion	<b>52.5</b>	49.2	<b>39.4</b>	37.2
Past	<b>80.8</b>	77.5	<b>54.6</b>	52.8
Future	<b>81.1</b>	78.0	<b>56.1</b>	54.3
Interesting	<b>78.2</b>	76.5	<b>54.2</b>	53.7
Obj. attr.	<b>62.4</b>	47.5	<b>45.7</b>	43.6
Obj. aff.	<b>83.3</b>	73.0	<b>63.6</b>	63.5
Obj. pos.	<b>77.5</b>	65.9	<b>56.3</b>	55.7
Per. attr.	<b>56.0</b>	48.0	<b>44.2</b>	38.6
Per. act.	<b>83.0</b>	80.7	<b>65.5</b>	65.4
Per. loc.	<b>84.3</b>	82.7	<b>65.2</b>	63.3
Pair’s rel.	<b>75.3</b>	63.0	<b>55.7</b>	54.3
Average	<b>75.0</b>	69.1	<b>55.8</b>	54.4

Table 3: Accuracies computed for different approaches on the easy and hard tasks. nCCA (ours) uses the representation with object proposals (NMS 0.75, and 100 proposals with mean-pooling). nCCA uses the whole image representation. Results in %.

	Easy Task		Hard Task	
	nCCA (ours)	nCCA (bbox) [32]	nCCA (ours)	nCCA (bbox) [32]
Obj. attr.	<b>62.4</b>	54.7	45.7	<b>49.8</b>
Obj. aff.	<b>83.3</b>	72.2	<b>63.6</b>	63.0
Obj. pos.	<b>77.5</b>	58.9	<b>56.3</b>	50.7
Per. attr.	<b>56.0</b>	53.1	44.2	<b>46.1</b>
Per. act.	<b>83.0</b>	75.6	<b>65.5</b>	65.1
Per. loc.	<b>84.3</b>	73.8	<b>65.2</b>	57.8
Pair’s rel.	<b>75.3</b>	64.2	55.7	<b>56.5</b>
Average	<b>74.5</b>	64.6	<b>56.6</b>	55.6

Table 4: Accuracies computed for different approaches on the easy and hard task. nCCA (ours) uses the representation with object proposals (NMS 0.75, and 100 proposals with mean-pooling). nCCA(bbox) mean-pools over the representations computed on the available ground-truth bounding boxes both at train and test time. The averages are computed only over 7 categories. Results in %.

improves over global, full frame nCCA by 5.9 percentage points on the easy task, and about 1.4 percentage points on the difficult task in average. However, our nCCA also consistently outperforms state-of-the-art on every category except the ‘Scenes’ category. This suggests that better localized object oriented representation is beneficial. However, Edge Boxes only roughly localize objects. This naturally leads to the following question if better localization helps. To see the limits, we compare nCCA (ours) against nCCA (bbox) [32] that crops over ground truth bounding boxes from MS COCO segmentations and next averages over their representations (Table 3 in [32] shows that ground truth bounding boxes outperforms automatically detected bounding boxes, and hence they can be seen as an upper bound for a detection method trained to detect objects on MS COCO). Surprisingly, nCCA (ours) outperforms nCCA (bbox) by a large margin as Table 4 shows. Arguably, object proposals have better recall and captures multi-scale, multi-parts phenomena.

**CNN+LSTM with comparison in the output embedding space.** On one hand nCCA tops the leaderboard on the Visual Madlibs task [32]. On the other hand, the largest body of work on the question answering about images [2, 4, 14, 20] combines a CNN with an LSTM [4, 13, 17, 20, 29, 30, 31]. We hypothesize that, likewise to nCCA, in order to choose a



	Easy Task			Hard Task	
	Embedded CNN+LSTM (ours)	Ask Your Neurons [17]	CNN+LSTM [32]	Embedded CNN+LSTM (ours)	CNN+LSTM [32]
Scenes	<b>74.7</b>	70.6	71.1	<b>62.1</b>	60.5
Emotion	<b>36.2</b>	35.7	34.0	<b>34.3</b>	32.7
Past	<b>46.8</b>	44.9	35.8	<b>42.5</b>	32.0
Future	<b>48.1</b>	41.2	40.0	<b>41.4</b>	34.3
Interesting	<b>49.9</b>	49.1	39.8	<b>40.1</b>	33.3
Obj. attr.	<b>46.5</b>	45.5	45.4	<b>40.6</b>	40.3
Obj. aff.	<b>68.5</b>	64.3	-	<b>86.4</b>	-
Obj. pos.	<b>53.3</b>	49.5	50.9	<b>45.0</b>	44.9
Per. attr.	<b>40.7</b>	39.9	37.3	<b>40.0</b>	35.1
Per. act.	<b>64.1</b>	62.6	63.7	<b>53.7</b>	53.6
Per. loc.	<b>61.5</b>	59.1	59.2	<b>51.4</b>	49.3
Pair’s rel.	<b>66.2</b>	60.1	-	<b>54.5</b>	-
Average	<b>54.7</b>	51.9	47.7	<b>49.3</b>	41.7

Table 5: Comparison between our Embedded CNN+LSTM approach that computes the similarity between input and candidate answers in the embedding space, and the plain CNN+LSTM original approach from [32]. Since the accuracies of CNN+LSTM [32] are unavailable for two categories, we report average over 10 categories in this case. Results in %.

completion of the prompt sentence out of four candidates, the comparison between the candidate completions should be directly done in the output embedding space. This contrasts to a post-hoc process used in [32] where an image description architecture (CNN+LSTM) first generates a completion that is next compared against the candidates in the word2vec space (see section 3 for more details). Moreover, since the “Ask Your Neurons” architecture [17] is more suitable for the question answering task, we extend that method to do comparisons directly in the embedding space (“Embedded CNN+LSTM” in Table 5). Note that, here we feed the sentence prompt to LSTM even though it is fixed per category. Table 5 shows the performance of different methods. Our “Embedded CNN+LSTM” outperforms other methods on both tasks confirming our hypothesis. “Ask Your Neurons” [17] is also slightly better than the original CNN+LSTM [32] (on the 10 categories that the results for CNN+LSTM are available it achieves 49.8% accuracy on the easy task, which is 2.1 percentage points higher than CNN+LSTM).

## 5 Conclusion

We study an image representation formed by averaging over representations of object proposals, and show its effectiveness through experimental evaluation on the Visual Madlibs dataset [32]. We achieve state of the art performance on the multi-choice “filling the blank” task. We have also shown and discussed effects of different parameters that affect how the proposals are obtained. Surprisingly, the larger number of proposals the better overall performance. Moreover, the model benefits even from highly overlapping proposals. Such model even outperforms the prior work that uses ground truth bounding boxes from the MS COCO dataset. The proposed representation can be considered as a valid alternative to ‘soft’ attention representations such as implemented in recent work of visual question answering using memory networks [31]. Due to its popularity on question answering about images tasks, we also investigate a CNN+LSTM approach that chooses a prompt completion candidate by doing comparisons directly in the embedding space. This approach contrasts with a post-hoc solution of the previous work allowing for a tighter integration of the model with the

multi-choice task.

**Acknowledgements:** This work is supported by the German Research Foundation (DFG) under the SFB/CRC 1223.

## References

- [1] Zeynep Akata, Mateusz Malinowski, Mario Fritz, and Bernt Schiele. Multi-cue zero-shot learning with strong supervision. In *CVPR*, 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] Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *CVPR*, 2015.
- [4] Haoyuan Gao, Junhua Mao, Jie Zhou, Zhiheng Huang, Lei Wang, and Wei Xu. Are you talking to a machine? dataset and methods for multilingual image question answering. In *NIPS*, 2015.
- [5] Donald Geman, Stuart Geman, Neil Hallonquist, and Laurent Younes. Visual turing test for computer vision systems. *Proceedings of the National Academy of Sciences*, 2015.
- [6] Yunchao Gong, Qifa Ke, Michael Isard, and Svetlana Lazebnik. A multi-view embedding space for modeling internet images, tags, and their semantics. *IJCV*, 2014.
- [7] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 1997.
- [8] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015.
- [9] Andrej Karpathy, Armand Joulin, and Li Fei-Fei. Deep fragment embeddings for bidirectional image sentence mapping. In *NIPS*, 2014.
- [10] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [11] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012.
- [12] Percy Liang, Michael I Jordan, and Dan Klein. Learning dependency-based compositional semantics. *Computational Linguistics*, 2013.
- [13] Lin Ma, Zhengdong Lu, and Hang Li. Learning to answer questions from image using convolutional neural network. In *AAAI*, 2016.
- [14] Mateusz Malinowski and Mario Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *NIPS*, 2014.

- 
- [15] Mateusz Malinowski and Mario Fritz. Towards a visual turing challenge. In *Learning Semantics (NIPS workshop)*, 2014.
- [16] Mateusz Malinowski and Mario Fritz. Hard to cheat: A turing test based on answering questions about images. *AAAI Workshop: Beyond the Turing Test*, 2015.
- [17] Mateusz Malinowski, Marcus Rohrbach, and Mario Fritz. Ask your neurons: A deep learning approach to visual question answering. *arXiv preprint arXiv:1605.02697*, 2016.
- [18] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013.
- [19] Seong Joon Oh, Rodrigo Benenson, Mario Fritz, and Bernt Schiele. Person recognition in personal photo collections. In *ICCV*, 2015.
- [20] Mengye Ren, Ryan Kiros, and Richard Zemel. Image question answering: A visual semantic embedding model and a new dataset. In *NIPS*, 2015.
- [21] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. Imagenet large scale visual recognition challenge. *arXiv:1409.0575*, 2014.
- [22] Kevin J Shih, Saurabh Singh, and Derek Hoiem. Where to look: Focus regions for visual question answering. In *CVPR*, 2016.
- [23] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.
- [24] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *CVPR*, 2015.
- [25] Alan M Turing. Computing machinery and intelligence. *Mind*, pages 433–460, 1950.
- [26] Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. Cider: Consensus-based image description evaluation. 2015.
- [27] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015.
- [28] Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. In *ICLR*, 2015.
- [29] Caiming Xiong, Stephen Merity, and Richard Socher. Dynamic memory networks for visual and textual question answering. *arXiv:1603.01417*, 2016.
- [30] Huijuan Xu and Kate Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. *arXiv:1511.05234*, 2015.
- [31] Zichao Yang, Xiaodong He, Jianfeng Gao, Li Deng, and Alex Smola. Stacked attention networks for image question answering. In *CVPR*, 2016.

- [32] Licheng Yu, Eunbyung Park, Alexander C. Berg, and Tamara L. Berg. Visual madlibs: Fill in the blank image generation and question answering. In *ICCV*, 2015.
- [33] Ning Zhang, Ryan Farrell, and Trever Darrell. Pose pooling kernels for sub-category recognition. In *CVPR*, 2012.
- [34] C Lawrence Zitnick and Piotr Dollár. Edge boxes: Locating object proposals from edges. In *ECCV*. 2014.