

Understanding Image and Text Simultaneously: a Dual Vision-Language Machine Comprehension Task

Nan Ding
Google

dingnan@google.com

Sebastian Goodman
Google

seabass@google.com

Fei Sha
Google

fsha@google.com

Radu Soricut
Google

rsoricut@google.com

Abstract

We introduce a new multi-modal task for computer systems, posed as a combined vision-language comprehension challenge: identifying the most suitable text describing a scene, given several similar options. Accomplishing the task entails demonstrating comprehension beyond just recognizing “keywords” (or key-phrases) and their corresponding visual concepts. Instead, it requires an alignment between the representations of the two modalities that achieves a visually-grounded “understanding” of various linguistic elements and their dependencies. This new task also admits an easy-to-compute and well-studied metric: the accuracy in detecting the true target among the decoys.

The paper makes several contributions: an effective and extensible mechanism for generating decoys from (human-created) image captions; an instance of applying this mechanism, yielding a large-scale machine comprehension dataset (based on the COCO images and captions) that we make publicly available; human evaluation results on this dataset, informing a performance upper-bound; and several baseline and competitive learning approaches that illustrate the utility of the proposed task and dataset in advancing both image and language comprehension. We also show that, in a multi-task learning setting, the performance on the proposed task is positively correlated with the end-to-end task of image captioning.

1. Introduction

There has been a great deal of interest in multi-modal artificial intelligence research recently, bringing together the fields of Computer Vision and Natural Language Processing. This interest has been fueled in part by the availability of many large-scale image datasets with textual annotations. Several vision+language tasks have been proposed around these datasets [15, 16, 21, 3]. Image Captioning [15, 10, 16, 11, 17, 33, 25, 37] and Visual Question Answering [23, 24, 31, 3, 40, 35, 28, 13, 38, 41, 22] have

in particular attracted a lot of attention. The performances on these tasks have been steadily improving, owing much to the wide use of deep learning architectures [6].

A central theme underlying these efforts is the use of natural language to identify how much visual information is perceived and understood by a computer system. Presumably, a system that understands a visual scene well enough ought to be able to describe what the scene is about (thus “captioning”) or provide correct and visually-grounded answers when queried (thus “question-answering”).

In this paper, we argue for directly measuring how well the semantic representations of the visual and linguistic modalities align (in some abstract semantic space). For instance, given an image and two captions – a correct one and an incorrect yet-cunningly-similar one – can we both qualitatively and quantitatively measure the extent to which humans can dismiss the incorrect one but computer systems blunder? Arguably, the degree of the modal alignment is a strong indicator of task-specific performance on any vision+language task. Consequentially, computer systems that can learn to maximize and exploit such alignment should outperform those that do not.

We take a two-pronged approach for addressing this issue. First, we introduce a new and challenging Dual Machine Comprehension (DMC) task, in which a computer system must identify the most suitable textual description from several options: one being the target and the others being “adversarialy”-chosen decoys. All options are free-form, coherent, and fluent sentences with *high degrees of semantic similarity* (hence, they are “cunningly similar”). A successful computer system has to demonstrate comprehension beyond just recognizing “keywords” (or key phrases) and their corresponding visual concepts; they must arrive at a coinciding and visually-grounded understanding of various linguistic elements and their dependencies. What makes the DMC task even more appealing is that it admits an easy-to-compute and well-studied performance metric: the accuracy in detecting the true target among the decoys. Second, we illustrate how solving the DMC task benefits related vision+language tasks. To this end, we render the

DMC task as a classification problem, and incorporate it in a multi-task learning framework for end-to-end training of joint objectives.

Our work makes the following contributions: (1) an effective and extensible algorithm for generating decoys from human-created image captions (Section 3.2); (2) an instantiation of applying this algorithm to the COCO dataset [21], resulting in a large-scale dual machine-comprehension dataset that we make publicly available (Section 3.3); (3) a human evaluation on this dataset, which provides an upper-bound on performance (Section 3.4); (4) a benchmark study of baseline and competitive learning approaches (Section 5), which underperform humans by a substantial gap (about 20% absolute); and (5) a novel multi-task learning model that simultaneously learns to solve the DMC task and the Image Captioning task (Sections 4.3 and 5.4).

Our empirical study shows that performance on the DMC task positively correlates with performance on the Image Captioning task. Therefore, besides acting as a standalone benchmark, the new DMC task can be useful in improving other complex vision+language tasks. Both suggest the DMC task as a fruitful direction for future research.

2. Related work

Image understanding is a long-standing challenge in computer vision. There has recently been a great deal of interest in bringing together vision and language understanding. Particularly relevant to our work are image captioning (IC) and visual question-answering (VQA). Both have instigated a large body of publications, a detailed exposition of which is beyond the scope of this paper. Interested readers should refer to two recent surveys [7, 34].

In IC tasks, systems attempt to generate a fluent and correct sentence describing an input image. IC systems are usually evaluated on how well the generated descriptions align with human-created captions (ground-truth). The language generation model of an IC system plays a crucial role; it is often trained such that the probabilities of the ground-truth captions are maximized (MLE training), though more advanced methods based on techniques borrowed from Reinforcement Learning have been proposed [27]. To provide visual grounding, image features are extracted and injected into the language model. Note that language generation models need to both decipher the information encoded in the visual features, and model natural language generation.

In VQA tasks, the aim is to answer an input question correctly with respect to a given input image. In many variations of this task, answers are limited to single words or a binary response (“yes” or “no”) [3]. The Visual7W dataset [41] contains answers in a richer format such as phrases, but limits questions to “wh-” style (what, where, who, etc). The Visual Genome dataset [18], on the other hand, can potentially define more complex questions and

answers due to its extensive textual annotations.

Our DMC task is related but significantly different. In our task, systems attempt to discriminate the best caption for an input image from a set of captions — all but one are decoys. Arguably, it is a form of VQA task, where the same default (thus uninformative) question is asked: *Which of the following sentences best describes this image?* However, unlike current VQA tasks, choosing the correct answer in our task entails a deeper “understanding” of the available answers. Thus, to perform well, a computer system needs to understand both complex scenes (visual understanding) and complex sentences (language understanding), and be able to reconcile them.

The DMC task admits a simple classification-based evaluation metric: the accuracy of selecting the true target. This is a clear advantage over the IC tasks, which often rely on imperfect metrics such as BLEU [26], ROUGE [20], METEOR [5], CIDEr [32], or SPICE [2].

Related to our proposal is the work in [15], which frames image captioning as a ranking problem. While both share the idea of selecting captions from a large set, our framework has some important and distinctive components. First, we devise an algorithm for smart selection of candidate decoys, with the goal of selecting those that are sufficiently similar to the true targets to be challenging, and yet still be reliably identifiable by human raters. Second, we have conducted a thorough human evaluation in order to establish a performance ceiling, while also quantifying the level to which current learning systems underperform. Lastly, we show that there exists a positive correlation between the performance on the DMC task and the performance on related vision+language tasks by proposing and experimenting with a multi-task learning model. Our work is also substantially different from their more recent work [14], where only one decoy is considered and its generation is either random, or focusing on visual concept similarity (“switching people or scenes”) instead of our focus on both linguistic surface and paragraph vector embedding similarity.

3. The Dual Machine Comprehension Task

3.1. Design overview

We propose a new multi-modal machine comprehension task to examine how well visual and textual semantic understanding are aligned. Given an image, human evaluators or machines must accurately identify the best sentence describing the scene from several decoy sentences. Accuracy on this task is defined as the percentage that the true targets are identified.

It seems straightforward to construct a dataset for this task, as there are several existing datasets which are composed of images and their (multiple) ground-truth captions, including the popular COCO dataset [21]. Thus, for any

given image, it appears that one just needs to use the captions corresponding to other images as decoys. However, this naïve approach could be overly simplistic as it provides no control over the properties of the decoys.

Specifically, our desideratum is to recruit *challenging* decoys that are sufficiently similar to the targets. However, for a small number of decoys, e.g. 4-5, randomly selected captions could be significantly different from the target. The resulting dataset would be too “easy” to shed any insight on the task. Since we are also interested in human performance on this task, it is thus impractical to increase the number of decoys to raise the difficulty level of the task at the expense of demanding humans to examine tediously and unreliably a large number of decoys. In short, we need an *automatic procedure to reliably create difficult sets of decoy captions* that are sufficiently similar to the targets.

We describe such a procedure in the following. While it focuses on identifying decoy captions, the main idea is potentially adaptable to other settings. The algorithm is flexible in that the “difficulty” of the dataset can be controlled to some extent through the algorithm’s parameters.

3.2. Algorithm to create an MC-IC dataset

The main idea behind our algorithm is to carefully define a “good decoy”. The algorithm exploits recent advances in paragraph vector (PV) models [19], while also using linguistic surface analysis to define similarity between two sentences. Due to space limits, we omit a detailed introduction of the PV model. It suffices to note that the model outputs a continuously-valued embedding for a sentence, a paragraph, or even a document.

The pseudo-code is given in Algorithm 1 (the name MC-IC stands for “Machine-Comprehension for Image & Captions”). As input, the algorithm takes a set C of $\langle \text{image}, \{\text{caption}(s)\} \rangle$ pairs¹, as those extracted from a variety of publicly-available corpora, including the COCO dataset [21]. The output of the algorithm is the MC_{IC} set.

Concretely, the MC-IC Algorithm has three main arguments: a dataset $C = \{\langle \mathbf{i}_i, \mathbf{c}_i \rangle | 1 \leq i \leq m\}$ where \mathbf{i}_i is an image and \mathbf{c}_i is its ground-truth caption²; an integer N which controls the size of \mathbf{c}_i ’s neighborhood in the embedding space defined by the paragraph vector model PV; and a function Score which is used to score the N items in each such neighborhood.

The first two steps of the algorithm tune several hyperparameters. The first step finds optimal settings for the PV model given the dataset C . The second finds a weight parameter λ given PV, dataset C , and the Score function.

¹On the order of at least hundreds of thousands of examples; smaller sets result in less challenging datasets.

²For an image with multiple ground-truth captions, we split it to multiple instances with the same image for each one of the ground-truth captions; the train/dev/test splits are done such that they contain disjoint *image* sets, as opposed to disjoint *instance* sets.

Algorithm 1: MC-IC(C, N, Score)

```

Result: Dataset  $\text{MC}_{\text{IC}}$ 
 $\text{PV} \leftarrow \text{OPTIMIZE-PV}(C)$ 
 $\lambda \leftarrow \text{OPTIMIZE-SCORE}(\text{PV}, C, \text{Score})$ 
 $\text{MC}_{\text{IC}} \leftarrow \emptyset$ 
 $nr\_decoys = 4$ 
for  $\langle \mathbf{i}_i, \mathbf{c}_i \rangle \in C$  do
   $A \leftarrow []$ 
   $T_{\mathbf{c}_i} \leftarrow \text{PV}(\mathbf{c}_i)[1..N]$ 
  for  $\mathbf{c}_d \in T_{\mathbf{c}_i}$  do
     $score \leftarrow \text{Score}(\text{PV}, \lambda, \mathbf{c}_d, \mathbf{c}_i)$ 
    if  $score > 0$  then
       $A.append(\langle score, \mathbf{c}_d \rangle)$ 
    end
  end
  if  $|A| \geq nr\_decoys$  then
     $R \leftarrow \text{descending-sort}(A)$ 
    for  $l \in [1..nr\_decoys]$  do
       $\langle score, \mathbf{c}_d \rangle \leftarrow R[l]$ 
       $\text{MC}_{\text{IC}} \leftarrow \text{MC}_{\text{IC}} \cup \{(\langle \mathbf{i}_i, \mathbf{c}_d \rangle, \text{false})\}$ 
    end
   $\text{MC}_{\text{IC}} \leftarrow \text{MC}_{\text{IC}} \cup \{(\langle \mathbf{i}_i, \mathbf{c}_i \rangle, \text{true})\}$ 
end
end

```

These hyperparameters are dataset-specific. Details are discussed in the next section.

The main body of the algorithm, the outer **for** loop, generates a set of nr_decoys (4 here) decoys for each ground-truth caption. It accomplishes this by first extracting N candidates from the PV neighborhood of the ground-truth caption, excluding those that belong to the same image. In the inner **for** loop, it computes the similarity of each candidate to the ground-truth and stores them in a list A . If enough candidates are generated, the list is sorted in descending order of score. The top nr_decoys captions are marked as “decoys” (*i.e.* **false**), while the ground-truth caption is marked as “target” (*i.e.* **true**).

The score function $\text{Score}(\text{PV}, \lambda, \mathbf{c}', \mathbf{c})$ is a crucial component of the decoy selection mechanism. Its definition leverages our linguistic intuition by combining linguistic surface similarity, $\text{sim}_{\text{SURF}}(\mathbf{c}', \mathbf{c})$, with the similarity suggested by the embedding model, $\text{sim}_{\text{PV}}(\mathbf{c}', \mathbf{c})$:

$$\text{Score} = \begin{cases} 0 & \text{if } \text{sim}_{\text{SURF}} \geq L \\ \lambda \text{sim}_{\text{PV}} + (1 - \lambda) \text{sim}_{\text{SURF}} & \text{otherwise} \end{cases} \quad (1)$$

where the common argument $(\mathbf{c}', \mathbf{c})$ is omitted. The higher the similarity score, the more likely that \mathbf{c}' is a good decoy for \mathbf{c} . Note that if the surface similarity is above the threshold L , the function returns 0, flagging that the two captions are too similar to be used as a pair of target and decoy.

Split	dev	test	train	total
#unique_images	2,000	2,000	110,800	114,800
# instances	9,999	10,253	554,063	574,315

Table 1. MCIC-COCO dataset descriptive statistics

In this work, sim_{SURF} is computed as the BLEU score between the inputs [26] (with the brevity penalty set to 1). The embedding similarity, sim_{PV} , is computed as the cosine similarity between the two in the PV embedding space.

3.3. The MCIC-COCO dataset

We applied the MC-IC Algorithm to the COCO dataset [21] to generate a dataset for the visual-language dual machine comprehension task. The dataset is called MCIC-COCO and it is made publicly available³. We describe the details of this dataset below.

We set the neighborhood size at $N = 500$, and the threshold at $L = 0.5$ (see Eq. 1). As the COCO dataset has a large body of images (thus captions) focusing on a few categories (such as sports activities), this threshold is important in discarding significantly similar captions to be decoys – otherwise, even human annotators will experience difficulty in selecting the ground-truth captions.

The hyperparameters of the PV model, dim (embedding dimension) and epochs (number of training epochs), are optimized in the OPTIMIZE-PV step of the MC-IC Algorithm. The main idea is to learn embeddings such that ground-truth captions from the same image have similar embeddings. Concretely, the optimization step is a grid-search over the hyper-parameters of the PV-DBOW model [19], which we train using a softmax loss. Since there are multiple ground-truth captions associated with each image, the dataset is denoted by $C = \{(\mathbf{i}_{r_c}, \mathbf{c}_{r_c}) | 1 \leq r \leq n, 1 \leq c \leq s_r\}$, where r is the index for each unique image ($\mathbf{i}_{r_c} \equiv \mathbf{i}_r$), n is the total number images and $s_r > 1$ is the number of unique captions for image r . The total number of data examples $m = \sum_{r=1}^n s_r$. Here the hyper-parameters are searched on a grid to minimize “multiple ground-truth score” rank (mgs-rank): the average rank (under the cosine-distance score) between \mathbf{c}_{r_c} and $\{\mathbf{c}_{r_l} | 1 \leq l \leq s_r, l \neq c\}$. The lower the mgs-rank, the better the resulting paragraph vector model is at modeling multiple ground-truths for a given image as being similar. As such, our grid-search over the MCIC-COCO dev dataset yields a minimum mgs-rank at $\text{dim}=1024$ and $\text{epochs}=5$.

Similarly, the OPTIMIZE-SCORE(PV, Score) step is a grid-search over the λ parameter of the Score function, given a paragraph vector embedding model PV and a dataset C of captions and images, as before. A well-chosen λ will ensure the multiple ground-truth captions for the

Correct responses	# instances	Accuracy%
3 out of 3	673	67.3
at least 2 out of 3	828	82.8
at least 1 out of 3	931	93.1
0 out of 3	69	0.0

Table 2. Human performance on the DMC task with the MCIC-COCO dataset. **Bold** denotes the performance ceiling.

same image will be measured with high degree of similarity with the Score function. The $\lambda \in [0, 1]$ parameter is searched on a grid to minimize the “weighted multiple ground-truths score” rank (wmgs-rank): the average rank (under the Score) between \mathbf{c}_{r_c} and $\{\mathbf{c}_{r_l} | 1 \leq l \leq s_r, l \neq c\}$, relative to the top N -best closest-cosine neighbors in PV. For example, if given five ground-truths for image \mathbf{i}_r , and when considering \mathbf{c}_{r_1} , ground-truths \mathbf{c}_{r_2} to \mathbf{c}_{r_5} are ranking at #4, #10, #16, and #22 (in top-500 closest-cosine neighbors in PV), then $\text{wmgs-rank}(\mathbf{c}_{r_1}) = 13$ (the average of these ranks). Our grid-search over the MCIC-COCO dev dataset yields a minimum wmgs-rank at $\lambda=0.3$.

The resulting MCIC-COCO dataset has 574,315 instances that are in the format of $\{i : (\langle \mathbf{i}_i, \mathbf{c}_i^j \rangle, \text{label}_i^j), j = 1 \dots 5\}$ where $\text{label}_i^j \in \{\text{true}, \text{false}\}$. For each such instance, there is one and only one j such that the label is **true**. We have created a train/dev/test split such that all of the instances for the same image occur in the same split. Table 1 reports the basic statistics for the dataset.

3.4. Human performance on MCIC-COCO

Setup To measure how well humans can perform on the DMC task, we randomly drew 1,000 instances from the MCIC-COCO dev set and submitted those instances to human “raters”⁴ via a crowd-sourcing platform.

Three independent responses from 3 different rates were gathered for each instance, for a total of 3,000 responses. To ensure diversity, raters were prohibited from evaluating more than six instances or from responding to the same task instance twice. In total, 807 distinct raters were employed.

Raters were shown one instance at a time. They were shown the image and the five caption choices (ground-truth and four decoys, in randomized order) and were instructed to choose the best caption for the image. Before starting evaluation, the raters were trained with sample instances from the *train* dataset, disjoint from the *dev* dataset on which their performance data were collected. The training process presents an image and five sentences, of which the ground-truth caption is highlighted. In addition, specific instructions and clarification were given to the raters on how to choose the best caption for the image. In Figure 1, we present three instances on how the rater instructions were

³<http://www.github.com/google/mcic-coco>

⁴Raters are vetted, screened and tested before working on any tasks; requirements include native-language proficiency level.



1. a herd of giraffe standing next to each other in a dirt field
2. a pack of elephants standing next to each other
3. animals are gathering next to each other in a dirt field
4. three giraffe standing next to each other on a grass field
5. **two elephants standing next to each other in a grass field**

Instructions:

Captions 1 and 4 are clearly incorrect. They do not match up with the image at all.

Caption 2 calls the elephants a "pack", which is vague and a bit subjective. It does not mention the grass at all.

Caption 3 only uses the word "animals" to describe what is in the picture, when the picture clearly shows elephants.

Caption 5 gives the exact count and correct animal type and even mentions the grass field. It is a more accurate, descriptive, and objective caption than the other options.



1. a meal covered with a lot of broccoli and tomatoes
2. a pan filled with a mixture of vegetables and meat
3. **a piece of bread covered in a meat and sauces**
4. a pizza smothered in cheese and meat with french fries
5. a plate of fries and a sandwich cut in half

Instructions:

Caption 3 does not mention the fork, the plate, or the eggs, but it is still the best option because the other captions are inaccurate. Captions 1, 2, 4, and 5 all describe items not present in the picture, such as broccoli, a pan, cheese, or fries.



1. the man in the picture is reaching toward a frisbee
2. a middle aged man in a field tossing a frisbee
3. a woman in stance to throw a frisbee
4. **a man dives for a catch in this ultimate frisbee match**
5. there is a male tennis player playing in a match

Instructions:

Caption 5 is clearly incorrect (the game being played here is not tennis). Captions 2 and 3 describe the act of tossing, but the picture shows the act of catching, so these captions are both inaccurate. Captions 1 and 4 seem close, but the phrase "dives for a catch" is more descriptive than the phrase "reaching toward a frisbee". In addition, Caption 4 mentions the name of the game that they are playing, so it better informs the reader about what is happening in the rest of the image than the other caption does.

Figure 1. Examples of instances from the MCIC-COCO dataset (the ground-truth is in **bold** face), together with rater instructions.

presented for rater training.

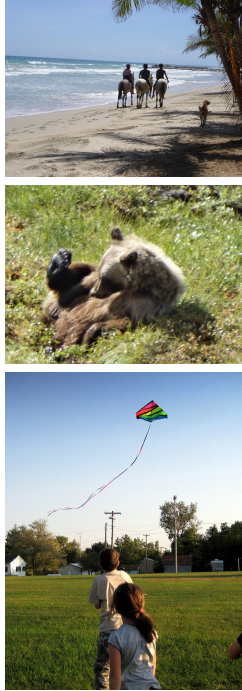
Quantitative results We assessed human performance in two metrics: (1) Percentage of correct rater responses (1-human system): **81.1%** (2432 out of 3000); (2) Percentage of instances with at least 50% (*i.e.* 2) correct responses (3-human system): **82.8%** (828 out of 1000).

Table 2 gives a detailed breakdown on the statistics related to the inter-rater (dis)agreement. The first row, with accuracy at 67.3%, suggests that this is the level at which the correct answer is obvious (*i.e.*, percentage of "easy" instances). The second row, at 82.8%, indicates that this is the performance ceiling in terms of accuracy that can be expected for the MCIC-COCO dataset; at the same time, it suggests that the difference between 67.3% and 82.8% (*i.e.*, about 15% of instances) is caused by "difficult" instances. Finally, the third row, at 93.1%, indicates that the level of

"unanswerable" instances is somewhere in the 10%-15% range (combining the increase from 82.8% to 93.1% and the remaining 6.9% that no one gets right).

We will investigate those instances in detail in the future. The COCO dataset has a significant number of captions that fit more than one image in the dataset, given the biased concentration on certain categories. Thus, we suspect that even with our threshold-check (*cf.* the introduction of L in Eq. 1), our procedure might have failed to filter out some impossible-to-distinguish decoys.

Qualitative examples We present in Figure 2 several example instances from the MCIC-COCO dataset. The first example illustrates how certain aspects of VQA are subsumed by the DMC task: in order to correctly choose answer 3, a system needs to implicitly answer questions like "how many people are in the image?" (answer: three, thus choices 4.



1. three bikes on the shore while people talk on a small boat
 2. three people and a dog are running on the beach
 3. **three people ride horses along the beach while a dog follows**
 4. two people on horseback ride along a beach
 5. two people on horses trot along a sandy beach
-
1. a brown bear is standing in the grass
 2. a brown bear standing in the water of a river
 3. a dark brown bear standing in the woods
 4. **a small brown bear rolling in the grass**
 5. a small brown bear standing in the dirt
-
1. a family is playing the wii in a house
 2. a man playing with a frisbee in a park
 3. a small boy playing with kites in a field
 4. three women play with frisbees in a shady park
 5. **boy flies a kite with family in the park**

Figure 2. Examples of instances from the MCIC-COCO dataset (the ground-truth is in **bold** face). The correct answers for the first two examples are relatively obvious to humans, but less so to computer systems. The third example illustrates one of the difficult cases in which the humans annotators did not agree (option 3. was also chosen).

and 5. are wrong), and “what are the people doing?” (answer: riding horses, thus choices 1. and 2. are wrong). The second example illustrates the extent to which a successful computer system needs to be able to differentiate between “standing” and “rolling” in a visually-grounded way, presumably via a pose model [39] combined with a translation model between poses and their verbal correspondents. Last but not least, the third examples illustrates a difficult case, which led to human annotator disagreement in our annotation process (both choice 3. and 5. were selected by different annotators).

4. Learning Methods

We describe several learning methods for the dual machine comprehension (DMC) task with the MC_{IC} dataset. We start with linear models which will be used as baselines. We then present several neural-network based models. In particular, we describe a novel, hybrid neural network model that combines the feedforward architecture and the seq2seq architecture [29] for multi-task learning of the DMC task and the image captioning task. This new model achieves the best performance in both tasks.

4.1. Linear models as baselines

Regression To examine how well the two embeddings are aligned in “semantic understanding space”, a simple ap-

proach is to assume that the learners do not have access to the decoys. Instead, by accessing the ground-truth captions only, the models learn a linear regressor from the image embeddings to the target captions’ embeddings (“forward regression”), or from the captions to the images (“backward regression”). With the former approach, referred as **Baseline-l2C**, we check whether the predicted caption for any given image is closest to its true caption. With the latter, referred as **Baseline-C2I**, we check whether the predicted image embedding by the ground-truth caption is the closest among predicted ones by decoy captions to the real image embeddings.

Linear classifier Our next approach **Baseline-LinM** is a linear classifier learned to discriminate true targets from the decoys. Specifically, we learn a linear discriminant function $f(\mathbf{i}, \mathbf{c}; \Theta) = \mathbf{i}^\top \Theta \mathbf{c}$ where Θ is a matrix measuring the compatibility between two types of embeddings, cf. [12]. The loss function is then given by

$$L(\Theta) = \sum_i [\max_{j \neq j^*} f(\mathbf{i}_i, \mathbf{c}_i^j; \Theta) - f(\mathbf{i}_i, \mathbf{c}_i^{j^*}; \Theta)]_+ \quad (2)$$

where $[\]_+$ is the hinge function and j indexes over all the available decoys and i indexes over all training instances. The optimization tries to increase the gap between the target $\mathbf{c}_i^{j^*}$ and the worst “offending” decoy. We use stochastic (sub)gradient methods to optimize Θ , and se-

lect the best model in terms of accuracy on the MCIC-COCO development set.

4.2. Feedforward Neural Network (FFNN) models

To present our neural-network-based models, we use the following notations. Each training instance pair is a tuple $\langle \mathbf{i}_i, \mathbf{c}_i^j \rangle$, where \mathbf{i}_i denotes the image, and \mathbf{c}_i^j denotes the caption options, which can either be the target or the decoys. We use a binary variable $y_{ijk} \in \{0, 1\}$ to denote whether j -th caption of the instance i is labeled as k , and $\sum_k y_{ijk} = 1$.

We first employ the standard feedforward neural-network models to solve the DMC task on the MCIC-COCO dataset. For each instance pair $\langle \mathbf{i}_i, \mathbf{c}_i^j \rangle$, the input to the neural network is an embedding tuple $\langle \text{DNN}(\mathbf{i}_i; \Gamma), \text{Emb}(\mathbf{c}_i^j; \Omega) \rangle$, where Γ denotes the parameters of a deep convolutional neural network DNN. DNN takes an image and outputs an image embedding vector. Ω is the embedding matrix, and $\text{Emb}(\cdot)$ denotes the mapping from a list of word IDs to a list of embedding vectors using Ω . The loss function for our FFNN is given by:

$$L(\Gamma, \Omega, \mathbf{u}) = \sum_{i,j,k} y_{ijk} \log \text{FN}_k(\text{DNN}(\mathbf{i}_i; \Gamma), \text{Emb}(\mathbf{c}_i^j; \Omega); \mathbf{u}) \quad (3)$$

where FN_k denotes the k -th output of a feedforward neural network, and $\sum_k \text{FN}_k(\cdot) = 1$. Our architecture uses a two hidden-layer fully connected network with Rectified Linear hidden units, and a softmax layer on top.

The formula in Eq. 3 is generic with respect to the number of classes. In particular, we consider a 2-class-classifier ($k \in \{0, 1\}$, 1 for ‘yes’, this is a correct answer; 0 for ‘no’, this is an incorrect answer), applied independently on all the $\langle \mathbf{i}_i, \mathbf{c}_i^j \rangle$ pairs and apply one FFNN-based binary classifier for each; the final prediction is the caption with the highest ‘yes’ probability among all instance pairs belonging to instance i .

4.3. Vec2seq + FFNN Model

We describe here a hybrid neural-network model that combines a recurrent neural-network with a feedforward one. We encode the image into a single-cell RNN encoder, and the caption into an RNN decoder. Because the first sequence only contains one cell, we call this model a vector-to-sequence (Vec2seq) model as a special case of Seq2seq model as in [29, 4]. The output of each unit cell of a Vec2seq model (both on the encoding side and the decoding side) can be fed into an FFNN architecture for binary classification. See Figure 3 for an illustration of the Vec2seq + FFNN model architecture.

Multi-task learning In addition to the classification loss (Eq. 3), we also include a loss for generating an output sequence \mathbf{c}_i^j based on an input \mathbf{i}_i image. We define a binary

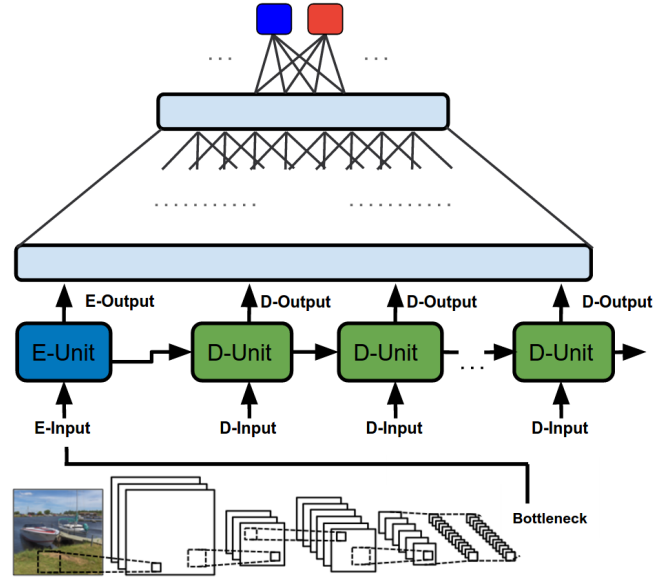


Figure 3. Vec2seq + FFNN model architecture.

variable $z_{ijlv} \in \{0, 1\}$ to indicate whether the l th word of \mathbf{c}_i^j is equal to word v . \mathbf{O}_{ijl}^d denotes the l -th output of the decoder of instance pair $\langle \mathbf{i}_i, \mathbf{c}_i^j \rangle$, \mathbf{O}_{ij}^e denotes the output of the encoder, and \mathbf{O}_{ij}^d denotes the concatenation of decoder outputs.

With these definitions, the loss function for the Vec2seq + FFNN model is:

$$\begin{aligned} L(\Theta, \mathbf{w}, \mathbf{u}) &= \sum_{i,j,k} y_{ijk} \log \text{FN}_k(\mathbf{O}_{ij}^e(\mathbf{i}_i, \mathbf{c}_i^j; \Theta), \mathbf{O}_{ij}^d(\mathbf{i}_i, \mathbf{c}_i^j; \Theta); \mathbf{u}) \\ &\quad + \lambda_{gen} \sum_{i,j,l,v} y_{ij1} z_{ijlv} \log \text{softmax}_v(\mathbf{O}_{ijl}^d(\mathbf{i}_i, \mathbf{c}_i^j; \Theta); \mathbf{w}) \end{aligned} \quad (4)$$

where $\sum_v \text{softmax}_v(\cdot) = 1$; Θ are the parameters of the Vec2seq model, which include the parameters within each unit cell, as well as the elements in the embedding matrices for images and target sequences; \mathbf{w} are the output projection parameters that transform the output space of the decoder to the vocabulary space. \mathbf{u} are the parameters of the FFNN model (Eq. 3); λ_{gen} is the weight assigned to the sequence-to-sequence generation loss. Only the true target candidates (the ones with $y_{ij1} = 1$) are included in this loss, as we do not want the decoy target options to affect this computation.

The Vec2seq model we use here is an instantiation of the attention-enhanced models proposed in [4, 8]. However, our current model does not support location-wise attention, as in the Show-Attend-and-Tell [36] model. In this sense, our model is an extension of the Show-and-Tell model with a single attention state representing the entire image, used as image memory representation for all decoder decisions.

We apply Gated Recurrent Unit (GRU) as the unit cell [9]. We also compare the influence on performance of the λ_{gen} parameter.

5. Experiments

5.1. Experimental Setup

Baseline models For the baseline models, we use the 2048-dimensional outputs of Google-Inception-v3 [30] (pre-trained on ImageNet ILSVRC 2012) to represent the images, and 1024-dimensional paragraph-vector embeddings (section 3.2) to represent captions. To reduce computation time, both are reduced to 256-dimensional vectors using random projections.

Neural-nets based models The experiments with these models are done using the Tensorflow package [1]. The hyper-parameter choices are decided using the hold-out development portion of the MCIC-COCO set. For modeling the input tokens, we use a vocabulary size of 8,855 types, selected as the most frequent tokens over the captions from the COCO training set (words occurring at least 5 times). The models are optimized using ADAGRAD with an initial learning rate of 0.01, and clipped gradients (maximum norm 4). We run the training procedures for 3,000,000 steps, with a mini-batch size of 20. We use 40 workers for computing the updates, and 10 parameter servers for model storing and (asynchronous and distributed) updating.

We use the following notations to refer to the neural network models: $\text{FFNN}_{2\text{-class}}^{\text{argmax } 1..5}$ refers to the version of feedforward neural network architecture with a 2-class-classifier ('yes' or 'no' for answer correctness), over which an argmax function computes a 5-way decision (i.e., the choice with the highest 'yes' probability); we henceforth refer to this model simply as FFNN.

The Vec2seq+FFNN refers to the hybrid model described in Section 4.3, combining Vec2seq and $\text{FFNN}_{2\text{-class}}^{\text{argmax } 1..5}$. The RNN part of the model uses a two-hidden-layer GRU unit-cell [9] configuration, while the FFNN part uses a two-hidden-layer architecture. The λ_{gen} hyper-parameter from the loss-function $L(\Theta, \mathbf{w}, \mathbf{u})$ (Eq. 4) is by default set to 1.0 (except for Section 5.4 where we directly measure its effect on performance).

Evaluation metrics The metrics we use to measure performance come in two flavors. First, the accuracy in detecting (the index of) the true target among the decoys provides a direct way of measuring the performance level on the comprehension task. We use this metric as the main indicator of comprehension performance. Second, because our Vec2seq+FFNN models are multi-task models, they can also generate new captions given the input image. The performance level for the generation task is measured using the standard scripts measuring ROUGE-L [20] and CIDEr [32], using as reference the available captions from the COCO

Model	dim	Dev	Test
Baseline-l2C	256	19.6 ±0.4	19.3±0.4
Baseline-C2I	256	32.8 ±0.5	32.0±0.5
Baseline-LinM	256	44.6 ±0.5	44.5±0.5
FFNN	256	56.3 ±0.5	55.1±0.5
Vec2seq+FFNN	256	60.5 ±0.5	59.0±0.5

Table 3. Performance on the DMC Task, in accuracies (and standard deviations) on MCIC-COCO for baselines and NN models.

data (around 5 for most of the images). Code for these metrics is available as part of the COCO evaluation toolkit⁵. As usual, both the hypothesis strings and the reference strings are preprocessed: remove all the non-alphabetic characters; transform all letters to lowercase, and tokenize using white space; replace all words occurring less than 5 times with an unknown token ⟨UNK⟩ (total vocabulary of 8,855 types); truncate to the first 30 tokens.

5.2. Results

Table 3 summarizes our main results on the comprehension task. We report the accuracies (and their standard deviations) for random choice, baselines, and neural network-based models.

Interestingly, the Baseline-l2C model performs at the level of random choice, and much worse than the Baseline-C2I model. This discrepancy reflects the inherent difficulty in vision-Language tasks: for each image, there are several possible equally good descriptions, thus a linear mapping from the image embeddings to the captions might not be enough – statistically, the *linear* model will just predict the mean of those captions. However, for the reverse direction where the captions are the independent variables, the learned model does not have to capture the variability in image embeddings corresponding to the different but equally good captions – there is only one such image embedding.

Nonlinear neural networks overcome these modeling limitations. The results clearly indicate their superiority over the baselines. The Vec2seq+FFNN model obtains the best results, with accuracies of 60.5% (dev) and 59.0% (test); the accuracy numbers indicate that the Vec2seq+FFNN architecture is superior to the non-recursive fully-connected FFNN architecture (at 55.1% accuracy on test). We show next the impact on performance of the embedding dimension and neural-network sizes, for both the feedforward and the recurrent architectures.

5.3. Analysis: embedding dimension and neural-network sizes

In this section, we compare neural networks models of different sizes. Specifically, we compare embedding dimensions of {64, 256, 512, 1024, 2048},

⁵<https://github.com/tylin/coco-caption>

dim	hidden-1	hidden-2	Dev	Test
FFNN				
64	64	16	56.5	53.9 ± 0.5
256	64	16	56.3	55.1 ± 0.5
256	256	64	55.8	54.3 ± 0.5
512	512	128	54.1	52.5 ± 0.5
1024	1024	256	52.2	51.3 ± 0.5
2048	2048	512	50.7	50.7 ± 0.5
Vec2seq+FFNN			(with default $\lambda_{gen} = 1.0$)	
64	64	16	55.3	54.0 ± 0.5
256	64	16	60.5	59.0 ± 0.5
256	256	64	61.2	58.8 ± 0.5
512	512	128	61.6	59.6 ± 0.5
1024	1024	256	62.5	60.8 ± 0.5
2048	2048	512	63.4	60.8 ± 0.5

Table 4. The impact of model sizes on MCIC-COCO accuracy for the FFNN model.

and two hidden-layer architectures with sizes of $\{(64, 16), (256, 64), (512, 128), (1024, 256), (2048, 512)\}$.

The results in Table 4 illustrate an interesting behavior for the neural-network architectures. For the FFNN models, contrary to expectations, bigger network sizes leads to decreasing accuracy. On the other hand, for Vec2seq+FFNN models, accuracy increases with increased size in model parameters, up until the embedding dimension of the RNN model matches the embedding dimension of the Inception model, at 2048.

At accuracy levels of 63.4% (dev) and 60.8% (test), this performance establishes a high-bar for a computer model performance on the DMC task using the MCIC-COCO dataset. According to the estimate from Table 2, this level of performance is still *significantly* below the 82.8% accuracy achievable by humans, which makes MCIC-COCO a challenging testbed for future models of Vision-Language machine comprehension.

5.4. Multi-task learning for DMC and Image Captioning

In this section, we compare models with different values of λ_{gen} in Eq. 4. This parameter allows for a natural progression from learning for the DMC task only ($\lambda_{gen} = 0$) to focusing on the image captioning loss ($\lambda_{gen} \rightarrow +\infty$). In between the two extremes, we have a multi-task learning objective for jointly learning related tasks.

The results in Table 5 illustrate one of the main points of this paper. That is, the ability to perform the comprehension task (as measured by the accuracy metric) positively correlates with the ability to perform other tasks that require machine comprehension, such as caption generation. At $\lambda_{gen} = 4$, the Vec2seq+FFNN model not only has a high accuracy of detecting the ground-truth option, but it also generates its own captions given the input image, with

λ_{gen}	Acc		ROUGE-L		CIDEr	
	Dev	Test	Dev	Test	Dev	Test
0.0	50.7	50.7 ± 0.5	-	-	-	-
0.1	61.1	59.0 ± 0.5	0.517	0.511	0.901	0.865
1.0	63.4	60.8 ± 0.5	0.528	0.518	0.972	0.903
2.0	63.4	61.3 ± 0.5	0.528	0.519	0.971	0.921
4.0	63.0	60.9 ± 0.5	0.533	0.524	0.989	0.938
8.0	62.1	60.1 ± 0.5	0.526	0.520	0.957	0.914
16.0	61.8	59.6 ± 0.5	0.530	0.519	0.965	0.912

Table 5. The impact of λ_{gen} on MCIC-COCO accuracy, together with caption-generation performance (ROUGE-L and CIDEr against 5 references). All results are obtained with a Vec2seq+FFNN model (embedding size 2048 and hidden-layer sizes of 2048 and 512).

an accuracy measured on MCIC-COCO at 0.9890 (dev) and 0.9380 (test) CIDEr scores. On the other hand, at an accuracy level of about 59% (on test, at $\lambda_{gen} = 0.1$), the generation performance is at only 0.9010 (dev) and 0.8650 (test) CIDEr scores.

We note that there is an inherent trade-off between prediction accuracy and generation performance, as seen for λ_{gen} values above 4.0. This agrees with the intuition that training a Vec2seq+FFNN model using a loss $L(\Theta, \mathbf{w}, \mathbf{u})$ with a larger λ_{gen} means that the ground-truth detection loss (the first term of the loss in Eq.4) may get overwhelmed by the word-generation loss (the second term). However, our empirical results suggest that there is value in training models with a multi-task setup, in which both the comprehension side as well as the generation side are carefully tuned to maximize performance.

6. Discussion

We have proposed and described in detail a new multi-modal machine comprehension task (DMC), combining the challenges of understanding visual scenes and complex language constructs simultaneously. The underlying hypothesis for this work is that computer systems that can be shown to perform increasingly well on this task will do so by constructing a visually-grounded understanding of various linguistic elements and their dependencies. This type of work can therefore benefit research in both machine visual understanding and language comprehension.

The Vec2seq+FFNN architecture that we propose for addressing this combined challenge is a generic multi-task model. It can be trained end-to-end to display both the ability to choose the most likely text associated with an image (thus enabling a direct measure of its ‘‘comprehension’’ performance), as well as the ability to generate a complex description of that image (thus enabling a direct measure of its performance in an end-to-end complex and meaningful task). The empirical results we present validate the underlying hypothesis of our work, by showing that we can measure

the decisions made by such a computer system and validate that improvements in comprehension and generation happen in tandem.

The experiments presented in this work are done training our systems in an end-to-end fashion, starting directly from raw pixels. We hypothesize that our framework can be fruitfully used to show that incorporating specialized vision systems (such as object detection, scene recognition, pose detection, etc.) is beneficial. More precisely, not only it can lead to a direct and measurable impact on a computer system's ability to perform image understanding, but it can express that understanding in an end-to-end complex task.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. Corrado, A. Davis, J. Dean, M. Devin, S. Ghemawat, I. Goodfellow, A. Harp, G. Irving, M. Isard, Y. Jia, R. Jozefowicz, L. Kaiser, M. Kudlur, J. Levenberg, D. Mané, R. Monga, S. Moore, D. Murray, C. Olah, M. Schuster, J. Shlens, B. Steiner, I. Sutskever, K. Talwar, P. Tucker, V. Vanhoucke, V. Vasudevan, F. Viégas, O. Vinyals, P. Warden, M. Wattenberg, M. Wicke, Y. Yu, and X. Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [2] P. Anderson, B. Fernando, M. Johnson, and S. Gould. SPICE: semantic propositional image caption evaluation. *CoRR*, abs/1607.08822, 2016.
- [3] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual question answering. In *International Conference on Computer Vision (ICCV)*, 2015.
- [4] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. In *Proceedings of ICLR*, 2015.
- [5] S. Banerjee and A. Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the ACL Workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*, 2005.
- [6] Y. Bengio. Learning deep architectures for ai. *Found. Trends Mach. Learn.*, 2(1):1–127, Jan. 2009.
- [7] R. Bernardi, R. Cakici, D. Elliott, A. Erdem, E. Erdem, N. Ikizler-Cinbis, F. Keller, A. Muscat, and B. Plank. Automatic description generation from images: A survey of models, datasets, and evaluation measures. *JAIR*, 55, 2016.
- [8] D. Chen, J. Bolton, and C. D. Manning. A Thorough Examination of the CNN/Daily Mail Reading Comprehension Task. In *Proceedings of ACL*, 2016.
- [9] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of EMNLP, October 25-29, 2014, Doha, Qatar*, pages 1724–1734, 2014.
- [10] J. Donahue, L. A. Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014.
- [11] H. Fang, S. Gupta, F. Iandola, R. Srivastava, L. Deng, P. Dollár, J. Gao, X. He, M. Mitchell, J. Platt, et al. From captions to visual concepts and back. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [12] A. Frome, G. S. Corrado, J. Shlens, S. Bengio, J. Dean, T. Mikolov, et al. Devise: A deep visual-semantic embedding model. In *Advances in Neural Information Processing Systems (NIPS)*, 2013.
- [13] H. Gao, J. Mao, J. Zhou, Z. Huang, L. Wang, and W. Xu. Are you talking to a machine? dataset and methods for multilingual image question answering. In *NIPS*, 2015.
- [14] M. Hodosh and J. Hockenmaier. Focused evaluation for image description with binary forced-choice tasks. In *Proc. 5th Vision and Language Workshop*, 2016.
- [15] M. Hodosh, P. Young, and J. Hockenmaier. Framing image description as a ranking task: Data, models and evaluation metrics. *JAIR*, 2013.
- [16] A. Karpathy and L. Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [17] R. Kiros, R. Salakhutdinov, and R. S. Zemel. Unifying visual-semantic embeddings with multimodal neural language models. *Transactions of the Association for Computational Linguistics*, 2015.
- [18] R. Krishna, Y. Zhu, O. Groth, J. Johnson, K. Hata, J. Kravitz, S. Chen, Y. Kalantidis, L.-J. Li, D. A. Shamma, M. Bernstein, and L. Fei-Fei. Visual Genome: Connecting language and vision using crowdsourced dense image annotations. 2016.

- [19] Q. Le and T. Mikolov. Distributed representations of sentences and documents. In *Proceedings of the 31st International Conference on Machine Learning*, Beijing, China, 2014.
- [20] C.-Y. Lin and F. J. Och. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of ACL*, 2004.
- [21] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014.
- [22] X. Lin and D. Parikh. Leveraging visual question answering for image-caption ranking. *CoRR*, abs/1605.01379, 2016.
- [23] M. Malinowski and M. Fritz. A multi-world approach to question answering about real-world scenes based on uncertain input. In *NIPS*, 2014.
- [24] M. Malinowski, M. Rohrbach, and M. Fritz. Ask your neurons: A neural-based approach to answering questions about images. In *ICCV*, 2015.
- [25] J. Mao, W. Xu, Y. Yang, J. Wang, and A. Yuille. Deep captioning with multimodal recurrent neural networks (mRNN). In *Proc. Int. Conf. Learn. Representations*, 2015.
- [26] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: A method for automatic evaluation of machine translation. In *Proceedings of ACL*, pages 311–318, 2002.
- [27] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba. Sequence level training with recurrent neural networks. *CoRR*, abs/1511.06732, 2015.
- [28] M. Ren, R. Kiros, and R. Zemel. Image question answering: A visual semantic embedding model and a new dataset. In *NIPS*, 2015.
- [29] I. Sutskever, O. Vinyals, and Q. V. V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems 27*, pages 3104–3112. Curran Associates, Inc., 2014.
- [30] C. Szegedy, V. Vanhoucke, S. Ioffe, J. Shlens, and Z. Wojna. Rethinking the inception architecture for computer vision. volume abs/1512.00567, 2015.
- [31] K. Tu, M. Meng, M. W. Lee, T. E. Choe, and S. C. Zhu. Joint video and text parsing for understanding events and answering queries. *IEEE MultiMedia*, 2014.
- [32] R. Vedantam, C. Lawrence Zitnick, and D. Parikh. Cider: Consensus-based image description evaluation. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2015.
- [33] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *Proc. of IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [34] Q. Wu, D. Teney, P. Wang, C. Shen, A. R. Dick, and A. van den Hengel. Visual question answering: A survey of methods and datasets. *CoRR*, abs/1607.05910, 2016.
- [35] Q. Wu, P. Wang, C. Shen, A. Dick, and A. van den Hengel. Ask me anything: Free-form visual question answering based on knowledge from external sources. In *CVPR*, 2016.
- [36] K. Xu, J. Ba, R. Kiros, K. Cho, A. C. Courville, R. Salakhutdinov, R. S. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. *CoRR*, abs/1502.03044, 2015.
- [37] K. Xu, J. Ba, R. Kiros, A. Courville, R. Salakhutdinov, R. Zemel, and Y. Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *Proc. of the 32nd International Conference on Machine Learning (ICML)*, 2015.
- [38] Z. Yang, X. He, J. Gao, L. Deng, and A. Smola. Stacked attention networks for image question answering. In *CVPR*, 2016.
- [39] B. Yao and F.-F. Li. Modeling mutual context of object and human pose in human-object interaction activities. In *Proceedings of the 2010 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2010.
- [40] L. Yu, E. Park, A. C. Berg, and T. L. Berg. Visual madlibs: Fill-in-the-blank description generation and question answering. In *ICCV*, 2015.
- [41] Y. Zhu, O. Groth, M. Bernstein, and L. Fei-Fei. Visual7w: Grounded question answering in images. In *CVPR*, 2016.