

Bottom-Up and Top-Down Attention for Image Captioning and Visual Question Answering

Peter Anderson^{1*}, Xiaodong He², Chris Buehler², Damien Teney³

Mark Johnson⁴, Stephen Gould¹, Lei Zhang²

¹Australian National University ²Microsoft Research

³University of Adelaide ⁴Macquarie University

Abstract

Top-down visual attention mechanisms have been used extensively in image captioning and visual question answering (VQA) to enable deeper image understanding through fine-grained analysis and even multiple steps of reasoning. In this work, we propose a combined bottom-up and top-down attention mechanism that enables attention to be calculated at the level of objects and other salient image regions. This is the natural basis for attention to be considered. Within our approach, the bottom-up mechanism (based on Faster R-CNN) proposes image regions, each with an associated feature vector, while the top-down mechanism determines feature weightings. Applying this approach to image captioning, our results on the MSCOCO test server establish a new state-of-the-art for the task, improving the best published result in terms of CIDEr score from 114.7 to 117.9 and BLEU-4 from 35.2 to 36.9. Demonstrating the broad applicability of the method, applying the same approach to VQA we obtain first place in the 2017 VQA Challenge.

1. Introduction

Problems combining image and language understanding such as image captioning [4] and visual question answering (VQA) [13] continue to inspire considerable research at the boundary of computer vision and natural language processing. In both these tasks it is often necessary to perform some fine-grained visual processing, or even multiple steps of reasoning to generate high quality outputs. As a result, visual attention mechanisms have been widely adopted in both image captioning [37, 29, 54, 52] and VQA [12, 30, 51, 53, 59]. These mechanisms improve performance by learning to focus on the regions of the image that are salient and are currently based on deep neural network architectures.

*Work performed while interning at Microsoft.

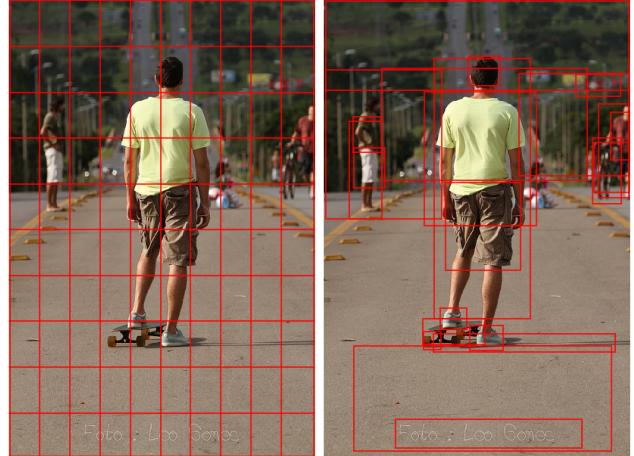


Figure 1. Typically, attention models operate on CNN features corresponding to a uniform grid of equally-sized image regions (left). Our approach enables attention to be calculated at the level of objects and other salient image regions (right).

In the human visual system, attention can be focused volitionally by top-down signals determined by the current task (e.g., looking for something), and automatically by bottom-up signals associated with unexpected, novel or salient stimuli [3, 6]. In this paper we adopt similar terminology and refer to attention mechanisms driven by non-visual or task-specific context as ‘top-down’, and purely visual feed-forward attention mechanisms as ‘bottom-up’.

Most conventional visual attention mechanisms used in image captioning and VQA are of the top-down variety. Taking as context a representation of a partially-completed caption output, or a question relating to the image, these mechanisms are typically trained to selectively attend to the output of one or more layers of a convolutional neural net (CNN). However, this approach gives little consideration to how the image regions that are subject to attention are determined. As illustrated conceptually in Figure 1, the resulting input regions correspond to a uniform grid of equally sized

and shaped receptive fields – irrespective of the content of the image. To generate more human-like captions and question answers, objects and other salient image regions are a much more natural basis for attention [10, 39].

In this paper we propose a combined bottom-up and top-down visual attention mechanism. The **bottom-up** mechanism proposes a set of salient image regions, with each region represented by a pooled convolutional feature vector. Practically, we implement bottom-up attention using Faster R-CNN [36], which represents a natural expression of a bottom-up attention mechanism. The **top-down** mechanism uses task-specific context to predict an attention distribution over the image regions. The feature glimpse is computed as a weighted average of image features over all regions.

We evaluate the impact of combining bottom-up and top-down attention on two tasks. We first present an image captioning model that takes multiple glimpses of salient image regions during caption generation. Empirically, we find that the inclusion of bottom-up attention has a significant positive benefit for image captioning. Our results on the MSCOCO test server establish a new state-of-the-art for the task, improving the best published result in terms of CIDEr score from 114.7 to 117.9 and BLEU-4 from 35.2 to 36.9. Demonstrating the broad applicability of the method, we additionally present a VQA model using the same bottom-up attention features that won first place in the 2017 VQA Challenge, achieving 70.2% overall accuracy on the VQA v2.0 test-standard server.

2. Related Work

A large number of attention-based deep neural networks have been proposed for image captioning and VQA. Typically, these models can be characterized as top-down approaches, with context provided by a representation of a partially-completed caption in the case of image captioning [37, 29, 54, 52], or a representation of the question in the case of VQA [12, 30, 51, 53, 59]. In each case attention is applied to the output of one or more layers of a CNN, by predicting a weighting for each spatial location in the CNN output. However, determining the optimal number of image regions invariably requires an unwinnable trade-off between coarse and fine levels of detail. Furthermore, the arbitrary positioning of the regions with respect to image content may make it more difficult to detect objects that are poorly aligned to regions and to bind visual concepts associated with the same object.

We are aware of a few models that seek to combine bottom-up and top-down attention processing. Semantic attention mechanisms [56, 57] apply top-down attention to a list of visual attributes detected in an image. Attribute detection could be regarded as a bottom-up attention process outputting words rather than image feature vectors. How-

ever, these and other non-attention attribute detection approaches [55, 50] do not retain spatial information, treating the detected attributes as a bag of words. In contrast, in our feature-based approach a single feature vector can be discriminative for several visual words, for example an adjective and a noun, providing a strong signal that the detected concepts belong to the same object.

The most similar works to ours are those that also propose attention over salient image regions. Jin et. al. [20] use selective search [45] to identify salient image regions, which are filtered with a classifier then resized and CNN-encoded as input to an image captioning model with attention. Similarly to DenseCap [21], Pedersoli et. al. [33] use spatial transformer networks [19] to generate image features, which are processed using a custom attention model based on three bi-linear pairwise interactions. The use of spatial transformer networks allows for end-to-end training of region proposal coordinates. In contrast, in our approach we solve the region supervision problem using pretraining. In addition, we apply our approach to both image captioning and VQA, and for better comparison with prior work we use more conventional captioning / VQA architectures.

3. Approach

Given an image I , both our image captioning model and our VQA model take as input a possibly variably-sized set of k image features, $V = \{v_1, \dots, v_k\}$, $v_i \in \mathbb{R}^D$, such that each image feature encodes a salient region of the image.

The spatial image features V can be variously defined as the output of our bottom-up attention model, or, following standard practice, as the spatial output layer of a CNN. We describe our approach to implementing a bottom-up attention model in Section 3.1. In Section 3.2 we outline the architecture of our image captioning model and in Section 3.3 we outline our VQA model. We note that for the top-down attention component, both models use simple one-pass attention mechanisms, as opposed to the more complex schemes of recent models such as stacked, multi-headed, or bidirectional attention [53, 18, 23, 30] that could also be applied.

3.1. Bottom-Up Attention Model

The definition of spatial image features V is generic. However, in this work we define spatial regions in terms of bounding boxes and implement bottom-up attention using Faster R-CNN [36]. Faster R-CNN is an object detection model designed to identify instances of objects belonging to certain classes and localize them with bounding boxes. Other region proposal networks could also be trained as an attentive mechanism [35, 27].

Faster R-CNN detects objects in two stages. The first stage, described as a Region Proposal Network (RPN), predicts object proposals. A small network is slid over features at an intermediate level of a CNN. At each spatial loca-

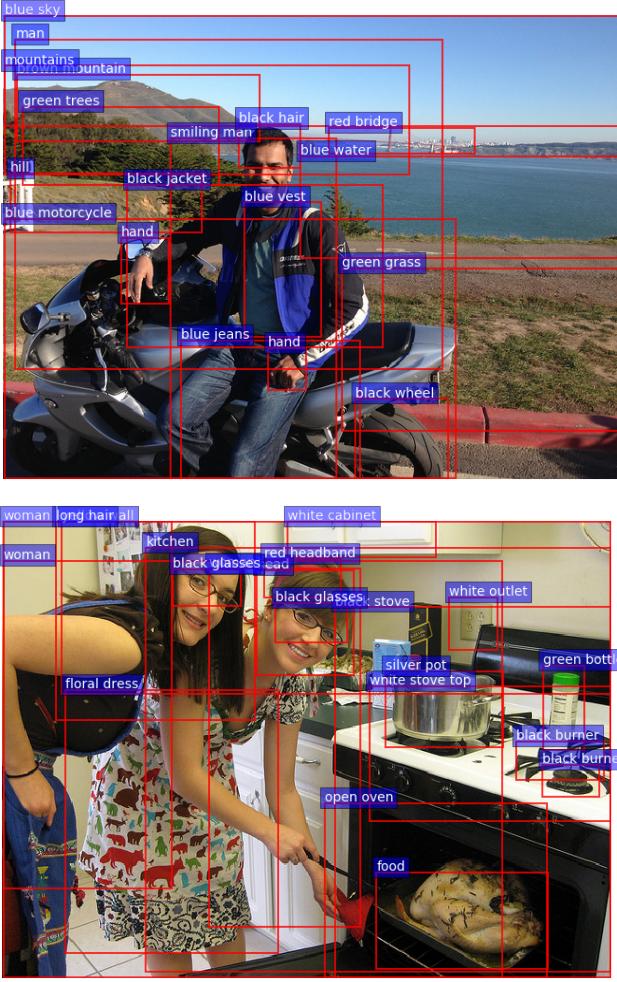


Figure 2. Example output from our Faster R-CNN bottom-up attention model. Each bounding box is labeled with an attribute class followed by an object class. Note however, that in captioning and VQA we utilize only the feature vectors – not the predicted labels.

tion the network predicts a class-agnostic objectness score and a bounding box refinement for anchor boxes of multiple scales and aspect ratios. Using greedy non-maximum suppression with an intersection-over-union (IoU) threshold, the top box proposals are selected as input to the second stage. In the second stage, region of interest (RoI) pooling is used to extract a small feature map (e.g. 14×14) for each box proposal. These feature maps are then batched together as input to the final layers of the CNN. The final output of the model consists of a softmax distribution over class labels and class-specific bounding box refinements for each box proposal.

In this work, we use Faster R-CNN in conjunction with the ResNet-101 [14] CNN. To generate an output set of image features V for use in image captioning or VQA, we take the final output of the model and perform non-maximum

suppression for each object class using an IoU threshold. We then select all regions where any class detection probability exceeds a confidence threshold. For each selected region i , v_i is defined as the mean-pooled convolutional feature from this region, such that the dimension D of the image feature vectors is 2048. Used in this fashion, Faster R-CNN effectively functions as a ‘hard’ attention mechanism, as only a relatively small number of image bounding box features are selected from a large number of possible configurations.

To pretrain the bottom-up attention model, we first initialize Faster R-CNN with ResNet-101 pretrained for classification on ImageNet [38]. We then train on Visual Genome [47] data. To aid the learning of good feature representations, we add an additional training output for predicting attribute classes (in addition to object classes). To predict attributes for region i , we concatenate the mean pooled convolutional feature v_i with a learned embedding of the ground-truth object class, and feed this into an additional output layer defining a softmax distribution over each attribute class plus a ‘no attributes’ class.

The original Faster R-CNN multi-task loss function contains four components, defined over the classification and bounding box regression outputs for both the RPN and the final object class proposals respectively. We retain these components and add an additional multi-class loss component to train the attribute predictor. In Figure 2 we provide some examples of model output. We make code, models and features available from the project page¹.

3.2. Captioning Model

Given a set of image features V , our proposed captioning model uses a ‘soft’ top-down attention mechanism to weight each feature during caption generation, using the existing partial output sequence as context. This approach is broadly similar to several previous works [37, 29, 52]. However, the particular design choices outlined below make for a relatively simple yet high-performing baseline model. Even without bottom-up attention, our captioning model achieves performance comparable to state-of-the-art on most evaluation metrics (refer Table 1).

At a high level, the captioning model is composed of two LSTM [17] layers using a standard implementation [9]. In the sections that follow we will refer to the operation of the LSTM over a single time step using the following notation:

$$h_t = \text{LSTM}(x_t, h_{t-1}) \quad (1)$$

where x_t is the LSTM input vector and h_t is the LSTM output vector. Here we have neglected the propagation of memory cells for notational convenience. We now describe the formulation of the LSTM input vector x_t and the output

¹<http://www.panderson.me/up-down-attention>

vector \mathbf{h}_t for each layer of the model. The overall captioning model is illustrated in Figure 3.

3.2.1 Top-Down Attention LSTM

Within the captioning model, we characterize the first LSTM layer as a top-down visual attention model, and the second LSTM layer as a language model, indicating each layer with superscripts in the equations that follow. Note that the bottom-up attention model is described in Section 3.1, and in this section it's outputs are simply considered as features V . The input vector to the attention LSTM at each time step consists of the previous output of the language LSTM, concatenated with the mean-pooled image feature $\bar{\mathbf{v}} = \frac{1}{k} \sum_i \mathbf{v}_i$ and an encoding of the previously generated word, given by:

$$\mathbf{x}_t^1 = [\mathbf{h}_{t-1}^2, \bar{\mathbf{v}}, W_e \Pi_t] \quad (2)$$

where $W_e \in \mathbb{R}^{E \times \Sigma}$ is a word embedding matrix for a vocabulary of size Σ , and Π_t is one-hot encoding of the input word at timestep t . These inputs provide the attention LSTM with maximum context regarding the state of the language LSTM, the overall content of the image, and the partial caption output generated so far, respectively. The word embedding is learned from random initialization without pretraining.

Given the output \mathbf{h}_t^1 of the attention LSTM, at each time step t we generate a normalized attention weight $\alpha_{i,t}$ for each of the k image features \mathbf{v}_i as follows:

$$a_{i,t} = \mathbf{w}_a^T \tanh(W_{va} \mathbf{v}_i + W_{ha} \mathbf{h}_t^1) \quad (3)$$

$$\alpha_t = \text{softmax}(\mathbf{a}_t) \quad (4)$$

where $W_{va} \in \mathbb{R}^{H \times V}$, $W_{ha} \in \mathbb{R}^{H \times M}$ and $\mathbf{w}_a \in \mathbb{R}^H$ are learned parameters. The attended image feature used as input to the language LSTM is calculated as a convex combination of all input features:

$$\hat{\mathbf{v}}_t = \sum_{i=1}^K \alpha_{i,t} \mathbf{v}_i \quad (5)$$

3.2.2 Language LSTM

The input to the language model LSTM consists of the attended image feature, concatenated with the output of the attention LSTM, given by:

$$\mathbf{x}_t^2 = [\hat{\mathbf{v}}_t, \mathbf{h}_t^1] \quad (6)$$

Using the notation $y_{1:T}$ to refer to a sequence of words (y_1, \dots, y_T) , at each time step t the conditional distribution over possible output words is given by:

$$p(y_t | y_{1:t-1}) = \text{softmax}(W_p \mathbf{h}_t^2 + \mathbf{b}_p) \quad (7)$$

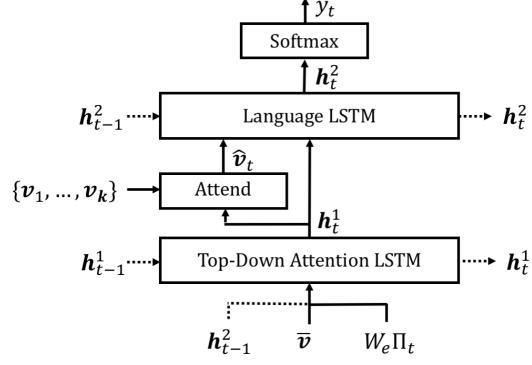


Figure 3. Overview of the proposed captioning model. Two LSTM layers are used to selectively attend to spatial image features $\{\mathbf{v}_1, \dots, \mathbf{v}_k\}$. These features can be defined as the spatial output of a CNN, or following our approach, generated using bottom-up attention.

where $W_p \in \mathbb{R}^{\Sigma \times M}$ and $\mathbf{b}_p \in \mathbb{R}^\Sigma$ are learned weights and biases. The distribution over complete output sequences is calculated as the product of conditional distributions:

$$p(y_{1:T}) = \prod_{t=1}^T p(y_t | y_{1:t-1}) \quad (8)$$

3.2.3 Objective

Given a target ground truth sequence $y_{1:T}^*$ and a captioning model with parameters θ , we minimize the following cross entropy loss:

$$L_{XE}(\theta) = - \sum_{t=1}^T \log(p_\theta(y_t^* | y_{1:t-1}^*)) \quad (9)$$

Although directly optimizing the CIDEr metric reportedly has little impact on perceived caption quality in human evaluations [26], for fair comparison with recent work [37] we also report results optimized for CIDEr [46]. In this stage of training we begin with the cross-entropy trained model and then seek to minimize the negative expected score:

$$L_R(\theta) = -\mathbf{E}_{y_{1:T} \sim p_\theta} [r(y_{1:T})] \quad (10)$$

where r is the score function (e.g., CIDEr). Similar loss function have been demonstrated effective in other tasks with structural output such as machine translation [16]. Following the REINFORCE [49] approach from Self-critical Sequence Training [37] (SCST), the gradient of this loss can be approximated:

$$\nabla_\theta L_R(\theta) \approx -(r(y_{1:T}^s) - r(\hat{y}_{1:T})) \nabla_\theta \log p_\theta(y_{1:T}^s) \quad (11)$$

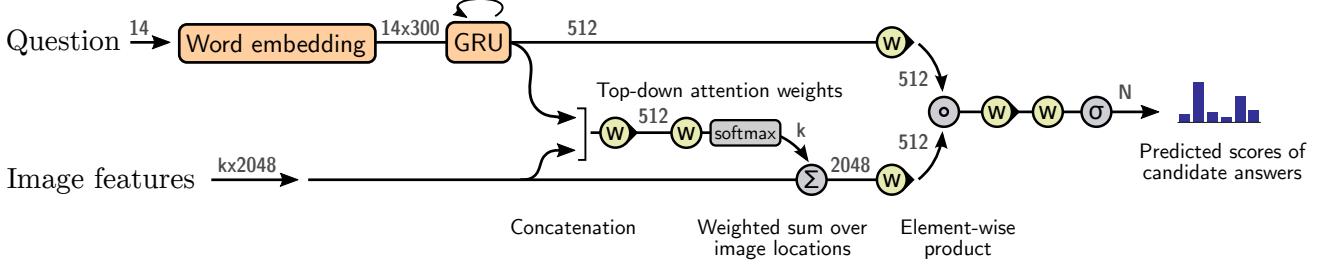


Figure 4. Overview of the proposed VQA model. A deep neural network implements a joint embedding of the question and image features $\{v_1, \dots, v_k\}$. These features can be defined as the spatial output of a CNN, or following our approach, generated using bottom-up attention. Output is generated by a multi-label classifier operating over a fixed set of candidate answers. Gray numbers indicate the dimensions of the vector representations between layers. Yellow elements use learned parameters.

where $y_{1:T}^s$ is a sampled caption and $r(\hat{y}_{1:T})$ defines the baseline score obtained by greedily decoding the current model. This gradient tends to increase the probability of sampled captions that score higher than the score from the current model.

We have observed when decoding using beam search that the resulting beam typically contains at least one very high quality (and high scoring) caption – although frequently the best caption does not have the highest log-probability of the set. Therefore, we make one additional approximation. Rather than sampling captions from the entire probability distribution, for more rapid training we take the captions in the decoded beam as a sample set. Using this approach, we complete CIDEr optimization in a single epoch.

3.3. VQA Model

Given a set of spatial image features V , our proposed VQA model also uses a ‘soft’ top-down attention mechanism to weight each feature, using the question representation as context. As illustrated in Figure 6, the proposed model implements the well-known joint multimodal embedding of the question and the image, followed by a prediction of regression of scores over a set of candidate answers. This approach has been the basis of numerous previous models [18, 23, 42]. However, as with our captioning model, implementation decisions are important to ensure that this relatively simple model delivers high performance.

The learned non-linear transformations within the network are implemented with gated hyperbolic tangent activations [7]. These are a special case of highway networks [40] that have shown a strong empirical advantage over traditional ReLU or tanh layers. Each of our ‘gated tanh’ layers implements a function $f_a : x \in \mathbb{R}^m \rightarrow y \in \mathbb{R}^n$ with parameters a defined as follows:

$$\tilde{y} = \tanh(Wx + b) \quad (12)$$

$$g = \sigma(W'x + b') \quad (13)$$

$$y = \tilde{y} \circ g \quad (14)$$

where σ is the sigmoid activation function, $W, W' \in \mathbb{R}^{n \times m}$ are learned weights, $b, b' \in \mathbb{R}^n$ are learned biases, and \circ is the Hadamard (element-wise) product. The vector g acts multiplicatively as a gate on the intermediate activation \tilde{y} .

Our proposed approach first encodes each question as the hidden state q of a gated recurrent unit [5] (GRU), with each input word represented using a learned word embedding. Similar to Equation 3, given the output q of the GRU, we generate an unnormalized attention weight a_i for each of the k image features v_i as follows:

$$a_i = w_a^T f_a([v_i, q]) \quad (15)$$

where w_a^T is a learned parameter vector. Equation 4 and Equation 5 (neglecting subscripts t) are used to calculate the normalized attention weight and the attended image feature \hat{v} . The distribution over possible output responses y is given by:

$$h = f_q(q) \circ f_v(\hat{v}) \quad (16)$$

$$p(y) = \sigma(W_o f_o(h)) \quad (17)$$

Where h is a joint representation of the question and the image, and $W_o \in \mathbb{R}^{\Sigma \times M}$ are learned weights. Due to space constraints, many important aspects of our VQA approach are not detailed here. For full specifics of the VQA model including a detailed exploration of architectures and hyper-parameters, refer to Teney et. al. [41].

3.4. Implementation Details

Our bottom-up attention Faster R-CNN implementation uses an IoU threshold of 0.7 for region proposal suppression, and 0.3 for object class suppression. To select salient image regions, a class detection confidence threshold of 0.2 is used, allowing the number of regions per image k to vary with the complexity of the image, up to a maximum of 100. However, in initial experiments we find that simply selecting the top 36 features in each image works almost as well in both downstream tasks. Since Visual Genome contains a

relatively large number of annotations per image, the model is relatively intensive to train. Using 8 Nvidia M40 GPUs, we take around 5 days to complete 380K training iterations, although we suspect that faster training regimes could also be effective.

In the captioning model, we set the number of hidden units M in each LSTM to 1,000, the number of hidden units H in the attention layer to 512, and the size of the input word embedding E to 1,000. In training, we use a simple learning rate schedule, beginning with a learning rate of 0.01 which is reduced to zero on a straight-line basis over 60K iterations using a batch size of 100 and a momentum parameter of 0.9. Training using two Nvidia Titan X GPUs takes around 9 hours (including less than one hour for CIDEr optimization). During optimization and decoding we use a beam size of 5. When decoding we also enforce the constraint that a single word cannot be predicted twice in a row. Note that in both our captioning and VQA models, image features are fixed and not finetuned.

In the VQA model, we use 300 dimension word embeddings, initialized with pretrained GloVe vectors [34], and we use hidden states of dimension 512. We train the VQA model using AdaDelta [58] and regularize with early stopping. The training of the model takes in the order of 12–18 hours on a single Nvidia K40 GPU. Refer to Teney et al. [41] for further details of the VQA model implementation.

4. Evaluation

4.1. Datasets

4.1.1 Visual Genome Dataset

We use the Visual Genome [47] dataset to pretrain our bottom-up attention model, and for data augmentation when training our VQA model. The dataset contains 108K images densely annotated with scene graphs containing objects, attributes and relationships, as well as 1.7M visual question answers.

For pretraining the bottom-up attention model, we use only the object and attribute data. We reserve 5K images for validation, and 5K images for future testing, treating the remaining 98K images as training data. As approximately 51K Visual Genome images are also found in the MSCOCO captions dataset [25], we are careful to avoid contamination of our MSCOCO validation and test sets. We ensure that any images found in both datasets are contained in the same split in both datasets.

As the object and attribute annotations consist of freely annotated strings, rather than classes, we perform extensive cleaning and filtering of the training data, including manually removing abstract concepts with low precision in initial experiments. Our final training set contains 1,600 object classes and 400 attribute classes. Note that we have not

filtered overlapping classes (e.g. ‘person’, ‘man’, ‘guy’), classes with both singular and plural versions (e.g. ‘tree’, ‘trees’) and classes that are difficult to precisely localize (e.g. ‘sky’, ‘grass’, ‘buildings’).

When training the VQA model, we augment the VQA v2.0 training data with Visual Genome question and answer pairs provided the correct answer is present in model’s answer vocabulary. This represents about 30% of the available data, or 485K questions.

4.1.2 Microsoft COCO Dataset

To evaluate our proposed captioning model, we use the MSCOCO 2014 captions dataset [25]. For validation of model hyperparameters and offline testing, we use the ‘Karpathy’ splits [22] that have been used extensively for reporting results in prior work. This split contains 113,287 training images with five captions each, and 5K images respectively for validation and testing. Our MSCOCO test server submission is trained on the entire MSCOCO 2014 training and validation set (123K images).

We follow standard practice and perform only minimal text pre-processing, converting all sentences to lower case, tokenizing on white space, and filtering words that do not occur at least five times, resulting in a model vocabulary of 10,010 words. To evaluate caption quality, we use the standard automatic evaluation metrics, namely SPICE [1], CIDEr [46], METEOR [8], ROUGE-L [24] and BLEU [32].

4.1.3 VQA v2.0 Dataset

To evaluate our proposed VQA model, we use the recently introduced VQA v2.0 dataset [13], which attempts to minimize the effectiveness of learning dataset priors by balancing the answers to each question. The dataset, which was used as the basis of the 2017 VQA Challenge², contains 1.1M questions with 11.1M answers relating to MSCOCO images.

We perform standard question text preprocessing and tokenization. Questions are trimmed to a maximum of 14 words for computational efficiency. The set of candidate answers is restricted to correct answers in the training set that appear more than 8 times, resulting in an output vocabulary size of 3,129. Our VQA test server submissions are trained on the training and validation sets plus additional questions and answers from Visual Genome. To evaluate answer quality, we report accuracies using the standard VQA metric [2], which takes into account the occasional disagreement between annotators for the ground truth answers.

²<http://www.visualqa.org/challenge.html>

	Cross-Entropy Loss							CIDEr Optimization						
	BLEU-1	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE		BLEU-1	BLEU-4	METEOR	ROUGE-L	CIDEr	SPICE	
SCST:Att2in [37]	-	31.3	26.0	54.3	101.3	-	-	33.3	26.3	55.3	111.4	-	-	
SCST:Att2all [37]	-	30.0	25.9	53.4	99.4	-	-	34.2	26.7	55.7	114.0	-	-	
Ours: ResNet	74.5	33.4	26.1	54.4	105.4	19.2	76.6	34.0	26.5	54.9	111.1	20.2		
Ours: Up-Down	77.2	36.2	27.0	56.4	113.5	20.3	79.8	36.3	27.7	56.9	120.1	21.4		
Relative Improvement	4%	8%	3%	4%	8%	6%	4%	7%	5%	4%	8%	6%		

Table 1. Single-model image captioning performance on the MSCOCO Karpathy test split. Our baseline ResNet model obtains similar results to SCST [37], the existing state-of-the-art. Our Up-Down model demonstrates the contribution of bottom-up attention, with significant gains across all metrics regardless of whether cross-entropy loss or CIDEr optimization is used.

	Cross-Entropy Loss							CIDEr Optimization						
	SPICE	Objects	Attributes	Relations	Color	Count	Size	SPICE	Objects	Attributes	Relations	Color	Count	Size
Ours: ResNet	19.2	35.4	8.6	5.3	12.2	4.1	3.9	20.2	37.0	9.2	6.1	10.6	12.0	4.3
Ours: Up-Down	20.3	37.1	9.2	5.8	12.7	6.5	4.5	21.4	39.1	10.0	6.5	11.4	18.4	3.2

Table 2. Breakdown of SPICE F-scores over various subcategories on the MSCOCO Karpathy test split. Our full Up-Down model is better than the baseline at identifying objects, as well as describing object attributes and the relations between objects.

4.2. ResNet Baseline

To quantify the impact of bottom-up attention, in both our captioning and VQA experiments we evaluate our full model (*Up-Down*) against prior work as well as an ablated baseline. In each case, the baseline (*ResNet*), uses a ResNet [14] CNN pretrained on ImageNet [38] to encode each image in place of the bottom-up attention mechanism.

In image captioning experiments, similarly to previous work [37] we encode the full-sized input image with the final convolutional layer of Resnet-101, and use bilinear interpolation to resize the output to a fixed size spatial representation of 10×10 . This is equivalent to the maximum number of spatial regions used in our full model. In VQA experiments, we encode the resized input image with ResNet-200 [15]. In separate experiments we evaluate the effect of varying the size of the spatial output from its original size of 14×14 , to 7×7 (using bilinear interpolation) and 1×1 (i.e., mean pooling without attention).

4.3. Image Captioning Results

In Table 1 we report the performance of our full model and the ResNet baseline in comparison to the existing state-of-the-art Self-critical Sequence Training [37] (SCST) approach on the test portion of the Karpathy splits. For fair comparison, results are reported for models trained with both standard cross-entropy loss, and models optimized for CIDEr score. Note that the SCST approach uses ResNet-101 encoding of full images, similar to our ResNet baseline. All results are reported for a single model with no fine-tuning of the input ResNet / R-CNN model. However,

the SCST results are selected from the best of four random initializations, while our results are outcomes from a single initialization.

Relative to the SCST models, our ResNet baseline obtains slightly better performance under cross-entropy loss, and slightly worse performance when optimized for CIDEr score. After incorporating bottom-up attention, our full Up-Down model shows significant improvements across all metrics regardless of whether cross-entropy loss or CIDEr optimization is used. Using just a single model, we obtain the best reported results for the Karpathy test split. As illustrated in Table 2, the contribution from bottom-up attention is broadly based, illustrated by improved performance in terms of identifying objects, object attributes and also the relationships between objects.

Table 3 reports the performance of 4 ensembled models trained with CIDEr optimization on the official MSCOCO evaluation server, along with the highest ranking previously published results. At the time of submission (18 July 2017), we outperform all other test server submissions on all reported evaluation metrics.

4.4. VQA Results

In Table 4 we report the single model performance of our full Up-Down VQA model relative to several ResNet baselines on the VQA v2.0 validation set. The addition of bottom-up attention provides a significant improvement over the best ResNet baseline across all question types, even though the ResNet baseline uses approximately twice as many convolutional layers. Table 5 reports the performance of 30 ensembled models on the official VQA 2.0

	BLEU-1		BLEU-2		BLEU-3		BLEU-4		METEOR		ROUGE-L		CIDEr		SPICE	
	c5	c40	c5	c40	c5	c40										
Google NIC [48]	71.3	89.5	54.2	80.2	40.7	69.4	30.9	58.7	25.4	34.6	53.0	68.2	94.3	94.6	18.2	63.6
MSR Captivator [11]	71.5	90.7	54.3	81.9	40.7	71.0	30.8	60.1	24.8	33.9	52.6	68.0	93.1	93.7	18.0	60.9
M-RNN [31]	71.6	89.0	54.5	79.8	40.4	68.7	29.9	57.5	24.2	32.5	52.1	66.6	91.7	93.5	17.4	60.0
LRCN [9]	71.8	89.5	54.8	80.4	40.9	69.5	30.6	58.5	24.7	33.5	52.8	67.8	92.1	93.4	17.7	59.9
Hard-Attention [52]	70.5	88.1	52.8	77.9	38.3	65.8	27.7	53.7	24.1	32.2	51.6	65.4	86.5	89.3	17.2	59.8
ATT-FCN [56]	73.1	90.0	56.5	81.5	42.4	70.9	31.6	59.9	25.0	33.5	53.5	68.2	94.3	95.8	18.2	63.1
Review Net [54]	72.0	90.0	55.0	81.2	41.4	70.5	31.3	59.7	25.6	34.7	53.3	68.6	96.5	96.9	18.5	64.9
MSM [55]	73.9	91.9	57.5	84.2	43.6	74.0	33.0	63.2	25.6	35.0	54.2	70.0	98.4	100.3	-	-
Adaptive[29]	74.8	92.0	58.4	84.5	44.4	74.4	33.6	63.7	26.4	35.9	55.0	70.5	104.2	105.9	19.7	67.3
PG-SPIDER-TAG [26]	75.1	91.6	59.1	84.2	44.5	73.8	33.1	62.4	25.5	33.9	55.1	69.4	104.2	107.1	19.2	63.7
SCST:Att2all [37]	78.1	93.7	61.9	86.0	47.0	75.9	35.2	64.5	27.0	35.5	56.3	70.7	114.7	116.7	-	-
Ours: Up-Down	80.2	95.2	64.1	88.8	49.1	79.4	36.9	68.5	27.6	36.7	57.1	72.4	117.9	120.5	-	-

Table 3. Leaderboard of published image captioning models on the online MSCOCO test server. Our submission, an ensemble of 4 models optimized for CIDEr with different initializations, outperforms previous work on all reported metrics.

	Yes/No	Number	Other	Overall
Ours: ResNet (1×1)	76.0	36.5	46.8	56.3
Ours: ResNet (14×14)	76.6	36.2	49.5	57.9
Ours: ResNet (7×7)	77.6	37.7	51.5	59.4
Ours: Up-Down	80.3	42.8	55.8	63.2
Relative Improvement	3%	14%	8%	6%

Table 4. Single-model performance on the VQA v2.0 validation set. The use of bottom-up attention in the Up-Down model provides a significant improvement over the best ResNet baseline across all question types, even though the ResNet baselines use almost twice as many convolutional layers.

	Yes/No	Number	Other	Overall
Prior [13]	61.20	0.36	1.17	25.98
Language-only [13]	67.01	31.55	27.37	44.26
d-LSTM+n-I [28, 13]	73.46	35.18	41.83	54.22
MCB [12, 13]	78.82	38.28	53.36	62.27
UPMC-LIP6	82.07	41.06	57.12	65.71
Athena	82.50	44.19	59.97	67.59
HDU-USYD-UNCC	84.50	45.39	59.01	68.09
Ours: Up-Down	86.60	48.64	61.15	70.34

Table 5. VQA v2.0 test-standard server accuracy, ranking our submission against recent published and unpublished work for each question type. Our approach, an ensemble of 30 models trained with different random seeds, outperforms all other leaderboard entries.

test-standard evaluation server, along with the previously published baseline results and the highest ranking other entries. At the time of submission (8 August 2017), we outperform all other test server submissions.

4.5. Qualitative Analysis

To qualitatively compare attention methodologies, we begin by visualizing the salient image regions for different words generated by the captioning model. It is common practice to directly visualize the attention weights α_t associated with word y_t emitted at the same time step [29, 52]. In Figure 5 we illustrate attention weights for the ResNet baseline and our full Up-Down model on the same image.

The selected image is unusual because it depicts a bathroom containing a couch but no toilet. Nevertheless, the baseline ResNet model hallucinates a toilet and therefore generates a poor quality caption. Attention visualizations suggest that, despite correctly identifying the sitting man, the baseline model is unable to successfully query the region of the image containing the couch, because there is no attention region capturing the man and the couch together. Presumably the model then relies on language priors to hallucinate a toilet as the most likely sitting affordance in a bathroom. In contrast, our Up-Down model correctly identifies the couch, despite the novel scene composition. Attention visualizations show that the Up-Down model is able to query the area of the image around the sitting man as there is a bounding box that encompasses this salient region. We include an example of VQA attention in Figure 6.

These and other examples, combined with the broadly-based improvements recorded in Table 2, suggest that considering attention at the level of objects and other salient image regions more effectively exposes the structure of the scene for subsequent processing. More specifically, when a candidate attention region corresponds to an object, or several related objects, all the visual concepts associated with those objects appear to be spatially co-located. In other words, our approach is able to consider all of the information pertaining to an object at once. This is also a nat-

Ours: Resnet – A man sitting on a *toilet* in a bathroom.



Ours: Up-Down – A man sitting on a *couch* in a bathroom.



Figure 5. Qualitative differences between attention methodologies in caption generation. For each generated word, we visualize the attended image region, outlining the region with the maximum attention weight in red. The ResNet baseline model is unable to query the image to find out where the man is sitting, because there is no attention region that captures the man and the couch together. Without a clear visual signal, the ResNet model presumably relies on language priors to hallucinate a toilet by looking at the sink. In contrast, our Up-Down model clearly identifies the out-of-context couch, generating a correct caption while also providing more interpretable attention weights.

Question: What room are they in? Answer: kitchen



Figure 6. Example VQA attention output, showing attention concentrated on the stovetop.

ural way for attention to be implemented. In the human visual system, the problem of integrating the separate features of objects in the correct combinations is known as the feature binding problem, and experiments suggest that attention plays a central role in the solution [44, 43].

5. Conclusion

In this paper we present a novel combined bottom-up and top-down visual attention mechanism. Qualitative evaluations suggest that the combination of bottom-up attention with the more traditional top-down approach more effectively exposes the structure of a scene, while also improving the interpretability of the resulting attention weights.

Applying this approach to image captioning, our model achieves state-of-the-art performance on the MSCOCO test server by a significant margin on all evaluation metrics. The broader applicability of the method is demonstrated with additional experiments on the VQA v2.0 dataset, which also achieve state-of-the-art performance. In both cases the underlying bottom-up attention features are the same. This suggests that the benefits of our approach may be captured by simply replacing pretrained CNN features with pretrained bottom-up attention features.

Acknowledgements

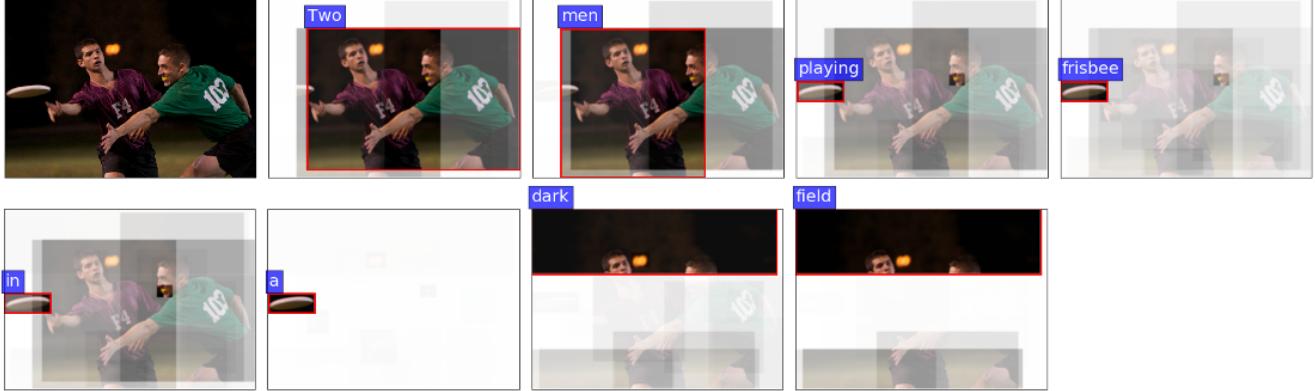
This research is partially supported by an Australian Government Research Training Program (RTP) Scholarship and by the Australian Research Council Centre of Excellence for Robotic Vision (project number CE140100016).

References

- [1] P. Anderson, B. Fernando, M. Johnson, and S. Gould. SPICE: Semantic propositional image caption evaluation. In *ECCV*, 2016. 6
- [2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh. VQA: Visual Question Answering. In *ICCV*, 2015. 6
- [3] T. J. Buschman and E. K. Miller. Top-down versus bottom-up control of attention in the prefrontal and posterior parietal cortices. *Science*, 315(5820):1860–1862, 2007. 1
- [4] X. Chen, T.-Y. L. Hao Fang, R. Vedantam, S. Gupta, P. Dollar, and C. L. Zitnick. Microsoft COCO captions: Data collection and evaluation server. *arXiv preprint arXiv:1504.00325*, 2015. 1
- [5] K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *EMNLP*, 2014. 5
- [6] M. Corbetta and G. L. Shulman. Control of goal-directed and stimulus-driven attention in the brain. *Nature Reviews Neuroscience*, 3(3):201–215, 2002. 1
- [7] Y. N. Dauphin, A. Fan, M. Auli, and D. Grangier. Language modeling with gated convolutional networks. *arXiv preprint arXiv:1612.08083*, 2016. 5
- [8] M. Denkowski and A. Lavie. Meteor Universal: Language Specific Translation Evaluation for Any Target Language. In *Proceedings of the EACL 2014 Workshop on Statistical Machine Translation*, 2014. 6
- [9] 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 *CVPR*, 2015. 3, 8
- [10] R. Egly, J. Driver, and R. D. Rafal. Shifting visual attention between objects and locations: evidence from normal and parietal lesion subjects. *Journal of Experimental Psychology: General*, 123(2):161, 1994. 2
- [11] H. Fang, S. Gupta, F. N. Iandola, R. Srivastava, L. Deng, P. Dollar, J. Gao, X. He, M. Mitchell, J. C. Platt, C. L. Zitnick, and G. Zweig. From captions to visual concepts and back. In *CVPR*, 2015. 8
- [12] A. Fukui, D. H. Park, D. Yang, A. Rohrbach, T. Darrell, and M. Rohrbach. Multimodal compact bilinear pooling for visual question answering and visual grounding. In *EMNLP*, 2016. 1, 2, 8
- [13] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh. Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering. In *CVPR*, 2017. 1, 6, 8
- [14] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 3, 7
- [15] K. He, X. Zhang, S. Ren, and J. Sun. Identity mappings in deep residual networks. *arXiv preprint arXiv:1603.05027*, 2016. 7
- [16] X. He and L. Deng. Maximum expected BLEU training of phrase and lexicon translation models. In *ACL*, 2012. 4
- [17] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 1997. 3
- [18] A. Jabri, A. Joulin, and L. van der Maaten. Revisiting visual question answering baselines. *arXiv preprint arXiv:1606.08390*, 2016. 2, 5
- [19] M. Jaderberg, K. Simonyan, A. Zisserman, and K. Kavukcuoglu. Spatial transformer networks. In *NIPS*, 2015. 2
- [20] J. Jin, K. Fu, R. Cui, F. Sha, and C. Zhang. Aligning where to see and what to tell: image caption with region-based attention and scene factorization. *arXiv preprint arXiv:1506.06272*, 2015. 2
- [21] J. Johnson, A. Karpathy, and L. Fei-Fei. Densecap: Fully convolutional localization networks for dense captioning. In *CVPR*, 2016. 2
- [22] A. Karpathy and F.-F. Li. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015. 6
- [23] V. Kazemi and A. Elqursh. Show, ask, attend, and answer: A strong baseline for visual question answering. *arXiv preprint arXiv:1704.03162*, 2017. 2, 5
- [24] C. Lin. Rouge: a package for automatic evaluation of summaries. In *ACL Workshop*, 2004. 6
- [25] T. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollar, and C. L. Zitnick. Microsoft COCO: Common objects in context. In *ECCV*, 2014. 6
- [26] S. Liu, Z. Zhu, N. Ye, S. Guadarrama, and K. Murphy. Optimization of image description metrics using policy gradient methods. *arXiv preprint arXiv:1612.00370v1*, 2016. 4, 8
- [27] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg. SSD: Single shot multibox detector. *arXiv preprint arXiv:1512.02325*, 2015. 2
- [28] J. Lu, X. Lin, D. Batra, and D. Parikh. Deeper lstm and normalized cnn visual question answering model. https://github.com/VT-vision-lab/VQA_LSTM_CNN, 2015. 8
- [29] J. Lu, C. Xiong, D. Parikh, and R. Socher. Knowing when to look: Adaptive attention via a visual sentinel for image captioning. In *CVPR*, 2017. 1, 2, 3, 8
- [30] J. Lu, J. Yang, D. Batra, and D. Parikh. Hierarchical question-image co-attention for visual question answering. In *NIPS*, 2016. 1, 2
- [31] J. Mao, W. Xu, Y. Yang, J. Wang, and A. L. Yuille. Deep captioning with multimodal recurrent neural networks (m-RNN). In *ICLR*, 2015. 8
- [32] K. Papineni, S. Roukos, T. Ward, and W. Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2002. 6
- [33] M. Pedersoli, T. Lucas, C. Schmid, and J. Verbeek. Areas of attention for image captioning. *arXiv preprint arXiv:xxxx*, 2015. 2

- [34] J. Pennington, R. Socher, and C. D. Manning. GloVe: Global Vectors for Word Representation. In *Proceedings of the Conference on Empirical Methods for Natural Language Processing (EMNLP)*, 2014. 6
- [35] J. Redmon, S. K. Divvala, R. B. Girshick, and A. Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 2
- [36] S. Ren, K. He, R. Girshick, and J. Sun. Faster R-CNN: Towards real-time object detection with region proposal networks. In *NIPS*, 2015. 2
- [37] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel. Self-critical sequence training for image captioning. In *CVPR*, 2017. 1, 2, 3, 4, 7, 8
- [38] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. Imagenet large scale visual recognition challenge. *IJCV*, 2015. 3, 7
- [39] B. J. Scholl. Objects and attention: The state of the art. *Cognition*, 80(1):1–46, 2001. 2
- [40] R. K. Srivastava, K. Greff, and J. Schmidhuber. Highway networks. *arXiv preprint arXiv:1505.00387v1*, 2015. 5
- [41] D. Teney, P. Anderson, X. He, and A. van den Hengel. Tips and tricks for visual question answering: Learnings from the 2017 challenge. *arXiv preprint arXiv:1708.02711*, 2017. 5, 6
- [42] D. Teney and A. van den Hengel. Zero-shot visual question answering. *arXiv preprint arXiv:1611.05546*, 2016. 5
- [43] A. Treisman. Perceptual grouping and attention in visual search for features and for objects. *Journal of Experimental Psychology: Human Perception and Performance*, 8(2):194, 1982. 8
- [44] A. M. Treisman and G. Gelade. A feature-integration theory of attention. *Cognitive Psychology*, 12:97–136, 1980. 8
- [45] J. R. Uijlings, K. E. Van De Sande, T. Gevers, and A. W. Smeulders. Selective search for object recognition. *IJCV*, 2013. 2
- [46] R. Vedantam, C. L. Zitnick, and D. Parikh. CIDEr: Consensus-based image description evaluation. In *CVPR*, 2015. 4, 6
- [47] I. Vendrov, R. Kiros, S. Fidler, and R. Urtasun. Order-embeddings of Images and Language. *arXiv preprint arXiv:1511.06361*, 2015. 3, 6
- [48] O. Vinyals, A. Toshev, S. Bengio, and D. Erhan. Show and tell: A neural image caption generator. In *CVPR*, 2015. 8
- [49] R. J. Williams. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Mach. Learn.*, 8(3-4):229–256, May 1992. 4
- [50] Q. Wu, C. Shen, L. Liu, A. Dick, and A. van den Hengel. What value do explicit high level concepts have in vision to language problems? In *CVPR*, 2016. 2
- [51] H. Xu and K. Saenko. Ask, attend and answer: Exploring question-guided spatial attention for visual question answering. In *ECCV*, 2016. 1, 2
- [52] 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. In *ICML*, 2015. 1, 2, 3, 8
- [53] Z. Yang, X. He, J. Gao, L. Deng, and A. J. Smola. Stacked attention networks for image question answering. In *CVPR*, 2016. 1, 2
- [54] Z. Yang, Y. Yuan, Y. Wu, R. Salakhutdinov, and W. W. Cohen. Review networks for caption generation. In *NIPS*, 2016. 1, 2, 8
- [55] T. Yao, Y. Pan, Y. Li, Z. Qiu, and T. Mei. Boosting image captioning with attributes. *arXiv preprint arXiv:1611.01646v1*, 2016. 2, 8
- [56] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo. Image captioning with semantic attention. In *CVPR*, 2016. 2, 8
- [57] Y. Yu, H. Ko, J. Choi, and G. Kim. End-to-end concept word detection for video captioning, retrieval, and question answering. In *CVPR*, 2017. 2
- [58] M. D. Zeiler. ADADELTA: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012. 6
- [59] Y. Zhu, O. Groth, M. Bernstein, and L. Fei-Fei. Visual7w: Grounded question answering in images. In *CVPR*, 2016. 1, 2

Two men playing frisbee in a dark field.



A group of people are playing a video game.



A brown sheep standing in a field of grass.

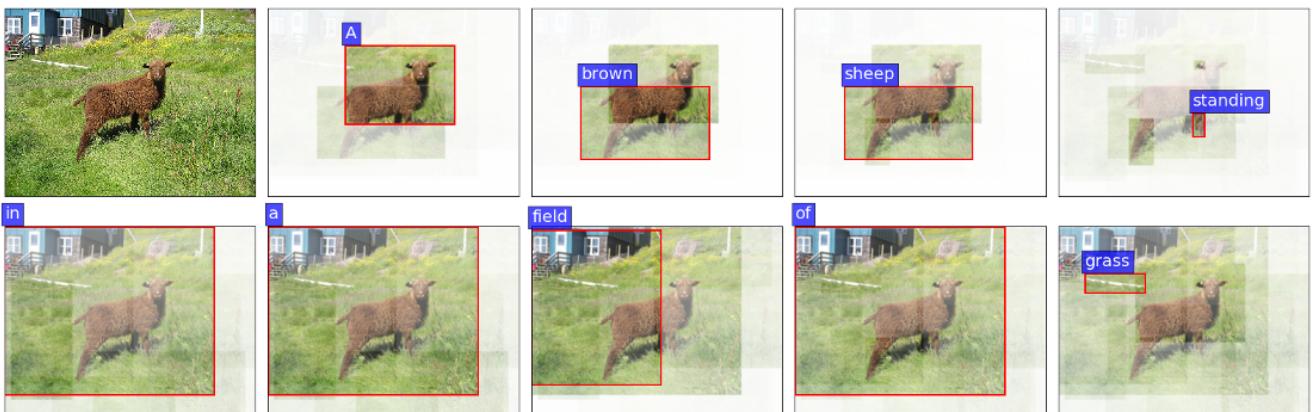
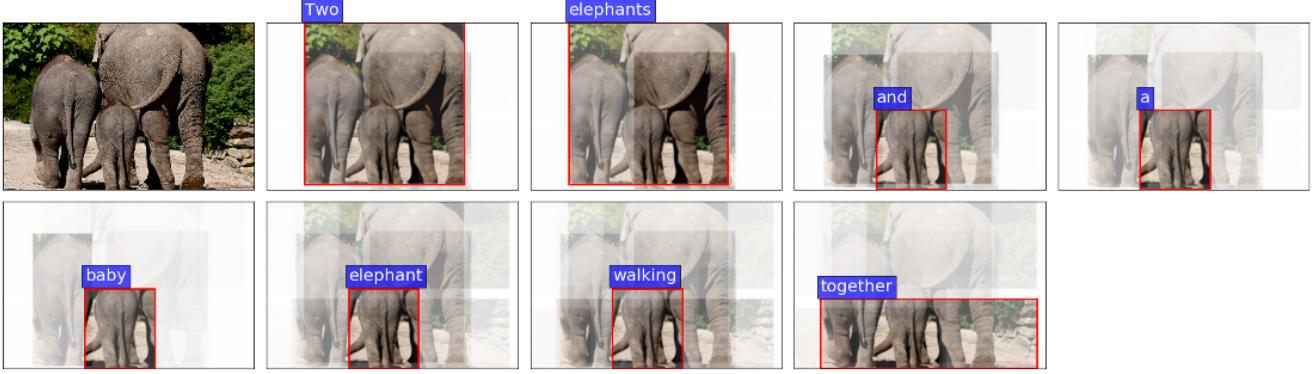


Figure 7. Examples of generated captions showing attended image regions. Notice the attention given to fine details, such as: (1) the frisbee and the green player's mouthguard when generating the word 'playing' in the first image, (2) the man's hands holding the game controllers in the middle image, and (3) the sheep's legs when generating the word 'standing' in the bottom image.

Two elephants and a baby elephant walking together.



A close up of a sandwich with a stuffed animal.



A dog laying in the grass with a frisbee.



Figure 8. Further examples of generated captions showing attended image regions. The first example suggests an understanding of spatial relationships when generating the word ‘together’. The middle image demonstrates the successful captioning of a compositionally novel scene. The bottom example is a failure case. The dog’s pose is mistaken for laying, rather than jumping – possibly due to poor salient region cropping that misses the dog’s head and feet.

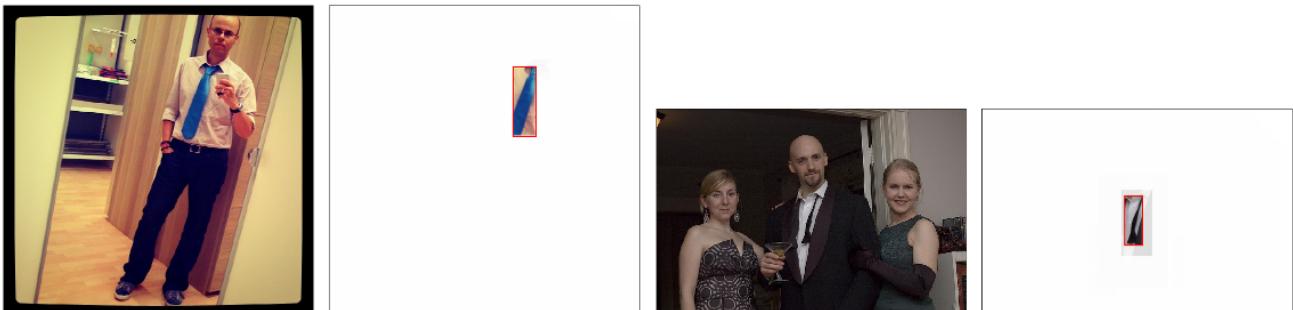
Question: What color is illuminated on the traffic light? Answer left: green. Answer right: red.



Question: What is the man holding? Answer left: phone. Answer right: controller.



Question: What color is his tie? Answer left: blue. Answer right: black.



Question: What sport is shown? Answer left: frisbee. Answer right: skateboarding.



Question: Is this the handlebar of a motorcycle? Answer left: yes. Answer right: no.



Figure 9. Further examples of successful visual question answering results, showing attended image regions.

Question: What is the name of the realty company? Answer left: none. Answer right: none.



Question: What is the bus number? Answer left: 2. Answer right: 23.



Question: How many cones have reflective tape? Answer left: 2. Answer right: 1.



Question: How many oranges are on pedestals? Answer left: 2. Answer right: 2.



Figure 10. Examples of visual question answering failure cases involving reading and counting.