

# Visual Coreference Resolution in Visual Dialog using Neural Module Networks

Satwik Kottur<sup>1,2\*</sup>, José M. F. Moura<sup>2</sup>, Devi Parikh<sup>1,3</sup>, Dhruv Batra<sup>1,3</sup>, and  
Marcus Rohrbach<sup>1</sup>

<sup>1</sup> Facebook AI Research, Menlo Park, USA

<sup>2</sup> Carnegie Mellon University, Pittsburgh, USA

<sup>3</sup> Georgia Institute of Technology, Atlanta, USA

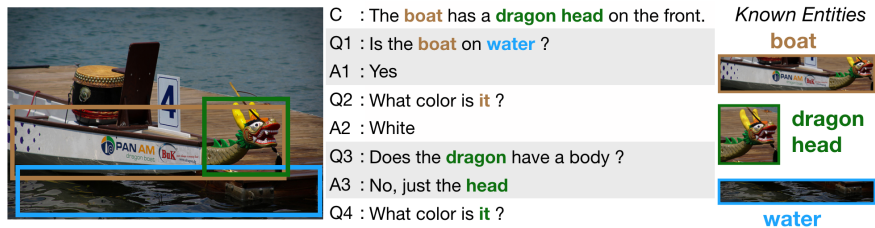
**Abstract.** Visual dialog entails answering a series of questions grounded in an image, using dialog history as context. In addition to the challenges found in visual question answering (VQA), which can be seen as one-round dialog, visual dialog encompasses several more. We focus on one such problem called *visual coreference resolution* that involves determining which words, typically noun phrases and pronouns, *co-refer* to the same entity/object instance in an image. This is crucial, especially for pronouns (e.g., ‘*it*’), as the dialog agent must first link it to a previous coreference (e.g., ‘*boat*’), and only then can rely on the visual grounding of the coreference ‘*boat*’ to reason about the pronoun ‘*it*’. Prior work (in visual dialog) models visual coreference resolution either (a) implicitly via a memory network over history, or (b) at a coarse level for the entire question; and not explicitly at a phrase level of granularity. In this work, we propose a neural module network architecture for visual dialog by introducing two novel modules—**Refer** and **Exclude**—that perform explicit, grounded, coreference resolution at a finer word level. We demonstrate the effectiveness of our model on MNIST Dialog, a visually simple yet coreference-wise complex dataset, by achieving near perfect accuracy, and on VisDial, a large and challenging visual dialog dataset on real images, where our model outperforms other approaches, and is more interpretable, grounded, and consistent qualitatively.

## 1 Introduction

The task of Visual Dialog [13, 44] involves building agents that ‘see’ (i.e. understand an image) and ‘talk’ (i.e. communicate this understanding in a dialog). Specifically, it requires an agent to answer a sequence of questions about an image, requiring it to reason about both the image and the past dialog history. For instance, in Fig. 1, to answer ‘*What color is it?*’, the agent needs to reason about the history to know what ‘*it*’ refers to and the image to find out the color. This generalization of visual question answering (VQA) [8] to dialog takes a step closer to real-world applications (aiding visually impaired users, intelligent home

---

\* Work partially done as an intern at Facebook AI Research



**Fig. 1:** Our model begins by grounding entities in the caption (C), *boat* (brown) and *dragon head* (green), and stores them in a pool for future coreference resolution in the dialog (right). When asked ‘Q1: Is the *boat* on *water*?’, it identifies that the *boat* (known entity) and *water* (unknown entity) are crucial to answer the question. It then grounds the novel entity *water* in the image (blue), but resolves *boat* by referring back to the pool and reusing the available grounding from the caption, before proceeding with further reasoning. Thus, our model explicitly resolves coreferences in visual dialog.

assistants, natural language interfaces for robots) but simultaneously introduces new modeling challenges at the intersection of vision and language. The particular challenge we focus on in this paper is that of *visual coreference resolution* in visual dialog. Specifically, we introduce a new model that performs explicit visual coreference resolution and interpretable entity tracking in visual dialog.

It has long been understood [18, 48, 34, 50] that humans use *coreferences*, different phrases and short-hands such as pronouns, to refer to the same entity or referent in a single text. In the context of visually grounded dialog, we are interested in referents which are in the image, e.g. an object or person. All phrases in the dialog which refer to the same entity or referent in the image are called visual coreferences. Such coreferences can be noun phrases such as ‘*a dragon head*’, ‘*the head*’, or pronouns such as ‘*it*’ (Fig. 1). Especially when trying to answer a question that contains an anaphora, for instance the pronoun ‘*it*’, which refers to its full form (the antecedent) ‘*a dragon head*’, it is necessary to *resolve* the coreference on the language side and ground it to the underlying visual referent. More specifically, to answer the question ‘*What color is it?*’ in Fig. 1, the model must correctly identify which object ‘*it*’ refers to, in the given context. Notice that a word or phrase can refer to different entities in different contexts, as is the case with ‘*it*’ in this example. Our approach to explicitly resolve visual coreferences is inspired from the functionality of variables or memory in a computer program. In the same spirit as how one can refer back to the contents of variables at a later time in a program without explicitly re-computing them, we propose a model which can refer back to entities from previous rounds of dialog and reuse the associated information; and in this way resolve coreferences.

Prior work on VQA [31, 15, 4] has (understandably) largely ignored the problem of visual coreference resolution since individual questions asked in isolation rarely contain coreferences. In fact, recent empirical studies [3, 22, 51, 17] suggest that today’s vision and language models seem to be exploiting dataset-level statistics and perform poorly at grounding entities into the correct pixels. In

contrast, our work aims to explicitly reason over past dialog interactions by referring back to previous references. This allows for increased interpretability of the model. As the dialog progresses (Fig. 1), we can inspect the pool of entities known to the model, and also visualize which entity a particular phrase in the question has been resolved to. Moreover, our explicit entity tracking model has benefits even in cases that may not strictly speaking require coreference resolution. For instance, by explicitly referring ‘*dragon*’ in Q3 (Fig. 1) back to a known entity, the model is consistent with itself and (correctly) grounds the phrase in the image. We believe such consistency in model outputs is a strongly desirable property as we move towards human-machine interaction in dialog systems.

Our main technical contribution is a neural module network architecture for visual dialog. Specifically, we propose two novel modules—**Refer** and **Exclude**—that perform explicit, grounded, coreference resolution in visual dialog. In addition, we propose a novel way to handle captions using neural module networks at a word-level granularity finer than a traditional sentence-level encoding. We show quantitative benefits of these modules on a reasoning-wise complicated but visually simple MNIST dialog dataset [41], where achieve near perfect accuracy. On the visually challenging VisDial dataset [13], our model not only outperforms other approaches but also is more interpretable by construction and enables word-level coreference resolution. Furthermore, we qualitatively show that our model is (a) more interpretable (a user can inspect which entities were detected and tracked as the dialog progresses, and which ones were referred to for answering a specific question), (b) more grounded (where the model looked to answer a question in the dialog), (c) more consistent (same entities are considered across rounds of dialog).

## 2 Related Work

We discuss: (a) existing approaches to visual dialog, (b) related tasks such as visual grounding and coreference resolution, and (c) neural module networks.

**Visual Dialog.** Though the origins of visual dialog can be traced back to [47, 16], it was largely formalized by [13, 44] who collected human annotated datasets for the same. Specifically, [13] paired annotators to collect free-form natural-language questions and answers, where the questioner was instructed to ask questions to help them imagine the hidden scene (image) better. On the other hand, dialogs from [44] are more goal driven and contain yes/no questions directed towards identifying a secret object in the image. The respective follow up works used reinforcement learning techniques to solve this problem [14, 43]. Other approaches to visual dialog include transferring knowledge from a discriminatively trained model to a generative dialog model [30], using attention networks to solve visual coreferences [41], and more recently, a probabilistic treatment of dialogs using conditional variational autoencoders [33]. Amongst these, [41] is the closest to this work, while [30, 33] are complementary. To solve visual coreferences, [41] relies on global visual attentions used to answer previous questions. They store these attention maps in a memory against keys based on

textual representations of the entire question and answer, along with the history. In contrast, operating at a finer word-level granularity within each question, our model can resolve different phrases of a question, and ground them to different parts of the image, a core component in correctly understanding and grounding coreferences. E.g., ‘*A man and woman in a car. Q: Is he or she driving?*’, which requires resolving ‘*he*’ and ‘*she*’ individually to answer the question.

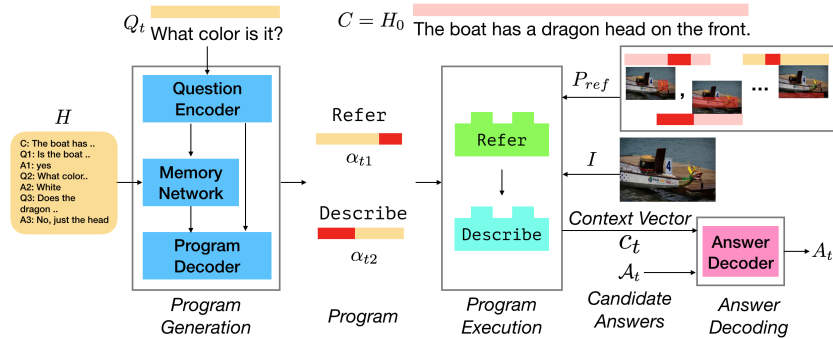
**Grounding language in images and video.** Most works in this area focus on the specific task of localizing a textual referential expression in the image [21, 25, 32, 36, 39, 45, 50] or video [38, 27, 49, 7]. Similar to these works, one component of our model aims to localize words and phrases in the image. However, the key difference is that if the phrase being grounded is an anaphora (e.g., ‘*it*’, ‘*he*’, ‘*she*’, etc.), our model first resolves it explicitly to a known entity, and then grounds it by borrowing the referent’s visual grounding.

**Coreference resolution.** The linguistic community defines coreference resolution as the task of clustering phrases, such as noun phrases and pronouns, which refer to the same entity in the world (see, for example, [10]). The task of visual coreference resolution links the coreferences to an entity in the visual data. For example, [37] links character mentions in TV show descriptions with their occurrence in the video, while [25] links text phrases to objects in a 3D scene. Different from these works, we predict a program for a given natural language question about an image, which then tries to resolve any existing coreferences, to then answer the question. An orthogonal direction is to generate language while jointly grounding and resolving coreferences – e.g., [40] explore this for movie descriptions. While out of scope for this work, it is an interesting direction for future work in visual dialog, especially when generating questions.

**Neural Module Networks** [6] are an elegant class of models where an instance-specific architecture is composed from neural ‘modules’ (or building blocks) that are shared across instances. The high-level idea is inspired by ‘options’ or sub-tasks in hierarchical RL. They have been shown to be successful for visual question answering in real images and linguistic databases [5] and for more complex reasoning tasks in synthetic datasets [23, 20]. For this, [23, 20] learn program prediction and module parameters jointly, end-to-end. Within this context, our work generalizes the formulation in [20] from VQA to visual dialog by introducing a novel module to perform explicit visual coreference resolution.

### 3 Approach

Recall that visual dialog [13] involves answering a question  $Q_t$  at the current round  $t$ , given an image  $I$ , and the dialog history (including the image caption)  $H = (\underbrace{C}_{H_0}, \underbrace{(Q_1, A_1)}_{H_1}, \dots, \underbrace{(Q_{t-1}, A_{t-1})}_{H_{t-1}})$ , by ranking a list of 100 candidate answers  $\mathcal{A}_t = \{A_t^{(1)}, \dots, A_t^{(100)}\}$ . As a key component for building better visual dialog agents, our model explicitly resolves visual coreferences in the current question, if any.



**Fig. 2:** Overview of our model architecture. The question  $Q_t$  (orange bar) is encoded along with the history  $H$  through a memory augmented question encoder, using which a program (Refer Describe) is decoded. For each module in the program, an attention  $\alpha_{ti}$  over  $Q_t$  is also predicted, used to compute the text feature  $x_{t,xt}$ . For  $Q_t$ , attention is over ‘it’ for Refer and ‘What color’ for Describe, respectively (orange bars with red attention). Refer module uses the coreference pool  $P_{ref}$ , a dictionary of all previously seen entities with their visual groundings, resolves ‘it’, and borrows the referent’s visual grounding (boat in this case). Finally, Describe extracts the ‘color’ to produce  $c_t$  used by a final decoder to pick the answer  $A_t$  from the candidate pool  $A_t$ .

Towards this end, our model first identifies relevant words or phrases in the current question that refer to entities in the image (typically objects and attributes). The model also predicts whether each of these has been mentioned in the dialog so far. Next, if these are novel entities (unseen in the dialog history), they are localized in the image before proceeding, and for seen entities, the model predicts the (first) relevant coreference in the conversation history, and retrieves its corresponding visual grounding. Therefore, as rounds of dialog progress, the model collects unique entities and their corresponding visual groundings, and uses this *reference pool* to resolve any coreferences in subsequent questions.

Our model has three broad components: (a) *Program Generation* (Sec. 3.3), where a reasoning pathway, as dictated by a *program*, is predicted for the current question  $Q_t$ , (b) *Program Execution* (Sec. 3.4), where the predicted program is executed by dynamically connecting neural modules [5, 6, 20] to produce a *context* vector summarizing the semantic information required to answer  $Q_t$  from the context  $(I, H)$ , and lastly, (c) *Answer Decoding* (Sec. 3.4), where the context vector  $c_t$  is used to obtain the final answer  $\hat{A}_t$ . We begin with a general characterization of neural modules used for VQA in Sec. 3.1 and then discuss our novel modules for coreference resolution (Sec. 3.2) with details of the reference pool. After describing the inner working of the modules, we explain each of the above three components of our model.

### 3.1 Neural Modules for Visual Question Answering

The main technical foundation of our model is the neural module network (NMN) [6]. In this section, we briefly recap NMNs and more specifically, the attentional modules from [20]. In the next section, we discuss novel modules we propose to handle additional challenges in visual dialog.

For a module  $m$ , let  $x_{vis}$  and  $x_{txt}$  be the input image and text embeddings, respectively. In particular, the image embeddings  $x_{vis}$  are spatial activation maps of the image  $I$  from a convolutional neural network. The text embedding  $x_{txt}$  is computed as a weighted sum of embeddings of words in the question  $Q_t$  using the soft attention weights  $\alpha$  predicted by a program generator for module  $m$  (more details in Sec. 3.3). Further, let  $\{a_i\}$  be the set of  $n_m$  single-channel spatial maps corresponding to the spatial image embeddings, where  $n_m$  is the number of attention inputs to  $m$ . Denoting the module parameters with  $\theta_m$ , a neural module  $m$  is essentially a parametric function  $y = f_m(x_{vis}, x_{txt}, \{a_i\}_{i=1}^{n_m}, \theta_m)$ . The output from the module  $y$  can either be a spatial image attention map (denoted by  $a$ ) or a context vector (denoted by  $c$ ), depending on the module. The output spatial attention map  $a$  feeds into next level modules while a context vector  $c$  is used to obtain the final answer  $A_t$ . The upper part of Tab. 1 lists modules we adopt from prior work, with their functional forms. We shortly summarize their behavior. A **Find** module localizes objects or attributes by producing an attention over the image. The **Relocate** module takes in an input image attention and performs necessary spatial relocations to handle relationships like ‘next to’, ‘in front of’, ‘beside’, etc. Intersection or union of attention maps can be obtained using **And** and **Or**, respectively. Finally, **Describe**, **Exist**, and **Count** input an attention map to produce the context vector by describing an attribute, checking for existence, or counting, respectively, in the given input attention map. As noted in [20], these modules are designed and named for a potential ‘atomic’ functionality. However, we do not enforce this explicitly and let the modules discover their expected behavior by training in an end-to-end manner.

### 3.2 Neural Modules for Coreference Resolution

We now introduce novel components and modules to handle visual dialog.

**Reference Pool** ( $P_{ref}$ ). The role of the reference pool is to keep track of entities seen so far in the dialog. Thus, we design  $P_{ref}$  to be a dictionary of key-value pairs  $(x_{txt}, a)$  for all the **Find** modules instantiated while answering previous questions  $(Q_i)_{i=1}^{t-1}$ . Recall that **Find** localizes objects/attributes specified by  $x_{txt}$ , and thus by storing each output attention map  $y$ , we now have access to all the entities mentioned so far in the dialog with their corresponding visual groundings. Interestingly, even though  $x_{txt}$  and  $y$  are intermediate outputs from our model, both are easily interpretable, making our reference pool a *semantic dictionary*. To the best of our knowledge, our model is the first to attempt explicit, interpretable coreference resolution in visual dialog. While [41] maintains a dictionary similar to  $P_{ref}$ , they do not consider word/entity level coreferences

Name	Inputs	Output	Function
<b>Neural Modules for VQA [20]</b>			
<b>Find</b>	$x_{vis}, x_{txt}$	attention	$y = \text{conv}_2(\text{conv}_1(x_{vis} \odot Wx_{txt}))$
<b>Relocate</b>	$a, x_{vis}, x_{txt}$	attention	$\tilde{y} = W_1 \text{sum}(a \odot x_{vis})$ $y = \text{conv}_2(\text{conv}_1(x_{vis}) \odot \tilde{y} \odot W_2 x_{txt})$
<b>And</b>	$a_1, a_2$	attention	$y = \min\{a_1, a_2\}$
<b>Or</b>	$a_1, a_2$	attention	$y = \max\{a_1, a_2\}$
<b>Exist</b>	$a, x_{vis}, x_{txt}$	context	$y = W^T \text{vec}(a)$
<b>Describe</b>	$a, x_{vis}, x_{txt}$	context	$y = W_1^T (W_2 \text{sum}(a \odot x_{vis}) \odot W_3 x_{txt})$
<b>Count</b>	$a, x_{vis}, x_{txt}$	context	$y = W_1^T ([\text{vec}(a), \max\{a\}, \min\{a\}])$
<b>Neural Modules for Coreference resolution (Ours)</b>			
<b>Not</b>	$a$	attention	$y = \text{norm}_{L_1}(1 - a)$
<b>Refer</b>	$x_{txt}, P_{ref}$	attention	(see text for details, (3))
<b>Exclude</b>	$a, x_{vis}, x_{txt}$	attention	$y = \text{And}[\text{Find}[x_{vis}, x_{txt}], \text{Not}[a]]$

**Table 1:** Neural modules used in our work for visual dialog, along with their inputs, outputs, and function formulations. The upper portion contains modules from prior work used for visual question answering, while the bottom portion lists our novel modules designed to handle additional challenges in visual dialog.

nor do their keys lend similar interpretability as ours. With  $P_{ref} = \{(x_p^{(i)}, a_p^{(i)})\}_i$  as input to **Refer**, we can now resolve references in  $Q_t$ .

**Refer Module.** This novel module is responsible for resolving references in the question  $Q_t$  and ground them in the conversation history  $H$ . To enable grounding in dialog history, we generalize the above formulation to give the module access to a pool of references  $P_{ref}$  of previously identified entities. Specifically, **Refer** only takes the text embedding  $x_{txt}$  and the reference pool  $P_{ref}$  as inputs, and resolves the entity represented by  $x_{txt}$  in the form of a soft attention  $\alpha$  over  $Q_t$ . In this section after introducing  $P_{ref}$ . For the example shown in Fig. 2,  $\alpha$  for **Refer** attends to ‘it’, indicating the phrase it is trying to resolve.

At a high level, **Refer** treats  $x_{txt}$  as a ‘query’ and retrieves the most likely match from  $P_{ref}$  as measured by some similarity with respect to keys  $\{x_p^{(i)}\}_i$  in  $P_{ref}$ . The associated image attention map of the best match is used as the visual grounding for the phrase that needed resolution (i.e. ‘it’). More concretely, we first learn a *scoring network* which when given a query  $x_{txt}$  and a possible candidate  $x_p^{(i)}$ , returns a scalar value  $s_i$  indicating how likely these text features refer to the same entity (1). To enable **Refer** to consider the sequential nature of dialog when assessing a potential candidate, we additionally provide  $\Delta_i t$ , a measure of the ‘distance’ of a candidate  $x_p^{(i)}$  from  $x_{txt}$  in the dialog history, as input to the scoring network.  $\Delta_i t$  is formulated as the absolute difference between the round of  $x_{txt}$  (current round  $t$ ) and the round when  $x_p^{(i)}$  was first mentioned. Collecting these scores from all the candidates, we apply a softmax function to compute contributions  $\tilde{s}_i$  from each entity in the pool (2). Finally,

we weigh the corresponding attention maps via these contributions to obtain the visual grounding  $a_{out}$  for  $x_{txt}$  (3).

$$s_i = \text{MLP}([x_{txt}, x_p^{(i)}, \Delta_i t]) \quad (1)$$

$$\tilde{s}_i = \text{Softmax}(s_i) \quad (2)$$

$$a_{out} = \sum_{i=1}^{|P_{ref}|} \tilde{s}_i a_p^{(i)} \quad (3)$$

**Not Module.** Designed to focus on regions of the image **not** attended by the input attention map  $a$ , it outputs  $y = \text{norm}_{L_1}(1-a)$ , where  $\text{norm}_{L_1}(\cdot)$  normalizes the entries to sum to one. This module is used in **Exclude**, described next.

**Exclude Module.** To handle questions like ‘*What other red things are present?*’, which seek other objects/attributes in the image than those specified by an input attention map  $a$ , we introduce yet another novel module – **Exclude**. It is constructed using **Find**, **Not**, and **And** modules as  $y = \text{And}[\text{Find}[x_{txt}, x_{vis}], \text{Not}[a]]$ , where  $x_{txt}$  is the text feature input to the **Exclude** module, for example, ‘*red things*’. More explicitly, **Find** first localizes all objects instances/attributes in the image. Next, we focus on regions of the image other than those specified by  $a$  using **Not**[ $a$ ]. Finally, the above two outputs are combined via **And** to obtain the output  $y$  of the **Exclude** module.

### 3.3 Program Generation

A *program* specifies the network layout for the neural modules for a given question  $Q_t$ . Following [20], it is serialized through the reverse polish notation (RPN) [11]. This serialization helps us convert a hard, structured prediction problem into a more tractable sequence prediction problem. In other words, we need a program predictor to output a series of module tokens in order, such that a valid layout can be retrieved from it. There are two primary design considerations for our predictor. First, in addition to the program, our predictor must also output a soft attention  $\alpha_{ti}$ , over the question  $Q_t$ , for every module  $m_i$  in the program. This attention is responsible for the *correct* module instantiation in the current context. For example, to answer the question ‘*What color is the cat sitting next to the dog?*’, a **Find** module instance attending to ‘*cat*’ qualitatively serves a different purpose than one attending to ‘*dog*’. This is implemented by using the attention over  $Q_t$  to compute the text embedding  $x_{txt}$  that is directly fed as an input to the module during execution. Second, to decide whether an entity in  $Q_t$  has been seen before in the conversation, it must be able to ‘peek’ into the history  $H$ . Note that this is unique to our current problem and does not exist in [20]. To this effect, we propose a novel augmentation of attentional recurrent neural networks [9] with memory [46] to address both the requirements (Fig. 2).

The program generation proceeds as follows. First, each of the words in  $Q_t$  are embedded to give  $\{w_{ti}\}_{i=1}^T$ , where  $T$  denotes the number of tokens in  $Q_t$ . We then use a *question encoder*, a multi-layer LSTM, to process  $w_{ti}$ ’s, resulting in a sequence of hidden states  $\{\hat{w}_{ti}\}_{i=1}^T$  (4). Notice that the last hidden state  $h_T$  is the question encoding, which we denote with  $q_t$ . Next, each piece of history  $(H_i)_{i=0}^{t-1}$  is processed in a similar way by a *history encoder*, which is a multi-layer LSTM



akin to the question encoder. This produces encodings  $(h_i)_{i=0}^{t-1}$  (5) that serve as memory units to help the program predictor ‘peek’ into the conversation history. Using the question encoding  $q_t$ , we attend over the history encodings  $(h_i)_{i=0}^{t-1}$ , and obtain the history vector  $\hat{h}_t$  (6). The history-agnostic question encoding  $q_t$  is then fused with the history vector  $\hat{h}_t$  via a fully connected layer to give a history-aware question encoding  $\hat{q}_t$  (7), which is fed into the *program decoder*.

#### Question Encoder

$$\{\hat{w}_{ti}\} = \text{LSTM}(\{w_{ti}\}) \quad (4)$$

$$q_t = \hat{w}_{tT}$$

#### History Memory

$$\hat{h}_i = \text{LSTM}(h_i) \quad (5)$$

$$\beta_{ti} = \text{Softmax}(q_t^T \hat{h}_i)$$

$$\hat{h}_t = \sum_{i=0}^{t-1} \beta_{ti} \hat{h}_i \quad (6)$$

$$\hat{q}_t = \text{MLP}([q_t, \hat{h}_t]) \quad (7)$$

#### Program Decoder

$$\tilde{u}_{ti}^{(j)} = \text{Linear}([\hat{w}_{tj}, d_{ti}])$$

$$u_{ti}^{(j)} = v^T \tanh(\tilde{u}_{ti}^{(j)})$$

$$\alpha_{ti}^{(j)} = \text{Softmax}(u_{ti}^{(j)})$$

$$e_{ti} = \sum_{j=1}^T \alpha_{ti}^{(j)} \hat{w}_{tj} \quad (8)$$

$$\tilde{e}_{ti} = \text{MLP}([e_{ti}, d_{ti}]) \quad (9)$$

$$p(m_i | \{m_k\}_{k=1}^{i-1}, Q_t, H) \\ = \text{Softmax}(\tilde{e}_{ti}) \quad (10)$$

The decoder is another multi-layer LSTM network (with hidden states  $\{d_{ti}\}$ ) which, at every time step  $i$ , produces a soft attention map  $\alpha_{ti}$  over the input sequence  $(Q_t)$  [9]. This soft attention map for each module is used to compute the corresponding text embedding,  $x_{t,xt} = \sum_j \alpha_{ti}^{(j)} w_{tj}$ . Finally, to predict a module token  $m_i$  at time step  $i$ , a weighted sum of encoder hidden states  $e_{ti}$  (8) and the history-aware question vector  $\hat{q}_t$  are combined via another fully-connected layer (9), followed by a softmax to give a distribution  $P(m_i | \{m_k\}_{k=1}^{i-1}, Q_t, H)$  over the module tokens (10). During training, we minimize the cross-entropy loss  $\mathcal{L}_Q^{prog}$  between this predicted distribution and the ground truth program tokens. Fig. 2 outlines the schematics of our program generator.

**Modules on captions.** As the image caption  $C$  is also a part of the dialog (history  $H_0$  at round 0), it is desirable to track entities from  $C$  via the coreference pool  $P_{ref}$ . To this effect, we propose a novel extension of neural module networks to captions by using an auxiliary task that checks the alignment of a (caption, image) pair. First, we learn to predict a program from  $C$ , different from those generated from  $Q_t$ , by minimizing the negative log-likelihood  $\mathcal{L}_C^{prog}$ , akin to  $\mathcal{L}_Q^{prog}$ , of the ground truth caption program. Next, we execute the caption program on two images  $I^+ = I$  and  $I^-$  (a random image from the dataset), to produce caption context vectors  $c_C^+$  and  $c_C^-$ , respectively. Note that  $c_C^+$  and  $c_C^-$  are different from the context vector  $c_t$  produced from execution of the question program. Finally, we learn a binary classifier on top to output classes  $+1/-1$  for  $c_C^+$  and  $c_C^-$ , respectively, by minimizing the binary cross entropy loss  $\mathcal{L}_C^{aux}$ . The intuition behind the auxiliary task is: to rightly classify aligned  $(C, I^+)$  from misaligned  $(C, I^-)$ , the modules will need to localize and focus on salient entities in the caption. These entities (specifically, outputs from **Find** in the caption program) are then collected in  $P_{ref}$  for explicit coreference resolution on  $Q_t$ .

**Entities in answers.** Using an analogous argument as above, answers from the previous rounds  $\{A_i\}_{i=1}^{t-1}$  could have entities necessary to resolve coreferences in  $Q_t$ . For example, ‘ $Q$ : What is the boy holding?  $A$ : A ball.  $Q$ : What color is it?’ requires resolving ‘it’ with the ‘ball’ mentioned in the earlier answer. To achieve this, at the end of round  $t - 1$ , we encode  $H_{t-1} = (Q_{t-1}, A_{t-1})$  as  $h_t^{ref}$  using a multi-layer LSTM, obtain the last image attention map  $a$  fed to the last module in the program that produced the context vector  $c_t$ , and add  $(h_t^{ref}, a)$  as an additional candidate to the reference pool  $P_{ref}$ . Notice that  $h_t^{ref}$  contains the information about the answer  $A_{t-1}$  in the context of the question  $Q_{t-1}$ , while  $a$  denotes the image attention which was the last crucial step in arriving at  $A_{t-1}$  in the earlier round. In resolving coreferences in  $Q_t$ , if any, all the answers from previous rounds now become potential candidates by virtue of being in  $P_{ref}$ .

### 3.4 Other Model Components

**Program Execution.** This component takes the generated program and associated text features  $x_{t,xt}$  for each participating module, and executes it. To do so, we first deserialize the given program from its RPN to a hierarchical module layout. Next, we arrange the modules dynamically according to the layout, giving us the network to answer  $Q_t$ . At this point, the network is a simple feed-forward neural network, where we start the computation from the leaf modules and feed outputs activations from modules at one layer as inputs into modules at the next layer (see Fig. 2). Finally, we feed a context vector  $c_t$  produced from the last module into the next answer decoding component.

**Answer Decoding.** This is the last component of our model that uses the context vector  $c_t$  to score answers from a pool of candidates  $\mathcal{A}_t$ , based on their correctness. The answer decoder: (a) encodes each candidate  $A_t^{(i)} \in \mathcal{A}_t$  with a multi-layer LSTM to obtain  $o_t^{(i)}$ , (b) computes a score via a dot product with the context vector, i.e.,  $c_t^T o_t^{(i)}$ , and (c) applies a softmax activation to get a distribution over the candidates. During training, we minimize the negative log-likelihood  $\mathcal{L}_A^{dec}$  of the ground truth answer  $A_t^{gt}$ . At test time, the candidate with the maximum score is picked as  $\mathcal{A}_t$ . Using nomenclature from [13], this is a *discriminative* decoder. Note that our approach is not limited to a discriminative decoder, but can also be used with a *generative* decoder (see supplement).

**Training Details.** Our model components have fully differentiable operations within them. Thus, to train our model, we combine the supervised loss terms from both program generation  $\{\mathcal{L}_Q^{prog}, \mathcal{L}_C^{prog}, \mathcal{L}_C^{aux}\}$  and answer decoding  $\{\mathcal{L}_A^{dec}\}$ , and minimize the sum total loss  $\mathcal{L}^{total}$ .

## 4 Experiments

We first show results on the synthetic MNIST Dialog dataset [41], designed to contain complex coreferences across rounds while being relatively easy textually

and visually. It is important to resolve these coreferences accurately in order to do well on this dataset, thus stress testing our model. We then experiment with a large visual dialog dataset on real images, VisDial [13], which offers both linguistic and perceptual challenge in resolving visual coreferences and grounding them in the image. Implementation details are in the supplement.

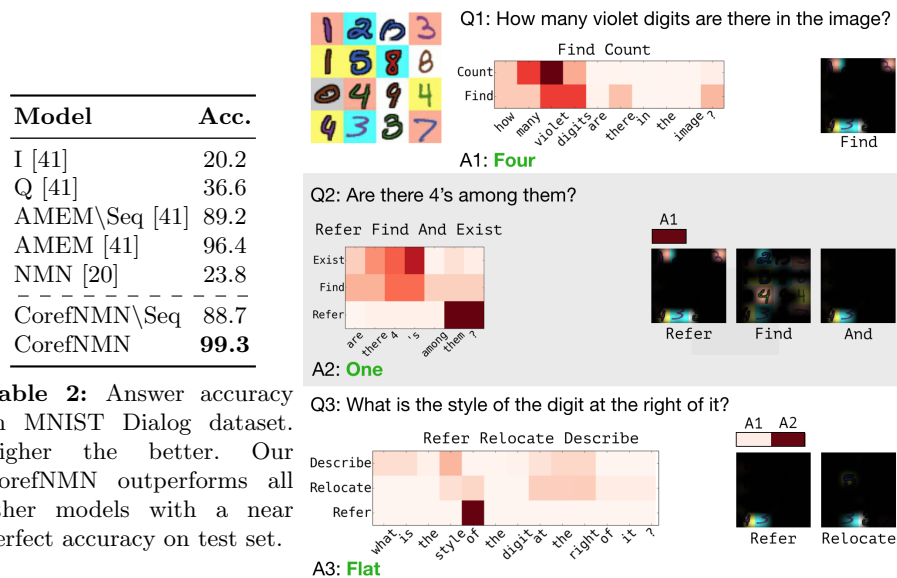
#### 4.1 MNIST Dialog Dataset

**Dataset.** The dialogs in the MNIST dialog dataset [41] are grounded in images composed from a  $4 \times 4$  grid of MNIST digits [26]. Digits in the grid have four attributes—digit class (0 – 9), color, stroke, and background color. Each dialog has 10 question-answer pairs, where the questions are generated through language templates, and the answers are single words. Further, the questions are designed to query attributes of target digit(s), count digits with similar attributes, etc., all of which need tracking of the target digits(s) by resolving references across dialog rounds. Thus, coreference resolution plays a crucial part in the reasoning required to answer the question, making the MNIST dataset both interesting and challenging (Fig. 3). The dataset contains  $30k$  training,  $10k$  validation, and  $10k$  test images, with three 10-round dialogs for each image.

**Models and baselines.** Taking advantage of single-word answers in this dataset, we simplify our answer decoder to be a  $N$ -way classifier, where  $N$  is the number of possible answers. Specifically, the context vector  $c_t$  now passes through a fully connected layer of size  $N$ , followed by softmax activations to give us a distribution over possible answer classes. At training time, we minimize the cross-entropy  $\mathcal{L}_A^{dec}$  of the predicted answer distribution with the ground truth answer, at every round. Note that single-word answers also simplify evaluation as answer accuracy can now be used to compare different models. We further simplify our model by removing the memory augmentation to the program generator, i.e.,  $\hat{q}_t = q_t$  (7), and denote it as CorefNMN. In addition to the full model, we also evaluate an ablation, CorefNMN\Seq, without  $\Delta_i t$  that additionally captured sequential nature of dialog (see **Refer** description). We compete against the explicit reasoning model (NMN) [20] and a comprehensive set of baselines AMEM, image-only (I), and question-only (Q), all from [41].

**Supervision.** In addition to the ground truth answer, we also need program supervision for questions to learn the program generation. For each of the 5 ‘types’ of questions, we manually create one program which we apply as supervision for all questions of the corresponding type. The type of question is provided with the question. Note that our model needs program supervision only while training, and uses predictions from program generator at test time.

**Results.** Tab. 2 shows the results on MNIST dataset. The following are the key observations: (a) The text-only Q (36.6%) and image-only I (20.2%) do not perform well, perhaps as expected as MNIST Dialog needs resolving strong coreferences to arrive at the correct answer. For the same reason, NMN [20] has a low accuracy of 23.8%. Interestingly, Q outperforms NMN by around 13% (both use question and image, but not history), possibly due to the explicit reasoning



**Table 2:** Answer accuracy on MNIST Dialog dataset. Higher the better. Our CorefNMN outperforms all other models with a near perfect accuracy on test set.

**Fig. 3:** Illustration of explicit coreference resolution reasoning of our model on the MNIST dialog dataset. For each question, a program and corresponding attentions ( $\alpha$ 's) over question words (hot matrix on the left) is predicted. A layout is unpacked from the program, and modules are connected to form a feed-forward network used to answer the question, shown in green to indicate correctness. We also visualize output attention maps (right) from each participating module. Specifically, in Q1 and Q2, **Find** localizes all violet digits and 4's, respectively (indicated by the corresponding  $\alpha$ ). In Q2, **Refer** resolves 'them' and borrows the visual grounding from previous question.

nature of NMN prohibiting it from capturing the statistic dataset priors. (b) Our CorefNMN outperforms all other models with near perfect accuracy of 99.3%. Examining the failure cases reveals that most of the mistakes made by CorefNMN was due to misclassifying qualitatively hard examples from the original MNIST dataset. (c) Factoring the sequential nature of the dialog additionally in the model is beneficial, as indicated by the 10.6% improvement in CorefNMN, and 7.2% in AMEM. Intuitively, phrases with multiple potential referents, more often than not, refer to the most recent referent, as seen in Fig. 1, where 'it' has to be resolved to the closest referent in history. Fig. 3 shows a qualitative example.

## 4.2 VisDial v0.9 Dataset

**Dataset.** The VisDial dataset [13] is a crowd-sourced dialog dataset on COCO images [28], with free-form answers. The publicly available VisDial v0.9 contains 10-round dialogs on around 83k training images, and 40k validation images. VisDial was collected from pairs of human workers, by instructing one of them to ask questions in a live chat interface to help them imagine the scene better.

Model	MRR	R@1	R@5	R@10	Mean
MN-QIH-D [13]	0.597	45.55	76.22	85.37	5.46
HCIAE-D-MLE [30]	0.614	47.73	77.50	86.35	5.15
AMEM+SEQ-QI [41]	0.623	48.53	78.66	87.43	4.86
NMN[20]	0.616	48.24	77.54	86.75	4.98
CorefNMN\Mem	0.618	48.56	77.76	86.95	4.92
CorefNMN\ $\mathcal{L}_C^{aux}$	<b>0.636</b>	<b>50.49</b>	79.56	88.30	4.60
CorefNMN\Mem\ $\mathcal{L}_C^{aux}$	0.617	48.47	77.54	86.77	4.99
CorefNMN	<b>0.636</b>	50.24	<b>79.81</b>	<b>88.51</b>	<b>4.53</b>
CorefNMN (ResNet-152)	0.641	50.92	80.18	88.81	4.45

**Table 3:** Retrieval performance on the validation set of VisDial v0.9 [13] (discriminative models) using VGG [42] features (except last row). Higher the better for mean reciprocal rank (MRR) and recall@ $k$  (R@1, R@5, R@10), while lower the better for mean rank. Our CorefNMN model outperforms all other models across all metrics.

Thus, the dialogs contain a lot of coreferences in natural language, which need to be resolved to answer the questions accurately.

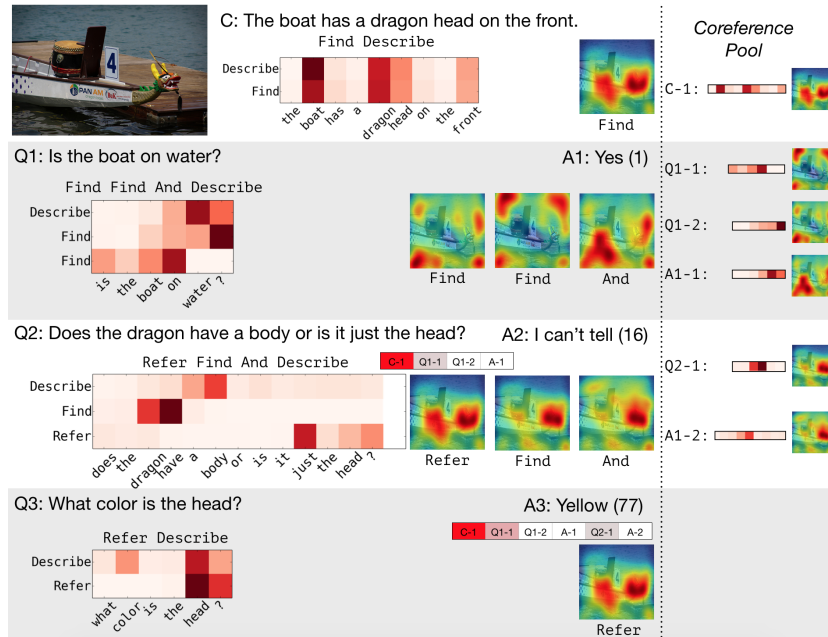
**Models and baselines.** In addition to the CorefNMN model described in Sec. 3, we also consider ablations without the memory network augmented program generator (CorefNMN\Mem) or the auxiliary loss  $\mathcal{L}_C^{aux}$  to train modules on captions (CorefNMN\ $\mathcal{L}_C^{aux}$ ), and without both (CorefNMN\Mem\ $\mathcal{L}_C^{aux}$ ). As strong baselines, we consider: (a) neural module network without history [20] with answer generation, (b) the best *discriminative* model based on memory networks MN-QIH-D from [13], (c) history-conditioned image attentive encoder (HCIAE-D-MLE) [29], and (d) Attention-based visual coreference model (AMEM+SEQ-QI) [41]. We use ImageNet pretrained VGG-16 [42] to extract  $x_{vis}$ , and also ResNet-152 [19] for CorefNMN. Further comparisons are in supplement.

**Evaluation.** Evaluation in visual dialog is via retrieval of the ground truth answer  $A_t^{gt}$  from a pool of 100 candidate answers  $\mathcal{A}_t = \{A_t^{(1)}, \dots, A_t^{(100)}\}$ . These candidates are ranked based the discriminative decoder scores. We report Recall@ $k$  for  $k = \{1, 5, 10\}$ , mean rank, and mean reciprocal rank (MRR), as suggested by [13], on the set of 40 $k$  validation images (there is not test available for v0.9).

**Supervision.** In addition to the ground truth answer  $A_t^{gt}$  at each round, our model gets program supervision for  $Q_t$ , to train the program generator. We automatically obtain (weak) program supervision from a language parser on questions (and captions) [21] and supervision to predict for **Refer** from an off-the-shelf text coreference resolution tool<sup>4</sup>, based on [12]. For questions that are a part of coreference chain, we replace **Find** with **Refer** in the parser supervised program. Our model predicts everything from the questions at test time.

**Results.** We summarize our observations from Tab. 3 below: (a) Our CorefNMN outperforms all other approaches across all the metrics, highlighting the impor-

<sup>4</sup> <https://github.com/huggingface/neuralcoref>



**Fig. 4:** Example to demonstrate explicit coreference resolution by our CorefNMN model. It begins by grounding ‘*dragon head*’ from the caption *C* (shown on top), and saves it in the coreference pool  $P_{ref}$  (right). At this point however, it does not consider the entity ‘*boat*’ important, and misses it. Next, to answer *Q1*, it localizes ‘*boat*’ and ‘*water*’, both of which are ‘unseen’, and rightly answers with *Yes*. The ground truth rank (1 for *Q1*) is shown in the brackets. Additionally, it also registers these two entities in  $P_{ref}$  for coreference resolution in future dialog. For *Q2*, it refers the phrase ‘*the head*’ to the referent registered as *C-1*, indicated by attention on the bar above *Refer*.

tance of explicitly resolving coreferences for visual dialog. Specifically, our  $R@k$  ( $k = 1, 2, 5$ ) is at least 1 point higher than the best prior work (AMEM+SEQ-QI), and almost 2 points higher than NMN. (b) Removing memory augmentation (CorefNMN\Mem) hurts performance uniformly over all metrics, as the model is unable to peek into history to decide when to resolve coreferences via the *Refer* module. Modules on captions seems to have varied effect on the full model, with decrease in  $R@1$ , but marginal increase or no effect in other metrics. (c) Fig. 4 illustrates the interpretable and grounded nature of our model.

### 4.3 VisDial v1.0 Dataset

**Dataset.** Das *et al.*[13] recently released VisDial v1.0 dataset. Specifically, VisDial v1.0 comprises of: (a) A re-organization of train and validation splits from v0.9 to form the new train v1.0. Thus, train v1.0 now contains 120k images with 10–round dialogs for each images, resulting in a total of 1.2 million question-answer pairs. (b) An additional 10k COCO-like images from Flickr, on which

Model	MRR	R@1	R@5	R@10	Mean	NDCG
LF-QIH-D [13] (VGG)	0.554	40.95	72.45	82.83	5.95	0.453
HRE-QIH-D [13] (VGG)	0.542	39.93	70.45	81.50	6.41	0.455
MN-QIH-D [13] (VGG)	0.555	40.98	72.30	83.30	5.92	0.475
NMN [20]	0.588	44.15	76.88	86.88	4.81	<b>0.581</b>
CorefNMN	<b>0.615</b>	<b>47.55</b>	<b>78.10</b>	<b>88.80</b>	<b>4.40</b>	0.547

**Table 4:** Retrieval performance on the test-standard split of VisDial v1.0 dataset [13] (discriminative models). Higher the better for mean reciprocal rank (MRR), recall@ $k$  (R@1, R@5, R@10), and normalized discounted cumulative gain (NDCG) while lower the better for mean rank. Our CorefNMN model outperforms all other models across all metrics, except neural module baseline (NMN) on NDCG.

crowd-sourced dialogs between pairs of humans were collected similar to v0.9. The 10k images are further split into 2k validation (val v1.0) and 8k test sets (test-std v1.0). Dense candidate option annotations, which indicate the correctness of each candidate in the pool, were also collected for these 10k images. Each image in the val v1.0 split is associated with a 10-round dialog, while an image in test-std v1.0 has a variable-round dialog.

**Additional Metrics and Models.** Just as in the previous version (v0.9), the performance on the VisDial v1.0 dataset is benchmarked using standard retrieval metrics like Recall@ $k$  ( $k = \{1, 5, 10\}$ ), mean reciprocal rank (MRR) and mean rank. Further, Das *et al.*[13] also propose to use normalized discounted cumulative gain (NDCG) to score the sorted pool of candidate answers, to evaluate on VisDial v1.0<sup>5</sup>. Intuitively, NDCG penalizes accurate answers that appear lower in the sorted pool based on a logarithmic weighting scheme, normalizing for the number of accurate answers across instances. We train our CorefNMN model on train v1.0 and report numbers on test-std v1.0 split. We compare against LF-QIH-D, HRE-QIH-D, and MN-QIH-D from [13], and out-of-the-box neural module network (NMN) [20]. The LF-QIH-D, HRE-QIH-D, and MN-QIH-D use VGG [42] image features, while the neural module based models (CorefNMN and NMN) use ResNet-152 [19] features.

### Results.

Performance on VisDial v1.0 is given in Tab. 4. Our CorefNMN outperforms all other approaches on all metrics, except the neural module baseline (NMN) on the NDCG metric. We note that recent state-of-the-art, as reported on the leaderboard<sup>6</sup>, has reached up to 0.578 (NDCG) from team DL-61, but the approach and other details (e.g. features, use of an ensemble) are not fully known and unpublished at this point in time.

<sup>5</sup> <https://visualdialog.org/challenge/2018>

<sup>6</sup> <https://evalai.cloudcv.org/web/challenges/challenge-page/103/leaderboard/298>

## 5 Conclusions

We introduced a novel model for visual dialog based on neural module networks that provides an introspective reasoning about visual coreferences. It explicitly links coreferences and grounds them in the image at a word-level, rather than implicitly or at a sentence-level, as in prior visual dialog work. Our CorefNMN outperforms prior work on both the MNIST dialog dataset (close to perfect accuracy), and on VisDial dataset, while being more interpretable, grounded, and consistent by construction.

**Acknowledgements.** This work was supported in part by NSF, AFRL, DARPA, Siemens, Google, Amazon, ONR YIPs and ONR Grants N00014-16-1-{2713,2793}, N000141210903. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of the U.S. Government, or any sponsor.



## Overview of Supplement

The supplement is organized as follows:

- Sec. A shows the results of our model using a discriminative decoder with image features extracted using ImageNet pretrained ResNet-152 [19], showing superior performance of our explicit coreference model CorefNMN,
- Sec. B details our experiments with a generative answer decoder,
- Implementation details for our experiments are given in Sec. C, and
- Schematics of our novel **Refer** module are in Fig. 5, Fig. 6 visualizes the auxiliary task used to run modules on captions, a novel way to handle captions at a fine word-level granularity, and Fig. 7 shows another qualitative example from VisDial.

## A Discriminative Decoder Experiments

**Comparisons with ResNet-152 features.** As mentioned in Sec. 4.2 of the main paper, the models trained on VisDial v0.9 used an ImageNet pretrained VGG [42] to extract the image features  $x_{vis}$ . In this section, we present results where a pretrained ResNet-152 [19] was used to obtain the image features for our CorefNMN model, in Tab. 5. For a fair comparison, we obtained performance metrics from the authors for few of these baselines (MNQIH-G, LF-QH-G), and retrain NMN with ResNet-152 features.

## B Generative Decoder Experiments

The main paper describes a discriminative answer decoder (Sec.3.4) and presents results for the same (Sec. 4.2) on VisDial v0.9 dataset. As a reminder, a discriminative decoder takes the context vector  $c_t$  as an input, and scores candidate answers according to their correctness. In other words, a discriminative decoder needs to be presented with a list of candidate answers and cannot ‘generate’ novel answers. We now introduce a generative answer decoder and present results on VisDial v0.9 dataset.

**Generative Answer Decoder** is a language model composed of a multi-layer LSTM with  $c_t$  as its initial state. During training, we minimize the negative log-likelihood  $\mathcal{L}_A^{dec}$  of the ground truth answer  $A_t^{gt}$  with respect to the model. At test time, we use the decoder to score all candidates in the answer pool  $\mathcal{A}_t$  by model log-likelihood, and rank them accordingly. Note that this ranking of candidate answers is done to comply with the evaluation protocol of VisDial v0.9, and is not a limitation of generative answer decoder, unlike the discriminative one. Thus, the generative answer decoder can potentially be used to generate novel answers to a given question in the visual dialog via language generation.

Model	MRR	R@1	R@5	R@10	Mean
LF-QIH-D* [13]	0.591	44.91	75.68	84.92	5.55
HRE-QIH-D* [13]	0.586	44.86	74.35	83.86	5.81
MN-QIH-D* [13]	0.601	46.04	76.78	85.93	5.29
NMN[20]	0.620	48.89	77.83	86.99	4.94
CorefNMN	<b>0.641</b>	<b>50.92</b>	<b>80.18</b>	<b>88.81</b>	<b>4.45</b>

**Table 5:** Retrieval performance on the validation set of VisDial dataset v0.9 [13] (discriminative models with ResNet-152 [19] features). Higher the better for mean reciprocal rank (MRR) and recall@ $k$  (R@1, R@5, R@10), while lower the better for mean rank. Our CorefNMN model outperforms all other models across all metrics. \*indicates numbers obtained from authors for models retrained on ResNet-152 features.

**Models and baselines.** We denote our model, as described in Sec. 3, with CorefNMN to indicate that the model uses history  $H$ . We also consider ablations which do not have the memory network augmented program generator (CorefNMN\Mem), or, the auxiliary loss  $\mathcal{L}_C^{aux}$  to train modules on captions (CorefNMN\mathcal{L}\_C^{aux}), and a combination of both (CorefNMN\Mem\mathcal{L}\_C^{aux}). As strong baselines, we consider: (a) neural module network without history [20] with answer generation, (b) the best *generative* model based on memory networks MN-QIH-G from [13], in addition to their LF-QIG-G and HRE-QIH-G models, and (c) history-conditioned image attentive encoder (HCIAE-G-MLE) [29]. We do not consider the HCIAE-G-DIS model with perceptual loss [29] as its contribution is complementary to our model. Our model uses ImageNet pre-trained ResNet-152 [19] to extract features for images  $x_{vis}$ , while some of these baseline models originally use VGG-16 [42] features. For a fair comparison, we obtained performance metrics from the authors for few of these baselines (MN-QIH-G, LF-QH-G), and report the others (HCIAE-G-MLE) as is.

**Results.** We summarize our observations from Tab. 6 below: (a) Our CorefNMN outperforms all other approaches according to R@5, R@10, and mean rank metrics, highlighting the importance of explicitly resolving coreferences for visual dialog. Specifically, our mean rank of 15.69 is a 4% improvement over the NMN baseline. (b) However, the NMN baseline has a higher R@1, and perhaps as a result, the best MRR. A possible reason could be due to the noisy supervision from the out-of-domain, automatic, text-based coreference tool, as we entirely rely on it to predict **Refer**, the module responsible for coreference resolution. (c) Our novel way of handling captions using modules (indicated  $\mathcal{L}_C^{aux}$ ) boosts the full model, while the hurting the ablation without the memory augmentation. That is, CorefNMN\Mem\mathcal{L}\_C^{aux} is better than CorefNMN\Mem across all metrics, while CorefNMN\mathcal{L}\_C^{aux} is uniformly worse than CorefNMN.

Model	MRR	R@1	R@5	R@10	Mean
MN-QIH-G [13] (VGG)	0.526	42.29	62.85	68.88	17.06
HCIAE-G-MLE[30](VGG)	0.539	44.06	<b>63.55</b>	69.24	16.01
LF-QIH-G* [13]	0.515	41.04	61.63	67.54	17.32
HRE-QIH-G* [13]	0.523	42.26	62.20	67.95	16.96
MN-QIH-G* [13]	0.527	42.60	62.58	68.52	17.21
NMN[20]	<b>0.542</b>	<b>45.05</b>	63.27	69.28	16.34
CorefNMN\Mem	0.531	43.67	62.31	68.27	16.73
CorefNMN\Mem\mathcal{L}_C^{aux}	0.537	44.26	63.02	69.01	16.47
CorefNMN\mathcal{L}_C^{aux}	0.533	43.62	63.08	69.12	16.39
CorefNMN	0.535	43.66	<b>63.54</b>	<b>69.93</b>	<b>15.69</b>

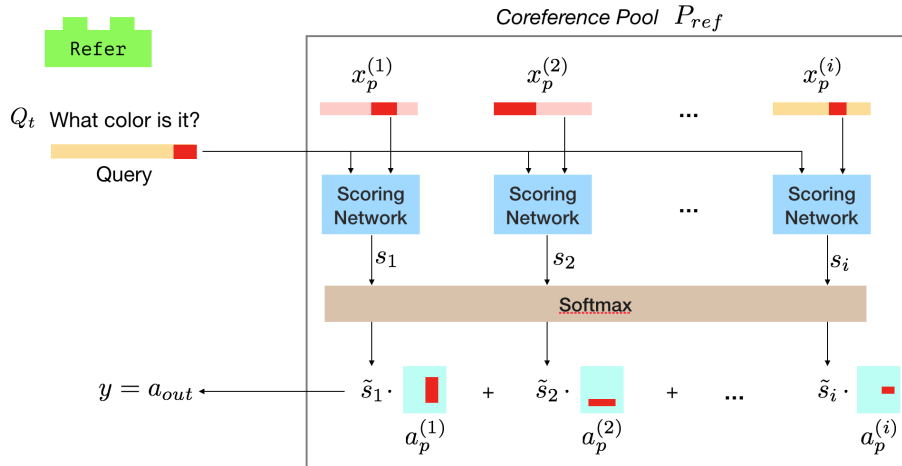
**Table 6:** Retrieval performance on the validation set of VisDial v0.9 [13] (generative models). Higher the better for mean reciprocal rank (MRR) and recall@ $k$  (R@1, R@5, R@10), while lower the better for mean rank. Our CorefNMN model outperforms all other models in R@5, R@10, and mean rank. However, the neural module network baseline (NMN) has the best R@1, and perhaps as a result, the best MRR as well. \*indicates numbers obtained from authors for models retrained on ResNet-152 features.

## C Implementation Details

All our models are implemented with Tensorflow v1.0 [2]. To optimize, we use Adam [24] with a learning rate of 0.0001. Gradients at each iteration are clamped to  $[-2.0, 2.0]$  to avoid gradient explosion. To preprocess text, we follow [13], i.e., we lowercase all questions and answers, and tokenize using the Python NLTK framework [1]. We then construct a dictionary of all words that appear at least five times in the training set. Specific model hyperparameters for each experiment are given below.

### C.1 MNIST Dialog Dataset

Due to the synthetic nature, the text in the dialog is of low variability and is made up of a small vocabulary. Thus, we only use a single-layered LSTM with a hidden size of 64 to encode both questions and history. For each word in our vocabulary of size 73, we learn embeddings with 32 dimensions. We learn a similar dimensional embedding for each of our modules. To extract image features, we design a convolutional neural network (CNN) with the same architecture as [41]. Specifically, our CNN has four  $3 \times 3$  convolutional layers, each followed by a batch norm, ReLU non-linearity, and a  $2 \times 2$  max pool layer. While the first two convolutional layers have 32 feature channels, the last two have 64 channels each. To pick the best model, we use early stoppage on the provided validation set of  $10k$  images.



**Fig. 5:** Visualization how our refer module accesses the coreference pool, to understand “it” in this example. For notation see main paper Section 3.  $x_p$  can be seen as the keys to retrieve attentions  $a_p$  from the coreference pool  $P_{ref}$ . The yellow and pink lines symbolize the sentences/questions, and red is the text attention. Cyan boxes symbolize the image, with red being the spatial attention.

## C.2 VisDial Dataset

Each LSTM used in our VisDial experiments has two layers and with a hidden size of 1000. To represent images, we use convolutional features before the final mean pooling from a ImageNet pre-trained ResNet-152 [19] model. Further, we also add two additional dimensions indicating the  $X$  (columns) and  $Y$  (rows) locations respectively, to facilitate the model in handling spatial reasoning. With a large vocabulary of around  $8k$  words, our word and module embeddings are 300 dimensional vectors. We also initialize our word embeddings with GloVe [35]. We pick the best model via early stopping using mean reciprocal rank metric on a subset of  $3k$  images, set aside from the  $83k$  training images of VisDial v0.9.

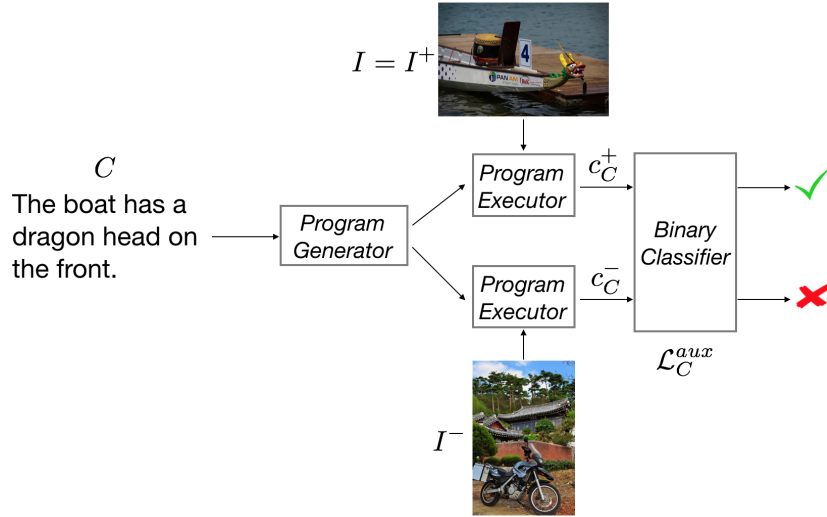
## C.3 Document Changelog

To help the readers track changes to this document, a brief changelog describing the revisions is provided below:

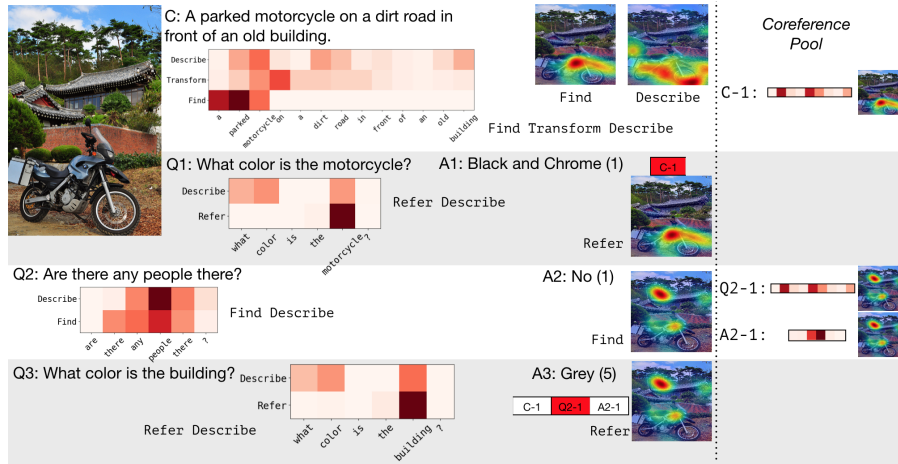
- v0:** ECCV 2018 camera-ready version (not on arXiv).
- v1:** Initial arXiv version: Added experiments on VisDial v1.0 dataset.

## References

1. NLTK. <http://www.nltk.org/>



**Fig. 6:** Flow diagram for the auxiliary task that measures the alignment between a caption and image as a binary classification task. By predicting a program for the caption, we propose a novel way to process captions at a finer word-level. Sec. 3.3 of the main paper motivates the use of modules on captions via the auxiliary task.



**Fig. 7:** Qualitative example on VisDial dataset.

2. Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G.S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., Zheng, X.: TensorFlow: Large-scale machine learning on heterogeneous systems (2015), <https://www.tensorflow.org/>, software available from tensorflow.org
3. Agrawal, A., Batra, D., Parikh, D.: Analyzing the behavior of visual question answering models. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP) (2016)
4. Anderson, P., He, X., Buehler, C., Teney, D., Johnson, M., Gould, S., Zhang, L.: Bottom-up and top-down attention for image captioning and vqa. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2018)
5. Andreas, J., Rohrbach, M., Darrell, T., Klein, D.: Learning to Compose Neural Networks for Question Answering (2016)
6. Andreas, J., Rohrbach, M., Darrell, T., Klein, D.: Neural module networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
7. Anne Hendricks, L., Wang, O., Shechtman, E., Sivic, J., Darrell, T., Russell, B.: Localizing moments in video with natural language. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2017)
8. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Zitnick, C.L., Parikh, D.: Vqa: Visual question answering. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2015)
9. Bahdanau, D., Cho, K., Bengio, Y.: Neural Machine Translation by Jointly Learning to Align and Translate. In: Proceedings of the International Conference on Learning Representations (ICLR) (2015)
10. Bergsma, S., Lin, D.: Bootstrapping path-based pronoun resolution. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL). Association for Computational Linguistics (2006)
11. Burks, A.W., Warren, D.W., Wright, J.B.: An analysis of a logical machine using parenthesis-free notation. *Mathematical Tables and Other Aids to Computation* **8**(46), 53–57 (1954), <http://www.jstor.org/stable/2001990>
12. Clark, K., Manning, C.D.: Deep reinforcement learning for mention-ranking coreference models. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP) (2016)
13. Das, A., Kottur, S., Gupta, K., Singh, A., Yadav, D., Moura, J.M., Parikh, D., Batra, D.: Visual Dialog. In: CVPR (2017)
14. Das, A., Kottur, S., Moura, J.M., Lee, S., Batra, D.: Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning. arXiv preprint arXiv:1703.06585 (2017)
15. Fukui, A., Park, D.H., Yang, D., Rohrbach, A., Darrell, T., Rohrbach, M.: Multi-modal compact bilinear pooling for visual question answering and visual grounding. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP) (2016)
16. Geman, D., Geman, S., Hallonquist, N., Younes, L.: Visual Turing Test for computer vision systems. In: Proceedings of the National Academy of Sciences (2015)
17. Goyal, Y., Khot, T., Summers-Stay, D., Batra, D., Parikh, D.: Making the v in vqa matter: Elevating the role of image understanding in visual question answering. In:

- Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
18. Grice, H.P.: Logic and conversation. In: Cole, P., Morgan, J.L. (eds.) *Syntax and Semantics: Vol. 3: Speech Acts*, pp. 41–58. Academic Press, New York (1975), <http://www.ucl.ac.uk/lis/studypacks/Grice-Logic.pdf>
  19. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016)
  20. Hu, R., Andreas, J., Rohrbach, M., Darrell, T., Saenko, K.: Learning to reason: End-to-end module networks for visual question answering. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)* (2017)
  21. Hu, R., Xu, H., Rohrbach, M., Feng, J., Saenko, K., Darrell, T.: Natural language object retrieval. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016)
  22. Johnson, J., Hariharan, B., van der Maaten, L., Fei-Fei, L., Zitnick, C.L., Girshick, R.: Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE (2017)
  23. Johnson, J., Hariharan, B., van der Maaten, L., Hoffman, J., Fei-Fei, L., Zitnick, C.L., Girshick, R.: Inferring and executing programs for visual reasoning. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)* (2017)
  24. Kingma, D., Ba, J.: Adam: A method for stochastic optimization. arXiv:1412.6980 (2014)
  25. Kong, C., Lin, D., Bansal, M., Urtasun, R., Fidler, S.: What are you talking about? text-to-image coreference. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014)
  26. LeCun, Y., Cortes, C.: MNIST handwritten digit database (2010), <http://yann.lecun.com/exdb/mnist/>
  27. Lin, D., Fidler, S., Kong, C., Urtasun, R.: Visual semantic search: Retrieving videos via complex textual queries. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2014)
  28. Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollr, P., Zitnick, C.L.: Microsoft COCO: Common Objects in Context. In: *Proceedings of the European Conference on Computer Vision (ECCV)* (2014)
  29. Lu, J., Kannan, A., Yang, J., Parikh, D., Batra, D.: Best of both worlds: Transferring knowledge from discriminative learning to a generative visual dialog model. In: *Advances in Neural Information Processing Systems (NIPS)* (2017)
  30. Lu, J., Yang, J., Batra, D., Parikh, D.: Hierarchical Question-Image Co-Attention for Visual Question Answering. In: *Advances in Neural Information Processing Systems (NIPS)* (2016)
  31. Malinowski, M., Rohrbach, M., Fritz, M.: Ask your neurons: A neural-based approach to answering questions about images. In: *Proceedings of the IEEE International Conference on Computer Vision (ICCV)* (2015)
  32. Mao, J., Huang, J., Toshev, A., Camburu, O., Yuille, A., Murphy, K.: Generation and comprehension of unambiguous object descriptions. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016)
  33. Massiceti, D., Siddharth, N., Dokania, P.K., Torr, P.H.S.: Flipdial: A generative model for two-way visual dialogue (2018)
  34. Mitchell, M., van Deemter, K., Reiter, E.: Generating expressions that refer to visible objects. In: *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.

- pp. 1174–1184. Association for Computational Linguistics, Atlanta, Georgia (June 2013), <http://www.aclweb.org/anthology/N13-1137>
35. Pennington, J., Socher, R., Manning, C.D.: Glove: Global vectors for word representation. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP) (2014)
  36. Plummer, B., Wang, L., Cervantes, C., Caicedo, J., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV) (2015)
  37. Ramanathan, V., Joulin, A., Liang, P., Fei-Fei, L.: Linking people in videos with "their" names using coreference resolution. In: Proceedings of the European Conference on Computer Vision (ECCV) (2014)
  38. Regneri, M., Rohrbach, M., Wetzel, D., Thater, S., Schiele, B., Pinkal, M.: Grounding Action Descriptions in Videos. *Transactions of the Association for Computational Linguistics (TACL)* **1**, 25–36 (2013)
  39. Rohrbach, A., Rohrbach, M., Hu, R., Darrell, T., Schiele, B.: Grounding of textual phrases in images by reconstruction. In: Proceedings of the European Conference on Computer Vision (ECCV) (2016)
  40. Rohrbach, A., Rohrbach, M., Tang, S., Oh, S.J., Schiele, B.: Generating descriptions with grounded and co-referenced people. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
  41. Seo, P.H., Lehrmann, A., Han, B., Sigal, L.: Visual reference resolution using attention memory for visual dialog. In: Advances in Neural Information Processing Systems (NIPS) (2017)
  42. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. In: Proceedings of the International Conference on Learning Representations (ICLR) (2015)
  43. Strub, F., de Vries, H., Mary, J., Piot, B., Courville, A.C., Pietquin, O.: End-to-end optimization of goal-driven and visually grounded dialogue systems. In: Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI) (2017)
  44. de Vries, H., Strub, F., Chandar, S., Pietquin, O., Larochelle, H., Courville, A.C.: Guesswhat?! visual object discovery through multi-modal dialogue. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2017)
  45. Wang, L., Li, Y., Lazebnik, S.: Learning deep structure-preserving image-text embeddings. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)
  46. Weston, J., Chopra, S., Bordes, A.: Memory networks. In: Proceedings of the International Conference on Learning Representations (ICLR) (2015)
  47. Winograd, T.: Procedures as a representation for data in a computer program for understanding natural language. Tech. rep., DTIC Document (1971)
  48. Winograd, T.: *Understanding Natural Language*. Academic Press, Inc., Orlando, FL, USA (1972)
  49. Yu, H., Siskind, J.M.: Grounded language learning from videos described with sentences. In: Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL) (2013)
  50. Yu, L., Poirson, P., Yang, S., Berg, A.C., Berg, T.L.: Modeling context in referring expressions. In: Proceedings of the European Conference on Computer Vision (ECCV) (2016)



51. Zhang, P., Goyal, Y., Summers-Stay, D., Batra, D., Parikh, D.: Yin and Yang: Balancing and Answering Binary Visual Questions. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (2016)