

# QUERY-REDUCTION NETWORKS FOR QUESTION ANSWERING

**Minjoon Seo**<sup>1</sup>   **Sewon Min**<sup>2</sup>   **Ali Farhadi**<sup>1,3</sup>   **Hannaneh Hajishirzi**<sup>1</sup>  
 University of Washington<sup>1</sup>, Seoul National University<sup>2</sup>, Allen Institute for Artificial Intelligence<sup>3</sup>  
 {minjoon, ali, hannaneh}@cs.washington.edu, shmsw25@snu.ac.kr

## ABSTRACT

In this paper, we study the problem of question answering when reasoning over multiple facts is required. We propose Query-Reduction Network (QRN), a variant of Recurrent Neural Network (RNN) that effectively handles both short-term (local) and long-term (global) sequential dependencies to reason over multiple facts. QRN considers the context sentences as a sequence of state-changing triggers, and *reduces* the original query to a more informed query as it observes each trigger (context sentence) through time. Our experiments show that QRN produces the state-of-the-art results in bAbI QA and dialog tasks, and in a real goal-oriented dialog dataset. In addition, QRN formulation allows parallelization on RNN’s time axis, saving an order of magnitude in time complexity for training and inference.

## 1 INTRODUCTION

In this paper, we address the problem of question answering (QA) when reasoning over multiple facts is required. For example, consider we know that `Frogs eat insects` and `Flies are insects`. Then answering `Do frogs eat flies?` requires reasoning over both of the above facts. Question answering, more specifically context-based QA, has been extensively studied in machine comprehension tasks (Richardson et al., 2013; Hermann et al., 2015; Hill et al., 2016; Rajpurkar et al., 2016). However, most of the datasets are primarily focused on lexical and syntactic understanding, and hardly concentrate on inference over multiple facts. Recently, several datasets aimed for testing multi-hop reasoning have emerged; among them are story-based QA (Weston et al., 2016) and the dialog task (Bordes and Weston, 2016).

Recurrent Neural Network (RNN) and its variants, such as Long Short-Term Memory (LSTM) (Hochreiter and Schmidhuber, 1997) and Gated Recurrent Unit (GRU) (Cho et al., 2014), are popular choices for modeling natural language. However, when used for multi-hop reasoning in question answering, purely RNN-based models have shown to perform poorly (Weston et al., 2016). This is largely due to the fact that RNN’s internal memory is inherently unstable over a long term. For this reason, most recent approaches in the literature have mainly relied on global attention mechanism and shared external memory (Sukhbaatar et al., 2015; Peng et al., 2015; Xiong et al., 2016; Graves et al., 2016). The attention mechanism allows these models to focus on a single sentence in each layer. They can sequentially read multiple relevant sentences from the memory with multiple layers to perform multi-hop reasoning. However, one major drawback of these standard attention mechanisms is that they are insensitive to the time step (memory address) of the sentences when accessing them.

Our proposed model, Query-Reduction Network<sup>1</sup>(QRN), is a single recurrent unit that addresses the long-term dependency problem of most RNN-based models by simplifying the recurrent update, while taking the advantage of RNN’s capability to model sequential data (Figure 1). QRN considers the context sentences as a sequence of state-changing triggers, and transforms (*reduces*) the original query to a more informed query as it observes each trigger through time. For instance in Figure 1b, the original question, `Where is the apple?`, cannot be directly answered by any single sentence from the story. After observing the first sentence, `Sandra got the apple there`, QRN transforms the original question to a reduced query `Where is Sandra?`, which is presumably

<sup>1</sup>Code is publicly available at: [seominjoon.github.io/qrn/](http://seominjoon.github.io/qrn/)

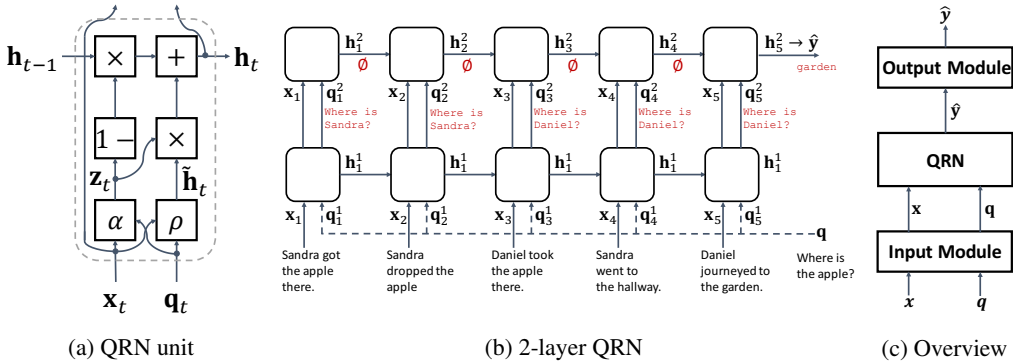


Figure 1: (1a) QRN unit, (1b) 2-layer QRN on 5-sentence story, and (1c) entire QA system (QRN and input / output modules).  $x$ ,  $q$ ,  $\hat{y}$  are the story, question and predicted answer in natural language, respectively.  $\mathbf{x} = \langle x_1, \dots, x_T \rangle$ ,  $\mathbf{q}$ ,  $\hat{\mathbf{y}}$  are their corresponding vector representations (upright font).  $\alpha$  and  $\rho$  are update gate and reduce functions, respectively.  $\hat{\mathbf{y}}$  is assigned to be  $h_5^2$ , the local query at the last time step in the last layer. Also, red-colored text is the inferred meanings of the vectors (see ‘Interpretations’ of Section 5.3).

easier to answer than the original question given the context provided by the first sentence.<sup>2</sup> Unlike RNN-based models, QRN’s candidate state ( $\tilde{h}_t$  in Figure 1a) does not depend on the previous hidden state ( $h_{t-1}$ ). Compared to memory-based approaches (Weston et al., 2015; Sukhbaatar et al., 2015; Peng et al., 2015; Kumar et al., 2016; Xiong et al., 2016), QRN can better encode *locality* information because it does not use a global memory access controller (circle nodes in Figure 2), and the query updates are performed locally.

In short, the main contribution of QRN is threefold. First, QRN is a simple variant of RNN that *reduces* the query given the context sentences in a differentiable manner. Second, QRN is situated between the attention mechanism and RNN, effectively handling time dependency and long-term dependency problems of each technique, respectively. Hence it is well-suited for sequential data with both local and global interactions (note that QRN is *not* the replacement of RNN, which is arguably better for modeling complex local interactions). Third, unlike most RNN-based models, QRN can be parallelized over time by computing candidate reduced queries ( $\tilde{h}_t$ ) directly from local input queries ( $q_t$ ) and context sentence vectors ( $x_t$ ). In fact, the parallelizability of QRN implies that QRN does not suffer from the vanishing gradient problem of RNN, hence effectively addressing the long-term dependency. We experimentally demonstrate these contributions by achieving the state-of-the-art results on story-based QA and interactive dialog datasets.

## 2 MODEL

In story-based QA (or dialog dataset), the input is the *context* as a sequence of sentences (story or past conversations) and a *question* in natural language (equivalent to the user’s last utterance in the dialog). The output is the predicted answer to the question in natural language (the system’s next utterance in the dialog). The only supervision provided during training is the answer to the question.

In this paper we particularly focus on end-to-end solutions, i.e., the only supervision comes from questions and answers, and we restrain from using manually defined rules or external language resources, such as lexicon or dependency parser. Let  $\langle x_1, \dots, x_T \rangle$  denote the sequence of sentences, where  $T$  is the number of sentences in the story, and let  $q$  denote the question. Let  $\hat{y}$  denote the predicted answer, and  $y$  denote the true answer. Our proposed system for end-to-end QA task is divided into three modules (Figure 1c): input module, QRN layers, and output module.

**Input module.** Input module maps each sentence  $x_t$  and the question  $q$  to  $d$ -dimensional vector space,  $x_t \in \mathbb{R}^d$  and  $q_t \in \mathbb{R}^d$ . We adopt a previous solution for the input module (details in Section 5).

<sup>2</sup>This mechanism is akin to logic regression in situation calculus (Reiter, 2001).

**QRN layers.** QRN layers use the sentence vectors and the question vector from the input module to obtain the predicted answer in vector space,  $\hat{y} \in \mathbb{R}^d$ . A QRN layer refers to the recurrent application of a QRN unit, which can be considered as a variant of RNN with two inputs, two outputs, and a hidden state (reduced query), all of which operate in vector space. The details of the QRN module is explained throughout this section (2.1, 2.2).

**Output module.** Output module maps  $\hat{y}$  obtained from QRN to a natural language answer  $\hat{g}$ . Similar to the input module, we adopt a standard solution for the output module (details in Section 5).

We first formally define the base model of a QRN unit, and then we explain how we connect the input and output modules to it (Section 2.1). We also present a few extensions to the network that can improve QRN’s performance (Section 2.2). Finally, we show that QRN can be parallelized over time, giving computational advantage over most RNN-based models by one order of magnitude (Section 3).

## 2.1 QRN UNIT

As an RNN-based model, QRN is a single recurrent unit that updates its hidden state (reduced query) through time and layers. Figure 1a depicts the schematic structure of a QRN unit, and Figure 1b demonstrates how layers are stacked. A QRN unit accepts two inputs (*local* query vector  $\mathbf{q}_t \in \mathbb{R}^d$  and sentence vector  $\mathbf{x}_t \in \mathbb{R}^d$ ), and two outputs (reduced query vector  $\mathbf{h}_t \in \mathbb{R}^d$ , which is similar to the hidden state in RNN, and the sentence vector  $\mathbf{x}_t$  from the input without modification). The local query vector is not necessarily identical to the original query (question) vector  $\mathbf{q}$ . In order to compute the outputs, we use *update gate* function  $\alpha : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, 1]$  and *reduce* function  $\rho : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}^d$ . Intuitively, the update gate function measures the relevance between the sentence and the local query and is used to update the hidden state. The reduce function transforms the local query input to a candidate state which is a new reduced (easier) query given the sentence. The outputs are calculated with the following equations:

$$z_t = \alpha(\mathbf{x}_t, \mathbf{q}_t) = \sigma(\mathbf{W}^{(z)}(\mathbf{x}_t \circ \mathbf{q}_t) + b^{(z)}) \quad (1)$$

$$\tilde{\mathbf{h}}_t = \rho(\mathbf{x}_t, \mathbf{q}_t) = \tanh(\mathbf{W}^{(h)}[\mathbf{x}_t; \mathbf{q}_t] + \mathbf{b}^{(h)}) \quad (2)$$

$$\mathbf{h}_t = z_t \tilde{\mathbf{h}}_t + (1 - z_t) \mathbf{h}_{t-1} \quad (3)$$

where  $z_t$  is the scalar update gate,  $\tilde{\mathbf{h}}_t$  is the candidate reduced query, and  $\mathbf{h}_t$  is the final reduced query at time step  $t$ ,  $\sigma(\cdot)$  is sigmoid activation,  $\tanh(\cdot)$  is hyperbolic tangent activation (applied element-wise),  $\mathbf{W}^{(z)} \in \mathbb{R}^{1 \times d}$ ,  $\mathbf{W}^{(h)} \in \mathbb{R}^{d \times 2d}$  are weight matrices,  $b^{(z)} \in \mathbb{R}$ ,  $\mathbf{b}^{(h)} \in \mathbb{R}^d$  are bias terms,  $\circ$  is element-wise vector multiplication, and  $[\cdot]$  is vector concatenation along the row. As a base case,  $\mathbf{h}_0 = \mathbf{0}$ . Here we have explicitly defined  $\alpha$  and  $\rho$ , but they can be any reasonable differentiable functions.

The update gate is similar to the global attention mechanism (Sukhbaatar et al., 2015; Xiong et al., 2016) in that it measures the similarity between the sentence (a memory slot) and the query. However, a significant difference is that the update gate is computed using sigmoid ( $\sigma$ ) function on the current memory slot only (hence internally embedded within the unit), whereas the global attention is computed using softmax function over the entire memory (hence globally defined). The update gate can be rather considered as *local sigmoid* attention.

**Stacking layers** We just showed the single-layer case of QRN, but QRN with multiple layers is able to perform reasoning over multiple facts more effectively, as shown in the example of Figure 1b. In order to stack several layers of QRN, the outputs of the current layer are used as the inputs to the next layer. That is, using superscript  $k$  to denote the current layer’s index (assuming 1-based indexing), we let  $\mathbf{q}_t^{k+1} = \mathbf{h}_t^k$ . Note that  $\mathbf{x}_t$  is passed to the next layer without any modification, so we do not put a layer index on it.

**Bi-direction.** So far we have assumed that QRN only needs to look at past sentences, whereas often times, query answers can depend on future sentences. For instance, consider a sentence “John dropped the football.” at time  $t$ . Then, even if there is no mention about the “football” in the past (at time  $i < t$ ), it can be implied that “John” has the “football” at the current time  $t$ . In order to incorporate the future dependency, we obtain  $\vec{\mathbf{h}}_t$  and  $\overleftarrow{\mathbf{h}}_t$  in both forward and backward directions,

respectively, using Equation 3. We then add them together to get  $\mathbf{q}_t$  for the next layer. That is,

$$\mathbf{q}_t^{k+1} = \overrightarrow{\mathbf{h}}_t^k + \overleftarrow{\mathbf{h}}_t^k \quad (4)$$

for layer indices  $1 \leq k \leq K - 1$ . Note that the variables  $\mathbf{W}^{(z)}, b^{(z)}, \mathbf{W}^{(h)}, \mathbf{b}^{(h)}$  are shared between the two directions.

**Connecting input and output modules.** Figure 1c depicts how QRN is connected with the input and output modules. In the first layer of QRN,  $\mathbf{q}_t^1 = \mathbf{q}$  for all  $t$ , where  $\mathbf{q}$  is obtained from the input module by processing the natural language question input  $\mathbf{q}$ .  $\mathbf{x}_t$  is also obtained from  $\mathbf{x}_t$  by the same input module. The output at the last time step in the last layer is passed to the output module. That is,  $\hat{\mathbf{y}} = \mathbf{h}_t^K$  where  $K$  represent the number of layers in the network. Then the output module gives the predicted answer  $\hat{\mathbf{y}}$  in natural language.

## 2.2 EXTENSIONS

Here we introduce a few extensions of QRN, and later in our experiments, we test QRN’s performance with and without each of these extensions.

**Reset gate.** Inspired by GRU (Cho et al., 2014), we found that it is useful to allow the QRN unit to reset (nullify) the candidate reduced query (i.e.,  $\tilde{\mathbf{h}}_t$ ) when necessary. For this we use a *reset gate* function  $\beta : \mathbb{R}^d \times \mathbb{R}^d \rightarrow [0, 1]$ , which can be defined similarly to the update gate function:

$$r_t = \beta(\mathbf{x}_t, \mathbf{q}_t) = \sigma(\mathbf{W}^{(r)}(\mathbf{x}_t \circ \mathbf{q}_t) + b^{(r)}) \quad (5)$$

where  $\mathbf{W}^{(r)} \in \mathbb{R}^{1 \times d}$  is a weight matrix, and  $b^{(r)} \in \mathbb{R}$  is a bias term. Equation 3 is rewritten as

$$\mathbf{h}_t = z_t r_t \tilde{\mathbf{h}}_t + (1 - z_t) \mathbf{h}_{t-1}. \quad (6)$$

Note that we do not use the reset gate in the last layer.

**Vector gates.** As in LSTM and GRU, update and reset gates can be vectors instead of scalar values for fine-controlled gating. For vector gates, we modify the row dimension of weights and biases in Equation 1 and 5 from 1 to  $d$ . Then we obtain  $\mathbf{z}_t, \mathbf{r}_t \in \mathbb{R}^d$  (instead of  $z_t, r_t \in \mathbb{R}$ ), and these can be element-wise multiplied ( $\circ$ ) instead of being broadcasted in Equation 3 and 6.

## 3 PARALLELIZATION

An important advantage of QRN is that the recurrent updates in Equation 3 and 5 can be computed in parallel across time. This is in contrast with most RNN-based models that cannot be parallelized, where computing the candidate hidden state at time  $t$  explicitly requires the previous hidden state. In QRN, the final reduced queries ( $\mathbf{h}_t$ ) can be decomposed into computing over candidate reduced queries ( $\tilde{\mathbf{h}}_t$ ), without looking at the previous reduced query. Here we primarily show that the query update in Equation 3 can be parallelized by rewriting the equation with matrix operations. The extension to Equation 5 is straightforward. The proof for QRN with vector gates is shown in Appendix B. The recursive definition of Equation 3 can be explicitly written as

$$\mathbf{h}_t = \sum_{i=1}^t \left[ \prod_{j=i+1}^t 1 - z_j \right] z_i \tilde{\mathbf{h}}_i = \sum_{i=1}^t \exp \left\{ \sum_{j=i+1}^t \log(1 - z_j) \right\} z_i \tilde{\mathbf{h}}_i. \quad (7)$$

Let  $b_i = \log(1 - z_i)$  for brevity. Then we can rewrite Equation 7 as the following equation:

$$\begin{pmatrix} \mathbf{h}_1^\top \\ \mathbf{h}_2^\top \\ \mathbf{h}_3^\top \\ \vdots \\ \mathbf{h}_T^\top \end{pmatrix} = \left[ \exp \left\{ \begin{pmatrix} 0 & -\infty & -\infty & \dots & -\infty \\ b_2 & 0 & -\infty & \dots & -\infty \\ b_2 + b_3 & b_3 & 0 & \dots & -\infty \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{j=2}^T b_j & \sum_{j=3}^T b_j & \sum_{j=4}^T b_j & \dots & 0 \end{pmatrix} \right\} \right] \begin{pmatrix} z_1 \tilde{\mathbf{h}}_1^\top \\ z_2 \tilde{\mathbf{h}}_2^\top \\ z_3 \tilde{\mathbf{h}}_3^\top \\ \vdots \\ z_T \tilde{\mathbf{h}}_T^\top \end{pmatrix} \quad (8)$$

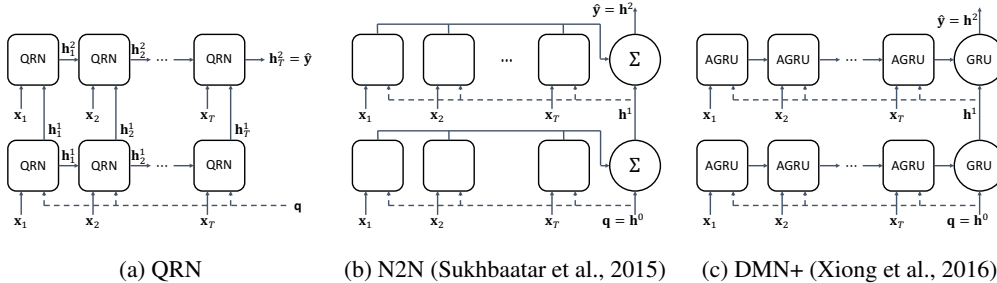


Figure 2: The schematics of QRN and the two state-of-the-art models, End-to-End Memory Networks (N2N) and Improved Dynamic Memory Networks (DMN+), simplified to emphasize the differences among the models. AGRU is a variant of GRU where the update gate is replaced with soft attention, proposed by Kumar et al. (2016). For QRN and DMN+, only forward direction arrows are shown.

Let  $\mathbf{H} = [\mathbf{h}_1^\top; \dots; \mathbf{h}_T^\top]$  be a  $T$ -by- $d$  matrix where the transposes ( $\top$ ) of the column vectors  $\mathbf{h}_t$  are concatenated across row. We similarly define  $\tilde{\mathbf{H}}$  from  $\tilde{\mathbf{h}}_t$ . Also, let  $\mathbf{z} = [z_1; \dots; z_T]$  and  $\mathbf{b} = [0; b_2; \dots; b_T]$  be column vectors (note that we use 0 instead of  $b_1$ ). Then Equation 8 is:

$$\mathbf{H} = [\mathbf{L} \circ \exp(\mathbf{L} [\mathbf{B} \circ \mathbf{L}'])] [\mathbf{Z} \circ \tilde{\mathbf{H}}] \quad (9)$$

where  $\mathbf{L}, \mathbf{L}' \in \mathbb{R}^{T \times T}$  are lower and *strictly* lower triangular matrices of 1's, respectively,  $\circ$  is element-wise multiplication, and  $\mathbf{B}$  is a matrix where  $T$   $\mathbf{b}$ 's are tiled across the column, i.e.  $\mathbf{B} = [\mathbf{b}, \dots, \mathbf{b}] \in \mathbb{R}^{T \times T}$ , and similarly  $\mathbf{Z} = [\mathbf{z}, \dots, \mathbf{z}] \in \mathbb{R}^{T \times d}$ . All implicit operations are matrix multiplications. With reasonable  $N$  (batch size),  $d$  and  $T$  (e.g.  $N, d, T = 100$ ), matrix operations in Equation 9 can be comfortably computed in most modern GPUs.

## 4 RELATED WORK

QRN is inspired by RNN-based models with gating mechanism, such as LSTM (Hochreiter and Schmidhuber, 1997) and GRU (Cho et al., 2014). While GRU and LSTM use the previous hidden state and the current input to obtain the candidate hidden state, QRN only uses the current two inputs to obtain the candidate reduced query (equivalent to candidate hidden state). We conjecture that this not only gives computational advantage via parallelization, but also makes training easier, i.e., avoiding vanishing gradient (which is critical for long-term dependency), overfitting (by simplifying the model), and converging to local minima.

The idea of structurally simplifying (constraining) RNNs for learning longer-term patterns has been explored in recent previous work, such as Structurally Constrained Recurrent Network (Mikolov et al., 2015) and Strongly-Typed Recurrent Neural Network (STRNN) (Balduzzi and Ghifary, 2016). QRN is similar to STRNN in that both architectures use gating mechanism, and the gates and the candidate hidden states do not depend on the previous hidden states, which simplifies the recurrent relation. However, QRN can be distinguished from STRNN in three ways. First, QRN's update gate simulates attention mechanism, measuring the relevance between the input sentence and query. On the other hand, the gates in STRNN can be considered as the simplification of LSTM/GRU by removing their dependency on previous hidden state. Second, QRN is an RNN that is natively compatible with context-based QA tasks, where the QRN unit accepts two inputs, i.e. each context sentence and query. This is distinct from STRNN which has only one input. Third, we show that QRN is timewise-parallelizable on GPUs. Our parallelization algorithm is also applicable to STRNN.

End-to-end Memory Network (N2N) (Sukhbaatar et al., 2015) uses external memory with multi-layer attention mechanism to focus on sentences that are relevant to the question. There are two key differences between N2N and our QRN. First, N2N summarizes the entire memory in each layer to control the attention in the next layer (circle nodes in Figure 2b). Instead, QRN does not have any controller node (Figure 2a) and is able to focus on relevant sentences through the update gate that is internally embodied within its unit. Second, N2N adds time-dependent trainable weights to the sentence representations to model the time dependency of the sentences (as discussed in Section 1). QRN does not need such additional weights as its inherent RNN architecture allows QRN

to effectively model the time dependency. Neural Reasoner (Peng et al., 2015) and Gated End-to-end Memory Network (Perez and Liu, 2016)) are variants of MemN2N that share its fundamental characteristics.

Improved Dynamic Memory Network (DMN+) (Xiong et al., 2016) uses the hybrid of the attention mechanism and the RNN architecture to model the sequence of sentences. It consists of two distinct GRUs, one for the time axis (rectangle nodes in Figure 2c) and one for the layer axis (circle nodes in Figure 2c). Note that the update gate of the GRU for the time axis is replaced with external softmax attention weights. DMN+ uses the time-axis GRU to summarize the entire memory in each layer, and then the layer-axis GRU controls the attention weights in each layer. In contrast, QRN is simply a single recurrent unit without any controller node.

## 5 EXPERIMENTS

### 5.1 DATA

**bAbI story-based QA dataset** bAbI story-based QA dataset (Weston et al., 2016) is composed of 20 different tasks (Appendix A), each of which has 1,000 (1k) synthetically-generated story-question pair. A story can be as short as two sentences and as long as 200+ sentences. A system is evaluated on the accuracy of getting the correct answers to the questions. The answers are single words or lists (e.g. “football, apple”). Answering questions in each task requires selecting a set of relevant sentences and applying different kinds of logical reasoning over them. The dataset also includes 10k training data (for each task), which allows training more complex models. Note that DMN+ (Xiong et al., 2016) only reports on the 10k dataset.

**bAbI dialog dataset** bAbI dialog dataset (Bordes and Weston, 2016) consists of 5 different tasks (Table 3), each of which has 1k synthetically-generated goal-oriented dialogs between a user and the system in the domain of restaurant reservation. Each dialog is as long as 96 utterances and comes with external knowledge base (KB) providing information of each restaurant. The authors also provide Out-Of-Vocabulary (OOV) version of the dataset, where many of the words and KB keywords in test data are not seen during training. A system is evaluated on the accuracy of its response to each utterance of the user, choosing from up to 2500 possible candidate responses. A system is required not only to understand the user’s request but also refer to previous conversations in order to obtain the context information of the current conversation.

**DSTC2 (Task 6) dialog dataset** Bordes and Weston (2016) transformed the Second Dialog State Tracking Challenge (DSTC2) dataset (Henderson et al., 2014) into the same format as the bAbI dialog dataset, for the measurement of performance on a real dataset. Each dialog can be as long as 800+ utterances, and a system needs to choose from 2407 possible candidate responses for each utterance of the user. Note that the evaluation metric of the original DSTC2 is different from that of the transformed DSTC2, so previous work on the original DSTC2 should not be directly compared to our work. We will refer to this transformed DSTC2 dataset by “Task 6” of dialog dataset.

### 5.2 MODEL DETAILS

**Input Module.** In the input module, we are given sentences (previous conversations in dialog)  $\mathbf{x}_t$  and a question (most recent user utterance)  $\mathbf{q}$ , and we want to obtain their vector representations,  $\mathbf{x}_t, \mathbf{q} \in \mathbb{R}^d$ . We use a trainable embedding matrix  $\mathbf{A} \in \mathbb{R}^{d \times V}$  to encode the one-hot vector of each word  $\mathbf{x}_{tj}$  in each sentence  $\mathbf{x}_t$  into a  $d$ -dimensional vector  $\mathbf{x}_{tj} \in \mathbb{R}^d$ . Then the sentence representation  $\mathbf{x}_t$  is obtained by Position Encoder (Weston et al., 2015). The same encoder with the same embedding matrix is also used to obtain the question vector  $\mathbf{q}$  from  $\mathbf{q}$ .

**Output Module for story-based QA.** In the output module, we are given the vector representation of the predicted answer  $\hat{\mathbf{y}}$  and we want to obtain the natural language form of the answer,  $\hat{\mathbf{y}}$ . We use a  $V$ -way single-layer softmax classifier to map  $\hat{\mathbf{y}}$  to a  $V$ -dimensional sparse vector,  $\hat{\mathbf{v}} = \text{softmax}(\mathbf{W}^{(y)}\hat{\mathbf{y}}) \in \mathbb{R}^V$ , where  $\mathbf{W}^{(y)} \in \mathbb{R}^{V \times d}$  is a weight matrix. Then the final answer  $\hat{\mathbf{y}}$  is simply the argmax word in  $\hat{\mathbf{v}}$ . To handle questions with multiple-word answers, we consider each

of them as a single word that contains punctuations such as space and comma, and put it in the vocabulary.

**Output Module for dialog.** We use a fixed number single-layer softmax classifiers, each of which is similar to that of the story-based QA model, to sequentially output each word of the system’s response. While it is similar in spirit to the RNN decoder (Cho et al., 2014), our output module does not have a recurrent hidden state or gating mechanism. Instead, it solely uses the final output of the QRN,  $\hat{y}$ , and the current word output to influence the prediction of the next word among possible candidates.

**Training.** We withhold 10% of the training for development. We use the hidden state size of 50 by default. Batch sizes of 32 for bAbI story-based QA 1k, bAbI dialog and DSTC2 dialog, and 128 for bAbI QA 10k are used. The weights in the input and output modules are initialized with zero mean and the standard deviation of  $1/\sqrt{d}$ . Weights in the QRN unit are initialized using techniques by Glorot and Bengio (2010), and are tied across the layers. Forget bias of 2.5 is used for update gates (no bias for reset gates). L2 weight decay of 0.001 (0.0005 for QA 10k) is used for all weights. The loss function is the cross entropy between  $\hat{v}$  and the one-hot vector of the true answer. The loss is minimized by stochastic gradient descent for maximally 500 epochs, but training is early stopped if the loss on the development data does not decrease for 50 epochs. The learning rate is controlled by AdaGrad (Duchi et al., 2011) with the initial learning rate of 0.5 (0.1 for QA 10k). Since the model is sensitive to the weight initialization, we repeat each training procedure 10 times (50 times for 10k) with the new random initialization of the weights and report the result on the test data with the lowest loss on the development data.

### 5.3 RESULTS.

We compare our model with baselines and previous state-of-the-art models on story-based and dialog tasks (Table 1). These include LSTM (Hochreiter and Schmidhuber, 1997), End-to-end Memory Networks (N2N) (Sukhbaatar et al., 2015), Dynamic Memory Networks (DMN+) (Xiong et al., 2016), Gated End-to-end Memory Networks (GMemN2N) (Perez and Liu, 2016), and Differentiable Neural Computer (DNC) (Graves et al., 2016).

**Story-based QA.** Table 1(top) reports the summary of results of our model (QRN) and previous work on bAbI QA (task-wise results are shown in Table 2 in Appendix). In 1k data, QRN’s ‘2r’ (2 layers + reset gate +  $d = 50$ ) outperforms all other models by a large margin (2.8+%). In 10k dataset, the average accuracy of QRN’s ‘6r200’ (6 layers + reset gate +  $d = 200$ ) model outperforms all previous models by a large margin (2.5+%), achieving a nearly perfect score of 99.7%.

**Dialog.** Table 1(bottom) reports the summary of the results of our model (QRN) and previous work on bAbI dialog and Task 6 dialog (task-wise results are shown in Table 3 in Appendix). As done in previous work (Bordes and Weston, 2016; Perez and Liu, 2016), we also report results when we use ‘Match’ for dialogs. ‘Match’ is the extension to the model which additionally takes as input whether each answer candidate matches with context (more details on Appendix). QRN outperforms previous work by a large margin (2.0+%) in every comparison.

**Ablations.** We test four types of ablations (also discussed in Section 2.2): number of layers (1, 2, 3, or 6), reset gate (r), and gate vectorization (v) and the dimension of the hidden vector (50, 100). We show a subset of combinations of the ablations for bAbI QA in Table 1 and Table 2; other combinations performed poorly and/or did not give interesting observations. According to the ablation results, we infer that: **(a)** When the number of layers is only one, the model lacks reasoning capability. In the case of 1k dataset, when there are too many layers (6), it seems correctly training the model becomes increasingly difficult. In the case of 10k dataset, many layers (6) and hidden dimensions (200) helps reasoning, most notably in difficult task such as task 16. **(b)** Adding the reset gate helps. **(c)** Including vector gates hurts in 1k datasets, as the model either overfits to the training data or converges to local minima. On the other hand, vector gates in bAbI story-based QA 10k dataset sometimes help. **(d)** Increasing the dimension of the hidden state to 100 in the dialog’s Task 6 (DSTC2) helps, while there is not much improvement in the dialog’s Task 1-5. It can be hypothesized that a larger hidden state is required for real data.

Task	1k						10k				
	Previous works				QRN		Previous works				QRN
	LSTM	N2N	DMN+ <sup>†</sup>	GMemN2N	2r	3r	N2N	DMN+	GMemN2N	DNC	6r200
# Failed	20	10	16	10	7	5	3	1	3	2	0
Average error rates	51.3	15.2	33.2	12.7	9.9	11.3	4.2	2.8	3.7	3.8	0.3

Task	Plain				With Match		
	Previous works		QRN		Previous works		QRN
	N2N	GMemN2N	2r	2r100	N2N+	GMemN2N+	2r+
bAbI dialog Average error rates	13.9	14.3	5.5	5.5	6.7	5.4	1.5
bAbI dialog (OOV) Average error rates	30.3	27.9	11.1	11.1	11.2	10.3	2.3
DSTC2 dialog Average error rates	58.9	52.6	49.5	48.9	59.0	51.3	49.3

Table 1: (top) bAbI QA dataset (Weston et al., 2016): number of failed tasks and average error rates (%). <sup>†</sup> is obtained from [github.com/therne/dmn-tensorflow](https://github.com/therne/dmn-tensorflow). (bottom) bAbI dialog and DSTC2 dialog dataset (Bordes and Weston, 2016) average error rates (%) of QRN and previous work (LSTM, N2N, DMN+, GMemN2N, and DNC). For QRN, the first number (1, 2, 3) indicates the number of layers, ‘r’ means the reset gate is used, and the last number (100, 200), if exists, indicates the dimension of the hidden state, where the default value is 50. ‘+’ indicates that ‘match’ (See Appendix for details) is used. The task-wise results are shown in Appendices: Table 2 (bAbI QA) and Table 3 (dialog datasets). See Section 5.3 for details.

Task 2: Two Supporting Facts	Layer 1			Layer 2	Task 15: Deduction	Layer 1			Layer 2
	$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$		$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$
Sandra picked up the apple there.	0.95	0.89	0.98	0.00	Mice are afraid of wolves.	0.11	0.99	0.13	0.78
Sandra dropped the apple.	0.83	0.05	0.92	0.01	Gertrude is a mouse.	0.77	0.99	0.96	0.00
Daniel grabbed the apple there.	0.88	0.93	0.98	0.00	Cats are afraid of sheep.	0.01	0.99	0.07	0.03
Sandra travelled to the bathroom.	0.01	0.18	0.63	0.02	Winona is a mouse.	0.14	0.85	0.77	0.05
Daniel went to the hallway.	0.01	0.24	0.62	0.83	Sheep are afraid of wolves.	0.02	0.98	0.27	0.05
Where is the apple?	hallway				What is Gertrude afraid of?	wolf			

Task 3: Displaying options	Layer 1			Layer 2	Task 6: DSTC2 dialog	Layer 1			Layer 2
	$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$		$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$
resto-paris-expen-frech-8stars?	0.00	1.00	0.96	0.91	Spanish food.	0.84	0.07	0.00	0.82
Do you have something else?	0.41	0.99	0.00	0.00	You are lookng for a spanish restaurant right?	0.98	0.02	0.49	0.75
Sure let me find another option.	1.00	0.00	0.00	0.12	Yes.	0.01	1.00	0.33	0.13
resto-paris-expen-frech-5stars?	0.00	1.00	0.96	0.91	What part of town do you have in mind?	0.20	0.73	0.41	0.11
No this does not work for me.	0.00	0.00	0.14	0.00	I don't care.	0.00	1.00	0.02	0.00
Sure let me find an other option.	1.00	0.00	0.00	0.12	What price range would you like?	0.72	0.46	0.52	0.72
What do you think of this? resto-paris-expen-french-4stars					I don't care.	API CALL spanish R-location R-price			

Figure 3: (top) bAbI QA dataset (Weston et al., 2016) visualization of update and reset gates in QRN ‘2r’ model (bottom two) bAbI dialog and DSTC2 dialog dataset (Bordes and Weston, 2016) visualization of update and reset gates in QRN ‘2r’ model. Note that the stories can have as many as 800+ sentences; we only show part of them here. More visualizations are shown in Figure 4 (bAbI QA) and Figure 5 (dialog datasets).

**Parallelization.** We implement QRN with and without parallelization in TensorFlow (Abadi et al., 2016) on a single Titan X GPU to quantify the computational gain of the parallelization. For QRN without parallelization, we use the RNN library provided by TensorFlow. QRN with parallelization gives 6.2 times faster training and inference than QRN without parallelization on average. We expect that the speedup can be even higher for datasets with larger context.

**Interpretations.** An advantage of QRN is that the intermediate query updates are interpretable. Figure 1 shows intermediate local queries ( $q_t^k$ ) interpreted in natural language, such as “Where is Sandra?”. In order to obtain these, we place a decoder on the input question embedding  $q$  and add its loss for recovering the question to the classification loss (similarly to Peng et al. (2015)). We then use the same decoder to decode the intermediate queries. This helps us understand the flow of information in the networks. In Figure 1, the question Where is apple? is transformed into Where is Sandra? at  $t = 1$ . At  $t = 2$ , as Sandra dropped the apple, the apple is no more relevant to Sandra. We obtain Where is Daniel? at time  $t = 3$ , and it is propagated until  $t = 5$ , where we observe a sentence (fact) that can be used to answer the query.

**Visualization.** Figure 3 shows visualization of the (scalar) magnitudes of update and reset gates on story sentences and dialog utterances. More visualizations are shown in Appendices: Figure 4 and Figure 5. In Figure 3, we observe high values on facts that provide information to answer question (the system’s next utterance for dialog). In QA Task 2 example (top left), we observe high update gate values in the first layer on facts that state who has the apple, and in the second layer, the high update gate values are on those that inform where that person went to. We also observe that the forward reset gate at  $t = 2$  in the first layer ( $\vec{r}_2^1$ ) is low, which is signifying that apple no more



belongs to Sandra. In dialog Task 3 (bottom left), the model is able to infer that three restaurants are already recommended so that it can recommend another one. In dialog Task 6 (bottom), the model focuses on the sentences containing Spanish, and does not concentrate much on other facts such as I don't care.

## 6 CONCLUSION

In this paper, we introduce Query-Reduction Network (QRN) to answer context-based questions and carry out conversations with users that require multi-hop reasoning. We show the state-of-the-art results in the three datasets of story-based QA and dialog. We model a story or a dialog as a sequence of state-changing triggers and compute the final answer to the question or the system's next utterance by recurrently updating (or *reducing*) the query. QRN is situated between the attention mechanism and RNN, effectively handling time dependency and long-term dependency problems of each technique, respectively. It addresses the long-term dependency problem of most RNNs by simplifying the recurrent update, in which the candidate hidden state (reduced query) does not depend on the previous state. Moreover, QRN can be parallelized and can address the well-known problem of RNN's vanishing gradients.

## ACKNOWLEDGMENTS

This research was supported by the NSF (IIS 1616112), Allen Institute for AI (66-9175), Allen Distinguished Investigator Award, Google Research Faculty Award, and Samsung GRO Award. We thank the anonymous reviewers for their helpful comments.

## REFERENCES

- Martin Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- David Balduzzi and Muhammad Ghifary. Strongly-typed recurrent neural networks. In *ICML*, 2016.
- Antoine Bordes and Jason Weston. Learning end-to-end goal-oriented dialog. *arXiv preprint arXiv:1605.07683*, 2016.
- Kyunghyun Cho, Bart Van Merriënboer, Caglar Gulcehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *EMNLP*, 2014.
- John Duchi, Elad Hazan, and Yoram Singer. Adaptive subgradient methods for online learning and stochastic optimization. *JMLR*, 12, 2011.
- Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *JMLR*, 2010.
- Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, et al. Hybrid computing using a neural network with dynamic external memory. *Nature*, 2016.
- Matthew Henderson, Blaise Thomson, and Jason Williams. The second dialog state tracking challenge. In *SIGdial*, 2014.
- Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. In *NIPS*, 2015.
- Felix Hill, Antoine Bordes, Sumit Chopra, and Jason Weston. The goldilocks principle: Reading children's books with explicit memory representations. In *ICLR*, 2016.
- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.

- Ankit Kumar, Ozan Irsoy, Jonathan Su, James Bradbury, Robert English, Brian Pierce, Peter Ondruska, Ishaan Gulrajani, and Richard Socher. Ask me anything: Dynamic memory networks for natural language processing. In *ICML*, 2016.
- Tomas Mikolov, Armand Joulin, Sumit Chopra, Michael Mathieu, and Marc’Aurelio Ranzato. Learning longer memory in recurrent neural networks. In *ICLR 2015 Workshop*, 2015.
- Baolin Peng, Zhengdong Lu, Hang Li, and Kam-Fai Wong. Towards neural network-based reasoning. *arXiv preprint arXiv:1508.05508*, 2015.
- Julien Perez and Fei Liu. Gated end-to-end memory networks. *arXiv preprint arXiv:1610.04211*, 2016.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. In *EMNLP*, 2016.
- Raymond Reiter. *Knowledge in Action*. MIT Press, 1st edition, 2001.
- Matthew Richardson, Christopher JC Burges, and Erin Renshaw. Mctest: A challenge dataset for the open-domain machine comprehension of text. In *EMNLP*, 2013.
- Sainbayar Sukhbaatar, Arthur Szlam, Jason Weston, and Rob Fergus. End-to-end memory networks. In *NIPS*, 2015.
- Jason Weston, Sumit Chopra, and Antoine Bordes. Memory networks. In *ICLR*, 2015.
- Jason Weston, Antoine Bordes, Sumit Chopra, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. In *ICLR*, 2016.
- Caiming Xiong, Stephen Merity, and Richard Socher. Dynamic memory networks for visual and textual question answering. In *ICML*, 2016.

## A TASK-WISE RESULTS

Here we provide detailed per-task breakdown of our results in QA(Table 2) and dialog datasets (Table 3).

Task	1k														10k					
	Previous works				QRN						Previous works			QRN						
	LSTM	N2N	DMN+	GMemN2N	1r	2	2r	3r	6r	6r200*	N2N	DMN+	GMemN2N	2r	2rv	3r	6r200			
1: Single supporting fact	50.0	0.1	1.3	0.0	0.0	0.0	0.0	0.0	13.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0			
2: Two supporting facts	80.0	18.8	72.3	8.1	65.7	1.2	0.7	0.5	1.5	15.3	0.3	0.3	0.0	0.4	0.8	0.4	0.0			
3: Three supporting facts	80.0	31.7	73.3	38.7	68.2	17.5	5.7	1.2	15.3	13.8	2.1	1.1	4.5	0.4	1.4	0.0	0.0			
4: Two arg relations	39.0	17.5	26.9	0.4	0.0	0.0	0.0	0.7	9.0	13.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0			
5: Three arg relations	30.0	12.9	25.6	1.0	1.0	1.1	1.1	1.2	1.3	12.5	0.8	0.5	0.2	0.5	0.2	0.3	0.0			
6: Yes/no questions	52.0	2.0	28.5	8.4	0.1	0.0	0.9	1.2	50.6	15.5	0.1	0.0	0.0	0.0	0.0	0.0	0.0			
7: Counting	51.0	10.1	21.9	17.8	10.9	11.1	9.6	9.4	13.1	15.3	2.0	2.4	1.8	1.0	0.7	0.7	0.0			
8: Lists/sets	55.0	6.1	21.9	12.5	6.8	5.7	5.6	3.7	7.8	15.1	0.9	0.0	0.3	1.4	0.6	0.8	0.4			
9: Simple negation	36.0	1.5	42.9	10.7	0.0	0.6	0.0	0.0	32.7	13.0	0.3	0.0	0.0	0.0	0.0	0.0	0.0			
10: Indefinite knowledge	56.0	2.6	23.1	16.5	0.8	0.6	0.0	0.0	3.5	12.9	0.0	0.0	0.2	0.0	0.0	0.0	0.0			
11: Basic coreference	38.0	3.3	4.3	0.0	11.3	0.5	0.0	0.0	0.9	14.7	0.1	0.0	0.0	0.0	0.0	0.0	0.0			
12: Conjunction	26.0	0.0	3.5	0.0	0.0	0.0	0.0	0.0	0.0	15.1	0.0	0.0	0.0	0.0	0.0	0.0	0.0			
13: Compound coreference	6.0	0.5	7.8	0.0	5.3	5.5	0.0	0.3	8.9	13.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0			
14: Time reasoning	73.0	2.0	61.9	1.2	20.2	1.3	0.8	3.8	18.2	14.5	0.1	0.0	0.0	0.2	0.0	0.0	0.1			
15: Basic deduction	79.0	1.8	47.6	0.0	39.4	0.0	0.0	0.0	0.1	14.7	0.0	0.0	0.0	0.0	0.0	0.0	0.0			
16: Basic induction	77.0	51.0	54.4	0.1	50.6	54.8	53.0	53.4	53.5	15.5	51.8	45.3	0.0	49.4	50.4	49.1	0.0			
17: Positional reasoning	49.0	42.6	44.1	41.7	40.6	36.5	34.4	51.8	52.0	13.0	18.6	4.2	27.8	0.9	0.0	5.8	4.1			
18: Size reasoning	48.0	9.2	9.1	9.2	8.2	8.6	7.9	8.8	47.5	14.9	5.3	2.1	8.5	1.6	8.4	1.8	0.7			
19: Path finding	92.0	90.6	90.8	88.5	88.8	89.8	78.7	90.7	88.6	13.6	2.3	0.0	31.0	36.1	1.0	27.9	0.1			
20: Agents motivations	9.0	0.2	2.2	0.0	0.0	0.0	0.2	0.3	5.5	14.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0			
# Failed	20	10	16	10	12	8	7	5	13	20	3	1	3	2	2	3	0			
Average error rates (%)	51.3	15.2	33.2	12.7	20.1	11.7	<b>9.9</b>	11.3	20.5	14.2	4.2	2.8	3.7	4.6	3.2	4.3	<b>0.3</b>			

Table 2: bAbI QA dataset (Weston et al., 2016) error rates (%) of QRN and previous work: LSTM (Weston et al., 2016), End-to-end Memory Networks (N2N) (Sukhbaatar et al., 2015), Dynamic Memory Networks (DMN+) (Xiong et al., 2016), Gated End-to-end Memory Networks(GMemN2N) (Perez and Liu, 2016). Results within each task of Differentiable Neural Computer(DNC) were not provided in its paper Graves et al. (2016)). For QRN, a number in the front (1, 2, 3, 6) indicates the number of layers. A number in the back (200) indicates the dimension of hidden vector, while the default value is 50. ‘r’ indicates that the reset gate is used, and ‘v’ indicates that the gates were vectorized. ‘\*’ indicates joint training.

Task	Plain						With Match		
	Previous works		QRN				Previous works		QRN
	N2N	GMemN2N	1r	2r	2r100	2rv	N2N+	GMemN2N+	2r+
1: Issuing API calls	0.1	<b>0.0</b>	0.02	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
2: Updating API calls	<b>0.0</b>	<b>0.0</b>	0.11	0.01	<b>0.0</b>	<b>0.0</b>	1.7	<b>0.0</b>	<b>0.0</b>
3: Displaying options	25.1	25.1	<b>12.6</b>	<b>12.6</b>	<b>12.6</b>	<b>12.6</b>	25.1	25.1	<b>7.6</b>
4: Providing extra information	40.5	42.8	<b>14.3</b>	<b>14.3</b>	<b>14.3</b>	<b>14.3</b>	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
5: Conducting full dialogs	3.9	3.7	9.4	<b>0.6</b>	0.7	0.9	6.6	2.0	<b>0.0</b>
Average error rates (%)	13.9	14.3	7.3	<b>5.5</b>	<b>5.5</b>	5.6	6.7	5.4	<b>1.5</b>
1 (OOV): Issuing API calls	27.7	17.6	<b>6.6</b>	<b>6.6</b>	<b>6.6</b>	<b>6.6</b>	3.5	<b>0.0</b>	<b>0.0</b>
2 (OOV): Updating API calls	21.1	21.1	8.5	8.4	<b>8.4</b>	8.5	5.5	5.8	<b>0.0</b>
3 (OOV): Displaying options	25.6	24.7	<b>12.4</b>	<b>12.4</b>	<b>12.4</b>	12.5	24.8	24.9	<b>7.7</b>
4 (OOV): Providing extra information	42.4	43.0	<b>14.3</b>	<b>14.3</b>	14.4	14.4	<b>0.0</b>	<b>0.0</b>	<b>0.0</b>
5 (OOV): Conducting full dialogs	34.5	33.3	19.5	14.0	13.71	<b>13.6</b>	22.3	20.6	<b>4.0</b>
Average error rates (%)	30.3	27.9	12.3	<b>11.1</b>	<b>11.1</b>	<b>11.1</b>	11.2	10.3	<b>2.3</b>
6: DSTC2 dialog	58.9	52.6	49.9	49.5	<b>48.9</b>	53.8	59.0	51.3	<b>49.3</b>

Table 3: bAbI dialog and DSTC2 dialog dataset (Bordes and Weston, 2016) average error rates (%) of QRN and previous work: End-to-end Memory Networks(N2N (Bordes and Weston, 2016)) and Gated End-to-end Memory Networks(GMemN2N (Perez and Liu, 2016)). For QRN, a number in the front (1, 2, 3, 6) indicates the number of layers and a number in the back (100) indicates the dimension of hidden vector, while the default value is 50. ‘r’ indicates that the reset gate is used, ‘v’ indicates that the gates were vectorized, and ‘+’ indicates that ‘match’ was used.

## B VECTOR GATE PARALLELIZATION

For vector gates, we have  $\mathbf{z}_t \in \mathbb{R}^d$  instead of  $z_t \in \mathbb{R}$ . Therefore the following equation replaces Equation 7:

$$\mathbf{h}_t = \sum_{i=1}^t \exp \left\{ \begin{pmatrix} \sum_{j=i+1}^t \log(1 - z_j^1) \\ \sum_{j=i+1}^t \log(1 - z_j^2) \\ \vdots \\ \sum_{j=i+1}^t \log(1 - z_j^d) \end{pmatrix} \right\} \circ \mathbf{z}_i \circ \tilde{\mathbf{h}}_i \quad (10)$$

where  $z_j^k$  is the  $k$ -th column vector of  $z_j$ . Let  $b_{ij} = \log(1 - z_i^j)$  for brevity. Then, we can rewrite Equation 8 as following:

$$\begin{pmatrix} \mathbf{h}_1^\top \\ \mathbf{h}_2^\top \\ \mathbf{h}_3^\top \\ \vdots \\ \mathbf{h}_T^\top \end{pmatrix}^j = \left[ \exp \left\{ \begin{pmatrix} 0 & -\infty & -\infty & \dots & -\infty \\ b_{2j} & 0 & -\infty & \dots & -\infty \\ b_{2j} + b_{3j} & b_{3j} & 0 & \dots & -\infty \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \sum_{k=2}^T b_{kj} & \sum_{k=3}^T b_{kj} & \sum_{k=4}^T b_{kj} & \dots & 0 \end{pmatrix} \right\} \right] \begin{pmatrix} \mathbf{z}_1 \circ \tilde{\mathbf{h}}_1^\top \\ \mathbf{z}_2 \circ \tilde{\mathbf{h}}_2^\top \\ \mathbf{z}_3 \circ \tilde{\mathbf{h}}_3^\top \\ \vdots \\ \mathbf{z}_T \circ \tilde{\mathbf{h}}_T^\top \end{pmatrix}^j \quad (11)$$

Let  $\mathbf{H} = [\mathbf{h}_1^\top; \dots; \mathbf{h}_T^\top]$  be a  $T$ -by- $d$  matrix where the transposes ( $\top$ ) of the column vectors  $\mathbf{h}_t$  are concatenated across row. We similarly define  $\tilde{\mathbf{H}}$  from  $\tilde{\mathbf{h}}_t$ . Also, let  $\mathbf{z} = [z_1; \dots; z_T]$ , and  $\mathbf{B}_d$  be a  $T$ -by- $T$  matrix where  $T$   $[0; b_{2d}; \dots; b_{Td}]$ 's are tiled across the column.

Then Equation 11 is:

$$\mathbf{H} = \begin{pmatrix} [\mathbf{L} \circ \exp(\mathbf{L} [\mathbf{B}_1 \circ \mathbf{L}'])] [\mathbf{Z} \circ \tilde{\mathbf{H}}]^1 \\ [\mathbf{L} \circ \exp(\mathbf{L} [\mathbf{B}_2 \circ \mathbf{L}'])] [\mathbf{Z} \circ \tilde{\mathbf{H}}]^2 \\ \vdots \\ [\mathbf{L} \circ \exp(\mathbf{L} [\mathbf{B}_d \circ \mathbf{L}'])] [\mathbf{Z} \circ \tilde{\mathbf{H}}]^d \end{pmatrix} \quad (12)$$

where  $\mathbf{L}, \mathbf{L}' \in \mathbb{R}^{T \times T}$  are lower and *strictly* lower triangular matrices of 1's are tiled across the column.  $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_d] \in \mathbb{R}^{T \times d}$ .

## C MODEL DETAILS

**Match.** While similar in spirit, our ‘Match’ model is slightly different from previous work (Bordes and Weston, 2016; Perez and Liu, 2016). We use answer candidate embedding matrix, and add 2 dimension of 0-1 matrix which expresses whether the answer candidate matches with any word in the paragraph and the question. In other words, the softmax is computed by  $\hat{\mathbf{v}} = \text{softmax}(\mathbf{W}[\mathbf{W}^{(y)}; \mathbf{M}^{(y)}]\hat{\mathbf{y}}) \in \mathbb{R}^V$ , where  $\mathbf{W} \in \mathbb{R}^{d \times d}$  and  $\mathbf{W}^{(y)} \in \mathbb{R}^{V \times (d-2)}$  are trainable weight matrices, and  $\mathbf{M}^{(y)} \in \mathbb{R}^{V \times 2}$  is the 0-1 match matrix.

## D VISUALIZATIONS

**Visualization of Story-based QA.** Figure 4 shows visualization of models for story-based QA tasks.

In the task 3 (left), the model focuses on the facts that contain ‘football’ in the first layer, and found out where Mary journeyed to before the bathroom in the second layer. In task 7 (right), the model focuses on the facts that provide information about the location of Sandra.

Task 3: Three Supporting Facts	Layer 1			Layer 2	Task 7: Counting	Layer 1			Layer 2
	$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$		$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$
Mary got the football there.	0.82	1.00	0.0	0.06	Mary journeyed to the garden.	0.67	0.08	0.58	0.12
John went back to the bedroom.	0.01	0.00	0.72	0.57	Mary journeyed to the office.	0.91	0.44	0.11	0.21
Mary journeyed to the office.	0.01	0.04	0.06	0.88	Sandra grabbed the apple there.	0.02	0.34	0.92	0.89
Mary journeyed to the bathroom.	0.44	0.00	0.89	0.05	Sandra discarded the apple.	0.26	0.61	0.95	0.97
Mary dropped the football.	0.62	0.01	0.00	0.03	Daniel went to the bedroom.	0.70	0.44	0.99	0.03
Where was the football before the bathroom?				office	How many objects is Sandra carrying?				none

Figure 4: Visualization of update and reset gates in QRN ‘2r’ model for on several tasks of bAbI QA (Table 2). We do not put reset gate in the last layer. Note that we only show some of recent sentences here, though the stories can have as many as 200+ sentences.

Task 1 Issuing API calls	Layer 1			Layer 2	Task 1 Issuing API calls	Layer 1			Layer 2
	$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$		$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$
Good morning.	0.12	0.34	0.98	0.20	Can you make a restaurant reservation for eight in a cheap price range in madrid	0.00	1.00	0.93	1.00
Hello what can i help you with today.	0.97	0.97	0.12	0.12	I’m on it.	0.00	1.00	0.74	0.00
Can you book a table in rome with italian cuisine.	0.00	0.87	1.00	1.00	Any preference on a type of cuisine.	0.00	0.11	1.00	0.01
I’m on it.	0.73	0.97	0.38	0.00	I love british food.	0.00	0.99	0.99	0.57
How many people would you in your party.	1.00	1.00	0.00	0.41	Okay let me look into some options for you.	1.00	0.00	0.00	0.02
For four people please.					<SILENCE>				
Which price range are you looking for.					API CALL british madrid eight cheap				
Task 4 Providing extra-information	Layer 1			Layer 2	Task 4 Providing extra-information	Layer 1			Layer 2
	$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$		$z^1$	$\vec{r}^1$	$\overleftarrow{r}^1$	$z^2$
resto-paris-expen-spanish-8stars R-phone resto-paris-expen-spanish-8stars-phone	0.71	0.84	0.99	0.36	resto-paris-expen-spanish-8stars R-address resto-paris-expen-spanish-8stars-address	1.00	0.99	1.00	1.00
resto-paris-expen-spanish-8stars R-location paris	0.05	0.01	1.00	0.00	resto-paris-expen-spanish-8stars R-number four	0.02	0.95	0.97	0.00
resto-paris-expen-spanish-8stars R-price expensive	0.00	0.05	0.92	0.00	resto-paris-expen-spanish-8stars R-rating 8	0.38	0.91	1.00	0.10
resto-paris-expen-spanish-8stars R-rating 8	0.38	0.91	1.00	0.10	What do you think of this option: resto-paris-expen-spanish-8stars	0.90	0.93	0.99	1.00
Let’s do it.	0.00	0.00	1.00	0.00	Great let me do the reservation.	0.98	0.99	0.97	0.00
Do you have its address.					Here it is: resto-paris-expen-spanish-8stars-address				

Figure 5: Visualization of update and reset gates in QRN ‘2r’ model for on several tasks of bAbI dialog and DSTC2 dialog (Table 3). We do not put reset gate in the last layer. Note that we only show some of recent sentences here, even the dialog has more sentences.

**Visualization of Dialog.** Figure 5 shows visualization of models for dialog tasks.

In the first dialog of task 1, the model focuses on the user utterance that mentions the user’s desired cuisine and location, and the current query (user’s last utterance) informs the system of the number of people, so the system is able to learn that it now needs to ask the user about the desired price range. In the second dialog of task 1, the model focuses on the facts that provide information about the requests of the user. In task 4 (third), the model focuses on what restaurant a user is talking about and the information about the restaurant.