

TOWARDS NEURAL NETWORK-BASED REASONING

Baolin Peng^{1*} Zhengdong Lu² Hang Li² Kam-Fai Wong¹
¹Dept. of Systems Engineering & Engineering Management,
 The Chinese University of Hong Kong
 {blpeng, kfwong}@se.cuhk.edu.hk
²Noah's Ark Lab, Huawei Technologies
 {Lu.Zhengdong, HangLi.HL}@huawei.com

Abstract

We propose NEURAL REASONER, a framework for neural network-based reasoning over natural language sentences. Given a question, NEURAL REASONER can infer over multiple supporting facts and find an answer to the question in specific forms. NEURAL REASONER has 1) a specific interaction-pooling mechanism, allowing it to examine multiple facts, and 2) a deep architecture, allowing it to model the complicated logical relations in reasoning tasks. Assuming no particular structure exists in the question and facts, NEURAL REASONER is able to accommodate different types of reasoning and different forms of language expressions. Despite the model complexity, NEURAL REASONER can still be trained effectively in an end-to-end manner. Our empirical studies show that NEURAL REASONER can outperform existing neural reasoning systems with remarkable margins on two difficult artificial tasks (Positional Reasoning and Path Finding) proposed in [8]. For example, it improves the accuracy on Path Finding(10K) from 33.4% [6] to over 98%.

1 Introduction

Reasoning is essential to natural language processing tasks, most obviously in examples like document summarization, question-answering, and dialogue. Previous efforts in this direction are built on rule-based models, requiring first mapping natural languages to logic forms and then inference over them. **The mapping (roughly corresponding to semantic parsing), and the inference, are by no means easy, given the variability and flexibility of natural language, the variety of the reasoning tasks, and the brittleness of a rule-based system.**

Just recently, there is some new effort, mainly represented by Memory Network and its dynamic variants [9, 5], trying to build a purely neural network-based reasoning system with fully distributed semantics that can infer over multiple facts to answer simple questions, all in natural language, e.g.,

Fact1: John travelled to the hallway.

Fact2: Mary journeyed to the bathroom.

Question: Where is Mary?

The Memory Nets perform fairly well on simple tasks like the examples above, but poorly

*The work is done when the first author worked as intern at Noah's Ark Lab, Huawei Technologies.

on more complicated ones due to their simple and rigid way of modeling the dynamics of question-fact interaction and the complex process of reasoning.

In this paper we give a more systematic treatment of the problem and propose a flexible neural reasoning system, named NEURAL REASONER. It is purely neural network based and can be trained in an end-to-end way [6], using only supervision from the final answer. Our contributions are mainly two-folds

- we propose a novel neural reasoning system NEURAL REASONER that can infer over multiple facts in a way insensitive to 1) the number of supporting facts, 2) the form of language, and 3) the type of reasoning;
- we give a particular instantiation of NEURAL REASONER and a multi-task training method for effectively fitting the model with relatively small amount of data, yielding significantly better results than existing neural models on two artificial reasoning task;

2 Overview of Neural Reasoner

NEURAL REASONER has a layered architecture to deal with the complicated logical relations in reasoning, as illustrated in Figure 1. It consists of one encoding layer and multiple reasoning layers. The encoder layer first converts the question and facts from natural language sentences to vectorial representations. More specifically,

$$Q \xrightarrow{\text{encode}} \mathbf{q}^{(0)}, \quad F_k \xrightarrow{\text{encode}} \mathbf{f}_k^{(0)}, \quad k = 1, 2, \dots, K.$$

where $\mathbf{q}^{(0)} \in \mathbb{R}^{d_Q}$ and $\mathbf{f}_k^{(0)} \in \mathbb{R}^{d_F}$. With the representations obtained from the encoding layer, the reasoning layer recursively updates the representations of questions and facts,

$$\{\mathbf{q}^{(\ell)} \quad \mathbf{f}_1^{(\ell)} \quad \dots \quad \mathbf{f}_K^{(\ell)}\} \xrightarrow{\text{reason}} \{\mathbf{q}^{(\ell+1)} \quad \mathbf{f}_1^{(\ell+1)} \quad \dots \quad \mathbf{f}_K^{(\ell+1)}\}$$

through the interaction between question representation and fact representations. Intuitively, this interaction models the reasoning, including examination of the facts and comparison of the facts and the questions. Finally at layer- L , the resulted question representation $\mathbf{q}^{(L)}$ is fed to an answerer, which layer can be a classifier for choosing between a number of pre-determined classes (e.g., {Yes, No}) or a text generator for create a sentence.

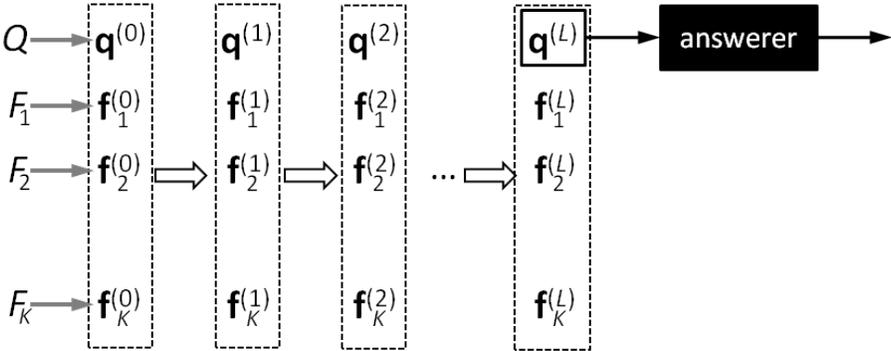


Figure 1: High level system diagram of NEURAL REASONER .

We argue that NEURAL REASONER has the following desired properties:

- it can handle varying number of facts, including irrelevant ones, and reach the final conclusion through repeated processing of filtering and combining;
- it makes no assumption about the form of language, as long as enough training examples are given.

3 Model

In this section we give an instantiation of NEURAL REASONER described in Section 2, as illustrated in Figure 2. In a nutshell, question and facts, as symbol sequences, are first converted to vectorial representations in the encoding layer via recurrent neural networks (RNNs). The vectorial representations are then fed to the reasoning layers, where the question and the facts get updated through an nonlinear transformation jointly controlled by deep neural networks (DNNs) and pooling. Finally at the answering layer, the resulted question representation is used to generate the final answer to the question. More specifically

- in the encoding layer (Layer-0) we use recurrent neural networks (RNNs) to convert question and facts to their vectorial representations, which are then forwarded to the first reasoning layer;
- in each reasoning layer (i.e., Layer- ℓ with $1 \leq \ell \leq L - 1$), we use a deep neural network (denoted as DNN_ℓ) to model the pairwise interaction between question representation $\mathbf{q}^{(\ell-1)}$ and each fact representation $\mathbf{f}_k^{(\ell-1)}$ from the previous layer, which yields updated fact representation $\mathbf{f}_k^{(\ell)}$ and updated (fact-dependent) question representation $\mathbf{q}_k^{(\ell)}$;
- we then fuse the individual updated fact representations $\{\mathbf{q}_1^{(\ell)}, \mathbf{q}_2^{(\ell)}, \dots, \mathbf{q}_K^{(\ell)}\}$ for the global updated representation $\mathbf{q}^{(\ell)}$ through a pooling operation (see Section 3.2 for more details)
- finally in Layer- L , the interaction net (DNN_L) returns only question update, which, after summarization by the pooling operation, will serve as input to the Answering Layer.

In the rest of this section, we will give details of different components of the model.

3.1 Encoding Layer

The encoding layer is designed to find semantic representations of question and facts. Suppose that we are given a fact or a question as word sequence $\{x_1, \dots, x_T\}$, the encoding module summarizes the word sequence with a vector with fixed length. We have different modeling choices for this purpose, e.g., CNN [4] and RNN [7], while in this paper we use GRU [2], a variant of RNN, as the encoding module. GRU is shown to be able to alleviate the gradient vanishing issue of RNN and have similar performance to the more complicated LSTM [3].

As shown in Figure 3, GRU takes as input a sequence of word vectors (for either question or facts)

$$\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_T\}, \mathbf{x}_i \in \mathbb{R}^{|\mathcal{V}|} \tag{1}$$

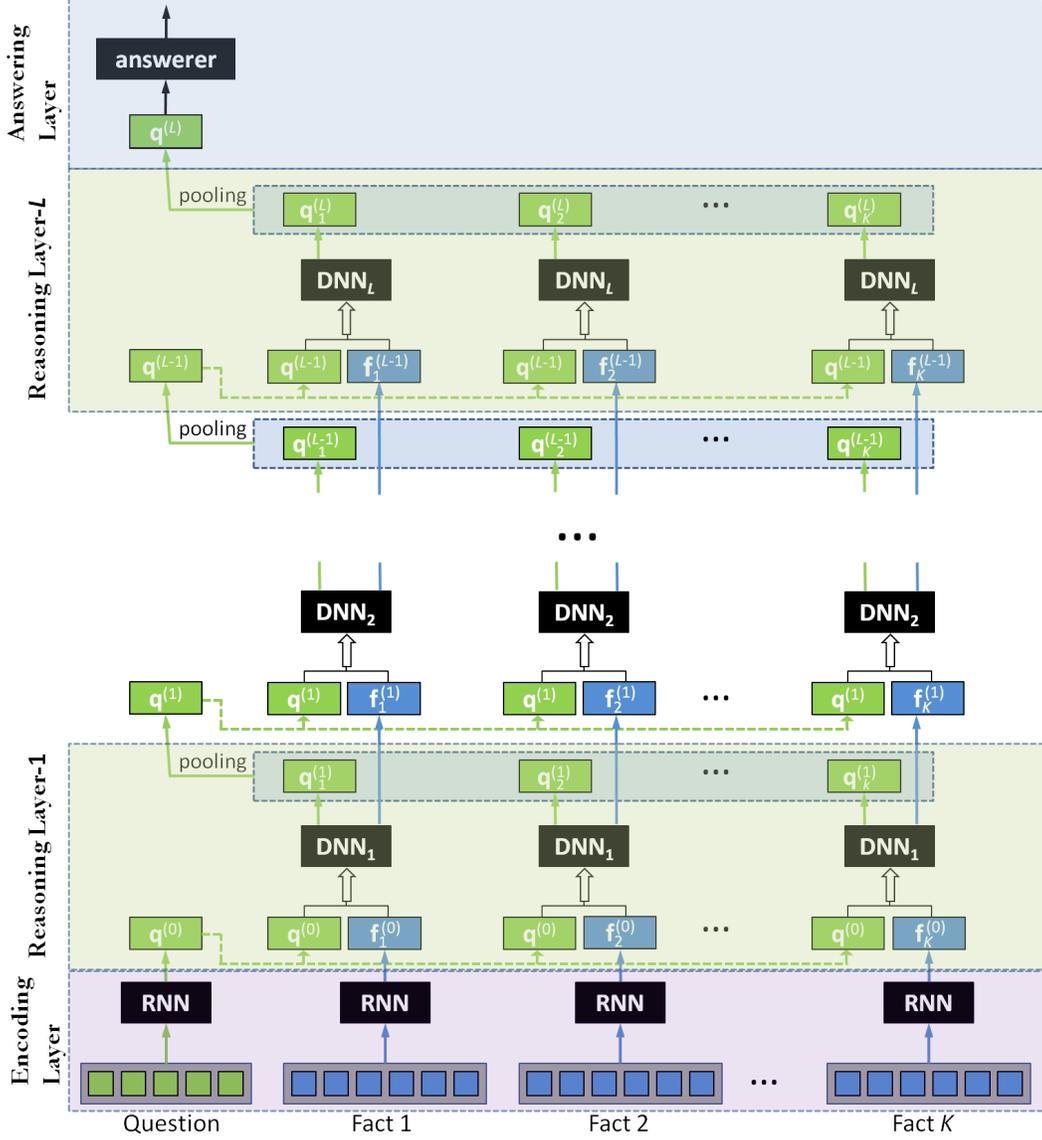


Figure 2: A diagram of our implementation of NEURAL REASONER with L reasoning layers, operating on one question and K facts.

where $|\mathcal{V}|$ stands for the size of vocabulary for input sentences. Detailed forward computations are as follows:

$$\mathbf{z}_t = \sigma(\mathbf{W}_{xz}\mathbf{E}x_t + \mathbf{W}_{hz}\mathbf{h}_{t-1}) \quad (2)$$

$$\mathbf{r}_t = \sigma(\mathbf{W}_{xr}\mathbf{E}x_t + \mathbf{W}_{hr}\mathbf{h}_{t-1}) \quad (3)$$

$$\hat{\mathbf{h}}_t = \tanh(\mathbf{W}_{xh}\mathbf{E}x_t + \mathbf{U}_{hh}(\mathbf{r}_t \odot \mathbf{h}_{t-1})) \quad (4)$$

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \hat{\mathbf{h}}_t \quad (5)$$

where $\mathbf{E} \in \mathbb{R}^{m \times k}$ is the word embedding and $\mathbf{W}_{xz}, \mathbf{W}_{xr}, \mathbf{W}_{xh}, \mathbf{W}_{hz}, \mathbf{W}_{hr}, \mathbf{U}_{hh}$ are weight matrices. We take the last hidden state \mathbf{h}_t as the representation of the word sequence.

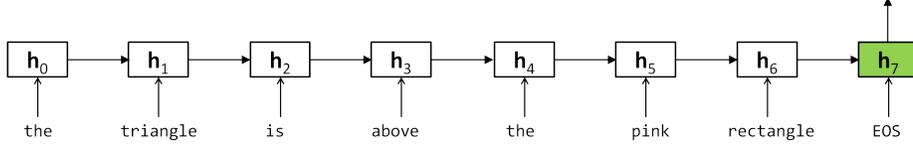


Figure 3: The RNN encoder. The last state is used to summarize the word sequence.

3.2 Reasoning Layers

The modules in the reasoning layers include those for question-fact interaction, pooling.

3.2.1 Question-Fact Interaction

On reasoning layer ℓ , the k^{th} interaction is between $\mathbf{q}^{(\ell-1)}$ and $\mathbf{f}_k^{(\ell-1)}$, resulting in updated representations $\mathbf{q}_k^{(\ell)}$ and $\mathbf{f}_k^{(\ell)}$

$$[\mathbf{q}_k^{(\ell)}, \mathbf{f}_k^{(\ell)}] \stackrel{\text{def}}{=} g_{\text{DNN}_\ell}([\mathbf{q}^{(\ell-1)\top}, \mathbf{f}_k^{(\ell-1)\top}]^\top; \Theta_\ell), \quad (6)$$

with Θ_ℓ being the parameters. In general, $\mathbf{q}_k^{(\ell)}$ and $\mathbf{f}_k^{(\ell)}$ can be of different dimensionality as those of the previous layers. In the simplest case with a single layer in DNN_ℓ , we have

$$\mathbf{q}_k^{(\ell)} \stackrel{\text{def}}{=} \sigma(\mathbf{W}_\ell^\top [(\mathbf{q}^{(\ell-1)})^\top, \mathbf{f}_k^{(\ell-1)\top}] + \mathbf{b}_\ell), \quad (7)$$

where $\sigma(\cdot)$ stands for the nonlinear activation function.

Roughly speaking, $\mathbf{q}_k^{(\ell)}$ contains the update of the system’s understanding on answering the question after its interaction with fact K , while $\mathbf{f}_k^{(\ell)}$ records the change of the K^{th} fact. Therefore, $\{(\mathbf{q}_k^{(\ell)}, \mathbf{f}_k^{(\ell)})\}$ constitute the “state” of the reasoning process.

3.2.2 Pooling

Pooling aims to fuse the understanding of the question right after its interaction with all the facts to form the current status of the question, through which we can enable the comparison between different facts. There are several strategies for this pooling

- **Average/Max Pooling:** To obtain the n^{th} element in $\mathbf{q}^{(\ell)}$, we can take the average or the maximum of the elements at the same location from $\{\mathbf{q}_1^{(\ell)}, \dots, \mathbf{q}_K^{(\ell)}\}$. For example, with max-pooling, we have

$$\mathbf{q}^{(\ell)}(d) = \max(\{\mathbf{q}_1^{(\ell)}(d), \mathbf{q}_2^{(\ell)}(d), \dots, \mathbf{q}_K^{(\ell)}(d)\}), \quad d = 1, 2, \dots, D_\ell$$

where $\mathbf{q}^{(\ell)}(d)$ stands for the d^{th} element of vector $\mathbf{q}^{(\ell)}$. Clearly this kind of pooling is the simplest, without any associated parameters;

- **Gating:** We can have an extra gating network $g^{(\ell)}(\cdot)$ to determine the certainty of the features in $\mathbf{q}_k^{(\ell)}$ based on $\{\mathbf{q}^{(\ell-1)}, \mathbf{f}_k^{(\ell-1)}\}$ (the input for getting $\mathbf{q}_k^{(\ell)}$). The output $g^{(\ell)}(\mathbf{q}^{(\ell-1)}, \mathbf{f}_k^{(\ell-1)})$ has the same dimension as $\mathbf{q}_k^{(\ell)}$, whose n^{th} element, after normalization, can be used as weight for the corresponding element in $\mathbf{q}_k^{(\ell)}$ in obtaining $\mathbf{q}^{(\ell)}$.

- **Model-based:** In the case of temporal-reasoning, there is crucial information in the sequential order of the facts. To account for this temporal structure, we can use a CNN or RNN to combine the information in $\{\mathbf{q}_1^{(\ell)}, \dots, \mathbf{q}_K^{(\ell)}\}$.

At layer- L , the query representation $\mathbf{q}^{(L)}$ after the pooling will serve as the features for the final decision.

3.3 Answering Layer

For simplicity, we focus on the reasoning tasks which can be formulated as classification with predetermined classes. More specifically, we apply NEURAL REASONER to deal with the following two types of questions

- **Type I:** General questions, i.e., questions with Yes-No answer;
- **Type II:** Special questions with a small set of candidate answers.

At reasoning Layer- L , it performs pooling over the intermediate results to select important information for further uses.

$$\mathbf{q} = \text{pool}(\{\mathbf{q}_1^{(L)}, \mathbf{q}_2^{(L)}, \dots, \mathbf{q}_K^{(L)}\}) \quad (8)$$

$$\mathbf{y} = \text{softmax}(\mathbf{W}_{\text{softmax}}^\top \mathbf{q}^{(L)}) \quad (9)$$

After reaching the last reasoning step, in this paper we take two steps, Q^2 is sent to a standard softmax layer to generate an answer which is formulated as a classification problem.

There is another type of prediction as classification where the *effective* classes dynamically change with instances, e.g., the **Single-Supporting-Fact** task in [9]. Those tasks cannot be directly solved with NEURAL REASONER. One simple way to circumvent this is to define the following score function

$$\text{score}_z = g_{\text{match}}(\mathbf{q}^{(L)}, \mathbf{w}_z; \theta)$$

where g_{match} is a function (e.g., a DNN) parameterized with θ , and \mathbf{w}_z is the embedding for class z , with z being dynamically determined for the task.

3.4 Training

The training of model tunes the parameters in $\{\text{RNN}_0, \text{DNN}_1, \dots, \text{DNN}_L\}$ and those in the softmax classifier. Similar to [6], we perform end-to-end training, taking the final answer as the only supervision. More specifically, We use the cross entropy for the cost of classification

$$E_{\text{reasoning}} = \sum_{n \in \mathcal{T}} D_{\text{CE}}(p(\mathbf{y}|r_n) || \mathbf{y}_n)$$

where n indexes the instances in the training set \mathcal{T} , and $r_n = \{Q_n, F_{n,1}, \dots, F_{n,K_n}\}$ stands for question and facts for the n^{th} instance.

Our end-to-end training is the same as [6], while the training in [9] and [5] use the step-by-step labels on the supporting facts for each instance (see Table 1 for examples) in addition to the answer. As described in [6], those extra labels brings much stronger supervision just the answer in the end-to-end learning setting, and typically yield significantly better result on relatively complicated tasks.

4 Auxiliary Training for Question/Fact Representation

We use auxiliary training to facilitate the learning of representations of question and facts. Basically, in addition to using the learned representations of question and facts in the reasoning process, we also use those representations to reconstruct the original questions or their more abstract forms with variables (elaborated later in Section 4.2).

In the auxiliary training, we intend to achieve the following two goals

- to compensate the lack of supervision in the learning task. In our experiments, the supervision can be fairly weak since for each instance it is merely a classification with no more than 12 classes, while the number of instances are 1K to 10K.
- to introduce beneficial bias for the representation learning task. Since the network is a complicated nonlinear function, the back-propagation from the answering layer to the encoding layer can easily fail to learn well.

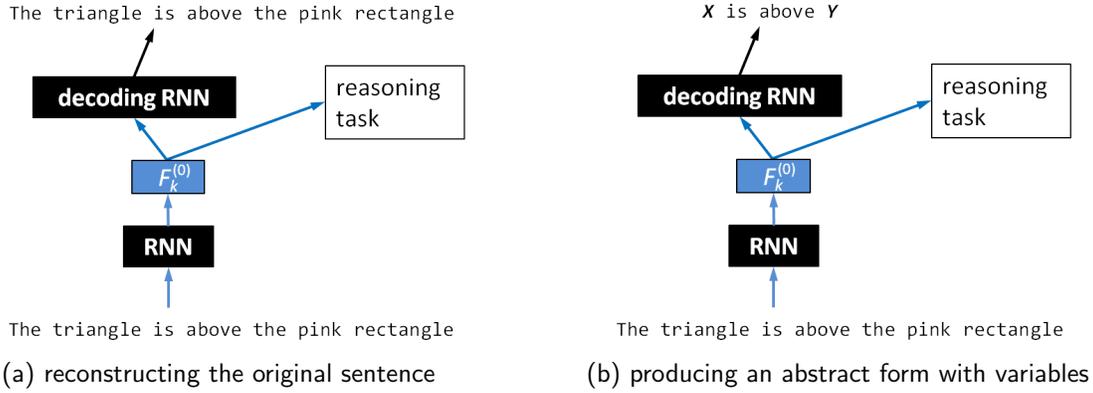


Figure 4: Auxiliary training for question representation. The training for fact representation is identical and therefore omitted.

4.1 Multi-task Learning Setting

As illustrated in Figure 4, we take the simplest way to fuse the auxiliary tasks (recovering) with the main task (reasoning) through linearly combining their costs with trade-off parameter α

$$E = \alpha E_{\text{recovering}} + (1 - \alpha) E_{\text{reasoning}} \quad (10)$$

where $E_{\text{reasoning}}$ is the cross entropy loss describing the discrepancy of model prediction from correct answer (see Section 3.4), and $E_{\text{recovering}}$ is the negative log-likelihood of the sequences (question or facts) to be recovered. More specifically,

$$E_{\text{recovering}} = \sum_{n \in \mathcal{T}} \left\{ \sum_{k=1}^{K_n} \log p(F_{n,k} | \mathbf{f}_{n,k}^{(0)}) + \log p(Q_n | \mathbf{q}_n^{(0)}) \right\}$$

where the likelihood is estimated as in the encoder-decoder framework proposed in [2]. On top of the encoding layer (RNN), we add another decoding layer (RNN) which is trained to sequentially predict words in the original sentence.

4.2 Abstract Forms with Variables

Instead of recovering the original sentence in question and facts, we also study the effect of producing a more abstract form in the auxiliary training task. **More specifically, we let the decoding RNN to recover a sentence with entities replaced with variables** (treated as particular symbols), e.g.,

$$\begin{array}{l} \text{The triangle is above the pink rectangle.} \xrightarrow{\text{recover}} x \text{ is above } y. \\ \text{The blue square is to the left of the triangle.} \xrightarrow{\text{recover}} z \text{ is to the left of } x. \\ \text{Is the pink rectangle to the right of the square?} \xrightarrow{\text{recover}} \text{Is } y \text{ to the right of the } z \text{ ?} \end{array}$$

Through this, we intend to teach the system a more abstract way of representing sentences (both question and facts) and their interactions. More specifically,

- all the entities are only meaningful only when they are compared with each other. In other words, the model (in the encoding and reasoning layers) should not consider specific entities, but their general notions.
- it helps the model to focus on the relations between the entities, the commonality of different facts, and the patterns shared between different instances.

5 Experiments

We report our empirical study on applying NEURAL REASONER to the Question Answer task defined in [8], and compare it against state-of-the-art neural models [9, 5].

5.1 Setup

bAbI is a synthetic question and answering dataset. It contains 20 tasks, and each of them is composed of a set of facts, a question and followed by an answer which is mostly a single word. For most of the time, only a subset of facts are relevant to the given question. Two versions of the data are available, one has 1K training instances per task and the other has 10K instances per task, while the testing set are the same for the two versions.

We select the two most challenging tasks (among the 20 tasks in [8]) Positional Reasoning and Path Finding, to test the reasoning ability of NEURAL REASONER. Positional Reasoning task tests model’s spatial reasoning ability, while Path Finding task, first proposed in [1] tests the ability to reason the correct path between objects based on natural language instructions. In Table 1, we give an instance of each task.

5.2 Implementation Details

In our experiments, we actually used a simplified version of NEURAL REASONER . In the version

- we choose to keep the representation un-updated on each layer, e.g.,

$$F_k \xrightarrow{\text{encode}} \mathbf{f}_k^{(0)} = \mathbf{f}_k^{(1)} = \dots = \mathbf{f}_k^{(L-1)}, \quad k = 1, 2, \dots, K.$$

This choice pushes the update $\mathbf{q}_k^{(\ell)}$ (and its summarization $\mathbf{q}^{(\ell)}$) to record all the information in the interaction between facts and question.

Task I: path finding

1.The office is east of the hallway.
2.The kitchen is north of the office.
3.The garden is west of the bedroom.
4.The office is west of the garden.
5.The bathroom is north of the garden.
How do you go from the kitchen to the garden? south, east, relies on 2 and 4
How do you go from the office to the bathroom? east, north, relies on 4 and 5

Task II: positional reasoning

1.The triangle is above the pink rectangle.
2.The blue square is to the left of the triangle.
Is the pink rectangle to the right of the blue square? Yes, relies on 1 and 2
Is the blue square below the pink rectangle? No, relies on 1 and 2

Table 1: Samples of the two tasks: path finding (upper panel) and positional reasoning (lower panel), with facts, questions and given answers (following each question). For each panel, we first list facts and then question that one needs to answer based on the given facts. On Task I, the answer to the first question is *south, east*, standing for going south first and then east, obtained based on fact 2 and 4.

- we use only two layers, i.e., $L = 2$, for the relatively simple task in the experiments.

Our model was trained with the standard back-propagation (BP) aiming to maximize the likelihood of correct answers. All the parameters including the word-embeddings were initialized by randomly sampling from a uniform distribution $[-0.1, 0.1]$. No momentum and weight decay was used. We trained all the tasks for 200 epochs with stochastic gradient descent and the gradients which had ℓ_2 norm larger than 40 were clipped, learning rate being controlled by AdaDelta [10]. For multi-task learning, different mixture ratios were tried, from 0.1 to 0.9.

5.3 Neural Reasoner vs. Competitor Models

We compare NEURAL REASONER with the following three neural reasoning models: 1)Memory Network, including the one with step-by-step supervision [9](denoted as MEMORY NET-STEP) and the end-to-end version [6] (denoted as MEMORY NET-N2N), and 2) DYNAMIC MEMORY NETWORK, proposed in [5], also with step-by-step supervision. In Table 2, we report the performance of a particular case of NEURAL REASONER with 1) two reasoning layers, 2) 2-layer DNNs as the interaction modules in each reasoning layer, and 3) auxiliary task of recovering the original question and facts. The results are compared against three neural competitors. We have the following observations.

- The proposed NEURAL REASONER performs significantly better than Memory Net-N2N, especially with more training data.
- Although not a fair comparison to our model, NEURAL REASONER is actually better than MEMORY NET-N2N and DYNAMIC MEMORY NET on Positional Reasoning (1K)

& (10K) as well as Path Finding (10K), with about 20% margin on both tasks with 10K training instances.

	Posi. Reason. (1K)	Posi. Reason. (10K)
Step-by-step Supervision		
MEMORY NET-STEP	65.0%	75.4%
DYNAMIC MEMORY NET	59.6%	-
End-to-End		
MEMORY NET-N2N	59.6%	60.3%
NEURAL REASONER	66.4%	97.9%

	Path Finding (1K)	Path Finding (10K)
Step-by-step Supervision		
MEMORY NET-STEP	36.0%	68.1%
DYNAMIC MEMORY NET	34.5%	-
End-to-End		
MEMORY NET-N2N	17.2%	33.4%
NEURAL REASONER	17.3%	87.0%

Table 2: Results on two reasoning tasks. The results of MEMORY NET-STEP, MEMORY NET-N2N, and DYNAMIC MEMORY NET are taken respectively from [9],[6] and [5].

Please note that the results of NEURAL REASONER reported in Table 2 are not based on architectures specifically tuned for the tasks. As a matter of fact, with more complicated models (more reasoning layers and deeper interaction modules), we can achieve even better results on large datasets (e.g., over 98% accuracy on Path Finding with 10K instances). We will however leave the discussion on different architectural variants to the next section.

5.4 Architectural Variations

This section is devoted to the study of architectural variants of NEURAL REASONER. More specifically, we consider the variations in 1)the number of reasoning layers, 2) the depth of the interaction DNN, and 3) the auxiliary tasks, with results summarized by Table 3. We have the following observations:

- Auxiliary tasks are essential to the efficacy of NEURAL REASONER, without which the performances of NEURAL REASONER drop dramatically. The reason, as we conjecture in Section 4, is that the reasoning task alone cannot give enough supervision for learning accurate word vectors and parameters of the RNN encoder. We note that NEURAL REASONER can still outperform Memory Net (N2N) with 10K data on both tasks.
- NEURAL REASONER with shallow architectures, more specifically two reasoning layers and 1-layer DNN, apparently can benefit from the auxiliary learning of recovering abstract forms on small datasets (1K on both tasks). However, with deeper architectures or more training data, the improvement over that of recovering original sentences become smaller, despite the extra information it utilizes.
- When larger training datasets are available, NEURAL REASONER appears to prefer relatively deeper architectures. More importantly, although both tasks require two reasoning steps, the performance does not deteriorate with three reasoning layers. On both

	Posi. Reason. (1K)	Posi. Reason. (10K)
No auxiliary task		
2-layer reasoning, 1-layer DNN	60.2%	72.1%
2-layer reasoning, 2-layer DNN	59.6%	69.3%
3-layer reasoning, 3-layer DNN	58.7%	59.7%
Auxiliary task: Original		
2-layer reasoning, 1-layer DNN	63.1%	93.8%
2-layer reasoning, 2-layer DNN	66.4%	97.9%
3-layer reasoning, 3-layer DNN	69.4%	99.1%
Auxiliary task: Abstract		
2-layer reasoning, 1-layer DNN	70.9%	95.2%
2-layer reasoning, 2-layer DNN	66.6%	95.6%
3-layer reasoning, 3-layer DNN	68.3%	97.4%

	Path Finding (1K)	Path Finding (10K)
No auxiliary task		
2-layer reasoning, 1-layer DNN	13.6%	52.2%
2-layer reasoning, 2-layer DNN	12.3%	54.2%
3-layer reasoning, 3-layer DNN	13.1%	51.7%
Auxiliary task: Original		
2-layer reasoning, 1-layer DNN	14.1%	57.0%
2-layer reasoning, 2-layer DNN	17.3%	87.0%
3-layer reasoning, 3-layer DNN	14.0%	98.4%
Auxiliary task: Abstract		
2-layer reasoning, 1-layer DNN	18.1%	55.8%
2-layer reasoning, 2-layer DNN	15.4%	87.8%
3-layer reasoning, 3-layer DNN	11.3%	98.6%

Table 3: Results on two reasoning tasks yielded by NEURAL REASONER with different architectural variations.

tasks, with 10K training instances, NEURAL REASONER with three reasoning layers and 3-layer DNN can achieve over 98% accuracy.

6 Conclusion and Future Work

We have proposed NEURAL REASONER, a framework for neural network-based reasoning over natural language sentences. NEURAL REASONER is flexible, powerful, and language independent. Our empirical studies show that NEURAL REASONER can dramatically improve upon existing neural reasoning systems on two difficult artificial tasks proposed in [9]. For future work, we will explore 1) tasks with higher difficulty and reasoning depth, e.g., tasks which require a large number of supporting facts and facts with complex intrinsic structures, 2) the common structure in different but similar reasoning tasks (e.g., multiple tasks all with general questions), and 3) automatic selection of the reasoning architecture, for example, determining when to stop the reasoning based on the data.

References

- [1] D. L. Chen and R. J. Mooney. Learning to interpret natural language navigation instructions from observations. In *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence*,

- AAAI 2011, San Francisco, California, USA, August 7-11, 2011*, 2011.
- [2] K. Cho, B. van Merriënboer, C. Gulcehre, F. Bougares, H. Schwenk, and Y. Bengio. Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of EMNLP*, pages 1724–1734, 2014.
 - [3] J. Chung, cC. Gülcehre, K. Cho, and Y. Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *CoRR*, abs/1412.3555, 2014.
 - [4] B. Hu, Z. Lu, H. Li, and Q. Chen. Convolutional neural network architectures for matching natural language sentences. In *Advances in Neural Information Processing Systems 27*, pages 2042–2050, 2014.
 - [5] A. Kumar, O. Irsoy, J. Su, J. Bradbury, R. English, B. Pierce, P. Ondruska, I. Gulrajani, and R. Socher. Ask me anything: Dynamic memory networks for natural language processing. *CoRR*, abs/1506.07285, 2015.
 - [6] S. Sukhbaatar, A. Szlam, J. Weston, and R. Fergus. Weakly supervised memory networks. *CoRR*, abs/1503.08895, 2015.
 - [7] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112, 2014.
 - [8] J. Weston, A. Bordes, S. Chopra, and T. Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. *CoRR*, abs/1502.05698, 2015.
 - [9] J. Weston, S. Chopra, and A. Bordes. Memory networks. *CoRR*, abs/1410.3916, 2014.
 - [10] M. D. Zeiler. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*, 2012.