

Neural Enquirer: Learning to Query Tables in Natural Language

Pengcheng Yin¹ Zhengdong Lu² Hang Li² Ben Kao¹

¹ Department of Computer Science, The University of Hong Kong

² Noah’s Ark Lab, Huawei Technologies

¹{pcyin,kao}@cs.hku.hk ²{Lu.Zhengdong, HangLi.HL}@huawei.com

Abstract

We propose NEURAL ENQUIRER — a neural network architecture for answering natural language (NL) questions based on a knowledge base (KB) table. Unlike existing work on end-to-end training of semantic parsers [Pasupat and Liang, 2015; Neelakantan *et al.*, 2015], NEURAL ENQUIRER is fully “neuralized”: it finds distributed representations of queries and KB tables, and executes queries through a series of neural network components called “executors”. Executors model query operations and compute intermediate execution results in the form of table annotations at different levels. NEURAL ENQUIRER can be trained with gradient descent, with which the representations of queries and the KB table are jointly optimized with the query execution logic. The training can be done in an end-to-end fashion, and it can also be carried out with stronger guidance, e.g., step-by-step supervision for complex queries. NEURAL ENQUIRER is one step towards building neural network systems that can understand natural language in real-world tasks. As a proof-of-concept, we conduct experiments on a synthetic QA task, and demonstrate that the model can learn to execute reasonably complex NL queries on small-scale KB tables.

1 Introduction

Natural language dialogue and question answering often involve querying a knowledge base [Wen *et al.*, 2015; Berant *et al.*, 2013]. The traditional approach involves two steps: First, a given query Q is semantically parsed into an “executable” representation, which is often expressed in certain logical form Z (e.g., SQL-like queries). Second, the representation is executed against a knowledge base from which an answer is obtained. For queries that involve complex semantic constraints and logic (e.g., “Which city hosted the longest Olympic Games before the Games in Beijing?”), semantic parsing and query execution become extremely complex. For example, carefully hand-crafted features and rules are needed to correctly parse a complex query into its logical form (see example shown in the lower-left corner of Figure 1). This

complexity often results in poor accuracy of the system. To partially overcome this difficulty, recent works [Clarke *et al.*, 2010; Liang *et al.*, 2011; Pasupat and Liang, 2015] attempt to “backpropagate” query execution results to revise the semantic representation of a query, which is an example of learning from grounding [Chen and Mooney, 2008; Kim and Mooney, 2012]. This approach, however, is greatly hindered by the fact that traditional semantic parsing mostly involves rule-based features and symbolic manipulation, and is subject to intractable search space incurred by the great flexibility of natural language.

Neural network-based models have enjoyed much successes in natural language processing, particularly in machine translation and syntactic parsing. These successes are attributable to direct and strong supervision. The recent work on learning to execute simple program codes with LSTM [Zaremba and Sutskever, 2014] pioneers in the direction on learning to parse structured objects through executing it in a purely neural way, while the more recent work on Neural Turing Machines (NTMs) [Graves *et al.*, 2014] introduces more modeling flexibility by equipping the LSTM with external memory and various means for interacting with it.

Inspired by the above-mentioned research, we aim to design a neural network system that learns to understand queries and execute them on a knowledge base table from examples of queries and answers. We propose NEURAL ENQUIRER, a fully neuralized, end-to-end differentiable system that jointly models semantic parsing and query execution. NEURAL ENQUIRER encodes queries and KB tables into distributed representations, and executes compositional queries against the KB through a series of differentiable executors. The model is trained using query-answer pairs, where the distributed representations of queries and the KB are optimized together with the query execution logic in an end-to-end fashion. As the first step along this line of research, we evaluate NEURAL ENQUIRER using a synthetic question-answering task as a proof-of-concept, and demonstrate that our proposed model is capable of learning to execute complex compositional natural language questions on small-scale KB tables.

2 Model

Following [Pasupat and Liang, 2015], we study the problem of question answering on a single KB table. Specifically, given an NL query Q and a KB table T , NEURAL ENQUIRER

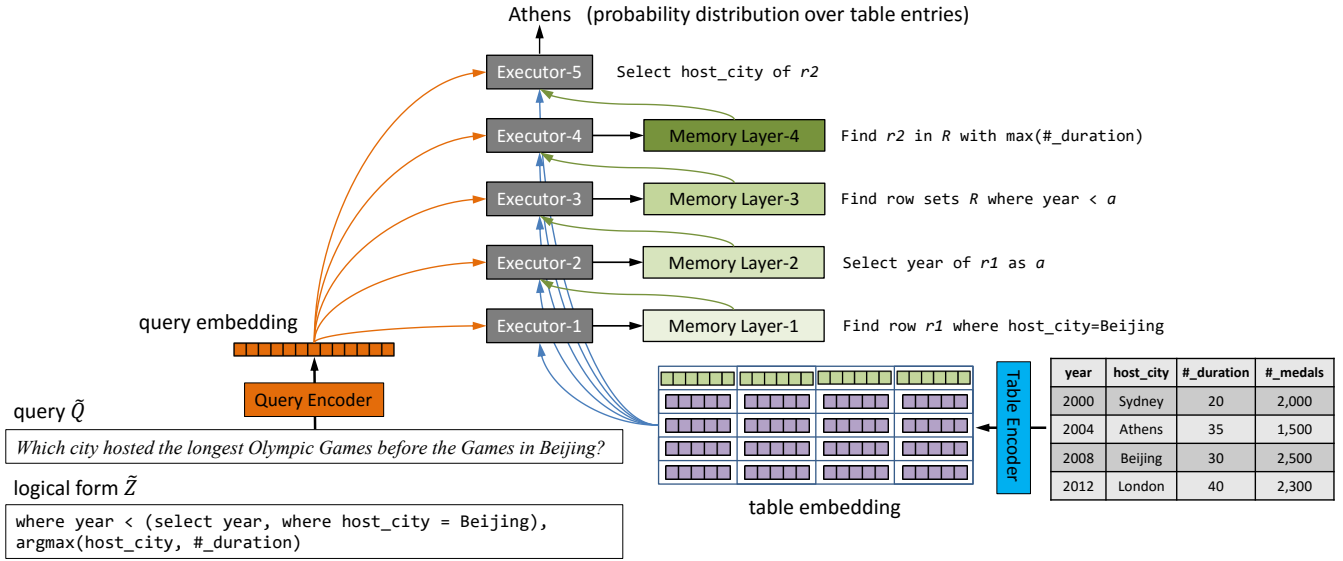


Figure 1: An overview of NEURAL ENQUIRER with five executors

executes Q against \mathcal{T} and outputs a ranked list of answers. The execution is done by first using *Encoders* to encode the query and table into distributed representations, which are then sent to a cascaded pipeline of *Executors* to derive the answer. Figure 1 gives an illustrative example (with five executors). It consists of the following components:

Query Encoder (Section 2.1), which abstracts the semantics of an NL query and encodes it into a query embedding.

Table Encoder (Section 2.2), which derives a table embedding by encoding entries in the table into distributed vectors.

Executor (Section 2.3), which executes the query against the table and outputs *annotations* that encode intermediate execution results. Annotations are stored in the memory of each layer to be accessed by the executor of the next layer. Since complex compositional queries can be answered in multiple steps of computation, each executor models a specific type of operation conditioned on the query. Figure 1 illustrates the operation each executor is assumed to perform in answering the example query \tilde{Q} . Different from classical semantic parsing approaches which require a predefined set of all possible logical operations, NEURAL ENQUIRER *learns* the logic of executors via end-to-end training using query-answer pairs. By stacking several executors, our model is able to answer complex queries that involve multiple steps of computation.

2.1 Query Encoder

Query Encoder converts an NL query Q into a query embedding $\mathbf{q} \in \mathbb{R}^{d_q}$. Let $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ be the embeddings of words in Q , where $\mathbf{x}_t \in \mathbb{R}^{d_w}$ is from an embedding matrix \mathbf{L} . We employ a bidirectional Gated Recurrent Unit (GRU) [Bahdanau *et al.*, 2015] to summarize the sequence $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$ in forward and reverse orders. \mathbf{q} is formed by concatenating the last hidden states in the two directions.

It is worth noting that Query Encoder can find the representation of a rather general class of symbol sequences, agnostic to the actual representation of the query (e.g., natural

language, SQL, etc). The model is able to learn the semantics of input queries through end-to-end training, making it a generic model for query understanding and query execution.

2.2 Table Encoder

Table Encoder converts a KB table \mathcal{T} into a distributed representation, which is used as an input to executors. Suppose \mathcal{T} has M rows and N columns. In our model, the n -th column is associated with a *field name* (e.g., *host_city*). Each cell value is a word (e.g., *Beijing*) in the vocabulary. We use w_{mn} to denote the cell value in row m column n , and \mathbf{w}_{mn} to denote its embedding. Let \mathbf{f}_n be the embedding of the field name for column n . For each entry (cell) w_{mn} , Table Encoder computes a (field, value) composite embedding $\mathbf{e}_{mn} \in \mathbb{R}^{d_e}$ by fusing \mathbf{f}_n and \mathbf{w}_{mn} using a single-layer Neural Network:

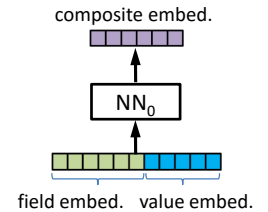
$$\mathbf{e}_{mn} = \text{NN}_0(\mathbf{f}_n, \mathbf{w}_{mn}) = \tanh(\mathbf{W} \cdot [\mathbf{f}_n; \mathbf{w}_{mn}] + \mathbf{b}),$$

where $[\cdot; \cdot]$ denotes vector concatenation. The output of Table Encoder is an $M \times N \times d_e$ tensor that consists of $M \times N$ embeddings, each of length d_e .

We remark that our Table Encoder is different from classical knowledge embedding models (e.g., TransE [Bordes *et al.*, 2013]). While traditional methods learn the embeddings of entities (cell values) and relations (field names) in an unsupervised fashion via minimizing certain reconstruction errors, embeddings in Table Encoder are optimized via supervised learning in end-to-end QA tasks.

2.3 Executor

NEURAL ENQUIRER executes an input query on a KB table through layers of execution. Each layer consists of an executor that, after learning, performs certain operation (e.g.,



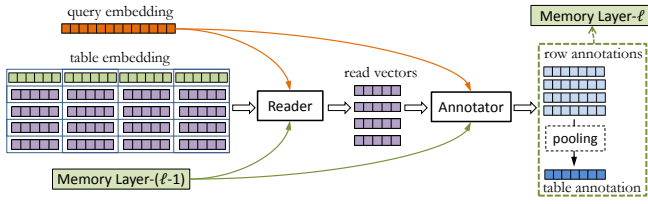


Figure 2: Overview of an Executor- ℓ

`select, max`) relevant to the input query. An executor outputs intermediate execution results, referred to as *annotations*, which are saved in the external memory of the executor. A query is executed sequentially through a stack of executors. Such a cascaded architecture enables the model to answer complex, compositional queries. An example is given in Figure 1 in which descriptions of the operation each executor is assumed to perform for the query \tilde{Q} are shown. We will demonstrate in Section 4 that the model is capable of learning the operation logic of each executor via end-to-end training.

As illustrated in Figure 2, an executor at Layer- ℓ (denoted as Executor- ℓ) consists of two major neural network components: a *Reader* and an *Annotator*. The executor processes a table row-by-row. The Reader reads in data from each row m in the form of a *read vector* \mathbf{r}_m^ℓ , which is then sent to the Annotator to perform the actual execution. The output of the Annotator is a *row annotation* \mathbf{a}_m^ℓ , which captures the row-wise local computation result. Once all row annotations are obtained, Executor- ℓ generates a *table annotation* \mathbf{g}^ℓ to summarize the global computation result on the whole table by pooling all row annotations. All the row and table annotations are saved in the memory of Layer- ℓ : $\mathcal{M}^\ell = \{\mathbf{a}_1^\ell, \mathbf{a}_2^\ell, \dots, \mathbf{a}_M^\ell, \mathbf{g}^\ell\}$. Intuitively, row annotations handle operations that require only row-wise, local information (e.g., `select`, `where`), while table annotations model superlative operations (e.g., `max`, `min`) by aggregating table-wise, global execution results. A combination of row and table annotations enables the model to perform a wide variety of query operations in real world scenarios.

Reader

As illustrated in Figure 3, for the m -th row with N (field, value) composite embeddings $\mathcal{R}_m = \{\mathbf{e}_{m1}, \mathbf{e}_{m2}, \dots, \mathbf{e}_{mN}\}$, the Reader fetches a read vector \mathbf{r}_m^ℓ from \mathcal{R}_m via an attentive reading operation:

$$\mathbf{r}_m^\ell = f_R^\ell(\mathcal{R}_m, \mathcal{F}_T, \mathbf{q}, \mathcal{M}^{\ell-1}) = \sum_{n=1}^N \tilde{\omega}(\mathbf{f}_n, \mathbf{q}, \mathbf{g}^{\ell-1}) \mathbf{e}_{mn}$$

where $\mathcal{M}^{\ell-1}$ denotes the content of memory Layer- $(\ell-1)$, and $\mathcal{F}_T = \{\mathbf{f}_1, \mathbf{f}_2, \dots, \mathbf{f}_N\}$ is the set of field name embeddings. $\tilde{\omega}(\cdot)$ is the normalized attention weights given by:

$$\tilde{\omega}(\mathbf{f}_n, \mathbf{q}, \mathbf{g}^{\ell-1}) = \frac{\exp(\omega(\mathbf{f}_n, \mathbf{q}, \mathbf{g}^{\ell-1}))}{\sum_{n'=1}^N \exp(\omega(\mathbf{f}_{n'}, \mathbf{q}, \mathbf{g}^{\ell-1}))} \quad (1)$$

where $\omega(\cdot)$ is modeled as a Deep Neural Network (denoted as $\text{DNN}_1^{(\ell)}$). Since each executor models a specific type of computation, it should only attend to a subset of entries that are pertinent to its execution. This is modeled by the Reader. Our

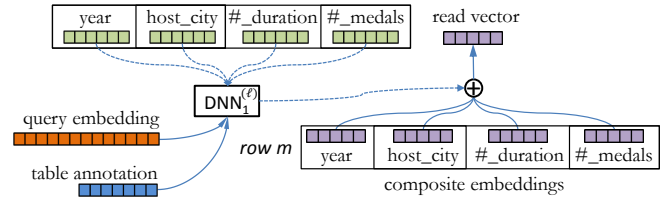


Figure 3: Illustration of the Reader in Executor- ℓ .

approach is related to the content-based addressing of Neural Turing Machines [Graves *et al.*, 2014] and the attention mechanism in neural machine translation models [Bahdanau *et al.*, 2015].

Annotator

The Annotator of Executor- ℓ computes row and table annotations based on read vectors fetched by the Reader. The results are stored in the ℓ -th memory layer \mathcal{M}^ℓ accessible to Executor- $(\ell+1)$. The last executor is the only exception, which outputs the final answer.

[Row annotations] Capturing row-wise execution result, the annotation \mathbf{a}_m^ℓ for row m in Executor- ℓ is given by

$$\mathbf{a}_m^\ell = f_A^\ell(\mathbf{r}_m^\ell, \mathbf{q}, \mathcal{M}^{\ell-1}) = \text{DNN}_2^{(\ell)}([\mathbf{r}_m^\ell; \mathbf{q}; \mathbf{a}_m^{\ell-1}; \mathbf{g}^{\ell-1}]). \quad (2)$$

$\text{DNN}_2^{(\ell)}$ fuses the corresponding read vector \mathbf{r}_m^ℓ , the results saved in the previous memory layer (row and table annotations $\mathbf{a}_m^{\ell-1}, \mathbf{g}^{\ell-1}$), and the query embedding \mathbf{q} . Intuitively, row annotation $\mathbf{a}_m^{\ell-1}$ and table annotation $\mathbf{g}^{\ell-1}$ summarize the local and global status of execution up to Layer- $(\ell-1)$, respectively. $\text{DNN}_2^{(\ell)}$ then performs the actual query execution by combining these annotations with the read vector and query embedding, and outputs a row annotation \mathbf{a}_m^ℓ that encodes the local execution result on row m .

[Table annotations] Capturing global execution state, a table annotation summarizes all row annotations via a global max pooling operation:

$$\mathbf{g}^\ell = f_{\text{MAXPOOL}}(\mathbf{a}_1^\ell, \mathbf{a}_2^\ell, \dots, \mathbf{a}_M^\ell) = [g_1, g_2, \dots, g_{d_g}]^\top \quad (3)$$

where $g_k = \max(\{\mathbf{a}_1^\ell(k), \mathbf{a}_2^\ell(k), \dots, \mathbf{a}_M^\ell(k)\})$ is the maximum value among the k -th elements of all row annotations.

Last Layer Executor

Instead of computing annotations based on read vectors, the last executor in NEURAL ENQUIRER directly outputs the probability of an entry w_{mn} in table \mathcal{T} being the answer a :

$$p(a=w_{mn}|Q, \mathcal{T}) = \frac{\exp(f_{\text{ANS}}^\ell(\mathbf{e}_{mn}, \mathbf{q}, \mathbf{a}_m^{\ell-1}, \mathbf{g}^{\ell-1}))}{\sum_{m'=1, n'=1}^{M, N} \exp(f_{\text{ANS}}^\ell(\mathbf{e}_{m'n'}, \mathbf{q}, \mathbf{a}_{m'}^{\ell-1}, \mathbf{g}^{\ell-1}))} \quad (4)$$

where $f_{\text{ANS}}^\ell(\cdot)$ is modeled as a DNN ($\text{DNN}_3^{(\ell)}$). Note that the last executor, which is devoted to returning answers, could still carry out execution in $\text{DNN}_3^{(\ell)}$.

3 Learning

NEURAL ENQUIRER can be trained in an *end-to-end* (N2N) fashion on QA tasks. During training, both the representations of queries and tables, as well as the execution logic captured by the weights of executors are learned. Given a set of

year	host_city	#_participants	#_medals	#_duration	#_audience	host_country	GDP	country_size	population
2008	Beijing	4,200	2,500	30	67,000	China	2,300	960	130

Figure 4: An example table in the synthetic QA task (only one row shown)

① SELECT-WHERE Queries [select F_a , where $F_b = w_b$]	③ WHERE-SUPERLATIVE Queries [where $F_a > < w_a, \text{argmax/min}(F_b, F_c)$]
▷ How many people participated in the Games in Beijing?	▷ How long was the Games with the most medals that had fewer than 3,000 participants?
▷ In which city was the Games hosted in 2012?	▷ How many medals were in the first Games after 2008?
② SUPERLATIVE Queries [argmax/min(F_a, F_b)]	④ NEST Queries [where $F_a > < (\text{select } F_a, \text{where } F_b = w_b), \text{argmax/min}(F_c, F_d)$]
▷ When was the latest Games hosted?	▷ Which country hosted the longest Games before the Games in Athens?
▷ How big is the country which hosted the shortest Games?	▷ How many people watched the earliest Games that lasts for more days than the Games in 1956?

Table 1: Example queries for each query type, with annotated SQL-like logical form templates

N_D query-table-answer triples $\mathcal{D} = \{(Q^{(i)}, \mathcal{T}^{(i)}, y^{(i)})\}$, we learn the model parameters by maximizing the log-likelihood of gold-standard answers:

$$\mathcal{L}_{N2N}(\mathcal{D}) = \sum_{i=1}^{N_D} \log p(a = y^{(i)} | Q^{(i)}, \mathcal{T}^{(i)}) \quad (5)$$

In end-to-end training, each executor discovers its operation logic from training data in a purely data-driven fashion, which could be difficult for complex queries requiring four or five sequential operations.

This can be alleviated by softly guiding the learning process via controlling the attention weights $\tilde{w}(\cdot)$ in Eq. (1). By enforcing $\tilde{w}(\cdot)$ to bias towards a field pertaining to a specific operation, we can “coerce” the executor to figure out the logic of this operation relative to the field. For example, for Executor-1 in Figure 1, by biasing the attention weight of the `host_city` field towards 1.0, only the value of `host_city` will be fetched and sent to the Annotator. In this way we can “force” the executor to learn the where operation to find the row whose `host_city` is *Beijing*. This method will be referred to as *step-by-step* (SbS) training. Formally, this is done by introducing additional supervision signal to Eq. (5):

$$\mathcal{L}_{SbS}(\mathcal{D}) = \sum_{i=1}^{N_D} [\log p(a = y^{(i)} | Q^{(i)}, \mathcal{T}^{(i)}) + \alpha \sum_{\ell=1}^{L-1} \log \tilde{w}(\mathbf{f}_{i,\ell}^*, \cdot, \cdot)] \quad (6)$$

where α is a tuning weight, and L is the number of executors. $\mathbf{f}_{i,\ell}^*$ is the embedding of the field known *a priori* to be used by Executor- ℓ in answering the i -th example.

4 Experiments

In this section we evaluate NEURAL ENQUIRER on synthetic QA tasks with NL queries of varying compositional depths.

4.1 Synthetic QA Task

We present a synthetic QA task with a large number of QA examples at various levels of complexity to evaluate the performance of NEURAL ENQUIRER. Starting with “artificial” tasks accelerates the development of novel deep models [Weston *et al.*, 2015], and has gained increasing popularity in recent research on modeling symbolic computation using DNNs [Graves *et al.*, 2014; Zaremba and Sutskever, 2014].

Our synthetic dataset consists of query-table-answer triples $\{(Q^{(i)}, \mathcal{T}^{(i)}, y^{(i)})\}$. To generate a triple, we first randomly sample a table $\mathcal{T}^{(i)}$ of size 10×10 from a synthetic schema

of Olympic Games. The cell values of $\mathcal{T}^{(i)}$ are drawn from a vocabulary of 120 location names and 120 numbers. Figure 4 gives an example table. Next, we sample a query $Q^{(i)}$ generated using NL templates, and obtain its gold-standard answer $y^{(i)}$ on $\mathcal{T}^{(i)}$. Our task consists of four types of NL queries, with examples given in Table 1. We also give the logical form template for each type of query. The templates define the semantics and compositionality of queries. We generate queries at various compositional depths, ranging from simple SELECT-WHERE queries to more complex NEST ones. This makes the dataset have similar complexity as a real-world one, except for the relatively small vocabulary. The queries are flexible enough to involve complex matching between NL phrases and logical constituents, which makes query understanding nontrivial: (1) the same field is described by different NL phrases (e.g., “*How big is the country ...*” and “*What is the size of the country ...*” for the `country_size` field); (2) different fields may be referred to by the same NL pattern (e.g., “*in China*” for `host_country` and “*in Beijing*” for `host_city`); (3) simple NL constituents may be grounded to complex logical operations (e.g., “*after the Games in Beijing*” implies comparing between the values of `year` fields).

To simulate the read-world scenario where queries of various types are issued to the model, we construct two MIXED datasets, with 25K and 100K training examples respectively, where four types of queries are sampled with the ratio 1 : 1 : 1 : 2. Both datasets share the same testing set of 20K examples, 5K for each type of query. We enforce that no tables and queries are shared between training/testing sets.

4.2 Setup

[Tuning] We adopt a model with five executors. The lengths of hidden states for GRU and DNNs are 150, 50. The numbers of layers for $\text{DNN}_1^{(\ell)}$, $\text{DNN}_2^{(\ell)}$ and $\text{DNN}_3^{(\ell)}$ are 2, 3, 3. The length of word embeddings and annotations is 20. α is 0.2. We train the model using ADADelta [Zeiler, 2012] on a Tesla K40 GPU. The training converges fast within 2 hours.

[Metric] We evaluate in terms of accuracy, defined as the fraction of correctly answered queries.

[Models] We compare the results of the following settings:

- **Sempre** [Pasupat and Liang, 2015] is a state-of-the-art semantic parser and serves as the baseline;
- **N2N**, our model trained using *end-to-end* setting (Sec 4.3);
- **SbS**, our model trained using *step-by-step* setting (Sec 4.4);

	MIXED-25K				MIXED-100K		
	SEMPRE	N2N	SbS	N2N-OOV	N2N	SbS	N2N-OOV
SELECT.WHERE	93.8%	96.2%	99.7%	90.3%	99.3%	100.0%	97.6%
SUPERLATIVE	97.8%	98.9%	99.5%	98.2%	99.9%	100.0%	99.7%
WHERE.SUPERLATIVE	34.8%	80.4%	94.3%	79.1%	98.5%	99.8%	98.0%
NEST	34.4%	60.5%	92.1%	57.7%	64.7%	99.7%	63.9%
Overall Accuracy	65.2%	84.0%	96.4%	81.3%	90.6%	99.9%	89.8%

Table 2: Accuracies on MIXED datasets

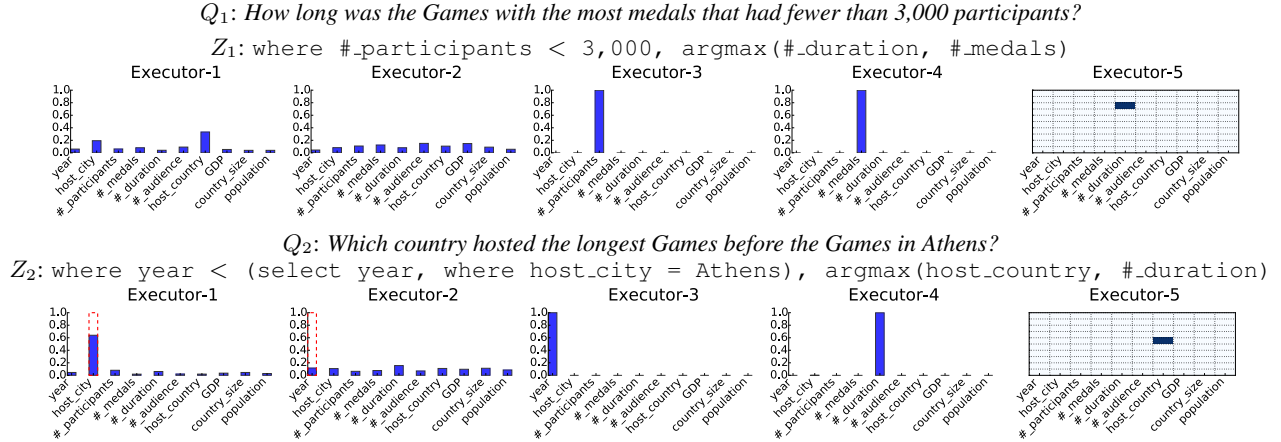


Figure 5: Weights visualization of queries Q_1 and Q_2

- **N2N-OOV**, a variant of the N2N model to deal with out-of-vocabulary words (Sec 4.5)

4.3 End-to-End Evaluation

Table 2 summarizes the results of SEMPRES and NEURAL ENQUIRER under different settings. We show both the individual performance for each query type and the overall accuracy. We evaluate SEMPRES only on MIXED-25K because of its long training time even on this small dataset (about 3 days).

In this section we discuss the results under end-to-end (N2N) training setting. On MIXED-25K, the relatively low performance of SEMPRES indicates that our QA task, although synthetic, is highly nontrivial. Surprisingly, NEURAL ENQUIRER outperforms SEMPRES on all query types, with a marginal gain on *simple* queries (SELECT.WHERE, SUPERLATIVE), and significant improvement on *complex* queries (WHERE.SUPERLATIVE, NEST). We posit that the low performance of SEMPRES on complex queries is likely due to the intractable search space incurred by the flexibility of its float parsing algorithm. On MIXED-100K, our model registers an overall accuracy of 90.6%. These results show that in our QA task, NEURAL ENQUIRER is very effective in answering compositional NL queries, especially those with complex semantics compared with the state-of-the-art system.

To further understand why our model is capable of answering compositional queries, we study the attention weights of Readers (Eq. 1) for intermediate executors, and the answer probability (Eq. 4) the last executor outputs for each table entry. These statistics are obtained on MIXED-100K. We sample two queries (Q_1 and Q_2) in the testing set that our model

answers correctly and visualize their corresponding values in Figure 5. To better understand the query execution process, we also give the logical forms (Z_1 and Z_2) of the two queries. Note that the logical forms are just for reference purpose and unknown by the model. We find that each executor actually *learns* its execution logic in N2N training, which is in accordance with our assumption. The model executes Q_1 in three steps, with each of the last three executors performs a specific type of operation. For each row, Executor-3 takes the value of the #_participants field as input, while Executor-4 attends to the #_medals field. Finally, Executor-5 outputs a high probability for the #_duration field in the 3-rd row. The attention weights for Executor-1 and Executor-2 appear to be meaningless because Q_1 requires only three steps of execution, and the model learns to defer the meaningful execution to the last three executors. Comparing with the logical form Z_1 of Q_1 , we can deduce that Executor-3 “executes” the where clause in Z_1 to find row sets R satisfying the condition, and Executor-4 performs the first part of argmax to find the row $r \in R$ with the maximum value of #_medals, while Executor-5 outputs the value of #_duration in r .

Compared with the relatively simple Q_1 , Q_2 is more complicated. According to Z_2 , Q_2 involves an additional nest sub-query to be solved by two extra executors, and requires a total of five steps of execution. The last three executors function similarly as in answering Q_1 , yet the execution logic for the first two executors (devoted to solving the sub-query) is a bit obscure, since their attention weights are scattered instead of being perfectly centered on the ideal fields as highlighted in red dashed rectangles. We posit that this is because dur-

ing the end-to-end training, the supervision signal propagated from the top layer has decayed along the long path down to the first two executors, which causes vanishing gradients.

4.4 With Additional Step-by-Step Supervision

To alleviate the vanishing gradient problem when training on complex queries like Q_2 , we train the model using step-by-step (SbS) setting (Eq. 6), where we encourage each intermediate executor to attend to the field that is known *a priori* to be relevant to its execution logic. Results are shown in Table 2 (column SbS). With stronger supervision signal, the model significantly outperforms the N2N setting, and achieves perfect accuracy on MIXED-100K. This shows that NEURAL ENQUIRER is capable of leveraging the additional supervision signal given to intermediate layers in SbS training. Let us revisit the query Q_2 in SbS setting. In contrast to the result in N2N setting (Figure 5) where the attention weights for the first two executors are obscure, now the weights are perfectly skewed towards each relevant field with a value of 1.0, which corresponds with the highlighted ideal weights.

4.5 Dealing with Out-Of-Vocabulary Words

One of the major challenges for applying neural network models to NLP applications is to deal with out-of-vocabulary (OOV) words (e.g., new entities for QA). Surprisingly, we find that a simple variant of NEURAL ENQUIRER is able to handle unseen entities almost without loss of accuracy.

Specifically, we divide words in the vocabulary into *entity* words and *operation* words. Embeddings of entity words (e.g., *Beijing*) function like indices to facilitate the matching between the entities in queries and tables during query execution, and therefore are not updated once randomly initialized; while those of operation words, i.e., all non-entity words (e.g., numbers, *longest*, *before*, etc), carry semantic meanings relevant to execution and will be optimized in training. Therefore, after randomly initializing the embedding matrix L , we only update the embeddings of operation words in training, while keeping those of entity words unchanged. To evaluate the model we modify the queries in the testing set to replace all entity words (i.e., all country and city names) with those unseen in training. Results obtained using N2N training, reported in Table 2 (column N2N-OOV), show that the model yields performance comparable with non-OOV settings.

An interesting question is how the model resolves the types of OOV entity words (i.e., cities vs. countries) in ambiguous queries, e.g., Q_3 : “How many people watched the Games in Macau?”, since the random embeddings of entity words (e.g., *Macau*) cannot link them to their corresponding fields. The model executes Q_3 using the last three executors, with the last executor attending to the `#_audience` field as expected. Interestingly, however, the model attends to the `host_city` field in Executor-3, and then `host_country` in Executor-4 (see Figure 6), indicating the model learns to scan all possible fields to figure out the correct field of an OOV entity.

4.6 Querying Expanded Knowledge Source

We simulate a test case to evaluate the model’s ability to generalize to an expanded knowledge source. We train a model on tables whose field sets are either $\mathcal{F}_1, \mathcal{F}_2, \dots, \mathcal{F}_5$, where \mathcal{F}_i

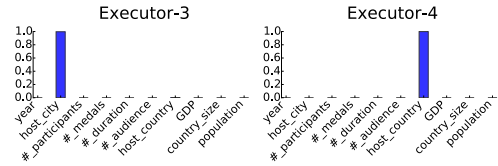


Figure 6: Weights visualization of query Q_3

Query Type	SELECT_WHERE	SUPERLATIVE	WHERE_SUPERLATIVE	Overall
Accuracy	68.2%	84.8%	80.2%	77.7%

Table 3: Accuracies for querying expanded knowledge source

is a subset of the entire field set $\mathcal{F}_{\mathcal{T}}$ and $|\mathcal{F}_i| = 5$. We then test the model on tables with all fields $\mathcal{F}_{\mathcal{T}}$ and queries whose fields span multiple subsets \mathcal{F}_i . Figure 7 illustrates the setting. Note that all testing queries exhibit field combinations unseen in training. This is to simulate the difficulty the model often encounters when scaling to large knowledge sources, which usually poses a great challenge on model’s generalization ability. We evaluate the N2N model on a dataset of the first three types of relatively simple queries. The sizes of training/testing splits are 75,000 and 30,000, with equal numbers for each query type. Table 3 lists the results. The model still maintains a reasonable performance even when the compositionality of testing queries is previously unseen, showing the model’s generalization ability in tackling unseen query patterns through the composition of familiar ones, and hence the potential to scale to larger and unseen knowledge sources.

5 Related Work

This work is related to semantic parsing, which aims to parse NL queries into logical forms executable on KBs [Zettlemoyer and Collins, 2005; Artzi *et al.*, 2015]. Recent studies take a semi-supervised learning approach, and adopt the results of query execution as indirect supervision to train a parser [Berant *et al.*, 2013; Berant and Liang, 2014; Yih *et al.*, 2015; Pasupat and Liang, 2015; Misra *et al.*, 2015]. Semantic parsers learned in this way can scale to large open domain KBs, but are inadequate for understanding complex queries because of the intractable search space incurred by the flexibility of parsing algorithms. Our work follows this approach in using query answers as indirect supervision, but jointly performs semantic parsing and query execution in distributional spaces, where the distributed representations of logical forms are implicitly learned in end-to-end QA tasks.

Our work is also related to the recent research of modeling symbolic computation using neural networks, pioneered by the development of Neural Turing Machines (NTMs) [Graves *et al.*, 2014] and the work of learning to execute (LTE) simple Python programs using LSTM [Zaremba and Sutskever, 2014]. Our work is related to both lines of research in using external memories like NTMs and learning by executing like LTE. As a highlight and difference, our work employs multiple layers of deep memories, with the neural network operations highly customized towards querying KB tables.

Perhaps the most related work is the recently proposed

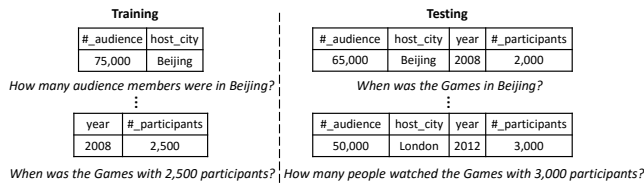


Figure 7: Expanded knowledge source querying simulation

NEURAL PROGRAMMER [Neelakantan *et al.*, 2015], which studies the same task of executing queries on tables using DNNs. While in NEURAL PROGRAMMER, the query planning is modeled using DNNs to determine which operation to execute at each step, the symbolic operations are predefined by users. **In contrast our model is fully neuralized:** it models both the query planning and query execution using DNNs, which are jointly optimized via end-to-end training. Our model learns symbolic operations using a data-driven approach. We also present results on NL queries and demonstrate that a fully neural system is capable of executing compositional logic operations up to a certain level of complexity.

6 Conclusion

We propose NEURAL ENQUIRER, a fully neural, end-to-end differentiable network that learns to execute compositional natural language queries on knowledge base tables.

References

- [Artzi *et al.*, 2015] Yoav Artzi, Kenton Lee, and Luke Zettlemoyer. Broad-coverage CCG semantic parsing with AMR. In *EMNLP*, pages 1699–1710, 2015.
- [Bahdanau *et al.*, 2015] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. In *ICLR*, 2015.
- [Berant and Liang, 2014] Jonathan Berant and Percy Liang. Semantic parsing via paraphrasing. In *ACL (1)*, pages 1415–1425, 2014.
- [Berant *et al.*, 2013] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. Semantic parsing on freebase from question-answer pairs. In *EMNLP*, pages 1533–1544, 2013.
- [Bordes *et al.*, 2013] Antoine Bordes, Nicolas Usunier, Alberto Garcia-Durn, Jason Weston, and Oksana Yakhnenko. Translating embeddings for modeling multi-relational data. In *NIPS*, pages 2787–2795, 2013.
- [Chen and Mooney, 2008] David L. Chen and Raymond J. Mooney. Learning to sportscast: A test of grounded language acquisition. In *ICML*, pages 128–135, 2008.
- [Clarke *et al.*, 2010] James Clarke, Dan Goldwasser, Ming-Wei Chang, and Dan Roth. Driving semantic parsing from the world’s response. In *CoNLL*, pages 18–27, 2010.
- [Graves *et al.*, 2014] Alex Graves, Greg Wayne, and Ivo Danihelka. Neural turing machines. *CoRR*, abs/1410.5401, 2014.
- [Kim and Mooney, 2012] Joohyun Kim and Raymond J. Mooney. Unsupervised pcfg induction for grounded language learning with highly ambiguous supervision. In *EMNLP-CoNLL*, pages 433–444, 2012.
- [Liang *et al.*, 2011] Percy Liang, Michael I. Jordan, and Dan Klein. Learning dependency-based compositional semantics. In *ACL (1)*, pages 590–599, 2011.
- [Misra *et al.*, 2015] Dipendra Kumar Misra, Kejia Tao, Percy Liang, and Ashutosh Saxena. Environment-driven lexicon induction for high-level instructions. In *ACL (1)*, pages 992–1002, 2015.
- [Neelakantan *et al.*, 2015] Arvind Neelakantan, Quoc V. Le, and Ilya Sutskever. Neural programmer: Inducing latent programs with gradient descent. *CoRR*, abs/1511.04834, 2015.
- [Pasupat and Liang, 2015] Panupong Pasupat and Percy Liang. Compositional semantic parsing on semi-structured tables. In *ACL (1)*, pages 1470–1480, 2015.
- [Wen *et al.*, 2015] Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Pei hao Su, David Vandyke, and Steve J. Young. Semantically conditioned lstm-based natural language generation for spoken dialogue systems. In *EMNLP*, pages 1711–1721, 2015.
- [Weston *et al.*, 2015] Jason Weston, Antoine Bordes, Sumit Chopra, and Tomas Mikolov. Towards ai-complete question answering: A set of prerequisite toy tasks. *CoRR*, abs/1502.05698, 2015.
- [Yih *et al.*, 2015] Wentau Yih, Mingwei Chang, Xiaodong He, and Jianfeng Gao. Semantic parsing via staged query graph generation: Question answering with knowledge base. In *ACL (1)*, pages 1321–1331, 2015.
- [Zaremba and Sutskever, 2014] Wojciech Zaremba and Ilya Sutskever. Learning to execute. *CoRR*, abs/1410.4615, 2014.
- [Zeiler, 2012] Matthew D. Zeiler. ADADELTA: an adaptive learning rate method. *CoRR*, abs/1212.5701, 2012.
- [Zettlemoyer and Collins, 2005] Luke S. Zettlemoyer and Michael Collins. Learning to map sentences to logical form: Structured classification with probabilistic categorical grammars. In *UAI*, pages 658–666, 2005.