

# Paper Reading

---

Jun Gao

June 26, 2018

Tencent AI Lab

# Neural Generative Question Answering [IJCAI2016]

---

# Introduction

This paper presents an end-to-end neural network model, named Neural Generative Question Answering (GENQA), that can generate answers to simple factoid questions, based on the facts in a knowledge-base.

- The model is built on the encoder-decoder framework for sequence-to-sequence learning, while equipped with the ability to enquire a knowledge-base
- Its decoder can switch between generating a common word and outputting a term ) retrieved from knowledge-base with a certain probability.
- The model is trained on a dataset composed of real world question-answer pairs associated with triples in the knowledge-base.

# The GENQA Model

The GENQA model consists of **Interpreter**, **Enquirer**, **Answerer**, and an external knowledgebase. **Answerer** further consists of **Attention Model** and **Generator**.

- **Interpreter** transforms the natural language question  $Q$  into a representation  $H_Q$  and saves it in the short-term memory.
- **Enquirer** takes  $H_Q$  as input to interact with the knowledge-base in the long-term memory, retrieves relevant facts (triples) from the knowledge-base, and summarizes the result in a vector  $r_Q$ .
- The **Answerer** feeds on the question representation  $r_Q$  as well as the vector  $r_Q$  and generates an answer with Generator.

Given the question represented as word sequence  $Q = (x_1, \dots, x_{T_Q})$ , Interpreter encodes it to an array of vector representations.

- In our implementation, we adopt a bi-directional recurrent neural network (GRU).
- By concatenating the hidden states (denoted as  $(\mathbf{h}_1, \dots, \mathbf{h}_{T_Q})$ ), the embeddings of words (denoted as  $(\mathbf{x}_1, \dots, \mathbf{x}_{T_Q})$ ), and the one-hot representations of words, we obtain an array of vectors  $\mathbf{H}_Q = (\tilde{\mathbf{h}}_1, \dots, \tilde{\mathbf{h}}_{T_Q})$ , where  $\tilde{\mathbf{h}}_t = [\mathbf{h}_t; \mathbf{x}_t; \mathbf{x}_t]$ .
- This array of vectors is saved in the short-term memory, allowing for further processing by Enquirer and Answerer.

- Enquirer first performs term-level matching to retrieve a list of relevant candidate triples, denoted as  $\tau_Q = \{\tau_k\}_{k=1}^{k_Q}$ .  $k_Q$  is the number of candidate triples.
- After obtaining  $\tau_Q$ , Enquirer calculates the relevance (matching) scores between the question and the  $K_Q$  triples. The  $k^{th}$  element of  $\mathbf{r}_Q$  is defined as the probability

$$r_{Q_k} = \frac{e^{S(Q, \tau_k)}}{\sum_{k'=1}^{K_Q} e^{S(Q, \tau_{k'})}}$$

- where  $S(Q, \tau_k)$  denotes the matching score between question  $Q$  and triple  $\tau_k$ . The probability in  $\mathbf{r}_Q$  will be further taken into the probabilistic model in Answerer for generating an answer.

In this work, we provide two implementations for Enquirer to calculate the matching scores between question and triples.

- Bilinear Model: simply takes the average of the word embedding vectors in  $\mathbf{H}_Q$  as the representation of the question (with the result denoted as  $\bar{\mathbf{x}}_Q$ ).

$$\bar{S}(Q, \tau) = \bar{\mathbf{x}}_Q^T \mathbf{M} \mathbf{u}_\tau$$

where  $\mathbf{M}$  is a matrix parameterizing the matching between the question and the triple.

- CNN-based Matching Model: the question is fed to a convolutional layer followed by a max-pooling layer, and summarized as a fixed-length vector  $\hat{\mathbf{h}}_Q$ .

$$\bar{S}(Q, \tau) = f_{MLP}([\hat{\mathbf{h}}_Q; \mathbf{u}_\tau])$$

## Answerer

Answerer uses an RNN to generate an answer based on the information of question saved in the short-term memory (represented as  $\mathbf{H}_Q$ ) and the relevant facts retrieved from the long-term memory (indexed by  $\mathbf{r}_Q$ ).

In generating the  $t^{th}$  word  $y_t$  in the answer, the probability is given by the following mixture model

$$\begin{aligned} p(y_t|y_{t-1}, s_t, \mathbf{H}_Q, \mathbf{r}_Q; \theta) = \\ p(z_t = 0|s_t; \theta)p(y_t|y_{t-1}, s_t, \mathbf{H}_Q, z_t = 0; \theta) + \\ p(z_t = 1|s_t; \theta)p(y_t|\mathbf{r}_Q, z_t = 1; \theta) \end{aligned}$$

which sums the contributions from the language part and the knowledge part, with the coefficient  $p(z_t|s_t; \theta)$  being realized by a logistic regression model with  $s_t$  as input.





# Conclusion

The model is built on the encoder-decoder framework for sequence-to-sequence learning, while equipped with the ability to query a knowledge-base.

# **A Knowledge-Grounded Neural Conversation Model [AAAI2018]**

---

# Introduction

This paper presents a novel, fully data-driven, and knowledge-grounded neural conversation model aimed at producing more contentful responses.

- It offers a framework that generalizes the SEQ2SEQ approach of most previous neural conversation models, as it naturally combines conversational and non-conversational data via multi-task learning.
- It not only conditions responses based on conversation history, but also on external facts that are relevant to the current context.
- Our approach only requires a way to ground external information based on conversation context (e.g., via simple entity name matching), which makes it highly versatile and applicable in an open-domain setting.

# Grounded Response Generation

In order to infuse the response with factual information relevant to the conversational context, we propose a knowledge-grounded model architecture.

- First, we have available a large collection of world facts, which is a large collection of raw text entries indexed by named entities as keys.
- Then, given a conversational history or source sequence  $S$ , we identify the focus in  $S$ , which is the text span based on which we form a query to link to the facts.
- Finally, both conversation history and relevant facts are fed into a neural architecture that features distinct encoders for conversation history and facts.

# Dialog Encoder and Decoder

- The dialog encoder and response decoder form together a sequence-to-sequence (SEQ2SEQ model)
- This part of our model is almost identical to prior conversational SEQ2SEQ models, except that we use gated recurrent units (GRU) instead of LSTM cells.

## Facts Encoder

Given an input sentence  $S = \{s_1, s_2, \dots, s_n\}$ , and a fact set  $F = \{f_1, f_2, \dots, f_k\}$  The RNN encoder reads the input string word by word and updates its hidden state.

- $u$  is the summary of the input sentence and  $r_i$  is the bag of words representation of  $f_i$ . The hidden state of the RNN is initialized with  $\hat{u}$  to predict the response sentence R word by word.

$$m_i = Ar_i$$

$$c_i = Cr_i$$

$$p_i = \text{softmax}(u^T m_i)$$

$$o = \sum_{i=1}^k p_i c_i$$

$$\hat{u} = o + u$$

# Multi-Task Learning

We train our system using multi-task learning as a way of combining conversational data that is naturally associated with external data and other businesses. We use multi-task learning with these tasks:

- NOFACTS task: We expose the model without fact encoder with  $(S, R)$  training examples, where  $S$  represents the conversation history and  $R$  is the response.
- FACTS task: We exposes the full model with  $(\{f_1, \dots, f_k, S\}, R)$  training examples.
- AUTOENCODER task: It is similar to the FACTS task, except that we replace the response with each of the facts.

The tasks FACTS and NOFACTS are representative of how our model is intended to work, but we found that the AUTOENCODER tasks helps inject more factual content into the response.



# Multi-Task Learning

The different variants of our multi-task learned system exploits these tasks as follows:

- SEQ2SEQ: This system is trained on task NOFACTS with the 23M general conversation dataset. Since there is only one task, it is not per se a multi-task setting.
- MTASK: This system is trained on two instances of the NOFACTS task, respectively with the 23M general dataset and 1M grounded dataset (but without the facts).
- MTASK-R: This system is trained on the NOFACTS task with the 23M dataset, and the FACTS task with the 1M grounded dataset.

- MTASK-F: This system is trained on the NOFACTS task with the 23M dataset, and the AUTOENCODER task with the 1M dataset.
- MTASK-RF: This system blends MTASK-F and MTASK-R, as it incorporates 3 tasks: NOFACTS with the 23M general dataset, FACTS with the 1M grounded dataset, and AUTOENCODER again with the 1M dataset.

We use the same learning technique as (Luong et al., 2015) for multi-task learning. In each batch, all training data is sampled from one task only. For task  $i$  we define its mixing ratio value of  $\alpha_i$ , and for each batch we select randomly a new task  $i$  with probability of  $\alpha_i / \sum_j \alpha_j$  and train the system by its training data.



- The model is a largescale, scalable, fully data-driven neural conversation model that effectively exploits external knowledge, and does so without explicit slot filling.
- It generalizes the SEQ2SEQ approach to neural conversation models by naturally combining conversational and non-conversational data through multi-task learning.