

# nlp\_enjoyers at TextGraphs-17 Shared Task: Text-Graph Representations for Knowledge Graph Question Answering using all-MPNet

Nikita Kurdiukov\*, Viktoriia Zinkovich\*, Sergey Karpukhin\*, and Pavel Tikhomirov

The Skolkovo Institute of Science and Technology, Moscow, Russia

{nikita.kurdiukov, viktoriia.zinkovich, sergey.karpukhin, pavel.tikhomirov}@skoltech.ru

## Abstract

This paper presents a model for solving the Multiple Choice Question Answering (MCQA) problem, focusing on the impact of subgraph extraction from a Knowledge Graph on model performance. The proposed method combines textual and graph information by adding linearized subgraphs directly into the main question prompt with separate tokens, enhancing the performance of models working with each modality separately. The study also includes an examination of Large Language Model (LLM) backbones and the benefits of linearized subgraphs and sequence length, with efficient training achieved through fine-tuning with LoRA. The top benchmark, using subgraphs and MPNet, achieved an F1 score of 0.3887. The main limitation of the experiments is the reliance on pre-generated subgraphs/triplets from the graph, and the lack of exploration of in-context learning and prompting strategies with decoder-based architectures.

## 1 Introduction

With the exponential growth of digital information, developing tools for prompt and efficient data retrieval has become a top priority in Natural Language Processing (NLP). Many state-of-the-art approaches have been proposed to solve such problems, especially encoder-only models, including BERT (Devlin et al., 2019) and its variants, such as RoBERTa (Liu et al., 2019) and ALBERT (Lan et al., 2020), which show good performance in retrieval tasks.

However, one important area of research focuses on solving Multiple Choice Question Answering problems (MCQA), where the model needs to select one correct answer among several options autonomously, without external context (Huang et al., 2022; Sakhovskiy et al., 2024). This task remains quite challenging in NLP, as in order to answer a

\* Equal contribution

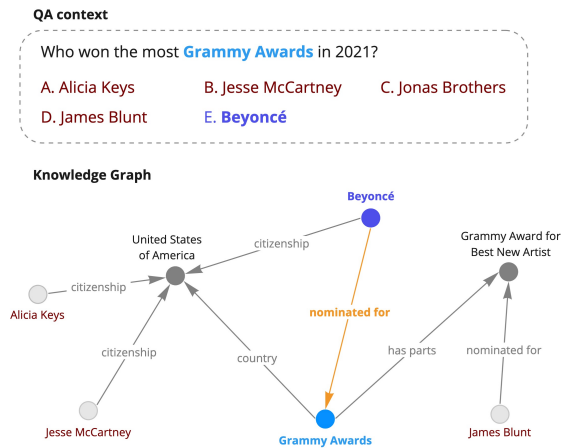


Figure 1: Example of a knowledge graph instance for a sample in a text dataset: the graph incorporates information about relations between the concept in question ("Grammy Awards") and candidate answer concepts

quiz question, the developed model should not only have a large knowledge base (Talmor et al., 2019), but also be able to make logical inferences (Li et al., 2022).

To solve such tasks, different LLMs can be applied, e.g., T5 (Raffel et al., 2020) and BART (Lewis et al., 2020), which are encoder-decoder models for natural language generation (NLG). However, even such SOTA models can generally fall short on MCQA. One common reason is that models try to predict the most likely answer in terms of grammatical construction without considering the logical coherence of the text (Robinson et al., 2023).

To enhance the performance of LLMs, in the following work, we incorporate structured knowledge graphs into the model training process (Fig. 1), as this method has been noted many times in earlier works (Salnikov et al., 2023). The graph is obtained by taking the shortest paths from all mentioned concepts in the corresponding questions to a candidate answer entity in the knowledge graph of Wikidata.

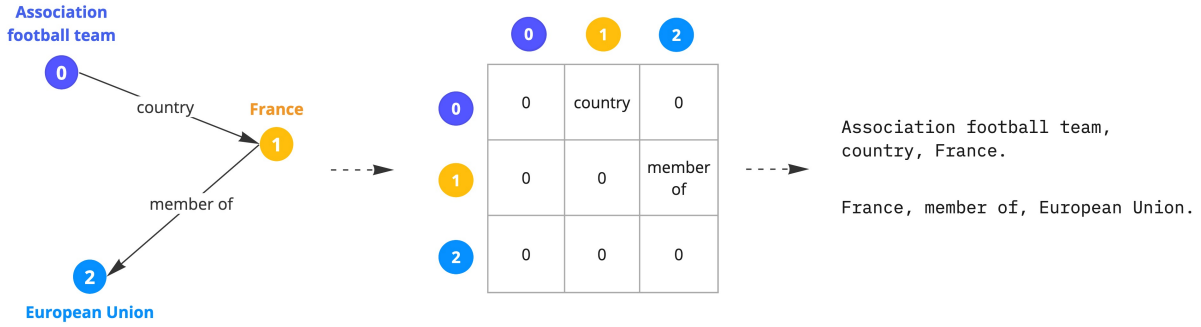


Figure 2: Example of the process of a subgraph linearization into text

Thus, the **main contributions of the following work** are as follows:

- We propose a method of combining textual and graph information. Adding linearized subgraphs directly into the main question prompt with additional separate tokens allows for improved performance of models working with each modality separately.
- We conducted a thorough study of LLM backbones and performed a wide hyper-parameter search. For efficient training, we applied fine-tuning with LoRA.

## 2 Method

We propose implementing the MPNet (Song et al., 2020) model and training it on question-answer pairs with incorporated linearized knowledge graphs. Additionally, we utilize the LoRA implementation from the peft library and apply an oversampling technique to address imbalance in the training dataset.

Our approach ultimately relies on tuning of LLM for binary classification task while also including information from the Wikidata graph domain in the LLM pipeline. The representations for target prediction on the question-answer pair are acquired by accessing the last hidden layer representation of the [CLS] token of the model.

Given the nature of the task, it is obvious that only one of the candidate answers is correct; however the number of candidate answers for a single question is not known beforehand. During inference, we utilize the knowledge that only one candidate answer is correct and select the most probable answer based on model scores. This naturally allows the use of a model trained for a classification target to rank the top-1 candidate answer.

### 2.1 Dataset

For our research, we utilized the [TextGraphs17-shared-task](#) dataset, which consists of 37,672 question-answer pairs annotated with Wikidata entities. This dataset includes 10 different types of data, notably entities from Wikidata mentioned in both the answer and the corresponding question, as well as a shortest-path graph for each <question, candidate answer> pair.

### 2.2 Evaluation metrics

During training and evaluation of our models, we use the same metrics as those present in the workshop leaderboard, which include **accuracy, precision, recall** and **F1-score**. It is important to note that accuracy is quite uninformative here due to the dataset’s imbalance, with incorrect answers constituting 90% of the data.

### 2.3 Input preprocessing

Since the subgraphs from the knowledge graph are already provided, we only need to preprocess them for the model. To incorporate information from the subgraphs, they are linearized into text according to [Salnikov et al. \(2023\)](#). The process is nearly identical, except that distinct triplets are separated with a semicolon. Specifically, subgraphs are converted to a binary adjacency matrix. If nodes indexed  $i$  and  $j$  are connected, their edge label is stored in the corresponding  $[i, j]$  matrix element. The matrix is then unraveled row by row to generate linearized sentences from corresponding triples (node\_from, edge, node\_to) in the adjacency matrix (Fig.2).

The resulting input text for the model has the following form: Question entities + ' : ' + Question + ' [SEP] ' + Linearized graph. Details of various backbones, processing pipelines and scores are reported in Sections 3 and 4.

### 3 Experiments

All fine-tuning experiments were conducted using the LoRA implementation from the peft library (Hu et al., 2021). The default LoRA parameters are as follows: a LoRA rank of 16, a LoRA alpha of 32, and a LoRA dropout of 0.1. The target modules of LoRA are the query and value weight matrices.

Our default model training is conducted for 50 epochs with best checkpoint saving, Binary Cross-entropy loss, a batch size of 64, a sequence length of 256, the AdamW optimizer, a learning rate of  $3 \cdot 10^{-4}$ , and a default weight decay of  $10^{-2}$ . Additionally, we apply oversampling during training by using a weighted sampler with probabilities inversely proportional to the labels in dataset.

We split the data into training and validation subsets by grouping samples with distinct questions in an 80:20 proportions, respectively.

#### 3.1 MiniLM experiments

The MiniLM employed is all-MiniLM<sup>1</sup>, a fine-tuned and diminished version of MiniLM by Wang et al. (2020). The training procedure is default.

#### 3.2 T5 experiments

We fine-tuned T5-Small<sup>2</sup> by Chepurova et al. (2023), which was trained on tail and entity prediction in a knowledge graph using the graph’s context represented by the node’s neighborhood. The result on the public test is presented in Table 1.

The classifier head utilizes the last hidden representation of the [EOS] token due to the encoder-decoder architecture. The model was fine-tuned for 30 epochs with the Adafactor optimizer, a learning rate of  $8 \cdot 10^{-5}$ , and a batch size of 32. LoRA alpha was set to 64 for this model.

The input format for this model was adjusted to match the original format the model was trained on. The resulting input format: 'predict [SEP] ' + Question + '[SEP]' + Linearized graph + '[SEP]' + Answer Entity

#### 3.3 MPNet experiments

Another BERT-like model we used is all-MPNet-base<sup>3</sup>. The model was trained for 20 epochs with a batch size of 32, a sequence

length of 200, the Adam optimizer, and a learning rate of  $1 \cdot 10^{-4}$ .

### 4 Additional Experiments

#### 4.1 Ablation study of sequence length and linearized graph usage

The impact of sequence length and linearized graph usage on performance was examined using the all-MiniLM model, see Table 2. We report the F1 score on the public test subset achieved by our best model checkpoints.

SL	Linearized Graph	F1 Score
256	No	0.2276
256	Yes	0.3279
<b>512</b>	<b>Yes</b>	<b>0.3463</b>

Table 2: Ablation of the Sequence Length (SL) and usage of linearized graph on all-MiniLM performance. Public test scores achieved by best model checkpoints.

#### 4.2 Usage of different backbones

Additionally, we experimented with Phrase-BERT. In brief, this model was pretrained with a contrastive objective to predict similarity between texts separated by the [SEP] token hidden state. In our pipeline, we attempted to predict the correct answer from the candidates as the 'closest' to the question. We fine-tuned this model with LoRA parametrization, as described in Section 3, structuring the input as Question entities + ' ' + Question + '[SEP] + Answer entities. Information from the graph was not used during experiments with this model. In our experiments this approach didn't provide significant quality improvements.

### 5 Conclusion

The encoder transformer architecture showed the best results in text comprehension tasks. The size of the model once again proved to have a positive influence on its performance, with the MPNet architecture outperforming MiniLM.

Despite the popularity of the T5 model for answer candidates generation, it underperformed in our experiments. Perhaps it is worth utilizing only the encoder part of the model or using a different training procedure.

Another valuable aspect that was confirmed is the benefit of incorporating graph knowledge into the model. The linearized graph indeed provided the model with valuable information, improving

<sup>1</sup><https://huggingface.co/sentence-transformers/all-MiniLM-L6-v2>

<sup>2</sup><https://huggingface.co/DeepPavlov/t5-wikidata5M-with-neighbors>

<sup>3</sup><https://huggingface.co/sentence-transformers/all-mpnet-base-v2>

Model	F1 Score
T5-Small-wikidata5M (Chepurova et al., 2023)	0.3180
all-MiniLM	0.3463
<b>all-MPNet</b>	<b>0.3887</b>

Table 1: Public test F1 scores. Best checkpoints’ scores are reported.

its ability to answer questions. More advanced subgraph/triplet sampling or generation strategies could further improve the model’s performance, making this a promising direction for future research.

## Limitations

The biggest constraint of our experiments is the reliance on pre-existing subgraphs or triplets derived from the graph. There remains a wide array of potential experiments to be conducted in this area.

Furthermore, we have not investigated the application of in-context learning and prompting techniques with decoder architectures, which could be of even more significant interest due to their current popularity and proven effectiveness.

## Acknowledgements

The authors gratefully acknowledge Professor Alexander Panchenko of Skoltech for his guidance during the "Deep Learning for Natural Language Processing" course, during which this work was conducted.

## References

- Alla Chepurova, Aydar Bulatov, Yuri Kuratov, and Mikhail Burtsev. 2023. [Better together: Enhancing generative knowledge graph completion with language models and neighborhood information](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 5306–5316, Singapore. Association for Computational Linguistics.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. 2021. [Lora: Low-rank adaptation of large language models](#).
- Zixian Huang, Ao Wu, Jiaying Zhou, Yu Gu, Yue Zhao, and Gong Cheng. 2022. [Clues before answers: Generation-enhanced multiple-choice qa](#).
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. [Albert: A lite bert for self-supervised learning of language representations](#). In *International Conference on Learning Representations*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Xiao Li, Gong Cheng, Ziheng Chen, Yawei Sun, and Yuzhong Qu. 2022. [Adalogn: Adaptive logic graph network for reasoning-based machine reading comprehension](#).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#).
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Joshua Robinson, Christopher Michael Rytting, and David Wingate. 2023. [Leveraging large language models for multiple choice question answering](#).
- Andrey Sakhovskiy, Mikhail Salnikov, Irina Nikishina, Aida Usmanova, Angelie Kraft, Cedric Möller, Debayan Banerjee, Junbo Huang, Longquan Jiang, Rana Abdullah, Xi Yan, Elena Tutubalina, Ricardo Usbeck, and Alexander Panchenko. 2024. [TextGraphs 2024 shared task on text-graph representations for knowledge graph question answering](#). In *Proceedings of the Seventeen Workshop on Graph-Based Methods for Natural Language Processing (TextGraphs-17)*, Bangkok, Thailand. Association for Computational Linguistics.
- Mikhail Salnikov, Hai Le, Prateek Rajput, Irina Nikishina, Pavel Braslavski, Valentin Malykh, and

- Alexander Panchenko. 2023. [Large language models meet knowledge graphs to answer factoid questions.](#)
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. [Mpnet: Masked and permuted pre-training for language understanding.](#)
- Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. [CommonsenseQA: A question answering challenge targeting commonsense knowledge.](#) In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4149–4158, Minneapolis, Minnesota. Association for Computational Linguistics.
- Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. 2020. [Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers.](#)