
EmbedKGQA: Improving Multi-hop Question Answering over Knowledge Graphs using Knowledge Base Embeddings

Anonymous Author(s)

Affiliation

Address

email

Reproducibility Summary

1

2 Scope of Reproducibility

3 Our work consists of four parts: (1) Reproducing results from Saxena et al. [2020] (2) Adding more experiments by
4 replacing the knowledge graph embedding method (3) and exploring the question embedding method using various
5 transformer models (4) Verifying the importance of Relation Matching (RM) module. Based on the code shared by
6 the authors, we have reproduced the results for the EmbedKGQA method. We have not performed relation matching
7 deliberately to validate point-4.

8 Methodology

9 We have used the code provided by Saxena et al. [2020] with some customization for reproducibility. In addition to
10 making the codebase more modular and easy to navigate, we have made changes to incorporate different transformers
11 in the question embedding module. Question-Answering models were trained from scratch as no pre-trained models
12 were available for our particular dataset. The code for this work is available on GitHub¹.

13 Results

14 We were able to reproduce the Hits@1 to be within $\pm 2.35\%$ of reported value (in most cases). Anomalies were ob-
15 served in 3 cases. [1] In MetaQA-KG-Full (3-hop) dataset [2] In WebQSP-KG-50 and [3] WebQSP-KG-Full datasets.
16 From our experiments on the QA model, we have found that a recent transformer architecture, SBERT (Reimers and
17 Gurevych [2019]) produced better accuracy than the original paper. Replacing RoBERTa (Liu et al. [2019]) with
18 SBERT (Reimers and Gurevych [2019]) increased the absolute accuracy by $\approx 3.4\%$ in half KG case and $\approx 0.6\%$ in the
19 full KG case. (KG: Knowledge Graph, " \approx ": Approximately)

20 What was easy

21 As the code was open-sourced, we didn't have to implement the paper giving us the liberty to customize the codebase
22 to focus on the author's claim validation, perform extended experiments and explore shared as well as new models. In
23 addition to this, pretrained KG embedding models were shared which helped in the reproduction experiment.

24 What was difficult

25 The lack of a comprehensive documentation alongwith missing comments defining functions/classes/attributes etc.,
26 made it laborious to review the code and modify it. In addition to large training times for question answering models,
27 the knowledge graph embeddings also required a significant amount of computing resources.

28 Communication with original authors

29 We had a couple of virtual meetings with Apoorv Saxena², the primary author of EmbedKGQA.

¹<https://github.com/jishnujayakumar/MLRC2020-EmbedKGQA>

²<https://apoorvumang.github.io>

30 1 Introduction

31 Knowledge is the key to question answering task. Knowledge Graph (KG) is a multi-relational graph consisting of
32 entities as nodes and relations among them as typed edges. KGs can accommodate a wide variety of facts, making
33 them one of the potential candidates for intelligent decision-making. Question Answering over KG (KGQA) task aims
34 to answer natural language queries posed over the KG. Multi-hop KGQA is a trending topic and has gained traction
35 from both academia and industry recently. Multi-hop KGQA task involves reasoning over multiple edges of the KG
36 to arrive at the correct answer.

37 Earlier works on KGs(e.g. Suchanek et al. [2007], Google [2013], Lehmann et al. [2015], Mitchell et al. [2018])
38 have some element of sparsity, i.e. they do not capture all the facts available in the real world. Recent research on
39 multi-hop KGQA has attempted to reduce this sparsity with the help of relevant external textual resources that are not
40 readily available. On the other side, KG embeddings have emerged as an effective tool to overcome the KG sparsity
41 by predicting missing links in the KG. Although effective, KG embeddings have not been explored for the multi-hop
42 KGQA task. Saxena et al. [2020] fills this gap with the proposed EmbedKGQA method.

43 This work intends to reproduce and perform an ablation (removing relation matching module) as well as extended
44 study on EmbedKGQA(Saxena et al.[2020]). EmbedKGQA claims to be the first of its kind to use KG embeddings
45 for multi-hop KGQA and improves over other state-of-the-art (SOTA) baselines.

46 2 Scope of reproducibility

47 According to Saxena et al. [2020], using ComplEx (Trouillon et al. [2016]) KG embeddings significantly improves
48 Hits@1 for multi-hop KGQA task and it has been proved with the help of the results on MetaQA (Zhang et al. [2018])
49 and WebQSP (Yih et al. [2016]) datasets. This reproducibility work tries to test this claim and conducts experiments
50 as mentioned in table:{2,3} of the original paper. Section 4.1 contains the corresponding results which support the
51 claim with some anomalies.

52 3 Methodology

53 The authors of the original paper have open-sourced the code along with the data and pre-trained ComplEx KG
54 embedding models. We have used the same codebase (commit:5d8fdbd4) and customized it for our purposes. In
55 addition to this, we have added a comprehensive documentation to make it more interpretable. Moreover, a command-
56 line functionality is also added to easily configure various transformers models in the training workflow.

57 3.1 Model descriptions

58 As shown in Figure:1, EmbedKGQA has three modules:

- 59 1. KG Embedding Module - This module contains a KG embedding model called ComplEx (Trouillon et al.
60 [2016]) to learn embeddings for all entities in the input KG. 4 pretrained models have been shared by the
61 author which contains 2 models for MetaQA-KG-`{Full, 50}` as well as 2 models WebQSP-KG-`{Full,50}`
62 dataset. Details about the dataset are mentioned in section:3.2.
- 63 2. Question Embedding Module: Given a question q , head entity h and set of answer entities A , this module
64 learns the question embeddings based on the score function defined by the KG embedding method used in 1.
- 65 3. Answer Selection Module: This module uses the outputs of module:1 and 2 to select the final answer by scor-
66 ing the (head entity and question) pair against all possible answers. The strategy is mentioned in section:`{4.4,`
67 `4.4.1}` of the original paper respectively.

68 3.2 Datasets

69 There are two datasets used in the original paper. MetaQA (Zhang et al. [2018]) and WebQSP (Yih et al. [2016]). Both
70 datasets have two portions. (1) KG data (2) QA data. KG data for both are further divided into two categories. (1)
71 Using the full KG (indicated by suffix KG-Full) and (2) Using only 50% of the facts in the respective KGs (indicated
72 by suffix KG-50). The details of generating custom KG datasets are discussed here³. Both datasets are taken from
73 here⁴. Statistics for table:`{1, 2}` have been taken from Saxena et al. [2020]. For generating question embeddings the

³<https://github.com/malllabiisc/EmbedKGQA>

⁴<https://drive.google.com/drive/folders/1RlqGBMo451TmWz9MUPTq-0KcjSd3ujxc> (As of 01/2021)

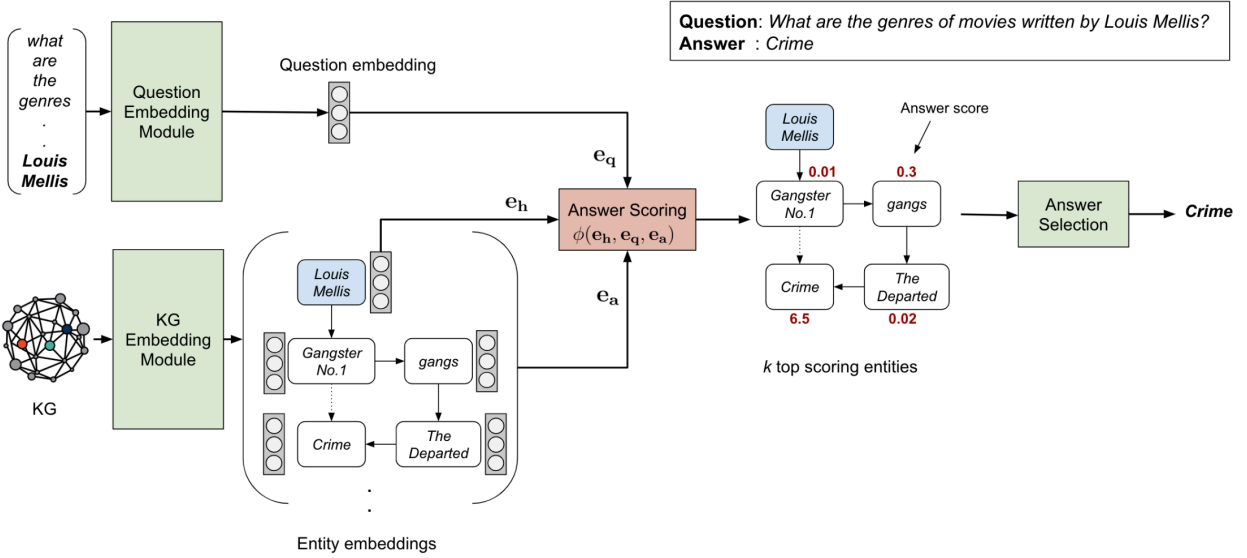


Figure 1: Overview of EmbedKGQA, the proposed method for Multi-hop KGQA.
Image source: Saxena et al. [2020].

Dataset	Train	Dev	Test
MetaQA 1-hop	96,106	9,992	9,947
MetaQA 2-hop	118,948	14,872	14,872
MetaQA 3-hop	114,196	14,274	14,274
WebQSP	2,998	100	1,639

Table 1: QA data statistics for each dataset according to Saxena et al. [2020]

Dataset	Triples	Entities	Relations	Experiment-Alias
MetaQA-KG-Full	135k	43k	9	MetaQA_full
WebQSP-KG-Full	5.7 million	1.8 million	γ	fbwq_full
MetaQA-KG-50	ϕ	-	-	MetaQA_half
WebQSP-KG-50	ψ	-	-	fbwq_half

Table 2: KG data statistics for each dataset. Refer³ for more details.

Experiment-Alias is the name used for the respective datasets in experiments.
 γ = Contains all facts that are within 2-hops of any entity mentioned in the questions of WebQSP.
 ϕ = Contains only 50% of the triples (randomly selected without replacement).
 ψ = Contains 50% of the edges sampled randomly from fbwq_full.

74 question is placed between <s> and </s> tags for all transformers except SentenceTransformer as it takes the input
75 sentence in its pure form. The preprocessing used in the original code has been used here. No additional preprocessing
76 has been performed from our end.

77 3.3 Hyper-Parameters

78 Hyper-parameters used to train the models aren't explicitly shared in the codebase or the paper, hence we decided to
79 use the default values provided in the codebase³ to compensate the lack of time. For reproduction, a pretrained model
80 shared along with the data was used; ComplEx(Trouillon et al. [2016]) was used as the knowledge graph embedding
81 method for all the KG types, i.e., full and half of both datasets types.

Type	MetaQA-KG-Full			MetaQA-KG-50		
	1-hop	2-hop	3-hop	1-hop	2-hop	3-hop
Train (t)	350 seconds	380 seconds	380 seconds	280 seconds	330 seconds	320 seconds
Validation (v)	42 seconds	95 seconds	147 seconds	47 seconds	108 seconds	182 seconds
T	10.89 hours	13.19 hours	14.64 hours	9.08 hours	12.16 hours	13.94 hours

Type	WebQSP-KG-Full	WebQSP-KG-50
Train (t)	280 seconds	300 seconds
Validation (v)	95 seconds	105 seconds
T	10.42 hours	10.92 hours

Table 3: Time for training/validation. Refer section:3.3 for hyper-parameters. (r=1, total_epochs=100)

For table:3, validate_every = The number of train routines before validation for a single epoch

Total runs (r) = Number of times the training has been performed for a particular task

Total train time (GPU hours) excluding early stopping, $T = (total_epochs \times (t + v)) \times r$

82 For the purpose of reproducibility, hyper-parameters for training MetaQA and WebQSP QA models have been taken
83 from section: {MetaQA⁵, WebQuestionsSP⁶} of the original codebase respectively. For RoBERTa (Liu et al. [2019]),
84 a pretrained model 'roberta-base' has been taken from HuggingFace transformers package (Wolf et al. [2020]). Other
85 hyper-parameters are populated by default values in the codebase³.

86 3.4 Experimental setup and code

87 Experiments have been performed on the NVIDIA DGX-1 server with 8xV100 GPUs, out of which 6 were used in
88 this work. The metric used for validating the claims is Hits@1. According to Wang et al. [2019b], Hits@k is the
89 proportion of test triples ranking in the top-k results. The code for this work is open-sourced on GitHub¹. In addition
90 to this, we have shared couple of Docker images⁷ for easy kick-starting of experiments without the hassle of setting
91 up the environment.

92 Following trained models are made available in our Docker image⁷, chosen on the basis of better performance in our
93 extended study.

- 94 • Tucker KG embedding model for Meta-QA- {Full, 50}
- 95 • QA models trained using ComplEx as KG embedding model and SBERT mentioned in table:4 as question
96 embedding model for WebQSP-KG- {Full, 50}

97 3.5 Computational requirements

98 This work has been performed on 6 V100-16GB GPUs connected via NVLink. NVLink reduced multi-GPU training
99 time by 1/4. The time required for various reproductions are mentioned in table: {5,6}.

100 3.6 Extended Experiments

101 Apart from reproducing the results mentioned in the original paper, a couple of extended experiments have been
102 performed to find answers to the following two questions:

- 103 1. Can recent KG embedding methods like Tucker (Balažević et al. [2019]) give higher accuracy on higher
104 levels of hops, i.e., 3-hop scenario to be specific compared to Trouillon et al. [2016] used in the original
105 paper?

⁵<https://github.com/malllabiisc/EmbedKGQA#metaqa>

⁶<https://github.com/malllabiisc/EmbedKGQA#webquestionssp>

⁷<https://github.com/jishnujayakumar/MLRC2020-EmbedKGQA#helpful-pointers>

Transformer	Pretrained-Model
RoBERTa	roberta-base
XLNet	xlnet-base-cased
ALBERT	albert-base-v2
SentenceTransformer (SBERT)	sentence-transformers/bert-base-nli-mean-tokens
Longformer	allenai/longformer-base-4096

Table 4: Pretrained models from HuggingFace transformers package (Wolf et al. [2020])

106 2. Can other transformer architectures like ALBERT (Lan et al. [2019]), XLNet (Yang et al. [2019]), Long-
107 former (Beltagy et al. [2020]) and SBERT (Reimers and Gurevych [2019]) improve the results on WebQSP
108 (Yih et al. [2016])?

109 Details of hyper-parameters used for these experiments are available in our GitHub repository¹. Various transformer
110 models used for experiment-2 are mentioned in table:4.

111 4 Results

112 We report results for reproducibility as well as our extended experiments. The results of reproduction have a mixed
113 nature while the ones for our extended experiments show positive signs to support claim-1, 2. Detailed discussion
114 about the results can be found in section:5. For all tables in this report, bold values indicate better performance.

115 4.1 Results reproducing original paper

116 We perform two experiments based on the two datasets introduced in section:3.2. These experiments provide vital
117 information about the results mentioned in table:{2,5} of the original paper. The results of the two are reported in
118 table:{5, 6} respectively. From the results of table:5 in Saxena et al. [2020] and table:6 in this report, it is evident
119 that relation matching(RM) is an important component in multi-hop KGQA when the given KG is considerably large,
120 i.e. {MetaQA, WebQSP} KG-Full; Definitely, WebQSP-KG-50 also shows improvement in presence of RM but the
121 performance significantly improves when applied to KG-Full setting. The author of Saxena et al. [2020] had also
122 expressed the same opinion in one of the virtual meetings.

123 4.2 Results beyond the original paper

124 We have conducted two additional experiments from our end to find an answer to claim:{1,2}. The results in table:{7,8}
125 support claim-1 but with a caveat. On the other hand, values in table:9 improve upon the results reported by the
126 original paper creating a new SOTA baseline. Additional experiments ingest custom hyper-parameters mentioned in
127 our codebase in absence of the original hyper-parameters. None of these experiments include the RM module.

128 5 Discussion

129 The reproducibility results from table:{5,6} corroborate the claims mentioned in section:2 to some extent. The repro-
130 duced version is within ± 2.4 range (positive value indicates better performance and vice-versa) except for MetaQA-
131 KG-Full dataset’s 3-hop and WebQSP-KG-Full scenario which has a significant drop of 22.5% and 18.5% respectively.
132 The absence of RM module has been reported and discussed here^{8,9,10}. For a given question, the RM module uses it’s
133 context to extract useful information from the available edges present in the KG. This information is further plugged
134 into the answer selection module to select more relevant answers. Thus, relation matching is a vital component in
135 multi-hop question answering, especially in KG-Full setting where more number of the edges are present w.r.t. KG-50
136 setting or any smaller KG w.r.t. the KG-Full setting. Results from table:{5,6} corroborates the previous statement.
137 Moreover, MetaQA-KG-50 3-hop outperforms the original model by a margin of +0.9% without using RM which is
138 an interesting observation. Apart from one reported anomaly, the reproduced results are pretty close to the original

⁸<https://github.com/malllabiisc/EmbedKGQA/issues/1>

⁹<https://github.com/malllabiisc/EmbedKGQA/issues/51>

¹⁰<https://github.com/malllabiisc/EmbedKGQA/issues/56>

Model	RM	MetaQA-KG-Full			MetaQA-KG-50		
		1-hop	2-hop	3-hop	1-hop	2-hop	3-hop
EmbedKGQA	✓	97.5	98.8	94.8	83.9	91.8	70.3
EmbedKGQA (Reproduced)	✗	95.4	96.4	<u>72.3</u>	83.2	91.6	71.2
Δ	-	-2.1	-2.4	<u>-22.5</u>	-0.7	-0.2	0.9

Table 5: Hits@1 results for original and reproduced experiments using MetaQA-KG- $\{Full, 50\}$ datasets. $\Delta = (\text{Reproduced Hits@1 without RM}) - (\text{Original Hits@1 with RM})$

Model	RM	WebQSP-KG-Full	WebQSP-KG-50
EmbedKGQA	✓	66.6	53.2
EmbedKGQA	✗	<u>48.1</u>	47.4
EmbedKGQA (Reproduced)	✗	54.9	41.3
$\Delta_{original}$	-	<u>18.5</u>	5.8
Δ	-	6.8	-6.1

Table 6: Hits@1 results for original and reproduced experiments using WebQSP-KG- $\{Full, 50\}$ datasets. $\Delta = (\text{Reproduced Hits@1 without RM}) - (\text{Original Hits@1 without RM})$, $\Delta_{original} = (\text{Original Hits@1 with RM}) - (\text{Original Hits@1 without RM})$.

For table: $\{5, 6\}$, KG-Embedding-Model=ComplEx. RM=Relation Matching, ✓= inclusion, ✗= exclusion. The original values for EmbedKGQA are taken from Saxena et al. [2020]. Underline indicates anomaly due to the absence of RM module.

KG-Model	MetaQA-KG-Full								
	1-hop			2-hop			3-hop		
	Hits@1	Hits@5	Hits@10	Hits@1	Hits@5	Hits@10	Hits@1	Hits@5	Hits@10
ComplEx	95.39	99.83	99.97	96.46	99.02	99.27	72.33	93.27	95.66
TuckER	95.51	99.81	99.97	93.13	98.7	99.28	73.81	93.6	96.09
Δ	0.12	-0.02	0	-3.33	-0.32	0.01	1.48	0.33	0.43

Table 7: Comparison of ComplEx with TuckER based on Hits@ k results for MetaQA-KG-Full dataset. $k \in \{1, 5, 10\}$.

KG-Model	MetaQA-KG-50								
	1-hop			2-hop			3-hop		
	Hits@1	Hits@5	Hits@10	Hits@1	Hits@5	Hits@10	Hits@1	Hits@5	Hits@10
ComplEx	83.24	89.83	91.22	91.63	97.08	98.04	71.2	90.77	93.72
TuckER	83	89.36	90.41	86.07	94.66	96.4	71.96	91.16	93.94
Δ	-0.24	-0.47	-0.81	-5.56	-2.42	-1.64	0.76	0.39	0.22

Table 8: Comparison of ComplEx with TuckER based on Hits@ k results for MetaQA-KG-50 dataset. $k \in \{1, 5, 10\}$.

For table: $\{7, 8\}$, $\Delta = (\text{TuckER Hits@}k) - (\text{ComplEx Hits@}k)$.

Question-Embedding-Method	WebQSP-KG-Full			WebQSP-KG-50		
	Hits@1	Hits@5	Hits@10	Hits@1	Hits@5	Hits@10
RoBERTa (Liu et al. [2019])	54.96	67.62	71.97	41.27	51.14	54.19
XLNet (Yang et al. [2019])	51.98	64.44	69.11	39.33	49.25	52.04
ALBERT (Lan et al. [2019])	47.31	59.83	63.98	31.15	42.31	45.68
Longformer (Beltagy et al. [2020])	54.9	66.77	70.47	41.92	51.98	54.83
SBERT (Reimers and Gurevych [2019])	55.55	68.98	72.74	44.65	53.86	56.13
Δ	0.59	1.36	0.77	3.38	2.72	1.94

Table 9: Hits@ k results for recent transformer models by Wolf et al. [2020] used for generating question embeddings. KG-Embedding-Method=Complex, $\Delta = (\text{SBERT_Hits@}k - \text{RoBERTa_Hits@}k)$, $k \in \{1, 5, 10\}$

139 results in case of MetaQA dataset. Default set of hyper-parameters mentioned in the original codebase(Refer section:
140 3.3) were used in the reproducibility study. The anomaly in WebQSP-KG-Full,i.e. 18.5% drop bolsters the importance
141 of RM in KG-Full setting. The reproduced results for WebQSP-KG-50 are within $\pm 7\%$ range. The use of different
142 hyper-parameters can be one of the possible answers to this variation. This value is significant but not w.r.t. WebQSP-
143 KG-Full’s drop of 18.5% which again strengthens the importance of RM in KG-Full setting. As mentioned in 4.1, RM
144 is highly useful when the KG is considerably large.

145 From table:{7, 8}, it is clear that TuckER (Balažević et al. [2019]) performs better than ComplEx (Trouillon et al.
146 [2016]) for the 3-hop scenario for both MetaQA-KG datasets, i.e., Full and 50. Though these results strengthen claim-
147 1, a more comprehensive set of tests may lead to a concrete conclusion. (e.g., experiments employing a broader set of
148 hyper-parameters).

149 According to table:9, in all the cases, SBERT (Reimers and Gurevych [2019]) outperforms RoBERTa (Liu et al. [2019])
150 used in the original paper creating a new SOTA benchmark which supports claim-2.

151 There were some experiments which didn’t work out due to the lack of time. E.g. Using RelationalTucker3 (Wang
152 et al. [2019a]) and Simple (Kazemi and Poole [2018]) to test claim-1. Furthermore, hyper-parameter search couldn’t
153 be done due to the same reason hence we had to pick the default ones mentioned in the codebase. All these create a
154 room for further experiments and improvements.

155 5.1 What was easy

156 The paper was straightforward to understand. The open-sourced codebase helped us get kick-started.

157 5.2 What was difficult

158 The structure of the codebase made it difficult to navigate it. Since the code relied upon different techniques for the two
159 datasets, the development of one function that trains different kinds of KG embeddings and another function that trains
160 different kinds of QA models for both datasets was difficult. MetaQA uses LSTM/GRU (Hochreiter and Schmidhuber
161 [1997] / Chung et al. [2014]) while WebQSP uses RoBERTa (Liu et al. [2019]) to perform the same task of generating
162 question embeddings. Also, training KG embeddings for MetaQA yields files in the form of NumPy (Harris et al.
163 [2020]) files while WebQSP uses LibKGE (Broscheit et al. [2020]) for the same purpose which produces LibKGE
164 specific KG embedding(KGE) models. Reproduction and the extensive study was a bit hard in the beginning as KGE
165 and question embedding methodology varied for both datasets. After having a couple of virtual meetings with the
166 author and code review, it became easier to conduct the planned experiments. The unavailability of hyper-parameters
167 used to train each module increased the experiment cycle by multi-fold.

168 5.3 Communication with original authors

169 We had a couple of virtual meetings with the primary author of Saxena et al. [2020]. Though it was daunting to
170 understand the codebase due to the reasons mentioned in section:5.2 with the help and support of the author, it became
171 easier to navigate the codebase.

172 5.4 Future Scope

173 We think that there is a wide range of empirical analysis and experimentation that can be performed for multi-hop QA
174 task, out of which we are sharing a few here:

- 175 1. KG embedding compression (Sachan [2020])
- 176 2. Using recent transformer models like Performer (Choromanski et al. [2020]), Reformer (Kitaev et al. [2020])
177 etc. for generating question embeddings.
- 178 3. Using low-dimensional hyperbolic KG embeddings (Chami et al. [2020]) in KG embedding module along
179 with hyperbolic word embeddings (Dhingra et al. [2018]) for question embedding module.
- 180 4. A new approach for sentence embedding, SBERT-WK (Wang and Kuo [2020]) instead of SBERT (Reimers
181 and Gurevych [2019]) can be tried out.

182 References

- 183 Ivana Balažević, Carl Allen, and Timothy M Hospedales. Tucker: Tensor factorization for knowledge graph comple-
184 tion. *arXiv preprint arXiv:1901.09590*, 2019.
- 185 Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer, 2020.
- 186 Samuel Broscheit, Daniel Ruffinelli, Adrian Kochsiek, Patrick Betz, and Rainer Gemulla. Libkge-a knowledge graph
187 embedding library for reproducible research. In *Proceedings of the 2020 Conference on Empirical Methods in*
188 *Natural Language Processing: System Demonstrations*, pages 165–174, 2020.
- 189 Ines Chami, Adva Wolf, Da-Cheng Juan, Frederic Sala, Sujith Ravi, and Christopher Ré. Low-dimensional hyperbolic
190 knowledge graph embeddings. In *Proceedings of the 58th Annual Meeting of the Association for Computational*
191 *Linguistics*, pages 6901–6914, Online, July 2020. Association for Computational Linguistics. doi: 10.18653/v1/
192 2020.acl-main.617. URL <https://www.aclweb.org/anthology/2020.acl-main.617>.
- 193 Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter
194 Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. Re-
195 thinking attention with performers, 2020.
- 196 Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent
197 neural networks on sequence modeling, 2014.
- 198 Bhuwan Dhingra, Christopher Shallue, Mohammad Norouzi, Andrew Dai, and George Dahl. Embedding text in
199 hyperbolic spaces. In *Proceedings of the Twelfth Workshop on Graph-Based Methods for Natural Language Pro-*
200 *cessing (TextGraphs-12)*, pages 59–69, New Orleans, Louisiana, USA, June 2018. Association for Computational
201 Linguistics. doi: 10.18653/v1/W18-1708. URL <https://www.aclweb.org/anthology/W18-1708>.
- 202 Google. Freebase data dumps. <https://developers.google.com/freebase/data>, 2013.
- 203 Charles R. Harris, K. Jarrod Millman, Stéfan J van der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric
204 Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti Picus, Stephan Hoyer, Marten H. van
205 Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-
206 Marchant, Kevin Sheppard, Tyler Reddy, Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E.
207 Oliphant. Array programming with NumPy. *Nature*, 585:357362, 2020. doi: 10.1038/s41586-020-2649-2.
- 208 Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- 209 Seyed Mehran Kazemi and David Poole. Simple embedding for link prediction in knowledge graphs, 2018.
- 210 Nikita Kitaev, ukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer, 2020.
- 211 Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite
212 bert for self-supervised learning of language representations. *arXiv preprint arXiv:1909.11942*, 2019.
- 213 Jens Lehmann, Robert Isele, Max Jakob, Anja Jentzsch, Dimitris Kontokostas, Pablo N Mendes, Sebastian Hellmann,
214 Mohamed Morsey, Patrick Van Kleef, Sören Auer, et al. Dbpedia—a large-scale, multilingual knowledge base
215 extracted from wikipedia. *Semantic web*, 6(2):167–195, 2015.

- 216 Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettle-
217 moyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach, 2019.
- 218 T. Mitchell, W. Cohen, E. Hruschka, P. Talukdar, B. Yang, J. Betteridge, A. Carlson, B. Dalvi, M. Gardner, B. Kisiel,
219 J. Krishnamurthy, N. Lao, K. Mazaitis, T. Mohamed, N. Nakashole, E. Platanios, A. Ritter, M. Samadi, B. Settles,
220 R. Wang, D. Wijaya, A. Gupta, X. Chen, A. Saparov, M. Greaves, and J. Welling. Never-ending learning. *Commun.*
221 *ACM*, 61(5):103115, April 2018. ISSN 0001-0782. doi: 10.1145/3191513. URL [https://doi.org/10.1145/](https://doi.org/10.1145/3191513)
222 3191513.
- 223 Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings*
224 *of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational
225 Linguistics, 11 2019. URL <https://arxiv.org/abs/1908.10084>.
- 226 Mrinmaya Sachan. Knowledge graph embedding compression. In *Proceedings of the 58th Annual Meeting of the*
227 *Association for Computational Linguistics*, pages 2681–2691, 2020.
- 228 Apoorv Saxena, Aditay Tripathi, and Partha Talukdar. Improving multi-hop question answering over knowledge
229 graphs using knowledge base embeddings. In *Proceedings of the 58th Annual Meeting of the Association for*
230 *Computational Linguistics*, pages 4498–4507, Online, July 2020. Association for Computational Linguistics. doi:
231 10.18653/v1/2020.acl-main.412. URL <https://www.aclweb.org/anthology/2020.acl-main.412>.
- 232 Fabian M. Suchanek, Gjergji Kasneci, and Gerhard Weikum. Yago: A core of semantic knowledge. In *Proceedings*
233 *of the 16th International Conference on World Wide Web*, WWW ’07, page 697706, New York, NY, USA, 2007.
234 Association for Computing Machinery. ISBN 9781595936547. doi: 10.1145/1242572.1242667. URL <https://doi.org/10.1145/1242572.1242667>.
235 [//doi.org/10.1145/1242572.1242667](https://doi.org/10.1145/1242572.1242667).
- 236 Théo Trouillon, Johannes Welbl, Sebastian Riedel, Éric Gaussier, and Guillaume Bouchard. Complex embeddings
237 for simple link prediction. In *Proceedings of the 33rd International Conference on International Conference on*
238 *Machine Learning - Volume 48*, ICML’16, page 20712080. JMLR.org, 2016.
- 239 B. Wang and C. J. Kuo. SBERT-WK: A sentence embedding method by dissecting BERT-based word models.
240 *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:2146–2157, 2020.
- 241 Yanjie Wang, Samuel Broscheit, and Rainer Gemulla. A relational tucker decomposition for multi-relational link
242 prediction, 2019a.
- 243 Yanjie Wang, Daniel Ruffinelli, Rainer Gemulla, Samuel Broscheit, and Christian Meilicke. On evaluating embedding
244 models for knowledge base completion, 2019b.
- 245 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim
246 Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite,
247 Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M.
248 Rush. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on*
249 *Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online, October 2020.
250 Association for Computational Linguistics. URL [https://www.aclweb.org/anthology/2020.emnlp-demos.](https://www.aclweb.org/anthology/2020.emnlp-demos.6)
251 6.
- 252 Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. Xlnet: Generalized
253 autoregressive pretraining for language understanding. In *Advances in neural information processing systems*, pages
254 5753–5763, 2019.
- 255 Wen-tau Yih, Matthew Richardson, Chris Meek, Ming-Wei Chang, and Jina Suh. The value of semantic parse labeling
256 for knowledge base question answering. In *Proceedings of the 54th Annual Meeting of the Association for Computa-*
257 *tional Linguistics (Volume 2: Short Papers)*, pages 201–206, Berlin, Germany, August 2016. Association for Com-
258 putational Linguistics. doi: 10.18653/v1/P16-2033. URL <https://www.aclweb.org/anthology/P16-2033>.
- 259 Yuyu Zhang, Hanjun Dai, Zornitsa Kozareva, Alexander J Smola, and Le Song. Variational reasoning for question
260 answering with knowledge graph. In *AAAI*, 2018.