

Unsupervised multiple-choice question generation for out-of-domain Q&A fine-tuning

Anonymous ACL submission

Abstract

Pre-trained models have shown very good performances on a number of question answering benchmarks especially when fine-tuned on multiple question answering datasets at once. In this work, we propose an approach for generating a fine-tuning dataset thanks to a rule-based algorithm that generates questions and answers from unannotated sentences. We show that the state-of-the-art model UnifiedQA can greatly benefit from such a system on a multiple-choice benchmark about physics, biology and chemistry it has never been trained on. We further show that improved performances may be obtained by selecting the most challenging distractors (wrong answers), with a dedicated ranker based on a pretrained RoBERTa model.

1 Introduction

In the past years, deep learning models have greatly improved their performances on a large range of question answering tasks, especially using pre-trained models such as BERT (Devlin et al., 2018), RoBERTa (Liu et al., 2019) and T5 (Raffel et al., 2020). More recently, these models have shown even better performances when fine-tuned on multiple question answering datasets at once. Such a model is UnifiedQA (Khashabi et al., 2020), which, starting from a T5 model, is trained on a large number of question answering datasets including multiple choices, yes/no, extractive and abstractive question answering. UnifiedQA is, at the time of writing, state-of-the-art on a large number of question answering datasets including multiple-choice datasets like OpenBookQA (Mihaylov et al., 2018) or ARC (Clark et al., 2018). However, even if UnifiedQA achieves good results on previously unseen datasets, it often fails to achieve optimal performances on these datasets until it is further fine-tuned on dedicated human annotated data. This tendency is increased when the target dataset deals with questions about a very specific domain.

One solution to this problem would be to fine-tune or retrain these models with additional human annotated data. However, this is expensive both in time and resources. Instead, a lot of work has been done lately on automatically generating training data for fine-tuning or even training completely unsupervised models for question answering. One commonly used dataset for unsupervised question answering is the extractive dataset SQUAD (Rajpurkar et al., 2016). Lewis et al. (2019) proposed a question generation method for SQUAD using an unsupervised neural based translation method. Fabbri et al. (2020) and Li et al. (2020) further gave improved unsupervised performances on SQUAD and showed that simple rule-based question generation could be as effective as the previously mentioned neural method. These approaches are rarely applied to multiple-choice questions answering in part due to the difficulty of selecting distractors. A few research papers however proposed distractor selection methods for multiple-choice questions using either supervised approaches (Sakaguchi et al., 2013; Liang et al., 2018) or general purpose knowledge bases (Ren and Zhu, 2020).

In this paper, we propose an unsupervised process to generate questions, answers and associated distractors in order to fine-tune and improve the performance of the state-of-the-art model UnifiedQA on unseen domains. This method, being unsupervised, needs no additional annotated domain specific data requiring only a set of unannotated sentences of the domain of interest from which the questions are created. Contrarily to most of the aforementioned works, our aim is not to train a new completely unsupervised model but rather to incorporate new information into an existing state-of-the-art model and thus to take advantage of the question-answering knowledge already learned.

We conduct our experiments on the SciQ dataset (Welbl et al., 2017). SciQ contains multiple-

041
042
043
044
045
046
047
048
049
050
051
052
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073
074
075
076
077
078
079
080
081

Question: What type of organism is commonly used in preparation of foods such as cheese and yogurt?

- (A) **mesophilic organisms** (B) protozoa
(C) gymnosperms (D) viruses

Support text: Mesophiles are often found living in or on the bodies of humans or other animals. The optimal growth temperature of many pathogenic mesophiles is 37°C (98°F), the normal human body temperature. Mesophilic organisms have important uses in food preparation, including cheese, yogurt, beer and wine.

Figure 1: Example of a question in SciQ. The answer in bold is the correct one.

082 choice questions (4 choices) featuring subjects centered
083 around physics, biology and chemistry. An
084 example of question can be found in Figure 1.
085 We focus on the SciQ dataset because it has not
086 yet been used for training UnifiedQA and it re-
087 quires precise scientific knowledge. Furthermore,
088 our experiments reveal that the direct application
089 of UnifiedQA on the SciQ benchmark leads to a
090 much lower performance than when fine-tuning it
091 on the SciQ training set (see Section 4). Our ob-
092 jective in this work is to solve this gap between
093 UnifiedQA and UnifiedQA fine-tuned on super-
094 vised data with the unsupervised question gener-
095 ation approach described in Section 2. We addition-
096 ally test our method on two commonly used multi-
097 ple choice question answering datasets: Common-
098 senseQA (Talmor et al., 2018) and QASC (Khot
099 et al., 2020). These datasets contain questions with
100 similar domains to SciQ even though the questions
101 are slightly less specific.

102 2 Question Generation Method

103 We propose a method for generating multiple-
104 choice questions in order to fine-tune and improve
105 UnifiedQA. This process is based on 3 steps. First,
106 a set of sentences is being selected (Section 2.1)
107 from which a generic question generation system is
108 applied (Section 2.2). Then a number of distractors
109 are added to each question (Section 2.3).

110 2.1 Sentence Selection

111 Our question generation method uses a set of un-
112 annotated sentences from which the questions will be
113 generated. We compare three selection methods.

114 First, we consider a scenario where the applica-
115 tion developer does not manually collect any sen-
116 tence, but simply gives the name (or topic) of the
117 target domain. In our case, the topics are “Physics”,
118 “Biology” and “Chemistry” since these are the main
119 domains in SciQ. A simple information retrieval

120 strategy is then applied to automatically mine sen-
121 tences from Wikipedia. We first compute a list
122 of Wikipedia categories by recursively visiting all
123 subcategories starting from the target topic names.
124 The maximum recursion number is limited to 4. We
125 then extract the summary (head paragraph of each
126 Wikipedia article) for each of the articles matching
127 the previously extracted categories and subcate-
128 gories. We only keep articles with more than 800
129 average visitors per day for the last ten days (on
130 April 27, 2021), resulting in 12 656 pages.

131 The two other selection methods extract sen-
132 tences from SciQ itself and therefore are not en-
133 tirely unsupervised but rather simulate a situation
134 where we have access to unannotated texts that
135 precisely describe the domains of interest such as
136 a school book for example. The SciQ dataset in-
137 cludes a support paragraph for each question (see
138 Figure 1). Pooled together, these support para-
139 graphs provide us with a large dataset of texts about
140 the domains of interest. We gather the paragraphs
141 corresponding to all questions and split them into
142 sentences to produce a large set of sentences that
143 are no longer associated with any particular ques-
144 tion but cover all the topics found in the questions.
145 We compare two different setups. In the first one,
146 we include all the sentences extracted from the
147 train, validation and test sets thus simulating a per-
148 fect selection of sentences that cover all the knowl-
149 edge expressed in the questions. Still, we only
150 use the support paragraphs and not the annotated
151 questions themselves. As compared to the classical
152 supervised paradigm, this setting removes all anno-
153 tation costs for the application developer, but it still
154 requires to gather sentences that are deemed useful
155 for the test set of interest. We then compare this
156 setup with another one, where only the sentences
157 from the train set are included. This scenario ar-
158 guably meets more practical needs since it would
159 suffice to gather sentences close to the domain of
160 interest. The number of sentences for each dataset
161 is presented in Table 1.

162 2.2 Questions Generation

163 The generation of questions from a sentence re-
164 lies on the jsRealB text realizer (Lapalme, 2021)
165 which generates an affirmative sentence from a con-
166 stituent structure. It can also be parameterized to
167 generate variations of the original sentence such
168 as its negation, its passive form and different types
169 of questions such as *who*, *what*, *when*, etc. The

	Sentences	Questions
SciQ data	53 270	77 873
SciQ data (train only)	45 526	66 552
Wikipedia data	45 327	62 848

Table 1: Number of sentences selected for each of the datasets considered as well as the number of questions automatically generated from these sentences.

constituency structure of a sentence is most often created by a user or by a program from data. In this work, it is instead built from a Universal Dependency (UD) structure using a technique developed for *SR'19* (Lapalme, 2019). The UD structure of a sentence is the result of a dependency parse with Stanza (Qi et al., 2020). We thus have a pipeline composed of a neural dependency parser, followed by a program to create a constituency structure used as input for a text realizer, both in JavaScript. Used without modification, this would create a *complex* echo program for the original affirmative sentence, but by changing parameters, its output can vary.

In order to create questions from a single constituency structure, jsRealB uses the *classical* grammar transformations: for a *who* question, it removes the subject (i.e. the first noun phrase before the verb phrase), for a *what* question, it removes the direct object (i.e. the first noun phrase within the verb phrase); for other types of questions (*when, where*) it removes the first prepositional phrase within the verb phrase. Depending on the preposition the question will be a *when* or a *where*. Note that the *removed* part becomes the answer to the question.

In order to determine which questions are appropriate for a given sentence, we examine the dependency structure of the original sentence and check if it contains the required part to be removed before parameterizing the realization. The generated questions are then filtered to remove any question for which the answer is composed of a single stopword. Table 1 shows the number of questions generated for each dataset. An example of a synthetic question is shown in Figure 2.

2.3 Distractors Selection

Since SciQ is a multiple-choice dataset, we must add distractors to each question we generate, to match the format of SciQ. A simple solution to this problem is to select random distractors among answers to other similar questions generated from the dataset of sentences we gathered. Obviously, selecting random distractors may lead to a fine-tuning dataset that is too easy to solve. Therefore,

Question: What often is found living in or on the bodies of humans or other animals?
Right answer: mesophile
Random distractors:
(A) the most magnetic material in nature
(B) this energy
(C) climate
Refined distractors:
(A) carbohydrates
(B) small cell fragments called platelet
(C) echinoderm

Figure 2: Example of a synthetic question generated from the second sentence of the support paragraph in Figure 1 with a set of random distractors and with the set of refined ones.

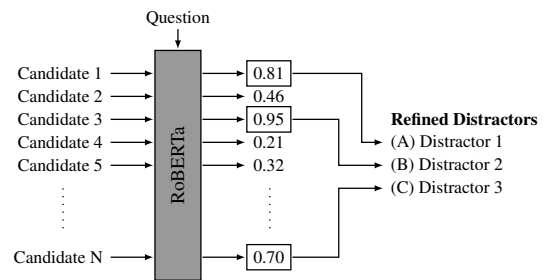


Figure 3: Description of the distractor refining method. RoBERTa scores each candidate distractor with regard to the question and the best 3 are selected to become the new refined distractors.

we propose another strategy that selects hard distractors for each question. To do so, starting from our synthetic dataset with random distractors, we fine-tune RoBERTa (Liu et al., 2019) using the standard method of training for multiple choices question answering. Each pair question/choice is fed to RoBERTa and the embedding corresponding to the first token (“[CLS]”) is given to a linear layer to produce a single scalar score for each choice. The scores corresponding to every choice for a given question are then compared to each other by a softmax and a cross-entropy loss. With this method, RoBERTa is trained to score a possible answer for a given question, based on whether or not it is a credible answer to that question. For each question, we then randomly select a number of candidate distractors from the answers to other questions and we use our trained RoBERTa to score each of these candidates. The 3 candidates with the highest scores (and thus the most credible answers) are selected. The idea is that during this first training, RoBERTa will learn a large amount of insignificant logic and the re-selection then minimizes the amount of trivial distractors. The process is better shown in Figure 3, and an example of

refined distractors can be found in Figure 2.

The number of candidate distractors is an hyperparameter. A small number of candidates may result in a situation where none of the candidates are credible enough, while a large number requires more computation time, since the score of each candidate for every question needs to be computed, and has a higher risk of proposing multiple valid answers. In our experiments, we use a number of 64 candidates in order to limit computation time.

3 Training and Implementation Details

To refine distractors, we use the “Large” version of RoBERTa and all models are trained for 4 epochs and a learning rate of 1×10^{-5} . These hyperparameters are chosen based on previous experiments with RoBERTa on other multiple-choice datasets. The final UnifiedQA fine-tuning is done using the same multiple choices question answering setup as the one used in the original UnifiedQA paper (Khashabi et al., 2020). We use the “Large” version of UnifiedQA and all the models are trained for 4 epochs using Adafactor and a learning rate of 1×10^{-5} . The learning rate is loosely tuned to get the best performance on the validation set during the supervised training of UnifiedQA. We use the Hugging Face pytorch-transformers (Wolf et al., 2020) library for model implementation. The datasets as well as the code used to create questions and evaluate the models are provided as supplementary materials and will be made available on GitHub once the anonymity requirement is lifted.

4 Results

Accuracy results in Table 2 have a 95% Wald confidence interval of $\pm 2.8\%$. The first row of Table 2 presents the accuracy results of a vanilla UnifiedQA large model on SciQ. The second line shows the accuracy when UnifiedQA is fine-tuned over the full training corpus. Our objective is thus to get as close as possible to this accuracy score using only unsupervised methods. The results using Wikipedia are the only ones that are unsupervised and therefore are the ones directly comparable to UnifiedQA with no fine-tuning or other unsupervised methods. The other results serve to illustrate what could be obtained with a tighter selection of sentences.

Fine-tuning UnifiedQA on synthetic questions with random distractors improves the results as compared to the baseline and, as expected, the closer the unlabeled sentences are to the topics

	Dev	Test
UnifiedQA (no fine-tuning)	64.6	63.4
UnifiedQA (supervised)	78.7	78.7
Unsupervised - Random distractors		
SciQ data	71.3	70.8
SciQ data (train only)	70.9	70.1
Wikipedia data	68.3	67.5
Unsupervised - Refined distractors		
SciQ data	75.4	74.2
SciQ data (train only)	73.1	72.4
Wikipedia data	70.6	69.4

Table 2: Accuracy on SciQ by UnifiedQA fine-tuned on our synthetic datasets. “SciQ data” refers to the questions generated using the support paragraphs in SciQ while “Wikipedia data” refers to questions generated using sentences harvested from Wikipedia. All scores are averaged over 3 independent runs (including the complete question generation process and the final UnifiedQA fine-tuning).

of the questions, the better is the accuracy. Hence, generating questions from only the train set of SciQ gives performances that are comparable but slightly lower to the ones obtained from the combined train, dev and test set of SciQ. Finally, questions selected from Wikipedia also improve the results, despite being loosely related to the target test corpus. Our distractor selection method further boosts the accuracy results in all setups. This suggests that a careful selection of distractors is important, and that the hard selection criterion used here seems adequate in our context.

	CQA	QASC
UnifiedQA (no fine-tuning)	60.9	44.5
UnifiedQA (supervised)	74.3	61.0
Wikipedia data (Random)	64.9	57.2
Wikipedia data (Refined)	65.1	59.4

Table 3: Accuracy results obtained on the dev set of CommonsenseQA and QASC when fine-tuning UnifiedQA using data from Wikipedia.

The results for CommonsenseQA and QASC using the same selection of sentences from Wikipedia are reported in table 3. Overall, we obtain similar results to SciQ with a large improvement of performances when generating questions and a further boost with refined distractors.

5 Conclusion

In this work, we proposed a multiple-choice question generation method that can be used to fine-tune the state-of-the-art UnifiedQA model and improve its performance on an unseen and out of domain dataset.

311
312
313
314
315
316

317
318
319
320

321
322
323
324
325
326
327

328
329
330
331

332
333
334
335
336

337
338
339
340
341
342
343

344
345

346
347
348
349
350
351

352
353
354
355
356
357

358
359
360
361
362
363
364

References

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. [Think you have solved question answering? try arc, the AI2 reasoning challenge](#). *CoRR*, abs/1803.05457.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.

Alexander Fabbri, Patrick Ng, Zhiguo Wang, Ramesh Nallapati, and Bing Xiang. 2020. [Template-based question generation from retrieved sentences for improved unsupervised question answering](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 4508–4513, Online. Association for Computational Linguistics.

D. Khashabi, S. Min, T. Khot, A. Sabharwal, O. Tafjord, P. Clark, and H. Hajishirzi. 2020. Unifiedqa: Crossing format boundaries with a single qa system. *EMNLP - findings*.

Tushar Khot, Peter Clark, Michal Guerquin, Peter Jansen, and Ashish Sabharwal. 2020. Qasc: A dataset for question answering via sentence composition. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8082–8090.

Guy Lapalme. 2019. Realizing Universal Dependencies structures using a symbolic approach. In *The Second Multilingual Surface Realisation Shared Task (SR'19): Overview and Evaluation Results*. In *Proceedings of the 2nd Workshop on Multilingual Surface Realisation (MSR), (EMNLP-2019)*, page 8 pages, Hong-Kong. ACL.

Guy Lapalme. 2021. [The jsRealB text realizer: Organization and use cases](#). (arXiv:2012.15425).

Patrick Lewis, Ludovic Denoyer, and Sebastian Riedel. 2019. [Unsupervised question answering by cloze translation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4896–4910, Florence, Italy. Association for Computational Linguistics.

Zhongli Li, Wenhui Wang, Li Dong, Furu Wei, and Ke Xu. 2020. [Harvesting and refining question-answer pairs for unsupervised QA](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6719–6728, Online. Association for Computational Linguistics.

Chen Liang, Xiao Yang, Neisarg Dave, Drew Wham, Bart Pursel, and C. Lee Giles. 2018. [Distractor generation for multiple choice questions using learning to rank](#). In *Proceedings of the Thirteenth Workshop on Innovative Use of NLP for Building Educational Applications*, pages 284–290, New Orleans, Louisiana. Association for Computational Linguistics.

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. *arXiv preprint arXiv:1907.11692*.

Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. In *EMNLP*.

Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. [Stanza: A Python natural language processing toolkit for many human languages](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*.

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.

Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [Squad: 100, 000+ questions for machine comprehension of text](#). *CoRR*, abs/1606.05250.

Siyu Ren and Kenny Q. Zhu. 2020. [Knowledge-driven distractor generation for cloze-style multiple choice questions](#). *CoRR*, abs/2004.09853.

Keisuke Sakaguchi, Yuki Arase, and Mamoru Komachi. 2013. [Discriminative approach to fill-in-the-blank quiz generation for language learners](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 238–242, Sofia, Bulgaria. Association for Computational Linguistics.

Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2018. [Commonsenseqa: A question answering challenge targeting commonsense knowledge](#). *CoRR*, abs/1811.00937.

Johannes Welbl, Nelson F. Liu, and Matt Gardner. 2017. [Crowdsourcing multiple choice science questions](#). In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 94–106, Copenhagen, Denmark. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.