

JEFF - Just Another EFFicient Reading Comprehension Test Generation

Anonymous ACL submission

Abstract

We introduce a method for generating vocabulary questions on reading comprehension of a given English article. In our approach, the method involves selecting target words in the given English article, finding synonyms as answer keys, and generating seemingly reasonable words in context as distractors. At runtime, some target words in the inputted article will be identified as questions, and automatically generating one answer key and three distractors. We present a AQG (automatic question generation) system, JEFF, that applies the method to generate questions automatically. Evaluation on a set of questions generated by JEFF shows that the method is close to the human-designed ones.

1 Introduction

Many reading comprehension vocabulary tests (e.g. *The atmosphere on Mars was thicker and liquid water widespread. Question: The word atmosphere here is A. ice B. air C. soil D. crust*) are widely used in text books and standardized tests. For example, Vocabulary Workshop¹ used this kind of practice to familiarize English learners with the meaning of target words, while TOEFL² used to assess language learners' proficiency in English.

In this paper, we generate the vocabulary type test questions of reading comprehension.³ This type of question is a four-choice multiple choice question. It have three main elements, which are target words, answer keys and distractors. Target words are the words in the article which is very

difficult to test takers, so test takers don't know their meanings and only can conjecture them by contexts. Answer keys are the words have the same meaning as corresponding target words, which are correct answers of the questions. As for distractors, they should be the wrong options and the main object of distractors is attracting test takers to choose them but they can't be possible to be answer keys.

Vocabulary test items of reading comprehension such as Vocabulary Workshop and TOEFL are typically manually designed. However, it is almost impossible to generate suitable hand-crafting test items on demand for user-selected online reading materials. Such test items could be generated automatically and on the fly if we generate paraphrases of the target word (e.g., atmosphere) in the stem sentence and use these paraphrases as an answer key (e.g., air) and distractors (e.g., ice, soil, and crust).

Consider the sentence "The atmosphere on Mars was thicker and liquid water widespread." of a given article. The best target word for this sentence is probably not "water" for B2 level learners, but rather "atmosphere" or "widespread". A good answer key might be "air" with the similar meaning to "atmosphere". Proper distractors might be "ice", "soil" or "crust", because they are reasonable in phrasal context, but impossible in sentential context. Intuitively, by paraphrasing the target word in sentential and phrasal contexts, these answer key and distractors the can be retrieved.

We present a new system, JEFF, that automatically generates items of vocabulary tests on reading comprehension. An example JEFF generates a test item is shown in Figure 1. JEFF has found a target word (e.g., atmosphere) as questions and determined an answer key and three distractors. JEFF generates these options automatically by using a pre-trained language model with contexts of various width to predict replacements. We describe the JEFF question-generating process in more detail in

¹It was published by William H. Sadlier. See <https://www.sadlier.com/>

²Test Of English as a Foreign Language, <http://www.ets.org/toefl/>

³According to The Official Guide to the TOEFL Test (<https://www.ets.org/toefl/test-takers/ibt/prepare/guides-books/>), the vocabulary type test questions of reading comprehension account for 25%-50% of TOEFL reading test, which is the highest proportion question type in the TOEFL reading test.

Running Water on Mars?

Photographic evidence suggests that liquid water once existed on the surface of Mars. Two types of flow features are seen: runoff channels and outflow channels. Geologists think that they are dried-up beds of long-gone rivers. Runoff channels on Mars speak of a time 4 billion years ago, when the **atmosphere** on Mars was thicker and liquid water widespread.

Question

The word **atmosphere** in the passage is closest in meaning to

- A. ice
- B. air
- C. soil
- D. crust

Figure 1: A sample question generated by JEFF.

Section 3.

At run-time, JEFF starts with an English article and the level⁴ of the test taker submitted by the user, and selects proper target words, and then generates answer keys and distractors for each target word. In our prototype, JEFF returns the questions to the user directly (see Figure 1), alternatively the test items can be presented to a teacher for making exam papers by identifying ideal ones.

The rest of the article is organized as follows. We review the related work in the next section. Then, we present our method for automatically finding proper target words and corresponding answer keys and distractors(Section 3). As part of our evaluation, we use indicators, including discrimination and redability, which are usually used in the tests evaluation(Section 4). We obtained astonishing result(Section 5), which shows the machine-generating questions almost meet the standard of ideal tests.

2 Related Work

Automatic question generation (AQG) has been an area of active research. Recently, the state-of-the-art in AQG research has been represented in [Ch and Saha \(2018\)](#) and [Kurdi et al. \(2020\)](#), which involves

⁴In our system, we use Common European Framework of Reference(CEFR) level (<https://www.eur.nl/en/education/language-training-centre/cefr-levels>) as input.

AGQ papers for educational purposes. In our work we address an aspect of multiple choice question generation, that was direct focus of some papers reviewed by [Kurdi et al. \(2020\)](#). We also consider levels of test takers, which tend to be different for test takers, and teachers might be interested in tailored levels of the test to the students.

More specifically, we focus on automatically generating multiple choice question, namely, deciding target words of a given article and generating answer keys and distractors. Generating distractors has long been an active topic of AQG research. An interesting approach presented by [Goto et al. \(2010\)](#) describes how to generate distractors by investigating existing human-designed tests. In general, this system generated distractors which is similar to existing tests based on the question patterns. For example, pattern “Interrogative” may generate “which”, “what”, “who”, “when”, or “where” as distractors. In contrast, we will show how to generate distractors not based on existing tests because some patterns like vocabulary will have many possibilities, and there will be few existing tests for most of vocabularies.

Most of the automatic vocabulary question generation systems are choosing distractors by existing database or dictionaries [[Sumita et al. \(2005\)](#)]. Similarly, [Wita et al. \(2018\)](#) describes a method for utilizing Japanese WordNet to generate distractors. This approach is not optimal because distractors generated by this method doesn’t refer to the contexts of the target words and this method doesn’t consider levels of vocabularies. However, the user (e.g., the teacher) is actually looking for distractors which almost make sense with the context. Generating context-related distractors is particularly important when making test takers hard to choose the answer key and help them force them to realize the given article.

Some of the traditional AQG systems have begun to generate distractors by large corpora. [Brown et al. \(2005\)](#) describes a system that generates distractors by pre-processing BNC⁵. The table used for generating distractors is built counting the frequencies of all words, because [Coniam \(1997\)](#) say distractors with similar frequency to the answer key are ideal. Additionally, the [Liu et al. \(2005\)](#) generated distractors based on the frequency and the collocation. However, this approach has not been

⁵British National Corpus, which is a 100-million-word text corpus of samples of written and spoken English

148 shown to significantly make test takers confused if
149 they don't understand the given article enough.

150 Recent work has been done on using word em-
151 bedding with the goal of selecting distractors simi-
152 lar to the target word. For example, an embedding-
153 based approach is described in Jiang and Lee
154 (2017). Aldabe and Maritxalar (2010) show how
155 to use a method almost identical to word embed-
156 ding to create distractors in Basque. These methods
157 considering meanings of words have a great break-
158 through on generating distractors. However, it's
159 hard to define how close a relationship is an appro-
160 priate distractor. If the meaning of the distractor is
161 too similar to the target word, there will be more
162 than two answer keys.

163 Recently, Sakaguchi et al. (2013) independently
164 presented a method for producing distractors based
165 on the context of the target word. This work took
166 four words beside the target word as features. In
167 contrast to our work, the length of the context this
168 work reflecting on is too short.

169 In a study more closely related to our work, Su-
170 santi et al. (2015) generated vocabulary questions
171 of reading comprehension but rather cloze tests.
172 The main difference from cloze tests is that the
173 cloze tests utilizing the target word as answer key.
174 Hence, cloze tests are unnecessary to produce an
175 extra answer key. In contrast, vocabulary ques-
176 tions of reading comprehension need to generate
177 an answer key whose meaning is the same as the
178 target word. However, Susanti et al. (2015) still
179 used WordNet to find answer keys by looking up
180 synonyms and distractors by words with the same
181 hypernym. Again, this method didn't consider the
182 context.

183 In contrast to the previous research in generat-
184 ing distractors, we present a system that generates
185 distractors related to the large scope of the context.
186 In addition, we have one more step to generate test
187 items. We need to produce an answer key. We
188 exploit the power of deep learning by paraphrasing
189 the target word of the given article.

190 3 THE JEFF system

191 Manually generating test items (e.g., The test items
192 for *The atmosphere on Mars was thicker and liq-
193 uid water widespread.*) for the user-selected ma-
194 terials is almost impossible. Human-designed test
195 items for one user is uneconomic. To generate test
196 items for the user-selected materials, a promising
197 approach is to automatically generate test items.

198 3.1 Selecting Target Words

199 In the first stage, we select a set of words in the
200 given article. These words are functioned as target
201 words which will be shown in the questions (e.g.,
202 In "The word atmosphere in the passage is closest
203 in meaning to", "atmosphere" is the target word).
204 One target word will generate one test item later.
205 The input of this stage is an arbitrary article and
206 the target word level⁶ inputted by the user, and the
207 output of this stage is a set of target words.

208 The method for selecting target words is based
209 on the level. That significantly influences the dif-
210 ficulty of test items. Hence, we consider the level
211 of test takers to select target words. We utilize the
212 Oxford 3000 and the Oxford 5000⁷ to pick the ideal
213 target words. The level of the selected target words
214 should be higher than the level of test takers. Then,
215 they can't directly realize the meaning of the target
216 words, but rather they will conjecture them by the
217 context. Additionally, based on the analysis of the
218 The Official Guide to the TOEFL Test⁸, we only
219 found noun, verb, adjective and adverb words as
220 the target words. Consequently, we follow this rule
221 and only take words of these four classifications as
222 the target words. We combine the two conditions
223 (i.e., the level and the POS) to select the target
224 words meeting conditions.

225 3.2 Paraphrasing Target Words

226 In the second stage, we find some words as the
227 replacement of the target word to paraphrase the
228 given article. These words also or almost make
229 sense in their context. Then, they will be taken as
230 answer keys or distractors (Section 3.3).

231 For this stage, we use a robust language model
232 (i.e., Devlin et al. (2018)) considering the context
233 to paraphrase the article. The input of this stage
234 is a set of target words, which is selected from the
235 last stage, and we use this language model with the
236 context of different widths to get the replacement
237 of these target words. These replacements will be
238 the output of this stage.

239 We address this problem by using a simple
240 heuristic that worked well in our experiments. In

⁶Users can choose A1, A2, B1, B2, C1, C2, or FL. They can choose one or more of them.

⁷They were presented by oxford dictionary (<https://www.oxfordlearnersdictionaries.com/about/wordlists/oxford3000-5000>). They classified most vocabularies into six levels, and we add one to put words which is not in these six levels.

⁸<https://www.ets.org/toefl/test-takers/ibt/prepare/guides-books/>

the question type of vocabulary test of reading comprehension, some test takers may choose the distractor as a result of the context of partial width. For example, if the test taker only considers the sentence “The **atmosphere** on Mars was thicker”, the distractor “crust” may be a possible answer here. By this idea, we input the context of various widths into the language model. We set widths from 1 to 49 words and use 3 words as a gap.

3.3 Generating Answer Keys and Distractors

In the third stage, we use some rules to select answer keys and distractors from the replacements from section 3.2. Thus, the output of this stage are an answer key and 3 distractors for each target word. If the number of replacements passing the rules meets the setting of multiple choice questions (i.e., an answer key and 3 distractors), it will be a test item in section 3.4.

The first rule is following the option level, which is another input from users. The replacements matching the option level are kept. Next, POS of replacements different from that of target words are eliminated. The third rule is used to distinguish answer keys and distractors. If the replacement is in the synonym dictionary⁹ of its target word, then this may be an answer key. Otherwise, if the replacement isn't in the synonym dictionary of its target word, this may be a distractor.

Unfortunately, some replacements, which are selected from the third rule and considered as distractors, are still reasonable in the given article and that will make the possible answer more than one. Therefore, we use an pre-trained natural language inference model¹⁰ to dismiss them. If the score of ENTAILMENT of the replacement is higher than 0.97, this replacement will be removed. Similarly, for potential answer key, if the score of ENTAILMENT of this is lower than 0.98, this will be knocked out.

3.4 Generating Test Items

In the final stage, we process the distractors and answer keys found in the former stage and then generate last test items.

The target words with more than 3 distractors and more than 1 answer key will be kept. The tense

⁹<https://www.lexico.com/synonyms>, which is presented by Oxford.

¹⁰The model output three scores for ENTAILMENT, NEUTRAL, and CONTRADICTION. Total of them are 1. See <https://huggingface.co/roberta-large-mnli>.

or the singular and plural of options of these target words will be adjusted to match corresponding target words. Afterwards, we randomly array the positions of options. Then, join the target word to the stem (e.g., “The word atmosphere in the passage is closest in meaning to”). Finally, we complete generating test items for a given article with given levels.

4 Experiments

JEFF was designed to generate questions for given articles. As such, JEFF will be evaluated by comparing with human-designed test items.

We collected the articles from *VOA Learning English*¹¹. For consistency, we only select the articles which are marked as advanced level and presented in 2021. We got 304 articles in the end. Then, we set target word level as B2, C1, C2, and FL and set option level as A2 to FL to implement on all collected articles. Hence, we had 304 articles with questions and kept 85 articles with more than 3 test items.

Next, we randomly reviewed 50 and selected the articles with more than three excellent test items. There were 22 articles selected.

On the other hand, we gathered the human-designed questions from The Official Guide to the TOEFL Test. There were 8 articles in Reading and 7 of them with more than 3 test items about vocabulary. Equally, we kept these 7.

Finally, we have three kinds of questions. One is machine-generated, another is machine-generated but selected by us and the other is human-designed. We randomly selected 2 articles with 3 randomly-selected test items for each kind and randomly arrayed the six articles to make an exam paper. By this way, we made 5 exam papers, in which two kinds of machine-generated test items were unrepeatable. Therefore, there were 30 test items for both machine-generated kinds and 21 test items for human-designed kind.

The two kinds of experimental setting simulate the two kinds of applications of JEFF. Directly machine-generated one simulates that learners read arbitrary online materials and want to test themselves. The kind which is machine-generated and selected by us simulates that teachers use JEFF to accelerate making exam papers.

We invited 10 English experts to evaluate these exam papers. One expert randomly reviewed one

¹¹<https://www.voanews.com/>

exam paper, so each paper was reviewed twice. The background of experts are 2 professors, 7 English teachers and 1 native speaker. The length of teaching years of 9 teachers is 8.5 years on average. They were asked to score test items by the instruction "Based on the quality of test items, score each of them (1-5 points)".

5 Results

In this section, we report the results of the experimental evaluation using the methodology described in the previous section. 74 questions are scored by 10 English experts at least 2 times.

The average score of human-designed test items is 4.06. The test items JEFF directly generated got 3.74, which is close to the human-designed ones. As for those generated by JEFF and selected by us, they got 4.11 which is even a little higher than human-designed ones.

Figure 2 shows the frequencies of three kinds of questions for each point. As we can see, the kind which is machine-generated and selected by us performs the best on 5 points, and the machine-generated kind is comparable to the human-designed kind. However, in the low point area (i.e., 1 and 2 points), there is only one review result of human-designed questions, which means almost all human-designed questions are almost above the level, while there are 9 and 14 ones for questions which is machine-generated and selected by us and machine-generated respectively. We think since we are not experts and we didn't totally select the acceptable questions for the kind which is machine-generated and selected by us, this kind doesn't perform as good as the human-designed kind in low score area.

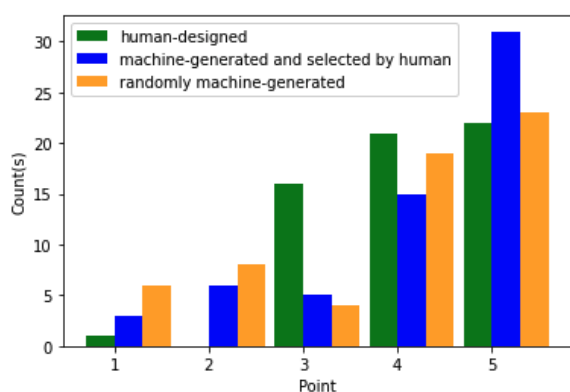


Figure 2: Frequency for each point

It is interesting to note that JEFF can almost

replace humans in the generating question task. If JEFF has a little assistant by human, it is even better than only human-designed questions.

6 Future Work and summary

Many avenues exist for future research and improvement of our system. For example, joining a mechanism for filtering out the questions may be marked as 1 or 2 points. Additionally, an interesting direction to explore is selecting target words or generating options with more than one word (e.g., moving forward). Yet another direction of research would be to consider the attributes of options. For example, in four choice "A. explorers", "B. traders", "C. immigrants" and "D. ships", the first three are identities of people and different from the option D. The test takers may easily eliminate this option.

In summary, we have introduced a method for automatically generating questions of vocabulary tests on reading comprehension with controllable level that is suitable for using them on user-selected materials for learners and making exam papers with the assistance of teachers for them. The method involves selecting target words by inputted level, paraphrasing the given articles, and generating answer keys and distractors. We have implemented and thoroughly evaluated the method as applied to practical tests. In evaluations of three kinds of test source, we have shown that the method is close to the human-designed one and even outperforms that with our help.

References

- Itziar Aldabe and Montse Maritxalar. 2010. Automatic distractor generation for domain specific texts. In *International Conference on Natural Language Processing*, pages 27–38. Springer.
- Jonathan Brown, Gwen Frishkoff, and Maxine Eskenazi. 2005. Automatic question generation for vocabulary assessment. In *Proceedings of Human Language Technology Conference and Conference on Empirical Methods in Natural Language Processing*, pages 819–826.
- Dhawaleswar Rao Ch and Sujana Kumar Saha. 2018. Automatic multiple choice question generation from text: A survey. *IEEE Transactions on Learning Technologies*, 13(1):14–25.
- David Coniam. 1997. A preliminary inquiry into using corpus word frequency data in the automatic generation of english language cloze tests. *Calico Journal*, pages 15–33.

- 421 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and
422 Kristina Toutanova. 2018. Bert: Pre-training of deep
423 bidirectional transformers for language understand-
424 ing. *arXiv preprint arXiv:1810.04805*.
- 425 Takuya Goto, Tomoko Kojiri, Toyohide Watanabe, To-
426 moharu Iwata, and Takeshi Yamada. 2010. Auto-
427 matic generation system of multiple-choice cloze
428 questions and its evaluation. *Knowledge Man-
429 agement & E-Learning: An International Journal*,
430 2(3):210–224.
- 431 Shu Jiang and John SY Lee. 2017. Distractor generation
432 for chinese fill-in-the-blank items. In *Proceedings
433 of the 12th Workshop on Innovative Use of NLP for
434 Building Educational Applications*, pages 143–148.
- 435 Ghader Kurdi, Jared Leo, Bijan Parsia, Uli Sattler, and
436 Salam Al-Emari. 2020. A systematic review of auto-
437 matic question generation for educational purposes.
438 *International Journal of Artificial Intelligence in Ed-
439 ucation*, 30(1):121–204.
- 440 Chao-Lin Liu, Chun-Hung Wang, Zhao Ming Gao, and
441 Shang-Ming Huang. 2005. Applications of lexical
442 information for algorithmically composing multiple-
443 choice cloze items. In *Proceedings of the second
444 workshop on Building Educational Applications Us-
445 ing NLP*, pages 1–8.
- 446 Keisuke Sakaguchi, Yuki Arase, and Mamoru Komachi.
447 2013. Discriminative approach to fill-in-the-blank
448 quiz generation for language learners. In *Proceed-
449 ings of the 51st Annual Meeting of the Association for
450 Computational Linguistics (Volume 2: Short Papers)*,
451 pages 238–242.
- 452 Eiichiro Sumita, Fumiaki Sugaya, and Seiichi Ya-
453 mamoto. 2005. Measuring non-native speakers’ pro-
454 ficiency of english by using a test with automati-
455 cally-generated fill-in-the-blank questions. In *Proceedings
456 of the second workshop on Building Educational Ap-
457 plications Using NLP*, pages 61–68.
- 458 Yuni Susanti, Ryu Iida, and Takenobu Tokunaga. 2015.
459 Automatic generation of english vocabulary tests. In
460 *CSEDU (1)*, pages 77–87.
- 461 Ratsameetip Wita, Sahussarin Oly, Sununta Choomok,
462 Thanabhorn Treeratsakulchai, and Surarat Wita. 2018.
463 A semantic graph-based japanese vocabulary learning
464 game. In *International Conference on Web-Based
465 Learning*, pages 140–145. Springer.