

Scalable Educational Question Generation with Pre-trained Language Models

Sahan Bulathwela, Hamze Muse and Emine Yilmaz

Centre for Artificial Intelligence, University College London, United Kingdom
{m.bulathwela, hamze.muse.20, emine.yilmaz}@ucl.ac.uk

Abstract. The automatic generation of educational questions will play a key role in scaling online education, enabling self-assessment at scale when a global population is manoeuvring their personalised learning journeys. We develop *EduQG*, a novel educational question generation model built by adapting a large language model. Our extensive experiments demonstrate that *EduQG* can produce superior educational questions by further pre-training and fine-tuning a pre-trained language model on the scientific text and science question data.

1 Introduction

Digital learning resources such as Massively Open Online Courses (MOOC) and Open Educational Resources (OER) are abundant, but they often lack associated questions that enable self-testing and skill verification [3,7,5] once the learning resources are consumed. Generating scalable educational questions is crucial for democratising education [6]. While existing language models are used for question generation, their utility in education has only been explored recently. This work demonstrates how a large language model can be adapted for educational question generation. The experiments validate the improvement of questions through additional pre-training with educational text. The study also explores the impact of pre-training data size on question generation and investigates the enhancement of educational questions through fine-tuning with a science question dataset. The experimental results show that pre-training and fine-tuning with domain-specific scientific text can outperform a state-of-the-art baseline, providing significant evidence for building an effective educational question-generation model.

2 Related Work

This work focuses on developing AI systems capable of generating educational questions for technology-enhanced learning. It involves two main sub-tasks: Question Generation (QG), where a model generates a question based on given information, and Question Answering (QA), where a model generates a response to a question. QG is essential for QA and both tasks are part of reading comprehension tasks. This paper focuses on QG specifically.

2.1 Automatic Question Generation (QG)

Automatic question generation involves creating valid and coherent questions based on given sentences and desired responses. Previous approaches have used rule-based and neural-based models, with neural models dominating in various applications [21]. Recent advancements in deep learning have led to the adoption of sequence-to-sequence models. By leveraging question-answering datasets, neural models can generate questions using both the context and expected response, ensuring high-quality questions. However, this approach often relies on an additional system to identify relevant responses [15], limiting its real-world applicability. The scarcity of public datasets also hinders the development of QG systems that generate both questions and answers. Alternatively, QG models can be trained to rely solely on the context, allowing the creation of questions that belong to a specific type [19] for the document, paragraph, or sentence level [9,8]. This work specifically focuses on the latter task setting, where only the context is used as input.

2.2 Pre-trained Language Models (PLMs) for Educational QG

In the field of educational neural question generation, state-of-the-art (SOTA) systems leverage pre-trained language models (PLMs) such as GPT-3 [2] and Google T5 [13]. These models, pre-trained on massive text corpora, enable zero-shot question generation without additional training. Recent research has demonstrated the potential for generating educational questions using GPT models [17,1].

Leaf, a cutting-edge question generation system, fine-tunes a large language model for the question and multiple-choice distracter generation [11]. It uses the SQuAD 1.1 dataset [14] to train its question generation component by fine-tuning a pre-trained T5 model [13]. This work diverges from SOTA approaches by employing pre-training to further enhance the PLM’s handling of scientific language in the educational context [12], a technique that has shown promise in domain-specific applications like medicine [20].

Our hypothesis is that pre-training with scientific text can lead to better educational question generation even when models are fine-tuned for general-purpose tasks. To evaluate the quality of generated questions, various metrics are utilized, such as BLEU, ROUGE, METEOR, F1-Score, Human Ratings, Perplexity, and Diversity [1,17,11]. This study selects a representative subset of these metrics to measure success in terms of linguistic validity and fluency.

2.3 Related Datasets

S2ORC is a corpus comprising 81.1 million English scholarly publications across various academic fields [10]. For question generation (QG) and question-answering (QA) datasets, [21] offers a comprehensive review. The Leaf system, our baseline, is designed for educational purposes by fine-tuning the T5 model using the

SQuAD 1.1 dataset, which focuses on reading comprehension [14]. However, this dataset is less suited for evaluating educational QG capabilities.

In contrast, SciQ [18] is a collection of 13,679 crowd-sourced scientific exam questions covering physics, chemistry, and other sciences. Although smaller than SQuAD, SciQ is more relevant for objectively evaluating educational QG models. Therefore, we use the SciQ dataset to assess the models developed in this work, aligning our evaluation with real-world scenarios.

3 Methodology

This study aims to study the effect of further pre-training and fine-tuning the Pre-trained Language Model (PLM) on Educational QG.

3.1 Research Questions

- **RQ1:** Can PLMs generate human-like educational questions?
- **RQ2:** Does pre-training PLMs with scientific text improve educational QG?
- **RQ3:** How does the training dataset size affect the pre-training?
- **RQ4:** Does fine-tuning the model with educational questions improve it?

3.2 Question Generations Models

Our experiments develop QG systems that utilise different PLMs trained using different task settings. It is important to note that we were not interested in training a neural model from scratch as this is impractical in real-world scenarios due to data scarcity and computational cost [2]. Instead, we used a PLM as the foundation of the different QG systems we developed for our experiments.

Baseline Leaf Model: Based on the relevant literature, we identified Leaf system [16] as the state-of-the-art educational question generation system to use as our baseline. In Leaf, the pre-trained language model, T5, a text-to-text transformer-based language model [13] (already trained on web-crawled data and Wikipedia articles) is fine-tuned for question generation using a reading comprehension dataset.

Proposed EduQG Models: The key differentiator between the baseline model and our proposal is that the EduQG model uses an additional pre-training step that further trains the PLM with scientific text documents before fine-tuning it for question generation. The expectation here is that the additional pre-training on scientific text is going to provide the PLM with more understanding of scientific language and knowledge that is relevant for generating good educational questions.

We also develop *Leaf+* and *EduQG+*, extending the Leaf model and the EduQG model that is further fine-tuned using an educational question dataset

that is more specialised than a reading comprehension dataset that only contains general-purpose questions. We hypothesise that further pre-training harnesses the model’s ability to generate educational questions.

3.3 Data

There are different types of datasets that are utilised in different stages of training the PLMs unto question generation models. These datasets allow us to:

1. Pre-train the PLM further with additional scientific language data
2. Fine-tune the PLM to carry out question generation, which is different from the initial task it was trained on
3. Objectively evaluate the performance of the question generation model

We incorporate a subset of datasets described in section 2.3 in our experiments. While the pre-training step is skipped when building the baseline Leaf model, the S2ORC corpus [10] is used for pre-training the EduQG models. The resultant language model is fine-tuned for question generation using the SQuAD 1.1 dataset [14]. Finally, we use the test set data from the SciQ question dataset [18] for evaluation. This is because the SciQ dataset exclusively contains science questions from examinations making it suitable for objectively evaluating the model’s suitability in *educational question generation*.

3.4 Evaluation Metrics

As identified in section 2, two aspects of quality are considered when evaluating the QG models, i) the prediction accuracy and ii) the linguistic quality of the generated questions. To measure the predictive accuracy of the questions, we use the BLEU score and the F1 score that is used in prior work [14,1,11]. To measure how human-like the generated questions are (i.e. linguistic quality), we use perplexity and diversity [17]. A lower perplexity score indicates better coherence. The diversity score indicates how diverse the vocabulary of the generated questions is. Larger diversity values coupled with low perplexity, indicate the use of a richer vocabulary with grammatical precision.

3.5 Experimental Setup

Our experiments are designed to answer the research questions that are outlined in section 3.1. To address RQ1, we calculate the linguistic quality-related metrics (specifically, perplexity and diversity) of the human-generated questions (the ground truth) in the SQuAD 1.1 and SciQ datasets. We hypothesise that the machine-generated questions are acceptable if they demonstrate superior or similar linguistic quality metrics in comparison to the metrics computed using the human-generated questions in the datasets (SQuAD and SciQ). The source code is available publicly ¹.

¹ https://github.com/hmuus01/Educational_QG

Fig. 1 illustrates the experiments we set out to answer RQs 2-4. The foundational language model to all the developed models (baselines and proposals) is the *T5-small* language model (hereafter referred to as T5 model). Altogether 5 models are developed (coloured boxes in the figure), all of which are evaluated using the SciQ test data.

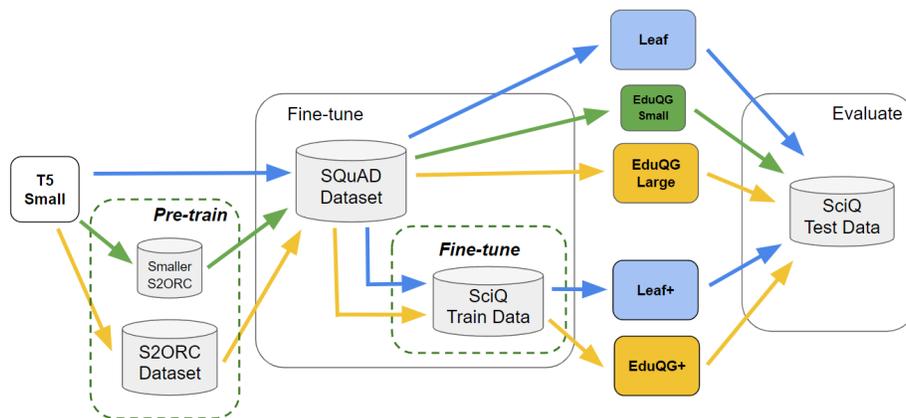


Fig. 1. Methodology for training and evaluating the baseline *Leaf* model (blue), novel *EduQG Small* (green) and *EduQG Large* (yellow) models (and their $\cdot +$ counterparts), introducing additional pre-training and fine-tuning steps (green dashed boxes) to address RQ 2,3 and 4.

To address RQ2, we develop *Leaf* and *EduQG Large* models as per Fig. 1. As the baseline, we develop the *Leaf* model by fine-tuning the T5 model on the SQuAD 1.1 dataset (blue flow of arrows in Fig. 1 through the *Leaf* model). Our proposal, *EduQG Large*, additionally pre-trains the T5 model with a down-sampled version of the S2ORC dataset that contains approx. 23.2M scientific abstracts related to Chemistry, Biology and Physics research papers (yellow flow of in Fig. 1 through the *EduQG Large* model). To answer RQ3, we use two models, i) *EduQG Large* from the previous experiment, and ii) *EduQG Small* (green flow of arrows through the *EduQG Small* model) using a smaller number of training examples from 23.5M data points. To answer RQ4, we develop *Leaf+* and *EduQG+* (blue and yellow flows of arrows passing through the $\cdot +$ models), extensions of the *Leaf* and *EduQG Large* models (baselines for RQ4 experiment) that are further fine-tuned using the training data from the SciQ dataset. While the SQuAD dataset will help the PLM to learn question generation in general, the SciQ training data is expected to teach the model *educational question generation*. We hypothesise this change will lead to superior performance.

4 Results

As per section 3.5, several experiments are executed. Table 1 shows the perplexity and diversity scores computed on the human-generated questions found in SQuAD 1.1 and SciQ datasets (RQ1). Table 2 presents the prediction accuracy and the linguistic quality metrics calculated for the models described in section 3.2 (RQ 2 and 3). Fig. 2 further elaborates the distribution of metric scores across the test data. Table 4 presents the improvement of predictive performance and the linguistic quality of the models *Leaf+* and *EduQG+* which are further fine-tuned using the SciQ training data (RQ4). Finally, Table 3 shows a handful of randomly selected test examples from the SciQ dataset where the baseline *Leaf* and the novel *EduQG Large* models have generated questions using the same context.

Table 1. Linguistic quality of the human-generated questions in the datasets.

Dataset	Perplexity ↓	Diversity ↑
SQuAD 1.1	84.16	0.779
SciQ	18.74	0.824

Table 2. Comparison of predictive performance and linguistic quality between Leaf (baseline) and EduQG (our proposals). The best and second best performance is indicated in **bold** and *italic* faces respectively. The proposed models that outperform the baseline counterpart ($p < 0.01$ in a one-tailed paired t-test) are marked with *.**.

Model	Predictive Performance					Linguistic Quality	
	BLEU-1 ↑	BLEU-2 ↑	BLEU-3 ↑	BLEU-4 ↑	F1-Score ↑	Perplexity ↓	Diversity ↑
Leaf	27.07	20.22	17.17	<i>16.46</i>	30.90	30.82	0.735
EduQG Small	<i>29.07</i> .*	<i>21.52</i> .*	<i>17.49</i> .*	15.94	<i>33.12</i> .*	34.51	<i>0.736</i>
EduQG Large	29.19 .*	21.69 .*	18.03 .*	16.76 .*	33.18 .*	<i>34.36</i>	0.749 .*

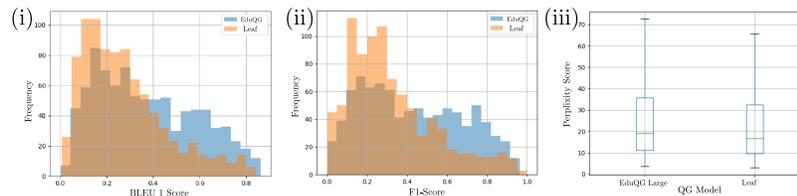


Fig. 2. The distribution of (i) BLEU 1, (ii) F1 and (iii) Perplexity Score between the Leaf and EduQG models.

Table 3. Randomly selected contexts From SciQ test data used to create questions using the Leaf and EduQG Large models.

Context	EduQG	Leaf
(1) Scientific models are useful tools for scientists. Most of Earth’s systems are extremely complex. Models allow scientists to work with systems that are nearly impossible to study as a whole. Models help scientists to understand these systems. They can analyze and make predictions about them using the models. There are different types of models.	What is used to analyze and make predictions about systems that are nearly impossible or easy to study as a whole?	What help scientists to understand the systems of Earth?
(2) Muscles That Move the Head The head, attached to the top of the vertebral column, is balanced, moved, and rotated by the neck muscles (Table 11.5). When these muscles act unilaterally, the head rotates. When they contract bilaterally, the head flexes or extends. The major muscle that laterally flexes and rotates the head is the sternocleidomastoid. In addition, both muscles working together are the flexors of the head. Place your fingers on both sides of the neck and turn your head to the left and to the right. You will feel the movement originate there. This muscle divides the neck into anterior and posterior triangles when viewed from the side (Figure 11.14).	What is the major muscle that laterally rotates?	What is the major muscle that laterally flexes and rotates the head?
(3) Biodiversity refers to the variety of life and its processes, including the variety of living organisms, the genetic differences among them, and the communities and ecosystems in which they occur. Scientists have identified about 1.9 million species alive today. They are divided into the six kingdoms of life shown in the Figure below. Scientists are still discovering new species. Thus, they do not know for sure how many species really exist today. Most estimates range from 5 to 30 million species.	What term refers to the variety of life and its processes?	How many species are identified today?
(4) Take-Home Experiment: The Pupil Look at the central transparent area of someone’s eye, the pupil, in normal room light. Estimate the diameter of the pupil. Now turn off the lights and darken the room. After a few minutes turn on the lights and promptly estimate the diameter of the pupil. What happens to the pupil as the eye adjusts to the room light? Explain your observations. The eye can detect an impressive amount of detail, considering how small the image is on the retina. To get some idea of how small the image can be, consider the following example.	What is the central transparent area of someone’s eye?	What is the name of a take-home Experiment?
(5) In both eukaryotes and prokaryotes, ribosomes are the non-membrane bound organelles where proteins are made. Ribosomes are like the machines in the factory that produce the factory’s main product. Proteins are the main product of the cell.	What are the non-membrane bound organelles where proteins are made?	What is the main product of a cell?

Table 4. Comparison of predictive performance and linguistic quality between Leaf and EduQG models in Table 2 to the new proposals further fine-tuned on SciQ training data, *Leaf+* and *EduQG+*. The best and second best performance is indicated in **bold** and *italic* faces respectively. The new models that outperform the baseline counterparts ($p < 0.01$ in a one-tailed paired t-test) are marked with $\cdot^{(*)}$.

Model	Predictive Performance					Linguistic Quality	
	BLEU-1 \uparrow	BLEU-2 \uparrow	BLEU-3 \uparrow	BLEU-4 \uparrow	F1-Score \uparrow	Perplexity \downarrow	Diversity \uparrow
Leaf	27.07	20.22	17.17	16.46	30.90	<i>30.82</i>	0.735
EduQG	29.19	21.69	18.03	16.76	33.18	34.36	0.749
Leaf+	<i>36.67^(*)</i>	<i>31.45^(*)</i>	<i>28.17^(*)</i>	24.26^(*)	<i>41.65^(*)</i>	26.43^(*)	<i>0.801^(*)</i>
EduQG+	37.20^(*)	33.86^(*)	28.49^(*)	<i>22.35^(*)</i>	43.04^(*)	33.88 ^(*)	0.812^(*)

5 Discussion

The results presented in section 4 provide sufficient information for us to answer the research questions pointed in section 3.1.

5.1 Ability of PLMs to Generate Educational Questions (RQ1)

The results presented in Tables 1 and 2 together allow us to answer RQ1. It is seen from the linguistic quality metrics in Table 2 that the perplexity score obtained by all the trained models (both baseline and novel) is acceptable. That is, the perplexity scores obtained by the model-generated questions are much lower compared to the perplexity score of the SQuAD 1.1 questions that are human-generated. The language used in academic texts can be highly advanced and rich. This is reflected by the very low perplexity score and the high vocabulary diversity score of the SciQ questions in Table 1. While the proposed models haven’t achieved a perplexity score close to the SciQ questions, having a superior perplexity in comparison to the SQuAD 1.1 question shows that the generated questions inherit coherent language and human readability. The random examples presented in Table 3 further reinforce this conclusion.

5.2 Effect of Pre-training with a Scientific Text Corpus (RQ2)

Table 2 demonstrates that the novel models, EduQG Small and EduQG Large, surpass the baseline Leaf model in nearly all evaluation metrics for predicting educational questions in the SciQ test dataset. This improvement highlights the impact of additional pre-training on scientific text for generating educational questions. With all models fine-tuned using the same question generation dataset (the SQuAD dataset), the only intervention in the proposed models is during the pre-training phase as per Fig. 1.

The T5 language model, the foundational PLM in this experiment, is trained primarily on web-crawled data and Wikipedia articles [13]. However, this training corpus lacks scientific texts, leading to a weaker understanding of scientific knowledge and language. The improvement in predicting educational questions

signifies that additional pre-training enhances the model’s grasp of scientific knowledge and language, even without specific training on educational questions during fine-tuning.

Table 2 shows higher mean perplexity scores for EduQG models, though the difference is not statistically significant. Fig. 2 (iii) indicates that the perplexity distribution between the two models is not statistically different. The observations in Table 3 further illustrate that the EduQG model generates more educational and pedagogically sound questions, as seen in rows 3 and 4.

5.3 Impact of the Training Size on the Question Quality (RQ3)

The results in Table 2 further points out the performance difference between models *EduQG Small* and *EduQG Large* where the only difference is the size of pre-training data (green vs. yellow arrows in Fig. 1). The *EduQG Large* model is superior in all evaluation metrics with the larger pre-training dataset of 23.2M data abstracts. The *EduQG Small* model outperforms the baseline *Leaf* model that uses fewer pre-training examples from the S2ORC dataset. This trend suggests that the increasing number of training examples used in the pre-training step leads to a better QG model. The increasing diversity values with the growing number of pre-training examples is also noticeable from Table 2. The improvement of BLEU and F1-Scores with diversity indicates that the validity of questions is not harmed by the diversity of the vocabulary used by the model.

5.4 Effect of Fine-tuning Using Educational Questions (RQ4)

The experimental setups of RQ 2 and 3 use a *zero-shot* evaluation where no observations from the SciQ dataset are used during the training phase. On the contrary, the experiments relating to RQ4 ($\cdot +$ models in Fig. 1) use the training data from the SciQ dataset that allows the newly proposed models, *Leaf+* and *EduQG+* to learn from educational question examples. Table 4 indicates that the additional fine-tuning significantly improves the predictive accuracy. It is noteworthy that fine-tuning is also improving the perplexity score of the generated questions which was absent in the previous experiments. We can see that the new models are outperforming the baselines. This improvement attributes to the low perplexity score of SciQ questions as per Table 1 that are exposed to the model during training.

5.5 Opportunities

The examples in Table 3 with all the above results indicate that educational QG systems are very close to becoming part of human-facing technology-enhanced learning systems (Such as X5Learn that leverages Open Educational Resources [4]). Many works in the past have shown how zero-shot question generation is operationally feasible using very large language models gated behind an API from a large corporation (Model-as-a-Service architecture) [17]. However, our

result contributes to this topic as we introduce methods to enhance openly-available PLMs (in our case, T5) to support educational QG. We intentionally use the *T5-Small* model that has 60M parameters in comparison models such as GPT-3 XL that has 1.3B parameters [2] to show that relatively small models can be trained with domestic hardware to create SOTA educational QG capabilities. Our method also gives the stakeholder full control and ownership, a critical feature for quality assurance of the downstream educational systems that rely on this model (contrary to having no control over a third party that can change their model over time). This work also informs the educational data mining community that domain-specific data can be used with language models to harness them to specific educational use cases (e.g. extend to other domains, different question types that support diverse pedagogy etc.). While the proposed systems are not perfect, the quality of AI-generated questions indicates that a teacher or an educator can re-purpose these questions with minimum effort and time. Human-in-the-loop systems can be built to support educators while their corrections will harvest more training data to improve the models over time. Educational questions can be generated at scale using the proposed model both for existing and newly created learning resources, adding more testing opportunities for learners/teachers to use when needed.

We see our work being foundational to building a series of tools that can support educators with scalable/personalised assessments. Ultimately, we have the opportunity to improve these models to the point where an intelligent tutor can rely on them to create on-demand questions to verify a learner’s knowledge state with no human intervention.

5.6 Limitations

We need to be cautious to avoid the obvious pitfalls of such automatic systems. Intelligent QG models we build tend to exhibit the patterns in the data that we feed them. We need to be mindful that we take rigorous steps to validate the datasets to be ethically and pedagogically sound. Putting emphasis on quality assurance of the training data will help us to build ethical, unbiased QG models that can benefit all learners equally.

Many intelligent learning systems exploit learner engagement signals to determine what characteristics of the system should sharpen and weaken [3]. In the context of question generation, it is important to distinguish between *bad* questions vs. *difficult* questions as the latter, although demanding, may positively impact a learner while the former will only hinder and diminish learning gains. The AI-generated questions should allow users to improve their learning gains over time.

Another gap in this work is the lack of human evaluation of the AI-generated questions. While offline evaluation on labelled datasets is useful, having teachers and learners evaluate and contrast between human vs. AI-generated questions will provide much more insightful findings that can improve this line of research in the future. Our subsequent work will focus on this aspect.

6 Conclusion

This work demonstrates the operational feasibility of adapting pre-trained language models for educational question generation. Specifically, we argue that a relatively small language model manageable with domestic hardware can be further trained and harnessed with low computational costs and produce a humanly-acceptable educational question generation model. We validate that a PLM fine-tuned with question generation data can generate questions that are linguistically valid and human-like. We show that the quality of the educational questions generated can be significantly improved by pre-training using domain-specific corpora alone. We use a corpus of scientific abstracts to empirically demonstrate this while we point out the relationship between the prediction quality and the amount of data. Going further, we improve the model’s question generation capabilities significantly by further fine-tuning it using a domain-specific question dataset, indicating fine-tuning can be used to further improve the model.

A few promising steps remain to take this work to the future. Validating the generalisability of our approach to other PLMs such as GPT [2] and extending evaluation to human experts [1,17] are the immediate next steps. Establishing methods to audit the ethical and pedagogical value of training datasets will improve the use of the downstream QG models. Identifying systematic approaches (e.g. using curriculum learning) to identify the most useful training examples would allow us to make QG models significantly better with less number of training examples leading to computational cost savings. Finally, formalising concepts such as question difficulty, and value for learning will allow us to evaluate the quality of generated questions much more pragmatically.

Acknowledgements This work is also partially supported by the European Commission-funded project "Humane AI: Toward AI Systems That Augment and Empower Humans by Understanding Us, our Society and the World Around Us" (grant 820437), EU Erasmus+ project 621586-EPP-1-2020-1-NO-EPPKA2-KA and the EPSRC Fellowship "Task Based Information Retrieval" (grant EP/P024289/1). This research is conducted as part of the X5GON project (www.x5gon.org) funded by the EU’s Horizon 2020 grant No 761758.

References

1. Bhat, S., Nguyen, H.A., Moore, S., Stamper, J., Sakr, M., Nyberg, E.: Towards automated generation and evaluation of questions in educational domains. In: Proc. of the 15th Int. Conf. on Educational Data Mining, 701. vol. 704 (2022)
2. Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J.D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al.: Language models are few-shot learners. *Advances in neural information processing systems* **33**, 1877–1901 (2020)
3. Bulathwela, S., Perez-Ortiz, M., Yilmaz, E., Shawe-Taylor, J.: Truelearn: A family of bayesian algorithms to match lifelong learners to open educational resources. In: *AAAI Conference on Artificial Intelligence* (2020)

4. Bulathwela, S., Kreitmayer, S., Pérez-Ortiz, M.: What's in It for Me? Augmenting Recommended Learning Resources with Navigable Annotations. In: Proc. of the Int. Conf. on Intelligent User Interfaces Companion (2020)
5. Bulathwela, S., Pérez-Ortiz, M., Yilmaz, E., Shawe-Taylor, J.: Semantic TrueLearn: Using Semantic Knowledge Graphs in Recommendation Systems. In: Proc. of First KGSWC International Workshop on Joint Use of Probabilistic Graphical Models and Ontology (PGMOnto) (2021), <https://arxiv.org/abs/2112.04368>
6. Bulathwela, S., Pérez-Ortiz, M., Holloway, C., Shawe-Taylor, J.: Could ai democratise education? socio-technical imaginaries of an edtech revolution. In: Proc. of NeurIPS Workshop on ML4D. arXiv (2021), <https://arxiv.org/abs/2112.02034>
7. Bulathwela, Sahan and Pérez-Ortiz, Maria and Yilmaz, Emine and Shawe-Taylor, John: Power to the Learner: Towards Human-Intuitive and Integrative Recommendations with Open Educational Resources. *Sustainability* **14**(18) (2022)
8. Du, X., Shao, J., Cardie, C.: Learning to ask: Neural question generation for reading comprehension. arXiv preprint arXiv:1705.00106 (2017)
9. Guo, H., Pasunuru, R., Bansal, M.: Soft layer-specific multi-task summarization with entailment and question generation. arXiv preprint arXiv:1805.11004 (2018)
10. Lo, K., Wang, L.L., Neumann, M., Kinney, R., Weld, D.: S2ORC: The semantic scholar open research corpus. In: Proc. of the Ann. Meet. of the ACL. Online (2020)
11. Lopez, L.E., Cruz, D.K., Cruz, J.C.B., Cheng, C.: Simplifying paragraph-level question generation via transformer language models. In: Pacific Rim International Conference on Artificial Intelligence. pp. 323–334. Springer (2021)
12. Muse, H., Bulathwela, S., Yilmaz, E.: Pre-training with scientific text improves educational question generation (student abstract). In: AAAI Conference on Artificial Intelligence (2023)
13. Raffel, C., Shazeer, N., Roberts, A., Lee, K., Narang, S., Matena, M., Zhou, Y., Li, W., Liu, P.J., et al.: Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.* **21**(140), 1–67 (2020)
14. Rajpurkar, P., Zhang, J., Lopyrev, K., Liang, P.: Squad: 100, 000+ questions for machine comprehension of text. *CoRR* **abs/1606.05250** (2016)
15. Tamang, L.J., Banjade, R., Chapagain, J., Rus, V.: Automatic question generation for scaffolding self-explanations for code comprehension. In: International Conference on Artificial Intelligence in Education. pp. 743–748. Springer (2022)
16. Vachev, K., Hardalov, M., Karadzhov, G., Georgiev, G., Koychev, I., Nakov, P.: Leaf: Multiple-choice question generation. In: Proc. of the European Conf. on Information Retrieval (2022)
17. Wang, Z., Valdez, J., Basu Mallick, D., Baraniuk, R.G.: Towards human-like educational question generation with large language models. In: Proc. of Int. Conf. on Artificial Intelligence in Education (2022)
18. Welbl, J., Liu, N.F., Gardner, M.: Crowdsourcing multiple choice science questions. In: Proc. of the 3rd Workshop on Noisy User-generated Text. ACL (Sep 2017). <https://doi.org/10.18653/v1/W17-4413>
19. Wu, X., Jiang, N., Wu, Y.: A question type driven and copy loss enhanced framework for answer-agnostic neural question generation. arXiv preprint arXiv:2005.11665 (2020)
20. Xu, H., Van Durme, B., Murray, K.: BERT, mBERT, or BiBERT? a study on contextualized embeddings for neural machine translation. In: Proc. of Conf. on Empirical Methods in Natural Language Processing (2021)
21. Zhang, R., Guo, J., Chen, L., Fan, Y., Cheng, X.: A review on question generation from natural language text. *Trans. on Information Systems* **40**(1), 1–43 (2021)