

WikiFlash: Generating Flashcards from Wikipedia Articles

Yuang Cheng¹, Yue Ding¹, Damián Pascual², Oliver Richter², Martin Volk¹, Roger Wattenhofer^{2*}

¹University of Zurich.

²ETH Zurich.

{yuang.cheng, yue.ding}@uzh.ch, volk@cl.uzh.ch, {dpascual, orichter, wattenhofer}@ethz.ch

Abstract

Flashcards, or any sort of question-answer pairs, are a fundamental tool in education. However, the creation of question-answer pairs is a tedious job which often defers independent learners from properly studying a topic. We seek to provide a tool to automatically generate flashcards from Wikipedia articles to make independent education more attractive to a broader audience. We investigate different state-of-the-art natural language processing models and propose a pipeline to generate flashcards with different levels of detail from any given article. We evaluate the proposed pipeline based on its computing time and the number of generated and filtered questions, given the proposed filtering method. In a user study, we find that the generated flashcards are evaluated as “helpful” on average. Further, users evaluated the quality of human created flashcards that are available open source as comparable to or only slightly better than the automatically generated cards.¹

1 Introduction

The recent advances in artificial intelligence make available a new set of tools that can be exploited to advance the field of personalized education. In the last years, we have seen how, thanks to new deep learning methods, machines have attained super-human performance in a large number of language-related tasks (Wang et al., 2019b,a). These methods can accelerate the development of personalized education by automatically generating instructional material.

Generating instructional materials manually is a costly task that requires instructors to select and cure large amounts of information. With a growing internet an ever-increasing (and overwhelming) amount of information and data is available. However, it is challenging for a person to learn in a systematic manner from this information. To improve human learning, it is necessary to structure the information into instructional materials that select the most relevant points and guide learning. Automatically generating these materials can widely accelerate human learning while giving each person the freedom to learn any arbitrary topic of her interest. The fast digitalization of education in response to the COVID-19 pandemic has brought many challenges.

However, a digital education system also opens opportunities thanks to the new tools at our disposal, specifically in terms of personalization.

A well-known and effective format for instructional materials are flashcards (Thalheimer, 2003). Flashcards are small cards (physical or virtual) with a question written on the front face and the answer to that question written on the back face. Flashcards stimulate learning by hiding the answer that the student is trying to learn. A big advantage of flashcards is that they are topic-independent, i.e., flashcards can be used to learn anything: languages, history, mathematics... Nevertheless, a large number of flashcards is necessary to cover a given topic or subtopic, and preparing good flashcards requires good summarization skills, all of which makes the process of manually producing flashcards challenging and time consuming.

In this work, we present a system for automatically generating flashcards about any arbitrary topic. We leverage recent advances in language processing, in particular transformer-based models (Vaswani et al., 2017), to extract questions and answers from input text. We implement our system as a web application that takes as input the title of a Wikipedia article and outputs flashcards for that article. We evaluate the application, profiling generation time and the number of flashcards produced. Furthermore, we run a user study to assess the quality of our automatically generated cards in comparison to human-created cards. The results show that the quality of our automatically generated cards is similar to the quality of cards generated by humans.

Our system has the flexibility of generating instructional materials (in the form of flashcards) for any topic the student is interested in, beyond standard curricula. We consider our web application to be both, a proof-of-concept of how current technologies allow automatic generation of materials for learning, as well as a first step towards a completely functional tool to enhance learning anywhere and about anything.

2 Related Work

Automatic question generation for educational purposes is a growing research area with many works focusing on assessment and template based question generation (Kurdi et al., 2020). In a recent trend, data driven approaches that use neural networks became more prominent in many natural lan-

* Authors in alphabetical order

¹Our application is available at: flashcard.ethz.ch

guage processing tasks, including question generation (Pan et al., 2019). These data driven approaches might struggle to extract questions that require several steps of reasoning as in the LearningQ dataset (Chen et al., 2018). However, for flashcard generation, simple factoid questions are often preferred. We therefore focus on models that perform well on the Wikipedia based SQuADv1 dataset (Rajpurkar et al., 2016), which was originally developed for question answering models but can be re-purposed for context based question generation.

On this dataset, transformer based approaches for question generation (Kriangchaivech and Wangperawong, 2019; Chan and Fan, 2019; Lopez et al., 2020; Dong et al., 2019; Bao et al., 2020) are currently performing best in terms of n -gram similarity metrics such as BLEU (Papineni et al., 2002), METEOR (Banerjee and Lavie, 2005) and ROUGE-L (Lin, 2004). This is likely due to the fact that these models benefit from large scale unsupervised pre-training. Our implementation is based on the publicly available code of (Patil, 2020), which follows ideas from (Chan and Fan, 2019; Lopez et al., 2020) and (Alberti et al., 2019) and achieves results not far behind the current state-of-the-art (Bao et al., 2020).

As a pre-processing step, text summarization can be used to reduce the text from which questions are to be generated. Automatically summarizing text is the focus of a large body of research and a number of datasets exist that are used to benchmark progress (Hermann et al., 2015; See, Liu, and Manning, 2017; Rush, Chopra, and Weston, 2015; Narayan, Cohen, and Lapata, 2018). There are two types of summarization: extractive (Zhong et al., 2019), the summary consist of sentences copied from the original text; and abstractive (Gupta and Gupta, 2019), the sentences do not coincide with the original text but the meaning does. Abstractive summarization is both, more natural and harder. Recently proposed models (Yan et al., 2020; Zhang et al., 2019; Lewis et al., 2020; Raffel et al., 2019) have achieved new state-of-the-art results as measured by ROUGE-L score. Here, we leverage this progress and use abstractive summarization as a content selection step before the question generation.

The general idea of filtering questions in a post-processing step has been explored in different settings (Kwankajornkiet, Suchato, and Punyabukkana, 2016; Blšták and Rozinajová, 2016; Liu, Rus, and Liu, 2017; Niraola and Rus, 2015; Alberti et al., 2019). Using a question-answering system to filter questions where the answers do not align was proposed by Alberti et al. to create a synthetic data corpus for pretraining a question-answering model. We use this approach in our system with slight adjustments. Compared to their approach of filtering all questions where answers do not align, we relax the filtering by allowing for questions where the extracted answers and the answers produced by the question-answer model yield a sufficient overlap.

The main contribution of this work is an end-to-end application that allows for flashcard generation based on a Wikipedia article freely chosen by the users. Our work thereby differs from the work of (Du and Cardie, 2018) that created a fixed size corpus for scientific investigation. Also,

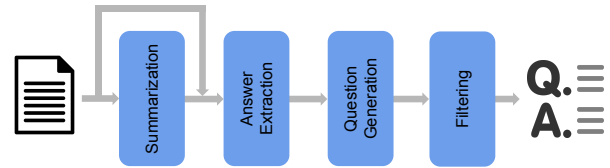


Figure 1: Pipeline of the flashcard generation system. The summarization step is optional.

despite the existence of many applications that allow for the design and/or studying of flashcards, we only encountered one working application which allows for automated flashcard creation (Examiners, 2020). This application uses a key-phrase based system for the creation of flashcards in the biological and medical domain. In contrast, our approach does not rely on any domain specific key phrases and is therefore applicable to a much wider range of topics.

3 Method

Generating meaningful flashcards from an arbitrary piece of text is not a trivial problem. Currently, there does not exist a single model that can alone perform this task. We therefore divide the flashcard generation process into four sub-tasks that cover more general and well-studied problems that can be individually addressed by state-of-the-art models. In particular, we build a pipeline consisting of four stages: summarization, answer identification, question generation and question answering.

Summarization By definition, a summary contains the most relevant information for a general understanding of the corresponding text. Thus, generating flashcards from a summary reduces the level of detail in the resulting flashcards, in comparison to using the original text as input. A summarization stage gives the user the freedom of deciding between two levels of detail for the information contained in the flashcards. If more detailed flashcards are preferred, the summarization step is skipped and the input text is passed directly to the next step of the pipeline. Otherwise, a summary is generated from the input text and fed into the next stage.

Answer extraction After the optional summarization step, we proceed to generate flashcards by identifying potential answers in the text. To this end, we use a model for answer extraction, which receives as input a piece of text and finds words or groups of words that can be answers to questions. These answers, together with the text they are extracted from, are passed as input to the next stage.

Question generation In this stage we use an answer-aware question generation model to generate answer specific questions. This way, the output of this stage is the set

of question-answer tuples that we need for flashcards. However, the question-answer tuples generated at this point tend to include some questions that either make no sense or are incorrect. Therefore, we include a final step in our pipeline to filter out unusable questions.

Filtering To filter out erroneous questions, we use a model for question answering. For each question-answer tuple we provide this model with the question and the paragraph where the answer can be found. If the answer provided by the question-answering model overlaps enough with the answer from which the question was generated, then the question-answer tuple is accepted, otherwise it is discarded.

A depiction of the complete pipeline can be seen in Figure 1. This pipeline represents a general approach to automatically generate and select flashcards that can be used for learning.

4 Implementation

We implement our flashcard generation pipeline as a web application. The interface of our application is simple and intuitive: the user should provide as input the title of a Wikipedia article she wants flashcards from. If the article name given contains a typo, correctly spelled alternatives are suggested to the user. Likewise, if the article name is redundant, i.e., there are more than one article with the same name, disambiguation results are suggested.

The user also has to specify the level of detail of the generated cards. Our current application allows for four levels, which differ in the text from which the cards are generated and whether summarization is used or not. To determine the detail level, we exploit the fact that Wikipedia articles always follow the same structure, with an introduction that contains high level information and a main text that is divided in different sections and subsections with more in-depth information. This way, we create the levels of detail offered in our application as follows:

- **Highlight Introduction:** card generation from a summary of the introduction.
- **Introduction:** card generation directly from the introduction text.
- **Highlight Full:** card generation from a summary of the entire article.
- **Full:** card generation directly from the text of the entire article.

Regarding, the implementation of each stage of our system, we use the following models to build the pipeline:

Summarization We use DistilBART for summarization (Shleifer, 2020; Lewis et al., 2020), pre-trained on the CNN/DailyMail summarization dataset (Hermann et al., 2015). The maximum input length of this model is 1024 tokens, which is less than a long Wikipedia article. To circumvent this issue, our summaries are generated paragraph-wise.

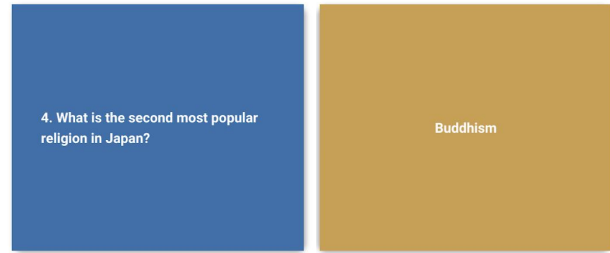


Figure 2: Example of a flashcard for the topic *Shinto*. The blue square (left) is the front face of the card, i.e., the question. After clicking on it, the card flips and shows the answer (right).

Answer extraction We use T5 fine-tuned on the SQuADv1 dataset for answer extraction (Patil, 2020; Raffel et al., 2019). At inference time, for each paragraph we highlight one sentence at a time and feed it together with the rest of the paragraph to the model. The model extracts answers from the highlighted sentence leveraging the additional context information contained in the rest of the paragraph. To stay within the admitted input size of the model, we clip the length of the paragraphs to 512 tokens.

Question generation Here we use T5 fine-tuned on the SQuADv1 dataset for answer-aware question generation (Patil, 2020; Raffel et al., 2019). For each extracted answer, we append the corresponding paragraph as context and feed it to the model. Again, to not exceed the maximum input size we clip the input to a length of 512 tokens.

Filtering For filtering we use DistilBERT fine-tuned using a second step of knowledge distillation on the SQuADv1 dataset for question answering (Transformers, 2020; Sanh et al., 2020). Similar to the previous steps, we feed the model at inference time with each of the generated questions together with their corresponding paragraphs. We calculate an overlap score between the answer extracted in the answer extraction step and the answer produced by this question-answering model. The overlap score we calculate here is the ratio of identical bigrams over the total number of bigrams. Questions with an overlap score below 0.75 are discarded. Duplicates and questions whose answer is the title of the article are also discarded.

For each of the stages of our system, many different models exist in the literature. We selected each specific model based on their fitness to the task (i.e., models that are trained on Wikipedia based data-sets) as well as their availability as open source implementation. Once the cards are generated, they are displayed in a user-friendly way: the question is shown on the front side of the card, and after clicking on it, the card flips and shows on the other side the answer to the question. An example of a generated card as displayed in the application can be seen in Figure 2. To make our cards usable beyond our web application, we provide the option of exporting the generated cards as text file that can be im-

ported into Anki. Anki is a popular framework for flashcard-based learning with a large community of users that share their own flashcard decks as well as a number of commercial applications for smart-phones and web to help learning. This way, our generated flashcards are compatible with existing commercial applications and the user can choose the learning platform she prefers.

5 Evaluation

In this section we evaluate objective parameters of our flashcard generation pipeline, such as compute time or number of questions generated. Conversely, in Section 6 we evaluate the subjective quality of the generated cards through a user study. We divide our objective evaluation in two parts: 1) summarization step and 2) question generation and filtering step.

Summarization

Since we do not have reference summaries of the pieces of text that we are aiming to summarize, we cannot rely on the ROUGE score, which is the most common metric for summary quality. Instead, we calculate two values, similarity and error rate, that do not require a reference summary. The similarity score gives us a notion of how faithful the summary is to the original text, while the error rate quantifies the linguistic correctness of the summary.

To calculate the similarity score we use SentenceBERT (Reimers and Gurevych, 2019) to compute an embedding of each sentence in the original text and in the summary. Then, we calculate a context vector for the original text by adding up all of its sentence embeddings. We do the same for the summarized text. This results in two context vectors, one representing the original text and one representing the summary. The similarity score is the cosine similarity of these two vectors.

The error rate is the percentage of erroneous tokens. To calculate it, we determine the number of wrong tokens using *LanguageTool* (LanguageTool, 2020) and divide it by the total number of tokens. If a sentence has no end-of-sentence token, it is considered incomplete and an error is added to the count.

To determine which model to use in the summarization step, we compare two state-of-the-art models, T5 and BART. In Table 1 we compare both models in terms of similarity and error rate scores over the introduction of 256 Wikipedia articles. These articles were randomly selected based on the requirement that their introductions have more than 200 tokens. BART presents higher similarity score and lower error rate. This result is in line with the fact that BART obtains higher ROUGE score than T5 in summarization benchmarks such as CNN/DailyMail (Raffel et al., 2019; Lewis et al., 2020).

Next, since we are implementing our system as a web application, we need to consider computation time. To improve user experience we are interested in reducing as much as possible the time needed for the system to generate the cards. To this end we compare BART to its distilled version, i.e., DistilBART. We use the same set of 256 Wikipedia

Model	Similarity	Error Rate
BART	0.947	0.057
T5	0.912	0.129

Table 1: Comparison of T5 and BART summarization.

articles as in the previous experiment and for each model we calculate the average time it takes to summarize their introductions, as well as similarity and error rate. We run this experiment on a 24GB Nvidia Titan RTX GPU. The results in Table 2 show that DistilBART is 1.64 times faster, while it performs equally well in terms of similarity and error rate. While the absolute difference in computation time might seem small, we note that the computation time scales linearly with the article length, as articles are fed one paragraph at a time. We therefore choose DistilBART for the summarization step of our system, as the total speed up is significant.

Model	Time	Similarity	Error Rate
BART (large)	6.10 s	0.947	0.057
DistilBART	3.72 s	0.937	0.052

Table 2: Comparison of BART and DistilBART summarization.

Question Generation and Filtering

We study the performance of the question generation and filtering stage of our pipeline in terms of computing time and questions generated. We use the same 1024 randomly selected articles from Wikipedia as in the previous experiment and analyse the number of questions generated. In Table 3, we report the average number of flashcards generated and the average number of flashcards kept after the filtering stage.

		Time	All Qs	Qs after filter
Orig.	Per section	14.3 s	10.4	8.7
	Per article	240.5 s	178.4	148.2
Sum.	Per section	9.3 s	8.6	7.2
	Per article	151.2 s	144.0	120.5

Table 3: Average number of Questions (Qs) generated and kept after filtering for generating from original text (Orig.) and summaries (Sum.).

We see that even after applying our filtering step the number of questions kept, i.e., questions that meet a minimal quality requirement, is relatively large. In particular, generating 148.2 questions on average for a Wikipedia article implies that a student can access a significant amount of information from the cards. Furthermore, from the results we see that summarization helps in reducing the number of questions that are discarded.

Question	Answer
<i>What political party was George Orwell hostile to?</i>	Stalinism
<i>Animal Farm was a great commercial success when international relations were transformed as the wartime alliance gave way to what?</i>	Cold War
<i>When was Animal Farm first published?</i>	17 August 1945
<i>What group of people rebel against the human farmer in Animal Farm?</i>	Farm animals
<i>What magazine named Animal one of its 100 best English-language novels?</i>	Time

Table 4: Question-answer examples from the Wikipedia article on the novel *Animal Farm* by George Orwell.

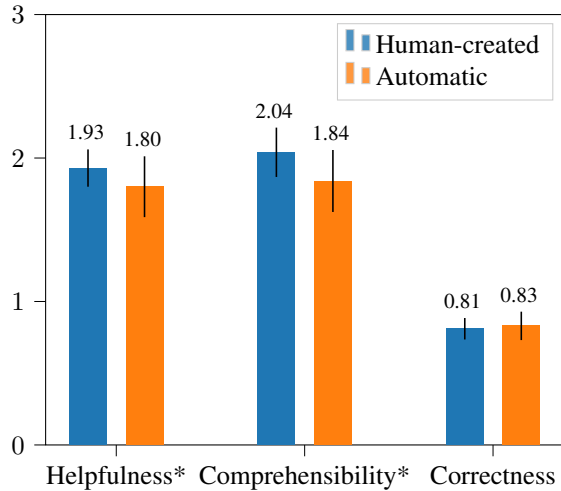


Figure 3: Results of the user study for the category history. Statistically significant differences ($p < 0.01$) are marked with an asterisk.

From the results presented in this section, we cannot assess the quality and usefulness of the generated cards, since this is a feature that depends on human perception. However, we can visually examine some examples of flashcards to have a notion of what kind of question-answer pairs our model generates. Table 4 shows the first five question-answer tuples generated for the article *Animal Farm* (novel by George Orwell) for the level “Highlight Introduction”. From the examples we see that generally, the generated cards are grammatically correct and contain meaningful information. However, to evaluate flashcard quality in a more rigorous manner, in the next section we conduct a user study.

6 User study

Given the strong perceptual component of flashcards, the best way of evaluating the quality of automatically generated cards is with a user study. In this study, we are interested in determining three aspects: usefulness for learning, linguistic comprehensibility and content correctness. In our user study, we ask about this three aspects and define a four-point scale for usefulness and comprehensibility (strongly disagree, disagree, agree and strongly agree), and a binary scale for content correctness (incorrect, correct/unknown). Table 5 dis-

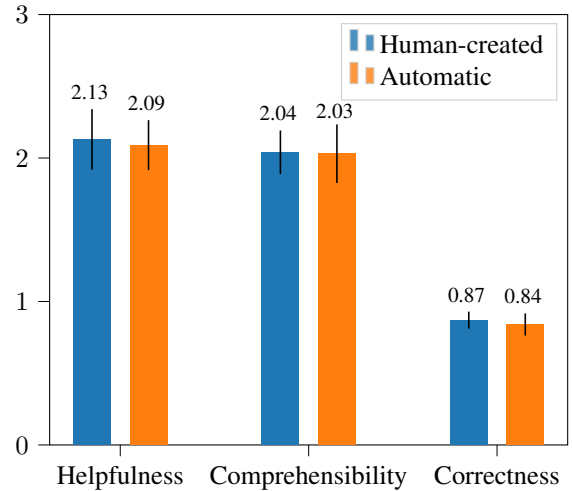


Figure 4: Results of the user study for the category geography. The difference between human-created and automatic cards is not statistically significant.

plays the detail of the questions asked in the study. During the study, the user is shown one card at a time and has to answer the three questions before the next card is displayed.

Question	Scale
1) <i>Is this card helpful for people who are studying this topic?</i>	0 – 3
2) <i>The text on the card makes sense to me.</i>	0 – 3
3) <i>Is the answer to this question correct?</i>	0 – 1

Table 5: Questions in the user study

The study consisted of 50 cards, from which 25 are generated by our automatic flashcard generation system, and the other 25 are created by humans. We obtain the human-created cards from flashcard decks that are freely available online in Anki format (AnkiWeb, 2020). At the beginning of the study, the user can choose between two topics, History and Geography; this gives the user the possibility of deciding the topic she is most familiar with, or interested in. We chose History and Geography as representative topics since they consist of factual knowledge, which is often

studied with flashcards. The human-created cards for History were taken from decks with titles: “Christianity” and “French Revolution”, while our cards were generated from the Wikipedia article “French Revolution” and the history section of the article “Germany”. For Geography, the topics were “India”, “Physical Geography” and “General Geography” for the human-created cards; and our cards were generated from the article “Atmosphere”, and the geography section of the articles “India” and “China”. For each category, we randomly chose 25 cards from the generated cards and mixed them with 25 randomly chosen cards from the human-created decks. The origin of the flashcards, i.e., whether they are automatically generated or human created, was not revealed to the participants.

50 participants from Amazon Mechanical Turks took part in the study. Data from participants which completed the study in less than 1,000 seconds was discarded. From the remaining, 21 participants selected the category history and 27 geography. Figures 3 and 4 show the results of our user study. The maximum score for helpfulness and comprehensibility is 3 and minimum is 0; and for correctness maximum is 1 and minimum is 0.

The results show that in the case of geography there is no statistically meaningful difference between human-created and our cards for either of the three aspects. For history, the difference for helpfulness and comprehensibility is statistically significant ($p < 0.01$), with human cards being marginally better than our cards. Neither category revealed a statistically significant difference in perceived correctness. Upon further investigation we found that the difference in the history category is mainly due to three automatically generated flashcards which are too ambiguous. We intend to improve our generation and filtering procedure in future work based on this insight.

Overall, this study demonstrates that the quality of our automatically generated cards is close to the quality of cards created manually by humans. This result validates our system and evidences its potential for enhancing personalized learning.

7 Discussion

In this work, we have presented a system for flashcard generation from raw text. Our system builds on recent advances in natural language processing, which have made available models for summarization, answer extraction, question generation and question answering. We thereby base our work on recent ideas on combining different models for question-answer generation and filtering. We have implemented our system as a web application that generates flashcards from Wikipedia articles with four different levels of detail. Our user study shows that the quality of the cards generated by our application is comparable, or only slightly worse, than human-created flashcards. Our work makes available a valuable tool for personalized education. By speeding up and automatizing flashcard generation, we give students the flexibility to decide which topics to learn, beyond standard curricula. Moreover, our work can be extended and combined with existing curricula by mapping course concepts to Wikipedia pages. A usage of knowledge graphs can also

be envisioned to link a user to adjacent topics for an automatically generated curriculum. We will explore these ideas in future work. We believe that in the near future tools and applications such as the one presented here will play a major role in enhancing autonomous and personalized learning. Although our application is already functional, there is still a lot of room for improvement and we plan to develop it further in order to improve computing efficiency and user experience.

References

- Alberti, C.; Andor, D.; Pitler, E.; Devlin, J.; and Collins, M. 2019. Synthetic QA corpora generation with roundtrip consistency. In Korhonen, A.; Traum, D. R.; and Màrquez, L., eds., *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, 6168–6173. Association for Computational Linguistics.
- AnkiWeb. 2020. Shared decks - ankiweb.
- Banerjee, S., and Lavie, A. 2005. METEOR: an automatic metric for MT evaluation with improved correlation with human judgments. In Goldstein, J.; Lavie, A.; Lin, C.; and Voss, C. R., eds., *Proceedings of the Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization@ACL 2005, Ann Arbor, Michigan, USA, June 29, 2005*, 65–72. Association for Computational Linguistics.
- Bao, H.; Dong, L.; Wei, F.; Wang, W.; Yang, N.; Liu, X.; Wang, Y.; Piao, S.; Gao, J.; Zhou, M.; and Hon, H. 2020. Unilmv2: Pseudo-masked language models for unified language model pre-training. *CoRR* abs/2002.12804.
- Blšák, M., and Rozinajová, V. 2016. Automatic question generation based on analysis of sentence structure. In Sojka, P.; Horák, A.; Kopeček, I.; and Pala, K., eds., *Text, Speech, and Dialogue*, 223–230. Cham: Springer International Publishing.
- Chan, Y., and Fan, Y. 2019. A recurrent bert-based model for question generation. In Fisch, A.; Talmor, A.; Jia, R.; Seo, M.; Choi, E.; and Chen, D., eds., *Proceedings of the 2nd Workshop on Machine Reading for Question Answering, MRQA@EMNLP 2019, Hong Kong, China, November 4, 2019*, 154–162. Association for Computational Linguistics.
- Chen, G.; Yang, J.; Hauff, C.; and Houben, G. 2018. Learningq: A large-scale dataset for educational question generation. In *Proceedings of the Twelfth International Conference on Web and Social Media, ICWSM 2018, Stanford, California, USA, June 25-28, 2018*, 481–490. AAAI Press.
- Dong, L.; Yang, N.; Wang, W.; Wei, F.; Liu, X.; Wang, Y.; Gao, J.; Zhou, M.; and Hon, H. 2019. Unified language model pre-training for natural language understanding and generation. In Wallach, H. M.; Larochelle, H.; Beygelzimer, A.; d’Alché-Buc, F.; Fox, E. B.; and

- Garnett, R., eds., *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, 8-14 December 2019, Vancouver, BC, Canada*, 13042–13054.
- Du, X., and Cardie, C. 2018. Harvesting paragraph-level question-answer pairs from Wikipedia. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1907–1917. Melbourne, Australia: Association for Computational Linguistics.
- Examiners, T. 2020. theexaminers.
- Gupta, S., and Gupta, S. 2019. Abstractive summarization: An overview of the state of the art. *Expert Systems with Applications* 121:49–65.
- Hermann, K. M.; Kocisky, T.; Grefenstette, E.; Espeholt, L.; Kay, W.; Suleyman, M.; and Blunsom, P. 2015. Teaching machines to read and comprehend. In *Advances in neural information processing systems*, 1693–1701.
- Kriangchaivech, K., and Wangperawong, A. 2019. Question generation by transformers.
- Kurdi, G.; Leo, J.; Parsia, B.; Sattler, U.; and Al-Emari, S. 2020. A systematic review of automatic question generation for educational purposes. *International Journal of Artificial Intelligence in Education* 30(1):121–204.
- Kwankajornkiet, C.; Suchato, A.; and Punyabukkana, P. 2016. Automatic multiple-choice question generation from thai text. In *2016 13th International Joint Conference on Computer Science and Software Engineering (JCSSE)*, 1–6.
- LanguageTool. 2020. LanguageTool.
- Lewis, M.; Liu, Y.; Goyal, N.; Ghazvininejad, M.; Mohamed, A.; Levy, O.; Stoyanov, V.; and Zettlemoyer, L. 2020. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Jurafsky, D.; Chai, J.; Schluter, N.; and Tetreault, J. R., eds., *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL 2020, Online, July 5-10, 2020*, 7871–7880. Association for Computational Linguistics.
- Lin, C.-Y. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, 74–81. Barcelona, Spain: Association for Computational Linguistics.
- Liu, M.; Rus, V.; and Liu, L. 2017. Automatic chinese factual question generation. *IEEE Transactions on Learning Technologies* 10(2):194–204.
- Lopez, L. E.; Cruz, D. K.; Cruz, J. C. B.; and Cheng, C. 2020. Transformer-based end-to-end question generation. *CoRR* abs/2005.01107.
- Narayan, S.; Cohen, S. B.; and Lapata, M. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. In Riloff, E.; Chiang, D.; Hockenmaier, J.; and Tsujii, J., eds., *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, 1797–1807. Association for Computational Linguistics.
- Niraula, N. B., and Rus, V. 2015. Judging the quality of automatically generated gap-fill question using active learning. In *Proceedings of the Tenth Workshop on Innovative Use of NLP for Building Educational Applications*, 196–206. Denver, Colorado: Association for Computational Linguistics.
- Pan, L.; Lei, W.; Chua, T.-S.; and Kan, M.-Y. 2019. Recent advances in neural question generation. *arXiv preprint arXiv:1905.08949*.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA*, 311–318. ACL.
- Patil, S. 2020. Question generation using transformers.
- Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.
- Rajpurkar, P.; Zhang, J.; Lopyrev, K.; and Liang, P. 2016. Squad: 100, 000+ questions for machine comprehension of text. In Su, J.; Carreras, X.; and Duh, K., eds., *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, 2383–2392. The Association for Computational Linguistics.
- Reimers, N., and Gurevych, I. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. In Inui, K.; Jiang, J.; Ng, V.; and Wan, X., eds., *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 3980–3990. Association for Computational Linguistics.
- Rush, A. M.; Chopra, S.; and Weston, J. 2015. A neural attention model for abstractive sentence summarization. In Márquez, L.; Callison-Burch, C.; Su, J.; Pighin, D.; and Marton, Y., eds., *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, 379–389. The Association for Computational Linguistics.
- Sanh, V.; Debut, L.; Chaumond, J.; and Wolf, T. 2020. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter.
- See, A.; Liu, P. J.; and Manning, C. D. 2017. Get to the point: Summarization with pointer-generator networks. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, 1073–1083. Vancouver, Canada: Association for Computational Linguistics.
- Shleifer, S. 2020. Distilbart model.

- Thalheimer, W. 2003. The learning benefits of questions. *Work Learning Research*.
- Transformers, H. 2020. Question answering using distilbert.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention is all you need. In Guyon, I.; von Luxburg, U.; Bengio, S.; Wallach, H. M.; Fergus, R.; Vishwanathan, S. V. N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, 5998–6008.
- Wang, A.; Pruksachatkun, Y.; Nangia, N.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. 2019a. Superglue: A stickier benchmark for general-purpose language understanding systems. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d'Alché-Buc, F.; Fox, E.; and Garnett, R., eds., *Advances in Neural Information Processing Systems*, volume 32, 3266–3280. Curran Associates, Inc.
- Wang, A.; Singh, A.; Michael, J.; Hill, F.; Levy, O.; and Bowman, S. R. 2019b. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net.
- Yan, Y.; Qi, W.; Gong, Y.; Liu, D.; Duan, N.; Chen, J.; Zhang, R.; and Zhou, M. 2020. Prophetnet: Predicting future n-gram for sequence-to-sequence pre-training. *CoRR* abs/2001.04063.
- Zhang, J.; Zhao, Y.; Saleh, M.; and Liu, P. J. 2019. PE-GASUS: pre-training with extracted gap-sentences for abstractive summarization. *CoRR* abs/1912.08777.
- Zhong, M.; Liu, P.; Wang, D.; Qiu, X.; and Huang, X. 2019. Searching for effective neural extractive summarization: What works and what's next. In Korhonen, A.; Traum, D. R.; and Màrquez, L., eds., *Proceedings of the 57th Conference of the Association for Computational Linguistics, ACL 2019, Florence, Italy, July 28- August 2, 2019, Volume 1: Long Papers*, 1049–1058. Association for Computational Linguistics.