

EXPLORING QUESTION ANSWERING AND GENERATION USING NEURAL MODELS AND SYNTACTIC TRANSFORMATIONS ON AN ENGLISH AND MULTILINGUAL CORPUS

Blessed Guda
Nebiyu Daniel Hailemariam
Oluwadara Adedeji
Chukwuemeka Malachi Ugwu

Abstract

This research explores the use of neural models for question answering (QA) and syntactic transformations and neural models for question generation (QG). For the QA system, the question type is first classified as either factoid or polar using Finetuned DistilBERT model trained on BookCorpus and English Wikipedia, which achieved an accuracy of 99.96%. For the factoid questions, the QA system uses a Finetuned BERT large model trained on BookCorpus and English Wikipedia and finetuned on the SQuAD dataset for question-answering. For polar questions, the context and question are passed to the Finetuned Roberta model to output a yes or no answer. The system achieves an f1 score of 91.3% and an Exact Match of 86.91% with the BERT large model for answering factoid questions. This research also explores the MT5 model on the MLQA dataset to answer questions from high-resource language (English) using content from low-resource languages. Furthermore, the research investigates the use of syntactic transformations and T5 neural models in the QG system. It is observed that the syntactic transformation system requires domain-specific linguistic knowledge to improve its performance. However, despite this limitation, the system can generate approximately 30 to 40% of the total generated questions fluently and answerably. To improve this, a T5 model trained on the SQuAD dataset was used for generating Wh-questions and a BoolQ dataset for polar questions. The manual evaluation shows that this question generation system produced over 95% of fluent and answerable questions from a pool of 20 questions generated.

Question Answering (QA) is a subfield of Natural Language Processing (NLP) that aims to build systems that can automatically answer questions posed by humans in natural language (Ostapov, 2011). QA systems can be classified into two categories: open-domain and closed-domain. Open-domain QA systems aim to answer any question that a user might ask, while closed-domain QA systems are designed to answer questions within a specific domain (Lende and Raghuvanshi, 2016).

Deep learning methods have been used extensively in recent years for QA tasks due to their ability to learn complex patterns in data and their ability to generalize well on unseen data (Sarker, 2021). Deep learning models such as Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformers have been used for QA tasks (Chaturvedi, 2018)(Hao et al., 2022)(Sarker, 2021). However, deep learning methods require large amounts of labeled data and computational resources (Hao et al., 2022)(Sarker, 2021).

Traditional methods such as Information Retrieval (IR) and Rule-based methods have also been used for QA tasks (Adnan and Akbar, 2019). IR-based methods retrieve relevant documents from a large corpus of text based on keyword matching or semantic similarity (Zhou et al., 2020). Rule-based methods use hand-crafted rules to extract answers from text (Rivas et al., 2019). Traditional methods are less computationally expensive than deep learning methods and require less labeled data (Adnan and Akbar, 2019). However, traditional methods may not perform

1 Introduction

1.1 Question Answering

well on complex questions or questions that require reasoning (Zhou et al., 2020).

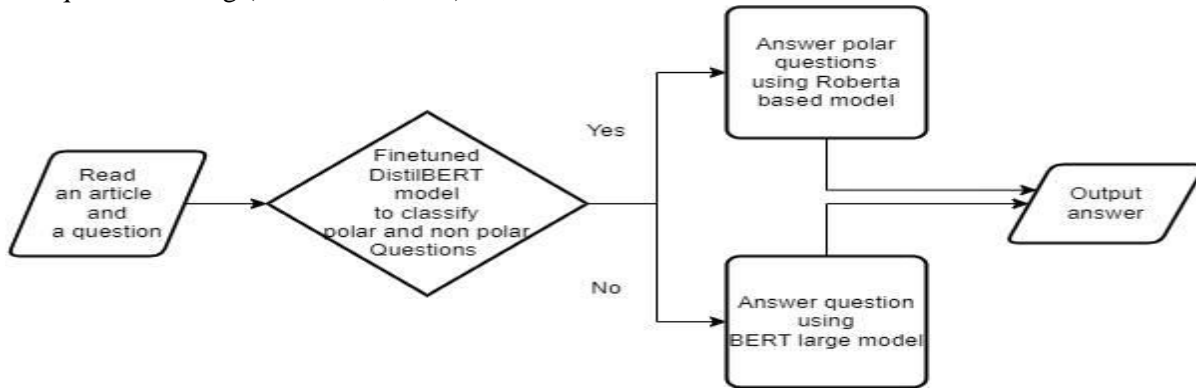


Figure 1: Architecture of Question Answering System

1.2 Question Generation

Question Generation (QG) is a field of study in natural language processing that focuses on developing systems that generate questions automatically using different sources such as unprocessed text, semantic representation, or databases (Pan et al., 2019). It is useful for educational purposes (Le et al., 2014) or to generate questions and answers for training and improving a Question Answering system (Mulla and Gharpure, 2023) and reduce the time needed for human labor to annotate question-answer datasets (Zhang et al., 2022). Different approaches have been explored to generate questions. These approaches to question generation can be template-based (Fabbri et al., 2020), syntactic transformations based or neural model-based (Mulla and Gharpure, 2023).

In template-based matching, a sentence similar to a context being considered is retrieved from the corpus, and a question is generated from this sentence (Fabbri et al., 2020). In syntactic transformation, sentences from a passage are transformed based on the constituency and dependency parse of the sentence (the syntax of the sentence) (Varga and Ha, 2010). Transformation rules are then used to generate questions based on the on the syntax of the sentence.

The other approach to question generation is the neural model approach which uses the encoder-decoder architecture to generate questions (Dwivedi and Marappan, 2009). This encoder-decoder model is trained on datasets such as SQuAD, which contain passages, questions and questions. Typically, the questions are concatenated with the corresponding passages and fed into the decoder. The encoding representation

of the question is used by the decoder in generating the answer. Some works have also extracted more features such as Style and Clue (Liu et al., 2020) to make the questions generated to be humanlike. This neural-based task of generating questions can be performed either from scratch or by fine-tuning previously pre-trained models. This is the contemporary way of question generation as it generates better questions which are harder to answer and not limited to the syntax sentences as in syntactic transformation (Du et al., 2017; however, they require a huge amount of training data and computational resources (Zohuri, 2020)).

2 System Architecture

2.1 Question Answering

Figure 1 depicts the architecture of the complete QA system. After reading the context and questions, each question is passed to the Finetuned DistilBERT model to classify if the question is polar or non-polar. The Finetuned DistilBERT model is a transformer model trained with BERT serving as a teacher to guide its training on BookCorpus and English Wikipedia, the same datasets the BERT base model was trained on (Sanh et al., 2019). It was trained on the None dataset (Question Classifier V2 (2023)) and has achieved an accuracy of 99.96%. If the question is classified as a polar question, it is passed to the Finetuned Roberta model - a model that enhances the BERT model using larger mini-batches and learning rates during training. This subsystem outputs a yes or no answer to the polar question.

On the other hand, if the Finetuned DistilBERT model classifies the question as non-polar, the question is passed to the BERT large model, which

is pre-trained on BookCorpus and English Wikipedia and finetuned on the SQuAD dataset for question-answering purposes (Devlin et al., 2019). The BERT large model is trained to read a given passage of text and provide an answer to a question along with its probability. The system applies a threshold score greater than 0.3 to the probability score to determine if an answer is correct. If the probability score is greater than or equal to 0.3, the system returns the answer as the final output. If the probability score is lower than 0.3, it returns an empty string as the answer. Returning an empty string means that the question is unanswerable given the input passage.

This research also explores multilingual QA systems. The research in (Asai et al., 2021) collected a multilingual dataset (XOR-TYDI QA) of questions from 7 languages. The questions were collected from native speakers of those languages rather than using other existing translation. This research investigates multilingual QA systems. The motivation for this is that when questions are asked of information not directly related to the culture of the speaker, they are less likely to have answers in the questioner's language (Information asymmetry) (Ewa S. Callahan, 2013). The approach in (Asai et al., 2021) used content from a high-resource language (English) to answer questions from a low-resource language.

However, research was not conducted to answer questions from a high-resource language using content from low-resource languages. This approach may help in answering questions in a language with high resources, such as English, using resources from low-resource languages. It may be helpful for questions with cultural inclinations tied to the low-resource language that may not be captured by the resources in English. To investigate this, the research explores multilingual QA systems using the MT5 model

and the MLQA dataset that contains English-to-Vietnamese, English-to-Arabic, and English-to-Spanish translation datasets. The aim is to answer questions from a high-resource language (English) using content from low-resource languages.

2.2 Question generation

The methodology explored in this project for question generation is syntactic transformations and neural models. In this section, we explain each methodology explored.

2.2.1 Syntactic Transformation

Figure 2 depicts the system architecture for the syntactic transformation process. The first step involves preprocessing the document, during which all pronouns are replaced with their references through coreference resolution. It is done using the Coreferee module within the spaCy Python library. Subsequently, the subject, verbs, and objects were extracted using the dependency tree technique in spaCy. It involved looping through all the tokens to obtain the root token. The root token is the verb on which other tokens depend.

After obtaining the root token, we iterate through all its children. Two lists of probable subject and object dependencies are used as a lookup for identifying if children of the root are subjects or objects. We loop through all the children of the root token. If the children are in the predefined subject dependencies list, they are marked as subjects. Likewise, if the children are in the predefined objects dependencies list, they are marked as objects. Finally, the subject, root token (verb), and object are returned.

The subjects and objects that are returned also have dependent tokens, which provide additional meaning to them. The subtrees are extracted and form the subject and object words or phrases. Subsequently, transformations were performed using the Named Entity Recognition (NER) tags.

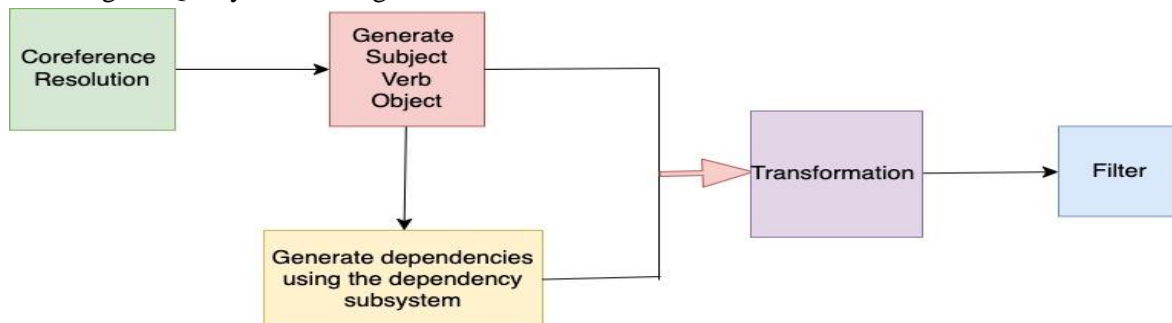


Figure 2: Syntactic transformation architecture.

For example, if the NER tag is DATE or TIME, the

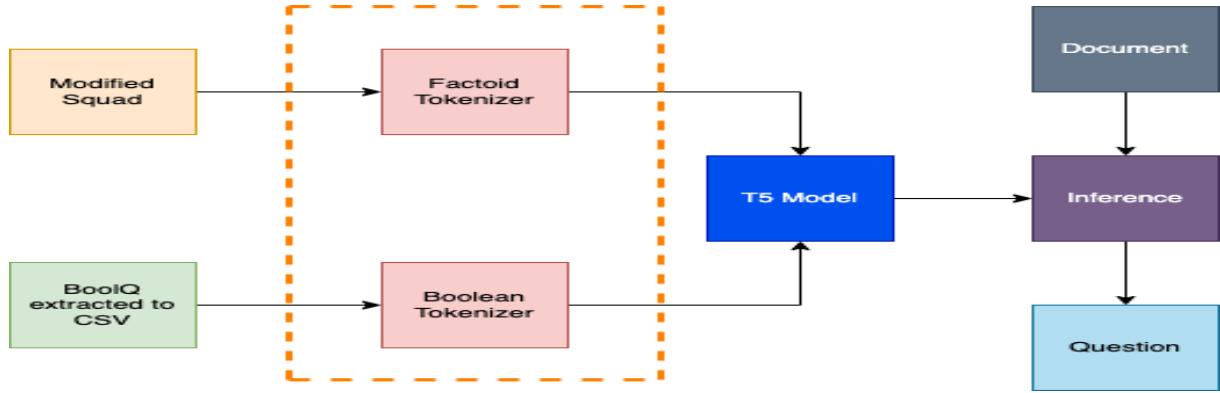


Figure 3: System diagram for the Neural model

question transformation will be “When did” + subject text + verb inflection. The verb inflection was applied. This was necessary to make the questions grammatical and fluent. However, the questions generated from the syntactic transformation were not good enough. A filtering technique was explored. However, due to the simplicity of the rules implemented the quality of the questions generated was low. Hence, a neural network approach was explored.

2.2.2 Neural network model

Figure 3 shows the system architecture for the neural model. The Text-to-Text Transfer Transformer (T5) is a model based on the encoder-decoder architecture (Raffel et al., 2020) and is suitable for text-to-text problems such as question generation. The T5 model was trained twice, one for generating Wh-questions and the other for polar questions. For the Wh-questions, the model was trained using the SQuAD dataset, while the polar QA were trained on the BoolQ dataset.

Each input question-passage pair was concatenated, tokenized, and terminated with an end-of-sequence token. Both models were trained and evaluated on the train and validation sets defined in the original datasets without any modifications. The factoid T5 model was trained for 5 epochs with a learning rate of 0.0001. Similarly, the polar T5 model was trained for 4 epochs, using a learning rate of 0.0001. After training, both models were saved and used for inference.

3. Experiments

3.1 Question Answering

The Text-to-Text Transfer Transformer (T5) is a model based on the encoder-decoder architecture. It was pre-trained using teacher forcing on the diverse C4 dataset, which has achieved state-of-

the-art results on benchmark NLP tasks, such as summarization, question-answering, and text classification (Raffel et al., 2020). This research evaluates the T5 large model both on question generation and question answering tasks.

Also, this research carried out experiments on developing multilingual QA systems. The traditional approaches to multilingual QA systems are highly dependent on the performance of the machine translation systems. This is because the query is first translated and then handled as a monolingual QA problem. This is susceptible to ambiguities arising from machine translation (Jiang et al., 2020). Therefore, we experimented with the MT5 model, which is a variation of the T5 model that supports multiple languages. It is pre-trained using the MC4 dataset, which contains text from 101 different languages.

BERT is a deep bi-directional Transformers-based language representation model. Unlike the T5 and mT5 model, which uses both the encoder and decoder components of the transformer model, BERT only uses the encoder component. It can be employed in different NLP tasks, such as question answering and language inference, using only one additional output layer without needing major modifications to the architecture for specific tasks.

The BERT (Devlin et al., 2019) Large uncased model is pre-trained on the BookCorpus (Zhu et al., 2015), and the English Wikipedia English Wikipedia (2023) and finetuned with the SQuAD dataset for question answering. During training BERT Large uncased model (Devlin et al., 2019), all texts are converted to lowercase, which is why the model is referred to as “uncased”. Subsequently, these case-insensitive texts are tokenized using WordPiece with a vocabulary size of 30,000. For a single input sample, the tokenized question is added after the beginning of the sentence token, after which a separation token is added to separate the sentence from the question. The model was

325 finetuned for 2 epochs for a learning rate of $3e-5$ on
326 the SQuAD dataset, which showed an f1 of 91.3%
327 and an Exact Match (EM) of 86.91%.

328 Similarly, the T5 (Ewa S. Callahan, 2013) Large
329 model was also finetuned on the SQuAD dataset
330 for 5 epochs, which yielded an F1 score of 81.32
331 and an EM score of 77.64%. Therefore, the
332 extractive approach using the BERT-uncased
333 model outperformed the generative approach on
334 the SQuAD dataset. Furthermore, during our
335 testing, we discovered that the T5 model performed
336 even more poorly on Wikipedia texts when
337 compared to the SQuAD dataset. It may be due to
338 the longer lengths of the contexts in the Wikipedia
339 corpus compared to the SQuAD data. Moreover,
340 when comparing the inference time of the two
341 models on the AWS EC2 g3sxlarge instance, the
342 BERT uncased takes 23.53% less time to answer a
343 question compared to the T5 model.

344 Both the T5 model and BERT Large uncased
345 model are not specifically trained to answer polar
346 questions, and as a result, they did not perform well
347 on polar questions, which require a simple "yes" or
348 "no" answer. For this, we used a separate model
349 for the boolean questions. To solve this, each
350 question was first passed to the Finetuned
351 DistilBERT model, trained on the BookCorpus and
352 English Wikipedia datasets to classify whether a
353 question is polar or non-polar. This model achieved
354 an accuracy of 99.96% on the None dataset. If the
355 question was classified as polar, it was then passed
356 to the Finetuned Roberta model, which enhances
357 the BERT model using larger mini-batches and
358 learning rates during training, to output a yes or no
359 answer. This subsystem effectively addressed the
360 limitations of the previous models in handling
361 polar questions, improving the overall performance
362 of the question generation system.

363 3.2 Question Generation

365 To generate questions, we first explored the
366 syntactic transformation of all the sentences in the
367 passage. Several experiments were conducted with
368 different transformations based on our limited
369 knowledge of English syntax. Both factoid and
370 polar questions were generated from the syntactic
371 transformation, as shown in Figure 4.

372 However, it was observed that some of the
373 questions were not grammatically fluent, with
374 punctuation errors and long, poorly constructed
375 questions that lacked clear meaning and
376 answerability. To generate more quality questions,

377 a simple filter was implemented, but it did not
378 significantly improve the quality of the questions
379 generated. The filter is constructed using a matcher
380 in spaCy Spacy API Documentation (2015). The
381 matcher uses a rule-based approach or pattern
382 defined by spaCy to filter sentences based on rules
383 and patterns passed. Although the questions
384 generated were not significantly better than the
385 ones generated without the filter, some
386 improvement was observed. Based on manual
387 testing, we realized that about 30-40% of the
388 questions generated were fluent and answerable.
389 Furthermore, after applying the filter, there were no
390 repeated questions generated. Subsequently, a T5
391 model trained on the SQuAD dataset and syntactic
392 transformation was used to generate factoid and
393 polar questions, respectively. However, questions
394 generated by the syntactic transformation approach
395 had numerous grammatical errors and were not
396 fluent.

397 This research also experimented with only
398 neural models to generate questions. The SQuAD
399 dataset was used to train the T5 model to generate
400 factoid questions, whereas BoolQ was used to train
401 the T5 model to generate polar questions. Initially,
402 the entire context was provided to the model while
403 performing inference. However, this approach
404 resulted in generating only a limited number of
405 questions. Despite adjusting the maximum length
406 parameter to tackle this issue, it did not lead to an
407 increase in the number of questions generated
408 during inference. To address this problem, the
409 context was divided into paragraphs and passed
410 separately to the model. This approach generated
411 more questions.

412 The performance of the model improved with a
413 higher number of epochs. However, due to time
414 constraints, we trained the factoid T5 model for the
415 or 5 epochs with a learning rate of 0.0001 and the
416 polar T5 model for 4 epochs using a learning rate
417 of 0.0001. Overall, the questions generated using
418 this technique were more fluent and answerable
419 than questions generated using the syntactic
420 transformation method we explored earlier. By
421 manually inspecting the generated questions, we
422 found that the system produced over 95% of fluent
423 and answerable questions from a pool of 20
424 questions generated. Some examples of the
425 questions generated are shown in Figure 5.

Did Messi win the 2005 FIFA World Youth Champi
 Did Messi relocate to Spain?
 Did Messi struggle with injury?
 Where did Messi relocate?
 Did Messi first uninterrupted campaign come in
 Did Messi personal best campaign to date am th
 Who did Messi regain?
 What did Messi won?

Figure 4: Snapshot of some sample questions generated by the syntactic transformation approach.

Did lionel messi play in the 2014 world cup?
 Did messi's father have growth hormone?
 Did messi ever win a cup with barcelona?
 When was Lionel Messi born?
 How many times has Messi won the FIFA Ballon
 What is the Guinness World Records for most
 When did Messi make his competitive debut?

Figure 5: Snapshot of some sample questions generated by the T5 neural model.

4. Discussion

4.1 Question Answering

The study compared the performance of two models, BERT-uncased and T5-Large, on the SQuAD dataset for extractive question-answering. The BERT-uncased model achieved higher F1 and EM scores than the T5-Large model, with scores of 91.3% and 86.91% compared to 81.32% and 77.64%, respectively. Additionally, the study found that the BERT-uncased model had a shorter inference time, taking 23.53% less time to answer a question compared to the T5-Large model. It is worth noting that the T5-Large model performed poorly on Wikipedia texts, which could be due to longer context lengths in the corpus.

Furthermore, the Multilingual Question Answering (MLQA) (Lewis et al., 2020) dataset was used to finetune the MT5 pre-trained model and investigate the performance and viability of this approach. In the subset of the dataset used, the question and answer are in English, while the context is in the low-resource language. The model is trained and evaluated on pairs of languages, between English and (Vietnamese, Arabic, and Spanish). Each model is trained for 5 epochs. A summary of the performance is summarized in Table 1. On each set, the model is trained on the multilingual MLQA and evaluated on the validation set.

Table 1: Result of training the MT5 model on the MLQA corpus

Translation	F1	EM	Size of Dataset
Eng-Vietnamese	22.50	6.50	6.01k
Eng-Arabic	25.14	14.51	5.85k
Eng-Spanish	39.41	26.75	5.75k
Eng-Vietnamese	22.50	6.50	6.01k

It can be observed, from Table 1, that the highest result was obtained when answering questions using Spanish contexts, while Vietnamese gave the lowest scores. The scores for these languages were low compared to the monolingual QA system. However, it shows that questions in English can be answered using passages from other languages, which is the opposite of the approach used in (Asai et al., 2021). These results show promising prospects considering the small training sets and few epochs used in training the models. However, this mode of question-answering should be further investigated to see how the model performs with more data and longer training times.

4.2 Question Generation

We could not improve the syntactic transformation because of insufficient linguistics knowledge. Furthermore, the team implemented a simple filter using matcher in spaCy and conditional filters to identify the Wh-and polar questions; however, this did not give the best results. High-precision transformation templates based on English lingual rules are required to improve the syntactic transformation system. In conjunction with the syntactic transformation approach, a T5 model, trained on the SQuAD dataset, was used to generate Wh-questions. The syntactic transformation complements the T5 model to generate polar questions. However, this approach had numerous grammatical errors and was not fluent. Therefore, a T5 model was trained on the BoolQ dataset to generate polar questions, and another T5 model was trained on the SQuAD dataset to generate factoid questions. This approach, which relied solely on T5 models, was more effective in producing accurate and fluent questions in comparison with the initial method.

Implementing this approach resulted in the desired outcome for factoid and polar questions.

5. Conclusion

This research explored neural models for QA and syntactic transformations and neural models for QG. The QA system leverages two different models for question answering based on the type of question. To classify whether a question is polar or non-polar, a Finetuned DistilBERT model trained on BookCorpus and English Wikipedia is used. If the question is polar, the system passes the question and context to a Finetuned Roberta model, which provides a yes or no answer. For non-polar questions, the BERT large model, which is pre-trained on BookCorpus and English Wikipedia and finetuned on the SQuAD dataset for question-answering, is utilized. We achieved an f1 score of 91.3% and an Exact Match (EM) of 86.91% with the BERT large model. An experiment was also conducted on the multilingual MT5 model using English-to-Vietnamese, English-to-Arabic, and English-to-Spanish translation datasets. The study showed promising results for answering questions in English using context in other languages. Although the performance was lower than that of monolingual QA systems, further investigation is needed to see the viability of this approach.

Moreover, the research attempted to use a syntactic transformation approach for generating factoid and polar questions. Challenges arose due to insufficient linguistics knowledge, leading to unsatisfactory results. A simple filter using matcher in spaCy and conditional filters were implemented to improve the quality of the questions, but it did not yield the desired outcome. A T5 model, trained on the SQuAD dataset for generating factoid questions, was incorporated into the syntactic transformation system to enhance the performance. However, this approach had numerous grammatical errors and lacked fluency. Two separate T5 models were trained: one on the BoolQ dataset for polar questions and another on the SQuAD dataset for factoid questions to address this problem. This approach, relying solely on T5 models, proved to be more effective in generating accurate and fluent questions compared to the initial method. It was manually evaluated and found to produce over 95% fluent and answerable questions out of a pool of 20 generated questions. The findings highlight the importance of linguistic

knowledge in improving syntactic transformations and demonstrate the effectiveness of T5 models trained on specialized datasets for generating different types of questions.

6. References

- Kiran Adnan and Rehan Akbar. 2019. Limitations of information extraction methods and techniques for heterogeneous unstructured big data. *International Journal of Engineering Business Management*, 11:1–23.
- Akari Asai, Jungo Kasai, Jonathan H. Clark, Kenton Lee, Eunsol Choi, and Hannaneh Hajishirzi. 2021. XOR QA: Cross-lingual Open-Retrieval Question Answering. *NAACL-HLT 2021 - 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Proceedings of the Conference*:547–564.
- Akshay Chaturvedi. 2018. CNN for Text-Based Multiple Choice Question Answering. *Computer Vision and Pattern Recognition Unit, Indian Statistical Institute*:272–277.
- Jacob Devlin, Ming Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. *NAACL HLT 2019 - 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies - Proceedings of the Conference*, 1(Mlm):4171–4186.
- Xinya Du, Junru Shao, and Claire Cardie. 2017. Learning to ask: Neural question generation for reading comprehension. *ACL 2017 - 55th Annual Meeting of the Association for Computational Linguistics, Proceedings of the Conference (Long Papers)*, 1:1342–1352.
- Ankit Dwivedi and Arunothia Marappan. 2009. A Comparative Study of Neural Question Generation Models.
- Susan C. Herring Ewa S. Callahan. 2013. Cultural Bias in Wikipedia Content on Famous Persons. *Journal of the American Society for Information Science and Technology*, 64(July):1852–1863.
- Alexander R. Fabbri, Patrick Ng, Zhiguo Wang, Ramesh Nallapati, and Bing Xiang. 2020. Template-based question generation from retrieved sentences for improved unsupervised question answering. *Proceedings of the Annual Meeting of the Association for Computational Linguistics*:4508–4513.
- Tianyong Hao, Xinxin Li, Yulan He, Fu Lee Wang, and Yingying Qu. 2022. Recent progress in leveraging deep learning methods for question answering.

612 *Neural Computing and Applications*, 34(4):2765–
613 2783

614 Zhuolin Jiang, Amro El-Jaroudi, William Hartmann,
615 Damianos Karakos, and Lingjun Zhao. 2020. Cross-
616 lingual Information Retrieval with BERT.

617 Nguyen-Thinh Le, Nhu-Phuong Nguyen, Kazuhisa
618 Seta, and Niels Pinkwart. 2014. Automatic question
619 generation for supporting argumentation. *Vietnam*
620 *Journal of Computer Science*, 1(2):117–127

621 Sweta P. Lende and M. M. Raghuwanshi. 2016.
622 Question answering system on education acts using
623 NLP techniques. *IEEE WCTFTR 2016 -*
624 *Proceedings of 2016 World Conference on*
625 *Futuristic Trends in Research and Innovation for*
626 *Social Welfare*:0–5

627 Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian
628 Riedel, and Holger Schwenk. 2020. MLQA:
629 Evaluating cross-lingual extractive question
630 answering. *Proceedings of the Annual Meeting of*
631 *the Association for Computational*
632 *Linguistics*:7315–7330.

633 Bang Liu, Haojie Wei, Di Niu, Haolan Chen, and
634 Yancheng He. 2020. Asking Questions the Human
635 Way: Scalable Question-Answer Generation from
636 Text Corpus. *The Web Conference 2020 -*
637 *Proceedings of the World Wide Web Conference,*
638 *WWW 2020*:2032–2043

639 Nikahat Mulla and Prachi Gharpure. 2023. Automatic
640 question generation: a review of methodologies,
641 datasets, evaluation metrics, and applications.
642 *Progress in Artificial Intelligence*, 12(1):1–32.

643 Yuriy Ostapov. 2011 Question Answering in a Natural
644 Language Understanding System Based on Object –
645 Oriented Semantics.

646 Liangming Pan, Wenqiang Lei, Tat-Seng Chua, and
647 Min-Yen Kan. 2019. Recent Advances in Neural
648 Question Generation. (3).

649 Colin Raffel, Noam Shazeer, Adam Roberts, Katherine
650 Lee, Sharan Narang, Michael Matena, Yanqi Zhou,
651 Wei Li, and Peter J. Liu. 2020. Exploring the limits
652 of transfer learning with a unified text-to-text
653 transformer. *Journal of Machine Learning*
654 *Research*, 21:1–67

655 Carol Rivas, Daria Tkacz, Laurence Antao, Emmanouil
656 Mentzakis, Margaret Gordon, Sydney Anstee, and
657 Richard Giordano. 2019. Automated analysis of
658 free-text comments and dashboard representations
659 in patient experience surveys: a multimethod co-
660 design study. *Health Services and Delivery*
661 *Research*, 7(23):1–160..

662 Victor Sanh, Lysandre Debut, Julien Chaumond, and
663 Thomas Wolf. 2019. DistilBERT, a distilled version
664 of BERT: smaller, faster, cheaper and lighter. :2–6.

665 Iqbal H. Sarker. 2021. Deep Learning: A
666 Comprehensive Overview on Techniques,
667 Taxonomy, Applications and Research Directions.
668 *SN Computer Science*, 2(6):1–20.

669 Andrea Varga and La Ha. 2010. A question generation
670 system for the qgstec 2010 task b. *Proceedings of*
671 *the Third Workshop on Question Generation*(June
672 2010):80–83.

673 Ruqing Zhang, Jiafeng Guo, Lu Chen, Yixing Fan, and
674 Xueqi Cheng. 2022. A Review on Question
675 Generation from Natural Language Text. *ACM*
676 *Transactions on Information Systems*, 40(1):1–43..

677 Ming Zhou, Nan Duan, Shujie Liu, and Heung Yeung
678 Shum. 2020. Progress in Neural NLP: Modeling,
679 Learning, and Reasoning. *Engineering*, 6(3):275–
680 290.

681 Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan
682 Salakhutdinov, Raquel Urtasun, Antonio Torralba,
683 and Sanja Fidler. 2015. Aligning books and movies:
684 Towards story-like visual explanations by watching
685 movies and reading books. *Proceedings of the IEEE*
686 *International Conference on Computer Vision, 2015*
687 *International Conference on Computer Vision,*
688 *ICCV 2015*:19–27.

689 English Wikipedia (2023) Wikipedia.
690 Wikimedia Foundation. Available at:
691 https://en.wikipedia.org/wiki/English_Wikipedia
692 (Accessed: May 2, 2023).

693 Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan
694 Salakhutdinov, Raquel Urtasun, Antonio Torralba,
695 and Sanja Fidler. 2015. Aligning books and movies:
696 Towards story-like visual explanations by watching
697 movies and reading books. *Proceedings of the IEEE*
698 *International Conference on Computer Vision, 2015*
699 *International Conference on Computer Vision,*
700 *ICCV 2015*:19–27.

701 Spacy API Documentation (2015) *Matcher*.
702 Available at: <https://spacy.io/api/matcher> (Accessed:
703 May 2, 2023).

704 *Question Classifier V2* (2023)
705 *alangpp255/Question_classifier_V2 · Hugging*
706 *Face*. Available at:
707 https://huggingface.co/alangpp255/Question_classifier_V2 (Accessed: May 2, 2023).