# Building Efficient and Effective OpenQA Systems for Low-Resource Languages

Emrah Budur[a,b,*,1], Rıza Özçelik[c], Dilara Soylu[d], Omar Khattab[d], Tunga Güngör[a] and Christopher Potts[d]

[a]*Boğaziçi University, Bebek, Istanbul, 34342, Turkey*

[b]*Amazon, Toronto, ON, Canada*

[c]*Stanford University, Stanford, 94305, CA, USA*

[d]*Eindhoven University of Technology, Eindhoven, 5612 AZ, The Netherlands*

## ABSTRACT

Question answering (QA) is the task of answering questions posed in natural language with free-form natural language answers extracted from a given passage. In the OpenQA variant, only a question text is given, and the system must retrieve relevant passages from an unstructured knowledge source and use them to provide answers, which is the case in the mainstream QA systems on the Web. QA systems currently are mostly limited to the English language due to the lack of large-scale labeled QA datasets in non-English languages. In this paper, we show that effective, low-cost OpenQA systems can be developed for low-resource languages. The key ingredients are (1) weak supervision using machine-translated labeled datasets and (2) a relevant unstructured knowledge source in the target language. Furthermore, we show that only a few hundred gold assessment examples are needed to reliably evaluate these systems. We apply our method to Turkish as a challenging case study, since English and Turkish are typologically very distinct. We present SQuAD-TR, a machine translation of SQuAD2.0, and we build our OpenQA system by adapting ColBERT-QA for Turkish. We obtain a performance improvement of 9-34% in the EM score and 13-33% in the F1 score compared to the BM25-based and DPR-based baseline QA reader models by using two versions of Wikipedia dumps spanning two years. Our results show that SQuAD-TR makes OpenQA feasible for Turkish, which we hope encourages researchers to build OpenQA systems in other low-resource languages. We make all the code, models, and the dataset publicly available.

## 1. Introduction

Question answering (QA) is the task of answering questions posed in natural language with free-form natural language answers. In its standard formulation, QA is posed in a highly constrained way. The system is given a passage and a question with a guarantee that the answer can be found in the passage [50, 32, 27]. The main component of standard QA systems is a *reader*, which takes a passage and a question as input and returns an answer. Present day systems are extremely successful at such tasks, often surpassing human performance [9]. However, they are of limited use, since real-world question answering scenarios mostly do not involve gold passages or provide answerability guarantees.

This observation has motivated a move towards Open Domain Question Answering (OpenQA), where only the question text is given as input without any passage. Related passages are retrieved from a large corpus by a *retriever* and then used by the reader to predict an answer. In this setting, there is no guarantee that the retrieved passages will contain the answer, and the success of the system thus depends on having a successful retriever module to provide appropriate passages to the reader.

Recent years have seen rapid improvements in such systems stemming from the use of neural retriever modules that can provide semantically rich representations of documents. We are approaching the point where OpenQA systems will be as effective as standard QA systems [30].

However, this rapid progress in both standard QA and OpenQA systems is largely confined to English. Progress in other languages is constrained by a scarcity of gold data. While there are some high-quality multilingual resources in this domain [36, 4], the amount and diversity of such data remain low. The cost of creating new datasets is the main obstacle to progress in this area.

In this paper, we address the following research question: *Can we build cost-effective QA systems for low-resource languages without having a gold training dataset?* As a positive answer to this question, we propose a cost-effective approach to remedying the data scarcity problem for the QA task in non-English languages. Our proposal extends previous work on standard QA in languages other than English [44, 20, 1], and we argue that adopting the OpenQA formulation of the problem is a key step. For OpenQA, we require only gold question–answer pairs, and only for assessment. In particular, passages need not be a component of the gold data, since they are retrieved by the system to use as (perhaps noisy) evidence. Our formulation still requires training data, but this can be created by automatic translation from English datasets. These translations may contain mistakes, but we show that they can still lead to robust QA systems. Whereas

---

*Corresponding author

✉ emrah.budur@boun.edu.tr (E. Budur); r.ozcelik@tue.nl (R. Özçelik); soylu@stanford.edu (D. Soylu); okhattab@stanford.edu (O. Khattab); gungort@boun.edu.tr (T. Güngör); cgpotts@stanford.edu (C. Potts)

[1]Work done outside of Amazon.

the cost of creating a dataset like SQuAD [50, 49] can be upwards of US$50,000, our costs are only around US$500, most of which is for machine translation services. The cost of creating a gold assessment set could in principle be very large, but we show that one can get robust assessments of OpenQA systems with only around 200 question–answer pairs. Such gold datasets can be created by a small team very quickly.

As a case study, we explore OpenQA for Turkish, beginning from English resources. We expect English/Turkish to be a challenging case for our approach, since the two languages are typologically diverse. This case can also serve as a representation of challenges found between other languages and English. In §3, we introduce SQuAD-TR, an automatic translation of SQuAD2.0 [49] into Turkish using Amazon Translate.[2] SQuAD-TR serves as our noisy training data. For gold assessments, we rely on the XQuAD dataset [4] that provides a human-translated version of the SQuAD1.0 [50] dev set into Turkish (and nine other languages). Our core OpenQA system is a variant of the ColBERT-QA model of Khattab et al. [30], but it is built entirely over Turkish resources and is trained on SQuAD-TR. We show that the proposed methodology and resources can lead to highly effective standard QA systems (§4) and OpenQA systems (§5) for Turkish, and using only a few hundreds gold question–answer pairs is sufficient for the robust evaluation of QA systems. We make our code, models, and the dataset publicly available.[3]

## 2. Related Work and Background

### 2.1. English Question Answering Datasets

There are many QA datasets for English used to address different challenges; see Cambazoglu et al. [10] for a thorough review. One class of QA datasets consists of multiple-choice questions. MCTest [51] is an early dataset built in this style (see also CBT [24]; Booktest [6]). MCTest contains 2640 human-generated questions associated with a correct answer from a set of candidate answers. The questions and answers are based on 660 short fictional stories at a grade-school level. The fictional nature of the stories limits the use of world knowledge to answer the questions, which is one of the special challenges of this dataset. The main drawback of MCTest is its small size.

SQuAD1.1 [50] was the first major extractive reading comprehension dataset. SQuAD1.1 contains over 100K examples, and each example is a question–passage–answer triple, where annotators selected a span of text from the passage as the answer to the question. SQuAD2.0 [49] is a follow-up that includes over 50K additional examples representing unanswerable questions. The goal here is to encourage the development of systems that detect whether a question is answerable based on the passage given and abstain from answering if necessary [28]. Although we did not use unanswerable questions in our experiments and they

are out of the scope of this paper, we built SQuAD-TR from SQuAD2.0 to facilitate future research on unanswerable questions in Turkish.

HotPotQA [60] extends the extractive reading comprehension paradigm to multi-hop questions, i.e., questions whose answers need to be pieced together from information in multiple passages. A closely related task is multi-hop claim verification, as in HoVer [25].

Another class of datasets leverages an existing set of human-generated question–answer pairs, and augments these with supporting passages from external knowledge sources. A prominent example of this type of dataset is TriviaQA [27], which contains 95K question–answer pairs that were prepared by trivia enthusiasts. The question–answer pairs are accompanied by documents retrieved from the Web and Wikipedia. In a similar manner, Dunn et al. [21] built SearchQA by using the Google search engine to retrieve context snippets relevant for question–answer pairs obtained from the Jeopardy! game show archive.[4]

Search engine query logs are also used as a source of examples. WikiQA [59] and Natural Questions [32] are the most commonly used datasets in this class. WikiQA is derived from the 3K most frequent user queries in the query logs of the Bing search engine. Each query is paired with a Wikipedia page clicked by at least five unique users. If a sentence in the summary part of the associated Wikipedia page includes the answer to the query, the sentence is labeled as *correct*, otherwise *incorrect*. This version of the QA task is referred as *answer-sentence selection*, as it only selects the target sentence answering the question without requiring extraction of the correct answer span from that sentence. WikiQA includes question–page pairs with no correct sentences, so the dataset can also be used to build *answer triggering* models, which predict whether the sentences include the answer or not and then select a sentence only if it answers the question.

Like WikiQA, Natural Questions (NQ) relies on queries to a real search engine. NQ contains a total of 320K examples with queries obtained from Google query logs. Each query is associated with a Wikipedia page, which may or may not contain the answer for the query. If the Wikipedia page has the answer, *a long answer* is included in the example to show the passage answering the question. The example may also contain *a short answer* denoting the short form of the target answer. If the example contains neither long nor short answer, then no answer span exists on the page. NQ is a challenging dataset with realistic queries supported by high-quality annotations for the long and short answers. Likewise WikiQA, NQ also provides an opportunity to build answer triggering models with its examples having no long and short answers.

All of these datasets can be re-cast in the OpenQA mould, assuming we can find a large collection of relevant unlabeled documents to be used as a knowledge source. SQuAD1.0, NQ, TriviaQA, HotPotQA have been extensively explored in these terms [30, 34].

---

[2]https://aws.amazon.com/tr/translate
[3]https://github.com/boun-tabi/SQuAD-TR

[4]http://j-archive.com

| Dataset | Language | Number of examples |
|---|---|---|
| KorQuAD [39] | Korean | 70,079 |
| FQuAD [20] | French | 62,003 |
| SberQuAD [22] | Russian | 50,364 |
| CMRC 2018 [18] | Chinese | 19,071 |
| GermanQuAD [43] | German | 13,722 |
| ARCD [44] | Arabic | 1,395 |

**Table 1**
QA datasets for non-English languages.

## 2.2. Multilingual Question Answering

Various methods have been used to address the dataset bottleneck for QA in non-English languages [12]. One approach is to curate in-language datasets from scratch. A number of datasets for different languages have been created in this way. We provide a summary in Table 1. Datasets created in this manner are likely to be of high quality, but they are expensive and time-consuming to create.

One way to reduce the cost of creating in-language QA systems is to try to rely on the zero-shot transfer capabilities of cross-lingual models. In this approach, multilingual language models (e.g., mBERT [19]; XLM [16]; XLM-RoBERTa [15]) are finetuned on English QA datasets and then used to answer questions in target non-English languages. Multilingual models are efficient and cost effective, especially in large-scale applications requiring multiple language support. However, their performance on the target languages is relatively lower than models that use in-language embeddings [36, 20, 39, 43, 8].

Another cost-effective approach is to rely on machine translation services. In this approach, in-language training datasets are automatically obtained by translating an existing English dataset using machine translation (MT). Previously, SQuAD1.1 [50] was translated into Arabic [44], French [20], and Spanish [11], and SQuAD2.0 [49] was translated into Persian [1]. Similar techniques have also been used in other areas of NLP [8, 57].

Using MT systems is undoubtedly productive, but relying on automatic translations for system assessment raises concerns about the validity of those assessments. To the extent that there are systematic errors in the MT output, assessment numbers are likely to be untrustworthy. To address this, Lewis, Oguz, et al. [36] proposed MLQA, which is a multi-way aligned QA dataset to be used for evaluation purposes in 7 languages, with over 5K examples for each language. Similarly, Artetxe et al. [4] developed XQuAD, which consists of a subset of the SQuAD1.0 development dataset with human translations into 10 languages, including Turkish. In what follows, we rely on the Turkish portion of XQuAD, namely XQuAD-TR, for evaluation as it is the only standard QA evaluation dataset that supports Turkish.

## 2.3. Open Domain Question Answering

Various methods are employed by researchers to develop OpenQA systems. A thorough investigation of these methods can be found in the comprehensive survey presented by Zhu et al. [62]. Traditionally, OpenQA systems involve two pipelined components: a *retriever* and a *reader*. Given a question, the retriever is expected to retrieve candidate passages, and the reader is supposed to extract the target answer span from those retrieved passages.

BM25 was a common choice for the retriever component in the earliest OpenQA systems [52] and it remains in wide use today [46, 38, 48, 47]. BM25 and other retrievers in its class rely on lexical matching. The guiding idea behind more recent neural retrievers is that lexical matching alone is not sufficiently semantic in nature to capture the nuanced ways in which passages can be relevant to user queries. Prominent recent examples of these neural retrievers include ORQA [34], REALM [23], DPR [29], RAG [37], and ColBERT [31]. The leaderboards for OpenQA systems are currently dominated by systems that employ neural retrievers, though BM25 remains a very strong baseline, especially where latency and cost are important additional considerations beyond accuracy-style metrics.

Neural retrievers, despite their advantages, suffer from the drawback of requiring a significant amount of storage space, especially for their indexes. To mitigate this limitation, S. Yang & Seo [58] proposed a solution that involves several techniques. These techniques include filtering out unnecessary passages prior to the retrieval step, consolidating retriever-reader models with a single encoder, and employing post-training compression.

The English-centric nature of research in this area is arguably holding back retriever development as well. The largest and most widely used dataset in this space is the MS MARCO Passage Ranking dataset [45, 17], and it contains only English texts and queries. However, Bonifacio et al. [7] translated MS MARCO into 13 different languages using automatic translation. The result is mMARCO [7], the first multilingual MS MARCO variant. mMARCO has enabled much new research on multilingual passage retrieval. However, mMARCO does not have any labels within the text to denote answer spans, and so it cannot by itself support the development of multilingual QA systems.

Neural information retrieval (IR) systems can begin from pretrained multilingual embeddings, and this can facilitate multilingual retrieval work. For example, Asai et al. [5] used DPR in the retriever step and proposed a cross-lingual transfer method (XOR QA) to obtain answers for unanswerable questions in the non-English languages from English Wikipedia. In order to do that, (1) they translate the questions

in non-English languages into English, (2) find relevant passages and answer spans from English Wikipedia, and (3) translate the English answer spans back to the original language. The leaderboard for their paper, XOR-TyDi [5], includes cross-lingual retrieval and OpenQA tasks.[5] XOR-TyDi has similar motivations to our work, in that it tackles issues around building OpenQA systems in non-English languages effectively, but it differs from our work in substantive ways. To achieve its goals, XOR-TyDi makes the English knowledge source available to non-English languages with the help of a cross-lingual retriever. In contrast, we propose a method to build an in-language retriever that benefits from an existing in-language knowledge source in a non-English language. Both methods are based on translation, but our method benefits from translated data at training time (TRANSLATE-TRAIN) even if the translation is noisy, whereas XOR-TyDi requires translations at test time (TRANSLATE-TEST), making the overall system highly vulnerable to translation errors.

For the most part, the reader component in OpenQA systems is an extractive reader: given a retrieved passage and a question, it is trained to extract a substring of the passage corresponding to the answer. Readers of this sort are clearly best aligned with standard QA datasets where the answer is guaranteed to be a substring of the passage provided. In datasets where the answer can be expressed more indirectly, extractive strategies will fail. Extractive readers are also potentially sub-optimal for OpenQA systems, for two reasons: we might be able to retrieve multiple relevant passages, and the passages themselves might indirectly express the answer. In response to these shortcomings, Lewis, Perez, et al. [37] explore readers that can consume multiple passages and generate original texts in response. Yu et al. [61] introduces KG-FiD, which incorporates knowledge graphs to rerank passages by utilizing Graph Neural Networks (GNN) before the reader generates the response. We leave exploration of generative readers for Turkish for future work.

Several end-to-end neural models have recently emerged in OpenQA (e.g., SOQAL [44]; DPR [29]; ColBERT-QA [30]; YONO [33]). Early examples predominantly relied on sparse vector representations in the retrieval component. For instance, Mozannar et al. [44] proposed SOQAL as an OpenQA system for Arabic using a hierarchical TF-IDF retriever pipelined with a BERT-based reader [19]. This was followed by an answer ranking component that assigns a score for each answer candidate obtained as a linear combination of the retriever and reader outputs. However, these retrievers, based on sparse representations, struggle to recognize similarity between synonyms and paraphrases that use different lexical terms.

To address the sparse vector representation problem, Karpukhin et al. [29] introduced DPR which is one of the early examples of dense retrievers in the OpenQA domain. DPR utilized dual-encoder architecture to encode dense and latent semantic representations of the questions and the contexts. Given a question, DPR was trained to distinguish the positive passages from the negative passages in the batch. One limitation of DPR was its use of a single vector for the question and the context, resulting in limited interactions between the terms in the two texts. Khattab et al. [30] recently developed ColBERT-QA as a novel end-to-end neural OpenQA model for English, offering more extensive and effective interaction between the question and context terms through a late-interaction mechanism. Alternatively, H. Lee et al. [33] proposed YONO, a single end-to-end architecture that jointly optimizes the retriever, reranking, and reader components. The fully end-to-end architecture of YONO contributes to its efficiency in terms of model size. However, there is a drawback to combining multiple components in a single architecture, as each component demonstrates different overfitting characteristics. This vulnerability becomes apparent especially when the training data is limited, which is often the case for low-resource languages.

In this paper, we focus on two advanced end-to-end neural models used in OpenQA, the DPR and ColBERT-QA models, for their ability to provide dense representations of queries and passages. Each model is explained briefly in the following sections.

## 2.4. DPR

DPR [29] employs a BERT-based dual-encoder architecture for the retriever component within an end-to-end OpenQA system. This architecture encodes the question and the passages through separate encoders, allowing it to provide relevant candidate passages based on their dense representations in response to the given question. A crucial aspect of DPR involves learning the similarity between questions and passages by employing in-batch negatives. These in-batch negatives are composed of relevant passages from other questions within the same training batch. Each training example in DPR benefits from hard negative passages, which are the top k negative results returned by BM25. The primary strength of utilizing in-batch negatives lies in expanding the number of training examples effectively while keeping memory footprint minimal. However, despite its capacity to encode questions and passages separately and index passages offline, its retrieval performance is considered to be sub-optimal due to the lack of interaction between the question and passage encoders.

## 2.5. ColBERT-QA

The ColBERT-QA system of Khattab et al. [30] is an OpenQA system built on top of the ColBERT retriever model [31]. In ColBERT-QA, the retriever is iteratively fine-tuned using weak supervision from the QA dataset so that it can perform task-specific retrieval. ColBERT-QA standardly uses an extractive reader, though its fine-tuned retriever is compatible with a wide range of reader designs.

The hallmark of the ColBERT model is its *late interaction* mechanism: both queries and passages are separately encoded into sequences of token-level vectors corresponding roughly to the output states of a BERT encoder [19]. Given a query $q$ encoded as a sequence of token-level vector representations $[q_1, \ldots, q_m]$ and a passage $p$ encoded as

---

$[p_1, \ldots p_n]$, ColBERT computes the similarity of every pair of vectors $q_i$ and $p_j$, and sums the scores only for the highest scoring $p_k$ for each $q_i$ ("MaxSim"). This is the basis for scoring documents with respect to queries. The architecture allows all passages in the knowledge source to be encoded off-line and indexed for fast comparisons with query representations. As a pure retriever, ColBERT achieves state-of-the-art results across a wide variety of IR benchmarks [54] and it can be implemented in a low-latency, space-efficient manner [53].

ColBERT-QA is a powerful example of recent general-purpose approaches to OpenQA and we base our models on this architecture. To adapt the model to Turkish, we made only language-specific adjustments (§5).

## 3. Datasets

### 3.1. SQuAD-TR

Inspired by previous work using machine translation as a stepping stone to obtain multilingual resources (§2.2), we translated SQuAD2.0 [49] to Turkish using Amazon Translate.[6] We translated the titles, context paragraphs, questions, and answer spans in the original dataset. As a natural consequence, we needed to remap the starting positions of the answer spans, since their positions were not maintained in the translated paragraphs. This is needed not only due to linguistic variation between the source and target languages [20] but also because the translation task is inherently context dependent [44]. A text span may have totally different translations depending on its context. This is a challenging issue for obtaining consistent translations, particularly for Turkish due to the context-dependent morphological variation of Turkish words, as exemplified in Table 2. The problem with mapping all of the answer spans after translation is that it requires a substantial amount of time and manual work. However, it is still possible to recover part of them automatically, so we mapped the answer spans automatically in the target translations, as in much related work in different languages [44, 20, 11, 1].

In this automatic post-processing step, we first looked for spans of text in the context paragraph that exactly matched the answer text. If we found such a span, we kept that answer text along with its starting position in the translated text following the previous works [44, 20, 11, 1]. For answer texts without matching spans, we searched for the spans of text that approximately matched with the target answer text using character-level edit distance [35].[7] We use different edit distance values based on the length of the answer text. For answer texts with lengths shorter than 4 characters, we try to match all spans that are 1-edit distance away from the answer text. For all other answer texts, we search for all spans that are *up to* 3-edit distance away from the answer text and

select all of the longest spans of texts that approximately match the target answer text. Table 2 shows examples of the answer spans that are recovered as a result of this post-processing.

This approximate matching is generally successful. However, for 25,528 question–answer pairs in SQuAD-TR-TRAIN, neither exact nor approximate matching returns a span in the translated paragraph. We removed these question–answer pairs from SQuAD-TR-TRAIN. This resulted in 259 paragraphs having no question–answer pairs. We removed those paragraphs from SQuAD-TR-TRAIN as well. Similarly, we removed 3,582 question–answer pairs from the SQuAD-TR-DEV dataset, but we did not need to remove any paragraphs from SQuAD-TR-DEV, as all paragraphs had at least one question–answer pair where the answer text has a matching span in the paragraph.

With this procedure, we obtained the training and evaluation splits of SQuAD-TR, namely SQuAD-TR-TRAIN and SQuAD-TR-DEV, respectively. We used SQuAD-TR-TRAIN as a training dataset but did not use SQuAD-TR-DEV for evaluation in our research. We share it for future work. For evaluation, we instead used the Turkish split of XQuAD [4], namely XQuAD-TR, which helped maximize the validity of our assessment results, since it is a high-quality, human-translated test set.

Table 3 provides basic statistics of SQuAD-TR and XQuAD-TR along with the training and dev splits of the original SQuAD2.0 dataset (SQuAD-EN), noted as SQuAD-EN-TRAIN and SQuAD-EN-DEV, respectively. The number of articles is identical for the SQuAD-EN and SQuAD-TR datasets, whereas the SQuAD-TR-TRAIN dataset has fewer paragraphs and answerable questions than SQuAD-EN-TRAIN due to the excluded paragraphs and questions. Similarly, the SQuAD-TR-DEV dataset has fewer answerable questions than SQuAD-EN-DEV, for the same reason. As a matter of course, the number of unanswerable questions did not change in any split of the SQuAD-TR dataset, as the original unanswerable questions remain unanswerable after translation. We release SQuAD-TR publicly.[8]

### 3.2. Knowledge Source

As we discussed above, in OpenQA, evidence passages are not given to the reader along with the questions, but rather are retrieved from a large corpus. Thus, we first need to prepare a knowledge source containing the passages to be retrieved. We used the Turkish Wikipedia as the main part of our knowledge source. We obtained the passages in our knowledge base by extracting *contents* and *titles* from Turkish Wikipedia articles. However, we observed that the majority of the target information available in SQuAD2.0 [49] was not actually available in Turkish Wikipedia due to two main issues.

One of these issues occurs when the article containing the target information in English Wikipedia is actually missing in Turkish Wikipedia. As an example, SQuAD2.0 has 50 question–answer pairs targeting 25 paragraphs about

---

[6]Amazon Translate was chosen thanks to the availability of AWS Cloud Credits for Research Grant for the authors, but it is possible to use other effective machine translation systems as well. Please refer to the disclaimer mentioned in the acknowledgements section for further information.

[7]We used the implementation in the Python regex package: https://pypi.org/project/regex/2021.4.4

[8]https://github.com/boun-tabi/SQuAD-TR

| | Language | Context span | Question | Answer Text (Before post processing) | Answer Text (After post processing) |
|---|---|---|---|---|---|
| **Example 1** *(Edit distance=1)* | Turkish | ...Görünüşü, o yılki MTV Video Müzik Ödülleri'nin MTV tarihinde en çok izlenen yayın haline gelmesine ve **12.4 milyon** izleyiciyi çekmesine yardımcı oldu;... | 2011 MTV Müzik Ödülleri'ni kaç kişi izledi? | 12,4 milyon *(12.4 million)* | 12,4 milyon *(12.4 million)* |
| | English | ...Her appearance helped that year's MTV Video Music Awards become the most-watched broadcast in MTV history, pulling in **12.4 million** viewers;... | How many people watched the 2011 MTV Music Awards? | 12.4 million | — |
| **Example 2** *(Edit distance=2)* | Turkish | ...Kariyerindeki en uzun süreli Hot 100 single'ı olma başarısına ulaşan "Halo"un ABD'deki başarısı, Beyoncé'nin **2000'li** yıllarda diğer kadınlardan daha fazla listede ilk on single elde etmesine yardımcı oldu.... | Hangi on yıl boyunca, Beyonce'ın diğer kadınlardan daha fazla şarkısı vardı? | 2000'ler *(2000s)* | 2000'li *(2000s)* |
| | English | ...The album featured the number-one song "Single Ladies (Put a Ring on It)" and the top-five songs "If I Were a Boy" and "Halo". Achieving the accomplishment of becoming her longest-running Hot 100 single in her career, "Halo"'s success in the US helped Beyoncé attain more top-ten singles on the list than any other woman during the **2000s**.... | For which decade, did Beyonce have more top ten songs than any other woman? | 2000s | — |
| **Example 3** *(Edit distance=3)* | Turkish | ...Amerika Kayıt Endüstrisi Birliği (RIAA), Beyoncé'yi 2000'lerin en iyi sertifikalı sanatçısı olarak toplamda **64 sertifikayla** listeledi.... | 2000'lerde kaç tane müzik sertifikası aldı? | 64 sertifikasyon *(64 certifications)* | 64 sertifikayla *(with 64 certificates)* |
| | English | ...The Recording Industry Association of America (RIAA) listed Beyoncé as the top certified artist of the 2000s, with a total of **64 certifications**.... | How many music certifications has she received in the 2000s? | 64 certifications | — |

**Table 2**
Examples for the answer spans that are recovered in SQuAD-TR-TRAIN after the automatic post-processing steps.

| | | | | Question Count | | |
|---|---|---|---|---|---|---|
| **Language** | **Dataset** | **Articles** | **Paragraphs** | **Answerable** | **Unanswerable** | **Total** |
| English | SQuAD-EN-TRAIN | 442 | 19035 | 86821 | 43498 | 130319 |
| | SQuAD-EN-DEV | 35 | 1204 | 5928 | 5945 | 11873 |
| Turkish | SQuAD-TR-TRAIN | 442 | 18776 | 61293 | 43498 | 104791 |
| | SQuAD-TR-DEV | 35 | 1204 | 2346 | 5945 | 8291 |
| | XQuAD-TR | 48 | 240 | 1190 | 0 | 1190 |

**Table 3**
Statistics for the SQuAD-EN, SQuAD-TR, and XQuAD-TR datasets.

Canada's national public broadcaster *CBC Television* [9] referenced as an article in the English Wikipedia. However, there is no corresponding article for the same entity in Turkish Wikipedia but rather on *TRT*,[10] which is Türkiye's national public broadcaster. Therefore, all the information required to answer the questions about the Canadian CBC Television is missing in Turkish Wikipedia. Another issue happens when the target article is actually available in Turkish Wikipedia with information-rich content but is missing the target information due to cultural bias. For example, SQuAD2.0 dataset has a question *When was the first known use of the word "computer"?* targeting a passage in the English Wikipedia article *Computer*.[11] The corresponding article *Bilgisayar*[12] in the Turkish Wikipedia does not have any information about the etymological origin of the English word '*computer*', but instead the origin of its Turkish translation '*bilgisayar*'. Asai et al. [5] succinctly describe the issues behind these two examples as *information scarcity* and *information asymmetry*, which can be commonly called the *missing information* in the knowledge source of the target language.

It should be noted that the missing information issues will resolve gradually as the Turkish Wikipedia grows over time in terms of the number of articles and their quality.

To quantify the effect of the improvement in the knowledge source to the success of the OpenQA models, we used two different dumps of the Turkish Wikipedia with the dates spanning about 2 years,[13] which we name as Wiki-TR-2021 and Wiki-TR-2023.

The missing information issues will understate the performance of the retriever models in the OpenQA systems if not mitigated properly. To mitigate these issues, we appended the target context passages of the SQuAD-TR-TRAIN and XQuAD-TR [4] datasets to the Turkish Wikipedia articles to complete our knowledge source. It should be noted that we do not append answer texts, but rather only the *contexts* and *titles*. In this way, we made the target passages in our knowledge source available to our models while ensuring the validity of our experimental protocol. As a result, the total number of passages in our knowledge source increased slightly with the addition of 19,117 unique passages in the SQuAD-TR-TRAIN and XQuAD-TR datasets to the existing articles in the Turkish Wikipedia dump used.

We split the combined passages of varying lengths in the knowledge source into equal chunks of passages using an enhanced whitespace tokenizer, as in the DPR model [29]. The original DPR model for English segments the passages into 100-word chunks resulting in 142 tokens (subwords) on average when the BERT [19] tokenizer is used. Turkish sentences produce about 1.3 times longer sequence of tokens

---

[9] https://en.wikipedia.org/wiki/CBC_Television
[10] https://tr.wikipedia.org/wiki/TRT
[11] https://en.wikipedia.org/wiki/Computer
[12] https://tr.wikipedia.org/wiki/Bilgisayar

[13] We used the data dumps of May 31st, 2021 and May 1st, 2023

| | Language | Short Name | Wikipedia Date | Passage Count | Passage Length | |
|---|---|---|---|---|---|---|
| | | | | | Word Count (Max) | Token Sequence Length (Avg) |
| DPR | English | Wiki-EN-2018 | Dec 20, 2018 | 21,015,324 | 100 | 142 |
| Ours | Turkish | Wiki-TR-2021 | May 31, 2021 | 1,719,277 | 75 | 136 |
| | | Wiki-TR-2023 | May 01, 2023 | 2,192,776 | 75 | 136 |

**Table 4**
Basic statistics for the knowledge sources used in our study and the one used in the DPR model of [29].

| Reader Model | Training Dataset | EM | F1 |
|---|---|---|---|
| mBERT [4] - *Baseline* | SQuAD-EN-TRAIN | – | 55.40 |
| mBERT | SQuAD-TR-TRAIN | 50.00 | 64.76 |
| BERTurk | SQuAD-TR-TRAIN | 51.17 | 67.78 |

**Table 5**
Reader results for the standard formulation of QA task tested on XQuAD-TR.

with the same number of words due to the very rich suffixing morphology of Turkish. For this reason, unlike Karpukhin et al. [29], we split the passages into 75 words instead of 100 words, as 100-word segments in Turkish run a high risk of being truncated by the ColBERT model [31], which accepts up to 180 tokens for documents by default. After splitting the combined passages into equal chunks of 75 words, we obtained a total of 1.7M and 2.1M passages for the Wikipedia dumps dated 2021 and 2023, respectively.

The resulting combined passages then served as the knowledge sources in our study. The basic statistics of these knowledge sources and the one used in the original DPR model for English are given in Table 4.

## 4. Standard QA Methods and Results

To help establish an upper-bound for OpenQA in Turkish, we first conducted a series of standard QA experiments. Artetxe et al. [4] established a baseline for these experiments with an mBERT model [19] trained on SQuAD-EN-TRAIN [49] and tested on XQuAD-TR [4] as a crosslingual QA application. We extended this experiment in two ways. First, we changed the training dataset to SQuAD-TR-TRAIN while keeping all other aspects of Artetxe et al.'s system fixed. The goal of this experiment is to begin to understand SQuAD-TR-TRAIN as a training resource. Second, we changed mBERT to BERTurk [55] to see the effects of pairing an in-language model with an in-language dataset.

For these experiments, we finetuned the BERTurk and mBERT models with the same hyperparameters for the standard QA formulation using SQuAD-TR-TRAIN and XQuAD-TR as the training and test datasets. We used a batch size of 16, without gradient accumulation, on a single NVIDIA Tesla V100 GPU. We applied $3 \times 10^{-5}$ as the learning rate, used a maximum length of 384, with a document stride of size 128,

and trained each model for 5 epochs using Huggingface's transformers library [56], Version 4.14.0.dev0.[14]
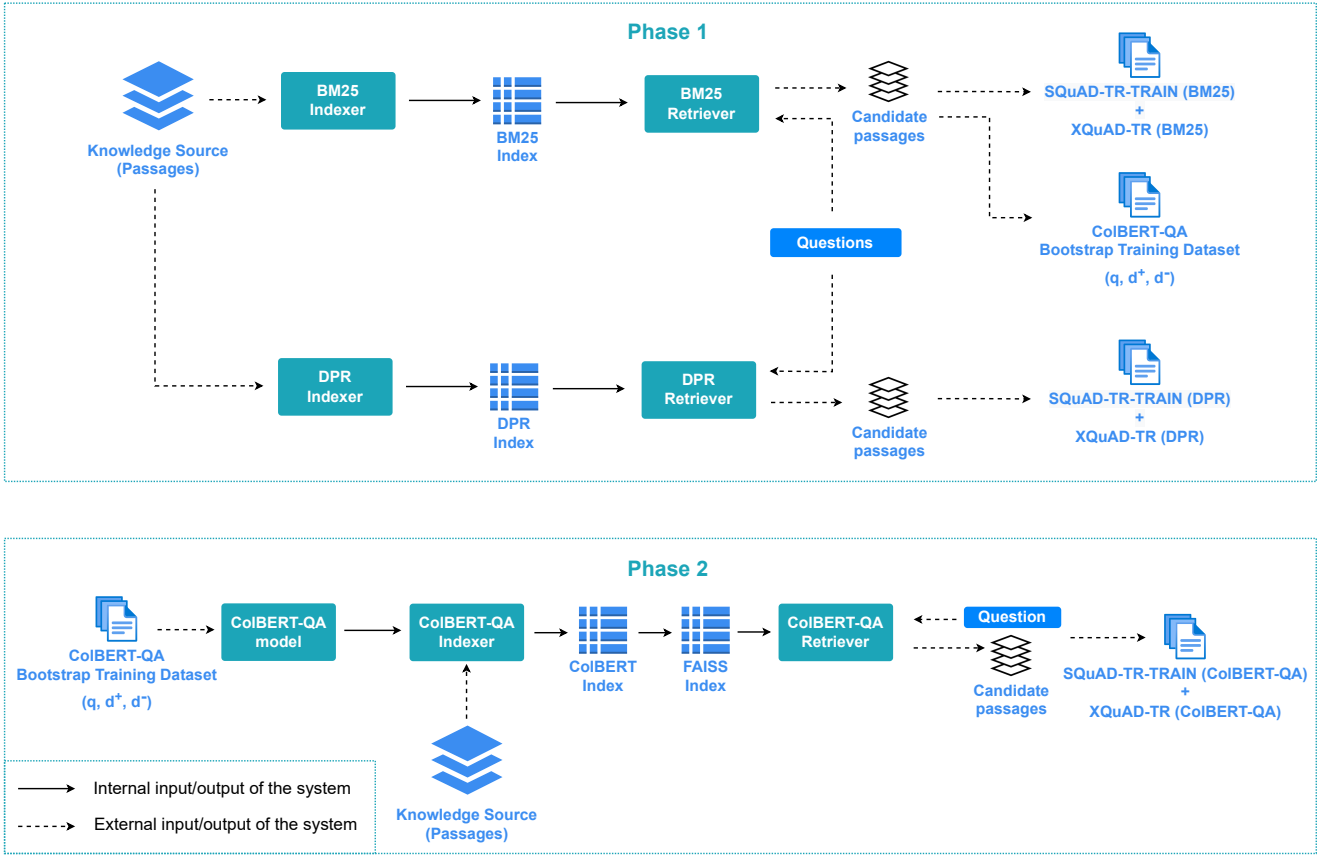
Our evaluation metrics are the standard ones from the literature [50]: we use Exact Match (EM) and F1 scores. EM is the percentage of the predicted answer texts matching at least one of the ground-truth answer texts in an exact manner. F1 is the average of the maximum overlap ratio between predicted answer tokens and ground truth answer tokens. While EM gives no credit to predictions that have no exact match in any of the ground truth answer texts, F1 gives partial credit to those predictions that have at least one partially matching ground truth answer token. We calculated the evaluation metrics on XQuAD-TR as our test set.

Table 5 summarizes the results of the experiments for standard QA where we experiment with mBERT and BERTurk as readers and SQuAD-EN-TRAIN and SQuAD-TR-TRAIN as training data. The results show that the use of an in-language model, BERTurk, and an in-language dataset, SQuAD-TR-TRAIN, yields the highest scores. The largest performance gap occurs when SQuAD-EN-TRAIN is replaced with SQuAD-TR-TRAIN, indicating that in-language datasets are essential for high-performing standard QA models, even if the datasets are machine-translated and potentially noisy. The use of an in-language model instead of a multi-lingual one has a smaller positive impact on the performance. This aligns with other recent findings in the literature [8, 42].

## 5. OpenQA Methods and Results

In this section, we turn to OpenQA for Turkish. We establish the baseline using BM25 and DPR [29] as early examples of sparse and dense retrievers. Then, we share the results of our proposed system based on ColBERT-QA [30]. We first review the main components of our system, the retriever and the reader, and then we report our experimental results. We conduct the experiments for each Wikipedia

---

[14]The choice of specific values for the hyperparameters in our study is primarily aimed at establishing an initial reference point for future studies within a constrained budget.

**Figure 1:** System overview diagram for the OpenQA retriever component. We assume the knowledge source is independent of the system for the sake of clarity in the figure.

dump as the knowledge source separately, which allows us to observe the contribution of the growth in the knowledge source.

### 5.1. Retriever

The first step in building our ColBERT-based [31] retriever involves a handful of steps that are specific to Turkish but that may have more general utility for cross-linguistic applications:

1. ColBERT-QA uses the WordPiece tokenizer of the original BERT-base model [19] in English. We replace this tokenizer with the WordPiece tokenizer of the BERTurk cased model [55], which was pretrained on a large Turkish corpus.

2. The original tokenizer of the ColBERT-QA model repurposes "unused" tokens in the tokenizer as query and document markers, which are available in the tokenizer of the original BERT-base model in English. As the BERTurk tokenizer did not have such unused tokens, we use alternative tokens for the document and question marker tokens. For the query marker we use a "blush" emoji (U+1F60A) and for the document marker we use a "smiley" emoji (U+1F603), as they are unlikely to occur in the Wikipedia articles yet likely to be present in various non-English BERT models.

3. ColBERT-QA initializes its weights using those of the original BERT model in English. We use the BERTurk weights to initialize the ColBERT weights before starting the training step. For languages without high-quality language-specific embeddings like BERTurk, one might use multilingual embeddings here instead.

In light of the above steps, our ColBERT-based retriever might more properly be called a ColBERTurk retriever. In the interest of clarity, we will continue to refer to it as a ColBERT model.

To train the retriever, we proceed in two phases, as outlined in Figure 1. In phase 1 (top row of the figure), we build our baseline retriever models. We rely on BM25 to index our knowledge sources, using `pyserini` [41][15] and `anserini` [40][16] wrappers for the Apache Solr search engine. We customize Apache Solr for Turkish by incorporating the Zemberek[17] plugin [3] as a morphological stemmer for Turkish. In addition to BM25 retriever, we use DPR model to index our knowledge sources and build our baseline for dense retriever.

In our experiments, the BM25 retriever provides the bootstrap dataset to train ColBERT-QA. With this

---

[15] https://github.com/castorini/pyserini
[16] https://github.com/castorini/anserini
[17] https://github.com/iorixxx/lucene-solr-analysis-turkish

lightweight BM25 retriever, we create a dataset of triples $(q, d^+, d^-)$, where $q$ is a question in SQuAD-TR-TRAIN, $d^+$ is a positive passage containing the target answer span for $q$, and $d^-$ is a negative passage that does not contain the target answer span for $q$. Both $d^+$ and $d^-$ are from the top $k$ results retrieved from the BM25 index. More specifically, we create the dataset of triples by pairing every $d^+$ with every other $d^-$ where $d^+$ is from the top $k^+$ results and $d^-$ is from the top $k^-$ results ($k^+ \leq k^-$) obtained for each question $q$.

For both of the Turkish Wikipedia dumps used as the knowledge source, the resulting bootstrap training dataset for ColBERT-QA model contains 6M triples for 86K questions, where $k^+ = 3$ and $k^- = 100$. It can be noted that we could use all question–answer pairs in SQuAD-TR that were originally labeled as answerable before translating SQuAD2.0. The reason for also including those question–answer pairs that we excluded from SQuAD-TR-TRAIN (§3.1) is that the retriever model, unlike the reader, does not require the location of the answer span in the context. Therefore, the retriever model can use all question–answer pairs in SQuAD-TR-TRAIN.

In phase 2 (bottom row of Figure 1), we use our BM25-derived dataset to train a ColBERT-QA model, and then we index our knowledge source using this retriever. This indexer computes the passage representations using ColBERT-QA to project them into an embedding space where the question and passage representations are close to each other if the passage has an answer for the question.

We trained our ColBERT-QA model in Turkish on a single NVIDIA Tesla V100 GPU with a maximum document length of 180 and batch size of 32 without gradient accumulation. Then, we indexed all the passages in the knowledge source once again, this time using ColBERT-QA. Following Khattab et al. [30], we further reindexed ColBERT-QA indexed document embeddings using FAISS [26] on Index-IVFPQ mode with a 16384 partition and a sample rate of 0.3 to speed up the retriever component.

Since there is no state-of-the-art model for OpenQA in Turkish yet, we compare the retriever and reader performance of our models with the performance of models based on the baseline BM25 and DPR retrievers. It is important to note that each retriever determines its own versions of the SQuAD-TR-TRAIN and XQuAD-TR [4] datasets specific to that retriever and to the knowledge source it uses. For each retriever, we retrieve the top $k$ passages for each question in SQuAD-TR-TRAIN and XQuAD-TR to set a context passage for that question. For the train sets, we use the first positive result out of top 5 retrieved results and remove those examples that have no positive result in the top 5 retrieved results.

## 5.2. Reader

In line with the three sets of OpenQA retrievers obtained in the retriever step for each knowledge source with different Wikipedia dumps, we build three sets of reader components in our OpenQA system in Turkish, depending on which retriever their training and test sets are based on. Figure 2 summarizes this process. We used the BERTurk cased model [55] along with its original

tokenizer for each reader, and we finetuned them using the retriever and knowledge source specific versions of SQuAD-TR-TRAIN, namely SQuAD-TR-TRAIN (BM25,YYYY) and SQuAD-TR-TRAIN (ColBERT-QA,YYYY), where YYYY $\in$ {2021, 2023} denotes the year of the Wikipedia dump.

We used the same hyperparameters as described in §4 to train and evaluate the BERTurk model on the retriever-specific datasets as shown in Figure 2. We also calculated EM and F1 scores on the open versions of XQuAD-TR [4] for each reader model.
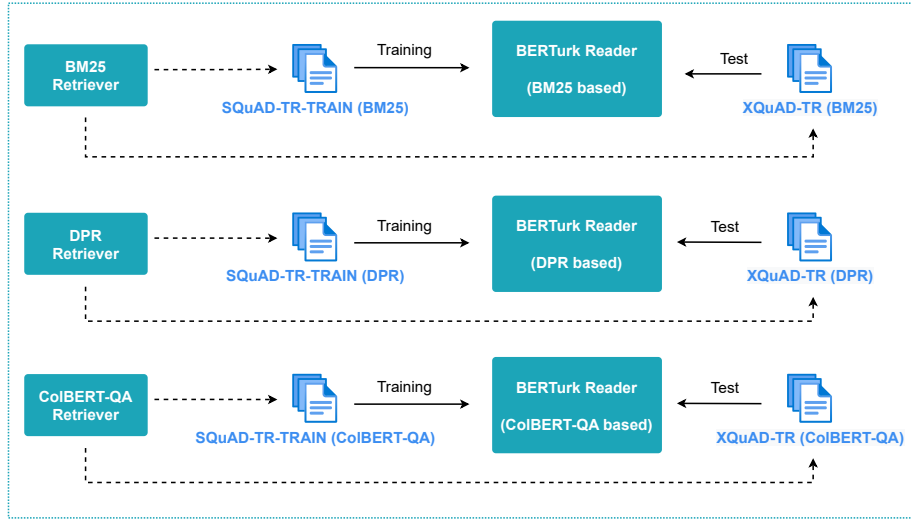
## 5.3. Retriever Results

The success of the retriever component sets an upper bound for the reader. Following previous works (DrQA [13]; DPR [29]; REALM [23]; ColBERT-QA [30]), we evaluated the success of the retriever component by means of $Success@k$, also noted as $S@k$, which is the ratio of the questions having a positive passage among the top $k$ results retrieved from the index. We evaluated $S@k$ values for $k \in \{1, 5, 20\}$. We supplemented $S@k$ with another metric, $Count@k$, also noted as $C@k$, which is the average number of positive passages among the top $k$ results retrieved for each question.

Turkish has different morphological characteristics from English, as it is an agglutinative language and has more morphological variants of each word. For this reason, the evaluation scores depend heavily on the tokenization scheme that is used when evaluating the results. Therefore, we used three different tokenization schemes: *whitespace*, *morphological*, and *enhanced whitespace*. Whitespace tokenization calculates the $S@k$ and $C@k$ values after splitting the retrieved passage and answer text into tokens whenever it finds a whitespace character. The morphological tokenization scheme segments the passage and answer texts into a list of stems by stripping all suffixes in all words before calculating the $S@k$ and $C@k$ values. The enhanced whitespace tokenization, which is the standard tokenizer of the DPR model, breaks the text into tokens not only when it encounters a whitespace character but also whenever it finds a list of predefined punctuations. We used the uncased version for all tokenization schemes to bring them in line with the output of our morphological tokenizer (Zemberek [2]), which was uncased out of the box.

Table 6 shows the results of each retriever. The results indicate that the contribution of the late interaction architecture of ColBERT [31] is markedly more effective than the baseline BM25 and DPR models, even though the BM25 retriever is empowered with a morphological stemmer as described in §5.2. We observe a performance improvement over the BM25 and DPR [29] models independent of the tokenization scheme. The results suggest that the ColBERT-QA retriever will give the reader module a better chance at finding correct answers.

Comparing the baseline retriever models, we noticed a substantial performance advantage of BM25 over the DPR retriever. This finding is consistent with the reported performance of DPR [29] on the English SQuAD 1.1 [50] dataset.

**Figure 2:** Overview diagram for the OpenQA reader component. For the sake of clarity in the figure, we assume the knowledge source that retrievers are based on is invariant of the system.

| Retriever Model | Knowledge Source | Whitespace | | | Morphological | | | Enhanced Whitespace | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | S@1/C@1 | S@5/C@5 | S@20/C@20 | S@1/C@1 | S@5/C@5 | S@20/C@20 | S@1/C@1 | S@5/C@5 | S@20/C@20 |
| BM25 - *Baseline - Sparse* | Wiki-TR-2021 | 42.79/0.43 | 58.91/<u>0.85</u> | 66.64/1.17 | 45.97/0.46 | 62.27/0.95 | 69.92/1.39 | 56.30/0.56 | 73.53/1.11 | 82.10/1.55 |
| | Wiki-TR-2023 | 41.68/0.42 | 58.15/0.85 | 66.22/1.18 | 44.62/0.45 | 61.43/<u>0.95</u> | 69.66/1.42 | 55.21/0.55 | 72.52/1.10 | 81.18/<u>1.55</u> |
| DPR - *Baseline - Dense* | Wiki-TR-2021 | 36.80/0.37 | 53.69/0.75 | 62.69/1.08 | 39.91/0.40 | 57.56/0.83 | 66.81/1.30 | 46.63/0.47 | 67.90/0.97 | 78.24/1.43 |
| | Wiki-TR-2023 | 37.90/0.38 | 53.53/<u>0.75</u> | 62.35/1.09 | 41.00/0.41 | 57.23/0.84 | 66.05/1.34 | 48.57/0.49 | 67.65/0.98 | 77.65/1.47 |
| ColBERT-QA | Wiki-TR-2021 | 58.91/0.59 | 69.58/1.06 | 74.54/1.44 | 62.10/0.62 | **73.28**/1.17 | **78.49**/**1.71** | 75.71/0.76 | **87.39**/1.39 | **92.86**/1.89 |
| | Wiki-TR-2023 | **59.66**/**0.60** | 69.58/1.07 | 74.20/1.43 | **62.61**/**0.63** | 72.77/**1.18** | 77.98/**1.71** | **76.05**/**0.76** | 86.81/<u>1.39</u> | 92.10/**1.89** |

**Table 6**
Retriever results for the OpenQA formulation of QA task tested on XQuAD-TR. All tokenizers are uncased. The highest values in each column are shown in **bold**. For equal pairs, the larger ones on more significant digits are <u>underlined</u>.

This is attributed to the fact that the annotators of the SQuAD datasets [50, 49] tended to formulate questions with significant lexical overlap with their passages, thereby providing an advantage to BM25. Karpukhin et al. [29] also point out the skewed distribution of the target Wikipedia passages compared to the vast number of Wikipedia articles in the knowledge source as another contributing factor. To address this issue, they suggested a hybrid approach that combines the outcomes of BM25 and DPR in order to achieve a result that surpasses each individually. When we compare the baseline models with the ColBERT retriever, we observe that ColBERT achieves markedly superior performance than the baselines on SQuAD-TR. This indicates the outstanding effectiveness of ColBERT retriever in handling examples characterized by lexical overlap as well as those requiring deep semantic understanding.

The results also indicate that the morphology-unaware *enhanced whitespace tokenizer* identifies the correct results better than the morphological tokenizer for all values of $S@k$ and $C@k$, suggesting that computationally-intensive morphological stemming can be avoided when evaluating QA systems in Turkish. Although the negative effect of morphological stemming may be surprising given the rich morphology of Turkish, this result is in line with previous literature [8].

As an additional perspective, we did not observe a significant improvement in the $S@k$ values when utilizing the newer Turkish Wikipedia dump, except for minor increases in S@1 scores for DPR and ColBERT-QA. This finding may imply that the expansion of the knowledge source resulted in more interference rather than providing valuable instances. However, we did observe a slight increase in the $C@k$ values even when the $S@k$ values remained unchanged or decreased. This suggests that the results contain a greater number of positive examples and presumably more challenging negative examples. Consequently, these outcomes indicate that the quality of the resulting bootstrap training dataset for the reader models may improve as the knowledge source continues to grow over time.

We also evaluated the retrievers in terms of the ranking of the answers returned, which is a widely used metric in the information retrieval domain. In order to assess this property, we computed mean reciprocal rank (MRR) for the questions correctly answered by each retriever by returning the relevant passage among their top k results at different ranks. Note that all retrievers can return the same passage as the highest-ranked relevant passage among their top k results, even though the passage's exact rank may differ in each retriever's output. For questions where the retrievers returned the same relevant passage with the highest rank

| | Turkish | English |
|---|---|---|
| **Question 165** | Bir öğretmenlik sertifikasının geçerli olduğu en uzun süre nedir? | What is the longest time that a teaching certificate is good for? |
| **Answer** | on yıla | ten years |
| **BM25** | ... için kullanılır. Genelde bu iptal bilgilerinin izlenmesinde kullanılır. Subject (Özne): Sertifikanın ait olduğu varlık: bir cihaz, birey, ya da kurum. Issuer (Sağlayıcı): Bilgileri doğrulayan ve sertifikayı imzalayan kuruluş. Not Before (Önce Değil): Sertifikanın geçerli olduğu en erken saat ve tarihi. Not After (Sonra Değil): Sertifikanın geçerli olduğu en geç saat ve tarihi. Key Usage (Anahtar Kullanımı): Sertifikanın açık anahtarındaki geçerli kriptografik kullanım. Ortak alanlar arasında dijital imza doğrulaması, anahtar şifreleme ve sertifika imzalama bulunur. Extended ... <br><br>Source: `Wiki-TR-2023` | It is used for monitoring the validity information. Subject: The entity to which the certificate belongs: a device, individual, or organization. Issuer: The organization that verifies the information and signs the certificate. Not Before: The earliest date and time at which the certificate is valid. Not After: The latest date and time at which the certificate is valid. Key Usage: The valid cryptographic usage in the public key of the certificate. Common fields include digital signature verification, key encryption, and certificate signing. Extended ... |
| **ColBERT-QA** | ... yönelik gereksinimler, genelde tam zamanlı profesyonellere yönelik gereksinimler kadar sert değildir. İş Gücü İstatistikleri Bürosu, ABD'de 1,4 milyon ilkokul öğretmeni, 674.000 ortaokul öğretmeni ve 1 milyon lise öğretmeni istihdam edildiğini tahmin etmektedir. Amerika Birleşik Devletleri'nde her eyalet devlet okullarında öğretmenlik yapma lisansı almak için gereksinimleri belirler. Öğretim sertifikasyonu genelde üç yıl devam eder, ama öğretmenler **on yıla** varan uzunlukta sertifikalar alabilirler. Devlet okulu öğretmenlerinin bir lisans derecesine sahip olması şart koşulmakta ve öğretmenlerinin çoğunun eğitim ... <br><br>Source: `XQuAD-TR` | ... are generally not as rigorous as those for full-time professionals. The Bureau of Labor Statistics estimates that there are 1.4 million elementary school teachers, 674,000 middle school teachers, and 1 million secondary school teachers employed in the U.S. In the United States, each state determines the requirements for getting a license to teach in public schools. Teaching certification generally lasts three years, but teachers can receive certificates that last as long as **ten years**. Public school teachers are required to have a bachelor's degree and the majority ... |
| **Question 484** | Siliya ne için kullanılır? | What are cilia used for? |
| **Answer** | hareket yöntemi | method of locomotion |
| **BM25** | 1 milimetreden (0,039 in) 1,5 metreye (4,9 ft) kadar değişen boyutlarıyla taraklılar, ana **hareket yöntemi** olarak siliya ("kıl") kullanan en büyük kolonyal olmayan hayvanlardır. Çoğu türün tarak dizisi denen ve vücutları boyunca devam eden, ktene adı verilen taraksı siliya grupları taşıyan sekiz dizisi vardır ve böylece siliya vurduğunda her tarak alttaki tarağa dokunur. "Ktenofor", Yunanca'da "tarak" anlamına gelen κτειζ (kök biçimi κτεν- ) ile "taşıyan" anlamına gelen Yunanca son ek -φοροζ 'tan gelir ve "tarak taşıyan" ... <br><br>Source: `XQuAD-TR` | Ranging from about 1 millimeter (0.039 in) to 1.5 meters (4.9 ft) in size, ctenophores are the largest non-colonial animals that use cilia ("hairs") as their main **method of locomotion**. Most species have eight strips, called comb rows, that run the length of their bodies and bear comb-like bands of cilia, called "ctenes" stacked along the comb rows so that when the cilia beat, those of each comb touch the comb below. The name "ctenophora" means "comb-bearing", from the Greek κτειζ (stem-form κτεν-) meaning "comb" and the Greek suffix -φοροζ meaning "carrying"... |
| **ColBERT-QA** | Silikon veya polisiloksan, siloksan'dan (-R2Si-O-SiR2-, burada R = organik grup) oluşan bir polimer'dir. Bunlar genellikle renksiz yağlar veya kauçuk benzeri maddelerdir. Silikonlar, dolgu macunlarında, yapıştırıcılarda, yağlayıcılarda, tıpta, pişirme kaplarında, ısı ve elektrik yalıtımında kullanılır. Bazı yaygın biçimler arasında silikon yağı, silikon gresi, silikon kauçuk, silikon reçine ve silikon kalafat bulunur. Daha kesin olarak polimerize edilmiş siloksan'lar veya polisiloksanlar olarak adlandırılan silikonlar, her silikon merkezine bağlı iki organik gruplu inorganik silikon-oksijen omurga zinciri'nden (····-Si-O-Si-O-Si-O-···· oluşur. Genellikle ... <br><br>Source: `Wiki-TR-2023` | Silicone or polysiloxane is a polymer composed of siloxane (-R2Si-O-SiR2-, where R = organic group). These are typically colorless oils or rubber-like substances. Silicones are used in caulks, adhesives, lubricants, medicine, cooking utensils, and for thermal and electrical insulation. Some common forms of silicones include silicone oil, silicone grease, silicone rubber, silicone resin, and silicone caulk. More precisely, silicones, also called polymerized siloxanes or polysiloxanes, consist of inorganic silicon-oxygen backbone chains with two organic groups attached to each silicon center (····-Si-O-Si-O-····). They are generally ... |

**Table 7**

The negative and positive effects of the TF-IDF approach used in BM25 are exemplified in Question 165 and Question 484, respectively. The content words in the questions and the corresponding correctly matched terms in the passages are shown underlined. The dashed lines represent the incorrectly matching terms (false positives) that adversely affect the results.

among their top k results, the ColBERT-QA, BM25, and DPR retrievers achieved MRR scores of 0.98, 0.89, and 0.85, respectively. These results show that ColBERT-QA performs significantly better than BM25, which also outperforms DPR in the task of ranking relevant passages.

## 5.4. Qualitative Analysis of the Retriever

In order to observe the strengths and weaknesses of the retrievers relative to each other, we manually analyzed the passages retrieved for the questions in the test set by the two top-performing retrievers, ColBERT-QA [30] and BM25. The analysis revealed a number of factors that help to explain the performance differences between the sparse and dense retriever models [14].

One important factor is the TF-IDF-based scoring mechanism used in BM25, which results in the retrieval of irrelevant passages that excessively mention the uncommon content words in the questions. While this approach proves advantageous when there are only a few relevant candidate passages, it comes with significant side effects when the model needs to suppress multiple related passages for the question to return the actual relevant passage. The decrease in the $S@k$ and $C@k$ values for the BM25 retriever as the knowledge source expands (Table 6) indicates that BM25 is

|  | Turkish | English |
|---|---------|---------|
| **Question 439** | Amazon Havzası'nda kaç ülke bulunmaktadır? | How many nations are within the Amazon Basin? |
| **Answer** | dokuz | nine |
| **BM25** | Amazon Havzası, Güney Amerika'nın Amazon Nehri ve kolları tarafından beslenen bölümüdür. Amazon drenaj havzası 6.300.000 kilometrekare (2.400.000 sq mi) bir alanı kaplamaktadır ve bu değer Güney Amerika kıtasının yaklaşık %35,5'ini oluşturmaktadır. Havza Bolivya, Brezilya, Kolombiya, Ekvador, Fransız Guyanası (Fransa), Guyana, Peru, Surinam ve Venezuela ülkeleri sınırları içinde yer almaktadır. Havzanın çoğu Amazon yağmur ormanları ile kaplıdır. Kapladığı 55 milyon kilometrekare $(21 \times 10^6$ sq mi) alan ile tropikal orman alanı, dünyanın en büyük yağmur ormanıdır. Dematteis, Lou; ... <br><br> Source: Wiki-TR-2023 | The Amazon basin is the part of South America drained by the Amazon River and its tributaries. The Amazon drainage basin covers an area of about 6,300,000 km2 (2,400,000 sq mi), or about 35.5 percent of the South American continent. It is located in the countries of Bolivia, Brazil, Colombia, Ecuador, Guyana, Peru, Suriname, and Venezuela, as well as the territory of French Guiana. <br> Most of the basin is covered by the Amazon rainforest, also known as Amazonia. With a 5.5 million km2 (2.1 million sq mi) area of dense tropical forest, it is the largest rainforest in the world. Dematteis, Lou; ... |
| **ColBERT-QA** | Amazon yağmur ormanı (Portekizce: Floresta Amazônica veya Amazônia; İspanyolca: Selva Amazónica, Amazonía veya genellikle Amazonia; Fransızca: Forêt amazonienne; Hollandaca: Amazoneregenwoud) İngilizce'de aynı zamanda Amazonia veya Amazon Jungle olarak da bilinir ve Güney Amerika'nın Amazon havzasının çoğunu kaplayan bir nemli geniş yapraklı ormandır. Bu havza 7.000.000 kilometre karelik alanı kaplamaktadır (2.700.000 mil kare) ve bunun 5.500.000 kilometre karesi (2.100.000 mil kare) yağmur ormanıyla kaplıdır. Bu bölge **dokuz** ulusa ait toprakları içermektedir. Ormanın çoğu yağmur ormanının %60'ı ... <br><br> Source: XQuAD-TR | The Amazon rainforest (Portuguese: Floresta Amazônica or Amazônia; Spanish: Selva Amazónica, Amazonía or usually Amazonia; French: Forêt amazonienne; Dutch: Amazoneregenwoud), also known in English as Amazonia or the Amazon Jungle, is a moist broadleaf forest that covers most of the Amazon basin of South America. This basin encompasses 7,000,000 square kilometres (2,700,000 sq mi), of which 5,500,000 square kilometres (2,100,000 sq mi) are covered by the rainforest. This region includes territory belonging to **nine nations**. The majority of the forest is contained within Brazil, with 60% ... |
| **Question 922** | Hangi Nobel Ekonomi Ödülü kazananı aynı zamanda bir üniversite mezun üyesidir? | What Nobel Memorial Prize in Economic Sciences winner is also a university alumni member? |
| **Answer** | Milton Friedman | Milton Friedman |
| **BM25** | Üniversitelerine göre Nobel Ödülü sahipleri listesi, Nobel Ödülü kazananların (öğrenci veya mezun oldukları üniversitelere göre) birinci derecede eğitim gördükleri Üniversitelere göre listelenmiş halidir. Üniversiteler Nobel Ödülü kazananların sayısına göre doğru orantılı şekilde sıralanmıştır. Üniversitelere göre listeleme işlemi oldukça kapsamlı bir çalışma gerektirmiştir. Bu nedenle çok çeşitli kaynaklardan yararlanılmıştır. Ödülü kazanan birçok kişi farklı Üniversitelere geçiş veya doktora yapmıştır. Bu nedenle bazı ödül sahipleri birden fazla Üniversitede eğitim görmüş veya doktora yapmış olabilir. Aşağıdaki listede ödül ... <br><br> Source: Wiki-TR-2023 | The list of Nobel Prize laureates according to their universities is a compilation that categorizes the winners (based on whether they were students or graduates) according to the universities where they received their primary education. The universities are ranked in proportion to the number of Nobel Prize recipients. The process of listing the universities required extensive and comprehensive research, utilizing various sources. Many prize winners have transitioned to different universities or pursued doctoral degrees, resulting in some recipients having received education or completed their doctorates at multiple universities. The list below features the recipients of the award. |
| **ColBERT-QA** | ... yer alır. Amerikalı ekonomist, sosyal kuramcı, politik filozof ve yazar Thomas Sowell de üniversitenin mezunları arasındadır. Ekonomide, tanınmış Nobel Ekonomi Ödüllü, ABD'nin Cumhuriyetçi Başkanı Ronald Reagan'ın Muhafazakar Britanya Başbakanı Margaret Thatcher'ın baş danışmanlarından biri olan **Milton Friedman**, Nobel ödüllü ve düzenleme tuzağı teorisini ileri süren George Stigler, ekonominin aile ekonomisi dalına önemli katkılar sunmuş olan Gary Becker, örgütsel karar verme konseptinin modern yorumundan sorumlu Herbert A. Simon, Nobel Ekonomi Ödüllü ilk Amerikalı olan Paul Samuelson ... <br><br> Source: XQuAD-TR | ... are. American economist, social theorist, political philosopher, and author Thomas Sowell is also an alumnus. In economics, notable Nobel Memorial Prize in Economic Sciences winners **Milton Friedman**, a major advisor to Republican U.S. President Ronald Reagan and Conservative British Prime Minister Margaret Thatcher, George Stigler, Nobel laureate and proponent of regulatory capture theory, Gary Becker, an important contributor to the family economics branch of economics, Herbert A. Simon, responsible for the modern interpretation of the concept of organizational decision-making, Paul Samuelson, the first American to win the Nobel Memorial Prize in Economic Sciences ... |

**Table 8**
Examples showcasing how ColBERT-QA can better capture the information needs implicit in questions. The content words in the questions and the corresponding correctly matched terms in the passages are shown underlined.

ineffective in inhibiting new irrelevant signals. Table 7 depicts two question–answer pairs[18] and passages returned by each retriever. The first example shows that TF-IDF-based scoring misleads the retriever and causes it to retrieve irrelevant passages containing the content words. Conversely, in the second example where there were limited relevant candidate passages available, BM25 identified the correct passage. This observation aligns with the qualitative analysis conducted by [29] for English OpenQA, which compares the results of BM25 and DPR.

Another factor is the ability of ColBERT-QA to represent questions better than BM25 and thus retrieve relevant passages more accurately. Two example questions are given in Table 8. In the first one, both models retrieved passages related to the Amazon (region) and its surrounding countries. However, only the passage retrieved by ColBERT-QA provided a specific numerical answer to the question "how many". In contrast, the passage retrieved by BM25 consisted of the correct list of the countries without an explicit count, posing a challenge for the reader in extracting the correct answer span. In the second example, ColBERT-QA successfully identified the information need as a Nobel Prize winner who is also a member of a university alumni, and retrieved the relevant passage. In contrast, BM25 retrieved a generic passage about universities and the Nobel Prize, completely

---

[18]Question numbers are 1-based indexes of the questions in XQuAD-TR.

|  | Turkish | English |
|---|---|---|
| **Question 430** | Kaç ülke isminde "Amazonas" bulunmaktadır? | How many nations contain "Amazonas" in their names? |
| **Answer** | Dört | four |
| **BM25** | Tarık el-Tayyib Muhammed Buazizi (29 Mart 1984 - 4 Ocak 2011), Tunuslu seyyar satıcı. 17 Aralık 2010'da kendisini yakarak intihar girişiminde bulundu. Bu olayın tesiri ile Tunus halkının ayaklanması üzerine 23 yıldır ülkeyi yöneten Zeynel Abidin Bin Ali ülkeden kaçmıştır. Bu olay aynı zamanda diğer Arap ülkelerindeki ayaklanmaları teşvik etmiştir. Ölümünden sonra Tunus'ta Yasemin Devrimi başlamıştır. 17 Ocak 2011'de başkent Tunus'un en ünlü caddesi olan 7 Kasım Caddesi'nin ismi (Zeynel ...<br><br>Source: `Wiki-TR-2023` | Tarık el-Tayyib Muhammed Buazizi (March 29, 1984 - January 4, 2011) was a Tunisian street vendor. On December 17, 2010, he set himself on fire in a suicide attempt. As a result of this incident, Zine El Abidine Ben Ali, who had been ruling the country for 23 years, fled from Tunisia. This event also inspired uprisings in other Arab countries. Following his death, the Jasmine Revolution began in Tunisia. On January 17, 2011, the name of 7 November Avenue, the most famous street in the capital city of Tunis (Zine El ... |
| **ColBERT-QA** | ...ile Brezilya sınırları içindedir, ardından %13 ile Peru, %10 ile Kolombiya, ve daha az oranlarla Venezuela, Ekvador, Bolivya, Guyana, Surinam ve Fransız Guyanası gelir. **Dört** ülkenin eyalet veya il isimlerinde "Amazonas" geçmektedir. Amazon gezegenin mevcut yağmur ormanlarının yarıdan fazlasını temsil etmektedir ve dünyadaki en büyük ve en çok biyoçeşitliliğe sahip tropik yağmur ormanı alanını içermektedir ve buna 16.000 türe ayrılan 390 milyar ağaç dahildir. Amazon yağmur ormanı (Portekizce: Floresta Amazônica veya Amazônia; İspanyolca: Selva Amazónica, ...<br><br>Source: `XQuAD-TR` | ...is contained within Brazil, followed by Peru with 13%, Colombia with 10%, and with minor amounts in Venezuela, Ecuador, Bolivia, Guyana, Suriname and French Guiana. States or departments in **four** nations contain "Amazonas" in their names. The Amazon represents over half of the planet's remaining rainforests, and comprises the largest and most biodiverse tract of tropical rainforest in the world, with an estimated 390 billion individual trees divided into 16,000 species. The Amazon rainforest (Portuguese: Floresta Amazônica or Amazônia; Spanish: Selva Amazónica ... |
| **Question 941** | Huihui neydi? | What was huihui? |
| **Answer** | Müslüman tıbbı | Muslim medicine |
| **BM25** | Batı tıbbı, bazen huihui ya da **Müslüman tıbbı** olarak adlandırıldığı Yuan meclisinin Nestûrî Hristiyanları tarafından Çin'de de uygulanmıştır. Nestûrî hekim Tercüman İsa, 1963 yılında, Kubilay'ın saltanatı döneminde Batı Tıbbı Ofisini kurmuştur. İki imparatorluk hastanesinde çalışan doktorlar imparatorluk ailesi ve meclisin üyelerini tedavi etmekten sorumluydu. Çinli hekimler, hümoral sistemi, geleneksel Çin tıbbının altında yatan yin-yang ve wuxing felsefesine karşı geldiği için Batı tıbbına karşı çıkıyorlardı. Batı tıbbı çalışmalarının bilinen bir Çin tercümesi yoktur ama Çinlilerin İbn-i ...<br><br>Source: `XQuAD-TR` | Western medicine was also practiced in China by the Nestorian Christians of the Yuan court, where it was sometimes labeled as huihui or **Muslim medicine**. The Nestorian physician Jesus the Interpreter founded the Office of Western Medicine in 1263 during the reign of Kublai. Huihui doctors staffed at two imperial hospitals were responsible for treating the imperial family and members of the court. Chinese physicians opposed Western medicine because its humoral system contradicted the yin-yang and wuxing philosophy underlying traditional Chinese medicine. No Chinese translation of Western medical works is known, but Chinese had Avicenna ... |
| **ColBERT-QA** | Tüzükleri, işbirliği yapmayan çocuk işçiler için hapis şartlarını öngörmüştür. Hong Kong gibi güneydoğu Asya kolonilerinde Mui Tsai () gibi çocuk işçiliği kültürel bir gelenek olarak rasyonelleştirildi ve İngiliz yetkililer tarafından göz ardı edildi. Hollanda Doğu Hindistan Şirketi yetkilileri, çocuklarının işçi tacizlerini "bu, bu çocukları daha kötü bir kaderden kurtarmanın bir yolu" ile mantıklı hale getirdiler. Zambiya'dan Nijerya'ya uzanan bölgelerdeki Hıristiyan misyon okulları da çocuklardan çalışma gerektirdi ve karşılığında laik eğitim değil din eğitimi sağladı. Başka ...<br><br>Source: `SQuAD-TR` | In southeast Asian colonies, such as Hong Kong, child labour such as the Mui Tsai (), was rationalised as a cultural tradition and ignored by British authorities. The Dutch East India Company officials rationalised their child labour abuses with, "it is a way to save these children from a worse fate." Christian mission schools in regions stretching from Zambia to Nigeria too required work from children, and in exchange provided religious education, not secular education. Elsewhere ... |

**Table 9**

Examples that illustrate how WordPiece tokenization can produce a mix of favorable and unfavorable outcomes, depending on its ability to resist the influence of lexical bias. The content words in the questions and the corresponding correctly matched terms in the passages are shown underlined. The dashed lines represent the incorrectly matching terms (false positives) that adversely affect the results.

oblivious to the specific person targeted as the information need.

During manual analysis, we also observed an intriguing aspect related to ColBERT-QA's WordPiece tokenization, which can have both positive and negative implications. Table 9 shows two example cases. In the first example, ColBERT-QA employed WordPiece tokenization to split the word "Amazonas" into the word pieces [`"Amazon"`, `"##as"`]. This split allowed ColBERT-QA to correctly associate the word "Amazonas" with the word "Amazon" and successfully retrieve the relevant passage. On the other hand, BM25 placed excessive emphasis on the term "Amazonas" and other content words in the question due to its lexical bias,

leading to the retrieval of an entirely unrelated passage that contained these content words extensively.

However, WordPiece tokenization can be a liability as well. In the second example, despite the word "Huihui" being a proper noun, ColBERT-QA tokenized it as [`"Hu"`, `"##ih"`, `"##u"`, `"##i"`], resulting in retrieving irrelevant passages. The same effect can also be seen in Question 484 in Table 7, where ColBERT-QA matched the words "Silikon" [`"Sili"`, `"##kon"`] and "Siliya" [`"Sili"`, `"##ya"`] due to their common prefix and incorrectly retrieved a passage on "Silikon" (Silicone) for a question about "Siliya" (Cilia), which are completely different concepts.

| Reader Model | Retriever Model | Training Dataset | Test Dataset | EM | F1 |
|---|---|---|---|---|---|
| BERTurk | BM25 - *Baseline - Sparse* | SQuAD-TR-TRAIN (BM25,2021) | XQuAD-TR (BM25,2021) | 30.00 | 40.78 |
| | | SQuAD-TR-TRAIN (BM25,2023) | XQuAD-TR (BM25,2023) | 38.15 | 50.48 |
| BERTurk | DPR - *Baseline - Dense* | SQuAD-TR-TRAIN (BM25,2021) | XQuAD-TR (BM25,2021) | 36.97 | 47.85 |
| | | SQuAD-TR-TRAIN (BM25,2023) | XQuAD-TR (BM25,2023) | 35.04 | 47.09 |
| BERTurk | ColBERT-QA | SQuAD-TR-TRAIN (ColBERT-QA,2021) | XQuAD-TR (ColBERT-QA,2021) | 40.34 | 54.47 |
| | | SQuAD-TR-TRAIN (ColBERT-QA,2023) | XQuAD-TR (ColBERT-QA,2023) | **44.87** | **59.19** |

**Table 10**
Reader results for the OpenQA formulation of QA task.

## 5.5. Reader Results

Table 10 shows the results of the reader step of the OpenQA formulation. The results demonstrate that the reader trained on the dataset obtained by the ColBERT-QA [30] retriever using `Wiki-TR-2021` achieves around 33-34% improvement and around 9-13% improvement, respectively, compared to the readers that use the baseline BM25-based and DPR retrievers. However, this improvement decreases up to 17% for both EM and F1 when the retrievers use `Wiki-TR-2023`. Based on these findings, we hypothesize that the substantial improvement in the quality of the bootstrap training dataset significantly contributes to narrowing the gap between the reader models towards the upper bound standard QA reader results shown in Table 5.

We notice that the reader scores improve when the retrievers are trained using the newer Wikipedia dump, except for the DPR-based reader. This is a counter-intuitive result, especially considering the fact that we observed a slight improvement in the S@1 score of the DPR retriever (Table 6) when using `Wiki-TR-2023` instead of `Wiki-TR-2021`. The declining reader scores for DPR indicates that the minor increase in the S@1 score of the DPR retriever may not necessarily lead to a higher quality training dataset for the reader, unlike the other retrievers in comparison.

Also, the OpenQA model for the ColBERT-QA based reader achieves almost 88% of the standard formulation QA reader results in terms of both EM and F1 scores. This result suggests that OpenQA formulation is productive for low-resource and resource-constrained languages, since we can rely on machine-generated noisy training data and unstructured knowledge sources.

## 5.6. OpenQA Results with Subsampled Test Sets

One of our central goals is to efficiently create OpenQA systems. Machine translation costs are already manageable and controlled. However, creating gold test sets can lead to unexpectedly high costs, especially if the goal is to have thousands or tens of thousands of examples. Thus, a key question for us is: *How small can our test sets be?*

To begin to address this question, we ran a series of experiments in which we subsampled the OpenQA test sets that we obtained using `Wiki-TR-2023`. Figure 3 summarizes the results of these experiments for our retrievers, and Figure 4 extends the protocols to our readers. In each panel, the x-axis tracks the number of assessment examples, and the y-axis shows our key metrics.
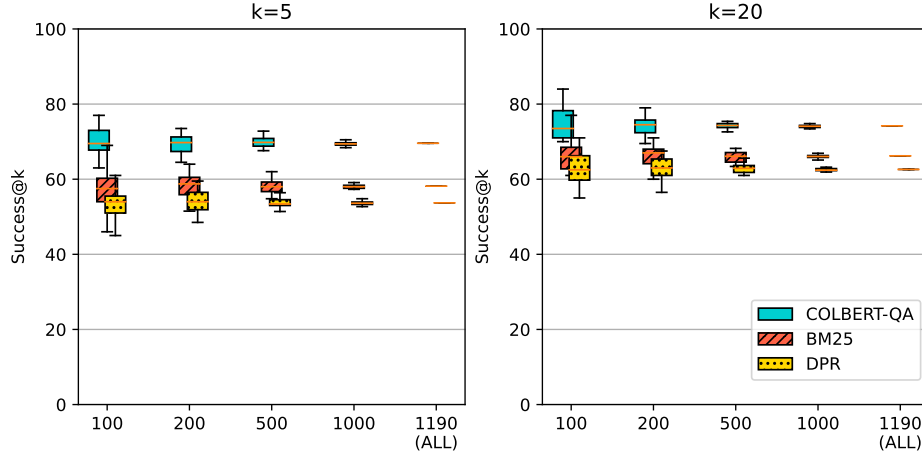
Strikingly, with only 100 examples we can already pretty clearly differentiate our BM25-based and DPR-based models from our ColBERT-QA-based models. By 200 examples, the systems are dramatically different on all metrics for BM25 and ColBERT-QA. As the test sets get larger, the variance of these measures gets tighter, as one would expect, but the core conclusions are unchanged beyond 200 examples for the models. The results of the BM25 and DPR models implicitly suggest that the number of examples needed to differentiate the benchmarked models would increase proportional to the competitiveness of the models with respect to each other.

In this setting, we are using the experiments to differentiate three systems, but the same logic would apply if we were seeking to determine whether a system had truly passed a lower-bound on performance that we set for a production system. Overall, these results show that OpenQA systems can be evaluated very efficiently, which opens the door to conducting multiple distinct evaluations of the same system, which could be crucial for piecing together a picture of how the system behaves overall.
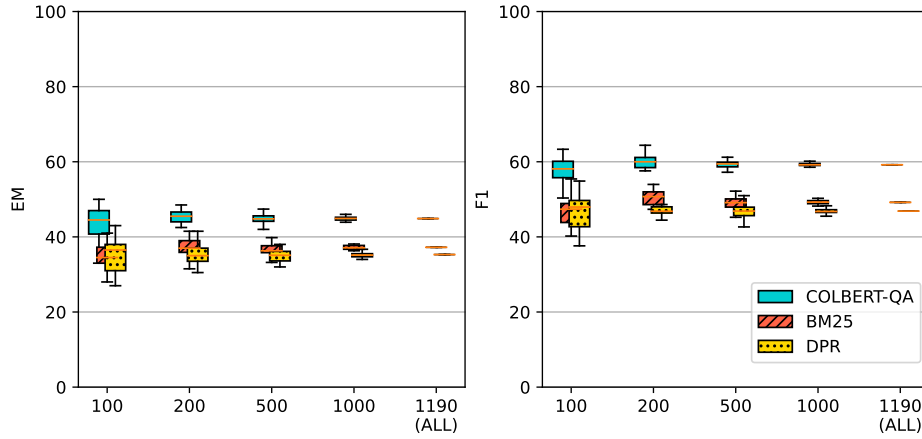
## 6. Discussion

In this section, we discuss the findings of our experiments to assess the effectiveness and limitations of the proposed approach. We begin by analyzing the results of the standard QA formulation in order to observe the potential performance for OpenQA. Then, we delve into the results for each component of the OpenQA formulation, namely, the retriever and reader modules. Finally, we discuss the required number of gold instances for a reliable evaluation of an OpenQA system.

Perhaps unsurprisingly, the best model in the standard QA results is the one that begins with BERTurk parameters and trains on our in-language dataset. However, using multilingual embeddings (mBERT) is competitive with this, and the only large gap in performance is between these two models and the mBERT model trained on English data. This aligns with other recent findings [8, 42] and shows that in-language training data is the real differentiator, even if it is potentially noisy MT data.

**Figure 3:** Retriever results for the OpenQA formulation for different sample portions of XQuAD-TR based on Wiki-TR-2023.



**Figure 4:** Reader results for the OpenQA formulation for different sample portions of XQuAD-TR based on Wiki-TR-2023.

The success of OpenQA systems highly depends on the performance of their retriever models. The baseline BM25 and DPR retrievers were able to capture at most, respectively, 56.30% and 48.57% of the target passages, whereas the ColBERT-QA retriever increased this to 76.05%, when evaluated using an enhanced whitespace tokenizer. As such, the ColBERT-QA retriever model increased the likelihood of the reader model finding more relevant passages to answer questions. These results also indicate that there is a room for further improvement within the retrieval module.

Additionally, we note that expanding the knowledge source has a potential to have a significant positive impact. We attribute this outcome primarily to the improved quality of negative examples acquired by the retrievers and utilized for training the readers. Consequently, we observe a notable increase in the overall success of BM25- and ColBERT-QA-based readers, towards their upper limit capped by the standard QA reader results, accompanied by a significant decrease in the gap between them. Despite the slight increase in the S@1 scores for the DPR retriever, we did not observe a similar increase between the DPR-based readers. This

highlights the need for validating the quality of the retriever scores with the reader scores.

In our OpenQA retriever results, we found that the morphological tokenizer does not identify correct answer spans more effectively than the enhanced whitespace tokenizer for Turkish. This counter-intuitive finding suggests that we should not necessarily favor computation-intensive morphological parsers over enhanced whitespace tokenizers when evaluating QA systems, even in morphologically-rich languages like Turkish. Furthermore, we observed that the baseline results of the BM25 retriever, which is supported by a morphological stemmer, were significantly outperformed by the morphology-unaware ColBERT-QA model.

In the standard QA formulation, the reader models achieved a maximum EM score of 51.17%, which sets the upper bound for the OpenQA reader models. The ColBERT-QA based OpenQA reader results demonstrated that we can preserve almost 88% of this score without requiring gold passages as input. This relaxation of the input requirement provides a significant advantage for

1. **Building a QA dataset in a new language.** Create an extractive reading comprehension dataset in the target language as follows:
   - **Training dataset:** Translate an existing extractive QA dataset using automatic translation.
   - **Evaluation dataset:** Obtain an evaluation dataset containing as little as 200 question–answer pairs.
2. **Compiling a knowledge source.** Create a knowledge source by compiling passages answering the questions in the extractive QA dataset obtained in the previous step.
3. **Training an OpenQA retriever.** Train a neural retriever weakly supervised by the training split of the extractive QA dataset and the knowledge source prepared in the previous steps.
4. **Creating an OpenQA reader.** Train an off-the-shelf reader using the extractive QA training dataset, where the contexts are now provided by the OpenQA retriever.

**Table 11**
Our general method for creating OpenQA systems in low-resource languages efficiently.

developing cost-effective QA systems in low-resource languages.

Although we do not require input gold passages at training time, we still need a labeled dataset at test time. Our experiments revealed that a few hundred evaluation examples may be enough to confidently differentiate the performance of models. This result not only helps limit the cost of obtaining QA systems but also paves the way for obtaining multiple evaluation datasets with different characteristics to better reflect the overall picture of the models in production.

In summary, our proposed system does not require gold datasets during training, but instead utilizes existing unstructured knowledge sources and MT systems to create a machine-translated labeled training dataset. The potential noise in the resulting training dataset can be overcome with the help of the weak supervision used in the OpenQA formulation, which is resilient to noisy data. As a result, we have shown that a cost-effective QA system is feasible for low-resource languages when we shift our focus to the OpenQA formulation.

## 7. Conclusion

In this paper, we obtained an affirmative answer to our core research question, *Can we develop cost-effective OpenQA systems for low-resource languages without requiring a gold training dataset?*. We further expanded our question to explore the minimum test set sizes required to reliably evaluate the performance of OpenQA models. Our findings help pave the way to transferring the rapid advancements made in English QA to non-English QA systems. Moreover, this new avenue not only allows one-way transfer of advancements, but also establishes a virtuous cycle between English OpenQA and non-English OpenQA systems, promoting mutual progress.

We presented a general method for creating efficient and effective OpenQA systems for low-resource languages, and we illustrated the method with a case study of Turkish. Our overall method is summarized in Table 11 as a recipe. As part of this, we introduced SQuAD-TR, a Turkish QA dataset derived by automatically translating SQuAD2.0 [49]. We

showed that SQuAD-TR can straightforwardly be used to train high quality OpenQA systems and benchmark different types of models, and we supported this assessment with detailed qualitative analysis. In addition, we provided evidence that the success of the OpenQA system is notably enhanced by expanding the knowledge source, and we showed that a relatively small number of gold test cases may be sufficient to obtain confident assessments of the quality of these systems.

The key to creating these systems in non-English languages so efficiently is the move from standard QA to OpenQA. In doing this, we greatly simplify the process of creating gold examples, which has been a barrier for the advancement in QA systems for low-resource languages. In OpenQA, these datasets are just question–answer pairs, completely eliminating the necessity for answer span annotation. Consequently, these datasets can now be acquired through automatic translation from the abundant resources available in English. The OpenQA task is also arguably more *relevant*, in that it comes much closer than standard QA to simulating the experience of searching a real-world knowledge store like the Web. Thus, we hope not only to have removed obstacles to creating QA systems for low-resource languages like Turkish, but we also hope to have helped motivate the OpenQA task more generally, as a step towards QA systems that can truly meet the information needs of real-world users. We publicly share our code, models, and data to encourage future research.

## CRediT authorship contribution statement

## Acknowledgments

## Declaration of AI and AI-assisted technologies in the writing process

During the preparation of this work the corresponding author used ChatGPT service for proofreading, spell checking, and grammar correction in specific sections of the manuscript. After using this service, the author carefully reviewed and edited the content as needed and the author takes full responsibility for the content of the publication.

## References

Abadani, N., Mozafari, J., Fatemi, A., Nematbakhsh, M., & Kazemi, A. (2021). ParSQuAD: Persian question answering dataset based on machine translation of SQuAD 2.0. *International Journal of Web Research*, *4*(1), 34–46.

Akın, A. A., & Akın, M. D. (2007). Zemberek, an open source NLP framework for Turkic languages. *Structure*, *10*, 1–5. Retrieved from http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.556.69 (https://github.com/ahmetaa/zemberek-nlp)

Arslan, A. (2016). DeASCIIfication approach to handle diacritics in Turkish information retrieval. *Information Processing & Management*, *52*(2), 326–339. Retrieved from http://www.sciencedirect.com/science/article/pii/S0306457315001053 doi: http://dx.doi.org/10.1016/j.ipm.2015.08.004

Artetxe, M., Ruder, S., & Yogatama, D. (2020, July). On the cross-lingual transferability of monolingual representations. In *Proceedings of the 58th annual meeting of the association for computational linguistics* (pp. 4623–4637). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.acl-main.421 doi: 10.18653/v1/2020.acl-main.421

Asai, A., Kasai, J., Clark, J., Lee, K., Choi, E., & Hajishirzi, H. (2021, June). XOR QA: Cross-lingual open-retrieval question answering. In *Proceedings of the 2021 conference of the north american chapter of the association for computational linguistics: Human language technologies* (pp. 547–564). Online: Association for Computational Linguistics.

---

Retrieved from https://aclanthology.org/2021.naacl-main.46 doi: 10.18653/v1/2021.naacl-main.46

Bajgar, O., Kadlec, R., & Kleindienst, J. (2017). Embracing data abundance. In *5th international conference on learning representations, ICLR 2017, Toulon, France, April 24-26, 2017, Workshop Track Proceedings*. ICLR. Retrieved from https://openreview.net/forum?id=H1U4mhVFe

Bonifacio, L. H., Campiotti, I., Jeronymo, V., Lotufo, R., & Nogueira, R. (2021). mMARCO: A multilingual version of MS MARCO passage ranking dataset. *arXiv preprint arXiv:2108.13897*.

Budur, E., Özçelik, R., Gungor, T., & Potts, C. (2020, November). Data and Representation for Turkish Natural Language Inference. In *Proceedings of the 2020 conference on empirical methods in natural language processing (emnlp)* (pp. 8253–8267). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.emnlp-main.662 doi: 10.18653/v1/2020.emnlp-main.662

Bui, M.-Q., Tran, V., Nguyen, H.-T., & Nguyen, L.-M. (2020, October). How state-of-the-art models can deal with long-form question answering. In *Proceedings of the 34th pacific asia conference on language, information and computation* (pp. 375–382). Hanoi, Vietnam: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.paclic-1.43

Cambazoglu, B. B., Sanderson, M., Scholer, F., & Croft, B. (2021, August). A review of public datasets in question answering research. *SIGIR Forum*, *54*(2). Retrieved from https://doi.org/10.1145/3483382.3483389 doi: 10.1145/3483382.3483389

Carrino, C. P., Costa-jussà, M. R., & Fonollosa, J. A. R. (2020, May). Automatic Spanish translation of SQuAD dataset for multi-lingual question answering. In *Proceedings of the twelfth language resources and evaluation conference* (pp. 5515–5523). Marseille, France: European Language Resources Association. Retrieved from https://aclanthology.org/2020.lrec-1.677

Chandra, A., Fahrizain, A., Ibrahim, & Laufried, S. W. (2021). *A survey on non-English question answering dataset.* arXiv. Retrieved from https://arxiv.org/abs/2112.13634 doi: 10.48550/ARXIV.2112.13634

Chen, D., Fisch, A., Weston, J., & Bordes, A. (2017, July). Reading Wikipedia to answer open-domain questions. In *Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: Long papers)* (pp. 1870–1879). Vancouver, Canada: Association for Computational Linguistics. Retrieved from https://aclanthology.org/P17-1171 doi: 10.18653/v1/P17-1171

Chen, X., Lakhotia, K., Oguz, B., Gupta, A., Lewis, P., Peshterliev, S., … Yih, W.-t. (2022, December). Salient phrase aware dense retrieval: Can a dense retriever imitate a sparse one? In *Findings of the association for computational linguistics: Emnlp 2022* (pp. 250–262). Abu Dhabi, United Arab Emirates: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2022.findings-emnlp.19

Conneau, A., Khandelwal, K., Goyal, N., Chaudhary, V., Wenzek, G., Guzmán, F., … Stoyanov, V. (2020, July). Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th annual meeting of the association for computational linguistics* (pp. 8440–8451). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.acl-main.747 doi: 10.18653/v1/2020.acl-main.747

Conneau, A., & Lample, G. (2019). Cross-lingual language model pretraining. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, & R. Garnett (Eds.), *Advances in neural information processing systems* (Vol. 32). Curran Associates, Inc. Retrieved from https://proceedings.neurips.cc/paper/2019/file/c04c19c2c2474dbf5f7ac4372c5b9af1-Paper.pdf

Craswell, N., Mitra, B., Yilmaz, E., Campos, D., & Voorhees, E. M. (2020). *Overview of the TREC 2019 deep learning track.* arXiv. Retrieved from https://arxiv.org/abs/2003.07820 doi: 10.48550/ARXIV.2003.07820

Cui, Y., Liu, T., Che, W., Xiao, L., Chen, Z., Ma, W., … Hu, G. (2019, November). A span-extraction dataset for Chinese machine reading comprehension. In *Proceedings of the 2019 conference on empirical methods in natural language processing and the 9th international joint*

---

conference on natural language processing (emnlp-ijcnlp) (pp. 5883–5889). Hong Kong, China: Association for Computational Linguistics. Retrieved from https://aclanthology.org/D19-1600 doi: 10.18653/v1/D19-1600

Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019, June). BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 conference of the north American chapter of the association for computational linguistics: Human language technologies, volume 1 (long and short papers) (pp. 4171–4186). Minneapolis, Minnesota: Association for Computational Linguistics. Retrieved from https://aclanthology.org/N19-1423 doi: 10.18653/v1/N19-1423

d'Hoffschmidt, M., Belblidia, W., Heinrich, Q., Brendlé, T., & Vidal, M. (2020, November). FQuAD: French question answering dataset. In Findings of the association for computational linguistics: Emnlp 2020 (pp. 1193–1208). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.findings-emnlp.107 doi: 10.18653/v1/2020.findings-emnlp.107

Dunn, M., Sagun, L., Higgins, M., Guney, V. U., Cirik, V., & Cho, K. (2017). SearchQA: A new Q&A dataset augmented with context from a search engine. arXiv. Retrieved from https://arxiv.org/abs/1704.05179 doi: 10.48550/ARXIV.1704.05179

Efimov, P., Chertok, A., Boytsov, L., & Braslavski, P. (2020). SberQuAD – Russian reading comprehension dataset: Description and analysis. In International conference of the cross-language evaluation forum for European languages (pp. 3–15).

Guu, K., Lee, K., Tung, Z., Pasupat, P., & Chang, M. (2020, 13–18 Jul). Retrieval augmented language model pre-training. In H. D. III & A. Singh (Eds.), Proceedings of the 37th international conference on machine learning (Vol. 119, pp. 3929–3938). PMLR. Retrieved from https://proceedings.mlr.press/v119/guu20a.html

Hill, F., Bordes, A., Chopra, S., & Weston, J. (2016). The Goldilocks Principle: Reading children's books with explicit memory representations.. (Publisher Copyright: © ICLR 2016: San Juan, Puerto Rico. All Rights Reserved.; 4th International Conference on Learning Representations, ICLR 2016 ; Conference date: 02-05-2016 Through 04-05-2016)

Jiang, Y., Bordia, S., Zhong, Z., Dognin, C., Singh, M., & Bansal, M. (2020, November). HoVer: A dataset for many-hop fact extraction and claim verification. In Findings of the association for computational linguistics: Emnlp 2020 (pp. 3441–3460). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.findings-emnlp.309 doi: 10.18653/v1/2020.findings-emnlp.309

Johnson, J., Douze, M., & Jégou, H. (2019). Billion-scale similarity search with GPUs. IEEE Transactions on Big Data, 7(3), 535–547.

Joshi, M., Choi, E., Weld, D., & Zettlemoyer, L. (2017, July). TriviaQA: A large scale distantly supervised challenge dataset for reading comprehension. In Proceedings of the 55th annual meeting of the association for computational linguistics (volume 1: Long papers) (pp. 1601–1611). Vancouver, Canada: Association for Computational Linguistics. Retrieved from https://aclanthology.org/P17-1147 doi: 10.18653/v1/P17-1147

Kamath, A., Jia, R., & Liang, P. (2020, July). Selective question answering under domain shift. In Proceedings of the 58th annual meeting of the association for computational linguistics (pp. 5684–5696). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.acl-main.503 doi: 10.18653/v1/2020.acl-main.503

Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., … Yih, W.-t. (2020, November). Dense passage retrieval for open-domain question answering. In Proceedings of the 2020 conference on empirical methods in natural language processing (emnlp) (pp. 6769–6781). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.emnlp-main.550 doi: 10.18653/v1/2020.emnlp-main.550

Khattab, O., Potts, C., & Zaharia, M. (2021). Relevance-guided supervision for OpenQA with ColBERT. Transactions of the Association for Computational Linguistics, 9, 929–944. Retrieved from https://aclanthology.org/2021.tacl-1.55 doi: 10.1162/tacl_a_00405

Khattab, O., & Zaharia, M. (2020). ColBERT: Efficient and effective passage search via contextualized late interaction over BERT. In Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval (p. 39–48). New York, NY, USA: Association for Computing Machinery. Retrieved from https://doi.org/10.1145/3397271.3401075 doi: 10.1145/3397271.3401075

Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., … Petrov, S. (2019). Natural questions: A benchmark for question answering research. Transactions of the Association for Computational Linguistics, 7, 452–466. Retrieved from https://aclanthology.org/Q19-1026 doi: 10.1162/tacl_a_00276

Lee, H., Kedia, A., Lee, J., Paranjape, A., Manning, C., & Woo, K.-G. (2022, December). You only need one model for open-domain question answering. In Proceedings of the 2022 conference on empirical methods in natural language processing (pp. 3047–3060). Abu Dhabi, United Arab Emirates: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2022.emnlp-main.198

Lee, K., Chang, M.-W., & Toutanova, K. (2019, July). Latent retrieval for weakly supervised open domain question answering. In Proceedings of the 57th annual meeting of the association for computational linguistics (pp. 6086–6096). Florence, Italy: Association for Computational Linguistics. Retrieved from https://aclanthology.org/P19-1612 doi: 10.18653/v1/P19-1612

Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions, and reversals. Soviet physics doklady, 10(8), 707–710.

Lewis, P., Oguz, B., Rinott, R., Riedel, S., & Schwenk, H. (2020, July). MLQA: Evaluating cross-lingual extractive question answering. In Proceedings of the 58th annual meeting of the association for computational linguistics (pp. 7315–7330). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.acl-main.653 doi: 10.18653/v1/2020.acl-main.653

Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., … Kiela, D. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. In Proceedings of the 34th international conference on neural information processing systems (pp. 9459–9474). Red Hook, NY, USA: Curran Associates Inc.

Li, C., Yates, A., MacAvaney, S., He, B., & Sun, Y. (2020). PARADE: Passage representation aggregation for document reranking. arXiv. Retrieved from https://arxiv.org/abs/2008.09093 doi: 10.48550/ARXIV.2008.09093

Lim, S., Kim, M., & Lee, J. (2019). KorQuAD1.0: Korean QA dataset for machine reading comprehension. arXiv. Retrieved from https://arxiv.org/abs/1909.07005 doi: 10.48550/ARXIV.1909.07005

Lin, J., Crane, M., Trotman, A., Callan, J., Chattopadhyaya, I., Foley, J., … Vigna, S. (2016). Toward reproducible baselines: The open-source IR reproducibility challenge. In European conference on information retrieval (pp. 408–420).

Lin, J., Ma, X., Lin, S.-C., Yang, J.-H., Pradeep, R., & Nogueira, R. (2021). Pyserini: A Python toolkit for reproducible information retrieval research with sparse and dense representations. In Proceedings of the 44th annual international ACM SIGIR conference on research and development in information retrieval (SIGIR 2021) (pp. 2356–2362).

Liu, J., Lin, Y., Liu, Z., & Sun, M. (2019, July). XQA: A cross-lingual open-domain question answering dataset. In Proceedings of the 57th annual meeting of the association for computational linguistics (pp. 2358–2368). Florence, Italy: Association for Computational Linguistics. Retrieved from https://aclanthology.org/P19-1227 doi: 10.18653/v1/P19-1227

Möller, T., Risch, J., & Pietsch, M. (2021, November). GermanQuAD and GermanDPR: Improving non-English question answering and passage retrieval. In Proceedings of the 3rd workshop on machine reading for question answering (pp. 42–50). Punta Cana, Dominican Republic: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2021.mrqa-1.4 doi: 10.18653/v1/2021.mrqa-1.4

Mozannar, H., Maamary, E., El Hajal, K., & Hajj, H. (2019, August). Neural Arabic question answering. In Proceedings of the fourth arabic natural language processing workshop (pp. 108–118). Florence, Italy: Association for Computational Linguistics. Retrieved from https://aclanthology.org/W19-4612 doi: 10.18653/v1/W19-4612

Nguyen, T., Rosenberg, M., Song, X., Gao, J., Tiwary, S., Majumder, R., & Deng, L. (2016). MS MARCO: A human generated MAchine Reading COmprehension dataset. In T. R. Besold, A. Bordes, A. S. d'Avila Garcez, & G. Wayne (Eds.), *Proceedings of the workshop on cognitive computation: Integrating neural and symbolic approaches 2016 co-located with the 30th annual conference on neural information processing systems (NIPS 2016), barcelona, spain, december 9, 2016* (Vol. 1773). CEUR-WS.org. Retrieved from http://ceur-ws.org/Vol-1773/CoCoNIPS_2016_paper9.pdf

Nogueira, R., Jiang, Z., Pradeep, R., & Lin, J. (2020, November). Document ranking with a pretrained sequence-to-sequence model. In *Findings of the association for computational linguistics: Emnlp 2020* (pp. 708–718). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.findings-emnlp.63 doi: 10.18653/v1/2020.findings-emnlp.63

Noraset, T., Lowphansirikul, L., & Tuarob, S. (2021). WabiQA: A Wikipedia-Based Thai question-answering system. *Information Processing & Management*, 58(1), 102431. Retrieved from https://www.sciencedirect.com/science/article/pii/S0306457320309249 doi: https://doi.org/10.1016/j.ipm.2020.102431

Pradeep, R., Nogueira, R., & Lin, J. (2021). *The Expando-Mono-Duo design pattern for text ranking with pretrained sequence-to-sequence models*. arXiv. Retrieved from https://arxiv.org/abs/2101.05667 doi: 10.48550/ARXIV.2101.05667

Rajpurkar, P., Jia, R., & Liang, P. (2018, July). Know what you don't know: Unanswerable questions for SQuAD. In *Proceedings of the 56th annual meeting of the association for computational linguistics (volume 2: Short papers)* (pp. 784–789). Melbourne, Australia: Association for Computational Linguistics. Retrieved from https://aclanthology.org/P18-2124 doi: 10.18653/v1/P18-2124

Rajpurkar, P., Zhang, J., Lopyrev, K., & Liang, P. (2016, November). SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 conference on empirical methods in natural language processing* (pp. 2383–2392). Austin, Texas: Association for Computational Linguistics. Retrieved from https://aclanthology.org/D16-1264 doi: 10.18653/v1/D16-1264

Richardson, M., Burges, C. J., & Renshaw, E. (2013, October). MCTest: A challenge dataset for the open-domain machine comprehension of text. In *Proceedings of the 2013 conference on empirical methods in natural language processing* (pp. 193–203). Seattle, Washington, USA: Association for Computational Linguistics. Retrieved from https://aclanthology.org/D13-1020

Robertson, S., Walker, S., Jones, S., Hancock-Beaulieu, M. M., & Gatford, M. (1995, January). Okapi at TREC-3. In *Overview of the third Text REtrieval Conference (trec-3)* (Overview of the Third Text REtrieval Conference (TREC–3) ed., pp. 109–126). Gaithersburg, MD: NIST. Retrieved from https://www.microsoft.com/en-us/research/publication/okapi-at-trec-3/

Santhanam, K., Khattab, O., Potts, C., & Zaharia, M. (2022). PLAID: An efficient engine for late interaction retrieval. In *Proceedings of the 31st ACM international conference on information & knowledge management* (p. 1747–1756). New York, NY, USA: Association for Computing Machinery. Retrieved from https://doi.org/10.1145/3511808.3557325 doi: 10.1145/3511808.3557325

Santhanam, K., Khattab, O., Saad-Falcon, J., Potts, C., & Zaharia, M. (2022, July). ColBERTv2: Effective and efficient retrieval via lightweight late interaction. In *Proceedings of the 2022 conference of the north american chapter of the association for computational linguistics: Human language technologies* (pp. 3715–3734). Seattle, United States: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2022.naacl-main.272 doi: 10.18653/v1/2022.naacl-main.272

Schweter, S. (2020, April). *BERTurk - BERT models for Turkish*. Zenodo. Retrieved from https://doi.org/10.5281/zenodo.3770924 doi: 10.5281/zenodo.3770924

Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., … Rush, A. (2020, October). Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: System demonstrations* (pp. 38–45). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2020.emnlp-demos.6 doi: 10.18653/v1/2020.emnlp-demos.6

Xue, L., Constant, N., Roberts, A., Kale, M., Al-Rfou, R., Siddhant, A., … Raffel, C. (2021, June). mT5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 conference of the north american chapter of the association for computational linguistics: Human language technologies* (pp. 483–498). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2021.naacl-main.41 doi: 10.18653/v1/2021.naacl-main.41

Yang, S., & Seo, M. (2021, June). Designing a minimal retrieve-and-read system for open-domain question answering. In *Proceedings of the 2021 conference of the north american chapter of the association for computational linguistics: Human language technologies* (pp. 5856–5865). Online: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2021.naacl-main.468 doi: 10.18653/v1/2021.naacl-main.468

Yang, Y., Yih, W.-t., & Meek, C. (2015, September). WikiQA: A challenge dataset for open-domain question answering. In *Proceedings of the 2015 conference on empirical methods in natural language processing* (pp. 2013–2018). Lisbon, Portugal: Association for Computational Linguistics. Retrieved from https://aclanthology.org/D15-1237 doi: 10.18653/v1/D15-1237

Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W., Salakhutdinov, R., & Manning, C. D. (2018, October-November). HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Proceedings of the 2018 conference on empirical methods in natural language processing* (pp. 2369–2380). Brussels, Belgium: Association for Computational Linguistics. Retrieved from https://aclanthology.org/D18-1259 doi: 10.18653/v1/D18-1259

Yu, D., Zhu, C., Fang, Y., Yu, W., Wang, S., Xu, Y., … Zeng, M. (2022, May). KG-FiD: Infusing knowledge graph in fusion-in-decoder for open-domain question answering. In *Proceedings of the 60th annual meeting of the association for computational linguistics (volume 1: Long papers)* (pp. 4961–4974). Dublin, Ireland: Association for Computational Linguistics. Retrieved from https://aclanthology.org/2022.acl-long.340 doi: 10.18653/v1/2022.acl-long.340

Zhu, F., Lei, W., Wang, C., Zheng, J., Poria, S., & Chua, T.-S. (2021). *Retrieving and reading: A comprehensive survey on open-domain question answering*.