

# WebFAQ: A Multilingual Collection of Natural Q&A Datasets for Dense Retrieval

Michael Dinzinger  
michael.dinzinger@uni-passau.de  
University of Passau  
Passau, Germany

Laura Caspari  
laura.caspari@uni-passau.de  
University of Passau  
Passau, Germany

Kanishka Ghosh Dastidar  
kanishka.ghoshdastidar@uni-passau.de  
University of Passau  
Passau, Germany

Jelena Mitrović  
jelena.mitrovic@uni-passau.de  
University of Passau  
Passau, Germany

Michael Granitzer  
michael.granitzer@uni-passau.de  
University of Passau  
Passau, Germany

## ABSTRACT

We present WebFAQ, a large-scale collection of open-domain question answering datasets derived from FAQ-style schema.org annotations. In total, the data collection consists of 96 million natural question-answer (QA) pairs across 75 languages, including 47 million (49%) non-English samples. WebFAQ further serves as the foundation for 20 monolingual retrieval benchmarks with a total size of 11.2 million QA pairs (5.9 million non-English). These datasets are carefully curated through refined filtering and near-duplicate detection, yielding high-quality resources for training and evaluating multilingual dense retrieval models. To empirically confirm WebFAQ’s efficacy, we use the collected QAs to fine-tune an in-domain pretrained XLM-RoBERTa model. Through this process of dataset-specific fine-tuning, the model achieves significant retrieval performance gains, which generalize – beyond WebFAQ – to other multilingual retrieval benchmarks evaluated in zero-shot setting. Last but not least, we utilize WebFAQ to construct a set of QA-aligned bilingual corpora spanning over 1000 language pairs using state-of-the-art bitext mining and automated LLM-assessed translation evaluation. Due to our advanced, automated method of bitext dataset generation, the resulting bilingual corpora demonstrate higher translation quality compared to similar datasets. WebFAQ and all associated resources are publicly available on GitHub<sup>1</sup> and HuggingFace.<sup>2</sup>

## CCS CONCEPTS

• Information systems → Retrieval models and ranking; Data mining.

## KEYWORDS

Question Answering, Dense Retrieval, Multilingual Text Embedding, Cross-Lingual Information Retrieval

## 1 INTRODUCTION

The rapid adoption of structured data annotations, such as Microdata and JSON-LD formats, have significantly transformed the way information is presented and consumed online. Among these structured resources, FAQ (Frequently Asked Questions) pages provide

<sup>1</sup><https://github.com/padas-lab-de/webfaq>

<sup>2</sup><https://huggingface.co/PaDaS-Lab>

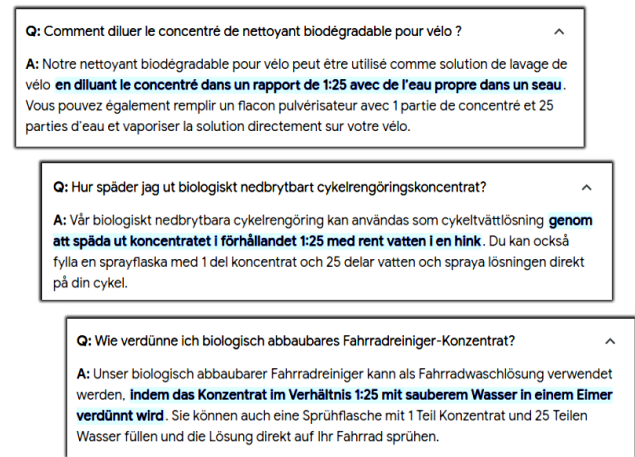


Figure 1: Exemplary FAQ entries across languages

an exceptionally rich collection of natural question-answer (QA) pairs across numerous topics and languages (see Figure 1). With the widespread use of schema.org annotations, search engines and web crawlers can extract these QA pairs efficiently, which provides a unique opportunity to leverage this data for research in the area of Natural Language Processing (NLP).

Furthermore, recent studies have demonstrated the effectiveness of Q&A datasets in training Open-Domain Question Answering (ODQA) systems and retrieval models [7, 15, 20, 32]. ODQA systems aim to accurately answer natural language questions across arbitrary domains, while retrieval models are designed to identify relevant documents or passages, and thus potential answers, from large-scale corpora. Particularly in this “open-retrieval” settings, various Q&A benchmarks, such as MS MARCO [3], HotpotQA [31] and Natural Questions (NQ) [19], have become essential cornerstones for the current broad scientific progress in Neural IR. For instance, these datasets are included in the training of numerous embedding models as bi-encoding dense retrievers [8, 27, 30, 34], and are commonly used in the evaluation of MTEB (Massive Text Embedding Benchmark) [24].

Recent initiatives such as the Massive Multilingual Text Embedding Benchmark (MMTEB) [9], an expansion of MTEB driven by the research community, demonstrate that there is an ongoing shift towards multilinguality and the diversification of evaluation tasks for text embedding models. The newly incorporated tasks include, among others, Cross-lingual STS (Semantic Text Similarity), Cross-lingual Retrieval as well as new forms of bitext mining, pointing towards the growth of relevance of diverse and multilingual resources. Despite these advancements, most large-scale Q&A datasets used for supervised fine-tuning and evaluation of Language Models (LMs) are restricted to English and/or Wikipedia-derived, limiting their applicability in truly open-domain multilingual scenarios.

Previous work by Huber et al. addresses these limitations by introducing CCQA [15], an open-domain question answering dataset derived from FAQ-style schema.org annotations. However, the dataset itself has not been publicly released, with only the data extraction code available. This limitation has hindered further research in leveraging FAQ-style QAs as a resource for fine-tuning or evaluation. To the best of our knowledge, no follow-up studies outside the original work have utilized CCQA, emphasizing the importance for publicly accessible datasets to support broader advancements in the field. In this regard, we introduce WebFAQ, a large-scale openly available resource for ODQA and open-domain Q&A retrieval, provided in a multitude of languages. WebFAQ encompasses 96 million QA pairs across 75 languages, including 47 million (49%) non-English samples. The extracted QAs are further labeled with topic and question type.

Building upon the WebFAQ base dataset, we introduce the following further key contributions:

### Monolingual retrieval datasets

We create 20 monolingual retrieval benchmarks with a total size of 11.2 million QA pairs (with 5.9 million non-English samples). To establish a clear notion of relevance within the retrieval datasets, refined filtering techniques have been applied, including near-duplicate detection and semantic consistency filtering for question-answer pairs. The result is a well curated resource to be used in the context of multilingual dense retrieval.

To provide initial baselines on our proposed benchmarks, we measure the performance of BM25 and three state-of-the-art embedding models in a zero-shot setting. These evaluations assess model performance without any dataset-specific adaptation. Separately, dataset-specific fine-tuning is applied to an in-domain pretrained XLM-RoBERTa model using WebFAQ data. The fine-tuned model achieves substantial performance gains, which generalize to other multilingual retrieval datasets. These improvements demonstrate that dense retrieval models benefit from exposure to WebFAQ data, leading to a concrete increase of model performance in open-domain Q&A retrieval.

### Bilingual datasets

This is an entirely novel contribution with respect to the work of Huber et al. Using state-of-the-art bitext mining techniques and automated LLM-based translation evaluation, we have constructed 1k bilingual datasets containing a total of 1.5 million aligned QAs (with each of the 1001 language pairs comprising at least 100 QA pairs). This effort takes advantage of a unique opportunity:

many FAQ pages exist in multiple languages, presenting the same questions and answers to users, but in translated form. Most notably, the aligned text sequences of our final bitext corpora exhibit high translation quality, even when compared to human-curated bitext datasets, demonstrating the effectiveness of our approach for automated bitext generation.

## 2 RELATED WORK

*QAs extracted from Common Crawl.* Our work builds upon the efforts of the Web Data Commons<sup>3</sup> (WDC) project, whose focus is the large-scale extraction of structured data from the Common Crawl<sup>4</sup> (CC) corpus. By systematically parsing and organizing schema.org annotations, embedded as JSON-LD, Microdata, RDF or Microformats, the WDC initiative provides the groundwork for various web data exploitations [6], such as the harvesting of FAQ pages. Accordingly, our work is indirectly inspired by CCQA [15], an open-domain question answering dataset from Meta AI, which also utilizes QA pairs extracted from Common Crawl. CCQA comprises approximately 55M unique QAs, including 24M English samples, gathered from 13 distinct web snapshots. In their paper, Huber et al. have demonstrated the effectiveness of CCQA for in-domain pre-training on tasks such as Closed-Book Question Answering (CBQA) and Passage Retrieval.

*Q&A datasets.* Beyond extracting QAs from Common Crawl, a variety of datasets from different sources have been developed to tackle distinct challenges in the context of question answering. WebQuestions [4], built using the Google Suggest API, focuses on entity-based queries, while ComplexWebQuestions [29] extends this scope by introducing broader and more complex questions requiring multi-step reasoning. ELI5 [10], sourced from Reddit, captures long-form explanatory Q&A. CoQA [25] shifts the focus to conversational Q&A, modeling multi-turn context progression. Quora [26] provides question pairs to evaluate semantic similarity rather than direct question answering. Additionally, the authors of PAQ (Probably Asked Questions) [20] generated questions for selected Wikipedia passages automatically, resulting in a large-scale collection of 65 million FAQ-style QA pairs.

Several large-scale Q&A datasets have also played a central role in advancing information retrieval research. For instance, HotpotQA [31] introduces multi-hop reasoning, requiring models to retrieve and integrate information from multiple documents. Natural Questions (NQ) [19] provides real user queries from Google Search along with corresponding Wikipedia passages, making it a widely used resource for training ODQA systems and Q&A retrievers. Similarly, MS MARCO [3], a large-scale dataset derived from real Bing search queries, has become a standard benchmark for passage ranking and retrieval tasks.

Existing multilingual Q&A datasets in the context of information retrieval include Mr. TyDi [32], MIRACL [33] and MLDR [8]. Mr. TyDi is a human-labeled retrieval dataset built on top of TyDi QA, which was sourced from Wikipedia and covers 11 typologically diverse languages. The limitations of Mr. TyDi, such as the methodology for annotating positive passages, were later addressed in

<sup>3</sup><https://webdatacommons.org/>

<sup>4</sup><https://commoncrawl.org/>

MIRACL. These two datasets as well as the recent MLDR (Multilingual Long-Document Retrieval), with a focus on more lengthy sample texts, were specifically crafted to facilitate the training and evaluation of multilingual retrieval systems.

*Cross-lingual datasets.* Significant IR datasets with queries and documents aligned across languages include mMARCO [5], the multilingual version of the MS MARCO passage retrieval dataset with the original English texts translated in 13 further languages. Further cross-lingual datasets are, e.g., MKQA [22], consisting of 10k QA pairs selected from NQ and translated from English into 25 additional languages, XQuAD [1], the cross-lingual adoption of SQuAD, and CLIRMatrix [28], currently the largest dataset in the context of Cross-Lingual IR (CLIR). CLIRMatrix includes 49M unique queries and 34 billion relevance labels across 139 languages with a massive bilingual corpus of 19,182 language combinations.

*Bitext mining.* The state-of-the-art approach for creating cross-lingual datasets is – beyond translation – automated sentence alignment via similarity search over text embeddings. Notable methods include LASER (Language-Agnostic SEntence Representations) [2] and its successor LaBSE (Language-Agnostic BERT Sentence Embeddings) [11]. LASER employs a sequence-to-sequence architecture for encoding, while LaBSE utilizes a Transformer-based model architecture. These advancements in multilingual sentence embeddings contribute to reduced error rates in cross-lingual similarity search and thus enable efficient bitext mining. Beyond bitext mining, recent research has explored the application of Large Language Models (LLMs) for automated translation evaluation. For example, Kocmi et al. introduce GEMBA [18], a GPT-based metric for translation evaluation, and demonstrate that LLMs can assess translation quality on par with human evaluators. These findings suggest that LLMs can play a critical role in both generating and validating cross-lingual datasets.

Notable datasets in the field of bitext mining include WMT 2019 [12], a massive dataset of 124M bitext pairs spanning nine language combinations, introduced as part of the reoccurring translation task at the Conference on Machine Translation (WMT). Another resource for bitext pairs is Tatoeba<sup>5</sup>, a community-driven collection of sentences and their translations provided in a multitude of languages, parts of which are easily accessible through HuggingFace<sup>6</sup>. Additionally, the BUCC 2018<sup>7</sup> dataset, originating from the 11th Workshop on Building and Using Comparable Corpora (BUCC2018), contains 35k bitext pairs in four language combinations.

### 3 DATA COLLECTION AND FILTERING

This section describes the methodology used to develop the WebFAQ Q&A dataset, including data collection, language detection, and topic and question type classification. Additionally, Section 3.3 outlines the refined filtering techniques employed to transform the raw QA corpus into a high-quality retrieval dataset with well-defined relevance relationships between queries and documents.

<sup>5</sup><https://tatoeba.org/>

<sup>6</sup><https://huggingface.co/datasets/Helsinki-NLP/tatoeba>

<sup>7</sup><https://comparable.limsi.fr/bucc2018/bucc2018-task.html>

### 3.1 Data Source

The raw data used to build WebFAQ originates from three Common Crawl snapshots, specifically the October dumps from 2022 to 2024. These web dumps were processed by the WDC initiative, which extracted structured schema.org annotations, including those marked with the FAQPage schema type. The extracted structured data is categorized into schema.org-specific subsets and made publicly available to support downstream research applications such as ours.

As noted by Huber et al. [15], the use of schema.org metadata requires additional effort from website creators, implying that the annotated QA pairs are intended for public use and are therefore more likely to be relevant, well-formed, and informative. Despite the inherent noisiness of web data, their findings confirm that the vast majority of web-mined QA pairs are sensible and answerable. This strongly supports their use as resource for large-scale dataset creation.

### 3.2 Processing

The structured FAQ data is first parsed to extract question-answer pairs while removing boilerplate text, quotation marks, and emojis. The extracted texts are then subjected to basic deduplication and filtering of incorrectly formatted samples.

*3.2.1 Language detection.* To classify the language of each QA pair, fastText [16, 17] is applied to the concatenated question-answer texts. In summary, the entire corpus comprises 75 languages with at least 1,000 samples each, while 49 languages appear in FAQ pages from at least 100 distinct websites<sup>8</sup> (see Table 1).

**Table 1: Distribution of languages**

Language	%	Language	%
eng (English)	51.2	ita (Italian)	2.7
deu (German)	6.9	jpn (Japanese)	2.6
spa (Spanish)	6.0	pol (Polish)	1.7
fra (French)	4.8	por (Portuguese)	1.7
rus (Russian)	3.8	tur (Turkish)	1.5
nld (Dutch)	2.8	Other	13.0

*3.2.2 Topic and Question Type Classification.* To further analyze and categorize the collected FAQ data, we fine-tune XLM-RoBERTa-base, a multilingual Transformer model, for two text classification tasks: (1) Topic Classification, (2) Question Type Classification. The resulting models are used to label the extracted QA pairs.

The respective training datasets of both classification tasks are composed of a sampled subset of QAs, which were automatically annotated using OpenAI’s GPT-4o-mini. To ensure sufficient training and validation data for each language, only those 49 languages were considered for labeling with at least 100 distinct websites contributing QA pairs. Additionally, to maximize diversity, the training sets include at most 1,000 QA pairs per language with no more than one sample per website, resulting in a final corpus of 37,383 samples.

<sup>8</sup>Websites are defined by their *origin*, which includes scheme, host and optionally port.

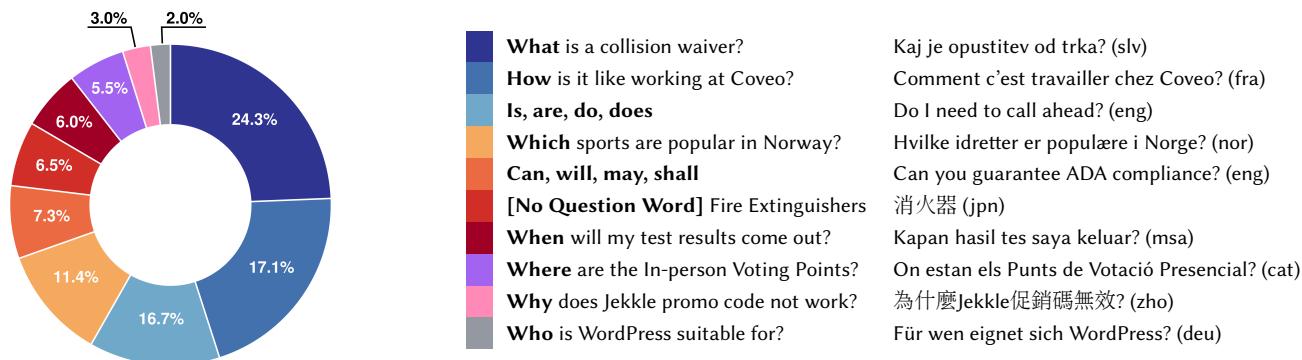


Figure 2: Distribution of Question types with examples

*Topic Classification.* The schema of tasks, provided to the classifier, is inspired by Curlie (formerly the Open Website Directory)<sup>9</sup>. However, because Curlie’s top-level categories focus more on commercial services and general information, we adapted its structure to better reflect the characteristics of FAQ-style question answering. Table 2 presents the final set of topics and their distribution across 96 million labeled QA pairs.

The dataset is split into 80% training, 10% validation, and 10% test sets. After two epochs of fine-tuning, the F1-score on the validation set is 81.29%.

Table 2: Distribution of topics

Topic	%
✈️ Traveling and Hospitality	34.1
🛒 Products and Commercial Services	19.8
🏥 Healthcare Services, Wellness, and Lifestyle	13.0
🎵 Entertainment, Recreation, and Leisure	9.7
🎓 Employment, Education, and Training	9.5
🏦 Banking, Financial Services, and Insurance	6.0
⚖️ Legal Services, Regulations, and Government	4.0
📍 General Information and Other	3.9

*Question Type Classification.* For question type classification, we adopt a schema similar to previous studies [10, 31]. While these works rely on string matching to determine question types, this approach is impractical in a vastly multilingual setting. Instead, we use XLM-RoBERTa to correctly map non-English questions to their English-equivalent question words. For example, the Slovenian question “Kaj je opustitev od trka?” is classified under “What”, without any custom definition of “Kaj” (*What*) as keyword to be matched, but merely relying on the capabilities of the pretrained model to map correctly across languages.

Figure 2 shows the distribution of question types among 96 million labeled QA pairs. Again, the dataset is split into 80% training, 10% validation, and 10% test, achieving an F1-score of 78.56% after three epochs of fine-tuning.

<sup>9</sup><https://curlie.org>

### 3.3 Refined Filtering

The initial processing pipeline produces a large-scale but relatively unclean dataset, containing ambiguous cases such as: (1) duplicate questions with different answers, (2) near-duplicate questions, and (3) generic QA pairs, describing those question or answer samples that require the context of their source webpage for meaningful interpretation. Additionally, some answers may be factually incorrect, even if they remain topically relevant to the question. Note that ensuring factual accuracy is however beyond the scope of this paper.

*Q&A and Retrieval datasets.* We publish the full Q&A dataset as an open resource<sup>10</sup>, allowing for custom filtering based on user needs. However, we find that datasets in the context of information retrieval have a more constrained notion of quality, as retrieval datasets must ensure that the concept of *relevance* between query and corresponding documents is well reflected and not blurred through insufficiencies during data preparation. This is particularly true for datasets like WebFAQ, which rely on sparse relevance judgments that can be obscured by near-duplicates and similar artifacts.

*Filtering.* We therefore propose three filtering techniques to eliminate the aforementioned ambiguities. The primary objective is to exclude QA pairs where a direct and unambiguous assignment between a question and its answer is not possible, coming at the cost of removing genuine QAs. This filtering involves:

- (1) Question-based Deduplication
- (2) Near-Duplicate Detection via Semantic Similarity Search
- (3) Question-Answer Semantic Consistency Filtering

Step (1) is rather straightforward, eliminating those duplicate questions with different answers, whereas steps (2) and (3) are more intricate and require the computation of text embeddings for both questions and answers. For this cause, we use the multilingual embedding model Jina (v3), as it performs well on Semantic Textual Similarity (STS) as well as retrieval tasks across multiple languages [27] and ranges in the top ranks of common leaderboards.

*Near-duplicate detection.* For (2), we establish a cosine similarity threshold  $\alpha$  among questions of same website origin. Figure 3a

<sup>10</sup><https://huggingface.co/datasets/PaDaS-Lab/webfaq>

presents exemplary filtered near-duplicate questions with cosine similarity  $\text{sim}_{Q_1, Q_2} > \alpha$ . For instance, question pairs that are merely reworded versions of each other, such as “Does Chrome have a free VPN?” and “Is there a free VPN for Chrome?”, are classified as near-duplicates. Since a *unique* relationship between one question and its corresponding answer cannot be established, these QA pairs are eliminated.

*Filtering based on semantic consistency between QAs.* For (3), we define a minimum cosine similarity threshold  $\beta$  between questions and their corresponding answers. QAs falling below this threshold are removed ( $\text{sim}_{Q,A} < \beta$ ). For instance, Figure 3b illustrates a case where the question lacks a specific entity (“Can I test *it* before purchasing?”). Without the context of the web page or the answer, it is unclear that “*it*” refers to an SMS service. As a result, the question by itself remains too generic, making it unrealistic to expect any state-of-the-art retriever to accurately determine its relevance to the corresponding answer.

Q <sub>1</sub> : Does Chrome have a free VPN? Q <sub>2</sub> : Is there a free VPN for Chrome?
(a) Near-duplicate detection ( $\alpha = 0.7$ )
Q: Can I test <i>it</i> before purchasing? A: Yes, you send 10 SMS campaigns to test our SMS service for free.
(b) QA Semantic Consistency Filtering ( $\beta = 0.5$ )

Figure 3: Examples of filtered questions/QA pairs

It is important to note that the goal of the proposed filtering steps is not to remove hard negatives, which can enhance a retrieval dataset – though this might be a side effect of the filtering process. Instead, the primary objective is to exclude QA pairs where state-of-the-art dense retrievers cannot reliably distinguish answers as either positive or negative.

The parameters  $\alpha = 0.7$  and  $\beta = 0.5$  were determined through manual inspection of 2,000 samples from the English subset, prioritizing dataset quality over sheer size. Applying those filtering steps to our collected QAs led to the creation of 20 large-scale retrieval corpora<sup>11</sup>, derived from the 20 largest language subsets, with a total size of 11.2 million QA pairs (5.9 million non-English). For the time being, we restricted the number to 20 subsets, such that each individual corpus contains at least around 80k distinct samples, ensuring a robust and diverse foundation for retrieval tasks.

## 4 EVALUATION

This section demonstrates the effectiveness of the WebFAQ retrieval datasets through experiments in an open-retrieval setting. First, we establish retrieval performance baselines by evaluating BM25 alongside three state-of-the-art embedding models and compare their results on WebFAQ to those obtained on other multilingual datasets.

Second, we fine-tune an in-domain pretrained language model using WebFAQ data. We hypothesize that this domain-specific fine-tuning step enhances the model’s understanding of relevance in retrieval tasks by leveraging supervised learning from WebFAQ’s QA pairs. As a result, substantial performance improvements across languages are expected, with potential benefits extending to other datasets in a true zero-shot setting.

All models are evaluated on the WebFAQ retrieval collection alongside two common multilingual retrieval datasets, namely Mr. TyDi [32] and MIRACL [33]. Both originate from Wikipedia-based sources and are widely used for training and evaluating multilingual retrievers. For MIRACL, we evaluate on the smaller hard negatives subset, which is part of the recent MMTEB leaderboard. Table 3 outlines retrieval performances on six languages – the intersection set of languages covered by WebFAQ, Mr. TyDi and MIRACL. The full results are available in Appendix A.

### 4.1 Initial Baselines

*BM25.* To begin with, the reported results in Table 3 include BM25, a strong traditional information retrieval method. Concretely, we use the Pyserini implementation [21] built on the Lucene search library to generate a BM25 index per language and dataset. The default settings of Pyserini were not changed, employing Lucene’s default whitespace analyzer and setting hyperparameter values of  $k_1 = 0.9$  and  $b = 0.4$ .

*SotA Embedding Models.* The baseline experiment further includes three state-of-the-art multilingual embedding models:

- *mGTE*: GTE-Multilingual-Base [34] (305M parameters)
- *mE5*: Multilingual-E5-Large-Instruct [30] (560M parameters)
- *Jina*: Jina (v3) [27] (572M parameters)

These models have demonstrated strong retrieval performance across various benchmarks and languages. They were selected as they had been – at the point of model selection – the three best ranked, public multilingual text embedding models on the MTEB leaderboard with less than 1B parameters.

*Interpretation.* The findings in Table 3 reveal inconsistent performance across datasets for the three aforementioned dense retrievers, while all of them perform relatively well compared to BM25. Jina, the largest model in terms of parameters, achieves the best results on WebFAQ but is outperformed on MIRACL and Mr. TyDi. This suggests that mGTE and mE5 likely benefit from prior exposure to the training splits of these datasets, indicating that their performance is influenced by the training data used during model development. According to the technical reports, mGTE and mE5 were trained using MIRACL and Mr. TyDi during supervised fine-tuning, whereas there is no explicit confirmation for Jina.

Beyond establishing baselines, the results further validate that WebFAQ’s retrieval dataset collection meets two key requirements. First, it is discriminative, effectively differentiating retrieval performance across models based on NDCG@10. Second, model rankings remain relatively stable across languages, demonstrating the dataset’s robustness and consistency in multilingual evaluation.

<sup>11</sup><https://huggingface.co/datasets/PaDaS-Lab/webfaq-retrieval>

**Table 3: Comparing retrieval performance on 3 multilingual datasets using NDCG@10 in %, including SotA embedding models and BM25 as baselines. Base and FT represent the pretrained XLM-RoBERTa model without/with fine-tuning on WebFAQ data. Hybrid combines BM25 and FT as described in Section 4.2. Bold font indicates top values w.r.t. the first 3 and last 4 rows.**

	WebFAQ						MIRACL (Hard Negatives)						Mr. Tydi					
	ara	eng	ind	jpn	kor	rus	ara	eng	ind	jpn	kor	rus	ara	eng	ind	jpn	kor	rus
mGTE	71.1	60.2	76.3	61.8	75.5	58.5	71.8	<b>54.9</b>	<b>50.5</b>	<b>66.4</b>	63.9	63.9	73.1	<b>57.1</b>	67.8	59.9	56.5	63.7
mE5	80.0	<b>67.1</b>	82.1	73.2	83.4	68.4	<b>75.6</b>	46.7	50.4	66.1	<b>65.8</b>	63.5	<b>76.5</b>	52.5	<b>71.3</b>	<b>61.9</b>	<b>59.5</b>	<b>65.4</b>
Jina	<b>85.5</b>	67.0	<b>85.2</b>	<b>75.2</b>	<b>87.4</b>	<b>72.6</b>	72.2	52.1	49.4	66.3	64.3	<b>65.4</b>	71.9	55.2	69.4	59.0	55.6	62.3
BM25	30.1	24.4	35.6	29.2	30.2	20.8	53.2	31.6	51.6	44.7	37.6	29.8	43.4	19.9	50.5	25.0	24.8	29.7
Base	61.0	49.7	74.6	56.5	73.2	50.1	36.1	30.6	28.1	27.6	39.3	29.6	30.5	18.1	33.2	16.0	27.7	22.5
FT	<b>74.1</b>	<b>59.0</b>	<b>82.7</b>	<b>68.8</b>	<b>81.6</b>	<b>61.5</b>	49.3	36.9	36.3	38.5	41.2	<b>38.5</b>	44.2	28.7	48.4	31.1	<b>34.8</b>	29.5
Hybrid	32.4	28.6	37.4	32.0	31.8	24.4	<b>61.7</b>	<b>50.8</b>	<b>57.8</b>	<b>51.3</b>	<b>43.2</b>	35.3	<b>54.4</b>	<b>31.7</b>	<b>58.5</b>	<b>37.8</b>	30.4	<b>36.0</b>

## 4.2 Model Fine-tuning

To empirically study its efficacy, we utilize WebFAQ to fine-tune a pretrained multilingual Transformer model. Specifically, we again employ XLM-RoBERTa-base as foundation model and apply task-specific pretraining with MS MARCO [3] data to create a base retrieval. While MS MARCO is a well-established training dataset for information retrieval, its corpus consists solely of English texts. Thus, this initial pretraining phase does not directly incorporate multilingual data but serves as an extensive in-domain warm-up before multilingual fine-tuning.

During this pretraining stage, Margin MSE [14] is employed as a contrastive loss function, using hard negatives scored by a MiniLM cross-encoder, provided by the Sentence-Transformers framework. The official Sentence-Transformers training script<sup>12</sup> is used for consistency and reproducibility. The model is trained for 30 epochs with 503k data samples and a batch size of 32.

After the task-specific pretraining on MS MARCO, we further fine-tune the model using the train splits of the WebFAQ retrieval datasets to better adapt it to the multilingual retrieval setting. The model undergoes an additional 20 epochs of training with 2,560,000 data samples and a batch size of 128, using Multiple Negatives Ranking Loss [13] with in-batch-negatives. The fine-tuning dataset is distributed relatively evenly across 20 languages, corresponding to the 20 retrieval corpora that we created. Each language contributes between 64,000 and 256,000 samples, depending on the size of the original train split. This ensures that the model is exposed to a diverse set of training samples.

## 4.3 Results

As shown in Table 3, fine-tuning with WebFAQ consistently improves retrieval performance. This effect is most evident in WebFAQ test splits, where the fine-tuned model achieves an average relative increase of 17% in NDCG@10 across six languages. Additionally, the model demonstrates improvements in zero-shot retrieval tasks, where evaluation datasets were not encountered during task-specific pretraining or fine-tuning. The observed gains on Mr. TyDi and MIRACL (Hard Negatives) suggest that exposure to WebFAQ’s diverse multilingual QA pairs enhances the model’s ability to judge

relevance more effectively. These findings confirm the hypothesis outlined in the first paragraph of Section 4.

*BM25.* When comparing the fine-tuned model with BM25, we find that our model outperforms the sparse retrieval approach on WebFAQ. However, on Mr. TyDi and MIRACL, the performance gap between BM25 and the fine-tuned dense retriever narrows. This outcome is expected, as these datasets are evaluated in a zero-shot setting. As noted by Zhang et al. [32], dense retrievers typically struggle in zero-shot retrieval scenarios due to their reliance on in-domain training for learning effective representations.

*Hybrid.* Lastly, we investigate whether our fine-tuned model retains useful relevance signals, even in cases where BM25 outperforms the dense retriever. To explore this, we employ a hybrid retrieval approach that integrates sparse (BM25) and dense (WebFAQ-trained XLM-RoBERTa) retrieval methods. The hybrid model first performs the regular retrieval step for both the sparse and dense model, retrieving the top 1000 documents, respectively. The retrieved sets are then merged by computing a combined similarity score for query  $q$  and document  $D_i$ , following Ma et al. [23]:

$$\lambda \text{sim}(q, D_i) + \text{BM25}(q, D_i) \quad (1)$$

where  $\text{sim}(q, D_i)$  represents the cosine similarity between  $q$  and  $D_i$  using XLM-RoBERTa, and BM25 provides the corresponding sparse retrieval score. Our implementation follows the original paper in setting  $\lambda = 1.1$ . If a document appears in only one retrieval set, its missing score is assigned as zero. Our findings indicate that the fine-tuned model continues to capture valuable relevance signals, enhancing retrieval performance beyond the BM25 baseline in several cases, even when BM25 alone outperforms the dense retriever.

## 5 BILINGUAL DATASETS

In addition to constructing a collection of monolingual retrieval corpora, we further utilize the collected QAs to create QA-aligned bitext datasets. Our approach exploits the fact that many websites provide FAQ pages in multiple languages, where the questions and answers are often literal translations of each other, since they aim to offer consistent services disregarding the users’ provenance.

<sup>12</sup>[https://sbert.net/examples/training/ms\\_marco/README.html](https://sbert.net/examples/training/ms_marco/README.html)

## 5.1 Methodology

The extraction of aligned text pairs employs state-of-the-art bitext mining techniques and utilizes cross-lingual similarity search based on sentence embeddings. Specifically, we use LaBSE (Language-Agnostic BERT Sentence Embeddings) [11] for bitext candidate generation. LaBSE generates vector representations of text sequences and enables cross-lingual similarity computation via cosine similarity scores. On top of this, the candidates’ translation quality is assessed using GEMBA, a GPT-based metric introduced by Kocmi et al [18]. Within our method, GEMBA is used to automatically annotate bitext candidates with a binary label (text pair is *accepted* / *rejected*) and find a reasonable cutoff value for cosine similarity between cross-lingual vector embeddings of candidate pairs. During alignment, the concatenation of a question and its corresponding answer is treated as a single unit, and bitexts are thus question-answer pairs.

The bitext mining process follows these key steps:

(1) **Elimination of near-duplicate QAs.**

The chosen approach of bitext mining includes cross-lingual similarity search, which is prone to similar text sequences within one language. Due to the large number of cross-lingual matchings of near-duplicates, the volume of generated bitext candidates had grown to an extremely large extent. To overcome this limitation, we applied the same elimination of near-duplicate QAs as described in Section 3.3.

(2) **Computation of LaBSE embeddings** for all text sequences (question + answer) of the 75 monolingual Q&A datasets.

(3) **For each language pair: Computation of cosine similarity values between cross-lingual text pairs.**

The matching of text pairs is conveyed across languages and restricted to website level. A text of language  $L_1$  is thus only matched with those texts of language  $L_2$  originating from the same website. This covers also those multilingual FAQs exposed under different language-specific URL paths<sup>13</sup>. Furthermore, to limit the volume of candidates to a reasonable size, all text pairs with cosine similarity below a tolerance threshold of  $s = 0.80$  are discarded.

(4) **For each language pair: Annotation of a sampled subset of text pairs as *accepted* or *rejected* translations.**

The remaining bitext candidates are randomly sampled with a probability of 0.1% and, if selected, automatically annotated with regards to their translation quality. This task is performed using GPT-4o-mini evaluating according to the GEMBA metric, which assigns a direct assessment (DA) score between 0 and 100. As observed by Kocmi et al, LLM-based DA scores are not uniformly distributed, as the model exhibits preference biases in assigning specific values.

Through manual inspection, we find that scores between 80 and 84 are rarely used (1%), whereas 85 (16%), 90 (18%) and 95 (51%) are frequently assigned by GPT-4o-mini to indicate high-confidence translations. Based on this observation, we classify bitexts with scores  $\geq 85$  as correct translations, while those below 85 are considered incorrect.

<sup>13</sup>Example of English FAQ page with German counterpart under a different URL path: <https://scubanana.es/faq/> and <https://scubanana.es/de/haeufig-gestellte-fragen/>

(5) **Definition of cosine similarity threshold.**

The LLM-generated annotations are used to define a final threshold with respect to a desired level of precision. For the concrete case of our dataset, we find that a threshold of  $s = 0.90$  provides a reasonable balance between quality and recall, yielding 95% correct translations according to the LLM-based assessment.

This method yields a total amount of 1.5 million aligned QA pairs across 1176 language pairs. For publishing our collection of bilingual datasets<sup>14</sup>, all language pairs featuring less than 100 samples were omitted, leaving 1,487,328 aligned texts across 1001 combinations. The most frequent language pairs are German-English (37,348 text pairs), followed by English-French (37,208) and English-Spanish (35,446).

**Table 4: Comparison of average translation quality scores for different bitext corpora. Pooled variance captures the variance of scores across language pairs.**

	Size	#Pairs	Translation quality Avg. score ( $\pm$ std)	Pooled variance
WMT 2019	124M	9	85.1 $\pm$ 13.4	173.1
BUCC 2018	35.0k	4	85.9 $\pm$ 12.6	158.7
Tatoeba	88.9k	113	77.3 $\pm$ 30.2	517.4
mMARCO	8.8M	182	86.6 $\pm$ 11.2	123.9
WebFAQ (Bitexts)	1.5M	1001	<b>91.0 <math>\pm</math> 8.6</b>	<b>73.1</b>

## 5.2 Quality Evaluation

To validate the presented bitext mining method, the resulting dataset is compared to existing parallel corpora in terms of average translation quality. Therefore, a random sample of 20,000 bitext pairs is selected from each corpus and evaluated based on the GEMBA metric, assigning scores on a 0 to 100 scale. The final average scores per bitext dataset are listed in Table 4.

Among the evaluated bitext datasets, WMT 2019 stands out as the largest, containing 124 million bitext pairs across nine language pairs. The corpus extracted from WebFAQ is two orders of magnitude smaller, but yet surpasses both BUCC 2018 and Tatoeba in terms of dataset size. Regarding translation quality, our dataset achieves the highest average GEMBA score and the lowest standard deviation. This suggests that the cosine similarity thresholding for cross-lingual sentence embeddings, described in the methodology, effectively produces a large-scale bitext corpus that remains comparable to existing resources in terms of translation accuracy. The six-point gap in scores compared to WMT 2019 and BUCC 2018 supports this claim. We believe that this holds true despite potential biases being introduced by using the same approach for LLM-based translation assessment for both dataset generation and evaluation. Additionally, our findings are corroborated by the comparison with mMARCO, an automatically translated retrieval dataset, which was created with a focus on translation quality, yet exhibits significantly

<sup>14</sup><https://huggingface.co/datasets/PaDaS-Lab/webfaq-bitexts>



lower GEMBA scores than our dataset. Overall, these findings confirm the effectiveness of our approach combining cross-lingual sentence embeddings with automated quality assessment in mining cross-lingual QA pairs.

## 6 DISCUSSION

The following section discusses several aspects that should be considered when using WebFAQ for research and evaluation. The first aspect is the prevalence of *sparse relevance judgments* in the retrieval datasets (Section 3.3). Generally speaking, this prevalence entails that no claim can be made about the exhaustiveness of the dataset’s relevance annotations. Even though this is a widely accepted practice in IR research (see e.g., MS MARCO [3] and Mr. Tydi [32]), one should keep in mind that the “qrels” of our dataset provide only one good answer to a question, and leave out many more answers that are also potentially relevant. This is particularly true for rather generic QAs, such as “What is the average time for my package to be delivered?”, which is asked and answered by different delivery service FAQs in a similar fashion.

Another discussion point arises with the the rigorous elimination of near-duplicate QA pairs, which primarily affects questions that revolve around entities. While we found this step helpful for removing many QAs one could call “spam” or low-quality, it may introduce an unintended bias by filtering out subtly different entity-based queries. For instance, the two questions “What is CBD?” and “What is CBDa?” are treated as near-duplicates and are thus removed, despite their subjective distinctions. This decision impacts analysis of retrieval models, as sparse methods (e.g., BM25) and dense models handle such minor lexical variations differently. The dataset, therefore, may not fully reflect scenarios where entity-level distinctions play a significant role in retrieval performance.

Finally, the quality of multilingual and cross-lingual datasets is inherently constrained by the accuracy of the language identification model. Through manual inspection, we observed systematic errors where the language detection system struggles with queries or answers containing long named entities from multiple languages. Additionally, in some edge cases, the question and answer are formulated in different languages, which the current identification method does not explicitly address, as it applies detection to the concatenated QA pair rather than handling them separately.

## 7 CONCLUSION

This paper introduces WebFAQ, a large-scale, multilingual resource for open-domain question answering and Q&A retrieval. Derived from FAQ-style schema.org annotations, WebFAQ encompasses 96 million natural QA pairs across 75 languages. With our work, we offer publicly available data that surpasses previous efforts in scale and linguistic diversity, and promises to surpass them in terms of usability and practical impact.

To demonstrate the utility of WebFAQ, we curated 20 monolingual retrieval benchmarks, applying advanced filtering techniques to ensure an unblurred notion of relevance reflected in the datasets. Our empirical results show that fine-tuning a text embedding model on WebFAQ’s QAs leads to significant retrieval performance improvements, with gains that generalize to other multilingual retrieval tasks. Additionally, we introduced a novel set of bilingual

Q&A datasets, constructed through state-of-the-art bitext mining and automated LLM-based translation evaluation. These bilingual corpora exhibit superior translation quality compared to existing datasets, which demonstrates the effectiveness of our automated approach for bitext generation.

With WebFAQ, we provide a publicly available, large-scale resource that facilitates the training and evaluation of multilingual retrieval models. We hope this work contributes to the diversification of evaluation tasks and training datasets and, in this regard, paves the way for broader advancements in open-domain multilingual Q&A retrieval.

## ACKNOWLEDGMENTS

This work has received funding from the European Union’s Horizon Europe research and innovation program under grant agreement No 101070014 (OpenWebSearch.EU, <https://doi.org/10.3030/101070014>).

## REFERENCES

- [1] Mikel Artetxe, Sebastian Ruder, and Dani Yogatama. 2020. On the Cross-lingual Transferability of Monolingual Representations. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 4623–4637. <https://doi.org/10.18653/v1/2020.acl-main.421>
- [2] Mikel Artetxe and Holger Schwenk. 2019. Massively Multilingual Sentence Embeddings for Zero-Shot Cross-Lingual Transfer and Beyond. *Transactions of the Association for Computational Linguistics* 7 (2019), 597–610. [https://doi.org/10.1162/tacl\\_a\\_00288](https://doi.org/10.1162/tacl_a_00288)
- [3] Payal Bajaj, Daniel Campos, Nick Craswell, Li Deng, Jianfeng Gao, Xiaodong Liu, Rangan Majumder, Andrew McNamara, Bhaskar Mitra, Tri Nguyen, Mir Rosenberg, Xia Song, Alina Stoica, Saurabh Tiwary, and Tong Wang. 2016. MS MARCO: A Human Generated Machine Reading Comprehension Dataset. <https://doi.org/10.48550/ARXIV.1611.09268>
- [4] Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. Semantic Parsing on Freebase from Question-Answer Pairs. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard (Eds.). Association for Computational Linguistics, Seattle, Washington, USA, 1533–1544. <https://aclanthology.org/D13-1160>
- [5] Luiz Bonifacio, Vitor Jeronymo, Hugo Queiroz Abonizio, Israel Campiotti, Marzieh Fadaee, Roberto Lotufo, and Rodrigo Nogueira. 2021. mMARCO: A Multilingual Version of the MS MARCO Passage Ranking Dataset. <https://doi.org/10.48550/ARXIV.2108.13897>
- [6] Alexander Brinkmann, Anna Primpeli, and Christian Bizer. 2023. The Web Data Commons Schema.org Data Set Series (*WWW ’23 Companion*). Association for Computing Machinery, New York, NY, USA, 136–139. <https://doi.org/10.1145/3543873.3587331>
- [7] Stefano Campese, Ivano Lauriola, and Alessandro Moschitti. 2023. QUADRO: Dataset and Models for Question-Answer Database Retrieval. In *Findings of the Association for Computational Linguistics: EMNLP 2023*. Association for Computational Linguistics, Singapore, 15573–15587. <https://doi.org/10.18653/v1/2023.findings-emnlp.1042>
- [8] Jianlyu Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. 2024. M3-Embedding: Multi-Linguality, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation. In *Findings of the Association for Computational Linguistics: ACL 2024*. Association for Computational Linguistics, Bangkok, Thailand, 2318–2335. <https://doi.org/10.18653/v1/2024.findings-acl.137>
- [9] Kenneth Enevoldsen et al. 2025. MMTEB: Massive Multilingual Text Embedding Benchmark. In *The 13th International Conference on Learning Representations*. <https://openreview.net/forum?id=z3pfz4VCV>
- [10] Angela Fan, Yacine Jernite, Ethan Perez, David Grangier, Jason Weston, and Michael Auli. 2019. ELI5: Long Form Question Answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, Anna Korhonen, David Traum, and Lluís Màrquez (Eds.). Association for Computational Linguistics, Florence, Italy, 3558–3567. <https://doi.org/10.18653/v1/P19-1346>
- [11] Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT Sentence Embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long*



- Papers*). Association for Computational Linguistics, Dublin, Ireland, 878–891. <https://doi.org/10.18653/v1/2022.acl-long.62>
- [12] Wikimedia Foundation. [n. d.]. *ACL 2019 Fourth Conference on Machine Translation (WMT19), Shared Task: Machine Translation of News*. <http://www.statmt.org/wmt19/translation-task.html>
- [13] Matthew Henderson, Rami Al-Rfou, Brian Strope, Yun-hsuan Sung, Laszlo Lukacs, Ruiqi Guo, Sanjiv Kumar, Balint Miklos, and Ray Kurzweil. 2017. Efficient Natural Language Response Suggestion for Smart Reply. <https://doi.org/10.48550/ARXIV.1705.00652>
- [14] Sebastian Hofstätter, Sophia Althammer, Michael Schröder, Mete Sertkan, and Allan Hanbury. 2020. Improving Efficient Neural Ranking Models with Cross-Architecture Knowledge Distillation. <https://doi.org/10.48550/ARXIV.2010.02666>
- [15] Patrick Huber, Armen Aghajanyan, Barlas Oguz, Dmytro Okhonko, Scott Yih, Sonal Gupta, and Xilun Chen. 2022. CCQA: A New Web-Scale Question Answering Dataset for Model Pre-Training. In *Findings of the Association for Computational Linguistics: NAACL 2022*, Marine Carpuat, Marie-Catherine de Marneffe, and Ivan Vladimir Meza Ruiz (Eds.). Association for Computational Linguistics, Seattle, United States, 2402–2420. <https://doi.org/10.18653/v1/2022.findings-naacl.184>
- [16] Armand Joulin, Edouard Grave, Piotr Bojanowski, Matthijs Douze, Herve Jégou, and Tomas Mikolov. 2016. FastText.zip: Compressing text classification models. <https://doi.org/10.48550/ARXIV.1612.03651>
- [17] Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. Bag of Tricks for Efficient Text Classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*. Association for Computational Linguistics, Valencia, Spain, 427–431. <https://aclanthology.org/E17-2068/>
- [18] Tom Kocmi and Christian Federmann. 2023. Large Language Models Are State-of-the-Art Evaluators of Translation Quality. In *Proceedings of the 24th Annual Conference of the European Association for Machine Translation*. European Association for Machine Translation, Tampere, Finland, 193–203. <https://aclanthology.org/2023.eamt-1.19/>
- [19] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, Kristina Toutanova, Llion Jones, Matthew Kelcey, Ming-Wei Chang, Andrew M. Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natural Questions: A Benchmark for Question Answering Research. *Transactions of the Association for Computational Linguistics* 7 (2019), 452–466. [https://doi.org/10.1162/tacl\\_a\\_00276](https://doi.org/10.1162/tacl_a_00276)
- [20] Patrick Lewis, Yuxiang Wu, Linqing Liu, Pasquale Minervini, Heinrich Küttler, Aleksandra Piktus, Pontus Stenetorp, and Sebastian Riedel. 2021. PAQ: 65 Million Probably-Asked Questions and What You Can Do With Them. *Transactions of the Association for Computational Linguistics* 9 (2021), 1098–1115. [https://doi.org/10.1162/tacl\\_a\\_00415](https://doi.org/10.1162/tacl_a_00415)
- [21] Jimmy Lin, Xueguang Ma, Sheng-Chieh Lin, Jheng-Hong Yang, Ronak Pradeep, and Rodrigo Nogueira. 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)*, 2356–2362.
- [22] Shayne Longpre, Yi Lu, and Joachim Daiber. 2021. MKQA: A Linguistically Diverse Benchmark for Multilingual Open Domain Question Answering. *Transactions of the Association for Computational Linguistics* 9 (2021), 1389–1406. [https://doi.org/10.1162/tacl\\_a\\_00433](https://doi.org/10.1162/tacl_a_00433)
- [23] Xueguang Ma, Kai Sun, Ronak Pradeep, and Jimmy Lin. 2021. A Replication Study of Dense Passage Retriever. arXiv:2104.05740 [cs.CL] <https://arxiv.org/abs/2104.05740>
- [24] Niklas Muennighoff, Nouamane Tazi, Loic Magne, and Nils Reimers. 2023. MTEB: Massive Text Embedding Benchmark. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics. <https://doi.org/10.18653/v1/2023.eacl-main.148>
- [25] Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A Conversational Question Answering Challenge. *Transactions of the Association for Computational Linguistics* 7 (2019), 249–266. [https://doi.org/10.1162/tacl\\_a\\_00266](https://doi.org/10.1162/tacl_a_00266)
- [26] Lakshay Sharma, Laura Graesser, Nikita Nangia, and Utku Evci. 2019. Natural Language Understanding with the Quora Question Pairs Dataset. <https://doi.org/10.48550/ARXIV.1907.01041>
- [27] Saba Sturua, Isabelle Mohr, Mohammad Kalim Akram, Michael Günther, Bo Wang, Markus Krimmel, Feng Wang, Georgios Mastrapas, Andreas Koukounas, Andreas Koukounas, Nan Wang, and Han Xiao. 2024. jina-embeddings-v3: Multilingual Embeddings With Task LoRA. arXiv:2409.10173 [cs.CL] <https://arxiv.org/abs/2409.10173>
- [28] Shuo Sun and Kevin Duh. 2020. CLIRMatrix: A massively large collection of bilingual and multilingual datasets for Cross-Lingual Information Retrieval. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online, 4160–4170. <https://doi.org/10.18653/v1/2020.emnlp-main.340>
- [29] Alon Talmor and Jonathan Berant. 2018. The Web as a Knowledge-Base for Answering Complex Questions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*. Association for Computational Linguistics, New Orleans, Louisiana, 641–651. <https://doi.org/10.18653/v1/N18-1059>
- [30] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. 2024. Multilingual E5 Text Embeddings: A Technical Report. <https://doi.org/10.48550/ARXIV.2402.05672>
- [31] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A Dataset for Diverse, Explainable Multi-hop Question Answering. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Brussels, Belgium, 2369–2380. <https://doi.org/10.18653/v1/D18-1259>
- [32] Xinyu Zhang, Xueguang Ma, Peng Shi, and Jimmy Lin. 2021. Mr. TyDi: A Multilingual Benchmark for Dense Retrieval. In *Proceedings of the 1st Workshop on Multilingual Representation Learning*. Association for Computational Linguistics, Punta Cana, Dominican Republic, 127–137. <https://doi.org/10.18653/v1/2021.mrl-1.12>
- [33] Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2023. MIRACL: A Multilingual Retrieval Dataset Covering 18 Diverse Languages. *Transactions of the Association for Computational Linguistics* 11 (2023), 1114–1131. [https://doi.org/10.1162/tacl\\_a\\_00595](https://doi.org/10.1162/tacl_a_00595)
- [34] Xin Zhang, Yanzhao Zhang, Dingkun Long, Wen Xie, Ziqi Dai, Jialong Tang, Huan Lin, Baosong Yang, Pengjun Xie, Fei Huang, Meishan Zhang, Wenjie Li, and Min Zhang. 2024. mGTE: Generalized Long-Context Text Representation and Reranking Models for Multilingual Text Retrieval. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*. Association for Computational Linguistics, Miami, Florida, US, 1393–1412. <https://doi.org/10.18653/v1/2024.emnlp-industry.103>

## A APPENDIX: FULL RESULTS

The below tables list the NDCG@10 in % results across models for all 20 WebFAQ retrieval corpora, the 18 languages included in MIRACL (Hard Negatives) and the 11 languages of Mr. Tydi. The top three rows show the performance of three state-of-the-art dense embedding models: GTE-Multilingual-Base (mGTE), Multilingual-E5-Large-Instruct (mE5) and Jina (v3). The bottom four rows of each table compare the results of BM25, a popular traditional sparse embedding model, which we use as a baseline, to our pretrained XML-RoBERTa models. Here, **Base** refers to the model with in-domain pretraining on MS MARCO, while **FT** refers to the model that was additionally fine-tuned on WebFAQ data. **Hybrid** combines BM25 and FT as described in Section 4.2. Bold font indicates top values w.r.t. the first 3 and last 4 rows.

**Table 5: Retrieval performance across models in NDCG@10 (%) on WebFAQ**

	ara	dan	deu	eng	fas	fra	hin	ind	ita	jpn	kor	nld	pol	por	rus	spa	swe	tur	vie	zho
mGTE	69.5	75.0	57.6	58.9	63.2	63.6	76.3	75.7	68.1	61.2	74.4	63.4	66.1	70.1	57.0	66.1	70.5	59.1	76.2	84.1
mE5	79.2	82.5	70.1	66.2	74.1	72.0	82.5	81.9	76.9	72.7	82.7	76.0	75.5	78.7	66.8	74.3	77.2	65.8	83.7	87.4
Jina	<b>84.6</b>	<b>87.5</b>	<b>74.2</b>	<b>66.4</b>	<b>78.7</b>	<b>76.5</b>	<b>86.6</b>	<b>84.8</b>	<b>82.1</b>	<b>74.5</b>	<b>85.8</b>	<b>80.1</b>	<b>81.5</b>	<b>83.7</b>	<b>71.4</b>	<b>78.3</b>	<b>82.8</b>	<b>70.6</b>	<b>86.9</b>	<b>88.9</b>
BM25	29.7	34.0	25.1	24.4	30.1	27.4	34.3	36.1	30.4	29.0	29.6	29.7	26.7	32.7	21.2	29.0	30.6	25.9	34.4	35.7
Base	61.0	74.7	53.1	49.7	63.2	54.9	70.1	74.6	60.0	56.5	73.2	61.2	61.5	65.3	50.1	57.1	67.6	49.8	70.2	73.5
FT	<b>74.1</b>	<b>83.7</b>	<b>66.3</b>	<b>59.0</b>	<b>74.4</b>	<b>69.1</b>	<b>79.6</b>	<b>82.7</b>	<b>73.8</b>	<b>68.8</b>	<b>81.6</b>	<b>74.6</b>	<b>74.2</b>	<b>78.3</b>	<b>61.5</b>	<b>69.6</b>	<b>78.8</b>	<b>66.1</b>	<b>80.7</b>	<b>84.8</b>
Hybrid	32.4	36.5	29.4	28.6	33.4	32.0	36.3	37.4	33.2	32.0	31.8	33.4	29.4	34.7	24.4	32.5	33.3	33.6	36.4	38.0

**Table 6: Retrieval performance across models in NDCG@10 (%) on MIRACL (Hard Negatives)**

	ara	ben	deu	eng	fas	fin	fra	hin	ind	jpn	kor	rus	spa	swa	tel	tha	yor	zho
mGTE	71.8	<b>72.7</b>	51.2	<b>54.9</b>	52.2	73.5	<b>55.2</b>	52.3	<b>50.5</b>	<b>66.4</b>	63.9	63.9	<b>53.0</b>	69.9	<b>83.1</b>	<b>74.5</b>	58.3	<b>62.3</b>
mE5	<b>75.6</b>	72.1	50.8	46.7	<b>56.4</b>	<b>74.6</b>	48.8	<b>58.2</b>	50.4	66.1	<b>65.8</b>	63.5	49.1	<b>71.7</b>	82.6	77.2	<b>59.9</b>	53.4
Jina	72.2	71.9	<b>53.3</b>	52.1	53.9	71.9	54.9	56.7	49.4	66.3	64.3	<b>65.4</b>	50.7	59.4	81.8	76.0	48.1	56.9
BM25	53.2	53.9	21.8	31.6	34.1	51.8	21.9	48.6	51.6	44.7	37.6	29.8	36.3	53.0	50.0	27.6	<b>64.4</b>	33.9
Base	36.1	32.0	32.8	30.6	36.2	42.1	25.9	31.9	28.1	27.6	39.3	29.6	30.2	23.1	34.6	41.0	11.3	26.9
FT	49.3	46.0	<b>35.8</b>	36.9	<b>43.0</b>	55.9	<b>39.5</b>	37.5	36.3	38.5	41.2	<b>38.5</b>	40.1	41.4	<b>59.6</b>	<b>52.1</b>	28.6	36.7
Hybrid	<b>61.7</b>	<b>61.1</b>	30.4	<b>50.8</b>	37.9	<b>62.0</b>	29.0	<b>54.8</b>	<b>57.8</b>	<b>51.3</b>	<b>43.2</b>	35.2	<b>43.8</b>	<b>57.9</b>	59.0	34.0	56.7	<b>40.3</b>

**Table 7: Retrieval performance across models in NDCG@10 (%) on Mr. Tydi**

	ara	ben	eng	fin	ind	jpn	kor	rus	swa	tel	tha
mGTE	73.1	<b>73.6</b>	<b>57.1</b>	63.9	67.8	59.9	56.5	63.7	69.6	<b>89.2</b>	72.7
mE5	<b>76.5</b>	72.6	52.5	<b>66.4</b>	<b>71.2</b>	<b>61.9</b>	<b>59.5</b>	<b>65.4</b>	<b>73.0</b>	87.4	<b>77.3</b>
Jina	71.9	71.3	55.2	63.2	69.4	59.0	55.6	62.3	61.6	87.0	74.4
BM25	43.4	50.2	19.9	34.0	50.5	25.0	24.8	29.7	53.9	51.7	27.9
Base	30.5	29.9	18.1	24.4	33.2	16.0	27.7	22.5	27.7	32.0	35.7
FT	44.2	45.2	28.7	38.7	48.4	31.1	<b>34.8</b>	29.5	46.7	62.2	<b>48.9</b>
Hybrid	<b>54.4</b>	<b>51.6</b>	<b>31.7</b>	<b>45.1</b>	<b>58.5</b>	<b>37.8</b>	30.4	<b>36.0</b>	<b>60.8</b>	<b>66.4</b>	35.7