

Leveraging LLMs for Synthesizing Training Data Across Many Languages in Multilingual Dense Retrieval

Anonymous ACL submission

Abstract

There has been limited success for dense retrieval models in multilingual retrieval, due to uneven and scarce training data available across multiple languages. Synthetic training data generation is promising (e.g., InPars or Promptagator), but has been investigated only for English. Therefore, to study model capabilities across both cross-lingual and monolingual retrieval tasks, we develop JUMP-IR,¹ a synthetic retrieval training dataset containing 33 (high to very-low resource) languages for fine-tuning multilingual dense retrievers without requiring any human supervision. To construct JUMP-IR, we propose SAP (*summarize-then-ask prompting*), where the large language model (LLM) generates a textual summary prior to the query generation step. SAP assists the LLM in generating informative queries in the target language. Using JUMP-IR, we explore synthetic fine-tuning of multilingual dense retrieval models and evaluate them robustly on three retrieval benchmarks: XOR-Retrieve (cross-lingual), XTREME-UP (cross-lingual) and MIRACL (monolingual). Our models, called JUMP-X, are competitive with human-supervised dense retrieval models, e.g., mContriever, finding that JUMP-IR can cheaply substitute for expensive human-labeled retrieval training data.²

1 Introduction

Dense retrieval models have demonstrated impressive performance in ad-hoc information retrieval (IR) tasks, e.g., web search, outperforming traditional retrieval systems such as BM25 (Karpukhin et al., 2020; Lin et al., 2021; Ni et al., 2022; Nee-lakantan et al., 2022, *inter alia*). A major reason for success lies in the availability of large-scale supervised training datasets in English, such

Dataset	Q Gen.	Cross.	Mono.	# L	Domain	# Train
NeuCLIR	Human	EN→L	L→L	3	News (hc4)	×
MKQA	Human	L→EN	×	26	Wikipedia	10K
mMARCO	Translate	×	L→L	13	MS MARCO	533K
Mr.TyDI	Human	×	L→L	11	Wikipedia	49K
MIRACL	Human	×	L→L	18	Wikipedia	726K
JH-POLO	GPT-3	EN→L	×	3	News (hc4)	78K
JUMP-IR	PaLM 2	L→EN	L→L	33	Wikipedia	28M

Table 1: Existing datasets only contain up to a few thousand training pairs, as scaling human annotations is both expensive and cumbersome. In our work, we construct JUMP-IR, a “synthetic” multilingual dataset with 28 million PaLM 2-generated training pairs across 33 languages; (Q Gen.) denotes the query generation task; (Cross. and Mono.) denotes the retrieval task and (query→document) language pair; (# L and # Train) denotes the language count and available training pairs.

as MS MARCO (Nguyen et al., 2016) or NQ (Kwiatkowski et al., 2019), and coupled with effective training strategies, such as custom hard-negative mining (Xiong et al., 2021; Lin et al., 2023), or teacher distillation (Hofstätter et al., 2021; Ren et al., 2021).

However, there is a limited exploration of dense retrieval models in multilingual retrieval,³ due to uneven and low distribution of human-supervised training data for other languages apart from English (Reimers and Gurevych, 2020; Ruder, 2022; Feng et al., 2022; Wieting et al., 2023). Collecting human annotations for training data generation is not scalable, as it is cumbersome to search and hire native speakers, check their language proficiency, and teach them. Additionally, human annotators are expensive, thereby requiring a large annotation budget for generating a sufficient amount of training pairs (cf. Figure 5).

Multilingual query generation is a complex task (Wang et al., 2021). It requires understanding of semantic mappings of words across languages,

¹Acronym blinded for review.

²Our dataset will be available in the supplementary material.

³Throughout the paper, we use “multilingual retrieval” to collectively denote both cross-language, i.e., cross-lingual and within language, i.e., monolingual retrieval tasks.

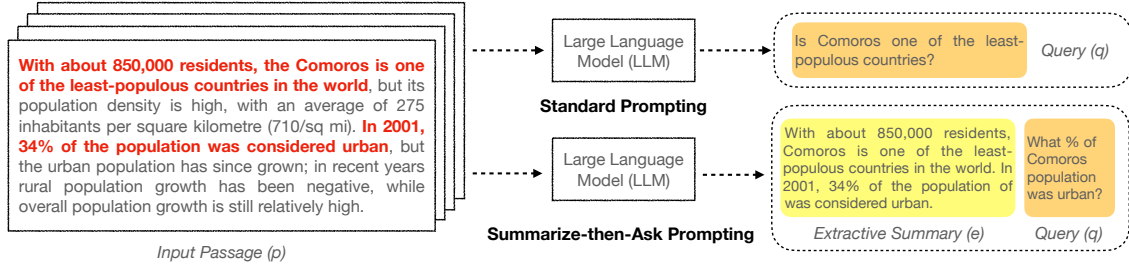


Figure 1: An illustration of SAP (*Summarize-then-Ask Prompting*) versus standard prompting for English query generation on English Wikipedia. SAP assists the large language model (LLM) in improving multilingual query generation (orange box) by identifying the relevant sections of the input passage (highlighted in red) using extractive summarization (yellow box) as an intermediate reasoning step.

similar to machine translation (Forcada, 2002; Tan et al., 2019; Zhu et al., 2023). Recently, using a Large Language Model (LLM) for query generation has been popular in English (Bonifacio et al., 2022; Dai et al., 2023). But as illustrated in Figure 1, standard prompt templates can lead the LLM to generate either extractive or uninformative⁴ queries across multiple languages.

To improve the quality of the generated query, we propose SAP (*Summarize-then-Ask Prompting*), where we prompt the LLM to break down the query generation in two stages: (i) *summary extraction*, which identifies the relevant information from the long input passage and extracts the best representative sentences as the summary, and (ii) *query generation*, which generates a multilingual query relevant for the input passage, using the extracted summary (first stage) as the intermediate step. SAP highlights the relevant information within the passage and produces difficult (i.e., informative) queries in the target language.

In our work, we utilize PaLM 2 (Anil et al., 2023) for multilingual query generation. The generated query paired with the original passage from Wikipedia is used to construct the JUMP-IR dataset. JUMP-IR spans 33 diverse languages, including both high and very-low resource-level languages. JUMP-IR provides synthetic training pairs for improving dense retrieval models without requiring any human supervision. JUMP-IR is one of largest multilingual synthetic training dataset with 28 million training pairs (cf. Table 1).

We develop synthetic multilingual (both monolingual and cross-lingual) dense retrieval models called JUMP-X, using mT5 (base) (Xue et al., 2021) as backbone and fine-tuning on JUMP-IR. We compare JUMP-X with models fine-tuned with

human supervision by changing only the training dataset while keeping other, i.e., both model parameters and training settings unchanged. We evaluate on three standard multilingual retrieval benchmarks (two cross-lingual and one monolingual). On XOR-Retrieve (Asai et al., 2021a), JUMP-X outperforms the best-supervised baseline (mContriever-X) by 7.1 points at Recall@5kt. On MIRACL (Zhang et al., 2023b), a monolingual retrieval benchmark, JUMP-X is inferior to the best-supervised baseline (mContriever) by 9.0 points at nDCG@10, which shows room for future improvement. On XTREME-UP (Ruder et al., 2023), a challenging benchmark containing 20 underrepresented Indo-European languages, JUMP-X outperforms the best-supervised baseline (mContriever-X) by 11.7 points at MRR@10.

2 JUMP-IR Dataset Overview

In our dataset overview, we describe the SAP design formulation for multilingual query generation (§2.1), data construction details (§2.2), and statistics and analysis (§2.3).

2.1 SAP Design Formulation

Multilingual query generation is not a trivial task as it requires a deep understanding of the passage content and its own translations across different languages (Wang et al., 2021). Passages can often be lengthy and contain information on multiple topics. Using the entire passage can potentially hallucinate models by generating non-meaningful queries, which affects the retrieval performance (Gospodinov et al., 2023).

To break down the task complexity of multilingual query generation and improve multilingual question quality, we implement summarize-then-ask prompting (SAP). As shown above in Figure 1, we identify the relevant information within a passage by asking the LLM to generate an extractive

⁴Uninformative denotes a query that can be easily answered using the first (or last) few words in the passage.

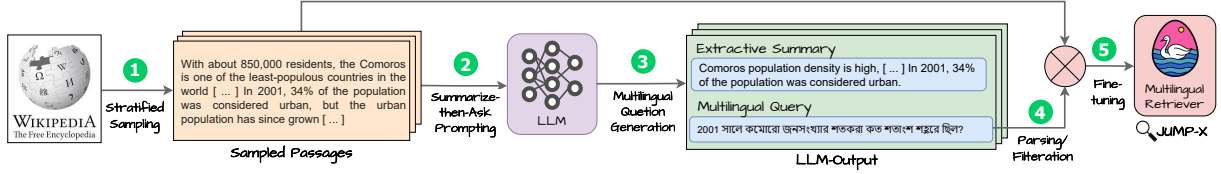


Figure 2: An illustration of cross-lingual JUMP-IR dataset construction procedure. (1) Sample N passages from English Wikipedia using stratified sampling for each target language out of a total of L languages; (2) Feed a single input passage along with few-shot exemplars to the LLM with SAP (summarize-then-ask prompting); (3 & 4) Parse the LLM output to receive the synthetic query in target language (above in Bengali); (5) Fine-tune a multilingual dense retriever model (JUMP-X) with training data combined for all languages, i.e., $N \times L$ pairs.

summary and use it as an intermediate step for generating informative queries. The procedure is described in more detail below:

(i) Summary extraction. The LLM constructs an extractive summary e_s of the input passage p_s , where s denotes the source language. The summary captures the most relevant information within the passage p_s (which occasionally may be long) acting as an useful intermediate signal for the LLM to generate a multilingual query in the later stage. We denote the first stage as $e_s = \text{LLM}(p_s; \theta^1, \dots, \theta^k)$, where $(\theta^1, \dots, \theta^k)$ denotes the k few-shot prompt exemplars⁵ containing the passage, summary in the source language s and the query in the target language t .⁶

(ii) Query Generation. Next, the LLM combines the summary e_s generated previously, with the original input passage p_s , highlighting the relevant information required for composing the query (q_t) in the target language t . We denote this stage as $q_t = \text{LLM}(e_s, p_s; \theta^1, \dots, \theta^k)$, where extractive summary e_s , input passage p_s and k -shot exemplars all appear from the first stage.

2.2 JUMP-IR Dataset Construction

For constructing JUMP-IR, we only require an unlabeled corpus of passages and generate multilingual training pairs. An overview of the cross-lingual generation procedure is shown in Figure 2. Prompt examples are shown in Appendix (§C.3).

Cross-lingual. The goal is to generate a query in the target language t using the input passage in English (source language s). We use a stratified sampling algorithm (for more details, refer

to §F.4 in Appendix) to sample a maximum of one million passages for each target language t from the English Wikipedia corpus used in XOR-Retrieve (Clark et al., 2020; Asai et al., 2021a) or XTREME-UP (Ruder et al., 2023). Next, we construct five prompt exemplars in English, where we generate both the exemplar summaries and queries in English. Further, we use Google Translate⁷ to translate the exemplar queries to other languages. Finally, we construct the prompt, where we explain our query generation task as an instruction, include the target language, and the 5-shot exemplars as an input to the LLM with SAP.

Monolingual. The goal is to generate a query in the same language as the input passage ($s = t$). We follow the setting similar to the cross-lingual task. We first sample one million passages (if available) for each language-specific Wikipedia corpus in MIRACL (Zhang et al., 2023b).⁸ Next, we carefully select three training pairs as our prompt exemplars.⁹ For languages with no training split, we manually construct our prompt exemplars. Further, we use Google Bard.¹⁰ to generate exemplar summaries in the target language. Finally, we construct the prompt, where we explain our query generation task, include the language, and the 3-shot exemplars with SAP.

2.3 Dataset Statistics and Human Validation

JUMP-IR synthetic training dataset spans 33 diverse languages, including both cross- and monolingual query-passage pairs. All queries in JUMP-IR are synthetic and LLM-generated using PaLM 2 (Anil et al., 2023) with small size (S). Detailed statistics can be found in the Appendix.

⁵Multilingual query generation requires few-shot prompt exemplars. As our experiments show in (§4), zero-shot prompting often generates unparseable outputs with PaLM 2.

⁶In our work, we did not use abstractive summarization, as LLMs have notoriously been shown to hallucinate and generate factual inconsistencies in their output generations (Maynez et al., 2020; Liu et al., 2023).

⁷Google Translate: translate.google.com

⁸For 16 / 18 languages, MIRACL contains a training split except for German (de) and Yoruba (yo).

⁹As language-specific passages consume more tokens, e.g., Telugu, to save computational budget, we rely only on 3-shot exemplars (instead of 5) for the monolingual task.

¹⁰Google Bard: bard.google.com

Model	PLM	PT	Finetune (Datasets)	Avg.	Ar	Bn	Fi	Ja	Ko	Ru	Te
<i>Existing Supervised Baselines (Prior work)</i>											
Dr. DECR (Li et al., 2022)	XLM-R	WikiM	NQ + XOR*	73.1	70.2	85.9	69.4	65.1	68.8	68.8	83.2
mDPR (Asai et al., 2021a)	mBERT	—	XOR	50.2	48.9	60.2	59.2	34.9	49.8	43.0	55.5
mBERT + xQG (Zhuang et al., 2023)	mBERT	—	XOR	53.5	42.4	54.9	54.1	33.6	52.3	33.8	52.5
Google MT + DPR (Asai et al., 2021a)	BERT	—	NQ	69.6	69.6	82.2	62.4	64.7	68.8	60.8	79.0
OPUS MT + DPR (Asai et al., 2021a)	BERT	—	NQ	50.6	52.4	62.8	61.8	48.1	58.6	37.8	32.4
<i>Zero-shot baselines (English-only supervision)</i>											
mContriever	mT5	mC4	—	38.9	35.9	33.9	43.6	34	35.1	45.1	44.5
mDPR (En)	mT5	—	MS MARCO	39.3	34.3	35.5	45.2	40.2	36.5	43.9	39.5
mContriever (En)	mT5	mC4	MS MARCO	44.0	37.5	38.2	50.6	41.1	37.2	49.8	53.8
<i>Supervised Baselines (Cross-lingual supervision)</i>											
mDPR-X	mT5	—	XOR	53.6	51.5	63.5	52.5	45.6	52.3	43.0	66.8
mContriever-X	mT5	mC4	XOR	55.3	52.1	68.1	54.5	47.7	50.5	50.2	64.3
mDPR-X	mT5	—	MS MARCO + XOR	58.2	55.3	70.1	56.7	49.8	55.8	50.6	69.3
mContriever-X	mT5	mC4	MS MARCO + XOR	59.6	54.7	73.4	57.0	53.1	56.5	51.5	71.0
<i>Synthetic Baselines (Our work)</i>											
JUMP-X (500K)	mT5	—	JUMP-IR	59.0	54.0	67.4	59.2	52.7	55.1	54.4	70.2
JUMP-X (500K)	mT5	mC4	JUMP-IR	63.0	57.0	71.1	61.8	56.8	60.7	63.3	70.2
JUMP-X (7M)	mT5	—	JUMP-IR	65.1	57.9	75.0	65.6	59.3	58.9	64.6	74.4
JUMP-X (7M)	mT5	mC4	JUMP-IR	66.7	61.2	77.0	65.0	62.2	62.8	65.4	73.5

Table 2: Experimental results showing Recall@5kt for cross-lingual retrieval on XOR-Retrieve dev (Asai et al., 2021a); (PLM) denotes the pretrained language model; (PT) denotes the pretraining dataset; (*) Dr.DECR is fine-tuned in a complex training setup across more datasets (§3.3); WikiM denotes WikiMatrix (Schwenk et al., 2021); XOR denotes XOR-Retrieve; JUMP-X (ours) is fine-tuned on 500K and 7M synthetic data.

Human validation. We conduct an validation study to evaluate the quality of generated queries in JUMP-IR for a subset of the languages.¹¹ We evaluate each query across a three-level rating scale measuring fluency, adequacy and language. From Appendix (Table 6), the generated query quality in English is found best. Around 86% of the generated queries are adequate and 88% are fluent (ratings 1 and 2) across five evaluated languages. For more details including results, refer to Appendix (§D).

Content Filtering. LLMs are shown to generate undesirable content, particularly under conditions that prime the model with material targeted at drawing out any negative patterns or associations in the training data (Gehman et al., 2020; Bender et al., 2021). We filter out training pairs in JUMP-IR with content classification of either /Adult or any of the /Sensitive Subjects labels. For more details on filtering, refer to Appendix (§D).

3 Experiments

3.1 Datasets and Metrics

We evaluate on three multilingual retrieval benchmarks: (i) **XOR-Retrieve** (Asai et al., 2021a), (ii) **MIRACL** (Zhang et al., 2023b) and (iii) **XTREME-UP** (Ruder et al., 2023). XOR-Retrieve and XTREME-UP are cross-lingual and MIRACL is monolingual. Following prior work,

¹¹Finding native speakers for all of the 33 languages, who are willing to annotate is both cumbersome and expensive.

we evaluate models at Recall@5kt on XOR-Retrieve, nDCG@10 on MIRACL and MRR@10 on XTREME-UP. For more details on evaluation benchmarks, refer to Appendix (§G).

3.2 Experimental Methods

Baselines. Following common practice across all datasets, we evaluate three broad range of baselines: (1) Zero-shot: where the model is fine-tuned only for human-labeled English training data such as MS MARCO (Nguyen et al., 2016) or NQ (Kwiatkowski et al., 2019). (2) Gold FT: where the model denoted by “X” (model-X) is fine-tuned on language-specific human labeled, i.e., gold training data. (3) Synthetic FT: where the model denoted by “JUMP-X” is fine-tuned without any gold training data, relying only on JUMP-IR training data. Additionally, we also report the amount of synthetic pairs used, e.g., 500K for fine-tuning a JUMP-X (500K) model.

Model Choices. For our dense retrieval models, we adapt DPR (Karpukhin et al., 2020) to the multilingual setting. Next, we include mContriever (Izacard et al., 2022) which adopts an additional pre-training stage with contrastive loss based on unsupervised data prepared from pairwise sentence cropping in mC4 (Xue et al., 2021).

Existing Baselines. For XOR-Retrieve, we include Dr. DECR (Li et al., 2022), a cross-lingual ColBERT (Khattab and Zaharia, 2020) fine-tuned

Model	Avg.	ar	bn	en	es	fa	fi	fr	hi	id	ja	ko	ru	sw	te	th	zh	de	yo
<i>Existing Supervised Baselines (Prior work)</i>																			
BM25	38.5	48.1	50.8	35.1	31.9	33.3	55.1	18.3	45.8	44.9	36.9	41.9	33.4	38.3	49.4	48.4	18.0	22.6	40.6
mDPR	41.8	49.9	44.3	39.4	47.8	48.0	47.2	43.5	38.3	27.2	43.9	41.9	40.7	29.9	35.6	35.8	51.2	49.0	39.6
Hybrid	56.6	67.3	65.4	54.9	64.1	59.4	67.2	52.3	61.6	44.3	57.6	60.9	53.2	44.6	60.2	59.9	52.6	56.5	37.4
Cohere-API	54.2	66.7	63.4	50.1	50.7	48.4	67.5	44.3	57.3	50.5	51.6	54.6	47.7	54.3	63.8	60.6	38.9	41.4	62.9
<i>Zero-shot baselines (English-only supervision)</i>																			
mDPR (En)	39.8	49.7	50.1	35.4	35.3	39.3	48.2	31.3	37.4	35.6	38.9	44.1	36.1	33.8	49.2	50.6	34.7	32.1	34.4
mContriever (En)	37.8	49.1	48.4	32.7	33.3	37.1	48.4	27.0	35.9	32.7	34.1	40.2	35.1	44.5	46.2	45.0	27.5	29.7	33.7
<i>Supervised Baselines (Monolingual supervision)</i>																			
mDPR-X	39.6	52.8	57.1	30.2	24.7	37.6	46.1	26.4	27.8	37.3	42.9	38.3	34.9	53.7	68.4	58.2	34.9	19.2	22.2
mContriever-X	55.4	66.4	68.4	44.2	42.8	48.9	65.2	46.2	45.0	45.8	56.8	58.8	51.2	67.7	79.0	70.7	49.4	42.3	48.4
<i>Synthetic Baselines (Our work)</i>																			
JUMP-X (180K)	46.4	60.2	57.1	34.7	33.4	36.3	40.6	64.3	33.0	39.5	40.8	43.3	49.7	40.0	55.9	56.3	63.3	50.2	36.5

Table 3: Experimental results for monolingual retrieval on MIRACL dev (Zhang et al., 2023b). All scores denote **nDCG@10**; (Hyb.) denotes Hybrid retriever with ranked fusion of three retrievers: mDPR, mColBERT and BM25; BM25, mDPR and Hybrid scores taken from (Zhang et al., 2023b); Cohere-API is used as a reranker on top of 100 BM25 results, taken from (Kamalloo et al., 2023). JUMP-X (ours) is fine-tuned on 180K synthetic data.

on large amounts of supervised data in a computationally expensive setup of knowledge distillation with English ColBERTv2 (Santhanam et al., 2022). xQG (Zhuang et al., 2023) involving cross-language query generation and concatenating the queries along with the passage representation. We also include two-stage translation baselines, Google Translate and Opus-MT from Asai et al. (2021a). For MIRACL, we include the official BM25, mDPR and Hybrid (combining BM25, mDPR and mColBERT) baseline available in Zhang et al. (2023b), and the Cohere-API is used as a reranker with top-100 BM25 retrieved results in Kamalloo et al. (2023).

3.3 Implementation Details

Supervised Baselines. We replicate mContriever and mDPR zero-shot baselines by initializing from a multilingual T5-base checkpoint (Xue et al., 2021) and fine-tune on MS MARCO, in a setup similar to Ni et al. (2022). Similarly, mContriever-X and mDPR-X have been additionally fine-tuned on training split available for each dataset. For additional technical details on supervised baselines, refer to Appendix (§F.2). mContriever includes an additional pre-training stage, we set the batch size to 8192, learning rate to $1e^{-3}$ and pre-train for 600K steps with mC4 (Xue et al., 2021). For more details on pretraining, refer to Appendix (§F.1).

Synthetic Baselines. For our synthetic baselines, we pre-train on mC4 and fine-tune on JUMP-IR with in-batch negatives with the contrastive loss function (van den Oord et al., 2018). During fine-tuning, we set the batch size to 4096, learning rate to $1e^{-3}$ and fine-tune JUMP-X for 5K to 50K steps, depending upon the size of the training dataset. In all our experiments, we use the PaLM

2 (S) (Anil et al., 2023) to generate the cross-language multilingual queries due to its rather low-cost and quick inference. For additional hyperparameter choices and fine-tuning details, refer to Appendix (§F.3). For all our experiments, we use T5X Retrieval (Ni et al., 2022) for pre-training, fine-tuning and evaluation.

3.4 Experimental Results

XOR-Retrieve. Table 2 shows that JUMP-X (7M) which is fine-tuned on 7M synthetic pairs (max. of 1M per language) outperforms the best FT model, mContriever-X, by 7.1 points on Recall@5kt. Without mC4 pre-training, our JUMP-X (7M) performance drops by only 1.6 points. We also evaluate JUMP-X (500k), a limited-budget baseline fine-tuned on 500k training pairs, that outperforms mContriever-X by 3.6 points. Few existing baselines outperform JUMP-X, however, the comparison is not fair, as Dr. DECR is a multilingual ColBERT (Khattab and Zaharia, 2020) model, which is computationally expensive at runtime (Thakur et al., 2021) and Google MT + DPR rely on a powerful Google Translate system for translation.

MIRACL. Table 3 shows that JUMP-X (180K) model is competitive on MIRACL. JUMP-X (180K) outperforms the best zero-shot model, by 6.6 points on nDCG@10. However, JUMP-X is unable to outperform mContriever-X, fine-tuned on around 90K human-labeled training pairs with up to four hard negatives available in MIRACL. However, JUMP-X have not been optimized with hard-negatives. Few existing baselines outperform JUMP-X, however the comparison is not fair, as the Hybrid baseline relies on information based on aggregation of three models, and for Cohere-API,

Model	Avg.	as	bho	brx	gbm	gom	gu	hi	hne	kn	mai	ml	mni	mr	mwr	or	pa	ps	sa	ta	ur
<i>Zero-shot baselines (English-only supervision)</i>																					
mDPR (En)	6.3	2.6	6.4	0.4	7.2	1.3	8.6	13.3	5.2	10.4	6.4	12.3	0.2	8.9	5.8	0.4	6.0	5.6	5.2	10.2	10.0
mContriever (En)	7.9	7.9	3.2	7.8	0.3	9.7	2.2	11.1	15.2	8.2	10.6	8.6	15.6	0.4	10.7	8.5	1.1	10.3	3.3	5.7	12.9
<i>Supervised Baselines (Cross-lingual supervision)</i>																					
mDPR-X	8.4	6.7	9.9	4.8	10.0	8.7	8.8	9.1	9.4	9.0	10.0	10.5	4.8	7.8	9.6	6.9	8.6	7.4	8.5	8.1	9.1
mContriever-X	12.4	9.8	15.7	6.7	14.0	11.7	13.3	15.5	13.9	13.6	13.9	16.9	6.5	12.0	13.8	7.5	13.4	9.8	12.4	13.0	14.1
mContriever-X [♡]	13.5	11.6	15.4	8.0	16.9	12.3	15.2	16.7	15.7	14.7	15.6	17.4	7.0	14.2	14.7	9.1	13.2	10.1	14.8	12.1	14.9
<i>Synthetic Baselines (Our work)</i>																					
JUMP-X (120K) ^{MT}	26.1	25.2	29.5	2.1	30.8	22.1	31.5	35.8	31.5	28.7	32.2	34.6	2.2	32.7	27.7	14.8	30.7	21.0	28.2	30.6	29.2
JUMP-X (120K)	25.2	24.4	27.7	4.3	28.3	25.4	29.4	32.4	28.8	30.1	31.8	34.4	5.1	30.7	25.7	15.8	29.6	20.6	26.1	27.9	26.1

Table 4: Experimental results for cross-lingual retrieval on XTREME-UP test (Ruder et al., 2023). (♡) denotes the mContriever-X model fine-tuned without MS MARCO (Nguyen et al., 2016); Two variants of JUMP-X considered, both fine-tuned on 120K synthetic data: (1) JUMP-X (120K)^{MT} fine-tuned using Google Translate, i.e., translated prompt exemplars for 15 languages, whereas (2) JUMP-X (120K) is fine-tuned using prompt exemplars sampled from XTREME-UP training split for all languages.

the underlying model information is unknown.

XTREME-UP. Table 4 shows the results on XTREME-UP. JUMP-X (120K) model is fine-tuned by randomly selecting 5 exemplars from the XTREME-UP training dataset (human-labeled queries) for all languages, whereas the MT variant reuses XOR-Retrieve prompt exemplars with translated summaries and queries for 15 languages.¹² JUMP-X (120K)^{MT} outperforms the best supervised baseline, mContriever-X[♡] (fine-tuned without MS MARCO) by a huge margin of 12.6 points on MRR@10. The JUMP-X (120K) model performs minimally worse than the MT version by 0.9 points. Interestingly, none of the models perform well on two extremely low-resource languages, Boro (brx) and Manipuri (mni).

3.5 SAP versus Standard Prompting

We evaluate whether the generated query quality using SAP against standard few-shot prompting affect the downstream retrieval performance on XOR-Retrieve. We additionally evaluate different LLM sizes to observe a correlation in retrieval model performance with change in LLM size. To ensure consistency, we adopt the experimental setup utilized in JUMP-X (500K). Our results are shown in Figure 3 (Left), we infer two insights: (i) Increase in the LLM size provides diminishing gains in JUMP-X performance on XOR-Retrieve and PaLM-2 (S) provides the best trade-off in terms of performance and query generation speed. (ii) SAP outperforms standard prompting by at least 0.6 points Recall@5kt for all PaLM-2 generators on XOR-Retrieve, where the max-

imum improvements are observed by up to 3.2 points Recall@5kt for models sizes (S) or smaller. We hypothesize that PaLM 2 (sizes > S) are inherently able to generate coherent questions, leading to diminishing improvements with SAP versus few-shot standard prompting.

3.6 How much Synthetic data to Generate?

We analyze the optimal value of synthetic training data for training JUMP-X models. Figure 5 depicts the relative improvement in JUMP-X on XOR-Retrieve, with the performance (gradually increasing) starting to saturate after 500K synthetic pairs. The first observation is that with only 2K pairs, the JUMP-X (2K) achieves 49.1 Recall@5kt on XOR-Retrieve, which outperforms the best zero-shot (English-only) baseline. The break-even point occurs at around 200K synthetic pairs, where the JUMP-X (250K) model achieves 60.5, outperforming the best supervised baseline of mContriever-X achieving 59.6 Recall@5kt.

3.7 Indo-European Language Transferability

We investigate language transfer capabilities of synthetic data generated with JUMP-IR on Indic (Indo-European language family). We fine-tune separate JUMP-X models individually for eight languages and evaluate on XTREME-UP. From Figure 4, we observe that models fine-tuned for Konkani (gom) or Hindi (hi) transfer best on all languages in XTREME-UP (rows 3&4), whereas Tamil (ta) transfers worst (row 8). Assamese (as), Konkani (gom), Odia (or), Pashto (pa) and Sanskrit (sa) have the lowest zero-shot capabilities with JUMP-X, where in-language synthetic data is found crucial for improvement in MRR@10. Hindi (hi), Kannada (kn), Malayalam (ml), Gujarati (gu) show good zero-shot transfer capabilities with all individual fine-tuned Indic languages.

¹²We were unable to translate our prompt exemplars for 5 languages due to language unavailability in Google Translate: Boro (brx), Garhwali (gbm), Chattisgarhi (hne) and Marwari (mwr). Manipuri (mni) is available in Google Translate in “Meitei” script instead of the “Bengali-Assamese” script present in the XTREME-UP dataset.

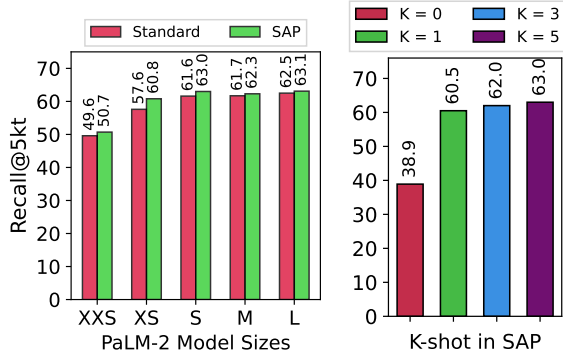


Figure 3: (Left) SAP (*Summarize-then-Ask Prompting*) (green) versus standard prompting (red) for various PaLM 2 model sizes. (Right) Varying K-shot prompt exemplars. All JUMP-X models are fine-tuned on 500K synthetic data and evaluated on XOR-Retrieve.

4 Ablation Studies

K-shot prompt exemplars. We investigate the number of K-shot prompt exemplars required by PaLM 2 and the variation in the cross-lingual performance with K on XOR-Retrieve.¹³ From Figure 3 (right), we observe a linear improvement in Recall@5kt with increase in K. Best Recall@5kt is observed with K = 5. Our SAP technique cannot perform well zero-shot (i.e., K = 0) due to the complex nature of the multilingual question query task which requires a few examples for PaLM 2 to understand the difficult task.

ByT5 tokenizer. We evaluate whether the poor performance of JUMP-X on low-resource languages in XTREME-UP can be attributed towards low-quality language tokenization. We reproduce JUMP-X, with a ByT5-base (Xue et al., 2022) model as backbone, which contains a language independent tokenizer extension. From our results in Table 5, ByT5 models underperform by up to 9.8 points MRR@10 on XTREME-UP, in contrast to mT5-base. Additionally, JUMP-X performance on both mni and brx do not improve with ByT5. We leave it as future work to investigate the low-quality performance of JUMP-X on mni and brx.

Training split query replacement. Next, we evaluate the impact of human-generated versus LLM-generated queries on retrieval performance on XTREME-UP. We replace all human-generated queries in the XTREME-UP training split with only synthetic queries generated using PaLM 2 (S). From Table 5, the performance drops by

¹³We limit K = 5, as it fits within the 4096 tokens in context length. Adding more exemplars require longer PaLM 2 contexts which increases the computational cost significantly.

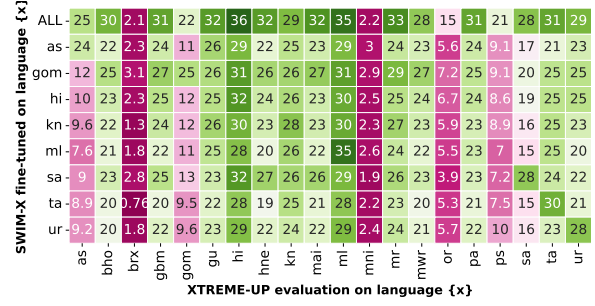


Figure 4: Heatmap showing MRR@10 denoting language-based transfer ability of JUMP-X (120K) across Indo-European languages available in XTREME-UP (Ruder et al., 2023). (ALL) denotes JUMP-X fine-tuned on all XTREME-UP languages.

2.0 points on MRR@10. This shows that better quality human-generated queries results in better MRR@10 in XTREME-UP. However, JUMP-X can be fine-tuned effectively with synthetic generated queries, by marginally dropping in retrieval performance.

5 Cost Comparison

Generating synthetic training data is relatively inexpensive however, not free. The cost is dependent upon the length of the prompt, input, and output generated from the LLM. The costs also linearly increase with each additional language pair. At this writing, PaLM 2 and similar LLMs cost about 0.0005 USD for 1000 characters in the input and output text.¹⁴ Our prompts on average contain about 8-9K characters in the prompt input and generate about 1-2K characters in the output. The relative performance improvement associated with annotation cost in XOR-Retrieve is shown in Figure 5. Generating 200K synthetic training pairs in JUMP-IR will roughly cost \$1K USD. JUMP-X (200K) performs comparably to the best supervised baseline (mContriever-X), trained on 15.2K human-annotated pairs, requiring roughly 14 times more, i.e., \$14.1K USD to annotate, if we pay an hourly rate of \$18.50 USD per hour for the annotator (local minimum wages is \$11.50 USD/hr) following (Zhang et al., 2023b), assuming an estimated annotation cost of 3.0 minutes per example (Ruder et al., 2023).

6 Background and Related Work

The development of pre-trained multilingual LMs has contributed toward recent progress in multilin-

¹⁴PaLM 2 pricing: cloud.google.com/vertex-ai/pricing

Model	Backbone	Query Gen.	brx	mni	MRR@10
1. Models with Byte-level (UTF-8) tokenizer					
mCon.-X [♡]	ByT5	Human	1.8	1.0	2.1
JUMP-X (120k) ^{MT}	ByT5	PaLM 2	2.1	4.9	13.3
JUMP-X (120k)	ByT5	PaLM 2	5.1	5.8	15.4
2. Human-generated query replacement in XTREME-UP					
mCon.-X [♡]	mT5	Human	-	-	13.5
JUMP-X (≈10K)	mT5	PaLM 2	-	-	11.5

Table 5: Ablations in XTREME-UP. First, we replace the mT5 backbone with ByT5. Next, we replace the human-generated queries in the XTREME-UP training dataset with PaLM-2 synthetic queries; MRR@10 scores are macro averaged across all 20 languages; brx denotes Boro and mni denotes the Manipuri language.

gual retrieval (Asai et al., 2021a; Izacard et al., 2022; Asai et al., 2021b; Li et al., 2022; Ruder et al., 2023; Zhang et al., 2023b,a). Notable baselines include mDPR and mContriever. mDPR (Asai et al., 2021a,b; Zhang et al., 2023a) extends English DPR (Karpukhin et al., 2020) to the multilingual setting. mContriever (Izacard et al., 2022) adopts an unsupervised pre-training objective using the contrastive loss function and data prepared from mC4 (Xue et al., 2021) and fine-tuned on MS MARCO (Nguyen et al., 2016).

Synthetic Data Generation. Traditionally, docT5query (Nogueira and Lin, 2019) for query generation has been prominent for generating synthetic training data in English (Ma et al., 2021; Thakur et al., 2021; Wang et al., 2022; Thakur et al., 2022). Recently, using LLMs for query generation has gained interest. Bonifacio et al. (2022) proposed InPars, where they few-shot prompt GPT-3 (Brown et al., 2020) to generate synthetic queries. Similarly, complementary works (Sachan et al., 2022; Jeronimo et al., 2023; Boytsov et al., 2023; Saad-Falcon et al., 2023; Dua et al., 2023) all follow a similar setup in Bonifacio et al. (2022). Dai et al. (2023) proposed Promptagator, which studied task-dependent few-shot LLM prompting and used the synthetic data for both retrieval and ranking models. Similarly, HyDE (Gao et al., 2023) and GenRead (Yu et al., 2023) generate synthetic documents instead of queries. However, prior work has focused on English, with the exception of HyDE. In our work, we robustly investigate how LLMs can be used for improving multilingual retrieval systems.

Multilingual Datasets. Prior work investigate techniques to build multilingual datasets for better fine-tuning or evaluation of dense retrieval models. Datasets such as NeuCLIR (Lawrie

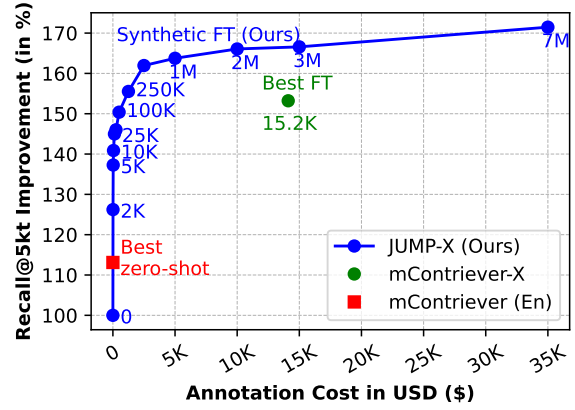


Figure 5: Recall@5kt improvement (in %) on XOR-Retrieve versus annotation cost in USD (\$) to create the training dataset. The amount of generated training pairs (human-generated marked in red and green; LLM-generated marked in blue) are mentioned with each marked datapoint in the graph.

et al., 2023), MKQA (Longpre et al., 2021) have been constructed using human annotators. Similarly, mMARCO (Bonifacio et al., 2021) has been generated using machine translation of MS MARCO (Nguyen et al., 2016). However, as translated documents are not written by a native speaker, mMARCO and similar datasets suffer from artifacts such as “Translationese” (Clark et al., 2020). A concurrent work (Mayfield et al., 2023) prompts GPT-3 to generate English queries from language specific passages in NeuCLIR.

7 Conclusion

In this work, we present JUMP-IR, a synthetic multilingual retrieval training dataset with 28 million training pairs across 33 diverse languages. JUMP-IR allows synthetic fine-tuning of multilingual dense retrieval models cheaply without human supervision. JUMP-IR is constructed using SAP (*summarize-then-ask prompting*) which assists the LLM to identify the relevant sections of the input passage, improving the quality of the generated multilingual query.

Our rigorous evaluation across three multilingual retrieval benchmarks assess our dataset quality. We find that JUMP-X, fine-tuned on JUMP-IR (keeping model and training parameters unchanged) outperform the best supervised cross-lingual baseline, mContriever-X by 7.1 points Recall@5kt on XOR-Retrieve and 11.7 points MRR@10 on XTREME-UP, while remaining competitive on monolingual retrieval in MIRACL.

8 Limitations of JUMP-IR dataset

JUMP-IR like any other dataset is not perfect and has limitations. These limitations do not directly affect the downstream multilingual retrieval task, where dense retrieval models learn how to match relevant passages to queries. JUMP-IR dataset has been created for the “sole” purpose of training multilingual retrieval models. We describe below few noted limitations:

1. Decontextualization. PaLM 2 captures the salient information from the paragraph, but can generate the query in a reduced context, which cannot be answered without the Wikipedia paragraph.

2. Code-Switching. PaLM 2 can occasionally generate a code-switched query with words combined for English and the target language. Code-switching is more frequently observed for cross-lingual generation in low-resource languages.

3. Passage Quality and Length. A good quality passage contains relevant information about a topic which PaLM 2 uses to generate a synthetic query. However, if the passage is really short with little or zero information, or contains noisy information, this likely can generate a subpar query.

4. Factual inconsistencies in LLM generation. LLMs have been found to generate text lacking sufficient grounding to knowledge sources (Dziri et al., 2022; Ji et al., 2023), thereby posing risks of misinformation and hallucination in their generated outputs (Maynez et al., 2020; Raunak et al., 2021; Muller et al., 2023). Queries in JUMP-IR are relevant for the input passage, but are not human-verified, thereby queries may contain factual inconsistencies.

References

Rohan Anil, Andrew M. Dai, Orhan Firat, Melvin Johnson, Dmitry Lepikhin, Alexandre Passos, Siamak Shakeri, Emanuel Taropa, Paige Bailey, Zhifeng Chen, Eric Chu, Jonathan H. Clark, Laurent El Shafey, Yanping Huang, Kathy Meier-Hellstern, Gaurav Mishra, Erica Moreira, Mark Omernick, Kevin Robinson, Sebastian Ruder, Yi Tay, Kefan Xiao, Yuanzhong Xu, Yujing Zhang, Gustavo Hernández Ábrego, Junwhan Ahn, Jacob Austin, Paul Barham, Jan A. Botha, James Bradbury, Siddhartha Brahma, Kevin Brooks, Michele Catasta, Yong Cheng, Colin Cherry, Christopher A. Choquette-Choo, Aakanksha Chowdhery, Clément Crepy, Shachi Dave, Mostafa Dehghani, Sunipa Dev, Jacob Devlin, Mark Díaz, Nan Du, Ethan Dyer, Vladimir Feinberg, Fangxiaoyu Feng, Vlad

Fienber, Markus Freitag, Xavier Garcia, Sebastian Gehrmann, Lucas Gonzalez, and et al. 2023. [PaLM 2 Technical Report](#). *CoRR*, abs/2305.10403.

Akari Asai, Jungo Kasai, Jonathan H. Clark, Kenton Lee, Eunsol Choi, and Hannaneh Hajishirzi. 2021a. [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, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 547–564. Association for Computational Linguistics.

Akari Asai, Xinyan Yu, Jungo Kasai, and Hanna Hajishirzi. 2021b. [One Question Answering Model for Many Languages with Cross-lingual Dense Passage Retrieval](#). In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 7547–7560.

Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. [On the Dangers of Stochastic Parrots: Can Language Models Be Too Big?](#) In *FAccT ’21: 2021 ACM Conference on Fairness, Accountability, and Transparency, Virtual Event / Toronto, Canada, March 3-10, 2021*, pages 610–623. ACM.

Luiz Henrique Bonifacio, Hugo Queiroz Abonizio, Marzieh Fadaee, and Rodrigo Frassetto Nogueira. 2022. [InPars: Unsupervised Dataset Generation for Information Retrieval](#). In *SIGIR ’22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, Madrid, Spain, July 11 - 15, 2022*, pages 2387–2392. ACM.

Luiz Henrique Bonifacio, Israel Campiotti, Roberto de Alencar Lotufo, and Rodrigo Frassetto Nogueira. 2021. [mMARCO: A Multilingual Version of MS MARCO Passage Ranking Dataset](#). *CoRR*, abs/2108.13897.

Leonid Boytsov, Preksha Patel, Vivek Sourabh, Ridhi Nisar, Sayani Kundu, Ramya Ramanathan, and Eric Nyberg. 2023. [InPars-Light: Cost-Effective Unsupervised Training of Efficient Rankers](#). *CoRR*, abs/2301.02998.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language Models are Few-Shot Learners](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.

Jonathan H. Clark, Eunsol Choi, Michael Collins, Dan Garrette, Tom Kwiatkowski, Vitaly Nikolaev, and Jennimaria Palomaki. 2020. TyDi QA: A Benchmark for Information-Seeking Question Answering in Typologically Diverse Languages . <i>Transactions of the Association for Computational Linguistics</i> , 8:454–470.	705
Zhuyun Dai, Vincent Y. Zhao, Ji Ma, Yi Luan, Jianmo Ni, Jing Lu, Anton Bakalov, Kelvin Guu, Keith B. Hall, and Ming-Wei Chang. 2023. Promptagator: Few-shot Dense Retrieval From 8 Examples . In <i>The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023</i> . OpenReview.net.	706
Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding . In <i>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)</i> , pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	707
Dheeru Dua, Emma Strubell, Sameer Singh, and Pat Verga. 2023. To Adapt or to Annotate: Challenges and Interventions for Domain Adaptation in Open-Domain Question Answering . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , <i>ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 14429–14446. Association for Computational Linguistics.	708
Nouha Dziri, Sivan Milton, Mo Yu, Osmar Zaiane, and Siva Reddy. 2022. On the Origin of Hallucinations in Conversational Models: Is it the Datasets or the Models? In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 5271–5285, Seattle, United States. Association for Computational Linguistics.	709
Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. 2022. Language-agnostic BERT Sentence Embedding . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , <i>ACL 2022, Dublin, Ireland, May 22-27, 2022</i> , pages 878–891. Association for Computational Linguistics.	710
Mikel L. Forcada. 2002. Explaining real MT to translators: between compositional semantics and word-for-word . In <i>Proceedings of the 6th EAMT Workshop: Teaching Machine Translation</i> , Manchester, England. European Association for Machine Translation.	711
Luyu Gao, Xueguang Ma, Jimmy Lin, and Jamie Callan. 2023. Precise Zero-Shot Dense Retrieval without Relevance Labels . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , <i>ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 1762–1777. Association for Computational Linguistics.	712
Samuel Gehman, Suchin Gururangan, Maarten Sap, Yejin Choi, and Noah A. Smith. 2020. RealToxicityPrompts: Evaluating Neural Toxic Degeneration in Language Models . In <i>Findings of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020</i> , volume EMNLP 2020 of <i>Findings of ACL</i> , pages 3356–3369. Association for Computational Linguistics.	713
Mitko Gospodinov, Sean MacAvaney, and Craig Macdonald. 2023. Doc2Query-: When Less is More . In <i>Advances in Information Retrieval - 45th European Conference on Information Retrieval, ECIR 2023, Dublin, Ireland, April 2-6, 2023, Proceedings, Part II</i> , volume 13981 of <i>Lecture Notes in Computer Science</i> , pages 414–422. Springer.	714
Sebastian Hofstätter, Sheng-Chieh Lin, Jheng-Hong Yang, Jimmy Lin, and Allan Hanbury. 2021. Efficiently Teaching an Effective Dense Retriever with Balanced Topic Aware Sampling . In <i>SIGIR '21: The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval, Virtual Event, Canada, July 11-15, 2021</i> , pages 113–122. ACM.	715
Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2022. Unsupervised Dense Information Retrieval with Contrastive Learning . <i>Transactions on Machine Learning Research</i> .	716
Vitor Jeronymo, Luiz Henrique Bonifacio, Hugo Queiroz Abonizio, Marzieh Fadaee, Roberto de Alencar Lotufo, Jakub Zavrel, and Rodrigo Frassetto Nogueira. 2023. InPars-v2: Large Language Models as Efficient Dataset Generators for Information Retrieval . <i>CoRR</i> , abs/2301.01820.	717
Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of Hallucination in Natural Language Generation . <i>ACM Comput. Surv.</i> , 55(12).	718
Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. 2020. The State and Fate of Linguistic Diversity and Inclusion in the NLP World . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 6282–6293, Online. Association for Computational Linguistics.	719
Ehsan Kamalloo, Xinyu Zhang, Odunayo Ogundepo, Nandan Thakur, David Alfonso-Hermelo, Mehdi Rezagholizadeh, and Jimmy Lin. 2023. Evaluating Embedding APIs for Information Retrieval . In <i>Proceedings of the The 61st Annual Meeting of the Association for Computational Linguistics: Industry Track, ACL 2023, Toronto, Canada, July 9-14, 2023</i> , pages 518–526. Association for Computational Linguistics.	720

763	Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick	Ji Ma, Ivan Korotkov, Yinfei Yang, Keith Hall, and	819
764	Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and	Ryan McDonald. 2021. Zero-shot Neural Passage	820
765	Wen-tau Yih. 2020. Dense Passage Retrieval for	Retrieval via Domain-targeted Synthetic Question	821
766	Open-Domain Question Answering . In <i>Proceed-</i>	Generation . In <i>Proceedings of the 16th Conference</i>	822
767	<i>ings of the 2020 Conference on Empirical Methods</i>	<i>of the European Chapter of the Association for Com-</i>	823
768	<i>in Natural Language Processing (EMNLP)</i> , pages	<i>putational Linguistics: Main Volume</i> , pages 1075–	824
769	6769–6781, Online. Association for Computational	1088, Online. Association for Computational Lin-	825
770	Linguistics.	guistics.	826
771	Omar Khattab and Matei Zaharia. 2020. ColBERT: Ef-	James Mayfield, Eugene Yang, Dawn J. Lawrie,	827
772	ficient and Effective Passage Search via Contextu-	Samuel Barham, Orion Weller, Marc Mason,	828
773	alized Late Interaction over BERT . In <i>Proceedings</i>	Suraj Nair, and Scott Miller. 2023. Synthetic	829
774	<i>of the 43rd International ACM SIGIR Conference on</i>	Cross-language Information Retrieval Training Data .	830
775	<i>Research and Development in Information Retrieval</i> ,	<i>CoRR</i> , abs/2305.00331.	831
776	SIGIR '20, page 3948, New York, NY, USA. Asso-		
777	ciation for Computing Machinery.	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and	832
		Ryan McDonald. 2020. On Faithfulness and Factu-	833
778	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	ality in Abstractive Summarization . In <i>Proceedings</i>	834
779	field, Michael Collins, Ankur Parikh, Chris Alberti,	<i>of the 58th Annual Meeting of the Association for</i>	835
780	Danielle Epstein, Illia Polosukhin, Matthew Kelcey,	<i>Computational Linguistics</i> , pages 1906–1919, On-	836
781	Jacob Devlin, Kenton Lee, Kristina N. Toutanova,	line. Association for Computational Linguistics.	837
782	Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob		
783	Uszkoreit, Quoc Le, and Slav Petrov. 2019. Natu-	Benjamin Muller, John Wieting, Jonathan H. Clark,	838
784	ral Questions: a Benchmark for Question Answering	Tom Kwiatkowski, Sebastian Ruder, Livio Baldini	839
785	Research . <i>Transactions of the Association of Com-</i>	Soares, Roei Aharoni, Jonathan Herzig, and Xinyi	840
786	<i>putational Linguistics</i> .	Wang. 2023. Evaluating and Modeling Attribu-	841
		tion for Cross-Lingual Question Answering . <i>CoRR</i> ,	842
787	Dawn J. Lawrie, Sean MacAvaney, James Mayfield,	abs/2305.14332.	843
788	Paul McNamee, Douglas W. Oard, Luca Soldaini,		
789	and Eugene Yang. 2023. Overview of the TREC	Arvind Neelakantan, Tao Xu, Raul Puri, Alec Radford,	844
790	2022 NeuCLIR Track . <i>CoRR</i> , abs/2304.12367.	Jesse Michael Han, Jerry Tworek, Qiming Yuan,	845
		Nikolas Tezak, Jong Wook Kim, Chris Hallacy,	846
791	Yulong Li, Martin Franz, Md Arafat Sultan, Bhavani	Johannes Heidecke, Pranav Shyam, Boris Power,	847
792	Iyer, Young-Suk Lee, and Avirup Sil. 2022. Learn-	Tyna Eloundou Nekoul, Girish Sastry, Gretchen	848
793	ing Cross-Lingual IR from an English Retriever . In	Krueger, David Schnurr, Felipe Petroski Such,	849
794	<i>Proceedings of the 2022 Conference of the North</i>	Kenny Hsu, Madeleine Thompson, Tabarak Khan,	850
795	<i>American Chapter of the Association for Computa-</i>	Toki Sherbakov, Joanne Jang, Peter Welinder, and	851
796	<i>tional Linguistics: Human Language Technologies</i> ,	Lilian Weng. 2022. Text and Code Embeddings by	852
797	pages 4428–4436, Seattle, United States. Associa-	Contrastive Pre-Training . <i>CoRR</i> , abs/2201.10005.	853
798	tion for Computational Linguistics.		
799	Jimmy Lin, Rodrigo Frassetto Nogueira, and Andrew	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao,	854
800	Yates. 2021. Pretrained Transformers for Text Rank-	Saurabh Tiwary, Rangan Majumder, and Li Deng.	855
801	ing: BERT and Beyond . Synthesis Lectures on Hu-	2016. MS MARCO: A Human Generated MA-	856
802	man Language Technologies. Morgan & Claypool	chine Reading Comprehension Dataset . <i>CoRR</i> ,	857
803	Publishers.	abs/1611.09268.	858
		Jianmo Ni, Chen Qu, Jing Lu, Zhuyun Dai, Gus-	859
804	Sheng-Chieh Lin, Akari Asai, Minghan Li, Barlas	tavo Hernández Ábrego, Ji Ma, Vincent Y. Zhao,	860
805	Oguz, Jimmy Lin, Yashar Mehdad, Wen-tau Yih,	Yi Luan, Keith B. Hall, Ming-Wei Chang, and Yin-	861
806	and Xilun Chen. 2023. How to Train Your	fei Yang. 2022. Large Dual Encoders Are General-	862
807	DRAGON: Diverse Augmentation Towards Gener-	izable Retrievers . In <i>Proceedings of the 2022 Con-</i>	863
808	alizable Dense Retrieval . <i>CoRR</i> , abs/2302.07452.	<i>ference on Empirical Methods in Natural Language</i>	864
		<i>Processing, EMNLP 2022, Abu Dhabi, United Arab</i>	865
809	Nelson F. Liu, Tianyi Zhang, and Percy Liang. 2023.	<i>Emirates, December 7-11, 2022</i> , pages 9844–9855.	866
810	Evaluating Verifiability in Generative Search En-	Association for Computational Linguistics.	867
811	gines . <i>CoRR</i> , abs/2304.09848.	Rodrigo Nogueira and Jimmy Lin. 2019. From	868
		doc2query to docTTTTTquery .	869
812	Shayne Longpre, Yi Lu, and Joachim Daiber. 2021.	Mahima Pushkarna, Andrew Zaldivar, and Oddur Kjar-	870
813	MKQA: A Linguistically Diverse Benchmark for	tansson. 2022. Data Cards: Purposeful and Trans-	871
814	Multilingual Open Domain Question Answering .	parent Dataset Documentation for Responsible AI .	872
815	<i>Trans. Assoc. Comput. Linguistics</i> , 9:1389–1406.	In <i>2022 ACM Conference on Fairness, Accountabil-</i>	873
		<i>ity, and Transparency, FAccT '22</i> , page 17761826,	874
816	Ilya Loshchilov and Frank Hutter. 2019. Decoupled	New York, NY, USA. Association for Computing	875
817	Weight Decay Regularization . In <i>International Con-</i>	Machinery.	876
818	<i>ference on Learning Representations</i> .		

877	Vikas Raunak, Arul Menezes, and Marcin Junczys-Dowmunt. 2021. The Curious Case of Hallucinations in Neural Machine Translation . In <i>Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1172–1183, Online. Association for Computational Linguistics.	935
878		936
879		937
880		
881		938
882		939
883		940
884		941
885	Nils Reimers and Iryna Gurevych. 2020. Making Monolingual Sentence Embeddings Multilingual using Knowledge Distillation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 4512–4525. Association for Computational Linguistics.	942
886		943
887		944
888		945
889		946
890		
891		
892	Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao, Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong Wen. 2021. RocketQAv2: A Joint Training Method for Dense Passage Retrieval and Passage Re-ranking . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021</i> , pages 2825–2835. Association for Computational Linguistics.	947
893		948
894		949
895		950
896		951
897		952
898		953
899		954
900		
901	Uma Roy, Noah Constant, Rami Al-Rfou, Aditya Barua, Aaron Phillips, and Yinfei Yang. 2020. LAREQA: Language-Agnostic Answer Retrieval from a Multilingual Pool . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP 2020, Online, November 16-20, 2020</i> , pages 5919–5930. Association for Computational Linguistics.	955
902		956
903		957
904		958
905		959
906		960
907		961
908		962
909	Sebastian Ruder. 2022. The State of Multilingual AI. http://ruder.io/state-of-multilingual-ai/ .	963
910		
911		
912	Sebastian Ruder, Jonathan H. Clark, Alexander Gutkin, Mihir Kale, Min Ma, Massimo Nicosia, Shruti Rijhwani, Parker Riley, Jean Michel A. Sarr, Xinyi Wang, John Wieting, Nitish Gupta, Anna Katanova, Christo Kirov, Dana L. Dickinson, Brian Roark, Bidisha Samanta, Connie Tao, David Ifeoluwa Adelani, Vera Axelrod, Isaac Caswell, Colin Cherry, Dan Garrette, R. Reeve Ingle, Melvin Johnson, Dmitry Panteleev, and Partha Talukdar. 2023. XTREME-UP: A User-Centric Scarce-Data Benchmark for Under-Represented Languages . <i>CoRR</i> , abs/2305.11938.	964
913		965
914		966
915		967
916		
917		968
918		969
919		970
920		971
921		972
922		973
923	Jon Saad-Falcon, Omar Khattab, Keshav Santhanam, Radu Florian, Martin Franz, Salim Roukos, Avirup Sil, Md. Arafat Sultan, and Christopher Potts. 2023. UDAPDR: Unsupervised Domain Adaptation via LLM Prompting and Distillation of Rerankers . <i>CoRR</i> , abs/2303.00807.	974
924		975
925		976
926		
927		
928		
929	Devendra Singh Sachan, Mike Lewis, Mandar Joshi, Armen Aghajanyan, Wen-tau Yih, Joelle Pineau, and Luke Zettlemoyer. 2022. Improving Passage Retrieval with Zero-Shot Question Generation . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7-11, 2022</i> , pages 3781–3797. Association for Computational Linguistics.	977
930		978
931		979
932		980
933		981
934		982
		983
		984
		985
		986
		987
		988
		989
		990
		991
		992

- Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Edouard Grave. 2020. [CCNet: Extracting High Quality Monolingual Datasets from Web Crawl Data](#). In *Proceedings of The 12th Language Resources and Evaluation Conference, LREC 2020, Marseille, France, May 11-16, 2020*, pages 4003–4012. European Language Resources Association.
- John Wieting, Jonathan Clark, William Cohen, Graham Neubig, and Taylor Berg-Kirkpatrick. 2023. [Beyond Contrastive Learning: A Variational Generative Model for Multilingual Retrieval](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12044–12066, Toronto, Canada. Association for Computational Linguistics.
- Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. [Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. [ByT5: Towards a Token-Free Future with Pre-trained Byte-to-Byte Models](#). *Transactions of the Association for Computational Linguistics*, 10:291–306.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [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, NAACL-HLT 2021, Online, June 6-11, 2021*, pages 483–498. Association for Computational Linguistics.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023. [Generate rather than Retrieve: Large Language Models are Strong Context Generators](#). In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Xinyu Zhang, Xueguang Ma, Peng Shi, and Jimmy Lin. 2021. [Mr. TyDi: A Multi-lingual Benchmark for Dense Retrieval](#). In *Proceedings of the 1st Workshop on Multilingual Representation Learning*, pages 127–137, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xinyu Zhang, Kelechi Ogueji, Xueguang Ma, and Jimmy Lin. 2023a. [Toward Best Practices for Training Multilingual Dense Retrieval Models](#). *ACM Trans. Inf. Syst.*, 42(2).
- Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. 2023b. [MIRACL: A Multilingual Retrieval Dataset Covering 18 Diverse Languages](#). *Transactions of the Association for Computational Linguistics*, 11:1114–1131.
- Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian Huang. 2023. [Multilingual Machine Translation with Large Language Models: Empirical Results and Analysis](#). *CoRR*, abs/2304.04675.
- Shengyao Zhuang, Linjun Shou, and Guido Zuccon. 2023. [Augmenting Passage Representations with Query Generation for Enhanced Cross-Lingual Dense Retrieval](#). In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR 2023, Taipei, Taiwan, July 23-27, 2023*, pages 1827–1832. ACM.

A Appendix

The following supplementary sections in JUMP-IR are arranged as follows:

- [Appendix B](#) provides information on the JUMP-IR dataset release.
- [Appendix C](#) provides extra material with JUMP-IR dataset: Datacard, Examples and Prompts. All the prompts for all languages will be provided as text files within our supplementary submission.
- [Appendix D](#) and [E](#) provides details on the human validation of JUMP-IR question quality and content filtering.
- [Appendix F](#) provides detailed information on hyperparameters and training settings for baselines, multilingual pre-training, synthetic finetuning, and sampling strategies.
- [Appendix G](#) provides statistics for three multilingual retrieval evaluation datasets: XOR-Retrieve, MIRACL and XTREME-UP.
- [Appendix H](#) contains additional results on the JUMP-IR dataset for XOR-Retrieve and MIRACL evaluation datasets.

B Details on JUMP-IR Dataset Release

Long Term Preservation. The dataset will be available for a longer time by continually updating the Tensorflow dataset (TFDS) and HuggingFace dataset. The authors will be responsible for maintaining the dataset and in future extension of the work for supporting more languages (Joshi et al., 2020) and other cross-language retrieval setting: English query retrieving across language specific corpora (En→L), inclusion of both would improve multilingual neural retrieval models on a wider variety of languages.

Licensing. The JUMP-IR dataset is based on language-specific Wikipedia. We follow the same license as Wikipedia for JUMP-IR: Creative Commons Attribution-ShareAlike 4.0 Unported License (CC BY-SA 4.0).¹⁵ Overall, the license allows both researchers and industry alike to access the dataset, and allow them to copy and redistribute the dataset for future work.

C JUMP-IR Extra Material

C.1 JUMP-IR Data Card

We provide the datacard associated with the JUMP-IR dataset along in the supplementary ma-

terial. The datacard generated using the template provided by the Data Cards Playbook (Pushkarna et al., 2022). The datacard has been generated using the Markdown format.¹⁶ The Datacard is provided along with our dataset release in the supplementary material.

C.2 JUMP-IR Dataset Statistics

The languages covered and the amount of training pairs available in JUMP-IR are provided in [Table 7](#). A majority of the training pairs (sampled a maximum of 1 million per language pair) are provided for 18 languages in MIRACL (Zhang et al., 2023b). The rest of 15 Indo-European languages from XTREME-UP contribute for 100K training pairs. We additionally, provide two examples from JUMP-IR dataset for each retrieval task, cross-lingual and monolingual in [Figure 7](#). The cross-lingual example is provided for Chinese (zh) and monolingual for Spanish (es).

There are six fields associated with every JUMP-IR training datapoint. We briefly describe each field available below: (i) `_id`: denotes the unique identifier of the training pair. (ii) `title`: denotes the title of the Wikipedia article. (iii) `text`: denotes the passage extracted from the Wikipedia article. (iv) `query`: denotes the synthetic multilingual query generated using PaLM 2 (Anil et al., 2023). (v) `lang`: denotes the language of the synthetic query. (v) `code`: denotes the ISO code of the synthetic query language.

C.3 JUMP-IR Prompts

All prompts and their templates (across all 33 languages) used for developing JUMP-IR have been provided in the supplementary material submission. We show individual prompt examples for a single language for the three datasets in the Appendix: (1) XOR-Retrieve (English passage; Synthetic Bengali query) in [Figure 8](#), (2) MIRACL (Chinese passage; Synthetic Chinese query) in [Figure 9](#), and (3) XTREME-UP (English Passage; Synthetic Hindi query) in [Figure 10](#). The rest of the prompts will be provided in the supplementary material.

D Human Validation

In this section, we evaluate the quality of the PaLM 2 generated questions available in the

¹⁵<https://creativecommons.org/licenses/by-sa/4.0>

¹⁶The Markdown format and the template of the datacard is available here: <https://github.com/pair-code/datacardsplaybook>

Lang. (ISO)	fluency (↑)			adequacy (↑)			language (↑)		
	0	1	2	0	1	2	0	1	2
English (en)	2%	3%	95%	2%	13%	85%	0%	0%	100%
Spanish (es)	1%	10%	89%	14%	12%	74%	1%	0%	99%
Chinese (zh)	7%	19%	74%	7%	30%	63%	0%	0%	100%
Hindi (hi)	12%	5%	83%	6%	19%	75%	0%	0%	100%
Bengali (bn)	6%	4%	90%	10%	14%	76%	1%	0%	99%

Table 6: Human validation statistics on JUMP-IR. Annotators (native speakers) evaluate the query quality on a three-level rating scale (0/1/2) measured for (i) fluency, (ii) adequacy and (iii) language.

JUMP-IR dataset using human annotators who are native speakers of different languages available in the dataset. For our annotation task, we evaluate five languages¹⁷ in total: English (en), Bengali (bn), Spanish (es), Chinese (zh) and Hindi (hi). Within the five languages, three are high-resource (en, es, zh), one is medium resource (hi) and low-resource (bn). For each language, we sample a fixed amount of question-passage pairs resulting in overall 500 question-passage pairs human evaluated. For English, Spanish and Chinese, we evaluate monolingual training pairs. For Hindi and Bengali, we mix and evaluate both cross-lingual and monolingual task-specific question-passage pairs.

We compute the question quality on a three-level rating scheme (0/1/2) based on three statistics, fluency, adequacy, and language. (i) Fluency measures the coherence of the generated question, i.e., whether the question can be perfectly understandable and readable by the user containing no spelling or grammatical mistakes. (ii) Adequacy measures the relevancy of the question with the Wikipedia passage (used for generation of the question), whether the question asked contains the answer within the passage. (iii) Language measures whether the generated question is in the correct language, or code-switching occurs in the generated question. We add these details in our annotation guidelines to teach the human annotator and attach it at the end of the Appendix section.

D.1 Human Validation Results

Table 6 shows the results of human validation across five languages on JUMP-IR. The human annotators get 99-100% for the language metric which denotes the PaLM 2 generated quality is always in the correct language. For Fluency, the major mistakes are observed in Hindi (12%), where few sampled passages in MIRACL can be too

short (2-3 words long), this confuses the PaLM 2 model which duplicates the exact text in the query. For Adequacy, we observe that in Chinese (30%) of the generated synthetic queries are not strongly related to the passage. Similar to fluency, a low adequacy is observed when the LLM-generated query is generated for a short sampled passage or when the query asks a question about a related topic which is not directly mentioned in the passage.

E Content Filtering

LLMs have been shown to generate undesirable content, particularly under conditions that prime the model with material targeted at drawing out any negative patterns or associations in the model’s training data (Gehman et al., 2020; Bender et al., 2021). We originally hoped that sampled Wikipedia passages would provide almost entirely safe material for prompting LLMs. However, for each combination of query-passage languages within JUMP-IR, we discovered that between 6–10% of the pairs contained sensitive subjects and adult content (i.e., weapons; violence and abuse; accidents and disasters; death and tragedy; war and conflict). We used the Google Cloud Natural Language content classification categories¹⁸ to identify and remove pairs when either the original sampled passage or the resulting LLM generated query has a content classification of either /Adult or any of the /Sensitive Subjects labels.

F Additional Technical Details

F.1 mContriever Pretraining

In the original implementation of mContriever (Izacard et al., 2022), the authors initialized the model using the mBERT (Devlin et al., 2019) pre-trained language model (PLM). Next, the model was jointly pre-trained on 29 languages covering the CCNet dataset (Wenzek et al., 2020) with a contrastive pre-training objective. In our implementation of mContriever, we initialize the model with the multilingual T5 (mT5) model (Xue et al., 2021). Next, we jointly pre-train the model on 101 languages¹⁹ available in mC4 (Xue et al., 2021). We sample two random non-overlapping texts from our document with a maximum size of 256 tokens. Similar to the mT5 pre-training ob-

¹⁷The authors in the paper are native speakers of the five languages chosen for evaluation: Bengali, Spanish, Chinese, Hindi and English.

¹⁸cloud.google.com/natural-language/docs/categories

¹⁹The list of all 101 languages in mC4 can be found at: www.tensorflow.org/datasets/catalog/c4

jective (Xue et al., 2021), examples are not uniformly sampled over languages, i.e., the probability that a training sample comes from a specific language is directly proportional to the amount of training data available in the language. We randomly sample a maximum of 20k samples per language and keep it as a validation subset. We optimize our mContriever model with the AdamW optimizer (Loshchilov and Hutter, 2019) with a learning rate of $1e^{-3}$, batch size of 8192, and for 600K training steps. For the first 500K steps, we pre-train with a language-mixed training objective, where a single training batch can contain examples across multiple languages. For the remaining 100k training steps, we pre-train with a language-unmixed training objective, where a single training batch contains all examples from a specific language, i.e., no mixing of different language pairs within a training batch. We internally conducted a quick evaluation of the mContriever pre-trained models with language-mixing (500k) and with both language-mixing and unmixing (600k) checkpoints. On XOR-Retrieve, we observe that the language-unmixed pre-training overall improves the model performance by 7.3 points on XOR-Retrieve.

F.2 Baseline FT Models

XOR-Retrieve. For the zero-shot baseline model, we fine-tune on the MSMARCO (Nguyen et al., 2016) dataset. Our base initialization model is mT5 (Xue et al., 2021). We use in-batch negatives, AdamW optimizer (Loshchilov and Hutter, 2019) and with a learning rate of $1e^{-3}$. The query sequence length contains a maximum sequence length of 64 tokens, whereas the document contains a maximum sequence length of 256 tokens. On MSMARCO, our models are fine-tuned with a batch size of 4096 and for 50k training steps. For our supervised fine-tuned baselines, we fine-tune on the XOR-Retrieve training dataset. The original dataset authors provide 1 hard negative per each training query in (Asai et al., 2021a). We fine-tune our baseline models on XOR-Retrieve on the triplets containing the query, positive passage and a hard negative, AdamW optimizer (Loshchilov and Hutter, 2019), learning rate of $1e^{-3}$ for a batch size of 4096 for 15K training steps.

MIRACL. For the zero-shot baseline model, we fine-tune on the MSMARCO (Nguyen et al., 2016) dataset. Details are shown above in XOR-Retrieve.

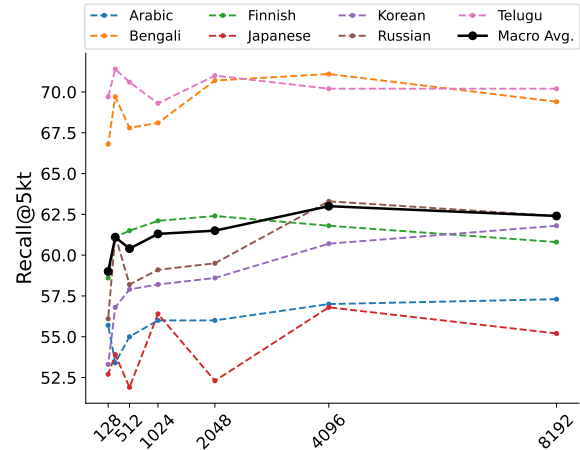


Figure 6: Training batch size ablation of JUMP-X (500K) model on XOR-Retrieve (Asai et al., 2021a). The best Recall@5kt (Macro Avg.) is achieved with batch size equal to 4096. To avoid overfitting, we fine-tune JUMP-X models with decreasing training steps of {40K, 40K, 30K, 30K, 20K, 20K, 15K} for increasing batch sizes of {128, 256, 512, 1024, 2048, 4096, 8192} respectively. We fine-tune all JUMP-X models on 500K synthetic JUMP-IR training pairs.

For the monolingual supervised models, we use the MIRACL training data for fine-tuning. The authors of MIRACL provided hard negatives for training samples. We sample up to a maximum of four hard negatives for each query and fine-tune our models on MIRACL for 15K training steps.

XTREME-UP. For the zero-shot baseline model, we fine-tune on the MSMARCO (Nguyen et al., 2016) dataset. For the supervised baselines, we use the XTREME-UP training data and fine-tune with in-batch negatives for a batch size of 1024 for 5K training steps.

F.3 Synthetic FT models

We fine-tune all JUMP-X models using in-batch negatives, AdamW optimizer (Loshchilov and Hutter, 2019) and with a learning rate of $1e^{-3}$. The pre-trained language model for JUMP-X is the mT5 Base model with 580M parameters (Xue et al., 2021). The batch size and the training steps varies for each retrieval setting. All training data is always split evenly across all languages present in the training data. For example, given 100K pairs with 5 different languages, each language includes 20K training pairs.

XOR-Retrieve. JUMP-X is fine-tuned with a batch size of 4096 and with a maximum of 50K steps on synthetic JUMP-IR cross-lingual pairs. For the 500K training pairs, we fine-tune for 20K

steps, and for the maximum of 7M pairs we fine-tune for 50K training steps. The training pairs within a single batch include language-mixing, i.e., one or more language-specific training pairs are sampled within a single training batch.

MIRACL. JUMP-X is fine-tuned for a batch-size of 4096 and for a maximum of 15K steps. As shown in (Roy et al., 2020; Zhang et al., 2023a), language-unmixed training setup is shown to work well for monolingual retrieval. Following prior work, our JUMP-X training pairs include language unmixing, i.e., all pairs are from a single language. The examples are uniformly sampled across all languages, i.e., probability that a training sample comes from a specific language is the same for all languages, unlike the previous experiment in mC4 pre-training.

XTREME-UP. JUMP-X has been fine-tuned for a batch size of 1024 and for a maximum of 15K training steps. Similar to XOR-Retrieve, training pairs include language-mixing with a single batch during fine-tuning.

F.4 Stratified Sampling in JUMP-IR

In our work, we use a stratified sampling technique to select a subset of passages from the Wikipedia corpus we use to generate questions for JUMP-IR. We ensure all languages have relatively an equal amount of training samples, wherever possible. Our Wikipedia corpus contains entities which are sorted alphabetically (A-Z). We then compute inclusion threshold I_{th} , which is defined as $I_{th} = D_{sample}/D_{total}$, where (D_{sample}) is number of passages required to sample and (D_{total}) is the total numbers of passages in corpus. Next, for each passage (p_i) in the corpus, we randomly generate an inclusion probability $\hat{p}_i \in [0, 1]$. We select the passage (p_i) if $p_i \leq I_{th}$. This ensures uniform sampling of passages with Wikipedia entities between all letters (A-Z).²⁰

G Evaluation Dataset Information

We evaluate on three multilingual retrieval benchmarks: (i) **XOR-Retrieve** (Asai et al., 2021a), (ii) **MIRACL** (Zhang et al., 2023b) and (iii) **XTREME-UP** (Ruder et al., 2023). NeuCLIR (Lawrie et al., 2023) was excluded from our evaluation as it contained a fewer subset of languages namely, Chinese (zh), Farsi (fa) and Russian (ru).

²⁰All Wikipedia entities starting with a non-alphabet are included in the beginning of the Wikipedia corpus.

Although MKQA (Longpre et al., 2021) contained a wider variety of languages, the dataset is commonly used for question-answering instead of multilingual retrieval. In Table 8, we provide an overview of the three evaluation datasets and provide statistics for each retrieval dataset.

Our three evaluation datasets contain a training split. Only XTREME-UP has released their test split publicly, as a result it was used for evaluation in the paper. However, for both XOR-Retrieve and MIRACL, we evaluate on the development split. The list of languages covered by each dataset and samples available for training and evaluation can be found in Table 8.

XOR-Retrieve (Asai et al., 2021a) is a cross-lingual open retrieval training and evaluation task within TyDi-QA (Clark et al., 2020). XOR-Retrieve contains 15K human annotated relevant passage-query pairs in the training set with one hard negative and 2K passage-answer pairs in the *dev* set. The corpus C contains 18.2M passages with a maximum of 100 word tokens from the English Wikipedia. The queries are multilingual and cover seven languages. We evaluate our models using recall at m kilo-tokens, i.e., Recall@mkt, which computes the fraction of queries for which the minimal answer is contained within the top m thousand tokens of the retrieved passages. Following prior work in Asai et al. (2021a), we evaluate our models at Recall@5kt and Recall@2kt.

MIRACL (Zhang et al., 2023b) is a monolingual open retrieval evaluation task containing 18 languages. MIRACL was developed on top of Mr. TyDi (Zhang et al., 2021), and covers more languages and provides denser judgments by human annotators. The test set is not publicly released, hence in this paper we evaluate using the *dev* set. The training set contains 88,288 pairs, with the exception of Yoruba (yo) and German (de) which do not have any training data available. The authors also provide labeled hard negatives for the training query-passage pairs. The *dev* set contains around 13,495 query-passage pairs. The corpus C in MIRACL are language-specific Wikipedia articles with various sizes starting from smallest, Yoruba (yo) with 49K passages, till the largest, English (en) with 39.2M passages. Following prior work in Zhang et al. (2023b) and Kamaloo et al. (2023), we evaluate our models at nDCG@10 and Recall@100.

XTREME-UP Ruder et al. (2023) contains di-

verse information-access and user-centric tasks focused on under-represented languages. In this paper, we evaluate cross-lingual retrieval task containing 5,280 query-passage pairs in the training set. The corpus C contains 112,426 passages sampled from TYDI-QA (Clark et al., 2020). The test set contains 10,705 query-passage pairs for evaluation. The cross-language retrieval for QA task contains 20 under-represented Indic languages. Following prior work in Ruder et al. (2023), we evaluate our models at MRR@10.

H Additional Results

XOR-Retrieve. In Table 9, we report the Recall@2kt scores across all multilingual retrievers on XOR-Retrieve. We find similar trends for improvement, where JUMP-X (7M) outperforms the best FT model on mContriever-X by 3.9 points on Recall@2kt. The JUMP-X (7M) without pre-training is also a strong baseline outperforming JUMP-X (7M) with pre-training on 4/7 languages in XOR-Retrieve.

MIRACL. In Table 10, we report the Recall@100 scores across all multilingual retrievers on MIRACL. We observe that the mContriever-X model overall achieves the highest Recall@100 score of 86.5, JUMP-X models achieve a recall of 78.9 which is competitive on MIRACL outperforming both the zero-shot mDPR and mContriever models. For Yoruba, Our JUMP-X outperforms mContriever which shows the importance of synthetic training data, as the model does not contain supervision for Yoruba (i.e., no human-labeled training pairs).

Cross-Lingual (33)		Monolingual (18)	
Q-P Lang.	# Train Pairs	Q-P Lang.	# Train Pairs
Languages available in MIRACL (Zhang et al., 2023b)			
ar-en	901,363	ar-ar	890,389
bn-en	909,748	bn-bn	257,327
de-en	909,145	de-de	943,546
en-en	-	en-en	936,481
es-en	905,771	es-es	947,340
fa-en	910,295	fa-fa	973,409
fi-en	906,429	fi-fi	967,139
fr-en	911,694	fr-fr	977,900
hi-en	919,729	hi-hi	466,272
id-en	907,826	id-id	837,459
ja-en	906,862	ja-ja	893,520
ko-en	905,669	ko-ko	941,459
ru-en	904,933	ru-ru	915,693
sw-en	905,242	sw-sw	123,099
te-en	902,190	te-te	220,431
th-en	914,610	th-th	451,540
yo-en	902,467	yo-yo	43,211
zh-en	921,701	zh-zh	946,757
Indo-European Languages in XTREME-UP (Ruder et al., 2023)			
as-en	5,899	as-as	-
bho-en	5,763	bho-bho	-
gom-en	5,755	gom-gom	-
gu-en	5,870	gu-gu	-
kn-en	5,763	kn-kn	-
mai-en	5,768	mai-mai	-
ml-en	5,907	ml-ml	-
mni-en	5,604	mni-mni	-
mr-en	5,977	mr-mr	-
or-en	5,837	or-or	-
pa-en	5,840	pa-pa	-
ps-en	5,694	ps-ps	-
sa-en	5,779	sa-sa	-
ta-en	5,930	ta-ta	-
ur-en	5,816	ur-ur	-
Total	15,532,876	Total	12,732,972
Overall Training Pairs = 28,265,848			

Table 7: Dataset Statistics of JUMP-IR for both cross-lingual and monolingual settings; (Q-P Lang.) denotes the language code of the query-passage training pair in JUMP-IR; (# Train Pairs) denotes the count of the relevant training pairs containing the synthetic query and original passage pair.

Benchmark	Retrieval Task	Query	Passage	#L	ISO	Languages	Train Split		Dev/Test Split		
							(#Queries)	(HNeg.)	(#Queries)	(#Passages)	(Metric)
XOR-Retrieve (Asai et al., 2021a)	Cross-lingual	<i>L</i>	English	7	ar, bn, fi, ja, ko, ru, te	Arabic, Bengali, Finnish, Japanese, Korean, Russian, Telugu	15,250	Yes (1 each)	2,110	18,003,200	Recall@5kt
MIRACL (Zhang et al., 2023b)	Monolingual	<i>L</i>	<i>L</i>	18	ar, bn, de, en, es, fa, fi, fr, hi, id, ja, ko, ru, sw, te, th, yo, zh	Arabic, Bengali, German, English, Spanish, Farsi, Finnish, French, Hindi, Indonesian, Japanese, Korean, Russian, Swahili, Telugu, Thai, Yoruba, Chinese	88,288	Yes (max 4)	13,495	106,332,152	nDCG@10
XTREME-UP (Ruder et al., 2023)	Cross-lingual	<i>L</i>	English	20	as, bho, brx, gbm, gom, gu, hi, hne, kn, mai, ml, mni, mr, mwr, or, pa, ps, sa, ta, ur	Assamese, Bhojpuri, Boro, Garhwali, Konkani, Gujarati, Hindi, Chhattisgarhi, Kannada, Maithili, Malayalam, Manipuri, Marathi, Marwari, Odia, Punjabi, Pashto, Sanskrit, Tamil, Urdu	13,270	No	5,300	112,426	MRR@10

Table 8: Statistics of multilingual retrieval evaluation benchmarks used in our work: XOR-Retrieve (Dev) (Asai et al., 2021a), MIRACL (Dev) (Zhang et al., 2023b) and XTREME-UP (Test) (Ruder et al., 2023). For each benchmark, we describe the retrieval task, language in which query and passage are available, train and dev/test split statistics and evaluation metric; (HNeg.) denotes availability of hard negatives for training multilingual models; (#L) denotes the number of languages covered by the benchmark.

(a) Cross-lingual Training Pair in JUMP-IR

Title: Menlo Park, New Jersey

Text: Menlo Park is an unincorporated community located within Edison Township in Middlesex County, New Jersey, United States. In 1876, Thomas Edison set up his home and research laboratory in Menlo Park, which at the time was the site of an unsuccessful real estate development named after the town of Menlo Park, California. While there, he earned the nickname "the Wizard of Menlo Park". The Menlo Park lab was significant in that it was one of the first laboratories to pursue practical commercial applications of research. It was in his Menlo Park laboratory that Thomas Edison invented the phonograph and developed it.

Passage (ID: 10770836) from English Wikipedia (en)

托马斯·爱迪生在哪里发明了留声机？

Translation: (Where did Thomas Edison invent the phonograph?)

LLM-generated Query in Chinese (zh)

(b) Monolingual Training Pair in JUMP-IR

Title: En la tierra del Guaraní

Text: Es considerada una de las primeras realizaciones sonoras de la región y uno de los primeros antecedentes de cooperación entre dos países de la zona (Paraguay y Argentina) para la realización de un filme.

Translation: (*In the land of Guaraní*: It is considered one of the first sound productions in the region and one of the first precedents of cooperation between two countries in the area (Paraguay and Argentina) for the making of a film.)

Passage (ID:spanish_5170543#3) from Spanish Wikipedia (es)

¿Qué película es una de las primeras realizaciones sonoras de la región?

Translation: (What film is one of the first sound films in the region?)

LLM-generated Query in Spanish (es)

Figure 7: Dataset examples showing both (a) cross-lingual and (b) monolingual training pairs in the JUMP-IR dataset. The passage is selected from English Wikipedia, and PaLM 2 generates the query. A detailed description of all the dataset column headers are provided in Appendix (§C.2). All translations in the figure above have been provided using Google Translate (translate.google.com) for illustration purposes.

Model	PLM	PT	Finetune (Datasets)	Recall@2kt							
				Avg.	Ar	Bn	Fi	Ja	Ko	Ru	Te
Existing Supervised Baselines (Prior work)											
Dr. DECR (Li et al., 2022)	XLM-R	WikiM	NQ + XOR*	66.0	—	—	—	—	—	—	—
mDPR (Asai et al., 2021a)	mBERT	—	XOR	40.5	38.8	48.4	52.5	26.6	44.2	33.3	39.9
mBERT + xQG (Zhuang et al., 2023)	mBERT	—	XOR	46.2	42.4	54.9	54.1	33.6	52.3	33.8	52.5
Google MT + DPR (Asai et al., 2021a)	BERT	—	NQ	62.2	62.5	74.7	57.3	55.6	60.0	52.7	72.3
OPUS MT + DPR (Asai et al., 2021a)	BERT	—	NQ	42.7	43.4	53.9	55.1	40.2	50.5	30.8	20.2
Zero-shot baselines (English-only supervision)											
mContriever	mT5	mC4	—	29.9	27.2	23.0	35.0	27.0	27.7	35.0	34.0
mDPR (En)	mT5	—	MS MARCO	30.6	26.2	26.0	37.9	32.8	24.6	34.6	32.4
mContriever (En)	mT5	mC4	MS MARCO	33.8	27.8	24.3	42.4	29.9	31.2	40.5	40.3
Supervised Baselines (Cross-lingual supervision)											
mDPR-X	mT5	—	XOR	43.6	43.7	50.0	44.6	36.1	41.1	35.9	54.2
mContriever-X	mT5	mC4	XOR	46.6	40.1	62.5	47.1	38.2	44.2	38.4	55.5
mDPR-X	mT5	—	MS MARCO + XOR	49.5	46.0	63.8	49.0	39.0	48.4	43.9	56.3
mContriever-X	mT5	mC4	MS MARCO + XOR	53.0	47.6	65.1	51.6	47.3	50.2	44.3	65.1
Synthetic Baselines (Our work)											
JUMP-X (500K)	mT5	—	JUMP-IR	49.2	46.3	57.2	49.0	42.7	45.6	44.7	58.8
JUMP-X (500K)	mT5	mC4	JUMP-IR	53.3	46.6	61.8	51.9	46.5	49.1	55.3	61.8
JUMP-X (7M)	mT5	—	JUMP-IR	56.6	50.8	65.1	56.1	48.1	54.0	55.7	66.4
JUMP-X (7M)	mT5	mC4	JUMP-IR	56.9	53.4	67.8	55.1	49.4	52.6	55.3	64.7

Table 9: Experimental results showing Recall@2kt for cross-lingual retrieval on XOR-Retrieve dev (Asai et al., 2021a); (PLM) denotes the pretrained language model; (PT) denotes the pretraining dataset; (*) Dr.DECR is fine-tuned in a complex training setup across more datasets (§3.3); WikiM denotes WikiMatrix (Schwenk et al., 2021); XOR denotes XOR-Retrieve; JUMP-X (ours) is fine-tuned on 500K and 7M synthetic data.

Model	Avg.	ar	bn	en	es	fa	fi	fr	hi	id	ja	ko	ru	sw	te	th	zh	de	yo
<i>Existing Supervised Baselines (Prior work)</i>																			
BM25	77.2	88.9	90.9	81.9	70.2	73.1	89.1	65.3	86.8	90.4	80.5	78.3	66.1	70.1	83.1	88.7	56.0	57.2	73.3
mDPR	79.0	84.1	81.9	76.8	86.4	89.8	78.8	91.5	77.6	57.3	82.5	73.7	79.7	61.6	76.2	67.8	94.4	89.8	79.5
Hybrid	88.0	94.1	93.2	88.2	94.8	93.7	89.5	96.5	91.2	76.8	90.4	90.0	87.4	72.5	85.7	82.3	95.9	88.9	80.7
Cohere-API	76.9	85.4	85.6	74.6	71.7	77.1	80.9	81.6	72.4	68.3	81.6	77.1	76.7	66.6	89.8	86.9	76.9	72.5	57.6
<i>Zero-shot baselines (English-only supervision)</i>																			
mDPR (En)	76.9	85.5	85.9	72.4	66.8	79.7	86.0	71.4	74.2	67.0	80.1	77.1	77.4	80.2	91.9	84.8	68.5	70.9	58.6
mContriever (En)	76.6	73.5	80.8	52.1	49.5	61.7	66.0	51.8	50.3	63.5	65.6	56.3	58.9	73.5	85.9	76.6	58.2	36.3	30.2
<i>Supervised Baselines (Monolingual supervision)</i>																			
mDPR-X	60.6	73.5	80.8	52.1	49.5	61.7	66.0	51.8	50.3	63.5	65.6	56.3	58.9	73.5	85.9	76.6	58.2	36.3	30.2
mContriever-X	86.5	92.0	95.3	80.6	78.8	84.0	93.1	86.0	82.1	83.7	89.5	87.7	86.7	93.3	96.7	94.3	85.9	79.3	68.8
<i>Synthetic Baselines (Our work)</i>																			
JUMP-X (180K)	78.9	89.2	87.8	72.9	70.0	76.3	91.6	75.8	72.5	74.3	77.6	76.8	77.9	87.8	84.9	92.9	69.9	72.4	69.3

Table 10: Experimental results for monolingual retrieval on MIRACL dev (Zhang et al., 2023b). All scores denote Recall@100; (Hyb.) denotes Hybrid retriever with ranked fusion of three retrievers: mDPR, mColBERT and BM25; BM25, mDPR and Hybrid scores taken from (Zhang et al., 2023b); Cohere-API is used as a reranker on top of 100 BM25 results, taken from (Kamalloo et al., 2023). JUMP-X is fine-tuned on 180K synthetic data.

5-shot Summarize-then-Ask Prompting for XOR-Retrieve

Read the following article and write a factual summary. Your summary will act as a surrogate for asking a question based on the article. Finally, translate the question to **Bengali**.

Article: Long Lost Family is a BAFTA award winning British television series that has aired on ITV since 21 April 2011. The programme, which is presented by Davina McCall and Nicky Campbell, aims to reunite close relatives after years of separation. It is made by the production company Wall to Wall. "Long Lost Family" is based on the Dutch series "Sporloos" (), airing on NPO 1 since February 1990 and it is made by KRO-NCRV. Presented by Davina McCall and Nicky Campbell, the series offers a last chance for people who are desperate to find long lost relatives.

Summary: Long Lost Family is a BAFTA award winning British television series aired since 2011. The series aim to reunite close relatives after years of separation which is presented by Davina McCall and Nicky Campbell.

Question [Bengali]: ব্রিটিশ টেলিভিশন সিরিজ লং লস্ট ফ্যামিলি কোন পুরস্কার জিতেছে?

Article: Muscular activity accounts for much of the body's energy consumption. All muscle cells produce adenosine triphosphate (ATP) molecules which are used to power the movement of the myosin heads. Muscles have a short-term store of energy in the form of creatine phosphate which is generated from ATP and can regenerate ATP when needed with creatine kinase. Muscles also keep a storage form of glucose in the form of glycogen. Glycogen can be rapidly converted to glucose when energy is required for sustained, powerful contractions. Within the voluntary skeletal muscles, the glucose molecule can be metabolized anaerobically in a process.

Summary: All muscle cells produce adenosine triphosphate (ATP) molecules for movement of myosin heads. A short term store of energy is generated from ATP in the form of creatine phosphate and can regenerate ATP when needed with creatine kinase.

Question [Bengali]: কীভাবে পেশী কোষগুলি মায়োসিন মাথার নড়াচড়ার জন্য শক্তিকে শক্তি দেয়?

Article: The 1960s brought anime to television and in America. The first anime film to be broadcast was "Three Tales" in 1960. The following year saw the premiere of Japan's first animated television series, "Instant History", although it did not consist entirely of animation. Osamu Tezuka's "Tetsuwan Atom" ("Astro Boy") is often miscredited as the first anime television series, premiering on January 1, 1963. "Astro Boy" was highly influential to other anime in the 1960s, and was followed by a large number of anime about robots or space.

Summary: First anime movie broadcast on TV was 'Three Tales' in 1960. First anime TV series was 'Instant History' in 1961. 'Astro Boy' first aired in 1963 was a highly influential anime about robots or space.

Question [Bengali]: ১৯৬০ সালে টিভিতে সম্প্রচারিত প্রথম অ্যানিমে ছবি কোনটি?

Article: Łęczna is a town in eastern Poland with 19,780 inhabitants (2014), situated in Lublin Voivodeship. It is the seat of Łęczna County and the smaller administrative district of Gmina Łęczna. The town is located in northeastern corner of historic province of Lesser Poland. Łęczna tops among the hills of the Lublin Upland, at the confluence of two rivers —the Wieprz, and the Świnka. On December 31, 2010, the population of the town was 20,706. Łęczna does not have a rail station, the town has been placed on a national Route 82 from Lublin to Włodawa. And shall be considered as a

Summary: Łęczna is a town in eastern Poland with 19,780 inhabitants. It is a hill located in the Lublin Upland, at the confluence of two rivers - Wieprz and Świnka. It is a road hub, and has no rail station.

Question [Bengali]: লিচেনা পোল্যান্ডের কোন দুটি নদীর সম্মুখস্থ অবস্থিত?

Article: The μ -law algorithm (sometimes written "mu-law", often approximated as "u-law") is a companding algorithm, primarily used in 8-bit PCM digital telecommunication systems in North America and Japan. It is one of two versions of the G.711 standard from ITU-T, the other version being the similar A-law, used in regions where digital telecommunication signals are carried on E-1 circuits, e.g. Europe. Companding algorithms reduce the dynamic range of an audio signal. In analog systems, this can increase the signal-to-noise ratio (SNR) achieved during transmission; in the digital domain, it can reduce the quantization error (hence increasing signal to quantization noise ratio).

Summary: The μ -law algorithm is a companding algorithm, which is used to reduce the dynamic range of audio signals. In analog systems, this can increase the signal-to-noise ratio (SNR) achieved during transmission.

Question [Bengali]: μ -আইন অ্যালগরিদম কীভাবে অ্যানালগ সিস্টেমে সংক্রমণকে প্রভাবিত করে?

Article: {Input Wikipedia Article in English}

Summary:

Figure 8: 5-shot SAP (*Summarize-then-Ask Prompting*) for XOR-Retrieve (Asai et al., 2021a) is shown for Bengali (bn). There are five exemplars (5-shot) in our cross-lingual question generation task. The passages are randomly selected from XOR-Retrieve. Summaries and questions are manually written in English by the authors. Finally, the questions in exemplars are translated to Bengali using Google Translate (translate.google.com).

3-shot Summarize-then-Ask Prompting for MIRACL

Read the following article in **Chinese** and write a factual summary in **Chinese**. Your summary will act as a surrogate for asking a question in **Chinese** based on the article.

Article: 四川各地小吃通常也被看作是川菜的组成部分。由于重庆地区小吃相对较少，除重庆麻辣小面外，川菜小吃主要以成都小吃为主。主要有担担面、川北凉粉、麻辣小面、酸辣麵、酸辣粉、叶儿粑、酸辣豆花、三合泥、红油抄手等以及用创始人姓氏命名的赖汤圆、龙抄手、钟水饺、吴抄手等。甜品方面，以原产四川眉山的冰粉和四川宜宾长宁县的凉糕最有名。

Summary: 四川美食种类繁多，小吃也非常有名，主要有担担面、川北凉粉、麻辣小面、酸辣粉、叶儿粑、酸辣豆花、三合泥、红油抄手、赖汤圆、龙抄手、钟水饺、吴抄手等。甜品方面，以原产四川眉山的冰粉和四川宜宾长宁县的凉糕最有名。

Question [Chinese]: 四川美食有哪些？

Article: 狮子座流星雨 (Leonids['li.ə.nɪdz] \ˈlee-uhˌnɪdz\')是與周期大約33年的坦普爾·塔特爾彗星有關的一個流星雨。狮子座流星雨的得名是因為這個流星雨輻射點的位置在獅子座。在2009年，這個流星雨的尖峰時間在11月17日（世界時），每小時的數量可能高達500顆，尚不足以成為流星暴（每小時超過1,000顆流星的大流星雨）。

Summary: 上一次狮子座流星雨发生在2009年11月17日。狮子座流星雨是与周期大约33年的坦普尔·塔特爾彗星有关的一个流星雨。狮子座流星雨的得名是因为这个流星雨辐射点的位置在狮子座。

Question [Chinese]: 上一次狮子座流星雨发生在什么时候？

Article: 清华大学（，縮寫：），简称清华，舊称清华学堂、游美肄业馆、清华学校、國立清華大學，是一所位于中华人民共和国北京市海淀区清华园的公立大学。始建于1911年，因北京西北郊清华园而得名。初为清政府利用美国退还的部分庚子赔款所建留美预备学校“游美学务处”及附设“肄业馆”，於1925年始设大学部。抗日战争爆发后，清华与北大、南开南迁长沙，组建国立长沙临时大学。1938年再迁昆明，易名国立西南联合大学。1946年迁回清华园复校，拥有文、法、理、工、农等5个学院。1949年中华人民共和国成立后，国立清华大学归属中央人民政府教育部，更名“清华大学”；而原国立清华大学校长梅贻琦于1955年在台湾新竹复校，仍沿用原名。

Summary: 清华大学始建于1911年，因北京西北郊清华园而得名。初为清政府利用美国退还的部分庚子赔款所建留美预备学校“游美学务处”及附设“肄业馆”。

Question [Chinese]: 清华大学什么时候成立的？

Article: {Input Wikipedia Article in Chinese}

Summary:

Figure 9: 3-shot SAP (*Summarize-then-Ask Prompting*) for MIRACL (Zhang et al., 2023b) is shown for Chinese (zh). There are three exemplars (3-shot) in our monolingual question generation task. The query-passage pairs are randomly selected from MIRACL training set. Finally, the summaries in exemplars are automatically generated using Google Bard (bard.google.com).

5-shot Summarize-then-Ask Prompting for XTREME-UP

Read the following article and write a factual summary. Your summary will act as a surrogate for asking a question based on the article. Finally, translate the question to Hindi.

Article: Long Lost Family is a BAFTA award winning British television series that has aired on ITV since 21 April 2011. The programme, which is presented by Davina McCall and Nicky Campbell, aims to reunite close relatives after years of separation. It is made by the production company Wall to Wall. "Long Lost Family" is based on the Dutch series "Sporloos" (), airing on NPO 1 since February 1990 and it is made by KRO-NCRV. Presented by Davina McCall and Nicky Campbell, the series offers a last chance for people who are desperate to find long lost relatives.

Summary: Long Lost Family is a BAFTA award winning British television series aired since 2011. The series aim to reunite close relatives after years of separation which is presented by Davina McCall and Nicky Campbell.

Question [Hindi]: ब्रिटिश टेलीविजन लॉन्ग लॉस्ट फैमिली ने कौन सा पुरस्कार जीता?

Article: Muscular activity accounts for much of the body's energy consumption. All muscle cells produce adenosine triphosphate (ATP) molecules which are used to power the movement of the myosin heads. Muscles have a short-term store of energy in the form of creatine phosphate which is generated from ATP and can regenerate ATP when needed with creatine kinase. Muscles also keep a storage form of glucose in the form of glycogen. Glycogen can be rapidly converted to glucose when energy is required for sustained, powerful contractions. Within the voluntary skeletal muscles, the glucose molecule can be metabolized anaerobically in a process.

Summary: All muscle cells produce adenosine triphosphate (ATP) molecules for movement of myosin heads. A short term store of energy is generated from ATP in the form of creatine phosphate and can regenerate ATP when needed with creatine kinase.

Question [Hindi]: मायोसिन हेड्स की गति के लिए मांसपेशियों की कोशिकाएं ऊर्जा को कैसे शक्ति देती हैं?

Article: The 1960s brought anime to television and in America. The first anime film to be broadcast was "Three Tales" in 1960. The following year saw the premiere of Japan's first animated television series, "Instant History", although it did not consist entirely of animation. Osamu Tezuka's "Tetsuwan Atom" ("Astro Boy") is often miscredited as the first anime television series, premiering on January 1, 1963. "Astro Boy" was highly influential to other anime in the 1960s, and was followed by a large number of anime about robots or space.

Summary: First anime movie broadcast on TV was 'Three Tales' in 1960. First anime TV series was 'Instant History' in 1961. 'Astro Boy' first aired in 1963 was a highly influential anime about robots or space.

Question [Hindi]: १९६० में टीवी पर प्रसारित होने वाली एनीमे फिल्म कौन सी थी?

Article: Łęczna is a town in eastern Poland with 19,780 inhabitants (2014), situated in Lublin Voivodeship. It is the seat of Łęczna County and the smaller administrative district of Gmina Łęczna. The town is located in northeastern corner of historic province of Lesser Poland. Łęczna tops among the hills of the Lublin Upland, at the confluence of two rivers —the Wieprz, and the Świnka. On December 31, 2010, the population of the town was 20,706. Łęczna does not have a rail station, the town has been placed on a national Route 82 from Lublin to Włodawa. And shall be considered as a

Summary: Łęczna is a town in eastern Poland with 19,780 inhabitants. It is a hill located in the Lublin Upland, at the confluence of two rivers - Wieprz and Świnka. It is a road hub, and has no rail station.

Question [Hindi]: लेक़ज़ना पोलैंड में किन दो नदियों के संगम पर स्थित है?

Article: The μ -law algorithm (sometimes written "mu-law", often approximated as "u-law") is a companding algorithm, primarily used in 8-bit PCM digital telecommunication systems in North America and Japan. It is one of two versions of the G.711 standard from ITU-T, the other version being the similar A-law, used in regions where digital telecommunication signals are carried on E-1 circuits, e.g. Europe. Companding algorithms reduce the dynamic range of an audio signal. In analog systems, this can increase the signal-to-noise ratio (SNR) achieved during transmission; in the digital domain, it can reduce the quantization error (hence increasing signal to quantization noise ratio).

Summary: The μ -law algorithm is a companding algorithm, which is used to reduce the dynamic range of audio signals. In analog systems, this can increase the signal-to-noise ratio (SNR) achieved during transmission.

Question [Hindi]: कैसे μ -नियम एल्गोरिथम एनालॉग सिस्टम में संचरण को प्रभावित करता है?

Article: {Input Wikipedia Article in English}

Summary:

Figure 10: 5-shot SAP (*Summarize-then-Ask Prompting* with Machine Translation (MT) for XTREME-UP (Ruder et al., 2023) is shown for Hindi (hi). There are five exemplars (5-shot) in our cross-lingual question generation. The passages are re-used from the XOR-Retrieve task. Summaries and questions are manually written in English by the authors. Finally, the questions in exemplars are translated to Hindi using Google Translate (translate.google.com).

Annotation Guidelines for JUMP-IR

- The goal of this task is to evaluate the quality of LLM-generated (PaLM 2-S) generated questions.
- Every annotator will receive a set of annotations containing the wikipedia paragraph and the question in the $\${target_language}$.
- Annotators should read each annotation carefully and provide feedback on the following:
 - The **fluency** of the question.
 - The **adequacy** of the question.
 - The **language** of the question.
- Annotators should be respectful and professional in their feedback.
- Annotators should complete all annotations within the allotted duration.

Here below we define the following terms:

Fluency

Rating Level	Explanation
2 (Flawless)	Perfect use of $\${target_language}$ with no mistakes at all.
1 (Good)	Few or minor spelling or grammar mistakes; the text is still mostly understandable and readable.
0 (Poor)	Many or serious spelling, grammar, or other mistakes, which make the text difficult to understand or hard to read.

Adequacy

Rating Level	Explanation
2 (Relevant)	Highly related to the wiki passage. The question can be answered using the wiki passage.

1 (Moderate)	The question is somewhat related to the wiki paragraph, the question cannot be answered using the passage.
0 (Not Relevant)	The question is not at all related to the wiki passage.

Language

Rating Level	Explanation
2 (Flawless)	The whole question is perfectly in the <code>\${target_language}</code> .
1 (Good)	Code-switching occurs with part of the question in the <code>\${target_language}</code> .
0 (Poor)	The whole question is not at all in <code>\${target_language}</code> .

Thank you for your participation in this task!