# WHAT DRIVES CROSS-LINGUAL RANKING? RETRIEVAL APPROACHES WITH MULTILINGUAL LANGUAGE MODELS

**Roksana Goworek**[*]
The Alan Turing Institute
Queen Mary University of London

**Olivia Macmillan-Scott**[*]
The Alan Turing Institute
University College London

**Eda B. Özyiğit**
The Alan Turing Institute

{rgowerek, omacmillan-scott, eozyigit}@turing.ac.uk

## ABSTRACT

Cross-lingual information retrieval (CLIR) enables access to multilingual knowledge but remains challenging due to disparities in resources, scripts, and weak cross-lingual semantic alignment in embedding models. Existing pipelines often rely on translation and monolingual retrieval heuristics, which add computational overhead and noise, degrading performance. This work systematically evaluates four intervention types, namely document translation, multilingual dense retrieval with pretrained encoders, contrastive learning at word, phrase, and query–document levels, and cross-encoder re-ranking, across three benchmark datasets. We find that dense retrieval models trained specifically for CLIR consistently outperform lexical matching methods and derive little benefit from document translation. Contrastive learning mitigates language biases and yields substantial improvements for encoders with weak initial alignment, and re-ranking can be effective, but depends on the quality of the cross-encoder training data. Although high-resource languages still dominate overall performance, gains over lexical and document-translated baselines are most pronounced for low-resource and cross-script pairs. These findings indicate that cross-lingual search systems should prioritise semantic multilingual embeddings and targeted learning-based alignment over translation-based pipelines, particularly for cross-script and under-resourced languages.

## 1 Introduction

Cross-lingual information retrieval (CLIR) aims to retrieve documents written in one language in response to queries expressed in another [1, 2]. As global information ecosystems expand, CLIR has become essential for equitable access to knowledge across linguistic boundaries [3, 4, 5]. Nevertheless, it remains a challenging problem because of disparities in resources, script variation, morphological complexity, and uneven pretraining coverage across languages [6, 7, 8]. Methods that perform well in monolingual information retrieval often degrade in multilingual settings, as lexical mismatch and typological distance hinder effective comparison of queries and documents.

A common strategy in practical CLIR systems is to translate queries or documents into a pivot language, typically English, and then apply monolingual retrieval models [1, 9]. Although straightforward, this approach introduces noise from machine translation (MT), increases latency in large-scale systems, and does not guarantee that the resulting representations are aligned with the semantic space used for ranking [10, 11, 12]. Multilingual encoders provide a compelling alternative by embedding texts from different languages into a shared vector space, yet prior work has shown that their cross-lingual alignment quality varies widely across languages and domains [13, 14].

Several linguistic and data-related factors compound the difficulty of CLIR. Lexical overlap is often non-existent across scripts; typological distance and morphological divergence reduce cross-lingual similarity; and pretraining corpora remain heavily skewed towards high-resource languages [13]. As a consequence, CLIR effectiveness can vary dramatically across language pairs, even for state-of-the-art multilingual encoders [15].

---

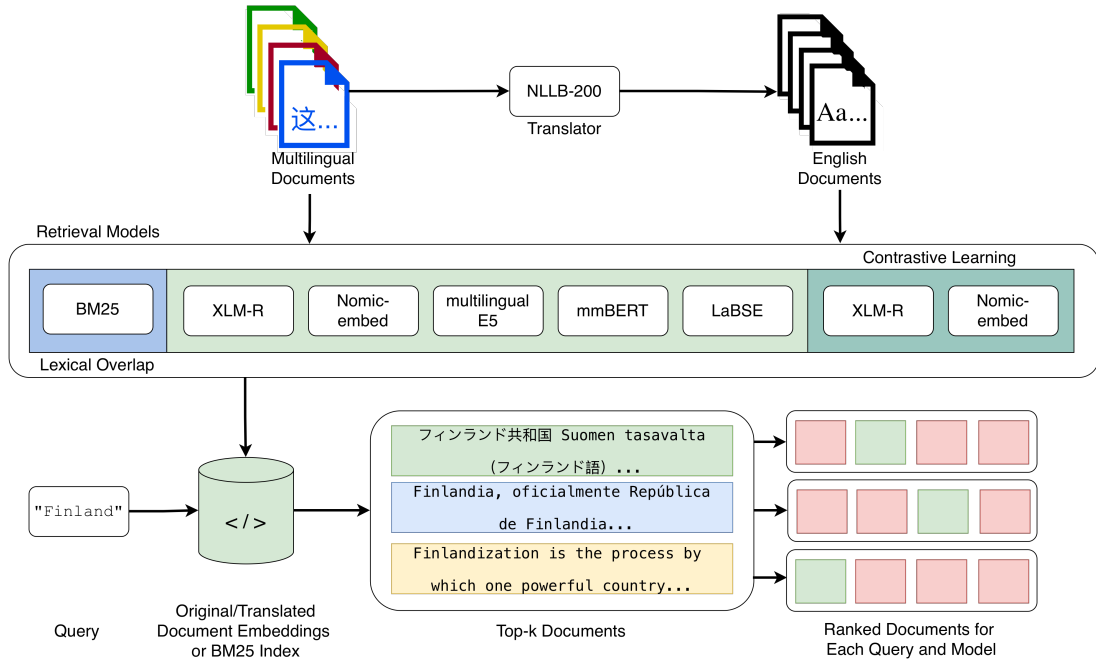[*]Work done during an internship at The Alan Turing Institute.

Figure 1: Overview of the CLIR ranking evaluation pipeline.

Prior research has examined individual components of CLIR pipelines in isolation, such as translation quality [16], dense retrieval [5], or cross-encoder re-ranking [17], and there is still limited understanding of how different interventions compare, especially when they are applied independently under controlled experimental conditions. In particular, it remains unclear when translation helps or harms retrieval, under which conditions contrastive alignment yields meaningful gains, whether re-ranking is effective for cross-script retrieval, and whether approximate nearest neighbour (ANN) [18] search offers efficiency benefits at the scale of multilingual document collections.

To address these gaps, this work evaluates several isolated interventions designed to improve encoder-based CLIR, focusing on four approaches: (i) embedding-based retrieval, comparing multilingual embedding models on native-language and English-translated corpora; (ii) contrastive alignment, applied at multiple linguistic granularities, namely, word, phrase, and alignment between queries and documents; and (iii) cross-encoder re-ranking with easy or hard negatives—where easy negatives are randomly sampled irrelevant documents, and hard negatives are documents that appear relevant to the query but are incorrect—used to train the cross-encoder to better discriminate fine-grained relevance; and a comparison of efficiency and effectiveness using ANN indexing. This experimental design makes it possible to determine when and why particular techniques yield improvements, and how these gains depend on linguistic factors, model alignment quality, and the underlying information retrieval architecture.

Figure 1 illustrates the cross-lingual information retrieval evaluation pipeline used in this study. Multilingual documents are first translated into English using a machine translation system, for example NLLB-200 [19]. For each retrieval model, separate document embeddings and BM25 indices are then constructed over both the original and the translated corpora. During retrieval, the system selects the top 100 candidate documents for each query using either cosine similarity in the embedding space or BM25 lexical matching. Retrieval effectiveness is assessed over 1,000 or 500 queries per language pair using standard information retrieval metrics. This pipeline provides a unified framework for comparing lexical and semantic methods and for analysing how multilingual semantic representations influence retrieval performance across diverse languages, scripts and resource conditions. The main contributions are as follows:

- **Evaluation of embedding-based CLIR methods.** This work presents a comparative evaluation of lexical, translation-based and multilingual dense retrieval, and re-ranking techniques across three diverse datasets, examining their behaviour under varying resource levels and script conditions.

- **Semantic alignment in CLIR.** It analyses the role of semantic alignment [5] in cross-lingual ranking by contrasting lexical overlap, translation and multilingual semantic modelling, and identifies the conditions under which contrastive alignment and cross-encoder re-ranking are most effective.

2

- **Linguistic and efficiency factors.** It quantifies retrieval biases and typological correlations across language pairs and investigates efficiency mechanisms in encoder-based models.

The remainder of this paper is structured as follows. Section 2 surveys the related work, beginning with lexical and translation-based approaches and outlining their limitations in cross-lingual settings, before introducing neural methods and semantic representations that currently underpin progress in the area. Section 3 describes the experimental setup, including datasets, retrieval models, training configurations and evaluation protocols, so as to support clarity and reproducibility. Section 4 reports the empirical results and examines the effect of each intervention. It first highlights the limitations of lexical BM25 retrieval, including under document translation, and then contrasts these findings with the substantially stronger performance of multilingual language models, followed by an analysis of the contributions of contrastive learning and cross-encoder re-ranking to retrieval effectiveness. Section 5 discusses the broader implications of the findings and presents an additional analysis of the linguistic factors that influence cross-lingual retrieval. Section 6 concludes the paper by summarising the main insights and indicating directions for future work.

## 2  Related Work

This section reviews prior research across three core dimensions of cross-lingual information retrieval (CLIR): (i) retrieval strategies including translation-based methods and modern sparse, dense, and hybrid architectures; (ii) multilingual semantic representations and alignment; and (iii) ranking strategies, including recent developments leveraging multilingual large language models (mLLMs).

**Ranking Strategies.**  Early approaches to CLIR relied on lexical matching, typically by translating queries into the document language using bilingual dictionaries or statistical machine translation, followed by monolingual ranking models such as BM25 or query likelihood formulations [3, 20, 21]. Document translation into a pivot language was also explored as an alternative, often yielding better disambiguation due to richer context [21, 22]. However, translation-based pipelines are hindered by ambiguity, limited vocabulary coverage, and poor robustness to morphological variation, particularly for low-resource languages [23, 24]. In addition, the translation step introduces latency and inconsistency in large-scale systems [11, 25], making explicit query translation less feasible in operational settings. Consequently, recent research has shifted towards language-agnostic retrieval strategies that minimise dependence on machine translation by using shared embedding spaces.

Modern CLIR systems increasingly adopt neural retrieval models. Sparse neural retrievers such as SPLADE [26] and SPLADE-v2 [27] enhance lexical matching via learned expansion, while dense retrieval approaches embed queries and documents into a shared vector space using multilingual encoders and rank via cosine similarity [5, 28, 29]. Hybrid strategies combine sparse and dense signals to balance precision and semantic recall [30, 31], with multilingual variants extending applicability across languages [32]. To support scalability, particularly in dense retrieval, ANN search is employed to reduce inference overhead [33, 34], though its behaviour under cross-lingual conditions is less thoroughly explored [35]. Multi-stage architectures are now prevalent, using sparse or dense retrieval for candidate generation, followed by more computationally intensive ranking models. Cross-encoders jointly encode query and document for deeper token interaction and typically achieve high effectiveness [36, 37], though at significant computational cost. Retrieval effectiveness, therefore, depends not only on the architecture but also on the quality of the underlying multilingual representation, which motivates the need for robust embedding alignment.

**Multilingual Representations and Alignment.**  Multilingual encoders such as XLM-R [6] and mBERT [38] enabled shared embedding spaces for cross-lingual Natural Language Processing (NLP); however, alignment quality varies considerably across languages and typologies [13, 14]. Retrieval-oriented encoders improve alignment using ranking objectives and supervision from parallel or translation-based data, as in LaBSE [39], LASER [40], and multilingual-E5 [41], achieving strong results on multilingual benchmarks such as MMTEB [42]. Contrastive learning is widely adopted in dense retrieval [43], with supervision sourced from sentence-level alignment, translation pairs, or relevance labels from datasets such as mMARCO [28], CLIRMatrix [29], and MIRACL [44]. More recent work explores fine-grained alignment at token level using word sense or semantic similarity resources [45, 46, 47].

Although such models are increasingly used as the basis for the first-stage retrieval, multi-stage architectures rely on more expressive re-ranking mechanisms to refine document order, using their output as the candidate pool for more computationally intensive models.

**Re-ranking with mLLMs.**  Recent work has applied mLLMs to CLIR re-ranking by exploiting their generative and reasoning capabilities. This includes query reformulation and expansion, relevance estimation via generative scoring or answer justification, and use as cross-lingual re-rankers that capture deeper semantic relations [48, 49]. Beyond

re-ranking, mLLMs have been used to reason directly over candidate passages, or to generate synthetic training data for retrievers [50, 51, 52]. mLLM-based re-ranking methods can be categorised as listwise, pairwise or pointwise. Listwise approaches prompt the model with multiple candidate passages simultaneously to estimate a ranked order [53, 54]. Open-source variants such as RankVicuna and RankZephyr [55, 56], and distillation-based methods like Rank-without-GPT [57], demonstrate competitive performance on English benchmarks. Pairwise reranking compares candidate documents in pairs [58]. Pointwise methods independently assess query–document relevance or generate query expansions [59, 60, 53]. While these approaches improve cross-lingual semantic matching, they incur high computational cost, latency and reproducibility challenges. As a result, mLLMs are primarily used at the re-ranking or post-retrieval stage rather than for first-stage retrieval. Our work focuses on pointwise re-ranking architectures that are more computationally scalable while preserving multilingual effectiveness.

# 3 Experimental Settings

Cross-lingual retrieval is evaluated within a unified encoder-based pipeline. Four families of interventions are examined: (i) the effect of document translation on lexical and semantic retrieval; (ii) the impact of contrastive learning on multilingual alignment; (iii) the contribution of re-ranking over the first-stage retrieval; and (iv) the trade-off between retrieval effectiveness and efficiency.

## 3.1 Datasets

Experiments are conducted on three established datasets: CLIRMatrix [29], mMARCO [28] and the Large-Scale CLIR dataset [1]. CLIRMatrix covers all combinations of eight query and document languages, excluding same-language pairs, with documents from multilingual Wikipedia. mMARCO provides translations of MS MARCO passages into 14 languages and supports all 14×14 language-pair combinations. The Large-Scale CLIR dataset consists of English queries paired with documents in 26 other languages.

| Dataset | AR | CA | CS | DE | EN | ES | FI | FR | HI | ID | IT | JA | KO | NL | NN | NO | PL | PT | RO | RU | SV | SW | TL | TR | UK | VI | ZH |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| CLIRMatrix | ✓ | | | ✓ | ✓ | ✓ | | ✓ | | | | ✓ | | | | | | | | ✓ | | | | | | | ✓ |
| mMARCO | ✓ | | | ✓ | ✓ | ✓ | | ✓ | ✓ | ✓ | ✓ | ✓ | | ✓ | | | | ✓ | | ✓ | | | | | | ✓ | ✓ |
| Large-Scale | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ |

Table 1: Language coverage of the three datasets. Columns list 2-letter ISO codes. A tick (✓) indicates inclusion. mMARCO and CLIRMatrix are cross-lingual (CLIRMatrix excludes same-language pairs), while Large-Scale provides only English queries, paired with documents in the listed languages.

Table 1 summarises the language coverage of the three datasets, with both queries and documents available in all languages with ✓ in Table 1 for CLIRMatrix and mMARCO, but only English queries, matched with documents in all the specified languages for Large-Scale.

This work uses the following ISO 639-1 language codes in the experiments: AR (Arabic), CA (Catalan), CS (Czech), DE (German), EN (English), ES (Spanish), FI (Finnish), FR (French), HI (Hindi), ID (Indonesian), IT (Italian), JA (Japanese), KO (Korean), NL (Dutch), NN (Norwegian Nynorsk), NO (Norwegian Bokmål), PL (Polish), PT (Portuguese), RO (Romanian), RU (Russian), SV (Swedish), SW (Swahili), TL (Tagalog), TR (Turkish), UK (Ukrainian), VI (Vietnamese), and ZH (Chinese).

To ensure comparability, all datasets are evaluated with a top-100 retrieval depth. In CLIRMatrix, 1,000 queries are sampled for each of the 56 language pairs, yielding 56,000 query instances and 9,055 unique documents (1,080–1,205 per language). In mMARCO, translations are deduplicated across languages and the resulting 7,433 unique documents are evenly redistributed across all 14 languages, assigning 530–531 queries per pair for a total of 97,720 query–document pairs. In the Large-Scale dataset, each English→X language pair includes 1,000 queries, giving 26,000 query pairs overall. For each dataset, top-100 document retrieval for each query is performed over the full multilingual document pool, and all documents are truncated to 512 tokens.

## 3.2 Retrieval Pipeline

This section introduces the pretrained multilingual encoders used as retrieval models in the cross-lingual experiments.

**Retrieval Models.** Five multilingual encoders are considered, spanning diverse pretraining and alignment paradigms and covering around 100 languages each. XLM-R [6] serves as a strong masked language modelling baseline, trained

on 100 languages with approximately 550 million parameters. LaBSE [39] is a dual encoder trained with a translation-ranking objective, supporting 109 languages with 471 million parameters. Multilingual-E5 [41] is an information retrieval-oriented extension of the E5 model with multilingual instruction tuning, covering more than 100 languages and comprising 560 million parameters. Nomic-embed-text-v2-moe [61] ("Nomic") is a sparse mixture-of-experts model with coverage of roughly 100 languages and an effective parameter budget of about 475 million. Finally, mmBERT [62] is a recent multilingual BERT variant with improved cross-lingual alignment and vocabulary coverage, supporting over 100 languages with approximately 307 million parameters.

**Document Translation.** To improve comparability across languages, all documents in CLIRMatrix and Large-Scale are translated into English using NLLB-200 [19]. For mMARCO, which already contains parallel documents in all languages, the corresponding English passages are used for the translated condition. The translation model is selected based on a pilot evaluation of four MT systems on a sample of 100 documents per language. Two criteria are used: COMET [63], which estimates translation adequacy, and the perplexity of LLaMA-3.1-8B [64], which reflects translation fluency. The systems considered are DeepL [65], a strong commercial MT service; NLLB-200, Meta's open multilingual model covering 200 languages; googletrans [66], a lightweight Python interface to Google Translate; and LibreTranslate [67], an open-source MT service that can be run locally. NLLB-200 is chosen as it consistently ranks near the top (typically second best across languages) on both adequacy and fluency. DeepL achieves the best scores overall, but rate limits and cost make it unsuitable for large-scale experiments, so it is treated as an approximate upper bound rather than a deployable option.

**First-stage Ranking.** Ranking is performed using a bi-encoder architecture. For each encoder, document embeddings are precomputed offline. At inference time, query embeddings are obtained and documents are ranked according to the cosine similarity between query and document embeddings, retaining the indices of the top 100 documents for each query. As a lexical baseline, BM25 [68] is included, with the document index constructed in advance and queried at retrieval time. This stage requires no additional training and serves as a control condition for assessing model- and data-level interventions. All encoders use mean pooling over the final hidden layer, and the resulting embeddings are L2-normalised. No adapters or projection layers are added in the base experiments, and retrieval relies on cosine similarity throughout.

**Contrastive Learning Settings.** Fine-tuning is applied only to the weakest (XLM-R) and strongest (Nomic) pretrained encoders in order to assess whether the effectiveness of contrastive interventions depends on baseline multilingual alignment. Contrastive learning is carried out using text pairs at three levels of granularity. At the word level, training uses multilingual word sense disambiguation datasets (XL-WiC [45], Am$^2$iCO [47], SenWiCh [69] restricted to languages that overlap with those in the retrieval datasets. At the phrase level, up to 1,000 Tatoeba [70] parallel sentences are used for each language pair present in the retrieval benchmarks. At the query–document level, fine-tuning is performed on relevance-annotated pairs from CLIRMatrix, mMARCO and the Large-Scale CLIR dataset.

**Cross-encoder Re-ranking Settings.** XLM-R and Nomic are further fine-tuned as cross-encoders that jointly encode each query–document pair and output a relevance score via a classification head. Re-ranking operates over the top-100 BM25 candidates; if the relevant document is absent from this candidate set, it replaces the item at rank 100. The classifier is trained with binary relevance labels, and two negative sampling strategies are compared: *easy negatives*, obtained by randomly sampling non-relevant documents, and *hard negatives*, obtained from non-relevant documents retrieved by the first-stage ranker and therefore harder to distinguish from true positives (available for CLIRMatrix and mMARCO).

Further details of the experimental settings and implementation are provided in Appendix A- F.

### 3.3 Evaluation

For each query, a ranked list of candidate documents is obtained using cosine similarity between L2-normalised bi-encoder embeddings, which serves as the query–document similarity score for all neural retrieval models. For BM25, document ranking is performed directly over the inverted index using the query tokens. In both cases, the top-100 documents are retained as the retrieval output for evaluation.

Retrieval effectiveness is measured over 1,000 (or 500) queries per language pair using standard IR metrics. The primary metric is Recall@100, defined as the proportion of queries for which the relevant document appears in the top-100 results. For re-ranking experiments, Recall@10 and nDCG@100 (normalized Discounted Cumulative Gain) [71] are additionally reported, where nDCG captures the quality of the ranking by assigning higher weight to relevant documents appearing near the top of the list and normalising scores across queries. Metrics are averaged over all language pairs within each dataset, and four representative pairs (en–es, en–zh, ja–en, ar–ru) are highlighted to illustrate

behaviour across high- and low-resource, same- and cross-script settings. Finally, post-hoc analyses examine (i) the relationship between ANN latency and retrieval accuracy; (ii) document-language retrieval bias; and (iii) correlations between retrieval quality and linguistic similarity.

A summary of the experiments, covering the models, training regimes, levels of comparison and types of retrieval can be seen in Table 2.

| Encoder(s) | Training | Level | Retrieval Type |
|---|---|---|---|
| BM25 | – | Original/Translated | Lexical |
| 5 pretrained multilingual encoders | – | Original/Translated | Cosine |
| Nomic, XLM-R | Contrastive | Word/Phrase/QD | Cosine |
| Nomic, XLM-R | Easy/Hard Negatives | QD | Cross-Encoder |
| 5 pretrained multilingual encoders | - | exact/ANN | Cosine |

Table 2: Summary of experimental settings across all intervention types. (QD = Query-Document)

ANN search is implemented using Hierarchical Navigable Small Worlds (HNSW) [34] indices in order to measure potential speed-ups over exact cosine-similarity search. All ANN experiments compare latency and Recall@100 against exhaustive nearest-neighbour retrieval on the same embeddings. Index construction details and parameter settings are reported in Appendix E.

# 4   Results

This section first selects the MT model used for document translation, then evaluates lexical BM25 retrieval with and without translation. It then examines the retrieval performance of pretrained multilingual encoders, before analysing the effects of contrastive fine-tuning, cross-encoder re-ranking and ANN search on effectiveness and efficiency.

## 4.1   Machine Translation Model

Candidate MT systems are evaluated to select the model used for document translation in the multilingual experiments. Translation quality is assessed on a subset of 100 documents per language from CLIRMatrix and Large-Scale dataset, after translating documents into English. Two criteria are used: COMET [63], which estimates translation adequacy, and the perplexity of LLaMA-3.1-8B [64], which reflects the fluency of the translated text. For the mMARCO dataset, the already available English documents are used in the translated condition, as the collection contains parallel queries and documents across all languages.

Figure 2 shows that DeepL consistently achieves the highest COMET scores and lowest perplexity. Because its free tier is limited to roughly 500k characters per month, the available quota is used twice to complete the evaluation. Among free and open-source systems, NLLB-200 emerges as the strongest option, outperforming googletrans and LibreTranslate on both adequacy and fluency. Googletrans is also less reliable in this setting, failing to translate two Arabic and one French document that contain extensive English code-switching, which triggers translation failures. On this basis, NLLB-200 is selected to translate the document corpora for subsequent experiments.

Across MT systems, Spanish and French consistently obtain low perplexity and high COMET scores, whereas Arabic, Japanese and Chinese show substantially lower COMET scores and highly variable perplexity. These differences reflect typological and script-related biases in both translation models and evaluation metrics. Script divergence in particular emerges as a persistent source of difficulty, with non-Latin scripts yielding less stable and generally poorer translations. This pattern aligns with broader observations that cross-script transfer remains a major bottleneck for multilingual NLP and contributes to uneven downstream cross-lingual retrieval performance.

## 4.2   Cross-Lingual Information Retrieval

After selecting the translation model for corpus normalisation, retrieval performance is examined under both translated and non-translated conditions. As a first step, BM25 is evaluated to characterise the limitations of purely lexical matching before moving to semantic encoders. Comparisons of BM25 with and without document translation confirm that lexical matching is largely ineffective for cross-lingual retrieval, particularly for non-English and cross-script language pairs.

In operational settings, where queries may arrive in unknown languages or MT coverage is limited, reliance on BM25 without translation reduces effectiveness to near zero for most non-English pairs. Figure 3 illustrates this limitation.
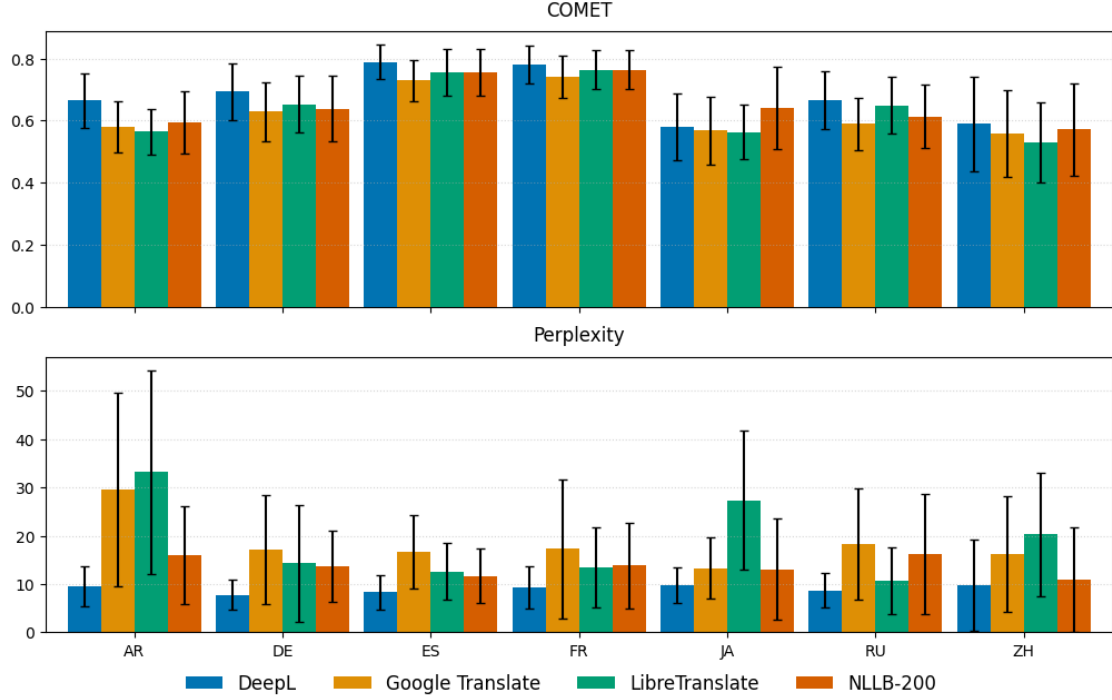
Figure 2: Translation quality of different models evaluated with COMET and Perplexity, by translating non-English documents from the CLIRMatrix dataset [29], sampling the same 100 documents from each language.
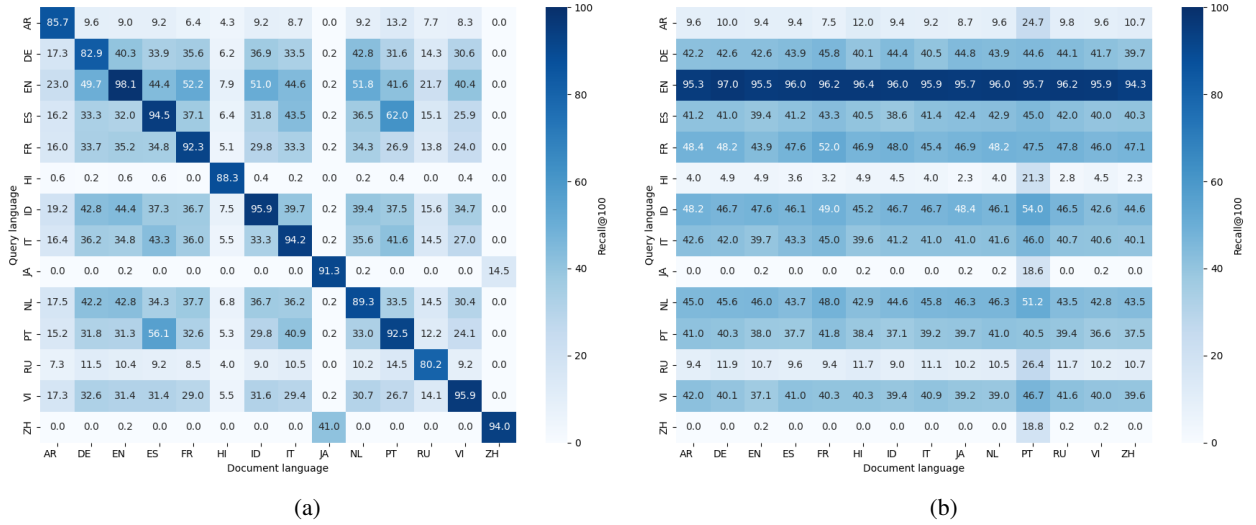


Figure 3: Comparison of lexical BM25 retrieval (Recall@100) on query–document language pairs for cross-lingual retrieval (a) without and (b) with document translation on the mMARCO dataset [28].

Without translation (Figure 3a), BM25 retrieves almost no relevant documents for non-English pairs, as expected for methods that depend on shared lexical forms. Translating all documents into English improves performance only for English queries (Figure 3b) by restoring lexical overlap between English queries and English documents, but it does not benefit queries in other languages and remains ineffective in cross-script settings. These findings motivate the move from lexical matching to semantic embedding-based ranking.

**Semantic Embedding-Based Ranking.** Given the weakness of lexical retrieval in multilingual settings, attention turns to the performance of pretrained multilingual encoders. Five encoders are compared across datasets and representative language pairs. This comparison shows that dense models dramatically outperform BM25 across all benchmarks. As displayed in Figure 4, Nomic consistently achieves the strongest performance, while XLM-R is the weakest. These two encoders therefore serve as representative extremes for analysing the effects of contrastive learning at different levels of granularity.

Performance patterns across language pairs reveal persistent difficulties for distant or low-resource combinations. English–Spanish is generally the easiest direction for most models (with the exception of LaBSE), whereas Arabic–Russian is consistently the most challenging, highlighting limitations in bridging typologically distant languages and cross-script retrieval. Evaluating the same models on a fully translated corpus (Appendix A) yields only minor changes: small improvements for some encoders and small declines for others.
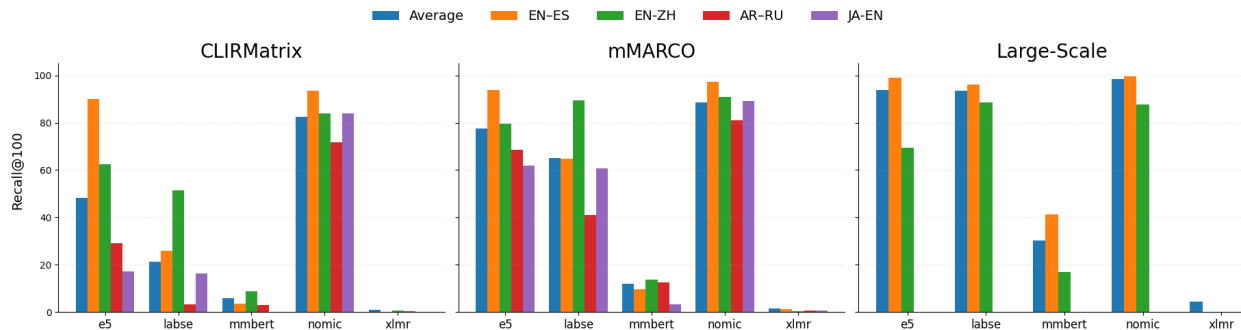


Figure 4: Performance (Recall@100) of five pretrained multilingual encoders using cosine similarity between query and document embeddings without translation. Results are averaged across all language pairs and shown for selected pairs.

| | Original Documents | Translated Documents |
|---|---|---|
| **CLIRMatrix** | | |
| multilingual-E5 | 48.3 | 19.5 |
| LaBSE | 21.4 | 19.0 |
| mmBERT | 6.0 | 4.2 |
| Nomic | 82.6 | 84.7 |
| XLM-R | 1.0 | 0.9 |
| **mMARCO** | | |
| multilingual-E5 | 77.5 | 94.7 |
| LaBSE | 65.2 | 69.6 |
| mmBERT | 12.1 | 10.2 |
| Nomic | 88.6 | 92.2 |
| XLM-R | 1.6 | 1.9 |
| **Large-Scale** | | |
| multilingual-E5 | 93.8 | 93.1 |
| LaBSE | 93.5 | 86.9 |
| mmBERT | 30.3 | 34.9 |
| Nomic | 98.3 | 93.9 |
| XLM-R | 4.3 | 6.9 |

Table 3: Impact of document translation on embedding-based retrieval. Average Recall@100 scores across all language pairs for each dataset when using document embeddings in their original language versus translated into English.

Document translation has a limited effect on embedding-based retrieval, with most models showing small or inconsistent changes in Recall@100 across datasets. In contrast to BM25, where translation is essential to restore lexical overlap, dense encoders already capture cross-lingual semantic similarity, making their performance largely invariant to whether documents are embedded in their original language or in English. Overall, these results indicate that multilingual encoders rely primarily on semantic rather than surface-level cues, and that document translation offers little additional benefit in dense retrieval pipelines. This motivates a closer examination of whether targeted contrastive alignment can further improve retrieval effectiveness.

8

**Contrastive Model Fine-Tuning.** Table 4 shows that contrastive fine-tuning affects XLM-R and Nomic in markedly different ways, depending on the level of supervision. For XLM-R, phrase-level alignment brings substantial gains, suggesting that its pretrained representations do not fully capture cross-lingual semantic relations and remain sensitive to language-specific noise. Query–document supervision yields by far the strongest improvements, increasing Recall@100 from 1.0% 4.3% and 1.6% in the baseline to 21.2%, 96.0% and 77.1% on CLIRMatrix, Large-Scale and mMARCO respectively, which indicates that retrieval-specific objectives are essential for general-purpose multilingual encoders. In contrast, word-level alignment has no meaningful effect, likely because query terms are rarely ambiguous in context and isolated word embeddings provide insufficient information for disambiguation.

For Nomic, which is already trained with cross-lingual alignment and retrieval objectives, additional contrastive fine-tuning has limited impact. Performance remains above 80% across most strategies, except with word-level alignment, and only query–document supervision on CLIRMatrix yields a modest improvement of around ~+4%. This suggests that Nomic is already well optimised for cross-lingual retrieval and that further supervision produces diminishing returns rather than systematic gains.

| Model | Dataset | None | Word | Phrase | Query–Doc |
|---|---|---|---|---|---|
| XLM-R | CLIRMatrix | 1.0 | 1.1 | 13.7 | 21.2 |
| | Large-Scale | 4.3 | 0.4 | 64.8 | 96.0 |
| | mMARCO | 1.6 | 1.4 | 68.0 | 77.1 |
| Nomic | CLIRMatrix | 82.6 | 31.1 | 79.9 | 88.3 |
| | Large-Scale | 98.3 | 98.3 | 98.3 | 97.2 |
| | mMARCO | 88.6 | 88.6 | 88.6 | 84.4 |

Table 4: Effect of contrastive fine-tuning (word, phrase, query–document level) on cosine-similarity retrieval (Recall@100 %) using XLM-R and Nomic, compared against non–fine-tuned baselines (None) across datasets, averaged across all language pairs.
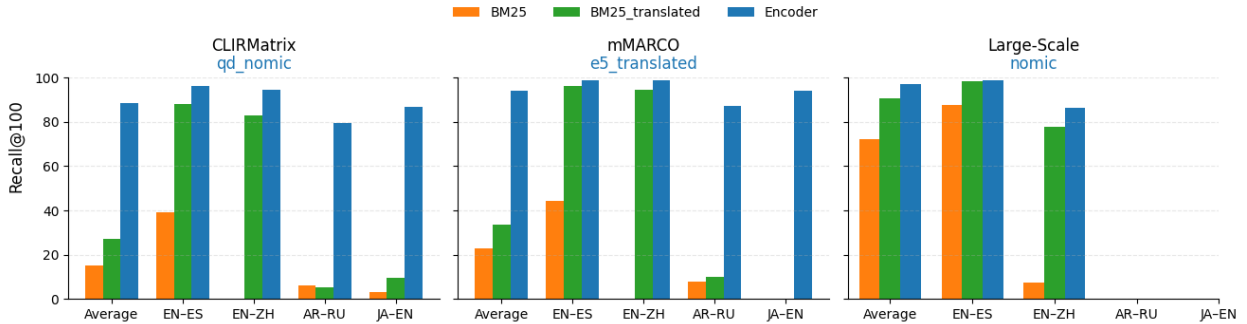


Figure 5: Performance of the best-performing embedding-based model on each dataset and BM25, with and without document translation, on average, and across selected query–document language pairs.

To contextualise the benefits of semantic encoders, the best-performing dense model on each dataset is contrasted with BM25. Figure 5 compares the strongest embedding-based configuration with BM25 applied to both original and English-translated corpora. Improvements over lexical retrieval are substantial across all datasets and are especially pronounced for low-resource and script-divergent language pairs, where BM25 fails almost entirely because of the lack of shared vocabulary. Dense retrievers also outperform BM25 when both languages are relatively well resourced, such as English–Spanish and English–Chinese, demonstrating that semantic representations provide benefits beyond vocabulary mismatch. These gains persist regardless of whether BM25 is applied to original or translated documents.

Table 5 summarises which model achieves the highest Recall@100 for each dataset, both on average and for selected language pairs. On CLIRMatrix, query–document contrastive learning substantially improves the already strong Nomic encoder, yielding the best results across all evaluated language pairs. The notably high performance observed for the Japanese–English pair in the translated condition does not follow from improved translation quality, since these documents are originally in English and remain untranslated. Instead, the improvement stems from reducing the influence of document embeddings in other languages that, in the non-translated setting, occasionally rank above the correct English documents. For mMARCO, both translation and strong cross-lingual alignment contribute to improved retrieval quality. This pattern is expected: the "multilingual" mMARCO variants are created by translating the original

| Dataset | Average | EN-ES | EN-ZH | AR-RU | JA-EN |
|---|---|---|---|---|---|
| CLIRMatrix | QD_N 88.3 | QD_N 96.2 | QD_N 94.6 | QD_N 79.6 | QD_N_T 87.1 |
| mMARCO | E5_T 94.7 | N_T 99.6 | N_T 99.6 | E5_T 88.5 | E5_T 95.3 |
| Large-Scale | N 98.3 | N 99.7 | N_T 91.3 | — | — |

Table 5: Best performing model on each dataset, and in the selected language pairs, and the Recall@100 score. QD = query-document contrastive learning, T = on translated documents, N = Nomic, E5 = multilingual-E5..

| Dataset | Model | Negatives | Recall@10 | | | | | nDCG@100 | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | Avg | EN-ES | EN-ZH | AR-RU | JA-EN | Avg | EN-ES | EN-ZH | AR-RU | JA-EN |
| CLIRMatrix | BM25 | – | 10.4 | 31.4 | 0.0 | 1.7 | 03.1 | 6.2 | 16.9 | 0.0 | 01.6 | 01.5 |
| | Nomic | easy | 11.8 | 10.0 | 60.0 | 0.0 | 10.0 | 21.5 | 23.0 | 41.1 | 17.0 | 23.2 |
| | Nomic | hard | 33.7 | 30.0 | 30.0 | 30.0 | 50.0 | 29.4 | 28.6 | 27.8 | 26.2 | 31.8 |
| | XLM-R | easy | 35.0 | 50.0 | 30.0 | 30.0 | 50.0 | 33.3 | 33.9 | 27.2 | 26.8 | 49.0 |
| | XLM-R | hard | 30.2 | 10.0 | 20.0 | 60.0 | 20.0 | 26.3 | 23.0 | 26.2 | 33.3 | 30.3 |
| mMARCO | BM25 | – | 17.3 | 29.8 | 0.0 | 05.3 | 0.2 | 14.9 | 26.5 | 0.0 | 04.3 | 0.2 |
| | Nomic | easy | 13.6 | 20.0 | 0.0 | 06.0 | 0.0 | 25.1 | 30.1 | 18.0 | 19.8 | 15.4 |
| | Nomic | hard | 08.7 | 18.0 | 04.0 | 03.0 | 0.0 | 23.3 | 30.8 | 19.5 | 20.4 | 16.9 |
| | XLM-R | easy | 16.1 | 32.0 | 0.0 | 04.0 | 0.0 | 21.9 | 29.0 | 15.0 | 17.8 | 15.0 |
| | XLM-R | hard | 83.4 | 90.0 | 91.0 | 82.0 | 78.0 | 62.5 | 69.9 | 73.3 | 58.8 | 49.5 |
| Large-Scale | BM25 | – | 63.5 | 76.3 | 07.1 | – | – | 57.9 | 69.7 | 06.5 | – | – |
| | Nomic | easy | 62.3 | 77.0 | 09.0 | – | – | 60.7 | 72.0 | 22.7 | – | – |
| | XLM-R | easy | 62.0 | 77.0 | 09.0 | – | – | 57.5 | 67.6 | 22.0 | – | – |

Table 6: Cross-encoder re-ranking with easy and hard negatives over the BM25 baseline. Entries report Recall@10 and nDCG@100 over the top 100 documents per query; improvements over the baseline are shown in **green** and degradations in **red** .

English MS MARCO passages, so retrieval over the original English corpus recovers the highest-quality texts, while Nomic and multilingual-E5 already provide strong cross-lingual alignment (as reflected in multilingual-E5's state-of-the-art performance on MMTEB [42]), leaving little room for further gains through fine-tuning. Finally, on the Large-Scale dataset, Nomic without any additional interventions achieves the best overall performance. Given that Nomic is explicitly optimised for cross-lingual retrieval and already attains very high effectiveness, further improvements are naturally hard to obtain.

## 4.3 Re-ranking Model

Cross-encoder re-ranking enables finer-grained modelling of interactions between query and document content, at the cost of higher computational load per inference. Across datasets, cross-encoders consistently improve the rank of the relevant document relative to BM25, although the magnitude of the improvement varies across language pairs. The largest gains occur where BM25 performs worst, namely in cross-script or lexically distant pairs such as English-Chinese (en–zh), Arabic-Russian (ar–ru) and Japanese-English (ja–en), where BM25 Recall@10 is close to zero but cross-encoders reach substantially higher values (e.g. CLIRMatrix en–zh improves from 0.0% to 30.0%; mMARCO en–zh from 0.0% to 91.0% with XLM-R trained on hard negatives). For language pairs with higher lexical overlap such as en–es, improvements are more modest, since BM25 already provides a reasonable initial ranking and the top-100 candidates often contain many plausible but non-relevant documents, which reduces the potential benefit of re-ranking.

Negative sampling emerges as a key factor in determining re-ranking effectiveness. Models trained with hard negatives generally outperform those trained with easy negatives, with XLM-R on mMARCO providing the clearest example: when trained on hard negatives, it reaches an average Recall@10 of 83.4%, nearly a five-fold improvement over BM25 at 17.3%, whereas the same architecture trained with easy negatives does not surpass the baseline. Although the available evidence is limited to two datasets and two encoders, this pattern is consistent with prior work that emphasises the importance of hard negatives in retrieval training more broadly, including in cross-lingual settings [72, 73, 74]. In contrast, little effect is observed on the Large-Scale dataset, likely because only easy negatives are available, preventing the cross-encoders from learning to sharpen semantic discrimination.

| | Average Time | | | Average Performance (Recall@100) | | |
|---|---|---|---|---|---|---|
| **Method** | CLIRMatrix | mMARCO | Large-Scale | CLIRMatrix | mMARCO | Large-Scale |
| DNN | **0.388** | **0.451** | **0.362** | **31.9** | **49.0** | **64.0** |
| ANN | 0.619 | 0.552 | 0.652 | 31.6 | 47.1 | 60.1 |

Table 7: Comparison of average inference time and Recall@100 performance for DNN and ANN retrieval across datasets. Averaged over the five pretrained encoders, over all language pairs in each dataset.

**Retrieval Efficiency and Approximate Indexing.** ANN indexing is typically used to improve retrieval efficiency for embedding-based systems over large corpora. However, in our setting, ANN does *not* yield computational benefits. As shown in Table 7, ANN search increases average inference time by 40 − 80% across datasets compared to direct retrieval, while only reducing Recall@100 by 0.3 − 3.9% on average. The document collections are relatively small, ranging from approximately 7k to 26k documents, and exact nearest neighbour retrieval, by comparing the query embedding to all the document embeddings, can be implemented as a single batched cosine-similarity computation between each query embedding and all document embeddings, which is highly optimised and vectorised. By contrast, the approximate setup, based on HNSW indices, requires index access for every query–document language pair and incurs additional overhead from graph traversal and index maintenance. For corpora of this size, these overheads outweigh any pruning benefit from approximate search, resulting in slower retrieval despite only minor reductions in recall.

These observations do not rule out the utility of approximate methods in larger-scale cross-lingual applications, particularly when embedding spaces span many languages and lexical shortcuts are limited. They do indicate, however, that approximate indexing is most effective beyond a certain corpus size. For the benchmarks considered here, approximate search yields negligible improvements in efficiency and slight losses in effectiveness, suggesting that exact dense retrieval remains the more appropriate choice. At larger scales (>1M vectors), ANN has been shown to outperform DNN substantially [33].

## 5 Analysis

The results indicate that the effectiveness of cross-lingual retrieval is shaped not only by the choice of retrieval architecture and fine-tuning strategy, but also by translation decisions and underlying linguistic factors. The following analysis unpacks these influences by examining how retrieval models encode cross-lingual meaning, the role of document translation and supervision, the conditions under which re-ranking is beneficial, and the extent to which linguistic distance affects performance. This provides a qualitative interpretation of the strengths and limitations of current methods beyond aggregate scores.

**Influence of Retrieval Models on Performance.** Lexical retrieval is fundamentally constrained by surface-form matching. Translating documents can partially mitigate this limitation, but only when the query language matches the translation language. Dense encoders, in contrast, consistently outperform BM25 on both original and translated corpora because they represent cross-lingual semantics rather than relying solely on shared vocabulary. Contrastive learning on parallel corpora or query–document pairs substantially improves retrieval when applied to base encoders with weak multilingual alignment. For XLM-R, recall increases from baseline values 1.0%, 4.3% and 1.6% on CLIRMatrix, Large-Scale and mMARCO datasets to 21.2%, 96.0% and 77.1% respectively, bringing performance close to that of Nomic, which is already optimised for cross-lingual retrieval. In contrast, encoders that are already strongly aligned, such as Nomic, exhibit only minor or negligible gains. Taken together, these patterns indicate that semantic alignment between languages and between queries and documents in the embedding space, rather than translation or lexical matching, is the primary driver of cross-lingual retrieval performance.

**Impact of Document Translation.** Because queries cannot reliably be translated at inference time in realistic systems, owing to uncertainty in language identification, variable translation quality and latency constraints, translation is applied only to documents. For BM25, document translation is essential to enable cross-lingual retrieval but is effective only when the query is in the same language as the translated collection. Translating documents into English substantially improves retrieval for English queries by restoring lexical overlap, yet yields little or no benefit for other query languages and fails entirely for cross-script scenarios.

Embedding-based retrieval behaves differently. Performance is largely unaffected by document translation and can even degrade slightly because of translation-induced noise (Table 3). Dense encoders already operate on semantic representations, which reduces dependence on surface form and minimises differences between original and translated conditions. Consequently, while translation supports shallow lexical methods, embedding-based systems derive meaning directly from multilingual representations. Combined with the observations above, this suggests that effective cross-

lingual retrieval depends on direct semantic encoding, and that introducing translation into embedding-based pipelines primarily adds computational overhead without improving retrieval quality.

**Role of Re-ranking.** Cross-encoder re-ranking can substantially improve the position of the relevant document when the initial candidate set is weak, especially for typologically distant or cross-script language pairs. In these cases, cross-encoders exploit detailed interactions between the query and candidate documents to recover relevant items that first-stage retrieval fails to rank highly. However, the magnitude and consistency of these gains vary across encoders and datasets, and depend strongly on the availability and quality of hard negatives. Models trained with hard negatives generally provide substantial improvements over BM25, whereas training with only easy negatives often fails to surpass the baseline. Re-ranking therefore remains a valuable but non-uniform refinement step. It is most beneficial when strong supervision and additional inference-time computation are available, but its training is complicated by limited lexical overlap and the difficulty of constructing informative cross-lingual negative examples.

**Analysis of Retrieval Efficiency.** Efficiency experiments show that approximate nearest neighbour search does not provide a speed advantage at the small to medium corpus scales considered here, with document collections between 7k and 26k items. Index traversal and graph-based overheads in the approximate method outweigh any benefit from pruning, making exact dense nearest neighbour search both faster and more reliable. Approximate indexing thus appears most useful as a scalability mechanism at substantially larger index sizes, rather than as a universal optimisation for cross-lingual retrieval.

**Linguistic Factors.** To understand why retrieval difficulty varies across language pairs, the quantitative evaluation is complemented by an analysis of how linguistic factors influence accuracy. Two questions are examined: whether linguistic similarity between query and document languages correlates with retrieval quality (Table 8) and whether models exhibit biases towards particular document languages (Figure 6).
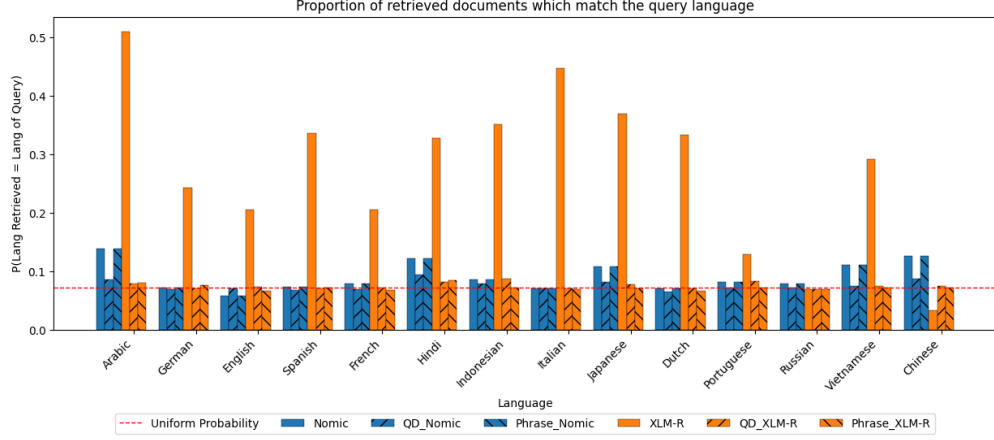
| Model | Dataset | geographic | syntax | phonology | inventory | genealogical |
|---|---|---|---|---|---|---|
| BM25 | clirmatrix | 80.5 | 74.3 | 36.4 | 6.3 | 63.7 |
| | mmarco | 48.9 | 70.4 | 41.5 | 38.5 | 54.8 |
| | large-scale | 62.9 | 57.2 | 46.7 | 33.6 | 51.3 |
| XLM-R | clirmatrix | -09.4 | -28.5 | -24.4 | 12.8 | -21.5 |
| | mmarco | 20.2 | 26.2 | 16.0 | 22.3 | 28.9 |
| | large-scale | 31.9 | 24.5 | 54.9 | 23.9 | 35.7 |
| Nomic | clirmatrix | 53.0 | 65.7 | 30.1 | -15.3 | 74.4 |
| | mmarco | 45.9 | 66.4 | 45.2 | 33.5 | 71.0 |
| | large-scale | 70.0 | 57.4 | 42.9 | 28.8 | 47.3 |
| **Average** | | 44.9 | 46.0 | 32.2 | 20.5 | 45.1 |

Table 8: Spearman correlations between linguistic similarity scores and retrieval performance across language pairs, models, and datasets.
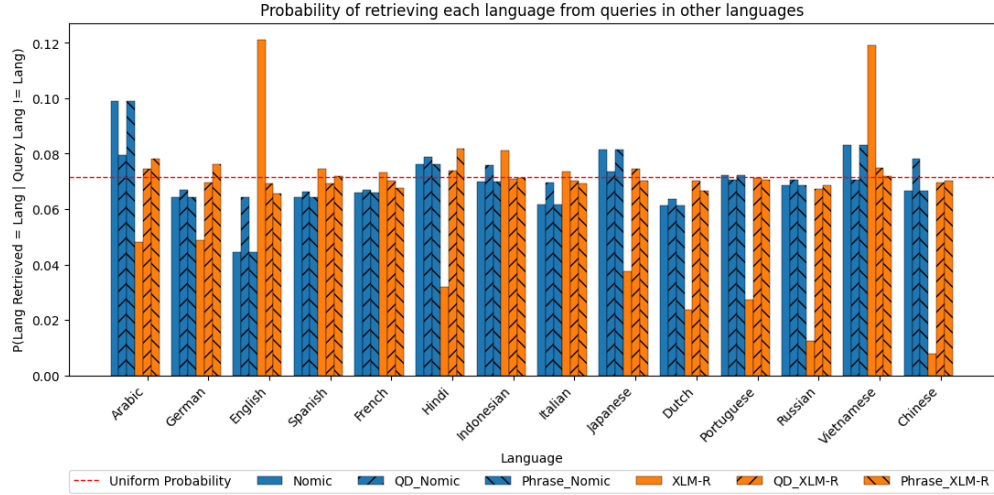
Table 8 shows that BM25 retrieval correlates strongly with features associated with typological proximity, particularly geographic, syntactic and genealogical similarity, reflecting its reliance on surface-form overlap. Dense encoders exhibit weaker correlations, indicating an ability to capture semantic alignment beyond lexical similarity. Nonetheless, some structural influence remains: Nomic, the strongest encoder, performs better between linguistically related languages, suggesting that cross-lingual alignment is not fully language agnostic. XLM-R shows very low correlation values, which are likely a consequence of its poor baseline performance rather than evidence of an absence of linguistic effects.

Figures 6a and 6b reveal fundamental differences in how models encode linguistic information during retrieval. XLM-R disproportionately retrieves documents written in the same language as the query, indicating limited cross-lingual generalisation and semantic abstraction. When the query and document languages diverge, XLM-R exhibits a marked preference for English and surprisingly Vietnamese, while under-retrieving Chinese, Russian, and other languages with distinct scripts. In contrast, Nomic demonstrates a much more uniform retrieval pattern, suggesting stronger cross-lingual alignment.

Critically, we find that contrastive learning, both at the phrase level and query-document level, substantially reduces retrieval bias, bringing XLM-R much closer to the uniform language distribution in the datasets. This suggests that contrastive alignment acts as a corrective mechanism, encouraging models to encode semantic equivalence rather than relying on language-specific features.

(a) Probability that the model retrieves a document written in the same language as the query. Nomic remains close to uniform, whereas XLM-R without fine-tuning shows strong query-language preference. Bars marked with (\\) and (//) denote phrase-level and query–document contrastive fine-tuning, respectively.



(b) Distribution of retrieved document languages when the correct document is not in the query language. Without fine-tuning, XLM-R displays language and script biases, whereas Nomic retrieves documents more uniformly. Contrastive fine-tuning settings follow the same notation as above.

Figure 6: Linguistic bias in document retrieval across Nomic and XLM-R under baseline and contrastive fine-tuning.

Further comparison with additional models (e.g. BM25, LaBSE, multilingual-E5, mmBERT) is provided in Appendix G. These results show that language bias is widespread across multilingual encoders and may be driven by their training objectives. Notably, multilingual-E5, despite achieving comparatively strong cross-lingual retrieval performance in Figure 4, exhibits very strong query-language retrieval bias, almost on par with BM25. When examining retrieval of documents in non-query languages, models show more similar behaviour; however, most still under-retrieve documents across several languages. This suggests that embeddings for those languages may be encoded in more separated, potentially language-specific clusters.

Overall, the linguistic analysis indicates that language similarity remains a core challenge for purely lexical systems and continues to significantly affect dense retrieval. While strong multilingual encoders substantially mitigate language and script biases, they do not eliminate them entirely. Nonetheless, targeted contrastive learning can further reduce such biases. The disparity between XLM-R and Nomic, as well as improvements observed through fine-tuning, demonstrates that cross-lingual generalisation relies heavily on encoder quality and training objective rather than being an inherent property of dense retrieval architectures.

In summary, the analyses indicate that performance in cross-lingual information retrieval is primarily driven by semantic alignment rather than surface-form similarity or translation effects. Linguistic structure still influences retrieval for weaker models, but its impact diminishes as encoders become more strongly grounded in multilingual semantics. These insights inform the concluding recommendations on the design of robust cross-lingual retrieval systems.

## 6 Conclusion

This work evaluated modelling strategies for CLIR, with a focus on semantic alignment, translation, supervision and the trade-off between efficiency and effectiveness. The study compared five pretrained multilingual encoders, contrastive fine-tuning at multiple levels of granularity, document translation, cross-encoder re-ranking and approximate nearest neighbour search, and examined how linguistic factors influence retrieval behaviour.

The results indicate that effective cross-lingual retrieval depends primarily on models that construct stable, language-agnostic semantic representations. Dense retrieval consistently outperforms lexical and translation-based baselines, confirming that semantic modelling, rather than lexical overlap, is the main driver of performance. Contrastive learning yields substantial gains only for encoders with weak initial alignment, whereas models already optimised for cross-lingual retrieval benefit little from additional supervision. Cross-encoder re-ranking offers inconsistent and often modest improvements, mainly in challenging cross-script or typologically distant language pairs, and its effectiveness is sensitive to the choice of negative documents used during training. Approximate nearest neighbour search does not accelerate inference at the scale of the evaluated benchmarks. Linguistic analysis further shows that structural similarity between languages affects retrieval for weaker models, but this influence diminishes as semantic alignment in the encoder improves. The findings suggest that progress in cross-lingual information retrieval depends on three main factors: (i) strong multilingual pretraining that yields robust language-agnostic embeddings; (ii) targeted contrastive supervision to correct specific alignment deficiencies; and (iii) selective use of translation or re-ranking only when justified by model limitations or typological characteristics of the language pair. Future work could extend these analyses to broader language coverage and larger collections, and explore alignment-aware retrieval strategies, particularly for low-resource and typologically distant languages.

**Limitations.** Although this study offers a structured analysis of cross-lingual retrieval across multiple architectures, supervision regimes and linguistic conditions, several limitations constrain the generality of the conclusions. First, the three datasets considered are modest in size, domain coverage and language diversity. This restricts the ability to draw firm conclusions about performance in truly large-scale or highly heterogeneous retrieval scenarios. In particular, the limited number of typologically distant and low-resource languages makes it difficult to fully characterise model behaviour in under-represented linguistic settings. Second, the cross-encoder re-ranking analysis is constrained by the availability of hard negatives in only two datasets, which prevents a systematic assessment of how negative sampling interacts with multilingual semantic alignment. The cross-lingual composition of the cross-encoder training data may also influence how effectively the models learn the task, but this factor is not exhaustively explored here. Third, the efficiency findings are tied to relatively small corpora, so the conclusion that approximate nearest neighbour search does not provide speed benefits is unlikely to hold at larger scales, where index traversal costs are amortised and approximate methods may become advantageous. Finally, the linguistic analyses rely on aggregate typological resources and broad similarity metrics, which may not capture all structural properties that are relevant for cross-lingual retrieval. These limitations highlight the need for larger and more typologically diverse datasets, broader architectural coverage and large-scale evaluations in order to fully characterise the behaviour of cross-lingual retrieval systems.

## References

[1] Shota Sasaki, Shuo Sun, Shigehiko Schamoni, Kevin Duh, and Kentaro Inui. Cross-lingual learning-to-rank with shared representations. In Marilyn Walker, Heng Ji, and Amanda Stent, editors, *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 458–463, New Orleans, Louisiana, June 2018. Association for Computational Linguistics.

[2] Roksana Goworek, Olivia Macmillan-Scott, and Eda B. Özyiğit. Bridging language gaps: Advances in cross-lingual information retrieval with multilingual llms. *arXiv preprint arXiv:2510.00908*, 2025.

[3] Douglas W. Oard and Anne R. Diekema. Cross-language information retrieval. In *Annual Review of Information Science and Technology (ARIST)*, volume 33, pages 223–256, 1998.

[4] Jian-Yun Nie. *Cross-Language Information Retrieval*, volume 8 of *Synthesis Lectures on Human Language Technologies*. Morgan & Claypool, 2010.

[5] Robert Litschko, Ivan Vulić, Simone Paolo Ponzetto, and Goran Glavaš. On cross-lingual retrieval with multilingual text encoders. *Information Retrieval Journal*, 25:149–183, 2022.

[6] Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. Unsupervised cross-lingual representation learning at scale. In Dan Jurafsky, Joyce Chai, Natalie Schluter, and Joel Tetreault, editors, *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online, July 2020. Association for Computational Linguistics.

[7] Pratik Joshi, Sebastin Santy, Amar Budhiraja, Kalika Bali, and Monojit Choudhury. The state and fate of linguistic diversity and inclusion in the nlp world. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6282–6293, 2020.

[8] Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. mt5: A massively multilingual pre-trained text-to-text transformer. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1002–1015, 2021.

[9] Victor Lavrenko and W. Bruce Croft. Cross-language information retrieval via neural network translation models. In *Proceedings of the 25th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 175–182. ACM, 2002.

[10] Atsushi Fujii and Tetsuya Ishikawa. Applying machine translation to two-stage cross-language information retrieval. In John S. White, editor, *Proceedings of the Fourth Conference of the Association for Machine Translation in the Americas: Technical Papers*, pages 13–24, Cuernavaca, Mexico, October 10-14 2000. Springer.

[11] Dong Zhou, Mark Truran, Tim Brailsford, Vincent Wade, and Helen Ashman. Translation techniques in cross-language information retrieval. *ACM Comput. Surv.*, 45(1), December 2012.

[12] Min Xiao and Yuhong Guo. Semi-supervised representation learning for cross-lingual text classification. In David Yarowsky, Timothy Baldwin, Anna Korhonen, Karen Livescu, and Steven Bethard, editors, *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1465–1475, Seattle, Washington, USA, October 2013. Association for Computational Linguistics.

[13] Telmo Pires, Éric Schlinger, and David Garrette. How multilingual is multilingual bert? In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, page 4996–5001, 2019.

[14] Jiayang Song, Yuheng Huang, Zhehua Zhou, and Lei Ma. Multilingual blending: Large language model safety alignment evaluation with language mixture. In Luis Chiruzzo, Alan Ritter, and Lu Wang, editors, *Findings of the Association for Computational Linguistics: NAACL 2025*, pages 3433–3449, Albuquerque, New Mexico, April 2025. Association for Computational Linguistics.

[15] Jimmy Lin, David Alfonso-Hermelo, Vitor Jeronymo, Ehsan Kamalloo, Carlos Lassance, Rodrigo Nogueira, Odunayo Ogundepo, Mehdi Rezagholizadeh, Nandan Thakur, Jheng-Hong Yang, and Xinyu Zhang. Simple yet effective neural ranking and reranking baselines for cross-lingual information retrieval, 2023.

[16] Shuo Sun, Suzanna Sia, and Kevin Duh. Clireval: Evaluating machine translation as a cross-lingual information retrieval task. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2020.

[17] Robert Litschko, Ivan Vulić, and Goran Glavaš. Parameter-efficient neural reranking for cross-lingual and multilingual retrieval. In Nicoletta Calzolari, Chu-Ren Huang, Hansaem Kim, James Pustejovsky, Leo Wanner, Key-Sun Choi, Pum-Mo Ryu, Hsin-Hsi Chen, Lucia Donatelli, Heng Ji, Sadao Kurohashi, Patrizia Paggio, Nianwen Xue, Seokhwan Kim, Younggyun Hahm, Zhong He, Tony Kyungil Lee, Enrico Santus, Francis Bond, and Seung-Hoon Na, editors, *Proceedings of the 29th International Conference on Computational Linguistics*, pages 1071–1082, Gyeongju, Republic of Korea, October 2022. International Committee on Computational Linguistics.

[18] Piotr Indyk and Rajeev Motwani. Approximate nearest neighbors: towards removing the curse of dimensionality. In *Proceedings of the Thirtieth Annual ACM Symposium on Theory of Computing*, STOC '98, page 604–613, New York, NY, USA, 1998. Association for Computing Machinery.

[19] NLLB Team, Marta R. Costa-jussà, James Cross, Onur Çelebi, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Maillard, Anna Sun, Skyler Wang, Guillaume Wenzek, Al Youngblood, Bapi Akula, Loic Barrault, Gabriel Mejia Gonzalez, Prangthip Hansanti, John Hoffman, Semarley Jarrett, Kaushik Ram Sadagopan, Dirk Rowe, Shannon Spruit, Chau Tran, Pierre Andrews, Necip Fazil Ayan, Shruti Bhosale, Sergey Edunov, Angela Fan, Cynthia Gao, Vedanuj Goswami, Francisco Guzmán, Philipp Koehn, Alexandre Mourachko, Christophe Ropers, Safiyyah Saleem, Holger Schwenk, and Jeff Wang. No language left behind: Scaling human-centered machine translation, 2022.

[20] Stephen E Robertson, Steve Walker, Susan Jones, Micheline M Hancock-Beaulieu, and Mike Gatford. Okapi at trec-3. In *Proceedings of the Third Text REtrieval Conference (TREC-3)*, pages 109–126. NIST, 1995.

[21] J. Scott McCarley. Should we translate the documents or the queries in cross-language information retrieval? In *Proceedings of the 37th Annual Meeting of the Association for Computational Linguistics*, pages 208–214. Association for Computational Linguistics, 1999.

[22] Kazuaki Kishida and Noriko Kando. Two-stage refinement of query translation in a pivot language approach to cross-lingual information retrieval: An experiment at clef 2003. In *Workshop of the Cross-Language Evaluation Forum for European Languages*, pages 253–262. Springer, 2003.

[23] Lisa Ballesteros and W. Bruce Croft. Resolving ambiguity for cross-language retrieval. In *Proceedings of the 21st Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, SIGIR '98, page 64–71, New York, NY, USA, 1998. Association for Computing Machinery.

[24] Javid Dadashkarimi, Azadeh Shakery, and Heshaam Faili. A probabilistic translation method for dictionary-based cross-lingual information retrieval in agglutinative languages. *arXiv preprint arXiv:1411.1006*, 2014.

[25] Liang Yao, Baosong Yang, Haibo Zhang, Weihua Luo, and Boxing Chen. Exploiting neural query translation into cross lingual information retrieval. *arXiv preprint arXiv:2010.13659*, 2020.

[26] Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. Splade: Sparse lexical and expansion model for first stage ranking. In *SIGIR*, 2021. arXiv preprint arXiv:2107.05720.

[27] Thibault Formal, Carlos Lassance, Benjamin Piwowarski, and Stéphane Clinchant. Splade v2: Sparse lexical and expansion model for information retrieval. In *SIGIR*, 2021. arXiv preprint arXiv:2109.10086.

[28] Luiz Bonifacio, Vitor Jeronymo, Hugo Queiroz Abonizio, Israel Campiotti, Marzieh Fadaee, Roberto Lotufo, and Rodrigo Nogueira. mmarco: A multilingual version of the ms marco passage ranking dataset, 2022.

[29] Shuo Sun and Kevin Duh. CLIRMatrix: A massively large collection of bilingual and multilingual datasets for cross-lingual information retrieval. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4160–4170, Online, November 2020. Association for Computational Linguistics.

[30] Shunyu Zhang, Yaobo Liang, Ming Gong, Daxin Jiang, and Nan Duan. Modeling sequential sentence relation to improve cross-lingual dense retrieval. *arXiv preprint arXiv:2302.01626*, 2023.

[31] Omar Khattab and Matei Zaharia. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *SIGIR*, 2020.

[32] Rohan Jha, Bo Wang, Michael Günther, Georgios Mastrapas, Saba Sturua, Isabelle Mohr, Andreas Koukounas, Mohammad Kalim Akram, Nan Wang, and Han Xiao. Jina-colbert-v2: A general-purpose multilingual late interaction retriever. In *Proceedings of the Fourth Workshop on Multilingual Representation Learning (MRL 2024)*, pages 159–166, Miami, Florida, USA, 2024. Association for Computational Linguistics.

[33] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with GPUs. In *IEEE Transactions on Big Data*, volume 7, pages 535–547. IEEE, 2019.

[34] Yury A. Malkov and Dmitry A. Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 42(4):824–836, 2020.

[35] M Wang, X Xu, Q Yue, and Y Wang. A comprehensive survey and experimental comparison of graph-based approximate nearest neighbor search. arxiv 2021. *arXiv preprint arXiv:2101.12631*.

[36] Rodrigo Nogueira and Kyunghyun Cho. Passage re-ranking with BERT. *arXiv preprint arXiv:1901.04085*, 2019.

[37] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense text retrieval. In *International Conference on Learning Representations (ICLR)*. OpenReview.net, 2021.

[38] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In Jill Burstein, Christy Doran, and Thamar Solorio, editors, *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)*, pages 4171–4186, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics.

[39] Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Arivazhagan, and Wei Wang. Language-agnostic BERT sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 878–891, Dublin, Ireland, 2022.

[40] Mikel Artetxe and Holger Schwenk. Massively multilingual sentence embeddings for zero-shot cross-lingual transfer and beyond. In *TACL*, 2019.

[41] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Multilingual e5 text embeddings: A technical report. `https://arxiv.org/abs/2402.05672`, 2024.

[42] Kenneth Enevoldsen, Isaac Chung, Imene Kerboua, Márton Kardos, Ashwin Mathur, David Stap, Jay Gala, Wissam Siblini, Dominik Krzemiński, Genta Indra Winata, Saba Sturua, Saiteja Utpala, Mathieu Ciancone, Marion Schaeffer, Gabriel Sequeira, Diganta Misra, Shreeya Dhakal, Jonathan Rystrøm, Roman Solomatin, Ömer Çağatan, Akash Kundu, Martin Bernstorff, Shitao Xiao, Akshita Sukhlecha, Bhavish Pahwa, Rafał Poświata, Kranthi Kiran GV, Shawon Ashraf, Daniel Auras, Björn Plüster, Jan Philipp Harries, Loïc Magne, Isabelle Mohr, Mariya Hendriksen, Dawei Zhu, Hippolyte Gisserot-Boukhlef, Tom Aarsen, Jan Kostkan, Konrad Wojtasik, Taemin Lee, Marek Šuppa, Crystina Zhang, Roberta Rocca, Mohammed Hamdy, Andrianos Michail, John Yang, Manuel Faysse, Aleksei Vatolin, Nandan Thakur, Manan Dey, Dipam Vasani, Pranjal Chitale, Simone Tedeschi, Nguyen Tai, Artem Snegirev, Michael Günther, Mengzhou Xia, Weijia Shi, Xing Han Lù, Jordan Clive, Gayatri Krishnakumar, Anna Maksimova, Silvan Wehrli, Maria Tikhonova, Henil Panchal, Aleksandr Abramov, Malte Ostendorff, Zheng Liu, Simon Clematide, Lester James Miranda, Alena Fenogenova, Guangyu Song, Ruqiya Bin Safi, Wen-Ding Li, Alessia Borghini, Federico Cassano, Hongjin Su, Jimmy Lin, Howard Yen, Lasse Hansen, Sara Hooker, Chenghao Xiao, Vaibhav Adlakha, Orion Weller, Siva Reddy, and Niklas Muennighoff. Mmteb: Massive multilingual text embedding benchmark. *arXiv preprint arXiv:2502.13595*, 2025.

[43] Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. Dense passage retrieval for open-domain question answering. In Bonnie Webber, Trevor Cohn, Yulan He, and Yang Liu, editors, *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online, November 2020. Association for Computational Linguistics.

[44] Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamalloo, David Alfonso-Hermelo, Xiaoguang Li, Qun Liu, Mehdi Rezagholizadeh, and Jimmy Lin. MIRACL: A multilingual retrieval dataset covering 18 diverse languages. *Transactions of the Association for Computational Linguistics*, 11:1114–1131, 2023.

[45] Alessandro Raganato, Tommaso Pasini, José Camacho-Collados, and Mohammad Taher Pilehvar. XL-WiC: A multilingual benchmark for evaluating semantic contextualization. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7193–7206. Association for Computational Linguistics, 2020. XL-WiC covers 12 languages.

[46] Mohammad Taher Pilehvar and Jose Camacho-Collados. Wic: The word-in-context dataset for evaluating context-sensitive meaning representations. In *Proceedings of NAACL-HLT*, pages 1267–1273, 2019.

[47] Qianchu Liu, Edoardo Maria Ponti, Diana McCarthy, Ivan Vulić, and Anna Korhonen. AM2iCo: Evaluating word meaning in context across low-resource languages with adversarial examples. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7151–7162, Online and Punta Cana, Dominican Republic, November 2021. Association for Computational Linguistics.

[48] Longfei Zuo, Pingjun Hong, Oliver Kraus, Barbara Plank, and Robert Litschko. Evaluating large language models for cross-lingual retrieval. In Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng, editors, *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 11415–11429, Suzhou, China, November 2025. Association for Computational Linguistics.

[49] Yuxin Huang, Simeng Wu, Ran Song, Yan Xiang, Yantuan Xian, Shengxiang Gao, and Zhengtao Yu. Multilingual generative retrieval via cross-lingual semantic compression. In Christos Christodoulopoulos, Tanmoy Chakraborty, Carolyn Rose, and Violet Peng, editors, *Findings of the Association for Computational Linguistics: EMNLP 2025*, pages 10855–10866, Suzhou, China, November 2025. Association for Computational Linguistics.

[50] Moritz Muennighoff. SGPT: Gpt sentence embeddings for semantic search. *arXiv preprint*, arXiv:2202.08904, 2023.

[51] Erxue Min, Hsiu-Yuan Huang, Xihong Yang, Min Yang, Xin Jia, Yunfang Wu, Hengyi Cai, Junfeng Wang, Shuaiqiang Wang, and Dawei Yin. From prompting to alignment: A generative framework for query recommendation. *arXiv preprint arXiv:2504.10208*, 2025.

[52] Zihan Zhang, Meng Fang, and Ling Chen. RetrievalQA: Assessing adaptive retrieval-augmented generation for short-form open-domain question answering. In *Findings of the Association for Computational Linguistics: ACL 2024*, pages 6963–6975, Bangkok, Thailand, August 2024. Association for Computational Linguistics.

[53] Xueguang Ma, Xinyu Zhang, Ronak Pradeep, and Jimmy Lin. Zero-shot listwise document reranking with a large language model. *arXiv preprint arXiv:2305.02156*, 2023.

[54] Weiwei Sun, Lingyong Yan, Xinyu Ma, Shuaiqiang Wang, Pengjie Ren, Zhumin Chen, Dawei Yin, and Zhaochun Ren. Is ChatGPT good at search? investigating large language models as re-ranking agents. In Houda Bouamor, Juan Pino, and Kalika Bali, editors, *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 14918–14937, Singapore, December 2023. Association for Computational Linguistics.

[55] Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. Rankvicuna: Zero-shot listwise document reranking with open-source large language models. arXiv preprint arXiv:2309.15088, 2023.

[56] Ronak Pradeep, Sahel Sharifymoghaddam, and Jimmy Lin. Rankzephyr: Effective and robust zero-shot listwise reranking is a breeze! arXiv preprint arXiv:2312.02724, 2023.

[57] Crystina Zhang, Sebastian Hofstätter, Patrick Lewis, Raphael Tang, and Jimmy Lin. Rank-without-gpt: Building gpt-independent listwise rerankers on open-source large language models. In *Advances in Information Retrieval : 47th European Conference on Information Retrieval, ECIR 2025, Lucca, Italy, April 6–10, 2025, Proceedings, Part II*, pages 233–247, Berlin, Heidelberg, 2025. Springer-Verlag.

[58] Qi Liu, Bo Wang, Nan Wang, and Jiaxin Mao. Leveraging passage embeddings for efficient listwise reranking with large language models. In *Proceedings of the ACM on Web Conference 2025 (WWW '25)*, pages 4274–4283, New York, NY, USA, 2025. Association for Computing Machinery.

[59] Rodrigo Nogueira, Zhiying Jiang, Ronak Pradeep, and Jimmy Lin. Document ranking with a pre-trained sequence-to-sequence model. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 708–718, Online, 2020. Association for Computational Linguistics.

[60] Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, Benjamin Newman, Binhang Yuan, Bobby Yan, Ce Zhang, Christian Cosgrove, Christopher D. Manning, Christopher Ré, Diana Acosta-Navas, Drew A. Hudson, et al. Holistic evaluation of language models. *arXiv preprint arXiv:2211.09110*, 2023. Preprint.

[61] Zach Nussbaum and Brandon Duderstadt. Training sparse mixture of experts text embedding models. *arXiv preprint arXiv:2502.07972*, 2025.

[62] Marc Marone, Orion Weller, William Fleshman, Eugene Yang, Dawn Lawrie, and Benjamin Van Durme. mmbert: A modern multilingual encoder with annealed language learning, 2025.

[63] Ricardo Rei, Craig Stewart, Ana C. Farinha, and Alon Lavie. COMET: A neural framework for MT evaluation. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2685–2702. Association for Computational Linguistics, November 2020.

[64] Meta. Introducing meta llama 3: The most capable openly available llm to date, 2024.

[65] DeepL SE. DeepL Translator, 2024. Accessed: 2025-10-22.

[66] Soomin Kim. googletrans: Free and Unlimited Google Translate API for Python, 2024. Accessed: 2025-10-22.

[67] Argos Open Technologies. LibreTranslate: Open-Source Machine Translation API, 2024. Accessed: 2025-10-22.

[68] Stephen Robertson and Hugo Zaragoza. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389, 2009.

[69] Roksana Goworek, Harpal Singh Karlcut, Hamza Shezad, Nijaguna Darshana, Abhishek Mane, Syam Bondada, Raghav Sikka, Ulvi Mammadov, Rauf Allahverdiyev, Sriram Satkirti Purighella, Paridhi Gupta, Muhinyia Ndegwa, Bao Khanh Tran, and Haim Dubossarsky. SenWiCh: Sense-annotation of low-resource languages for WiC using hybrid methods. In *Proceedings of the 7th Workshop on Research in Computational Linguistic Typology and Multilingual NLP*, pages 61–74, Vienna, Austria, August 2025. Association for Computational Linguistics.

[70] Jörg Tiedemann. The tatoeba translation challenge – realistic data sets for low resource and multilingual MT. In *Proceedings of the Fifth Conference on Machine Translation (WMT)*, pages 1174–1182, Online, November 2020. Association for Computational Linguistics.

[71] Kalervo J"arvelin and Jaana Kek"al"ainen. Cumulated gain-based evaluation of ir techniques. *ACM Transactions on Information Systems (TOIS)*, 20(4):422–446, 2002.

[72] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul N. Bennett, Junaid Ahmed, and Arnold Overwijk. Approximate nearest neighbor negative contrastive learning for dense text retrieval. *arXiv preprint arXiv:2007.00808*, 2020.

[73] Tianhao Shen, Mingtong Liu, Ming Zhou, and Deyi Xiong. Recovering gold from black sand: Multilingual dense passage retrieval with hard and false negative samples. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing (EMNLP '22)*, pages 10659–10670, 2022.

[74] Thilina Chaturanga Rajapakse, Andrew Yates, and Maarten de Rijke. Negative sampling techniques for dense passage retrieval in a multilingual setting. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 575–584, 2024.

[75] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *Introduction to Information Retrieval*. Cambridge University Press, Cambridge, UK, 2008.

[76] Patrick Littell, David R. Mortensen, Ke Lin, Katherine Kairis, Carlisle Turner, and Lori Levin. URIEL and lang2vec: Representing languages as typological, geographical, and phylogenetic vectors. In Mirella Lapata, Phil Blunsom, and Alexander Koller, editors, *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 8–14, Valencia, Spain, April 2017. Association for Computational Linguistics.

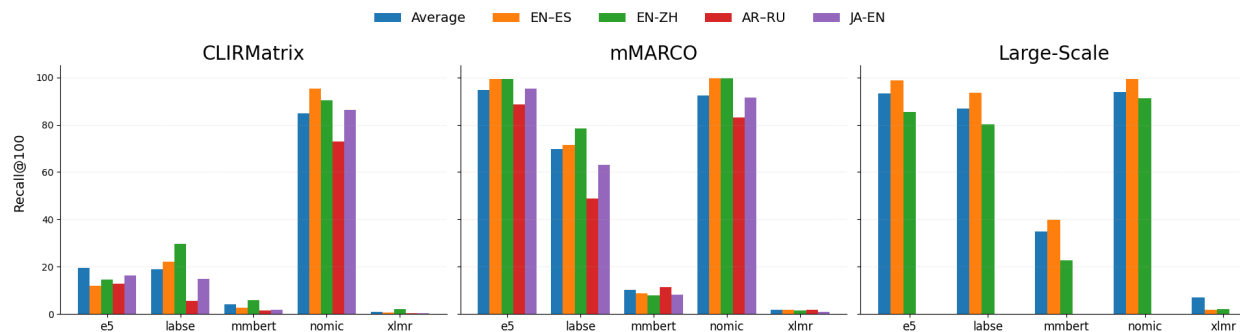## A  Embedding-Based Retrieval on Translated Corpora Performance



Figure 7: Performance (Recall@100) of five pretrained multilingual encoders using cosine similarity between query and document embeddings on translated document corpora. Results are averaged across all language pairs and shown for selected language pairs.

## B  Example Query-Document Pairs

| Dataset | Query | (Translated) Document | Description |
|---|---|---|---|
| **CLIRMatrix** [29] | Finland | *Republic of Finland (Finnish: Suomen tasavalta; Swedish: Republiken Finland) is a country in Northern Europe and a member of the European Union since 1995. Finland is bordered by the Baltic Sea...* | Queries: Wikipedia titles. Documents: Wikipedia articles. Relevance generated via monolingual IR in the query language and propagated across languages using Wikidata interlanguage links. |
| **Large-Scale** [1] | *Woolmer — is a place in hampshire, england.* | *Woolmer — Woolmer is a village in East Hampshire, England. It lies 36 km east of Winchester and 71 km southwest of London. In 2010, the village had a population of 550.* | Queries: first sentences of English Wikipedia articles. Documents: first 200 words of non-English Wikipedia articles (25 languages), with quality filtering. Relevance derived from interlanguage links. |
| **mMARCO** [28] | where can i find aztec healing clay | *Secret Aztec healing clay is bentonite clay from Death Valley, California, where it is sun-dried for up to six months at temperatures that sometimes reach 134 degrees. Use it for facials, acne, body wraps, clay baths, foot baths, and chilled clay for knee pads and insect bites...* | Queries: anonymised Bing search questions. Documents: MS MARCO web passages. Queries and passages machine-translated into multiple languages using Helsinki-NLP models and Google Translate. |

Table 9: Example queries, documents, and dataset construction details for the three datasets used in this study.

## C  Dataset Preprocessing

**De-duplication and Balancing (mMARCO).**  For the mMARCO dataset, we first deduplicated the translated passages across languages, resulting in a pool of 7,433 unique documents. These documents were then evenly distributed across the 14 target languages, assigning exactly one unique passage to each language to ensure balanced representation. For retrieval experiments, the full set of 7,433 queries was used for each query language, thereby standardizing the evaluation conditions across all language pairs.

**Query Selection.** For CLIRMatrix and the Large-Scale dataset, we sampled 1,000 queries per language pair. In mMARCO, which contains fewer available query–document pairs due to the deduplication step, we used 530–531 queries per language pair (depending on rounding). Sampling was performed uniformly at random from the available query pool for each language combination to avoid potential bias toward any specific topic or language domain.

# D Retrieval Pipeline Settings

## D.1 Encoder Model Details

For all retrieval experiments:

- Tokenization: model-native, maximum 512 tokens
- Pooling: mean pooling over last hidden layer
- Embedding dimension:
    - XLM-R: 768
    - LaBSE: 768
    - multilingual-E5: 1024
    - Nomic: 768
    - mmBERT: 768
- Normalization: L2 norm applied before retrieval
- Similarity: cosine similarity
- Caching: all document embeddings precomputed and cached

## D.2 Machine Translation: Evaluation and Settings

To standardize document representations across languages, all non-English documents from CLIRMatrix and Large-Scale were translated into English. For mMARCO, the provided English passages were used directly. We evaluated four MT systems on 100 randomly sampled documents per language using both semantic adequacy and fluency metrics: **DeepL** (2024 API), **NLLB-200**, **googletrans**, and **LibreTranslate** (self-hosted v1.5).

**Translation configuration.** Full-scale translation was performed using the `facebook/nllb-200-distilled-1.3B` checkpoint via the HuggingFace `translation` pipeline. We used a maximum input length of 1,200 tokens, beam search with `beam=5`, default length penalty, and batch size adapted to available device memory. English and Simple English source documents were copied without modification.

For comparison only, additional translation methods were implemented: DeepL (batch size 10), googletrans with retry-based fallbacks, and LibreTranslate (batch size 20 with per-document fallback). These were used solely for evaluation due to rate limits and instability at scale.

**Quality evaluation.** Semantic adequacy was assessed using COMET-QE (WMT22 COMET-Kiwi) in reference-free mode, comparing each original document with its translation. Fluency was estimated using the negative log-likelihood and perplexity of LLaMA-3.1-8B, computed over translated English text. Across languages, NLLB-200 achieved competitive adequacy and fluency scores (second only to DeepL) while remaining fully scalable and license-permissible. It was therefore selected as the default translation system for all primary experiments.

## D.3 Embedding Setup

For all cosine-based retrieval experiments, document embeddings were computed once per encoder (batch inference using model-native tokenisation, max 512 tokens), mean-pooled over the final hidden layer, L2-normalized, and cached. Query embeddings were generated on-the-fly using the same configuration. When multiple GPUs were available, documents were distributed across devices and the resulting embeddings concatenated; otherwise, encoding was performed in batches on a single GPU. Retrieval scores were computed via cosine similarity,

$$\text{sim}(q, d) = \frac{q \cdot d}{\|q\| \|d\|},$$

applied over normalised embeddings using matrix multiplication [75]. Documents were ranked by sorting scores in descending order. Optional query-side projection functions were used only in experiments involving alignment intervention.

### D.4 Contrastive Fine-Tuning

We fine-tuned Nomic and XLM-R at three levels: (i) word-level; (ii) phrase-level; and (iii) query–document (QD). All models were fine-tuned end-to-end using a contrastive objective with temperature scaling. Word-level training used a symmetric contrastive loss with dataset-provided negatives, while phrase- and QD-level training used an InfoNCE loss and in-batch negatives. No adapters or projection layers were added on top of the base encoders.

**Hyperparameters.** Table 10 summarises the training configuration used for all contrastive setups. Each set up was trained for a maximum of 10 epochs, with post-training best-epoch selection, based on performance on validation data. Batch size was fixed at 16. The learning rate was set to $1 \times 10^{-5}$ using the Adam optimizer. The loss temperature was $\tau = 0.1$. No scheduler, warmup, weight decay, or gradient clipping was applied.

| Setting | Value |
|---|---|
| Optimizer | Adam |
| Learning rate | 1e–5 |
| Batch size | 16 |
| Loss type | ContrastiveLoss/InfoNCE |
| Temperature ($\tau$) | 0.1 |
| Max length | 128 (word/phrase), 512 (QD) |
| Epochs | 10 |
| In-batch negatives | phrase/QD |
| Dataset-provided negatives | word |

Table 10: Contrastive learning hyperparameters for Nomic and XLM-R.

**Training data.**

- **Word-level:** Trained on XL-WiC, Am$^2$iCO, and SenWiCh, retaining only language subsets overlapping with our CLIR datasets. Evaluation was performed using accuracy-based threshold selection.
- **Phrase-level:** Parallel sentence pairs from Tatoeba, sampled up to 1,000 examples per language pair. Best checkpoint selected using the sum of A→B and B→A Recall@1.
- **Query–Document:** CLIRMatrix, mMARCO, and Large-Scale training pairs (one model per dataset). In-batch negatives were used. Validation used macro nDCG@10 across language pairs.

**Implementation.** All experiments used `PyTorch` with custom training loops. Checkpoints were selected using validation performance and non-optimal checkpoints were discarded.

### D.5 Re-Ranking (Cross-Encoder) Settings

| Setting | Value |
|---|---|
| Model architectures | XLM-R-base, Nomic (BERT-style) |
| Max sequence length | 512 (query + document) |
| Per-device train batch size | 8 |
| Per-device eval batch size | 8 |
| Optimizer | AdamW (default in `sentence-transformers`) |
| Learning rate | 2e–5 |
| LR scheduler | linear with warmup ratio 0.1 |
| Num train epochs | 10 |
| Loss | Binary cross-entropy over sigmoid scores |
| Train/validation split | 90/10, stratified by label |
| Seed | 42 |
| Precision | fp16 on CUDA, full precision otherwise |
| Model selection | best `eval_loss`, loaded at end |

Table 11: Cross-encoder fine-tuning hyperparameters.

We fine-tune XLM-R-base and Nomic as cross-encoders that jointly encode each query–document pair and output a scalar relevance score. Training is implemented with `sentence-transformers` using binary labels (1 for relevant, 0 for non-relevant) and the `BinaryCrossEntropyLoss` objective. For each dataset and negative sampling condition

(easy vs. hard), we train a separate cross-encoder on CSV files containing `text1` (query), `text2` (document), and `label`. The training/validation split is 90/10, stratified by label. Hyperparameters are summarised in Table 11.

Each positive query–document pair is paired with non-relevant documents according to the difficulty setting. *Easy negatives* are randomly sampled non-relevant documents. *Hard negatives* are derived from dataset-provided candidate rankings and correspond to seemingly relevant but non-match documents. At inference time, we re-rank the top-100 candidates from the BM25 first-stage ranker. For each query, we form (*query*, *document*) pairs and score them using the trained cross-encoder (batch size 32 during prediction). If the gold document is not present in the BM25 top-100, it is injected by replacing the document at rank 100, ensuring that re-ranking always operates over a candidate set containing the relevant document. Final rankings are obtained by sorting documents by cross-encoder score in descending order, and evaluated using Recall@10 and nDCG@100.

## E   Approximate Nearest Neighbour Search and Efficiency Calculations

For efficient retrieval, we used HNSW (`hnswlib`) over L2-normalized document embeddings. The index was created with `M=16` and `ef_construction=200`, storing one entry per document (ID-aligned with embeddings). At inference, queries were encoded using the same encoder and searched with `ef=50` (or `max(50, 2×top_k)` where adaptive). No projection was applied in main experiments.

Cosine similarity was computed as:
$$\text{sim}(q, d) = 1 - \text{dist}_{\text{cosine}}(q, d),$$
and candidates were ranked by descending similarity. Index metadata (model name, embedding dimension, HNSW parameters, and doc IDs) was saved for reproducibility.

### E.1   Latency Measurement Setup

To compare the efficiency of exact (dense) and approximate (ANN) retrieval, we measured end-to-end inference latency for both methods across all query–document language pairs in each dataset. For each of the five encoder models, retrieval was performed sequentially on each language pair: first using direct cosine similarity over precomputed embeddings (dense retrieval), followed immediately by ANN-based retrieval under identical conditions. Each method produced a timestamp upon completion.

Since different encoder models have varied processing speeds, raw latency values are not directly comparable. To account for this, timestamps were first interleaved in execution order (dense $\rightarrow$ ANN $\rightarrow$ dense $\rightarrow$ ANN, etc.) and min–max normalised to the $[0, 1]$ range. Normalised values were then scaled by the number of evaluated language pairs in each dataset (*CLIRMatrix*: 56, *mMARCO*: 196, *Large-Scale*: 26).

For each dataset, we computed the mean latency difference between:

- dense $\rightarrow$ ANN retrieval (i.e., ANN completion time relative to dense), and
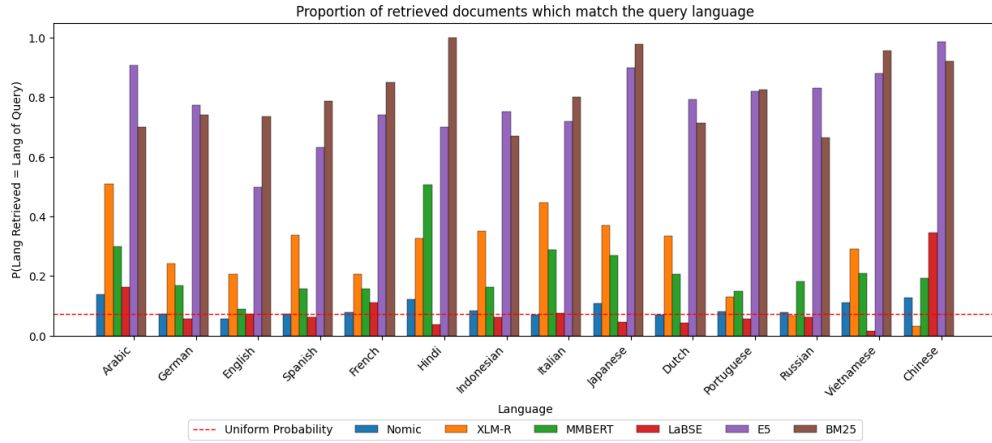- ANN $\rightarrow$ dense retrieval (reverse direction in the alternating sequence).

The reported latency values correspond to the average of these normalised differences across all models.
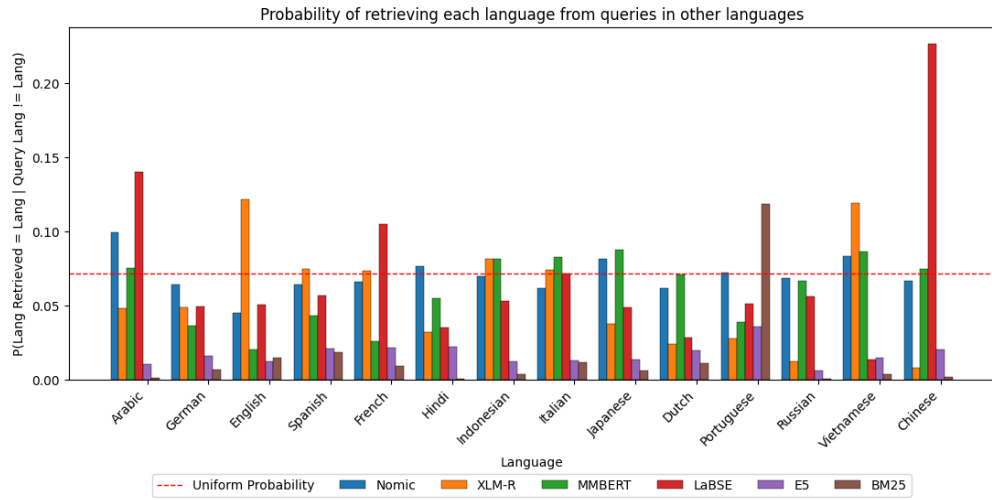
## F   Linguistic Similarity Analysis

To assess whether retrieval performance correlates with linguistic proximity, we computed cosine similarity between query and document languages using typological vectors from `lang2vec` [76]. Feature sets included geographical, syntactic, phonological, phonetic inventory, and family-based attributes; vectors were masked for missing values prior to similarity calculation.

Spearman correlations between language similarity and retrieval *Recall@100* were computed per evaluation file (dataset–model). For comparison across feature sets, we generated a per-file correlation matrix. Results were aggregated to quantify the influence of linguistic similarity on cross-lingual retrieval performance.

## G   Document Retrieval Language Bias of Evaluated Models

(a) Probability that models retrieve a document written in the same language as the query. BM25 exhibits the strongest bias due to reliance on lexical overlap, while Nomic is the most balanced. Notably, multilingual-E5 also shows a strong preference towards the query language despite strong retrieval performance in Figure 4



(b) Distribution of retrieved document languages when the relevant document is not in the query language. Most models retrieve non-query languages less frequently than the uniform probability.

Figure 8: Comparative linguistic retrieval bias across all evaluated models (Nomic, XLM-R, mmBERT, LaBSE, multilingual-E5 and BM25).