

# M3DR: Towards Universal Multilingual Multimodal Document Retrieval

Adithya S Kolavi  
CognitiveLab  
Bengaluru, India

adithyaskolavi@gmail.com

Vyoman Jain  
CognitiveLab  
Bengaluru, India

vyomanjain@gmail.com

## Abstract

Multimodal document retrieval systems have shown strong progress in aligning visual and textual content for semantic search. However, most existing approaches remain heavily English-centric, limiting their effectiveness in multilingual contexts. In this work, we present M3DR (Multilingual Multimodal Document Retrieval), a framework designed to bridge this gap across languages, enabling applicability across diverse linguistic and cultural contexts. M3DR leverages synthetic multilingual document data and generalizes across different vision-language architectures and model sizes, enabling robust cross-lingual and cross-modal alignment. Using contrastive training, our models learn unified representations for text and document images that transfer effectively across languages. We validate this capability on 22 typologically diverse languages, demonstrating consistent performance and adaptability across linguistic and script variations. We further introduce a comprehensive benchmark that captures real-world multilingual scenarios, evaluating models under monolingual, multilingual, and mixed-language settings. M3DR generalizes across both single dense vector and ColBERT-style token-level multi-vector retrieval paradigms. Our models, NetraEmbed and ColNetraEmbed achieve state-of-the-art performance with  $\sim 150\%$  relative improvements on cross-lingual retrieval.

## 1. Introduction

The exponential growth of digital documents across global enterprises, research institutions, and digital libraries has created an urgent need for effective multilingual document retrieval systems. While text-based retrieval has achieved remarkable success across languages [20, 24], retrieval based on traditional OCR-based pipelines face significant challenges: information loss from discarding visual elements (charts, diagrams, layout), brittleness with diverse fonts and scripts, and cascading errors particularly severe in low-resource languages.

Recent vision-based approaches like ColPali [11] have demonstrated promising results by directly encoding document images using vision language models (VLMs), thereby eliminating OCR dependencies and preserving rich visual-textual information. However, these systems remain predominantly English-centric, and our preliminary analysis shows that they perform very poorly on multilingual content, a critical limitation in a world where documents span hundreds of languages.

This work addresses a fundamental research question: *Can we develop universal document retrievers that maintain high performance across typologically diverse languages without sacrificing English competitiveness?*

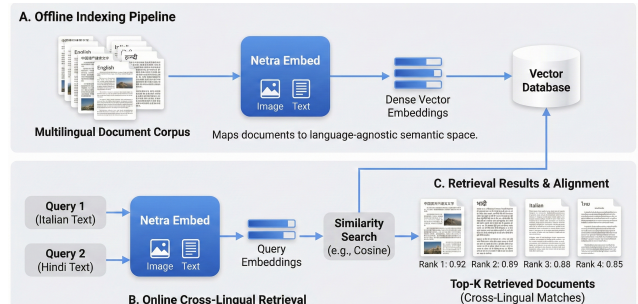


Figure 1. Overview of NetraEmbed, our multilingual multimodal document embedding model. (A) Offline indexing encodes documents into dense vectors in a shared semantic space, (B) online retrieval processes cross-lingual queries, and (C) results show effective matching across diverse scripts and languages.

We present M3DR, a scalable framework for training multilingual multimodal document retrievers across 22 typologically diverse languages spanning Latin, Devanagari, Dravidian, CJK, and other script families. M3DR generalizes across both single dense vector and ColBERT-style multi-vector retrieval paradigms. With our single dense vector model, we achieve 152% relative improvement over baselines on cross-lingual retrieval (0.716 vs 0.284 NDCG@5) establishing new state-of-the-art performance for multilingual multimodal document retrieval.

We make **three key contributions**:

- We develop a synthetic data generation pipeline for generating multilingual multimodal document retrieval data at scale. Our approach combines layout-aware document translation, authentic typography rendering with language-specific fonts, and VLM-based query synthesis (Llama 3.1 90B Vision [9], Llama 4 Scout [1]), producing  $\sim 1\text{M}$  parallel document images across 22 languages
- We introduce a comprehensive multilingual multimodal document retrieval benchmark **Nayana IR Bench**<sup>1</sup> with 23 datasets (1 cross lingual plus 22 monolingual),  $\sim 28\text{K}$  document images, and  $\sim 5.4\text{K}$  queries in BEIR compatible format [57]. This benchmark supports standardized evaluation across diverse script families for both cross lingual and monolingual retrieval capabilities.
- We release two 4B parameter models demonstrating M3DR’s generalization across retrieval paradigms. **NetraEmbed**<sup>2</sup> (Netra is Sanskrit for eye or vision), a single dense vector model using Matryoshka representation learning that supports multiple output dimensions (768, 1536, 2560), and **ColNetraEmbed**<sup>3</sup>, a ColBERT style multi vector variant. NetraEmbed reaches state of the art results on cross lingual and mono lingual retrieval while keeping strong English performance.

## 2. Related Work

Our work draws upon and extends several research threads. We organize related work into five key areas: visual document understanding, vision-based document retrieval, multimodal retrieval systems, multilingual embeddings, and training strategies.

**Visual Document Understanding.** Traditional document understanding relies on OCR followed by text processing [3, 18], suffering from information loss and error propagation. Recent vision-language models enable end-to-end understanding. mPLUG-DocOwl [16] introduced structure-aware learning for OCR-free document understanding across documents, tables, and charts. mPLUG-DocOwl2 [17] proposed high-resolution compression, reducing visual tokens from thousands to 324 while maintaining performance. HRVDA [35] addressed high-resolution challenges through content and instruction filtering. For long documents, LongDocURL [6] introduced comprehensive benchmarks spanning 33K+ pages, while MMLongBench-Doc [40] revealed performance gaps in current VLMs on lengthy PDFs (avg. 49.4 pages). SlideVQA [51] addressed multi-image document VQA with slide decks. While these works demonstrate impressive capabilities, they primarily focus on question-answering rather than retrieval, and evaluations center on English.

<sup>1</sup>NayanaIR bench

<sup>2</sup>NetraEmbed

<sup>3</sup>ColNetraEmbed

M3DR complements this by enabling multilingual document retrieval as a critical first step for document-centric RAG.

**Vision-Based Document Retrieval.** ColPali [11] pioneered vision-based document retrieval by adapting late-interaction mechanisms from ColBERT [25] to visual documents, achieving substantial improvements over OCR-based pipelines on the ViDoRe [41] benchmark. A reproducibility study [47] confirmed ColPali’s effectiveness and provided insights into query-patch matching. ModernVBERT [56] proposed a compact 250M-parameter encoder optimized for document retrieval, competitive with 10 $\times$  larger models through principled design. Document Screenshot Embedding [37] explored using screenshots as unified input, demonstrating versatility across document types. Guided Query Refinement [58] proposed test-time optimization for query embeddings, while EDJE [52] introduced efficient discriminative joint encoders for large-scale reranking. While these methods demonstrate strong English performance, none systematically address multilingual scenarios. M3DR extends vision-based retrieval to 22 languages, validating that these approaches generalize across linguistic and script boundaries.

**Multimodal Retrieval and RAG.** Universal multimodal retrieval aims to handle diverse modalities. UniIR [61] introduced instruction-guided retrieval with the M-BEIR benchmark. MM-Embed [34] trained multimodal embedders using MLLMs with modality-aware hard negative mining. GME [64] improved upon MM-Embed through high-quality fused-modal data synthesis. U-MARVEL [32] provided comprehensive study of key factors including progressive transition and hard negative mining. LamRA [36] re-purposed generative LMMs for retrieval through unified structure learning. For document-specific RAG, M3DocRAG [21] combined ColPali with MLMs for multi-page, multi-document QA. VisRAG [63] established vision-based RAG pipelines, demonstrating 20-40% gains over text-based RAG. VISA [38] enhanced RAG with visual source attribution using bounding boxes. Benchmarking efforts like MMDocIR [8], REAL-MM-RAG [60], and UDA [19] have been crucial. M3DR fills the gap by combining document-specific understanding with comprehensive multilingual support.

**Multilingual and Cross-Lingual Retrieval.** Text-based multilingual embeddings like mBERT [7], XLM-R [5], and LaBSE [12] achieve cross-lingual transfer but discard visual information. The GTE series [33] employed multi-stage contrastive learning on large-scale multilingual data. Jina-embeddings-v3 [50] introduced task-specific LoRA adapters, while Jina-v4 [15] extended to multimodal multilingual retrieval with separate text and image encoders, achieving strong performance on JVDR benchmark. However, separate encoders may miss fine-



grained visual-textual interactions that vision-based approaches capture. VLM2Vec [22] proposed converting VLMs into embedding models through contrastive training on MMEB. xVLM2Vec [45] extended to multilingual scenarios through self-knowledge distillation, though limited to European languages. BGE-M3 [4] supports 100+ languages over text. Unlike these works which focus on text-only multilingual retrieval or limited-language multimodal scenarios, M3DR systematically addresses multilingual visual document retrieval across diverse scripts.

**Training Strategies.** Modern dense retrievers rely on contrastive learning with InfoNCE loss [46] and hard negative mining. NV-Retriever [44] introduced positive-aware mining for effective false negative removal. LLaRA [30] adapted LLMs for retrieval through embedding-based auto-encoding. Llama2Vec [31] demonstrated unsupervised adaptation of Llama-2 for dense retrieval. DRAMA [39] leveraged LLMs for diverse data augmentation. Matryoshka representation learning [29] enables flexible embedding dimensions, critical for deployment efficiency. M3DR builds upon these strategies, conducting systematic ablations on loss functions, negative sampling, and Matryoshka training, identifying optimal configurations for multilingual multimodal document retrieval.

### 3. Training Data and Benchmark

A critical challenge in developing multilingual multimodal document retrieval systems is the scarcity of high-quality training data and comprehensive evaluation benchmarks. In this section, we describe our approach to addressing this challenge through: (1) a scalable synthetic data generation pipeline that creates parallel multilingual document corpora for training, (2) query synthesis using large vision-language models, and (3) construction of the **Nayana-IR Benchmark** for standardized evaluation.

#### 3.1. Training Dataset Construction

##### 3.1.1. Synthetic Parallel Corpus Generation

Building on the Nayana framework [26, 27], which demonstrated the effectiveness of layout-aware synthetic data generation for adapting vision-language models to low-resource languages, we extend this approach to create a large-scale multilingual parallel corpus for document retrieval. While the original Nayana work focused on OCR and document understanding tasks, we adapt their synthesis pipeline for retrieval by generating document images paired with diverse queries across 22 languages.

**Source Documents.** We curate 50,000 diverse English documents images spanning scientific papers, technical reports, educational materials, business documents, and forms from publicly available sources.

**Layout-Aware Translation.** Following the Nayana approach, our pipeline (Figure 2) employs: (1) *Layout De-*

*tection*: DocLayout-YOLO [62], Docling [2] extract text regions, visual elements, tables, and layout metadata. (2) *Neural Translation*: Context-aware translation using NLLB-200 [55] and language-specific models [14] preserves document semantics across 22 target languages. (3) *Visual Rendering*: Authentic typography with Noto Sans fonts [13] for universal script coverage, language-specific layout rules (character spacing, line breaking, text direction), and high-resolution rendering (1024-2048px height) maintaining visual elements from source documents.

##### 3.1.2. Query Synthesis

To support diverse retrieval scenarios, we generate multiple types of queries for each document image using state-of-the-art large vision-language models (Llama 3.1 90B Vision [9] and Llama 4 Scout [1]). We generate 5 diverse query types per document: (1) Basic factual questions (2×), (2) Long-answer questions (1×), (3) Multiple choice/short answer (1×), (4) Cross-paragraph reasoning (1×). LLM-based filtering ensures quality. This yields a training corpus of 1M document images with corresponding queries across 22 languages. **Query Distribution.** Training queries: 70% questions, 15% answers (reverse lookup), 15% text snippets (keyword search).

### 3.2. Nayana-IR Benchmark

While the synthetic parallel corpus enables model training, comprehensive evaluation requires carefully curated test sets that are completely separate from training data. We introduce **Nayana-IR**, the first comprehensive benchmark for multilingual multimodal document retrieval, comprising 23 datasets across 22 languages with both cross-lingual and monolingual evaluation protocols.

#### 3.2.1. Benchmark Structure

**Cross-Lingual Dataset.** 5,870 parallel document images across all 22 languages with 1,000 queries distributed uniformly (~45 per language). Evaluates: *Can models retrieve documents in any language given queries in any language?* Binary relevance (score 2: exact match, score 0: non-relevant).

**Monolingual Datasets (22).** Per-language datasets with ~1,000 documents and ~200 queries each. Evaluates: *Can models retrieve documents written in the same language given queries in that language?* Graded relevance (score 2: exact match, score 1: same document partial match, score 0: non-relevant). **Format and Metrics.** BEIR compatible structure [57] enables standardized evaluation. Metrics: NDCG@5 and NDCG@10, Recall@5 and Recall@10, MAP@10, and MRR@10. The benchmark provides comprehensive coverage across script families with balanced dataset sizes: Latin script (English, Spanish, French, German, Italian), Devanagari (Hindi, Marathi, Sanskrit), Dravidian (Kannada, Telugu, Tamil, Malayalam), CJK (Chi-

## M3DR: End-to-End Multilingual Multimodal Document Retrieval

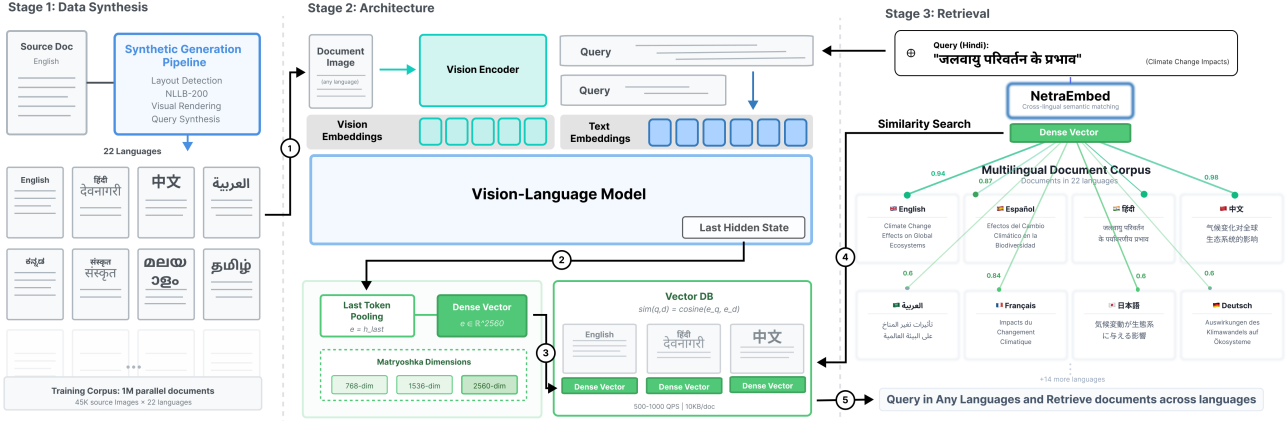


Figure 2. **M3DR Framework Overview.** Our complete pipeline encompasses synthetic data generation (layout detection, neural translation to 22 languages, visual rendering with authentic typography), query synthesis using large VLMs, dense embedding model training with Matryoshka representation learning, and multilingual document retrieval across diverse script families.

nese, Japanese, Korean), and other scripts (Arabic, Bengali, Gujarati, Odia, Punjabi, Russian, Thai). Each monolingual dataset contains approximately 1000 documents and 200 queries, while the cross lingual dataset spans all 22 languages with 5870 parallel documents and 1000 queries. In total this benchmark comprises **23 datasets, 27870 images, and 5400 queries.**

## 4. Methodology

Our design prioritizes: (1) *architectural flexibility* to accommodate VLM backbones from 256M to 4B parameters, (2) *efficiency and scalability* through single dense vector embeddings with flexible dimensionality and ColBERT-Style multi-vector model, (3) *multilingual generalization* across 22 languages with diverse scripts.

### 4.1. Model Architecture

#### 4.1.1. Single Dense Vector Model

Single dense vector models produce fixed-size vectors enabling efficient retrieval via approximate nearest neighbor search [23, 42]. Given text query  $q$  and document image  $d$ , we process each through a VLM backbone to obtain token-level hidden states  $\mathbf{H} \in \mathbb{R}^{n \times h}$ . We primarily use Gemma 3 4B-it [53], which directly outputs 2560-dimensional representations. We apply last token pooling to obtain document-level embeddings. Embeddings are L2-normalized and similarity computed via:

$$\text{sim}(q, d) = \frac{\mathbf{e}_q \cdot \mathbf{e}_d}{\|\mathbf{e}_q\|_2 \|\mathbf{e}_d\|_2} \quad (1)$$

We also evaluate with alternative backbones (Qwen2-VL 2B [59], SmolVLM [43]), employing the same approach with their native output dimensions.

**Matryoshka Representation Learning.** To enable flexible accuracy-efficiency trade-offs, we incorporate Matryoshka learning [29], training embeddings to be truncatable at dimensions 768, 1536, and 2560. During training, for each embedding  $\mathbf{e} \in \mathbb{R}^{2560}$ , we compute losses at three granularities:

$$\mathcal{L}_{\text{Matryoshka}} = \sum_{d \in \{768, 1536, 2560\}} w_d \cdot \mathcal{L}_{\text{base}}(\text{L2-norm}(\mathbf{e}[:d])) \quad (2)$$

where  $w_d = 1/3$ . This enables post-deployment dimension selection: 768-dim (70% storage reduction), 1536-dim (40% reduction), or 2560-dim (maximum accuracy), requiring only embedding truncation without retraining.

#### 4.1.2. ColBERT-Style Multi-Vector Model

To demonstrate M3DR’s generalizability across retrieval paradigms, we also develop a ColBERT-style multi-vector variant following ColPali [11]. Unlike single dense vector models that pool token representations into a single vector, ColBERT-style multi-vector models retain per-token embeddings to enable fine-grained matching.

Given query  $q$  and document  $d$ , we process each through the VLM backbone (Gemma 3 4B-it [53]) to obtain token-level representations  $\mathbf{Q} \in \mathbb{R}^{n_q \times h}$  and  $\mathbf{D} \in \mathbb{R}^{n_d \times h}$ . For Gemma 3 4B-it, each document image produces 256 visual tokens while queries vary based on query length. Similarity is computed via late interaction using MaxSim:

$$\text{sim}_{\text{late}}(q, d) = \sum_{i=1}^{n_q} \max_{j=1}^{n_d} \cos(\mathbf{q}_i, \mathbf{d}_j) \quad (3)$$

where  $\mathbf{q}_i$  and  $\mathbf{d}_j$  are L2-normalized token embeddings. This enables each query token to match its most

similar document token, capturing fine-grained semantic correspondences. Each document image produces  $256 \times 128$ -dimensional embeddings compared to a single 2560-dimensional embedding for single dense vector retrieval.

## 4.2. Training Objectives

### 4.2.1. Loss Functions for Single Dense Vector Models

**BiEncoderLoss (InfoNCE).** Our baseline employs InfoNCE [46] with in-batch negatives. Given batch  $B$  of query-document pairs  $\{(q_i, d_i^+)\}_{i=1}^B$ :

$$\mathcal{L}_{\text{BiEncoder}} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(s_{ii}/\tau)}{\sum_{j=1}^B \exp(s_{ij}/\tau)} \quad (4)$$

where  $s_{ij} = \text{sim}(q_i, d_j)$  and  $\tau = 0.02$ . This ranks positive  $d_i^+$  higher than in-batch documents  $\{d_j\}_{j \neq i}$ .

**BiNegativeCELoss (Hybrid Loss).** When explicit hard negatives  $\{d_i^{-,k}\}_{k=1}^K$  are available, we combine pairwise ranking with InfoNCE:

$$\mathcal{L}_{\text{BiNegCE}} = (1 - \lambda) \cdot \mathcal{L}_{\text{pairwise}} + \lambda \cdot \mathcal{L}_{\text{InfoNCE}} \quad (5)$$

$$\mathcal{L}_{\text{pairwise}} = \frac{1}{BK} \sum_{i=1}^B \sum_{k=1}^K \text{softplus} \left( \frac{s(q_i, d_i^{-,k}) - s(q_i, d_i^+)}{\tau} \right) \quad (6)$$

where  $\lambda = 0.5$  balances discrimination against hard negatives with in-batch diversity, essential for cross-lingual retrieval.

**MatryoshkaBiEncoderLoss.** We wrap base losses with Matryoshka mechanism:

$$\mathcal{L}_{\text{Matryoshka}} = \frac{1}{3} \sum_{d \in \{768, 1536, 2560\}} \mathcal{L}_{\text{base}}(\text{L2-norm}(\mathbf{E}[:, : d])) \quad (7)$$

### 4.2.2. Loss Functions for ColBERT-Style Multi-Vector Models

For the ColBERT-style multi-vector variant, we employ InfoNCE loss with in-batch negatives following ColPali [11]. Given batch  $B$  of query-document pairs, the similarity is first computed via MaxSim (Equation 3), optionally normalized by query length:

$$\text{sim}_{\text{norm}}(q_i, d_j) = \frac{\text{sim}_{\text{late}}(q_i, d_j)}{\text{len}(q_i)} \quad (8)$$

The InfoNCE loss with temperature scaling is then applied:

$$\mathcal{L}_{\text{late}} = -\frac{1}{B} \sum_{i=1}^B \log \frac{\exp(\text{sim}_{\text{norm}}(q_i, d_i^+)/\tau)}{\sum_{j=1}^B \exp(\text{sim}_{\text{norm}}(q_i, d_j)/\tau)} \quad (9)$$

where  $\tau = 0.02$  is the temperature parameter. We found that the standard ColBERT loss (without explicit hard negatives) worked best for our multilingual setting, using in-batch negatives for contrastive learning. This encourages query tokens to find strong matches with relevant document regions while maintaining discrimination against other documents in the batch.

## 4.3. Training Strategy

We investigate three progressively sophisticated strategies for our single dense vector models:

**Strategy 1: Positive-Only.** Uses only in-batch negatives  $\{d_j\}_{j \neq i}$  with BiEncoderLoss. Computationally efficient baseline but most negatives are trivially distinguishable.

**Strategy 2: Document-Level Negatives.** Samples  $K = 3$  hard negatives from nearby pages ( $p \pm \{1, 2, 3\}$ ) within the same document. These share thematic content and layout but lack specific query answers. Trained with BiNegativeCELoss ( $\lambda = 0.3$ ), encouraging fine-grained content discrimination.

**Strategy 3: Mined Hard Negatives.** Mines corpus-wide negatives via: (1) *Text similarity* using BM25 [49] and BGE-M3 [4] on OCR text, (2) *Visual similarity* using CLIP [48] and Jina-CLIP [28] on images, (3) *Fusion* via reciprocal rank fusion where  $k = 60$ . We sample  $K = 3$  negatives from top-20 combined ranking (excluding ground truth). This produces challenging multimodal negatives matching keywords and visual appearance without containing correct answers.

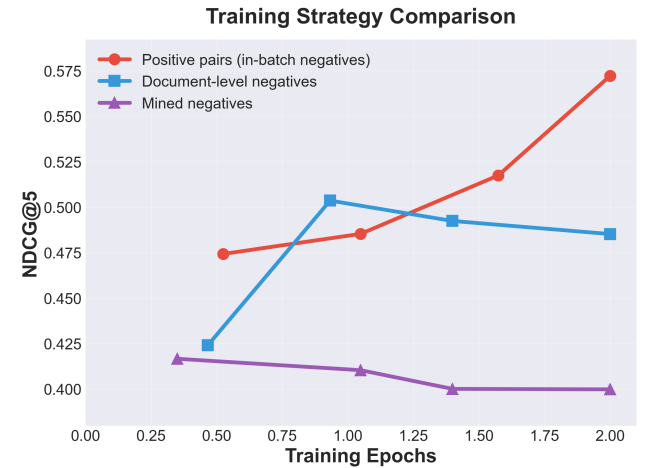


Figure 3. **Training Strategy Comparison.** Positive-only (in-batch negatives) training strategy substantially outperforms document-level negative and hard negative mining (combined text+visual) strategies, with consistent improvements throughout training.



## 5. Experiments

### 5.1. Experimental Setup

**Baselines.** We compare against state-of-the-art vision-language document retrieval models: ColPali-v1.2 [11], ColPali-v1.3, ColQwen2.5-v0.2, ColQwen2-v1.0, ColSmol-500M, GME-Qwen2-VL-2B [64] jina-embeddings-v4 [15], and ColNomic-Embed-Multimodal-3B [54].

**M3DR Models.** We evaluate two 4B-parameter models built on Gemma 3 4B-IT [53]: **NetraEmbed** (single dense vector with Matryoshka representation learning) and **ColNetraEmbed** (ColBERT-style multi-vector following ColPali [11]), both trained on 22 languages.

**Evaluation Protocol.** We evaluate on two benchmarks: (1) *Nayana-IR* covering 22 languages with both cross-lingual and monolingual retrieval tasks, and (2) *ViDoRe v2* [10] for English document retrieval. Primary metric: NDCG@5; secondary metrics: Recall@10, MAP@10, MRR@10.

### 5.2. Main Results

Table 1 presents comprehensive evaluation across multilingual and English benchmarks. **NetraEmbed** reaches 0.716 NDCG@5 on cross-lingual retrieval and 0.738 on monolingual tasks, representing 152% and 80% relative improvements over ColPali-v1.3 baseline (0.284 and 0.410 respectively), achieving state-of-the-art multilingual performance. **ColNetraEmbed** achieves 0.637 NDCG@5 on cross-lingual and 0.670 on monolingual tasks, demonstrating that both paradigms can be effectively trained for multilingual document retrieval, with single dense vectors providing superior efficiency-accuracy trade-offs.

Existing VLM-based retrieval models show substantial performance degradation on multilingual content, with most baselines achieving 0.00-0.32 NDCG@5 on cross-lingual tasks, validating the necessity of explicit multilingual training (§B presents detailed baseline comparisons). On ViDoRe v2 benchmark, NetraEmbed achieves 0.554 NDCG@5, demonstrating competitive performance on English content while prioritizing multilingual capabilities.

### 5.3. Per-Language Performance Analysis

Figure 4 illustrates per-language performance across all 22 languages in our benchmark. NetraEmbed achieves consistent high performance across diverse script families, validating robust cross-lingual transfer. Scaling analysis from 6 to 22 languages (§F) demonstrates progressive improvements as linguistic diversity increases.

### 5.4. Matryoshka Embeddings: Efficiency-Accuracy Trade-offs

Table 2. **Matryoshka Embedding Dimension Trade offs.** Storage and performance for different truncation levels of NetraEmbed on Nayana IR cross lingual benchmark.

Dimensions	Storage/Doc	NDCG@5	Rel. Perf.
768	~3 KB	0.680	95.0%
1536	~6 KB	0.706	98.6%
2560 (full)	~10 KB	0.716	100.0%

Matryoshka representation learning enables flexible post-deployment dimension selection without retraining, critical for adapting to diverse computational budgets. Table 2 quantifies storage-accuracy trade-offs for NetraEmbed. Comprehensive ablations across training configurations and language scales are presented in §D.

The 768-dimensional truncation achieves 95.0% of full model performance (0.680 vs 0.716 NDCG@5) while reducing storage by 70% (~3 KB vs ~10 KB per document), making it ideal for billion-scale deployments or edge devices with strict memory constraints. The 1536-dimensional representation offers a balanced middle ground at 98.6% performance with 40% storage reduction, suitable for production systems balancing accuracy and cost. The full 2560-dimensional representation maximizes retrieval accuracy while maintaining efficient storage requirements compared to token-level approaches, making NetraEmbed suitable for large-scale multilingual document retrieval deployments.

### 5.5. Single Dense Vector vs. ColBERT-Style Multi-Vector Retrieval

M3DR successfully generalizes across both retrieval paradigms. We compare the two approaches to understand architectural trade-offs for multilingual document retrieval.

**Performance.** NetraEmbed consistently outperforms ColNetraEmbed across multilingual tasks: +12.4% on cross-lingual NDCG@5 (0.716 vs 0.637) and +10.1% on monolingual tasks (0.738 vs 0.670). Both models achieve similar English performance on ViDoRe v2 (0.554 vs 0.551), suggesting that the single dense vector approach particularly benefits cross-lingual transfer learning.

**Efficiency.** Single dense vectors offer substantial computational advantages: (1) *Storage*: ~10 KB per document vs. ~2.5 MB for ColBERT-style multi-vector with 256 image tokens (250× reduction), (2) *Retrieval Speed*: Single-vector cosine similarity via HNSW vs. expensive MaxSim computation over token matrices, (3) *Scalability*: Single dense vectors scale efficiently to billion-document corpora with standard vector databases.

Table 1. **Main Results on Nayana-IR and ViDoRe v2.** M3DR models achieve state-of-the-art multilingual performance while maintaining strong English competitiveness. N@5: NDCG@5, R@10: Recall@10, M@10: MAP@10, MRR: MRR@10.

Model	Nayana-IR Cross-Lingual				Nayana-IR Monolingual				ViDoRe v2			
	N@5	R@10	M@10	MRR	N@5	R@10	M@10	MRR	N@5	R@10	M@10	MRR
<i>Baselines</i>												
ColPali-v1.2	0.224	0.237	0.198	0.328	0.402	0.474	0.383	0.412	0.506	0.591	0.411	0.611
ColPali-v1.3	0.284	0.347	0.249	0.403	0.410	0.484	0.393	0.422	0.538	0.627	0.436	0.644
ColQwen2-v1.0	0.050	0.065	0.038	0.109	0.413	0.466	0.398	0.422	0.545	0.640	0.438	0.653
ColQwen2.5-v0.2	0.143	0.160	0.127	0.220	0.453	0.513	0.437	0.464	0.592	0.664	0.484	0.711
ColSmol-500M	0.000	0.000	0.000	0.000	0.224	0.263	0.214	0.229	0.435	0.544	0.347	0.535
GME-Qwen2-VL-2B	0.235	0.308	0.209	0.314	0.444	0.525	0.426	0.452	0.574	0.630	0.466	0.690
ColNomic-Embed-3B	0.315	0.320	0.267	0.444	0.534	0.603	0.515	0.546	0.556	0.633	0.451	0.672
Jina-Embeddings-v4	0.435	0.435	0.390	0.548	X	X	X	X	0.576	0.686	X	X
<i>Our Models</i>												
NetraEmbed	0.716	0.871	0.703	0.775	0.738	0.844	0.709	0.751	0.554	0.637	0.437	0.647
ColNetraEmbed	0.637	0.700	0.610	0.610	0.670	0.764	0.645	0.686	0.551	0.664	0.445	0.445

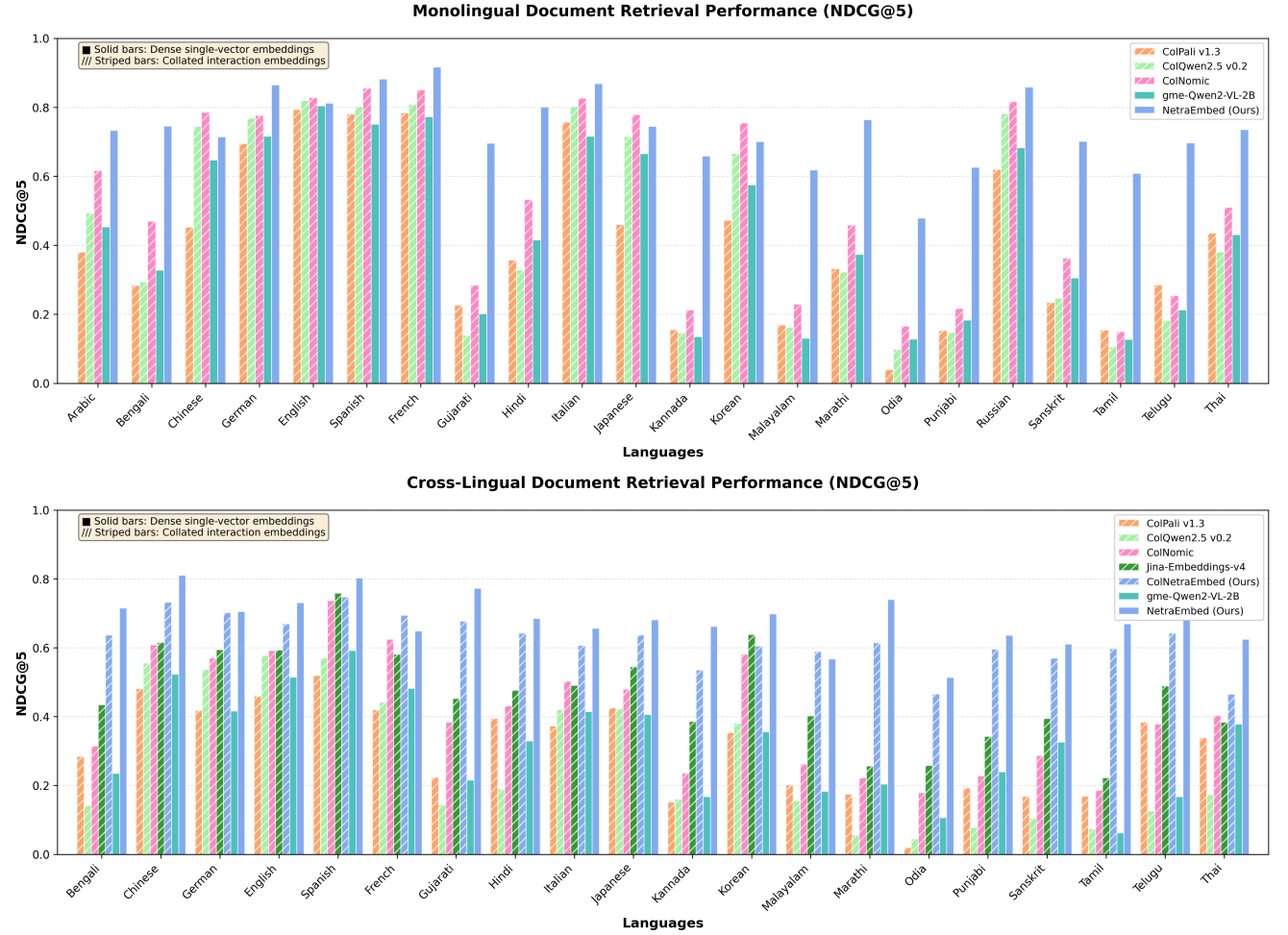


Figure 4. **Per Language Performance Across 22 Languages.** NetraEmbed achieves consistent high performance across all languages and script families such as Latin, Devanagari, CJK, Arabic, and others, while English centric baselines show significant drops on non English content.

Our results demonstrate that M3DR successfully generalizes across both paradigms, with single dense vectors providing the optimal balance of accuracy, efficiency, and multilingual generalization for large-scale document retrieval

deployments. Detailed architectural comparisons and deployment considerations are provided in §G.

## 6. Conclusion

This work addresses a fundamental gap in vision-based document retrieval: the inability of existing systems to handle multilingual content effectively. We presented M3DR, a comprehensive framework that achieves robust multilingual multimodal document retrieval across 22 typologically diverse languages spanning Latin, Devanagari, Dravidian, CJK, Arabic, and other script families.

**Key Results.** NetraEmbed achieves 0.716 NDCG@5 on cross-lingual retrieval and 0.738 on monolingual tasks, representing 152% and 80% relative improvements over the strongest baseline (ColPali-v1.3: 0.284 and 0.410 respectively). These gains demonstrate that explicit multilingual training is essential; existing English-centric VLMs catastrophically fail on non-English content despite strong English performance. NetraEmbed maintains competitive English performance (0.554 NDCG@5 on ViDoRe v2) while dramatically improving multilingual capabilities, validating that cross-lingual optimization does not require sacrificing monolingual competitiveness.

**Critical Design Insights.** Through extensive ablations (§A–J), we identified several key principles: (1) *Base model selection is decisive*: Gemma 3 4B-IT [53]’s multilingual pretraining enabled significant point gains over English-centric models (ColQwen2, ColPali) on cross-lingual tasks. (2) *Training simplicity wins*: BiEncoderLoss with in-batch negatives outperformed complex hard negative mining strategies, as multilingual batch diversity provides sufficient contrastive signal. (3) *Last token pooling dominates*: outperforming mean pooling by 13+ NDCG@5 points for decoder-based VLMs. (4) *Matryoshka enables flexibility*: 768-dimensional truncation retains 95% performance with 70% storage reduction.

**Resources and Reproducibility.** We release the Nayana-IR benchmark (23 datasets, ~28K images, ~5.4K queries in BEIR format) covering cross-lingual and monolingual retrieval across 22 languages. Our models NetraEmbed and ColNetraEmbed are publicly available, alongside training code and evaluation scripts. With LoRA fine-tuning on 4×A100 GPUs requiring only ~12 hours for SOTA results, we aim to democratize multilingual document retrieval research for groups with modest computational resources.

**Broader Impact.** M3DR democratizes document intelligence across linguistic communities, enabling equitable access to information retrieval technologies regardless of language. Applications span multilingual enterprise knowledge management, cross-border research collaboration, educational resource discovery for underserved languages, and cultural heritage preservation for digitized manuscripts. By achieving state-of-the-art multilingual performance without sacrificing English competitiveness,

M3DR demonstrates that inclusive AI systems need not compromise on quality. However, practitioners must monitor performance disparities across languages in deployment and implement continuous evaluation and mitigation strategies to ensure fairness and prevent marginalization of lower-resource languages.

**Limitations and Future Directions.** While our results demonstrate strong generalization across 22 languages, several limitations warrant further investigation: (1) Performance on rare language pairs (e.g., Tamil→Russian) shows 10-12% degradation compared to high-resource pairs, suggesting opportunities for improved cross-lingual alignment techniques. (2) Complex tabular content with language-specific number formatting (Hindi numerals, Arabic-Indic digits) remains challenging, indicating potential benefits from table-aware training objectives. (3) Our evaluation focuses on document-level retrieval; extending to passage-level or region-level retrieval within documents presents interesting future work. (4) Scaling beyond 22 languages to truly low-resource languages and investigating zero-shot transfer to unseen scripts would further validate the framework’s generalization capacity.

## Acknowledgments

This work benefited from compute credits for training, inference, and evaluation provided by [Modal](#), acknowledged as a compute sponsor. Dataset curation and synthesis were supported by the [Meta LLaMA Impact Grant](#) through our [Nayana initiative](#). We appreciate Meta for continued support of our research efforts at [CognitiveLab](#).



## References

- [1] Meta AI. The llama 4 herd: The beginning of a new era of natively multimodal ai innovation. <https://ai.meta.com/blog/llama-4-multimodal-intelligence/>, 2025. Meta AI Blog, April 5, 2025. 2, 3
- [2] Christoph Auer, Michele Dolfi, Joel Carvalho, Cesar Berrospi Ramis Fernandez, Valery Furrer, Marc Marone, Maksym Ahmed, Mikhail Podkorytov, and Peter Staar. Docling: A robust and extensible pipeline for pdf document conversion. *arXiv preprint arXiv:2408.09869*, 2024. 3
- [3] Lukas Blecher, Guillem Cucurull, Thomas Scialom, and Robert Stojnic. Nougat: Neural optical understanding for academic documents. In *arXiv preprint arXiv:2308.13418*, 2023. 2
- [4] Jianlv Chen, Shitao Xiao, Peitian Zhang, Kun Luo, Defu Lian, and Zheng Liu. Bge m3-embedding: Multilingual, multi-functionality, multi-granularity text embeddings through self-knowledge distillation. In *Findings of the Association for Computational Linguistics (ACL)*, pages 2304–2318, 2024. 3, 5
- [5] 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 *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 8440–8451, 2020. 2
- [6] Chao Deng, Jiale Yuan, Pi Bu, Peijie Wang, Zhongzhi Li, Jian Xu, Xiaohui Li, Yuan Gao, Jun Song, Bo Zheng, and Chenglin Liu. Longdocurl: A comprehensive multimodal long document benchmark integrating understanding, reasoning, and locating. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 1135–1159, 2025. 2
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Multilingual bert. In *NAACL*, pages 4171–4186, 2019. 2
- [8] Kuicai Dong, Yujing Chang, Derrick-Goh-Xin Deik, Dexun Li, Ruiming Tang, and Yong Liu. Mmdocir: Benchmarking multi-modal retrieval for long documents. *arXiv preprint arXiv:2501.08828*, 2025. 2
- [9] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 2024. 2, 3
- [10] Manuel Faysse, Hugues Sibille, Gautier Viaud, Hervé Le Borgne, and Céline Hudelot. Vidore: A visual document retrieval benchmark. In *Conference on Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2024. Part of ColPali paper at ICLR 2025. 6
- [11] Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. Colpali: Efficient document retrieval with vision language models. In *International Conference on Learning Representations (ICLR)*, 2025. 1, 2, 4, 5, 6
- [12] Fangxiaoyu Feng, Yinfei Yang, Daniel Cer, Naveen Ariavazhagan, and Wei Wang. Language-agnostic bert sentence embedding. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 878–891, 2022. 2
- [13] Google. Noto fonts: A free font family for all languages. *Google Fonts*, 2024. 3
- [14] Varun Gumma, Pranjal A Chitale, and Kalika Bali. Towards inducing long-context abilities in multilingual neural machine translation models. In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7158–7170, Albuquerque, New Mexico, 2025. Association for Computational Linguistics. 3
- [15] Michael Gunther, Saba Sturua, Mohammad Kalim Akram, Isabelle Mohr, Andrei Ungureanu, Bo Wang, Sedigheh Eslami, Scott Martens, Maximilian Werk, Nan Wang, and Han Xiao. jina-embeddings-v4: Universal embeddings for multimodal multilingual retrieval. In *Proceedings of the 5th Workshop on Multilingual Representation Learning (MRL 2025)*, pages 531–550, 2025. *arXiv preprint arXiv:2506.18902*. 2, 6
- [16] Anwen Hu, Haiyang Xu, Jiabo Ye, Mingshi Yan, Liang Zhang, Bo Zhang, Chen Li, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. mplug-docowl 1.5: Unified structure learning for ocr-free document understanding. In *Findings of the Association for Computational Linguistics: EMNLP 2024*, pages 3096–3120, 2024. 2
- [17] Anwen Hu, Haiyang Xu, Liang Zhang, Jiabo Ye, Mingshi Yan, Ji Zhang, Qin Jin, Fei Huang, and Jingren Zhou. mplug-docowl2: High-resolution compressing for ocr-free multi-page document understanding. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 5817–5834, 2025. 2
- [18] Yupan Huang, Tengchao Lv, Lei Cui, Yutong Lu, and Furu Wei. Layoutlmv3: Pre-training for document ai with unified text and image masking. In *Proceedings of the 30th ACM International Conference on Multimedia*, pages 4083–4091, 2022. 2
- [19] Yulong Hui, Yao Lu, and Huanchen Zhang. Uda: A benchmark suite for retrieval augmented generation in real-world document analysis. *arXiv preprint arXiv:2406.15187*, 2024. 2
- [20] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning. *Transactions on Machine Learning Research (TMLR)*, 2022. 1
- [21] Jaemin Jang, Jaehee Kang, and Nojun Kwak. M3docrag: Multi-modal retrieval is what you need for multi-page multi-document understanding. In *arXiv preprint arXiv:2411.04952*, 2024. 2
- [22] Ziyang Jiang, Rui Meng, Xinyi Yang, Semih Yavuz, Yingbo Zhou, and Wenhui Chen. Vlm2vec: Training vision-language models for massive multimodal embedding tasks. In *International Conference on Learning Representations (ICLR)*, 2025. 3

- [23] Jeff Johnson, Matthijs Douze, and Hervé Jégou. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547, 2021. 4
- [24] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wentaoh Yih. Dense passage retrieval for open-domain question answering. In *EMNLP*, pages 6769–6781, 2020. 1
- [25] Omar Khattab and Matei Zaharia. Colbert: Efficient and effective passage search via contextualized late interaction over bert. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 39–48, 2020. 2
- [26] Adithya Kolavi, Samarth P, and Vyoman Jain. Nayana OCR: A scalable framework for document OCR in low-resource languages. In *Proceedings of the 1st Workshop on Language Models for Underserved Communities (LM4UC 2025)*, pages 86–103, Albuquerque, New Mexico, 2025. Association for Computational Linguistics. 3
- [27] Adithya S Kolavi, Samarth P, and Vyoman Jain. Nayana: A foundation for document-centric vision-language models via multi-task, multimodal, and multilingual data synthesis. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 1678–1687, 2025. 3
- [28] Andreas Koukounas, Georgios Mastrapas, Michael Gunther, Bo Wang, Scott Roberts, Isabelle Mohr, Joan Fontanals Kacprzak, Saba Bhatt, Mohammad Kalim Akram Sturua, Nan Feng, et al. Jina clip: Your clip model is also your text retriever. In *arXiv preprint arXiv:2405.20204*, 2024. 5
- [29] Aditya Kusupati, Gantavya Bhatt, Aniket Rege, Matthew Wallingford, Aditya Sinha, Vivek Ramanujan, William Howard-Snyder, Kaifeng Chen, Sham Kakade, Prateek Jain, et al. Matryoshka representation learning. *Advances in Neural Information Processing Systems*, 35:30233–30249, 2022. 3, 4
- [30] Chaofan Li, Zheng Liu, Shitao Xiao, and Yingxia Shao. Making large language models a better foundation for dense retrieval. *arXiv preprint arXiv:2312.15503*, 2023. 3
- [31] Chaofan Li, Zheng Liu, Shitao Xiao, Yingxia Shao, and Defu Lian. Llama2vec: Unsupervised adaptation of large language models for dense retrieval. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3490–3500, Bangkok, Thailand, 2024. Association for Computational Linguistics. 3
- [32] Xiaojie Li, Chu Li, Shi-Zhe Chen, and Xi Chen. U-marvel: Unveiling key factors for universal multimodal retrieval via embedding learning with mllms. *arXiv preprint arXiv:2507.14902*, 2025. 2
- [33] Zehan Li, Xin Zhang, Yanzhao Zhang, Dingkun Long, Pengjun Xie, and Meishan Zhang. Towards general text embeddings with multi-stage contrastive learning. *arXiv preprint arXiv:2308.03281*, 2023. 2
- [34] Sheng-Chieh Lin, Chankyu Lee, Mohammad Shoeybi, Jimmy Lin, Bryan Catanzaro, and Wei Ping. Mm-embed: Universal multimodal retrieval with multimodal llms. *arXiv preprint arXiv:2411.02571*, 2024. 2
- [35] Chaohu Liu, Kun Yin, Haoyu Cao, Xinghua Jiang, Xin Li, Yinsong Liu, Deqiang Jiang, Xing Sun, and Linli Xu. Hrvda: High-resolution visual document assistant. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 15534–15545, 2024. 2
- [36] Yikun Liu, Pingan Chen, Jiayin Cai, Xiaolong Jiang, Yao Hu, Jiangchao Yao, Yanfeng Wang, and Weidi Xie. Lamra: Large multimodal model as your advanced retrieval assistant. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4015–4025, 2025. arXiv preprint arXiv:2412.01720. 2
- [37] Xueguang Ma, Sheng-Chieh Lin, Minghan Li, Wenhui Chen, and Jimmy Lin. Unifying multimodal retrieval via document screenshot embedding. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2024. 2
- [38] Xueguang Ma, Shengyao Zhuang, Bevan Koopman, Guido Zuccon, Wenhui Chen, and Jimmy Lin. Visa: Retrieval augmented generation with visual source attribution. *arXiv preprint arXiv:2412.14457*, 2024. 2
- [39] Xueguang Ma, Xi Victoria Lin, Barlas Oguz, Jimmy Lin, Wentaoh Yih, and Xilun Chen. Drama: Diverse augmentation from large language models to smaller dense retrievers. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (ACL)*, pages 30170–30186, 2025. arXiv:2502.18460. 3
- [40] Yubo Ma, Yuhang Zang, Liangyu Chen, Meiqi Chen, Yizhu Jiao, Kinze Li, Xinyuan Lu, Ziyu Liu, Yan Ma, Xiaoyi Dong, Pan Zhang, Liangming Pan, Yu-Gang Jiang, Jiaqi Wang, Yixin Cao, and Aixin Sun. Mmlongbench-doc: Benchmarking long-context document understanding with visualizations. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. 2
- [41] Quentin Macé, António Loison, and Manuel Faysse. Vidore benchmark v2: Raising the bar for visual retrieval, 2025. 2
- [42] 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. 4
- [43] Andrés Marafioti, Orr Zohar, Miquel Farré, Merve Noyan, Elie Bakouch, Pedro Cuenca, Cyril Zakka, Loubna Ben Allal, Anton Lozhkov, Nouamane Tazi, et al. Smolvlm: Redefining small and efficient multimodal models. *arXiv preprint arXiv:2504.05299*, 2025. 4
- [44] Gabriel Moreira, Radek Osmulski, Mengyao Xu, Ronay Ak, Benedikt Schifferer, and Even Oldridge. Improving text embedding models with positive-aware hard-negative mining. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management (CIKM)*, 2025. arXiv:2407.15831. 3
- [45] Elio Musacchio, Lucia Siciliani, Pierpaolo Basile, and Giovanni Semeraro. xvlm2vec: Adapting lvlm-based embedding models to multilinguality using self-knowledge distillation. *arXiv preprint arXiv:2503.09313*, 2025. 3
- [46] Aäron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018. 3, 5

- [47] Jingfen Qiao, Jia-Huei Ju, Xinyu Ma, Evangelos Kanoulas, and Andrew Yates. Reproducibility, replicability, and insights into visual document retrieval with late interaction. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR)*, pages 3335–3345, 2025. [2](#)
- [48] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Ilya Sutskever, and Prafulla Dhariwal. Learning transferable visual models from natural language supervision. In *International Conference on Machine Learning (ICML)*, pages 8748–8763, 2021. [5](#)
- [49] Stephen E. Robertson and Hugo Zaragoza. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends in Information Retrieval*, 3(4):333–389, 2009. [5](#)
- [50] Saba Sturua, Isabelle Mohr, Mohammad Kalim Akram, Michael Gunther, Bo Wang, Markus Krimmel, Feng Wang, Georgios Mastrapas, Andreas Koukounas, Nan Wang, and Han Xiao. jina-embeddings-v3: Multilingual embeddings with task lora. *arXiv preprint arXiv:2409.10173*, 2024. [2](#)
- [51] Ryota Tanaka, Kyosuke Nishida, Kosuke Nishida, Taku Hasegawa, Itsumi Saito, and Kuniko Saito. Slidevqa: A dataset for document visual question answering on multiple images. *arXiv preprint arXiv:2301.04883*, 2023. [2](#)
- [52] Mitchell Keren Taraday, Shahaf Wagner, and Chaim Baskin. Efficient discriminative joint encoders for large scale vision-language reranking. *arXiv preprint arXiv:2510.06820*, 2025. [2](#)
- [53] Gemma Team, Aishwarya Kamath, Johan Ferret, Shreya Pathak, Nino Vieillard, Ramona Merhej, Sarah Perrin, Tatiana Matejovicova, Alexandre Ramé, Morgane Rivière, et al. Gemma 3 technical report. *arXiv preprint arXiv:2503.19786*, 2025. [4](#), [6](#), [8](#), [12](#), [16](#)
- [54] Nomic Team. Nomic embed multimodal: Interleaved text, image, and screenshots for visual document retrieval, 2025. [6](#)
- [55] NLLB Team, Marta R Costa-jussà, James Cross, Çağrı Çöltekin, Maha Elbayad, Kenneth Heafield, Kevin Heffernan, Elahe Kalbassi, Janice Lam, Daniel Licht, Jean Mailhard, et al. No language left behind: Scaling human-centered machine translation. *arXiv preprint arXiv:2207.04672*, 2022. [3](#)
- [56] Paul Teietche, Quentin Macé, Max Conti, António Loison, Gautier Viaud, Pierre Colombo, and Manuel Faysse. Modernvbert: Towards smaller visual document retrievers. *arXiv preprint arXiv:2510.01149*, 2025. [2](#)
- [57] Nandan Thakur, Nils Reimers, Andreas Rücklé, Abhishek Srivastava, and Iryna Gurevych. Beir: A heterogeneous benchmark for zero-shot evaluation of information retrieval models. In *Conference on Neural Information Processing Systems (NeurIPS) Datasets and Benchmarks Track*, 2021. [2](#), [3](#)
- [58] Omri Uzan, Asaf Yehudai, Roi Pony, Eyal Shnarch, and Ariel Gera. Guided query refinement: Multimodal hybrid retrieval with test-time optimization. *arXiv preprint arXiv:2510.05038*, 2025. [2](#)
- [59] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024. [4](#)
- [60] Navve Wasserman, Roi Pony, Omer Naparstek, Adi Raz Goldfarb, Eli Schwartz, Udi Barzelay, and Leonid Karlinsky. Real-mm-rag: A real-world multi-modal retrieval benchmark. *arXiv preprint arXiv:2502.12342*, 2025. [2](#)
- [61] Cong Wei, Yang Chen, Haonan Chen, Hexiang Hu, Ge Zhang, Jie Fu, Alan Ritter, and Wenhui Chen. Uniir: Training and benchmarking universal multimodal information retrievers. *arXiv preprint arXiv:2311.17136*, 2023. [2](#)
- [62] Zhiyuan Xu, Qinghao Zhang, Xiaochen Wang, Liang Zhao, Beihan Wen, Zhaohui Dai, and Hao Tang. Doclayout-yolo: Enhancing document layout analysis through diverse synthetic data and global-to-local adaptive perception. *arXiv preprint arXiv:2410.12628*, 2024. [3](#)
- [63] Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Jianzong Ran, Yukun Yan, Zhenghao Liu, Shuo Wang, Xu Han, Zhiyuan Liu, and Maosong Sun. Visrag: Vision-based retrieval-augmented generation on multi-modality documents. In *International Conference on Learning Representations (ICLR)*, 2025. [arXiv:2410.10594](#). [2](#)
- [64] Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Liu, and Yicheng Chen. Generative multi-modal entity representation learning. *arXiv preprint arXiv:2409.14327*, 2024. [2](#), [6](#)



## Appendix

### A. Overview

This appendix presents a comprehensive analysis of the ablation studies conducted to develop the M3DR (Multilingual Multimodal Dense Retrieval) framework. We initially ran small ablations for each method with a dataset of 45k images across 6 languages to determine if it made sense to run larger scale experiments. Based on these preliminary results, we progressively scaled up to 22 languages with approximately 250k image text pairs.

The ablations cover critical design decisions including:

- Base model selection
- Loss function comparisons
- Pooling strategies
- Matryoshka embedding dimensions
- Model merging techniques
- Dense vs. late interaction architectures
- Cross-lingual scaling behavior

All experiments were evaluated on three benchmarks:

- **ViDoRe v2**: Document retrieval across 4 languages (English, French, German, Spanish)
- **NayanaIR-CrossBench**: Cross-lingual retrieval across 20 languages
- **NayanaIR-Bench-Monolingual**: Monolingual retrieval evaluation (22 datasets, 1 for each language)

**Metrics Reported**: NDCG@5 (N@5), Recall@10 (R@10), MAP@10 (M@10), MRR@10 (MRR) – consistent with main paper results.

#### Model Naming Conventions

##### Architecture Types:

- **SV (Single-Vector)**: Models producing one dense embedding per image/query for efficient similarity search
- **MV (Multi-Vector)**: ColBERT-style late interaction models producing multiple embeddings per image/query with MaxSim matching

##### Base Models:

- **Gemma3**: Our models built on Gemma 3 4B [53] backbone (finetuned versions)
- **colpali/colqwen/etc.**: External baseline models (original names preserved)

### B. Base Model Selection

This ablation study helped us decide which base model to use for our framework. We evaluated existing pre-trained models using two different retrieval architectures:

- **Col (Late Interaction)**: Multi-vector ColBERT-style models that produce multiple embeddings per image/query and use late interaction (MaxSim) for retrieval (e.g., ColPali, ColQwen, ColSmol)
- **Dense**: Single-vector embedding models that produce one dense embedding per image/query for efficient vector search (e.g., Gemma3, GME-Qwen, Colnomic)

While existing Col models performed well on ViDoRe v2, they were significantly behind on the NayanaIR-CrossBench, indicating poor cross-lingual generalization. We also tested an initial dense embedding model based on Gemma3 4B (Gemma3-InBatch) trained on the ColPali training set with in-batch negative loss.

Figure 5 visualizes the stark trade-off between ViDoRe and cross-lingual performance across baseline models. The scatter plot reveals a concerning pattern: models that achieve top performance on English-centric ViDoRe v2 (such as colqwen2.5-v0.2 at 59.20% NDCG@5) experience catastrophic failure on cross-lingual tasks, plummeting to just 14.26% NDCG@5, a 45%-point drop. This visualization demonstrates that despite using single-vector dense embeddings, our initial Gemma3-InBatch model shows promise due to Gemma3 4B’s robust multilingual vocabulary and pretraining. The figure clearly illustrates why we selected Gemma3 4B as our foundation: its multilingual capacity enables cross-lingual transfer that English-centric ColPali and ColQwen2 architectures cannot match, even when these models dominate on document retrieval benchmarks.

Table 3. **Base Model Selection: Performance Across Benchmarks.** Comparison reveals catastrophic cross-lingual failure of models pretrained primarily on English data. SV = Single-Vector; MV = Multi-Vector.

Model	Type	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10	Mono N@5	Mono R@10
colqwen2.5-v0.2	MV	59.20	66.44	14.26	15.98	45.27	51.30
ColQwen2-6langs	MV	58.63	67.62	37.15	39.85	N/A	N/A
gme-Qwen2-VL-2B	SV	57.39	63.01	23.53	30.83	44.36	52.51
ColPali-6langs	MV	56.41	64.94	42.22	50.05	53.38	62.34
colnomic-embed-3b	SV	55.55	63.26	31.47	31.96	53.37	60.28
colqwen2-v1.0	MV	54.54	64.04	4.97	6.47	41.27	46.58
colpali-v1.3	MV	53.85	62.69	28.45	34.71	41.03	48.41
Gemma3-InBatch	SV	52.44	60.34	20.65	23.14	37.95	46.91
colpali-v1.2	MV	50.55	59.07	22.41	23.68	40.21	47.45
colSmol-500M	MV	43.49	54.41	0.00	0.00	22.41	26.26

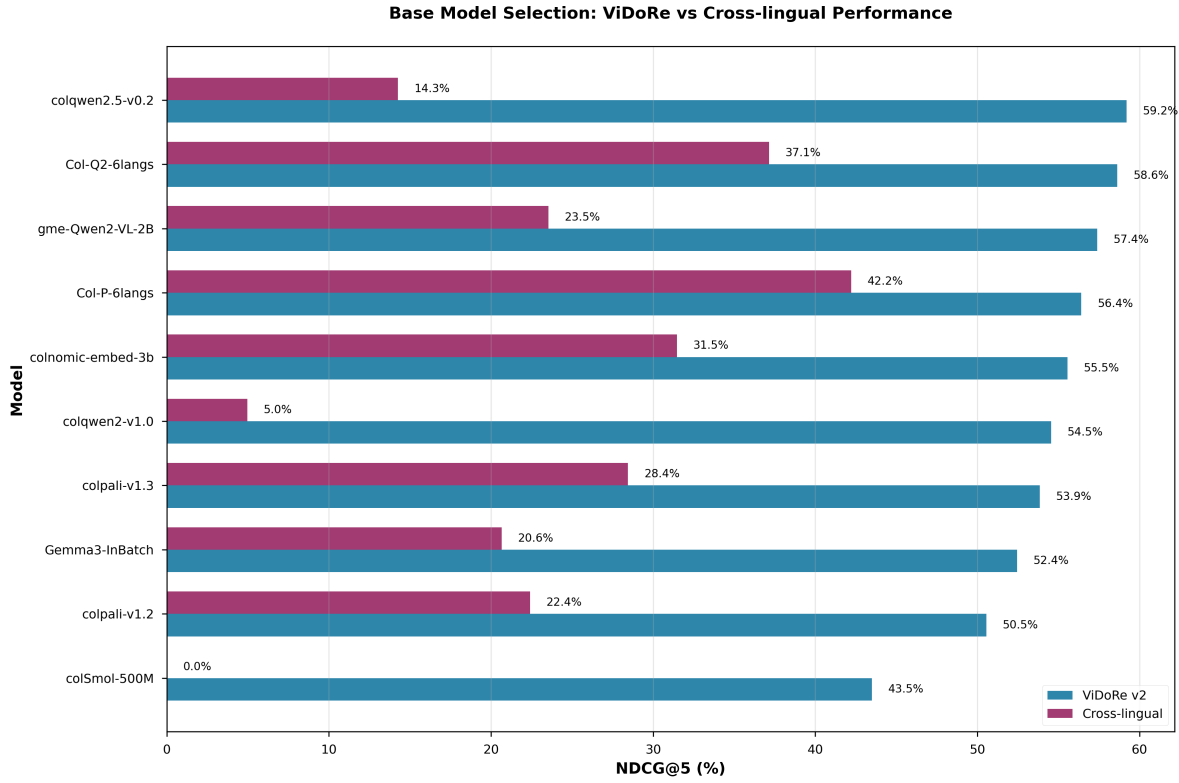


Figure 5. **Base Model Comparison: ViDoRe vs Cross-lingual NDCG@5 for all baseline models.** Models achieving high ViDoRe performance (English-dominated) often fail catastrophically on cross-lingual tasks.

## C. Preliminary Ablations on 6 Languages

Before committing substantial compute to full-scale experiments, we conducted targeted ablations on a smaller 6-language subset (Hindi, Kannada, Tamil, Telugu, Chinese, Japanese) comprising  $\sim 45k$  image-text pairs.

### C.1. Loss Function Ablations

We compared different loss functions for both dense and Col models to determine the optimal training objective. For dense models, we tested BiEncoderLoss, BiNegativeCELoss, BiPairwiseCELoss, and BiPairwiseNegativeCELoss. For Col models, we compared ColBERT loss against ColBERT with pairwise loss.

### C.1.1. Dense Model Loss Functions

We evaluated BiEncoderLoss (InBatch) as the baseline and compared it against hard negative mining at different training checkpoints to understand how hard negative mining affects model convergence over time.

Table 4. Dense Model Loss Function Ablation on ViDoRe v2.

Model	Loss Type	NDCG@5	Recall@10	MAP@10	MRR@10
Gemma3-InBatch (ckpt-2000)	BiEncoderLoss	49.31	58.14	38.78	59.03
Gemma3-HardNeg (ckpt-750)	Hard Negative	50.48	60.22	39.72	61.24
Gemma3-HardNeg (ckpt-1500)	Hard Negative	49.23	55.72	38.26	60.54
Gemma3-HardNeg (ckpt-1694)	Hard Negative	50.12	62.03	39.26	59.86
Gemma3-HardNeg (ckpt-1950)	Hard Negative	49.31	58.80	38.49	59.91
Gemma3-HardNeg (ckpt-2145)	Hard Negative	49.27	56.84	39.11	60.49
Gemma3-HardNeg (ckpt-2300)	Hard Negative	46.73	59.22	37.68	56.25

Figure 6 traces the evolution of four key metrics (NDCG@5, Recall@10, MAP@10, MRR@10) across training checkpoints when using hard negative mining. The line chart reveals erratic behavior: NDCG@5 fluctuates between 46.73% and 50.48% without clear convergence, while Recall@10 peaks at checkpoint 1694 (62.03%) before declining. This volatility contrasts sharply with the stable convergence of standard BiEncoderLoss, suggesting that hard negative mining introduces training instability in our multilingual multimodal setting. The visualization demonstrates that while hard negative mining provided marginal peak performance gains (1-3 points), the instability and careful checkpoint selection required make it impractical compared to the simpler in-batch negative approach.

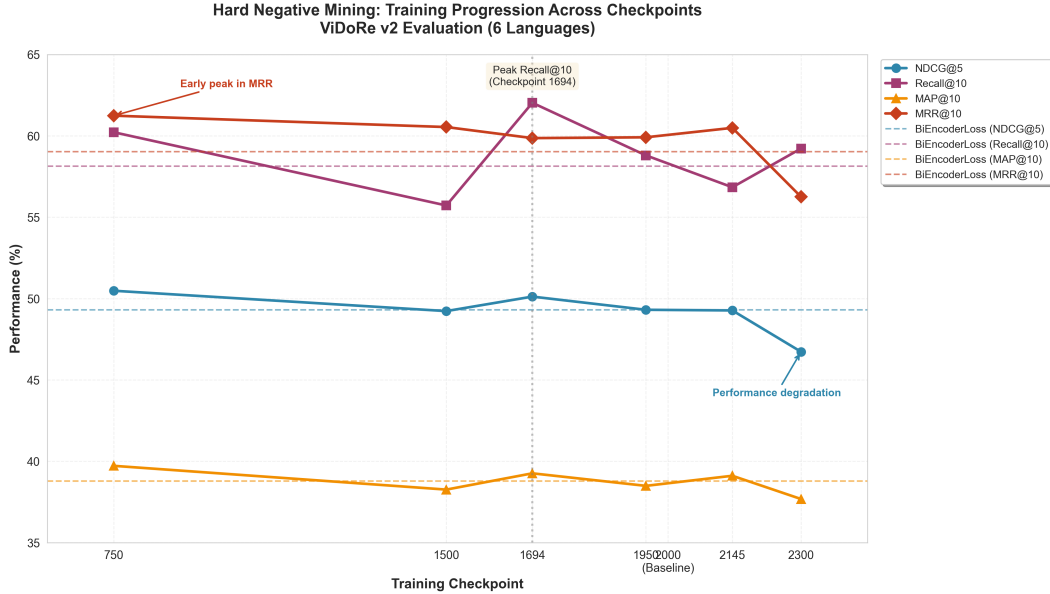


Figure 6. **Hard Negative Mining Training Progression.** Line chart showing metric evolution across training steps (750-2300).

### C.1.2. Col Model Loss Functions

Table 5. Col Model Loss Function Ablation.

Model	Loss	ViDoRe NDCG@5	ViDoRe Recall@10	Cross NDCG@5
ColPali-6langs-ColBERT	ColBERT	56.41	64.94	42.22
ColPali-6langs-Pairwise	ColBERT+Pairwise	39.53	51.98	33.94



The comparison reveals that for single-vector dense models, BiEncoderLoss emerged as the optimal choice, while for multi-vector Col models, standard ColBERT loss significantly outperformed ColBERT with pairwise loss (56.41% vs 39.53% NDCG@5 on ViDoRe), a 17 percentage point degradation. This is primarily because both the ColPali train set and Nayana IR dataset contained positive pairs without explicitly annotated hard negatives, making the diversity provided by in-batch negatives across languages and visual content sufficient for learning discriminative representations. Figure 7 summarizes these findings, illustrating that simpler objectives outperform complex negative sampling in multilingual settings.

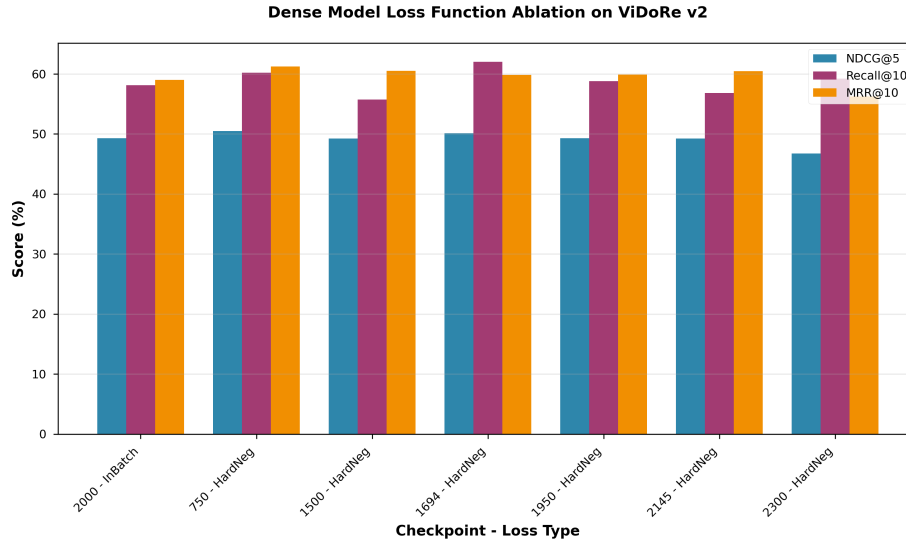


Figure 7. **Dense Model Loss Function Ablation.** Comparing BiEncoderLoss and Hard Negative mining variants.

## C.2. Pooling Strategy Ablations: Last Token vs Mean Pooling

For single-vector dense embedding models, we need to convert the final hidden states into a single embedding vector. Note: This ablation is specific to dense models only, as Col models inherently use all patch embeddings without pooling.

Table 6. **Pooling Strategy Ablation on ViDoRe v2.**

Model	Pooling	NDCG@5	Recall@10	MRR@10
Gemma3-LastToken	Last Token	<b>49.31</b>	<b>58.14</b>	<b>59.03</b>
Gemma3-MeanPool	Mean Pooling	35.85	45.04	45.07

Figure 8 demonstrates the dramatic superiority of last token pooling over mean pooling, with performance differences of 13.5 percentage points in NDCG@5 (49.31% vs 35.85%). This suggests that decoder-only VLMs like Gemma encode critical summary information in the final token representation, similar to how GPT-style models aggregate sequence information. Additionally, last token pooling proves computationally cheaper during inference as it requires processing only a single token representation rather than averaging across all tokens, avoiding the dilution of salient information through averaging. This became our default pooling strategy for all dense embedding models.

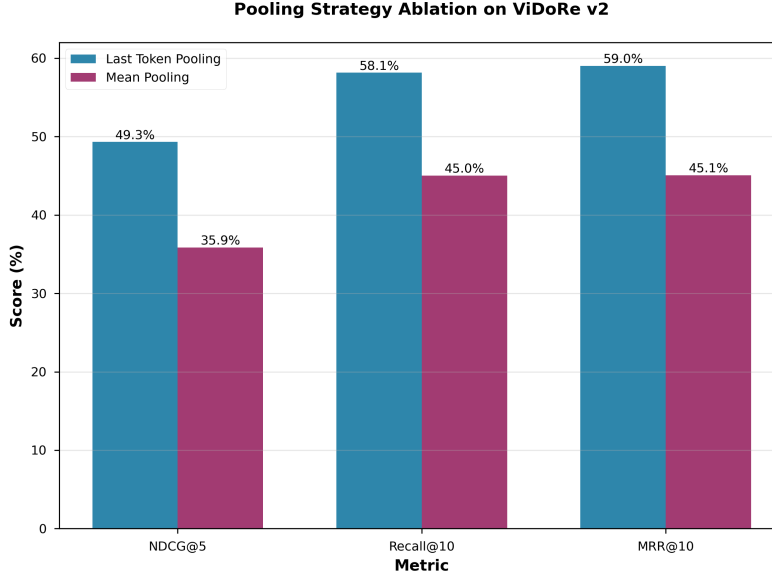


Figure 8. **Pooling Strategy Ablation.** Last Token vs Mean Pooling performance comparison.

### C.3. OCR Model Ablations

We hypothesized that existing VLMs pre-trained for OCR tasks (such as Docling and Granite OCR) might generalize well to text-heavy document retrieval. We finetuned these models for both dense and Col-style retrieval to test this hypothesis.

Table 7. **OCR Model Transfer Ablation on ViDoRe v2.** SV = Single-Vector; MV = Multi-Vector.

Model	Type	NDCG@5	Recall@10	MRR@10
Granite-Docling-258M	OCR baseline	0.97	1.16	2.40
Docling-Finetuned	SV	4.17	6.92	5.72
Gemma3-OCRInit	MV	1.37	1.05	1.94
Gemma3-Pretrained	VLM (no FT)	4.48	3.43	8.52

OCR-pretrained models failed to generalize to retrieval tasks even after finetuning. The pretrained Granite-Docling achieved near-zero retrieval performance (0.97% NDCG@5), and finetuning only recovered to 4.17%, still catastrophically poor. Remarkably, the base Gemma 3 4B [53] model without any finetuning matched this performance. This negative result demonstrates that OCR pretraining optimizes for character-level text recognition and spatial layout understanding, creating representations specialized for transcription rather than semantic retrieval. These representations are orthogonal to the semantic similarity space required for retrieval, validating our choice of general-purpose VLMs over task-specific OCR models.

### C.4. Training on 6 Languages: Initial Results

We finetuned existing multi-vector Col models (ColQwen2 and ColPali) and trained single-vector dense Gemma3 [53] models on our 6-language dataset to evaluate cross-lingual generalization capabilities.

Table 8. **6-Language Training Results.** SV = Single-Vector; MV = Multi-Vector.

Model	Type	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10	Mono N@5	Mono R@10
ColQwen2-6langs	MV	58.63	67.62	37.15	39.85	N/A	N/A
ColPali-6langs	MV	56.41	64.94	42.22	50.05	53.38	62.34
Gemma3-6langs	SV	49.08	59.76	60.39	72.94	62.48	74.70
Gemma3-6langs-v2	SV	49.06	60.44	59.62	76.47	62.42	73.88

Figure 9 presents a pivotal finding in our research journey. While multi-vector ColQwen2 and ColPali variants maintained

their ViDoRe advantage (56-59% vs 49% NDCG@5), the visualization reveals that dense Gemma3 models achieved dramatically superior cross-lingual performance: 60.39% vs 42.22% for the best late interaction model, an 18 % point gap. This massive difference validates our base model hypothesis: Gemma’s multilingual pretraining enabled cross-lingual transfer that English-centric ColPali and ColQwen2 foundations could not match, even after multilingual finetuning. The figure demonstrates the clear superiority of Gemma3-based models across cross-lingual and monolingual benchmarks, convincing us to pursue both single-vector Gemma3 as the primary model and multi-vector ColGemma3 for a comprehensive architectural comparison.

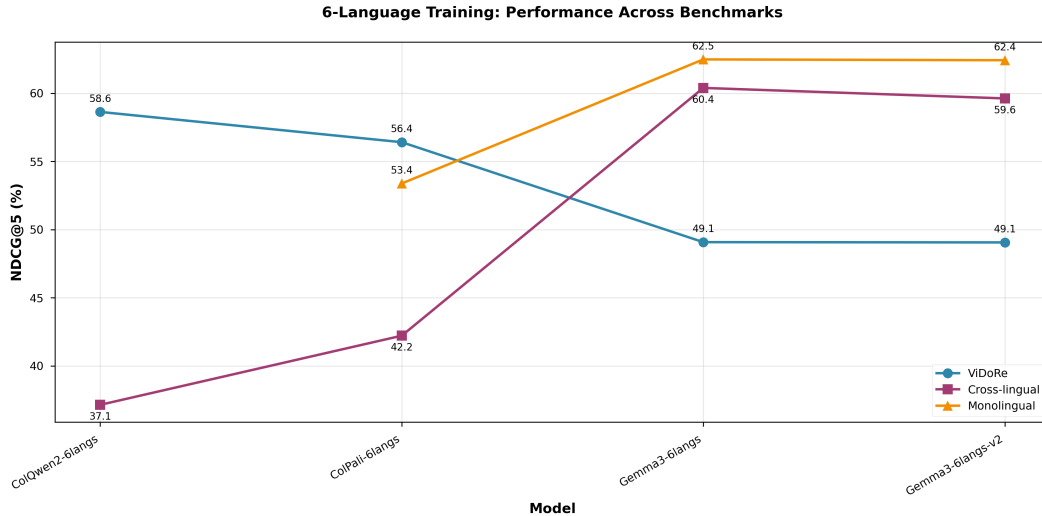


Figure 9. **6-Language Training: Performance across benchmarks.** Dense Gemma3 models excel at cross-lingual retrieval.

## D. Matryoshka Embedding Ablations

Matryoshka Representation Learning allows single-vector dense models to produce embeddings that maintain performance across different dimensionalities. This technique is specific to dense embedding models and does not apply to multi-vector Col models.

### D.1. Matryoshka on 6 Languages

Table 9. **Matryoshka Ablation on 6 Languages.**

Model	Dimensions	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10
Gemma3-Matryoshka-6langs	2056	45.81	60.32	57.22	71.32
Gemma3-6langs (baseline)	2056	49.08	59.76	60.39	72.94

The 6-language Matryoshka model showed 3-point degradation on ViDoRe and Cross benchmarks, with slight improvement on monolingual tasks. While concerning, we hypothesized that scaling to more languages might provide sufficient training signal to overcome this gap.

### D.2. Matryoshka on 22 Languages (Final Scale)

These models were trained on the ColPali train set and then finetuned with Matryoshka loss [768, 1536, 2056] and BiNegative CE loss on the Nayana IR dataset.

Figure 10 illustrates the graceful degradation of performance as embedding dimensionality decreases. The visualization shows that 768-dimensional embeddings retained 94.6% of full performance on cross-lingual tasks (73.15% vs 77.31% NDCG@5), a mere 4.2-point drop, while providing 70% storage reduction. The 1536-dimensional embeddings achieved 99.4% of full performance (76.87% vs 77.31%), effectively trading 1 NDCG@5 point for 40% storage savings. This flexibility is critical for production systems, where index size directly impacts cost and latency. Both the primary and v2 training runs

Table 10. **Primary Matryoshka Models on 22 Languages.**

Model	Dimensions	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10	Mono N@5	Mono R@10
Gemma3-Matryoshka (2056)	2056	49.88	62.00	77.31	88.38	74.10	85.03
Gemma3-Matryoshka (1536)	1536	48.21	60.47	76.87	86.91	73.25	84.07
Gemma3-Matryoshka (768)	768	45.10	56.52	73.15	79.12	70.77	82.42
Gemma3-Matryoshka-v2 (2056)	2056	45.63	57.33	72.60	83.09	73.77	84.29
Gemma3-Matryoshka-v2 (1536)	1536	45.00	56.42	72.16	82.50	73.14	84.02
Gemma3-Matryoshka-v2 (768)	768	42.86	53.91	70.45	84.12	72.05	82.68

showed consistent behavior across dimensions, confirming the robustness of Matryoshka training and enabling deployment teams to choose the optimal performance-storage trade-off for their specific use case.

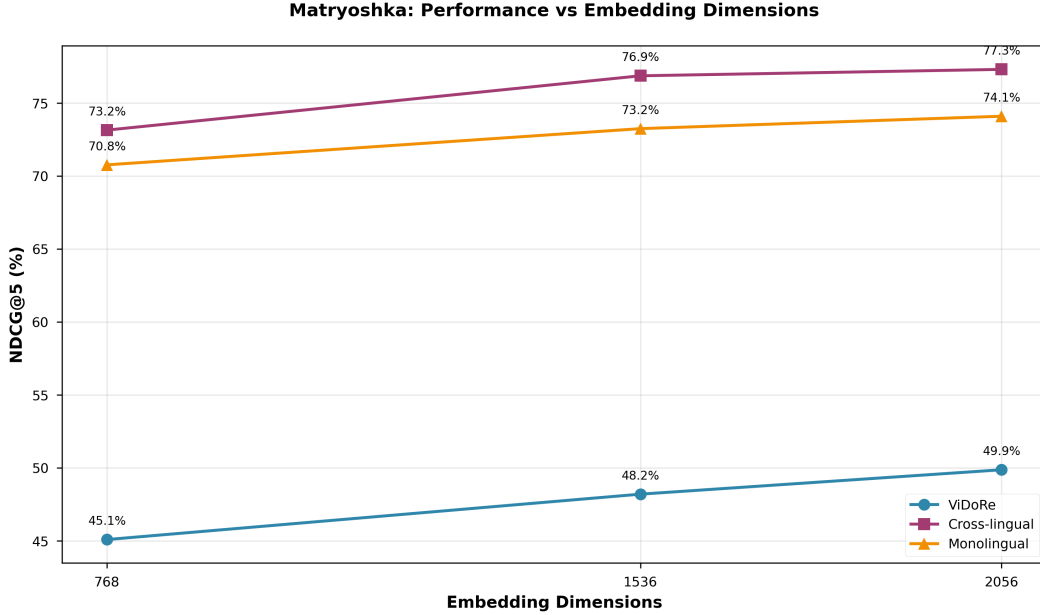


Figure 10. **Matryoshka Embedding Dimensions.** NDCG@5 performance across 768, 1536, and 2056 dimensions showing graceful degradation.

## E. Model Merging Strategies

To create a more balanced model that performs well across both ViDoRe and Nayana cross-lingual benchmarks, we explored model merging techniques. We merged two complementary checkpoints: Parent Model A (Gemma3-CrossLingual) with excellent cross-lingual performance (77.31% NDCG@5) but moderate ViDoRe performance (49.88%), and Parent Model B (Gemma3-ViDoRe) with top ViDoRe performance (55.40%) but decent cross-lingual performance (71.57%).

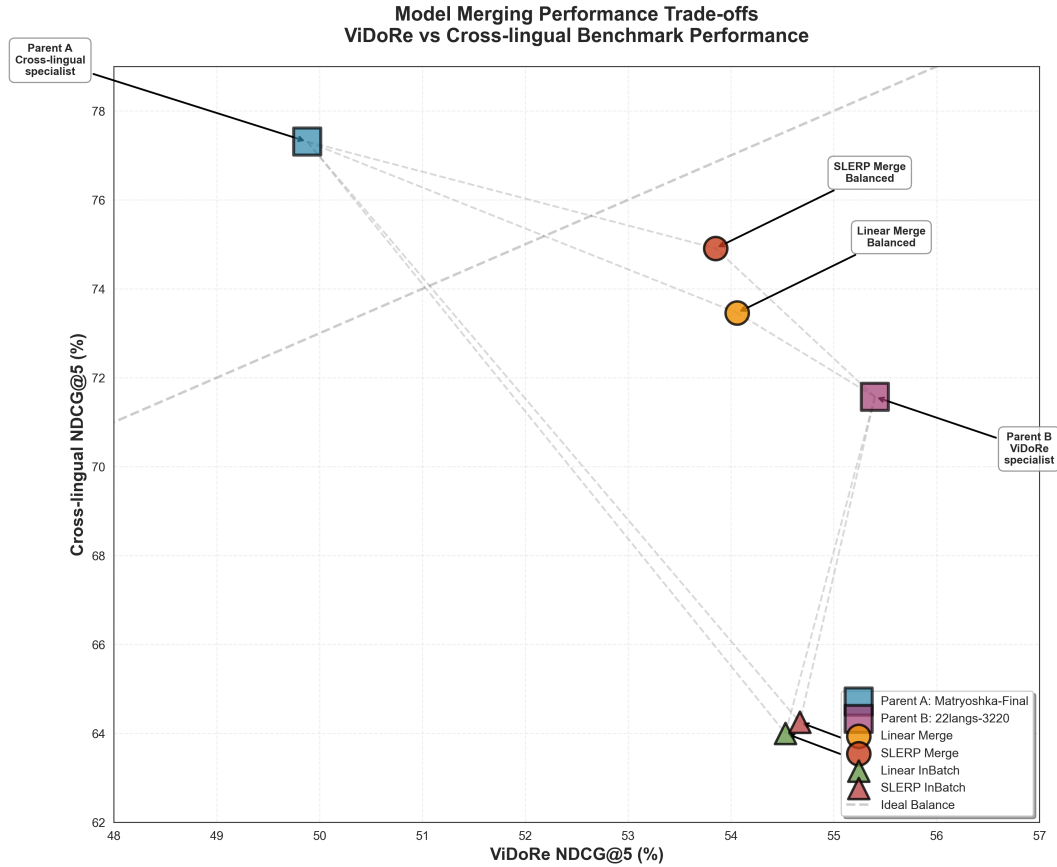
We explored two interpolation methods: Linear merging (weighted average  $\alpha \cdot \theta_1 + (1 - \alpha) \cdot \theta_2$ ) and SLERP (Spherical Linear Interpolation, interpolating along the geodesic on the hypersphere).

Figure 11 visualizes the performance trade-offs through a scatter plot of ViDoRe NDCG@5 versus Cross-lingual NDCG@5. The parent models occupy opposite corners of the performance space: Parent A (Gemma3-CrossLingual) in the upper-left with high cross-lingual but moderate ViDoRe performance, and Parent B (Gemma3-ViDoRe) in the lower-right with high ViDoRe but moderate cross-lingual performance. The merged models successfully populate the balanced middle ground, with Gemma3-SLERP achieving 53.85% ViDoRe and 74.91% cross-lingual, and Gemma3-LinearMerge reaching 54.06% ViDoRe and 73.46% cross-lingual. The visualization demonstrates that model merging successfully balances objectives without additional training: merged models achieve 95-98% of the best parent’s performance on each benchmark, with SLERP better preserving cross-lingual capabilities and Linear better retaining ViDoRe accuracy.



Table 11. Model Merging Results.

Model	Method	Dimensions	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10
Gemma3-LinearMerge (2056)	Linear	2056	54.06	63.43	73.46	89.27
Gemma3-LinearMerge (1536)	Linear	1536	52.04	62.01	74.77	88.53
Gemma3-LinearMerge (768)	Linear	768	48.29	58.33	70.84	80.00
Gemma3-SLERP (2056)	SLERP	2056	53.85	63.23	74.91	89.27
Gemma3-SLERP (1536)	SLERP	1536	51.03	61.93	74.40	88.53
Gemma3-SLERP (768)	SLERP	768	47.42	58.38	70.86	78.53

Figure 11. **Model Merging Performance Trade-offs.** Parent models at opposite corners, merged models achieving balanced intermediate performance.

## F. Scaling to 22 Languages

### F.1. Dense Models at Full Scale

Table 12. Dense Model Scaling from 6 to 22 Languages. SV = Single-Vector.

Model	Type	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10
Gemma3-22langs	SV	55.40	63.75	71.57	87.06
Gemma3-6langs	SV	49.08	59.76	60.39	72.94

Figure 12 demonstrates the benefits of linguistic diversity and training data scale. Scaling from 6 to 22 languages yielded substantial improvements: +6.3 points on ViDoRe (55.40% vs 49.08%), +11.2 points on cross-lingual tasks (71.57% vs 60.39%), and +11.3 points on monolingual tasks. The larger relative gains on cross-lingual and monolingual benchmarks (18%) versus ViDoRe (13%) suggest that linguistic diversity particularly benefits non-English retrieval, validating the scaling hypothesis.

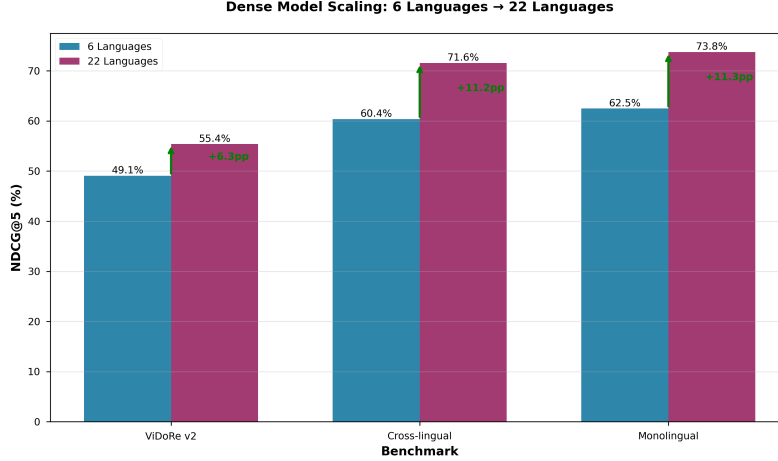


Figure 12. Language Scaling: Performance improvements when scaling from 6 to 22 languages.

## F.2. Multi-Vector Col Models at Full Scale

Table 13. Col Model Comparison at 22 Languages. MV = Multi-Vector.

Model	Type	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10	Mono N@5
ColGemma3-Merged-22langs	MV	55.31	67.29	63.75	69.95	67.02
ColGemma3-22langs	MV	52.91	62.13	64.74	72.35	71.21
ColGemma3-ColPaliOnly	MV	54.37	63.60	14.15	14.85	N/A
ColQwen2-22langs	MV	53.07	62.72	42.86	52.01	59.08
ColQwen2-22langs-Scaled	MV	N/A	N/A	51.08	60.45	64.31
ColPali-22langs	MV	50.37	61.83	42.64	53.63	53.38
<i>For comparison (6-language baselines):</i>						
ColQwen2-6langs	MV	58.63	67.62	37.15	39.85	N/A
ColPali-6langs	MV	56.41	64.94	42.22	50.05	53.38

Multi-vector ColGemma3 trained from scratch showed significantly better cross-lingual generalization (63.75% NDCG@5) compared to finetuned ColPali (42.64%) and ColQwen2 (42.86%). The final model followed a three-stage process: Stage 1 trained on ColPali data (54.37% ViDoRe, 14.15% cross-lingual), Stage 2 finetuned on Nayana IR 22 languages (52.91% ViDoRe, 64.74% cross-lingual), and Stage 3 merged the checkpoints (55.31% ViDoRe, 63.75% cross-lingual), creating the most balanced late interaction model. Interestingly, the 6-language ColQwen2 model outperformed the 22-language version on ViDoRe (58.63% vs 53.07%), suggesting potential overfitting or training instability.

## G. Dense vs Col Architecture Comparison

We trained models using both retrieval architectures to evaluate trade-offs: single-vector models enabling fast vector similarity search, and multi-vector late interaction models using multiple embeddings with MaxSim matching.

The architectural comparison reveals complementary strengths. Single-vector models achieved 73-77% NDCG@5 on cross-lingual tasks versus 64% for late interaction, a 10-13 point advantage. Late interaction models showed marginal superiority on ViDoRe (55.31% vs 53-55%), a 0-2 point difference. Beyond retrieval accuracy, single-vector embeddings dominate on practical deployment metrics: 8-10× smaller storage footprint enabling larger indices and lower costs; 10×

Table 14. **Dense vs. Late Interaction Architecture Performance.** SV = Single-Vector; MV = Multi-Vector.

Type	Model	ViDoRe N@5	ViDoRe R@10	Cross N@5	Cross R@10	Mono N@5
SV	Gemma3-LinearMerge	54.06	63.43	73.46	89.27	N/A
SV	Gemma3-SLERP	53.85	63.23	74.91	89.27	N/A
SV	Gemma3-Matryoshka	49.88	62.00	77.31	88.38	74.10
MV	ColGemma3-Merged	55.31	67.29	63.75	69.95	67.02
MV	ColGemma3-22langs	52.91	62.13	64.74	72.35	71.21

Table 15. **Dense vs. Late Interaction: Deployment Trade-offs.**

Aspect	Dense (Single-Vector)	Col (Multi-Vector)
Storage	2056 dims $\approx$ 8KB	128 tokens $\times$ 128 dims $\approx$ 64KB
Storage advantage	<b>8-10<math>\times</math> smaller</b>	-
Retrieval speed	500-1000 QPS	50-100 QPS
Speed advantage	<b>10<math>\times</math> faster</b>	-
Cross-lingual perf	<b>73-77% NDCG@5</b>	64% NDCG@5
ViDoRe perf	50-54% NDCG@5	<b>53-55% (marginal)</b>
Interpretability	Low (black-box vector)	<b>High (attention maps)</b>
Long-tail queries	Moderate	<b>Better (multi-vector)</b>
Deployment	<b>Simple (standard database)</b>	Complex (MaxSim)

faster retrieval (standard vector similarity vs expensive MaxSim across 128 token vectors); and compatibility with standard vector databases (Faiss, Milvus, Pinecone) simplifying infrastructure. Late interaction’s advantages lie in interpretability through token-level attention heatmaps and potentially better handling of long-tail queries through fine-grained multi-vector matching. Both architectures have merit depending on deployment priorities: single-vector for cost-efficiency and scale, late interaction for interpretability and explainability.

## H. Cross-Lingual Embedding Convergence: PCA Analysis

This section provides visual evidence that the Gemma3 model progressively learns cross-lingual alignment during training on the Nayana IR dataset. We visualize how embeddings for semantically equivalent content across different languages evolve from language-separated clusters to semantically-aligned representations.

### H.1. Methodology

We extract model checkpoints at regular intervals (500, 1500, 2500, 3500, 4500, 5066 steps), generate embeddings for document images paired with queries in multiple languages, apply PCA to project 2560-dimensional embeddings to 2D, and evaluate convergence behavior across different multilingual configurations (2, 6, and 15 languages).

### H.2. 2-Language Alignment: Kannada $\leftrightarrow$ English

Figure 13 visualizes the simplest case of cross-lingual alignment between Kannada queries and English documents across six training checkpoints. At checkpoint 500 (leftmost panel), the PCA projection shows clear separation between blue triangles (Kannada queries) and orange triangles (document embeddings), forming distinct clusters in opposite regions of the embedding space. As training progresses to checkpoint 2500 (middle panels), the visualization reveals boundaries beginning to blur as query and document embeddings start overlapping in the center of the plot. By checkpoint 5066 (rightmost panel), the separation has largely dissolved: blue and orange triangles intermingle throughout the embedding space, indicating the model learned to produce similar embeddings for semantically equivalent content regardless of language. This progressive convergence from language-separated to semantically-aligned representations demonstrates the model’s acquisition of cross-lingual understanding.

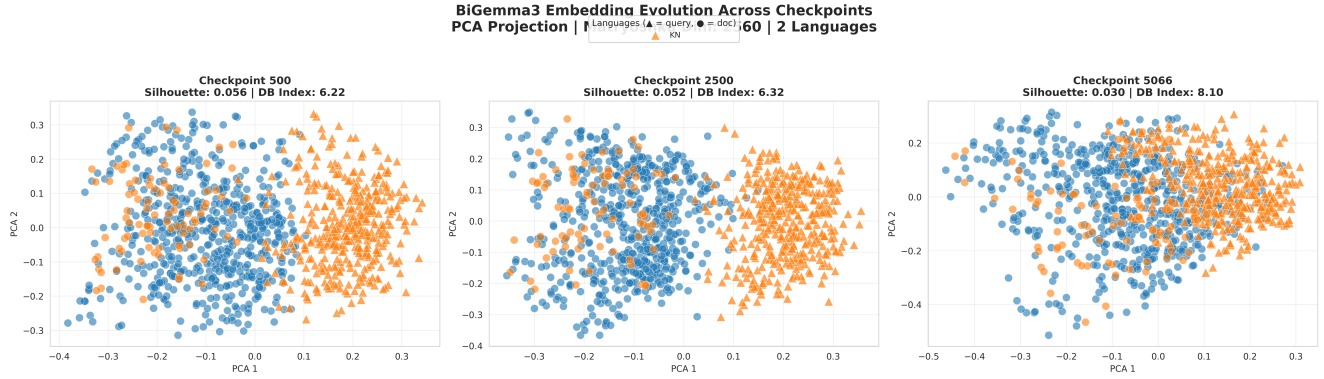


Figure 13. Gemma3 Embedding Evolution: 2 Languages (KN Query ↔ Doc).

### H.3. 6-Language Cross-Lingual Alignment

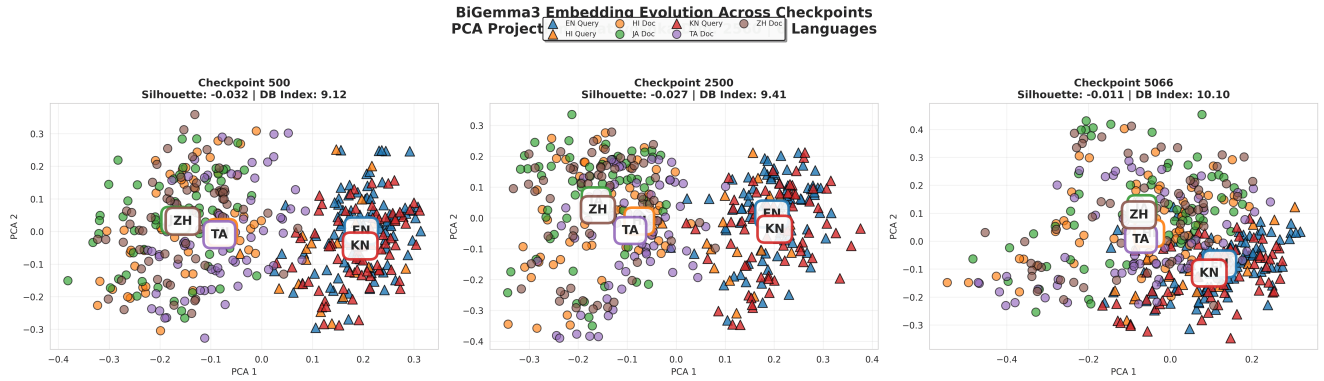


Figure 14. Gemma3 Embedding Evolution: 6 Languages.

Scaling to six languages (English, Hindi, Chinese, Tamil, Japanese, Kannada) reveals richer dynamics of cross-lingual convergence. Figure 14 shows six distinct color-coded language clusters clearly visible at checkpoint 500, with each language occupying its own region in the 2D PCA space. Queries (triangles) and documents (circles) form separate but overlapping zones. As training advances through checkpoints 1500 and 2500, the language-specific boundaries progressively dissolve i.e., the distinct color clusters merge toward the center of the plot. By checkpoint 5066, the visualization displays near-complete cross-lingual mixing: embeddings from different languages for the same semantic content cluster together regardless of script, with Hindi (Devanagari), Chinese (Hanzi), Tamil, Kannada, Japanese (Kanji/Kana), and English (Latin) all converging to similar regions. The transition from distinct language-separated clusters to a unified semantic space provides direct visual evidence that the model learns to align representations across diverse writing systems.

Figure 15 confirms the same convergence pattern using Matryoshka embeddings with 2560 dimensions. The visualization mirrors the behavior observed in the standard model: early checkpoints show language-clustered distributions, which progressively merge into semantically-organized representations. This demonstrates that Matryoshka’s multi-scale training objective preserves (and potentially enhances) cross-lingual alignment quality, validating that flexible dimensionality does not compromise the model’s ability to learn language-agnostic representations.

### H.4. 15-Language Scaling

Figure 16 tests the limits of cross-lingual alignment with 15 diverse languages spanning multiple language families and writing systems: Arabic, Bengali, German, Spanish, French, Hindi, Kannada, Russian, Tamil, Telugu, Thai, and Chinese. At checkpoint 500, the PCA projection shows significant dispersion across the entire embedding space, with language families forming loose clusters (e.g., Indic languages grouping together, Romance languages clustering separately). The visualization tracks progressive convergence: checkpoint 1500 shows early signals of alignment with document embeddings beginning to



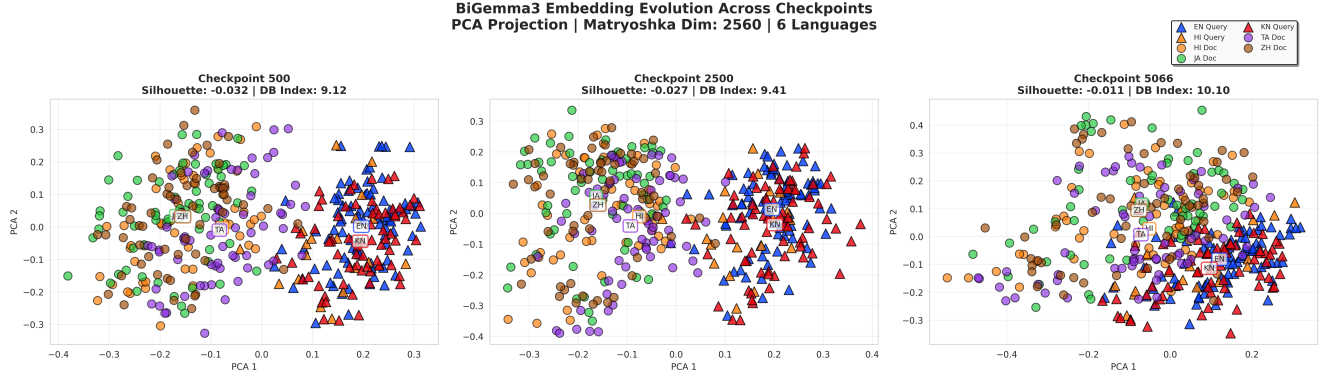


Figure 15. Gemma3 Embedding Evolution: 6 Languages (Matryoshka Dim: 2560).

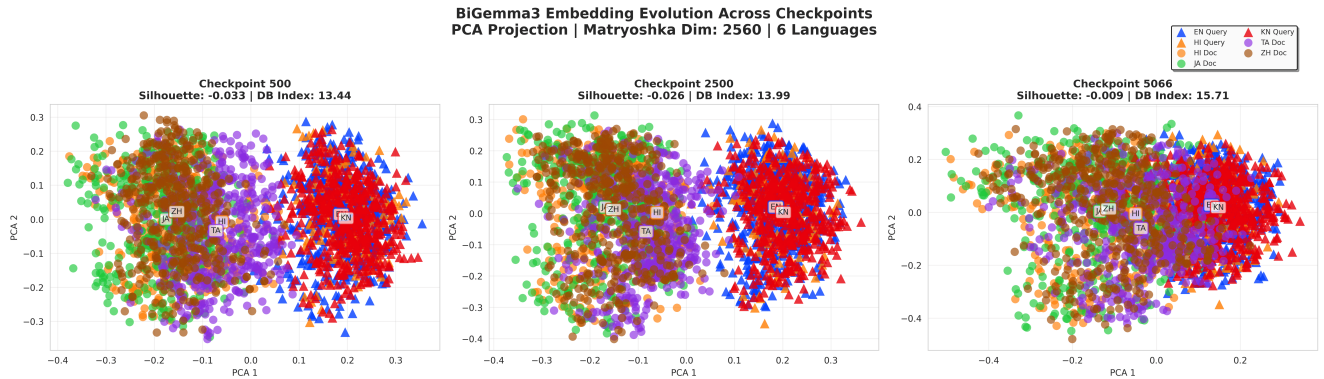


Figure 16. Gemma3 Embedding Evolution: 15 Languages.

form denser central clusters; checkpoint 2500 reveals clear movement toward the center with query embeddings from different languages increasingly overlapping; and checkpoint 5066 demonstrates strong cross-lingual mixing across all 15 languages despite the complexity. The model successfully learns to align semantically equivalent content even when confronted with high linguistic diversity, demonstrating the robustness of the M3DR framework to scale beyond language pairs to truly multilingual scenarios.

### H.5. Hindi ↔ Kannada Detailed Analysis

Figure 17 provides focused analysis of two Indic languages with distinct scripts to demonstrate fine-grained cross-lingual alignment. At checkpoint 500, Hindi documents (green circles) and Kannada documents (orange circles) are clearly separated in the PCA space, occupying opposite regions of the plot. The progression through checkpoints 1500 and 2500 shows gradual overlap beginning between the clusters, with the color boundaries becoming less distinct. By checkpoint 5066, the visualization reveals near-perfect alignment: green and orange circles are thoroughly intermixed throughout the embedding space, demonstrating that the model learned to produce nearly identical embeddings for the same visual content whether queried in Hindi (Devanagari script) or Kannada (Kannada script). This alignment between non-Latin scripts provides compelling evidence that the model develops script-agnostic semantic representations rather than superficial pattern matching.

### H.6. Implications for Cross-Lingual Retrieval

These visualizations collectively demonstrate that multilingual training with the Nayana IR dataset enables models to learn progressive cross-lingual alignment in the embedding space. The convergence happens gradually across training rather than as a sudden phase transition, with the model learning to project different languages into a shared semantic space. Despite using single-vector representations, these embeddings successfully capture cross-lingual semantics, as evidenced by the visual collapse of language-specific clustering into semantically-organized distributions. This validates our core hypothesis: the Nayana IR dataset contains sufficient cross-lingual signal to enable robust multilingual multimodal retrieval across diverse



Figure 17. **Gemma3 Embedding Convergence: Hindi ↔ Kannada.**

scripts and language families.

## I. Col Model Attention Visualization

Col models provide interpretability through token-level attention heatmaps via maximum similarity (MaxSim) maps between query tokens and image patches. We demonstrate that cross-lingual text queries attend to the same visual regions, validating language-agnostic visual grounding.

### I.1. Experimental Setup

We compare models trained solely on English-centric data (ColPali base versions) against those fine-tuned on the multilingual Nayana corpus. The test document contains the scientific term “*Acremonium coephenophialum*” in English, queried in three languages - English, Hindi and Kannada.

Figures 18 through 21 reveal the impact of multilingual fine-tuning on attention consistency. The base models (Figures 18 and 20) show attention heatmaps where English queries correctly focus on the target text with moderate confidence (red intensity), but Hindi and Kannada queries either diffuse attention across the document or fail to activate the correct region, ColGemma3-ColPaliOnly drops from 0.329 similarity in English to just 0.109 in Kannada, indicating near-complete failure to associate the Kannada script with the English visual text. In contrast, the Nayana-tuned models (Figures 19 and 21) demonstrate query consistency: all three language queries activate precisely the same visual regions with comparable confidence levels (red heatmap intensity), proving the models learned robust cross-lingual representations. Remarkably, ColGemma3-Finetuned achieves higher Hindi similarity (0.381) than English (0.335), demonstrating script-agnostic semantic understanding rather than relying on visual script matching.

Table 16 quantifies the dramatic improvement in cross-lingual consistency. ColGemma3-Finetuned achieves remarkable query consistency with only a +0.015 gap between English and the lowest language score, even showing higher Hindi scores (0.381) than English (0.329). This contrasts starkly with ColGemma3-ColPaliOnly which exhibits catastrophic cross-lingual failure with a -0.220 gap and Kannada performance plummeting to 0.109. The Nayana-tuned models demonstrate that multilingual fine-tuning does not degrade English capabilities while dramatically improving non-English query performance, with

ColQwen2 v1.0 - Base - Multilingual Query Consistency



Figure 18. ColQwen2-Base (MV). English Max Sim: 0.335 — Hindi: 0.247 — Kannada: 0.286

ColQwen2 v1.0 - Nayana 200k - Multilingual Query Consistency



Figure 19. ColQwen2-Finetuned (MV) - Ours. English Max Sim: 0.529 — Hindi: 0.436 — Kannada: 0.269

ColGemma - ColPali Train Set - Multilingual Query Consistency



Figure 20. ColGemma3-ColPaliOnly (MV). English Max Sim: 0.329 — Hindi: 0.203 — Kannada: 0.109

all language queries activating identical visual regions which is clear evidence of true cross-lingual semantic understanding rather than script-matching heuristics.



ColGemma - Nayana 200k - Multilingual Query Consistency



Figure 21. ColGemma3-Finetuned (MV) - Ours.  
English Max Sim: 0.335 — Hindi: 0.381 — Kannada: 0.349

Table 16. Cross-Lingual Query Consistency Analysis. MV = Multi-Vector.

Model	English MaxSim	Hindi MaxSim	Kannada MaxSim	Cross-Lingual Gap
ColGemma3-ColPaliOnly	0.329	0.203	0.109	-0.220
ColGemma3-Finetuned	0.335	<b>0.381</b>	<b>0.349</b>	+0.015
ColQwen2-Base	0.335	0.247	0.286	-0.088
ColQwen2-Finetuned	<b>0.529</b>	<b>0.436</b>	0.269	-0.260

## J. Training Configuration and Computational Cost

### J.1. LoRA Setup

Models were finetuned with LoRA (rank 32, alpha 32, dropout 0.1) applied to projection layers (down\_proj, gate\_proj, up\_proj, k\_proj, q\_proj, v\_proj, o\_proj).

### J.2. Hyperparameters

Training for 2 epochs gave the best results. Settings: per device batch size 32, learning rate 2e-4, gradient checkpointing, warmup 100 steps, logging 10 steps, and saving every 500 steps.

### J.3. Compute

Small runs used one A100 80 GB for 2 to 4 hours on roughly 45k pairs. Final multilingual training used 4 to 8 A100 80 GB GPUs with DDP and mixed precision. Two datasets were trained for 2 epochs each which took 6 to 8 hours in total (about 64 GPU hours, roughly 100 to 150 USD).

### J.4. Data

Table 17. Dataset Summary

Dataset	Pairs	Notes
ColPali	~150k	English document retrieval
Nayana IR	~250k	22 languages, cross lingual
<b>Total</b>	<b>~400k</b>	Multilingual retrieval

Nayana IR spans Indic scripts (11), major European languages (6), Asian languages (4), and Arabic. This coverage supports both multilingual and cross lingual document retrieval tasks.