

Llama Nemoretriever Colembed: Top-Performing Text-Image Retrieval Model

Core Contributors

Mengyao Xu, Gabriel Moreira, Ronay Ak, Radek Osmulski, Yauhen Babakhin, Zhiding Yu, Benedikt Schifferer, Even Oldridge
NVIDIA

Contributors¹

Adam Laiacano, Alex Richards, Andrew Tao, Ben Jarmak, Bo Liu, Charles Blackmon-Luca, Derek Whatley, Edward Kim, Fei Yu, Jeremy Dyer, Jeremy Jordan, Joey Conway, John Zedlewski, Julio Perez, Kalpesh Sutaria, Kam Mitchell, Kari Briski, Karan Sapra, Loan Luong, Maximilian Jeblick, Meghana Shrotri, Nave Algarici, Oliver Holworthy, Padmavathy Subramanian, Randy Gelhausen, Salik Siddiqui, Sean Sodha, Shizhe Diao, Sohail Sahi, Steven Baughman, Theo Viel, Tom Balough, Tom O’Brien.
NVIDIA

Abstract

Motivated by the growing demand for retrieval systems that operate across modalities, we introduce llama-nemoretriever-colembed, a unified text-image retrieval model that delivers state-of-the-art performance across multiple benchmarks. We release two model variants, 1B and 3B². The 3B model achieves state of the art performance, scoring NDCG@5 91.0 on ViDoRe V1 and 63.5 on ViDoRe V2, placing first on both leaderboards as of June 27, 2025.

Our approach leverages the NVIDIA Eagle2 Vision-Language model (VLM), modifies its architecture by replacing causal attention with bidirectional attention, and integrates a ColBERT-style late interaction mechanism to enable fine-grained multimodal retrieval in a shared embedding space. While this mechanism delivers superior retrieval accuracy, it introduces trade-offs in storage and efficiency. We provide a comprehensive analysis of these trade-offs. Additionally, we adopt a two-stage training strategy to enhance the model’s retrieval capabilities.

1 Introduction

Retrieval-Augmented Generation (RAG) has become a widely adopted paradigm for enhancing language models with external knowledge, enabling them to retrieve and reason on relevant content from large-scale corpora. Numerous high-performing text retrieval models including NV-Embed [1], NV-Retriever [2], Qwen3-Embedding [3], and e5-mistral [4] have been proposed, and evaluated on benchmarks such as MTEB [5, 6]. While these models and benchmarks focus primarily on text-only retrieval, they often assume clean, well-formatted textual inputs. In contrast, real-world use cases typically involve documents stored in formats like PDFs, PowerPoint slides, or Word documents, requiring preprocessing pipelines to extract textual content. This process often results in the loss of critical visual information for modalities like tables, charts, and infographics.

¹Sorted alphabetically

²We released our models at <https://huggingface.co/nvidia/llama-nemoretriever-colembed-3b-v1> and <https://huggingface.co/nvidia/llama-nemoretriever-colembed-1b-v1>.

To address this, an alternative approach proposed by ColPali [7] converts documents into images, enabling retrieval systems to handle both textual and visual modalities effectively.

Recent Vision-Language models (VLMs) aim to bridge the gap between text and image understanding by learning joint representations across modalities. Models such as Qwen-VL [8], LLaMA-3.1-Nemotron-Nano-VL [9], and NVIDIA’s Eagle2 models [10, 11] have demonstrated strong performance across a range of vision-language tasks by leveraging vision encoders like CLIP [12], SigLIP [13] and C-RADIO [14]. CLIP and SigLIP are trained on image-text pairs using contrastive learning. C-RADIO is trained through multi-teacher distillation. These limitations highlight the need for retrieval systems that can process documents in their native visual format, preserving both textual and visual information. This challenge has motivated the development of multimodal retrieval approaches that can understand and retrieve from documents as images. In order to evaluate multimodal retrieval models, several benchmarks have been introduced. The most popular benchmarks on visual document retrieval are ViDoRe V1 [7] and ViDoRe V2 [15], which encompass various domains, including academic, artificial intelligence, government reports, healthcare industry, etc.

In this report, we introduce llama-nemoretriever-colembed, a family of state-of-the-art text-image retrieval models designed for scalable and accurate multimodal retrieval. Our best-performing model, llama-nemoretriever-colembed-3b, achieves an NDCG@5 of 91.0 on ViDoRe V1 and 63.5 on ViDoRe V2, ranking No.1 on both benchmarks (as of June 27, 2025). We initialized our models from NVIDIA’s Eagle2 vision-language model [10, 11], replaced the causal attention with bidirectional attention, and fine-tuned the models through contrastive training on curated multimodal datasets. Our training datasets contain both text-only and text-image pairs, and we apply hard negative mining following the methods proposed in NV-Retriever [2] to improve retrieval performance. Our contributions include:

- We release two state-of-the-art text-image retrieval models: llama-nemoretriever-colembed-1B and 3B. The 3B model achieves top-1 performance on both ViDoRe V1 and ViDoRe V2 benchmarks, while the 1B model outperforms several leading 3B and 7B models.
- We explore two-stage training strategy, where the first stage leverages large-scale text-only data while the second stage uses text-image data. Our results demonstrate that pretraining on large-scale text-only retrieval data significantly enhances the model’s performance on downstream text-image retrieval tasks, highlighting the transferability of textual retrieval capabilities to multimodal settings.
- While the ColBERT-style late interaction mechanism enables fine-grained retrieval, it introduces additional overhead in terms of throughput and storage compared to simpler pooling strategies such as average or last-token pooling. We analyze and discuss the performance trade-offs of these approaches in production settings and compare an alternative method that incorporates a vision-language reranker model into the retrieval pipeline.

2 Model

2.1 Model Architecture

Our model adopts a bi-encoder retrieval framework inspired by prior dense retrieval approaches such as NV-Retriever [2], NV-Embed [1] and E5 [16]. In this setting, both the query and corpus items (text or image) are independently passed through a shared multimodal encoder, which projects them into a common embedding space. The relevance between a query and corpus entry is computed via a similarity function (e.g., cosine similarity or dot product), enabling fast and scalable retrieval across large corpora.

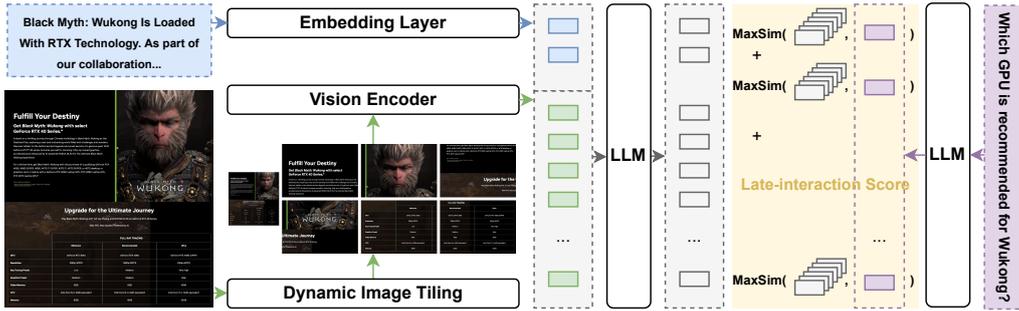


Figure 1: Multimodal Retrieval Architecture with Dynamic Image Tiling and Late-Interaction Scoring

We build our retrieval model on top of the NVIDIA Eagle 2 Vision-Language Models [10, 11]. These models adopt dynamic image tiling to support inputs of varying resolutions, and employ a carefully curated data strategy that improves multimodal learning. These design choices enable Eagle 2 models to achieve state-of-the-art results on several multimodal benchmarks, providing a solid foundation for retrieval tasks.

We further adapt the Eagle architecture by changing the causal attention to bidirectional attention and fine-tuning it under a retrieval-specific contrastive learning objective. As part of the dynamic tiling mechanism, `max_input_tiles` and `min_input_tiles` parameters are used to control the number of tiles produced from each image. For training, we set `max_input_tiles` = 2 to maintain memory efficiency, as increasing it to 4 showed no performance gains. During inference, we set `max_input_tiles` = 6 to allow finer visual granularity.

To support different deployment scenarios, we develop variants of the model at multiple scales (e.g., 1B³, 3B) as shown in Table 1, allowing trade-offs between performance and retrieval accuracy.

Model (Huggingface ID)	Parameters (B)	Embedding Dimension
nvidia/llama-nemoretriever-colembd-1b-v1	2.42	2048
nvidia/llama-nemoretriever-colembd-3b-v1	4.41	3072

Table 1: Overview of model architecture

2.2 Late-interaction

The late interaction mechanism introduced by ColBERT [17] enables fine-grained interactions between query and document tokens. As shown in Figure 2b, for a query, each token embedding interacts with all document token embeddings using a MaxSim operator, which selects the maximum similarity per query token and sums these scores to produce the final relevance score. This requires storing all token embeddings of the document corpus (text or images). At inference time, query token embeddings are computed and interact with the stored document embeddings through this MaxSim process. We adopt this mechanism in our models to enable fine-grained retrieval. While this approach offers the expressiveness of token-level matching, compared to simpler pooling methods such as average or last-token pooling, as shown in Figure 2a, the late-interaction

³We refer to it as the "1b" model because it leverages a 1B Llama model as the language backbone, with the complete architecture containing 2.42B parameters from both the vision encoder and language base model

method introduces latency and storage overhead that may need to be assessed, as they become a concern for real-world applications. We will discuss and compare these trade-offs in Section 5.

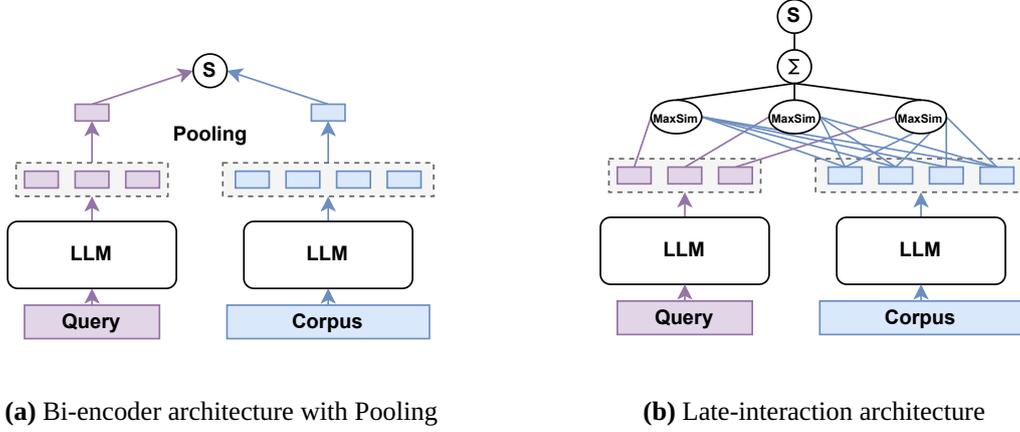


Figure 2: Visualization of bi-encoder and late-interaction architecture

3 Training

3.1 Contrastive learning

We leverage contrastive learning to maximize the embedding similarity between the query and positive passage, while minimizing the similarity between the query and negative corpus. We adopt the InfoNCE contrastive loss [18] to train the model to distinguish between positive and negative pairs in a shared embedding space:

$$\mathcal{L}(q, d^+, D_N) = -\log \frac{\exp(\text{sim}(q, d^+)/\tau)}{\sum_{d_i \in \{d^+\} \cup D_N} \exp(\text{sim}(q, d_i)/\tau)}, \quad (1)$$

where q is the embedding of a query, and d^+ are embeddings positive documents. D_N denotes the set of negative passages. $\text{sim}(\cdot)$ represents the similarity function (e.g., cosine similarity or dot product). τ is the temperature parameter.

To improve the effectiveness of contrastive learning, we incorporate the *top-k with percentage to positive threshold* strategy from NV-Retriever [2] for hard negative mining. We set the threshold as 0.95, meaning we select the K most relevant negative samples whose similarity to the query is less than 95% of the query–positive similarity score, we set $K = 2$. This encourages the model to learn from challenging negatives, while removing potential false negatives that have high similarity scores.

3.2 Two-stage training

Stage 1: Text-Only Pretraining In the first stage, we modify the model architecture by replacing causal attention with bidirectional attention and train it on a large-scale text-only retrieval corpora using a contrastive loss. This stage helps establish strong foundational retrieval models for textual queries and documents, allowing the model to learn semantic similarity in the embedding space.

Stage 2: Text-Image Fine-Tuning In the second stage, we fine-tune the model using text–image pairs. This stage aligns the learned text representations with visual inputs, enabling alignment across different modalities.

3.3 Datasets

In the first stage, we follow the NV-Retriever methodology and use large-scale text-text pairs to establish strong language representations. Including HotpotQA [19], MIRACL [20], Natural Questions (NQ) [21], Stack Exchange [22], SQuAD [23] and Tiger Math/Stack [24]. In the second stage, we fine-tune the model on a mixture of text-image retrieval datasets, including ColPali train set [7], Wiki-SS-NQ [25, 26], VDR [27], VisRAG-Ret-Train-Synthetic-data [28], VisRAG-Ret-Train-In-domain-data [28] and Docmatix [29]. These datasets provide diverse and challenging multimodal examples that help the model generalize to visual retrieval tasks. We provide a full list of training datasets in the model card on HuggingFace⁴.

4 Results

We evaluate our llama-nemoretriever-colembd models on two benchmarks: ViDoRe V1 and ViDoRe V2^{5 6}. The MTEB Visual Document Retrieval (VDR) leaderboard⁷ combines results from both ViDoRe benchmarks.⁸ Both model variants achieve state-of-the-art performance while maintaining superior parameter efficiency. As shown in Table 3 and 4, as of June 27, our llama-nemoretriever-colembd-3b-v1 model achieves top performance across both benchmarks: 91.0 on ViDoRe V1 (vs 89.9 best baseline) and 63.5 on ViDoRe V2 (vs 60.7 best baseline), resulting in a leading MTEB VDR score of 83.1 (vs 81.3 best baseline), as shown in Table 2. The 1B variant also consistently outperforms all baselines with scores of 90.5 and 62.1 on ViDoRe V1 and V2 respectively, achieving 82.63 on MTEB VDR. Our models demonstrate exceptional efficiency: the 1B variant (2.42B parameters) outperforms 3B and 7B baseline models, while the 3B variant (4.41B parameters) achieves state-of-the-art results using fewer parameters than 7B competitors. We provide the results for the deprecated ViDoRe V1 and V2 leaderboards in Appendix A.

Model	Parameters	Embedding Dim.	Max Tokens	MTEB VDR
nommic-ai/colnomic-embed-multimodal-7b	7B	128	128000	81.30
vidore/colqwen2.5-v0.2	3B	128	128000	81.23
nommic-ai/colnomic-embed-multimodal-3b	3B	128	128000	80.02
vidore/colqwen2-v1.0	2B	128	32768	79.74
vidore/colpali-v1.3	2B	128	16384	76.34
vidore/colpali-v1.2	2B	128	16384	74.72
vidore/colSmol-500M	500M	128	8192	71.36
Ours				
nvidia/llama-nemoretriever-colembd-1b-v1	2B	2048	8192	82.63
nvidia/llama-nemoretriever-colembd-3b-v1	4B	3072	8192	83.10

Table 2: Evaluation of baseline models and our models on MTEB: Visual Document Retrieval. Results are presented using nDCG@5 metrics.

MIRACL-VISION [30] is a multilingual visual document retrieval benchmark covering 18 languages, and serves as an extension of the original MIRACL dataset. We evaluate our models on this benchmark and compare them against text-only retrieval models on the MIRACL-VISION

⁴<https://huggingface.co/nvidia/llama-nemoretriever-colembd-3b-v1>

⁵ViDoRe leaderboard: <https://huggingface.co/spaces/vidore/vidore-leaderboard>

⁶The results can slightly change based on the code base. ViDoRe V1 and V2 were changed to use the results based on the MTEB evaluation code. In addition, ViDoRe V2 was changed from 7 datasets to 4 datasets. Our results are based on the new ViDoRe V1 and V2 protocols, if not stated otherwise.

⁷MTEB VDR: http://mteb-leaderboard.hf.space/?benchmark_name=VisualDocumentRetrieval

⁸MTEB ranks models by Borda Count [6], we use Avg NDCG@5 for simplicity because the Rank Borda scores are not visible in the MTEB leaderboards.

Model	Size (M)	Avg.	ArxivQA	DocVQA	InfoVQA	Shift Project	AI	Energy	Gov. Reports	Healthcare	TabFQuad	TAT-DQA
nomic-ai/colnomic-embed-multimodal-3b	3000	89.9	88.2	61.3	92.8	90.2	96.3	97.3	96.6	98.3	94.5	83.1
nomic-ai/colnomic-embed-multimodal-7b	7000	89.8	88.4	60.1	92.3	89.3	99.3	96.6	95.4	99.3	96.1	81.2
vidore/colqwen2.5-v0.2	3000	89.6	89.1	63.5	92.6	88.0	99.6	95.8	96.6	98	90.8	82.1
vidore/colqwen2-v1.0	2210	89.2	88	61.5	92.5	89.9	99.0	95.9	95.5	98.8	89	82.2
vidore/colpali-v1.3	2920	84.7	83.7	58.7	85.7	76.5	96.6	94.6	95.9	97.4	86.7	70.7
vidore/colpali-v1.2	2920	83.4	77.9	56.5	82.4	78.3	97.5	94.4	94.9	95.4	88.4	68.1
Ours												
nvidia/llama-nemoretriever-colembed-1b-v1	2418	90.5	87.6	64.5	93.6	92.3	100	96.6	96.7	99.6	94.3	79.8
nvidia/llama-nemoretriever-colembed-3b-v1	4407	91.0	88.4	66.2	94.9	90.7	99.6	96.6	97.8	99.3	95.9	80.6

Table 3: Evaluation of baseline models and our models on ViDoRe V1 (as of June 27th). Results are presented using nDCG@5 metrics

Model	Size (M)	Avg.	MIT	Economics	ESG Restaurant	ESG Restaurant
			Biomedical Multilingual	Macro Multilingual	Human English	Synthetic Multilingual
nomic-ai/colnomic-embed-multimodal-7b	7000	60.7	63.4	57	68	54.4
vidore/colqwen2.5-v0.2	3000	59.5	59.3	53	67.1	58.5
nomic-ai/colnomic-embed-multimodal-3b	3000	55.2	62.8	53.8	56.3	48
vidore/colpali-v1.3	2920	55.1	53.2	51	60.4	55.9
vidore/colqwen2-v1.0	2210	55	56.3	50.6	60.4	52.5
vidore/colpali-v1.2	2920	52.2	54.8	45.9	56.8	51.2
Ours						
nvidia/llama-nemoretriever-colembed-1b-v1	2418	62.1	62.9	53.2	76.4	55.9
nvidia/llama-nemoretriever-colembed-3b-v1	4407	63.5	64.3	55.9	75.4	58.6

Table 4: Evaluation of baseline models and our models on ViDoRe V2 (as of June 30). Results are presented using nDCG@5 metrics.

(text) subset. As shown in Table 5, our models consistently outperform prior visual retrieval models across the MIRACL-VISION benchmark. The 3B variant achieves the highest overall mean score (0.5841), demonstrating strong multilingual retrieval capabilities.

5 Real-World Applications

Leaderboards and benchmarks typically evaluate performance based on accuracy metrics. Some also include proxy indicators of computational efficiency, such as model size or embedding dimensionality. Ultimately, rankings are determined by accuracy, which may not reflect the broader needs of real-world applications. No solution fits all use-cases. In this section, we discuss the trade-offs of llama-nemoretriever-colembed in the context of production deployment.

5.1 Review Characteristics

Deploying a production system involves balancing accuracy, latency/throughput, and cost. A typical production system can be broken down into the following components:

- **Embedding:** All embeddings of documents are generated by retrieval model. This step can be performed in batches, with support for continuous updates as new documents arrive. Throughput and cost are key considerations, and the overall retrieval performance is primarily influenced by the size of the model.
- **Storage:** The embeddings are stored for retrieval, and storage requirements are primarily determined by the embedding dimension and precision, i.e., how many bytes are needed per document.
- **Serving:** Latency measures how quickly documents can be retrieved in response to a user query. Since queries are typically short (around 50–100 tokens), the size of the embedding

Language	MIRACL-VISION (Text)				MIRACL-VISION (Image)					
	multilingual e5-large	Snowflake/snowflake arctic-embed l-v2.0	Alibaba-NLP/gte multilingual base	BAAI/bge-m3	MrLight/dse qwen2-2b mrl-v1	Alibaba-NLP/gme Qwen2-VL-2B Instruct	llamaindex/vdr 2b-multi-v1	vidore/colqwen2 v1.0	nvidia/llama nemoretriever colembed-1b-v1	nvidia/llama nemoretriever colembed-3b-v1
Arabic	0.8557	0.8754	0.8503	0.8883	0.3893	0.4888	0.4379	0.4129	0.3596	0.4247
Bengali	0.8421	0.8325	0.8211	0.8585	0.2352	0.3755	0.2473	0.2888	0.3715	0.4878
Chinese	0.6900	0.7179	0.7167	0.7458	0.5962	0.6314	0.5963	0.4926	0.3869	0.4355
English	0.7029	0.7437	0.7345	0.7348	0.6605	0.6784	0.6784	0.6417	0.7165	0.7363
Farsi	0.6793	0.7001	0.6984	0.7297	0.2250	0.3085	0.2398	0.2616	0.2803	0.3109
Finnish	0.8974	0.9014	0.8957	0.9071	0.4162	0.6863	0.5283	0.6604	0.8278	0.8513
French	0.7208	0.8236	0.7771	0.8158	0.7160	0.6851	0.7194	0.6876	0.7959	0.7988
German	0.7622	0.7774	0.7498	0.7695	0.6267	0.6345	0.6205	0.5995	0.6515	0.6831
Hindi	0.7595	0.7255	0.6916	0.7581	0.1740	0.3127	0.2058	0.2209	0.4670	0.4867
Indonesian	0.6793	0.6906	0.6757	0.7049	0.4866	0.5416	0.5254	0.5320	0.6295	0.6428
Japanese	0.8378	0.8484	0.8442	0.8720	0.6232	0.7305	0.6553	0.6970	0.6790	0.7260
Korean	0.7327	0.7545	0.7397	0.7934	0.4446	0.6202	0.4952	0.4419	0.4430	0.5158
Russian	0.7857	0.8242	0.8023	0.8363	0.6505	0.7202	0.6995	0.6811	0.7227	0.7670
Spanish	0.6596	0.7250	0.7029	0.7268	0.5927	0.6277	0.6274	0.6224	0.7036	0.7109
Swahili	0.8157	0.8089	0.7987	0.8337	0.4156	0.5348	0.4509	0.4931	0.7326	0.7767
Telugu	0.8948	0.9201	0.9076	0.9090	0.0274	0.0893	0.0318	0.0264	0.0853	0.1669
Thai	0.8424	0.8485	0.8509	0.8682	0.2692	0.3563	0.3177	0.2389	0.3738	0.4035
Yoruba	0.5655	0.5332	0.5698	0.5842	0.4178	0.4884	0.4577	0.5120	0.5250	0.5888
AVG.	0.7624	0.7806	0.7682	0.7964	0.4426	0.5283	0.4741	0.4728	0.5414	0.5841

Table 5: Evaluation results on MIRACL-VISION benchmark comparing text-based and image-based retrieval models across multiple languages.

model plays a smaller role in this stage. Incorporating a reranker in the retrieval pipeline, such as a cross-encoder, can improve accuracy, but at the cost of increasing the latency to serve another model.

Each stage of the retrieval pipeline involves trade-offs that should be carefully aligned with the specific use case. For instance, in scenarios with a small corpus but a high volume of user queries, a larger embedding model without a reranker may offer better performance. On the other hand, for a large corpus with a moderate number of queries, a smaller embedding model combined with a reranker can be more cost-efficient. In [31], we explored these trade-offs by analyzing the effects of embedding model size, inference throughput, and accuracy, both individually and with a reranker integrated into the retrieval pipeline.

5.2 Retrieval Architecture Comparison

ColBERT introduced the late-interaction paradigm, which demonstrated significant performance improvements in retrieval tasks by preserving fine-grained token-level interactions between queries and documents. Unlike traditional pooling strategies that compress entire sequences into single vectors, late-interaction models leverage all token-level representations. However, this approach introduces a fundamental trade-off between accuracy and storage cost, as each document requires multiple token embeddings, leading to significantly increased storage requirements.

Table 6 summarizes the comprehensive trade-offs between different retrieval approaches, reporting storage requirements in gigabytes for embedding one million images. The storage footprint of late-interaction models depends on three key factors: token count (sequence length), embedding dimension, and numerical precision (e.g., float32, float16, int8). Our analysis reveals that ColEmbed models require more storage than bi-encoder alternatives. The nvidia/llama-nemoretriever-colembed-3b-v1 model with full dimensionality (3072) requires 10,311.1 GB for one million images, representing over 2,700 times more storage than comparable bi-encoder models.

Several techniques can reduce storage requirements for both paradigms. Linear projection layers can substantially reduce embedding dimensions. Following the approach used in vidore/colqwen2-v1.0 models [7], we applied linear projection layers to reduce the output dimension from 3072 to 512 and using smaller resolutions via dynamic image tiling, it reduces storage by approximately 88% with only modest accuracy degradation (ViDoRe V1: 0.9106 to 0.9064). While this significantly decreases storage requirements to 1,230.2 GB, it still remains over 300 times larger than bi-encoder approaches. The vidore/colqwen2-v1.0 model [7], with its more compact 128-dimension

embeddings, requires 179.1 GB, demonstrating how dimensionality reduction can significantly impact storage. However, bi-encoder models represent each document with a single embedding vector, requiring only 2.9-3.8 GB for one million images—at least 47 times less storage than the most compact late-interaction model in our comparison.

Additionally, binary quantization can store embeddings using only 1-bit per element, potentially reducing storage by 16x for both architectures. However, our experience with bi-encoders indicates that binary quantization performs poorly when the embedding dimensionality is too small, and these techniques require further testing with late-interaction embedding size of 128. AnswerAI’s late-pooling approach [32] can reduce token vectors by factors of 3-5, while MU-VERA [33] proposes converting multi-vector embeddings into single Fixed Dimensional Encodings (FDEs) whose inner product approximates multi-vector similarity, enabling the use of standard single-vector retrieval with smaller total embedding size.

Beyond storage, production systems must balance accuracy, latency, throughput, and inference costs. Late-interaction models face additional inference constraints due to late-interaction calculations, which require specialized vector database support and introduce latency overhead. A practical strategy is to enhance bi-encoder models with rerankers to improve retrieval accuracy while maintaining storage efficiency. Our experiments with the llama-3_2-nemoretriever-1b-vlm-embed-v1 model⁹ demonstrate this approach’s effectiveness: the base model achieves ViDoRe V1: 0.8313 and V2: 0.5178 with 3.8 GB storage, while adding reranking with 25 candidates improves performance to ViDoRe V1: 0.9064 and V2: 0.6214 at the cost of approximately 2,368 ms additional latency per query.

This hybrid approach achieves performance comparable to the llama-nemoretriever-colembed-3b-v1 model while maintaining the storage advantages of bi-encoder architectures. The results demonstrate a clear trade-off between the number of reranked candidates and both inference latency and retrieval accuracy. The choice between late-interaction and bi-encoder paradigms ultimately depends on specific use case requirements and system constraints.

Architecture	Model	Avg. # of embeddings (Sequence Length)	Embedding Dimension	# of floating points numbers per document	Embedding Storage for 1M images (GB)	Number Candidates Reranked	Additional Inference Time (ms/query)	ViDoRe V1	ViDoRe V2
ColEmbed	nvidia/llama-nemoretriever-colembed-3b-v1	1802	3072	5535744	10311.1	N/A	N/A	0.9106	0.6357
ColEmbed	nvidia/llama-nemoretriever-colembed-3b-v1	1290	512	660480	1230.2	N/A	N/A	0.9064	0.6109
ColEmbed	vidore/colqwen2-v1.0	751	128	96128	179.1	N/A	N/A	0.8906	0.5290
Bi-Encoder	MrLight/dse-qwen2-2b-mrl-v1	1	1536	1536	2.9	N/A	N/A	0.8510	0.5590
Bi-Encoder	lama-3_2-nemoretriever-1b-vlm-embed-v1	1	2048	2048	3.8	N/A	N/A	0.8313	0.5178
Bi-Encoder	lama-3_2-nemoretriever-1b-vlm-embed-v1 + reranker	1	2048	2048	3.8	10	960	0.8931	0.6025
Bi-Encoder	lama-3_2-nemoretriever-1b-vlm-embed-v1 + reranker	1	2048	2048	3.8	25	2368	0.9064	0.6214
Bi-Encoder	lama-3_2-nemoretriever-1b-vlm-embed-v1 + reranker	1	2048	2048	3.8	100	9392	0.9101	0.6182

Table 6: Comparison of system-relevant characteristics of different retrieval pipelines and models on ViDoRe benchmarks. The numbers slightly differ as we used another code base to evaluate the datasets. Sequence length is calculated by the median across ViDoRe V1. Note: There is a newer ColQwen model vidore/colqwen2.5-v0.2.

6 Conclusion

We present llama-nemoretriever-colembed, a family of scalable and high-performing text-image retrieval models that achieve state-of-the-art results on ViDoRe V1, ViDoRe V2, and MIRACL-VISION benchmarks. By modifying the Eagle2 VLM architecture with bidirectional attention and integrating a ColBERT-style late interaction mechanism, our models support fine-grained multimodal retrieval in a shared embedding space. Trained using a two-stage training pipeline combining large-scale text and image data, our models demonstrate strong generalization and multilingual

⁹A commercial multimodal retrieval model representing user queries as text and documents as images, https://build.nvidia.com/nvidia/llama-3_2-nemoretriever-1b-vlm-embed-v1

retrieval capabilities. We also analyze the trade-offs introduced by token-level late interaction and highlight key considerations for real-world deployment. Our release of both 1B and 3B model variants provides a strong foundation for future research and practical applications in multimodal retrieval.

Core Contributors

NeMo Retriever Applied Research – Embedding and Ranking Models: Mengyao Xu, Gabriel Moreira, Ronay Ak, Radek Osmulski, Yauhen Babakhin, Benedikt Schifferer

Contributors

NeMo Retriever Applied Research – OCR Models: Bo Liu, Theo Viel, Maximilian Jeblick

Nemotron VLM: Zhiding Yu ¹⁰, Padmavathy Subramanian, Karan Sapra, Andrew Tao

NeMo Product: Nave Algarici, Sean Sodha, Ben Jarmak

Data: Shizhe Diao, Tom Balough

NeMo Retriever Services: Kalpesh Sutaria, Loan Luong, Oliver Holworthy, Jeremy Jordan, Alex Richards, Fei Yu, Salik Siddiqui, Charles Blackmon-Luca, Derek Whatley, Adam Laiacano, Tom O’Brien, Randy Gelhausen, Jeremy Dyer, Edward Kim, Sohail Sahi, Julio Perez, Steven Baughman, Kam Mitchell, Meghana Shrotri

Management: Even Oldridge, Joey Conway, John Zedlewski, Kari Briski

A ViDore Legacy Results

While ViDore has released updated versions of their benchmarks, we include results from the previous evaluation framework for completeness. Tables 7 and 8 present our models’ performance on the legacy ViDore V1 and V2 benchmarks, calculated using the original codebase methodology.

On the legacy ViDore V1 benchmark (Table 8), our models maintain their competitive advantage, llama-nemoretriever-colembed-3b-v1 achieves the highest average score of 91.1, followed by our 1b variant at 90.5. Both models outperform all baseline models, with the next best performers achieving 90.4. Similarly, on the legacy ViDore V2 benchmark (Table 8), our models demonstrate superior performance with scores of 63.4 (3b model) and 62.6 (1b model), outperforming the best baseline of 62.1. These results are consistent with the updated benchmark evaluations.

	Avg.	ArxivQA	DocVQA	InfoVQA	Shift Project	AI	Energy	Gov. Reports	Healthcare	TabFQuad	TAT-DQA
tsystems/colqwen2.5-3b-multilingual-v1.0	90.4	93.4	64	93	88.3	100	95.9	96.1	97.7	95.1	80.7
Metric-AI/ColQwen2.5-7b-multilingual-v1.0	90.4	91.7	65.1	93.9	87.7	99.3	95.4	95.2	97.8	96.7	80.9
Metric-AI/ColQwen2.5-3b-multilingual-v1.0	90.3	92.2	64.4	93.5	88.4	98.9	96.5	96.4	98.4	93.7	80.7
nomie-ai/colnomie-embed-multimodal-7b	90.2	88.7	61.3	93.4	91.8	99.3	96.5	95.1	99.3	96	81.1
yydxlv/colqwen2.5-7b-v0.1	90.2	91.1	63.1	93.5	89.1	98.9	95.6	96.2	98.5	93.9	81.9
tsystems/colqwen2-7b-v1.0	90.1	90.7	64.5	92	89.3	99.3	96.3	96.3	99.3	95	78.6
Ours											
nvidia/llama-nemoretriever-colembed-1b-v1	90.5	87.6	64.2	93.6	92.3	100	96.6	96.7	99.6	94.3	79.9
nvidia/llama-nemoretriever-colembed-3b-v1	91.1	88.4	65.9	94.9	90.7	99.6	96.6	97.8	99.3	95.9	80.6

Table 7: Evaluation of baseline models and our models on legacy ViDore V1 (as of June 27). Results are presented using nDCG@5 metrics

References

- [1] Chankyu Lee, Rajarshi Roy, Mengyao Xu, Jonathan Raiman, Mohammad Shoeybi, Bryan Catanzaro, and Wei Ping. Nv-embed: Improved techniques for training llms as generalist embedding models. *arXiv preprint arXiv:2405.17428*, 2024.

¹⁰Core Contributor

Model	Avg.	ESG Restaurant Human	Economics Macro Multilingual	MIT Biomedical	ESG Restaurant Synthetic	ESG Restaurant Synthetic Multilingual	MIT Biomedical Multilingual	Economics Macro
nomic-ai/colnomic-embed-multimodal-7b	62.1	73.9	54.7	66.1	57.3	56.7	64.2	61.6
vidore/colqwen2.5-v0.2	60.6	68.4	56.5	63.6	57.4	57.4	61.1	59.8
nomic-ai/colnomic-embed-multimodal-3b	60.2	65.8	55.5	63.5	56.6	57.2	62.5	60.2
Alibaba-NLP/gme-Qwen2-VL-7B-Instruct	59.3	65.8	56.2	64.0	54.3	56.7	55.1	62.9
nomic-ai/nomic-embed-multimodal-7b	59.0	65.7	57.7	64.0	49.2	51.9	61.2	63.1
tsystems/colqwen2.5-3b-multilingual-v1.0	58.6	72.1	51.2	65.3	51.7	53.3	61.7	54.8
Ours								
nvidia/llama-nemoretriever-colembed-1b-v1	62.6	76.9	56.4	64.7	57.1	56.8	62.3	64.1
nvidia/llama-nemoretriever-colembed-3b-v1	63.4	74.7	58.0	65.7	58.8	57.6	63.2	66.0

Table 8: Evaluation of baseline models and our models on legacy ViDoRe V2 (as of June 27). Results are presented using nDCG@5 metrics

- [2] Gabriel de Souza P Moreira, Radek Osmulski, Mengyao Xu, Ronay Ak, Benedikt Schifferer, and Even Oldridge. Nv-retriever: Improving text embedding models with effective hard-negative mining. *arXiv preprint arXiv:2407.15831*, 2024.
- [3] Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, Fei Huang, and Jingren Zhou. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*, 2025.
- [4] Liang Wang, Nan Yang, Xiaolong Huang, Linjun Yang, Rangan Majumder, and Furu Wei. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*, 2023.
- [5] Niklas Muennighoff, Nouamane Tazi, Loïc Magne, and Nils Reimers. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*, 2022.
- [6] Isaac Chung, Imene Kerboua, Marton Kardos, Roman Solomatin, and Kenneth Enevoldsen. Maintaining mteb: Towards long term usability and reproducibility of embedding benchmarks. *arXiv preprint arXiv:2506.21182*, 2025.
- [7] Manuel Faysse, Hugues Sibille, Tony Wu, Bilel Omrani, Gautier Viaud, Céline Hudelot, and Pierre Colombo. Colpali: Efficient document retrieval with vision language models, 2024.
- [8] Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, et al. Qwen2-vl: Enhancing vision-language model’s perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*, 2024.
- [9] Akhiad Bercovich, Itay Levy, Izik Golan, Mohammad Dabbah, Ran El-Yaniv, Omri Puny, Ido Galil, Zach Moshe, Tomer Ronen, Najeeb Nabwani, et al. Llama-nemotron: Efficient reasoning models. *arXiv preprint arXiv:2505.00949*, 2025.
- [10] Zhiqi Li, Guo Chen, Shilong Liu, Shihao Wang, Vibashan VS, Yishen Ji, Shiyi Lan, Hao Zhang, Yilin Zhao, Subhashree Radhakrishnan, et al. Eagle 2: Building post-training data strategies from scratch for frontier vision-language models. *arXiv preprint arXiv:2501.14818*, 2025.
- [11] Guo Chen, Zhiqi Li, Shihao Wang, Jindong Jiang, Yicheng Liu, Lidong Lu, De-An Huang, Wonmin Byeon, Matthieu Le, Tuomas Rintamaki, et al. Eagle 2.5: Boosting long-context post-training for frontier vision-language models. *arXiv preprint arXiv:2504.15271*, 2025.

- [12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.
- [13] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muhammad Ferjad Naeem, Ibrahim Al-abdulmohsin, Nikhil Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil Mustafa, et al. Siglip 2: Multilingual vision-language encoders with improved semantic understanding, localization, and dense features. *arXiv preprint arXiv:2502.14786*, 2025.
- [14] Mike Ranzinger, Greg Heinrich, Jan Kautz, and Pavlo Molchanov. Am-radio: Agglomerative vision foundation model reduce all domains into one. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12490–12500, 2024.
- [15] Quentin Macé, António Loison, and Manuel Faysse. Vidore benchmark v2: Raising the bar for visual retrieval. *arXiv preprint arXiv:2505.17166*, 2025.
- [16] Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. Text embeddings by weakly-supervised contrastive pre-training. *arXiv preprint arXiv:2212.03533*, 2022.
- [17] 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*, pages 39–48, 2020.
- [18] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PmLR, 2020.
- [19] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.
- [20] Xinyu Zhang, Nandan Thakur, Odunayo Ogundepo, Ehsan Kamaloo, 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.
- [21] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Jacob Devlin, Kenton Lee, et al. Natural questions: a benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:453–466, 2019.
- [22] Stack Exchange, Inc. Stack Exchange Community Data Dump, 2023. <https://archive.org/details/stack-exchange-data-dump-2023-09-12>, 2023. Accessed: 2025-06-30.
- [23] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016.
- [24] Xiang Yue, Tuney Zheng, Ge Zhang, and Wenhui Chen. Mammoth2: Scaling instructions from the web. *Advances in Neural Information Processing Systems*, 2024.
- [25] et al. Xueguang Ma. Tevatron/wiki-ss-nq, 2024.

- [26] Xueguang Ma, Luyu Gao, Shengyao Zhuang, Jiaqi Samantha Zhan, Jamie Callan, and Jimmy Lin. Tevatron 2.0: Unified document retrieval toolkit across scale, language, and modality. *arXiv preprint arXiv:2505.02466*, 2025.
- [27] Marco Cimolai and Logan Markewich. llamaindex/vdr-multilingual-train, 2025.
- [28] Shi Yu, Chaoyue Tang, Bokai Xu, Junbo Cui, Junhao Ran, Yukun Yan, Zhenghao Liu, Shuo Wang, Xu Han, Zhiyuan Liu, and Maosong Sun. Visrag: Vision-based retrieval-augmented generation on multi-modality documents, 2024.
- [29] Hugo Laurençon, Andrés Marafioti, Victor Sanh, and Léo Tronchon. Building and better understanding vision-language models: insights and future directions., 2024.
- [30] Radek Osmulsk, Gabriel de Souza P Moreira, Ronay Ak, Mengyao Xu, Benedikt Schifferer, and Even Oldridge. Miracl-vision: A large, multilingual, visual document retrieval benchmark. *arXiv preprint arXiv:2505.11651*, 2025.
- [31] Gabriel de Souza P Moreira, Ronay Ak, Benedikt Schifferer, Mengyao Xu, Radek Osmulski, and Even Oldridge. Enhancing q&a text retrieval with ranking models: Benchmarking, fine-tuning and deploying rerankers for rag. In *Proceedings of the 1st Workshop on GenAI and RAG Systems for Enterprises, co-located with CIKM*, 2024.
- [32] Benjamin Clavié. A little pooling goes a long way for multi-vector representations, 2024.
- [33] Laxman Dhulipala, Majid Hadian, Rajesh Jayaram, Jason Lee, and Vahab Mirrokni. Muvera: multi-vector retrieval via fixed dimensional encodings. *arXiv preprint arXiv:2405.19504*, 2024.