

Decoding Latent Spaces: LLM-based Embeddings vs Encoder-Only Architectures

David N. Juboor
Dr. John MacLaren Walsh
Dr. Christopher Weyant

December 9, 2024

Abstract

Text embeddings are foundational to retrieval tasks, enabling applications such as search and recommendation systems. While encoder-only architectures, such as BERT-based models, are well-established for generating robust embeddings, decoder-only architectures, optimized for autoregressive tasks, remain underexplored in this context. This study evaluates the comparative performance of decoder-only and encoder-only models in creating text embeddings for retrieval tasks, leveraging the Massive Text Embedding Benchmark (MTEB) for evaluation across diverse datasets and domains. We empirically validate that decoder-only models, such as `mistralai/Mistral-7B-v0.1`, are poorly suited for retrieval, achieving significantly lower performance compared to encoder-only architectures. Notably, embeddings from intermediate and penultimate layers in decoder models outperformed final-layer embeddings, highlighting potential layer-wise trends in feature representation. Additionally, task-specific fine-tuning of encoder-only models resulted in performance trade-offs, with slight improvements in some domains but overall degradation across general retrieval tasks, underscoring the complexities of fine-tuning for domain-specific use cases. These findings reinforce the superiority of encoder-only models for retrieval and contribute to the understanding of architectural trade-offs in text embedding models, offering actionable insights for practitioners optimizing embeddings for retrieval-based applications.

1 Introduction

Text embeddings are a cornerstone of modern natural language processing (NLP), enabling tasks such as information retrieval, semantic search, and recommendation systems. Encoder-only architectures, such as BERT-based models [1], are well-established for generating robust embeddings due to their ability to effectively model contextual relationships in text. However, the potential utility of decoder-only architectures, which are optimized for autoregressive tasks like next-word prediction, remains largely unexplored in the context of embedding-based retrieval apart from training data generation tasks [2] [3] [4] and "LLM as a Judge" evaluation frameworks [5]. This distinction between encoder- and decoder-based embeddings presents an intriguing question: how do these latent spaces differ when applied to the same retrieval tasks?

This work investigates the comparative performance of decoder-only and encoder-only architectures in creating text embeddings for retrieval tasks. Specifically, decoder-only models, such as `mistralai/Mistral-7B-v0.1` [6], are evaluated alongside encoder-only models, including both task-tuned and untuned variants, to examine whether the autoregressive training objectives of decoder models can offer any utility in retrieval applications. By leveraging the Massive Text Embedding Benchmark (MTEB) [7], a widely adopted framework for evaluating embedding models, we provide a robust comparison of these architectures across diverse domains and datasets.

The motivation for this study lies in understanding the philosophical and practical overlap between the tasks these models are optimized for. Encoder-only models excel at representing semantic relationships in text, while decoder-only models focus on sequential coherence. By comparing their embeddings, we aim to uncover whether the "understanding" encoded by these architectures exhibits any meaningful overlap or divergence in retrieval tasks. In addition, this work evaluates the effects of task-specific fine-tuning on encoder-only models, shedding light on its impact on both domain-specific and general retrieval performance.

This study makes two primary contributions:

- **Empirical Validation of Decoder-Only Limitations:** We confirm that decoder-only models, such as `Mistral-7B`, are poorly suited for retrieval tasks, achieving significantly lower performance compared to encoder-based architectures. Interestingly, we observe layer-wise trends in decoder embeddings, with intermediate and penultimate layers outperforming final-layer embeddings.
- **Insights into Task-Specific Fine-Tuning:** Our results demonstrate that task-specific fine-tuning of encoder-only models can degrade general retrieval performance, emphasizing the need for careful ablation studies to balance domain-specific optimization and task generalizability.

To evaluate these questions, we compare models across a wide range of retrieval

datasets from the MTEB framework. These datasets span diverse domains, including scientific, medical, and general QA, enabling a robust analysis of model performance across varying contexts. By synthesizing these results, we aim to provide actionable insights for researchers and practitioners seeking to optimize text embedding models for retrieval tasks.

The remainder of this paper is organized as follows: Section 2 reviews related work on embedding-based retrieval and model fine-tuning. Section 3 describes the experimental methodology, including model selection, datasets, and evaluation metrics. Section 4 presents the experimental results, and Section 5 discusses the key findings and their implications. Finally, the paper concludes with a summary of contributions and future directions.

2 Related Work

The effectiveness of text embeddings in retrieval tasks has been extensively studied, with a primary focus on encoder-only architectures. Decoder-only models, in contrast, are optimized for autoregressive generation and remain underexplored in retrieval contexts. This section situates the present study within the broader landscape of embedding research, task-specific fine-tuning, and evaluation frameworks.

2.1 Embedding Models for Retrieval

Encoder-only architectures, such as BERT [1] and Sentence-BERT [8], are well-established for retrieval tasks due to their ability to capture semantic relationships in text. Recent advances, such as task-tuned LLMs like E5 [3], have further demonstrated the potential of scaling and domain-specific optimization for embedding quality. These models leverage bidirectional context, making them ideal for generating robust embeddings.

In contrast, decoder-only architectures like Mistral-7B [6] are primarily designed for autoregressive tasks, such as next-word prediction. While decoder embeddings have been explored in tasks such as summarization or text generation, their utility in retrieval remains unclear. This study aims to bridge this gap by empirically evaluating decoder embeddings across diverse retrieval tasks.

2.2 Task-Specific Fine-Tuning

Fine-tuning pre-trained models for specific tasks or domains is a common strategy for improving embedding quality. Methods such as triplet-loss-based training have been shown to enhance retrieval performance by optimizing query-positive-negative relationships [4]. However, fine-tuning may introduce trade-offs, such as reduced generalizability to other tasks or domains. This study contributes to this area by highlighting the risks of task-specific fine-tuning,

particularly its potential to degrade performance across both target and non-target domains.

2.3 Evaluation Frameworks

Evaluation benchmarks like the Massive Text Embedding Benchmark (MTEB) [7] have standardized the assessment of embedding models across diverse tasks and datasets. MTEB includes retrieval-focused datasets spanning domains such as science, healthcare, and general QA, providing a robust framework for comparing models. By leveraging MTEB, this study ensures consistency and comparability with existing literature.

2.4 Addressing Gaps in Existing Research

While encoder-only architectures dominate retrieval-focused research, decoder-only models remain under explored in this context. This study addresses this gap by empirically validating the limitations of decoder embeddings. Additionally, the results highlight subtleties in fine-tuning strategies, emphasizing the need for ablation studies to balance domain-specific optimization with general retrieval performance.

3 Methodology

This study evaluates the effectiveness of embeddings generated by a range of models across diverse retrieval tasks. As a structured approach to compare decoder-only and encoder-only architecture embedding performance, our experimental design relies heavily on the use of consistent, 3rd party evaluation frameworks across a consistent set of datasets and specific evaluation metrics [9] in the Retrieval space. Leveraging well accepted evaluation frameworks in our experimental design, systemically addresses our core experimental questions around embedding and retrieval performance between varying embedding architectures in a reasonably reproducible fashion.

3.1 Models Evaluated

The models analyzed are categorized into two primary architectural groups: decoder-only and encoder-only. As decoder-only architectures are primarily utilized in for auto-regressive text generation, we choose this architecture as a base-line for unidirectional encoding not tuned for any particular retrieval task. In addition, we select three models to represent reasonable variants of encoder-only architectures. These distinctions provide a foundation for analyzing performance across retrieval tasks.

Decoder-Only Model The `mistralai/Mistral-7B-v0.1` model [6] represents a decoder-only architecture optimized for text generation. While not

inherently designed for retrieval tasks, embeddings were extracted from the model’s hidden states to explore its potential for such applications. The following strategies were employed to extract embeddings:

- **Fused Average:** Combines hidden states from layers 15, 16, and the penultimate layer by averaging their representations.
- **Intermediate Layer:** Extracts embeddings directly from the 15th hidden state.
- **Penultimate Layer:** Uses the second-to-last hidden state for embeddings.
- **Final Layer Norm:** Relies on the normalized final hidden state for generating representations.

Encoder-Only Models Encoder-based architectures, optimized for high-quality text embeddings, were analyzed through three configurations each serving a unique embedding optimization strategy:

- **E5-Mistral-7B-Instruct:** This bi-encoder model [3] represents a more modern retrieval-optimized encoding architecture that has been fine-tuned on multi-domain data utilizing decoder-only model(s) to generate significant portions of its training data. Its embeddings were evaluated in their pre-trained form without additional modifications.
- **DistilBERT Baseline:** The `sentence-transformers/msmarco-distilbert-base-v4` model [1], widely used for retrieval tasks, represents a computationally efficient baseline in its pre-trained configuration. While used heavily in industry for its latency-performance tradeoff, this represents an "older" retrieval-optimized encoding architecture that has not been optimized explicitly for any of these particular retrieval domains, nor was it trained with LLM-generated data. It was selected primarily due to its retrieval-specific metric performance on general tasks with short queries and long answer context.
- **Task-Tuned DistilBERT:** A fine-tuned version of `msmarco-distilbert-base-v4`, trained over five epochs utilizing triplet loss and query-positive-negative pairs derived from selected in-domain yet out-of-distribution datasets. Training datasets included SciFact, `toughdata/quora-question-answer-dataset`, `kroshan/BioASQ`, and `deepset/covid_qa_deepset`. This architecture represents an "off-the-shelf retrieval optimized encoder" like our previous model, with the addition of being finetuned on data in the scientific domain. With this we aim to observe model performance gains or losses on a specific tasked domain while assessing generalization across other datasets both in and out of domain.

3.2 Datasets

The evaluation leveraged datasets from the Massive Text Embedding Benchmark (MTEB) Retrieval task [7], a widely accepted framework for benchmarking embedding models. The datasets span a broad range of domains and task types, ensuring a comprehensive evaluation. The full list includes:

- **Scientific and Academic Domains:** SciFact, SCIDOCS
- **General QA:** ArguAna
- **Finance:** FiQA2018
- **Healthcare and Environmental Domains:** NFCorpus
- **Community QA (CQADupStack):** CQADupstackAndroidRetrieval, CQADupstackEnglishRetrieval, CQADupstackGamingRetrieval, CQADupstackGisRetrieval, CQADupstackMathematicaRetrieval, CQADupstackPhysicsRetrieval, CQADupstackProgrammersRetrieval, CQADupstackStatsRetrieval, CQADupstackTexRetrieval, CQADupstackUnixRetrieval, CQADupstackWebmastersRetrieval, CQADupstackWordpressRetrieval

Additionally, these datasets have good variation in overall relevant queries per document and are consistent with our "short queries long documents" retrieval paradigm as seen in our dataset breakdowns in Appendix A.4 .

The inclusion of datasets from scientific, finance, healthcare, general QA, and community QA contexts provides a robust evaluation set that captures both domain-specific and generalized retrieval challenges.

3.3 Evaluation Metrics

All datasets were evaluated on the "test set" as described by their respective dataset. In addition, by using the MTEB framework for evaluation, all models and datasets were evaluated with identical code and roughly 100 individual retrieval specific metrics were accumulated for all models on all datasets. That said, all datasets contained a "main score" that described what specific evaluation metric is most relevant for the dataset. As our datasets were Retrieval specific, each dataset noted a "main score" of Normalized Discounted Cumulative Gain at top 10 (NDCG@10), a widely used metric that considers both relevance and ranking position.

NDCG@k, a metric that relies heavily on Discounted Cumulative Gain, rewards systems that rank highly relevant documents closer to the top of the list, aligning with practical retrieval goals on a search results page.

The Discounted Cumulative Gain (DCG) is computed as:

$$\text{DCG@k} = \sum_{i=1}^k \frac{2^{\text{rel}_i} - 1}{\log_2(i + 1)} \quad (1)$$

where rel_i represents the relevance score of the retrieved item at position i . To normalize DCG, the Ideal Discounted Cumulative Gain (IDCG) is calculated for the best possible ranking:

$$\text{IDCG@k} = \sum_{i=1}^k \frac{2^{\text{rel}_i^{\text{ideal}}} - 1}{\log_2(i + 1)} \quad (2)$$

Finally, NDCG is obtained as:

$$\text{NDCG@k} = \frac{\text{DCG@k}}{\text{IDCG@k}} \quad (3)$$

This normalization ensures that NDCG ranges between 0 and 1, with higher scores indicating better performance [9]. While NDCG@10 was prioritized due to its alignment with both our evaluation datasets and ranking-focused retrieval tasks; auxiliary metrics, such as Mean Average Precision (MAP@k), Mean Reciprocal Rank (MRR@k), Precision@k, and Recall@k, were also tracked for additional insights. Additionally, k was evaluated at 20 different positions including 1, 2, 3, 5, 10, 20, 30, 100, etc and all metrics were evaluated as Area Under Curve (AUC) of a typical Receiver Operating Characteristic (ROC) curve. Where applicable, embeddings were precomputed or cached and evaluation metrics were computed against these values for operational efficiency.

3.4 Experimental Design

The experiments were structured to address three central research questions:

- How does task-specific fine-tuning improve retrieval performance compared to pre-trained models?
- How do embeddings from decoder-only models compare to those of encoder-only models in terms of retrieval effectiveness?
- Does fine-tuning on domain-specific datasets enhance retrieval within the domain, and how does it impact general performance across other tasks?

Embeddings for encoder-only models were extracted using standard library functions, while embeddings from the decoder-only model were derived from intermediate and final hidden states, as described earlier. The task-tuned DistilBERT model was fine-tuned using triplet loss to optimize retrieval for query-positive-negative pairs.

Baseline comparisons included:

- **Task-Tuned Baseline:** The E5-Mistral model served as a task-tuned baseline for retrieval tasks.
- **Non-Tuned Baseline:** The pre-trained DistilBERT model provided a computationally efficient, non-domain-specific baseline for broader comparisons.

4 Experiments

The experiments were designed to evaluate the performance of a decoder-only and encoder-only architectures of various classes across a wide range of retrieval datasets for the retrieval task. Each experiment addresses specific research questions, leveraging metrics such as NDCG@k, MAP@k, and MRR@k for evaluation. The experimental setups ensure consistency across datasets and metrics, while focusing on the nuances of model architectures and task-specific fine-tuning to allow a direct comparison between various embedding implementations for retrieval.

4.1 Experiment 1: Decoder-Only Model Evaluation

This experiment investigates the suitability of decoder-only architectures for retrieval tasks. The `mistralai/Mistral-7B-v0.1` model [6] was evaluated by extracting embeddings from four different encoding stages: fused average, intermediate layer, penultimate layer, and final layer normalization. Decoder-only models, such as Mistral-7B, are primarily optimized for autoregressive tasks, and their embeddings lack the bidirectional context present in encoder-based architectures. While atypical for usage in retrieval, this experiment serves as a baseline for utilizing embeddings outside of their intended task (generation vs retrieval).

The primary objectives were:

- Quantify the performance gap between decoder-only and encoder-only models in retrieval tasks.
- Explore whether embeddings from specific layers exhibit better retrieval performance than others.

All evaluations were conducted across the full suite of evaluation datasets and metrics (described above), providing a robust assessment of decoder-only performance.

4.2 Experiment 2: Encoder-Only SLM vs LLM Comparison

This experiment compares an untuned encoder-only small language model (SLM) (`msmarco-distilbert-base-v4`) with the task-tuned large language model (LLM) encoder (`e5-mistral-7b-instruct`) [3]. The evaluation focuses on the trade-offs between embedding performance between relatively small and large encoder-only architectures. While the E5-Mistral model represents a high-performing large encoder-only baseline, DistilBERT provides a smaller, more lightweight, alternative for computationally constrained scenarios.

The goals of this experiment include:

- Evaluate the retrieval performance of the simpler SLM architecture relative to the more complex LLM encoder.
- Assess whether the increased performance of the LLM encoder justifies its higher computational cost in practical applications.

Metrics such as NDCG@1, 3, 5, and 10, as well as MAP, MRR, and others were used to quantify general performance while NDCG@10 was used as the primary score for all evaluation datasets.

4.3 Experiment 3: Task-Tuned SLM vs Non-Task-Tuned SLM

This experiment examines the impact of task-specific fine-tuning on an encoder only SLM. The `msmarco-distilbert-base-v4` model was fine tuned using out of distribution domain-specific training datasets from the scientific domain, including SciFact and BioASQ [5]. It is important to note that we chose training data in a domain within our evaluation set without utilizing training sets from the evaluation datasets themselves. This was a strategic choice intended to steer our finetuning in favor of a specific domain within our task without explicitly biasing our model towards in-distribution data. Fine-tuning was performed with triplet loss, which optimizes embeddings for query-positive-negative relationships, improving retrieval performance.

Dataset Formatting for Fine-Tuning For fine-tuning, training datasets were transformed into a query-positive-negative format:

- Queries and positives were directly derived from the datasets.
- Hard negatives were sampled randomly using in-distribution data (data within the particular dataset).

This in-distribution domain-specific hard-negative sampling method ensures that negative examples are challenging yet relevant, which is critical for effective triplet-loss training [8].

The experiment evaluated the fine-tuned SLM on the full set of retrieval datasets, focusing on:

- Quantifying improvements in task performance resulting from fine-tuning in the scientific domain.
- Assessing whether fine-tuning impacts the model’s generalization performance outside the scientific domain.

4.4 Experiment 4: Task-Tuned SLM vs Task-Tuned LLM

The final experiment compares the fine-tuned SLM from Experiment 3 with the task-tuned LLM encoder (`e5-mistral-7b-instruct`) [3]. Both models were

evaluated on all retrieval datasets, with particular attention to their performance on scientific datasets [5].

The objectives of this experiment are to determine whether task-specific fine-tuning allows the SLM to match or exceed the performance of the task-tuned LLM encoder within the scientific domain.

Metrics such as NDCG@k, MAP, and MRR were prioritized, and latency measurements were used to contextualize deployment considerations but not analyzed in depth for this study. For all experiments, NDCG@10 was used as the primary score under evaluation.

5 Key Results and Observations

The experiments provided insights into the performance of decoder-only and encoder-only architectures across a variety of retrieval tasks. This section summarizes the most critical findings, focusing on trends and anomalies observed during evaluation. Detailed numerical results and dataset-specific analyses are provided in the Appendix.

5.1 Decoder-Only Model Evaluation

The decoder-only architecture, represented by `mistralai/Mistral-7B-v0.1` [6], underperformed significantly across all retrieval tasks. This aligns with common industry knowledge that decoder-only models, optimized for autoregressive generation, struggle with retrieval tasks due to their unidirectional encoding and alternative training objectives. Key observations include:

- **Overall Performance:** Across all datasets, the decoder-only model achieved an average NDCG@10 score below 2%. Notable exceptions include SciFact and ArguAna, where performance increased to approximately 14–16%, likely due to alignment with specific pretraining data or objectives. Overall, these results show abysmal performance across all metrics and datasets when utilizing generative decoder-only architectures for retrieval tasks.
- **Layer-Wise Performance:** Embeddings from the penultimate layer consistently outperformed other layers, with the final layer performing worse than intermediate layers. This suggests that later layers, while still unsuitable for retrieval tasks, may contain more complex encodings better suited for transfer learning tasks.
- **Domain-Specific Trends:** The model performed relatively better on medical and scientific datasets compared to other domains (e.g., QA, programming). This trend may reflect pretraining biases, though further analysis of the model’s (closed source) training data is needed.

5.2 Encoder-Only SLM vs LLM Comparisons

The task-tuned LLM encoder (`e5-mistral-7b-instruct`) [3] consistently outperformed the untuned SLM (`msmarco-distilbert-base-v4`) across all retrieval tasks. These findings highlight the advantages of task-specific tuning in large models (even with generated data), particularly for complex retrieval tasks:

- **Performance Gap:** The LLM encoder demonstrated a significant performance advantage, achieving an average main-score improvement of approximately 22.3% over the untuned SLM across all dataset domains. Scientific and medical datasets showed the largest margins.
- **Fine-Tuned SLM vs Task-Tuned LLM:** Despite fine-tuning the SLM on domain-specific datasets, the LLM encoder maintained a consistent advantage, with an average performance improvement of 28.5% across all domains. In scientific datasets, the LLM achieved an NDCG@10 score of 37.16%, compared to 24.99% for the fine-tuned SLM.
- **[Note] Computational Trade-Offs:** While the LLM encoder demonstrated higher performance, latency measurements indicate that the SLM may be a more reasonable option in scenarios where embedding latency is a critical factor.

5.3 Effects of Task-Specific Fine-Tuning

While finetuning hyperparameter optimization is an active area of research, finetuning the encoder-only SLM (`msmarco-distilbert-base-v4`) yielded mixed results. Using standard tuning hyperparameters for this task [8], our fine-tuned SLM showed a notable decrease in overall performance across most domains. These findings challenge the assumption that domain-specific fine-tuning universally improves retrieval performance:

- **Overall Performance:** Fine-tuning resulted in an average NDCG@10 decrease of 6.2% across all datasets. Surprisingly, the scientific domain, which was the target of the fine-tuning, performance decreased by 5.3%.
- **Domain-Specific Trends:** General QA and programming datasets experienced the largest decreases (-9.8% and -6.7%, respectively). Medical datasets showed a slight improvement (+0.38%), though this result may fall within the bounds of statistical variability and warrants further investigation.
- **Dataset-Level Insights:** SciFact, a dataset included in the fine-tuning process, experienced a performance decrease of 4.3%. Conversely, ArguAna, which was not included in the fine-tuning data, showed an unexpected performance increase of 3.3%, suggesting potential generalization effects specific to this dataset.

These results highlight the complexity of fine-tuning strategies and underscore the need for further investigation into optimizing domain-specific training processes for retrieval tasks. Additionally, our results confirm that, while computationally less efficient, more modern LLM-based encoder-only embeddings significantly outperform their alternative counterparts for both in and out of domain retrieval tasks.

6 Conclusion

This study evaluated the performance of decoder-only and encoder-only architectures for text retrieval tasks, providing practical insights into their effectiveness across diverse datasets and task-specific fine-tuning scenarios. Several key takeaways emerged from the experiments:

- Decoder-only models, such as `mistralai/Mistral-7B-v0.1`, are poorly suited for retrieval tasks, achieving an average NDCG@10 score of less than 2% across most datasets. While intermediate and penultimate layers outperformed the final layer, our results suggest that embeddings at any stage of decoder-only models are not suited for retrieval objectives.
- Task-tuned LLM encoders (`e5-mistral-7b-instruct`) consistently outperformed both untuned and fine-tuned SLMs (`msmarco-distilbert-base-v4`). The LLM encoder demonstrated an average performance improvement of 28.5% over the fine-tuned SLM, reinforcing the value of task-specific pre-training at scale.
- Domain-specific fine-tuning of SLMs showed mixed results, with a 6.2% overall performance decrease and degradation even in the fine-tuned domain (scientific datasets). These findings highlight the importance of careful ablation studies when fine-tuning for specific retrieval tasks, as unoptimized fine-tuning can substantially degrade both generalizability and task performance.

While the findings largely confirm existing assumptions about decoder-only and encoder-only architectures, they also offer empirical validation across a wide range of datasets. This study serves as a practical guide for researchers and engineers in choosing models and tuning strategies for retrieval tasks, emphasizing the importance of balancing task performance with computational constraints.

Several limitations of this work warrant consideration. First, small performance changes, such as the slight improvement observed in medical datasets, may fall within statistical variability and should be interpreted cautiously. Second, the experiments focused on a subset of models and datasets, and other architectures or fine-tuning methods may yield different results. Finally, biases in pretraining data were not explicitly analyzed, which could explain some domain-specific trends observed during evaluation.

Future research could explore improved fine-tuning strategies, such as con-

trastive learning, alternative sampling methods, or more structured curriculum-based approaches, to enhance generalizability while improving domain-specific performance. Additionally, expanding the scope of evaluated models to include other encoder-only architectures or LLM variants could provide a more comprehensive understanding of trade-offs in retrieval-focused embeddings.

In conclusion, this study underscores the importance of architectural choice and tuning strategies in embedding models for retrieval tasks. Decoder-only models remain unsuitable for retrieval, and task-tuning must be approached with care to ensure alignment with desired task objectives. These findings provide actionable guidance for the development of efficient and effective retrieval systems.

References

- [1] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, “Roberta: A robustly optimized bert pretraining approach,” 2019. [Online]. Available: <https://arxiv.org/abs/1907.11692>
- [2] C. Li, M. Qin, S. Xiao, J. Chen, K. Luo, Y. Shao, D. Lian, and Z. Liu, “Making text embedders few-shot learners,” 2024. [Online]. Available: <https://arxiv.org/abs/2409.15700>
- [3] L. Wang, N. Yang, X. Huang, L. Yang, R. Majumder, and F. Wei, “Improving text embeddings with large language models,” 2024. [Online]. Available: <https://arxiv.org/abs/2401.00368>
- [4] N. Thakur, N. Reimers, J. Daxenberger, and I. Gurevych, “Augmented sbert: Data augmentation method for improving bi-encoders for pairwise sentence scoring tasks,” 2021. [Online]. Available: <https://arxiv.org/abs/2010.08240>
- [5] R. Friel, M. Belyi, and A. Sanyal, “Ragbench: Explainable benchmark for retrieval-augmented generation systems,” 2024. [Online]. Available: <https://arxiv.org/abs/2407.11005>
- [6] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. de las Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, L. R. Lavaud, M.-A. Lachaux, P. Stock, T. L. Scao, T. Lavril, T. Wang, T. Lacroix, and W. E. Sayed, “Mistral 7b,” 2023. [Online]. Available: <https://arxiv.org/abs/2310.06825>
- [7] N. Muennighoff, N. Tazi, L. Magne, and N. Reimers, “Mteb: Massive text embedding benchmark,” 2023. [Online]. Available: <https://arxiv.org/abs/2210.07316>
- [8] N. Reimers and I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” 2019. [Online]. Available: <https://arxiv.org/abs/1908.10084>
- [9] O. Jeunen, I. Potapov, and A. Ustimenko, “On (normalised) discounted cumulative gain as an off-policy evaluation metric for top- n recommendation,” 2024. [Online]. Available: <https://arxiv.org/abs/2307.15053>

A Appendix

A.1 Detailed Results (Aggregated)

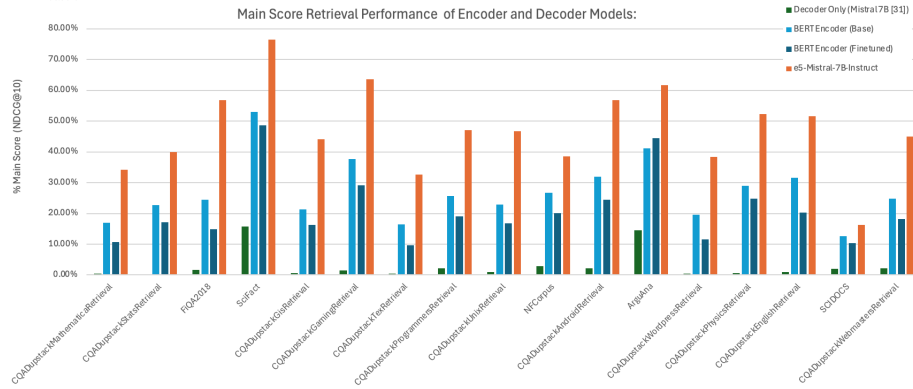


Figure 1: Dataset Specific Performance Metrics for All Models

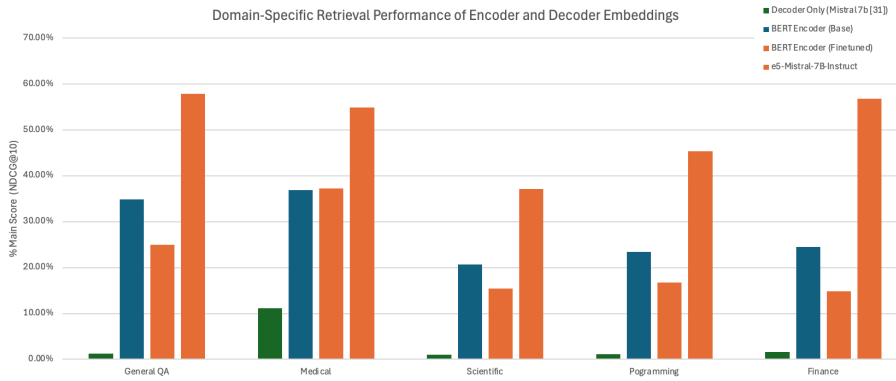


Figure 2: Domain Specific Performance Metrics for All Models

A.2 Decoder-Only Layer-Wise Performance

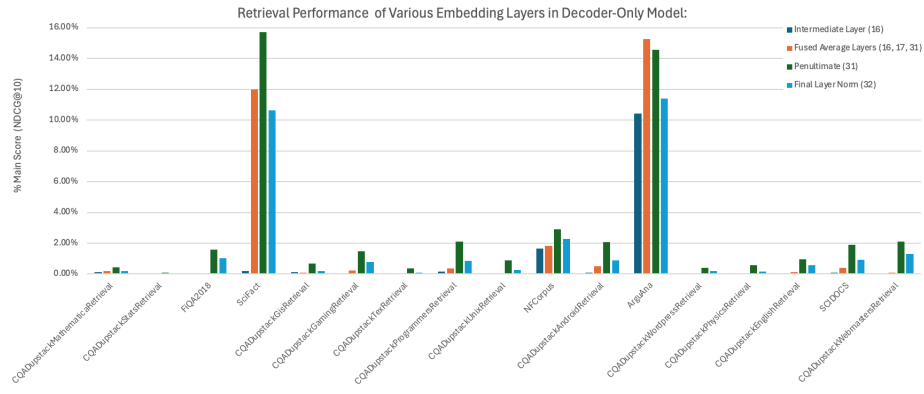


Figure 3: Dataset Specific Decoder-Only Embedding Performance by Layer

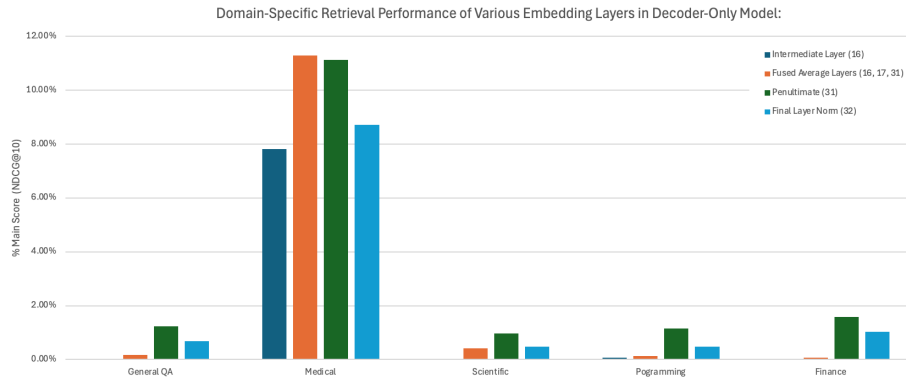


Figure 4: Domain Specific Decoder-Only Embedding Performance by Layer

A.3 Tabulated Results Across Various Metrics

A.3.1 Decoder-Only Intermediate Embedding Layer Performance

Main_NDC@10	Decoder-Only Intermediate Embedding Layer Performance										Main_NDC@10	Decoder-Only Intermediate Embedding Layer Performance								
	0.1%	0.2%	0.3%	0.4%	0.5%	0.6%	0.7%	0.8%	0.9%	1.0%		0.1%	0.2%	0.3%	0.4%	0.5%	0.6%	0.7%	0.8%	0.9%
MAP	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
MRR	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
NDCG	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Precision	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Recall	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 5: Decoder-Only Intermediate Layer [16] Performance

Main_NDC@10	Decoder-Only Fused Average Layers [16, 17, and 31] Performance										Main_NDC@10	Decoder-Only Fused Average Layers [16, 17, and 31] Performance								
	0.1%	0.2%	0.3%	0.4%	0.5%	0.6%	0.7%	0.8%	0.9%	1.0%		0.1%	0.2%	0.3%	0.4%	0.5%	0.6%	0.7%	0.8%	0.9%
MAP	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
MRR	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
NDCG	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Precision	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Recall	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 6: Decoder-Only Fused Average Layers [16, 17, and 31] Performance

Main_NDC@10	Decoder-Only Penultimate Layer [31] Performance										Main_NDC@10	Decoder-Only Penultimate Layer [31] Performance								
	0.1%	0.2%	0.3%	0.4%	0.5%	0.6%	0.7%	0.8%	0.9%	1.0%		0.1%	0.2%	0.3%	0.4%	0.5%	0.6%	0.7%	0.8%	0.9%
MAP	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
MRR	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
NDCG	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Precision	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00
Recall	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00	0.00

Figure 7: Decoder-Only Penultimate Layer [31] Performance

This figure displays a comprehensive table of performance metrics for various models across different tasks. The metrics include MAP, MRR, NDCG, Precision, and Recall. The models are grouped into Main-NDCG@10, MAP, MRR, NDCG, Precision, and Recall. The tasks are categorized into Medical, Academic, and Performance. The table shows that GPT-4 consistently achieves the highest performance across most metrics and tasks.

Figure 8: Decoder-Only Final Layer Norm [32] Performance

A.3.2 Encoder-Only Embedding Performance

This figure displays a comprehensive table of performance metrics for various models across different tasks. The metrics include MAP, MRR, NDCG, Precision, and Recall. The models are grouped into Main-NDCG@10, MAP, MRR, NDCG, Precision, and Recall. The tasks are categorized into Medical, Academic, and Performance. The table shows that GPT-4 consistently achieves the highest performance across most metrics and tasks.

Figure 9: Base DistilBERT Performance

This figure displays a comprehensive table of performance metrics for various models across different tasks. The metrics include MAP, MRR, NDCG, Precision, and Recall. The models are grouped into Main-NDCG@10, MAP, MRR, NDCG, Precision, and Recall. The tasks are categorized into Medical, Academic, and Performance. The table shows that GPT-4 consistently achieves the highest performance across most metrics and tasks.

Figure 10: Finetuned DistilBERT Performance

A.5 Finetuning Triplet Sample

Scientific Domain Sample:

Query: *Which hormone abnormalities are characteristic to Pendred syndrome?*

Positive Example:

A coherent organization of differentiation proteins is required to maintain an appropriate thyroid function in the Pendred thyroid. CONTEXT: Pendred syndrome is caused by mutations in the gene coding for pendrin, an apical Cl⁻/I⁻ exchanger. OBJECTIVE: To analyze intrathyroidal compensatory mechanisms when pendrin is lacking, we investigated the thyroid of a patient with Pendred syndrome. The expression of proteins involved in thyroid hormone synthesis, markers of oxidative stress (OS), cell proliferation, apoptosis, and antioxidant enzymes were analyzed. RESULTS: Three morphological zones were identified: nearly normal follicles with iodine-rich thyroglobulin

Negative Example:

A meta-analysis of prospective randomized controlled trials evaluating endovascular therapies for acute ischemic stroke. INTRODUCTION: A recent randomized controlled trial (RCT), the Multicenter Randomized CLinical trial of Endovascular treatment for Acute ischemic stroke in the Netherlands (MR CLEAN), demonstrated better outcomes with endovascular treatment compared with medical therapy for acute ischemic stroke (AIS). However, previous trials have provided mixed results regarding the efficacy of endovascular treatment for AIS. A meta-analysis of all available trial data was performed to summarize the available evidence. METHODS: A literature search was performed to identify all prospective RCTs comparing endovascular therapies with medical

General QA Domain Sample:

Query: *Why am I always rejected by the men I am interested in?*

Positive Example:

I am a man, just not competitively and I will tell you why I reject most of the time. I think this will actually mean something because I am homeless. I reject because I watch a lot of porn. I am 40, no kids. No money. I reject because I have no money. No car. Can't take you someplace anywhere. I would like to. Because I can talk to you. I don't have ANY gold. That's not me rejecting you. I know stuff other than rejection, projection, protection. How am I going to take care of myself?

Negative Example:

The use of microwave ovens to cook food has been a major source of debate for decades. While some people believe that microwaving food is dangerous and can cause health problems, others believe that microwaves are a safe and convenient way to quickly prepare meals. While both sides can agree that microwaving food does not provide the same nutritional value as traditional cooking methods, there is no scientific evidence that food prepared in a microwave oven is bad for your health. To understand the potential health risks associated with microwaving food, it is important to understand how microwaves work. Microwave ovens use high frequency radio waves to heat food. These waves penetrate food and cause the water molecules in it to vibrate, generating heat and cooking the food. Microwave ovens are designed to contain these waves, and the radiation emitted from them is non-ionizing, meaning it does not have enough energy to cause chemical changes in the food or the environment.

A.6 Alternative Evaluation Metrics

Supplemental Metric Definitions:

- **NDCG@k**: Normalized Discounted Cumulative Gain, measuring ranking quality by relevance and position.
- **MAP@k**: Mean Average Precision at rank k , capturing precision across top- k results.
- **MRR@k**: Mean Reciprocal Rank at k , reflecting the rank of the first relevant document.

Evaluation Framework: All models were evaluated on independent hardware using evaluation scripts sourced from the MTEB leaderboard framework [7], ensuring consistency and comparability across diverse datasets and domains.

A.7 Supplementary Operational Details

Caching Mechanism: LLM embeddings were precomputed for efficiency. Custom vector caching was created to efficiently retrieve precomputed vectors for evaluation in a "compute once, evaluate many" fashion.

Compute Resources: All experiments were conducted on dual Nvidia RTX A6000 GPUs, with total runtime for evaluations estimated at 48 hours for both LLM Encoder and Decoders (Mistral and e5), and estimated at 2 and a half hours for all SLM Encoders (BERT-based models).