

Vendi-RAG Diversity-Aware Retrieval: Efficiency and Overhead in Out-of-Domain ELI5 Queries

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 6 peer-reviewed papers addressing the following research question: What is the impact of Vendi-RAG's diversity-aware retrieval on inference efficiency and computational overhead compared to traditional BM25 and dense retrieval methods when processing out-of-domain. The advent of contextualised language models has brought gains in search effectiveness, not just when applied for re-ranking the output of classical weighting models such as BM25, but also when used directly for passage indexing and retrieval, a technique which is called dense. 18 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.7/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: On Single and Multiple Representations in Dense Passage Retrieval. Research question: What is the impact of Vendi-RAG's diversity-aware retrieval on inference efficiency and computational overhead compared to traditional BM25 and dense retrieval methods when processing out-of-domain questions in ELI5 benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 6 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.7/10.

3 Results

6 papers retrieved. 18 claims extracted; 0 independently verified. Quality review score: 3.7/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
The MSMARCO passage ranking dataset contains 8.8M passages.	×	0.07
The ANCE implementation is adapted from the provided code by its authors.	×	0.02
The ColBERT implementation is adapted from the provided code by its authors.	×	0.01
ANCE fine-tunes a RoBERTa model (specifically roberta-base).	×	0.01
ColBERT fine-tunes the bert-base-uncased BERT model.	×	0.03
Training ColBERT using bert-base-uncased for 44,500 batches achieves better performance than training ColBERT using robe	×	0.04
The RoBERTa-based ColBERT model had relative performance 25% less than the BERT-based ColBERT model (around NDCG@10 of 0	×	0.08
The ANCE document index is stored in FAISS using the uncompressed IndexFlatIP format.	×	0.02
The ColBERT document index is stored in FAISS using the compressed and quantised IndexIVFPQ format.	×	0.02
The ColBERT document index is trained on a random 5% sample of the document embeddings.	×	0.03
Mean response time for ANCE is 211ms.	×	0.06
Mean response time for ColBERT is 635ms.	×	0.06
The publicly available query sets with relevance assessments include 5000 queries sampled from the MSMARCO Dev set and t	×	0.05
The MSMARCO Dev set contains on average 1.1 judgements per query.	×	0.03
The TREC 2019 query set contains 43 queries with an average of 215.3 judgements per query.	×	0.05
There are three major types of dysarthria in cerebral palsy: spastic, dyskinetic (athetosis), and ataxic.	×	0.01
The main event that led the US to entering WW2 was Japan bombing Pearl Harbor.	×	0.00
The U.S entered WW1 for several reasons, including unlimited German submarine warfare and the Zimmermann note.	×	0.00

References

- <http://arxiv.org/abs/2205.02303v1>
- <http://arxiv.org/abs/2108.06279v2>
- <http://arxiv.org/abs/2210.05512v1>