

Manifold-Aware Distance Metrics in Dense Retrieval Across Extended Context Windows

Assignee Research

May 30, 2026

Abstract

This report synthesises findings from 16 peer-reviewed papers addressing the following research question: How does the performance of manifold-aware distance metrics in dense passage retrieval scale with increasing context window sizes beyond 512 tokens on Natural Questions and HotpotQA benchmarks. Dense Passage Retrieval (DPR) typically relies on Euclidean or cosine distance to measure query-passage relevance in embedding space, which is effective when embeddings lie on a linear manifold. However, our experiments across DPR benchmarks suggest that embeddings often lie on. 13 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: MA-DPR: Manifold-aware Distance Metrics for Dense Passage Retrieval. Research question: How does the performance of manifold-aware distance metrics in dense passage retrieval scale with increasing context window sizes beyond 512 tokens on Natural Questions and HotpotQA benchmarks?.

2 Methodology

Systematic literature search across multiple databases yielded 16 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

3 Results

16 papers retrieved. 13 claims extracted; 1 independently verified. Quality review score: 4.5/10.

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 system specifications include CPU—Intel(R) Core(TM) i7-14700HX and GPU—NVIDIA GeForce RTX 4070 Laptop GPU.	×	0.01
Average CPU utilization during measurement is $\sim 5\%$.	×	0.00
All codes and results are available online at github.com/QianfengWen/Manifold_Distance_Retrieval.git .	×	0.02
The experiments evaluate MA-DPR dManifold against several baselines including DPR with dEuclidean, DPR with dEuclidean +	×	0.08
The DPR benchmarks used are MS MARCO, NFCorpus, SciDocs, and ANTIQUE.	×	0.06
Two embedding models are used: msmarco-distilbert-base-tas-b (tas-b) and SciNCL.	×	0.04
MS MARCO is the in-distribution dataset for tas-b and SciDocs is the in-distribution dataset for SciNCL.	×	0.05
All embeddings are 2-normalized.	×	0.05
The evaluation metrics used are Recall, Mean Average Precision (MAP), and Normalized Discounted Cumulative Gain (nDCG) f	×	0.03
The research questions addressed are RQ1: Manifold Hypothesis Validation, RQ2: MA-DPR vs Baseline, RQ3: Design Choice Co	×	0.06
In Figure 2, for each ground truth relevant query q and passage p pair (orange dots) and irrelevant pair (blue dots), dE	×	0.06
In a perfectly linear embedding space, the manifold-aware distance induced by dKNN_Euclidean+cDC should closely align wi	✓	0.20
In the presence of non-linear structure, the two distances (dEuclidean and dManifold) are expected to diverge.	×	0.06

References

- <http://arxiv.org/abs/2509.13562v1>
- <http://arxiv.org/abs/2507.23334v2>
- <http://arxiv.org/abs/2605.02623v1>