

Comparison of Spell-Checking Algorithms in Dual-Encoder Retrieval Accuracy on NaturalQuestions

Assignee Research

June 11, 2026

Abstract

Dense retrieval is becoming one of the standard approaches for document and passage ranking. The dual-encoder architecture is widely adopted for scoring question-passage pairs due to its efficiency and high performance. Typically, dense retrieval models are evaluated on clean and curated datasets. However, when deployed in real-life applications, these models encounter noisy user-generated text. That said, the performance of state-of-the-art dense retrievers can substantially deteriorate when exposed to noisy text. In this work, we study the robustness of dense retrievers against typos in the

1 Introduction

This paper examines: Analysing the Robustness of Dual Encoders for Dense Retrieval Against Misspellings. Research question: How does the integration of different spell-checking algorithms (e.g., Hunspell, BERT-based, or phonetic) compare in terms of their impact on the retrieval accuracy of dual-encoder models on the NaturalQuestions benchmark?.

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 7.5/10.

3 Results

16 papers retrieved. 7 claims extracted; 5 independently verified. Quality review score: 7.5/10.

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
On clean questions, data augmentation, contrastive learning, and their combination do not harm the retrieval performance	✓	0.17
All robustification approaches (Data augm., CL, Data augm. + CL) perform significantly better than the original Dual Enc	×	0.12
The combined approach of data augmentation and contrastive learning achieves the highest performance among all tested me	×	0.11
Robustness of dual encoder models deteriorates when typos are restricted to non-stopwords or discriminative utterances c	✓	0.18
The most significant performance losses occur when typos appear on discriminative utterances (words overlapping between	✓	0.21
The proposed data augmentation combined with contrastive learning approach remains the best performing method across all	✓	0.24
There is a strong positive correlation between the frequency of typoed words in the training set and retrieval performan	✓	0.25

References

- <http://arxiv.org/abs/2108.06279v2>
- <http://arxiv.org/abs/2308.00480v1>
- <http://arxiv.org/abs/2205.02303v1>