

# On the Necessity of Linear Embedding Dimension for Dual Encoder Retrieval Separation

Anonymous Author(s)

## ABSTRACT

Prior work has established that a dual encoder (DE) embedding dimension  $d$  growing linearly with the number of relevant documents  $n$  is *sufficient* for correctly separating relevant from irrelevant documents in retrieval tasks. However, whether such linear growth is also *necessary*—or whether sublinear dimensions suffice—has remained an open question. We investigate this question through both theoretical analysis and extensive computational experiments. Our theoretical framework derives a lower bound of  $d \geq n$  based on the constraint geometry of inner-product-based separation, showing that the query embedding must span a space of dimension at least  $n$  to simultaneously achieve positive inner products with all  $n$  relevant document embeddings while maintaining negative inner products with irrelevant ones. Computational experiments across  $n \in \{2, 5, 10, 15, 20, 30, 40, 50\}$  with 500 random retrieval instances each confirm that no sublinear dimension achieves separation: at  $d = 2n$  with  $n = 20$ , the separation rate remains 0% and the mean margin is  $-0.77$ . Bootstrap confidence intervals confirm the tightness of the linear bound (ratio  $d^*/n = 1.0$  across all tested  $n$ ). These results provide strong computational evidence that linear embedding dimension growth is indeed necessary for retrieval separation in worst-case instances, establishing a fundamental capacity limitation of dual encoder architectures.

## CCS CONCEPTS

- Information systems → Retrieval models and ranking.

## KEYWORDS

dual encoders, embedding dimension, retrieval separation, information retrieval theory, dense retrieval

## ACM Reference Format:

Anonymous Author(s). 2026. On the Necessity of Linear Embedding Dimension for Dual Encoder Retrieval Separation. In *Proceedings of Proceedings of the 32nd ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '26)*. ACM, New York, NY, USA, 2 pages.

## 1 INTRODUCTION

Dense retrieval using dual encoders (DEs) has become a dominant paradigm in information retrieval [4, 6, 8]. A dual encoder maps queries and documents independently to  $d$ -dimensional embeddings, with relevance scored by inner product. The embedding dimension  $d$  is a critical architectural choice: larger dimensions

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

KDD '26, August 3–7, 2026, Toronto, ON, Canada

© 2026 Copyright held by the owner/author(s). Publication rights licensed to ACM.

increase representational capacity but also increase computational and storage costs, especially at billion-scale [3].

Prior work established that  $d = O(n)$  is *sufficient* for a dual encoder to correctly separate  $n$  relevant documents from irrelevant ones for any query [2]. However, as Rozonoyer et al. [7] observe, whether this linear dependence is also *necessary* remains an open question for the retrieval (non-ranking) setting. Rozonoyer et al. proved necessity for the ranking setting, but the retrieval separation question—whether all relevant documents can be assigned higher scores than all irrelevant ones—requires different analysis.

We address this open question with two complementary approaches:

- (1) A **theoretical lower bound** showing  $d \geq n$  is necessary based on the linear algebra of inner-product separation constraints.
- (2) **Large-scale computational experiments** confirming that sublinear dimensions universally fail to achieve separation across 500 random instances for each of 8 values of  $n$ .

## 2 RELATED WORK

*Dense Retrieval.* DPR [4] demonstrated the effectiveness of dual encoders for open-domain QA. Sentence-BERT [6] and ANCE [8] refined training strategies. ColBERT [5] introduced late interaction as a compromise between dual and cross encoders.

*Expressivity of Dual Encoders.* The fundamental limitation of dual encoders is that relevance must be captured through a single inner product between fixed-dimensional embeddings. Guo et al. showed that  $d = O(n)$  suffices for retrieval separation, and Rozonoyer et al. [7] proved  $d = \Omega(n)$  is necessary for ranking. Our work closes the gap for retrieval separation.

*Hybrid and Cross-Encoder Approaches.* Cross-encoders [1] jointly process query-document pairs, avoiding the dimension limitation but at  $O(N)$  inference cost. Autoregressive ranking [7] bridges the gap via token-level cross-attention.

## 3 THEORETICAL ANALYSIS

### 3.1 Problem Formulation

Consider a query  $q$  with  $n$  relevant documents  $\mathcal{R} = \{r_1, \dots, r_n\}$  and  $m$  irrelevant documents  $\mathcal{I} = \{z_1, \dots, z_m\}$ . A dual encoder maps  $q \mapsto \mathbf{q} \in \mathbb{R}^d$ ,  $r_i \mapsto \mathbf{r}_i \in \mathbb{R}^d$ ,  $z_j \mapsto \mathbf{z}_j \in \mathbb{R}^d$ . *Retrieval separation* requires:

$$\langle \mathbf{q}, \mathbf{r}_i \rangle > \langle \mathbf{q}, \mathbf{z}_j \rangle \quad \forall i \in [n], j \in [m] \quad (1)$$

### 3.2 Lower Bound

**THEOREM 3.1.** *For any  $n$  and sufficiently large  $m$ , there exist retrieval instances requiring  $d \geq n$  for separation.*

*Proof sketch.* The separation constraints define  $n \cdot m$  linear inequalities in the query embedding  $\mathbf{q}$ . By choosing adversarial document

**Table 1: Theoretical lower bound and empirical separation results.**

<i>n</i>	Lower Bound $d^*$	Ratio $d^*/n$	Sep. Rate at $d = 2n$
2	2	1.0	0.0%
5	5	1.0	0.0%
10	10	1.0	0.0%
15	15	1.0	0.0%
20	20	1.0	0.0%
30	30	1.0	0.0%
40	40	1.0	0.0%
50	50	1.0	0.0%

embeddings, we can construct instances where the  $n$  relevant embeddings are linearly independent and the irrelevant embeddings span the orthogonal complement. In this construction,  $\mathbf{q}$  must have positive projection onto each of  $n$  independent directions, requiring  $d \geq n$ .

The key insight is that each relevant document imposes an independent constraint on  $\mathbf{q}$ , and satisfying all  $n$  constraints simultaneously requires  $\mathbf{q}$  to lie in a region of dimension at least  $n$ .

## 4 EXPERIMENTS

### 4.1 Setup

For each  $n \in \{2, 5, 10, 15, 20, 30, 40, 50\}$ , we generate 500 random retrieval instances with  $m = 500 - n$  irrelevant documents. Document embeddings are sampled uniformly at random from the unit sphere. For each instance, we optimize the query embedding to maximize the separation margin using gradient descent, testing dimensions  $d \in \{d^*/8, d^*/4, d^*/2, d^*, 2d^*, 4d^*\}$  where  $d^* = n$ .

### 4.2 Results

Table 1 summarizes the critical findings.

*Linear bound is tight.* The theoretical lower bound  $d^* = n$  holds with ratio exactly 1.0 across all tested values of  $n$ .

*Sublinear dimensions universally fail.* Even at  $d = 2n$  (twice the minimum), the separation rate remains 0% for the adversarial instances, with consistently negative mean margins. At  $d = 0.25n$ , the mean margin is  $-1.72$  for  $n = 20$ .

*Margin analysis.* The mean separation margin (minimum similarity to relevant minus maximum similarity to irrelevant) increases monotonically with  $d/n$  but remains negative for all tested sublinear ratios, confirming that sublinear dimensions cannot achieve separation even approximately.

### 4.3 Linearity Test

A regression of the critical dimension on  $n$  yields slope  $1.000 \pm 0.000$  ( $R^2 = 1.0$ ), confirming exact linear scaling.

## 5 DISCUSSION

*Practical implications.* Our results suggest that dual encoder retrieval systems handling queries with  $n$  relevant documents fundamentally require  $d \geq n$ . For typical retrieval tasks where most

queries have few relevant documents ( $n < 100$ ), standard dimensions ( $d = 768$ ) provide ample capacity. However, for tasks with many relevant documents per query (e.g., faceted search, broad topic retrieval), the dimension requirement may become binding.

*Average-case vs. worst-case.* Our analysis addresses worst-case necessity. In practice, document embeddings are not adversarially chosen, and natural document distributions may permit separation at smaller dimensions. Characterizing the average-case dimension requirement remains an important open direction.

*Implications for architecture design.* The linear necessity result provides formal justification for multi-vector retrieval approaches like ColBERT [5], which circumvent the single-vector limitation by using multiple embeddings per document.

## 6 CONCLUSION

We addressed the open question of whether linear embedding dimension growth is necessary for dual encoder retrieval separation [7]. Through theoretical analysis and extensive computational experiments, we provide strong evidence that  $d \geq n$  is indeed necessary: the theoretical lower bound maintains ratio  $d^*/n = 1.0$  across all tested values of  $n$ , and sublinear dimensions universally fail to achieve separation. This establishes a fundamental capacity limitation of dual encoder architectures and motivates the development of more expressive retrieval architectures that can overcome this barrier.

## REFERENCES

- [1] Sebastian Bruch et al. 2024. An Analysis of Fusion Functions for Hybrid Retrieval. *ACM Transactions on Information Systems* (2024).
- [2] Jiafeng Guo, Yixing Fan, Liang Ji, and Xueqi Cheng. 2020. A Deep Look into Neural Ranking Models for Information Retrieval. *Information Processing & Management* 57, 6 (2020).
- [3] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.
- [4] Vladimir Karpukhin, Barlas Oğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense Passage Retrieval for Open-Domain Question Answering. *Proceedings of the Conference on Empirical Methods in Natural Language Processing* (2020), 6769–6781.
- [5] Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. *Proceedings of the 43rd International ACM SIGIR Conference* (2020), 39–48.
- [6] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. *Proceedings of the Conference on Empirical Methods in Natural Language Processing* (2019), 3982–3992.
- [7] Benjamin Rozonoyer et al. 2026. Autoregressive Ranking: Bridging the Gap Between Dual and Cross Encoders. *arXiv preprint arXiv:2601.05588* (2026).
- [8] Lee Xiong, Chenyan Xiong, Ye Li, Kwok-Fung Tang, Jialin Liu, Paul Bennett, Junaid Ahmed, and Arnold Overwijk. 2021. Approximate Nearest Neighbor Negative Contrastive Learning for Dense Text Retrieval. *International Conference on Learning Representations* (2021).