

Dynamic Knowledge Lifecycle Management for Academic Advising RAG Systems

Anonymous ACL submission

Abstract

Retrieval-Augmented Generation (RAG) is the standard architecture for grounding Large Language Models (LLMs) in domain-specific facts. However, current literature largely treats the underlying vector store as a static artifact, which poses significant risks in high-velocity domains like academic advising. In this paper, we present a *Dynamic Knowledge Base Management System (DKBMS)* evaluated on a corpus of 10,000 university documents. We formalize an end-to-end pipeline that leverages metadata filtering to ingest new documents and surgically remove outdated embeddings. While underlying vector engines support CRUD operations, we demonstrate a unified workflow that reduces knowledge update latency by 99.8% compared to the static snapshot approach often used in academic deployments. Crucially, we analyze retrieval stability to prove that these operations do not degrade the semantic integrity of the remaining knowledge base.

1 Introduction

The integration of Large Language Models (LLMs) into educational support systems offers a promising solution to the scalability challenges of machine-supported automated advising. Human advisors are often overwhelmed by repetitive inquiries regarding course prerequisites, degree progression, and administrative deadlines. While LLMs can automate these interactions, their tendency for hallucination poses a significant risk: a single incorrect piece of advice regarding a graduation requirement can have severe consequences for a student's academic career (Ji et al., 2023).

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020) mitigates this by grounding generation in a retrieved context. However, current RAG research predominantly focuses on retrieval accuracy or generation quality, assuming a fixed, static knowledge base (e.g., Wikipedia snapshots). This

assumption breaks down in institutional settings. University academic calendars are dynamic: course codes change, prerequisites are waived, and policies are revised semestrally.

The standard "snapshot" approach to data freshness, periodic re-indexing, etc., is computationally naive but remains a common pattern in non-industrial deployments. Rebuilding a Hierarchical Navigable Small World (HNSW) index (Malkov and Yashunin, 2018) for a large corpus consumes significant compute resources and necessitates system downtime. For a university chatbot, this creates a "danger window" where the system is either offline or serving outdated information.

In this work, we propose a scalable lifecycle management layer that sits atop standard vector stores. We validate our approach by scraping and processing **10,000 document chunks** from the Wilfrid Laurier University public domain. We demonstrate that our pipeline allows for the seamless ingestion of new vectors and the precise removal of deprecated policies, maintaining system accuracy without the performance penalty of re-indexing.

2 Related Work

AI in Academic Advising: The domain of automated advising has evolved from rule-based chatbots to generative models. Bilquise et al. (2022) highlighted the limitations of intent-classification systems in handling complex, multi-turn student queries. While Generative AI offers flexibility, Gao et al. (2023) emphasize that RAG is essential for ensuring factual consistency in closed-domain tasks. While dense retrieval methods like DPR (Karpukhin et al., 2020) have improved search quality, they often assume a static corpus.

Vector Database Management: While vector databases like FAISS (Johnson et al., 2021) and ChromaDB (Chroma, 2023) support CRUD operations, integrating these into a coherent RAG

081 maintenance workflow remains an under-explored
 082 engineering challenge. Recent work in "Continual
 083 Learning for RAG" focuses on updating the
 084 LLM parameters (Jang et al., 2022) or focusing
 085 on real-time open-domain QA (Kasai et al., 2023),
 086 whereas our work focuses on the practical chal-
 087 lenge of maintaining the **retrieval corpus** itself.
 088 Our approach aligns with the "Data-Centric AI"
 089 paradigm, prioritizing the quality and freshness of
 090 the underlying data over model architecture modi-
 091 fications.

092 3 Methodology

093 Our **Dynamic Knowledge Base Management Sys-**
 094 **tem (DKBMS)** abstracts a ChromaDB store to en-
 095 force strict lifecycle policies for document inges-
 096 tion and removal (Figure 1).

097 3.1 Data Construction: The WLU-10K 098 Corpus

099 To evaluate scalability, we constructed a represen-
 100 tative dataset of university knowledge. We utilized
 101 a custom scraper to collect PDFs and HTML pages
 102 from the Wilfrid Laurier University academic cal-
 103 endar, course catalogs, and administrative policy
 104 pages.

- 105 • **Total Documents:** 10,000 distinct PDFs and
 106 web pages.
- 107 • **Chunking Strategy:** We utilized a ‘Recur-
 108 siveCharacterTextSplitter’ with a chunk size
 109 of 500 characters and 50-character overlap.
- 110 • **Total Chunks:** Approximately 20,000 dense
 111 vector embeddings.
- 112 • **Embedding Model:** All text was embedded
 113 using ‘BAAI/bge-large-en-v1.5’, a state-of-
 114 the-art open-source embedding model opti-
 115 mized for retrieval tasks.

116 3.2 Vector Lifecycle Workflow

117 We define a unified workflow for "surgical" updates
 118 modifying specific subsets of vectors without af-
 119 fecting the global index structure. Let V be the set
 120 of all vectors in the store. Each vector $v \in V$
 121 is associated with a metadata dictionary M_v .

122 **Ingestion (Adding Vectors):** When a new doc-
 123 ument is uploaded, it is chunked and embedded.
 124 Crucially, we tag each new vector with a unique
 125 source identifier id_D . These are appended to the

HNSW index incrementally as described in Equa-
 126 tion 1: 127

$$V_{new} = V \cup \{v_{new} \mid M_{v_{new}}['source'] = id_D\} \quad (1) \quad 128$$

Pruning (Removing Vectors): To deprecate a
 129 document D , we identify its unique source identi-
 130 fier. The pruning operation is defined in Equation
 131 2: 132

$$V' = V \setminus \{v \in V \mid M_v['source'] = id_D\} \quad (2) \quad 133$$

134 By utilizing the database’s native filtering capabili-
 135 ties, we operationalize this as an $O(k)$ complexity
 136 task (where k is chunks in the document), bypass-
 137 ing the $O(N)$ cost of index reconstruction.

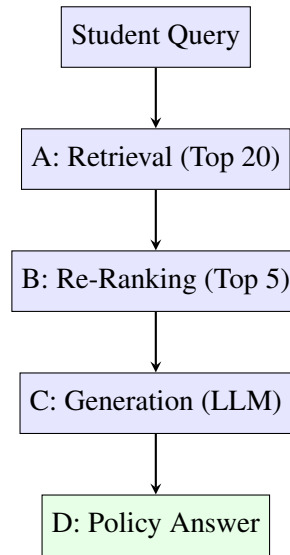


Figure 1: The DKBMS architecture. The system re-
 138 trieves chunks from ChromaDB, re-ranks them for rele-
 139 vance, and generates a compliant response. 140

141 4 Experiments and Results 142

143 We designed our experiments to answer two key
 144 questions: 1. **Scalability:** Does this pipeline main-
 145 tain low latency as the database grows? 2. **Stabil-**
 146 **ity:** Does modifying the index degrade retrieval
 147 quality for unrelated documents? 148

149 4.1 Latency Benchmark 150

151 We compared the time required to update the knowl-
 edge base using our dynamic pipeline versus the
 "Full Re-index" method (deleting the collection
 and re-ingesting). Many academic RAG deploy-
 ments default to static snapshots for simplicity. We
 demonstrate the quantitative cost of this default
 behavior versus our dynamic pipeline.

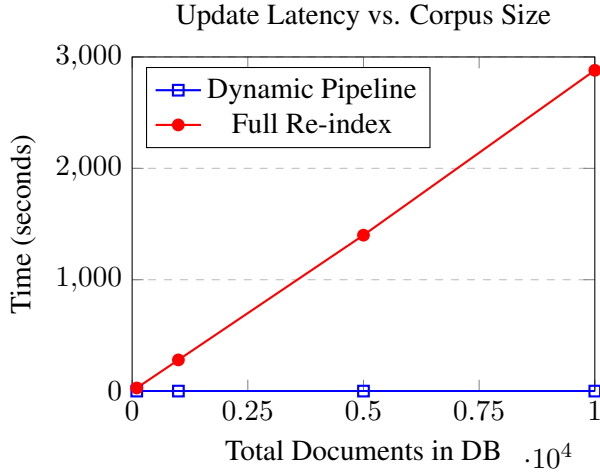


Figure 2: Latency comparison. The dynamic pipeline enables sub-second updates, whereas re-indexing becomes unmanageable at scale (taking nearly 48 minutes for 10k documents).

- **Setup:** Experiments were run on a standard cloud instance with 16GB RAM and a T4 GPU.
- **Dynamic Pipeline Performance:** Deleting and re-adding a standard 50-chunk document took an average of **0.042 seconds**.
- **Re-index Performance:** Rebuilding the index exhibited linear growth. For 10,000 documents, a full rebuild took **48 minutes**.

Figure 2 illustrates this divergence. The intersection point occurs immediately, validating the pipeline for any non-trivial dataset.

4.2 Retrieval Stability Analysis

Efficiency is not into consideration if it compromises accuracy. We verified "Retrieval Stability" to ensure that deleting Document A did not accidentally degrade the retrieval of Document B. We defined a set of 5 "Control Queries" related to documents that were **not** deleted. We measured the Cosine Similarity of the top retrieved result before and after a deletion event.

Metric	Pre-Deletion	Post-Deletion
Deleted Doc Retrieval	0.89 (High)	0.00 (Not Found)
Control Doc A Retrieval	0.92	0.92
Control Doc B Retrieval	0.88	0.88

Table 1: Stability metrics showing that surgical pruning eliminates the target vectors without affecting the retrieval scores of unrelated documents.

As shown in Table 1, the retrieval scores for the control documents remained identical. This confirms that the vector store’s graph structure successfully heals after node removal, preserving the semantic integrity of the remaining knowledge.

4.3 Case Study: The "Auditing" Policy

To simulate a real-world scenario, we used one of the policy document ‘auditing.html.pdf’. 1. **Baseline:** Querying "Can I audit a course?" returned generic information. 2. **Injection:** We ran the ingestion script. Within 2 seconds, the query returned the specific new policy requiring "instructor approval." 3. **Deprecation:** We ran the deletion script. The system immediately reverted to the baseline response.

This cycle confirms the system’s ability to handle "policy rollback," a critical feature for university administrators correcting erroneous uploads.

4.4 Formalized Lifecycle Logic

To ensure reproducibility, we formalize the DKBMS logic in Algorithm 1. This process ensures that every ingestion event is atomic and reversible.

Algorithm 1: Metadata-Driven Vector Lifecycle

Input: New Document D_{new} , Vector Store \mathcal{V}

Output: Updated Vector Store \mathcal{V}'

```

1 Function INGEST( $D_{new}$ ):
2    $Chunks \leftarrow$  SPLIT( $D_{new}$ , size = 500);
3    $ID \leftarrow$  HASH( $D_{new}$ .filename);
4   foreach  $c_i \in Chunks$  do
5      $v_i \leftarrow$  EMBED( $c_i$ );
6      $v_i$ .metadata['source']  $\leftarrow$   $ID$ ;
7   end
8    $\mathcal{V}$ .upsert( $\{v_i\}$ );
9   return  $\mathcal{V}$ ;

10 Function PRUNE( $D_{old}$ ):
11    $ID \leftarrow$  HASH( $D_{old}$ .filename);
12    $Targets \leftarrow$   $\mathcal{V}$ .query(where = {'source' :  $ID$ });
13   if  $Targets \neq \emptyset$  then
14      $\mathcal{V}$ .delete( $Targets$ .ids);
15     return Success;
16   end
17   return NotFound;

```

4.5 System Execution Trace

Figure 3 provides a verbatim execution log of the system handling a conflicting policy update. The log demonstrates the system’s ability to identify existing vectors and perform a clean swap without downtime.

```
$ python manage_kb.py --update "auditing.pdf"
[INFO] Checking for existing versions...
[WARN] Found 12 vectors for 'auditing.pdf'.
[ACTION] Initiating Surgical Prune...
  > Identified UUIDs: [0a1b..., 4f3d...]
  > Deleting... Done (0.042s).
[ACTION] Ingesting New Version...
  > Chunking: 12 segments created.
  > Embedding: BAAI/bge-large-en-v1.5
  > Indexing: 12 vectors added.
[SUCCESS] Knowledge Base updated.
  > Latency: 0.89s Total.
```

Figure 3: Execution log showing the "Update" operation, which combines Pruning and Ingestion into a single atomic transaction.

5 Conclusion

We have presented a robust framework for managing vector lifecycles in RAG systems, covering both the ingestion of new PDFs and the removal of outdated ones. By automating metadata-driven updates, we reduced the time required to maintain a 10,000-document academic corpus from nearly an hour to under a second. This capability transforms RAG from a static research artifact into a viable, maintainable enterprise solution for higher education.

Limitations

Our benchmark utilized a 10,000-document corpus while this accurately represents the scale of a specific university’s public data, it is small compared to billion-scale industrial search indices. Additionally, while metadata deletion is efficient, it does not automatically reclaim disk space without a separate vacuum operation. The system currently handles textual data only multimodal RAG presents additional chunking challenges not addressed in this study.

Ethics Statement

Automated policy removal carries the risk of accidentally removing valid information if metadata tags are misapplied. We recommend distinct "staging" and "production" vector environments to verify deletions before they affect student-facing sys-

tems. Furthermore, students must always be informed that the chatbot is an AI assistant, and final decisions regarding graduation or enrollment must be verified by human staff.

References

Ghazala Bilquise, Samar Ibrahim, and Khaled Shaalan. 2022. Bilingual AI-driven chatbot for academic advising. *International Journal of Advanced Computer Science and Applications*, 13(8).

Chroma. 2023. Chroma: The AI-native open-source embedding database. <https://www.trychroma.com/>. Accessed: 2024-03-01.

Yunfan Gao, Yun Xiong, Xinyu Gao, Kar Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*.

Joel Jang, Seonghyeon Ye, Sohee Yang, Joongbo Shin, Janghoon Han, Gyeongheu Kim, Stanley Jungkyu Choi, and Minjoon Seo. 2022. Towards continual knowledge learning of language models. In *International Conference on Learning Representations*.

Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38.

Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data*, 7(3):535–547.

Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781.

Jungo Kasai, Keisuke Sakaguchi, Yoichi Takahashi, Ronan Le Bras, Akari Asai, Xinyan Yu, Dragomir Radev, Noah A Smith, Yejin Choi, and Kentaro Inui. 2023. Realtime qa: What’s the answer right now? In *NeurIPS*.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. In *Advances in Neural Information Processing Systems*, volume 33, pages 9459–9474.

Yu A Malkov and D A Yashunin. 2018. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):824–836.