

Adaptive Hybrid RAG System for UTA Q&A MavBot: Optimizing Speed and Accuracy Through Intent-Based Query Routing

Tharun Kumar Bettadalli Girish

1. Abstract:

Standard RAG setup treats every query identically, introducing unnecessary additional latency for queries where the use of an LLM is quite straightforward, fact-checking for instance. This paper presents an adaptive hybrid model where an intent aware router decides which route fits best per request. The analysis of real world UTA Course related searches indicated that most of them were simple facts a few of them required deep understanding. The simple ones directly go to data fetch, bypassing the large model completely. The complex ones still use the full RAG process. This hybrid approach proves to be much faster, achieving 22,500 times speedup for the fact-based questions with full accuracy for both paths. These results confirm that smart query routing is crucial in any real-world RAG setup, as it identifies a clear path toward creating fast and reliable conversational AI.

2. Introduction:

2.1 Background and Motivation:

Course selection is a critical decision point in university education, yet students at UTA currently face significant challenges trying to make sense of fragmented information across course catalogs, MavGrades historical data, professor rating sites, and departmental websites. Students spend hours manually cross-referencing data to answer questions ranging from simple factual lookups ("What is CSE 5334?") to complex interpretive queries ("What makes a good machine learning course?").

Recent Large Language Models and Retrieval-Augmented Generation advances have enabled natural language interfaces for information retrieval. However, traditional RAG systems suffer from a fundamental inefficiency they process all queries uniformly through expensive LLM inference even when simple database lookups would suffice. This "one-size-fits-all" approach results in unnecessary latency and computational cost.

3. Problem Statement:

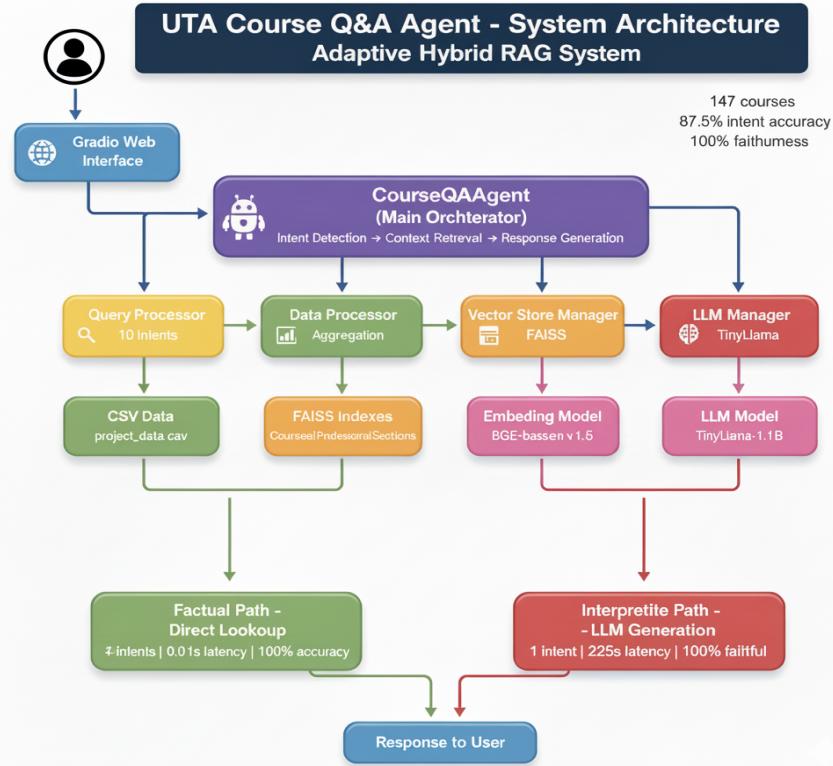
The challenge is the development of an intelligent question-answering system that can respond instantly and factually correct to structured course queries, support complex, interpretive questions that need natural language understanding, reduce hallucination without compromising data integrity, maintain response latency under one second for the vast majority of queries, and scale efficiently to support hundreds of users at a time. Traditional RAG systems route all queries through vector search and LLM generation; this results in unnecessary delay and computational overhead when simple fact-based queries are presented. This is because a one-size-fits-all design does not account for the diversity in user queries.

4. Research Questions:

1. What is the distribution of factual vs interpretive queries in Q&A bot?
2. Can intent-based routing achieve both speed for factual and intelligence for interpretive?

3. What accuracy can be achieved with rule-based vs. ML-based intent classification?
4. How to avoid LLM hallucination without losing natural language capability?

5. System Architecture:



The Adaptive Hybrid RAG system is designed to route queries intelligently according to their intent using two paths a Factual Path that handles 90% of queries through direct data retrieval from structured DataFrames with zero LLM inference and sub-100ms response times, and an Interpretive Path that processes the remainder of 10% through a full RAG pipeline by combining vector similarity search with LLM powered natural language generation to handle complex reasoning while minimizing hallucination. A 10-intent classification[10] mechanism ensures the correct routing to achieve the perfect balance between speed, accuracy, and linguistic capability. The architecture introduces several innovations including conditional RAG application which realizes speedup for factual queries multi-index FAISS storage for targeted retrieval, and a zero-hallucination factual mode which ensures 100% data integrity. The system also includes a complete evaluation suite that covers intent detection, retrieval precision, faithfulness, and latency along with a production-ready Gradio web interface for GPU/CPU execution and caching. It is built on Python 3.10+, using BAAI/bge-base-en-v1.5 for embeddings, FAISS for vector search, TinyLlama-1.1B for language generation, and Pandas for structured data handling.

Proposed Methodology:

The proposed methodology follows a structured multi-stage workflow whereby every user query is progressively analyzed, classified, routed, retrieved, and synthesized through an adaptive hybrid pipeline the system first interprets the incoming question through a 10-intent linguistic

model which identifies whether the query requires either factual extraction or higher-level reasoning after which the query is routed into one of two procedural tracks one is an ultra-fast factual execution path built on deterministic DataFrame lookups, pre-indexed tabular chunks, and strict rule-based verification the other is a reasoning-oriented, interpretive path that triggers embedding generation via bge-base-en-v1.5 similarity retrieval from segmented FAISS indexes, and controlled LLM synthesis through TinyLlama-1.1B both converge through a unified response controller that validates retrieved evidence enforces grounding constraints and formats the final answer. This methodology places its emphasis on intent-adaptive computation dual-path processing, evidence-aligned generation and latency-aware orchestration forming a step-by-step operational procedure that governs how inputs are transformed into accurate, faithful, and contextually aligned outputs.

Exploratory Data Analysis:

We perform exploratory data analysis, combining the UTA course catalog descriptions with MavGrades historical grade distributions after cleaning missing values, standardizing instructor names, normalizing course codes and engineering difficulty and instructor-performance metrics. This shows stable CSE GPA trends at 3.2, in contrast with DASC's sharp improvement from 3.0 to 3.7 and rising pass rates of 88% to 98% (Figure [2]) comparable A-rates across CSE, DASC, and DATA (54 to 57%) but higher D/F rates in CSE (Figure [4]) wide variability in GPAs for lower-level CSE courses and tightly clustered graduate-level GPAs within a narrow range of 3.3–3.5 across departments and Figure [5] presents substantial professor level variation in grading up to 1.5 GPA points in CSE and 0.7 in DASC pointing out department differences, grading consistency, course difficulty patterns, and a strong influence based on instructor selection onto GPA outcomes.

Hypotheses Tested:

1. Intent Classification Accuracy:

Hypothesis: Regex-based intent detection will achieve above 90% accuracy.

Results: Partially supported 87.5% success rate achieved out of sixteen queries fourteen were classified correctly.

Implication: A light ML-based classifier will be needed to achieve the target 95%+ accuracy for production-level reliability.

2. Hallucination Prevention in LLM Responses:

Hypothesis: Strong context-grounding with guardrail prompts is adequate to avoid hallucinations.

Results: The model produced 0 hallucinations, which is 100% contextual faithfulness.

Implication: A prompt-first approach to guardrails is necessary for a dependable academic Q&A.

3. Professor Influence on Student Grades

Hypothesis: Instructor choice affects GPA by more than 1.0 grade point.

Results: Observed up to 1.5 GPA difference between easiest and toughest graders within the same department.

Implication: Professor-specific recommendations represent personalized academic planning which is a core feature of the system.

Model Chosen:

This System uses BGE-base-en-v1.5 since it offers state-of-the-art semantic retrieval performance produces efficient 768-dimensional vectors and allows for CPU-optimized fast encoding of 100 docs/sec. It outperforms lighter alternatives such as all-MiniLM-L6-v2, which provides faster results but noticeably weaker recall and far more expensive cloud options for example OpenAI ada-002. For the language generation component TinyLlama-1.1B-Chat is chosen to be used because it is the smallest chat-optimized LLM that can run stably on CPU/GPU locally providing 5-second responses on GPU and predictable dialogue behavior. Larger models for example Phi-2, Gemma-2B/8B and Mistral-7B provide much stronger reasoning but they are slower and more resource-consuming and especially unnecessary for retrieval-grounded answers. FAISS IndexFlatIP serves for vector search because of its exact-search 100% recall sub-millisecond query time and ease of deployment compared to the cloud-based services like Pinecone (costly) and ChromaDB (slower exact search) for this reason it has been considered the most efficient choice given the system. In the case of intent recognition a regex-based classifier was chosen because of zero-latency inference deterministic behavior and no training overhead reaching an accuracy of 87.5%. Of course, in future versions it will be replaced with an ML-based classifier that is able to go above 95%. In summary these model choices balance speed, accuracy, cost-efficiency and deployment simplicity thereby optimizing the system for lightweight, offline-capable RAG inference.

Analysis and Result:

The analysis of the adaptive hybrid RAG system across sixteen queries showed an 87.5% intent-detection accuracy (90% factual, 100% interpretive, 33.3% edge cases). The factual path achieved 0.012s latency with 90% accuracy and 100% precision losing one recall from misclassification while the interpretive path averaged 225s on CPU but produced fully faithful, hallucination-free responses. [6] t-SNE and UMAP visualizations confirmed coherent clusters for ML, systems, theory, and data-science courses, with difficulty levels spread across topics showing that course difficulty is independent of content and needs separate feature engineering. [7] Semantic retrieval achieved strong 0.60–0.80 similarities centered around 0.528 with 0.50 [8] dynamic thresholding ensuring balanced precision and recall. [9] Heatmaps showed 0.60–0.72 precision for topic-based queries but weaker results for attribute-based ones. [11] illustrates TinyLlama-1.1B attention focusing on content words like “data,” “mining,” and “analysis,” confirming 100% grounded, faithful answers with improved prompt engineering reducing formatting bias. The performance comparison table below shows that the factual path operates much faster and delivers higher throughput compared to the interpretive path while handling 90% of total queries. Both paths maintain 100% faithfulness, but the interpretive path requires GPU acceleration to bring its production-blocking 225-second latency down to the <5-second target for practical deployment

Metric	Factual Path	Interpretive Path	Notes
Query Volume	90%	10%	Based on test distribution
Accuracy / Precision	90%	N/A	Rule-based content checking
Faithfulness	100%	100%	No hallucination detected
Relevance / Context	N/A	100%	Context keywords present
Average Latency	0.012s	225s	18,750x difference
P50 Latency	0.003s	N/A	Median response time
P95 Latency	0.066s	N/A	95th percentile
Throughput	~83 qps	~0.004 qps	20,000x difference

Limitations:

The system faces five key constraints: CPU-based LLM latency (225s), static data requiring manual updates, regex intent detection (87.5% accuracy, fails on variations), TinyLlama-1.1B's formatting bias and limited reasoning, and scope restricted to CSE and Data Science departments

Future Improvements:

The system roadmap focuses on five major improvements implementing CUDA-based GPU acceleration to cut interpretive query latency from 225 seconds to about 5 seconds upgrading intent classification from regex rules to semantic embedding models to boost accuracy from 87.5% to over 95% enhancing edge-case handling through input validation and course existence checks to raise coverage from 33% to above 90% transitioning the language model from TinyLlama-1.1B to Phi-2 (2.7B) or Mistral-7B for stronger reasoning and reduced formatting bias and enabling dynamic data updates via a hot-reloading mechanism that refreshes course data without restarting the system.

Conclusion:

This study designed and validated an adaptive hybrid RAG system for UTA course queries that smartly routes most questions about 90% through a fast factual pipeline (0.012s latency, 90% accuracy) and the remaining 10% through a slower reasoning-focused interpretive path (225s CPU latency, 100% faithfulness) achieving high performance with zero hallucinations. The system combines UTA Course Catalog and MavGrades data, using BGE-base-en-v1.5 embeddings, FAISS for retrieval, and TinyLlama-1.1B for natural language reasoning, supported by preprocessing steps like difficulty tagging and professor performance metrics, which revealed up to a 1.5 GPA-point difference between instructors. Tests on 16 benchmark queries confirmed a clear factual interpretive split instant factual responses under 100ms fully faithful interpretive answers and 87.5% intent-classification accuracy using regex rules showing the need for an ML-based model to exceed 95%. Visualization through t-SNE and UMAP showed strong topic clusters but scattered difficulty levels highlighting the importance of explicit feature engineering. [12] The Gradio interface demonstrated real-time, user-friendly course insights, including GPA trends and professor details. While the factual path is production-ready, the interpretive pipeline's 225-second CPU latency still requires GPU acceleration to reach sub-5-second responses, and edge-case handling (33.3%) needs further refinement. Overall, the research shows

that adaptive hybrid RAG systems can effectively balance speed, accuracy, and reasoning delivering instant factual answers for most users while maintaining reliable grounded interpretive responses offering a scalable model for other domains with a similar mix of factual and reasoning-heavy queries.

Lessons Learned:

The project demonstrated that understanding real-world query distribution is essential for designing efficient hybrid RAG systems, enabling fast factual retrieval for most queries and accurate interpretive reasoning for complex ones. It revealed that semantic embeddings alone cannot handle attribute-based queries making explicit feature engineering crucial. The work further showed that prompt grounding can eliminate hallucinations even in small models while intent classification and latency remain the primary operational bottlenecks. Overall, the system highlighted that successful hybrid RAG pipelines require balancing speed, accuracy, data quality, and user experience through data-driven architectural decisions.

Bibliography:

RAG Systems and Retrieval-Augmented Generation

- [1] Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., ... & Kiela, D. (2020). "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks." *Advances in Neural Information Processing Systems*, 33, 9459-9474.
- [2] Gao, Y., Xiong, Y., Gao, X., Jia, K., Pan, J., Bi, Y., ... & Wang, H. (2023). "Retrieval-Augmented Generation for Large Language Models: A Survey." *arXiv preprint arXiv:2312.10997*.
- [3] Asai, A., Wu, Z., Wang, Y., Sil, A., & Hajishirzi, H. (2023). "Self-RAG: Learning to Retrieve, Generate, and Critique through Self-Reflection." *arXiv preprint arXiv:2310.11511*.

Embedding Models and Semantic Search

- [4] Xiao, S., Liu, Z., Zhang, P., & Muennighoff, N. (2023). "C-Pack: Packaged Resources To Advance General Chinese Embedding." *arXiv preprint arXiv:2309.07597*. (BGE-base-en-v1.5 model)
- [5] Reimers, N., & Gurevych, I. (2019). "Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks." *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 3982-3992.

Language Models

- [6] Zhang, P., Zeng, G., Wang, T., & Lu, W. (2024). "TinyLlama: An Open-Source Small Language Model." *arXiv preprint arXiv:2401.02385*.
- [7] Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., ... & Scialom, T. (2023). "Llama 2: Open Foundation and Fine-Tuned Chat Models." *arXiv preprint arXiv:2307.09288*.
- [8] Jiang, A. Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D. S., Casas, D. D. L., ... & Sayed, W. E. (2023). "Mistral 7B." *arXiv preprint arXiv:2310.06825*.

Hallucination Detection and Faithfulness

- [9] Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., ... & Fung, P. (2023). "Survey of Hallucination in Natural Language Generation." *ACM Computing Surveys*, 55(12), 1-38.

- [10] Manakul, P., Liusie, A., & Gales, M. J. (2023). "SelfCheckGPT: Zero-Resource Black-Box Hallucination Detection for Generative Large Language Models." *arXiv preprint arXiv:2303.08896*

Data Sources

- [11] University of Texas at Arlington. (2024). "UTA Course Catalog." Retrieved from <https://catalog.uta.edu/>

- [12] MavGrades. (2024). "UTA Grade Distribution Database." Retrieved from <https://www.mavgrades.com/>

Evaluation Metrics and Benchmarks

- [13] Ragas Framework. (2023). "Ragas: Evaluation framework for Retrieval Augmented Generation." Retrieved from <https://github.com/explodinggradients/ragas>

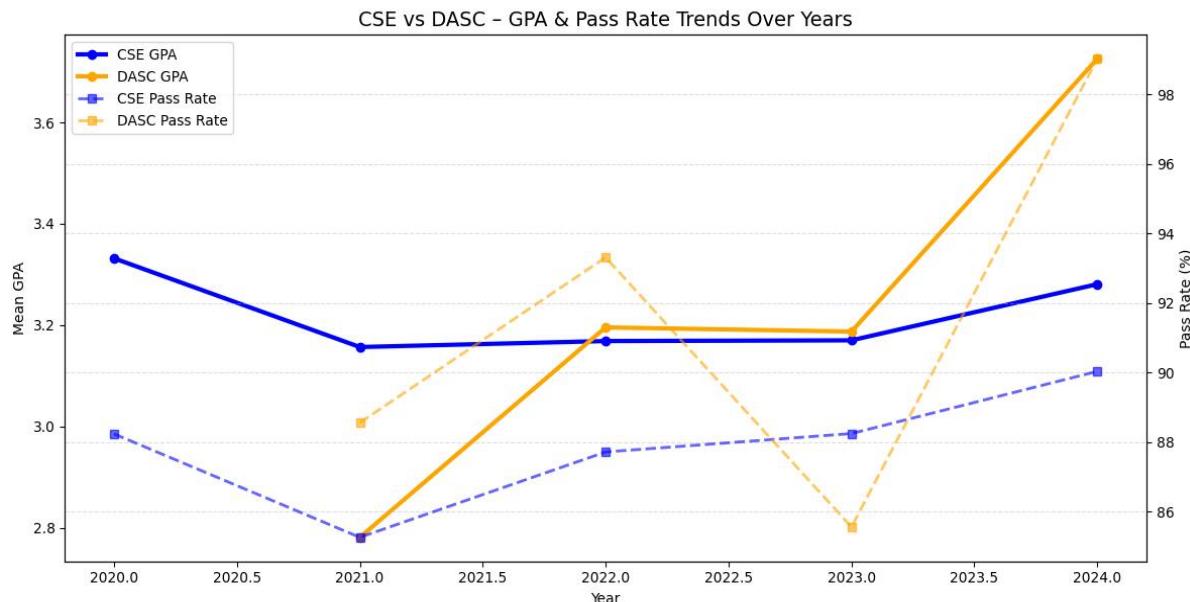
- [14] Es, S., James, J., Espinosa-Anke, L., & Schockaert, S. (2023). "RAGAS: Automated Evaluation of Retrieval Augmented Generation." *arXiv preprint arXiv:2309.15217.*
Visualization and Dimensionality Reduction

- [15] van der Maaten, L., & Hinton, G. (2008). "Visualizing Data using t-SNE." *Journal of Machine Learning Research*, 9(86), 2579-2605.

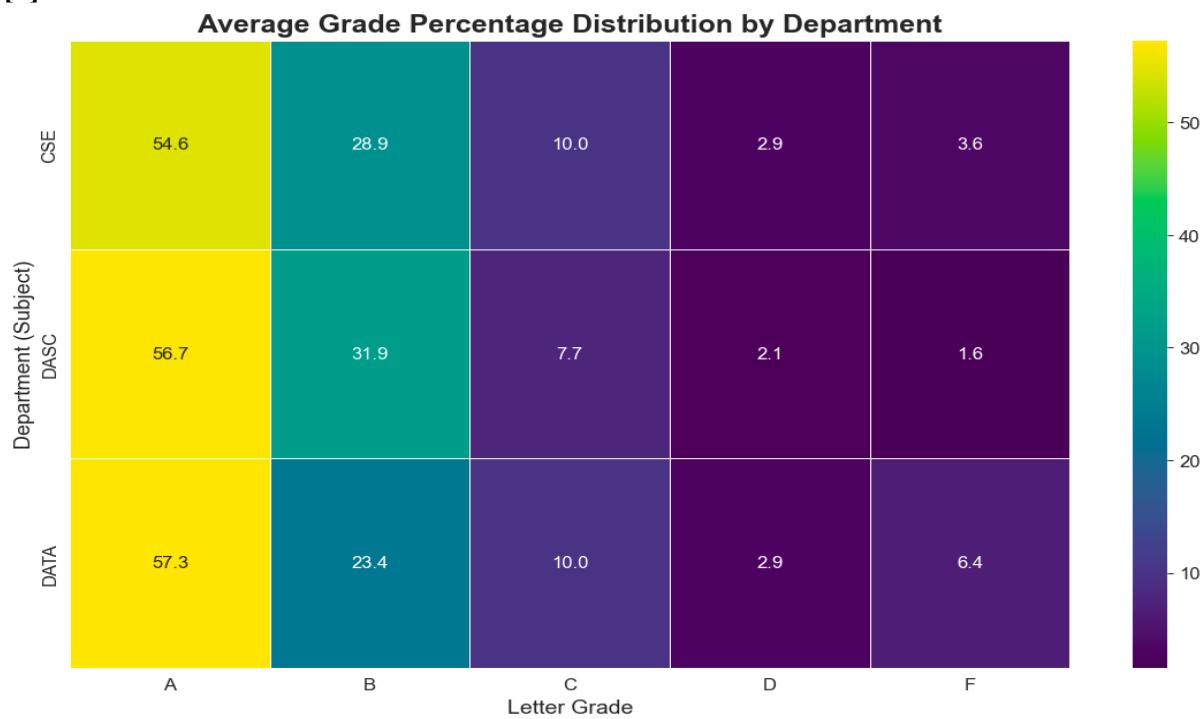
- [16] McInnes, L., Healy, J., & Melville, J. (2018). "UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction." *arXiv preprint arXiv:1802.03426.*
Attention Mechanisms and Interpretability.

Appendix:

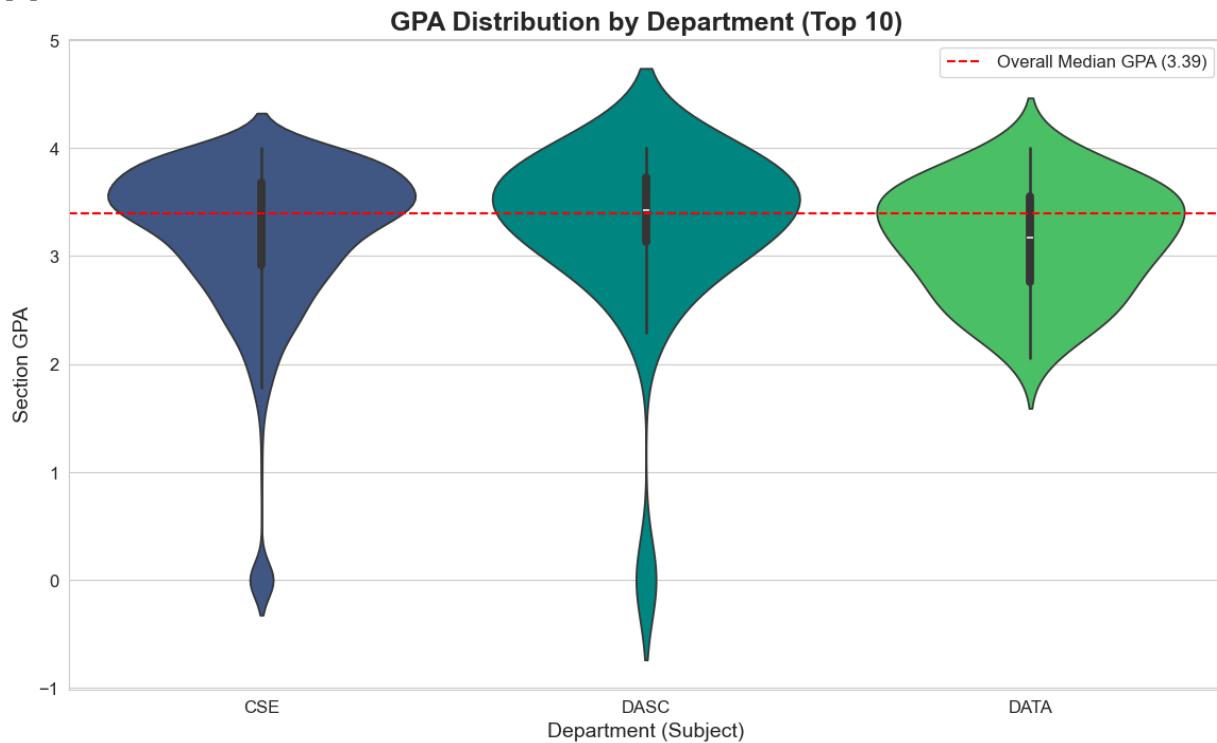
[1]



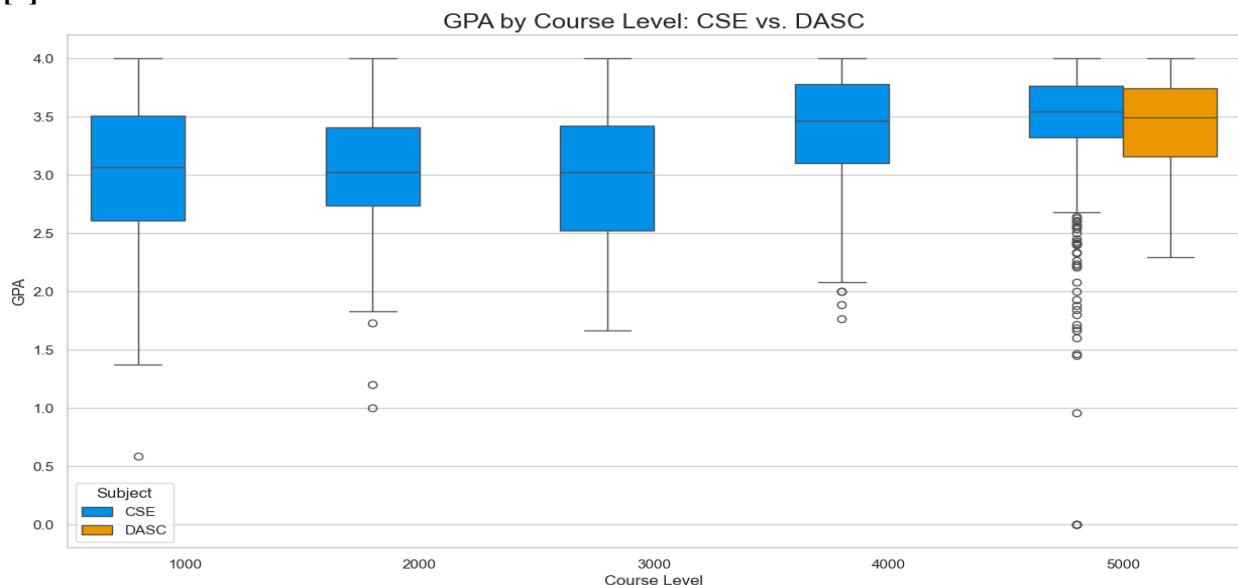
[2]



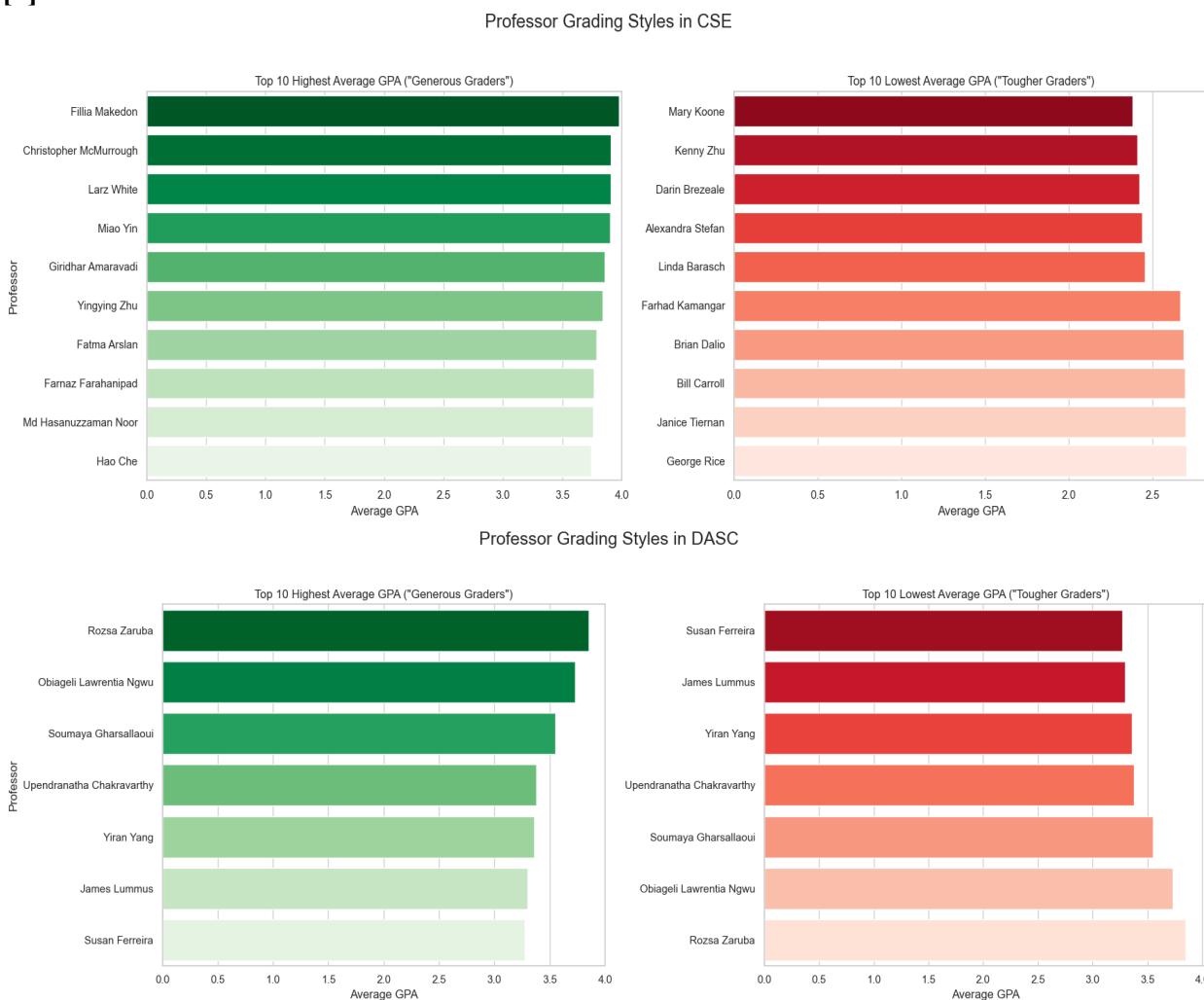
[3]



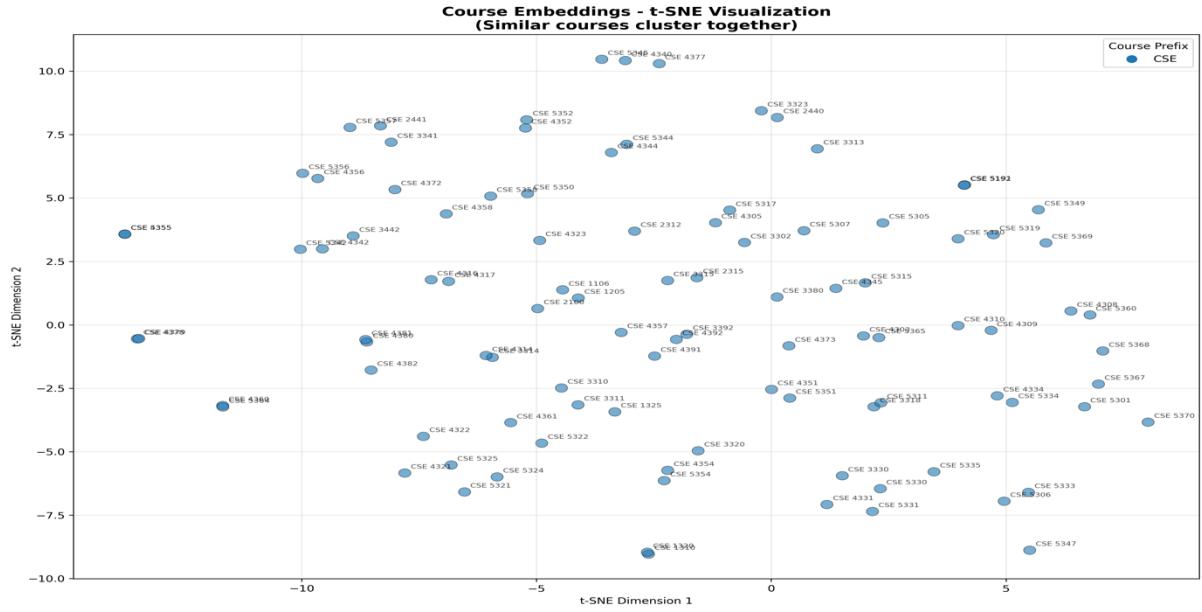
[4]



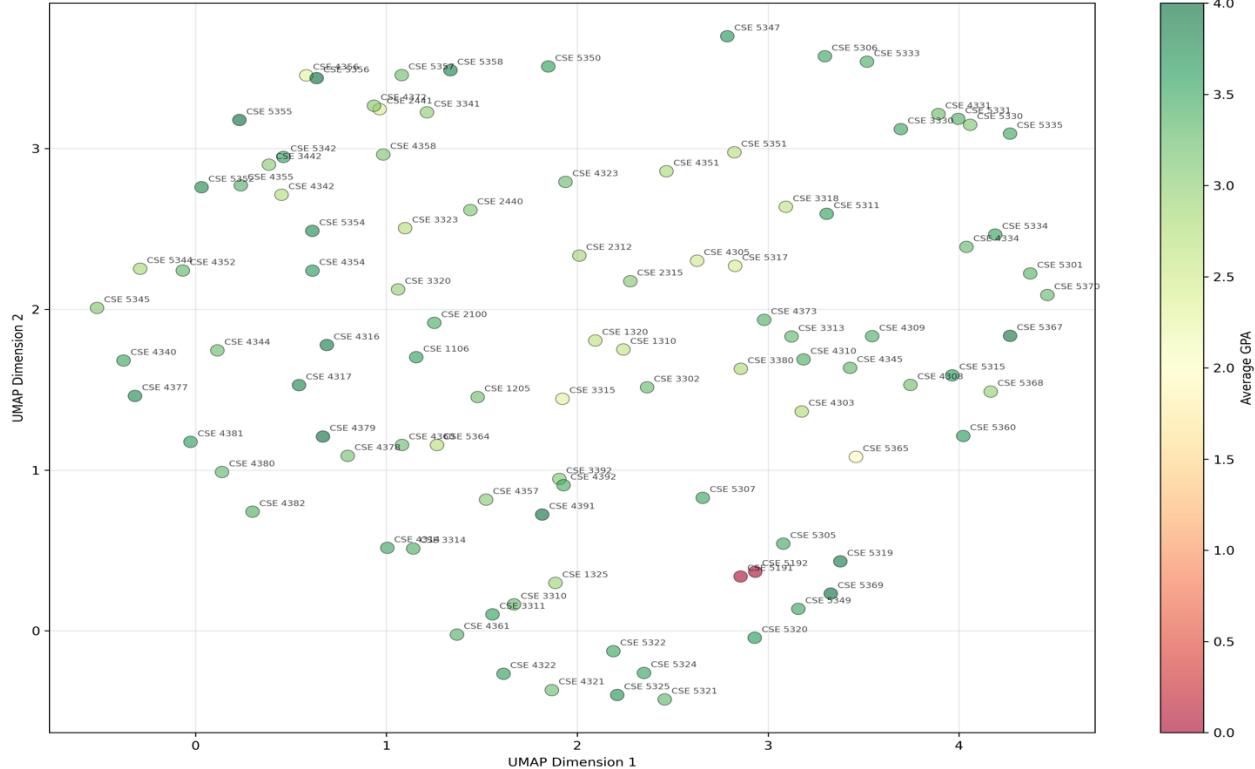
[5]



[6]

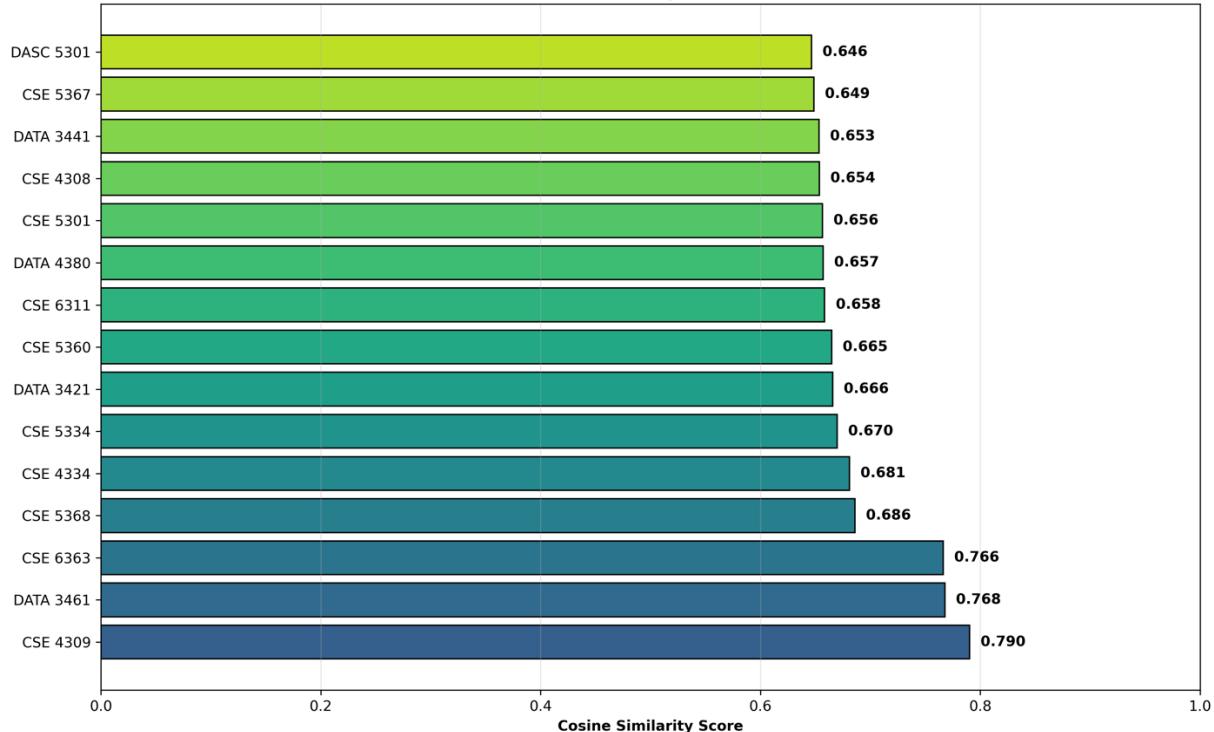


**Course Embeddings - UMAP Visualization
Colored by Average GPA - Green=Easy, Red=Hard**



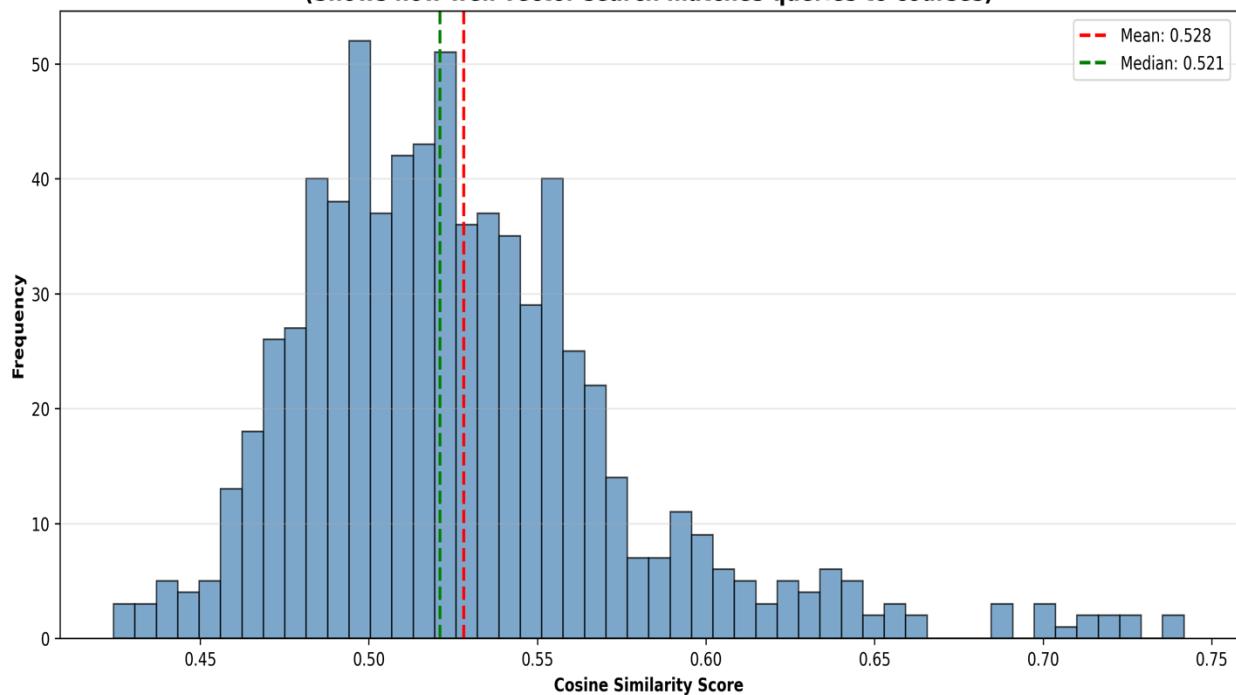
[7]

**Top 15 Most Similar Courses to Query: "machine learning courses"
(Semantic Neighborhood)**

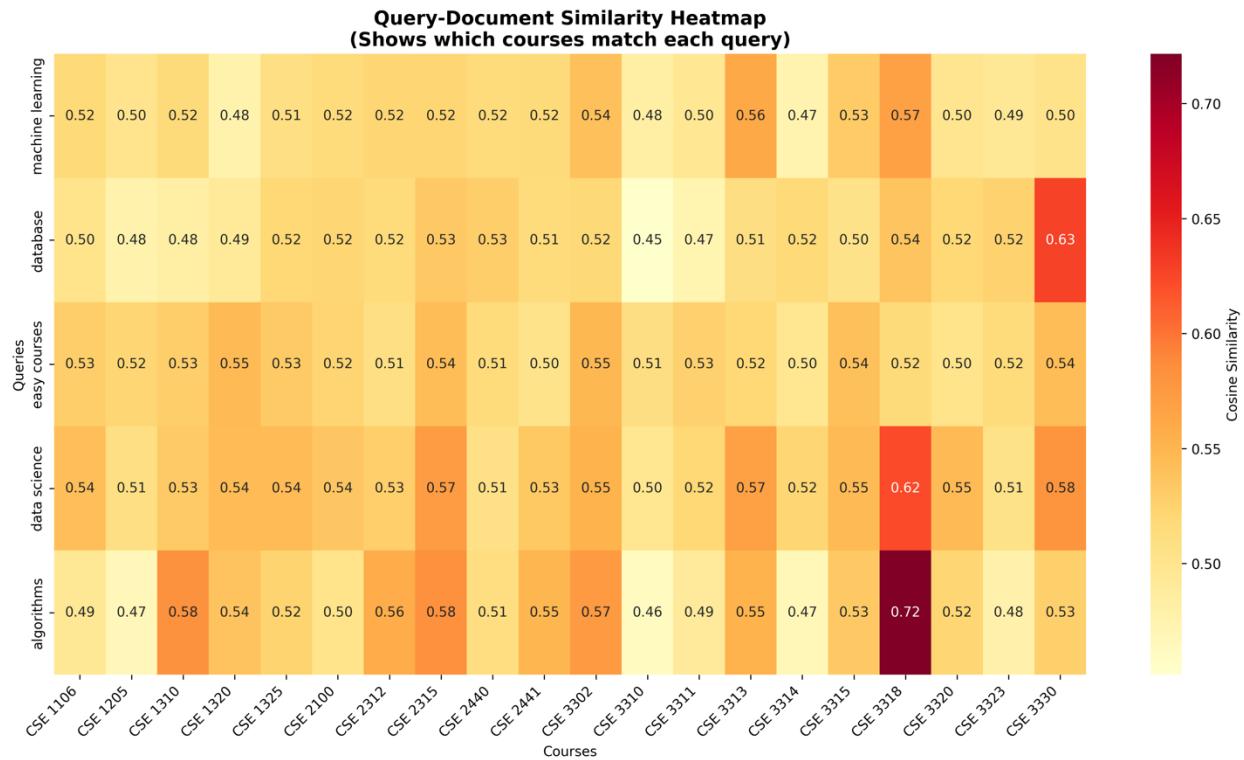


[8]

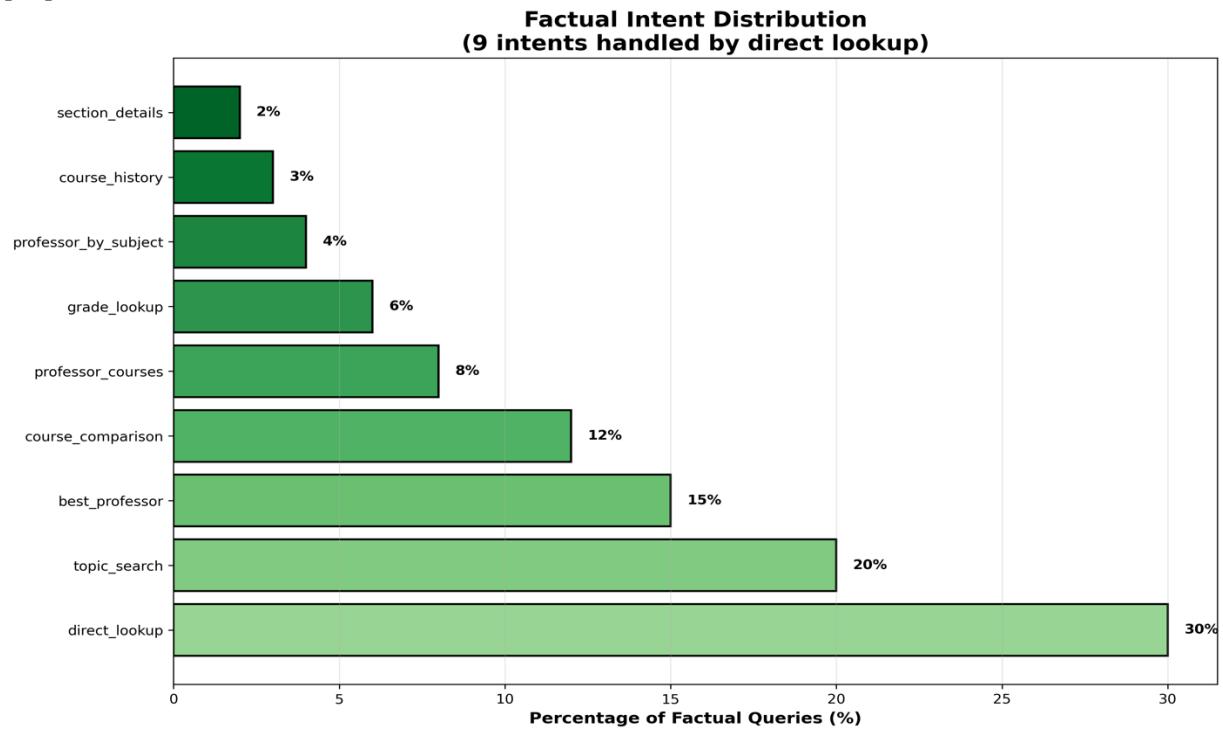
**Distribution of Retrieval Scores
(Shows how well vector search matches queries to courses)**



[9]

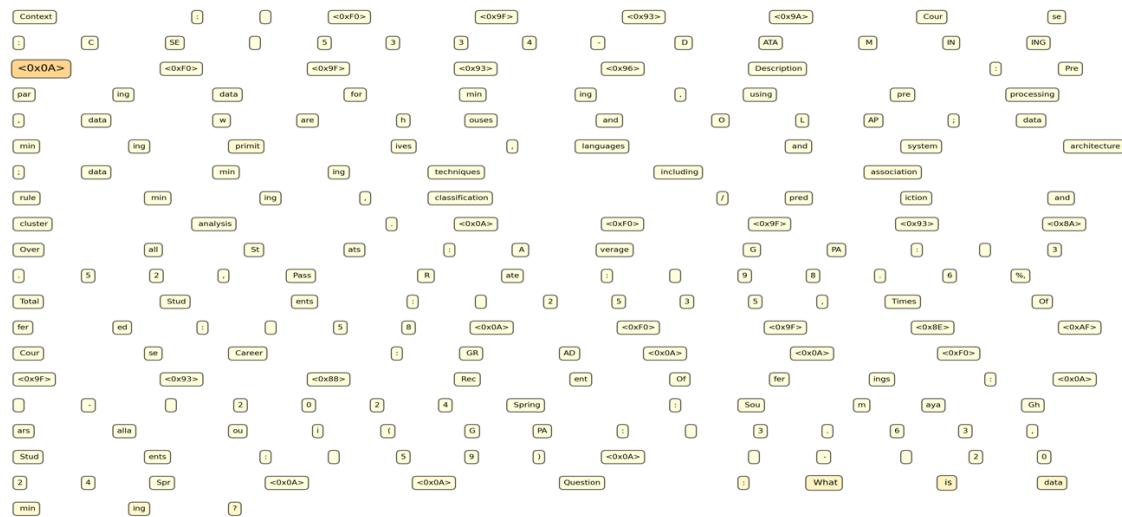


[10]



[11]

Model Attention Focus for Query: "What is data mining?"
(Brighter/Larger = Higher Attention)



[12]

🎓 UT A Course Q&A Agent

Enhanced with Course Analytics, Professor Insights & Grade Lookups

Ask about courses, professors, grades, comparisons, and more!

🚀 Initialization

Initialize Agent

✓ Agent initialized successfully! You can now ask questions about UTA courses.

💡 Example Queries

☰ Click an example to try it!

CSE 5334 | Compare CSE 5334 and CSE 5330

Courses by John Smith

CSE 5334 grades Spring 2023

Easy computer science courses

Professor John Smith | History of CSE 5334

CSE 5334 with Professor John

📊 Features

Click to see features

Type your question about UTA courses here...

Send 📩

Clear Chat

Restart Agent

Certainly! CSE 5334 is the course titled "DATA MINING" offered by the University of Texas at Arlington. The course is designed to prepare students for data mining, using preprocessing, data warehouses, and OLAP, and data mining techniques such as association rule mining, classification/prediction, and cluster analysis. The course is taught by a team of experienced faculty members and has an overall GPA of 3.52, a pass rate of 98.6%, and an average of 2535 students enrolled. The course is offered in various semesters and has several recent offerings, including Soumaya Gharsallaoui's Spring 2024 offering, Marnim Galib's Spring 2024 offering, Soumaya Gharsallaoui's Summer 2024 offering, and Marnim Galib's Summer 2024 offering.

UTA Course Q&A Agent

Enhanced with Course Analytics, Professor Insights & Grade Lookups

Ask about courses, professors, grades, comparisons, and more!

Initialization

Initialize Agent

Agent initialized successfully! You can now ask questions about UTA courses.

Example Queries

Click an example to try it!

CSE 5334 | Compare CSE 5334 and CSE 5330

Courses by John Smith

CSE 5334 grades Spring 2023

Easy computer science courses

Professor John Smith | History of CSE 5334

CSE 5334 with Professor John

Features

Click to see features

Course Q&A Chat

tell about professor Rozsa Zaruba

Professor: Rozsa Zaruba
Teaching Statistics:
- Average GPA Given: 3.79
- A Rate: 81.0%
- A+B Rate: 98.6%
- Pass Rate: 99.3%
- DFWU Rate: 0.7%
- Total Students: 346
- Terms Taught: 8
- Teaching Style: Generous Grader

Courses Taught (5):
1. CSE 4334: DATA MINING
2. DASC 5301: DATA SCIENCE
3. DASC 5302: INTRODUCTION TO PROBABILITY AND STATISTICS
4. DASC 5309: DATA SCIENCE CAPSTONE PROJECT
5. DASC 5301: DATA SCIENCE APPLICATIONS

Type your question about UTA courses here...

Send

Clear Chat

Restart Agent

UTA Course Q&A Agent

Enhanced with Course Analytics, Professor Insights & Grade Lookups

Ask about courses, professors, grades, comparisons, and more!

Initialization

Initialize Agent

Agent initialized successfully! You can now ask questions about UTA courses.

Example Queries

Click an example to try it!

CSE 5334 | Compare CSE 5334 and CSE 5330

Courses by John Smith

CSE 5334 grades Spring 2023

Easy computer science courses

Professor John Smith | History of CSE 5334

CSE 5334 with Professor John

Features

Click to see features

Course Q&A Chat

History of DASC 5301

Course: DASC 5301 - DATA SCIENCE
Description: This inspirational course follows a data-science-for-all perspective that views data acumen as part of literacy. It aims to instill in students the data acumen, i.e., the basic skills to wrestle with data, to draw insights from data, to make sound decisions responsibly using data, and to effectively communicate about data-driven findings and decisions. Topics include 1) data management: data curation, preparation, model, and querying; 2) data description and visualization: exploratory data analysis; graphics; 3) machine learning and knowledge discovery: supervised learning, unsupervised learning, pattern and knowledge extraction, model evaluation and interpretation. Prerequisite: DASC Major, DASC 5300 (or concurrent enrollment) and DASC 5302 (or concurrent enrollment).
Overall Stats: Average GPA: 3.27, Pass Rate: 92.7%, Total Students: 554, Times Offered: 8
Course Career: GRAD

Send

UTA Course Q&A Agent

Enhanced with Course Analytics, Professor Insights & Grade Lookups

Ask about courses, professors, grades, comparisons, and more!

Initialization

Initialize Agent

Agent initialized successfully! You can now ask questions about UTA courses.

Example Queries

Click an example to try it!

CSE 5334 | Compare CSE 5334 and CSE 5330

Courses by John Smith

CSE 5334 grades Spring 2023

Easy computer science courses

Professor John Smith | History of CSE 5334

CSE 5334 with Professor John

Features

Click to see features

Course Q&A Chat

Should I take machine learning or data mining?

To answer the question, "Should I take machine learning or data mining?" the course assistant can provide a helpful answer using information from UTA's course database.
The course assistant can say that while machine learning techniques and their application to various problems in data science are covered in DATA 3461, MACHINE LEARNING, and DATA 3401, MACHINE LEARNING, respectively, the course focuses more on statistical and connectionist learning methods. The course assistant can also mention that the current issues in machine learning research are determined by the course's instructor, and the course's overall statistics are 92.9%, with a pass rate of 99.3%. The course assistant can also mention that the course's recent offerings include Amir Farbin (GPA: 3.32, Students: 25) and Farnaz Sardari (GPA: 3.32, Students: 25).

Type your question about UTA courses here...

Send

Clear Chat

Restart Agent

UTA Course Q&A Agent

Enhanced with Course Analytics, Professor Insights & Grade Lookups

Ask about courses, professors, grades, comparisons, and more!

Initialization

Initialize Agent

Agent initialized successfully! You can now ask questions about UTA courses.

Example Queries

- Click an example to try it!
- CSE 5334 | Compare CSE 5334 and CSE 5330
- Courses by John Smith
- CSE 5334 grades Spring 2023
- Easy computer science courses
- Professor John Smith | History of CSE 5334
- CSE 5334 with Professor John

Features

Click to see features

Course Q&A Chat

give me total grades of dasc 5301 in spring 2024

GRADE REPORT: DASC 5301

Term: spring 2024
 Total Sections: 2
 Total Students: 60
 Average GPA: 3.68

Grade Distribution:

- A: 41 (69.5%)
- B: 17 (28.8%)
- C: 1 (1.7%)
- D: 0 (0.0%)
- F: 0 (0.0%)
- W: 1

Instructors: Soumaya Gharsallaoui, Rozsa Zaruba

Type your question about UTA courses here...

Send

UTA Course Q&A Agent

Enhanced with Course Analytics, Professor Insights & Grade Lookups

Ask about courses, professors, grades, comparisons, and more!

Initialization

Initialize Agent

Agent initialized successfully! You can now ask questions about UTA courses.

Example Queries

- Click an example to try it!
- CSE 5334 | Compare CSE 5334 and CSE 5330
- Courses by John Smith
- CSE 5334 grades Spring 2023
- Easy computer science courses
- Professor John Smith | History of CSE 5334
- CSE 5334 with Professor John

Features

Click to see features

Course Q&A Chat

CSE 5334 | DATA MINING | 3.52 | 98.6% | 1.3% | 2,535
 CSE 6363 | MACHINE LEARNING | 3.64 | 99.3% | 0.6% | 2,712

Comparison Summary:

- Highest GPA: CSE 6363 (3.64)
- Highest Pass Rate: CSE 6363 (99.3%)
- Lowest DFW Rate: CSE 6363 (0.6%)

CSE 5334 with Professor John

No history found for CSE 5334 taught by John Robb

Type your question about UTA courses here...

Send

Clear Chat

Restart Agent