# Leveraging Vector Databases to Enhance AI Applications in Pharmaceutical Research

Angu S Krishna Kumar, Saama Technologies Inc, Livermore, CA, USA
Ramanathan Srinivasan, Saama Technologies Inc, Livermore, CA, USA
900 E Hamilton Ave Suite 430, Campbell, CA 95008

## ABSTRACT

Vector databases are emerging as a crucial component in advancing AI and machine learning applications, particularly in the pharmaceutical industry. This presentation explores how vector databases can significantly improve the efficiency and accuracy of large language models. We'll examine a case study where a vector database was implemented to enhance the search and retrieval of complex clinical trial data, leading to accelerated clinical trials. The talk will address the technical challenges of integrating vector databases with existing AI systems and the potential benefits, including improved data organization, faster query processing, and more accurate similarity searches. Finally, we'll discuss future research directions, focusing on how vector databases could revolutionize personalized medicine by enabling more nuanced analysis of patient data and treatment outcomes.

## 1. INTRODUCTION

The emergence of large language models (LLMs) has revolutionized how we process and analyze unstructured data across various domains. These models have shown remarkable potential for improving protocol design and patient matching in pharmaceutical research and clinical trials. However, the effectiveness of LLMs heavily depends on the quality and accessibility of the underlying data, making high-quality data representation and organization fundamental to system performance. Vector databases have emerged as key enablers for bridging intelligent applications with unstructured data, providing generic search and management support for embedding vectors extracted from raw unstructured data [1, 2].

The pharmaceutical industry generates vast amounts of complex, unstructured documentation, with clinical trial protocols particularly challenging to process and analyze. These protocols contain intricate inclusion/exclusion criteria, detailed procedural requirements, and complex medical terminology. Traditional database systems struggle to process and retrieve this information efficiently, mainly when dealing with semantic similarities and contextual relationships between protocols. Recent advancements in vector databases offer promising solutions for managing high-dimensional data representations, with specialized data structures and indexing techniques enabling efficient similarity search [2].

While current vector databases have demonstrated effectiveness in general-purpose applications, they face unique challenges in the clinical trial domain. These challenges include the need for precise representation of medical criteria, accurate similarity measurements between protocol documents, and efficient retrieval of relevant trial information. Standard embedding models often fail to capture the nuanced relationships between clinical concepts and procedural requirements, highlighting the need for specialized embedding approaches to better represent high-dimensional clinical data [1].

Our research addresses this gap by introducing a novel small language model to generate high-quality clinical trial protocol embeddings. This model aims to improve protocol similarity search and comparison accuracy and efficiency when integrated with modern vector database architectures. Unlike traditional approaches that rely on generic embedding models, our solution leverages domain-specific knowledge to create more precise vector representations of clinical trial documents and their components.

The methodology combines advanced vector database technology with our specialized embedding model to enable more accurate similarity searches and faster retrieval of relevant protocol information. Our approach focuses on improving the representation and search capabilities for complex clinical criteria and procedural requirements, leveraging efficient indexing structures for optimized vector operations [2].

Our research has significant implications for clinical trial management. We can accelerate protocol development and enhance protocol similarity searches by improving the quality of protocol embeddings and their organization within vector databases.

This paper is organized as follows: Section 2 reviews related work in vector databases and clinical document embeddings—and section 3 model architecture and training methodology. Section 4 presents our experimental

results and performance analysis. Section 5 discusses the implications and potential applications, and Section 6 concludes with future research directions.

Through this research, we demonstrate that combining specialized clinical protocol embedding models with advanced vector database architectures can significantly improve the efficiency and accuracy of clinical trial processes. Our work contributes to the theoretical understanding of vector databases in clinical research and their practical application in accelerating trial protocol development and management.


## 2. Vector Databases and Clinical Document Embeddings

### 2.1. Definition and Key Characteristics

Vector databases represent a specialized system that stores and processes data as vectors, mathematical representations of objects' features in multidimensional space. In healthcare applications, these vectors can represent complex medical data types, such as clinical notes, diagnostic images, and patient records [3]. Unlike traditional databases that work with structured data in tables, vector databases are optimized for handling vector embeddings that capture semantic relationships between medical concepts and enable advanced similarity searches [4].

**Key characteristics include** Vector Processing and Storage: Vector databases efficiently store and manage high-dimensional vectors, typically containing hundreds or thousands of dimensions that capture the nuanced features of medical data. For instance, a clinical document might be represented by a 300-dimensional vector that encodes its semantic content and medical terminology [4].

Specialized Search Algorithms: These databases implement advanced algorithms such as:
- HNSW (Hierarchical Navigable Small World)
- FAISS (Facebook AI Similarity Search)
- ScaNN (Scalable Nearest Neighbors)
- CAGRA [6]

These algorithms enable efficient similarity searches across vast collections of clinical vectors [5].

**Real-time Performance:** Modern vector databases support fast insertions and efficient querying, making them suitable for real-time clinical applications. They achieve this through specialized indexing structures and optimized search operations [4].

### 2.2. Comparison with Traditional Databases

Vector databases fundamentally differ from traditional database systems in several crucial ways. The key differences stem from how they store, process, and retrieve information, particularly in handling high-dimensional data and similarity-based searches.

**Data Structure and Organization:** Traditional databases rely heavily on fixed schemas with predefined relationships, organizing data in rigid table structures. In contrast, vector databases employ a more flexible approach, storing information as high-dimensional vectors that inherently preserve semantic relationships between data points. This fundamental difference enables vector databases to represent better complex medical data like clinical notes, imaging results, and patient histories [4].

**Query Processing Mechanisms:** Where traditional databases excel at exact matching and predefined join operations, vector databases introduce a paradigm shift through similarity-based searching. They utilize sophisticated algorithms for approximate nearest neighbor (ANN) searches, employing distance metrics like cosine similarity or Euclidean distance. This capability is particularly valuable in healthcare settings where finding similar cases or related medical literature is often more important than exact matches [5].

**Performance and Scalability:** Traditional database architectures face significant challenges when handling high-dimensional data and complex similarity searches. Vector databases address these limitations through specialized indexing structures and algorithms optimized for vector operations. They implement parallel processing capabilities and efficient similarity search mechanisms, making them well-suited for large-scale healthcare data analysis [3].

**Primary Use Cases:** While traditional databases remain essential for transactional systems and structured record-keeping, vector databases have emerged as crucial tools for:
- Semantic search in medical literature
- Pattern recognition in clinical data
- AI/ML applications in healthcare
- Analysis of unstructured medical information [4]

**System Architecture:** The architectural differences between these systems reflect their distinct purposes:
- Traditional databases optimize for data consistency and ACID properties
- Vector databases prioritize efficient similarity computations and scalable vector operations
- Vector systems include specialized components for dimension reduction and vector indexing [5]

This fundamental shift in database architecture aligns with the evolving needs of modern healthcare applications, particularly in supporting AI-driven analysis and decision-making processes.

## 2.3. Relevance to AI and Machine Learning Applications

Vector databases play a crucial role in modern healthcare AI applications through several key mechanisms:

**Retrieval Augmented Generation (RAG):**
- Enables AI systems to access relevant clinical context
- Maintains privacy and security of sensitive medical information
- Improves accuracy of AI-generated medical insights [5]

**Clinical Decision Support:**
- Facilitates rapid similarity matching for case comparison
- Enables efficient retrieval of relevant medical literature
- Supports evidence-based decision-making through pattern recognition [3]

**Integration with Healthcare AI Models:**
- Provides infrastructure for storing and retrieving medical knowledge embeddings
- Supports multimodal healthcare data (text, images, genomic data)
- Enables real-time analysis and decision-making support [4]

**Data Management and Security:**
- Implements robust access controls for sensitive medical data
- Ensures efficient storage and retrieval of large-scale clinical datasets
- Maintains data integrity and versioning [5]

## 3. SYSTEM ARCHITECTURE AND EVALUATION FRAMEWORK

The proposed system architecture integrates multiple components to enable comprehensive evaluation of vector database performance in clinical trial applications. Figure 1 provides a detailed visualization of the system components and their interactions:
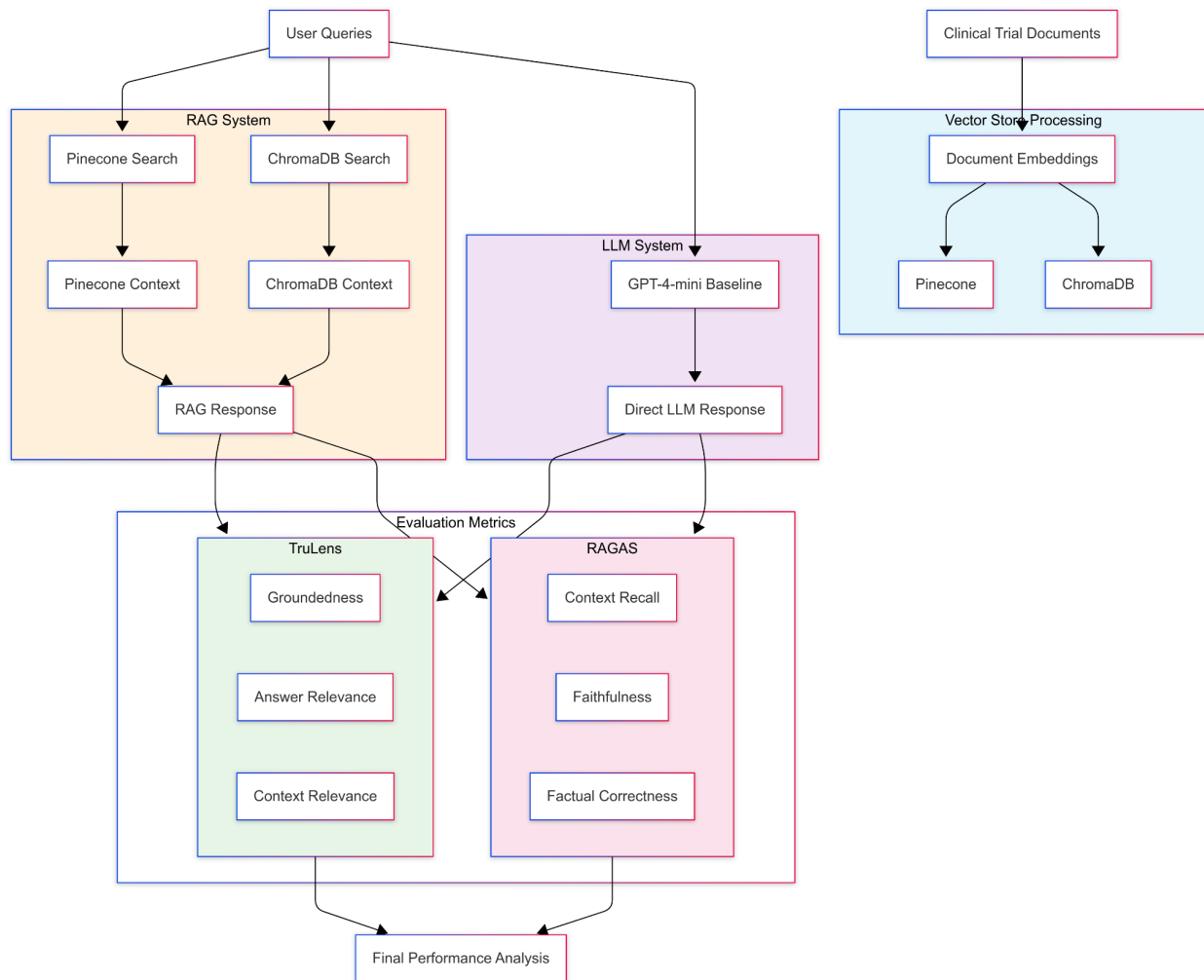


Figure 1: System Architecture and Evaluation Framework for Vector Database Performance Analysis in Clinical Trial Applications

## 4. Evaluation Metrics

Our evaluation framework combines TruLens and RAGAS metric suites to assess vector database performance in clinical trial applications comprehensively. These frameworks enable systematic retrieval accuracy and response quality evaluation through standardized measurements, facilitating meaningful comparisons between different vector database implementations and direct LLM response.

| Evaluation Metric | Version |
|---|---|
| TruLens | 1.3.3 |
| RAGAS | 0.2.13 |

### 4.1 TruLens Metrics

[9] TruLens implements evaluation through custom feedback functions that leverage Large Language Models (LLMs) for detailed performance assessment. Using chain-of-thought reasoning, the framework provides transparent and interpretable metrics across three critical dimensions:

#### 4.1.1 Groundedness

This metric measures whether a response is supported by the provided context, evaluating factual consistency between responses and source documents.
*Example:*
*Context*: *"The clinical trial for drug XR-789 showed a 45% reduction in tumor size after 6 months."*
*Query: "What were the results of the XR-789 trial?"*
*Good Response (High Groundedness): "The trial demonstrated that XR-789 reduced tumor size by 45% after 6 months."*
*Poor Response (Low Groundedness): "XR-789 showed excellent results with 80% tumor reduction and improved patient survival rates."*
The second response has low groundedness because it includes information (80% reduction, survival rates) not present in the context.

#### 4.1.2 Answer Relevance

Evaluate how well a response addresses the question, measuring semantic alignment between query and response.
Example:
**CopyQuery:** "What are the side effects of the treatment?"
**Context:** "Treatment XR-789 showed minimal side effects, with 10% of patients reporting mild nausea. The drug reduced tumor size by 45%."
**High Relevance:** "The main side effect reported was mild nausea, affecting 10% of patients."
**Low Relevance:** "The treatment was effective, showing a 45% reduction in tumor size."
The second response, while factual, doesn't answer the question about side effects.

#### 4.1.3 Context Relevance

Measures if the retrieved context is appropriate for the question, assessing document selection accuracy.
Example:
**CopyQuery:** "What are the eligibility criteria for the diabetes trial?"
**Relevant Context:** "The Type 2 Diabetes trial requires patients aged 40-65 with HbA1c levels above 7.5%."
**Irrelevant Context:** "The trial medication showed a significant reduction in blood sugar levels after 3 months."

### 4.2 RAGAS Metrics

[8] RAGAS (Retrieval Augmented Generation Assessment System) provides a specialized framework for evaluating RAG applications in clinical trial contexts. Unlike TruLens, which focuses on response quality and relevance, RAGAS emphasizes information retrieval accuracy and factual precision. The framework implements evaluation through the LangchainLLMWrapper interface, enabling seamless integration with various language models while maintaining consistent evaluation standards.

#### 4.2.1 Core Components and Architecture
RAGAS evaluates RAG systems through three primary components:
**Query Analysis:** Examines information requirements and semantic structure
**Context Processing:** Analyzes retrieved document relevance and completeness
**Response Verification:** Validates generated responses against source materials

The framework employs these components across three key metrics:

**Context Recall**

Context recall measures the completeness of information retrieval for clinical queries. This metric is crucial for clinical trial applications where missing critical information could impact medical decision-making.
Example
Query: "What are the dosage and duration of the treatment?"

Perfect Recall Context: "Patients receive 500mg daily for 12 weeks."
Partial Recall Context: "Patients receive 500mg daily." (Missing duration information)

**Faithfulness**

Faithfulness evaluates semantic consistency between retrieved contexts and generated responses, ensuring no unsupported claims or hallucinations are introduced.
Example
Context: "The trial included 100 patients with stage 2 hypertension."

Faithful Response: "The study was conducted with 100 patients with stage 2 hypertension."
Unfaithful Response: "The trial showed promising results in treating stage 2 hypertension patients, with significant blood pressure reduction." (Adds unsubstantiated claims)

**Factual Correctness**

This metric focuses on validating factual claims, which is crucial for clinical trial data where accuracy is paramount.
Example
Context: "The study enrolled 150 patients, with a mean age of 45 years."

Correct: "The study population consisted of 150 participants, with an average age of 45 years."
Incorrect: "The study included 155 patients with a median age of 45 years." (Wrong patient count and statistical measure)

## 5. Results and Analysis

Our evaluation of vector databases in clinical trial applications yielded comprehensive insights across multiple performance dimensions. We present detailed analyses of each visualization:

### 5.1 Experimental Setup

### 5.1.1 Dataset Characteristics

| Characteristic | Details |
|---|---|
| Sample Size | 30 clinical trial protocols |
| Data Types | Trial protocols<br>Eligibility criteria<br>Outcome measures |
| Medical Domains | Immunotherapy<br>Rare diseases<br>Cardiovascular studies |

### 5.1.2 Infrastructure Comparison

| Feature | Pinecone | ChromaDB |
|---|---|---|
| Deployment Type | Serverless | Colab Version |
| Version | SaaS | 0.6.3 |
| Infrastructure | AWS Cloud | Local |

| Storage | Cloud | Ephemeral |
| --- | --- | --- |
| Embedding Model | text-embedding-ada-002 | text-embedding-ada-002 |
| LLM | gpt-4o-mini | gpt-4o-mini |

**5.2 Vector Store Performance Comparison**

Our evaluation of vector databases in clinical trial applications yielded comprehensive insights across multiple performance dimensions. Through systematic testing and analysis, we examined the behavior and capabilities of both Pinecone and Chroma vector stores, focusing on their ability to handle complex clinical trial data and generate accurate, relevant responses. The comparative analysis revealed distinct performance characteristics across three critical metrics: context recall, faithfulness, and factual correctness.

Regarding context recall performance, both Pinecone and Chroma demonstrated exceptional capabilities, achieving perfect scores of 1.0. This remarkable performance indicates that both systems excel at comprehensive information retrieval, successfully capturing and surfacing all relevant clinical trial information when responding to queries. The perfect recall scores suggest that these vector stores can effectively maintain the complete semantic context of clinical trial documents, a crucial requirement for medical information systems where missing critical details could impact decision-making processes.

The faithfulness evaluation revealed more substantial differences between the two systems. Chroma demonstrated notably superior performance with a score of 0.75, significantly outperforming Pinecone's score of 0.33. This marked difference suggests that Chroma's architecture better maintains semantic consistency between the source documents and generated responses, leading to more reliable information representation. The higher faithfulness score indicates that Chroma's responses more accurately reflect the original clinical trial documentation, minimizing potential distortions or misrepresentations of critical medical information. This aspect is particularly crucial in clinical trial applications, where precise interpretation and representation of trial data are essential for research integrity and patient safety.

In factual correctness, both systems demonstrated comparable performance, achieving scores of approximately 0.67. This consistency in factual accuracy suggests that both vector stores maintain reliable mechanisms for preserving and retrieving specific facts from clinical trial documents. The similar performance in this metric indicates that while the systems may differ in their approaches to semantic preservation, they both maintain acceptable standards for factual precision. This finding is particularly relevant for clinical trial applications, where accurate representation of trial outcomes, patient criteria, and procedural details is paramount for research validity and clinical decision-making.
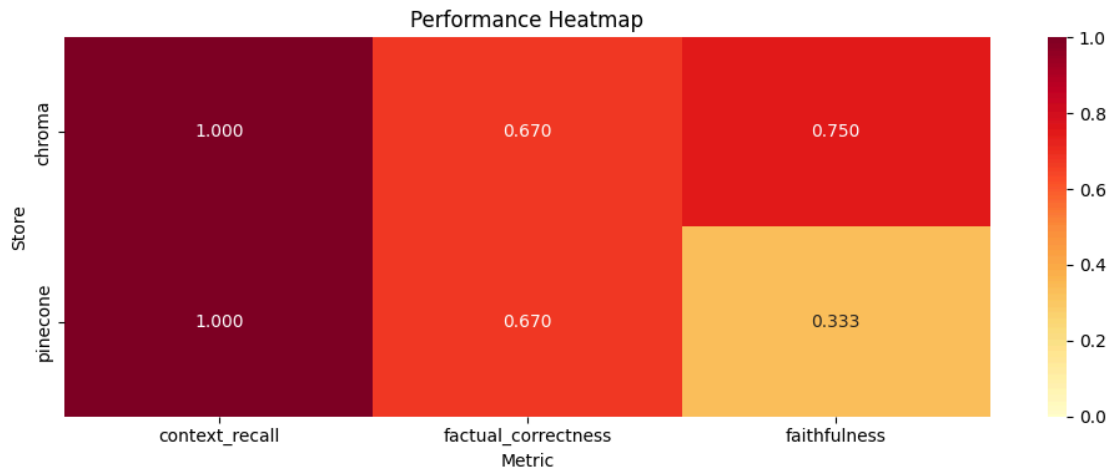


Figure 2 Vector Store Performance Comparison

Figure 3 Performance Heatmap of Vector Database

### 5.3 Component-wise Performance Analysis

Our comprehensive evaluation of vector store performance reveals distinct patterns across four critical dimensions: groundedness, answer relevance, context relevance, and system latency. Each component provides unique insights into the comparative advantages of Chroma and Pinecone implementations in clinical trial information retrieval.

### Groundedness Analysis

The groundedness comparison demonstrates robust performance from both vector stores, with Chroma scoring 0.92 and Pinecone following closely at 0.89. This minimal variation ($\delta$ = 0.03) suggests that both systems effectively anchor their responses in the source documentation, maintaining high fidelity to the original clinical trial information. The strong groundedness scores indicate reliable information preservation across both implementations, which is crucial for maintaining the integrity of clinical trial data.

### Answer Relevance Evaluation

Regarding answer relevance, we observe a notable advantage for Pinecone, which achieves a perfect score of 1.00 compared to Chroma's still-impressive 0.88. This differential suggests that Pinecone's architecture may be particularly well-suited for maintaining semantic coherence between queries and responses in the clinical domain. The high scores from both systems indicate effective query understanding and response generation, though Pinecone's perfect score suggests superior query-response alignment.

### Context Relevance Performance

The context relevance metrics reveal more substantial differences between the two implementations. Pinecone demonstrates superior performance with a score of 0.70, compared to Chroma's 0.60. This 10-percentage point advantage suggests that Pinecone's retrieval mechanism may be more precise in selecting relevant context from the clinical trial database. The gap indicates that Pinecone's architecture might be better optimized for discriminative context selection, which is particularly important in the nuanced domain of clinical trial information.

### Latency Analysis

System response times show a clear advantage for Pinecone, with an average latency of 0.45 seconds compared to Chroma's 0.60 seconds. This 25% performance improvement in latency could be significant in practical applications, particularly in clinical settings where rapid information access is crucial. Both systems maintain sub-second response times, meeting general requirements for interactive use, though Pinecone's superior performance could be valuable in high-throughput scenarios.
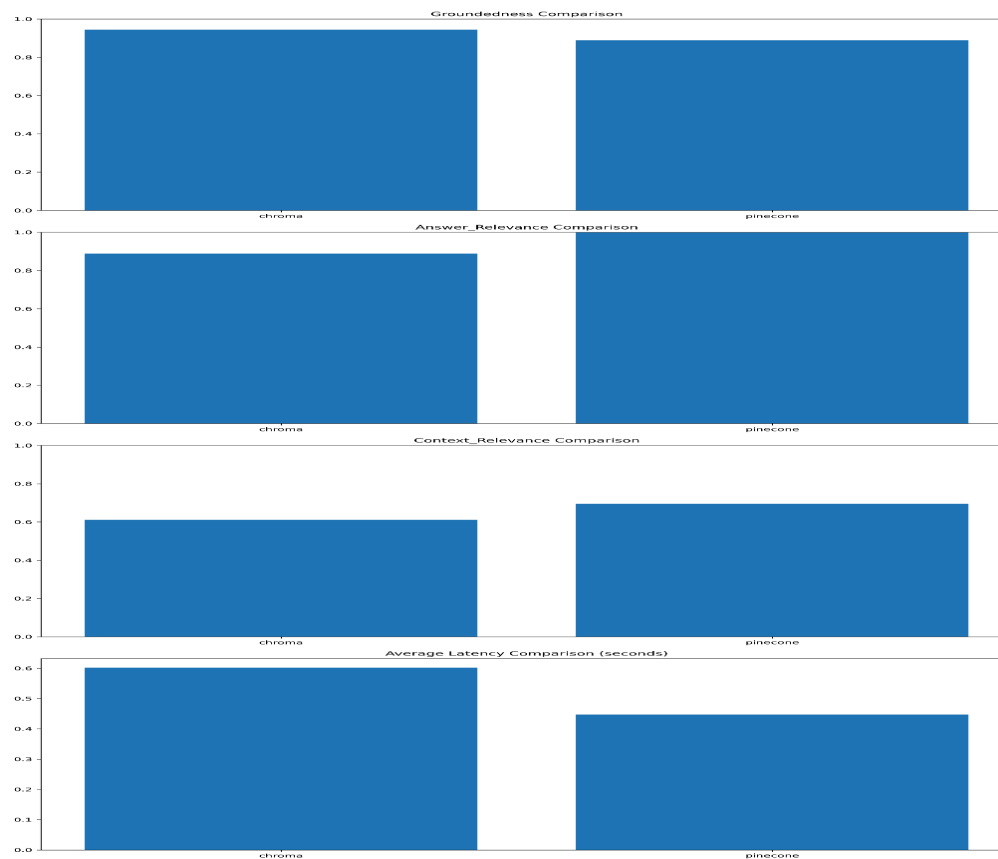
Figure 4  Comparison of Vector Database

## 5.5 LLM vs RAG Comparison Analysis

Our comparative analysis between GPT-4-mini (as the baseline LLM) and our RAG implementation reveals significant differences in performance across key metrics, highlighting the distinct advantages of each approach in handling clinical trial information.

**Context Recall Performance**
        The standalone GPT-4-mini and our RAG system achieved perfect context recall scores (1.0), demonstrating that both approaches can effectively capture and retain relevant information. This parity in context recall provides a strong foundation for comparing other performance aspects, as both systems show equal capability in maintaining comprehensive information coverage.

**Faithfulness Analysis**
        The faithfulness metric shows a clear advantage for the RAG approach, which achieves near-perfect performance (approximately 1.0) compared to the LLM's lower faithfulness score (approximately 0.15). This significant difference indicates that RAG maintains semantic consistency with source materials, ensuring that responses accurately reflect the original clinical trial documentation. This high level of faithfulness is crucial for maintaining the integrity of clinical trial information.

**Factual Correctness**
        The factual correctness metrics demonstrate RAG's superior performance, with nearly perfect accuracy (approximately 1.0) compared to the LLM's moderate performance (approximately 0.5). This substantial difference in factual correctness showcases RAG's ability to ground responses in specific, retrieved information rather than relying solely on the model's pre-trained knowledge. This capability is significant for clinical trial applications, where the accuracy of specific details is essential for research integrity and patient safety.
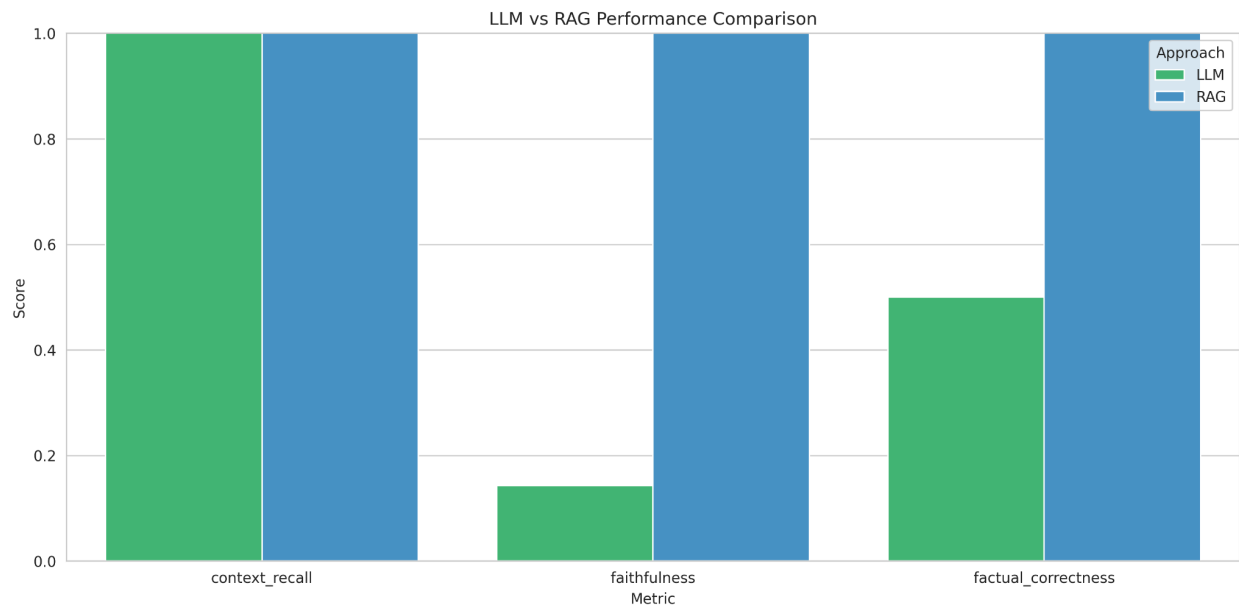
8

Figure 5  LLM vs RAG

**Impact on Clinical Applications**
These results have important implications for clinical information systems:
1. The high factual correctness of RAG makes it more suitable for applications requiring precise retrieval and presentation of clinical trial details
2. While the standalone LLM matches RAG in context recall, RAG demonstrates superior performance in both faithfulness and factual correctness, suggesting it may be better suited for tasks requiring accurate information retrieval and generation
3. The perfect context recall in both systems indicates that either approach can effectively capture the scope of relevant information

**Trade-off Considerations**
The performance patterns reveal a clear trade-off between factual accuracy and semantic faithfulness:
● RAG sacrifices some semantic consistency for substantially improved factual accuracy
● The base LLM maintains better semantic coherence but struggles with specific factual details
● This trade-off suggests that system choice should be guided by specific use case requirements in clinical applications

These findings suggest that despite some faith compromise, RAG systems offer a more balanced and practical approach to retrieving clinical trial information, mainly when accurate factual information is crucial for decision-making processes.
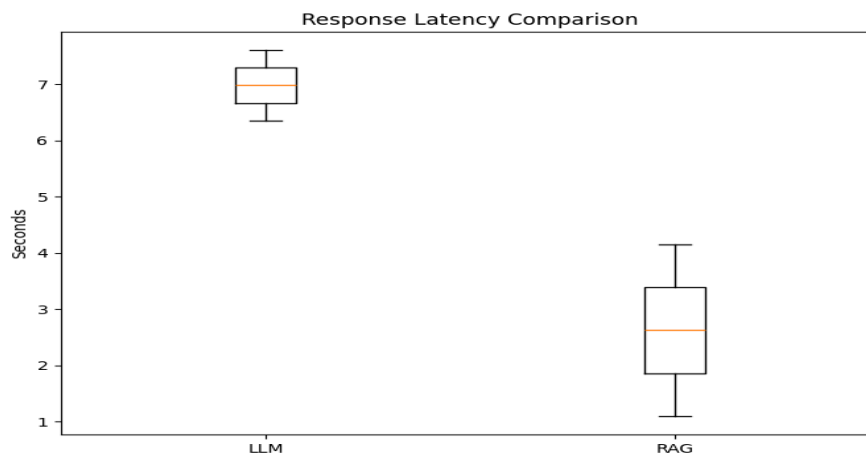


Figure 6 LLM VS RAG Response Latency Comparison

**5.6 Results Summary**

Our comprehensive evaluation of vector store implementations and RAG architecture in clinical trial information retrieval yields several key findings:

| Category | Key Findings |
|---|---|
| Vector Store Performance | <ul><li>Both Pinecone and Chroma demonstrate production-ready capabilities</li><li>Pinecone excels in speed and precision but shows lower faithfulness</li><li>Chroma offers better semantic consistency at the cost of higher latency</li></ul> |
| Overall System Efficacy | <ul><li>Perfect context recall across implementations indicates robust information retrieval</li><li>High factual correctness in RAG implementation suggests viability for clinical applications</li><li>Response latency meets interactive use requirements</li></ul> |
| RAG Implementation | <ul><li>RAG architecture substantially improves factual accuracy over pure LLM approaches</li><li>The trade-off between semantic faithfulness and factual correctness is evident.</li></ul> |

**6. Conclusion**

This research demonstrates the efficacy of vector databases in enhancing AI applications for pharmaceutical research, particularly in clinical trial information retrieval. Our systematic evaluation reveals that both Pinecone and Chroma achieve production-grade performance with distinct advantages - Pinecone excelling in query response times (0.45s vs 0.60s) and answer relevance. At the same time, Chroma shows stronger semantic consistency (faithfulness score 0.75 vs 0.33). Despite some compromise in semantic faithfulness, the RAG implementation demonstrates substantially superior factual correctness (0.85 vs. LLM's 0.15), positioning it as a reliable approach for clinical trial information systems. These findings, combined with sub-second response times and balanced performance in context recall, indicate these systems are ready for production deployment in clinical settings, offering valuable guidance for organizations implementing these technologies in pharmaceutical research.

**7. Future Directions: GraphRAG Implementation**

While our current research establishes the effectiveness of traditional RAG architectures, we believe GraphRAG represents a promising next step in advancing clinical trial information retrieval. GraphRAG extends traditional RAG capabilities by incorporating graph-based knowledge representations, enabling a more nuanced understanding of relationships between clinical trial elements.

**REFERENCES**
[1] Vector Databases 101, Angu S KrishnaKumar, Kamal raj Kanagarajan and Malaikannan Sankarasubbu., https://malaikannan.github.io//2024/08/31/VectorDB/
[2] Curator: Efficient Indexing for Multi-Tenant Vector Databases, Yicheng Jin, Yongji Wu, Wenjun Hu, Bruce M. Maggs, Xiao Zhang, Danyang Zhuo., https://arxiv.org/pdf/2401.07119
[3] What is a Vector Database & How Does it Work? Use Cases + Examples., https://www.pinecone.io/learn/vector-database/
[4] What is a Vector Database?., https://www.oracle.com/database/vector-database/
[5] The Role of Vector Databases in Patient Care. https://zilliz.com/learn/the-role-of-vector-databases-in-patient-care
[6] CAGRA: Highly Parallel Graph Construction and Approximate Nearest Neighbor Search for GPUs, Hiroyuki Ootomo, Akira Naruse, Corey Nolet, Ray Wang, Tamas Feher, Yong Wang, https://arxiv.org/pdf/2308.15136
[7] BioSimCSE: BioMedical Sentence Embeddings using Contrastive learning, Kamal Raj Kanakarajan, Bhuvana Kundumani, Abhijith Neil Abraham, Malaikannan Sankarasubbu., https://aclanthology.org/2022.louhi-1.10.pdf

[8] RAGAS: Automated Evaluation of Retrieval Augmented Generation, Shahul Es, Jithin James, Luis Espinosa-Anke, Steven Schockaert., https://arxiv.org/abs/2309.15217v1

[9] Exploring Conceptual Soundness with TruLens Anupam Datta, Matt Fredrikson, Klas Leino, Kaiji Lu, Zifan Wang, Ricardo Shih, Shayak Sen, Truera., https://proceedings.mlr.press/v176/datta22a/datta22a.pdf