# Advanced Data Management and Retrieval Techniques: A Survey

www.surveyx.cn

## Abstract

In the evolving landscape of data management, advanced techniques such as vector databases, Approximate Nearest Neighbor (ANN) search, and embedding indexing are pivotal in handling high-dimensional data efficiently. This survey explores these techniques, emphasizing their significance in applications ranging from e-commerce to healthcare. Vector databases, exemplified by systems like FAISS and Milvus, are optimized for storing and querying high-dimensional vectors, addressing the 'curse of dimensionality' through innovative indexing and retrieval methods. ANN search algorithms, including Hierarchical Navigable Small World (HNSW), enhance retrieval speed and accuracy, crucial for large-scale datasets. Embedding indexing transforms data into dense vectors, facilitating semantic search and improving information retrieval systems' contextual understanding. Retrieval-Augmented Generation (RAG) integrates retrieval with generative models, enhancing information synthesis by leveraging external knowledge sources. The integration of knowledge graphs further augments this process, enabling structured representation and reasoning. Despite challenges such as data integration and scalability, these technologies offer transformative potential across domains. Future research directions include optimizing index management, exploring quantum vector databases, and enhancing retrieval-augmented generation frameworks. By addressing these challenges, advanced data management techniques promise to significantly enhance the efficiency and accuracy of modern data-driven environments, ensuring robust performance across diverse applications.

## 1 Introduction

### 1.1 Importance of Advanced Data Management

Advanced data management techniques are essential for effectively handling complex, large-scale data, particularly as modern applications increasingly depend on high-dimensional vector representations. Traditional database systems frequently struggle with high-dimensional data, prompting the development of specialized systems, such as vector database management systems (VDBMS), tailored for managing rich, unstructured numerical vectors common in applications like recommender systems and similarity searches [1].

The evolution of AI applications further emphasizes the need for advanced data management, as the synthesis of quantum vector databases has been proposed to enhance data processing and organization efficiency [2]. This necessity is evident in domains such as drug discovery, where molecular similarity searches illustrate the demand for rapid identification of structurally similar compounds from extensive molecular databases, necessitating efficient management and retrieval methods for high-dimensional data [3].
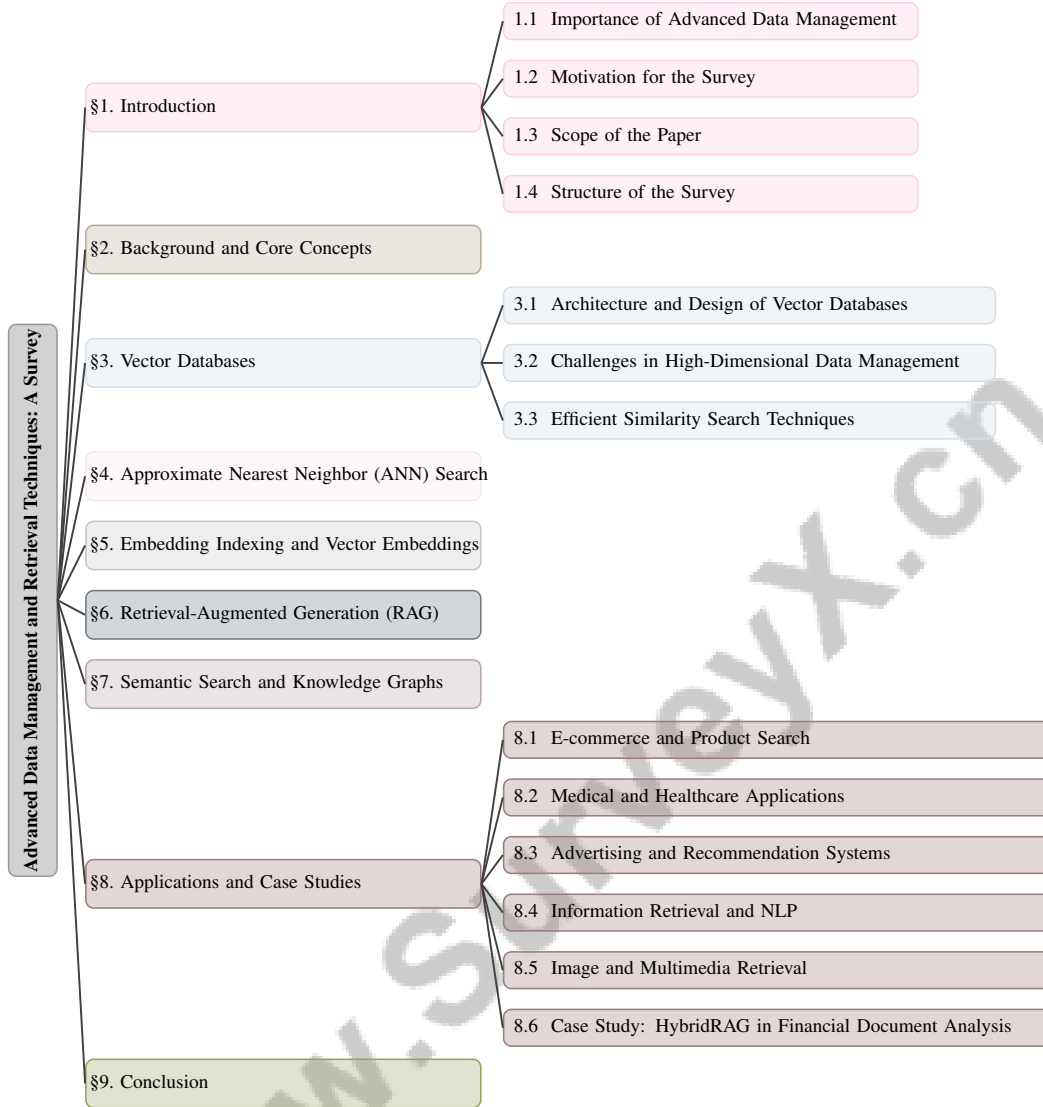
Figure 1: chapter structure

## 1.2 Motivation for the Survey

This survey is motivated by the significant challenges encountered by current data management systems, particularly in the efficiency and accuracy of information retrieval [4]. Traditional frameworks often fail to leverage large language models (LLMs) effectively in scenarios with limited domain-specific training data, leading to suboptimal retrieval outcomes [5]. This work seeks to explore advanced techniques to enhance LLM performance by addressing their inherent knowledge, memory, and capability limitations [6].

In investment management, the integration of Environmental, Social, and Governance (ESG) factors necessitates computationally efficient models for responsible investing [7]. This survey aims to identify and discuss advanced data management solutions that meet these computational demands.

Furthermore, enhancing retrieval methods for LLMs is crucial, as existing limitations hinder their effectiveness [6]. By providing a comprehensive overview of advanced data management and retrieval techniques, this survey aims to highlight their relevance and potential in addressing pressing challenges within modern data-driven environments.

## 1.3 Scope of the Paper

This survey thoroughly explores advanced data management and retrieval techniques, focusing on vector databases, embedding models, and retrieval-augmented generation (RAG) applications. It examines the development of vector-based knowledge bases from substantial unstructured environmental data, which is vital for precise data retrieval in RAG contexts [8]. Empirical benchmarks comparing various vector search implementations provide insights into trade-offs related to indexing time, query evaluation performance, and retrieval quality [9].

The scope includes an examination of retrieval-based language models (LMs) and their architectures, learning approaches, and applications, while excluding unrelated LMs that do not incorporate retrieval mechanisms [10]. The integration of vector databases and generative language models (GLMs) in educational assessments, particularly for scoring short answers, is also explored [11]. Additionally, the survey addresses user interactions in retrieval systems, focusing on embedding-based methods and their constraints in e-commerce search scenarios [12].

Recent advancements in VDBMSs are covered, including query processing, storage, indexing, and optimization techniques, while excluding traditional database management systems not specialized in vector data [1]. The survey also investigates the architectures, training strategies, and applications of retrieval-augmented language models (RA-LLMs), emphasizing the integration of retrieval mechanisms with LLMs [13]. Various indexing methods and configurations within the Faiss library are evaluated, aiding the selection of optimal approaches for diverse use cases in vector similarity search [14].

Moreover, the survey addresses the limitations of previous tuning-based and retriever-based methods, proposing innovations to enhance tool selection and retrieval processes [15]. It explores a hybrid context retrieval augmented generation pipeline that combines a vector database with a knowledge graph to assist institutions in creating accreditation reports [16]. Outlier detection techniques are examined to improve context retrieval for complex queries [17].

The automation and enhancement of bibliometric analysis using Generative AI models, particularly transformers and RAG, are also discussed [18]. Challenges and limitations in handling LLMs across specific application domains are highlighted [4]. This comprehensive examination aims to provide a detailed understanding of the current landscape and future directions in advanced data management and retrieval techniques.

## 1.4 Structure of the Survey

This survey is systematically organized to examine advanced data management and retrieval techniques comprehensively. The introduction outlines the significance of these techniques in modern data processing environments, followed by discussions on the motivation and scope of the survey. The background section elucidates core concepts, including vector databases, FAISS, Milvus, ANN search, embedding indexing, RAG, HNSW, semantic search, vector embeddings, and knowledge graphs, establishing a foundational understanding.

Subsequent sections delve into specific aspects of the topic. Section 3 explores vector databases, focusing on architecture, design, and challenges in managing high-dimensional data, highlighting efficient similarity search techniques employed by systems like FAISS and Milvus. Section 4 addresses Approximate Nearest Neighbor (ANN) Search, discussing its importance and key algorithms, such as HNSW, along with a comparative analysis of different approaches.

In Section 5, the survey examines embedding indexing and vector embeddings, detailing techniques for creating embeddings and optimizing indexing for improved retrieval speed and accuracy. Section 6 introduces RAG, detailing its conceptual framework, including the integration of vector databases for context retrieval, and discusses quality assurance mechanisms, such as a Bayesian approach, to evaluate the relevance of text chunks. Practical applications of RAG in domain-specific question answering are exemplified by a case study on Pittsburgh and Carnegie Mellon University, showcasing significant improvements in answer precision and relevance through enhanced document retrieval techniques [19, 20].

Section 7 delves into advancements in semantic search and the role of knowledge graphs, detailing how vector embeddings significantly enhance semantic search capabilities. The integration of knowledge graphs with vector retrieval techniques, such as HybridRAG, is highlighted, emphasizing

3

the importance of leveraging dense semantic representations and multi-categorization semantic analysis to optimize information retrieval across various applications, including e-commerce and financial data analysis [21, 12, 22, 23, 24]. Section 8 presents real-world applications and case studies, illustrating the effectiveness of these techniques across various domains such as e-commerce, healthcare, advertising, information retrieval, and multimedia.

The survey concludes with key findings and a discussion on future research directions and innovations in advanced data management. It seeks to deliver critical insights into the current research landscape and identify potential advancements, particularly by leveraging methodologies such as Generative AI models and RAG systems. These innovative approaches enhance bibliometric analysis by facilitating semantic searches and contextual understanding, uncovering valuable insights that traditional methods often overlook [19, 18].The following sections are organized as shown in Figure 1.

## 2 Background and Core Concepts

### 2.1 Vector Databases and Dimensionality Challenges

Vector databases are pivotal for handling high-dimensional data in applications like NLP, image recognition, and recommendation systems, optimizing the storage and querying of vector embeddings that numerically encapsulate semantic information. Managing high-dimensional data poses significant challenges, especially in RAG for LLMs, where effective data parsing and retrieval are crucial for context-specific responses. Techniques such as outlier detection and learned indexes like LIDER enhance retrieval by identifying semantically relevant documents and improving search efficiency, addressing limitations in traditional methods that struggle with high-dimensional embeddings [8, 17, 25, 26].

The 'curse of dimensionality' adversely affects distance-based similarity searches as dimensionality increases, necessitating dimensionality reduction techniques for improved storage efficiency and retrieval speed [1, 27]. Traditional database systems often fail to meet these requirements, especially as applications scale to billion-scale vector collections [28].

Quantum computing principles, such as Grover's algorithm, have been proposed to enhance querying efficiency in high-dimensional data management [2]. Despite advancements, misconceptions persist that enterprises need to overhaul existing systems for vector search capabilities, while frameworks like Lucene can effectively manage vectorized data [24].

Proximity graph-based methods excel in ANN searches but lack efficient mechanisms for updating graph indices during vertex deletions, highlighting the need for innovative approaches in dynamic environments [29]. Traditional semantic search methods often overlook the multi-category nature of documents, necessitating novel strategies for enhancing vector database efficacy [21].

In enterprise settings, managing diverse data sources and ensuring contextual data retrieval is crucial, emphasizing the integration of external knowledge bases with LLMs for complex data tasks. Implementing kNN retrieval augmentation can improve generalization from over-specified training data, enhancing retrieval system flexibility and efficiency [30]. Managing large datasets, such as medical images, underscores the need for advanced retrieval methods to support timely decision-making [31]. Additionally, integrating vector databases with Sentence Transformers and Retrieval Agents enhances contextual search and topic ranking, showcasing potential improvements in data management capabilities [18].

### 2.2 Approximate Nearest Neighbor (ANN) Search Techniques

ANN search techniques are vital for efficiently managing large-scale, high-dimensional datasets where traditional exhaustive search methods are computationally prohibitive. These techniques are essential in applications like image and text retrieval, recommendation systems, and RAG, where rapid identification of data points closest to a query vector is crucial. Recent advancements, including the Speed-ANN algorithm, utilize intra-query parallelism to enhance search efficiency across large datasets, achieving improved query latency compared to traditional methods like NSG and HNSW. Competitions such as the 2023 Big ANN Challenge have highlighted innovative approaches to ANN search that address diverse workloads and constraints, yielding substantial improvements in accuracy and efficiency. Solutions tailored for resource-constrained edge devices emphasize adapting ANN

4

techniques to meet real-world requirements while maintaining low latency and operational efficiency [32, 33, 34, 35].

Graph-based approaches, notably HNSW networks, have emerged as leading solutions due to their ability to construct efficient nearest neighbor graphs that enhance search speed and accuracy. However, these methods often face challenges related to memory consumption and scalability as datasets expand, necessitating continuous advancements to maintain optimal performance [36].

Hashing-based strategies, such as Locality-Sensitive Hashing (LSH), project high-dimensional data into lower-dimensional spaces to facilitate faster approximate searches. The DB-LSH framework exemplifies this by employing a dynamic bucketing strategy for efficient ANN searches in high-dimensional datasets [37]. Nonetheless, achieving high recall rates without significant computational costs remains challenging, particularly in scenarios demanding small approximation factors [33].

Innovations like Speed-ANN leverage parallel algorithms to enhance nearest neighbor search efficiency by utilizing multi-core architectures and optimizing search processes, addressing latency and accuracy issues [35]. These advancements underscore the ongoing evolution of ANN search techniques, increasingly tailored to the demands of contemporary data environments.

Despite these innovations, limitations persist in current ANN search methods, such as spill trees and vector quantization approaches, which often struggle with scalability and efficiency in large datasets [38]. Addressing these limitations is critical for advancing ANN search capabilities, ensuring robustness and effectiveness in handling modern data retrieval complexities.

## 2.3 Embedding Indexing and Vector Embeddings

Embedding indexing and vector embeddings are crucial for contemporary data retrieval systems, enabling efficient and semantically meaningful access to large-scale datasets. These techniques transform raw data into dense vector representations that encapsulate semantic relationships, facilitating similarity search, classification, and clustering. Reliable vector embeddings are particularly important in domains like biomedical literature, where existing models often fail to capture domain-specific nuances [39].

Creating vector embeddings typically involves training models to map input data—such as text or images—into a continuous vector space where semantically similar inputs are positioned closely together, enhancing retrieval accuracy. The Bi-Granular Document Representation, which uses both lightweight sparse embeddings for candidate search and heavyweight dense embeddings for fine-grained verification, illustrates a layered approach to embedding creation [40].

Vector databases are vital for storing these embeddings, allowing for efficient querying and retrieval based on proximity measures like cosine similarity. The integration of vector databases with RAG systems, such as RAGLog, exemplifies their role in enhancing retrieval accuracy through zero-shot semantic analysis of log entries for effective anomaly detection [41]. Hybrid systems like the Hybrid Context Retrieval Augmented Generation (HCRAG) pipeline combine vector embeddings with knowledge graphs to provide contextual information for generating responses, showcasing embedding indexing's versatility in complex data environments [16].

Dimensionality reduction techniques are often employed to manage the high-dimensional nature of embedding vectors, improving storage efficiency and retrieval speed. The use of vector databases in cloud-native applications, such as Shotit, an image-to-video search engine, highlights embeddings' role in enhancing search efficiency and reducing computational overhead [42]. Additionally, embedding-based approaches in recommendation systems that fuse embeddings from various data sources have been shown to improve job-to-candidate matching [43].

Advanced embedding methods continue to evolve, with libraries like string2string emphasizing embedding indexing and semantic search for efficient string-to-string problem-solving [44]. These innovations underscore embedding indexing and vector embeddings' foundational role in modern data retrieval systems, offering robust frameworks for managing and accessing high-dimensional data across diverse applications. Continuous refinement and innovation in these techniques ensure efficient and effective retrieval in increasingly complex and varied data environments.

# 3  Vector Databases

Understanding vector databases' architecture and design is crucial for appreciating their operational efficiencies and capabilities. This section delves into the structural intricacies of these databases, highlighting how innovative designs optimize high-dimensional data management. As illustrated in Figure 2, the hierarchical structure of vector databases emphasizes key aspects of their architecture and design, including the challenges associated with managing high-dimensional data and the efficient similarity search techniques employed. The architecture and design section specifically examines advanced systems such as Milvus and FAISS, while the challenges section outlines the 'curse of dimensionality' and presents innovative solutions like compact hash codes. Furthermore, the efficient similarity search techniques section details methods such as Approximate Nearest Neighbors (ANN) and advancements like SymphonyQG, showcasing the complexity and innovation inherent in vector database management. By examining their architecture, we can comprehend the mechanisms enabling rapid similarity searches and effective handling of extensive datasets, which are pivotal for various applications.
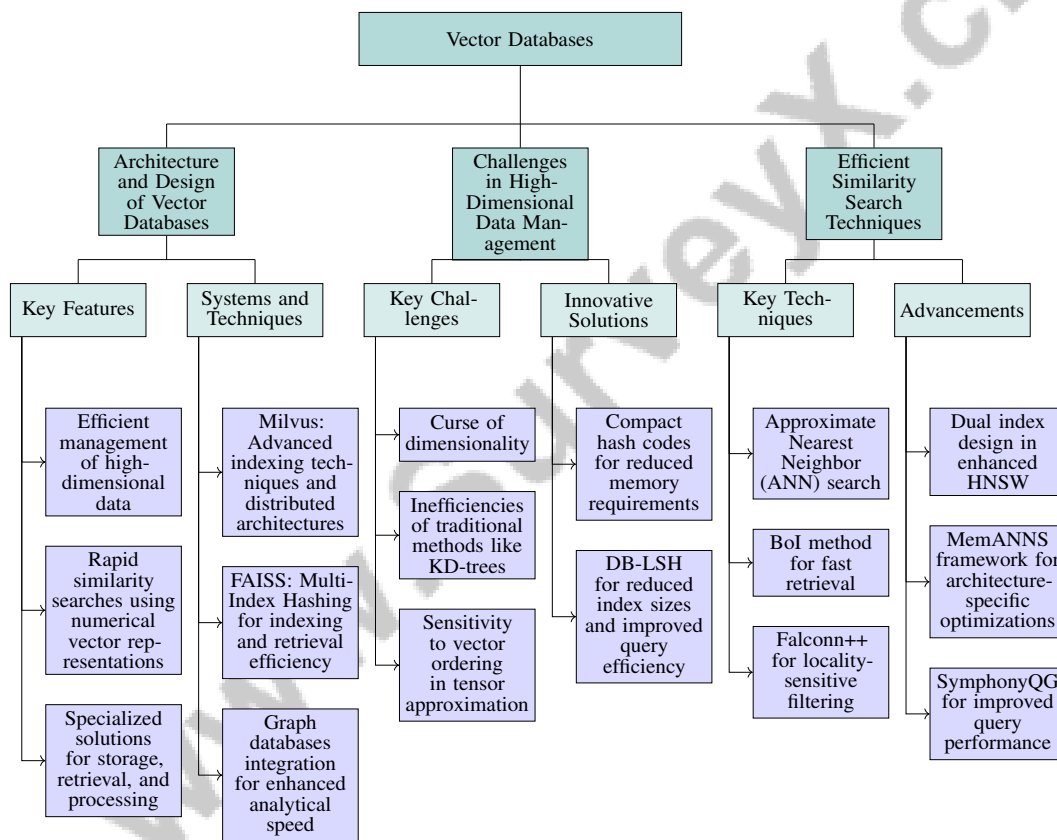
Figure 2: This figure illustrates the hierarchical structure of vector databases, highlighting key aspects of their architecture and design, challenges in managing high-dimensional data, and efficient similarity search techniques. The architecture and design section emphasizes advanced systems like Milvus and FAISS, while the challenges section outlines the 'curse of dimensionality' and innovative solutions like compact hash codes. The efficient similarity search techniques section details methods such as ANN and advancements like SymphonyQG, showcasing the complexity and innovation in vector database management.

## 3.1  Architecture and Design of Vector Databases

The architecture and design of vector databases are essential for efficiently managing high-dimensional data, particularly in applications such as LLMs, recommender systems, and chatbots. These databases facilitate rapid similarity searches by utilizing numerical vector representations

6

of complex data types, including text, images, and video. The high dimensionality and sparsity of vectorized data necessitate specialized solutions for efficient storage, retrieval, and processing. Recent advancements in vector database management systems (VDBMSs) have introduced innovative techniques for query processing, storage optimization, and hybrid query execution, addressing challenges like semantic similarity vagueness and high computational costs associated with similarity comparisons [45, 1].

Systems like FAISS and Milvus exemplify advanced structural designs necessary for vector data management. Milvus handles large-scale vector data efficiently through advanced indexing techniques and distributed architectures, allowing for independent scaling and optimization of system components [28]. FAISS, developed by Facebook AI, utilizes Multi-Index Hashing (MIH) to enhance indexing and retrieval efficiency with binary features, making it effective for applications such as job recommender systems [43]. Its architecture supports optimized designs for exhaustive and approximate search algorithms, facilitating rapid indexing and retrieval [3].

Innovative indexing structures are crucial for balancing memory efficiency with search accuracy. For example, SOAR improves approximate nearest neighbor (ANN) search efficiency by allowing data points to be assigned to multiple partitions, which is beneficial in GPU implementations [38]. Additionally, integrating graph databases with vector databases enhances analytical speed and supports complex relationship handling, especially in financial contexts [7].

The architecture of vector databases, such as FAISS and Milvus, is characterized by specialized capabilities for managing high-dimensional data, including advanced indexing techniques like HNSW and memory optimization through quantization methods. These features facilitate seamless interoperability with traditional databases, enhancing performance in similarity search, recommendation systems, and natural language processing [46, 47, 14, 45, 1]. As these systems evolve, they continue to incorporate innovative techniques that enhance their performance across diverse applications and environments.

## 3.2 Challenges in High-Dimensional Data Management

Managing high-dimensional data in vector databases presents significant challenges impacting efficiency and scalability. The 'curse of dimensionality' complicates the identification of exact nearest neighbors using brute-force methods. As dimensionality increases, the volume of space expands exponentially, leading to sparsity and difficulty in recognizing meaningful patterns or clusters within the data [38].

Traditional methods, such as KD-trees and brute-force searches, are inefficient in high-dimensional spaces due to excessive computational time and memory usage. The computational complexity is exacerbated by the need for exact distance calculations and inefficiencies associated with random memory accesses during querying and indexing [48]. These challenges are particularly pronounced in real-time data processing applications, such as mobile robotics, where efficient processing of high-dimensional visual data is critical [49].

Conventional tensor approximation methods also struggle with high-dimensional data representation and retrieval, often yielding poor results due to sensitivity to vector ordering [50]. This sensitivity complicates data representation and retrieval system effectiveness, underscoring the need for more robust approaches.

Long convergence steps in best-first search algorithms further limit speedup through multi-core processing [35]. The reliance on data-oriented partitioning trees, which are time-consuming to partition multi-dimensional spaces, results in decreased performance as dimensionality increases [51].

Innovative solutions, such as compact hash codes, offer reduced memory requirements, lower computational costs, and high retrieval performance with compact representations [52]. Techniques like DB-LSH, employing dynamic bucketing strategies, promise reduced index sizes and improved query efficiency, highlighting the ongoing evolution necessary to optimize vector databases for high-dimensional data management [37].

Addressing these multifaceted challenges is crucial for enhancing the performance and applicability of vector databases in modern data-intensive applications, such as recommender systems, natural language processing, and image retrieval. By developing specialized solutions that incorporate

7

advanced indexing, quantization, and hybrid query optimization, vector databases can better meet contemporary demands, ultimately improving reliability, speed, and scalability in handling vast amounts of unstructured data [45, 46, 1, 8]. Continuous refinement and innovation in these systems are essential for overcoming the inherent limitations of high-dimensional data management.

## 3.3 Efficient Similarity Search Techniques

Efficient similarity search techniques are crucial for enhancing vector databases' performance, especially when managing large-scale, high-dimensional datasets. Techniques such as Approximate Nearest Neighbor (ANN) search enable rapid retrieval of the k closest vectors to a query vector, essential for applications like eCommerce and recommendation systems. The effectiveness of these methods is significantly influenced by advanced indexing strategies, including graph-based approaches that optimize proximity graphs and novel relevance filtering mechanisms that enhance retrieval precision. As the demand for processing rich, unstructured data rises, developing sophisticated algorithms for efficient similarity search becomes increasingly important [1, 53, 45, 22, 54].

Notable advancements include the BoI method, which employs multi-index hashing to create multiple hash tables for fast retrieval without exhaustive distance calculations [55]. This significantly reduces computational overhead, making it suitable for large-scale applications. Similarly, the Falconn++ method utilizes locality-sensitive filtering to retain only relevant points during querying, optimizing both query time and indexing space [56].

The dual index design and MN-RU algorithm in enhanced HNSW index structures improve update speed and manage unreachable points, significantly enhancing the original HNSW approach [57]. Graph-based methods, such as reordering HNSW indices, improve cache efficiency and reduce query times by minimizing cache misses [58].

Architecture-specific optimizations, such as the MemANNS framework, incorporate architecture-aware data placement, efficient thread scheduling, and novel encoding methods to enhance the performance of the IVFPQ algorithm on PIM hardware [59]. This approach highlights the importance of tailoring similarity search techniques to specific hardware architectures for maximum efficiency.

The SymphonyQG method establishes a new benchmark in the time-accuracy trade-off for ANN search, achieving improved query performance and faster indexing through quantization integration [48]. Additionally, the Rii method utilizes a linear data layout for O(1) access to items by identifier, facilitating efficient subset searches through linear scans [60].

Efficient similarity search techniques in vector databases balance speed, accuracy, and resource utilization. Innovations in data structures, hybrid encoding methods, and optimized indexing strategies drive the evolution of Retrieval-Augmented Generation (RAG) techniques. These advancements enhance the performance of LLMs by improving the quality and relevance of retrieved information, reducing processing complexity, and enabling efficient parsing and vectorization of semi-structured data. Integrating vector databases with traditional full-text search engines allows for scalable and rapid semantic searching, while methodologies like HybridRAG combine knowledge graphs with vector retrieval to improve information extraction accuracy in specialized domains such as finance and environmental management. These continuous improvements ensure that RAG systems remain robust and effective in generating contextually rich and technically accurate responses [23, 8, 22, 61, 17].



(a) A Flowchart of a Query-Driven Search System[53]

(b) Comparison of Average Difference and Precision@10 for Different Page Sizes and Best Feature Selection[22]
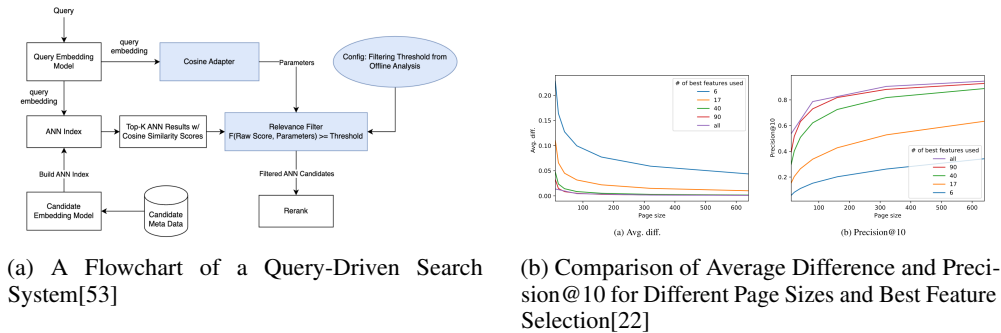
Figure 3: Examples of Efficient Similarity Search Techniques

As shown in Figure 3, two illustrative examples in vector databases and efficient similarity search techniques provide insights into optimizing search systems. The first example, a flowchart, outlines a query-driven search system that begins with embedding a query and adjusting its parameters using a cosine adapter to produce a relevance score. This score is filtered and reranked to enhance search results, utilizing a candidate embedding model and an approximate nearest neighbor (ANN) index. The systematic approach ensures efficient retrieval of the most relevant results. The second example offers a comparative analysis of average difference and Precision@10 across varying page sizes, highlighting how different page sizes and feature selection methods impact search result precision and accuracy. Together, these examples underscore the importance of systematic query processing and performance evaluation in developing robust similarity search techniques within vector databases [53, 22].

# 4 Approximate Nearest Neighbor (ANN) Search

| Category | Feature | Method |
|---|---|---|
| **Overview of ANN Search** | Parallel Processing Techniques | SANN[35] |
| **Key Algorithms and Methods** | Hierarchical and Layered Structures | HNSW[36], FPGA-MSS[3] |
| **Comparative Analysis and Trade-offs** | Residual and Encoding Techniques<br>Decomposition and Compression<br>Cluster and Assignment Strategies<br>Embedding and Optimization Methods | SOAR[38], STC[62]<br>TT-PCC[50]<br>m-k-means[52]<br>N/A[6] |

Table 1: This table provides a comprehensive summary of various methods and techniques employed in Approximate Nearest Neighbor (ANN) search. It categorizes the methods into three main areas: an overview of ANN search, key algorithms and methods, and comparative analysis with trade-offs, highlighting specific features and corresponding methods. The table serves as a reference for understanding the diverse approaches and their applications in high-dimensional data retrieval.

| Category | Feature | Method |
|---|---|---|
| **Overview of ANN Search** | Parallel Processing Techniques | SANN[35] |
| **Key Algorithms and Methods** | Hierarchical and Layered Structures | HNSW[36], FPGA-MSS[3] |
| **Comparative Analysis and Trade-offs** | Residual and Encoding Techniques<br>Decomposition and Compression<br>Cluster and Assignment Strategies<br>Embedding and Optimization Methods | SOAR[38], STC[62]<br>TT-PCC[50]<br>m-k-means[52]<br>N/A[6] |

Table 2: This table provides a comprehensive summary of various methods and techniques employed in Approximate Nearest Neighbor (ANN) search. It categorizes the methods into three main areas: an overview of ANN search, key algorithms and methods, and comparative analysis with trade-offs, highlighting specific features and corresponding methods. The table serves as a reference for understanding the diverse approaches and their applications in high-dimensional data retrieval.

Approximate Nearest Neighbor (ANN) search is vital for efficient data retrieval in high-dimensional datasets. Table 2 presents a detailed summary of the methodologies and techniques used in Approximate Nearest Neighbor (ANN) search, categorizing them into distinct areas for a clearer understanding of their applications and trade-offs. Additionally, Table 4 offers a comprehensive comparison of various methods employed in Approximate Nearest Neighbor (ANN) search, elucidating their distinct features and operational efficiencies. This section explores ANN search's significance and principles, establishing a basis for key algorithms and methods that drive advancements in this domain. ANN search addresses large-scale dataset challenges.

## 4.1 Overview of ANN Search

ANN search is integral to modern data retrieval systems, identifying data points near a query vector. Critical in image and text retrieval, recommendation systems, and retrieval-augmented generation, ANN search overcomes the computational challenges of exact nearest neighbor searches in high-dimensional vector databases. Techniques like graph-based indexing enhance accuracy and efficiency, improving user experiences in platforms like search engines and eCommerce [53, 54].

Traditional methods face inefficiencies in large datasets, particularly at the billion-scale level, where memory and latency constraints impede performance [35]. Techniques like SOAR improve ANN search indexing, enabling efficient retrieval of k nearest neighbors from d-dimensional vectors [38].

Graph-based methods advance ANN search by enhancing scalability and adaptability to dynamic data changes. These methods leverage Bayesian inference and optimized RAG modules to improve retrieved information quality [20, 18, 63, 61, 17]. Integrating ANN search with quantum computing exemplifies ongoing advancements to meet data demands.

Modern ANN search techniques demonstrate versatility across technological landscapes. For example, executing ANN search on edge devices for entity resolution in digital voice assistants shows adaptability and robustness. Innovations from the NeurIPS 2023 Big ANN Challenge highlight improvements in indexing structures and search algorithms for diverse workloads, including filtered searches and out-of-distribution data. Frameworks like LIDER enhance search efficiency with high-dimensional learned indices, while Speed-ANN leverages multi-core processing for low-latency, high-accuracy searches, underscoring ANN techniques' critical role in data management [25, 35, 34, 40, 54].

## 4.2 Key Algorithms and Methods

Advanced algorithms have transformed ANN search, crucial for managing large-scale, high-dimensional datasets. The Hierarchical Navigable Small World (HNSW) algorithm exemplifies this transformation, using a multi-layered graph structure for rapid, scalable nearest neighbor searches. This architecture allows dynamic updates and efficient high-dimensional space navigation, suitable for real-time applications [36]. HNSW's hierarchical structure outperforms flat graph-based approaches, especially on computational storage platforms that minimize data movement [3].

Graph-based methods, like Incremental Proximity Graph Maintenance (IPGM), enhance online operations by supporting vertex insertions and deletions in proximity graphs, essential for adapting to dynamic datasets in real-time applications [64, 29, 65, 66]. The NDSEARCH framework addresses ANNS method scaling inefficiencies by processing workloads within SSDs, utilizing LUN-level parallelism for improved data locality and bandwidth utilization.

Quantization strategies optimize ANN search efficiency. SymphonyQG integrates advanced quantization techniques with graph-based indexing, setting new benchmarks in time-accuracy trade-offs. Locality-sensitive hashing (LSH) remains key, with advancements like DB-LSH employing dynamic bucketing to enhance efficiency and reduce index sizes. DB-LSH organizes projected spaces with multi-dimensional indexes, allowing rapid high-quality candidate point generation and reduced query costs of $O(n^* d \log n)$ [51, 37].

Partitioning strategies optimize ANN search, addressing imbalanced partitioning challenges. Framing candidate set selection as a multilabel classification problem enhances retrieval performance, effectively grouping nearest neighbors of a query point within the same partition [67, 63]. The kNN-Embed method improves diversity in dense ANN-based retrieval by using a mixture distribution of user preferences.

Innovative methods like SOAR (Spilling with Orthogonality-Amplified Residuals) improve indexing quality and search efficiency [38]. The LLM-Embedder model unifies diverse retrieval tasks, optimizing training methodology to improve performance across scenarios [6]. These advancements highlight efforts to balance speed, accuracy, and scalability in ANN search systems, ensuring applicability in diverse data-intensive environments.

## 4.3 Comparative Analysis and Trade-offs

| Benchmark | Size | Domain | Task Format | Metric |
|---|---|---|---|---|
| BEIR[9] | 2,681,468 | Information Retrieval | Document Retrieval | nDCG@10 |
| EduBot[68] | 20 | Education | Question Answering | FAITH, $ANS_R EL$ |
| RAMAP[69] | 60,000 | Image Retrieval | Hashing Evaluation | RAMAP |
| RAG-Coding[70] | 150,523 | Program Synthesis | Commit Message Generation | RougeL |
| Thistle[71] | 1,000,000 | Information Retrieval | Query Result Retrieval | Accuracy, Runtime |
| Anserini-HNSW[72] | 8,800,000 | Information Retrieval | Passage Retrieval | MRR@10, Recall@1k |
| CoTAB[73] | 420,000 | Conversational AI | Question Answering | Recall, F2 |
| RGB[74] | 600 | Question Answering | Question Answering | Accuracy, Rejection Rate |

Table 3: This table summarizes various benchmarks relevant to the evaluation of ANN search algorithms across diverse domains. It includes details on the size, domain, task format, and evaluation metrics for each benchmark, providing a comprehensive overview of the datasets used in recent research.

Selecting an optimal ANN search algorithm requires understanding trade-offs between computational efficiency, accuracy, and resource utilization. Table 3 presents a detailed comparison of benchmarks used to assess the performance and trade-offs of different ANN search algorithms, highlighting the diversity in size, domain, task format, and evaluation metrics. Different methodologies offer distinct advantages and limitations critical for optimizing performance across application contexts. Advancements in the 2023 Big ANN Challenge emphasize the need for tailored indexing structures and search algorithms to accommodate complex workloads, including filtered searches and out-of-distribution data. Research on relevance filtering and operational choices between dense and sparse retrieval methods significantly impacts search accuracy, efficiency, and user experience, crucial for enhancing retrieval-augmented generation systems and AI applications [53, 34, 17, 9].

Graph-based approaches, like HNSW, are recognized for dynamic update capabilities and rapid search performance in high-dimensional spaces, though they may incur substantial memory overheads, posing constraints in resource-limited environments. The IPGM algorithm boosts query processing speed and efficiency by maintaining proximity graphs for online ANN searches [38]. Despite improvements, the performance gap between hierarchical and flat models in high-dimensional datasets can be minimal, necessitating further optimization exploration [50].

Quantization techniques minimize memory usage and accelerate query times. SymphonyQG exemplifies this integration, eliminating the explicit re-ranking step found in previous methods, while DB-LSH employs dynamic bucketing for enhanced indexing and querying efficiency [52].

Advanced algorithms like SOAR provide optimized assignments for data points, enhancing search efficiency and accuracy compared to traditional methods. Hybrid search methods, using sparse ternary codes, offer superior performance in the complexity-memory trade-off compared to conventional dense binary hashing schemes [62].

The LLM-Embedder model illustrates embedding-based approaches' potential, significantly outperforming general-purpose and task-specific models, highlighting embedding strategies' importance in optimizing ANN search performance across diverse environments [6].

Selecting an ANN search algorithm depends on application-specific demands, including retrieval speed, result accuracy, memory consumption, and computational efficiency. As workloads grow more complex, especially in real-world machine learning scenarios, careful consideration of these parameters is essential. Recent advancements in constrained optimization techniques for tuning ANN algorithms emphasize optimizing the speed-recall trade-off for superior performance. Challenges highlighted in the 2023 Big ANN Challenge stress the need for innovative solutions capable of handling filtered searches and out-of-distribution data while maintaining efficiency and accuracy [75, 17, 34]. Understanding each approach's inherent trade-offs allows practitioners to select the most suitable method for their use case, ensuring optimal performance across diverse data environments.

| Feature | Hierarchical Navigable Small World (HNSW) | Incremental Proximity Graph Maintenance (IPGM) | SymphonyQG |
|---|---|---|---|
| Memory Usage | High Memory Overhead | Not Specified | Minimized Memory Usage |
| Search Efficiency | Rapid Search Performance | Boosts Query Processing | Accelerates Query Times |
| Dynamic Update Capability | Supports Dynamic Updates | Supports Vertex Operations | Not Specified |

Table 4: This table provides a comparative analysis of three methods used in Approximate Nearest Neighbor (ANN) search: Hierarchical Navigable Small World (HNSW), Incremental Proximity Graph Maintenance (IPGM), and SymphonyQG. It highlights key features such as memory usage, search efficiency, and dynamic update capability, offering insights into the trade-offs and performance characteristics of each method.

# 5 Embedding Indexing and Vector Embeddings

Embedding indexing is crucial for modern data retrieval systems, as it enables efficient semantic relationship representation within datasets. Techniques such as Approximate Nearest Neighbor (ANN) search facilitate the retrieval of similar items from extensive datasets while ensuring high precision. Methods like the 'Cosine Adapter' optimize cosine similarity scores for better relevance filtering, and outlier detection enhances context quality in question-answering systems, improving user experience and retrieval accuracy across various applications, including e-commerce and large language models [53, 17].

## 5.1 Techniques for Embedding Creation

Vector embedding creation underpins modern data retrieval by enabling similarity searches, classification, and clustering through dense vector representations that encapsulate semantic relationships. Advanced techniques, such as HybridRAG, which integrates knowledge graphs with vector retrieval, and Bi-Granular Document Representation, optimize memory usage for large-scale corpus retrieval. These methods address challenges such as accuracy and relevance enhancement, risk reduction of irrelevant data, and efficiency improvement in question-answering systems across domains like finance and urban research [19, 18, 40, 23, 17].

Neural network architectures play a pivotal role in embedding creation. The MoE extension pipeline enhances pretrained models like BERT by incorporating multiple distinct experts for each transformer block, allowing for effective domain adaptation and improved embedding quality [76]. Fine-tuning models with proprietary documents tailors vector embeddings to specific enterprise needs [43]. ChatSOS employs methods to enhance information retrieval accuracy and reliability [4].

The HNSW algorithm, with its fully graph-based incremental structure, builds multi-layer graphs with proximity graphs for nested data subsets, significantly improving approximate K-Nearest Neighbor Search efficiency [36]. This underscores the importance of structured embedding creation in complex data environments.

Quantization techniques like Poeem's product quantization of residual vectors reduce quantization distortion and enhance retrieval performance while preserving essential information and minimizing dimensionality [27]. FPGA-based HNSW implementation exemplifies hardware-accelerated embedding creation, achieving high query per second (QPS) rates with substantial recall on large databases [3].

Innovative training methods like kNN-Embed transform single user dense embeddings into mixtures, improving user interest capture and candidate retrieval diversity [77]. In multi-modal retrieval systems, adversarial hubs enhance relevance across diverse queries by perturbing benign inputs [21].

The integration of classical and quantum computing principles in embedding creation, illustrated by the synthesis of quantum vector databases, enhances embedding indexing and storage processes [2]. This approach signifies potential advancements in embedding efficiency and retrieval performance.

The evolving landscape of vector embedding techniques reflects efforts to optimize efficiency, accuracy, and scalability in data retrieval systems. Recent advancements include integrating semantic vector encoding with traditional full-text search engines, enhancing querying performance in extensive datasets like Wikipedia. Innovations such as contrastive learning and the Mixture of Experts (MoE) framework improve complex document representation, particularly in specialized domains like biomedical literature, facilitating precise vector embeddings across scientific fields. The Bi-Granular Document Representation addresses embedding-based retrieval challenges by utilizing lightweight sparse embeddings for initial candidate searches and heavyweight dense embeddings for fine-grained verification, significantly improving recall rates in large-scale corpora [22, 39, 40]. These techniques ensure robust and effective embeddings across diverse applications and data environments.



(a) A diagram illustrating a machine learning model for generating code snippets[70]

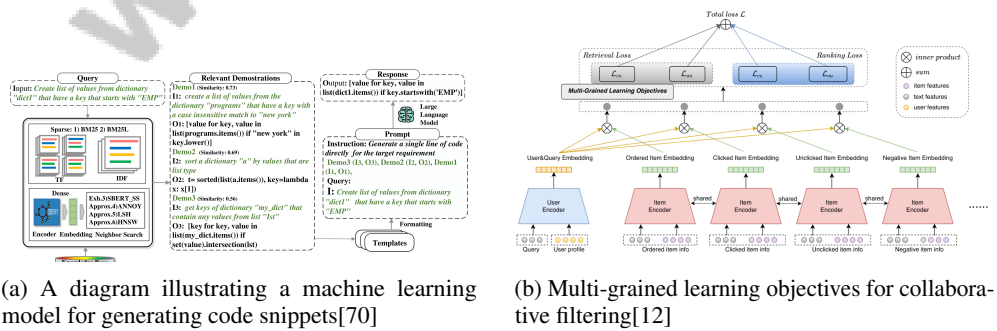(b) Multi-grained learning objectives for collaborative filtering[12]

Figure 4: Examples of Techniques for Embedding Creation

As shown in Figure 4, embedding indexing and vector embeddings enhance the efficiency and accuracy of information retrieval systems in machine learning. Transforming high-dimensional data

into lower-dimensional spaces facilitates effective data comparison and retrieval. The first example demonstrates a machine learning model generating code snippets based on user queries, retrieving relevant demonstrations from a database using a similarity scoring system. The second example illustrates a multi-grained learning objective framework for collaborative filtering, integrating user and query embeddings with item embeddings to optimize recommendations and refine prediction accuracy through a ranking loss component. These examples underscore the versatility and impact of embedding techniques in modern machine learning applications [70, 12].

## 5.2 Optimizing Embedding Indexing

Optimizing embedding indexing is essential for enhancing retrieval speed and accuracy in large-scale data systems. Effective strategies integrate advanced indexing techniques and algorithms to manage high-dimensional vector spaces efficiently. Dense hybrid representations (DHRs) combine dense lexical and semantic representations to optimize embedding indexing, improving retrieval speed and accuracy [27].

The GUITAR framework exemplifies the potential of gradient pruning techniques to reduce neural network evaluations, enhancing search speed and adaptability across various graph indices, making it suitable for diverse applications [78].

Incorporating adversarial techniques, as demonstrated by AdvBCT, plays a significant role in optimizing embedding indexing by aligning embedding distributions. This alignment allows new models to learn more discriminative features, improving retrieval accuracy and robustness, crucial for refining indexing strategies [79].

The integration of CNN and SURF features with modified HNSW graph structures, alongside novel geometrical verification strategies, exemplifies potential for optimizing memory usage and retrieval performance. This method employs binary hash codes to minimize memory consumption while maintaining high retrieval accuracy [49].

The NDSEARCH framework enhances embedding indexing by offloading graph traversal and distance computation tasks to an in-storage accelerator, SEARSSD, optimizing data access patterns. This significantly improves retrieval efficiency by streamlining data processing and reducing computational overhead [80].

Future research directions in embedding indexing optimization include exploring dynamic updates for graph structures and extending methods to support additional similarity metrics beyond Euclidean distance and cosine similarity. These advancements aim to optimize memory consumption and enhance the adaptability of indexing strategies [48].

## 6 Retrieval-Augmented Generation (RAG)

Retrieval-Augmented Generation (RAG) is a pivotal advancement in natural language processing, merging retrieval techniques with generative models to enhance their performance. This section explores the foundational principles and theoretical framework of RAG, emphasizing its mechanisms and operational strategies. By integrating retrieval capabilities with generative functions, RAG overcomes traditional model limitations, enabling innovative applications across diverse fields. The following subsection will delve into the specific concepts and framework of RAG, elucidating its structure and functionality.

### 6.1 Concept and Framework of RAG

RAG represents a sophisticated paradigm that combines retrieval mechanisms with generative models, significantly enhancing the contextual relevance and precision of language model outputs. Traditional models often lack the capacity to maintain extensive domain-specific knowledge, but RAG systems incorporate external knowledge sources, facilitating the generation of precise and contextually aware responses, thus improving performance on knowledge-intensive tasks [6].

The RAG framework involves retrieving relevant information from structured and unstructured data repositories, which guides the generative process of language models. This hybrid approach, integrating RAG techniques with Knowledge Graphs and vector retrieval methods, markedly improves

the precision of models in responding to domain-specific queries by providing pertinent contextual information. The use of multiple retrieval strategies and diverse annotated question-answer pairs enhances document retrieval accuracy, ensuring generated answers are both accurate and contextually relevant, demonstrated by significant performance improvements in evaluations [19, 23]. For example, systems like RAGLog utilize semantic relationships to improve log anomaly detection. Moreover, RAG's integration with multimodal models in tasks such as image geolocalization shows task accuracy improvements without further model training.

Central to the RAG framework is the retriever component, which identifies and retrieves relevant passages from large corpora. Advanced methodologies, such as the LLM-Embedder, consolidate diverse retrieval tasks into a single framework, optimizing retrieval processes for enhanced performance across various scenarios [6]. This adaptability is exemplified by frameworks like Hybrid Context Retrieval Augmented Generation (HCRAG), which synergistically integrates natural language queries with optimized retrieval methods and large language models (LLMs) to tackle complex tasks, such as accreditation reporting. HCRAG combines Knowledge Graph-based and vector retrieval techniques to enhance information extraction from complex documents, demonstrating superior performance in financial applications and broader domains [19, 81, 20, 23, 61]. Furthermore, outlier detection methods optimize context by filtering irrelevant documents, improving the quality and relevance of generated context.

RAG extends beyond traditional text-based applications, offering a comprehensive methodology for knowledge integration across multiple modalities, including structured data, unstructured text, and domain-specific information. This innovative approach enhances semantic and contextual understanding, enabling effective retrieval and analysis in diverse fields such as bibliometric studies, financial information extraction, and natural language processing. By leveraging advanced techniques like knowledge graphs and vector databases, RAG improves accuracy and relevance in information retrieval, broadening its applicability beyond conventional uses [18, 23, 81]. For instance, in bibliometric analysis, generative AI techniques automate contextual search and literature characterization, showcasing RAG's transformative potential in information synthesis.

## 6.2 Quality and Relevance in RAG Systems

RAG systems integrated with language models significantly enhance the quality and relevance of generated information by utilizing advanced retrieval techniques and external knowledge sources. A major advantage is the incorporation of structured relationships from knowledge graphs alongside broad contextual information from vector retrieval, ensuring comprehensive and accurate information extraction, critical for applications demanding high precision and context [23].

RAG systems like CRAG notably reduce token usage while maintaining high-quality responses, making them cost-effective for large-scale applications [82]. This efficiency is further bolstered by systems like RAGCache, which enhance cache efficiency and reduce computational overhead [83]. The ODCO method exemplifies retrieval process optimization by refining context through the elimination of semantically distant documents, thereby enhancing the relevance of generated responses [17].

The WRAG method improves retrieval accuracy by leveraging contextual information, emphasizing the most pertinent data sources based on query specifics [84]. This focus on contextual relevance is crucial for generating accurate outputs. Additionally, the AutoRAG framework systematically evaluates RAG setups to identify optimal modules for specific tasks and datasets, further enhancing RAG effectiveness [61].

Metrics in RAG systems evaluate the accuracy and relevance of generated responses, focusing on truthfulness and contextual appropriateness, ensuring reliable content production [63]. RAG systems thus represent a significant advancement in language model capabilities, enabling the generation of accurate and contextually relevant content across diverse applications.

## 6.3 Applications and Innovations in RAG

RAG systems have significantly advanced various fields by effectively integrating retrieval techniques, such as BM25 and FAISS, with generative models. This integration enhances contextual accuracy and relevance, as evidenced by improved performance metrics like F1 scores and recall in applications

14

such as domain-specific question answering and bibliometric analysis. RAG systems address common limitations of large language models, including hallucination and knowledge update challenges, by utilizing external knowledge databases to provide pertinent context, thus improving response quality and enabling nuanced answers in complex queries [19, 20, 18, 81].

In natural language processing (NLP), RAG systems have advanced conversational agents by effectively utilizing long-term memory and managing complex queries, enabling the generation of contextually relevant responses over extended interactions [73]. The development of benchmarks and retrieval models has further facilitated these applications, providing a robust foundation for enhancing conversational agents.

Recent innovations in RAG systems focus on optimizing retrieval mechanisms to enhance interpretability and broaden applications. The Atlas model exemplifies this trend by integrating advanced retrieval techniques to improve RAG performance across diverse NLP tasks [85]. Moreover, benchmarking large language models (LLMs) in RAG contexts has provided structured approaches for evaluating their capabilities, yielding valuable insights for future research and development [74].

Multimodal applications of RAG systems have expanded beyond traditional text and image modalities, incorporating innovative approaches to multimodal prompting and in-context learning. These advancements enable RAG systems to seamlessly process and generate outputs that integrate diverse data types [86]. The integration of Bayesian inference techniques represents a promising innovation in RAG systems, with ongoing research aimed at refining prior probability calculations and exploring additional features to enhance Bayesian filtering [20].

RAG systems have demonstrated significant improvements in retrieval-augmented generation by utilizing self-memory to achieve state-of-the-art results across multiple tasks, underscoring the importance of memory mechanisms in enhancing performance [87]. Future research is poised to explore the optimization of retrieval mechanisms, expansion of entity databases, and enhancement of query generation processes, further refining RAG capabilities for diverse applications and data environments.

The MoRSE model has shown superior performance, achieving over 10

The advancements and applications of RAG systems underscore their potential to transform information synthesis across various fields. By integrating sophisticated techniques such as Sentence Transformers and Bayesian inference, RAG enhances the accuracy and contextual relevance of content generated by large language models (LLMs). For instance, a pilot study demonstrated RAG's effectiveness in automating bibliometric analysis, enabling nuanced semantic searches and topic characterization in urban research. This innovative approach not only improves the retrieval of relevant information but also addresses traditional methods' limitations by ensuring high-quality and relevant text chunks, facilitating deeper insights and informed conclusions across diverse domains [20, 18]. Ongoing research and development position RAG systems to play an increasingly pivotal role in advancing modern data retrieval and generation frameworks.

## 7   Semantic Search and Knowledge Graphs

The integration of semantic search techniques with knowledge graph frameworks represents a significant advancement in information retrieval, necessitating a comprehensive understanding of their interplay. This section delves into the foundational aspects of this convergence, beginning with vector embeddings, which enhance semantic search by transforming data into dense numerical representations. These embeddings capture semantic relationships, facilitating the integration of knowledge graphs into information retrieval systems, thereby improving search accuracy and contextual relevance.

### 7.1   Leveraging Vector Embeddings for Semantic Search

Vector embeddings play a crucial role in advancing semantic search by providing a sophisticated framework for understanding complex semantic relationships, surpassing traditional keyword matching. By employing dense semantic representations and advanced indexing techniques, such as those in Elasticsearch and Lucene, vector embeddings enhance search performance and scalability, optimizing

15

applications across various domains [22, 24]. These embeddings convert data into dense numerical forms, capturing contextual nuances essential for effective retrieval.

The application of vector embeddings allows for the capture of intricate semantic relationships, yielding more accurate and contextually relevant search results. This transformation into a continuous vector space positions semantically similar items in proximity, enhancing retrieval through measures like cosine similarity. Integration with traditional systems, such as Elasticsearch, improves the speed and scalability of vector databases, surpassing conventional discrete indexing methods [22, 88]. Vector embeddings also enhance the interpretability of user queries, retrieving pertinent information even when query terms do not directly match indexed data.

Advanced techniques, including dense hybrid representations (DHRs), further optimize semantic search by merging dense lexical and semantic representations, improving retrieval accuracy and speed [27]. This highlights vector embeddings' potential to bridge lexical and semantic search paradigms, fostering a comprehensive understanding of user intent and data context. Their application extends to domains like natural language processing and information retrieval, where they improve the interpretability and relevance of search results, exemplified by their role in retrieval-augmented generation (RAG) systems [6].

## 7.2 Integration of Knowledge Graphs in Information Retrieval

The incorporation of knowledge graphs into information retrieval systems marks a substantial advancement in achieving accurate and contextually aware search results. Knowledge graphs provide a structured representation of entities and their interrelations, enabling retrieval systems to leverage rich semantic information beyond traditional keyword-based methods. This structure facilitates meaningful insights extraction and complex reasoning tasks, enhancing the overall efficacy of information retrieval systems [4].

Knowledge graphs improve retrieval by enabling systems to understand relationships between entities, allowing for precise query interpretation and result ranking. This capability is particularly valuable for ambiguous queries or those involving complex relationships that conventional methods struggle to capture. By utilizing the semantic structure of knowledge graphs, retrieval systems can infer additional context, delivering more relevant and comprehensive results [18].

Combining knowledge graphs with vector embeddings further augments the retrieval process by adding a semantic layer that enhances interpretability. This integration enables more effective semantic search, as knowledge graphs aid in understanding the broader context of queries, improving the relevance of retrieved information [2]. Additionally, knowledge graphs facilitate entity disambiguation in search queries, ensuring accurate identification and retrieval of intended information.

In RAG systems, knowledge graphs significantly enhance the quality and relevance of generated content by providing a structured framework for incorporating external knowledge, resulting in more accurate and contextually aware responses [6]. Approaches like HybridRAG, which merge vector databases and knowledge graph techniques, exemplify the enhanced search capabilities achieved by this integration, improving accuracy and contextual relevance in extracting information from complex unstructured data, such as financial documents and accreditation reports [16, 23]. By leveraging the rich semantic information in knowledge graphs, retrieval systems achieve a deeper understanding of user queries, delivering more relevant, contextually accurate results.

## 7.3 Challenges and Opportunities in Semantic Search and Knowledge Graphs

Despite their transformative potential, semantic search and knowledge graphs face challenges that must be addressed to fully realize their capabilities. A significant challenge is integrating heterogeneous data sources, which often differ in structure, format, and quality [18]. This diversity complicates the construction and maintenance of knowledge graphs, necessitating sophisticated data harmonization techniques to ensure consistency and accuracy in representing entities and relationships.

Scalability presents another challenge, particularly as data volumes grow exponentially. Efficiently managing and querying large-scale knowledge graphs requires advanced indexing and retrieval algorithms capable of handling their complexity and size without compromising performance [2]. Additionally, the dynamic nature of knowledge graphs, which are continuously updated, poses challenges in maintaining their accuracy and relevance over time.

Nevertheless, numerous opportunities exist for advancing semantic search and knowledge graph technologies. Integrating machine learning and artificial intelligence techniques offers promising avenues for automating the construction and updating of knowledge graphs, reducing manual effort and enhancing scalability [6]. Using vector embeddings alongside knowledge graphs can further enrich the semantic understanding of queries, yielding more accurate and contextually relevant search results.

The application of semantic search and knowledge graphs extends beyond traditional information retrieval, offering potential benefits in fields such as natural language processing, recommendation systems, and data analytics. By leveraging the rich semantic relationships encoded in knowledge graphs, these systems can provide nuanced insights and support complex reasoning tasks, enhancing decision-making processes across various domains [4].

# 8 Applications and Case Studies

Advanced data techniques are transforming various sectors, notably in e-commerce and product search, by enhancing user experiences and operational efficiencies. This section discusses case studies demonstrating these methodologies' real-world impact.

## 8.1 E-commerce and Product Search

Innovations such as Multi-level Multi-Grained Semantic Embeddings (MMSE) and Retrieval-Augmented Generation (RAG) have revolutionized e-commerce, improving retrieval efficiency and search accuracy for personalized user experiences. These techniques enhance understanding of user behavior, boost conversion rates, and increase shopping satisfaction [21, 12, 53, 18, 8]. Hybrid encoder approaches in native ad recommendations exemplify advanced data techniques' potential to refine ad targeting [89]. Systems like SnapMode efficiently process large-scale fashion data, enhancing product search and personalization [90].

Deep retrieval learning (DR) methods excel in e-commerce by learning retrievable representations from vast item catalogs, addressing industrial-scale data challenges [91]. Vector databases facilitate efficient similarity searches, enabling functionalities like image-based product searches that enhance user interaction [1]. Cloud-native vector databases like Manu support e-commerce applications by demonstrating versatility in recommendation systems and multimedia search [28].

## 8.2 Medical and Healthcare Applications

In medical and healthcare, advanced data techniques enhance data retrieval and analysis precision. Vector databases and embedding models improve medical text classification and diagnostic accuracy [92]. The DenseNet-FAISS method offers superior precision in medical image retrieval, particularly in BIRADS classification [93]. Retrieval-augmented methods enhance few-shot image classification, supporting informed healthcare decisions [94].

PLMs-based protein retrieval frameworks significantly impact bioinformatics by advancing protein analysis and retrieval, crucial for medical research [95]. Embedding spaces validated by ImageBind improve personalized medical services [96]. Fine-tuning models like LLaMA enhances healthcare applications by generating accurate, contextually relevant responses [26]. Integrating vector databases with large language models (LLMs) drives innovation and improves research outcomes [92, 8].

## 8.3 Advertising and Recommendation Systems

Advanced data management enhances advertising and recommendation systems, improving personalization and effectiveness. The NANN method on Taobao optimizes ad targeting, increasing advertising revenue by 3.1% [97]. ANN search methods are crucial for managing large-scale, high-dimensional datasets, enhancing user-item preference matching [53, 54, 97].

Vector embeddings and RAG techniques improve recommendation relevance by capturing semantic relationships between user behavior and product attributes. Methods like Multi-Categorization Semantic Analysis and kNN-Embed ensure diverse recommendations [21, 77, 17]. Vector database

systems and multi-categorization semantic analysis enhance targeting strategies, processing high-dimensional data efficiently [21, 1, 8].

## 8.4 Information Retrieval and NLP

Sophisticated algorithms in advanced data techniques improve information retrieval and NLP systems. The eCIL-MU framework enhances adaptability via a shift-to-the-nearest-class strategy [98]. Adaptive frame sampling methods in video moment search achieve superior retrieval performance, crucial for precise temporal localization [99]. The MS MARCO passage dataset benchmarks vector search effectiveness [24].

kNN-Embed improves recall and diversity in candidate retrieval tasks, emphasizing embedding-based approaches' effectiveness [77]. Adversarial hubs in multimodal retrieval systems craft hubs relevant to numerous queries, improving retrieval accuracy [100]. Advanced RAG-Tool Fusion enhances retrieval accuracy and agent performance without extensive model fine-tuning [15]. RAG and LLMs improve information retrieval and NLP systems, addressing traditional keyword search limitations [101, 18, 8, 19].

## 8.5 Image and Multimedia Retrieval

Advanced data techniques significantly enhance image and multimedia retrieval, focusing on efficiency and accuracy. Techniques like ANN search and RAG ensure rapid, precise retrieval of relevant multimedia content [19, 31, 53, 20, 17]. The Global-Local Image Retrieval competition method addresses image similarity retrieval challenges [102]. Shotit image-to-video search engine improves search speed and efficiency, achieving a 100x speedup [42].

Incorporating state-of-the-art retrieval techniques, such as deep hashing and vector databases, enhances content-based image retrieval (CBIR). RAG and frameworks like AutoRAG optimize retrieval processes, improving precision in specialized domains [5, 61, 31, 8]. Ongoing research continues to evolve image and multimedia retrieval, offering sophisticated solutions for managing large-scale multimedia datasets.

## 8.6 Case Study: HybridRAG in Financial Document Analysis

HybridRAG in financial document analysis exemplifies integrating retrieval-augmented generation with knowledge graphs to enhance information extraction from complex documents. Leveraging knowledge graphs' structured knowledge improves response accuracy and relevance [16]. HybridRAG efficiently manages large-scale financial datasets, crucial for timely decision-making processes, with performance comparable to advanced computational storage platforms [103].

Integrating HybridRAG with financial systems enhances data retrieval and processing speeds, providing a competitive edge. Its indexing methods facilitate rapid financial document processing, improving throughput and energy efficiency over traditional systems [80]. HybridRAG aligns with trends in advanced retrieval techniques for large-scale data environments, with DR methods achieving performance comparable to brute-force methods [91]. Future research could explore optimizations in memory usage and compressed representations for data vectors [104].

# 9  Conclusion

## 9.1 Future Directions and Innovations

Advancements in data management are poised to transform system capabilities, with future research focusing on optimizing index management through hierarchical storage-aware indexing and hardware integration. This approach promises to enhance retrieval efficiency and accuracy in high-dimensional data contexts. Further exploration of advanced embedding techniques and contextual factors in recommendation systems may lead to more personalized user experiences. Additionally, the exploration of alternative tree structures and their integration with Approximate Nearest Neighbor (ANN) methods could enhance database adaptability and update efficiency. Graph reordering techniques, coupled with hub-highway structures, offer promising avenues for improving graph traversal and edge pruning in ANN searches. Refinements in knowledge graph construction and document classification against

multiple standards are anticipated to enhance information retrieval systems. Future research will also likely focus on the deployment of local language models, optimization of token usage, and the transition to vector databases, aiming to improve system performance and alignment with user expectations. Real-time learning capabilities and multi-turn conversational workflows represent promising research areas, alongside search optimization algorithms on edge devices. In retrieval-augmented generation (RAG), optimizing frameworks like AutoRAG, exploring diverse datasets, and enhancing quantum vector databases for complex queries are critical for future advancements. These efforts are set to drive significant innovations in data management, addressing existing challenges and expanding the potential of data-driven technologies.

## 9.2 Large-Scale Data Challenges and Solutions

The management of large-scale data presents considerable challenges in retrieval, storage, and processing. Ensuring high performance in vector searches while integrating with SQL operations, as demonstrated by frameworks like SingleStore, is vital for applications handling both structured and unstructured data. Innovative solutions are required to balance performance and flexibility. Effective indexing structures, such as the inverted multi-index, offer potential improvements in grouping techniques, enhancing indexing efficiency and retrieval accuracy. Machine learning applications can benefit from optimized tensor-train decomposition processes, particularly for larger datasets and complex tasks, thus enhancing scalability and efficiency. Structured evaluation frameworks, essential for assessing the performance of large-scale data systems, ensure accuracy and relevance in applications like education. Addressing these challenges necessitates a multifaceted approach that combines advanced indexing strategies, seamless integration of diverse data operations, and machine learning optimization. By implementing these solutions, the effectiveness and efficiency of large-scale data management systems can be significantly enhanced, ensuring their adaptability and robustness across diverse domains.

# References

[1] Toni Taipalus. Vector database management systems: Fundamental concepts, use-cases, and current challenges, 2024.

[2] Cesar Borisovich Pronin and Andrey Vladimirovich Ostroukh. Synthesis of quantum vector databases based on grovers algorithm, 2023.

[3] Hongwu Peng, Shiyang Chen, Zhepeng Wang, Junhuan Yang, Scott A. Weitze, Tong Geng, Ang Li, Jinbo Bi, Minghu Song, Weiwen Jiang, Hang Liu, and Caiwen Ding. Optimizing fpga-based accelerator design for large-scale molecular similarity search, 2021.

[4] Haiyang Tang, Zhenyi Liu, Dongping Chen, and Qingzhao Chu. Chatsos: Llm-based knowledge qa system for safety engineering, 2023.

[5] Luiz Bonifacio, Hugo Abonizio, Marzieh Fadaee, and Rodrigo Nogueira. Inpars: Data augmentation for information retrieval using large language models. *arXiv preprint arXiv:2202.05144*, 2022.

[6] Peitian Zhang, Shitao Xiao, Zheng Liu, Zhicheng Dou, and Jian-Yun Nie. Retrieve anything to augment large language models. *arXiv preprint arXiv:2310.07554*, 2023.

[7] Partha Sen and Sumana Sen. Graph database while computationally efficient filters out quickly the esg integrated equities in investment management, 2024.

[8] Hang Yang, Jing Guo, Jianchuan Qi, Jinliang Xie, Si Zhang, Siqi Yang, Nan Li, and Ming Xu. A method for parsing and vectorization of semi-structured data used in retrieval augmented generation, 2024.

[9] Jimmy Lin. Operational advice for dense and sparse retrievers: Hnsw, flat, or inverted indexes?, 2024.

[10] Akari Asai, Sewon Min, Zexuan Zhong, and Danqi Chen. Retrieval-based language models and applications. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 6: Tutorial Abstracts)*, pages 41–46, 2023.

[11] Zifan Wang and Christopher Ormerod. Generative language models with retrieval augmented generation for automated short answer scoring, 2024.

[12] Binbin Wang, Mingming Li, Zhixiong Zeng, Jingwei Zhuo, Songlin Wang, Sulong Xu, Bo Long, and Weipeng Yan. Learning multi-stage multi-grained semantic embeddings for e-commerce search, 2023.

[13] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. Replug: Retrieval-augmented black-box language models. *arXiv preprint arXiv:2301.12652*, 2023.

[14] Matthijs Douze, Alexandr Guzhva, Chengqi Deng, Jeff Johnson, Gergely Szilvasy, Pierre-Emmanuel Mazaré, Maria Lomeli, Lucas Hosseini, and Hervé Jégou. The faiss library, 2025.

[15] Elias Lumer, Vamse Kumar Subbiah, James A. Burke, Pradeep Honaganahalli Basavaraju, and Austin Huber. Toolshed: Scale tool-equipped agents with advanced rag-tool fusion and tool knowledge bases, 2024.

[16] Candace Edwards. Hybrid context retrieval augmented generation pipeline: Llm-augmented knowledge graphs and vector database for accreditation reporting assistance, 2024.

[17] Vitaly Bulgakov. Optimization of retrieval-augmented generation context with outlier detection, 2024.

[18] Haowen Xu, Xueping Li, Jose Tupayachi, Jianming, Lian, and Femi Omitaomu. Automating bibliometric analysis with sentence transformers and retrieval-augmented generation (rag): A pilot study in semantic and contextual search for customized literature characterization for high-impact urban research, 2024.

[19] Haojia Sun, Yaqi Wang, and Shuting Zhang. Retrieval-augmented generation for domain-specific question answering: A case study on pittsburgh and cmu, 2024.

[20] Dattaraj Rao. Bayesian inference to improve quality of retrieval augmented generation, 2024.

[21] Yinglong Ma and Moyi Shi. Using multi-categorization semantic analysis and personalization for semantic search, 2014.

[22] Jan Rygl, Jan Pomikálek, Radim Řehůřek, Michal Růžička, Vít Novotný, and Petr Sojka. Semantic vector encoding and similarity search using fulltext search engines, 2017.

[23] Bhaskarjit Sarmah, Benika Hall, Rohan Rao, Sunil Patel, Stefano Pasquali, and Dhagash Mehta. Hybridrag: Integrating knowledge graphs and vector retrieval augmented generation for efficient information extraction, 2024.

[24] Jimmy Lin, Ronak Pradeep, Tommaso Teofili, and Jasper Xian. Vector search with openai embeddings: Lucene is all you need, 2023.

[25] Yifan Wang, Haodi Ma, and Daisy Zhe Wang. Lider: An efficient high-dimensional learned index for large-scale dense passage retrieval, 2022.

[26] Mathav Raj J, Kushala VM, Harikrishna Warrier, and Yogesh Gupta. Fine tuning llm for enterprise: Practical guidelines and recommendations, 2024.

[27] Sheng-Chieh Lin and Jimmy Lin. A dense representation framework for lexical and semantic matching, 2023.

[28] Rentong Guo, Xiaofan Luan, Long Xiang, Xiao Yan, Xiaomeng Yi, Jigao Luo, Qianya Cheng, Weizhi Xu, Jiarui Luo, Frank Liu, Zhenshan Cao, Yanliang Qiao, Ting Wang, Bo Tang, and Charles Xie. Manu: A cloud native vector database management system, 2022.

[29] Zhaozhuo Xu, Weijie Zhao, Shulong Tan, Zhixin Zhou, and Ping Li. Proximity graph maintenance for fast online nearest neighbor search, 2022.

[30] Ting-Rui Chiang, Xinyan Velocity Yu, Joshua Robinson, Ollie Liu, Isabelle Lee, and Dani Yogatama. On retrieval augmentation and the limitations of language model training. *arXiv preprint arXiv:2311.09615*, 2023.

[31] Deng Cai, Xiuye Gu, and Chaoqi Wang. A revisit on deep hashings for large-scale content based image retrieval, 2017.

[32] MD Shaikh Rahman, Syed Maudud E Rabbi, and Muhammad Mahbubur Rashid. Optimizing domain-specific image retrieval: A benchmark of faiss and annoy with fine-tuned features, 2024.

[33] Jianwei Zhang, Helian Feng, Xin He, Grant P. Strimel, Farhad Ghassemi, and Ali Kebarighotbi. Search optimization with query likelihood boosting and two-level approximate search for edge devices, 2023.

[34] Harsha Vardhan Simhadri, Martin Aumüller, Amir Ingber, Matthijs Douze, George Williams, Magdalen Dobson Manohar, Dmitry Baranchuk, Edo Liberty, Frank Liu, Ben Landrum, Mazin Karjikar, Laxman Dhulipala, Meng Chen, Yue Chen, Rui Ma, Kai Zhang, Yuzheng Cai, Jiayang Shi, Yizhuo Chen, Weiguo Zheng, Zihao Wan, Jie Yin, and Ben Huang. Results of the big ann: Neurips'23 competition, 2024.

[35] Zhen Peng, Minjia Zhang, Kai Li, Ruoming Jin, and Bin Ren. Speed-ann: Low-latency and high-accuracy nearest neighbor search via intra-query parallelism, 2022.

[36] Yu. A. Malkov and D. A. Yashunin. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs, 2018.

[37] Yao Tian, Xi Zhao, and Xiaofang Zhou. Db-lsh: Locality-sensitive hashing with query-based dynamic bucketing, 2022.

[38] Philip Sun, David Simcha, Dave Dopson, Ruiqi Guo, and Sanjiv Kumar. Soar: Improved indexing for approximate nearest neighbor search, 2024.

[39] Logan Hallee, Rohan Kapur, Arjun Patel, Jason P. Gleghorn, and Bohdan Khomtchouk. Contrastive learning and mixture of experts enables precise vector embeddings, 2024.

[40] Shitao Xiao, Zheng Liu, Weihao Han, Jianjin Zhang, Yingxia Shao, Defu Lian, Chaozhuo Li, Hao Sun, Denvy Deng, Liangjie Zhang, Qi Zhang, and Xing Xie. Progressively optimized bi-granular document representation for scalable embedding based retrieval, 2022.

[41] Jonathan Pan, Swee Liang Wong, and Yidi Yuan. Raglog: Log anomaly detection using retrieval augmented generation, 2023.

[42] Leslie Wong. Shotit: compute-efficient image-to-video search engine for the cloud, 2024.

[43] Jing Zhao, Jingya Wang, Madhav Sigdel, Bopeng Zhang, Phuong Hoang, Mengshu Liu, and Mohammed Korayem. Embedding-based recommender system for job to candidate matching on scale, 2021.

[44] Mirac Suzgun, Stuart M. Shieber, and Dan Jurafsky. string2string: A modern python library for string-to-string algorithms, 2023.

[45] Survey of vector database manage.

[46] Gulshan Yadav, RahulKumar Yadav, Mansi Viramgama, Mayank Viramgama, and Apeksha Mohite. Quantixar: High-performance vector data management system, 2024.

[47] Cheng Chen, Chenzhe Jin, Yunan Zhang, Sasha Podolsky, Chun Wu, Szu-Po Wang, Eric Hanson, Zhou Sun, Robert Walzer, and Jianguo Wang. Singlestore-v: An integrated vector database system in singlestore. *Proceedings of the VLDB Endowment*, 17(12):3772–3785, 2024.

[48] Yutong Gou, Jianyang Gao, Yuexuan Xu, and Cheng Long. Symphonyqg: Towards symphonious integration of quantization and graph for approximate nearest neighbor search, 2024.

[49] Shan An, Guangfu Che, Fangru Zhou, Xianglong Liu, Xin Ma, and Yu Chen. Fast and incremental loop closure detection using proximity graphs, 2019.

[50] Georgii Novikov, Alexander Gneushev, Alexey Kadeishvili, and Ivan Oseledets. Tensor-train point cloud compression and efficient approximate nearest-neighbor search, 2024.

[51] Jiuqi Wei, Botao Peng, Xiaodong Lee, and Themis Palpanas. Det-lsh: A locality-sensitive hashing scheme with dynamic encoding tree for approximate nearest neighbor search, 2024.

[52] Simone Ercoli, Marco Bertini, and Alberto Del Bimbo. Compact hash codes for efficient visual descriptors retrieval in large scale databases, 2016.

[53] Nicholas Rossi, Juexin Lin, Feng Liu, Zhen Yang, Tony Lee, Alessandro Magnani, and Ciya Liao. Relevance filtering for embedding-based retrieval, 2024.

[54] Yuting Qin et al. *Understanding Indexing Efficiency for Approximate Nearest Neighbor Search in High-dimensional Vector Databases*. PhD thesis, Massachusetts Institute of Technology, 2024.

[55] Federico Magliani, Tomaso Fontanini, and Andrea Prati. Efficient nearest neighbors search for large-scale landmark recognition, 2018.

[56] Ninh Pham and Tao Liu. Falconn++: A locality-sensitive filtering approach for approximate nearest neighbor search, 2022.

[57] Wentao Xiao, Yueyang Zhan, Rui Xi, Mengshu Hou, and Jianming Liao. Enhancing hnsw index for real-time updates: Addressing unreachable points and performance degradation, 2024.

[58] Benjamin Coleman, Santiago Segarra, Anshumali Shrivastava, and Alex Smola. Graph reordering for cache-efficient near neighbor search, 2021.

[59] Sitian Chen, Amelie Chi Zhou, Yucheng Shi, Yusen Li, and Xin Yao. Memanns: Enhancing billion-scale anns efficiency with practical pim hardware, 2024.

[60] Yusuke Matsui, Ryota Hinami, and Shin'ichi Satoh. Reconfigurable inverted index, 2018.

[61] Dongkyu Kim, Byoungwook Kim, Donggeon Han, and Matouš Eibich. Autorag: Automated framework for optimization of retrieval augmented generation pipeline, 2024.

[62] Sohrab Ferdowsi, Slava Voloshynovskiy, Dimche Kostadinov, and Taras Holotyak. Sparse ternary codes for similarity search have higher coding gain than dense binary codes, 2017.

[63] Tristan Kenneweg, Philip Kenneweg, and Barbara Hammer. Retrieval augmented generation systems: Automatic dataset creation, evaluation and boolean agent setup, 2024.

[64] Shuo Yang, Jiadong Xie, Yingfan Liu, Jeffrey Xu Yu, Xiyue Gao, Qianru Wang, Yanguo Peng, and Jiangtao Cui. Revisiting the index construction of proximity graph-based approximate nearest neighbor search, 2024.

[65] Kejing Lu, Chuan Xiao, and Yoshiharu Ishikawa. Probabilistic routing for graph-based approximate nearest neighbor search, 2024.

[66] Dantong Zhu and Minjia Zhang. Understanding and generalizing monotonic proximity graphs for approximate nearest neighbor search, 2021.

[67] Ville Hyvönen, Elias Jääsaari, and Teemu Roos. A multilabel classification framework for approximate nearest neighbor search, 2022.

[68] ementacija konverzacijskog agent.

[69] Qing-Yuan Jiang, Ming-Wei Li, and Wu-Jun Li. On the evaluation metric for hashing, 2024.

[70] Pengfei He, Shaowei Wang, Shaiful Chowdhury, and Tse-Hsun Chen. Exploring demonstration retrievers in rag for coding tasks: Yeas and nays!, 2024.

[71] Brad Windsor and Kevin Choi. Thistle: A vector database in rust, 2023.

[72] Xueguang Ma, Tommaso Teofili, and Jimmy Lin. Anserini gets dense retrieval: Integration of lucene's hnsw indexes, 2023.

[73] Nick Alonso, Tomás Figliolia, Anthony Ndirango, and Beren Millidge. Toward conversational agents with context and time sensitive long-term memory, 2024.

[74] Jiawei Chen, Hongyu Lin, Xianpei Han, and Le Sun. Benchmarking large language models in retrieval-augmented generation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 38, pages 17754–17762, 2024.

[75] Philip Sun, Ruiqi Guo, and Sanjiv Kumar. Automating nearest neighbor search configuration with constrained optimization, 2023.

[76] Kaixiang Yang, Hongya Wang, Bo Xu, Wei Wang, Yingyuan Xiao, Ming Du, and Junfeng Zhou. Tao: A learning framework for adaptive nearest neighbor search using static features only, 2021.

[77] Ahmed El-Kishky, Thomas Markovich, Kenny Leung, Frank Portman, Aria Haghighi, and Ying Xiao. knn-embed: Locally smoothed embedding mixtures for multi-interest candidate retrieval, 2023.

[78] Weijie Zhao, Shulong Tan, and Ping Li. Guitar: Gradient pruning toward fast neural ranking, 2023.

[79] Tan Pan, Furong Xu, Xudong Yang, Sifeng He, Chen Jiang, Qingpei Guo, Feng Qian Xiaobo Zhang, Yuan Cheng, Lei Yang, and Wei Chu. Boundary-aware backward-compatible representation via adversarial learning in image retrieval, 2023.

[80] Yitu Wang, Shiyu Li, Qilin Zheng, Linghao Song, Zongwang Li, Andrew Chang, Hai "Helen" Li, and Yiran Chen. Ndsearch: Accelerating graph-traversal-based approximate nearest neighbor search through near data processing, 2024.

[81] Shangyu Wu, Ying Xiong, Yufei Cui, Haolun Wu, Can Chen, Ye Yuan, Lianming Huang, Xue Liu, Tei-Wei Kuo, Nan Guan, et al. Retrieval-augmented generation for natural language processing: A survey. *arXiv preprint arXiv:2407.13193*, 2024.

[82] Simon Akesson and Frances A. Santos. Clustered retrieved augmented generation (crag), 2024.

[83] Chao Jin, Zili Zhang, Xuanlin Jiang, Fangyue Liu, Xin Liu, Xuanzhe Liu, and Xin Jin. Ragcache: Efficient knowledge caching for retrieval-augmented generation, 2024.

[84] Rajat Khanda. Agentic ai-driven technical troubleshooting for enterprise systems: A novel weighted retrieval-augmented generation paradigm, 2024.

[85] Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*, 1(2):4, 2022.

[86] Michihiro Yasunaga, Armen Aghajanyan, Weijia Shi, Rich James, Jure Leskovec, Percy Liang, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. Retrieval-augmented multimodal language modeling. *arXiv preprint arXiv:2211.12561*, 2022.

[87] Xin Cheng, Di Luo, Xiuying Chen, Lemao Liu, Dongyan Zhao, and Rui Yan. Lift yourself up: Retrieval-augmented text generation with self-memory. *Advances in Neural Information Processing Systems*, 36:43780–43799, 2023.

[88] Daniel Gillick, Alessandro Presta, and Gaurav Singh Tomar. End-to-end retrieval in continuous space, 2018.

[89] Junhan Yang, Zheng Liu, Bowen Jin, Jianxun Lian, Defu Lian, Akshay Soni, Eun Yong Kang, Yajun Wang, Guangzhong Sun, and Xing Xie. Hybrid encoder: Towards efficient and precise native adsrecommendation via hybrid transformer encoding networks, 2021.

[90] Narges Norouzi, Reza Azmi, Sara Saberi Tehrani Moghadam, and Maral Zarvani. Snapmode: An intelligent and distributed large-scale fashion image retrieval platform based on big data and deep generative adversarial network technologies, 2022.

[91] Weihao Gao, Xiangjun Fan, Chong Wang, Jiankai Sun, Kai Jia, Wenzhi Xiao, Ruofan Ding, Xingyan Bin, Hui Yang, and Xiaobing Liu. Deep retrieval: Learning a retrievable structure for large-scale recommendations, 2021.

[92] Rishabh Goel. Using text embedding models as text classifiers with medical data, 2024.

[93] MD Shaikh Rahman, Feiroz Humayara, Syed Maudud E Rabbi, and Muhammad Mahbubur Rashid. Efficient medical image retrieval using densenet and faiss for birads classification, 2024.

[94] Hangfei Lin, Li Miao, and Amir Ziai. Rafic: Retrieval-augmented few-shot image classification, 2023.

[95] Yuxuan Wu, Xiao Yi, Yang Tan, Huiqun Yu, Guisheng Fan, and Gaowei Zheng. A plms based protein retrieval framework, 2025.

[96] Andrew Hamara and Pablo Rivas. From latent to engine manifolds: Analyzing imagebind's multimodal embedding space, 2024.

[97] Rihan Chen, Bin Liu, Han Zhu, Yaoxuan Wang, Qi Li, Buting Ma, Qingbo Hua, Jun Jiang, Yunlong Xu, Hongbo Deng, and Bo Zheng. Approximate nearest neighbor search under neural similarity metric for large-scale recommendation, 2022.

[98] Zhiwei Zuo, Zhuo Tang, Bin Wang, Kenli Li, and Anwitaman Datta. ecil-mu: Embedding based class incremental learning and machine unlearning, 2024.

[99] Mahesh Kandhare and Thibault Gisselbrecht. An empirical comparison of video frame sampling methods for multi-modal rag retrieval, 2024.

[100] Tingwei Zhang, Fnu Suya, Rishi Jha, Collin Zhang, and Vitaly Shmatikov. Adversarial hubness in multi-modal retrieval, 2024.

[101] Sonal Prabhune and Donald J Berndt. Deploying large language models with retrieval augmented generation. *arXiv preprint arXiv:2411.11895*, 2024.

[102] Xinlong Sun, Yangyang Qin, Xuyuan Xu, Guoping Gong, Yang Fang, and Yexin Wang. 3rd place: A global and local dual retrieval solution to facebook ai image similarity challenge, 2021.

[103] Ji-Hoon Kim, Yeo-Reum Park, Jaeyoung Do, Soo-Young Ji, and Joo-Young Kim. Accelerating large-scale graph-based nearest neighbor search on a computational storage platform, 2022.

[104] Fabian Groh, Lukas Ruppert, Patrick Wieschollek, and Hendrik P. A. Lensch. Ggnn: Graph-based gpu nearest neighbor search, 2022.

**Disclaimer:**

SurveyX is an AI-powered system designed to automate the generation of surveys. While it aims to produce high-quality, coherent, and comprehensive surveys with accurate citations, the final output is derived from the AI's synthesis of pre-processed materials, which may contain limitations or inaccuracies. As such, the generated content should not be used for academic publication or formal submissions and must be independently reviewed and verified. The developers of SurveyX do not assume responsibility for any errors or consequences arising from the use of the generated surveys.