



ULTIMATE

Milvus Vector Database for AI Apps

Master Vector Search,
Distributed Data Management,
and GPU-Accelerated AI
Systems with Milvus

Prashanth Raghu

Copyright © 2026 Orange Education Pvt Ltd, AVA®

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author nor **Orange Education Pvt Ltd** or its dealers and distributors, will be held liable for any damages caused or alleged to have been caused directly or indirectly by this book.

Orange Education Pvt Ltd has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capital. However, **Orange Education Pvt Ltd** cannot guarantee the accuracy of this information. The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

First Published: April 2026

Published by: Orange Education Pvt Ltd, AVA®

Address: 9, Daryaganj, Delhi, 110002, India

275 New North Road Islington Suite 1314 London,
N1 7AA, United Kingdom

ISBN (PBK): 978-93-49887-18-3

ISBN (E-BOOK): 978-93-49887-58-9

Scan the QR code to explore our entire catalogue



www.orangeava.com

Table of Contents

1. Introduction to Vector Databases

Introduction

Structure

Introduction to Vector Databases

Differences between Vector and Relational and NoSQL DBMS

Exploring Use Cases

Search and Information Retrieval

Semantic Search: What is Semantic Search

Image and Video Search

Recommendation Systems

Anomaly Detection

Natural Language Processing (NLP)

Drug Discovery and Drug Prediction

Molecular Similarity Search

Financial Services

Content Moderation

Knowledge Graphs and Ontologies

Understanding the Architecture

Indexing Techniques

CPU-Based Indexing Techniques in Milvus

FLAT (Brute Force)

IVF (Inverted File Index)

IVF_PQ (Inverted File + Product Quantization)

HNSW (Hierarchical Navigable Small World Graph)

SCANN (Selective Search, Limited Support)

GPU-Based Indexing Techniques in Milvus

GPU FLAT

GPU IVF FLAT

GPU IVF PQ

CPU versus GPU Indexing: A Strategic View

Scalability and Performance

Indexing Techniques

Caching Strategies

Efficient Query Processing
Real-Time Data Handling
Fault Tolerance and High Availability

Conclusion

2. Fundamentals of Vectors and AI

Introduction

Structure

Background: Representing Language in AI Systems

Early Methods of Text Representation

Bag of Words (BoW)

Term Frequency-Inverse Document Frequency (TF-IDF)

Limitations of Traditional Methods

Embeddings: A Paradigm Shift in Language Representation

Word Embeddings

Sentence Embeddings

Why Embeddings are Essential in Modern AI

Vector Space and Similarity

Cosine Similarity

Euclidean Distance

Manhattan Distance (L1 Distance)

Jaccard Similarity

Hamming Distance

Summary of Similarity Measures

Applications of Vector Representation and Similarity Measures

Background on N-gram Models

Definition

How It Works

Next Word Prediction Using N-gram Model

Limitations of N-gram Models

Contextual Limitations

Data Sparsity

Limited Generalization

High Dimensionality

Lack of Semantic Understanding

Inflexibility

Background on Recurrent Neural Networks (RNNs)

[*Fundamental Concepts*](#)
[*Mathematical Formulation*](#)
[*Training RNNs*](#)
[*Advantages of RNNs*](#)
[Next Word Prediction Using Recurrent Neural Networks \(RNNs\)](#)
[*RNN Architecture*](#)
[*Advantages of RNNs*](#)
[*Limitations of RNNs*](#)
[Background on Transformers](#)
[*Key Concepts*](#)
[*Transformer Architecture*](#)
[*Training Transformers*](#)
[*Advantages of Transformers*](#)
[Next Word Prediction in Transformers](#)
[*Mechanism of Next Word Prediction*](#)
[*Benefits of Using Transformers for Next Word Prediction*](#)
[*Implementation Steps*](#)
[Conclusion](#)

3. Components of Milvus

[Introduction](#)
[Structure](#)
[Components of Milvus](#)
[*Micro-Service Components*](#)
[*Third-Party Dependencies*](#)
[Components of the Messaging Layer](#)
[*Key Roles of Kafka in Milvus*](#)
[*Key Roles of Pulsar in Milvus*](#)
[Components of the Coordinators](#)
[*Root Coordinator*](#)
[*Query Coordinator*](#)
[*Data Coordinator*](#)
[*Index Coordinator*](#)
[*Proxy \(Coordination Role\)*](#)
[Components of the Worker Nodes](#)
[*Types of Worker Nodes in Milvus*](#)
[*Query Node*](#)

[*Data Node*](#)

[*Index Node*](#)

[*Proxy \(Worker Role\)*](#)

[*Proxy Nodes*](#)

[*Query Coordinator Nodes*](#)

[Components of the Storage Layer](#)

[*Workflow Incorporating Log Snapshots, Index Files, and Delta Files*](#)

[*Summary of Relationships*](#)

[Components of the Storage Layer](#)

[*Workflow across Layers*](#)

[*Benefits of Milvus Storage Layering*](#)

[Conclusion](#)

4. Data, Storage and Cluster Management

[Introduction](#)

[Structure](#)

[Compute-Storage Disaggregation and Its Implementation in Kubernetes in Milvus](#)

[Understanding Compute-Storage Disaggregation](#)

[*Implementation in Kubernetes*](#)

[*Data Coordinator \(data-coord\) Deployment Configuration*](#)

[*Root Coordinator \(root-coord\) Deployment Configuration*](#)

[*Index Coordinator \(index-coord\) Deployment Configuration*](#)

[*Data Worker \(data-worker\) Deployment Configuration*](#)

[*Index Worker \(index-worker\) Deployment Configuration*](#)

[*Query Worker \(query-worker\) Deployment Configuration*](#)

[*Deploying the Configuration Files*](#)

[Cluster Management in Milvus: The Role of Data Coordinator](#)

[*The Cluster Interface*](#)

[*ClusterImpl: Implementing the Cluster Interface*](#)

[*Concurrency and Resource Management*](#)

[*Milvus Compaction Management*](#)

[*Milvus Garbage Collection for Data Coordinator*](#)

[*Milvus Index Service Management*](#)

[*Milvus Segment Management*](#)

[Implementation of the Root Coordinator in Milvus](#)

[*Understanding the Create Partition Task in Root Coordinator*](#)

[Overview of the `create_partition_task.go` File](#)

[Key Components of the Code](#)

[Practical Implications](#)

[Overview of Creating a Database Task in Root Coordinator](#)

[Overview of the `create_db_task.go` File](#)

[Key Components of the Code](#)

[Practical Implications](#)

[Overview of the Index Coordinator](#)

[Understanding Index Coordinator](#)

[Key Functions of the Index Coordinator](#)

[Benefits of the Index Coordinator in Milvus](#)

[Overview of the Indexing Operations](#)

[Important Functions in `indexnode.go`](#)

[Overview of the Index Creation Manager](#)

[Task Management Flow](#)

[Conclusion](#)

[5. Indexing Schemes](#)

[Introduction](#)

[Structure](#)

[Overview of the Indexing Types](#)

[Floating-Type Indexes](#)

[Sparse Indexes](#)

[GPU-Based Indexes](#)

[Binary Indexes](#)

[Summary of Binary Indexing Types in Milvus](#)

[Implementation of the IVF Indexing Scheme](#)

[IVF Implementation in FAISS - Key Code Snippets Summary](#)

[Additional Utility Functions](#)

[Key Code Summary](#)

[Implementation of the Index IVF Addition and Search Functions](#)

[Key Takeaways](#)

[`IndexIVF::search\(\)`](#)

[Implementation of the Index PQ Addition and Search Functions](#)

[`pq_estimators_from_tables` Function Summary](#)

[Implementation of the Index Flat Addition and Search Functions](#)

[Flat Index: The Foundation of FAISS \(Facebook AI Similarity Search\)](#)

[How It Works](#)

[Types of Flat Indexes](#)

[Search Operation in the Flat Index Scheme](#)

[Explanation](#)

[Other Similarity Metrics](#)

[Other Distance Metrics](#)

[Heaps Recap](#)

[Helper Functions Called](#)

[Index Addition Operation in the Flat Index Scheme](#)

[Detailed Explanation](#)

[Updating the Metadata](#)

[What Happens Internally](#)

[Implementation of the HNSW Flat Addition and Search Functions](#)

[Core Idea](#)

[Construction](#)

[Search](#)

[Pros](#)

[Cons](#)

[Structure Overview of `IndexHNSW`](#)

[Components](#)

[Constructors](#)

[Key Methods](#)

[Advanced Graph Functions](#)

[Addition of Vectors to the HNSW Index](#)

[Search Nearest Vectors Using the HNSW Index](#)

[Explanation of the Function](#)

[Conclusion](#)

[6. Indexing Schemes Binary, Sparse, and GPU](#)

[Introduction](#)

[Structure](#)

[Implementation of the Hierarchical Navigable Small World \(HNSW\)](#)

[Indexing Scheme](#)

[Overview](#)

[Key Components](#)

[Purpose](#)

[Applications](#)

[Implementation of the Hierarchical Navigable Small World
\(HNSW\) Search](#)

[Input Parameters](#)

[Key Operations](#)

[Detailed Steps in the Function](#)

[Key Structures Involved](#)

[Optimization Techniques](#)

[Implementation of the Binary IVF Index](#)

[Key Components of the `IndexBinaryIVF` Class](#)

[Main Functions](#)

[Searching within the Binary IVF Index](#)

[Function Parameters](#)

[Code Overview](#)

[Implementation of the GPU IVF Index](#)

[Class Overview](#)

[Constructor\(s\)](#)

[Destructor](#)

[Private Initialization Method](#)

[Public Methods](#)

[Protected Members](#)

[Searching the GPU IVF Index](#)

[Function Overview](#)

[Steps and Explanation](#)

[Conclusion](#)

[7. Distributed Data Management](#)

[Introduction](#)

[Structure](#)

[Implementation of the Data Coordinator](#)

[Operations of the Root Data Coordinator: The `DataNodeManager`
Interface](#)

[Session Management](#)

[Data Operations](#)

[Compaction and Health Management](#)

[Channel Operations](#)

Query and Slot Management

System Shutdown

Data Management by the Cluster

Function Overview

Detailed Steps

Implementation of the Data Nodes

Core Structure and Context Management

Node Role and State Management

Managers for Data Operations

Data Structures for Caching and Execution

External System Interfaces

Initialization and Lifecycle Control

Additional Components and Metrics

MultiSegmentWriter Struct

IO and Resource Allocation

Writer and Segment Tracking

Segment Size Management

Schema and Metadata

Compaction Data and BM25 Support

Summary and Final Remarks

Conclusion

8. Distributed Index and Query Management

Introduction

Structure

Implementation of the Index Co-ordinator

Implementation of the Index Nodes

Implementation of the Query Coordinator

Implementation of the Query Management

Key Concepts and Flow

Summary

Conclusion

9. Design of the Proxy Server

Introduction

Structure

The Role of Proxy in Milvus

Architectural Design of the Proxy.

Proxy Components

Communication between Components

Responsibilities of the Proxy.

Request Handling and Routing

Session Management

Load Balancing

Query Aggregation

Fault Tolerance and High Availability.

Proxy and Scalability.

Proxy and Query Execution

Proxy for Multi-Tenancy.

Conclusion

10. GPU Based Indexes and Optimizations

Introduction

Structure

Implementation of the GPU Based Index GPU IVF Flat

Key Concepts

Inverted File (IVF) Indexing

Flat Indexing within Clusters

GPU Acceleration

How GPU Index IVF Flat Works

Step 1: Clustering and Building the Inverted Index

Step 2: Query Processing and Cluster Selection

Step 3: Flat Search within Selected Clusters on GPU

Step 4: Results Aggregation and Return

Benefits of GPU IVF Flat

Scalability

Speed

Accuracy

Limitations and Considerations

Applications

Summary

Structure of the GPU Index IVF Flat

GPU Based Quantization

Function Overview

[Steps](#)

[Key Logic](#)

[Implementation of the GPU Based Index GPU_IVF_PQ](#)

[GPU-Based Index GPU_PQ: An In-depth Overview](#)

[Core Concepts of Product Quantization](#)

[Compression and Memory Efficiency](#)

[How GPU Index GPU_IVF_PQ Works](#)

[Training Phase on GPU](#)

[Query and Search Phase](#)

[Advantages of GPU IVF PQ](#)

[Memory Efficiency and Compression](#)

[High Speed and Throughput](#)

[Accuracy versus Speed Trade-off](#)

[Limitations and Considerations](#)

[Use Cases](#)

[Structure of the Index GPU_IVF_PQ](#)

[Implementation of the GPU Based Index GPU Scalar Quantization](#)

[Implementation of Scalar Quantization](#)

[Detailed Explanation](#)

[Key Functions](#)

[Conclusion](#)

[11. Auxiliary Components](#)

[Introduction](#)

[Structure](#)

[The Role of Etcd in Milvus](#)

[Cluster Coordination and Node Discovery](#)

[Configuration Management](#)

[Metadata Storage and Synchronization](#)

[Fault Tolerance and Recovery](#)

[Global State Management](#)

[How Etcd Integrates with Other Components in Milvus](#)

[The Role of Kafka in Milvus](#)

[Efficient Data Ingestion and Asynchronous Processing](#)

[Event-Driven Architecture](#)

[Real-Time Query Processing and Parallelization](#)

[Ensuring Consistency and Reliability](#)

[*Handling Large-Scale Indexing and Distributed Computation*](#)
[*The Role of Pulsar in Milvus*](#)
[*Real-Time Stream Processing and Data Ingestion*](#)
[*Distributed Messaging and Scalability*](#)
[*Data Persistence and Durability*](#)
[*Fault Tolerance and Recovery*](#)
[*Event-Driven Architecture and Task Management*](#)
[*Seamless Integration with Other Components*](#)
[*Conclusion*](#)

12. Network Management

[*Introduction*](#)

[*Structure*](#)

[*Client Management*](#)

[*Milvus Client SDK Overview*](#)

[*Connection Protocols*](#)

[*gRPC in Milvus: A Comprehensive Overview*](#)

[*Introduction*](#)

[*Understanding gRPC*](#)

[*Role of gRPC in Milvus*](#)

[*Primary Communication Protocol*](#)

[*Protocol Buffers for Data Serialization*](#)

[*Benefits of Using gRPC in Milvus*](#)

[*Performance*](#)

[*Scalability and Concurrency*](#)

[*Cross-language Support*](#)

[*How gRPC Works in Milvus*](#)

[*Connection Establishment*](#)

[*RPC Lifecycle*](#)

[*Streaming*](#)

[*Best Practices for Using gRPC with Milvus*](#)

[*Connection and Retry Configuration*](#)

[*Optimize Data Transfer*](#)

[*Security*](#)

[*Monitoring and Debugging*](#)

[*Limitations and Considerations*](#)

[*RESTful API in Milvus*](#)

[*Advantages of RESTful API in Milvus*](#)
[*How RESTful API Works in Milvus*](#)
[*Establishing a Connection to Milvus*](#)
[*Connection Basics*](#)
[*Connecting via SDK*](#)
[*Load Balancing and High Availability*](#)
[*Connection Pooling*](#)
[*Handling Client Queries in Milvus*](#)
[*Vector Search*](#)
[*Data Insertion*](#)
[*Updates and Deletions*](#)
[*Best Practices for Client Connections*](#)
[*Milvus Python SDK \(pymilvus\).*](#)
[*Installation*](#)
[*Establishing a Connection*](#)
[*Managing Collections*](#)
[*Inserting Data*](#)
[*Indexing and Searching*](#)
[*Querying Metadata*](#)
[*Best Practices*](#)
[*Use Cases*](#)
[*Milvus Java SDK*](#)
[*Installation \(Maven\)*](#)
[*Connecting to Milvus*](#)
[*Collection Management*](#)
[*Creating Index and Performing Search*](#)
[*Best Practices*](#)
[*Ideal Use Cases*](#)
[*Conclusion*](#)

[**Index**](#)

CHAPTER 1

Introduction to Vector Databases

Introduction

As data escaped the neat rows and columns of classical databases and began to live as text, images, audio, code, and graphs of relationships, a new geometry became necessary. Vector databases emerged to meet this shift. They store meaning rather than labels by representing data as high-dimensional vectors produced by embeddings. In this space, similarity becomes a distance problem, not a keyword match. This makes vector databases indispensable for modern AI systems like semantic search, recommendation engines, RAG pipelines, fraud detection, and multimodal assistants, where relevance depends on context and intent rather than exact matches. Their rise mirrors the transition from rule-based software to representation-based learning, where understanding emerges from proximity in latent space.

Mastery of vector databases is mastery of AI at scale because they sit at the junction of models, data, and systems. Knowing how vectors are generated, indexed, sharded, compressed, and queried allows engineers to reason about latency, recall, drift, and cost, not just accuracy. At scale, naive similarity search collapses under its own weight, and only a deep understanding of approximate nearest neighbor algorithms, embedding lifecycle management, and feedback loops keeps systems reliable and adaptive. In short, models think, but vector databases remember and retrieve meaning efficiently. Understanding these turns AI from a clever demo into an industrial, evolving intelligence. In this chapter, we aim to provide a comprehensive introduction to vector databases, an emerging technology in the field of data management and machine learning. Our objective is to elucidate the fundamental concepts, architectures, and applications of vector databases, making this chapter a valuable resource for both novices and experienced professionals.

Structure

In this chapter, we will cover the following topics:

- Introduction to Vector Databases
- Exploring Use Cases
- Understanding the Architecture
- Indexing Techniques
- Scalability and Performance

Introduction to Vector Databases

A vector database is designed to store and manage high-dimensional data, which traditional Database Management Systems (DBMS) struggle to handle effectively. This innovative type of database utilizes mathematical representations known as vectors, which can consist of numerous dimensions, typically ranging from tens to thousands, depending on the complexity of data. These vectors are generated from various forms of data, including text, images, audio, and video, allowing for sophisticated analysis and retrieval capabilities that go beyond simple keyword matching. Vector databases facilitate fast and precise similarity searches, making them essential for applications in machine learning, natural language processing, and computer vision. They also support the integration of Large Language Models (LLMs), enhancing the capabilities for tasks such as semantic search, real-time knowledge retrieval, and reducing inaccuracies often encountered in LLM outputs, commonly referred to as "hallucination".

Differences between Vector and Relational and NoSQL DBMS

Let us understand the differences between Vector and Relational Databases with an example. Consider a scenario where a user wants to find images similar to a specific photo. In a vector database, the user can input the photo, and the database will retrieve images based on their visual features represented as high-dimensional vectors. In contrast, a traditional database would require the user to search using keywords or exact matches, which may not yield the desired results based on similarity.

The following table demonstrates the key differences:

Vector Databases	Relational Databases
Store data as high-dimensional vectors that represent features or attributes of complex data types (like images, text).	Use structured data formats that fit rigid schemas, typically handling rows and columns.
Vector databases use a special methodology called as Semantic Search to perform the Search Capability. By using, Semantic Search, Vector Databases find the items based on the similarity and not just the exact matches. This is particularly useful for applications like recommendation systems or natural language processing.	Traditional Databases typically perform searches based on exact terms or defined criteria, limiting their ability to understand context.
Vector Databases can handle unstructured and complex data types such as images, audio and video, which are transformed into vectors for storage and retrieval.	Traditional Databases mainly support structured data, making them less effective for handling diverse data types.
Vector Databases are commonly used in applications requiring real-time knowledge search, semantic search, and machine learning integration.	Traditional Databases are often struggle with high-dimensionality due to dimensionality catastrophe.

Table 1.1: RDBMS versus Vector Databases

Exploring Use Cases

Vector databases have become increasingly important due to the rise of machine learning, AI, and other advanced computational tasks that require the processing and retrieval of high-dimensional data. These databases are optimized for storing and querying vector representations of data, typically generated by models such as embeddings from deep learning models. The detailed use cases of the Vector Databases are as following;

Search and Information Retrieval

Search and information retrieval form one of the core essence of the needs of the data management systems where users data is sliced and managed according to the needs of the individual applications.

In this section, we delve into the fundamentals of such systems and how vectors play such an important role in the storage and management of such data retrieval systems.

Semantic Search: What is Semantic Search

Semantic search retrieves information based on meaning and intent rather than exact keywords, using embeddings to understand context. It enables systems to find what you mean, not just what you type, even when the words differ.

For example, if a user searches for “how to cool a room without AC,” a keyword-based system might prioritize documents containing “cool,” “room,” and “AC,” even if they are irrelevant. A semantic search system, however, recognizes the goal of the user—finding alternative cooling methods—and could surface results about fans, cross-ventilation, or thermal curtains, even if those terms are not explicitly mentioned in the query.

Instead of relying solely on keyword matching, vector databases enable semantic search by allowing text embeddings to be stored. This enables search engines to better understand the context and meaning of queries and documents, resulting in more relevant search results. Instead of relying solely on keyword matching, vector databases enable semantic search by allowing text embeddings to be stored. This enables search engines to better understand the context and meaning of queries and documents, resulting in more relevant search results.

On an e-commerce platform, various product descriptions are embedded into vectors representing their semantic meaning. When a user searches for "comfortable running shoes," the query is also converted into a vector. The vector database then retrieves products with similar vector representations, ensuring that items such as "Lightweight sneakers for jogging" are included in the search results, even if they do not explicitly mention "running shoes."

This approach ensures that the search is more context-aware, improving the relevance of the results and enhancing the user experience.

Image and Video Search

By having visual embeddings from images and videos stored in vector databases, searches based on visual similarity are enabled. Users are allowed to search for images or videos that are visually similar to a given input, even if the metadata does not explicitly match.

On a media platform, embeddings of images and videos are created based on their visual content and stored as vectors in a vector database. When a

user provides an image of a "red sports car" as an input, the query is converted into a vector that captures the visual characteristics of the car. The vector database then retrieves images or videos with similar visual embeddings, ensuring that other visually similar "red sports cars" are included in the search results, even if they are labeled differently.

This approach allows for a more intuitive and accurate search experience, as users can visualize the related content without relying solely on metadata.

Recommendation Systems

A subset of search and retrieval systems, recommendation engines slice the data by understanding repeated patterns according to multiple cross sections of filters enabling the system to predict and suggest patterns to users or to use business management teams.

- **Content-Based Recommendations:** Embeddings of user preferences and item attributes can be stored in vector databases. Recommendations are generated by identifying the nearest neighbors (similar vectors) in this high-dimensional space, resulting in personalized content delivery.

On a streaming platform, user preferences and content attributes are embedded into vectors and stored in a vector database. When a user interacts with a particular movie or series, the platform retrieves similar content by finding vectors that are close to the preferred vector of the user. For instance, if a user frequently watches "action thrillers," the platform recommends other action-packed titles with similar thematic elements, even if the titles are not explicitly labeled as "action thrillers."

This approach ensures that recommendations are more aligned with the tastes of the user, enhancing the overall user experience through personalized content suggestions.

- **Collaborative Filtering:** By embedding users and items into the same vector space, vector databases enable the efficient computation of similarity scores. Recommendations are made based on the preferences of other users with similar tastes.

In an online shopping platform, user behavior and product attributes are embedded into a shared vector space. When a user interacts with certain products, the system identifies other users with similar behavior patterns by finding vectors close to the vector of the user in the space. The platform then recommends products that these similar users have liked or purchased.

This approach allows for personalized recommendations based on collective user behavior, ensuring that users discover products they are likely to enjoy based on the preferences of others with similar tastes.

Anomaly Detection

Systems which understand undesirable usage patterns leading to abuse of either the systems, business or users and try to prevent them are called Anomaly Detection and Management Systems.

- **Fraud Detection:** In financial transactions, network security, and other domains, embeddings of typical behavior can be stored in vector databases. Anomalous activities are detected by measuring the distance of new data points from these embeddings, allowing potential fraud or irregularities to be identified.

In a banking system, typical transaction behaviors are embedded into vectors and stored in a vector database. When a new transaction occurs, its vector is compared against the stored vectors representing normal behavior. If the vector of the new transaction is significantly distant from the typical behavior vectors, it is flagged as potentially fraudulent.

This approach enables the early detection of suspicious activities, helping to prevent fraud by identifying irregular patterns that deviate from the norm.

- **Industrial Monitoring:** In manufacturing or other industries, sensor data can be embedded and stored in a vector database. Anomalies in machine operations or production processes are detected by identifying deviations from the norm in the vector space.

In a manufacturing plant, sensor readings from machines are continuously embedded into vectors and stored in a vector database.

These vectors represent normal operating conditions. When new sensor data is received, its vector is compared to the stored vectors. If a significant deviation from the normal pattern is detected, the system flags it as an anomaly, indicating potential issues in the operation of the machine.

This approach allows for proactive maintenance by identifying irregularities early, reducing downtime and improving the overall efficiency of industrial operations.

Natural Language Processing (NLP)

Bringing computation for all users has the implicit need for computers to understand language data. However, simple the problem statement might look has captured the interest of the global community.

- **Text Classification:** By having text embeddings stored in vector databases, the classification of documents, emails, or other text data is facilitated. This is particularly useful in applications such as spam detection, sentiment analysis, and topic categorization.

In an email system, emails are embedded into vectors based on their content and stored in a vector database. When a new email arrives, its vector is compared against pre-labeled vectors representing categories such as "spam," "important," or "promotional." The system then classifies the email accordingly. For instance, if the vector closely aligns with the "spam" vectors, the email is flagged as spam. This approach improves the accuracy of text classification tasks, ensuring that documents or messages are correctly categorized based on their content.

- **Large Language Models (LLMs):** By having text and other data embeddings stored in vector databases, the retrieval of relevant information during text generation or query handling is facilitated. This enhances the ability of LLMs to produce accurate, contextually aware, and coherent outputs.

In a chatbot system, user queries are embedded into vectors and stored in a vector database. When a new query is received, its vector is compared against stored vectors representing similar queries or relevant information. The LLM then generates a response based on the

most relevant retrieved content, improving the accuracy and relevance of its answers.

This approach allows LLMs to handle complex queries more effectively by leveraging stored embeddings that capture the semantic meaning of the data.

- **Retrieval-Augmented Generation (RAG):** By storing embeddings of external knowledge in vector databases, the retrieval of relevant information to augment text generation is enabled. This allows RAG systems to generate more informed and contextually accurate content.

In a research assistant application, external documents and datasets are embedded into vectors and stored in a vector database. When the system is tasked with generating a report on a specific topic, it retrieves the most relevant information by comparing vectors in the database. The retrieved data is then used to inform the generative process, ensuring the output is up-to-date and contextually rich.

This approach improves the quality and relevance of content generated by RAG systems, particularly when specific or recent information is required.

[Drug Discovery and Drug Prediction](#)

Drug discovery is the scientific process of identifying new compounds that can treat or prevent diseases. Researchers study biological targets such as proteins or genes and screen thousands of molecules to find those that interact effectively with them. Promising compounds are then optimized and tested through laboratory studies and clinical trials to ensure safety and effectiveness before becoming medicines.

[Molecular Similarity Search](#)

By representing molecules as high-dimensional vectors based on their properties and storing these embeddings in vector databases, the search for similar molecules is facilitated. This capability is particularly valuable in drug discovery, where finding molecules with similar properties can accelerate the development of new treatments. In a drug discovery project, the chemical structures of various compounds are embedded into vectors and stored in a vector database. When researchers identify a molecule with

desirable properties, they can query the database to find other molecules with similar vector representations. This enables them to quickly identify potential drug candidates that share similar characteristics with the molecule of interest.

This approach significantly speeds up the drug discovery process by allowing researchers to efficiently explore chemical space and identify promising compounds.

Financial Services

AI in financial services enables institutions to analyze vast amounts of financial data to detect patterns, risks, and opportunities in real time. It powers applications such as fraud detection, credit scoring, algorithmic trading, and personalized financial recommendations.

- **Credit Scoring:** By storing and processing embeddings of financial histories, behaviors, and transactions in vector databases, more accurate and real-time credit scoring is enabled. This enhances the ability to evaluate creditworthiness based on a comprehensive analysis of the financial data of an individual.

In a credit scoring system, embeddings representing various aspects of the financial history of a person, such as transaction patterns, credit usage, and repayment behavior, is stored in a vector database. When a new credit application is evaluated, the system retrieves and compares the financial embeddings of the applicant with those of similar profiles stored in the database. This allows for a nuanced assessment of the credit risk of the applicant, leading to more accurate and timely credit scores.

This approach improves the precision of credit scoring by leveraging detailed embeddings, allowing for better-informed lending decisions and enhanced risk management.

- **Portfolio Management:** By representing financial portfolios as vectors based on asset characteristics and storing these embeddings in vector databases, quick comparisons and optimizations are facilitated. This approach helps in improving returns and reducing risks associated with portfolio management.

In a portfolio management system, different portfolios are embedded into vectors that capture asset characteristics such as risk levels, returns, and market correlations. When evaluating or optimizing a portfolio, the system queries the vector database to find similar portfolios or benchmark against optimal vectors. This allows for efficient analysis and adjustment of asset allocations to enhance portfolio performance and mitigate risks.

This approach streamlines the process of portfolio optimization, enabling more informed decision-making and better alignment with investment goals.

Content Moderation

- **Detection of Harmful Content:** By storing embeddings of flagged content, such as hate speech or adult content, in vector databases, real-time detection and moderation of harmful content are enabled. This allows for the effective identification and management of inappropriate content by comparing embeddings of new content against stored flagged embeddings.

In a content moderation system, various types of harmful content are embedded into vectors and stored in a vector database. When new user-generated content is submitted, its vector is compared to the stored vectors representing harmful content. If the vector of the new content is found to be similar to flagged embeddings, it is flagged for review or automatically moderated to prevent its dissemination.

This approach enhances content moderation by enabling prompt and accurate detection of harmful material, ensuring a safer and more compliant online environment.

- **Automatic Tagging and Categorization:** By embedding content metadata into vectors and storing these embeddings in vector databases, the process of automatically tagging and categorizing user-generated content on social platforms is facilitated. This approach streamlines content organization and enhances discoverability.

On a social media platform, content such as posts, images, and videos are embedded into vectors based on their metadata, like themes, topics, and user interactions. When new content is uploaded, its vector

is compared to stored vectors representing various tags and categories. The system then automatically assigns relevant tags or categories to the new content based on the closest matches, improving the efficiency of content management and user experience.

This approach ensures that user-generated content is systematically organized, making it easier for users to discover and engage with relevant material.

Knowledge Graphs and Ontologies

Knowledge graphs and ontologies use AI to organize information into interconnected entities, relationships, and concepts, creating structured representations of knowledge. AI techniques help automatically extract relationships from large datasets and continuously enrich these graphs.

- **Entity Linking:** By storing entity embeddings in vector databases, linking entities in knowledge graphs is facilitated through comparison of these embeddings. This aids in integrating data from various sources, ensuring consistency and enriching the knowledge graph.

In a knowledge graph, entities such as people, organizations, or concepts are represented by embeddings and stored in a vector database. When integrating data from different sources, the system compares the embeddings of new entities with those in the database to identify and link equivalent entities. For instance, if two data sources refer to the same person with slightly different names, their embeddings are compared to determine they represent the same entity, allowing for unified and comprehensive data integration.

This approach enhances the accuracy and completeness of knowledge graphs by effectively linking related entities across diverse datasets.

- **Relationship Discovery:** By storing and analyzing vector representations of entities and their relationships, vector databases assist in discovering hidden relationships and insights within knowledge graphs. This enables the identification of previously unknown connections and patterns.

In a knowledge graph, entities and their relationships are represented as vectors and stored in a vector database. By analyzing these vectors,

the system can uncover hidden relationships between entities that were not explicitly documented. For instance, if the vectors of two entities are found to be closely related in the vector space, the system might reveal an implicit connection between them, such as a shared affiliation or influence. This discovery can lead to new insights and enhance the depth of the knowledge graph.

This approach improves the ability to explore and understand complex networks of relationships, uncovering valuable insights and enhancing data analysis.

Understanding the Architecture

At its core, Milvus combines vector indexes with metadata filtering and coordination services, allowing it to scale horizontally while maintaining consistency and performance. This architecture makes it well-suited for production-grade semantic search, recommendations, and retrieval-augmented generation systems.

- **Data Representation:** Data is converted into high-dimensional vectors using embedding functions, such as machine learning models or feature extraction algorithms. This transformation allows the database to handle diverse data types, including text, images, audio, and video.
- **Storage Techniques:**
 - **Sharding:** Data is distributed across multiple machines or clusters (shards) to enhance scalability and performance. Sharding can be done through hash-based or range-based methods to balance load and avoid hotspots.
 - **Partitioning:** The database is divided into smaller, manageable segments based on specific criteria (like geographic location or category) to improve usability and query performance.
- **Caching:** Caching strategies, such as Least Recently Used (LRU) or partitioned caching, are implemented to optimize resource utilization by storing frequently accessed vector data.
- **Data Replication:** Multiple copies of data are created across different nodes to ensure availability and durability. Replication methods

include leaderless and leader-follower approaches, each balancing consistency and performance differently.

- **Indexing and Search:** Efficient indexing is crucial for handling high-dimensional vector data. Techniques like Approximate Nearest Neighbor (ANN) search and graph-based search methods facilitate quick and accurate similarity searches.
- **Support for Diverse Vector Types:** Vector databases accommodate various vector types (for example, dense, sparse, binary) and their unique characteristics, necessitating a flexible indexing system.
- **Distributed Processing:** To manage large-scale data and queries, vector databases support distributed parallel processing, addressing challenges like load balancing and fault tolerance.
- **Integration with Machine Learning Frameworks:** Seamless integration with popular machine learning frameworks (like TensorFlow, PyTorch) is essential for generating and utilizing vector embeddings, requiring robust APIs and connectors.
- **Cross-Modal Support:** Vector databases enable cross-modal searches, allowing users to retrieve relevant data across different formats (text, images, and the like.), thus enhancing the richness of search results.
- **Real-Time Knowledge Retrieval:** Real-time search capabilities are supported, providing access to the latest information and minimizing the chances of generating inaccurate results (hallucination).

[Indexing Techniques](#)

Vector similarity search is computationally expensive by nature. A naive brute-force search compares a query vector against every stored vector, an approach that collapses under scale. Milvus addresses this by supporting **Approximate Nearest Neighbor (ANN)** indexing techniques that trade a small amount of accuracy for massive gains in speed, memory efficiency, and scalability.

Milvus is architected to support **multiple index types**, allowing practitioners to choose the right geometric shortcut based on data size, dimensionality, latency requirements, recall targets, and hardware

availability. These index types broadly fall into **CPU-optimized** and **GPU-accelerated** families.

CPU-Based Indexing Techniques in Milvus

CPU indexes are the most widely used in production due to their maturity, flexibility, and lower infrastructure cost. They are implemented primarily using FAISS and custom Milvus logic.

FLAT (Brute Force)

Category:CPU

Type: Exact search

FLAT indexing performs a full linear scan over all vectors, computing exact distances. There is no approximation, no compression, and no clever geometry.

Characteristics:

- Highest recall (100%)
- No preprocessing or training
- Extremely high latency at scale
- Linear growth in compute cost

Use cases:

- Small datasets
- Ground truth evaluation
- Debugging or benchmarking other indexes

FLAT is the philosophical baseline, the Platonic ideal of similarity, but impractical beyond modest data sizes.

IVF (Inverted File Index)

Category:CPU

Type: Coarse quantization + ANN

IVF partitions the vector space into clusters using k-means. Each vector is assigned to a centroid, and searches only probe a subset of clusters (nprobe)

instead of the entire dataset.

Key parameters:

- nlist: number of clusters
- nprobe: number of clusters searched at query time

Characteristics:

- Significant speedup over FLAT
- Tunable recall versus latency
- Requires training phase
- Performance degrades if clusters are poorly balanced

Use cases:

- Medium to large datasets
- When latency matters but recall can be tuned
- General-purpose semantic search

IVF turns the search problem from “search everywhere” into “search where it makes sense.”

IVF_PQ (Inverted File + Product Quantization)

Category:CPU

Type: ANN + vector compression

IVF_PQ builds on IVF by compressing vectors using **Product Quantization (PQ)**. Vectors are split into subspaces, each quantized separately, dramatically reducing memory usage.

Characteristics:

- Much lower memory footprint
- Faster cache-friendly search
- Lower recall than IVF_FLAT
- Training complexity is higher

Use cases:

- Very large datasets (hundreds of millions to billions)

- Memory-constrained environments
- Recommendation systems at scale

IVF_PQ is where geometry meets thrift, squeezing vectors into compact symbolic representations.

HNSW (Hierarchical Navigable Small World Graph)

Category:CPU

Type: Graph-based ANN

HNSW builds a multi-layer graph where vectors are nodes connected by proximity edges. Search becomes a greedy graph traversal rather than a partitioned lookup.

Key parameters:

- M: graph connectivity
- efConstruction: build-time accuracy
- efSearch: query-time accuracy

Characteristics:

- Very high recall with low latency
- No clustering step
- High memory usage
- Slower index build time

Use cases:

- High-recall semantic search
- RAG pipelines
- Interactive AI systems

HNSW behaves like a well-designed city map: local roads, highways, and shortcuts—all working together.

SCANN (Selective Search, Limited Support)

Category:CPU

Type: Hybrid ANN

SCANN combines tree partitioning with asymmetric hashing. Despite being powerful, its support in Milvus is more limited compared to IVF and HNSW.

Characteristics:

- Good speed-recall tradeoff
- Complex tuning
- Less commonly deployed

[GPU-Based Indexing Techniques in Milvus](#)

GPU indexing is designed for **extreme throughput** and **real-time workloads**, trading hardware cost for massive parallelism.

[GPU_FLAT](#)

Category:GPU

Type: Exact search

GPU_FLAT performs brute-force search using CUDA cores. While still $O(N)$, it exploits GPU parallelism to reduce latency dramatically.

Characteristics:

- Exact results
- Very high throughput
- High GPU memory usage
- Limited by GPU VRAM

Use cases:

- Real-time exact search
- Small to medium datasets
- Latency-critical applications

[GPU_IVF_FLAT](#)

Category: GPU

Type: ANN

This is the GPU-accelerated version of IVF_FLAT. Clustering and probing are executed on the GPU.

Characteristics:

- Orders-of-magnitude faster than CPU IVF
- Requires careful memory management
- GPU memory bound

Use cases:

- Large-scale semantic search
- Real-time recommendation systems
- Streaming inference pipelines

[GPU IVF PQ](#)

Category: GPU

Type: ANN + compression

It combines GPU parallelism with PQ compression.

Characteristics:

- Best for massive datasets
- Excellent throughput
- Slightly lower recall
- More complex deployment

Use cases:

- Billion-scale vector databases
- Enterprise-grade AI platforms
- Multitenant AI systems

[CPU versus GPU Indexing: A Strategic View](#)

CPU indexing processes data sequentially or with limited parallelism, making it efficient for complex logic and smaller datasets. GPU indexing uses massive parallel cores to build and query indexes across large datasets much faster. As a result, GPUs excel in high-throughput tasks like vector search, large-scale analytics, and AI workloads

Dimensions	CPU Indexes	GPU Indexes
Cost	Lower	Higher
Latency	Moderate	Very low
Throughput	Limited	Extremely high
Memory	RAM-based	VRAM-limited
Deployment	Easier	More complex
Use Case	General production	Real-time, hyperscale

Table 1.2: CPU versus GPU Indexing

In practice, many Milvus deployments adopt **hybrid architectures**, using CPU indexes for long-term storage and GPU indexes for hot or frequently queried data.

Scalability and Performance

Scalability is the ability of a system to handle increasing workloads by efficiently adding resources such as compute, memory, or storage. Performance measures how quickly and efficiently the system processes tasks under a given load. Well-designed systems balance both, ensuring that performance remains stable even as scale grows.

Indexing Techniques

- **Description:** Efficient indexing is crucial for fast retrieval of high-dimensional vectors. Techniques such as Approximate Nearest Neighbor (ANN) search and graph-based indexing improve search speed and accuracy.
- **Mechanism:** Indexing structures like HNSW (Hierarchical Navigable Small World) graphs or LSH (Locality Sensitive Hashing) are used to quickly locate nearest neighbors in high-dimensional space.

Caching Strategies

- **Description:** Caching frequently accessed vectors or query results reduces the need for repetitive computations, improving response times.
- **Mechanism:** Implementations like Least Recently Used (LRU) or partitioned caching store recent queries and vector results in memory, enabling faster access.

Efficient Query Processing

- **Description:** Vector databases optimize query processing through techniques such as vector quantization and dimensionality reduction, which enhance performance without compromising accuracy.
- **Mechanism:** Methods like Product Quantization (PQ) or Principal Component Analysis (PCA) reduce the complexity of similarity searches by approximating high-dimensional vectors with lower-dimensional representations.

Real-Time Data Handling

- **Description:** For applications requiring real-time updates and queries, vector databases are optimized to handle rapid data ingestion and querying.
- **Mechanism:** Stream processing frameworks and in-memory data structures enable real-time analytics and quick retrieval of newly added vectors.

Fault Tolerance and High Availability

- **Description:** To ensure continuous operation and reliability, vector databases implement fault tolerance and high availability strategies.
- **Mechanism:** Data replication and redundancy mechanisms ensure that data is preserved across failures, and failover strategies allow for uninterrupted service.

Conclusion

Vector databases have emerged as a foundational component of modern AI systems, reshaping how machines store, search, and reason over unstructured data. As intelligence shifted from rule-driven logic to representation learning, data began to be encoded as high-dimensional vectors that capture semantic meaning rather than explicit structure. Traditional relational and keyword-based databases struggle in this regime, as they are optimized for exact matches and fixed schemas. Vector databases address this gap by enabling efficient similarity search over embeddings, making them indispensable in an era dominated by large language models, multimodal AI, and context-aware applications.

From an application standpoint, vector databases serve as the backbone for a wide range of intelligent systems. They power semantic search engines that understand intent beyond keywords, recommendation systems that learn taste and behavior, and retrieval-augmented generation pipelines that ground language models in external knowledge. Their relevance spans domains such as enterprise search, customer support automation, healthcare, finance, robotics, defense, and scientific discovery. Any system that must reason over meaning, relationships, or perception rather than rigid identifiers naturally benefits from vector-based storage and retrieval. As embeddings become the universal interface between data and models, vector databases act as the connective tissue that binds learning to memory.

Performance and scalability are where vector databases distinguish themselves as production-grade infrastructure rather than experimental tools. By employing approximate nearest neighbor indexing techniques, intelligent partitioning, compression, and distributed execution, they can handle millions to billions of vectors while maintaining low latency and high throughput. Modern vector databases are designed to scale horizontally, decouple compute from storage, and adapt to evolving data distributions without retraining entire systems. This ability to sustain semantic performance under massive data growth transforms AI from a laboratory artifact into a reliable, industrial capability. In essence, vector databases are not just faster storage engines; they are enablers of scalable intelligence in the data-rich world.

The upcoming chapter will focus on the basics of the mathematical principles of Vectors and how it forms the building block of AI systems to later investigate how Vector systems handle this data form which is the

interjection of the development of the next generation of computing applications.

You've Just Finished your Free Sample

Enjoyed the preview?

Buy: <http://www.ebooks2go.com>