

# Hybrid Search

Tim Spann @ Zilliz



# Slides



X



# Tim Spann

Principal Developer  
Advocate, Zilliz

[tim.spann@zilliz.com](mailto:tim.spann@zilliz.com)

<https://www.linkedin.com/in/timothyspann/>

<https://x.com/PaaSDev>



# Show Me A Demo



A screenshot of the Milvus UI interface. At the top, there's a navigation bar with "Database: default" and "airquality". Below it is a sidebar with icons for "Overview", "Vector Search", "Data", "Partitions", and "Segments". The main area shows an "Overview" of the "airquality" index, which has 34 partitions. It includes fields like "Created Time" (2024-04-24 03 PM) and "Consistency" (Bounded). A "Schema" table lists fields such as "id", "DateObserved", "HourObserved", "Latitude", "Longitude", "ParameterName", "ZipCode", "AQI", "vector", "details", and "location", each with its type (e.g., Int64, VarChar), index name, index type (e.g., id, vector), and parameters.



Lots  
of Slides

Cool Demo

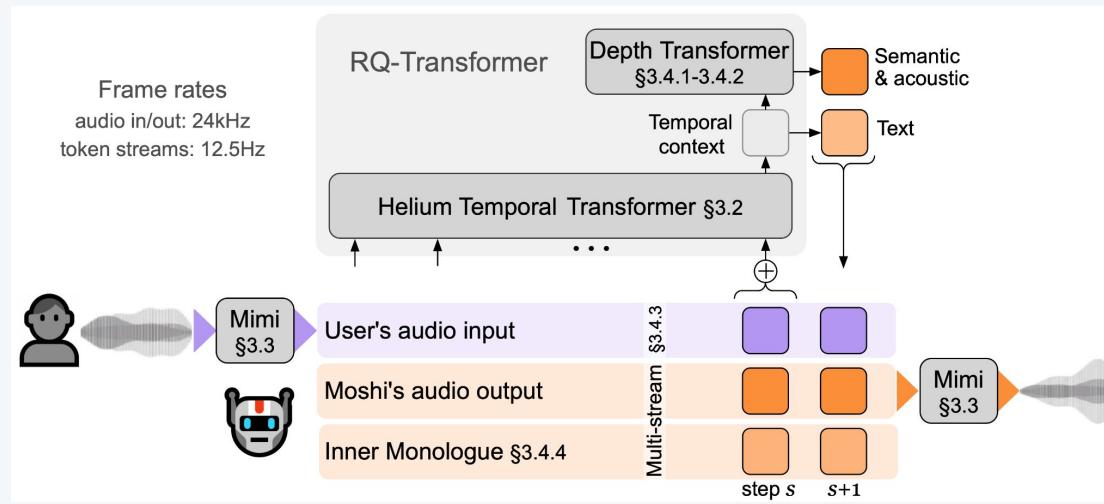


<https://multimodal-demo.milvus.io/>

<https://milvus.io/milvus-demos/reverse-image-search>

# What's New?

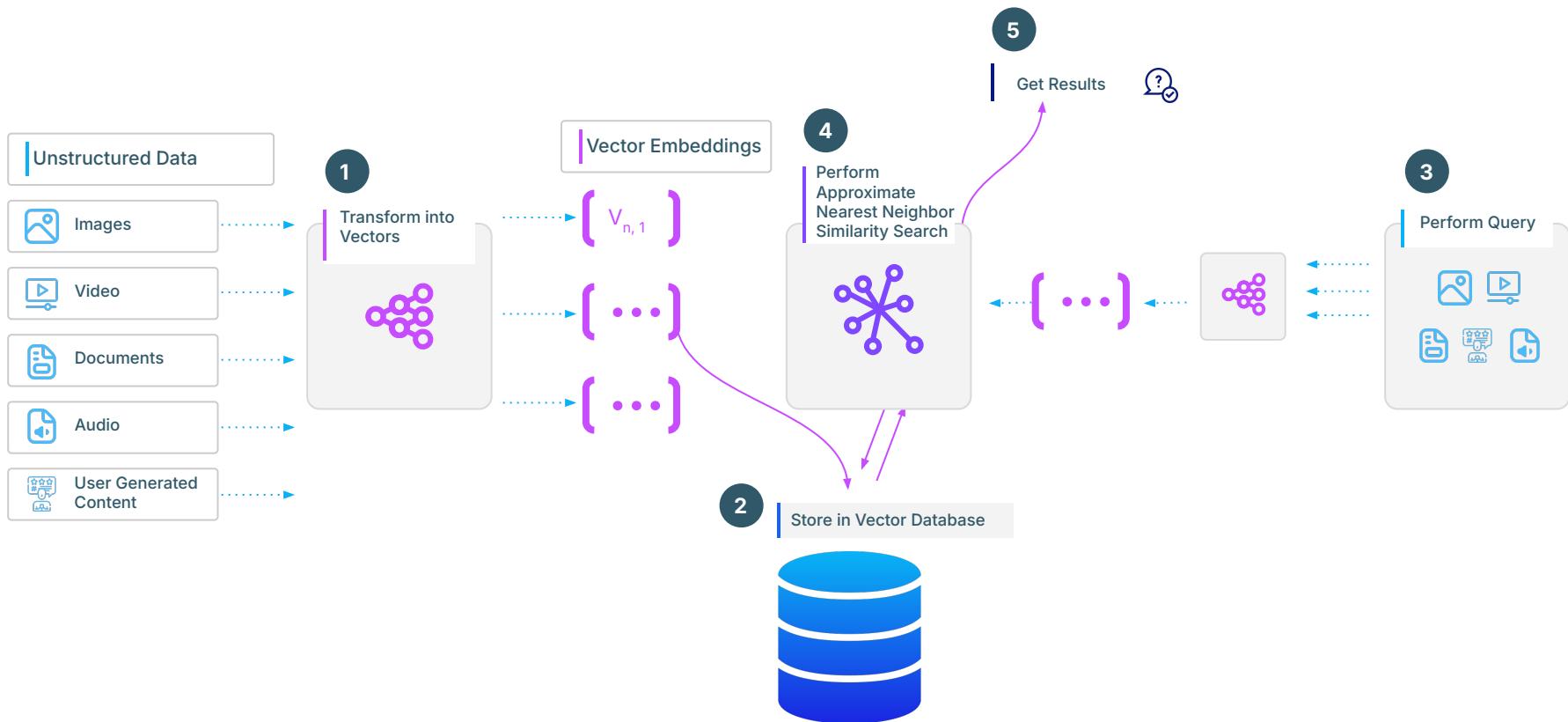
- <https://huggingface.co/mistralai/Mistral-Small-Instruct-2409>
- <https://github.com/kyutai-labs/moshi>



# 02

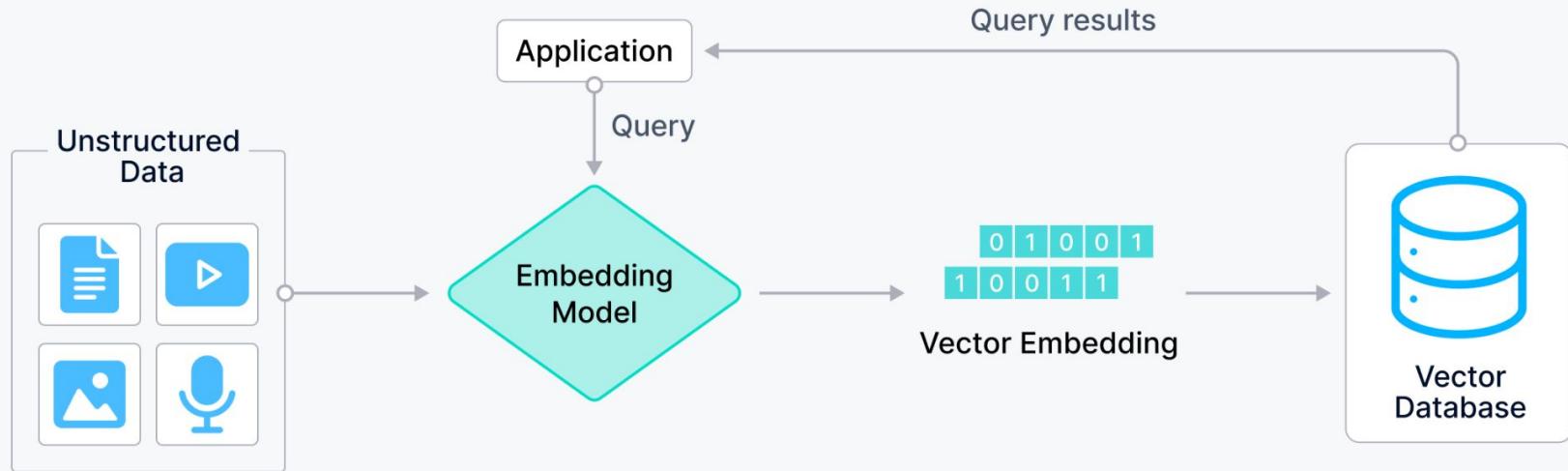
## Overview of Vector Databases

# How Similarity Search Works

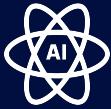


# Vector Database : making sense of unstructured data

A vector database stores embedding vectors and allows for semantic retrieval of various types of unstructured data.



# Do you really need a Vector Database?



## ANN Libraries

- FAISS, ANNOY, HNSW
- Supports 1M vectors
- Good for prototyping

## Existing Solutions

- 50M~100M vectors
- PostgreSQL, ElasticSearch, BigQuery, MongoDB, etc with ANNS plug-ins

## Vector Databases

- Purpose-built for vectors to support the requirements and lifecycle of vectors
- Billion+ scale
- CRUD, real-time search, top-k/range/hybrid search, multi-modal, multi-vector query, distributed
- Semantic Search is core to your business

Vector Databases are **purpose-built** to handle **indexing, storing, and querying** vector data.

# Vector Search + Indexing + Filtering =





# 03

## A Quick Introduction to Milvus



### Easy Setup

Pip-install to start coding in a notebook within seconds.



### Reusable Code

Write once, and deploy with one line of code into the production environment



### Integration

Plug into OpenAI, Langchain, LlmalIndex, and many more



### Feature-rich

Dense & sparse embeddings, filtering, reranking and beyond



**Milvus is an open-source vector database for GenAI projects.** pip install on your laptop, plug into popular AI dev tools, and push to production with a single line of code.



**29.4K+**

GitHub Stars



**2,800**

Folks



**25M+**

Downloads



**250+**

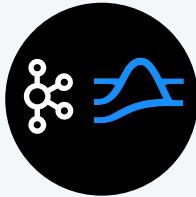
Contributors

# Milvus 🤝 Open-Source



## MINIO

Store Vectors and Indexes  
Enables Milvus' stateless  
architecture



## Kafka/ Pulsar

Handles Data Insertion  
stream  
Internal Component  
Communications  
Real-time updates to  
Milvus



## Prometheus / Grafana

Collects metrics from  
Milvus  
Provides real-time  
monitoring dashboards



## Kubernetes

Milvus Operator CRDs

# Distributed Architecture

## Microservice components

Query Coord

Data Coord

Root Coord

Query Node

Data Node

Index Node

Proxy

## Reliable States

Log Broker  
(Kafka/Pulsar)

Object Storage  
(S3/MinIO)

Key-Value-Meta-Store  
(etcd)



# Dynamic Scaling



## Stateless components for Easy Scaling

- **Query, Index, and Data Nodes** can be scaled **independently**
- Allows for **optimized resource allocation** based on workload characteristics

## Data sharding across multiple nodes

- **Distributes large datasets** across multiple **Data Nodes**
- Enables **parallel processing** for improved query performance

## Horizontal Pod Autoscaler (HPA)

- **Automatically scales** up and down
- **Custom metrics** can be used (e.g., query latency, throughput)

# Stateless Architecture



## Stateless Components

All Milvus components are deployed **Stateless**.



## Object Storage

Milvus relies on **Object Storage** (MinIO, S3, etc) for data **persistence**.

**Vectors** are stored in **Object Storage**, **Metadata** is in **etcd**.



## Scaling and Failover

**Scaling** and failover **don't** involve traditional **data rebalancing**.

When **new pods** are added or existing ones fail, they can immediately start handling requests by **accessing data** from the **shared object storage**.

# Different Consistency levels

Ensures every node or replica has the same view of data at a given time.

- **Strong**: Guaranteed up-to-date reads, highest latency
- **Bounded**: Reads may be slightly stale, but within a time bound
- **Session**: Consistent reads within a session, may be stale across sessions
- **Eventually**: Lowest latency, reads may be stale

## Trade Offs

- **Strong** consistency for **critical applications** requiring accurate results
- **Eventually** consistency for **high-throughput, latency-sensitive apps**

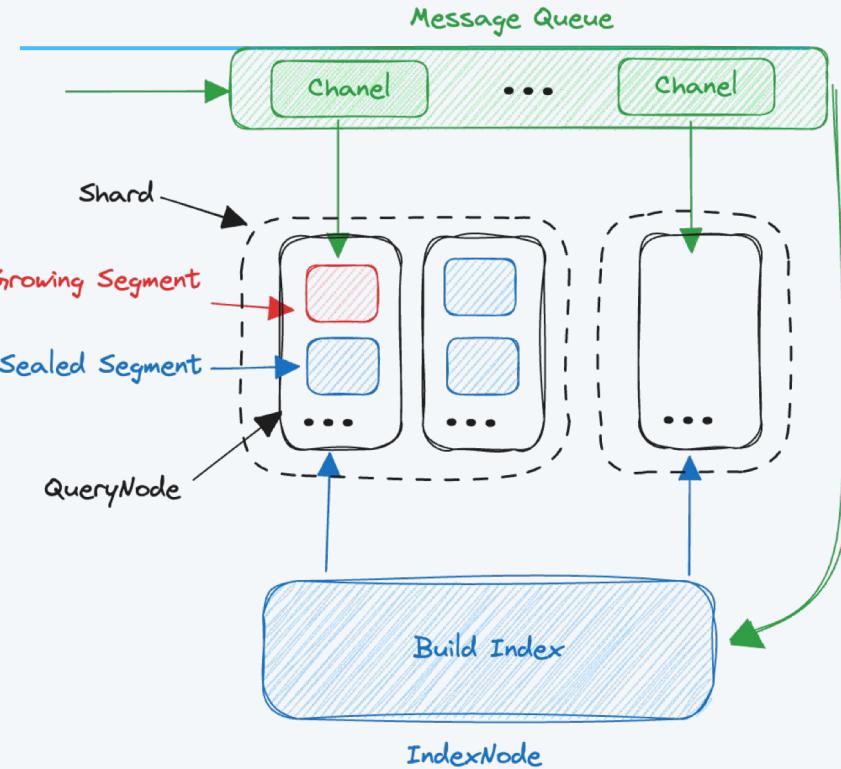
# Milvus Data Layout - Segments

## Growing Segment:

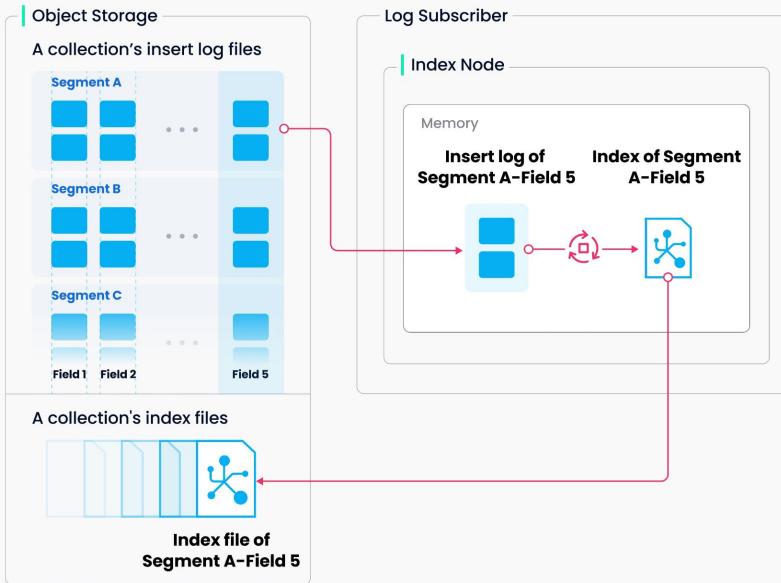
- In-memory segment replaying data from the Log Broker.
- Uses a FLAT index to ensure data is fresh and appendable.

## Sealed Segment:

- Immutable segment using alternative indexing methods for efficiency.



# Index Building



To **avoid frequent index building** for data updates.

A **collection** in Milvus is **divided** further **into segments**, each with its own index.

# Picking an Index

- 100% Recall – Use FLAT search if you need 100% accuracy
- 10MB < `index_size` < 2GB – Standard IVF
- 2GB < `index_size` < 20GB – Consider PQ and HNSW
- 20GB < `index_size` < 200GB – Composite Index, IVF\_PQ or HNSW\_SQ
- Disk-based indexes

# Indexes

Most of the vector index types supported by Milvus use approximate nearest neighbors search (ANNS),

- **HNSW**: HNSW is a graph-based index and is best suited for scenarios that have a high demand for search efficiency. There is also a GPU version **GPU\_CAGRA**, thanks to Nvidia's contribution.
- **FLAT**: FLAT is best suited for scenarios that seek perfectly accurate and exact search results on a small, million-scale dataset. There is also a GPU version **GPU\_BRUTE\_FORCE**.
- **IVF\_FLAT**: IVF\_FLAT is a quantization-based index and is best suited for scenarios that seek an ideal balance between accuracy and query speed. There is also a GPU version **GPU\_IVF\_FLAT**.
- **IVF\_SQ8**: IVF\_SQ8 is a quantization-based index and is best suited for scenarios that seek a significant reduction on disk, CPU, and GPU memory consumption as these resources are very limited.
- **IVF\_PQ**: IVF\_PQ is a quantization-based index and is best suited for scenarios that seek high query speed even at the cost of accuracy. There is also a GPU version **GPU\_IVF\_PQ**.

# Indexes Continued.

- **SCANN:** SCANN is similar to IVF\_PQ in terms of vector clustering and product quantization. What makes them different lies in the implementation details of product quantization and the use of SIMD (Single-Instruction / Multi-data) for efficient calculation.
- **DiskANN:** Based on Vamana graphs, DiskANN powers efficient searches within large datasets.

# New Stuff

## The Zilliz Sep '24 Launch

Supercharge Your GenAI with Production-Ready Data Infrastructure



### Data Sovereignty

#### Migration Services

Secure, complete and easy data migration to Zilliz from Milvus, PGVector, Elastic or between Zilliz clusters

#### Fivetran Source Connector

Enable unstructured data retrievals from 500+ system through Fivetran



### High-performance

#### Auto-Scale (Private Preview)

Dynamically increase cluster capacity based on workload to reduce operational burden and ensure continuous service

#### Multi-replica (Public Preview)

Replicated collections for increased throughput and fault tolerance



### Security & Reliability

#### Metrics & Alerts

18 core metrics for proactive issue identification and 39 alerts to ensure optimal performance

#### Auth0-based SSO

(Private Preview)  
Simplify IT management and improve security w/ Google/Github/SSO/Email login

99.95% Uptime SLAs

# New Stuff

<https://github.com/milvus-io/milvus-sdk-java/releases/tag/v2.4.4>

Milvus 2.4 introduces several new features and improvements:

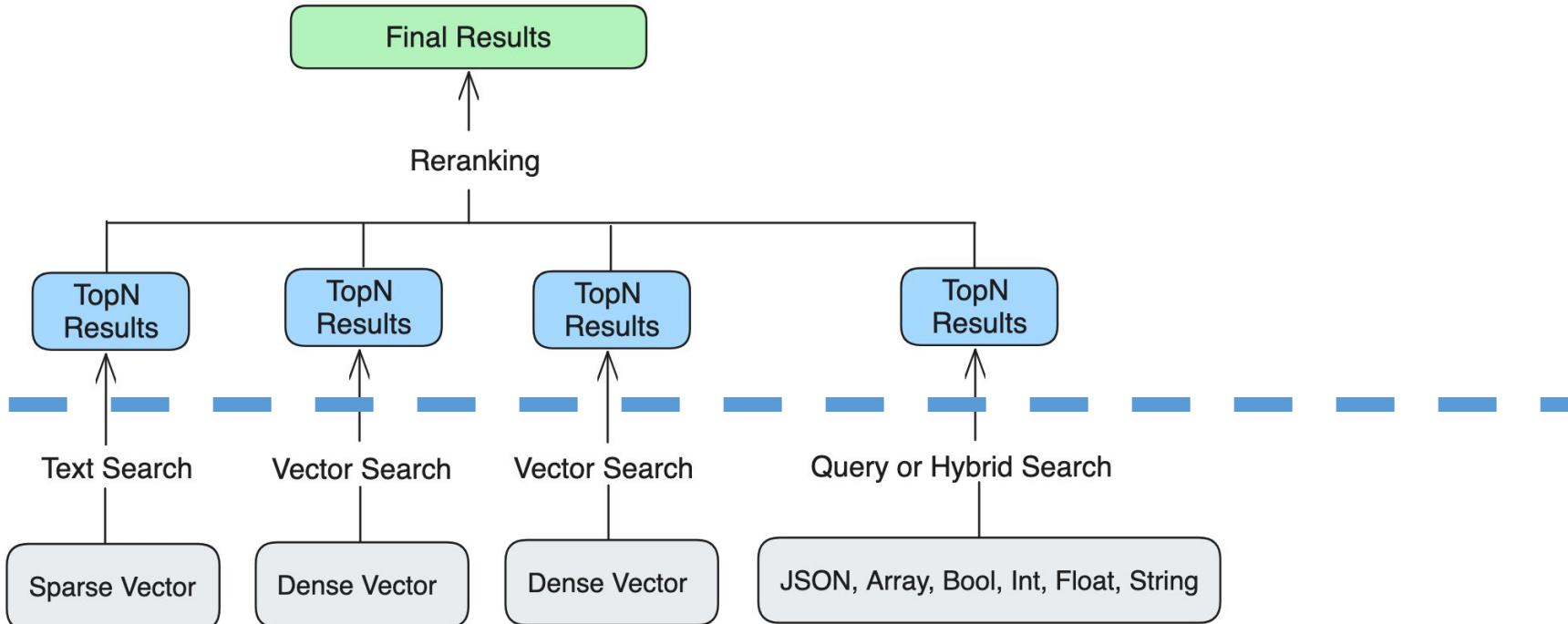
1. **New GPU Index - CAGRA:** This GPU-based index offers significant performance improvements, especially for batch searches
2. **Multi-vector and Hybrid Search:** This feature allows storing vector embeddings from multiple models and conducting hybrid searches.
3. **Sparse Vectors Support (Beta):** Milvus now supports sparse vectors for processing in collections, which is particularly useful for keyword interpretation and analysis
4. **Grouping Search:** This feature enhances document-level recall for Retrieval-Augmented Generation (RAG) applications by providing categorical aggregation
5. **Inverted Index and Fuzzy Matching:** These capabilities improve keyword retrieval for scalar fields
6. **Float16 and BF16 Vector Data Type Support:** Milvus now supports these half-precision data types for vector fields, which can improve query efficiency and reduce memory usage.
7. **L0 Segment:** This new segment is designed to record deleted data, enhancing the performance of delete and upsert operations.
8. **Refactored BulkInsert:** The bulk-insert logic has been improved, allowing for importing multiple files in a single bulk-insert request.

# MILVUS HAS VECTORS

I SEE ONLY CAT PICTURES



# Hybrid Search





The Landscape of GenAI Ecosystem: Beyond LLMs and Vector Databases

Read Blog



Workflow Orchestration / Optimization:  
LangChain, LlamaIndex, Haystack, DSPy, Semantic Kernel

Quality Evaluation / Observability:  
Ragas, Arize, Langfuse, Relari AI, Giskard, DeepEval

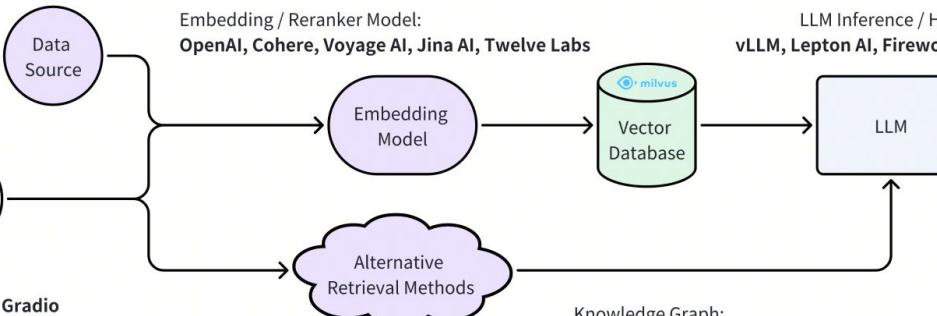
Data Connector / Web Scrapping:  
Airbyte, Fivetran, Apify, Zyte

Agent and Memory Management:  
MemGPT, MemO, Camel, AutoGPT, CrewAI, AutoGen

Embedding / Reranker Model:  
OpenAI, Cohere, Voyage AI, Jina AI, Twelve Labs

LLM Inference / Hosting:  
vLLM, Lepton AI, Fireworks AI, Octo AI

Frontend:  
Streamlit, Vercel, Gradio



Out-of-the-box / Low-code RAG Service:  
AnythingLLM, PrivateGPT, Dify, Shakudo, Vectara, Epsilla, FlowiseAI

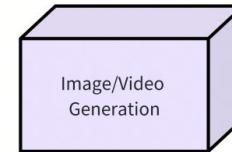


Image / Video Generation, Digital Twin:  
Midjourney, Stability AI, Runway, Pika, HeyGen



# 04

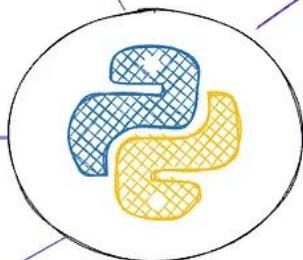
## APPS



<https://forecast.weather.gov>



**511NY**  
Traffic Cameras



Insert



Slack

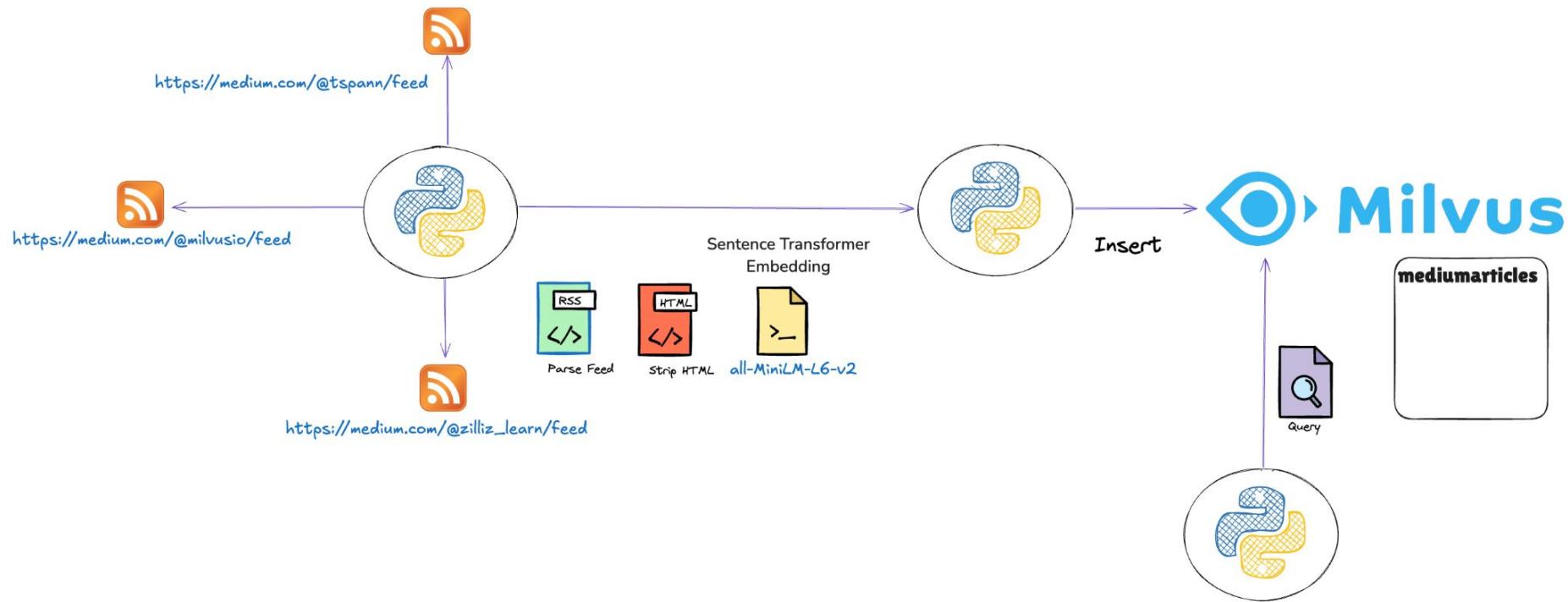


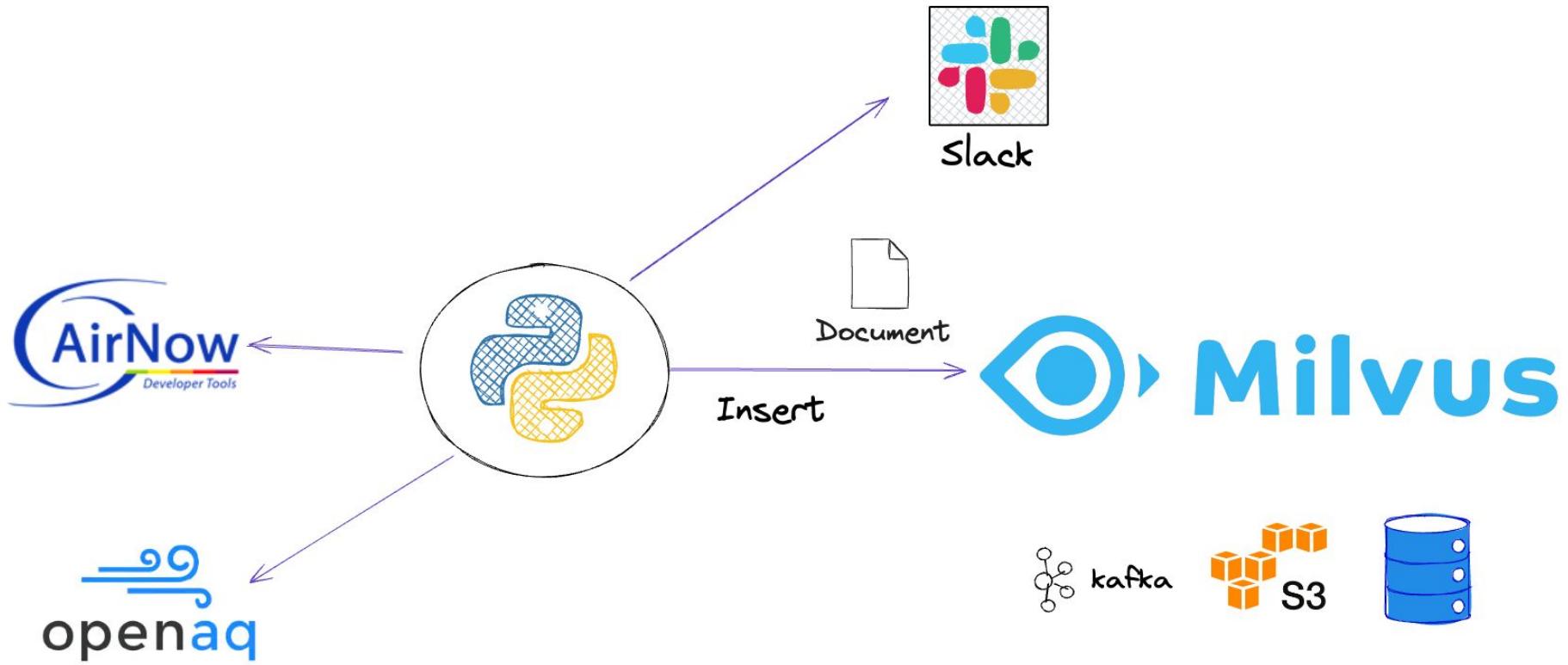
Document



**Milvus**







# 05

## Q & A

# RESOURCES



# Vector Database Resources

Give Milvus a Star!

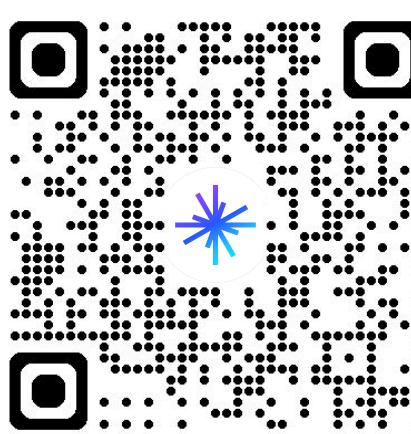


<https://github.com/milvus-io/milvus>

Chat with me on Discord!



# Unstructured Data Meetup



<https://www.meetup.com/unstructured-data-meetup-new-york/>

This meetup is for people working in unstructured data. Speakers will come present about related topics such as vector databases, LLMs, and managing data at scale. The intended audience of this group includes roles like machine learning engineers, data scientists, data engineers, software engineers, and PMs.

This meetup was formerly Milvus Meetup, and is sponsored by [Zilliz](#) maintainers of [Milvus](#).

# Generative AI Resource Hub

Tutorials, Code Examples, and Best Practices for Developing and Deploying GenAI Applications.



Learn



Build



Explore

<https://zilliz.com/learn/generative-ai>



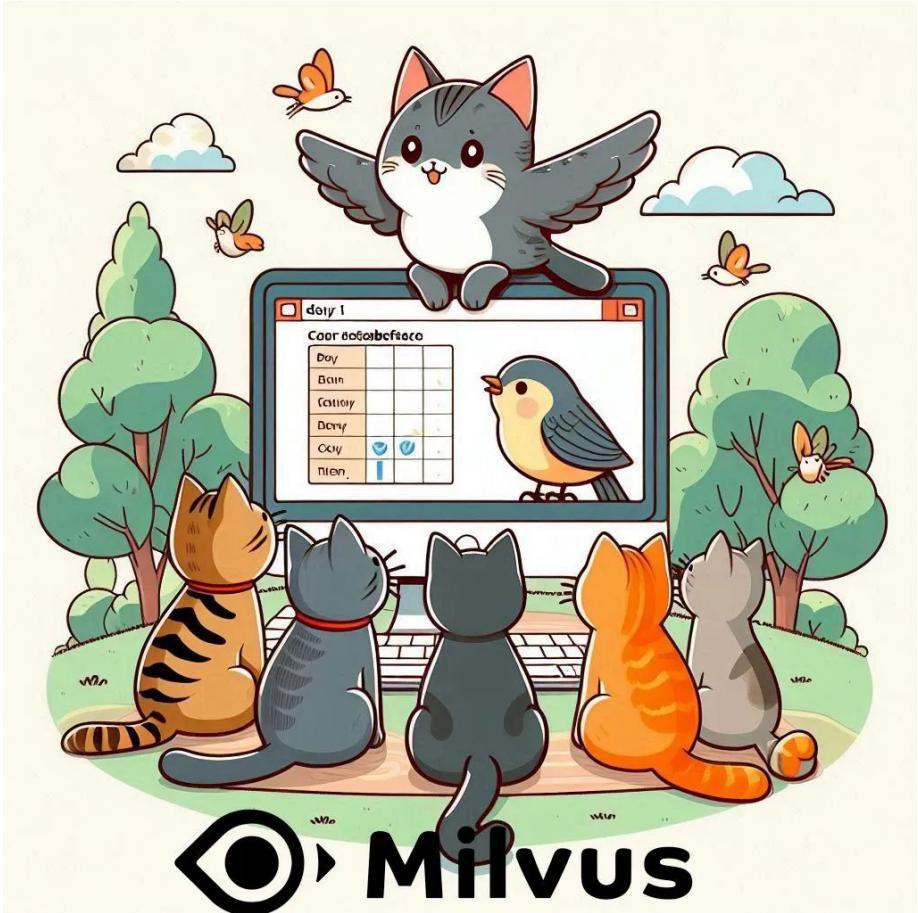




<https://medium.com/@tspann/unstructured-street-data-in-new-york-8d3cde0a1e5b>



<https://medium.com/@tspann/not-every-field-is-just-text-numbers-or-vectors-976231e90e4d>



 Milvus

<https://medium.com/@tspann/shining-some-light-on-the-new-milvus-lite-5a0565eb5dd9>



Raspberry Pi AI Kit - Hailo  
Edge AI



Milvus



<https://medium.com/@tspann/unstructured-data-processing-with-a-raspberry-pi-ai-kit-c959dd7fff47>

# AIM Weekly by Tim Spann



<https://bit.ly/32dAJft>

<https://github.com/milvus-io/milvus>

This week in Milvus, Towhee, Attu, GPT Cache, Gen AI, LLM, Apache NiFi, Apache Flink, Apache Kafka, ML, AI, Apache Spark, Apache Iceberg, Python, Java, Vector DB and Open Source friends.

# Thank you!

---



[milvus.io](https://milvus.io)



[github.com/milvus-io/](https://github.com/milvus-io/)



[@milvusio](https://twitter.com/milvusio)

# Connect with me!

---



[@paasDev](https://twitter.com/paasDev)



[/in/timothyspann](https://www.linkedin.com/in/timothyspann/)



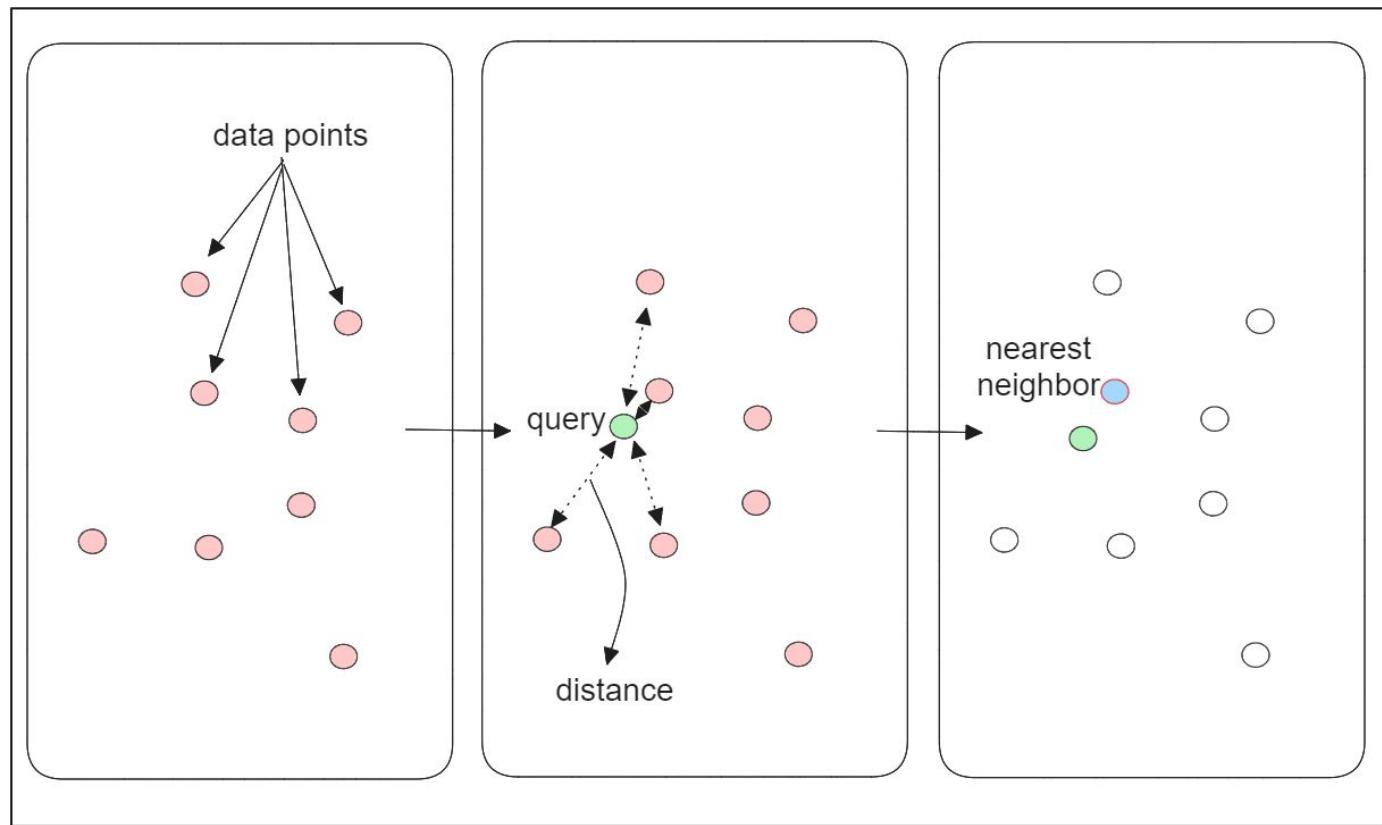
# Vector Search at Scale

# Indexing Strategies

- Tree based
- Graph based
- Hash based
- Cluster based

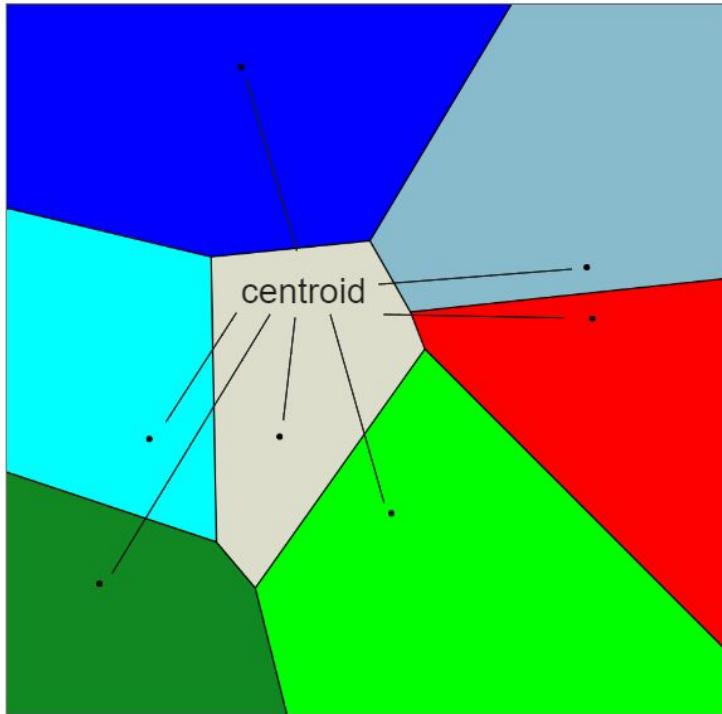
# FLAT

# FLAT Index

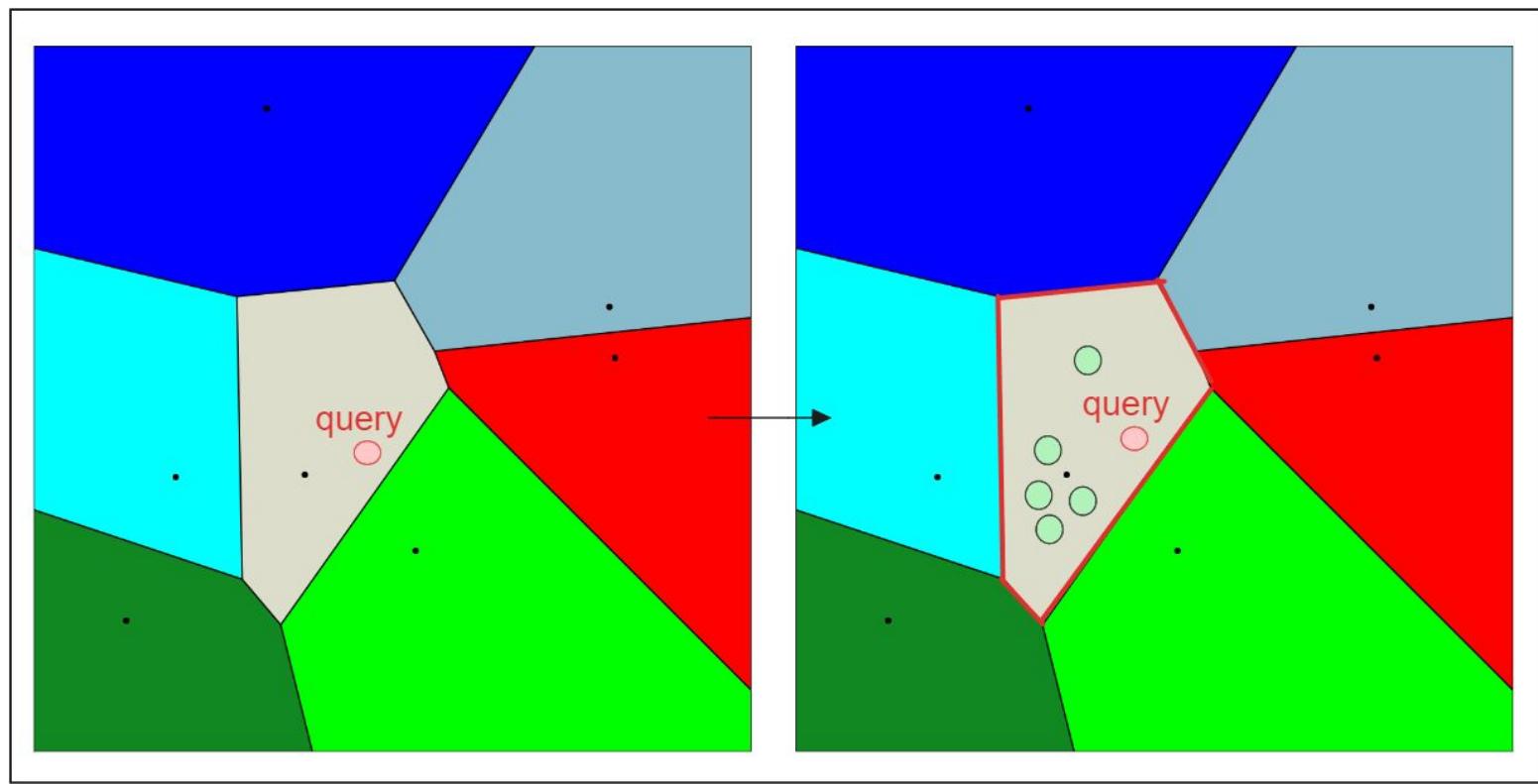


# Inverted File FLAT (IVF-FLAT)

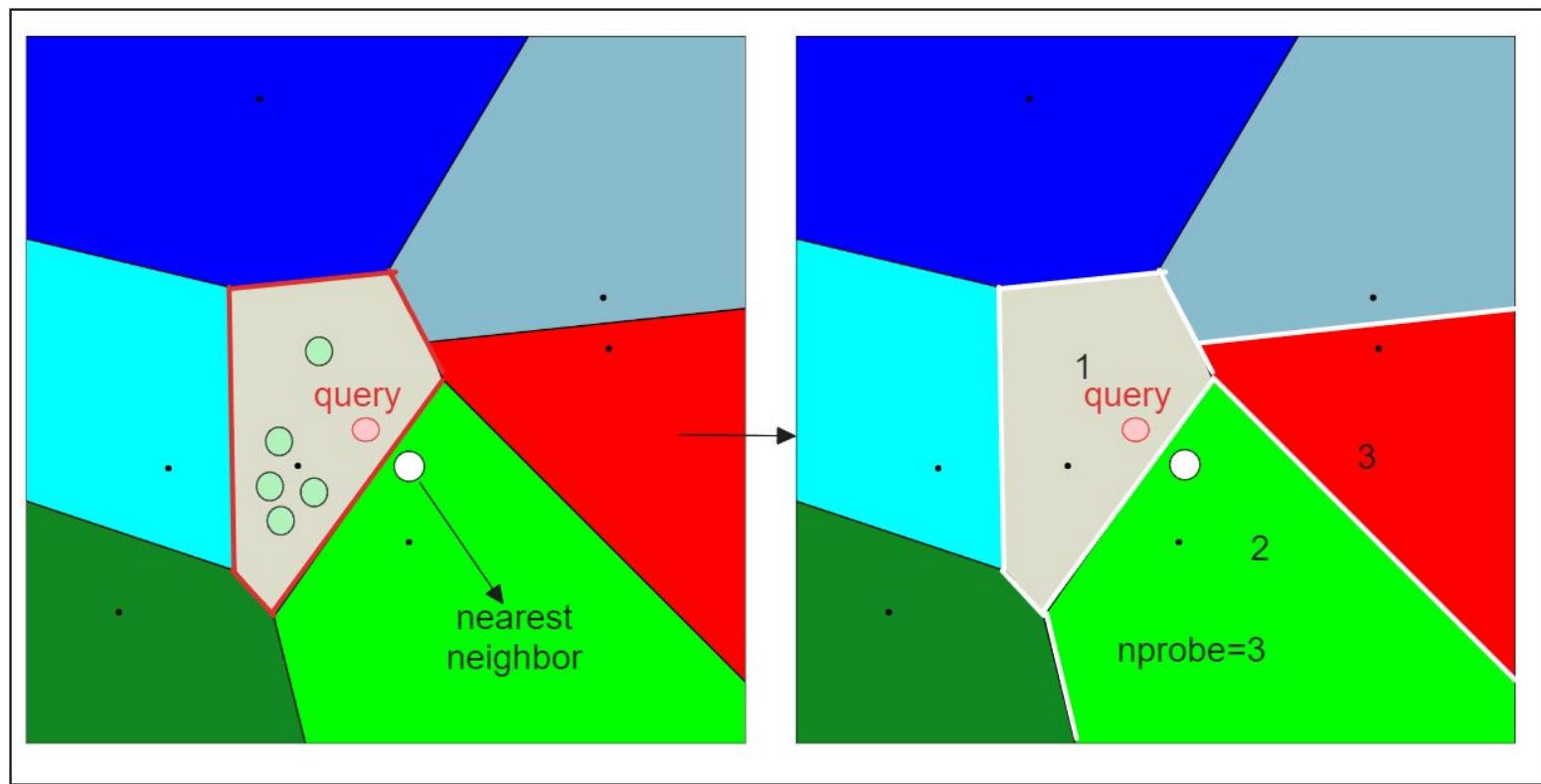
# IVF-FLAT Index



# IVF-FLAT Index



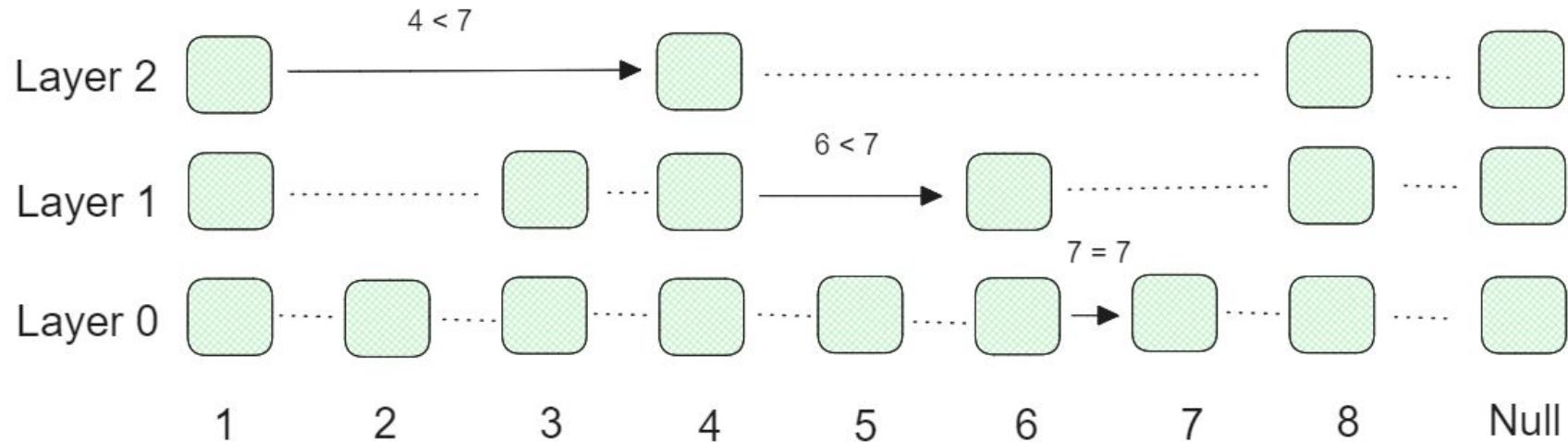
# IVF-FLAT Index



# Hierarchical Navigable Small World (HNSW)

# HNSW - Skip List

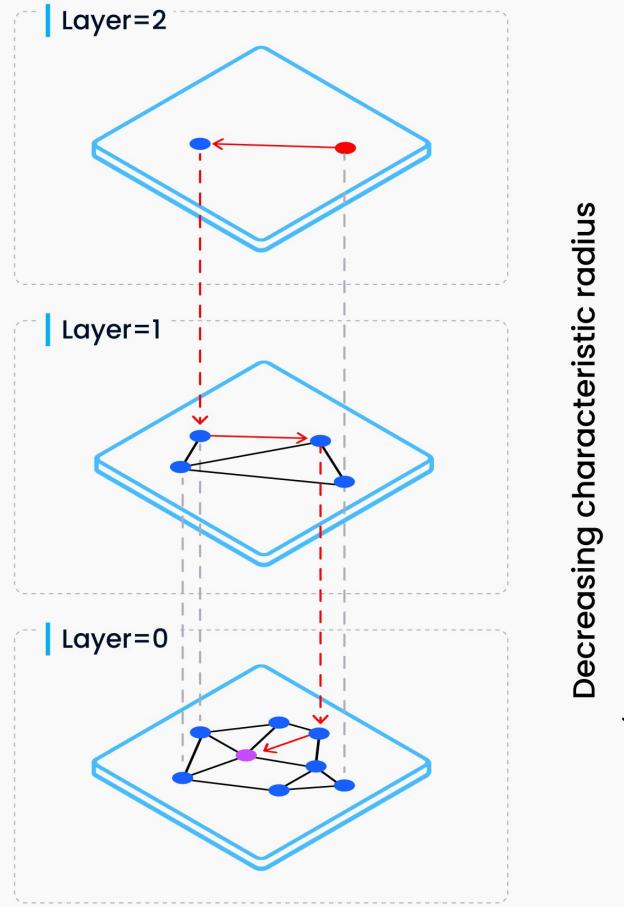
Finding element 7



# HNSW - NSW Graph

- Built by randomly shuffling data points and inserting them one by one, with each point connected to a predefined number of edges ( $M$ ).
  - ⇒ Creates a graph structure that exhibits the "small world".
  - ⇒ Any two points are connected through a relatively short path.

# HNSW



Category	Index	Accuracy	Latency	Throughput	Index Time	Cost
Graph-based	Cagra (GPU)	High	Low	Very High	Fast	Very High
	HNSW	High	Low	High	Slow	High
	DiskANN	High	High	Mid	Very Slow	Low
Quantization-based or cluster-based	ScaNN	Mid	Mid	High	Mid	Mid
	IVF_FLAT	Mid	Mid	Low	Fast	Mid
	IVF + Quantization	Low	Mid	Mid	Mid	Low

# Picking an Index

- 100% Recall – Use FLAT search if you need 100% accuracy
- 10MB < `index_size` < 2GB – Standard IVF
- 2GB < `index_size` < 20GB – Consider PQ and HNSW
- 20GB < `index_size` < 200GB – Composite Index, IVF\_PQ or HNSW\_SQ
- Disk-based indexes

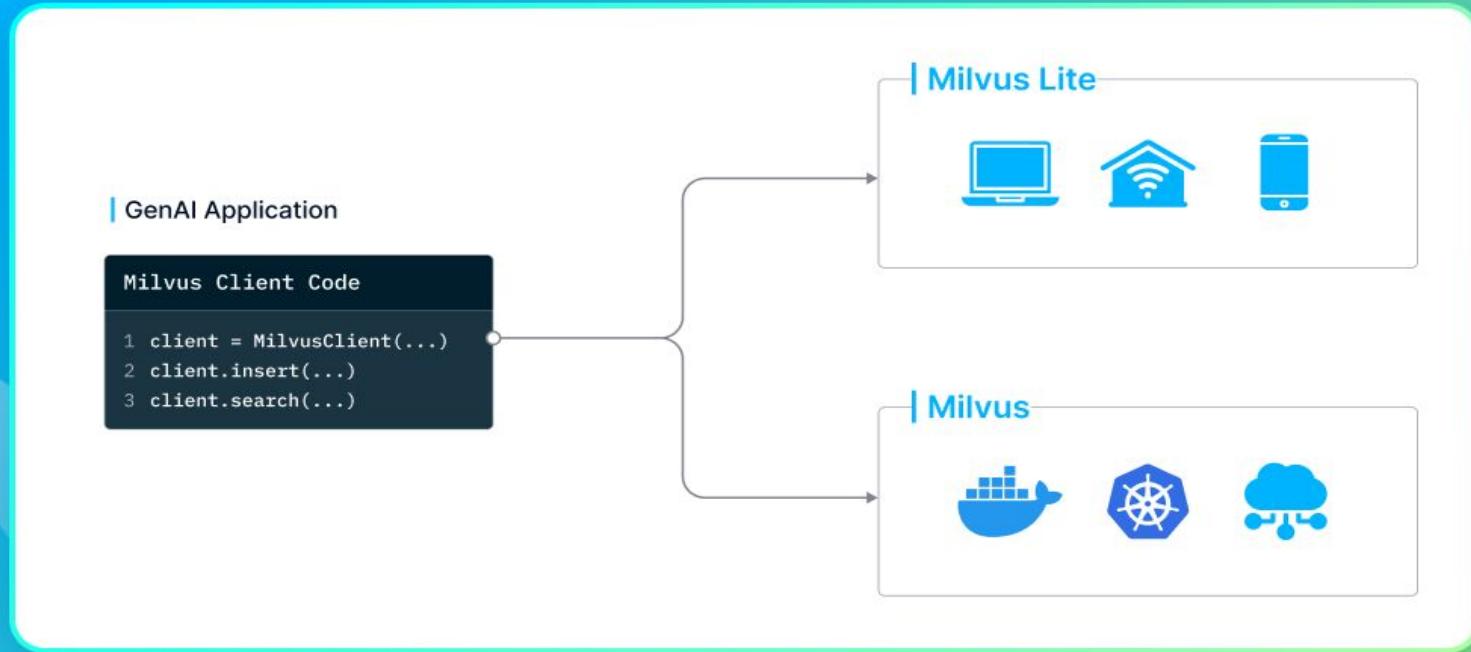
[zilliz.com/learn/choosing-right-vector-index-for-your-project](https://zilliz.com/learn/choosing-right-vector-index-for-your-project)

# Filtering

# Filtering on Metadata

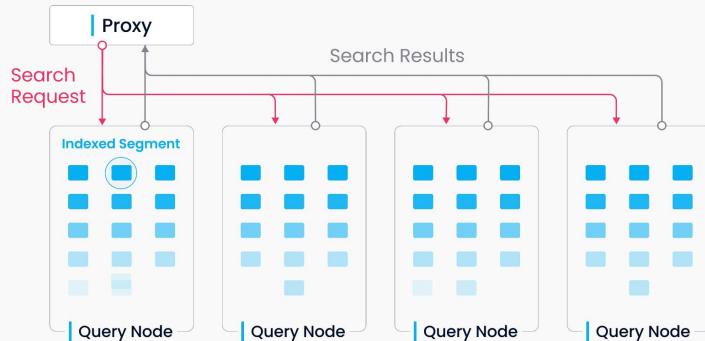
- **Search Space Reduction w/ Pre-Filtering**
- **Bitset Wizardry** 
  - Use Compact Bitsets to represent Filter Matches
  - Low-level CPU operations for speed
- **Scalar Indexing**
  - Bloom Filter
  - Hash
  - Tree-based

# Build Once Deploy Anywhere



# Scalable Search

---



- **Distributed Search** across shards
- **Parallel Processing**
- **Query Optimization**

# Why a Vector Database?

Purpose-built to store, index and query vector embeddings from unstructured data.

- Vector database

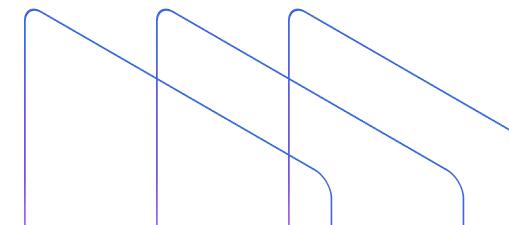
- Advanced filtering (filtered vector search, chained filters)
- Hybrid search (e.g. full text + dense vector)
- Durability (any write in a db is durable, a library typically only supports snapshotting)
- Replication / High Availability
- Sharding
- Aggregations or faceted search
- Backups
- Lifecycle management (CRUD, Batch delete, dropping whole indexes, reindexing)
- Multi-tenancy

- Vector search library

- High-performance vector search

- How do I support different applications?

- High query load
- High insertion/deletion
- Full precision/recall
- Accelerator support (GPU, FPGA)
- Billion-scale storage





## Comparison on functionalities

Feature	Milvus Lite	Milvus Standalone	Milvus Distributed
SDK / Client Library	Python gRPC	Python Go Java Node.js C# RESTful	Python Java Go Node.js C# RESTful
Data types	Dense Vector Sparse Vector Binary Vector Boolean Integer Floating Point VarChar Array JSON	Dense Vector Sparse Vector Binary Vector Boolean Integer Floating Point VarChar Array JSON	Dense Vector Sparse Vector Binary Vector Boolean Integer Floating Point VarChar Array JSON
Search capabilities	Vector Search (ANN Search) Metadata Filtering Range Search Scalar Query Get Entities by Primary Key Hybrid Search	Vector Search (ANN Search) Metadata Filtering Range Search Scalar Query Get Entities by Primary Key Hybrid Search	Vector Search (ANN Search) Metadata Filtering Range Search Scalar Query Get Entities by Primary Key Hybrid Search
CRUD operations	✓	✓	✓
Advanced data management	N/A	Access Control Partition Partition Key	Access Control Partition Partition Key Physical Resource Grouping
Consistency Levels	Strong	Strong Bounded Staleness Session Eventual	Strong Bounded Staleness Session Eventual

# Icons

ICON STYLE:



Database	Databases	Data Lake	DB Warehouse	Data Center	Cloud	Cloud to Cloud	Hybrid	Cloud	Cloud Dev	Equal Cloud	Cloud Management
Server	Payday	Workday	Docker	Operator	Kafka	KSQL Rocket	ksqldb	KSQL Circle	Connector	Microservices	Schema Registry
Streams	Event Streams	Central Nervous System	Early Production Streaming	Stream Designer	Balance	Rest	Trophy	Cluster	Certificate	Calendar	Stream Processing Cookbook
Apps	Service Apps	Coming Soon	Logs	Data Stacks	Stack Overflow	Storage	Platform	Data In	Data Out	Data Governance	Data Add
Processing	Real-time	Aggregate Data	Frameworks	Time / Money	Dev	Scale	Combine	Join	Transfer	Expand / Shrink	Add

Current

For the complete, most updated collection of Icons please go to: <https://cnfl.io/Icons>

# Icons

ICON STYLE:



Globe



Infinity



Settings



Monitoring



Anomaly  
Detection



Analytics



Real-time  
Analytics



Real-time



Processing



Process Data



Upload



Download



Computer



Devices



Computer /  
DB / Cloud



Speed



Time



Web Confirmed



RSS



ROI



Message



Quotes



Interview



# of Topics



Person



People



Webinar



Developer



Onboard



Offboard



People  
Manager



Career  
Enablement



Roadmap



Filter



Search



Solution



Features



Company  
Policies



Docs



Invoice



Blog



Podcast



Video



Book



Table



Email



Question



Check



Lock



Key



Warning



Hacker



Bug



GDPR



CCPA



Shield



Shield Open



Machine  
Learning



Continuous  
Learning



Current

For the complete, most updated collection of Icons please go to: <https://cnfl.io/Icons>

# Icons

ICON STYLE:



# of Events  
Per Day



Venue



Government



Business



Marketplace



Ecommerce



Sale



Money



Telecom



Support



Gaming



Healthcare



Badge



Love



Partner



Hand



Arm



Benefit



Thumbs Up



Swipe



Select



Promote



Awareness



Target



Car



Truck



Shirt



Food



Catalyst



Box



Puzzle



Lightening



Star



Sparkly New

Current

For the complete, most updated collection of Icons please go to: <https://cnfl.io/Icons>