



# Exploring Vector Search at Scale

Stephen Batifol

Haystack EU 24



# The Happy Beginnings



# The Tech Stack in Q4 2024



**CURSOR**



**Claude**

# Unstructured Data? No Problem!

- Docker + Vector DB
- Use some Embedding Model
- Store your Vectors in your DB, you're happy! :D



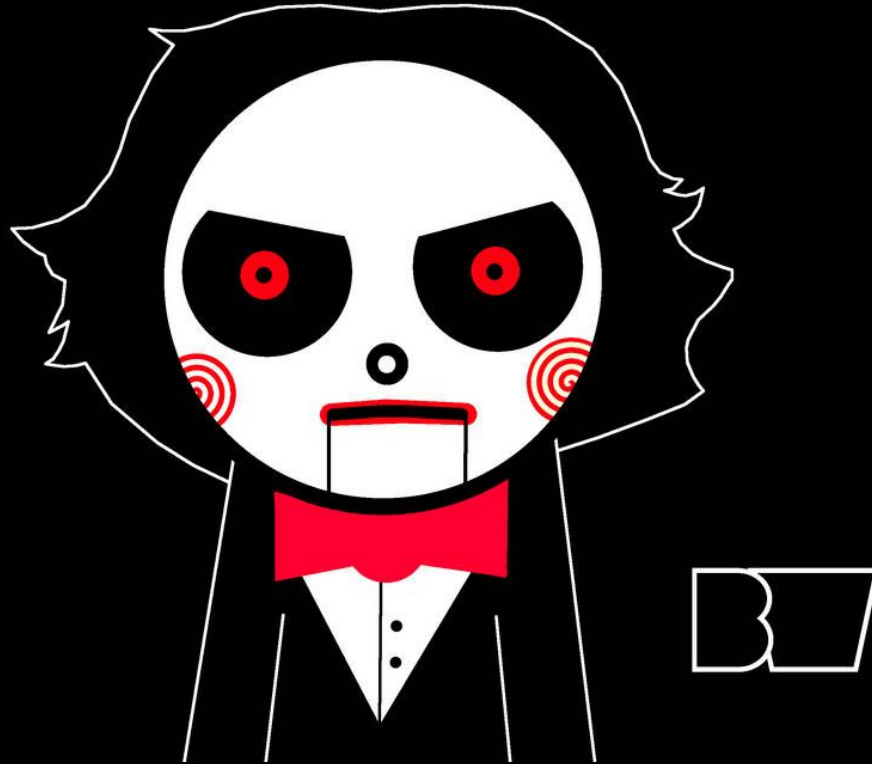
# The Growing Pains

You have more and more data!

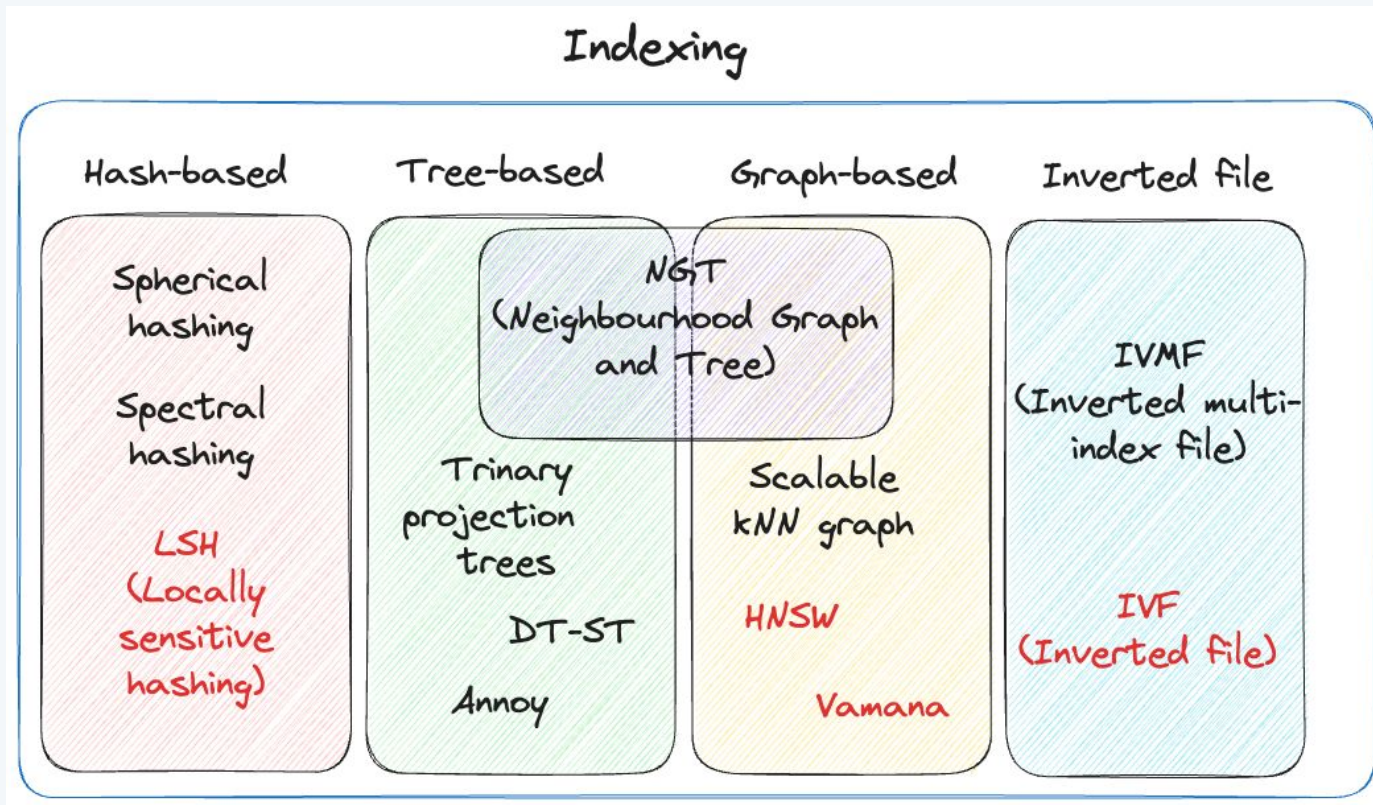
Maybe even some multimodal data 🎩

- Search Quality is declining
- Index updates take forever
- Your frustration rises!

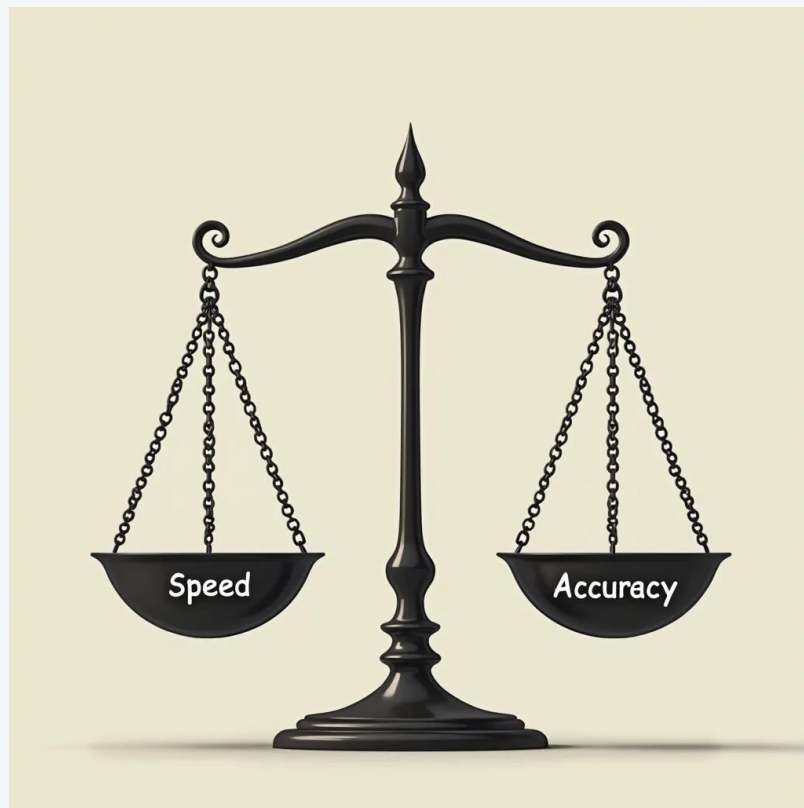
# LET THE GAME BEGIN



# Index Types Overview



# Balancing Speed & Accuracy



# Compression

## Compression

IVF Flat

HNSW Flat

Flat

IVF SQ

Vamana PQ

HNSW PQ

IVF PQ

Quantized



A night landscape with a starry sky. A bright star is visible in the upper left. In the foreground, there are dark, rocky terrain and silhouettes of evergreen trees on the left. On the right, three camels are silhouetted against a colorful, glowing horizon, standing on a rocky cliff. The sky transitions from deep blue to purple and orange near the horizon.

# *The Journey Continues*

# The Distributed Dilemma

- ✓ You've optimized your Index!
- ✓ You've chosen the right data structure
- ✓ You've dabbled in compression techniques!

More data → One instance isn't enough anymore!

⇒ It's time to go **Distributed!!**

# Sharding: Divide and Conquer

You brainstorm different sharding approaches:

- **Random sharding:** Simple, but might lead to uneven distribution.
- **Hash-based sharding:** More even distribution, but potentially tricky with updates.
- **Semantic sharding:** Grouping similar vectors together, which could speed up certain types of queries.



# Partitioning!

You sketch out a plan:

1. Incoming data gets routed to the appropriate language-specific partition.
2. Queries first determine the relevant language(s).
3. The search is executed only on the relevant partition.

This could reduce the search space for many queries, leading to faster results!

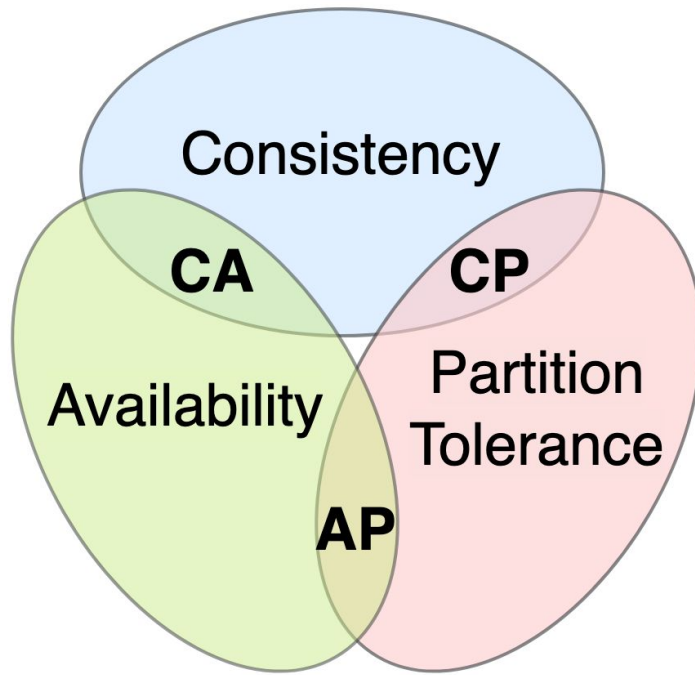


# Scaling Horizontally: The Replica Game

- You decide to implement a replica system:
  - Each shard has multiple replicas.
  - Read requests are load-balanced across replicas.
  - Write operations are synchronized across all replicas of a shard.

This setup allows you to scale out your read capacity simply by adding more replicas.

# Consistency vs Availability



# What about Real-Time?!

## Users want it now!

- Write Ahead Logs
- Two-tier Systems
- Incremental Indexing

# And Monitoring?! Keeping the Beast in Check

You set up dashboards to track:

- Query latency across different shards and replicas
- Index update times
- Resource utilization (CPU, memory, disk I/O)
- Shard balance and data distribution

You also implement automated processes for:

- Rebalancing shards when they become uneven
- Adding or removing replicas based on load
- Performing rolling updates to minimize downtime

# To Host or Not to Host?

As your system grows, you start weighing the pros and cons of cloud-hosted solutions versus managing your own infrastructure. You consider factors like:

- Scalability and elasticity
- Operational overhead
- Cost predictability
- Data privacy and compliance requirements



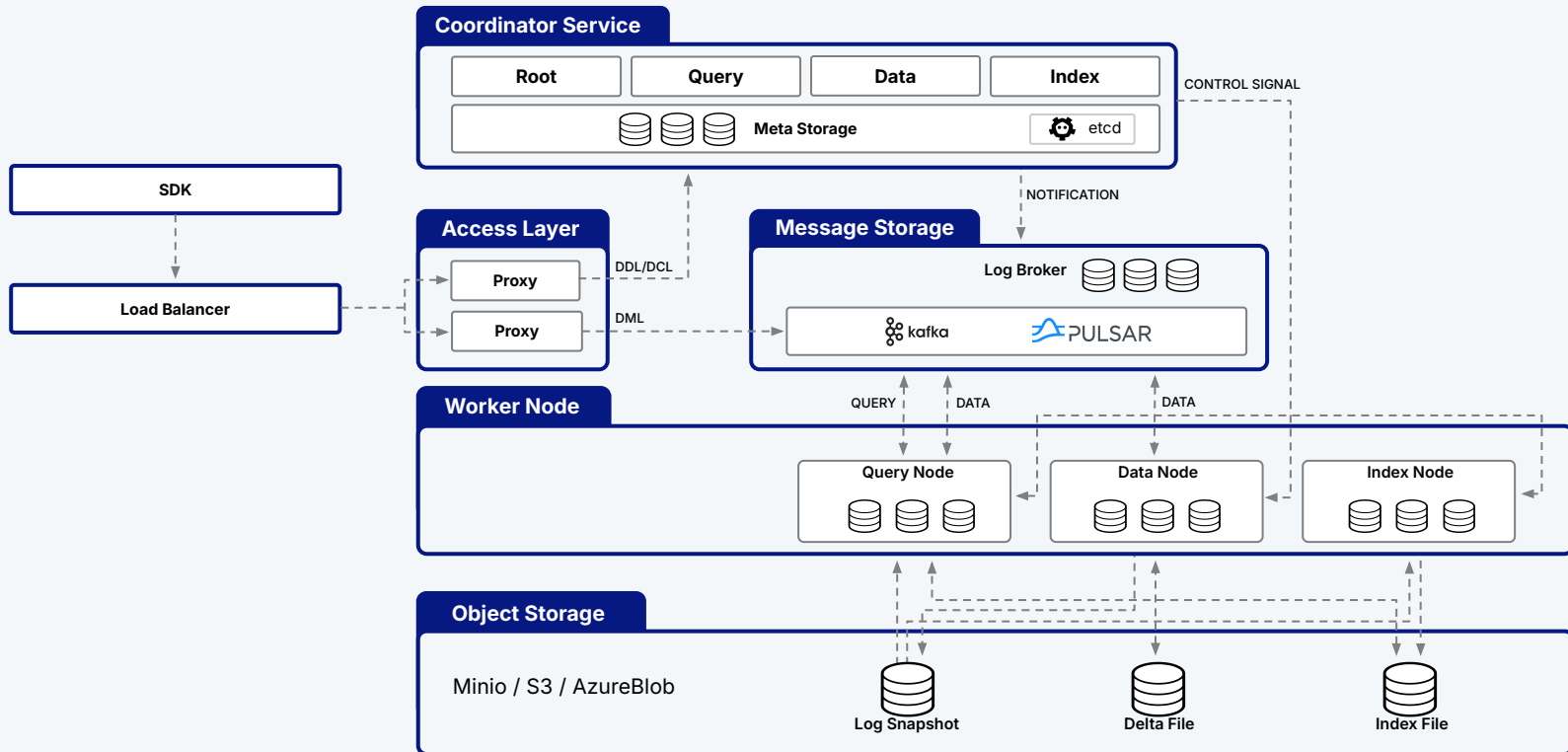
# A Never Ending Journey



**Congratulations, you just  
built Milvus!**



# Fully Distributed Architecture





## Achieving Billion+ Scale vector Search with K8s

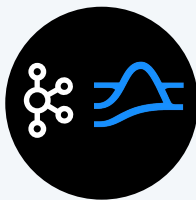
# Milvus 🤝 Open-Source



## 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

# Milvus Data Structures

## Shard

- Boost the ingestion rate

## Segment

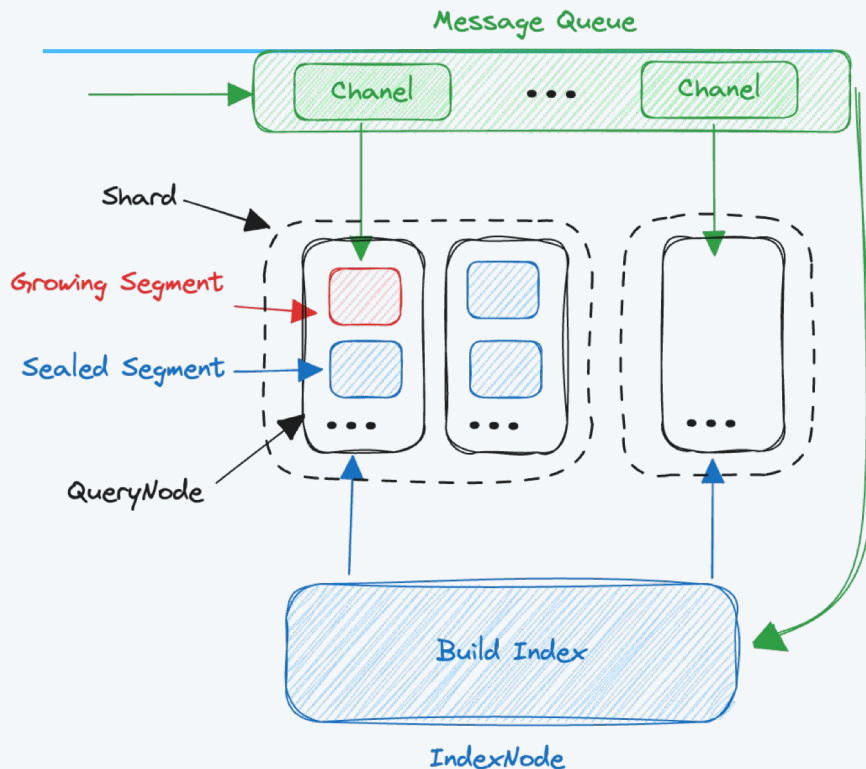
- A single unit of Data in Milvus.  
Segment < Partition < Collection

## Growing Segment

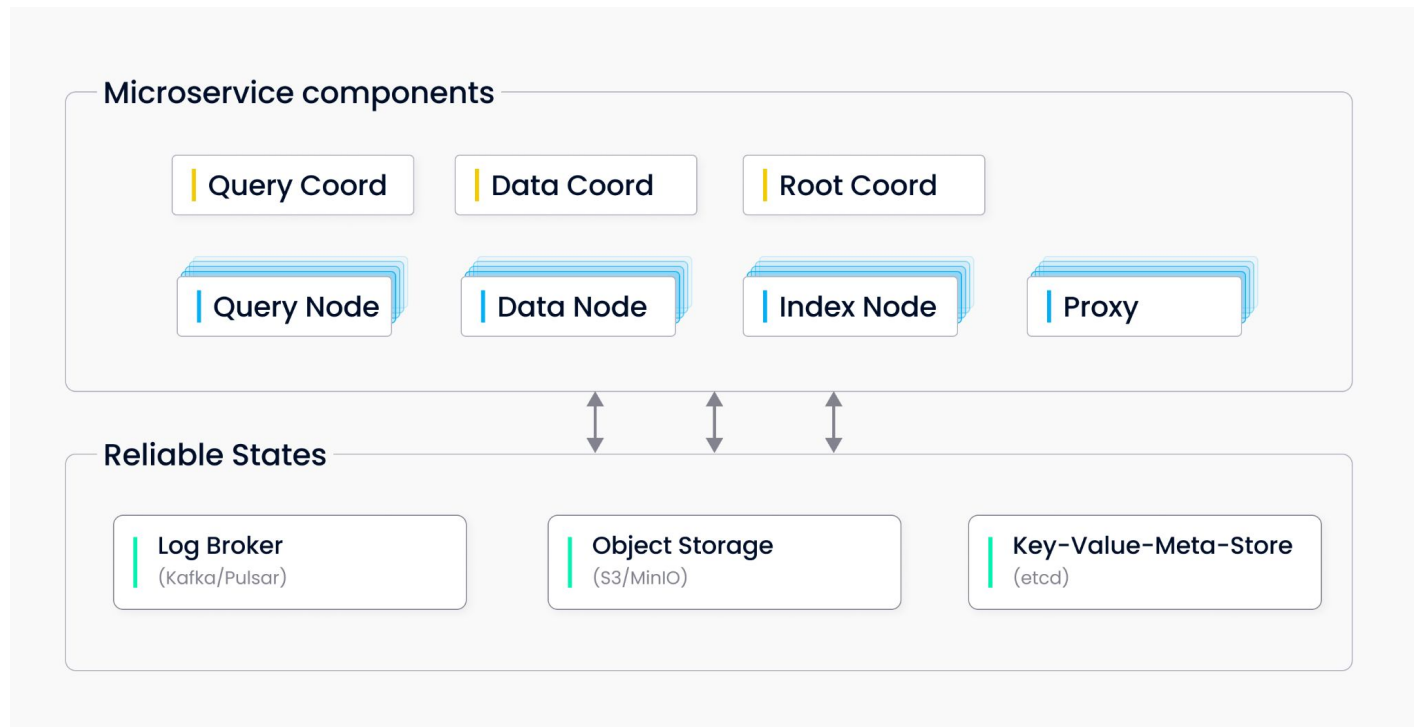
- Directly retrieves **data** from the **message queue** for rapid service. Utilizes a **brute-force index** and prioritizes **data freshness**.

## Sealed Segment

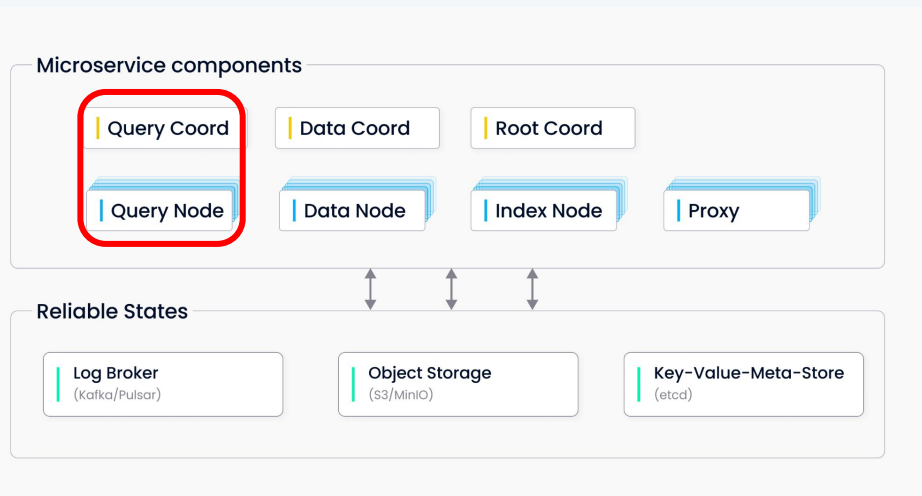
- An immutable segment uses **indexing** methods to guarantee **efficiency**.



# Distributed Architecture

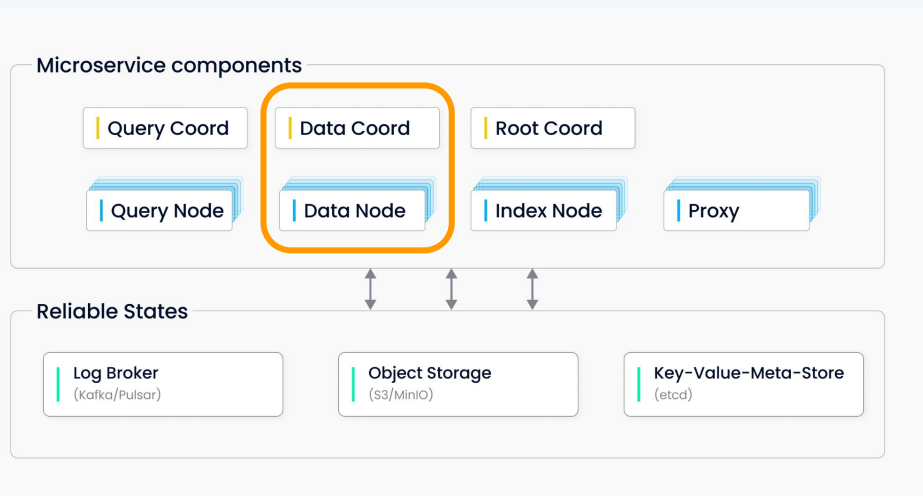


# Query Node: Serving Search Requests



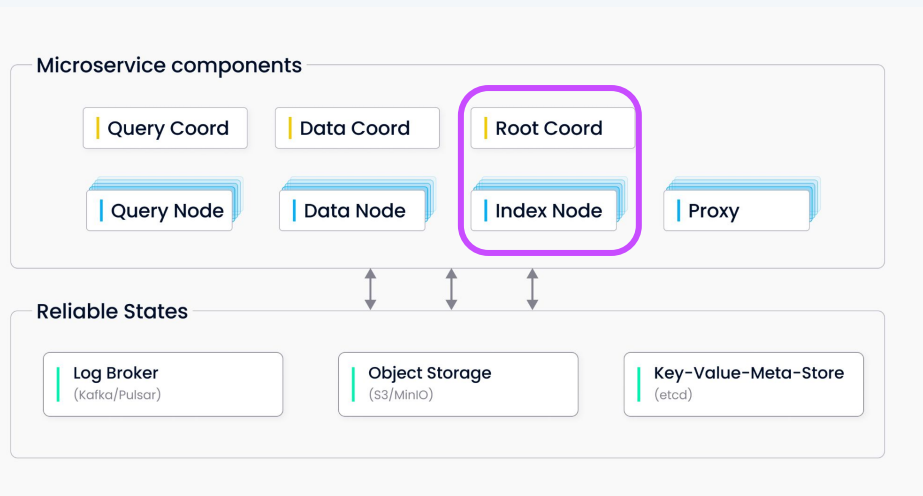
- Subscribe to the **log broker** for real-time querying
- Convert new data into **Growing Segments** - temporary in-memory structures for the latest information.
- Access **Sealed Segments** from object storage for comprehensive searches.
- Perform **hybrid searches** combining vector and scalar data for accurate retrieval.

# Data Node: Processing Data Updates



- Subscribe to the **log broker** for real-time updates.
- Process **mutation requests** for data changes or updates.
- Pack log data into **log snapshots** - compressed bundles of updates.
- Store log snapshots in **object storage** for persistence and scalability

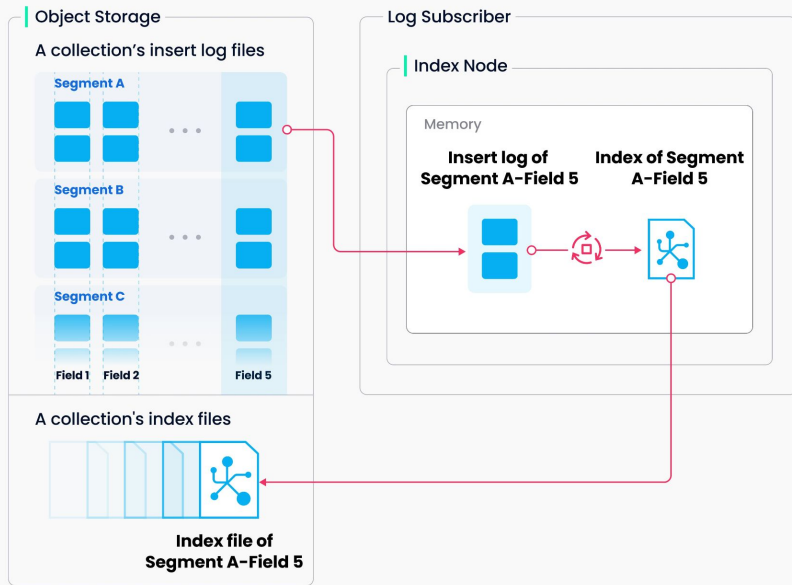
# Index Node: Building Search Indexes



- Build indexes on the data to facilitate faster search operations.
- Can be implemented using a **serverless framework** for cost-efficiency and scalability.



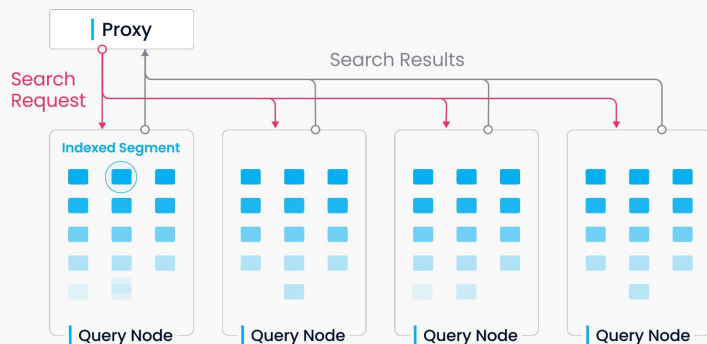
# Index Building



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

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

# Scalable Search



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



# Nearest Neighbors

## Brute-force

- Exhaustively search query vectors against index dataset

## IVF-Flat

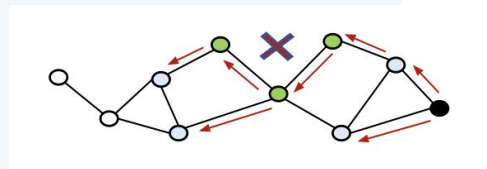
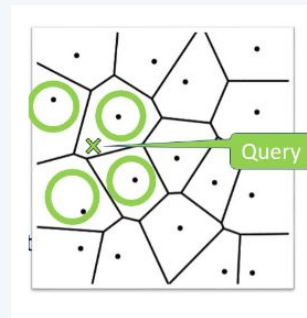
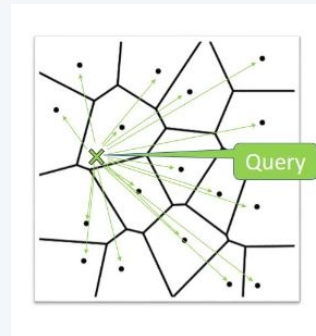
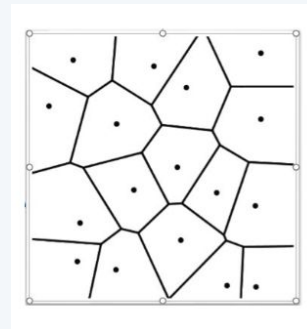
- Train clusters on Index dataset
- Partition Vectors by their closest clusters
- Search Smaller number closest lists.

## IVF-PQ

- Train clusters on Index dataset
- Partition and compress vectors by their closest clusters
- Search smaller number of closest lists

## CAGRA

- Build a graph from vector neighborhoods
- Reduce distances during search by traversing graph



# CAGRA

GPU accelerated Graph-Based ANN from NVIDIA

- Similar to HNSW for CPU
- Individual queries parallelized during Search
- Higher throughput than usual GPU Graph ANNs and lower latency than CPU Graph ANNs.

# Benchmarking setup

	Instance type	vCPU	Memory	GPU memory	Price(\$/h)	
T4	g4dn.2xlarge	8	32G	16G	0.752	1.58x expensive
A10G	g5.2xlarge	8	32G	24G	1.212	2.55x expensive
CPU	m6id.2xlarge	8	32G	–	0.4746	

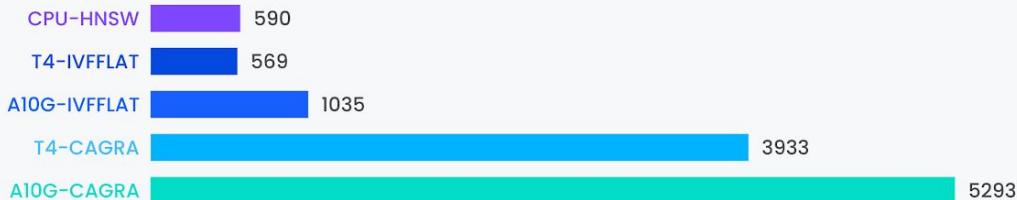
# Small Batch Performance

## Milvus-CAGRA vs Milvus GPU IVF vs Milvus-HNSW

(Search Batch Size: 1)

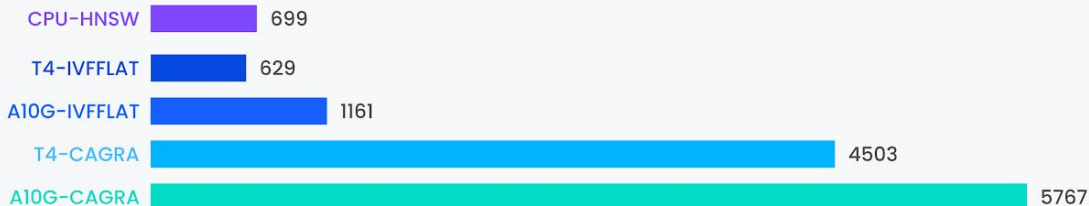
### Queries Per Second

Cohere:  
1M 768-dim vectors



### Queries Per Second

OpenAI:  
500K 1536-dim vectors





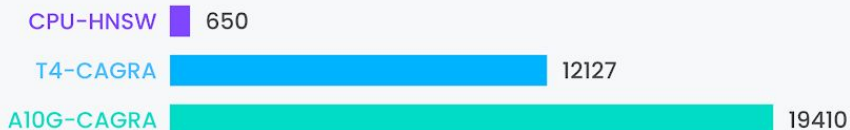
# Large Batch Performance

## Milvus-CAGRA vs Milvus-HNSW

(Cohere: 1M 768-dim vectors)

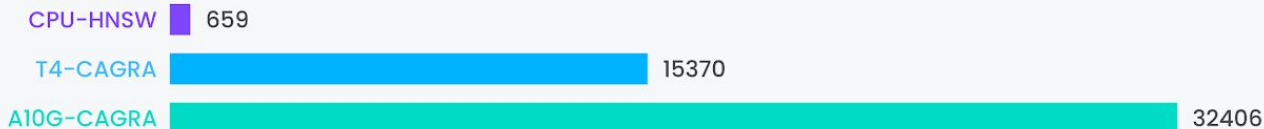
### Queries Per Second

Batch Size: 10



### Queries Per Second

Batch Size: 100



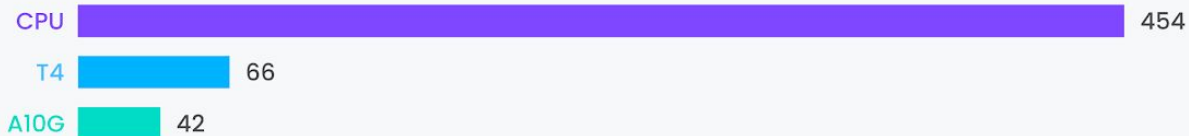


# Index Building (in seconds)

## Milvus-CAGRA vs Milvus-HNSW

### Index Building Time

Cohere:  
1M 768-dim vectors

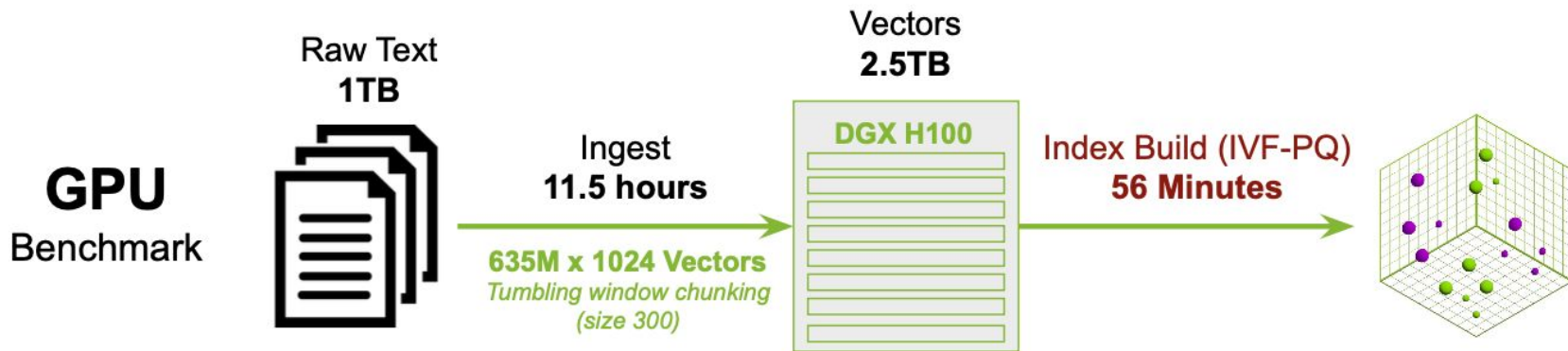


### Index Building Time

OpenAI:  
500K 1536-dim vectors



# Scale Indexing on DGX H100 (8x H100)



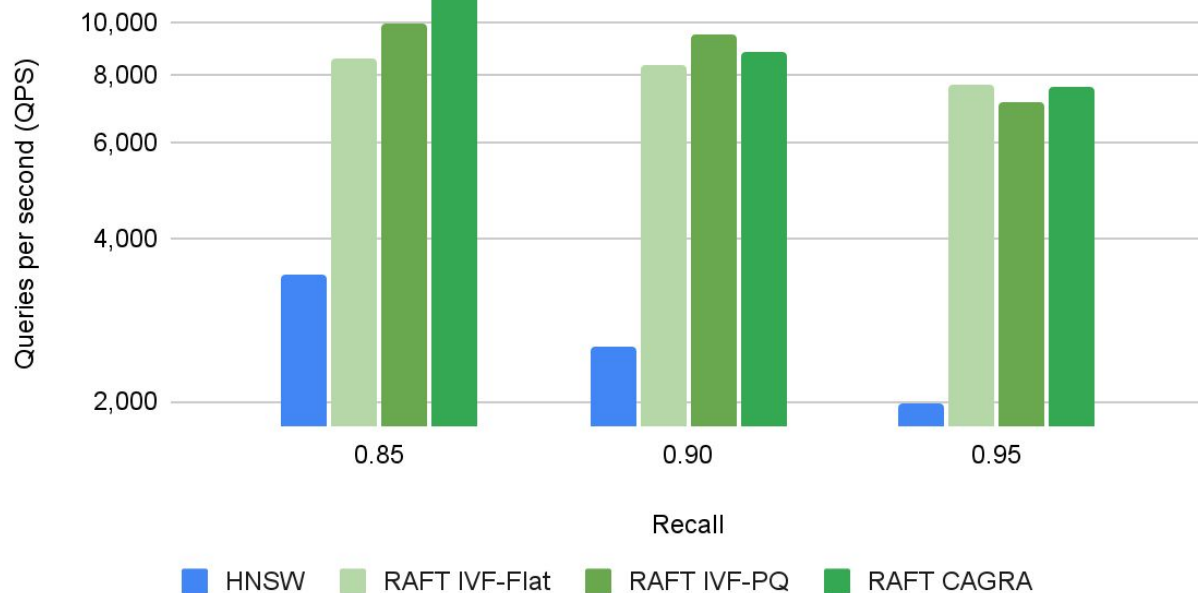
## CPU Estimate

GPU Index Build Time	56 Min
Number of GPUs	× 8
GPU / CPU Performance	× 20
Minutes per day	÷ 1440 Min/Day
<b>CPU Estimated Index Build</b>	<b>6.22 Days</b>

← *linear speedup*

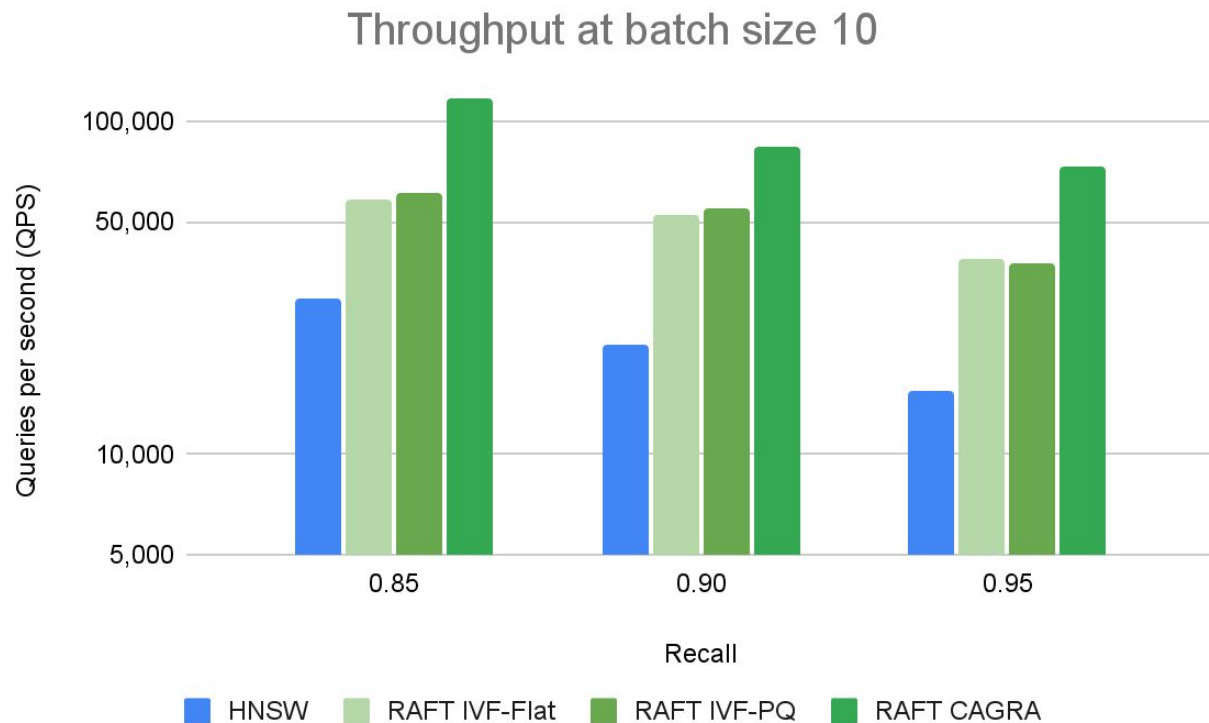
# Benchmarks (Batch size 1)

Throughput at batch size 1



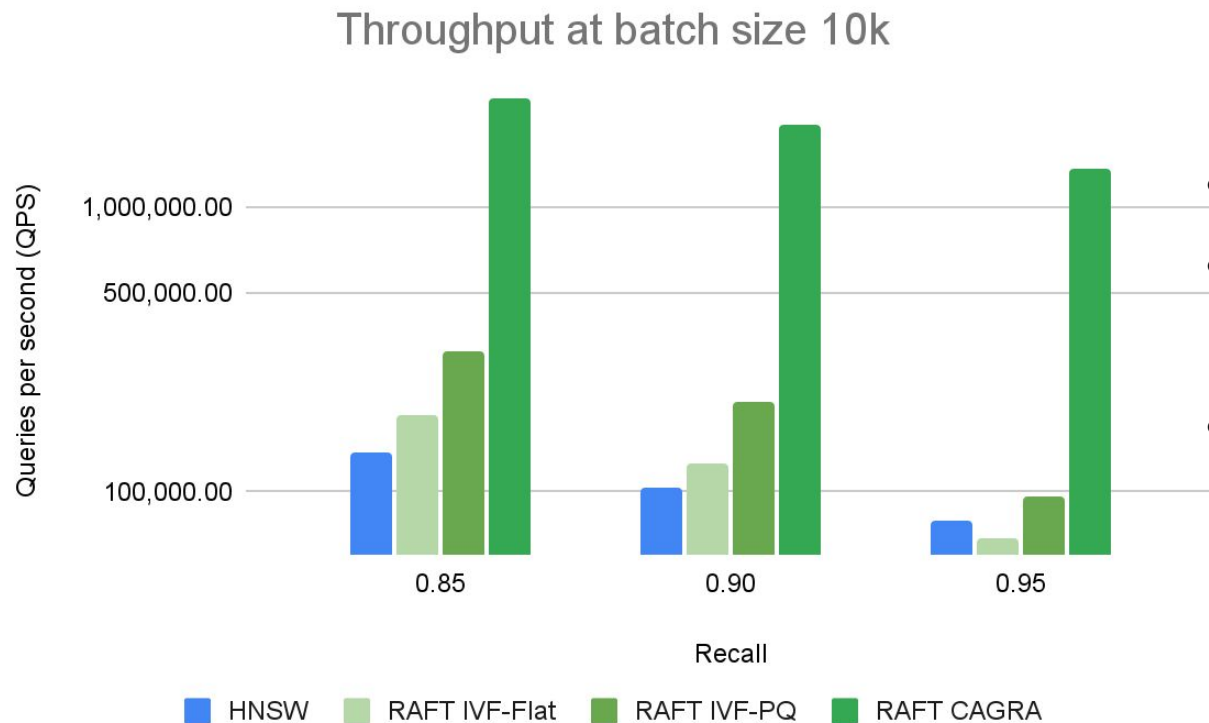
- **DEEP-100M** dataset
- **H100** GPU for RAFT indexes
- Intel **Xeon Platinum** 8480CL CPU for HNSW

# Benchmarks (Batch size 10)



- **DEEP-100M** dataset
- **H100** GPU for RAFT indexes
- Intel **Xeon Platinum** 8480CL CPU for HNSW

# Benchmarks (Batch size 10K)

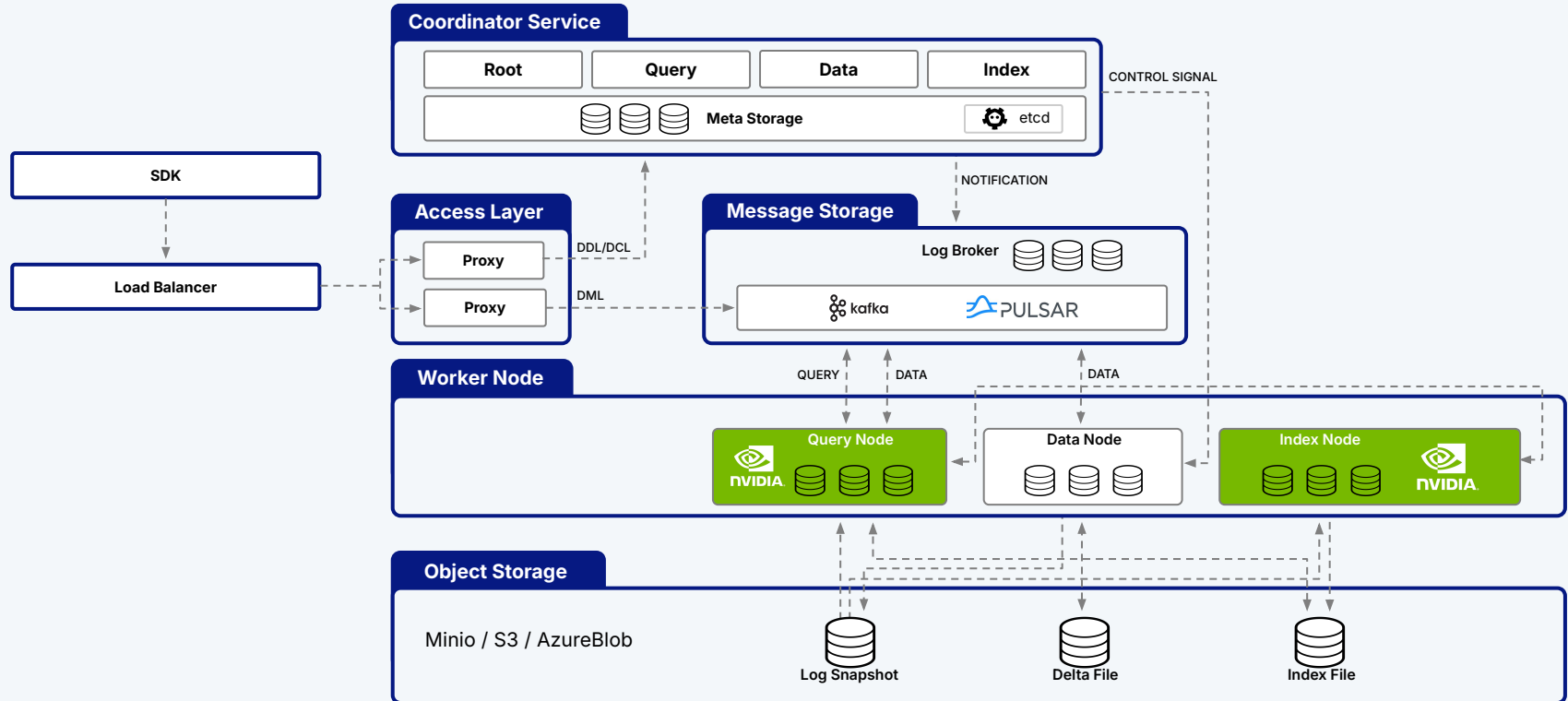


- **DEEP-100M** dataset
- **H100** GPU for RAFT indexes
- Intel **Xeon Platinum** 8480CL CPU for HNSW

# Tips for more performance

- Increase the number of **Data Nodes** ⇒ Improve **Streaming ingestion** performance
- Place **Index** and **Query nodes** on **separate GPUs**
- **Increase** max **segment** size can **improve search latency and throughput**

# Fully Distributed Architecture



# But at what Scale?

## 10B vectors

of 1536 dimensions  
in a single Milvus/Zilliz Cloud  
instance

## 100B vectors

in one of the largest deployment running  
on K8s.



# Speaker



## Stephen Batifol

Developer Advocate, Zilliz/ Milvus

✉ [stephen.batifol@zilliz.com](mailto:stephen.batifol@zilliz.com)

[in linkedin.com/in/stephen-batifol/](https://www.linkedin.com/in/stephen-batifol/)

✂ [@stephenbti](https://twitter.com/stephenbti)



# Thank you

---



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



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



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



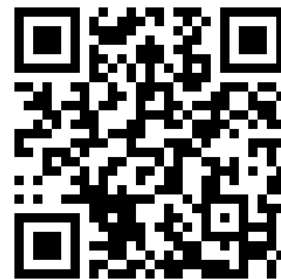
Notebooks



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



[in/stephen-batifol](https://www.linkedin.com/company/stephen-batifol/)



# Ready to scale

Write your code **once**, and run it **everywhere**, at **scale**!

- API and SDK are the same

## Milvus Lite

- Ideal for **prototyping**, **small scale** experiments.
- Easy to set up and use, `pip install pymilvus`
- Scale to **~1M** vectors

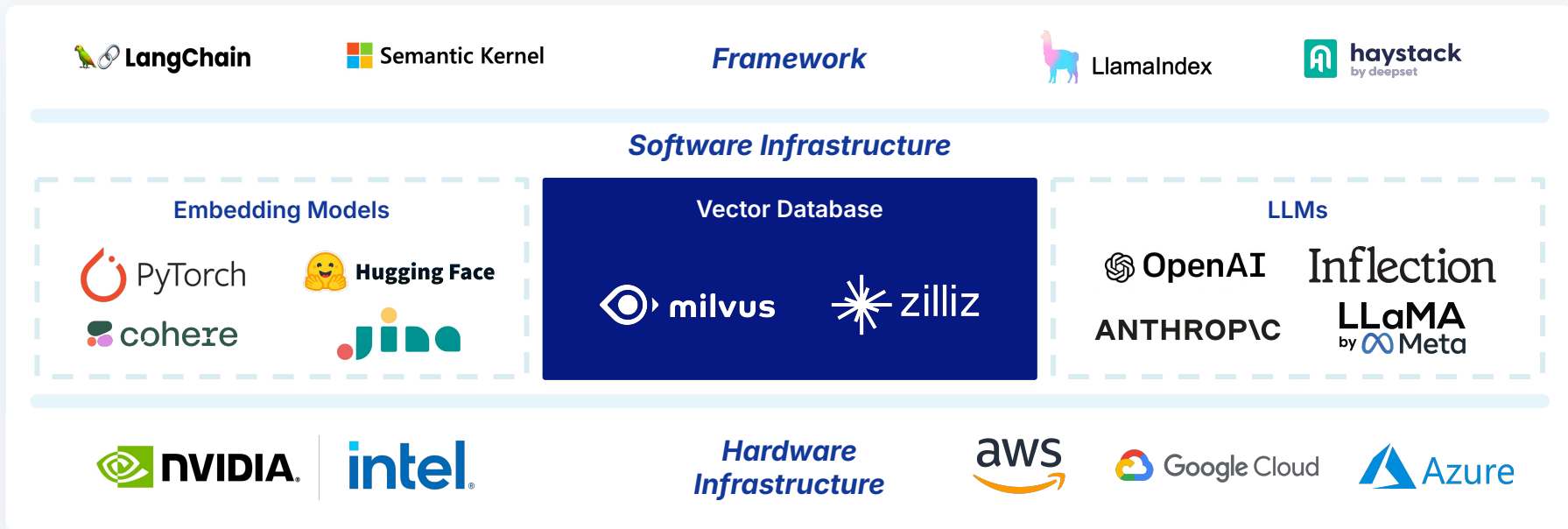
## Milvus Standalone

- Single-Node Deployment
- Bundled in a **single Docker** Image
- Supports Primary/Secondary
- Scale up to **100M vectors**

## Milvus Distributed

- Run on **K8s**
- **Load balancer** and **Multi-Node** Management
- **Scaling** of each component **independently**
- Scale to **100B** vectors

# Well-connected in LLM infrastructure to enable RAG use cases



# Thank you

---



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



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



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



Notebooks



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



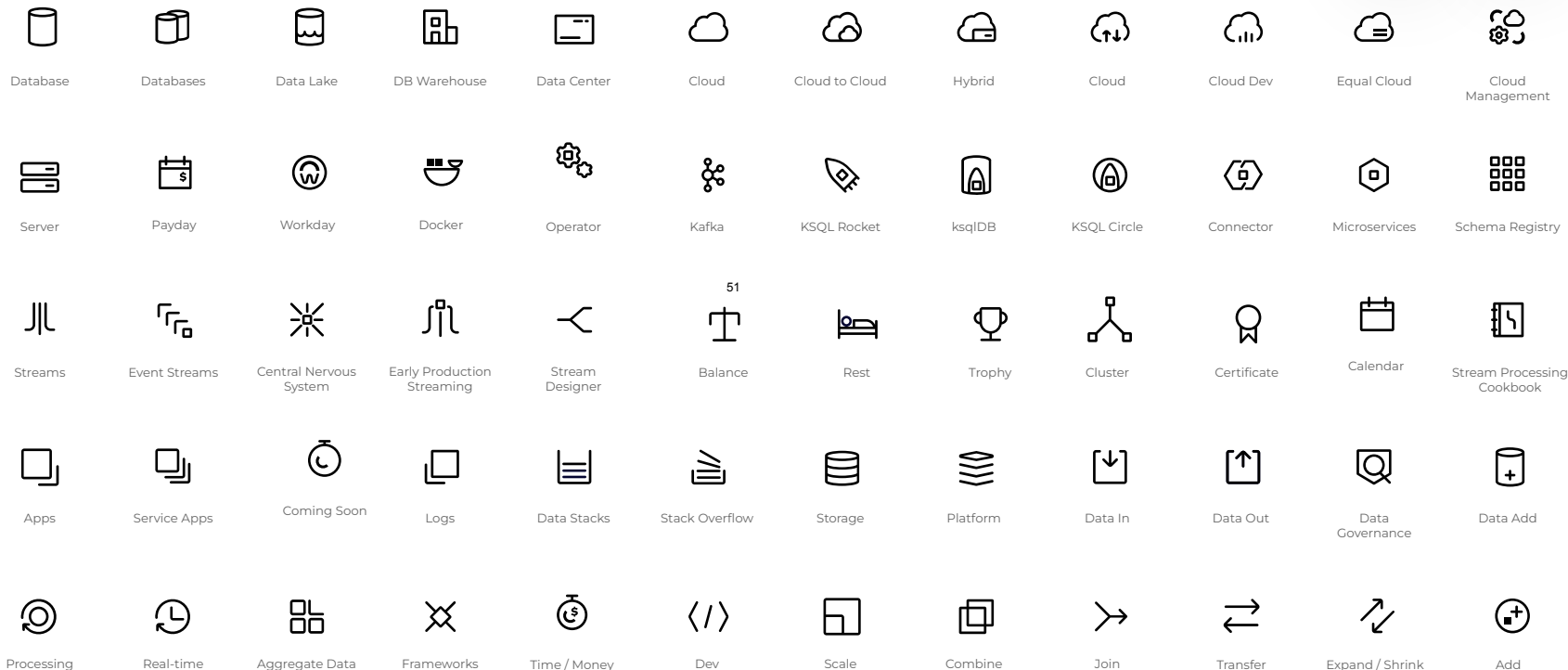
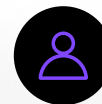
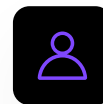
[/in/stephen-batifol](https://in.linkedin.com/in/stephen-batifol)





# Icons

ICON STYLE:

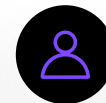
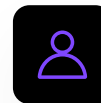


*Current*

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



Eye

*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>