



Zilliz Cloud



\$whoami

- Solutions Architect with Zilliz
- Previously
 - SA with Databricks, Confluent
 - Data scientist @ GIC
 - Data and AI practitioner for ~10 years



Mission:

Helping organizations make sense of unstructured data.



2017
Founded



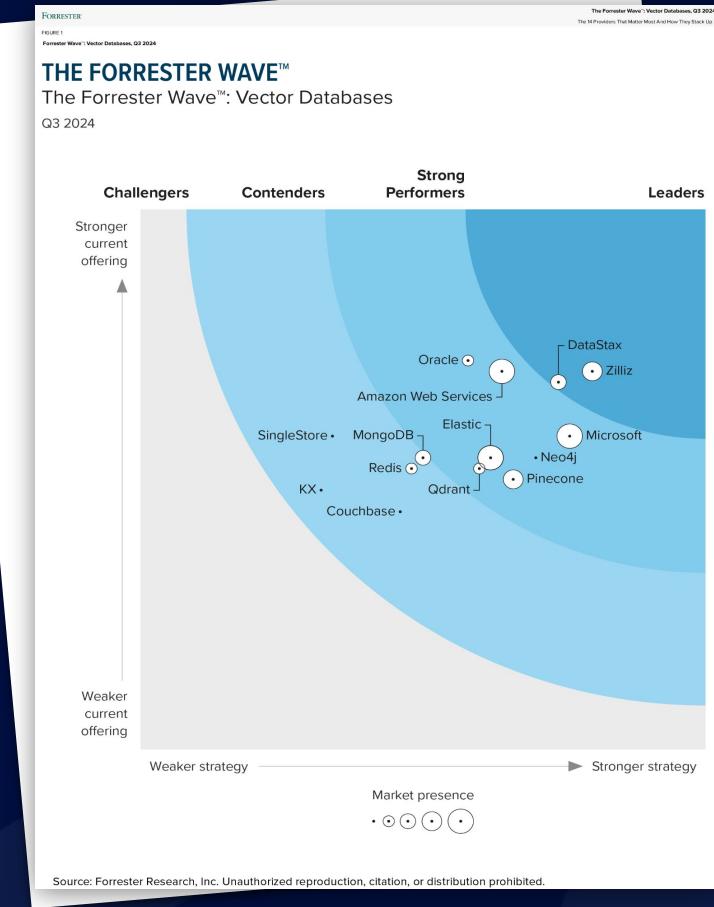
140+
Employees



\$113M
Raised



Redwood City, CA
Headquarters



The top industry analysts endorse our leading position

Specifically, we perform top at the following core vector database competency categories

- **Vector Index** - the only vendor in the market leverage multiple types of index and benefit from the unique strengths of each
- **Performance** - the only vendor provides advanced performance acceleration through various hardware optimization, query tuning
- **Scalability** - the only vendor offer compute storage separation and auto tiered storage

Milvus: The most widely-adopted vector database

Milvus is an **Open-Source Vector Database** to **store, index, manage, and use** the massive number of **embedding vectors** generated by deep neural networks and LLMs.



283+
contributors

30K+
stars

67M+
docker pulls

2.8K
+
forks

Built by database & AI experts

Milvus: A Purpose-Built Vector Data Management System
Jiaojiao Wu^{1*}, Xiaobin Yu, Ruihan Liang, Ling Li, Peng Xu, Qianqian Li, Xiangyu Wang,
Xianglong Guo, Shengming Xu, Jun Li, Tianqi Chen, Tianqi Zhou, Bojun Li, Jingbo Long,
Yadong Cai, Zhengming Li, Zhiqiang Zhang, Yihua Mo, Jun Gu, Gary Li, Yi Wei, Charles Xie
¹Zilliz & Peking University
^{*}Corresponding author

ABSTRACT

From the birth of vector search to the rise of large-scale distributed systems, Milvus has been at the forefront of innovation and application. This paper introduces Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

INTRODUCTION

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

CCS CONCEPTS
Information systems → Database management systems → Data access methods.

arXiv:2206.13843v1 [cs.DB] 28 Jun 2022

Manu: A Cloud Native Vector Database Management System
Bingtong Guo^{1*}, Xiaofan Luan², Long Xiong¹, Xiao Yan², Xiaoxiang Yi², Jiguo Lin¹,
Quanya Cheng², Weidu Xie², Jieren Luo², Frank Xu², Jianhua Cao², Yanhang Qiao², Ting Wang²,
Bo Tang², Charles Xie²

ACM Reference Format

Bingtong Guo, Xiaofan Luan, Long Xiong, Xiao Yan, Xiaoxiang Yi, Jiguo Lin, Quanya Cheng, Weidu Xie, Jieren Luo, Frank Xu, Jianhua Cao, Yanhang Qiao, Ting Wang, Bo Tang, and Charles Xie. 2022. Manu: A Cloud Native Vector Database Management System. In *Proceedings of the 2022 VLDB Endowment (VLDB '22)*, September 26–October 1, 2022, Virtual Event, China. ACM, New York, NY, USA, 16 pages. <https://doi.org/10.1145/3547921.3548431>.

KEYWORD

Vector database, High-dimensional, machine learning, Deep learning, Cloud native, Vector search, Vector data management system.

ABSTRACT

From the birth of vector search to the rise of large-scale distributed systems, Milvus has been at the forefront of innovation and application. This paper introduces Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

INTRODUCTION

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

1. INTRODUCTION

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

2. RELATED WORK

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

3. DESIGN

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

4. IMPLEMENTATION

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

5. CONCLUSION

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

6. ACKNOWLEDGMENTS

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

7. REFERENCES

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

8. FUTURE WORK

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

9. CONCLUSION

As Zilliz, we have experienced a growing need from various companies to manage large-scale high-dimensional vector data (long vectors) for their products. For example, our recommendation system needs to store millions of vectors for each user, and each user has tens of millions of vectors to achieve scalability and availability. In this paper, we introduce Milvus, a purpose-built vector data management system. It is designed to handle billions of vectors for data analysis, e.g., product recommendation, knowledge graph, and search. Milvus is built on top of a distributed storage system and two backends [1]. They store persistent data in memory and disk respectively. Milvus provides a unified interface for users to interact with both backends. It supports various data types, including binary vectors, sparse vectors, and structured data. Milvus also provides a powerful search engine that can meet the requirements of different applications. Milvus has been widely adopted in various industries, such as e-commerce, finance, and healthcare. In this paper, we introduce Milvus's architecture, design, and implementation. We also discuss how Milvus can be used to solve real-world problems.

Zilliz was built by a top-tier team of **algorithm and database engineers** with a strong pedigree in developing **high-performance, scalable, and highly available** distributed systems, uniquely tailored for **vector search**.



This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

© 2024 Zilliz. All rights reserved. This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

http://creativecommons.org/licenses/by-nc-sa/4.0/

6 | © Copyright 2024 Zilliz



Our Customers



AT&T



BOSCH

Chegg



CISION®

COMPASS

Deloitte.

ebay

FARFETCH

Grab



Inflection

intuit

Microsoft

new relic

NVIDIA®

OMERS

OII Otter.ai

PayPal

paloalto
NETWORKS

POSHMARK

RABLOX

salesforce

Shell

shutterstock

T

TREND
MICRO

Walmart

ZipRecruiter

zomato

zilliz

What is a vector database?

What is a vector database?

What are vector embeddings?

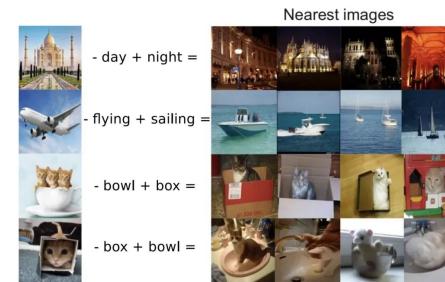
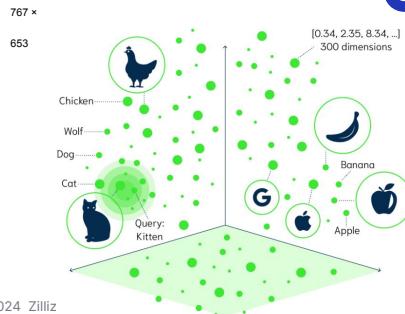
Numerical representations of unstructured data (text, images, audio, video, etc) that captures their semantics and relationships.

Vector database supports efficient storage, search and management of vector embeddings

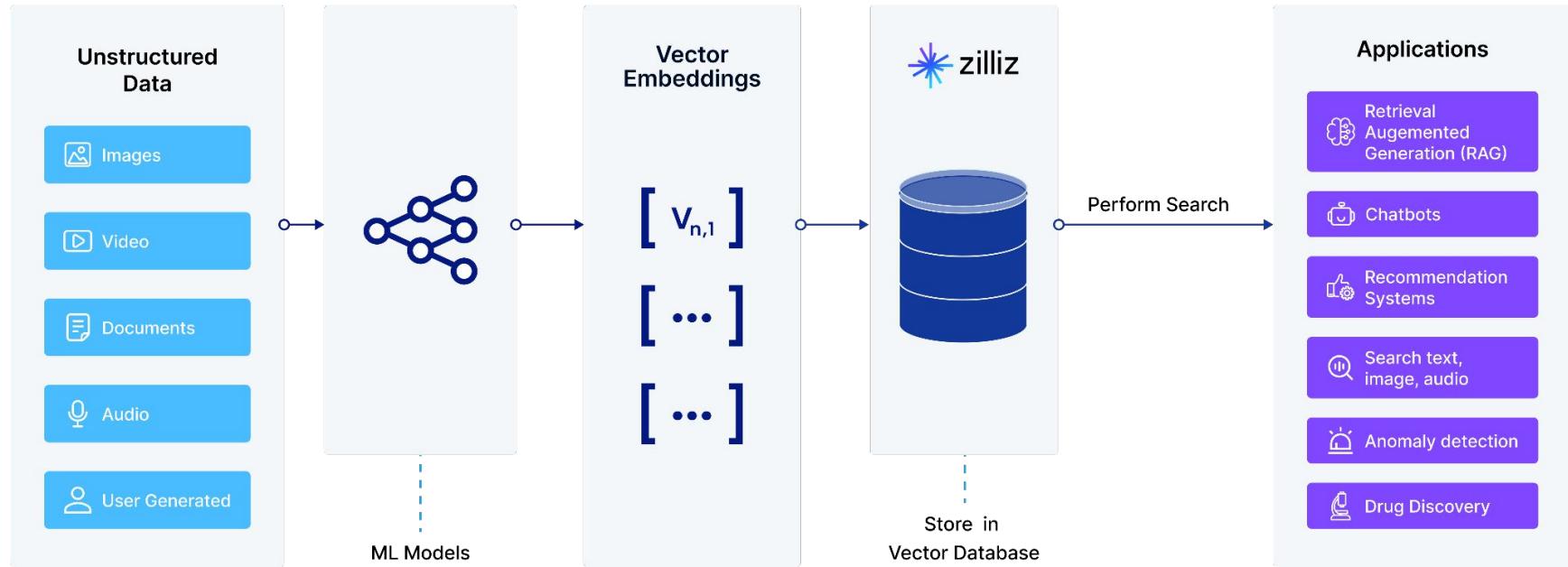
What is a vector database?

How to use vector db

extract → insert → index → search → operations



A New tool emerged. The Vector Database

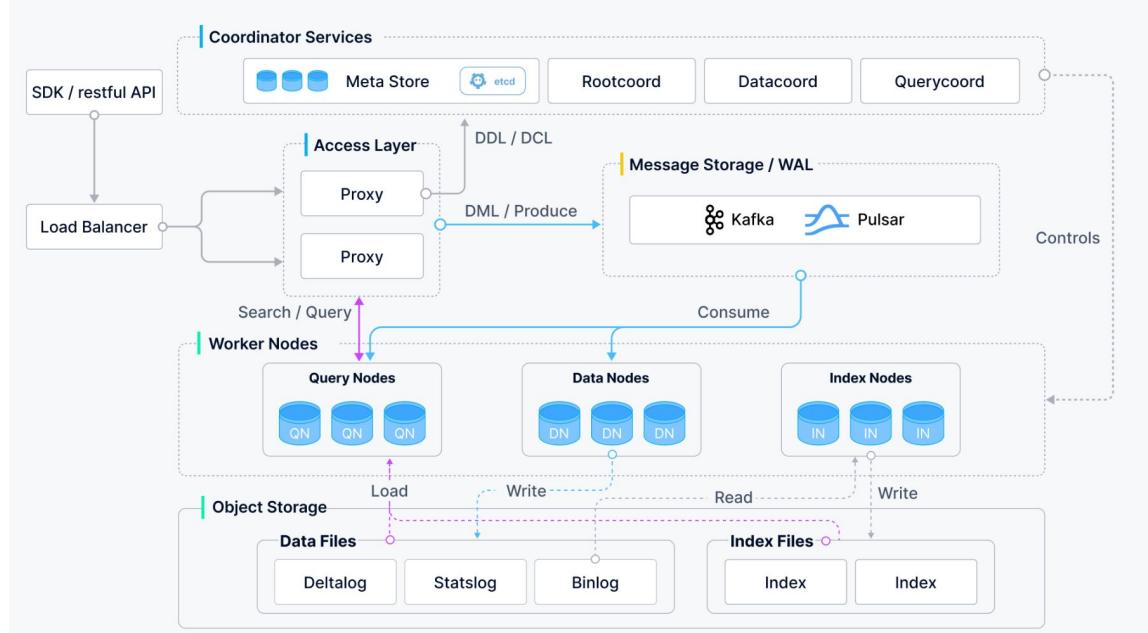


Milvus Architecture

Design Principles

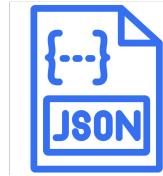
- Separation of storage and compute
- Fully depend on mature storage systems
- Microservice - scale by functionality
- Separate streaming and historical data
- Pluggable engine, storage and index
- Log as data

Fully distributed, designed for scalability

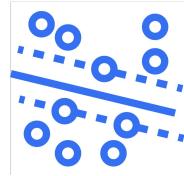


Milvus v2.4.x architecture overview

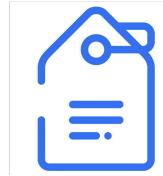
Rich functionality



Dynamic Schema



Float, Binary, &
Sparse Vector



Tag+Vector
Optimized Filtering



Hybrid Search
Dense & Sparse



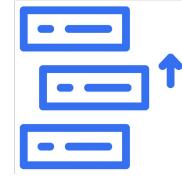
RBAC, TLS,
Encryption



Million+ level
tenant support



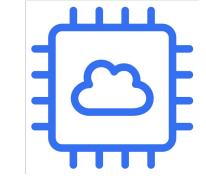
Disk Based
Index



Tiered Storage

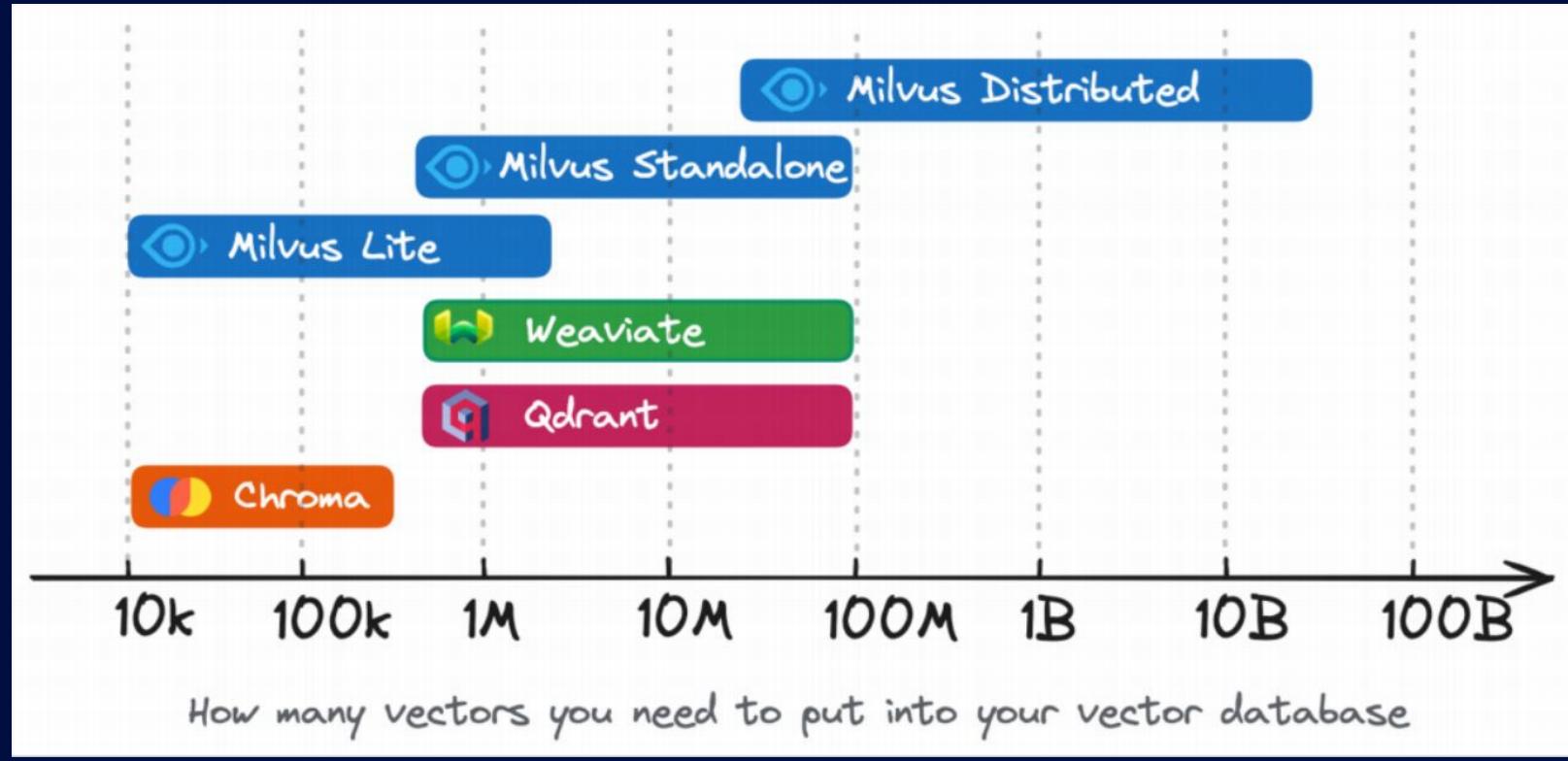


Bulk Import

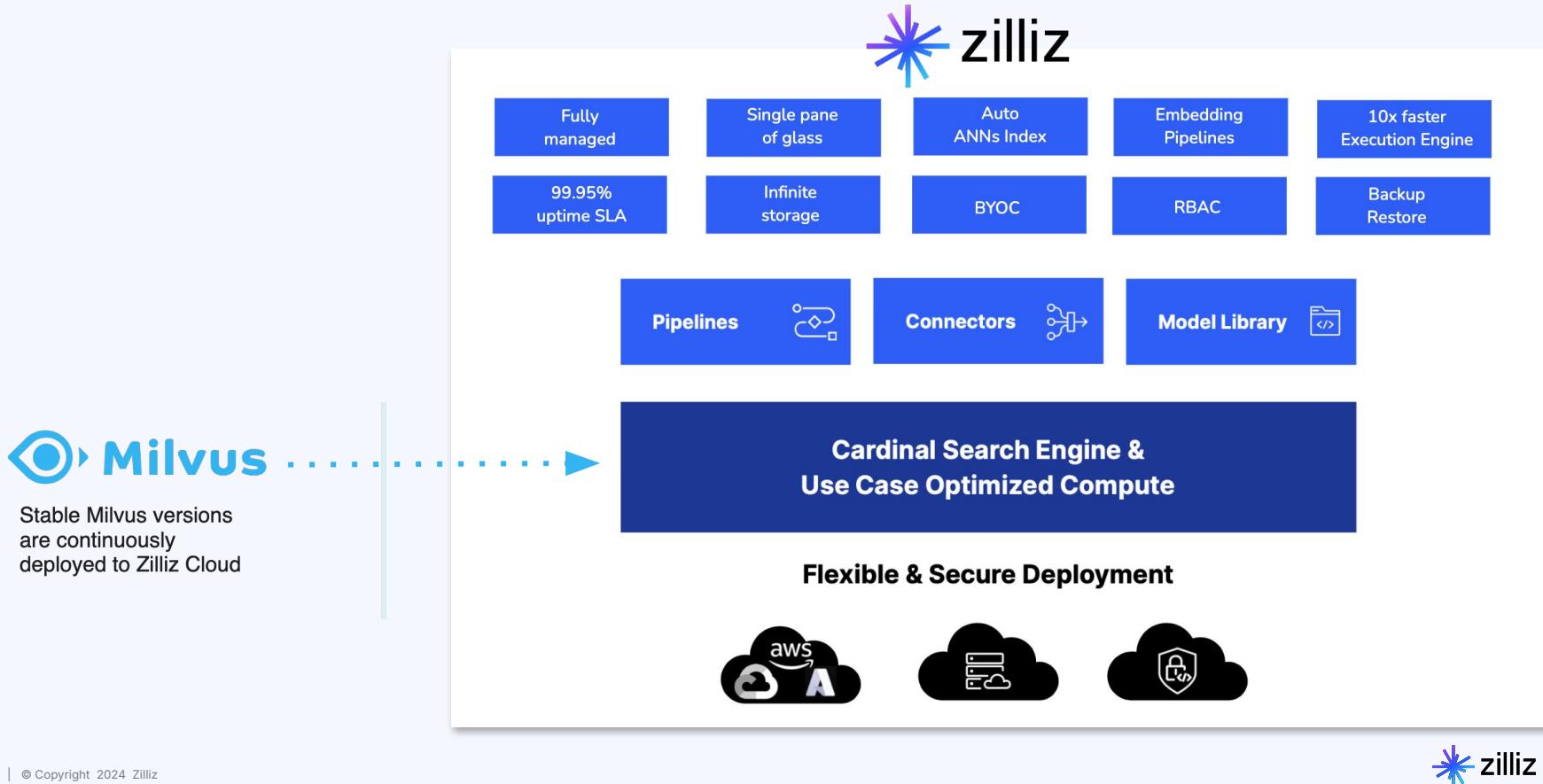


GPU, Intel & ARM
CPU support

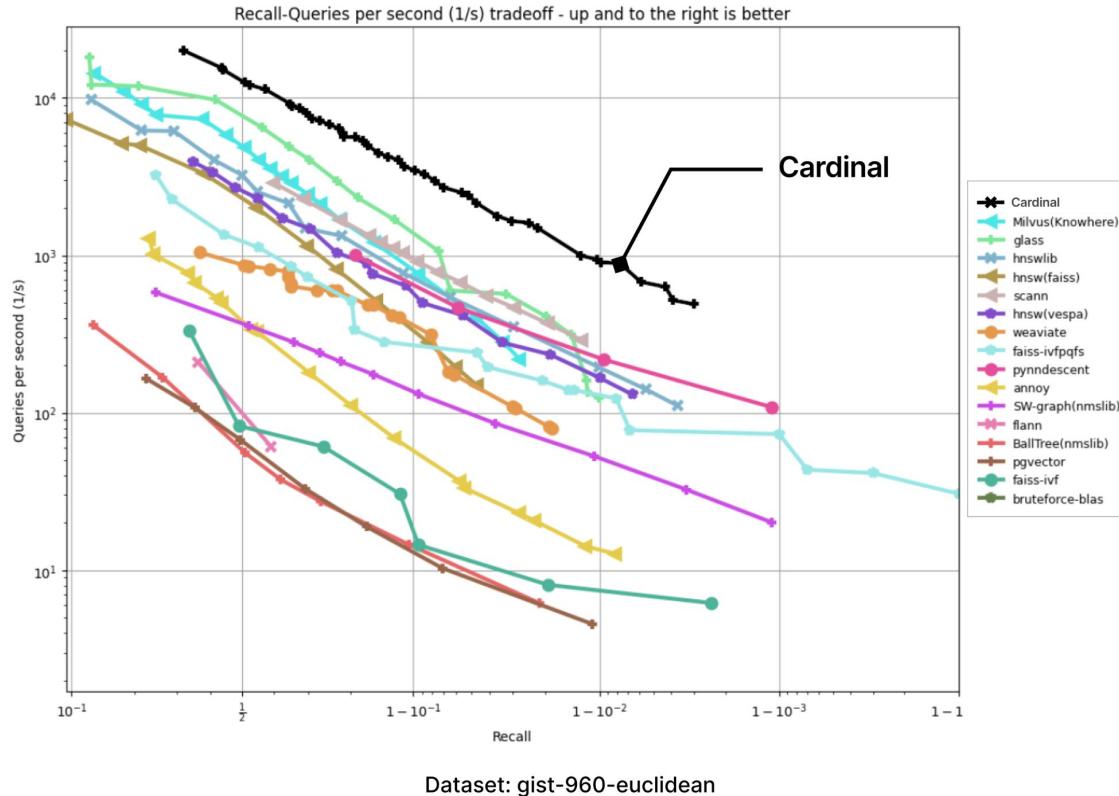
Milvus vs other open source vector databases



..but Zilliz Cloud is much more than OSS Milvus



ANN Benchmarks



Zilliz Offerings

SELF MANAGED SOFTWARE



Milvus

Most widely-adopted open source vector database



FULLY MANAGED SERVICE



Zilliz Cloud

AI Powered Search that is performant and scales



Google Cloud



Azure

BRING YOUR OWN CLOUD



Zilliz BYOC

For Private VPCs



Google Cloud



Azure
Coming Soon!



Set up Once: Common API across all products regardless of architecture

Getting started with Milvus Standalone

```
wget https://github.com/milvus-  
io/milvus/releases/download/v2.4.15/milvus-standalone-docker-  
compose.yml -O docker-compose.yml  
  
docker-compose up -d  
  
Creating milvus-etcd ... done  
Creating milvus-minio ... done  
Creating milvus-standalone ... done
```

Integrations

The Vector Database Integration Hub

Amplify your unstructured data potential and access a variety of vector database integration options for Zilliz Cloud.



Search



All Integrations

AI Models

Client Libraries

Data Sources

Observability

AI Frameworks

Orchestration

Other



CLIENT LIBRARIES

.NET

Use the .NET SDK with
Milvus or Zilliz Cloud



CLIENT LIBRARIES

Go

Use the Go SDK with
Milvus or Zilliz Cloud



CLIENT LIBRARIES

JavaScript

Use the JavaScript SDK
with Milvus or Zilliz Cloud



CLIENT LIBRARIES

Ruby SDK

Use the Ruby SDK with
Milvus or Zilliz Cloud



CLIENT LIBRARIES

Python

Use the Python SDK to
build Semantic Similarity
Search with Zilliz Cloud



Zilliz Use Cases

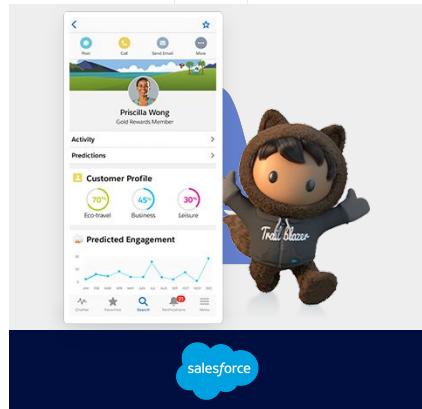
Industry leaders already use vector search in their apps

Use Case: Data Search

Vectors: 2 Billion

Req'ts: 200 ms, Cost mgmt

Index: DiskANN for cost savings



Use Case: Drug Discovery

Vectors: 12 Billion

Req'ts: High Recall

Index: BIN_FLAT



Use Case: Image Search

Vectors: 20 Billion

Req'ts: High Insertion, Cost

Index: Disk Based Index



Use Case: Recommender System

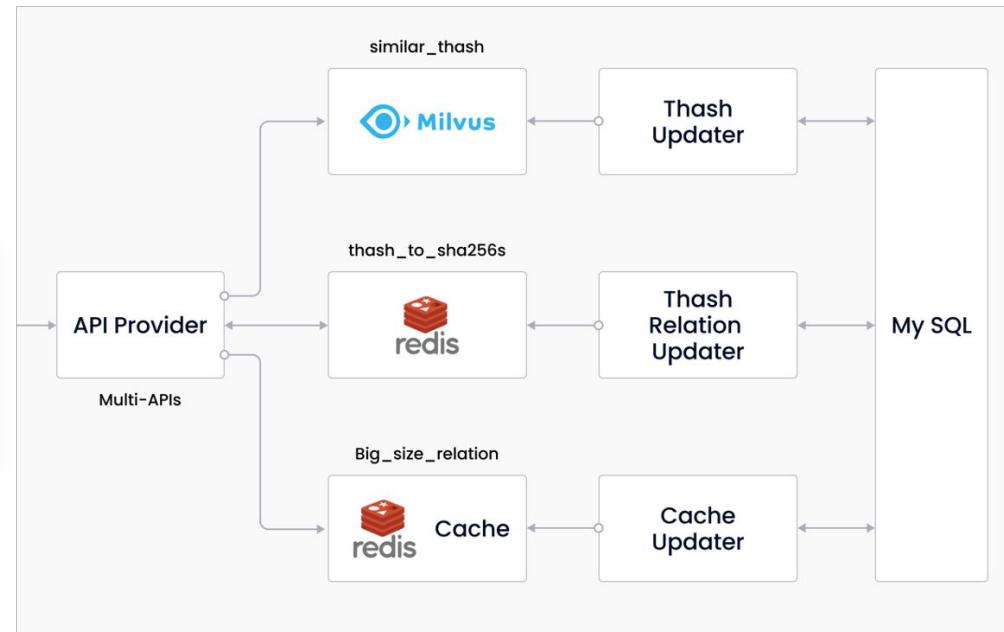
Vectors: 20 Billion

Req'ts: 5,000 QPS

Index: HNSW & CAGRA



Virus scan for Android app packages



e-Commerce: Similarity Search Engine



10x Smarter

search experience

Improved

user experience

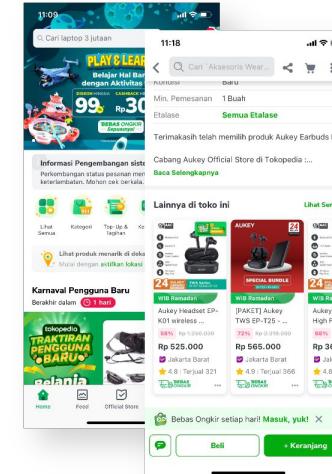
Enhanced

scalability and reliability

Moving from Elasticsearch to Milvus

Compared FAISS, Vearch and Milvus

- **Milvus** proved remarkably user-friendly. They found that you only need to pull its Docker image and adjust the parameters to suit your specific use cases.
- **Milvus** offers a broader range of supported indexes. Besides FAISS, HSNW, DISK_ANN, and ScaNN, there are 11 indexes to choose from.
- **Milvus** provides comprehensive documentation to aid users in their implementation.



e-Commerce: Similarity Search Engine

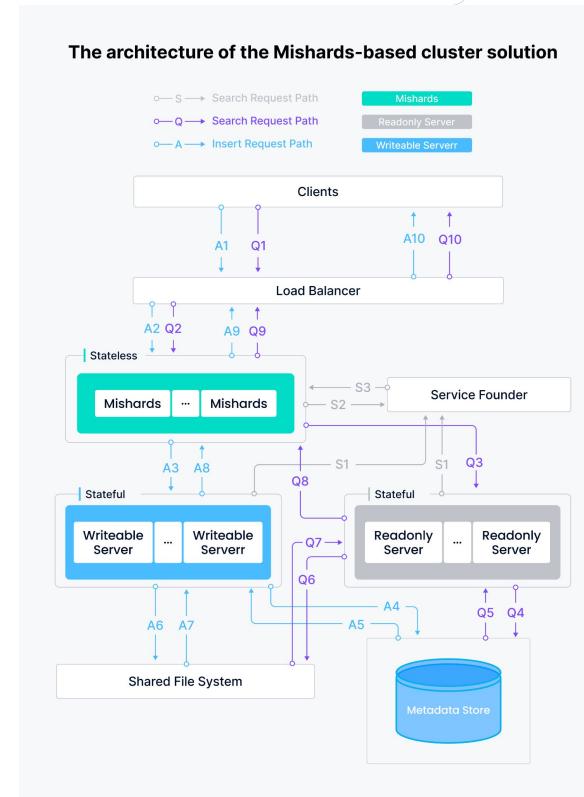


Ads service to match low-fill rate keywords with high-fill rate keywords

Compared Elasticsearch (Incumbent solution)

- 10x higher click-through-rate (CTR) and Conversion rate (CVR)
- Required HA implementation

Milvus offers two tools: **Mishards**, a cluster sharding middleware, and **Milvus-Helm** for streamlined configuration. At Tokopedia, they use Ansible playbooks for infrastructure setup, prompting them to create a playbook to orchestrate the infrastructure.



Useful resources

- [Forrester Wave Report: Vector Databases, Q3 2024](#)
- [Interactive Vector Database comparisons](#)
- [Zilliz Integration Hub](#)
- [Milvus demo notebooks](#)



@MILVUSAP



Discord



GitHub



Telegram
Milvus Asia Pacific



Backup slides

Multimedia Understanding



> 100M

embedding vector storage and searching

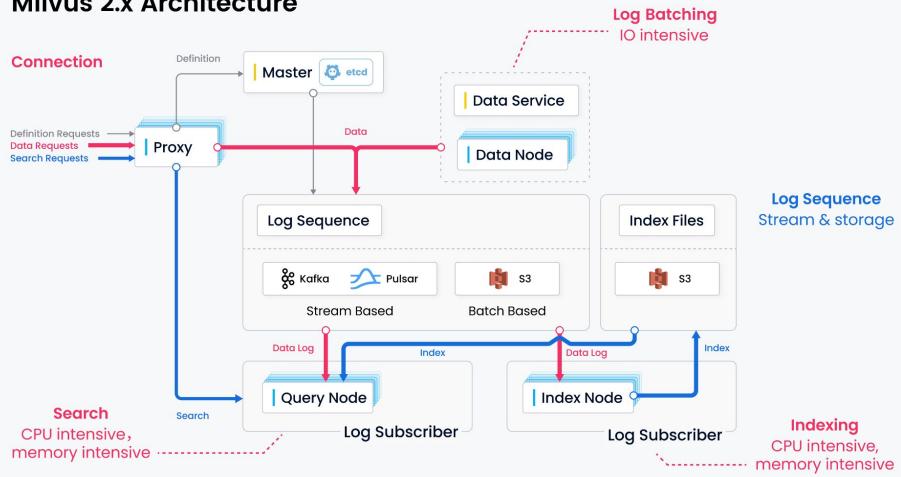
Seamlessly integration

with various internal systems and tech stacks

Enhanced real-time data retrieval

with reduced latency and increased system availability

Milvus 2.x Architecture



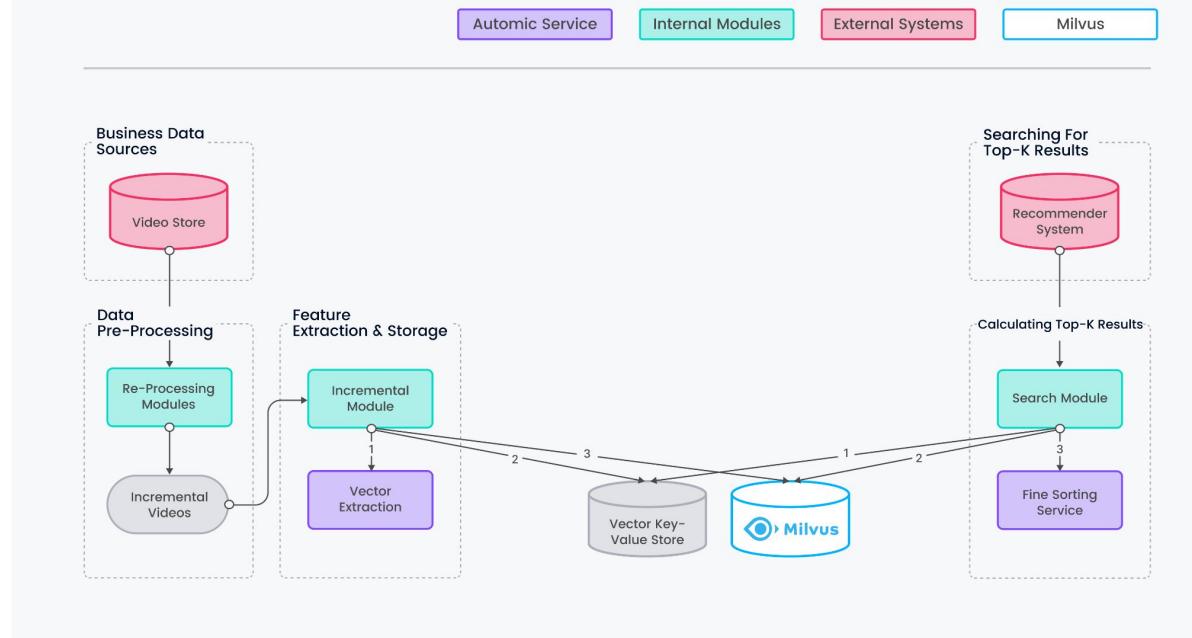
The Solution: Building a Vector Search Engine Using Milvus

After extensive research, Milvus emerged as the perfect fit. Milvus can handle billions of vectors and quickly scale out as data volume rises. Milvus' cloud-native architecture seamlessly integrated with Shopee's internal ecosystem, enabling the rapid setup of vector retrieval systems from scratch. Its feature-rich offerings, including distributed processing, GPU support, incremental updates, and scalar support, comprehensively addressed Shopee's multifaceted requirements. After careful consideration, the team selected Milvus as the foundation for their vector search engine to construct their vector search systems from scratch.

- Video Recall System
- Copyright Match System
- Video Deduplication System

Video Recall System: Improving Video Recommendation

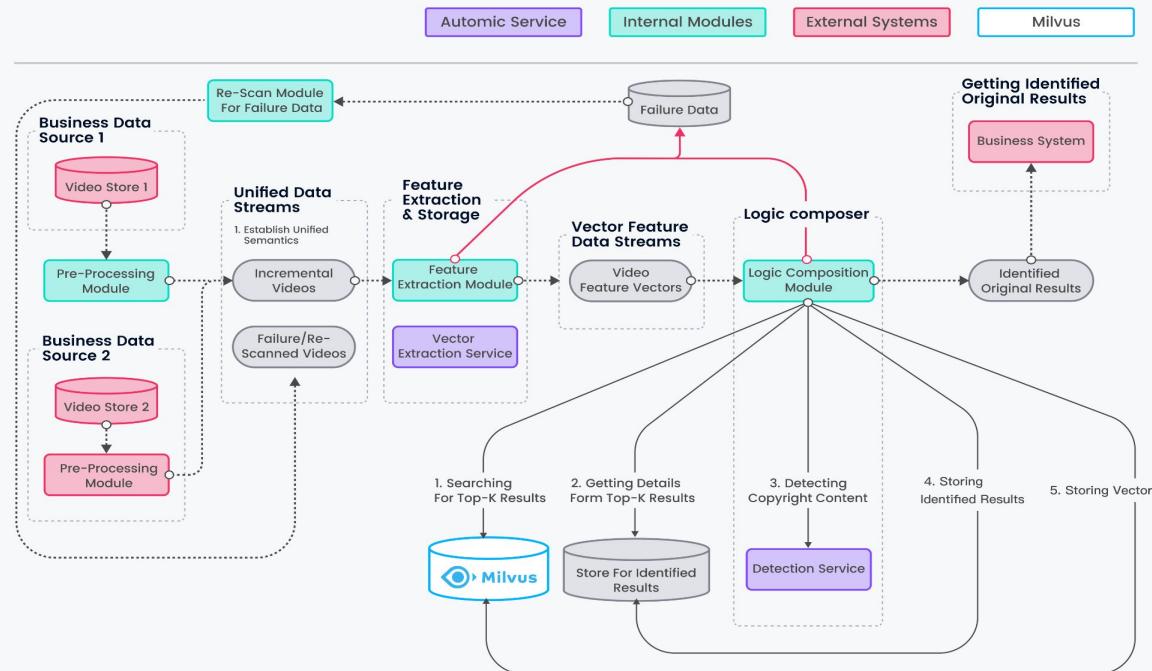
Video Recall System Architecture Using Milvus 2.x



Shopee's video recall system uses Milvus as a cornerstone in the process of recommending videos. When a user searches for a video, the business requests access to Milvus to retrieve the most similar Top-K candidates. These results undergo refinement through post-ranking algorithms before being returned to the user.

Copyright Match System

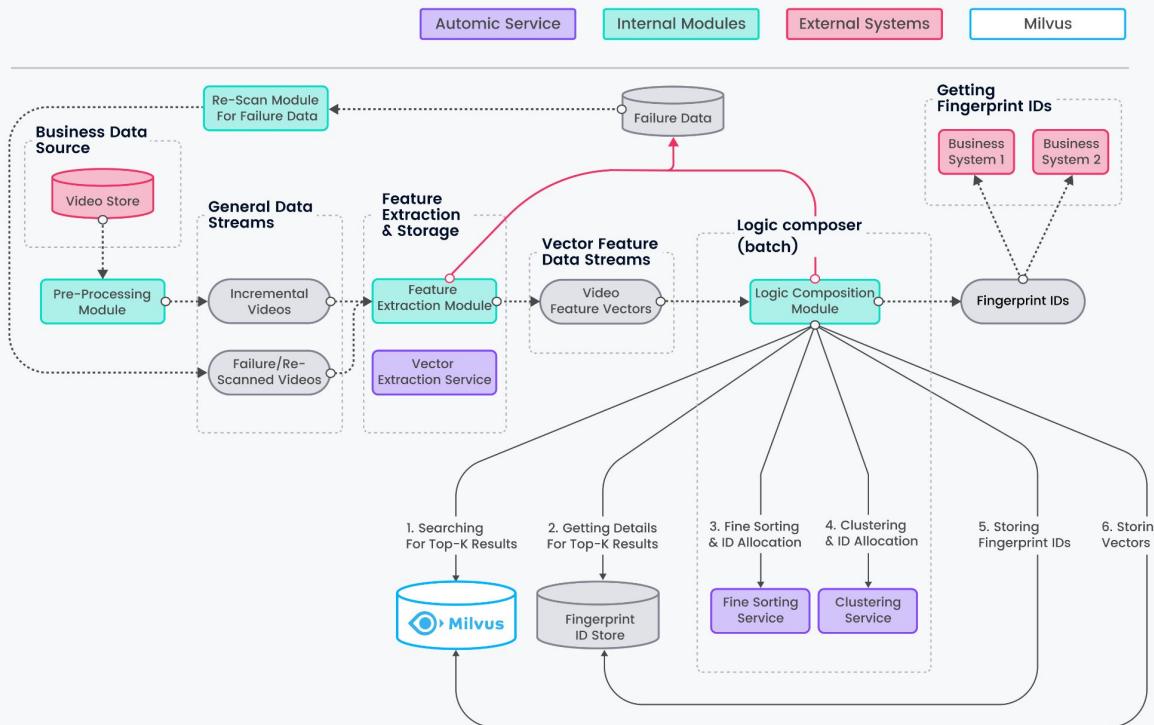
Copyright Match System Architecture



To maintain an excellent user experience and protect the copyrights of video creators, Shopee has implemented a copyright match system using Milvus. All released video features are transformed into vectors and stored in Milvus, and every newly uploaded video is matched with those held in Milvus by using similarity searches

Video Deduplication System

Video Deduplication System Architecture



The video deduplication system is designed to eliminate redundant content from Shopee's video platform. Like Shopee's copyright match system, the deduplication system uses Milvus to store embedding vectors transformed from video features. The system efficiently identifies and eliminates duplicate videos by searching for Top-K results in Milvus that are most similar to a specific part. Apart from the Top-K similarity search, the system involves other processing techniques such as batch data searching, post-ranking, clustering, and fingerprint assignment. In the end, Milvus stores all these results, providing valuable insights to various business units.

Compliance and Privacy



SOC 2 Type II

Zilliz Cloud's SOC2 Type II report offers important, third-party validation of our security practices, upheld consistently throughout the reporting period. This report provides a robust, evidence-based evaluation of our commitment to maintaining the highest security standards. By eliminating discrepancies, we aim to offer you enhanced confidence in the security posture of Zilliz Cloud.



ISO/IEC 27001

The ISO/IEC 27001 certification is an international benchmark for Information Security Management Systems (ISMS). Zilliz Cloud's adherence to this standard underscores a systematic approach to managing sensitive data, aligning with global best practices. By meeting this standard, we provide a stronger assurance that your information assets are well-protected.



GDPR Readiness

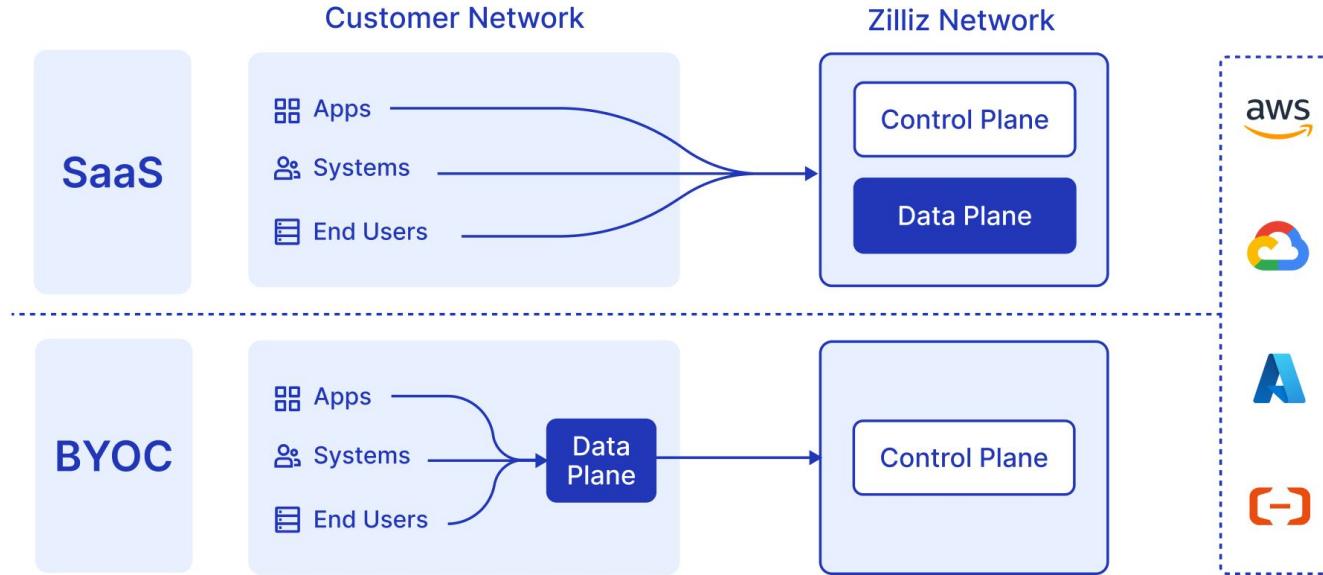
The General Data Protection Regulation (GDPR) sets standards for handling and protecting personal data from the European Economic Area (EEA) and ensures individuals' data rights. Zilliz is GDPR-ready and committed to supporting our customers' compliance efforts.



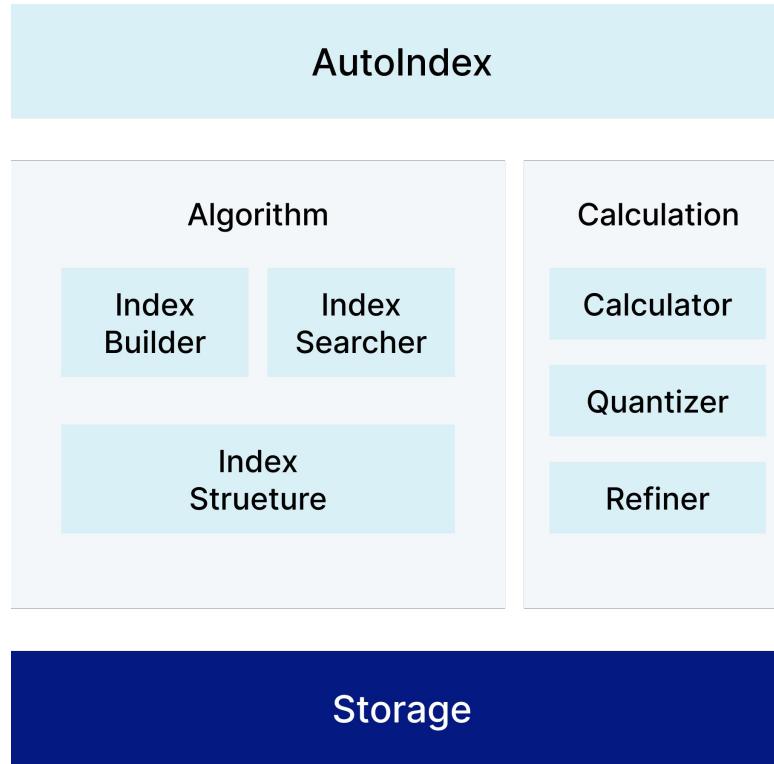
HIPAA Readiness

The Health Insurance Portability and Accountability Act of 1996 (HIPAA) regulates the protection of the privacy and security of health information. Zilliz is HIPAA-ready, enabling covered entities and their associates to use our secure cloud database to process, maintain, and store protected health information (PHI).

BYOC Versus SaaS

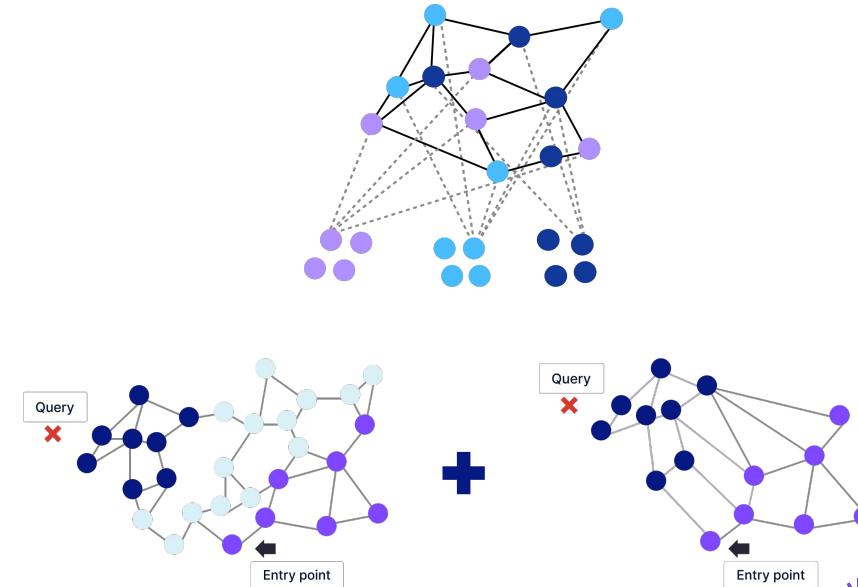


Cardinal Search Engine



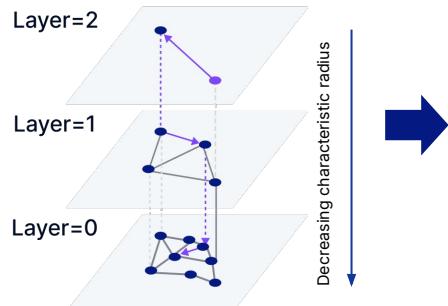
Graph-IVF hybrid solution helps:

Configurable storage media for all modules brings flexible performance & capacity tradeoff

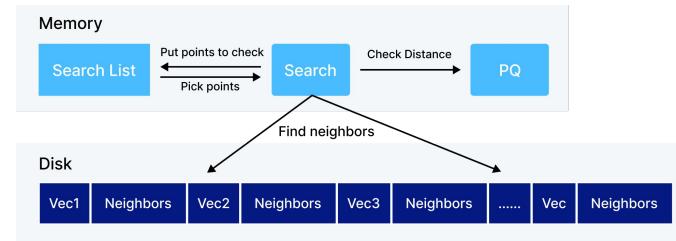


Intelligent Tiered Storage

Memory Based



Disk Based



Smart Tiers

