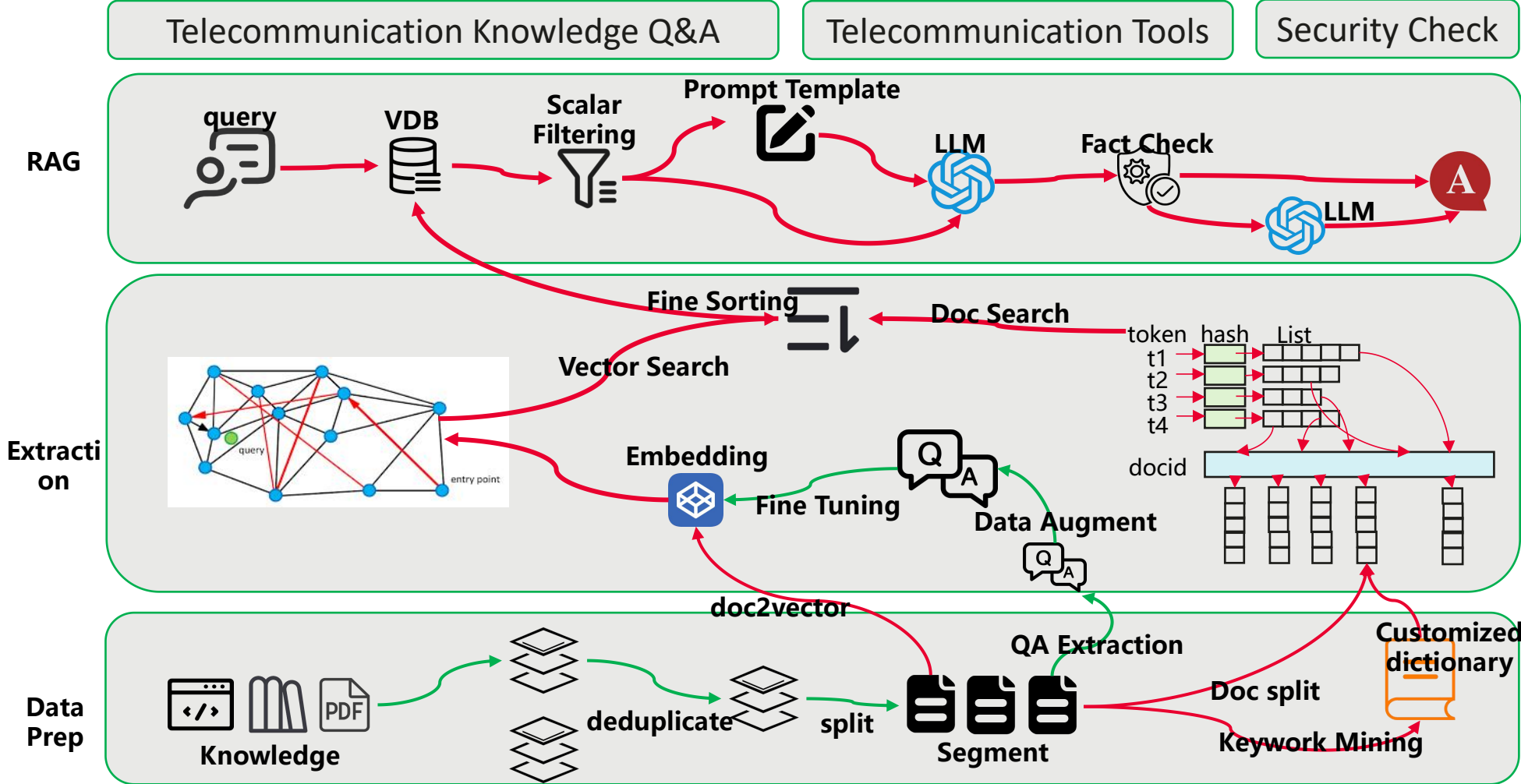# What Do We Really Need for Vector Database
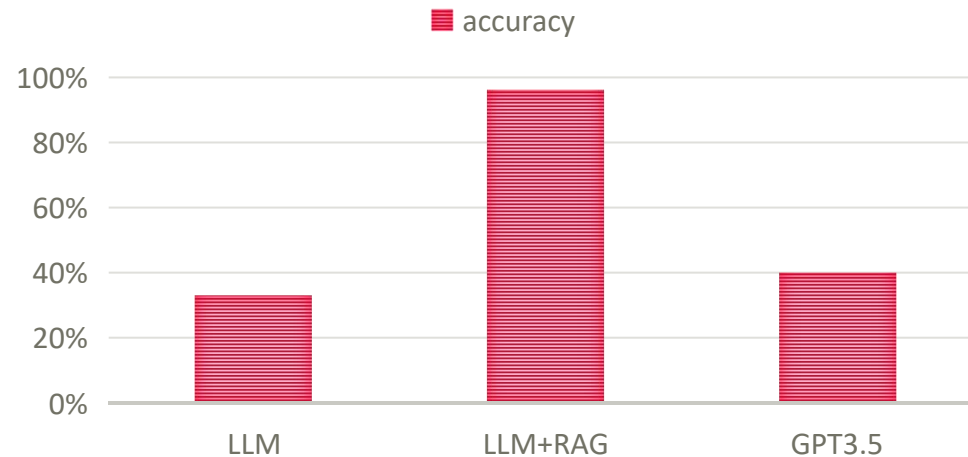
Ji Sun

Huawei

**VectorDB is widely used in the LLM applications, and the apps are highly dependent on the performance, reliability, security and scalability of VectorDB**

# A typical pipeline of LLM application in the Telecommunication Field

Telecommunication Knowledge Q&A

Telecommunication Tools

Security Check

**RAG**

query

VDB

Scalar Filtering

Prompt Template

LLM

Fact Check

LLM

A

**Extraction**

Fine Sorting

Doc Search

Vector Search

query

entry point

Embedding

Q A

Fine Tuning

Data Augment

Q A

token   hash   List

t1
t2
t3
t4

docid

doc2vector

QA Extraction

Customized dictionary

Doc split

Keywork Mining

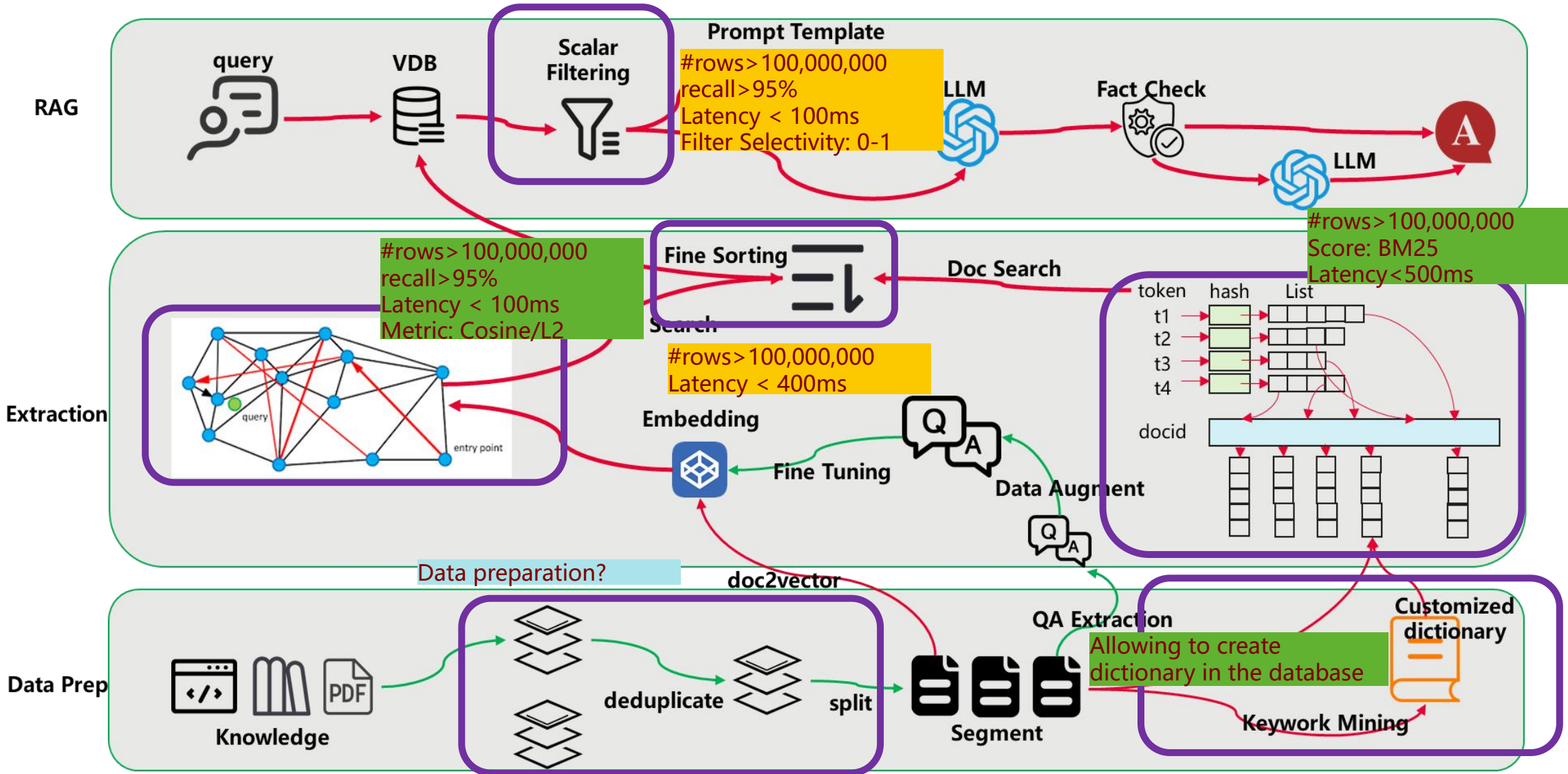**Data Prep**

Knowledge

deduplicate

split

Segment

# Details on this application

➢ **Pain Point:** LLM is not sensitive to the meaning of numeric data, and the conclusion is not trustworthy. For serious

questions in the telecommunication field, LLM may produce nonsense words (even if they seem to be logical).

➢ **Data Source:** The source documents of telecommunication knowledges come from real product, including the product

document, maintanence cases, external papers&journals

➢ **RAG:** users query the knowledge database by natural language, then the related knowledges are fed into LLM by prompt,

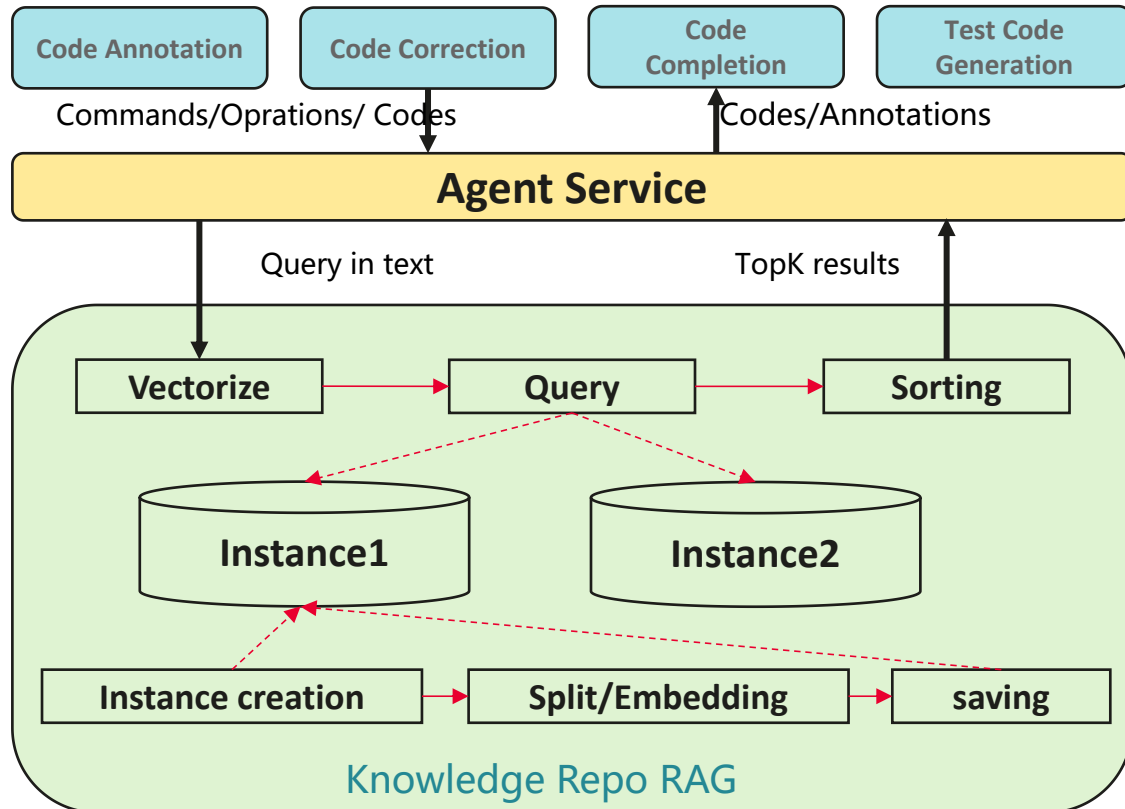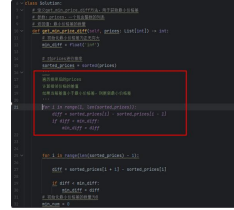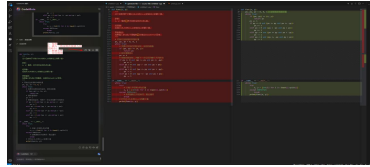and the LLM generate logical, smooth, correct and helpful answers.

# Vector Database Requirements



**RAG**

query → VDB → Scalar Filtering → Prompt Template → LLM → Fact Check → LLM → A

Prompt Template
#rows>100,000,000
recall>95%
Latency < 100ms
Filter Selectivity: 0-1

#rows>100,000,000
Score: BM25
Latency<500ms

**Extraction**

Fine Sorting ← Doc Search

#rows>100,000,000
recall>95%
Latency < 100ms
Metric: Cosine/L2

Search

#rows>100,000,000
Latency < 400ms

Embedding
Fine Tuning ← QA ← Data Augment

token  hash  List
t1
t2
t3
t4

docid

entry point
query

Data preparation?

doc2vector

QA Extraction

**Data Prep**

Knowledge → deduplicate → split → Segment

Allowing to create dictionary in the database

Customized dictionary

Keywork Mining

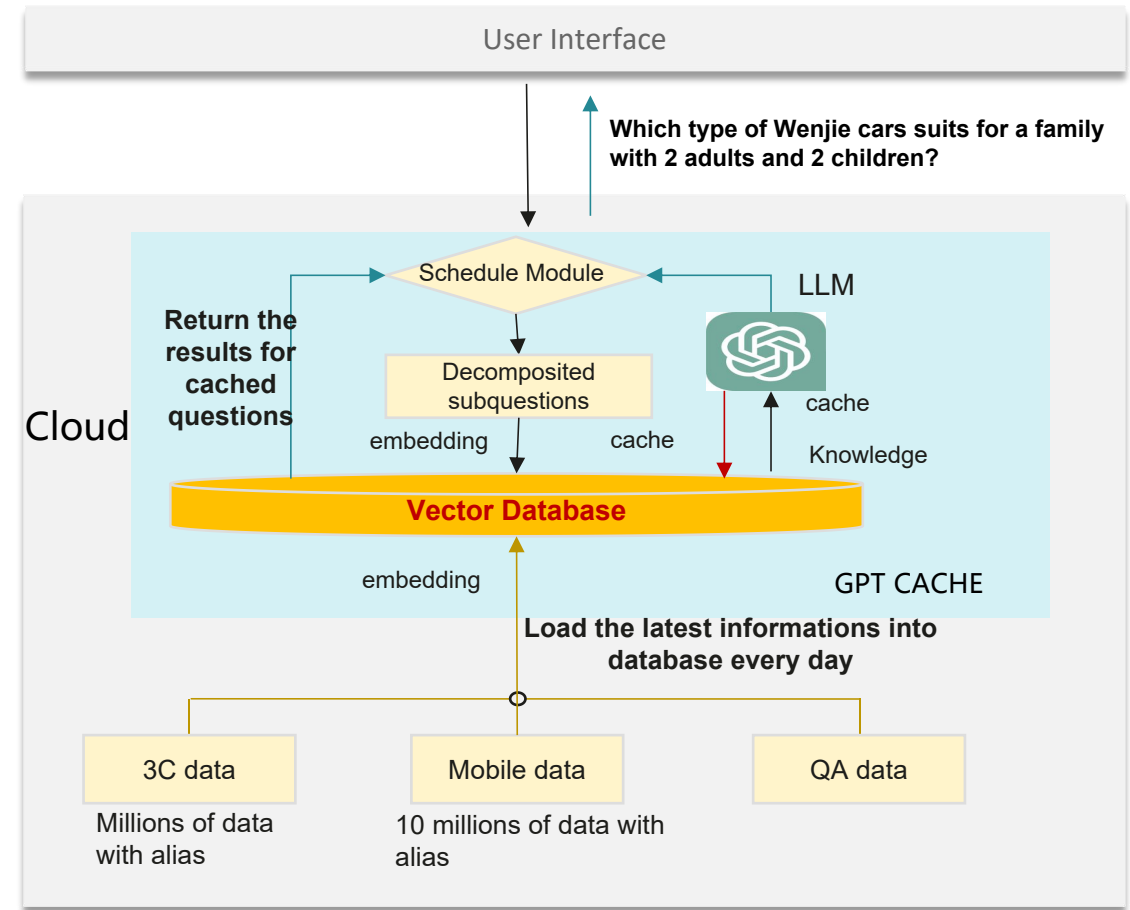# Applications with vector databases deployed on the cloud

## Developer Copilot

- **Millions of data**
- **accurate search(Top5)**
- **multiple tasks, different requirements**
- **heterogeneous inputs/outputs**



| Code Annotation | Code Correction | Code Completion | Test Code Generation |
|---|---|---|---|

Commands/Oprations/ Codes    Codes/Annotations

**Agent Service**

Query in text    TopK results

**Vectorize** → **Query** → **Sorting**

**Instance1**    **Instance2**

**Instance creation** → **Split/Embedding** → **saving**
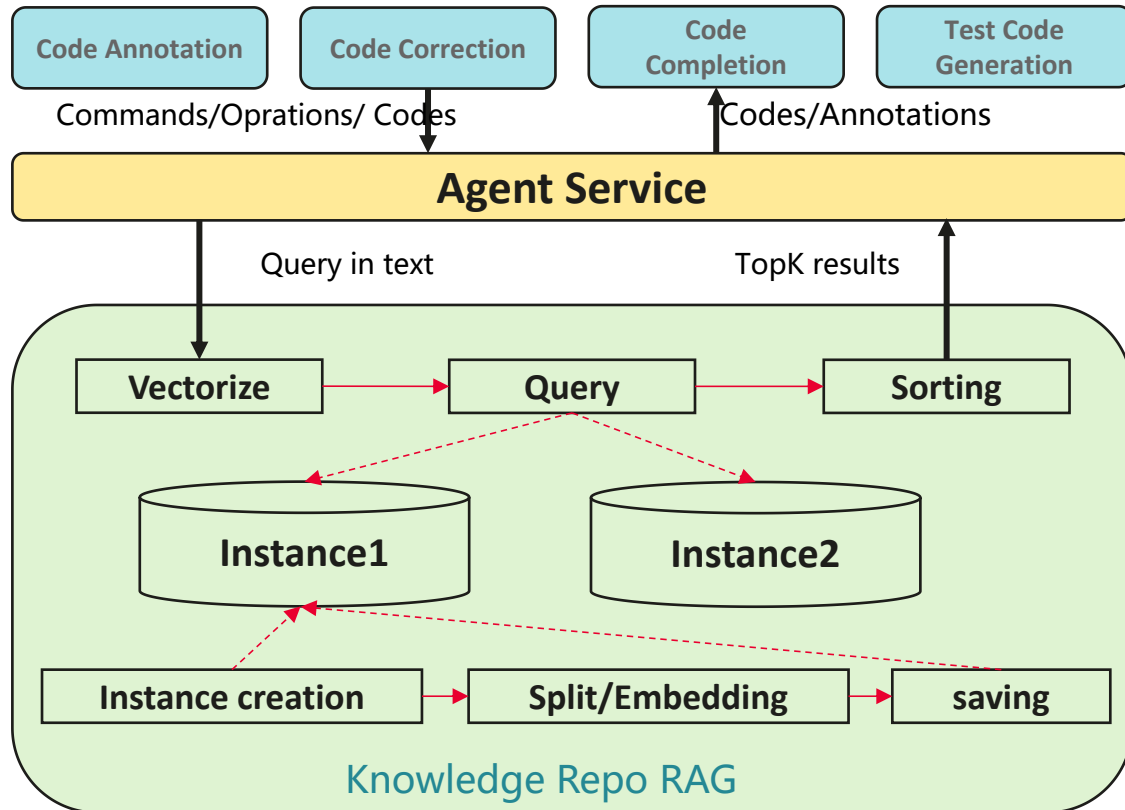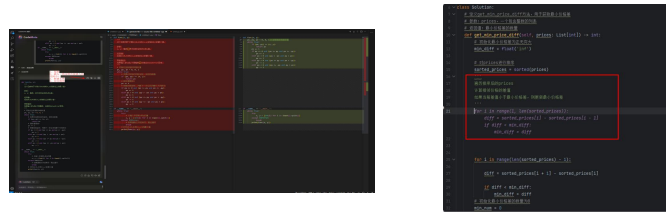
*Knowledge Repo RAG*

## Voice Assistant

- **Public service on cloud, privacy matters a lot**
- **High concurrency, high throughput is required**
- **Complicated questions, decomposition is required**
- **Relies on Time-to-live to manage the cache**



User Interface

Which type of Wenjie cars suits for a family with 2 adults and 2 children?

Cloud

Schedule Module    LLM

**Return the results for cached questions**

Decomposited subquestions    cache

embedding    cache    Knowledge

**Vector Database**

embedding    GPT CACHE

**Load the latest informations into database every day**

| 3C data | Mobile data | QA data |
|---|---|---|

Millions of data with alias    10 millions of data with alias

# Applications with vector databases deployed on the cloud

## Developer Copilot

- **Light-weighted index**
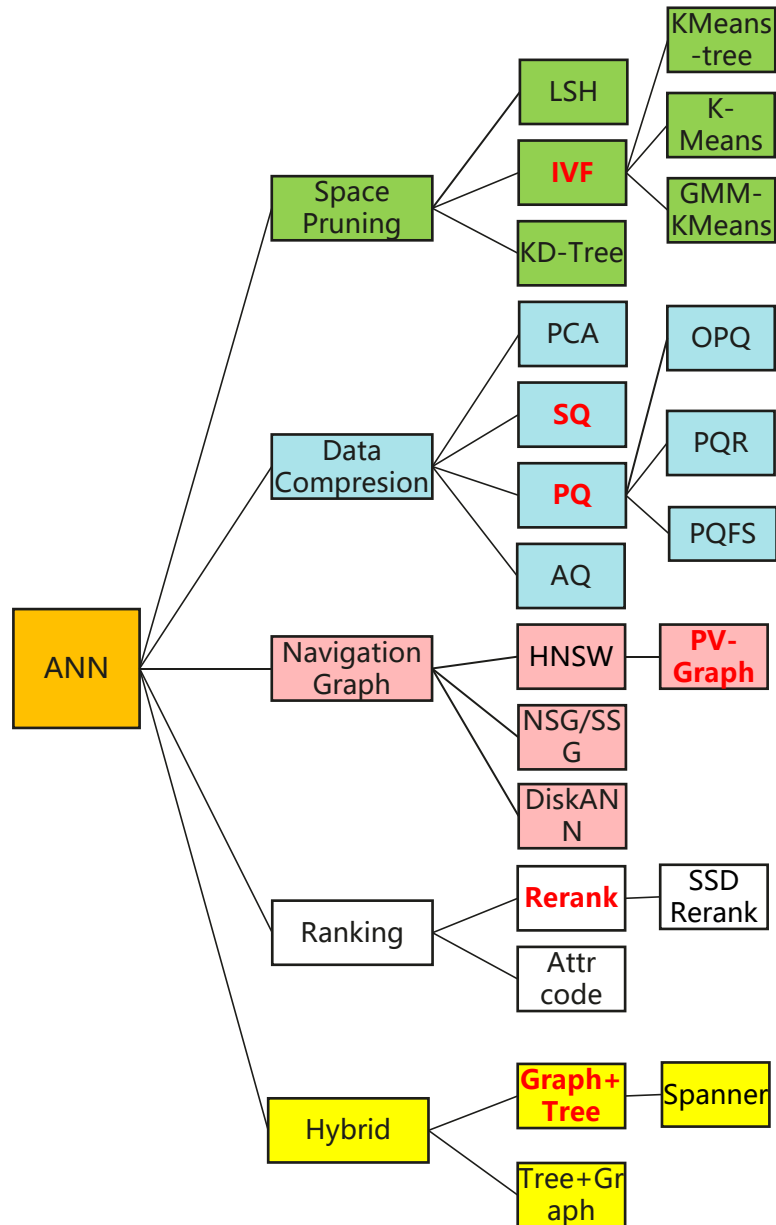- **Fast I/O, Computing**
- **Disaggregation**
- **Resource Scheduling/Scaling**



| Code Annotation | Code Correction | Code Completion | Test Code Generation |
|---|---|---|---|

Commands/Oprations/ Codes                    Codes/Annotations

**Agent Service**

Query in text                    TopK results

### Knowledge Repo RAG

Vectorize → Query → Sorting

Instance1          Instance2

Instance creation → Split/Embedding → saving

## Voice Assistant

- **Privacy protection**
- **Multi-tenant/Multi-read**
- **Query in batch**
- **Updatable**



User Interface

Which type of Wenjie cars suits for a family with 2 adults and 2 children?

Cloud

Schedule Module          LLM

**Return the results for cached questions**

Decomposed subquestions

embedding          cache          cache          Knowledge

**Vector Database**

embedding                    GPT CACHE

**Load the latest informations into database every day**

| 3C data | Mobile data | QA data |
|---|---|---|

Millions of data with alias

10 millions of data with alias

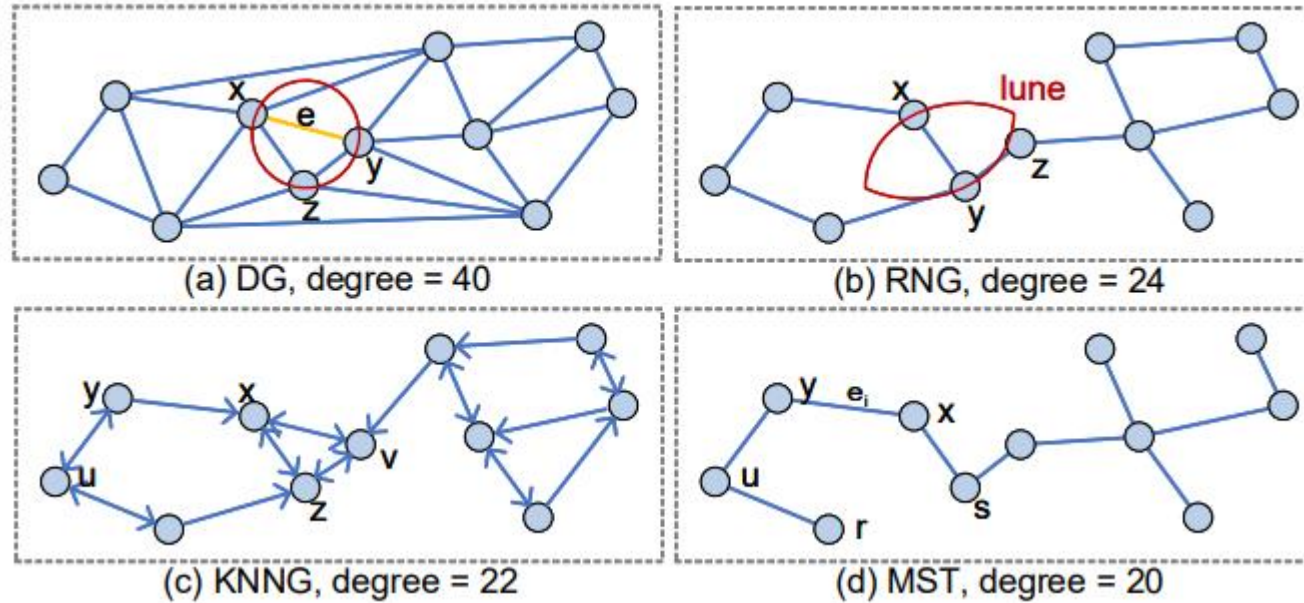# An algorithm view of vector similarity search

# Algorithm optimization for approximate nearest neighbors



Brute-force KNN takes over 2 minutes for 100 million vectors.

1. The recall is not necessarily 100%, the second nearest neighbor maybe better because of the limitation of embedding methods.

2. Nearest neighbors have locality naturally, if close, continue to search; otherwise, jump out.
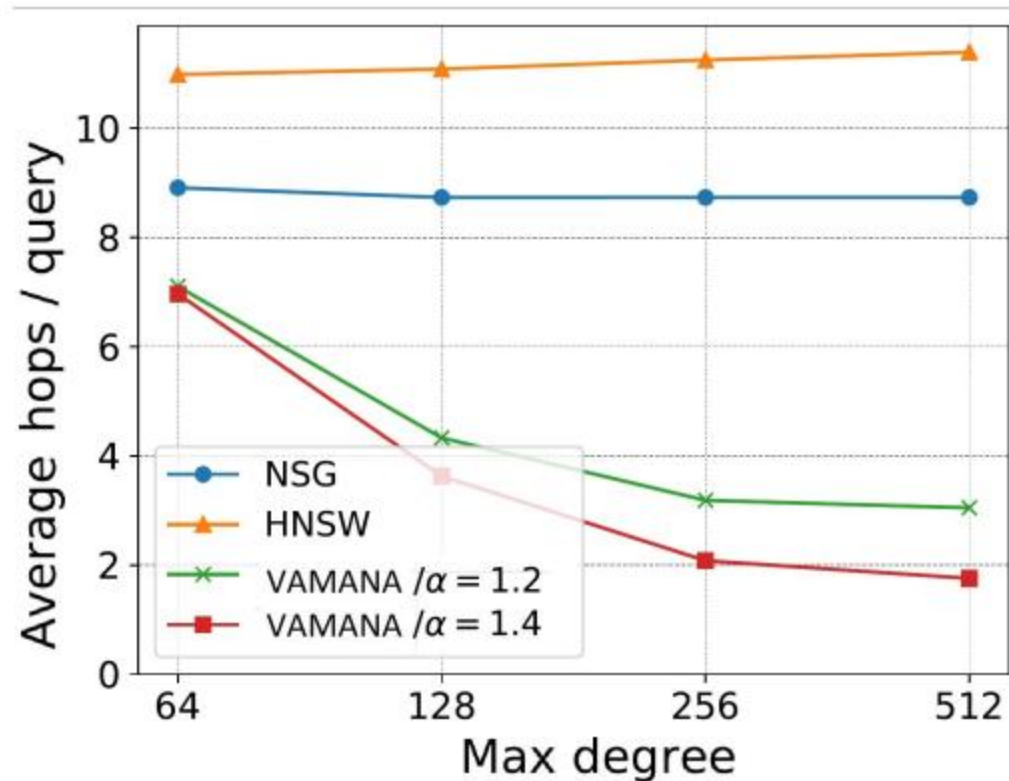
# Graph-based Algorithms



(a) DG, degree = 40

(b) RNG, degree = 24

(c) KNNG, degree = 22

(d) MST, degree = 20

|  | KGraph | NSW | HNSW | Vamana | NSG |
|---|---|---|---|---|---|
| Global Connectivity | Middle | Middle | Middle | Good | Good |
| Search Complexity (in memory, experimental) | $O(|S|^{0.54})$ | $O(\log^2(|S|))$ | $O(\log(|S|))$ | $O(|S|^{0.75})$ | $O(\log(|S|))$ |
| Space Consumption | Middle | Middle | Middle | Good | Good |
| Construction Complexity (in memory, experimental) | $O(|S|^{1.14})$ | $O(|S|\cdot\log^2(|S|))$ | $O(|S|\cdot\log(|S|))$ | $O(|S|^{1.16})$ | $O(|S|^{((1+c)/c}\cdot\log(|S|) + |S|^{1.14})$ |

A Comprehensive Survey and Experimental Comparison of Graph-Based Approximate Nearest Neighbor Search. Mengzhao Wang, Xiaoliang Xu, Qiang Yue, Yuxiang Wang. VLDB, 2021

# DiskANN shows the best #hops when searching

Average number of hops vs maximum graph degree
for achieving 98% 5-recall@5 on ANN_SIFT1M



DiskANN: Fast Accurate Billion-point Nearest Neighbor Search on a Single Node.
Suhas Jayaram Subramanya, Devvrit, Rohan Kadekodi, Ravishankar Krishaswamy, Harsha Vardhan Simhadri. NIPS, 2019
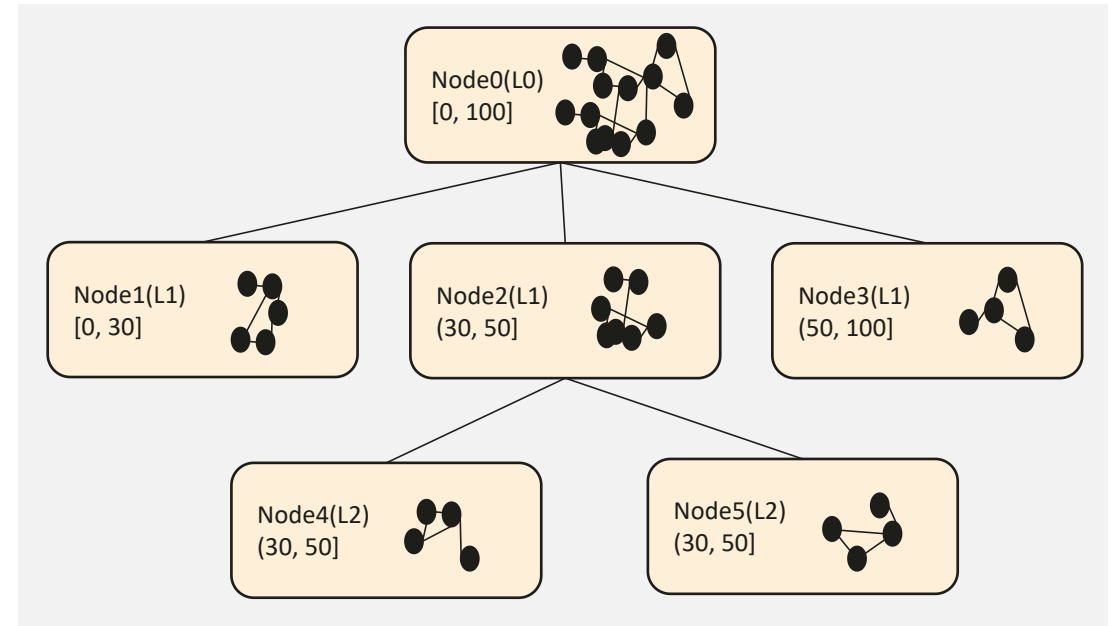
# Multi-column scalar/vector hybrid index, support attribute filtering

**Background**: Vector is suitable for semantic search, scalar attribute is also important for access control, accuracy improve, and performance improve.



Prior Filtering | Post Filtering | Hybrid Filtering



**Scalar-vector hybrid query is supported by following 3 mainstream methods:**

1. **Prior Filtering**. Sequence scan on selected vector is very slow. We can build index for each scalar label, but too many small index causes the range scan/ Full scan slow.

2. **Post Filtering.** Efficient when the selectivity of scalar condition is high. When the selectivity is low, the query will extend the candidate list constantly for enough result.

3. **Hybrid Filtering.** Consider scalar label when computing vector distances to find tuples whose vector is close to the query vector, and scalar matches the query condition. However, this method is not accurate enough.

1. Full scan can be executed in the root node0, which has the same performance with the full graph;

2. Query with equal condition can search on the leaf node with higher performance;

3. Query with range condition can search on a small part of nodes for better performance
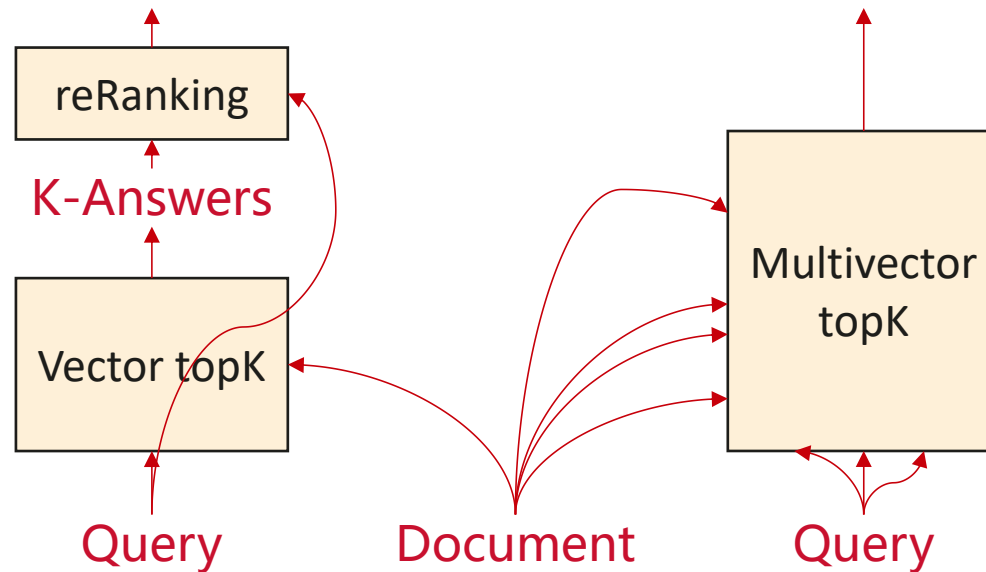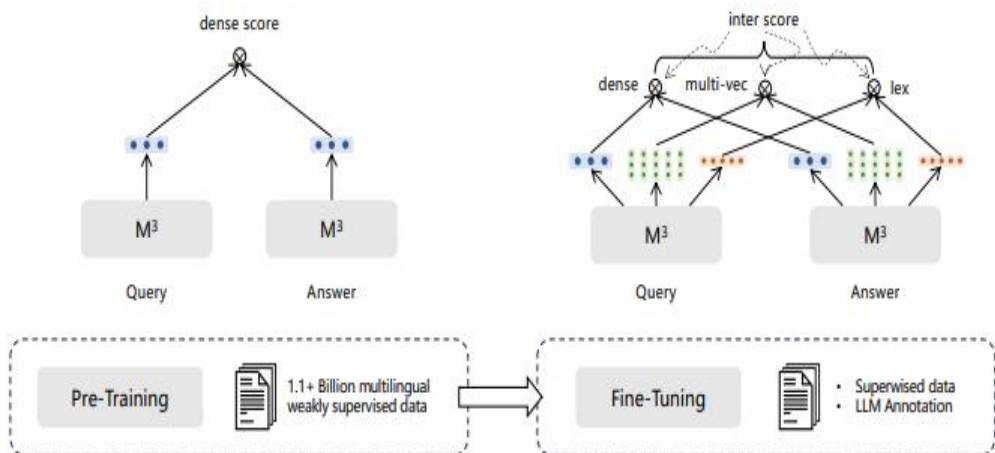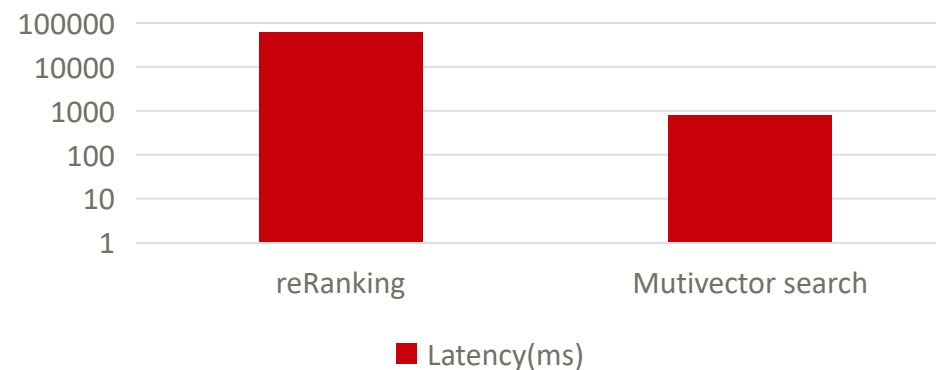
# Light-weighted multi-vector similarity query for ranking



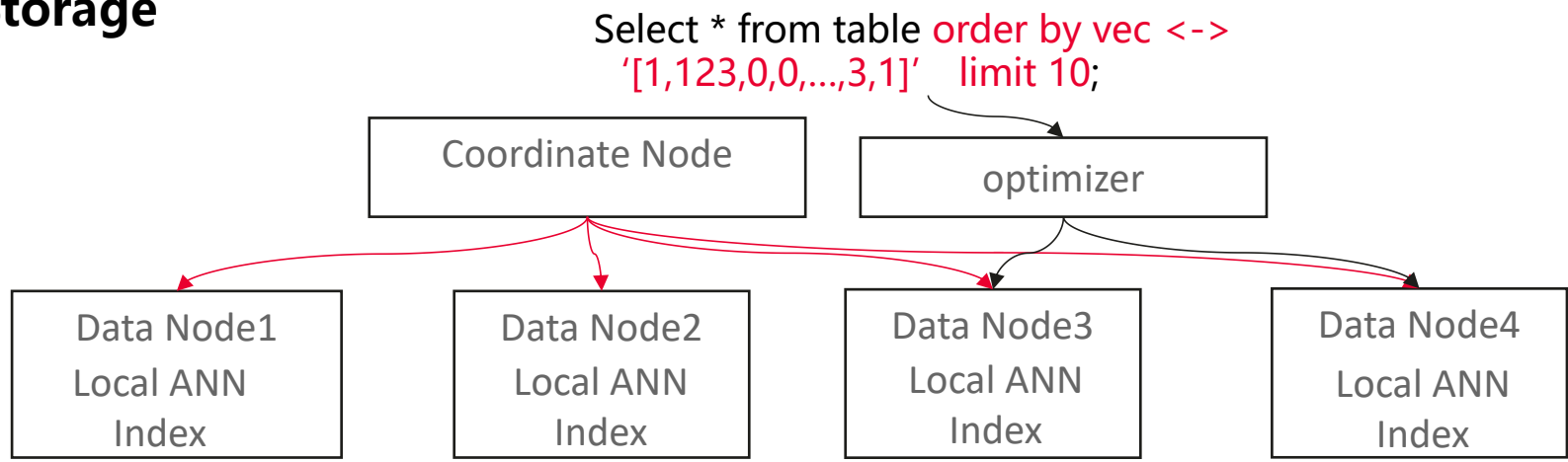ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT



BGE M3-Embedding: Multi-Lingual, Multi-Functionality, Multi-Granularity Text Embeddings Through Self-Knowledge Distillation
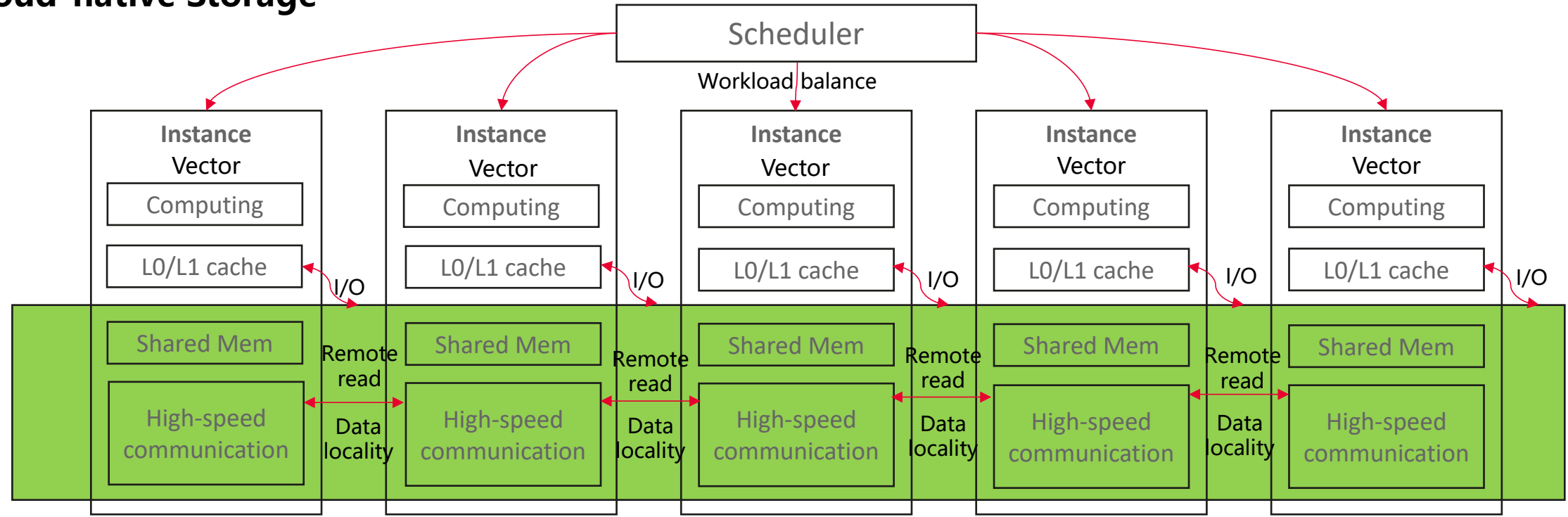




Performance Improvement

# A system view of vector database

# Vector Index on distributed database （ 10^9 < $N$ ）

## Shared-nothing Storage

Select * from table order by vec <->
'[1,123,0,0,...,3,1]'    limit 10;

Coordinate Node

optimizer

| Data Node1 Local ANN Index | Data Node2 Local ANN Index | Data Node3 Local ANN Index | Data Node4 Local ANN Index |
|---|---|---|---|

## Cloud-native Storage

Scheduler

Workload balance

| Instance Vector | Instance Vector | Instance Vector | Instance Vector | Instance Vector |
|---|---|---|---|---|
| Computing | Computing | Computing | Computing | Computing |
| L0/L1 cache | L0/L1 cache | L0/L1 cache | L0/L1 cache | L0/L1 cache |

I/O    I/O    I/O    I/O    I/O

| Shared Mem | Shared Mem | Shared Mem | Shared Mem | Shared Mem |
|---|---|---|---|---|
| High-speed communication | High-speed communication | High-speed communication | High-speed communication | High-speed communication |

Remote read / Data locality

Remote read / Data locality
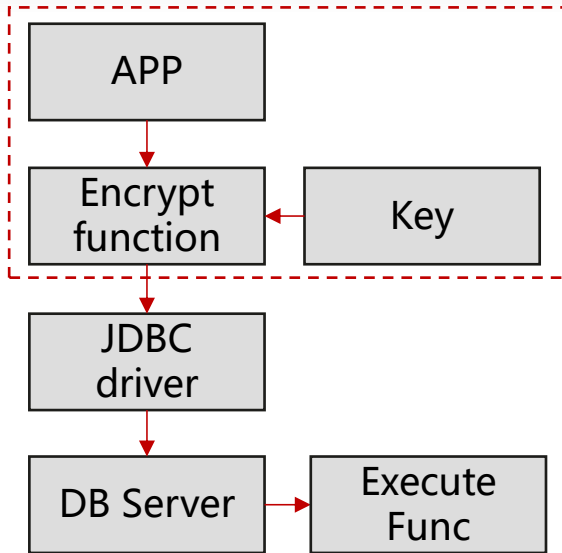
Remote read / Data locality

Remote read / Data locality

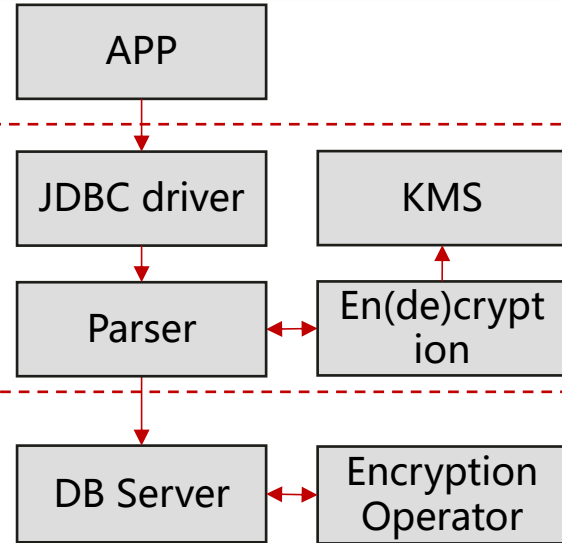# Data encryption for privacy protection

## Functional Encryption



- column-level, user-side key, server-side encryption
- APP knows encryption, not support encrypted condition
- Data is decrypted in session, database cannot decrypted automatically, avoid data being stolen by other users.
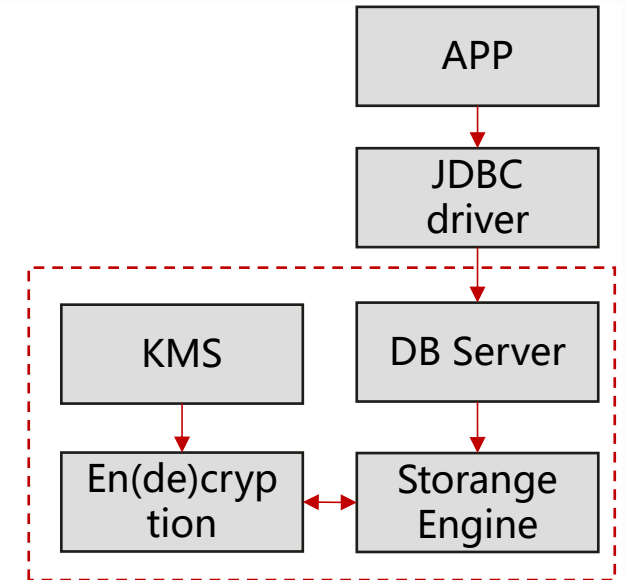
**VS**

## Fully Encryption



- Column-level, user-side key, driver embedded
- APP is not aware of encryption, support encrypted equal conditions.
- DB cannot decrypt data, avoid data leaking during OM activities.
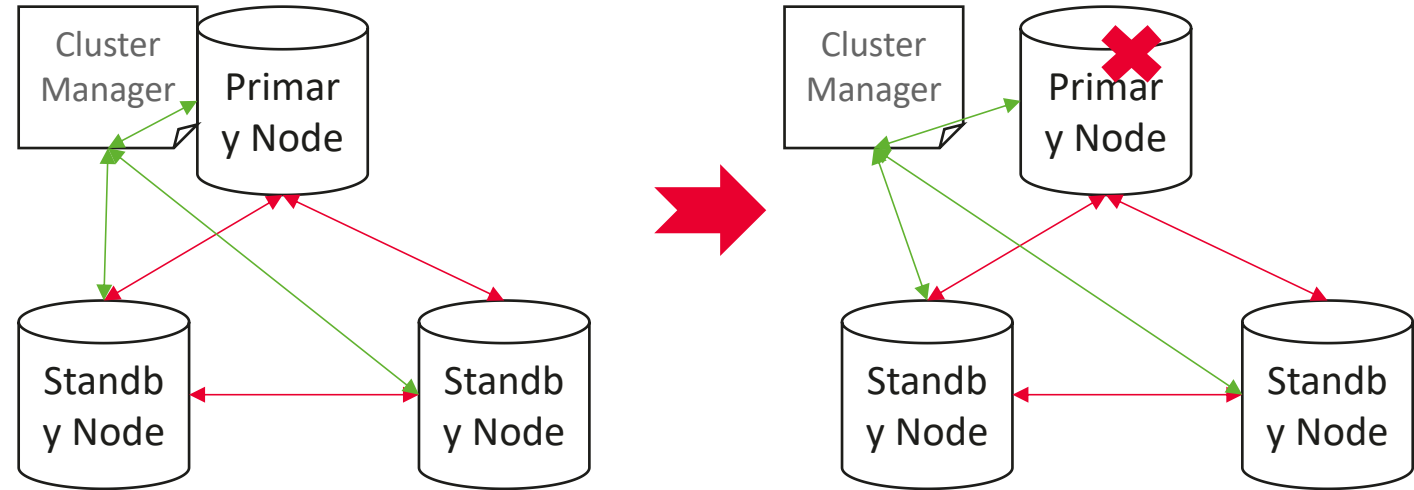
**VS**

## Transparent Data Encryption



- Table-level, server-side key, server-side encryption
- APP is not aware of encryption, plaintext in the memory.
- Prevent data from being stolen on the disk.

# High availability by data replication

| XLOG Types | XLOG Description |
|---|---|
| XLOG_DISKANN_OPERATIONS | Index Operations |
| ... | |

Failure Recovery



1. Customizes xlog collections for vector index, supports synchronize between primary node and standby nodes；

2. Supports fast RTO < 10s for service recovery；

3. Supports data recovery, users can recover the data to any version when misoperation or disk damage happen (PITR).

# Thanks!