

# GaussDB-Vector: A Large-Scale Persistent Real-Time Vector Database for LLM Applications

**Ji Sun<sup>1</sup>, Guoliang Li<sup>1</sup>, James Pan<sup>1</sup>, Jiang Wang<sup>2</sup>,  
Yongqing Xie<sup>2</sup>, Ruicheng Liu<sup>2</sup>, Wen Nie<sup>2</sup>**

<sup>1</sup>Department of Computer Science, Tsinghua University

<sup>2</sup>Huawei

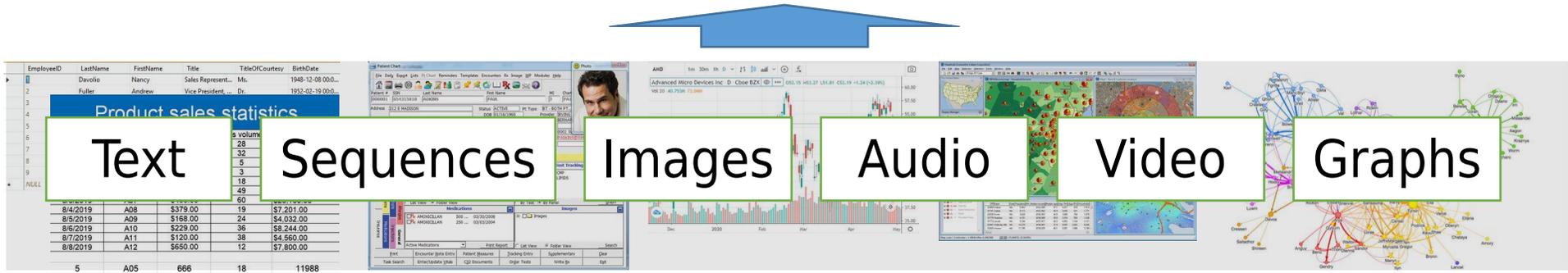
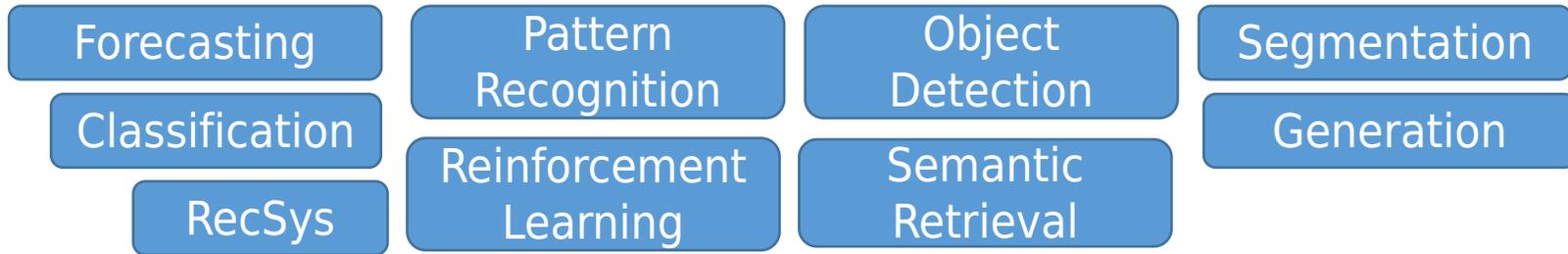


清华大学  
Tsinghua University

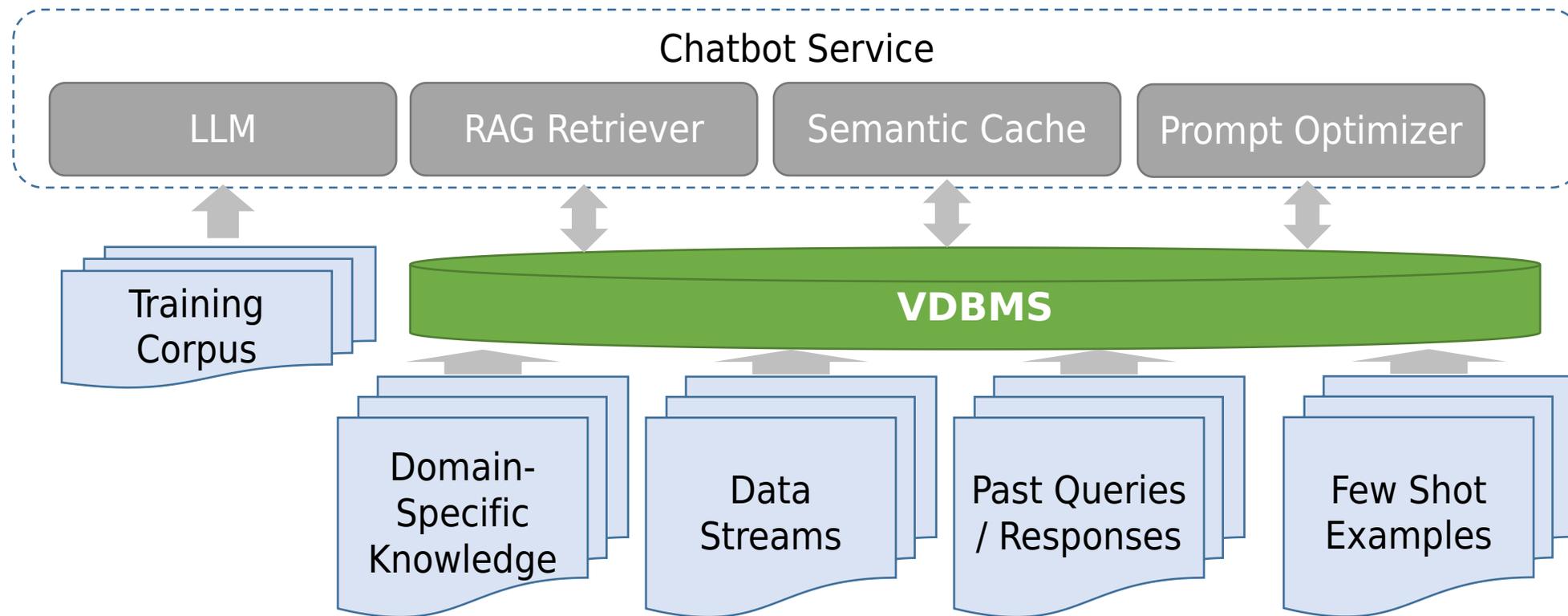


# Embeddings: Building Blocks of Future

*More and more applications rely on embedding vectors that are retrieved via similarity search*



# Vector Databases and LLMs



## Search Capabilities

- Support various queries, e.g. k-NN, Range NN, Predicated (Filtered) Search

## Performance

- Low latency queries
- High throughput
- Memory efficient

## Quality

- High recall
- High semantic relevance

## System-Level Features

- Scalable
- Available
- Elastic
- Consistent
- Persistent

# Vector Databases: Key Challenges

*More and more applications rely on embedding vectors that are retrieved via similarity search*

**Goal: Store embeddings in a DB and retrieve for whatever downstream task**



1) **Size:** Existing hardware not for large data units

2) **Structure:** Lack of structure (e.g. attrs) makes embeddings hard to index

## Performance

- **Latency**

- Fundamentally high I/O due to more frequent paging
- Naive disk-resident HNSW is not I/O optimal
- High latency of filtered search

## Quality

- **Recall**

- Accuracy degradation of IVF due to data drift following updates

## System-Level Features

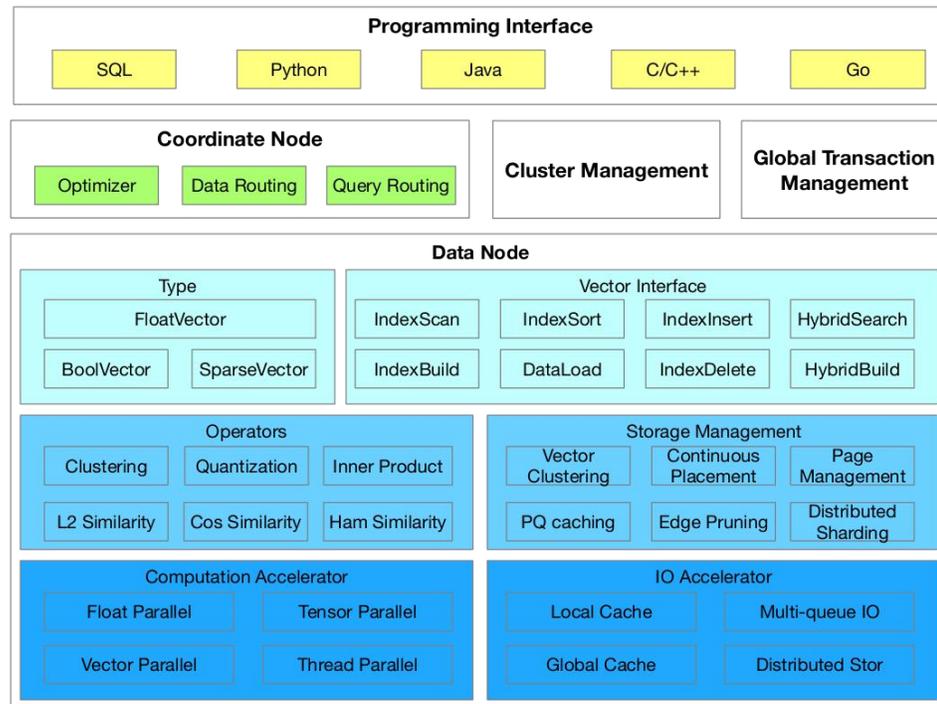
- **Scalability**

- Unclear sharding and query routing strategy for distributed deployments

# GaussDB-Vector

**Build a distributed persistent real-time database based on IVF & DiskANN**

- **Architecture:** Two-Tiered Distributed Database to provide scalability / availability
- **Techniques:** IVF & DiskANN but make it easier to build, update, and search



## Index Storage & Construction

- Two-Phase Neighborhood Pruning
- **Tailored Storage Structures to reduce I/O**
- **Data Sharding with Drift Detection & Auto Rebalancing**
- **Balanced Tree Hybrid Index**

## Search

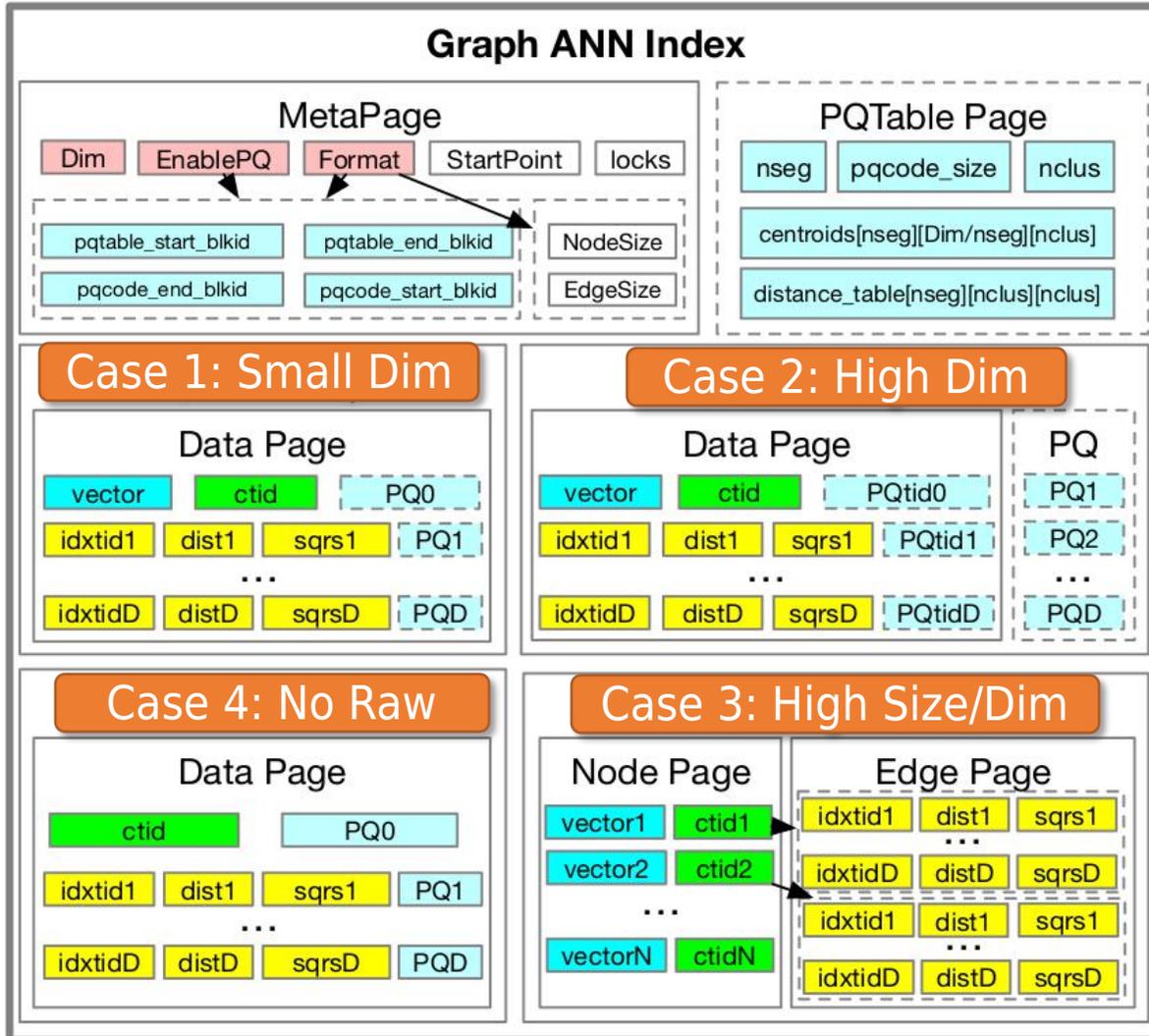
- Hot Nodes Buffering
- **Cardinality-Based Query Routing for Distributed Search**
- I/O Efficient Search Algorithms

## Update

- Asynchronous Delete
- **Incremental Index Updates for data freshness**

# Tailored Storage Structures

**Use Tailored Storage Structures based on dataset size/dimensions to reduce I/O**



## Case 1: Small Dimensionality

- Raw Vectors + PQ Codes stored together
- Each node is stored with all its neighbors on a single page to reduce I/O during neighborhood scan

## Case 2: High Dimensionality

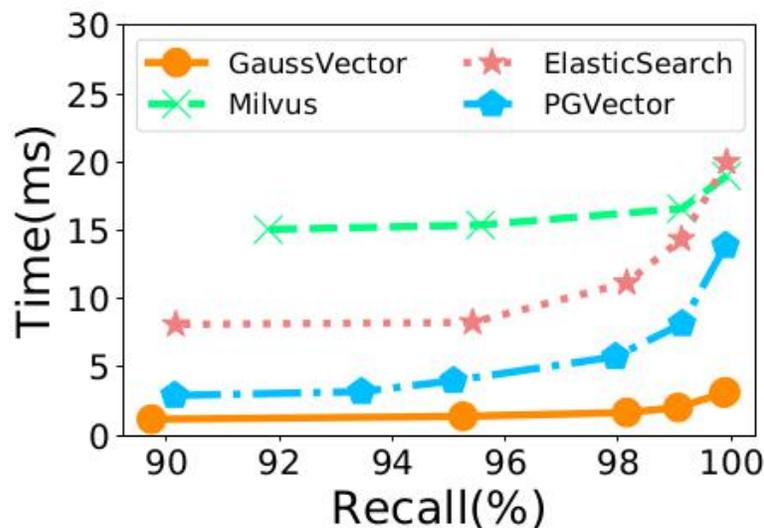
- Raw Vectors are too large and cannot be stored together with PQ Codes
- PQ Codes stored in separately pages and cached in memory

## Case 3: High Size/Dimensionality

- Raw Vectors are too large and cannot be stored with PQ Codes nor Neighbors
- Store Nodes and Neighborhoods separately

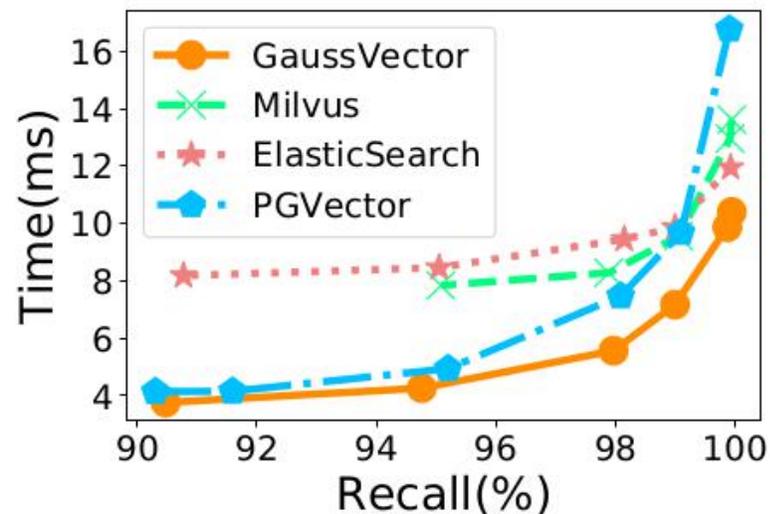
# Tailored Storage Structures

**Use Tailored Storage Structures based on dataset size/dimensions to reduce I/O**

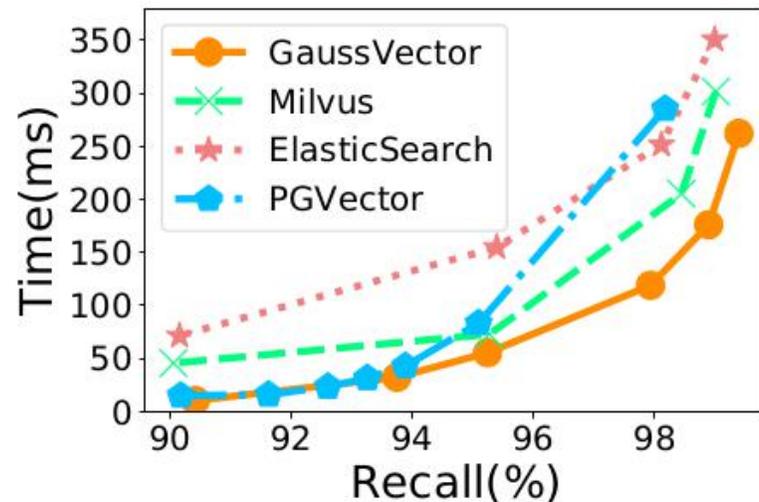


(a) SIFT dataset.

Case 1: Small Dim



(b) GIST dataset.



(c) HUWEINet dataset.

Case 2: High Dim

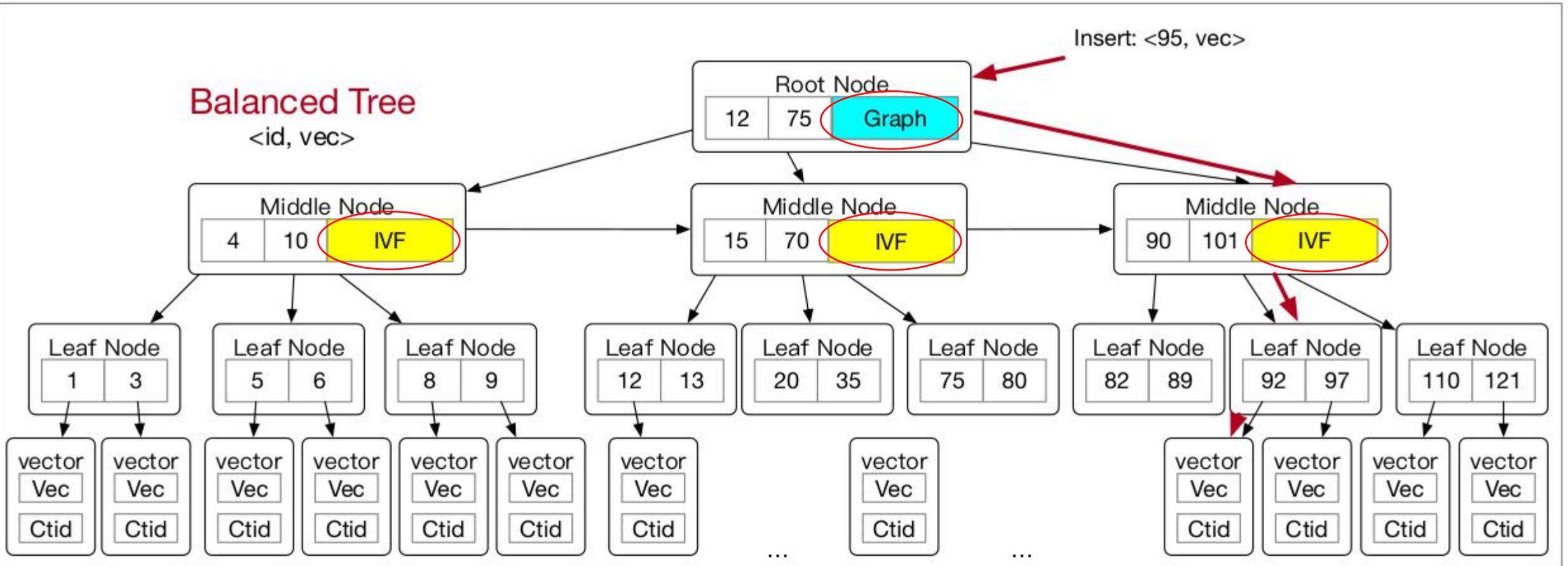
ElasticSearch & Milvus require loading multiple disk segments into memory

GaussDB-Vector adapts to high-dim scenarios to maintain high search performance

# Balanced Tree Hybrid Index

Use **Balanced Tree** with intermediate sub-indexes to support fast pre- / single-stage filtering

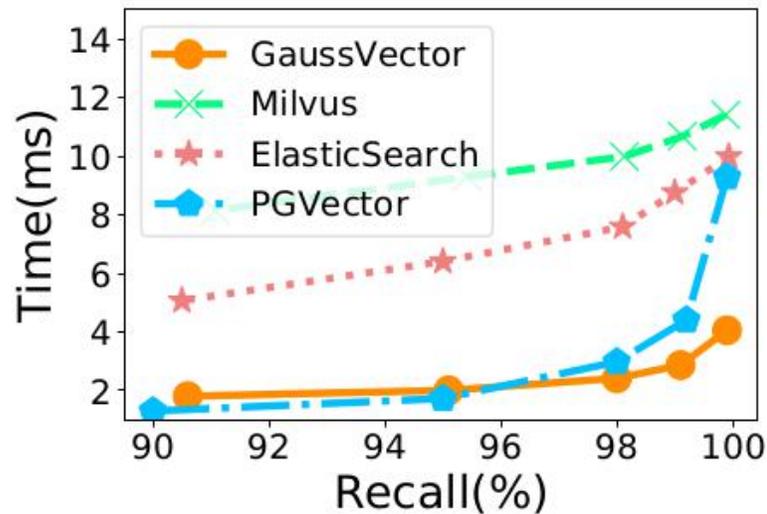
- **High Selectivity** (few satisfying records): Sequential Scan over leaf nodes
- **Low Selectivity** (many satisfying records): Single-Stage Scan over sub-indexes



# Balanced Tree Hybrid Index

Use **Balanced Tree** with intermediate sub-indexes to support fast pre- / single-stage filtering

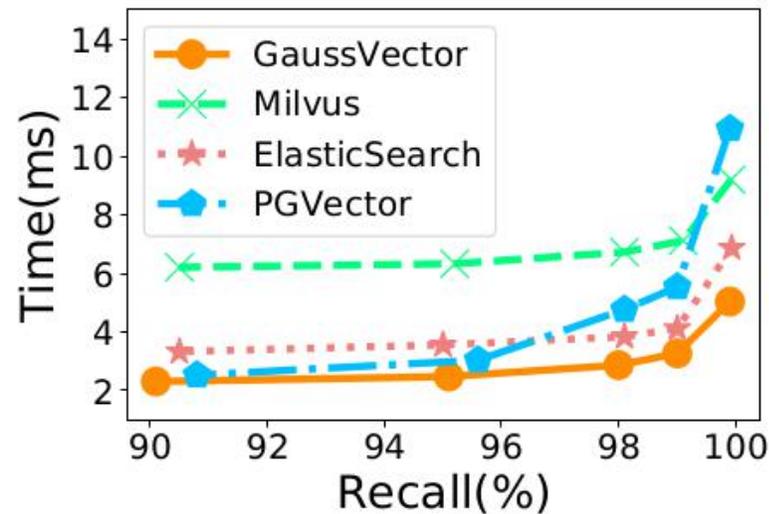
- **High Selectivity** (few satisfying records): Sequential Scan over leaf nodes
- **Low Selectivity** (many satisfying records): Single-Stage Scan over sub-indexes



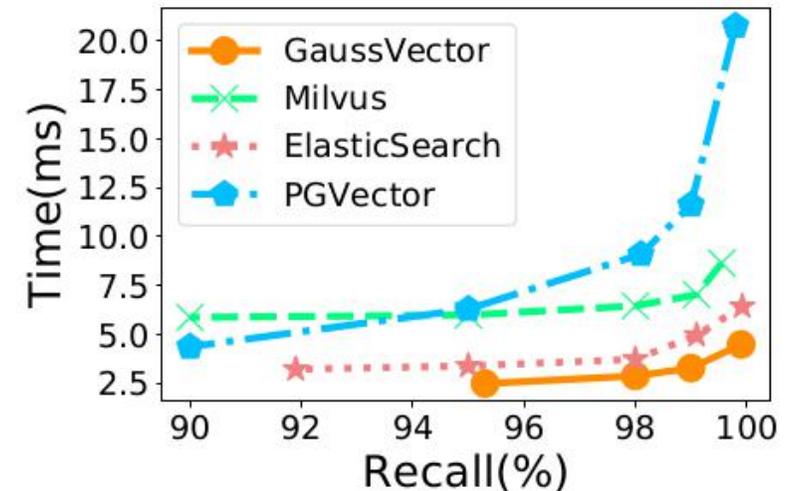
(a) SIFT dataset, 90% scalar selectivity.

Many Satisfying Records

ElasticSearch & Milvus require merging across many segments



(b) SIFT dataset, 50% scalar selectivity.



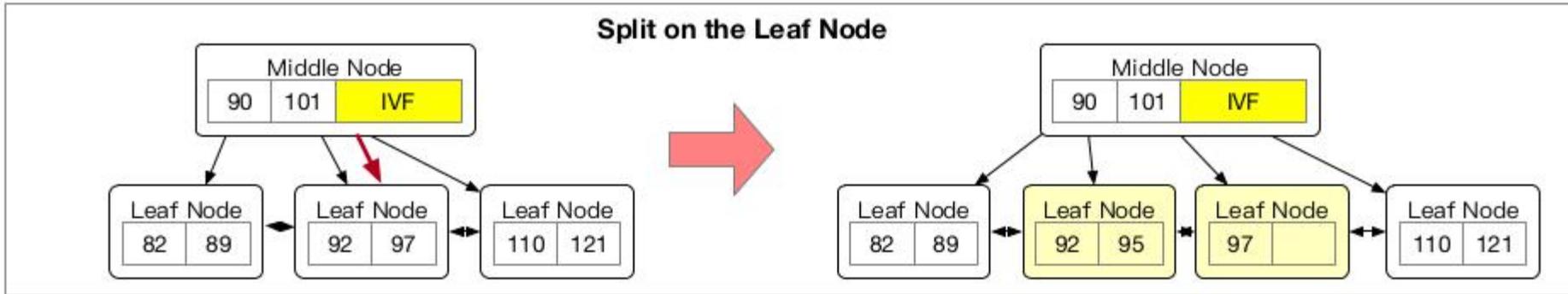
(c) SIFT dataset, 10% scalar selectivity.

Few Satisfying Records

Pre-Filtering via index followed by Sequential Scan beats single-stage scan over segments

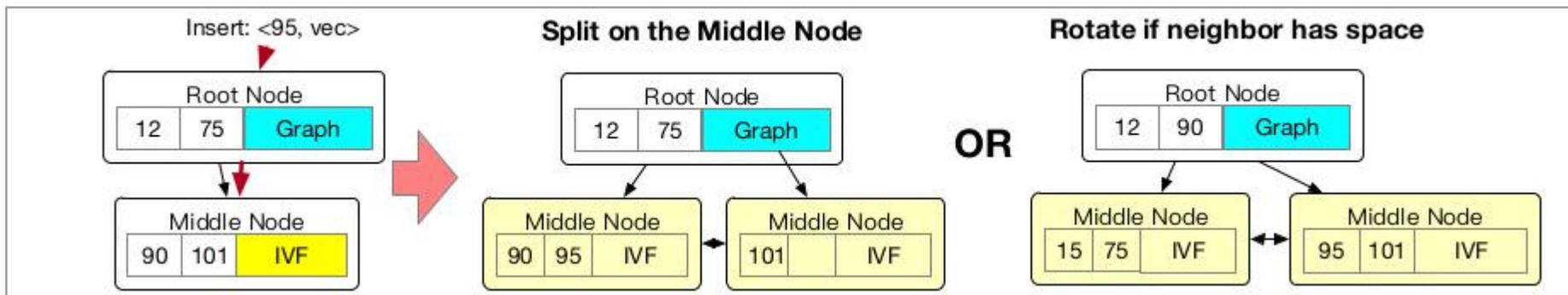
# Incremental Index Updates

**Balanced Tree with intermediate sub-indexes enables incremental updates for data freshness**



Leaf node split  
requires no index  
update

**(a) Splitting a leaf node following an insert**

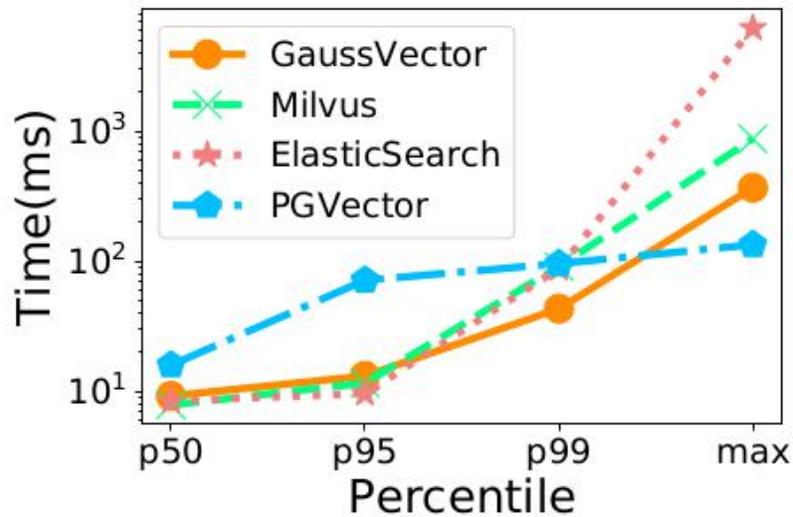


Inner node with IVF  
index is easy to  
update

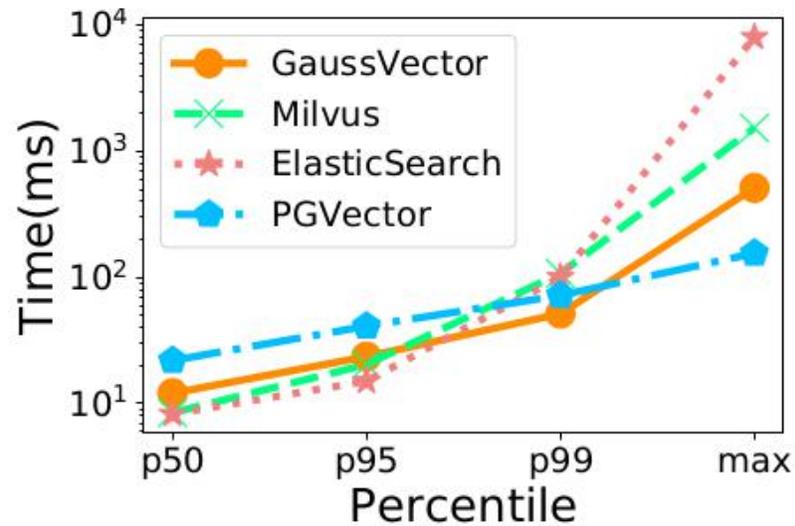
**(b) Splitting an inner node following an insert**

# Incremental Index Updates

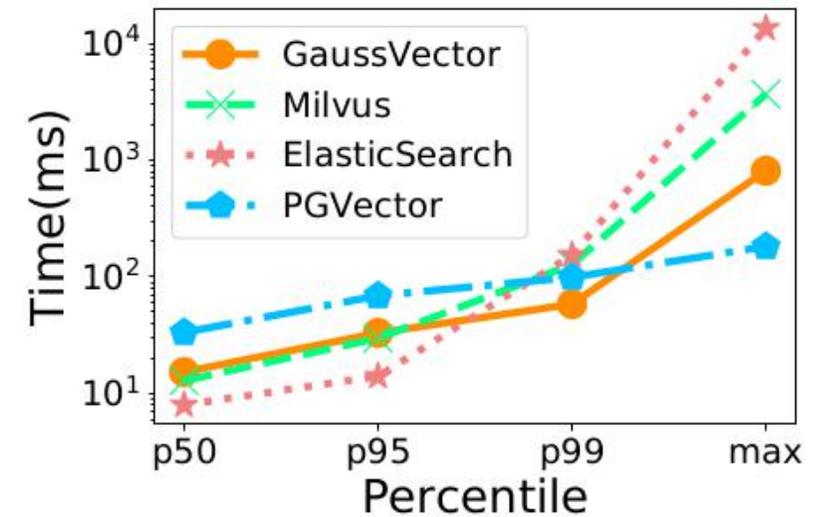
**Balanced Tree with intermediate sub-indexes enables incremental updates for data freshness**



(a) SIFT dataset.



(b) GIST dataset.

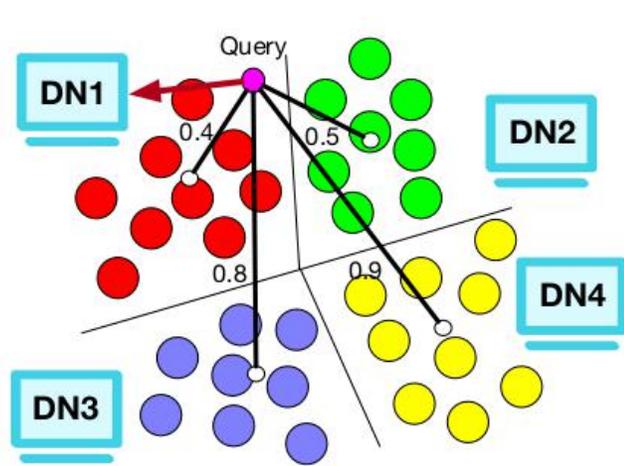


(c) HUAWEINet dataset.

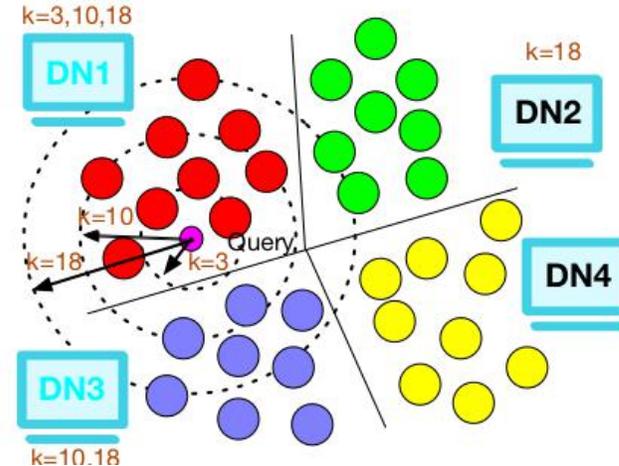
- Large max latency due to inner node / root node splits but otherwise competitive with baselines
- Insert / search latency balanced by adjusting indexing granularity

# Data Sharding & Distributed Search

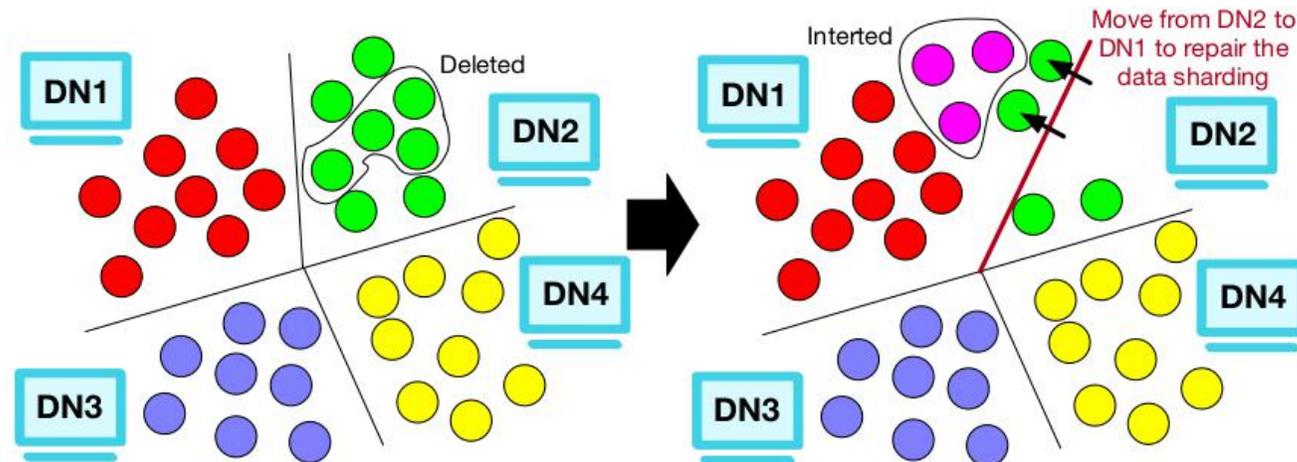
*Balanced Tree with intermediate sub-indexes enables incremental updates for data freshness*



**(a) Insert into Nearest Shard**



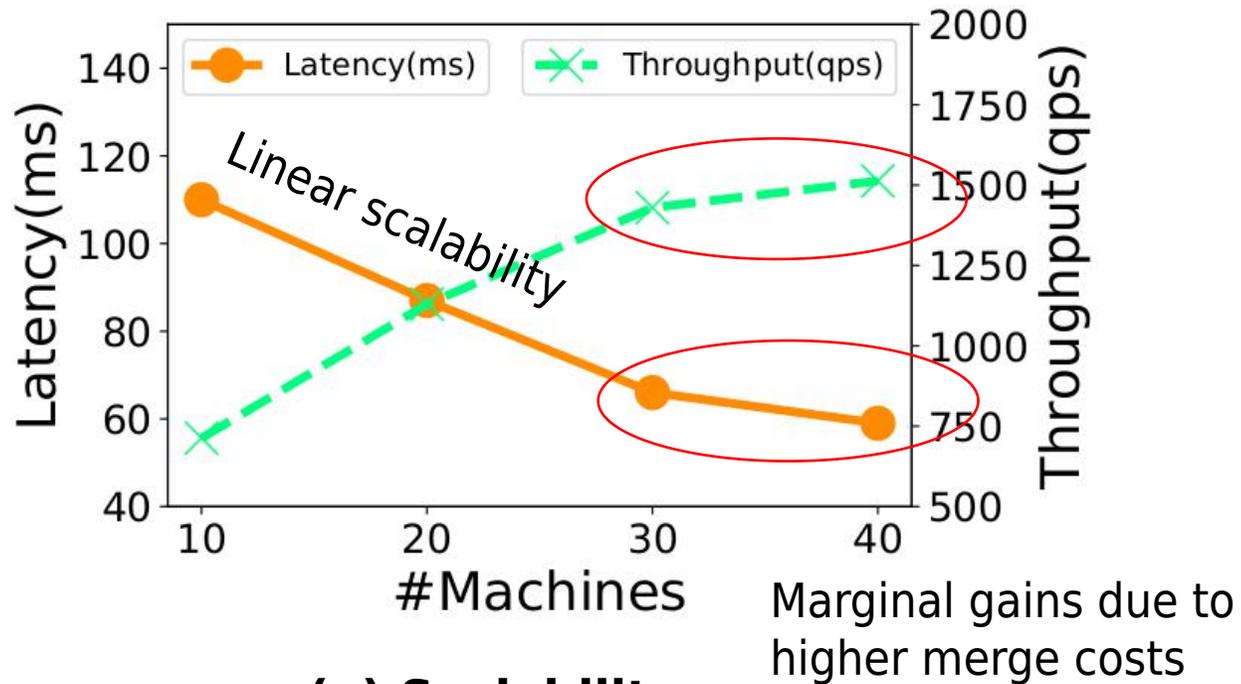
**(b) Search Shards based on Estimated Cardinality**



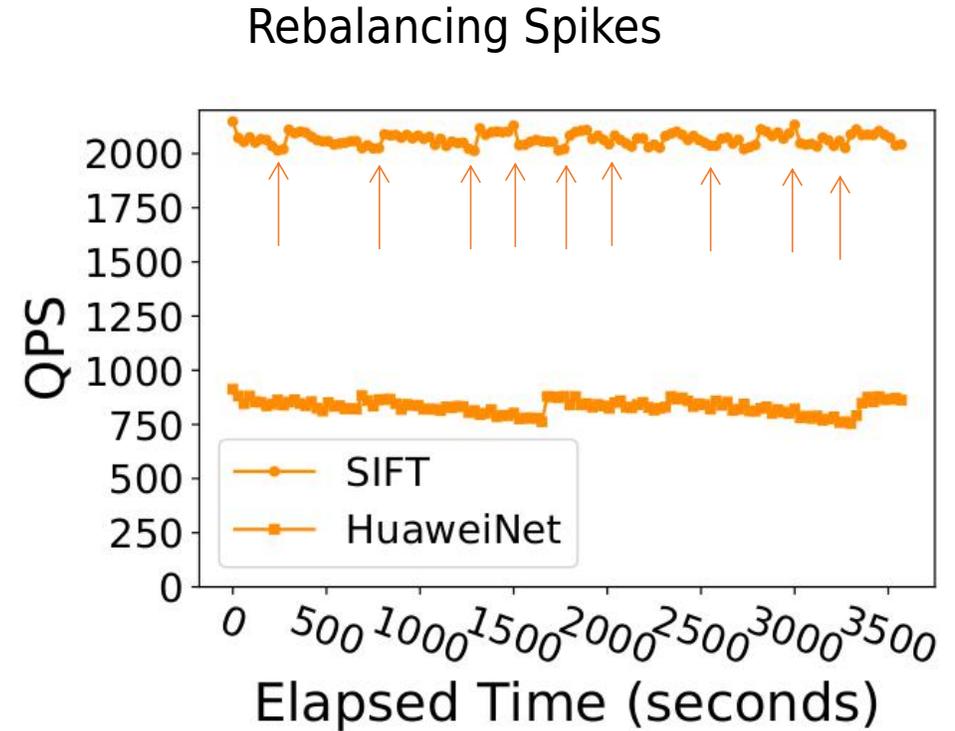
**(c) Background Monitoring & Remigration**

# Data Sharding & Distributed Search

**Balanced Tree with intermediate sub-indexes enables incremental updates for data freshness**



**(a) Scalability**



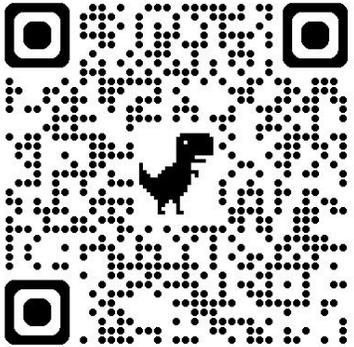
**(b) Update Throughput**

# Summary & Conclusion

**GaussDB-Vector integrates vector search capabilities within existing data management system**

	Query Processor	Storage & Indexing	Query Optimizer	Query Executor
 <p>Milvus</p>	<ul style="list-style-type: none"> <li>• k-NN</li> <li>• Range NN</li> <li>• Predicated k-NN</li> <li>• Multi-vector search</li> <li>• Grouping Search</li> <li>• Full Text Search</li> <li>• Reranking</li> </ul>	<ul style="list-style-type: none"> <li>• Replicas</li> <li>• Data Partitions</li> <li>• Multiple index types</li> <li>• Blob storage (backup &amp; persistence)</li> </ul>	<ul style="list-style-type: none"> <li>• Cost-based predicated query planning</li> </ul>	<ul style="list-style-type: none"> <li>• Scatter-gather</li> <li>• R/W Disagg.</li> <li>• CPU/Store Disagg.</li> <li>• Log-backed durable writes</li> <li>• Cache &amp; SIMD acceleration</li> <li>• GPU acceleration</li> <li>• Tunable consistency</li> </ul>
 <p>PostgreSQL (pgvector)</p>	<ul style="list-style-type: none"> <li>• k-NN</li> <li>• Range NN</li> <li>• Predicated k-NN</li> <li>• Full Text Search</li> <li>• <b>Relational queries</b></li> </ul>	<ul style="list-style-type: none"> <li>• Page-Based Storage</li> <li>• Shards &amp; Replicas</li> <li>• HNSW, IVF</li> <li>• Quantization</li> </ul>	<ul style="list-style-type: none"> <li>• Cost-based predicated query planning</li> </ul>	<ul style="list-style-type: none"> <li>• MVCC</li> <li>• Scatter-gather</li> <li>• Log-backed durable writes</li> <li>• etc.</li> </ul>
 <p>GaussDB-Vector</p>	<ul style="list-style-type: none"> <li>• k-NN</li> <li>• Range NN</li> <li>• Predicated k-NN</li> <li>• <b>Relational queries</b></li> </ul>	<ul style="list-style-type: none"> <li>• <b>Page-Based Storage</b></li> <li>• <b>Shards &amp; Replicas</b></li> <li>• HNSW, IVF, <b>Hybrid Index</b></li> <li>• Quantization</li> </ul>	<ul style="list-style-type: none"> <li>• Cost-based predicated query planning</li> </ul>	<ul style="list-style-type: none"> <li>• GaussDB Transaction Manager</li> <li>• <b>Scatter-gather</b></li> <li>• Log-backed durable writes</li> <li>• Cache &amp; SIMD acceleration</li> <li>• NPU/GPU acceleration</li> </ul>

# Thanks!



GaussDB-Vector:  
A Large-Scale Persistent Real-Time Vector  
Database for LLM Applications PVLDB (2025)

**Slides:**

<https://dbgrouop.cs.tsinghua.edu.cn/ligl/activities.html>



James Pan

[jpan@tsinghua.edu.cn](mailto:jpan@tsinghua.edu.cn)