

# Milvus: A Purpose-Built Vector Data Management System

Jianguo Wang etc.  
SIGMOD 2021

# Background: Unstructured Data are Dominating

Int, float,  
String, ...

01234  
56789

$e$   $\pi$

ABCDEFGH

2021.04.10

Text

## Abstract

Bigtable is a distributed storage system for managing structured data that is designed to scale to a very large size: petabytes of data across thousands of commodity servers. Many projects at Google store data in Bigtable, including web indexing, Google Earth, and Google Finance. These applications place very different demands on Bigtable, both in terms of data size (from URLs to web pages to satellite imagery) and latency requirements (from backend bulk processing to real-time data serving). Despite these varied demands, Bigtable has successfully provided a flexible, high-performance solution for all of these Google products. In this paper we describe the simple data model provided by Bigtable, which gives clients dynamic control over data layout and format, and we describe the design and implementation of Bigtable.

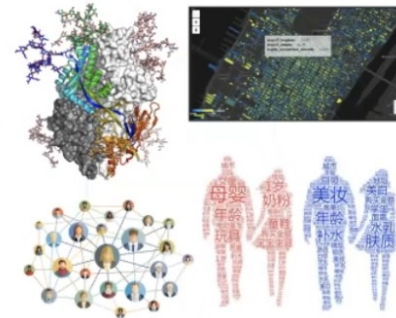
Json

```
{
  "firstName": "John",
  "lastName": "Smith",
  "isActive": true,
  "age": 12,
  "address": {
    "streetAddress": "21 2nd Street",
    "city": "New York",
    "state": "NY",
    "postalCode": "10021-3100"
  },
  "phoneNumbers": [
    {
      "type": "home",
      "number": "212 555-1234"
    },
    {
      "type": "office",
      "number": "646 555-4567"
    }
  ],
  "children": [],
  "spouse": null
}
```

Image  
video  
audio



Domain specific

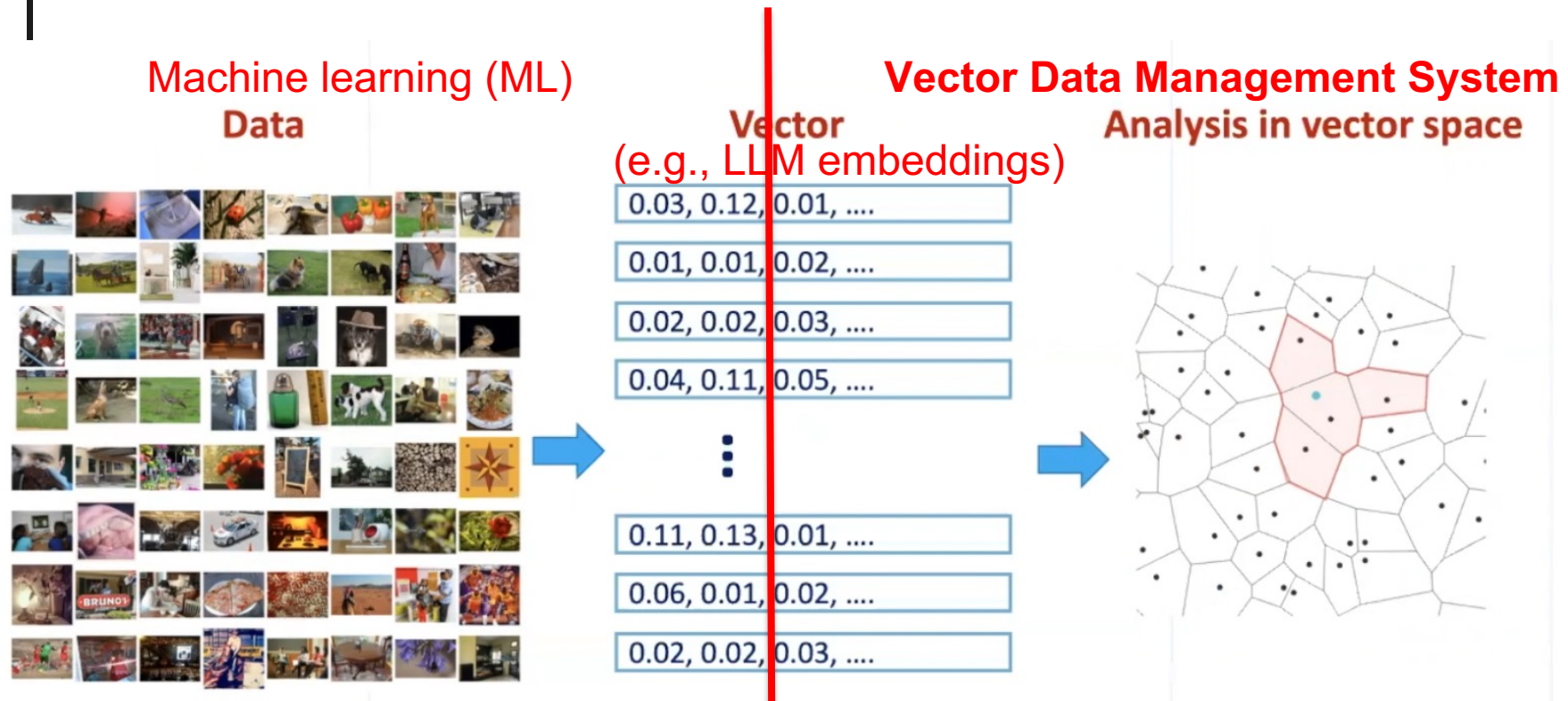


Structured data

Unstructured data

IDC: 80% data will be unstructured by 2025

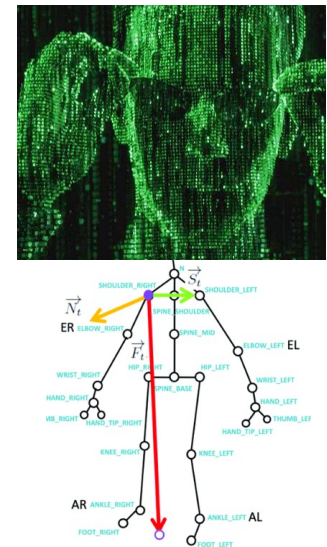
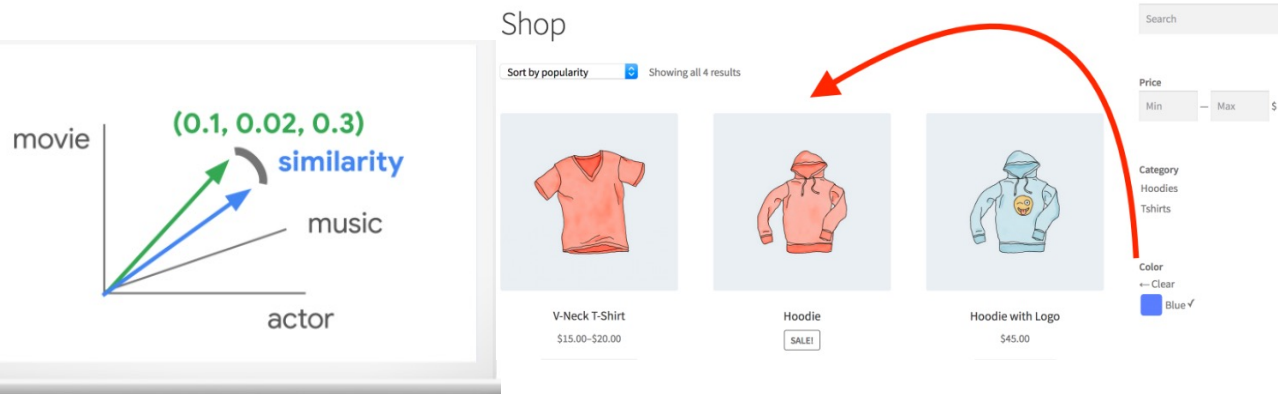
# Background: Vector Analysis Work Flow



ML is soaring and bringing increasingly large volume of vectors, and how about the Vector Data Management System?

# Requirements and Challenges for Scalable Vector Data Management System

- **Fast** query processing and **dynamic** vector manipulation
- **Advanced** query semantics:
  - Vector similarity search
  - Attribute filtering
  - Multi-vector search

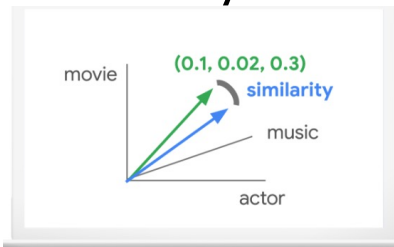


However, existing work fail  
for poor performance  
and limited functionalities...

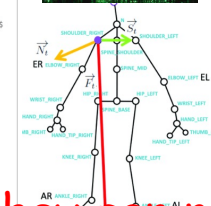
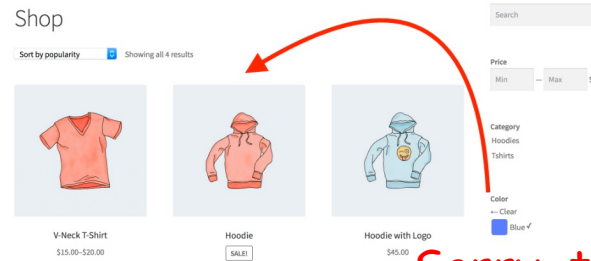
# Limitations of Existing Work

	Billion-Scale Data	Dynamic Data	GPU	Attribute Filtering	Multi-Vector Query	Distributed System
Facebook Faiss [3, 35]	✓	✗	✓	✗	✗	✗
Microsoft SPTAG [14]	✓	✗	✗	✗	✗	✗
ElasticSearch [2]	✗	✓	✗	✓	✗	✓
Jingdong Vearch [4, 39]	✗	✓	✓	✓	✗	✓
Alibaba AnalyticDB-V [65]	✓	✓	✗	✓	✗	✓
Alibaba PASE (PostgreSQL) [68]	✗	✓	✗	✓	✗	✗
Milvus (this paper)	✓	✓	✓	✓	✓	✓

- Existing algorithms (e.g., [14]) don't support large data storage at scale or dynamic data with fast updates.
- Existing systems (e.g., [3]) simply extend relational DB and fail to fully optimize query processing on vector data
- Both lack the support for more advanced query processing beyond simple **vector similarity search**.



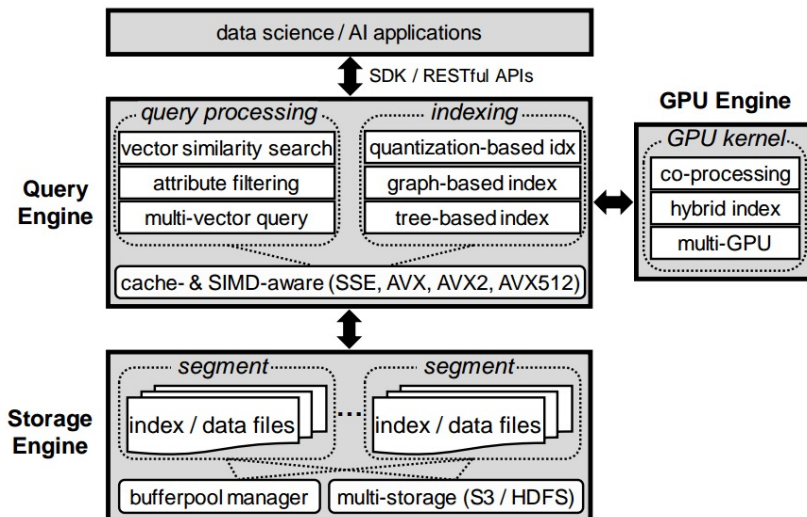
Yes, they can.



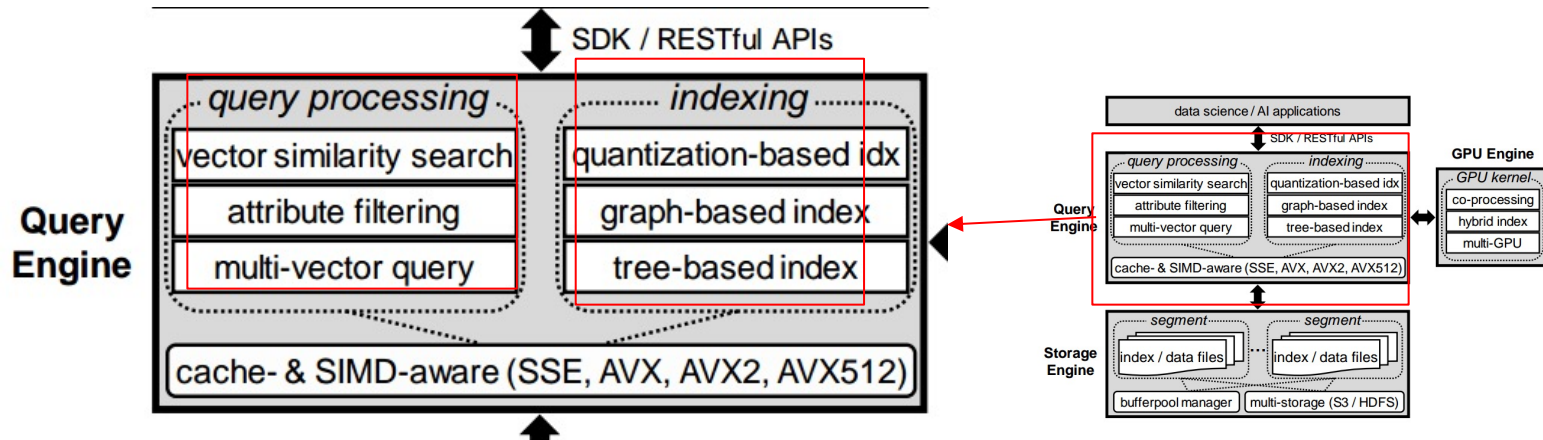
Sorry, they can not!!

# Milvus Overview

- All-in-one solution
  - Functional complete system + optimized algos
  - Hardware conscious design on both CPU and GPU
  - Easy SDK in Python/Java/Go/C++ and RESTful APIs
- Three major components
  - Query Engine
  - GPU Engine
  - Storage Engine



# The Query Engine

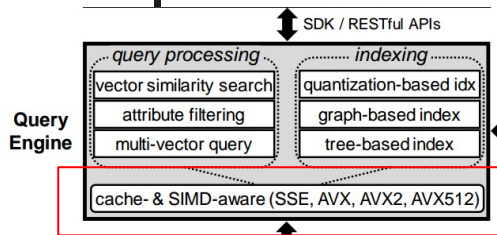


- Support vector(s) and vector(s)+non-vector attributes, therefore 3 primitive advanced queries\*
- Cover common similarity and index functions
- Adopts LSM-tree [47] for efficient insertions and deletions
  - First in MemTable
  - Next flush into disk by rules (e.g., by size or time)

\*Recall: Existing vectorDB only support similarity search



# The Query Engine



## Cache-aware and SIMD optimizations

### Cache-aware

- Reuse accessed data for L3 opt.
  - Hence reduce L3 miss
- Fine-grained parallelism of core-datablock assigns
  - Outperform OpenMP

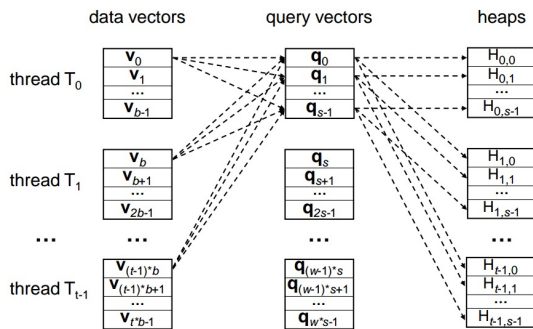
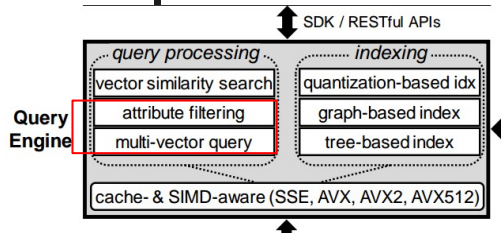


Figure 3: Cache-aware design in Milvus

### SIMD

- AVX512 support and automatic instruction choose in runtime

# The Query Engine



- Further opt. on attribute filtering and multi-vector query\*
  - Four adopted + one newly proposed attribute filtering
    - Partition-based is up to 13.7x faster than existing algos
  - Two new approaches instead of naïve solution for multi-vector query
    - i.e., vector fusion and iterative merging

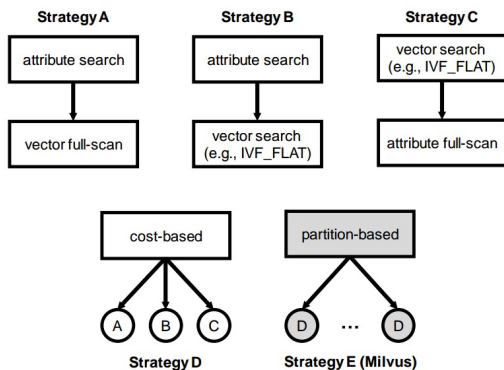
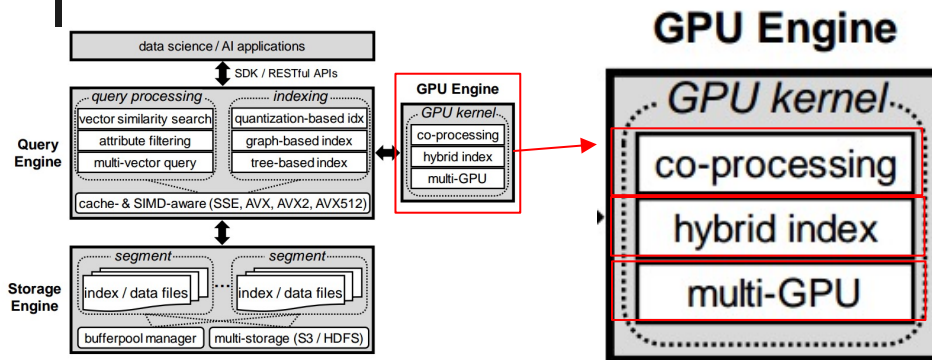


Figure 4: Different strategies for attribute filtering

\*Recall: Not supported by existing vector database

# The GPU Engine



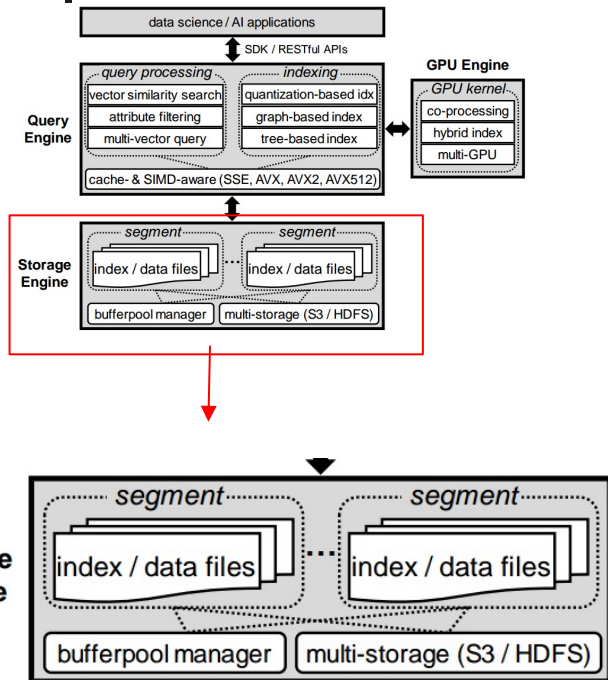
## Algorithm 1: SQ8H

```

1 let  $n_q$  be the batch size;
2 if  $n_q \geq \text{threshold}$  then
3   run all the queries entirely in GPU (load multiple buckets
   to GPU memory on the fly);
4 else
5   execute the step 1 of SQ8 in GPU: finding  $n_{probe}$  buckets;
6   execute the step 2 of SQ8 in CPU: scanning every relevant
   bucket;
  
```

- Optimized SQ8 algo [3] for CPU-GPU co-processing
  - Communication opt. by multi-bucket copying
  - Computation opt. by fine-grained decomposition
- Optimized GPU kernel to supported larger query within GPU memory
  - Progressively, round-by-round flushing results, instead of loading the entire and flushing only once
- Multi-GPU support by adjusting device number in runtime

# The Storage Engine



- Support various storage for different data structures and storage hardware
  - Vector storage (mem)
    - Continuous row for single vector
    - Columnar fashion for multiple vectors
  - Attribute storage (mem)
    - <key,value> pairs for each attribute column
    - Further optimize by building skip pointers[16]
  - Buffer pool to bridge mem and disk
  - Various file system support for in-disk and in-network

# Implementation Optimizations

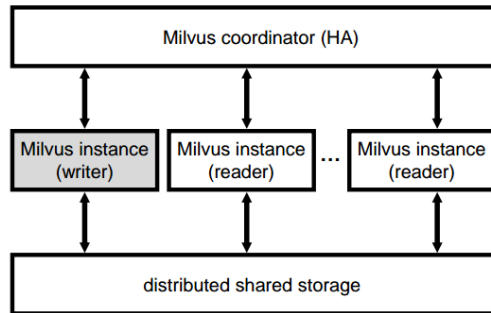
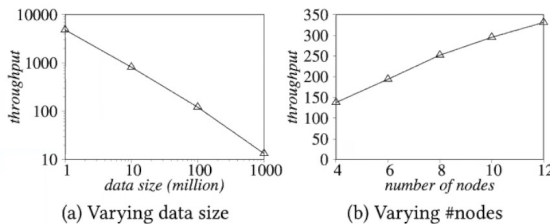


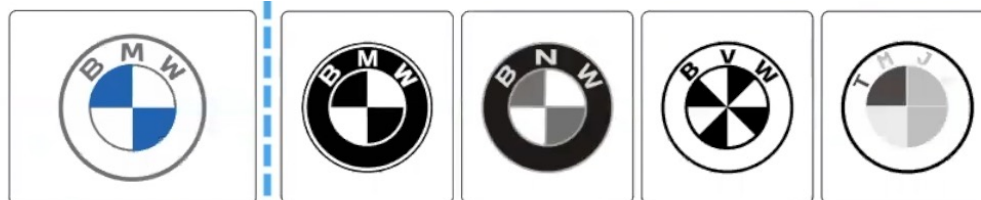
Figure 5: Milvus distributed system

- Asynchronous processing
  - To minimize foreground processing (e.g., receiving request)
  - Write-ahead logs + operating at back ground
- Snapshot isolation
  - To make sure of the R/W consistency
  - Using LSM structure[47]
- Distributed system
  - For scalability and availability
  - Shared-storage and separated compute/storage
  - Stateless computing instances
    - Single writer multiple reader as Milvus is read-heavy



# Application Examples

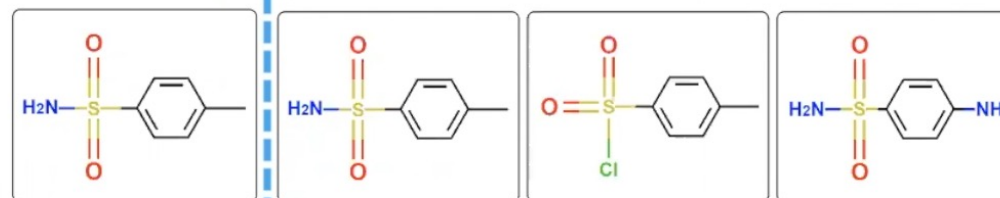
## Similar trademark search



## Floor plan search



## Chemical Structure Search



Cc1ccc(cc1)S(=O)(=O)N

Cc1ccc(cc1)S(=O)(=O)N

Cc1ccc(cc1)S(=O)(=O)Cl

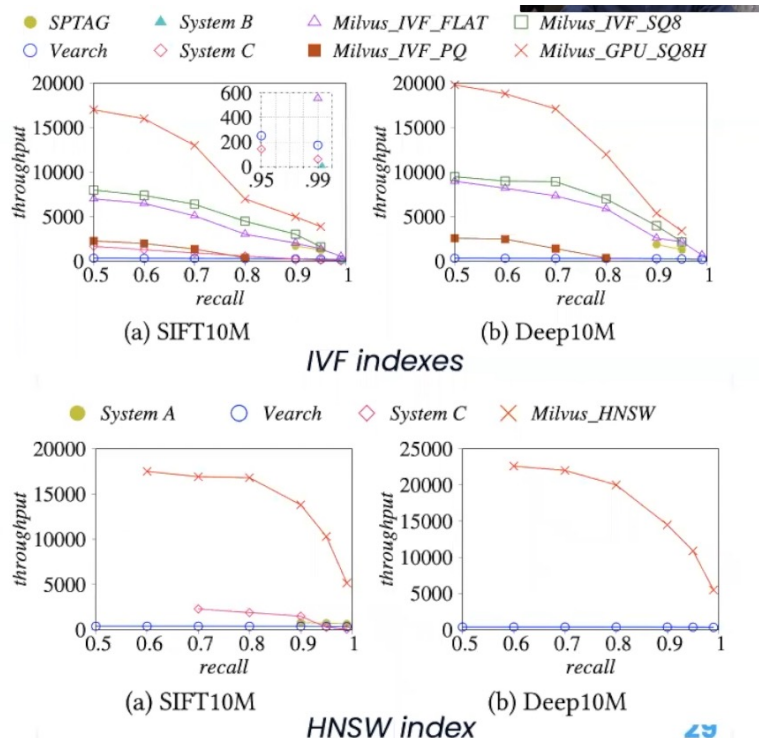
Nc1ccc(cc1)S(=O)(=O)N

Input

Output

Video search, chatbots,  
recommendation engine, biological  
multi-factor authentication...

# Outperform SOTA Significantly



- Platform: Alibaba cloud
- Dataset: SIFT1B and Deep1B
  - Both are 1 billion, 128-dimension
- Open source SOTA: SPTAG, Vearch
- Commercial system: System A, B, C

# Conclusion and Future Work

- Milvus satisfies the requirements of versatile applications by heuristic optimization in system and algorithm

System:

- SIMD support
- GPU accelerations
- Snapshot isolation
- Cache-aware optimization
- LSM structure
- Distributed processing

Algorithm:

- Hybrid search algorithm
- Cost base optimization
- Vector fusion
- Iterative merging

- Future work can be following:
  - Cloud native architecture
  - FPGA acceleration



# What Can We Learn and Do?

- This is an extremely complete system research.
  - Inclusively enough for wide community use

System:

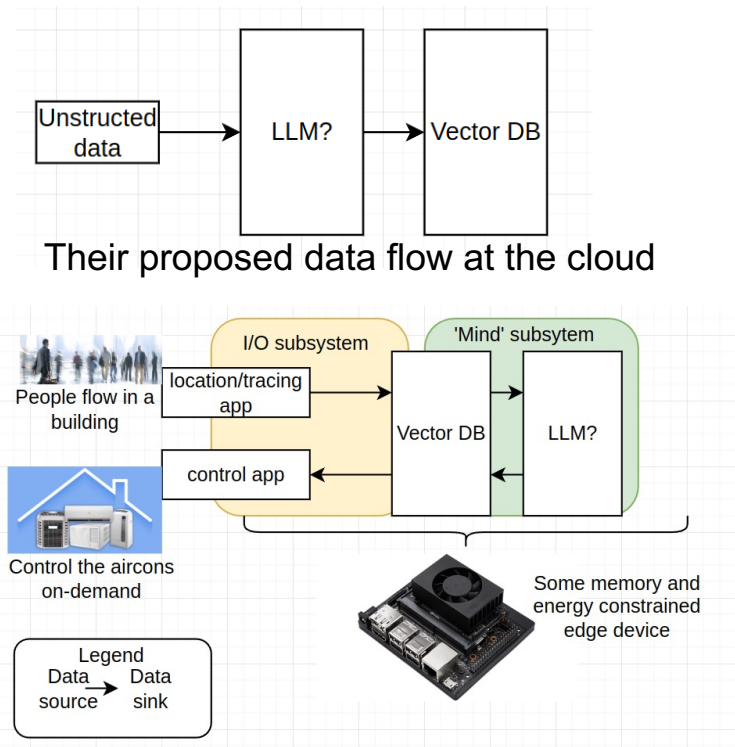
- SIMD support
- GPU accelerations
- Snapshot isolation
- Cache-aware optimization
- LSM structure
- Distributed processing

Algorithm:

- Hybrid search algorithm
- Cost base optimization
- Vector fusion
- Iterative merging

But it is not a good idea to apply Milvus on Edges...

# What if On Edge?



One possible edge application

- Our toy edge application
  - Continuously tracing people flow and control the aircons on demand
  - VectorDB bridges two subsystems on a device
- Problem 1: The indexing and storage overhead is still too heavy for constrained edges
  - Is there any cure by using approximate computation, or even new indexing ways?
- Problem 2: Two subsystems may raise high-speed and concurrent write (streams)
  - Will stateless, single writer still work well?