

Open Source AceCon

2021 智能云边开源峰会

AI x Cloud Native x Edge Computing

人工智能 × 云原生 × 边缘计算

Milvus: 探索云原生的向量数据库

郭人通

Partner at Zilliz

About Me

郭人通

兴趣领域：

分布式系统、数据库、异构计算

Milvus 系统架构师

CCF 分布式计算与系统专委会委员



计算机软件与理论博士

合伙人 & 架构师



- 01 What is Milvus
- 02 Real-world cases
- 03 Milvus2.0 Architecture Overview
- 04 RoadMap & Work in progress

01

What is Milvus?

Milvus is an open source vector database for Unstructured Data Search and Analytics

Unstructured Data Trend

Int, float,
string, ...

text

json

image
video
audio

domain specific

0 1 2 3 4
5 6 7 8 9

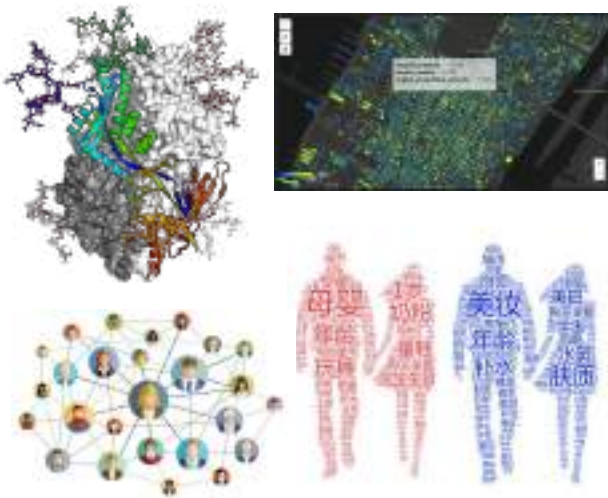
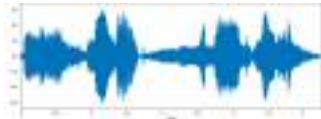
e π

ABCDEFG

2021.04.10

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 bulked 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.

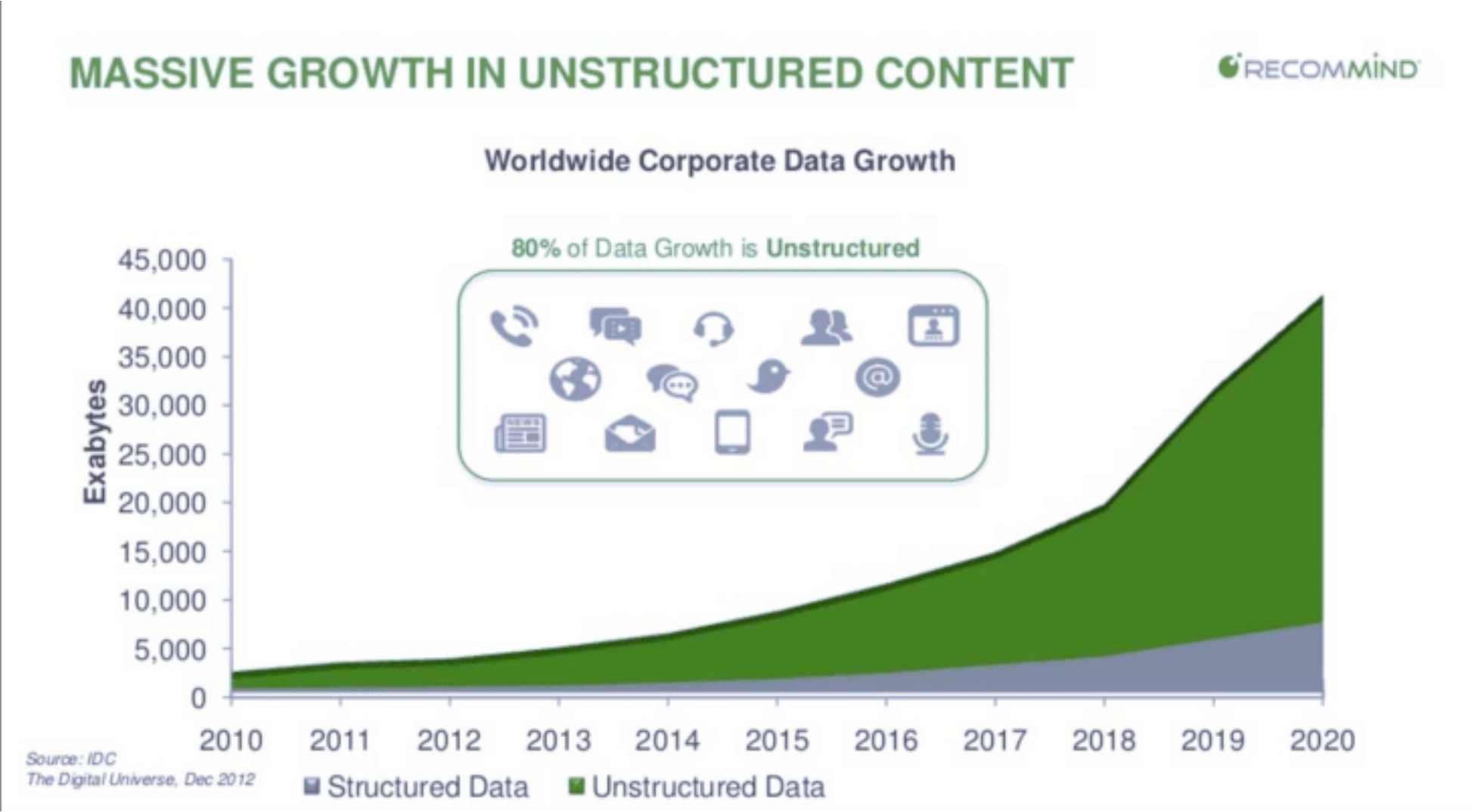
```
{  
  "firstName": "John",  
  "lastName": "Smith",  
  "isActive": true,  
  "age": 37,  
  "address": {  
    "streetAddress": "21 2nd Street",  
    "city": "New York",  
    "state": "NY",  
    "postalCode": "10011-3298"  
  },  
  "phoneNumbers": [ ]  
  {  
    "type": "home",  
    "number": "212 555-1234"  
  },  
  {  
    "type": "office",  
    "number": "212 555-4321"  
  },  
  "children": [ ],  
  "spouse": null  
}
```



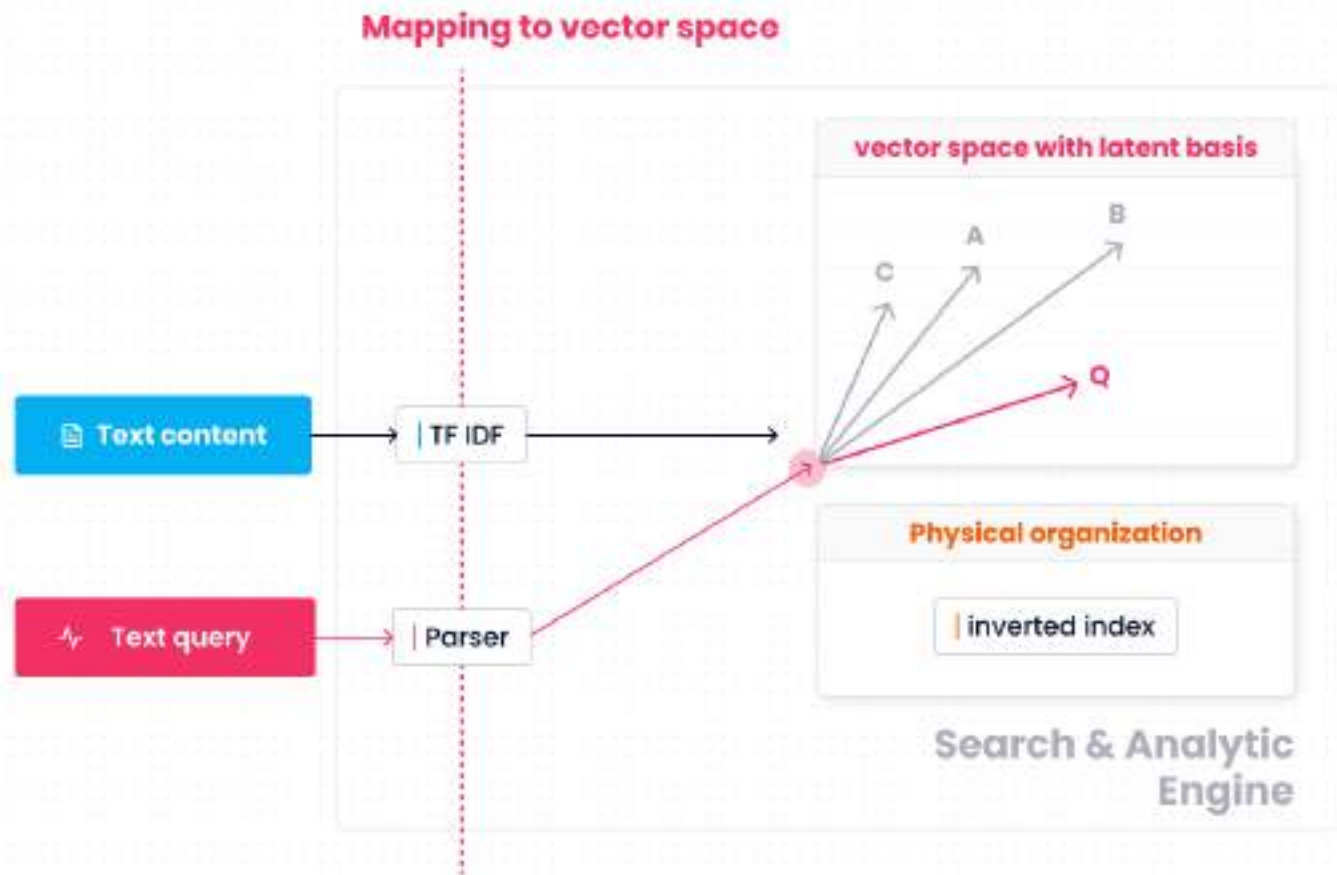
Structured data

Unstructured data

Unstructured Data Growth

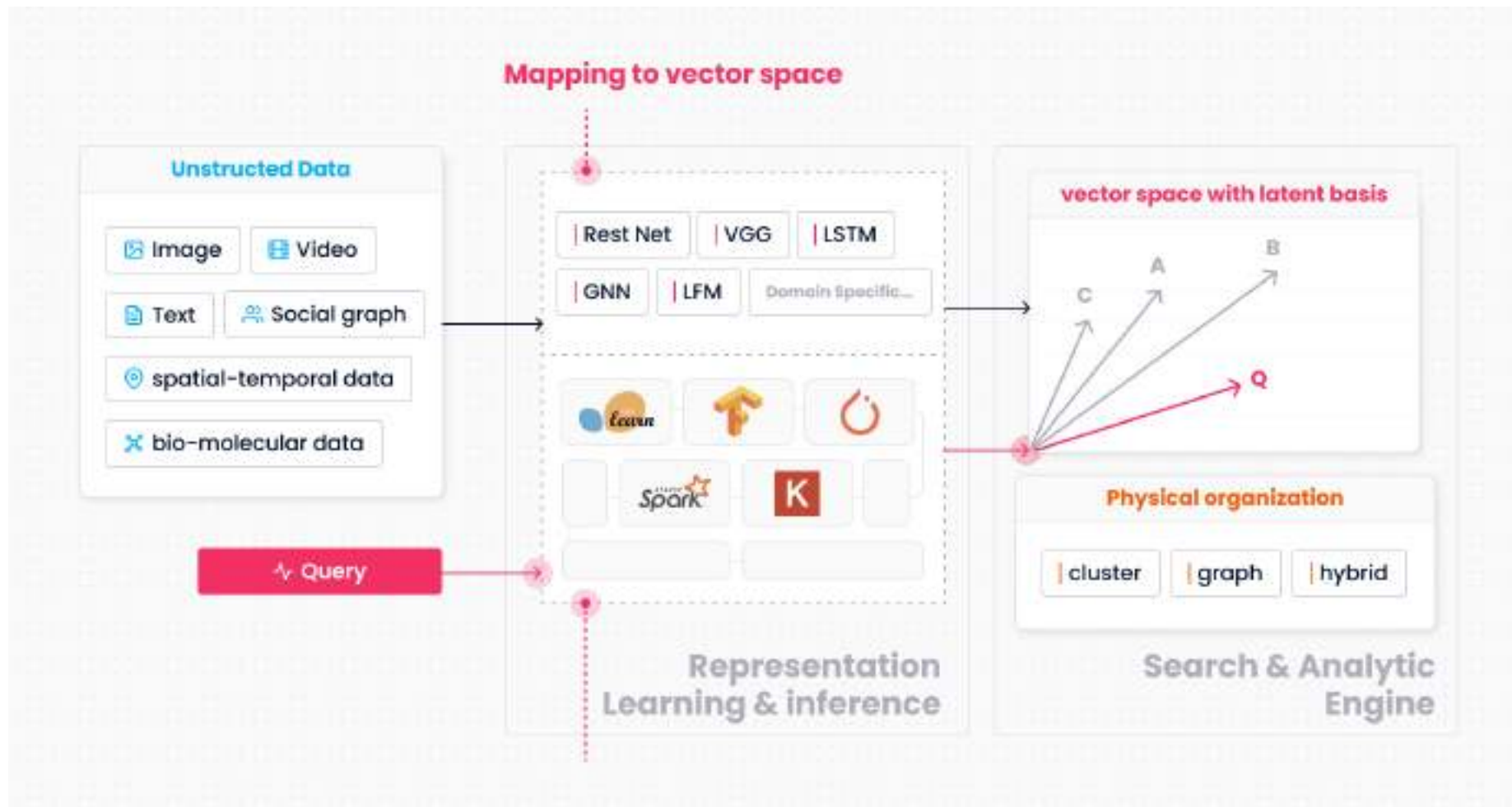


AI-driven Data Search and Analytics



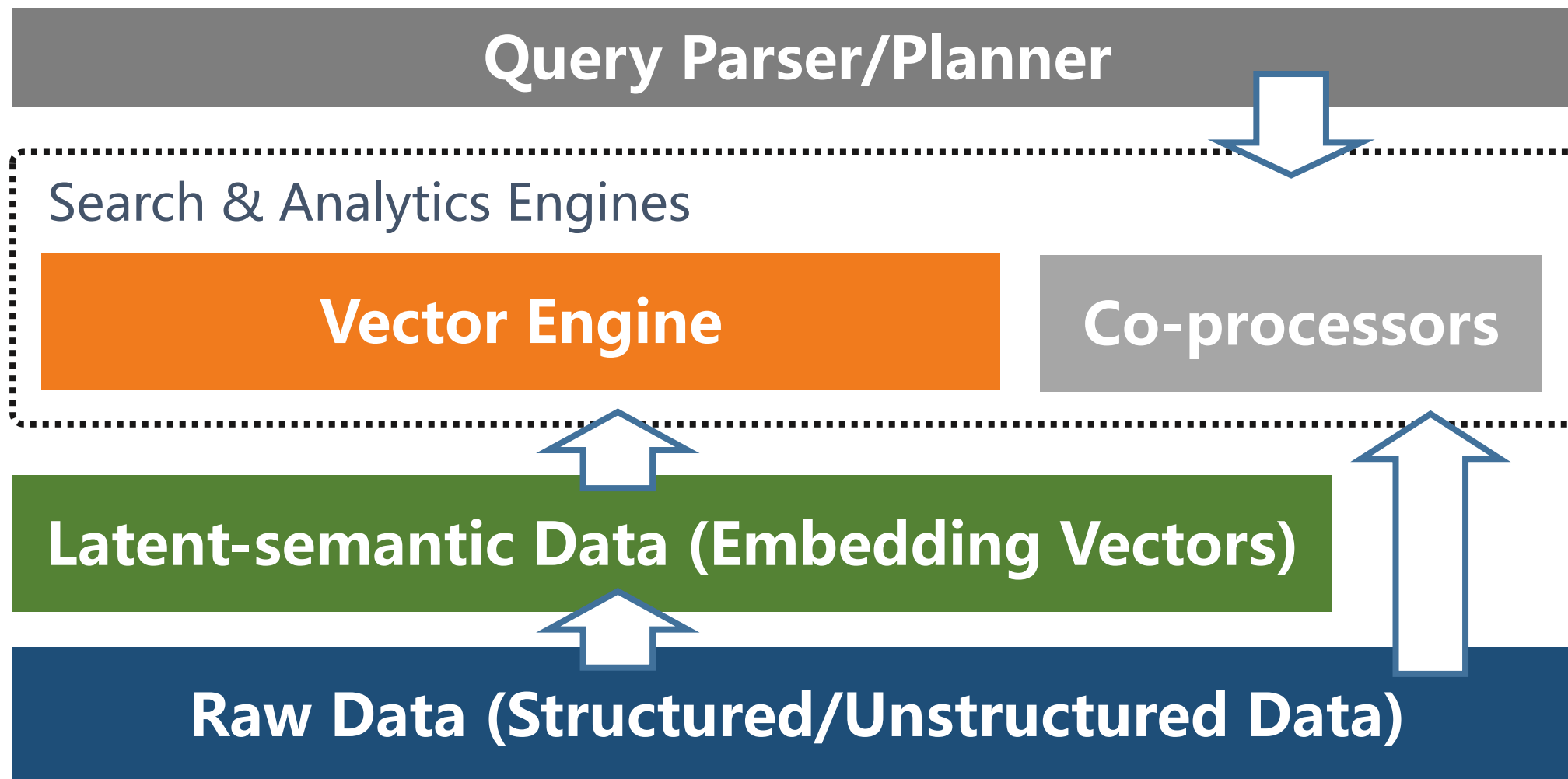
ElasticSearch
Perspective

AI-driven Data Search and Analytics



Milvus
Perspective

The Big Picture



Milvus is an LF AI & Data Graduated Project

LF AI & Data Foundation Interactive Landscape


The LF AI & Data Foundation landscape (png, pdf) is dynamically generated below. It is modeled after the CNCF landscape and based on the same open source code. Please [open](#) a pull request to correct any issues. Greyed logos are not open source. Last Updated: 2021-09-07 23:04:31Z

You are viewing 7 cards with a total of 42,422 stars.

[Landscape](#) [Card Mode](#) [LF AI & Data Members](#) [Companies Hosting Projects](#)

[Tweet](#) 151


Graduated LF AI & Data Projects (7)



Acumos

★ 17

LF AI & Data Foundation




Angel

Angel-ML

★ 5,279

LF AI & Data Foundation




EGERIA

Egeria

★ 565


LF AI & Data Foundation



Horovod

★ 11,333


LF AI & Data Foundation



Milvus

★ 6,445


LF AI & Data Foundation



ONNX

★ 10,808

LF AI & Data Foundation



Pyro

★ 6,975

LF AI & Data Foundation

Crunchbase data is used under license from Crunchbase to CNCF. For more information, please see the [license](#) info.

02

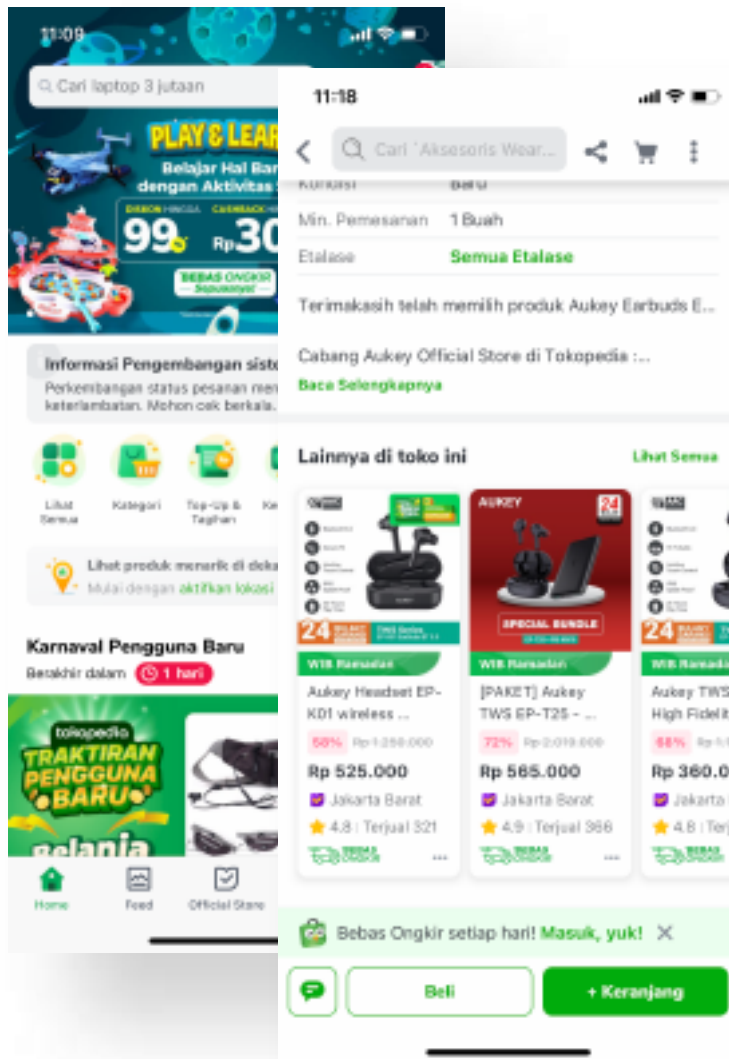
Real-World Cases

Users

1000+ Enterprise users around the global



Related product search and recommendations



ASCII representation

bread	b	r	e	a	d
	098	114	101	097	100
toast	t	o	a	s	t
	116	111	097	115	116

Elasticsearch for keyword search with ASCII codes.

The code we know about these two arrays of numbers is that bread not equal to toast.

We assume that similar contexts represent similar things, and try to compare them using mathematical methods. We could even find a way to encode whole sentences by their meaning.

Vector representation

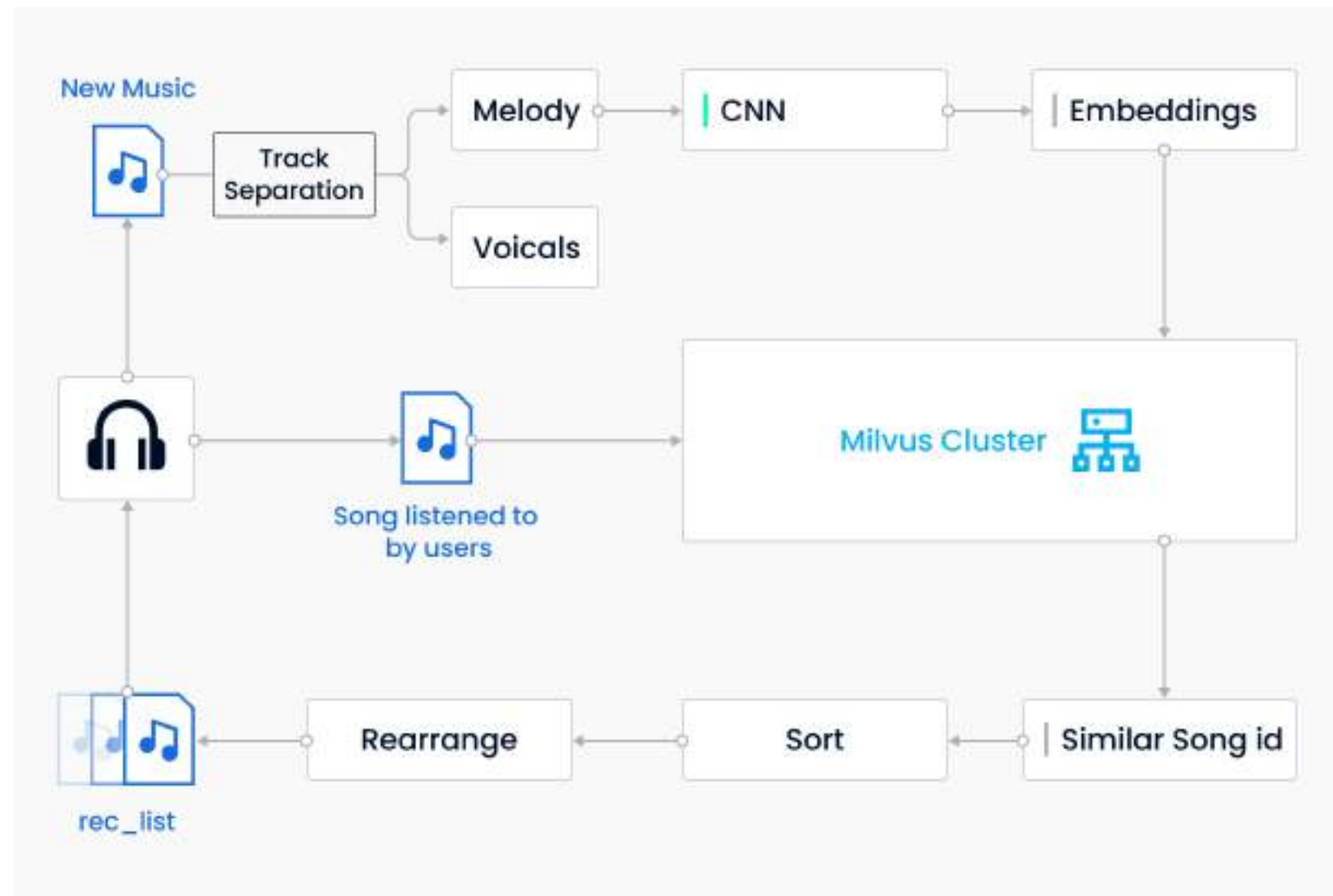
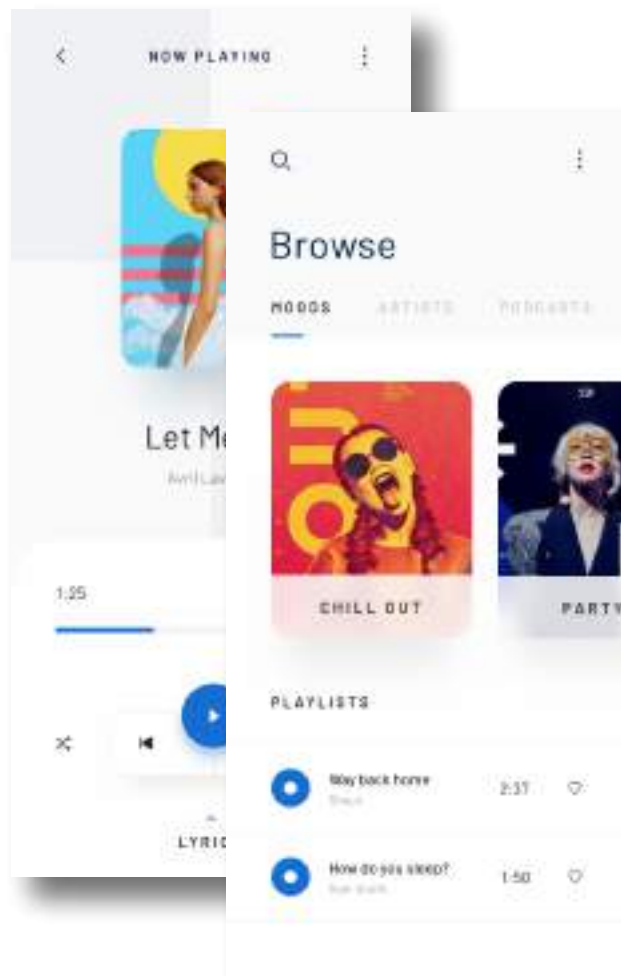
	butter	car	sun	brekfast	cat
bread	0.80	0.11	0.05	0.93	0.20...
toast	0.76	0.22	0.15	0.95	0.12...

The word bread is often used together with words butter, breakfast, and rarely with words car, sun and cat. So, is the toast.

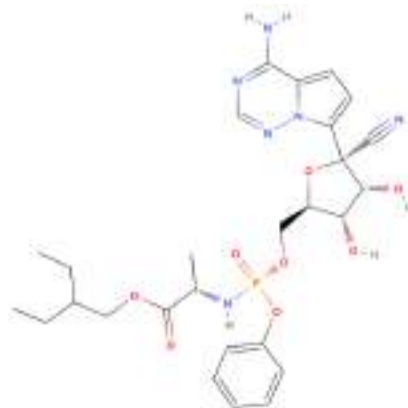
Reverse Image Search



Music Recommendations

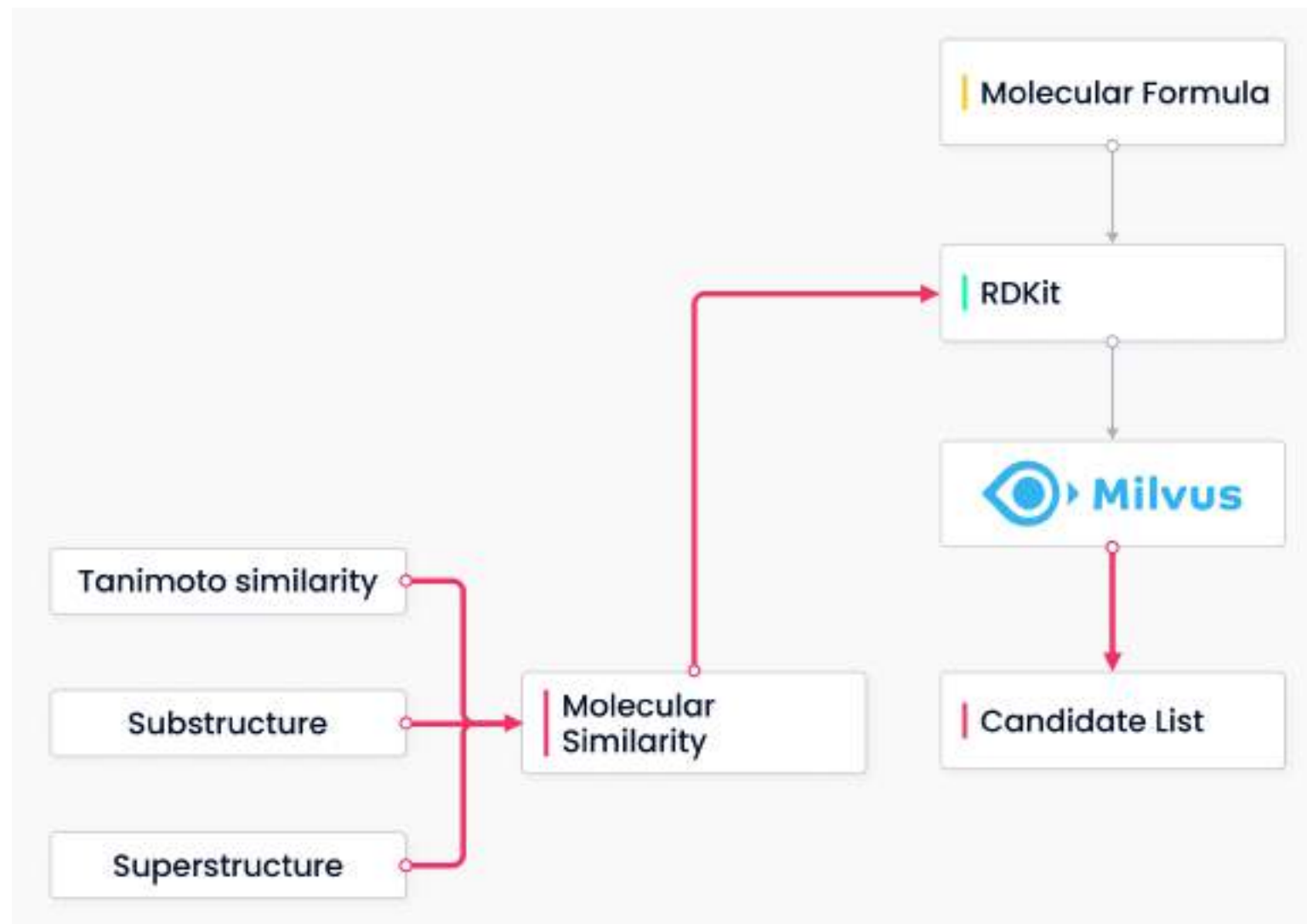


Pharmaceutical Molecular Analysis



CC(=O)Nc1ccc(S(=O)(=O)NCC(=O)N2CCS(=O)CC2cc1

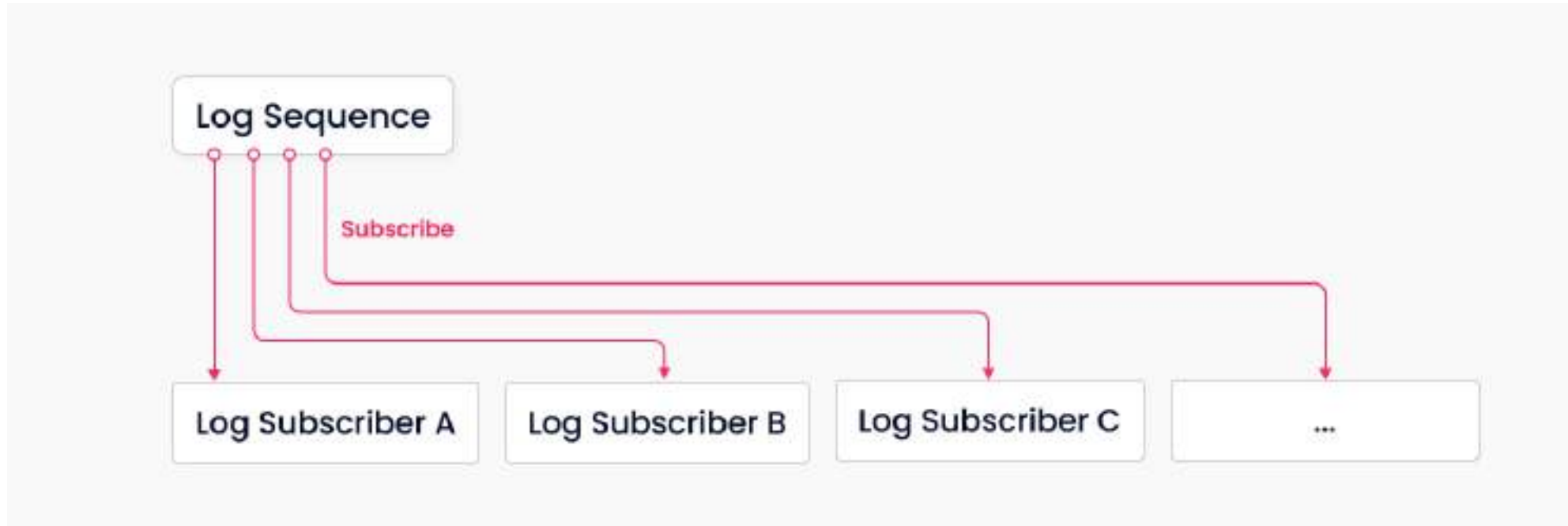
Molecular fingerprint: 1024 bits
00001100...10000000



03

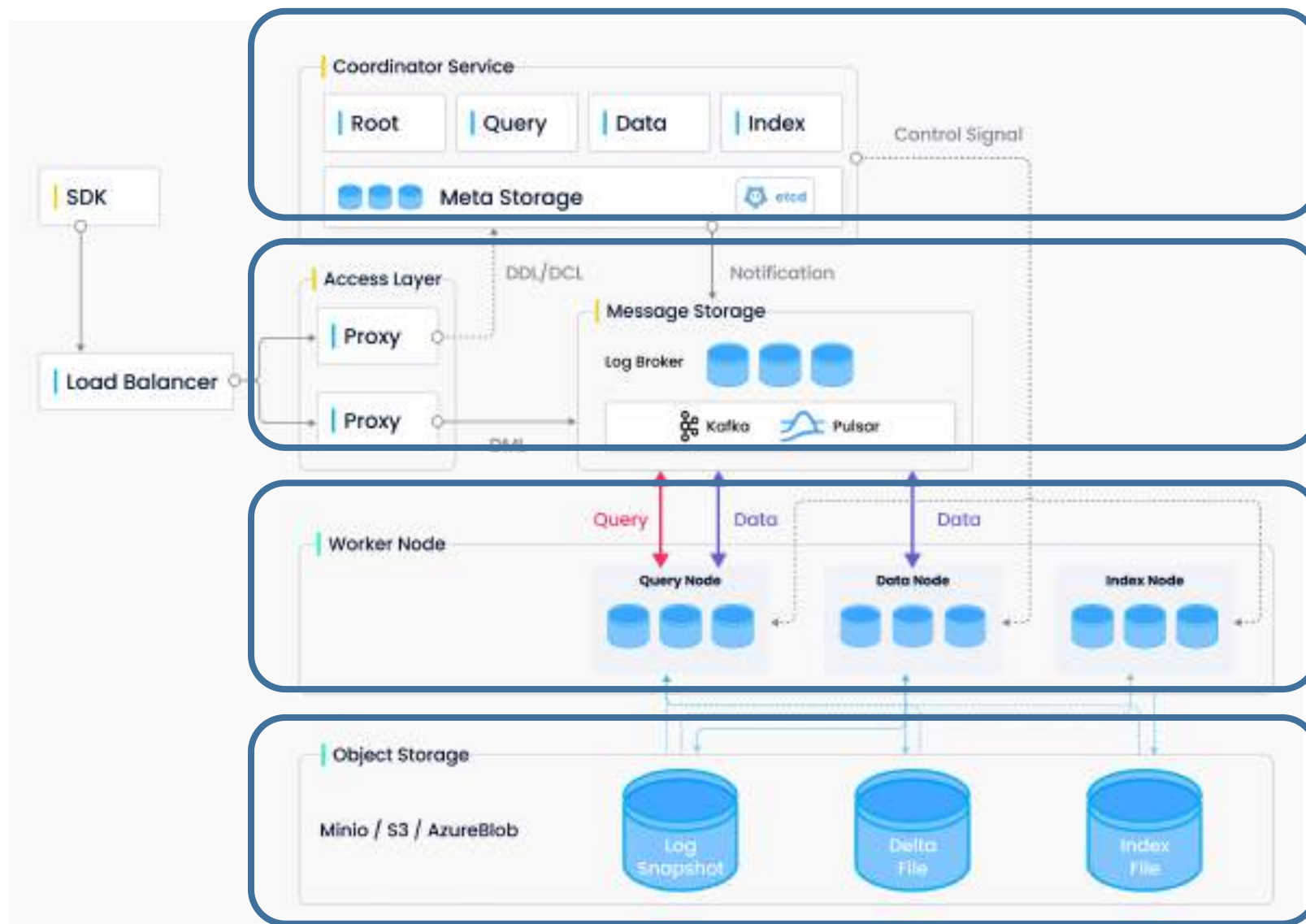
Milvus2.0 Architecture Overview

Log Sequence Pub-sub as System Backbone



- Disaggregate Log and database, make failure recovery easy and fast
- Guarantee data durability
- Make System extendable
- Reduce system complexity

Architecture



Coordinators

Log Broker

Log Subscribers

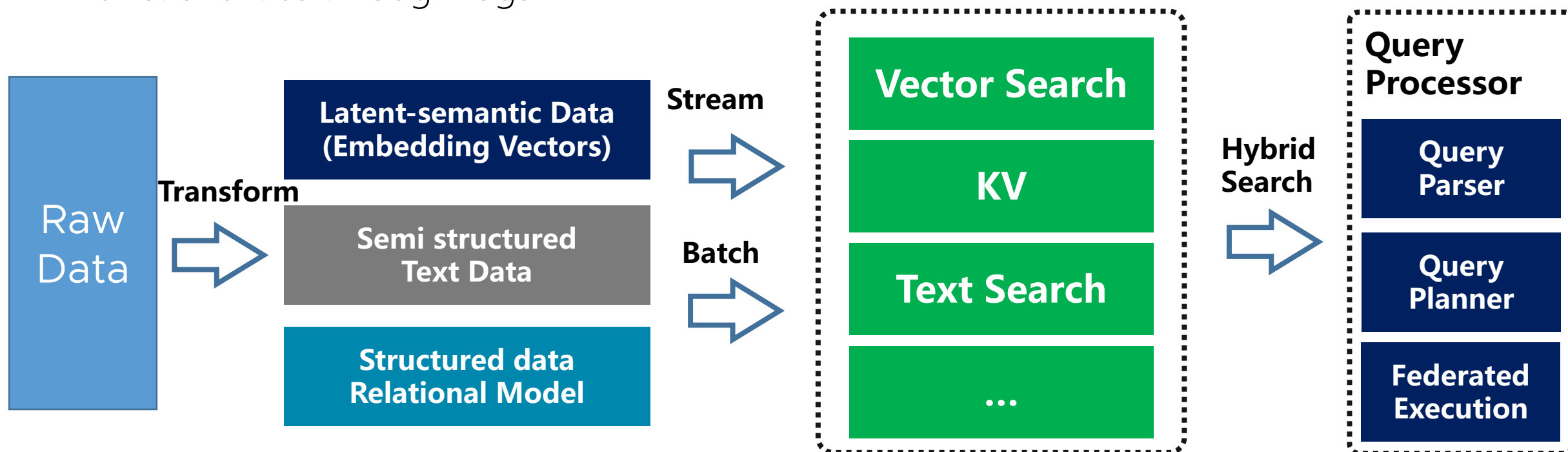
Storage

The Bazaar Architecture

How to further **extend** the system?

Solution

The 'bazaar' architecture. Loosely coupling multiple execution engines with different functionalities through logs.

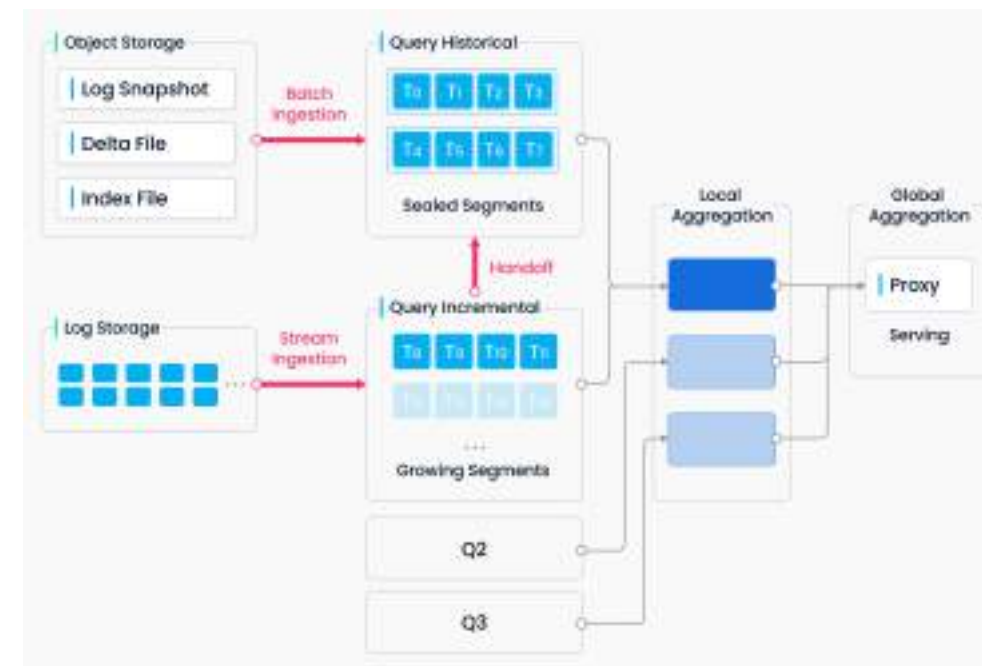
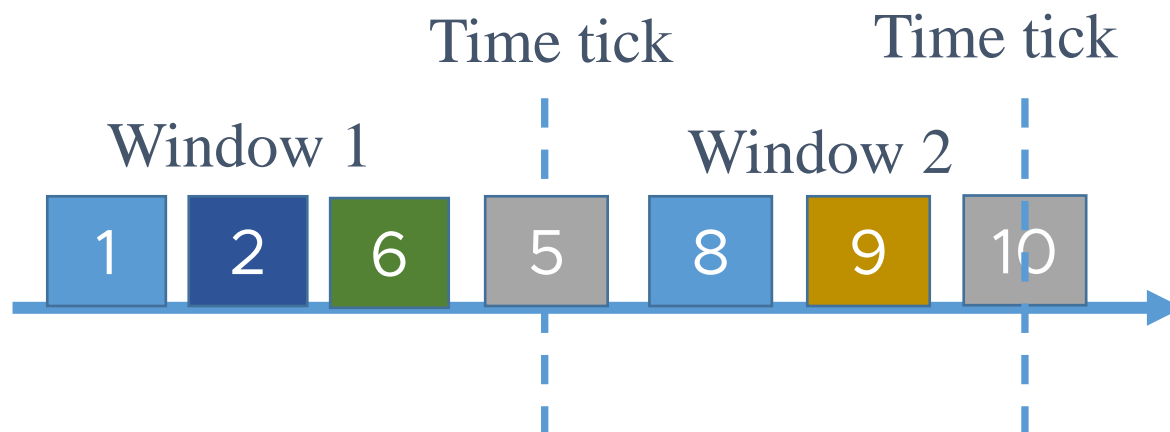


Combine Streaming and Batching

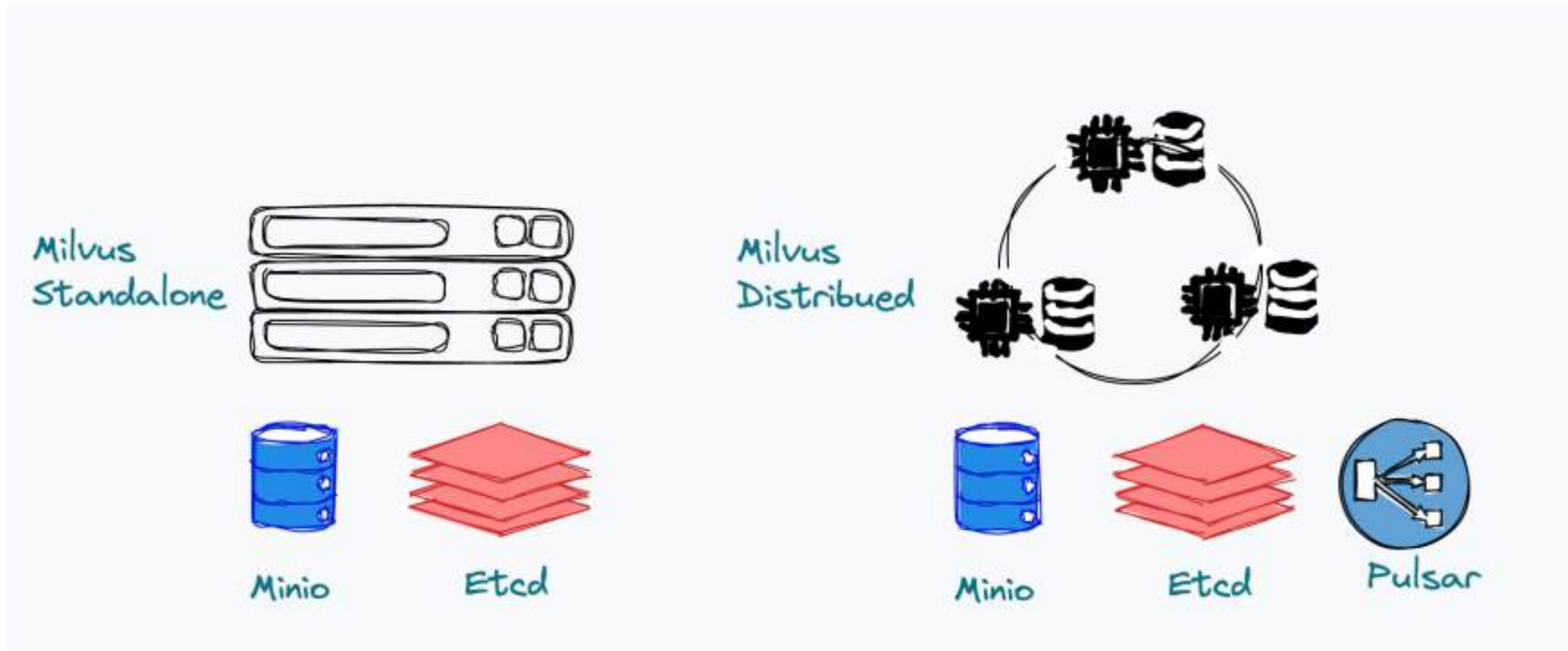
Relying only on log stream for reads is not practical (too slow)

Solution

Periodically backfill history data to segments, just like flush in LSM tree, and handoff growing segments to historical. Merge incremental and historical on Read to maintain data completeness



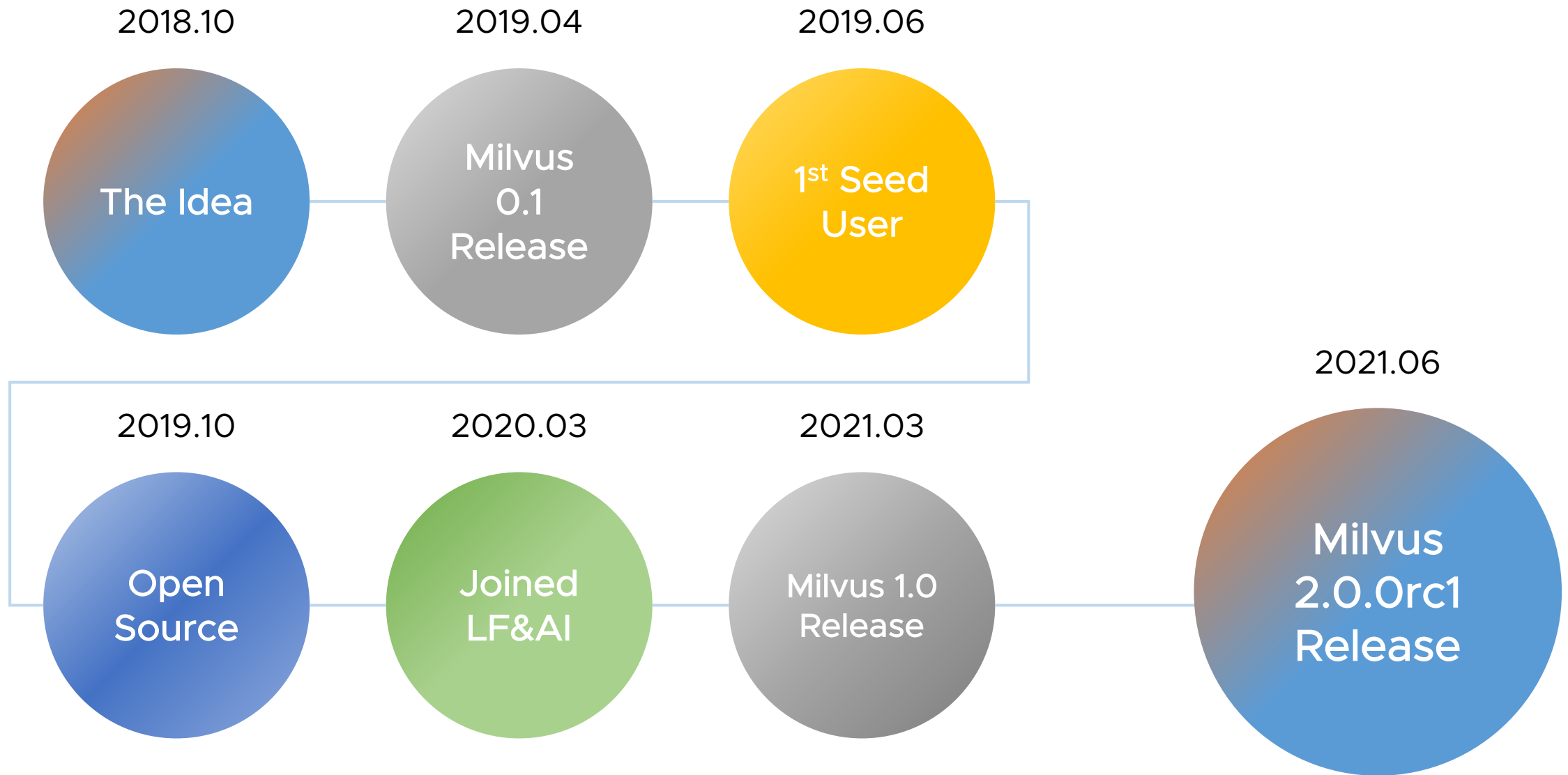
Milvus Deployment



04

RoadMap

History of Milvus



Milvus 2.0 GA will coming in Oct. 2021

Support String Data Type

Urgent needs from users

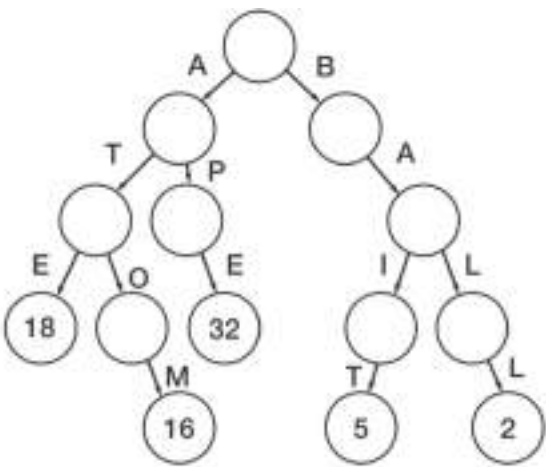
Scalar filtering on string field

Retrieve origin string

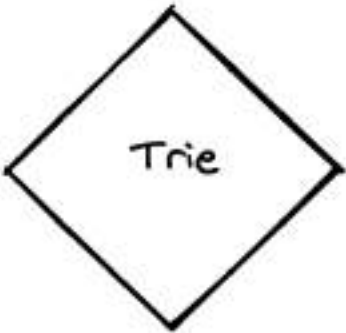
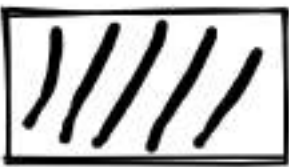
Memory consumption matters

3 million unicode russian words

DataStructure	MemoryUsage
Python-Dict	600MB
Python-List	300MB
PAT-Trie	242 MB
HAT-Trie	125 MB
DA-Trie	101 MB
Marisa-Trie	11 MB



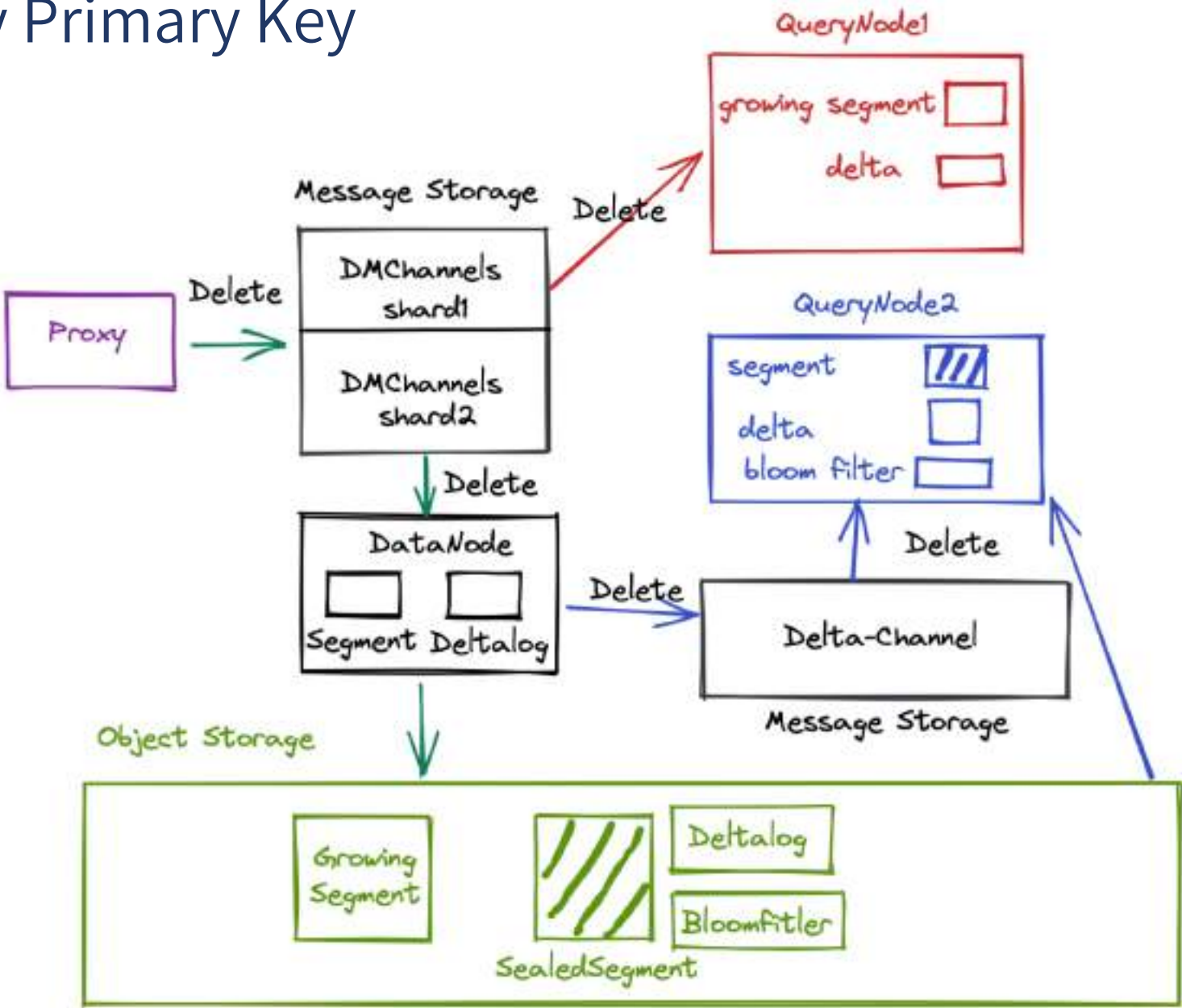
Sealed Segment



Growing Segment



Support Delete by Primary Key

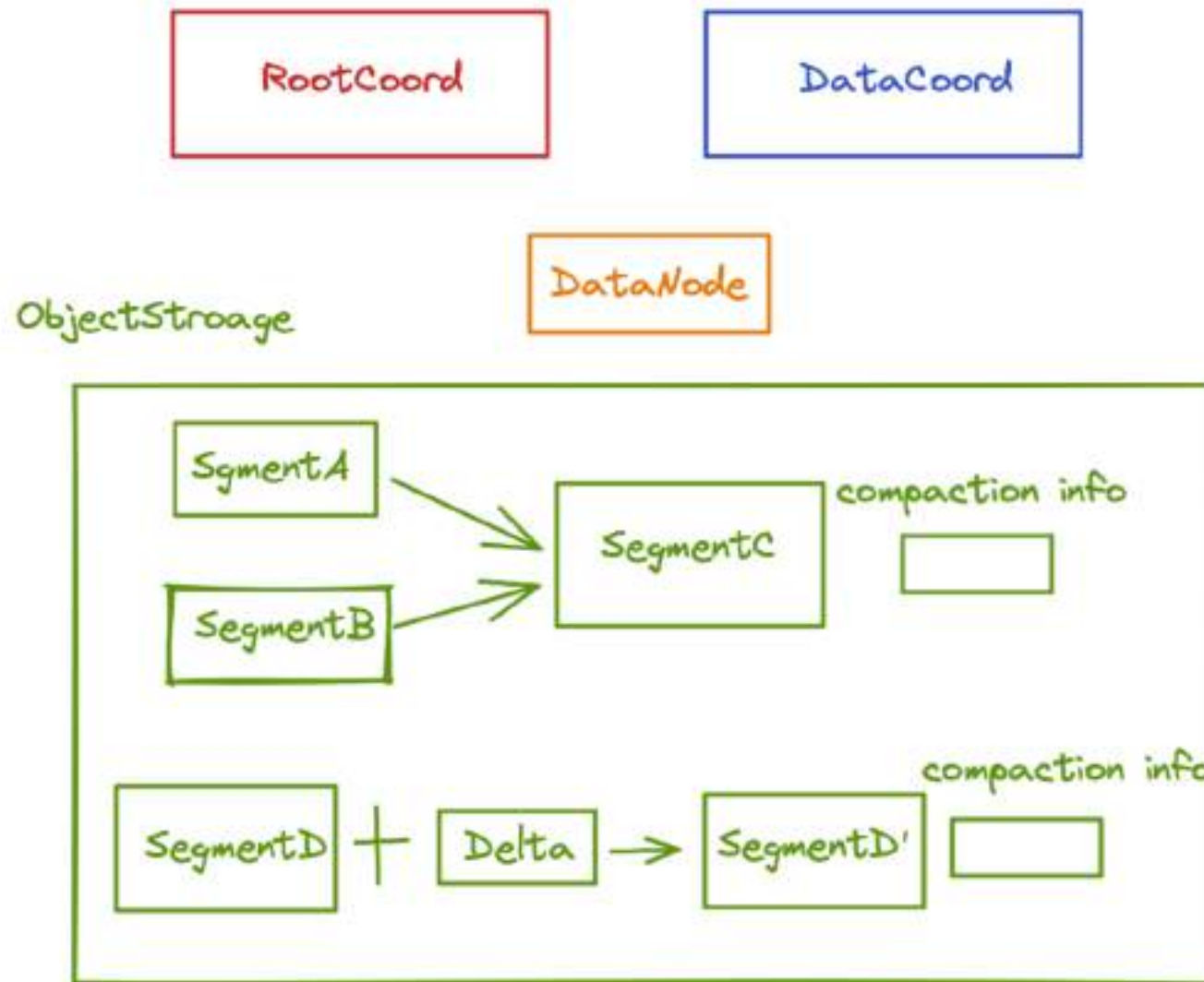


Segment Compaction

Segment size varies

Delete operation make
segment and index sparse

Indexes of large segments
are more efficient



Search/Query with Expression

Search

A set of criteria that results in a relevancy-ordered list that match the query.

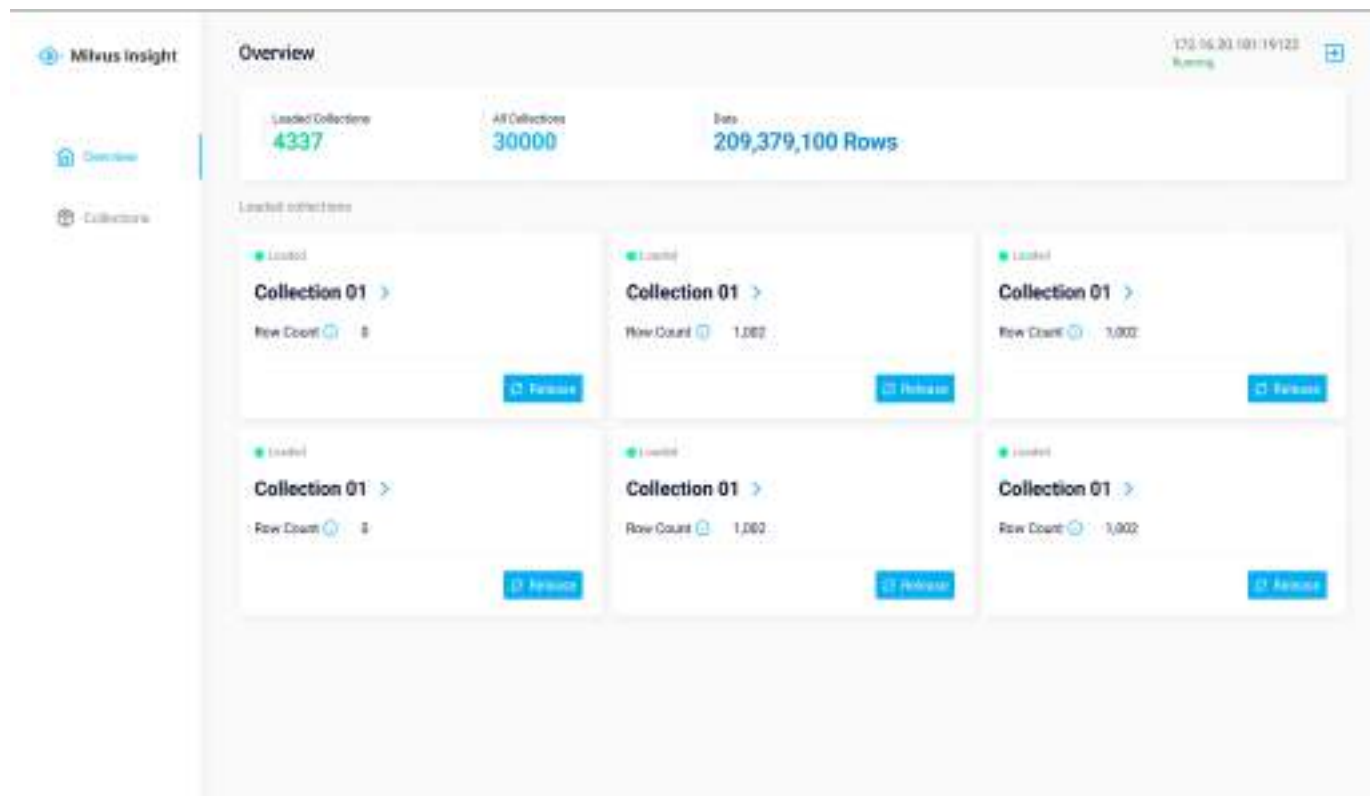
Query

A set of criteria that results in a list of records that match the query exactly, returned in order of particular field values

Operator	Description	Examples
Relational operators	Relational operators use symbols to check for equality, inequality, or relative order between two expressions. Relational operators include <code>></code> , <code>>=</code> , <code><</code> , <code><=</code> , <code>==</code> , and <code>!=</code> .	<ul style="list-style-type: none">• <code>A > 1</code>• <code>B >= 2</code>• <code>C < 3</code>• <code>D <= 4</code>• <code>E == 5</code>• <code>F != 6</code>
Logical operators	Logical operators perform a comparison between two expressions. The supported logical operators are: AND, && OR, , and NOT.	
IN operator	The IN condition is satisfied when the expression to the left of the keyword IN is included in the list of items.	<ul style="list-style-type: none">• <code>FloatCol in [1.0, 2, 3.0]</code>• <code>Int64Col in [1, 2, 3]</code>

Milvus combines scalar and vector search, such as “Find top 10 drama films similar to Forrest Gump”

- Support scalar datatypes columnar storage
- Filter scalar data by arithmetic and bool expressions
- Retrieve field data on query/search



- Cluster state visualization
- Meta Management
- Data Query
- Health Diagnosis
- Open source, Please join us



<https://github.com/milvus-io/milvus-insight>

Milvus CLI

```
milvus_cli > help
Usage: [OPTIONS] COMMAND [ARGS]...

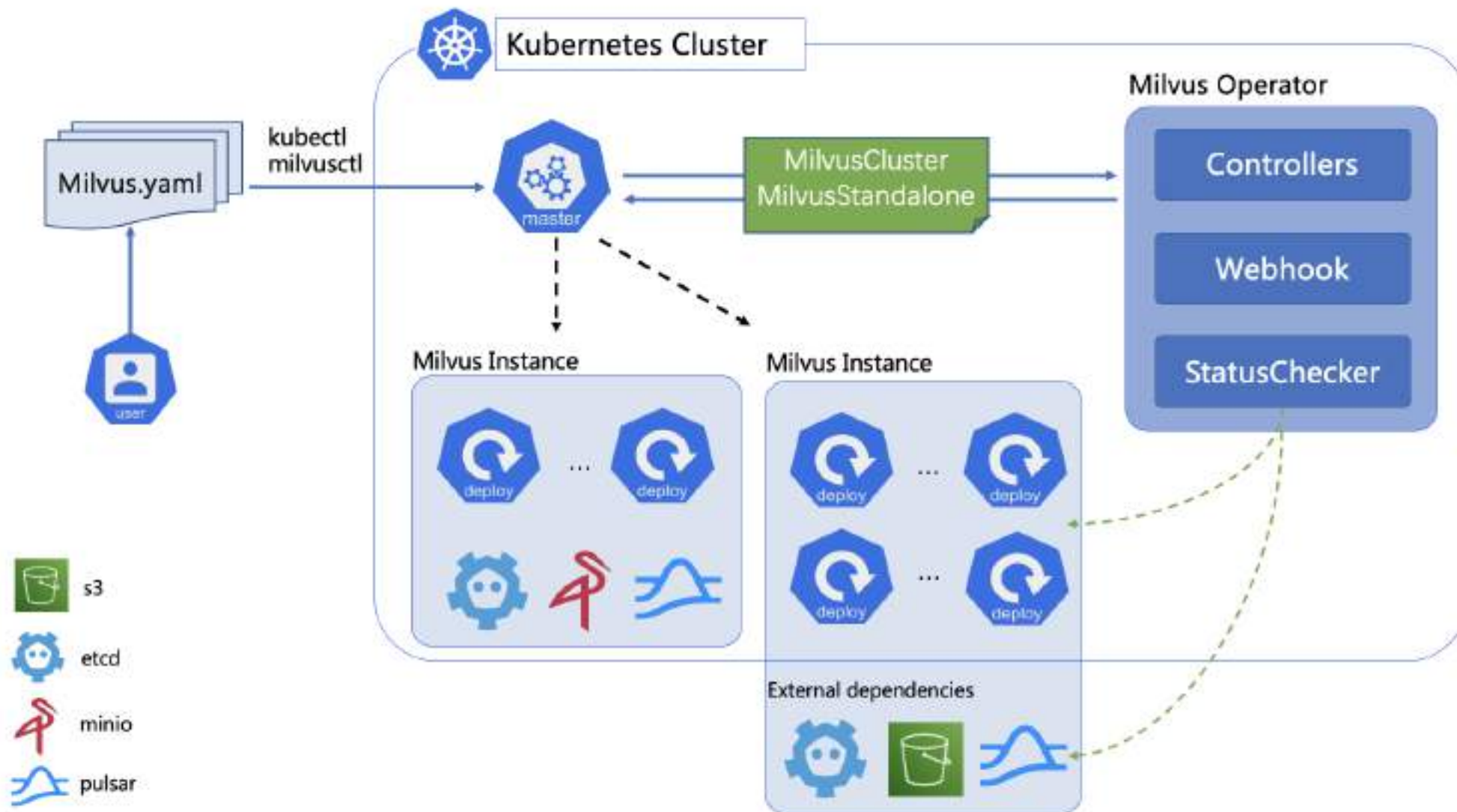
Milvus CLI

Commands:
  clear      Clear screen.
  connect    Connect to Milvus.
  create     Create collection, partition and index.
  delete     Delete specified collection, partition and index.
  describe   Describe collection or partition.
  exit       Exit the CLI.
  help       Show help messages.
  import     Import data.
  list       List collections, partitions and indexes.
  load       Load specified collection.
  query      Query with a set of criteria, and results in a list of...
  release    Release specified collection.
  search     Conducts a vector similarity search with an optional boolean...
  show       Show connection, loading_progress and index_progress.
  version    Get Milvus CLI version.
```



https://github.com/milvus-io/milvus_cli

K8s Operator



Open Source AceCon
2021 智能云边开源峰会
AI x Cloud Native x Edge Computing
人工智能 × 云原生 × 边缘计算

Thank You

