

Vearch --高性能、高可用的分布式向量数据库

孙延好

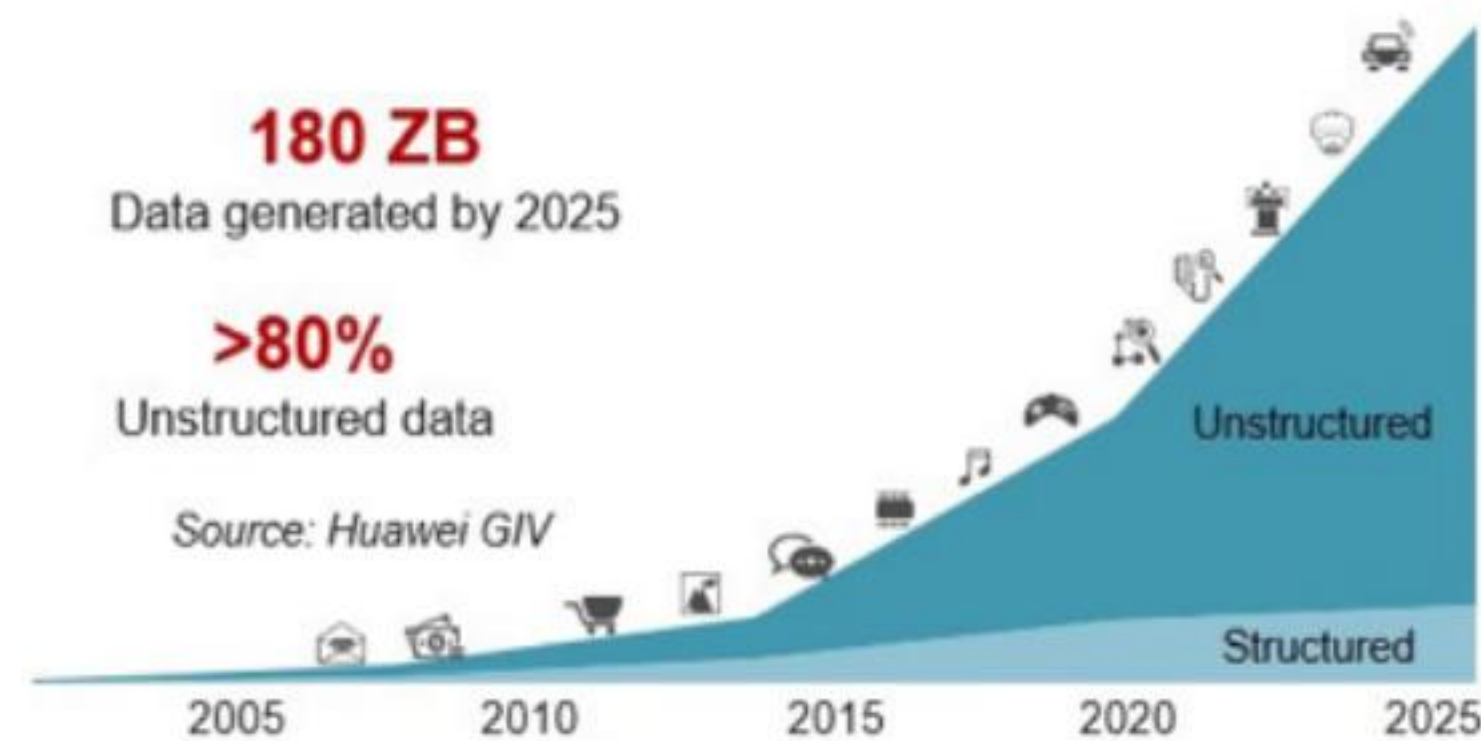


孙延好
京东后端技术部
架构师

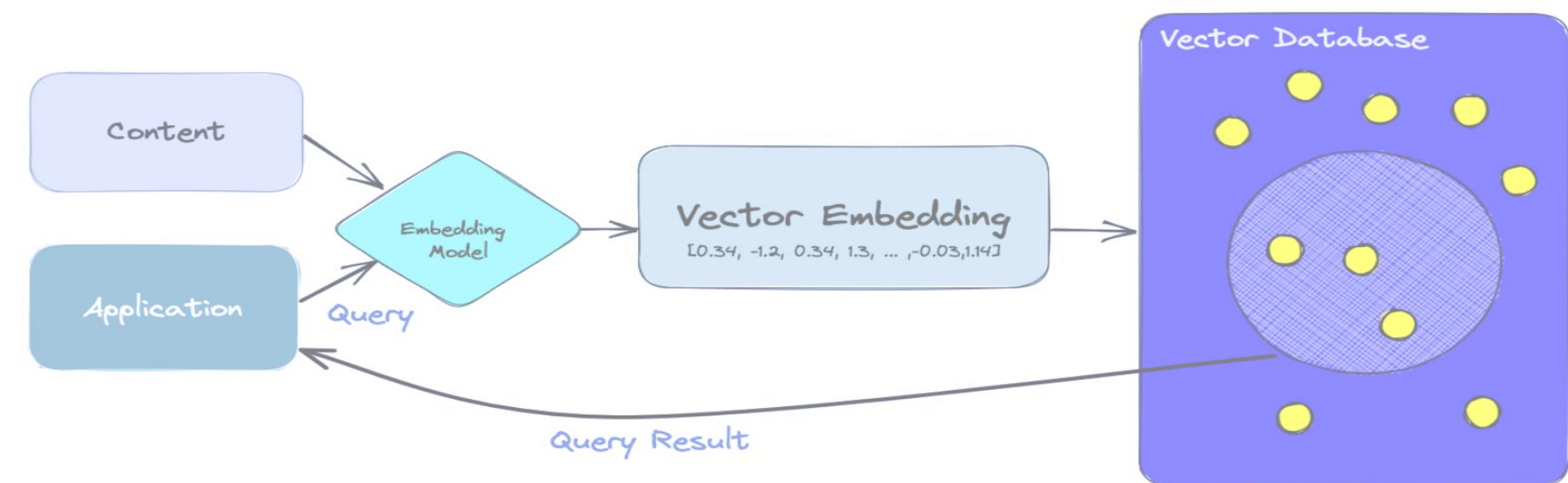
目录

- 向量数据库概述
- 向量数据库的发展历程
- 向量数据库的商业价值及在LLM中的应用
- Vearch的定位
- Vearch的整体架构
- Vearch的异构混合结构
- Vearch的应用情况

向量数据库概述



数据发展的趋势



非结构数据的使用情形

- **向量数据库**：是面向向量嵌入(vector embedding)而设计数据库系统,通过比较值查找比较相似的向量进行索引，从而进行搜索和分析

向量数据库的发展历程

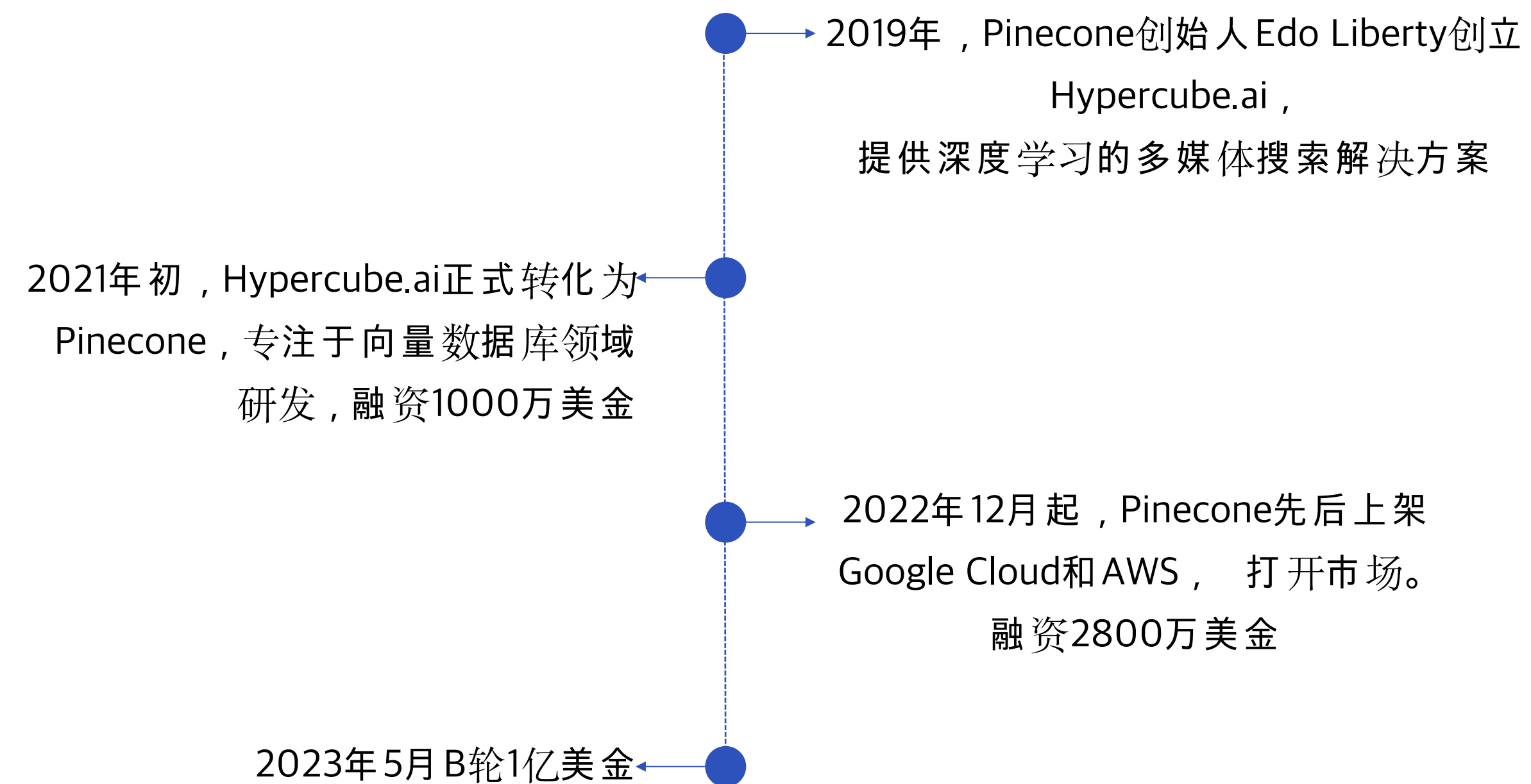


向量数据库的商业价值

Pinecone

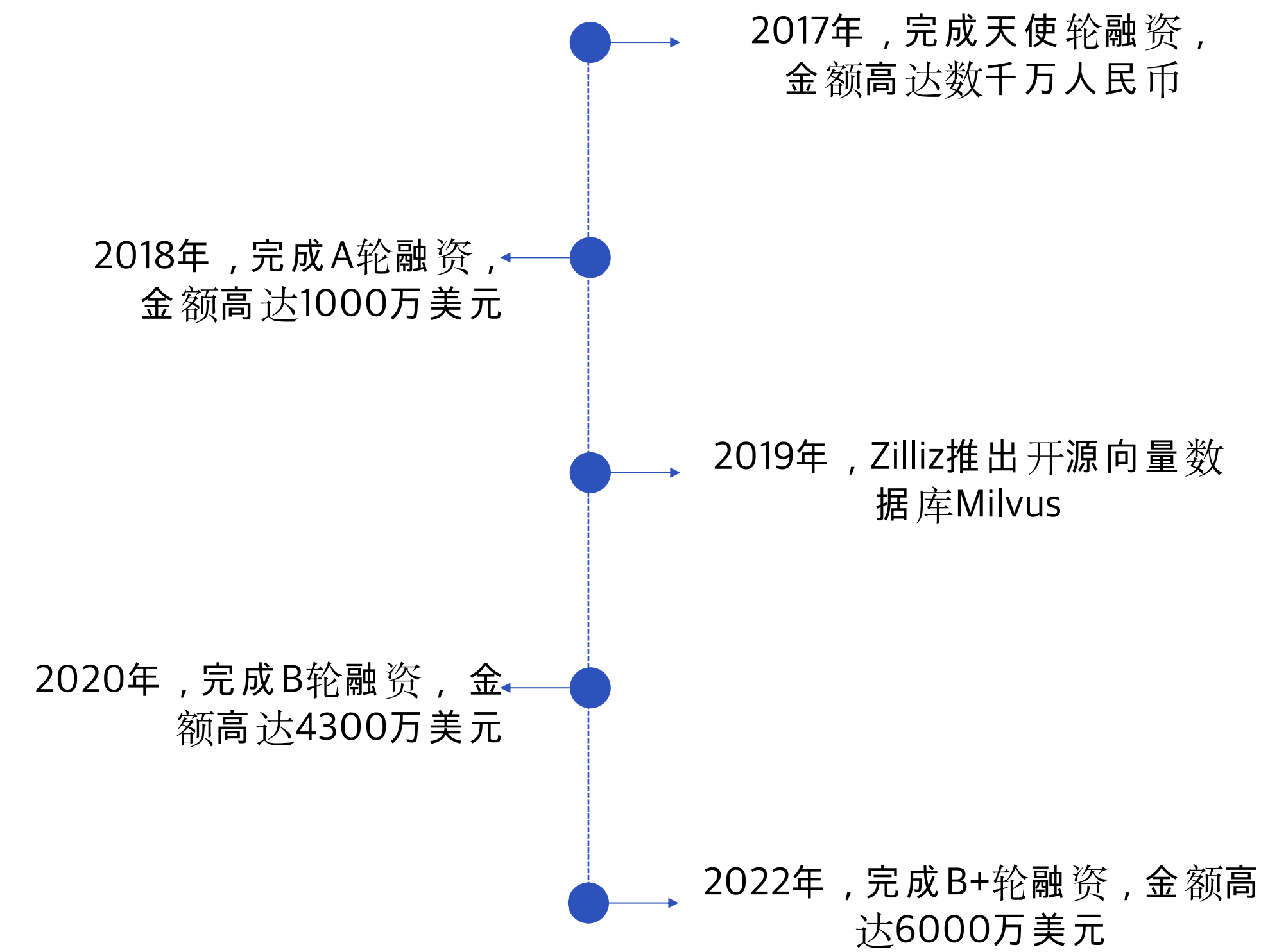
Pinecone总部位于纽约，专为OpenAI的GPT-4等大型语言模型(LLMs)提供长期记忆服务。Vector Search专注于通过AI生成的内容表示进行存储和搜索。

Pinecone是OpenAI、Cohere等LLM生成商的合作方



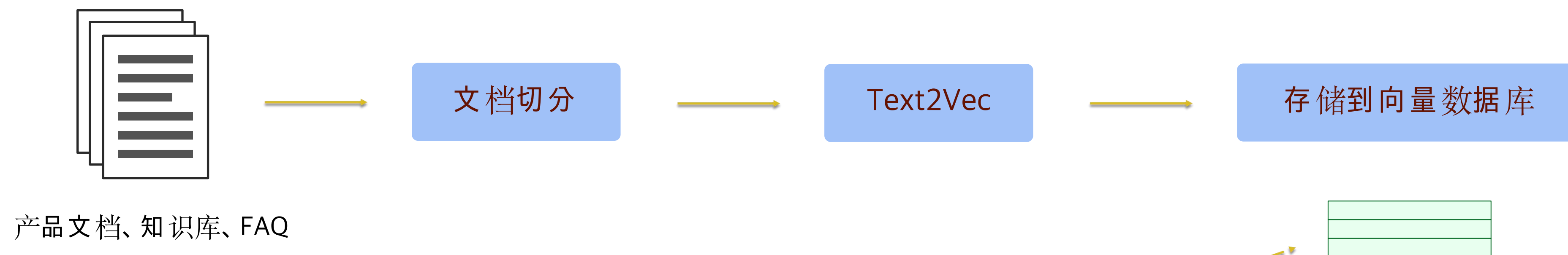
Milvus

Milvus的开发公司Zilliz成立于2017年，全球向量数据库领域的开拓者，以技术为核心，专注于在非结构化数据的分析中挖掘其价值，研发为人工智能服务的向量数据库系统，使更多的企业、组织、个人用更低的成本开发人工智能应用并从中获得便利。

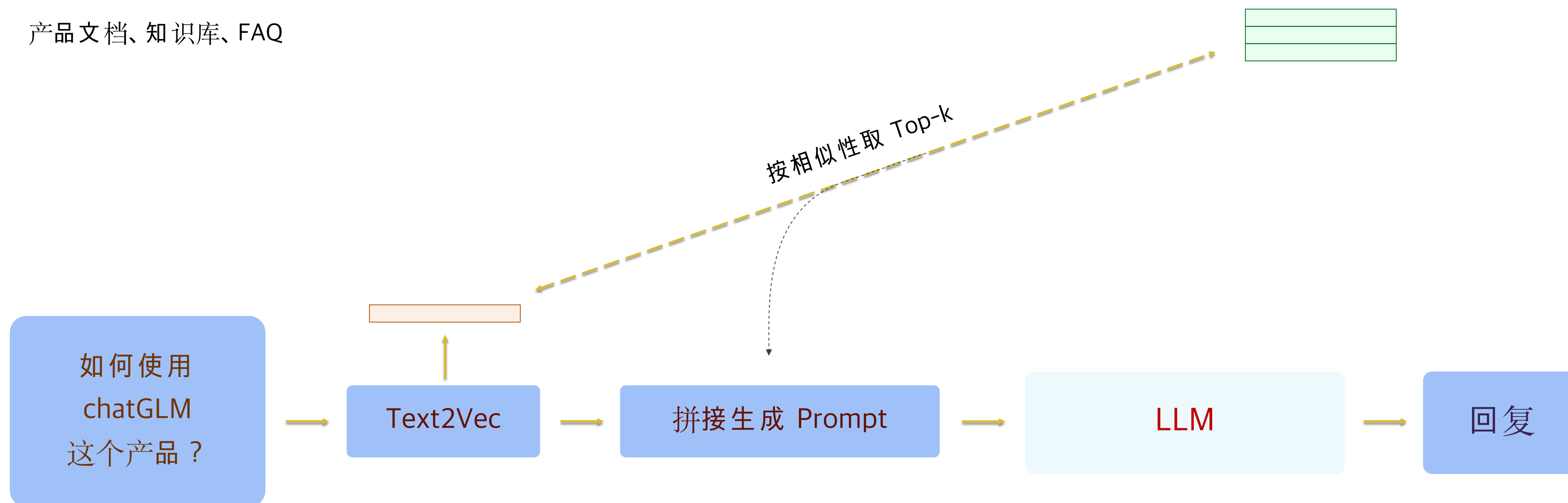


向量数据库在LLM中的应用

构建阶段



应用阶段



Vearch定位：海量向量高效存储和检索的一站式解决方案



高性能：千亿规模的向量检索毫秒级响应

高可用：经过多次大促流量洪峰锤炼

高可扩展：支持在线水平扩容或迁移

数据新鲜：数据写入实时可见

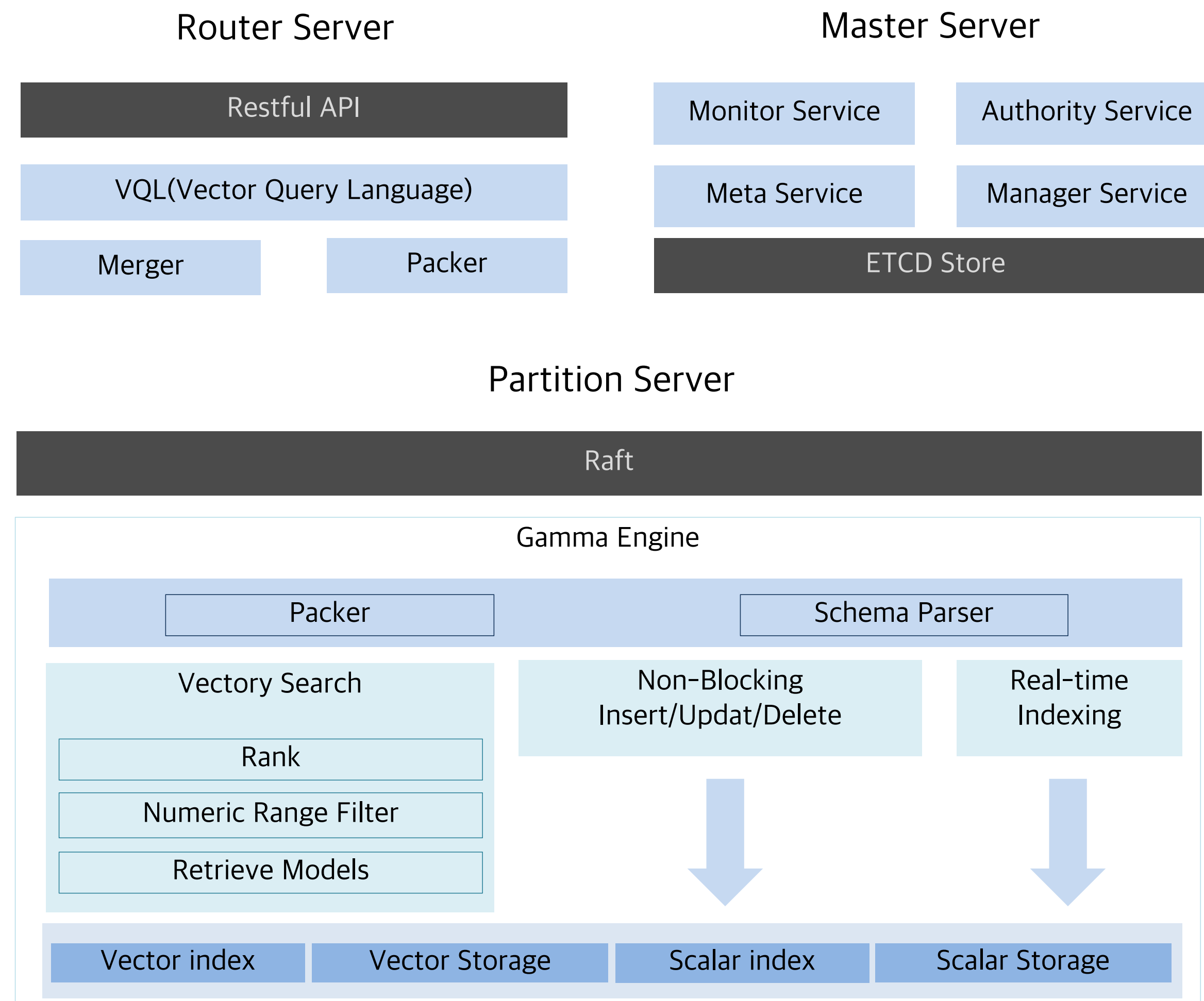
属性过滤：结合元数据过滤，查询更快，结果更准

定义数据逻辑视图，支撑功能丰富、高效灵活的查询接口

ID	Scalar-Column1	Vector-Column1	Vector-Column2
1	1.0	[1.0, 2.0]	[1.01, 2.87,..., 3.09]
2	2.0	[1.1, 2.1]	[1.10, 2.12,..., 3.61]
3	3.0	[1.2, 2.2]	[1.22, 2.52,..., 3.72]
4	4.0	[1.3, 2.3]	[1.34, 2.43,..., 3.34]
5	5.0	[1.5, 2.5]	[1.54, 2.45,..., 3.51]
6	6.0	[1.6, 2.6]	[1.56, 2.46,..., 3.46]
7	7.0	[1.7, 2.7]	[1.27, 2.57,..., 3.76]

- 以Table的形式进行数据组织，支持多Table逻辑隔离不同的业务数据
- 支持标量或者矢量，支持标量索引和向量索引加速检索
- 支持类ES的 Restful查询接口，支持 Attribute-Filtering和 Multi-Vector混合查询
- 支持 python sdk，能够更好与AI生态衔接

Vearch定位：构建分布式集群方案，实现高可用、可扩展



集群化

- 百亿级数据规模支撑

可扩展

- 支持横向在线扩容
- 支持数据在线迁移

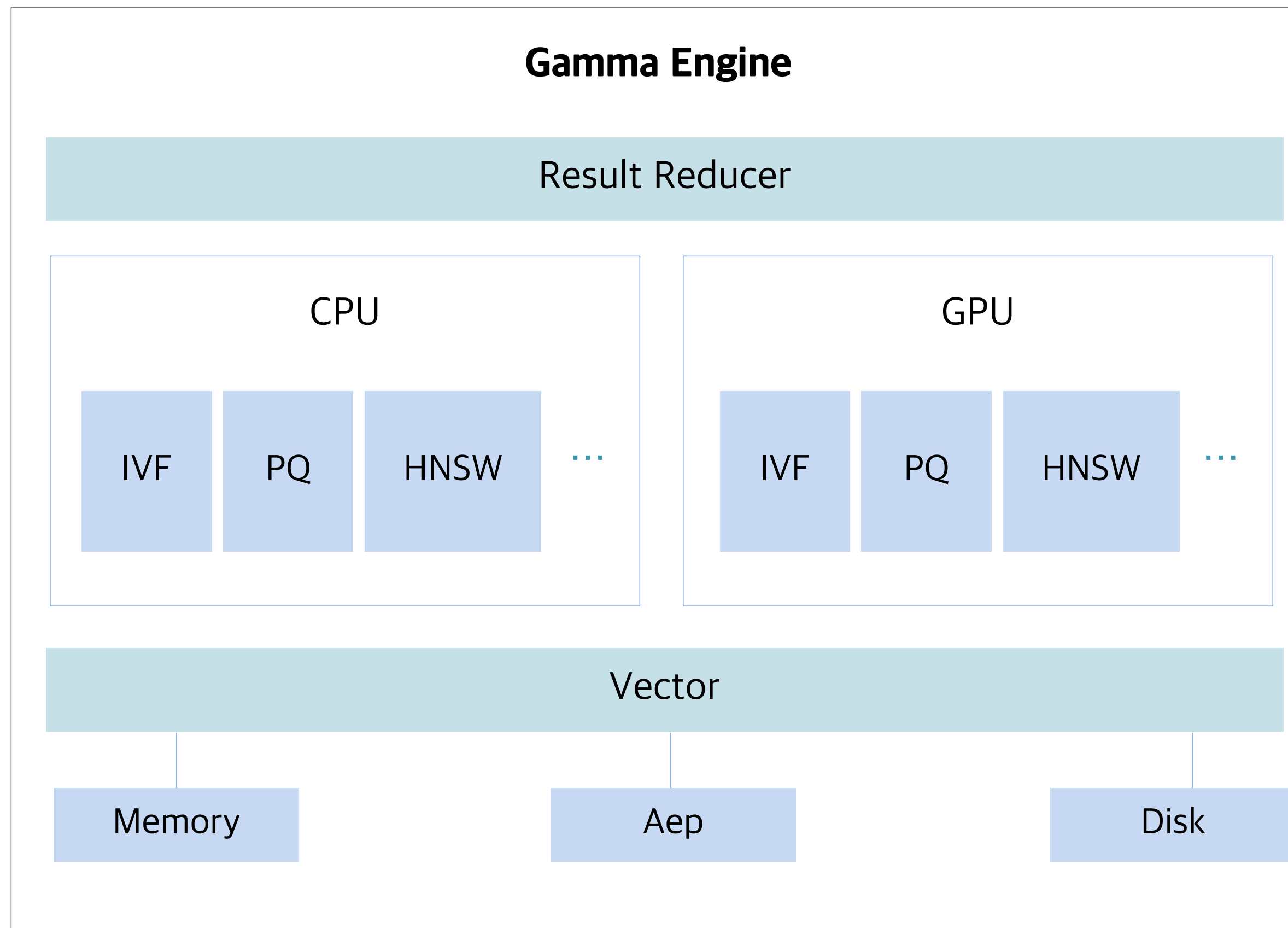
高可用

- 基于Raft算法，支持数据多副本可靠存储
- 支持自动故障发现及恢复
- 支持数据及索引的持久化存储

云原生

- 支持一键K8S部署
- 存储与计算分离

Vearch的异构混合架构，实现低成本、高性能



适应不同数据规模的最优召回匹配模型

- IVFPQ速度快，适用于亿级、十亿级数据规模
- HNSW查询精度高，适用于千万级数据规模
- 二进制索引用于处理unit8数据场景

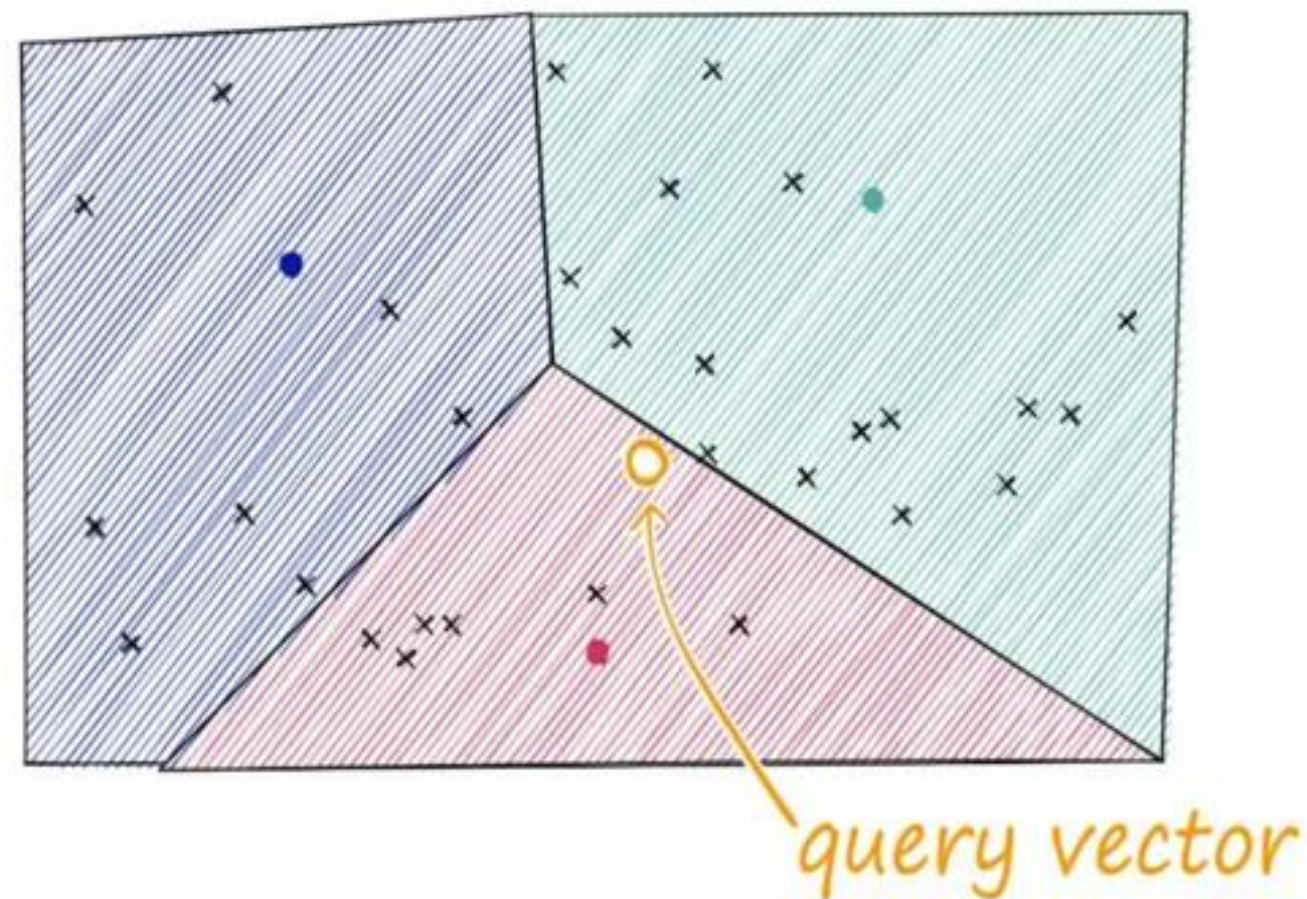
基于GPU和CPU异构平台的优化

- 指令重排，提升cache命中率
- 计算负载精细化调度，充分发挥GPU大数据量处理能力

分层存储体系构建

- 针对不同业务场景使用不同存储来平衡性能和成本消耗
- 支持在线热切换存储，自适应对性能和容量的诉求变更

Vearch支持实时索引，实现数据写入可实时检索技术

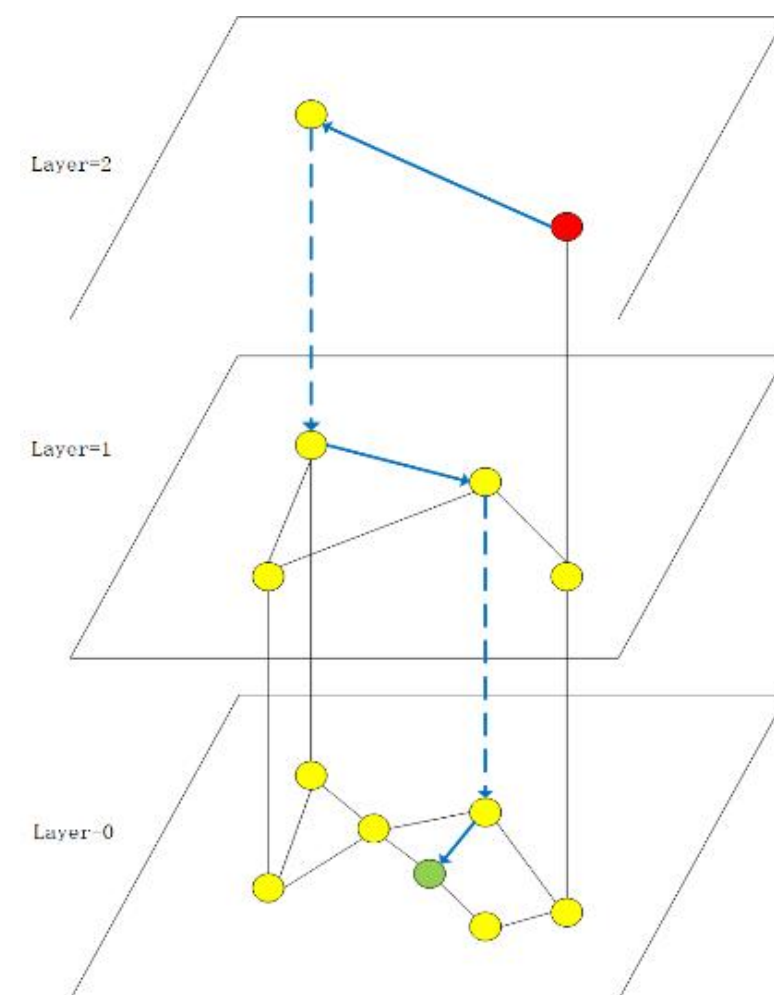


IVF

特点：构建索引快，召回非常高

算法原理：聚类进行分类

关键技术：原子指针、延迟释放

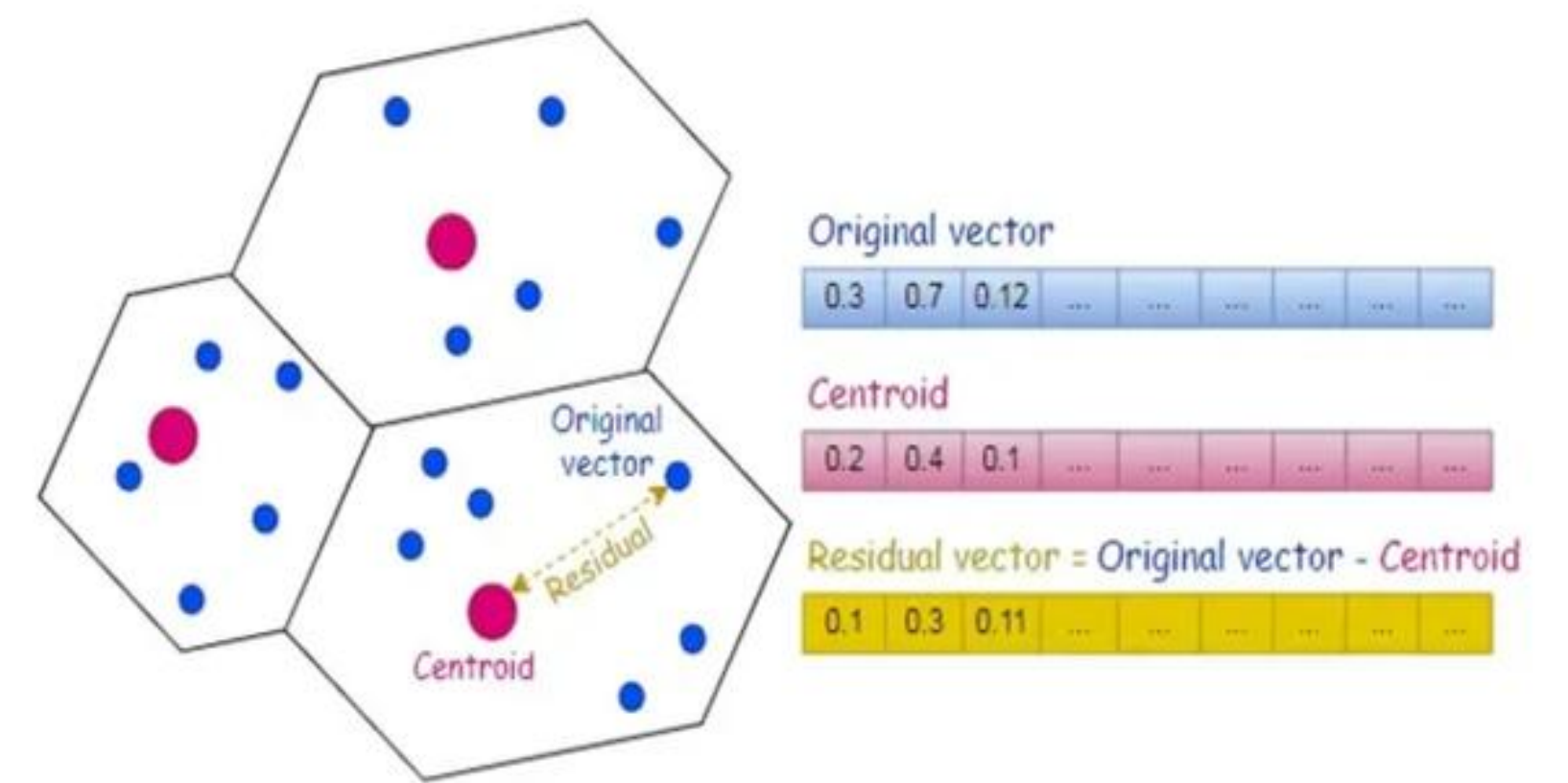


HNSW

特点：检索快、精度高

算法原理：图节点邻接加速检索路径

关键技术：图节点预分配、动态节点扩展



PQ

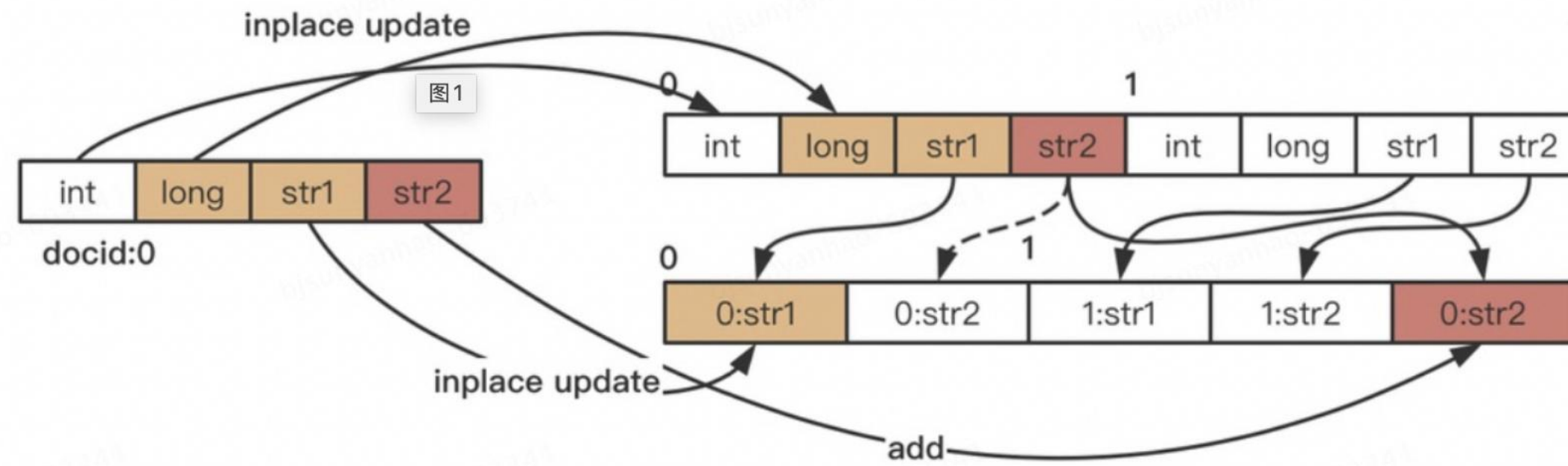
特点：检索快、内存占用小

算法原理：压缩量化

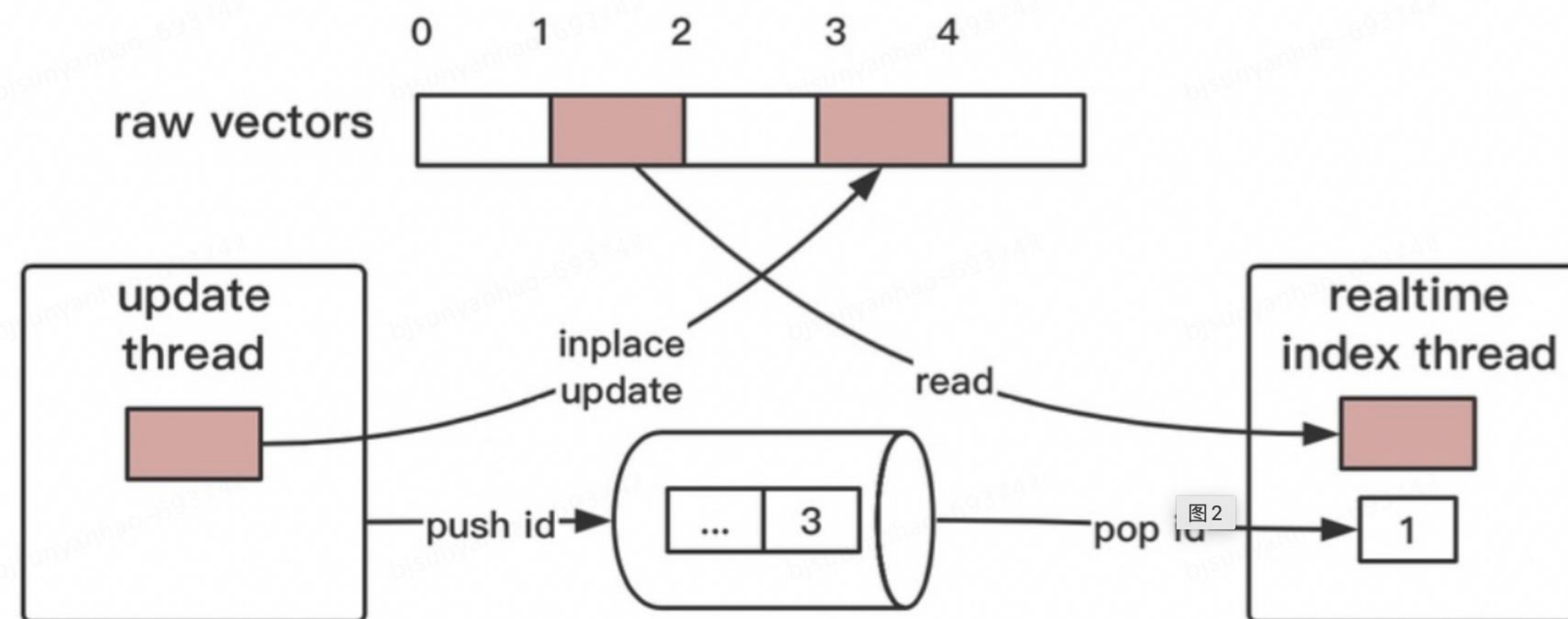
关键技术：动态索引切换

Vearch 的索引 Inplace update和 Compaction方案

- Table中的实现

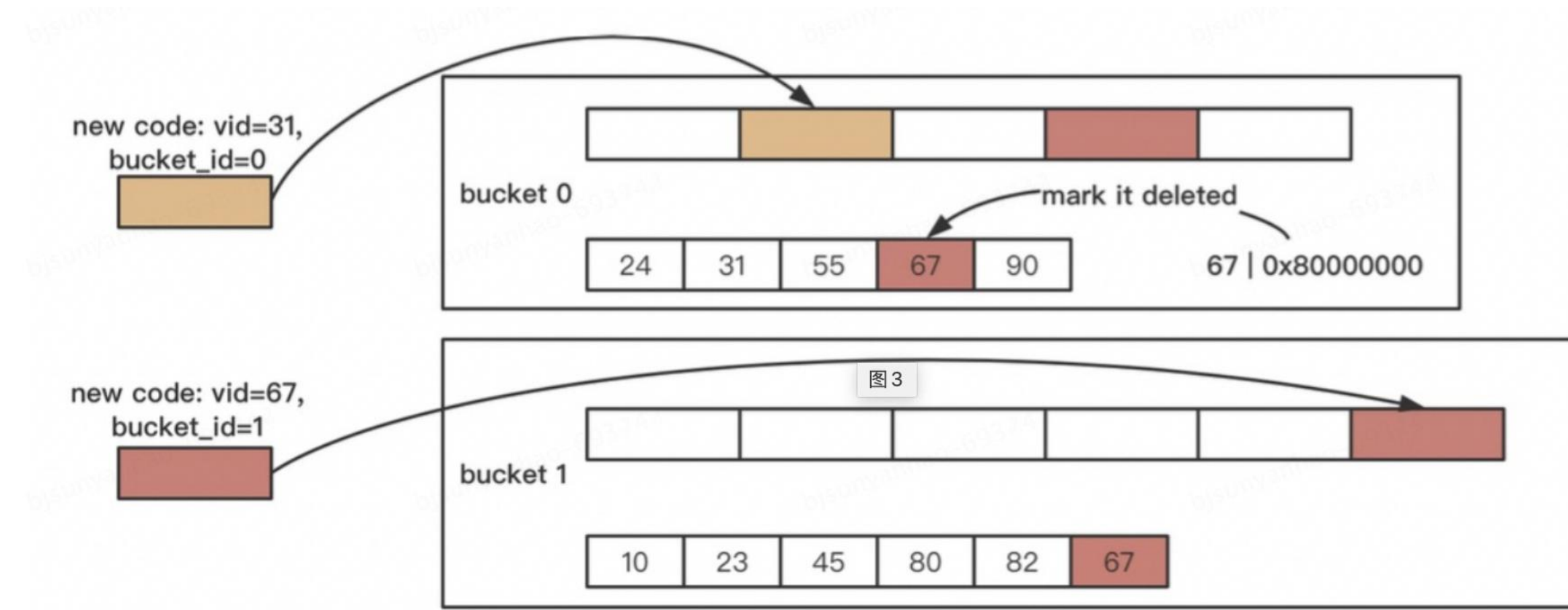


- Raw vector中的实现

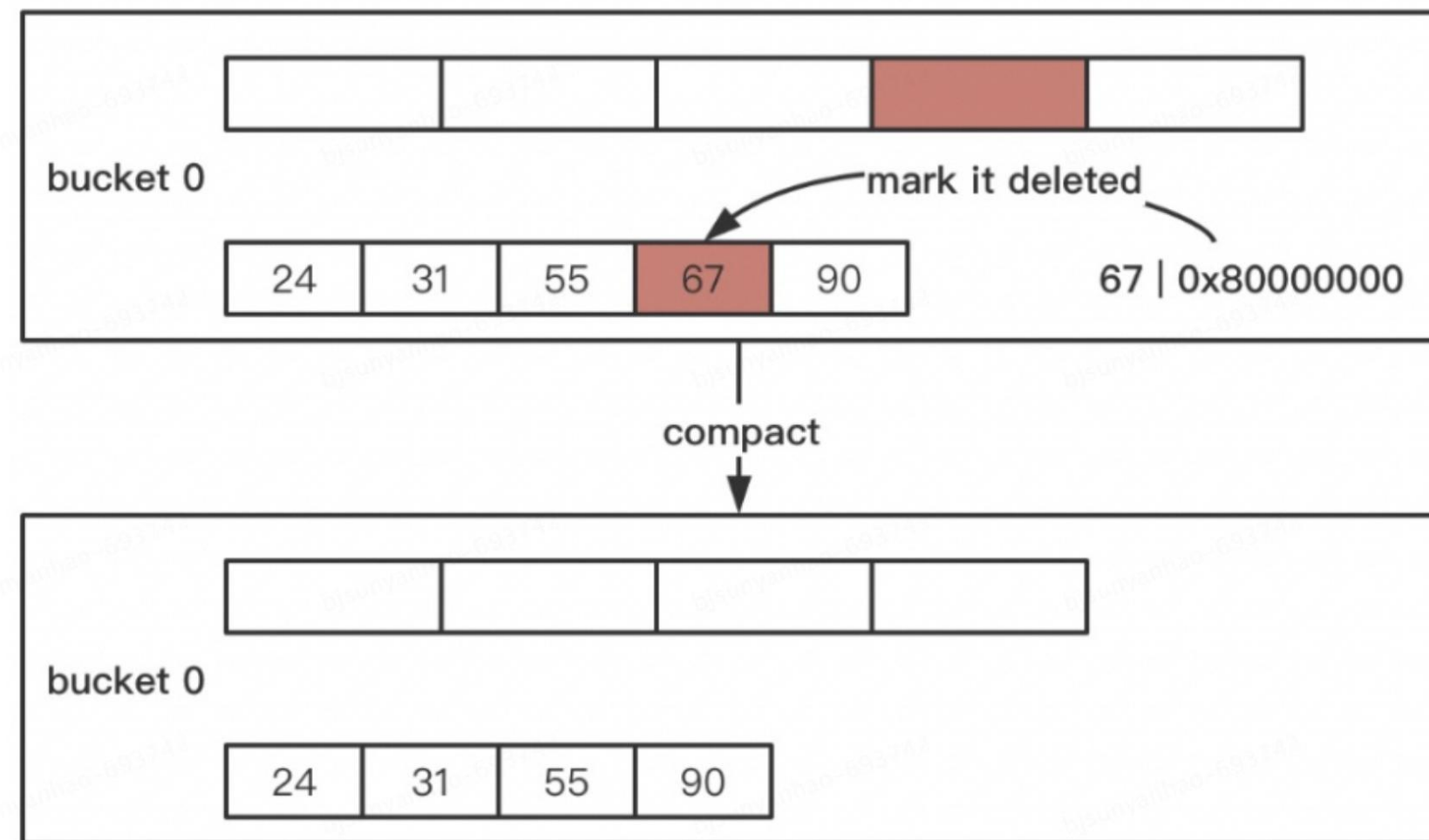


Vearch 的索引 Inplace update和 Compaction方案

- 倒排索引中的实现



- compaction



Vearch数据库：在内外多个领域被广泛应用

向量数据库的一些落地场景

使用企业：100+，总数据量：1000亿+
Star 1.7k，Issues 547

图片场景

重复铺货 | 同款推荐 | 图片侵权治理 | 敏感头像检测 | 新品识别 | ...

文本场景

智能问答词向量 | 用户意图分析 | 涅槃词向量 | 风控团伙作弊治理 | 本地知识库构建 | ...

视频场景

视频去重 | 人脸识别 | ...

推荐

直播种草推荐 | 商品推荐 | 优惠券推荐 | ...



【论文发表：The Design and Implementation of a Real Time Visual Search System on JD E-commerce Platform. In the 19th International ACM Middleware Conference, December 10-14, 2018, Rennes, France.】

THANKS



—
软件正在重新定义世界

Software Is Redefining The World