

KINGBASE | 金仓社区
金仓数据库



杭州站

KING 大 @ 加面对面

—— 大模型时代的数据库平替新范式 ——

基于数据库的 LLM体系化落地

分享嘉宾：张昊

King大咖面对面-杭州站

引子 – 印巴空战的体系化启示

KING BASE | 金仓社区



战机?

导弹?

飞行员?

体系化

数据链



算力



算法



数据



挑战一

模型训练与推理

挑战二

数据获取与利用

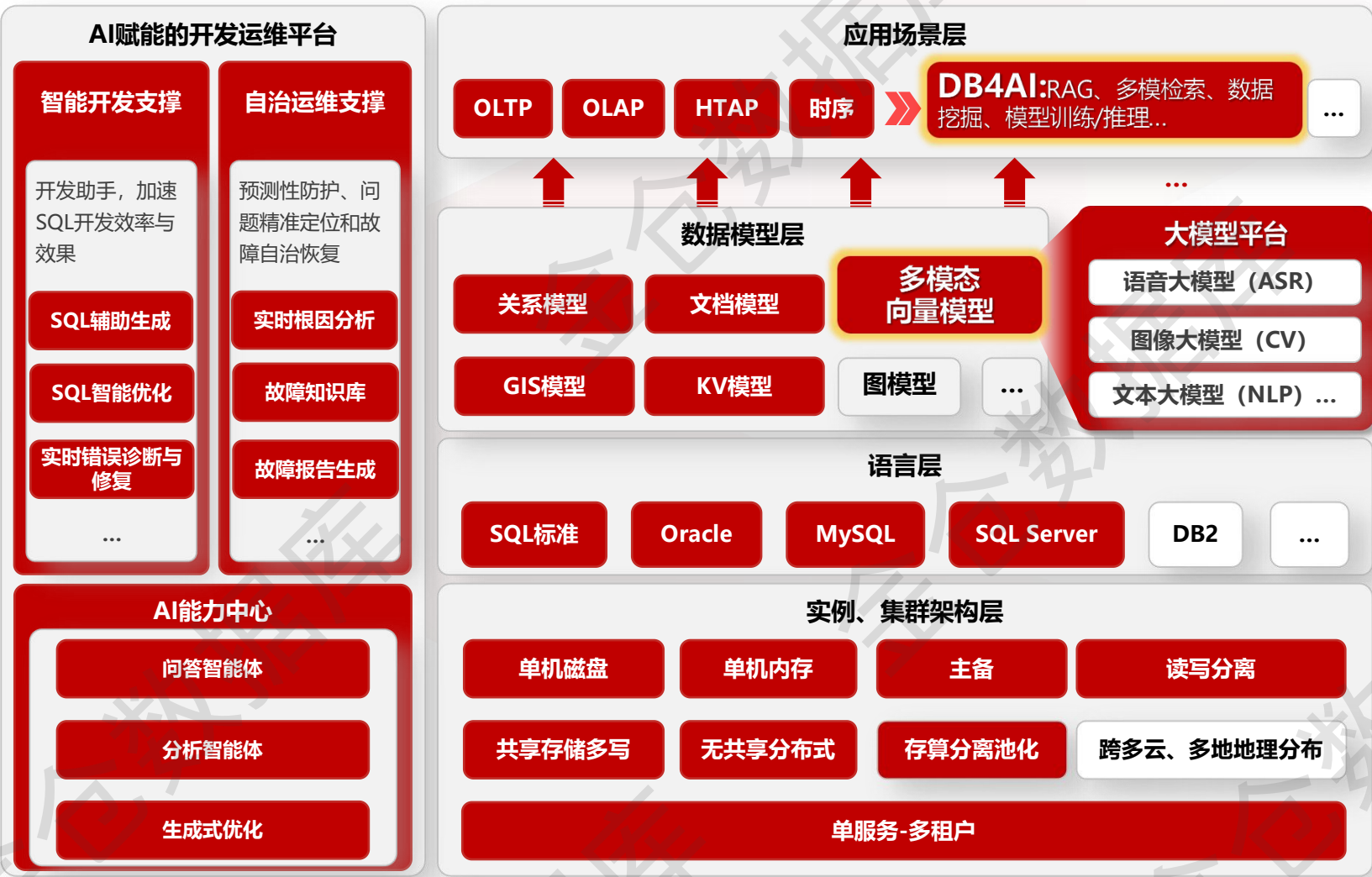
挑战三

隐私与数据安全



KES – AI时代的融合数据库架构

KING BASE | 金仓社区



一体化层次	一体化目标
多语法体系一体化兼容	平替: “0” 代码修改完成应用迁移
集中分布一体化架构	多集群架构 满足不同级别的可用/业务连续性、性能扩展性、成本需求 , 最大化投资价值
多应用场景一体化处理	满足 企业级应用与AI创新等不同场景 的数据存算需求 <ul style="list-style-type: none">结构化数据的存算问题, 在关系库框架内解决库内计算, 减少数据传递
多模数据一体化存储	收敛技术栈 , 降低应用复杂度和成本 减少库间数据同步, 降低同步开销
开发运维一体化管理	AI赋能开发运维: 基于大模型的开发与运维解决方案, 显著提升工作效率

打造支撑构建人工智能系统的数据底座

AI时代的融合数据库 - 支持人工智能的数据管理技术

KING BASE | 金仓社区

面向AI的数据治理

数据发现和清洗：可扩展的元数据管理
统一数据处理
自动数据流水线生成

数据血缘和标注：
基于主外键的价值分析
基于模糊集的数据补全

面向AI的模型训练与推理

模型训练平台

模型推理框架

异构执行引擎

异构AI计算引擎

并行加速技术

AI执行优化技术

智能的数据分析技术

面向AI的数据存储引擎

向量数据库

多模态数据混合检索

大模型缓存

业务领域		DB发挥的作用
计算	内置AI算法	数据挖掘算法
		机器学习算法
	内嵌大模型	模型训练
		模型推理
存储	向量数据库	向量数据、向量索引
		相似度计算、基于相似度聚类、向量化模型比较

多模态向量数据库 – LLM体系化落地的基石

KING BASE | 金仓社区

金仓向量数据库可解决大模型容易产生幻觉、知识更新不及时、数据隐私风险等问题，突破大模型在时间和空间上的限制，加速行业场景大模型落地

Deepseek从爆火到应用面临诸多挑战

> AI “幻觉”：LLM生成的内容与输入、上下文或已知世界不符，产生“幻觉”；

> 知识匮乏：训练数据的范围局限、时效性等问题，无法应对时效性知识要求场景；

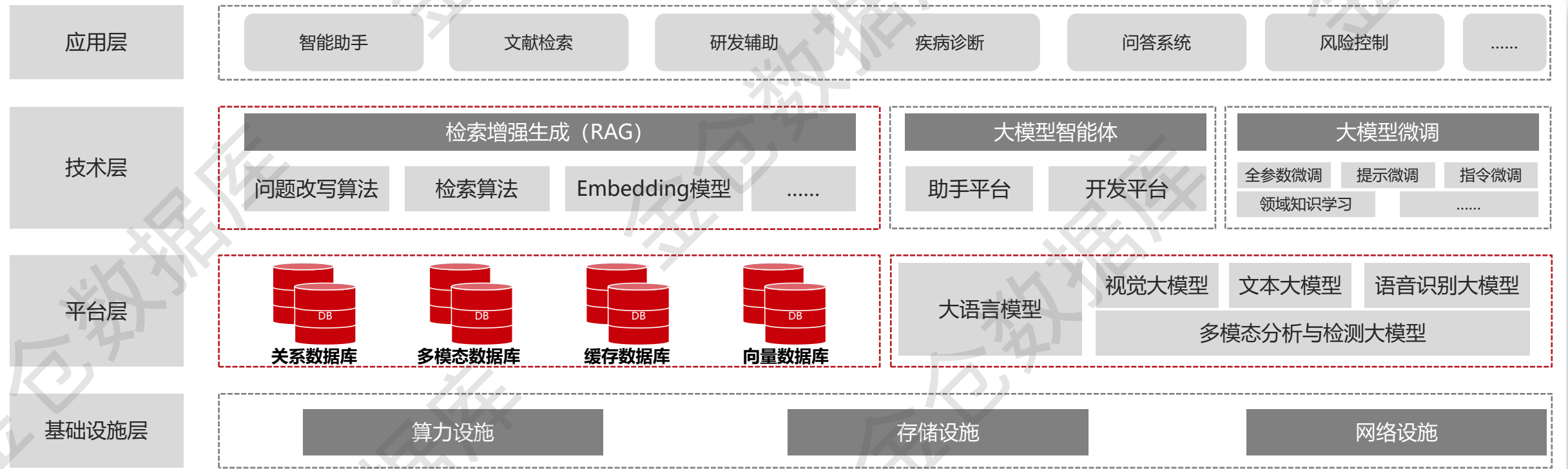
> 数据安全：私域数据、问答数据、用户信息暴露在公开模型，面临严重的安全风险。

向量数据库助力AI应用实践落地

> 数据赋能：数据赋能AI，解决大模型“幻觉”问题，让AI更靠谱；

> 实时导入：实时更新的私域数据，确保私域知识库完备性与时效性；

> 数据安全：纵深防御的数据库安全能力，确保AI应用不会出现数据泄露风险。



金仓向量数据库 – 关键特性

KING BASE | 金仓社区

金仓向量数据库是在关系型数据库KES之上的向量功能增强，除向量数据外也提供完整的关系、文档、GIS以及时序数据的一体化处理能力，天然继承了KES的高性能、高可用、高安全、易使用、易管理等特点，**提供一站式的大模型应用开发与落地解决方案，大幅提高AI技术在企业应用落地效果，降低使用成本。**

高性能

- 支持高效存储和检索高维向量数据，实现语义搜索和相似性匹配，支持亿级向量数据毫秒级召回，索引性能领先友商20%

高安全

- 从事前的身份鉴别、访问控制，到访问过程的传输加密，到数据存储的加密、脱敏，至事后的安全审计，全栈纵深防御。

高可用

- 从实例、集群到多中心的高可用保障，确保数据零丢失、服务持续在线

高兼容

- 支持多种异构数据库的应用API级兼容，以及开源数据库的原生协议兼容，最大程度降低迁移难度

全场景支持

- 支持向量数据与关系数据、文档数据、GIS以及时序数据的混合检索，一站式全场景支持

全周期自治

- 提供数据库部署、迁移、开发、测试、调优、运维等全生命周期智能自治管理

执行算法优化

- > **并行执行**: Scan / Join / Append /..., 多个进程并行处理一条SQL, 充分利用CPU资源
- > **共享执行计划缓存**, 全局复用, 生成更快, 内存更高效
- > **编译执行**, 将SQL表达式, 将PLSQL函数及存储过程, 并编译成动态库

数据存储优化

- > **基于逻辑时钟的快照优化**, 有效降低MVCC版本判断时间, 并支持更高并发响应。
- > **热点页面固定缓冲池技术**, 降低IO和封锁开销
- > **WAL并行回放技术**

多核优化

- > **进程绑核**, 避免进程在核间飘移
- > **NUMA化数据结构改造**, 减少跨核内存访问
- > **数据分区**, 减少进程访问冲突
- > **原子指令**, 减少计算开销

量化压缩优化

- > **标量量化**: 高精度浮点向量映射到低比特整数表示, 降低空间占用
- > **乘积化向量**: 高维向量分解为多个低维子空间
- > **量化索引结构**: 在HNSW中融入量化思想, 提高检索效率、降低内存占用

内外协同检索优化

- > **内存与磁盘数据协同设计**, 将压缩后的索引信息固化到内存, 结合磁盘索引将查询目标快速定位到磁盘中的局部搜索空间, 从而降低整体的内存占用需求、提高检索效率

高安全 – 纵深防御，资质完备

KING BASE | 金仓社区

信息技术产品安全分级评估证书 评估保障级4增强级 (EAL4+)

国家网络与信息系统安全产品质量
监督检验中心

中国国家信息安全产品认证 网络关键设备和网络安全专用产品安全认证

中国网络安全审查技术与认证中心

商用密码产品型号证书 (二级) 数据库管理密码模块

国家密码管理局

传输安全

传输加密: 提供SSL支持, 服务器端证书验证、客户端证书验证
基于国密算法的数据完整性保护

身份鉴别

口令复杂度控制

口令有效期控制

异常登录控制

客户端认证: 加密口令认证、Kerberos认证、LDAP认证、Radius认证、证书认证等

安全审计

DDL/DML/查询 生成审计结果
入侵检测
审计查询

审计存储
审计转储

数据中心

数据库服务器

数据库实例

数据库存储

三权分立

DAC自主访问控制: 角色启用禁用、受限DBA、备份恢复权限控制、系统ANY权限控制

MAC强制访问控制: 数据访问控制

加密: 透明存储加密、软加密设备、硬加密设备、主密钥管理、对象密钥管理、密态计算

动态脱敏: 完整脱敏、部分脱敏、随机数、邮件

动态脱敏: 备份数据脱敏

客体重用: 内存残留销毁

客体重用: 磁盘残留销毁

数据完整性保护: 数据页面一致性校验

数据完整性保护: 基于区块链防篡改

加密函数: MD5、SHA1 Blowfish、SHA224/256/384/512、AES、DES/3DES/CAST5、SM2、SM3、SM4、RC4

访问控制

存储安全

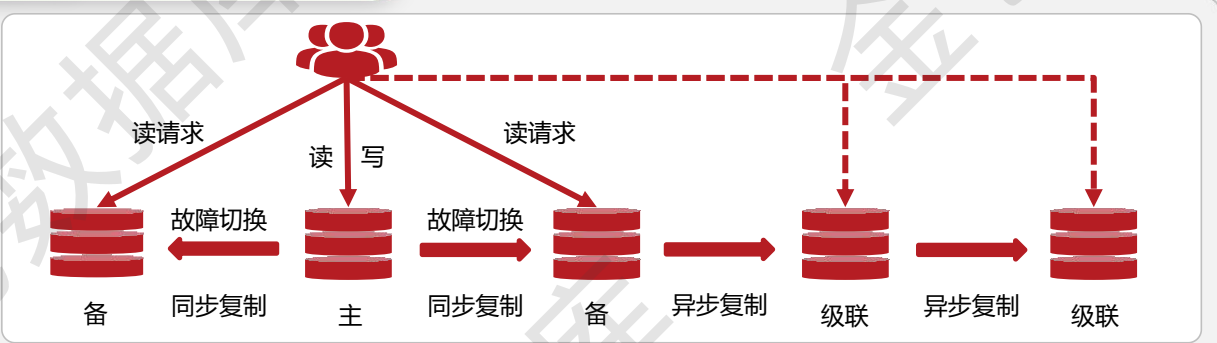
高可用 - 从实例、集群到多中心的高可用保障

实例级

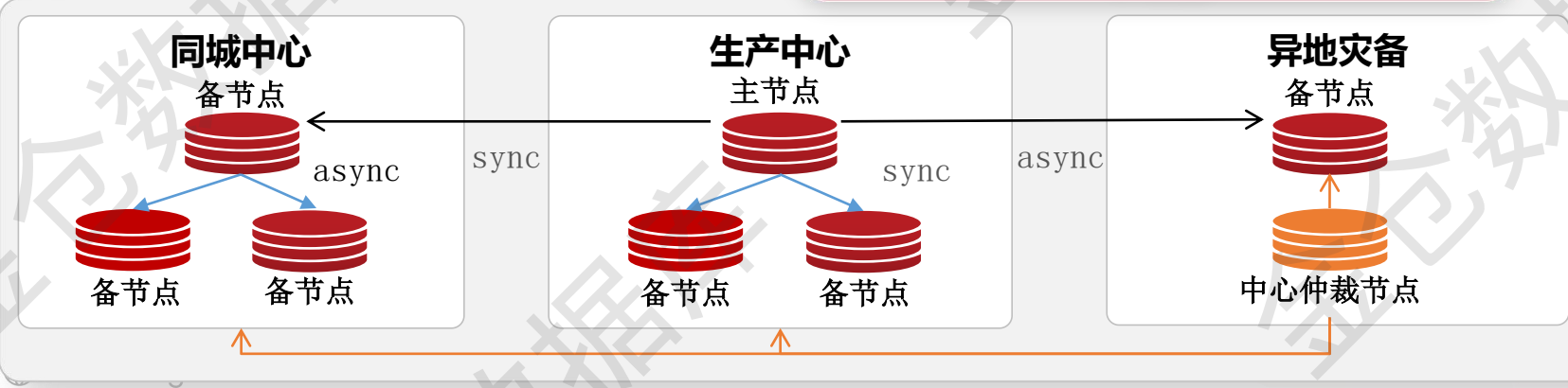
全量/增量/差异备份	全库/用户/对象备份	文本/二进制/加密输出	远程备份	备份管理
全库/事务/时刻恢复	全库/用户/对象还原	增量检查点	指定对象/时刻闪回	第三方备份平台集成
日志管理	数据一致性校验	控制文件多副本	进程级运行监控	故障进程自动重启

集群级

- 基于物理日志的全同步复制，数据不丢失
- 备节点支持异步复制，不阻塞事务，业务持续可用
- 备机日志回放并行加速，故障快速切换
- 增强型一致性协议选主，避免脑裂
- 节点故障恢复后自动回归，保持集群规模稳定



中心级



- 同城双中心物理全同步，确保数据“零”丢失，故障“秒”切换
- 跨城市物理日志高速复制，最优可实现秒级 RPO
- 中心级多数派一致性切换，避免中心级双主

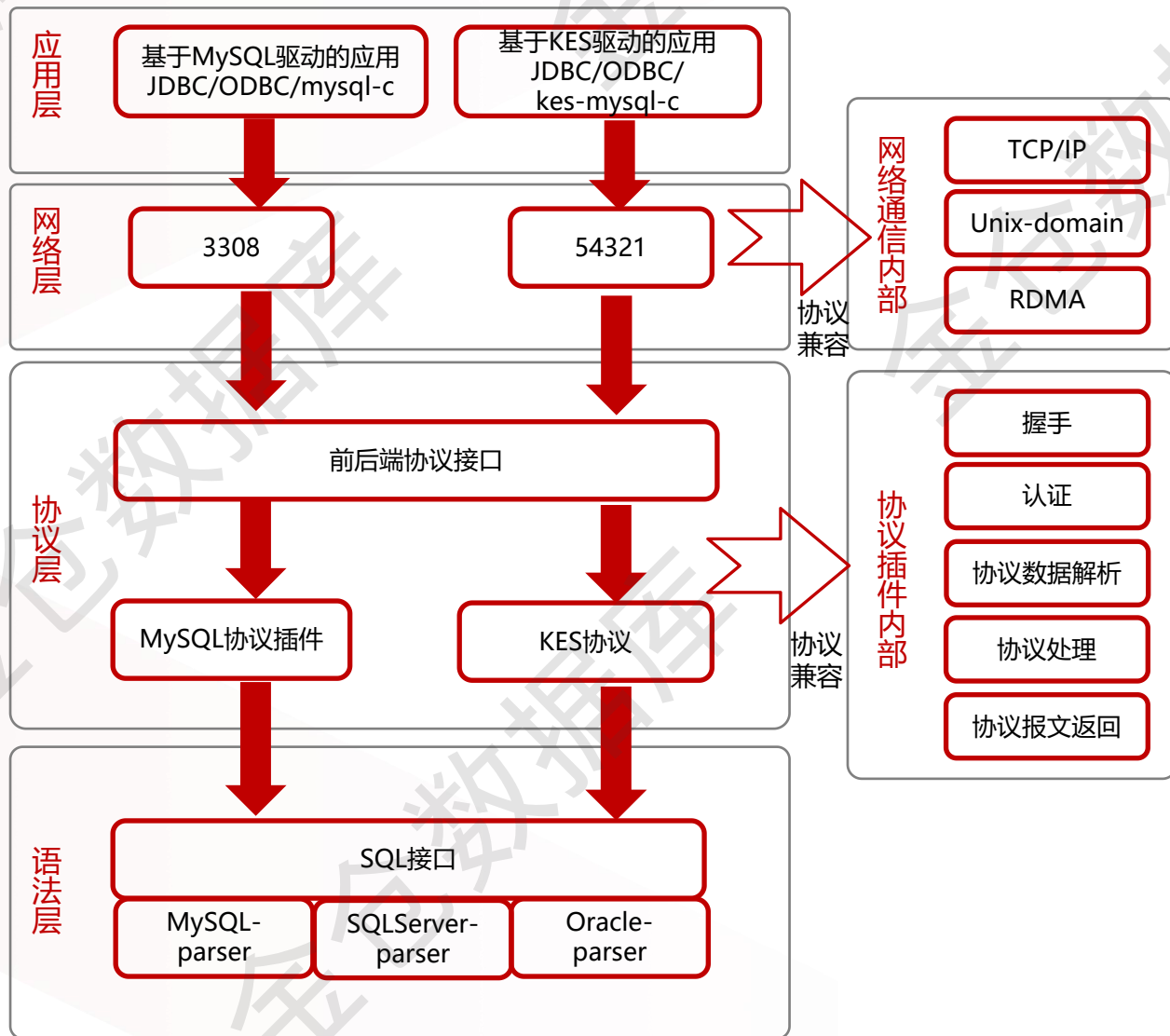
应用API级兼容

- 优势
 - ✓ 方案成熟
 - ✓ 适合协议不开源的产品，如兼容Oracle 应用API
- 缺点
 - ✓ 不同语言框架的API种类众多，尤其各个数据库都有非标准API
 - ✓ 需要应用替换驱动库和连接串等。

原生协议兼容

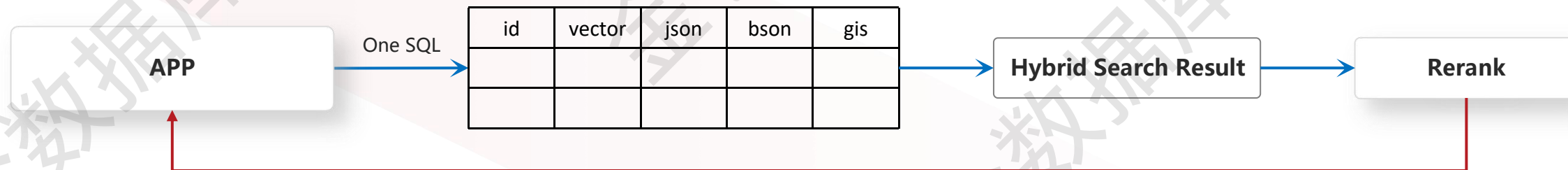
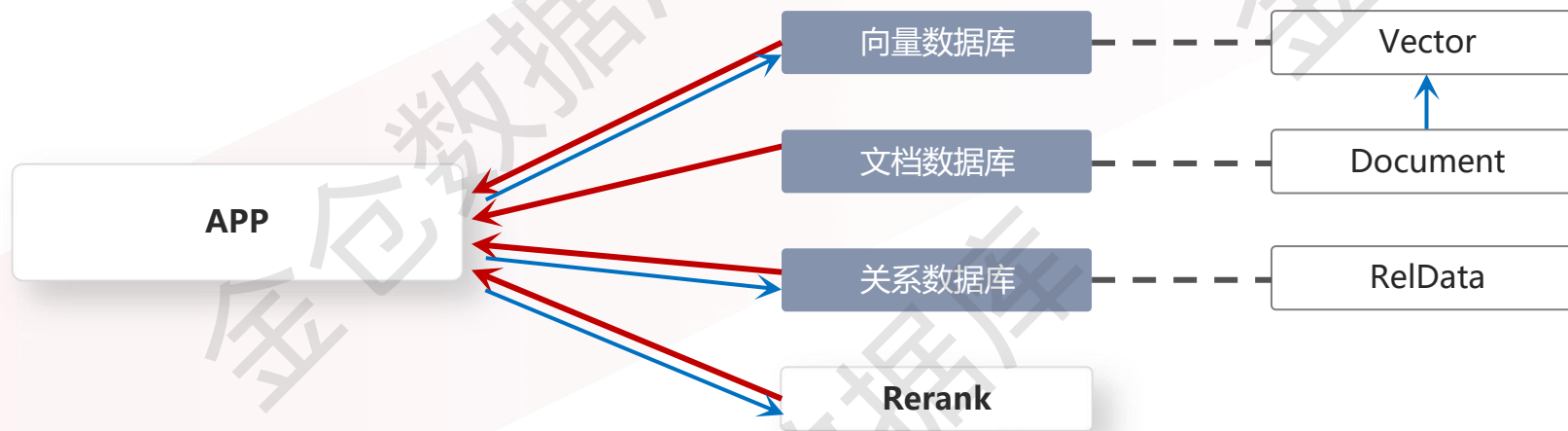
- 上层使用统一的协议接口，内部协议插件提供不同兼容协议的转换，达到一个库多种协议的效果。
- MySQL协议插件：
 - ✓ 使用原生MySQL客户端驱动直接连接KES
 - ✓ 使用原生MySQL客户端程序直接连接KES

保留所有MySQL生态



全场景 – 强大的混合检索能力

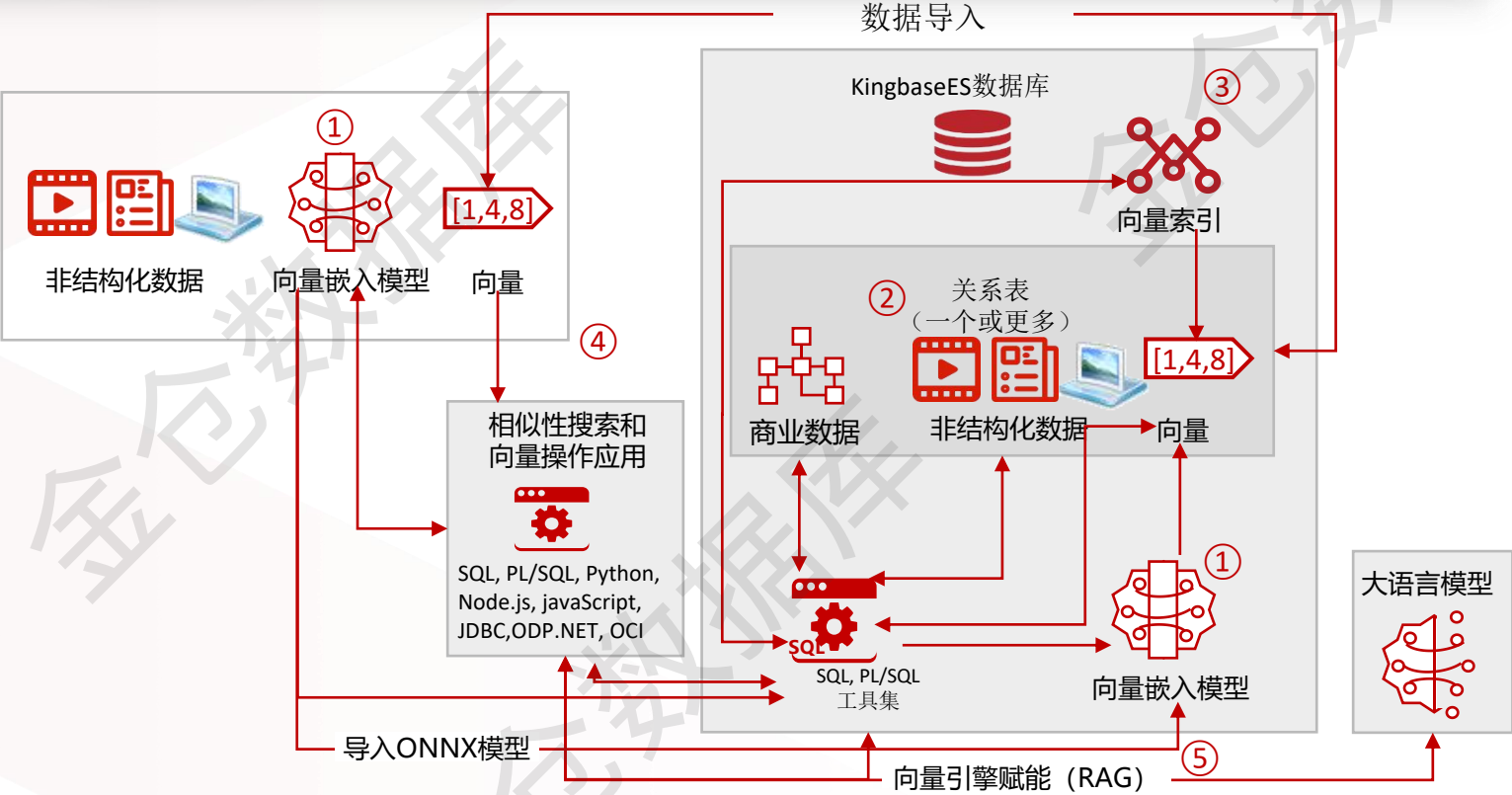
KING BASE | 金仓社区





RAG将传统信息检索系统的优势与生成式大语言模型 (LLM) 的功能相结合。通过将已知的数据和知识与LLM语言技能相结合，使能AI应用的输出更准确、更及时。

步骤	完成的任务		
检索	从知识库中检索到的问答对	问题向量化	将用户输入查询问题转换成向量
		相似度检索	在向量数据库中检索与问题向量最相似的知识库片段
		结果排序	按相似度得分对检索到结果排序，选择最相关片段作为后续生成的输入
增强	增强LLM的提示词 (prompt)		
生成	LLM拿着增强后的Prompt生成问题答案		

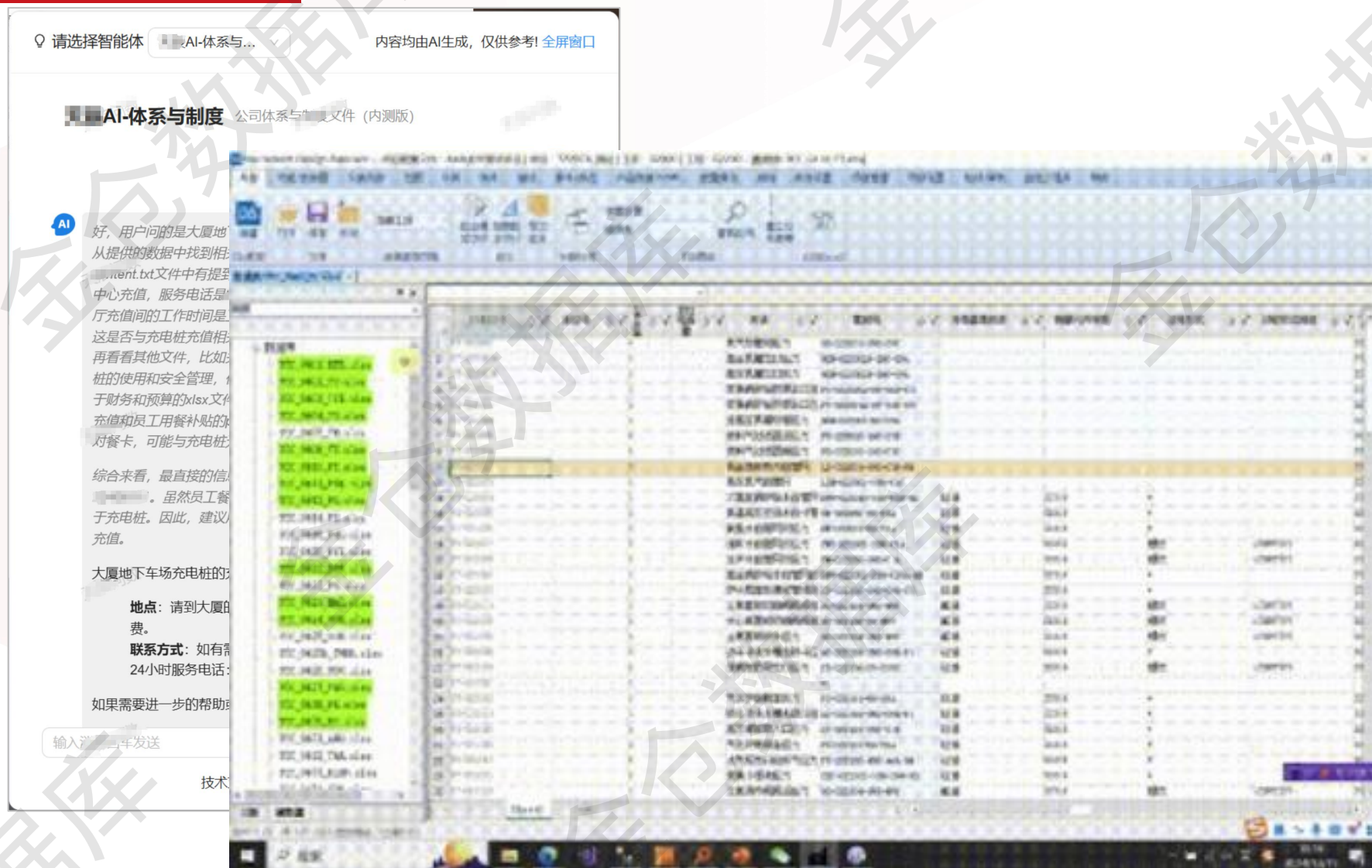


RAG场景 – 支撑某央企智能体全面落地

KING BASE | 金仓社区

落地实践

- 某央企通过KES向量数据库构建知识库，支持多种格式文档快速导入，助力OA智能问答、AI开发平台等智能体上线，为员工提供高效、准确的信息检索、知识推荐和设计开发辅助。
- AI智能体的上线显著提高该企业的数据资产利用率，践行AI赋能企业高质量发展。



典型场景需求

查询某市过去30天内与雨天+分支道路+追尾形式类似的交通事故频发（至少3次）的5个热点区域

返回结果包含

- 事故数量
- 事故详情
- 最早和最后发生时间
- 与预期特征的平均相似度

关键数据模型

- 某市热点区域：固定地理范围，以及细化区块/网格 -> GIS数据模型
- 雨天+分支道路+追尾形式：相似度检索 -> 向量数据模型
- 区域内各事故详情：文档聚合 -> 文档数据模型

```
WITH city_boundary AS (  
    SELECT geom AS city_geom  
    FROM areas  
    WHERE name = '某市'  
    LIMIT 1  
) ,  
filtered_accidents AS (  
    SELECT  
        a.accident_id,  
        a.accident_time,  
        a.location,  
        a.details,  
        a.embedding,  
        ar.area_id,  
        ar.name AS area_name,  
        a.embedding <=> '[1,2,3] '::vector AS vector_distance  
    FROM accidents a  
    CROSS JOIN city_boundary  
    JOIN areas ar ON ST_Contains(ar.geom, a.location)  
    WHERE ST_Within(a.location, city_boundary.city_geom)  
    AND a.accident_time >= NOW() - INTERVAL '30 days'  
    AND a.embedding <=> '[1,2,3] '::vector < 0.5  
)  
SELECT  
    fa.area_id,  
    fa.area_name,  
    COUNT(*) AS accident_count,  
    json_agg(fa.details) AS aggregated_details,  
    MIN(fa.accident_time) AS first_accident,  
    MAX(fa.accident_time) AS last_accident,  
    AVG(fa.vector_distance) AS avg_vector_distance  
FROM filtered_accidents fa  
GROUP BY fa.area_id, fa.area_name  
HAVING COUNT(*) > 3  
ORDER BY accident_count DESC  
LIMIT 5;
```

向量类型为 vector，3 维表示

- 事故类型：1（追尾） 2（侧撞） 3（翻车） 4（...）
- 天气情况：1（晴天） 2（雨天） 3（雪天） 4（...）
- 道路类型：1（高速） 2（城市干道） 3（分支道路） 4（...）

GIS计算

明确事故所在区域

向量计算

事故特征与[1,2,3] 的欧氏距离

文档处理

分区域聚合事故详情清单

分组排序

分组计算热点区域排序



KING BASE | 金仓社区

THANKS

成为世界卓越的数据库产品与服务提供商

