



COSCon'25

第十届中国开源年会

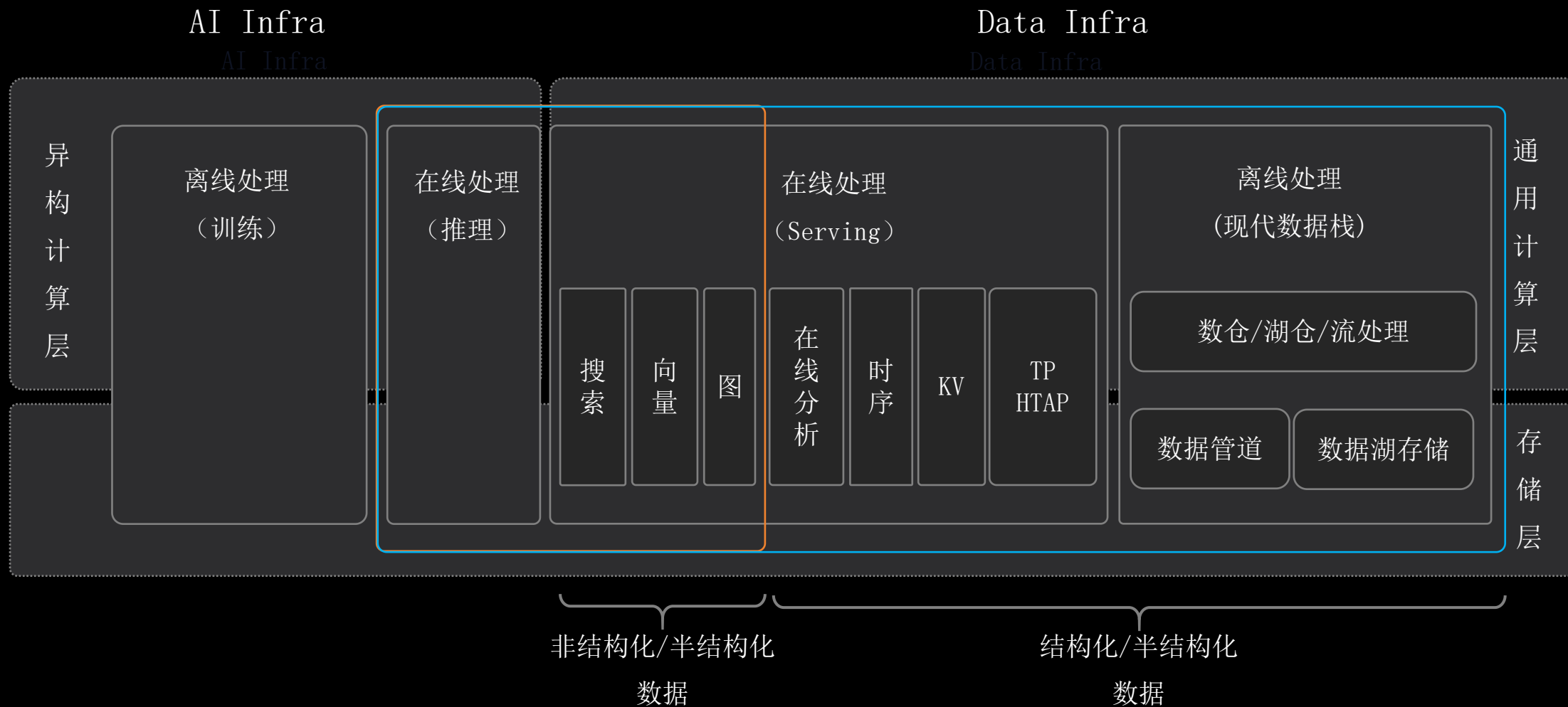
众智开源 | Open Source, Open Intelligence

回顾与展望
——构建大模型时代的核心数据基础设施

张颖峰@RAGFlow



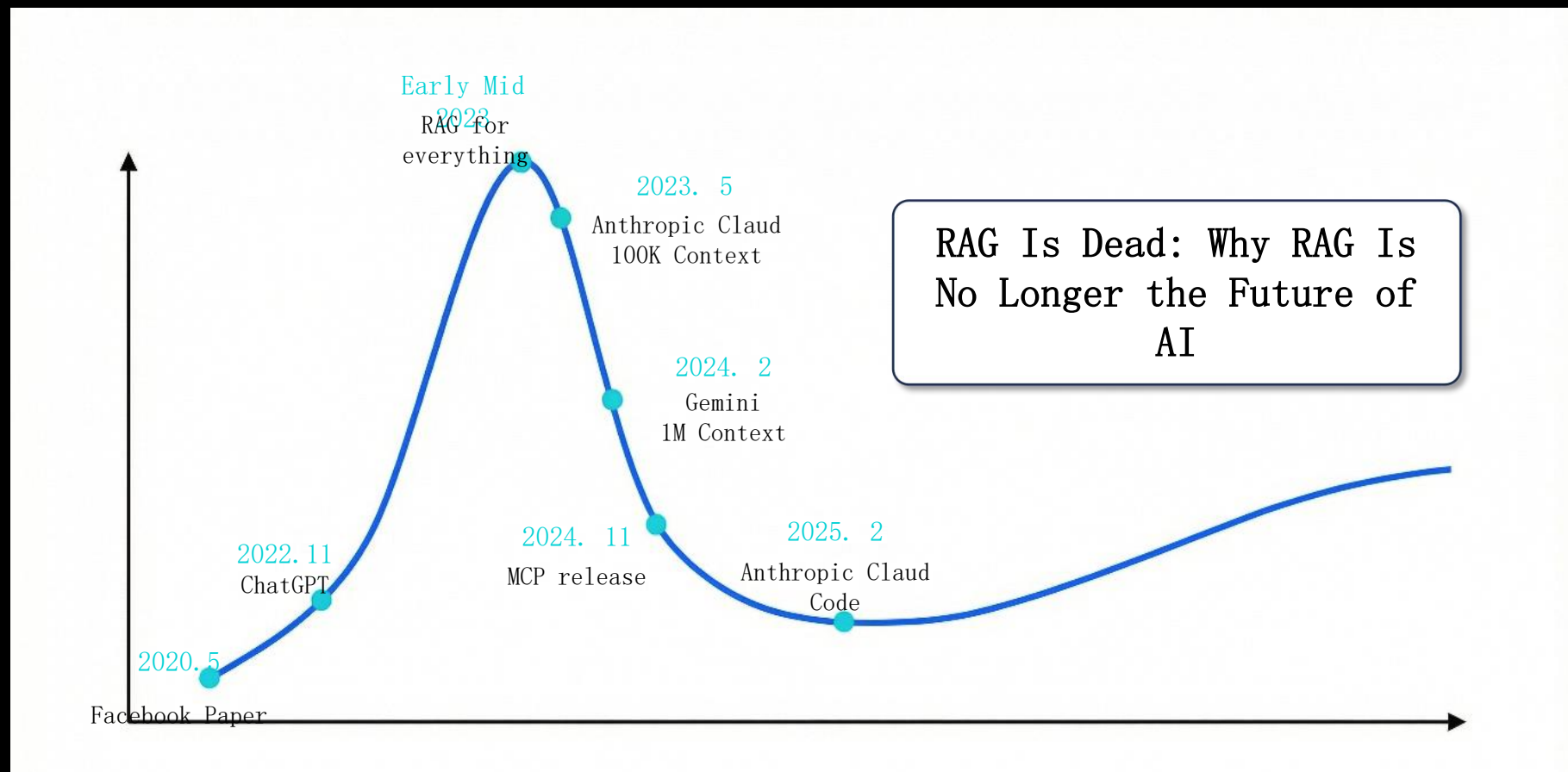
AI Infra 和 Data Infra



RAG Is Dead




- RAG 效果不好
- Long Context 足够了
- 有 Agents 不需要 RAG
- Memory 比 RAG 效果好



向 LLM 提供输入的四种模式





➤ 只依靠 LLM 的 Context

82 x 

➤ 借助于 KV Cache

12 x

➤ 无需索引的搜索 => Grep


7 x 

➤ 基于搜索引擎 => RAG

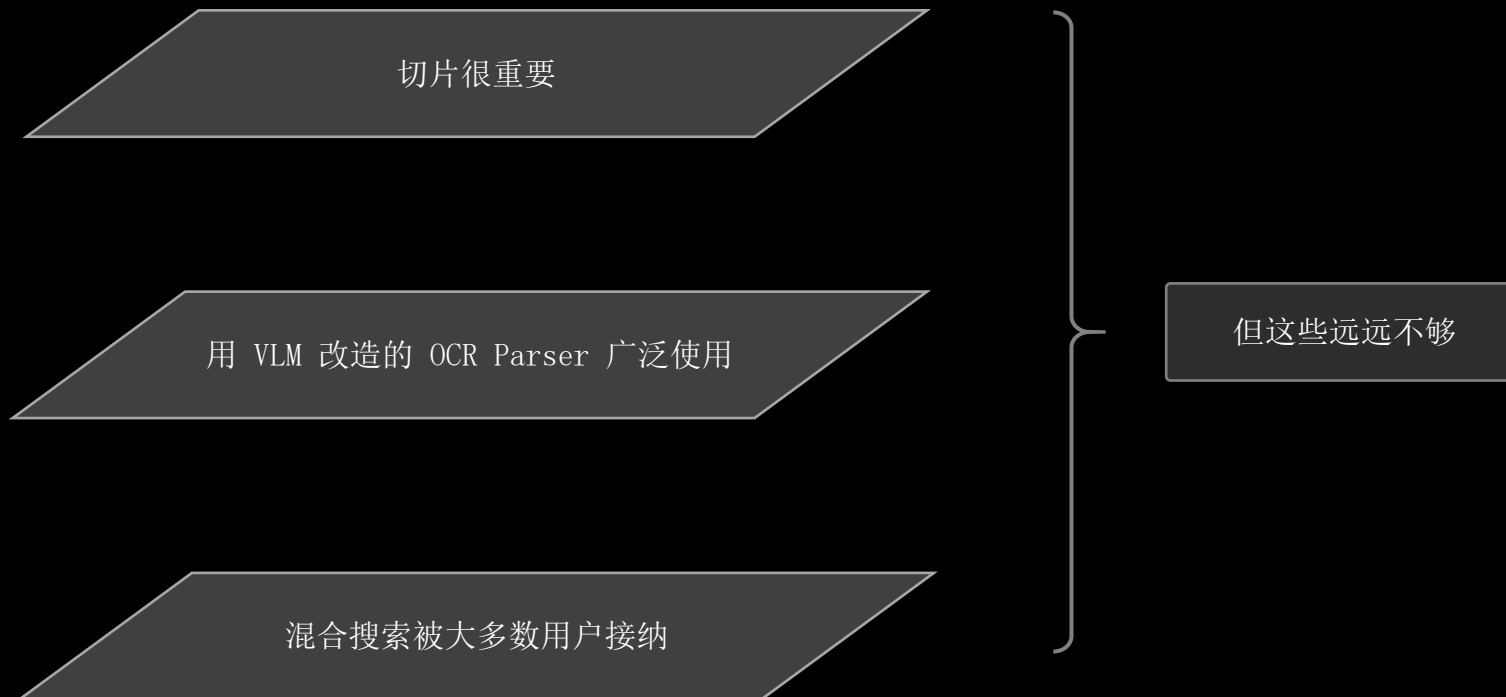
➤ 过长上下文导致成本激增和理解力下降

数据类型	Token消耗
单张技术图表	~1,000 tokens
典型代码库	数千万 tokens
企业级仓库	数亿 tokens
1,000页文档	600,000 tokens

➤ 为什么 Grep 不能满足要求

➤ RAG 就是 Context Engineering 1.0

过去一年 RAG 发展的共识



RAG 的痛点还有解么？

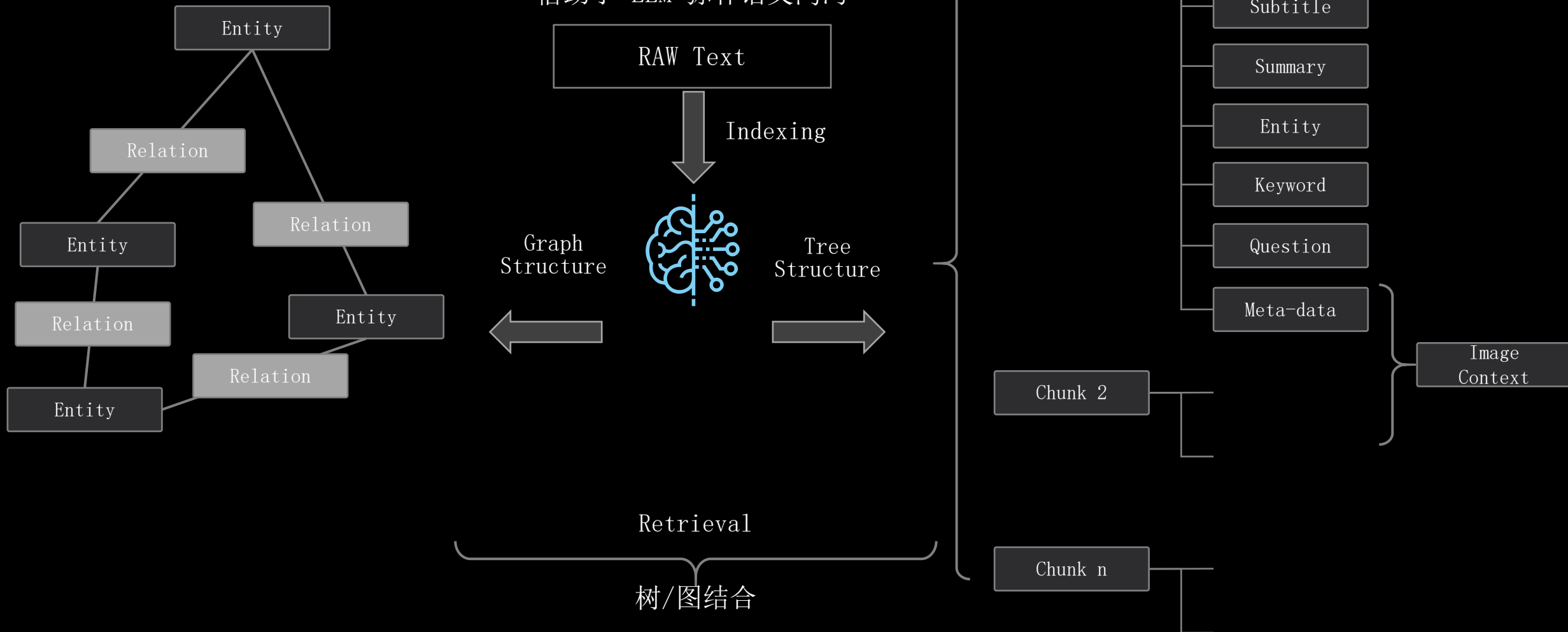


RAG 的痛点应该如何解



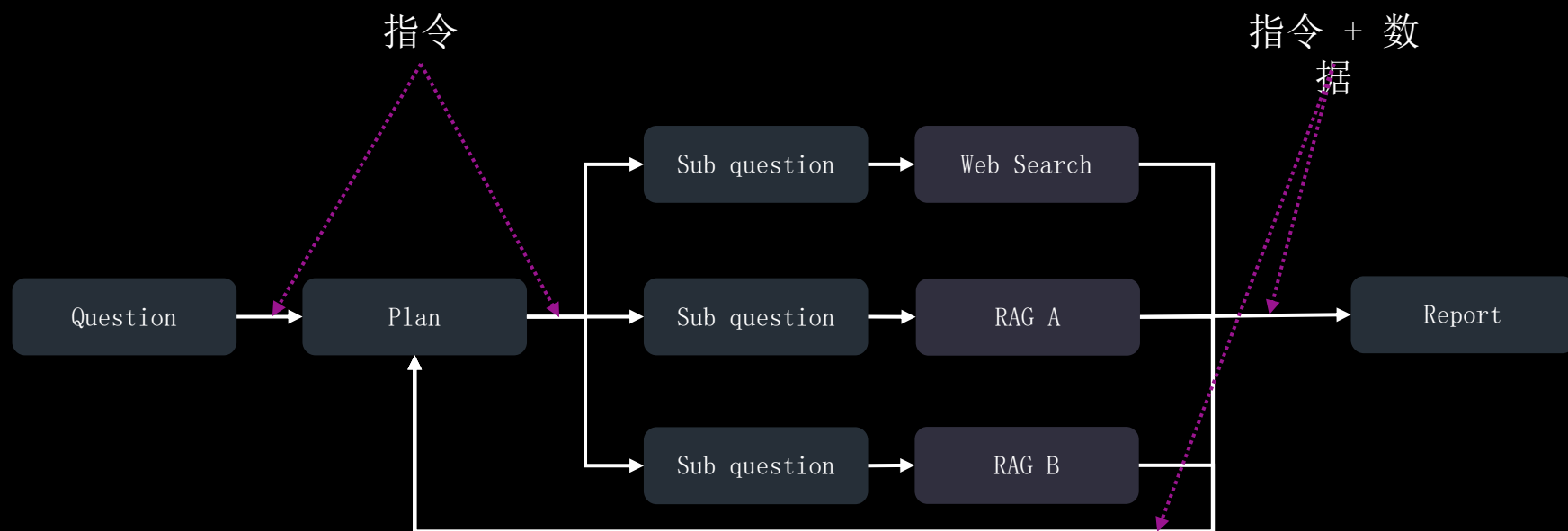
Long Context RAG

借助于 LLM 弥补语义鸿沟

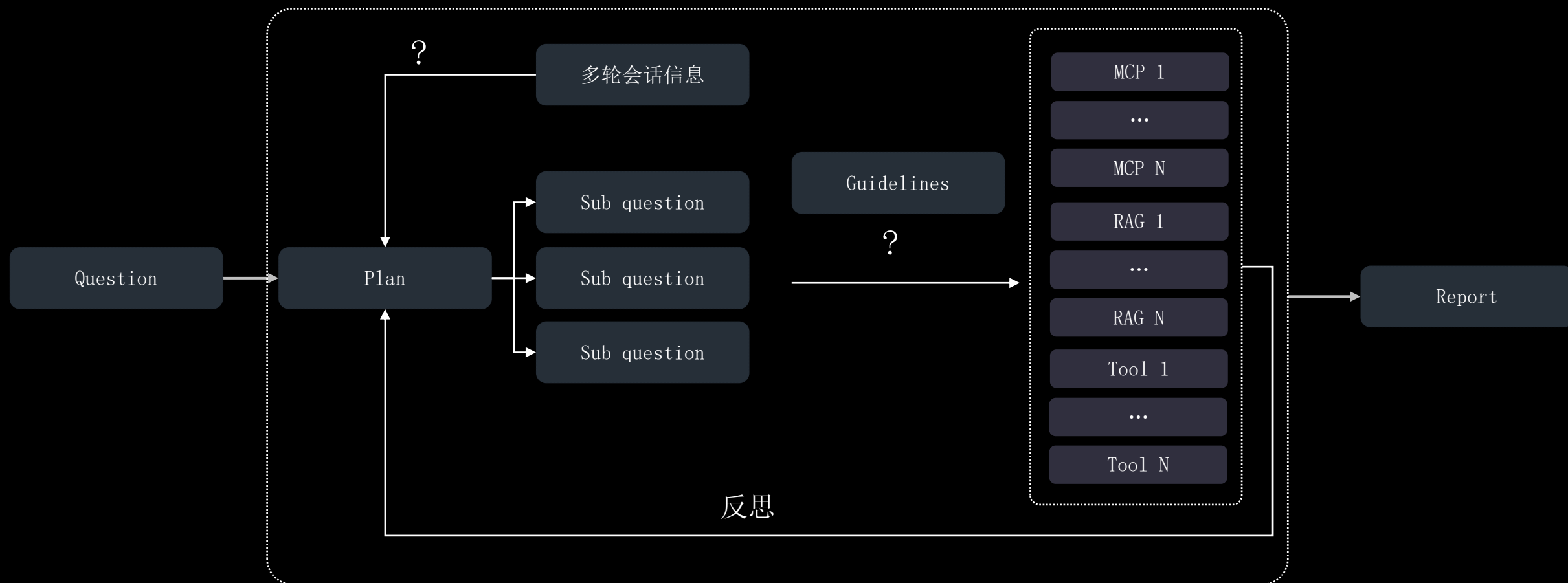


Context = 指令 + 数据

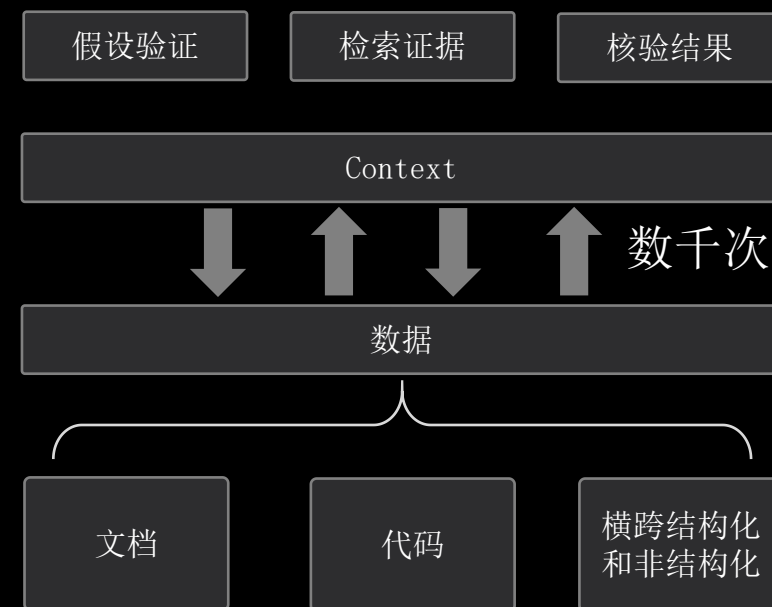
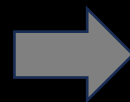
- 为什么需要 Agent
因为 LLM 一次只能干一件事
- 为什么 RAG/Agent 密不可分
反思无处不在
- 一切都是为了 Context



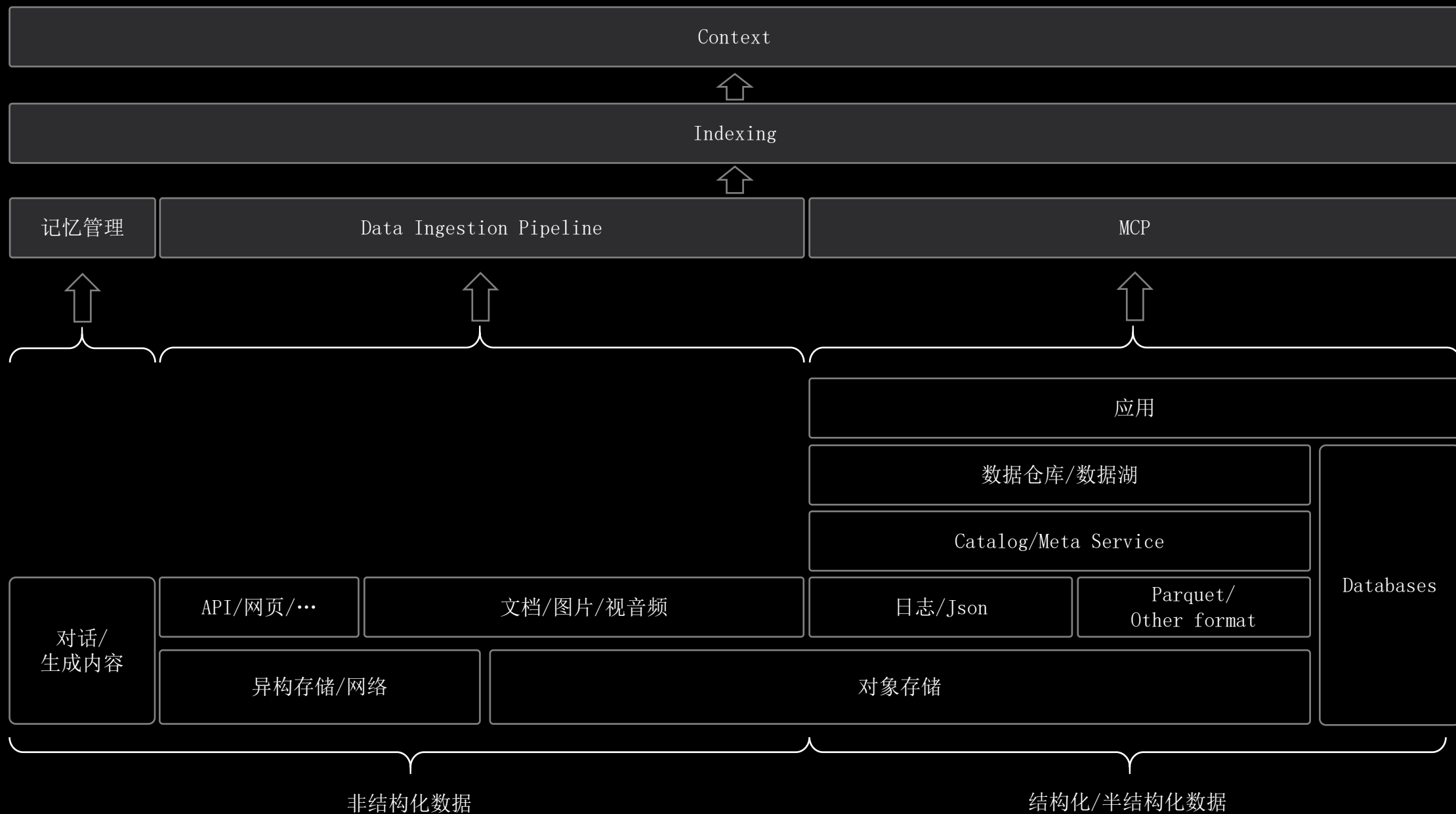
指令也不能无脑填充



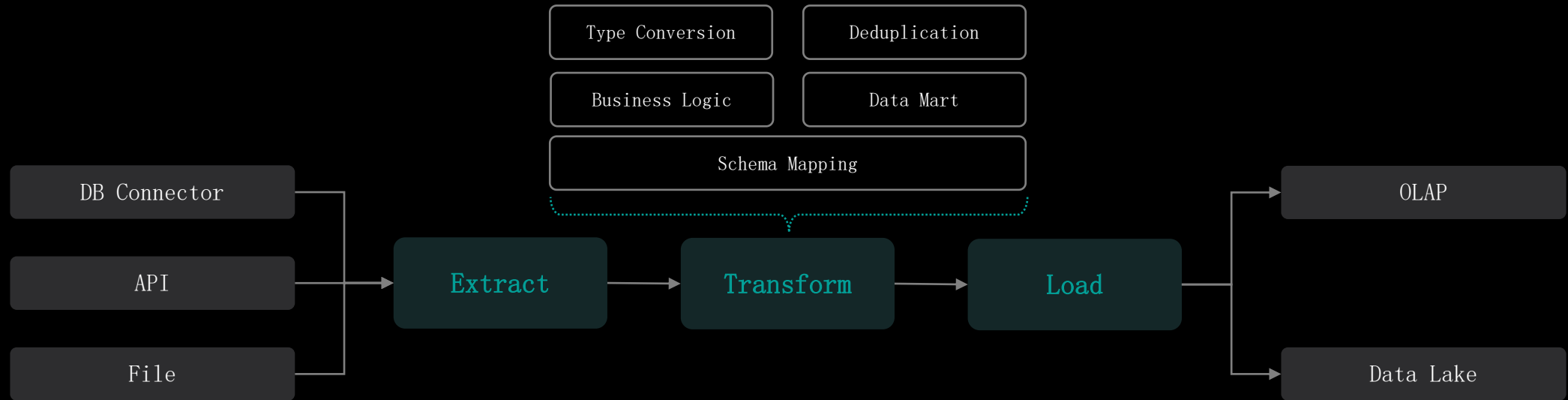
为什么 Retrieval 对 Agents 如此重要



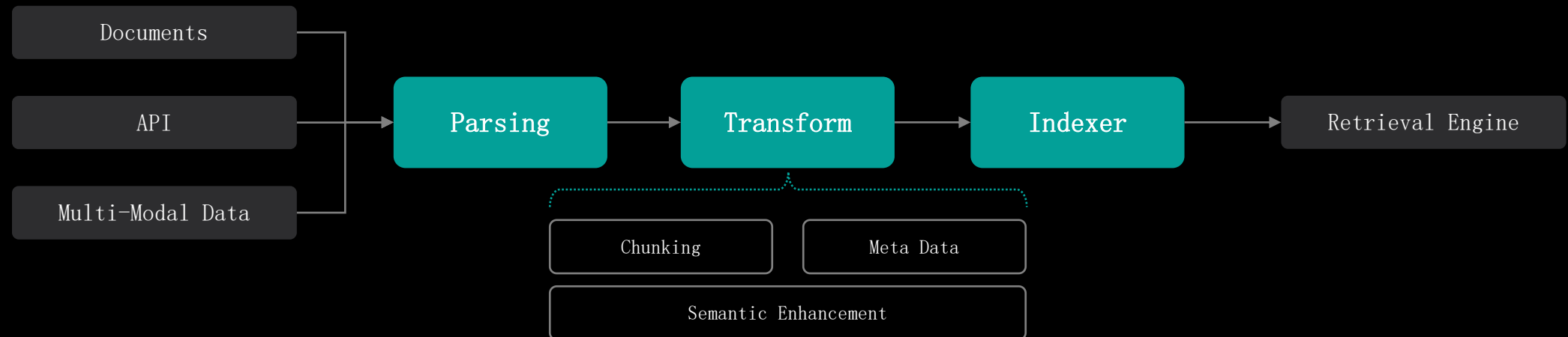
Agent Context 需要什么样的数据底座



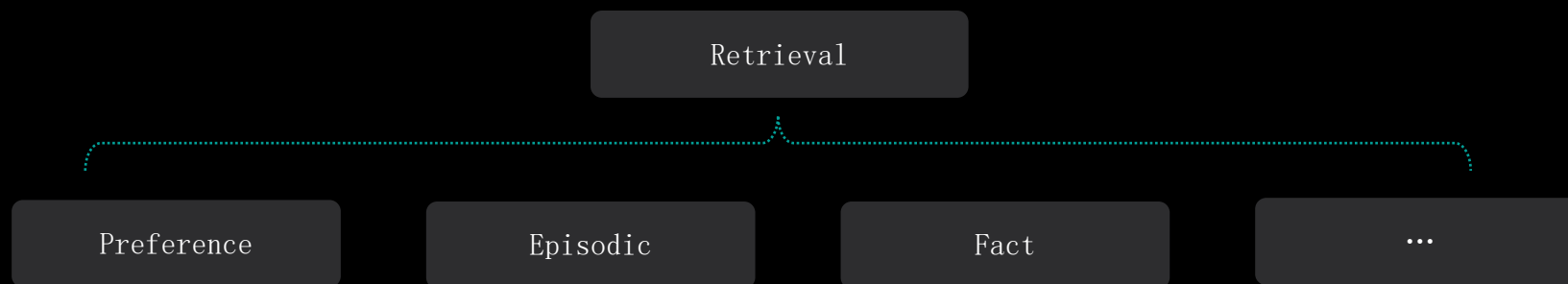
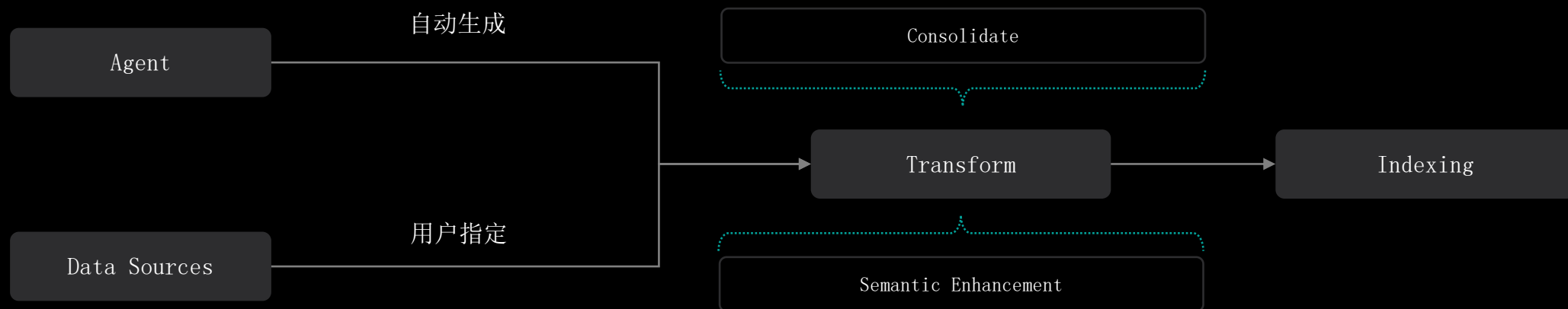
从 ETL 到 PTI



VS



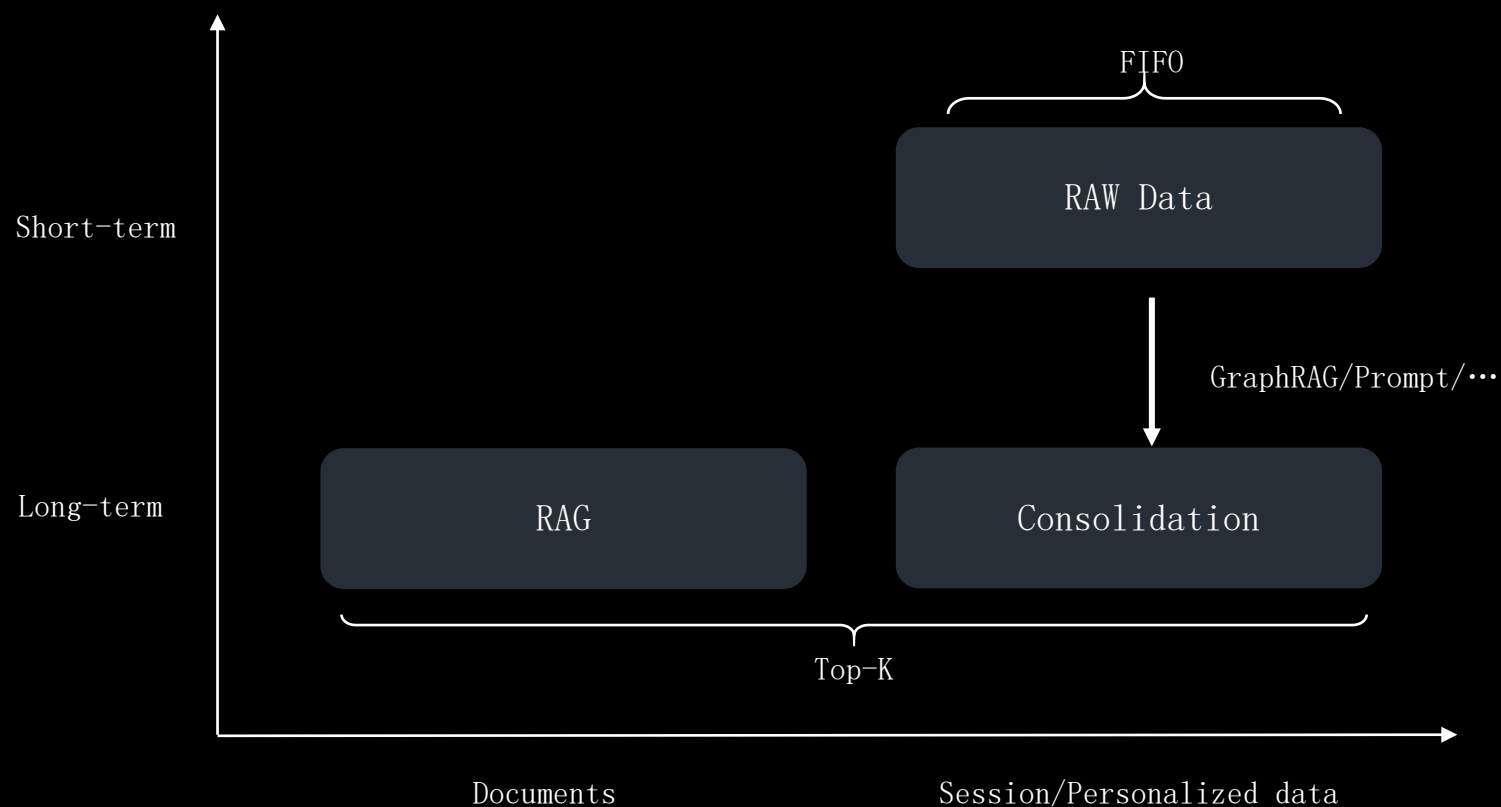
Memory 和 RAG——区别仅在于保存的数据



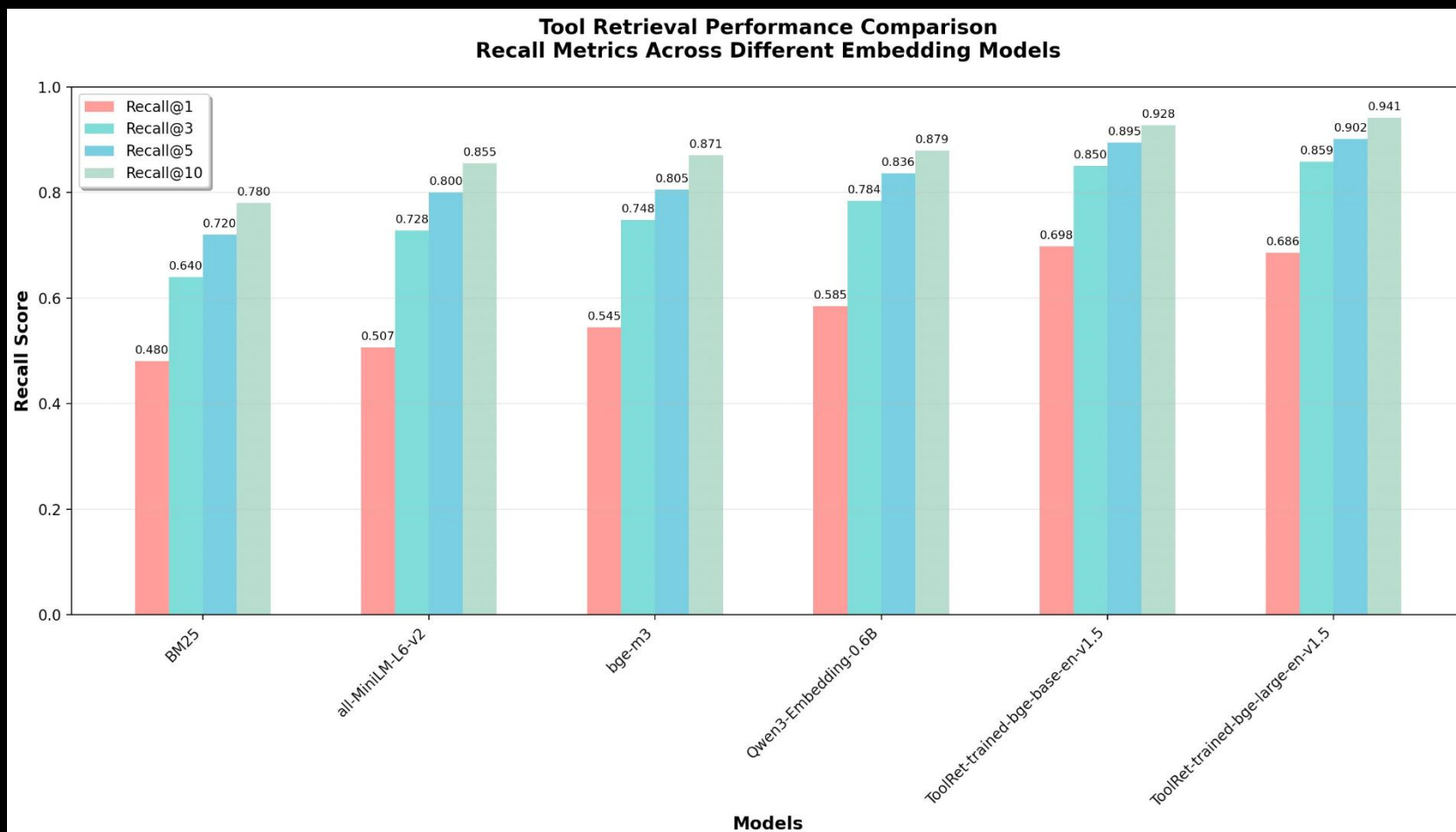
记忆可以成为独立平台么？



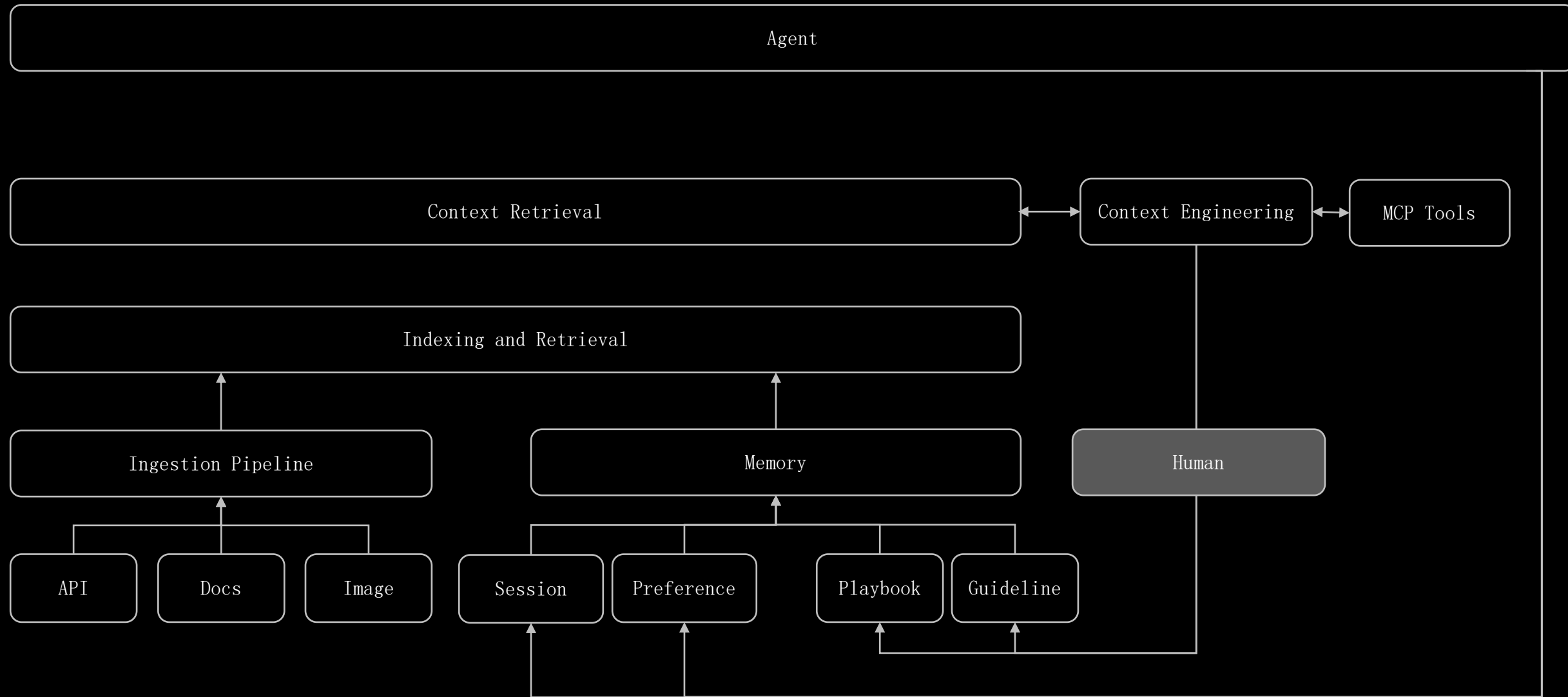
当下的 Memory 定义



Tools 也是需要 Retrieval 的



Tools 和 Guideline

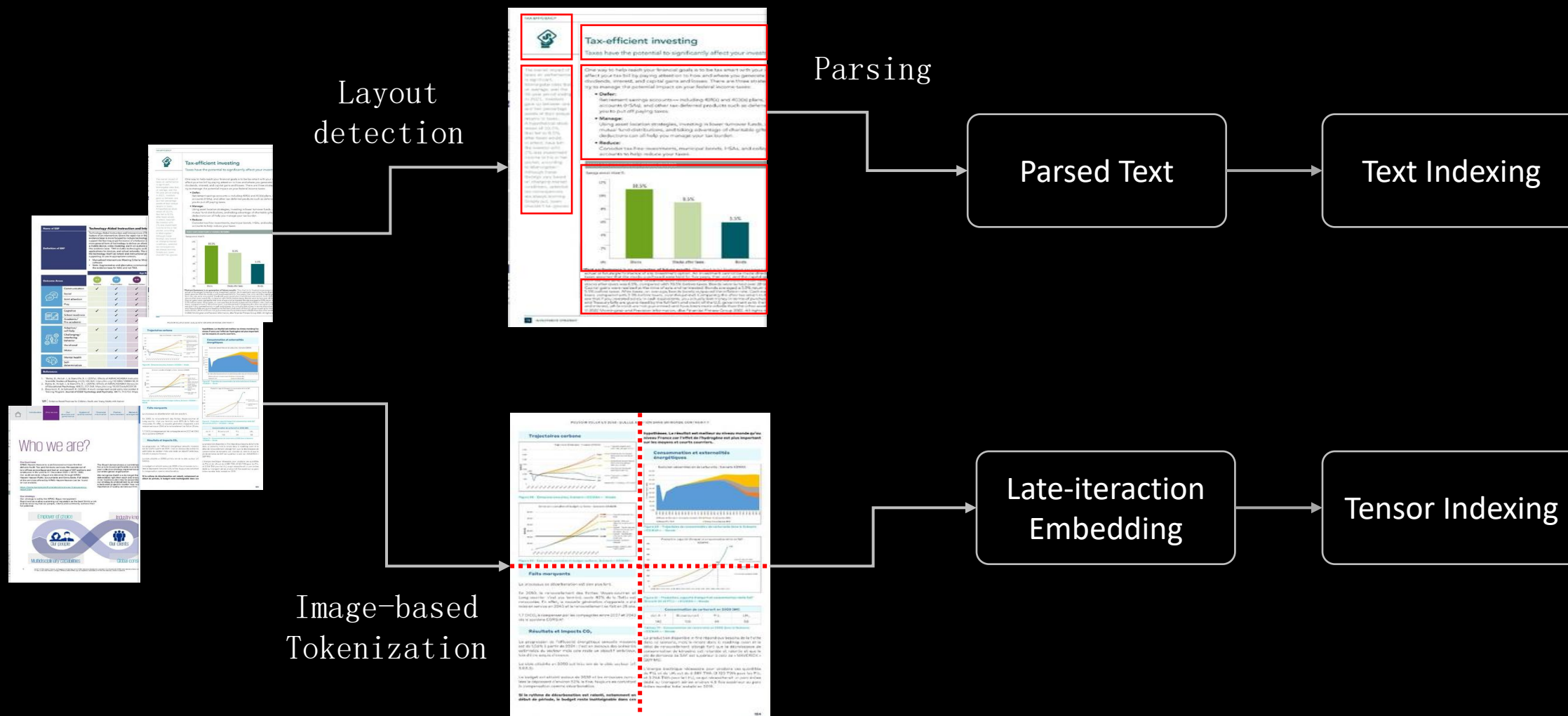




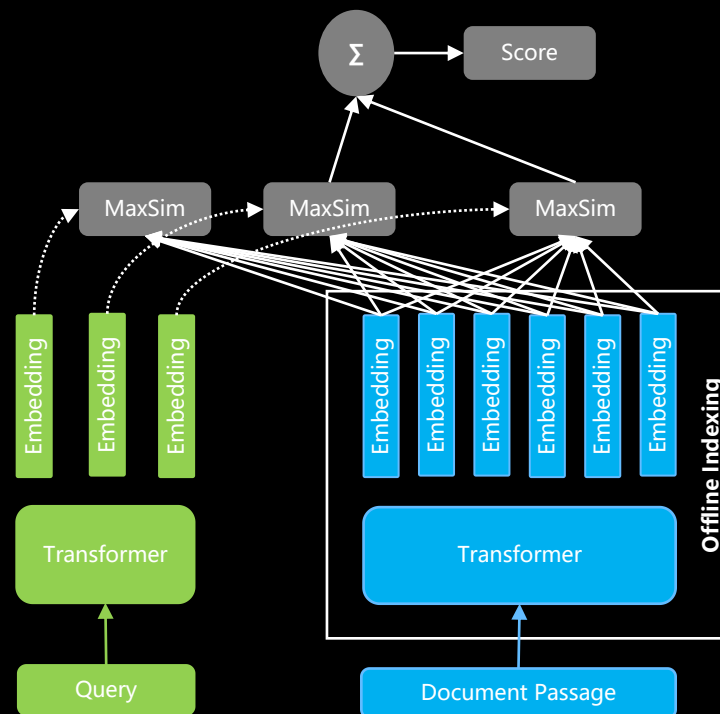
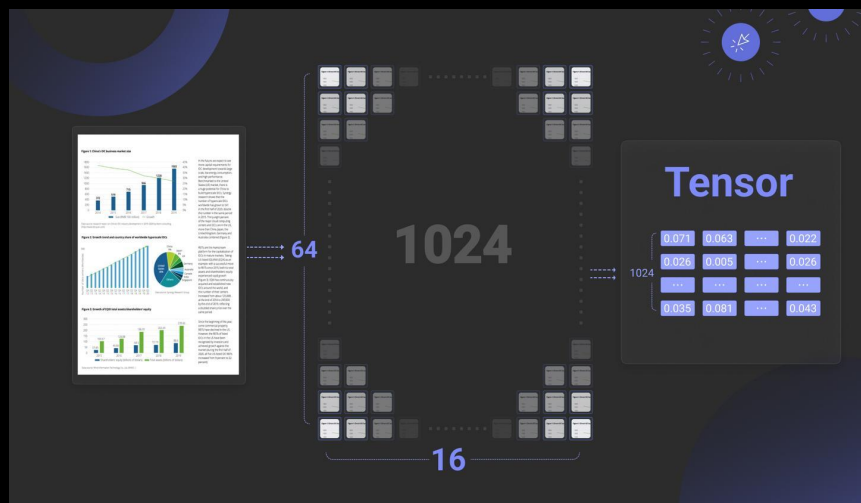
维度	Context Engineering（现状）	Context Platform（未来）
Context 创建	FDE 完成	大多自动化
Context 组装	在 Workflow 当中进行 Prompt 硬编码	Context Retrieval
Context 维护	Vendor 手动维护	客户自主维护

- Context 的“产品化”将是解锁下一代AI应用的关键
- 从 RAG 到 Memory 到 Context，代表了 LLM 和 Agent 应用范式的成熟

关于多模态和跨模态 Retrieval 的那些事情



单向量还是多向量 (Late Interaction)



Google DeepMind

On the Theoretical Limitations of Embedding-Based Retrieval

Orion Weller^{*,1,2}, Michael Boratko¹, Iftekhar Naim¹ and Jinhyuk Lee¹

¹Google DeepMind, ²Johns Hopkins University

- 图文结合任务优势显著
- 单向量检索存在瓶颈

➤ 工程化挑战

➤ 召回单元

➤ 存储膨胀/计算复杂度

➤ 训练对二值量化友好的 Embedding

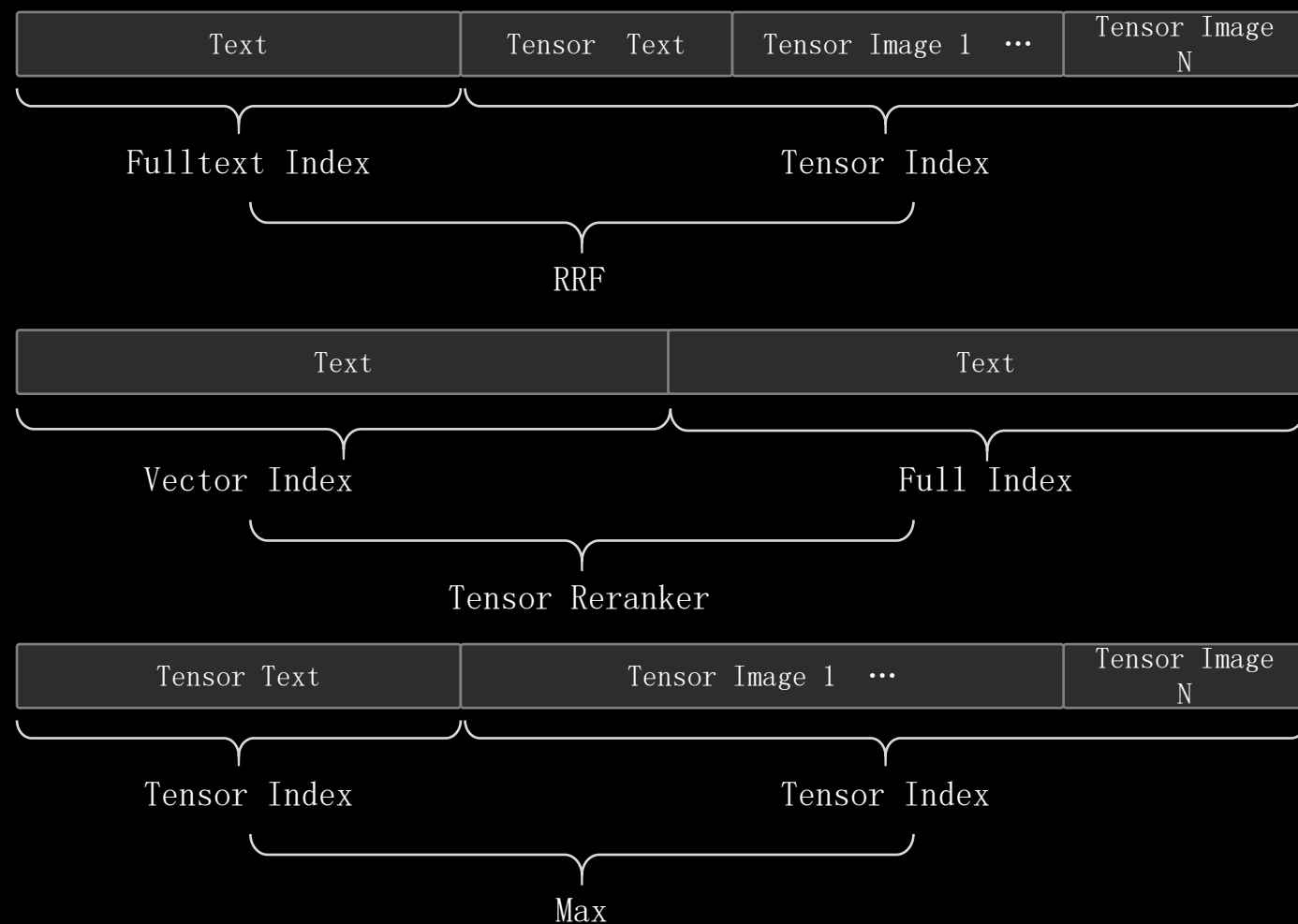
➤ Token 剪枝

➤ MUVERA

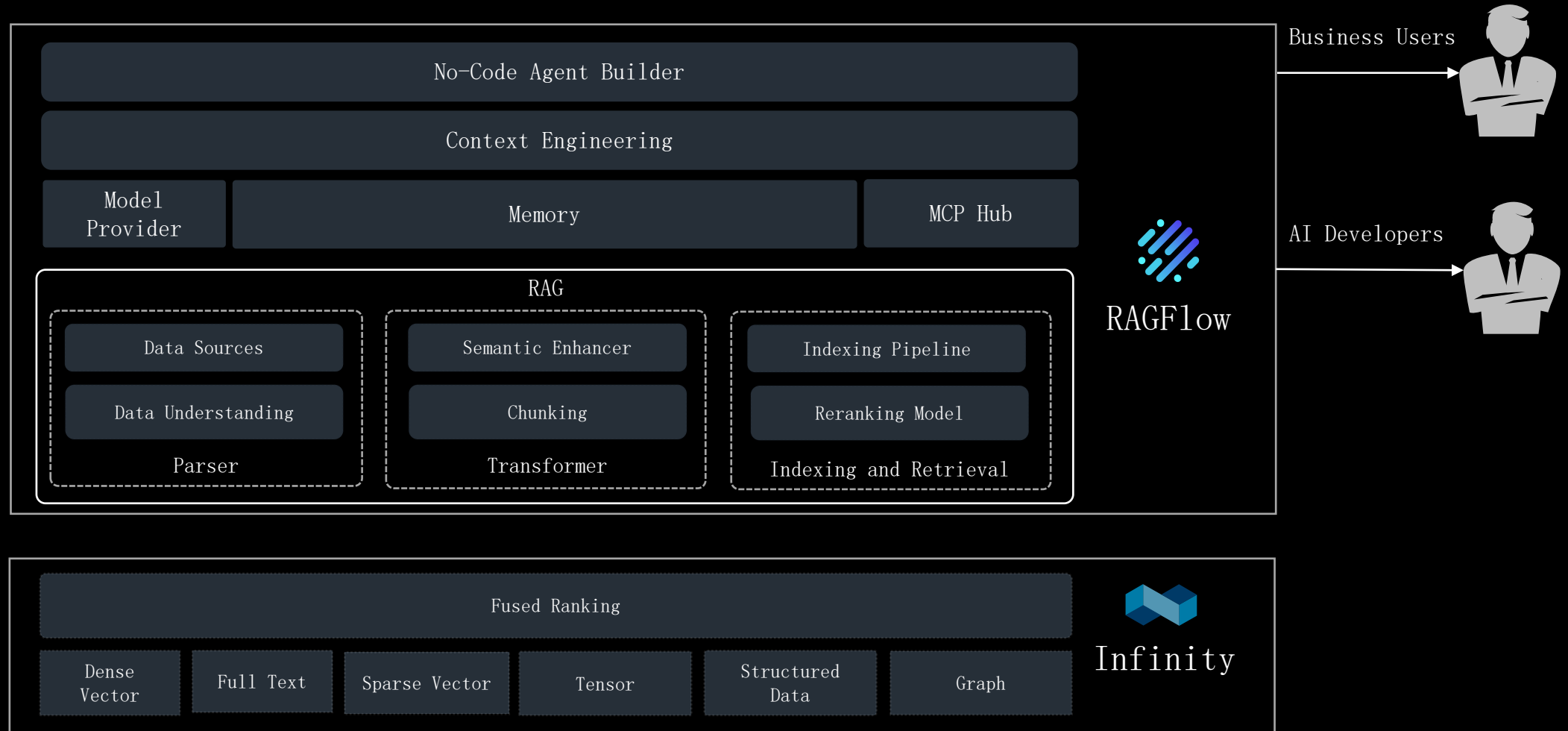
➤ Token Clustering

➤ 改造VLM，根据注意力分布输出剪枝结果

➤ 重新训练，直接在Tokenizer输出少的Token (ModernVBERT)



RAG is Dead? Long Live RAG !





COSCon'25 第十届中国开源年会

众智开源 | Open Source, Open Intelligence

Thanks

