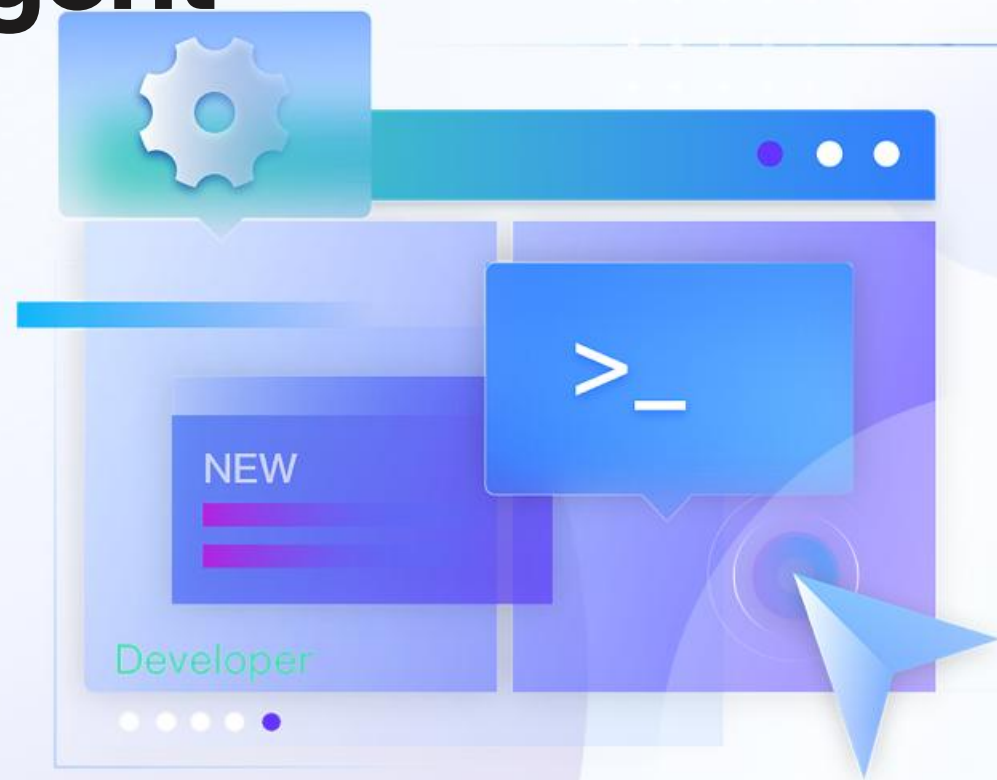


面向数据分析领域的Agent 思考和探索

演讲人

屿你数智AI负责人：eason

2024/07/06





Contents 目录

01 什么是Agent?

02 企业落地面临的挑战

03 数据分析领域的Agent探索-Tigi

04 Agent驱动数据生产变革

05 后续迭代方向与思考

Agent定义和特性

在人工智能领域，人们所指的Agent可以实现自主的意图理解、规划拆解问题、并调用合适的工具完成目标。

自主性

Agent智能体能独立运行，不需要人为干预，可以根据预设规则或学习算法自主决策

社会性

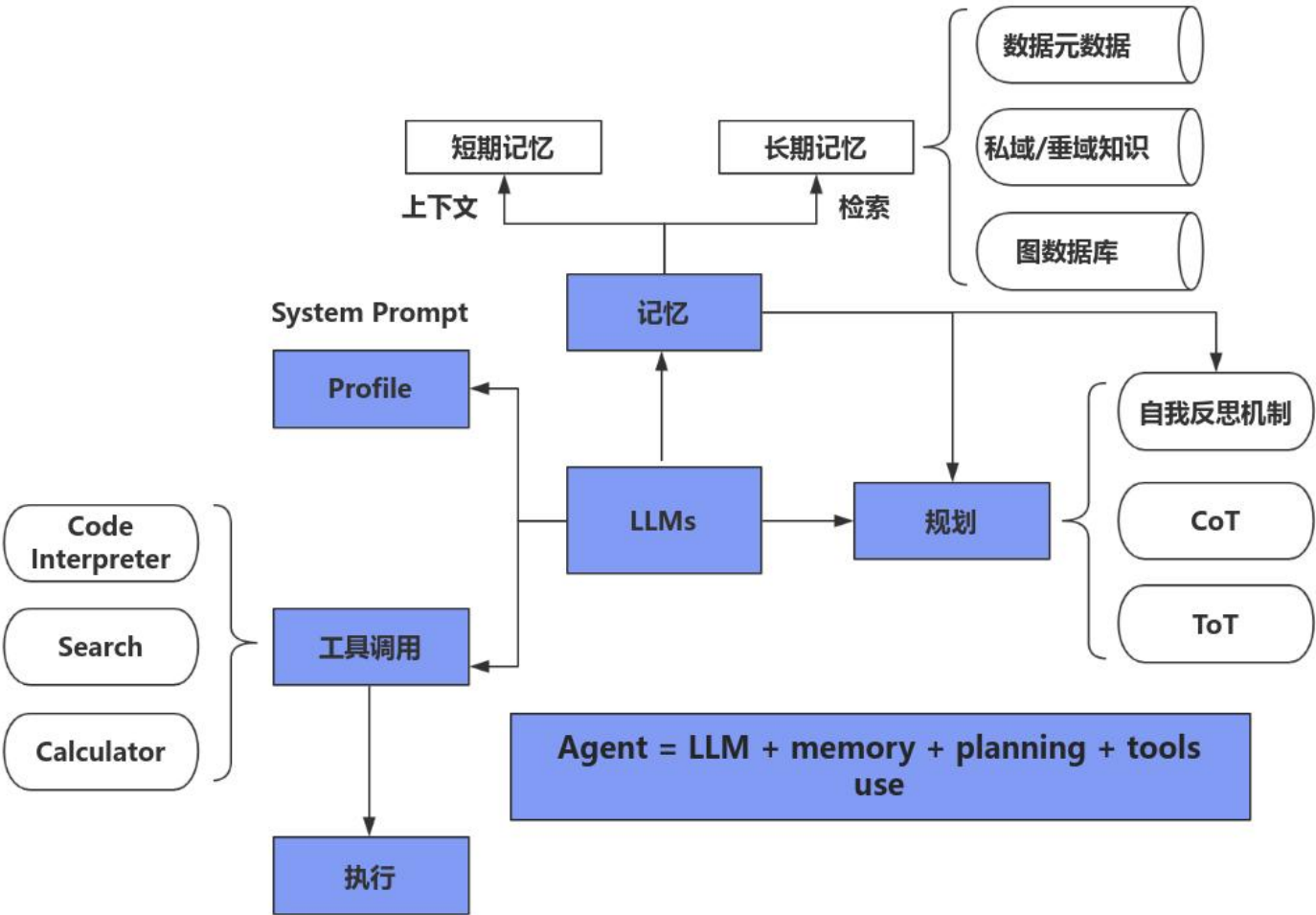
多个Agent智能体能够相互通信、协调与合作，共同完成任务。

反应性

Agent智能体能感知环境变化，并对其作出反应，通过不断调整自身行为以适应环境。

适应性

Agent智能体能够通过学习、进化等方式不断提升自身性能，更好地完成任务。





Contents 目录

01 什么是Agent?

02 企业落地面临的挑战

03 数据分析领域的Agent探索-Tigi

04 Agent驱动数据生产变革

05 后续迭代方向与思考

toC vs. toB

通用场景 toC

领域

解决通用领域的问题，**常识性问题**，比如旅游推荐、写作辅助等，一般采用基座大模型

架构

一般采用Prompt+大模型**端到端**解决问题

严谨性

追求输出的创意性、多样性，**严谨性要求低**

安全性

大多是**公开数据**用来训练基座模型，但是需要关注**种族歧视**和**伦理道德**

企业内场景 toB

解决某领域的**专业性问题**，需要结合垂域和私域内容解决问题，一般采用微调的手段

需对**大模型**进行微调、结合传统的**搜推算法**、现有**工具**串联组合调度以支撑复杂业务场景

对**准确性**要求很高（85线是企业可尝试的基准线）、同时要求**输出稳定**、**完整**和结果**可解释**

核心关注**数据安全和隐私问题**

数据、知识、场景是企业落地Agent的关键要素

企业生产往往都是复杂场景

2024W13, XX产品在IOS端的销售
额是多少?

难点: 时间识别、维度识别、指标
识别、过滤条件识别

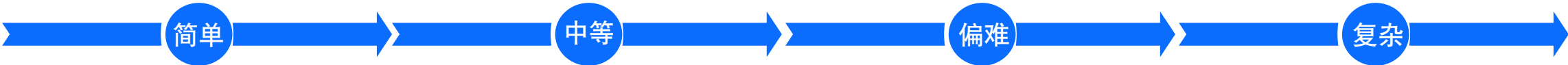
2024W13, XX产品在IOS端卖的怎
么样? 同比上一个周是降了还是升
了?

难点: 时间识别、维度识别、指标
识别、语义模糊、多任务拆解

简单场景

偏难场景

场景复杂度, 由易到难



中等场景

复杂场景

2024W13, XX产品在IOS端卖了多
少钱?

难点: 时间识别、维度识别、指标
识别、语义模糊

2024W13, 大盘的xx指标是否正常?
如果有异动, 是由什么导致的? 可
以用连环替代法来分析。

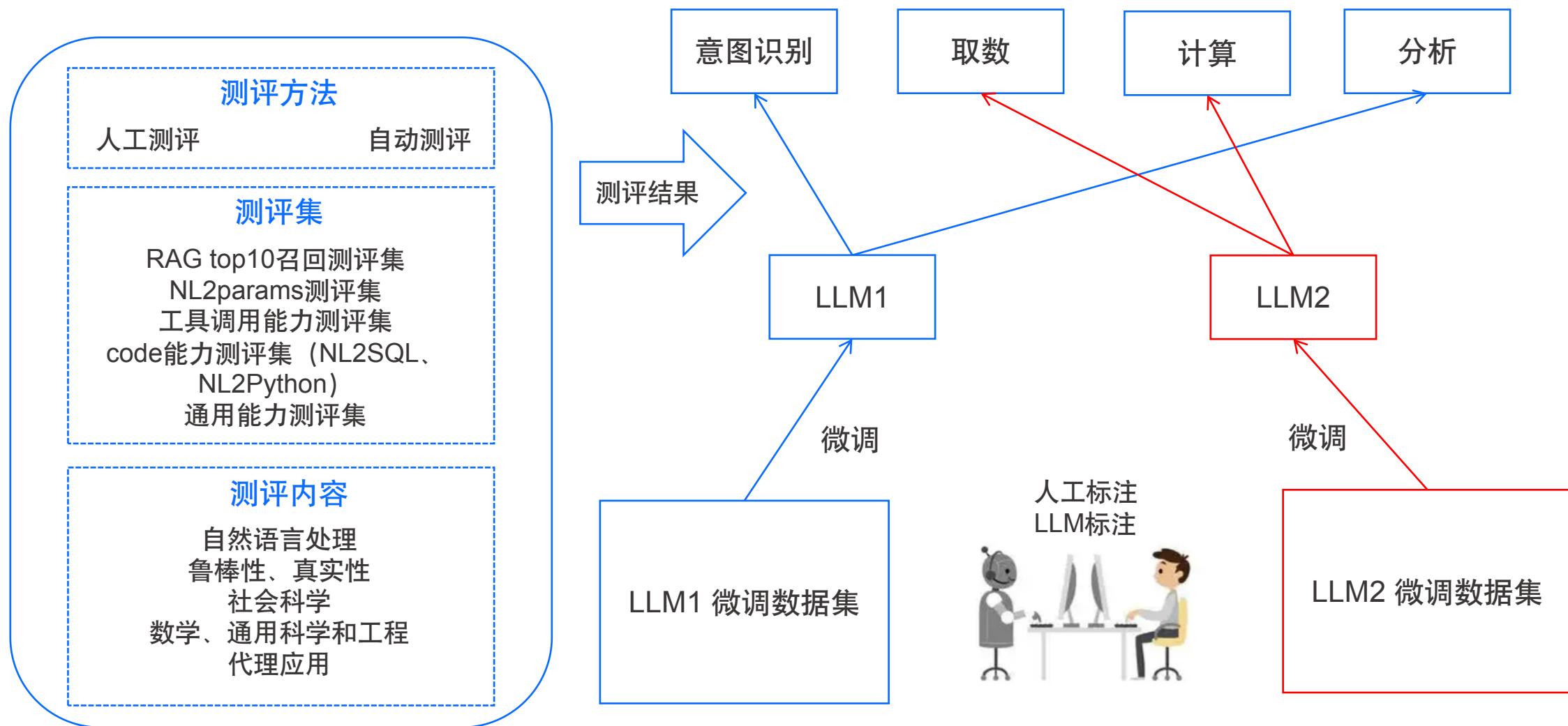
难点: 时间识别、维度识别、指标
识别、语义模糊、多任务拆解、连
环替代法、计算、归因分析



Contents 目录

- 01** 什么是Agent?
- 02** 企业落地面临的挑战
- 03** 数据分析领域的Agent探索-Tigi
- 04** Agent驱动数据生产变革
- 05** 后续迭代方向与思考

各个场景用什么大模型合适？



知识可信加工

Garbage in, Garbage out

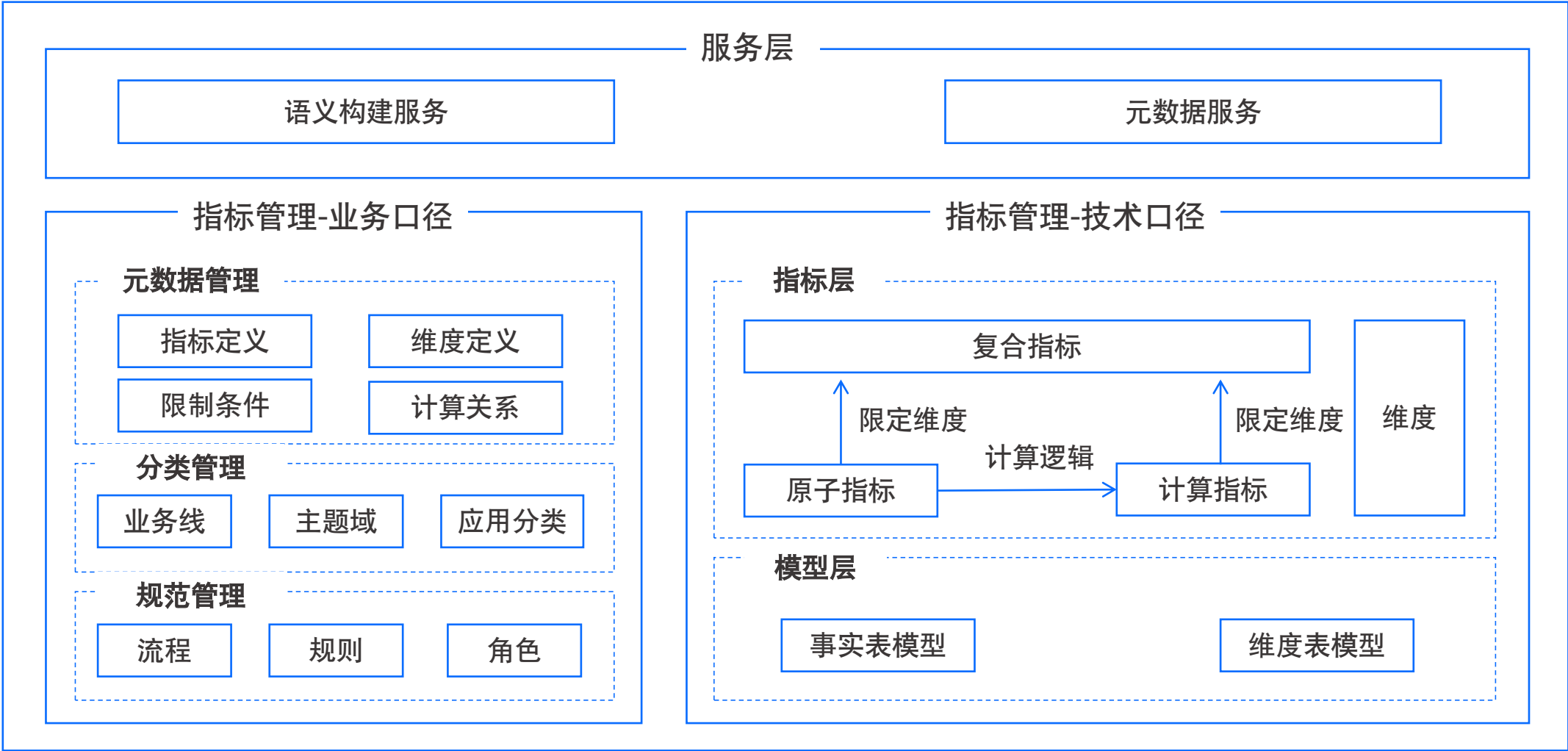


Garbage in, Quantity out



以上是通用的清洗模块和流程，可以保障知识库入库的**基线**，若要更好的效果，还需要针对公司内、领域内流行格式、排版做对应的清洗规则。

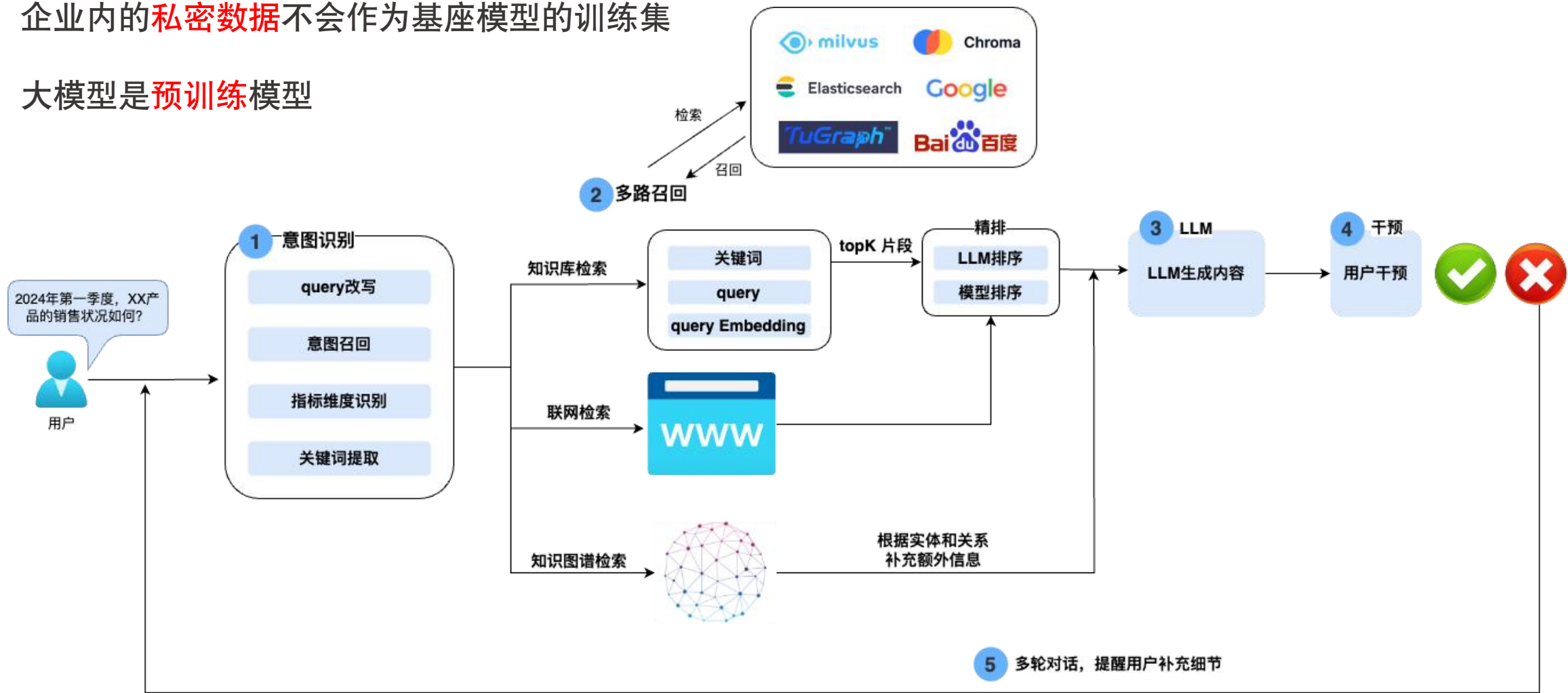
指标管理平台



基于知识图谱的RAGs

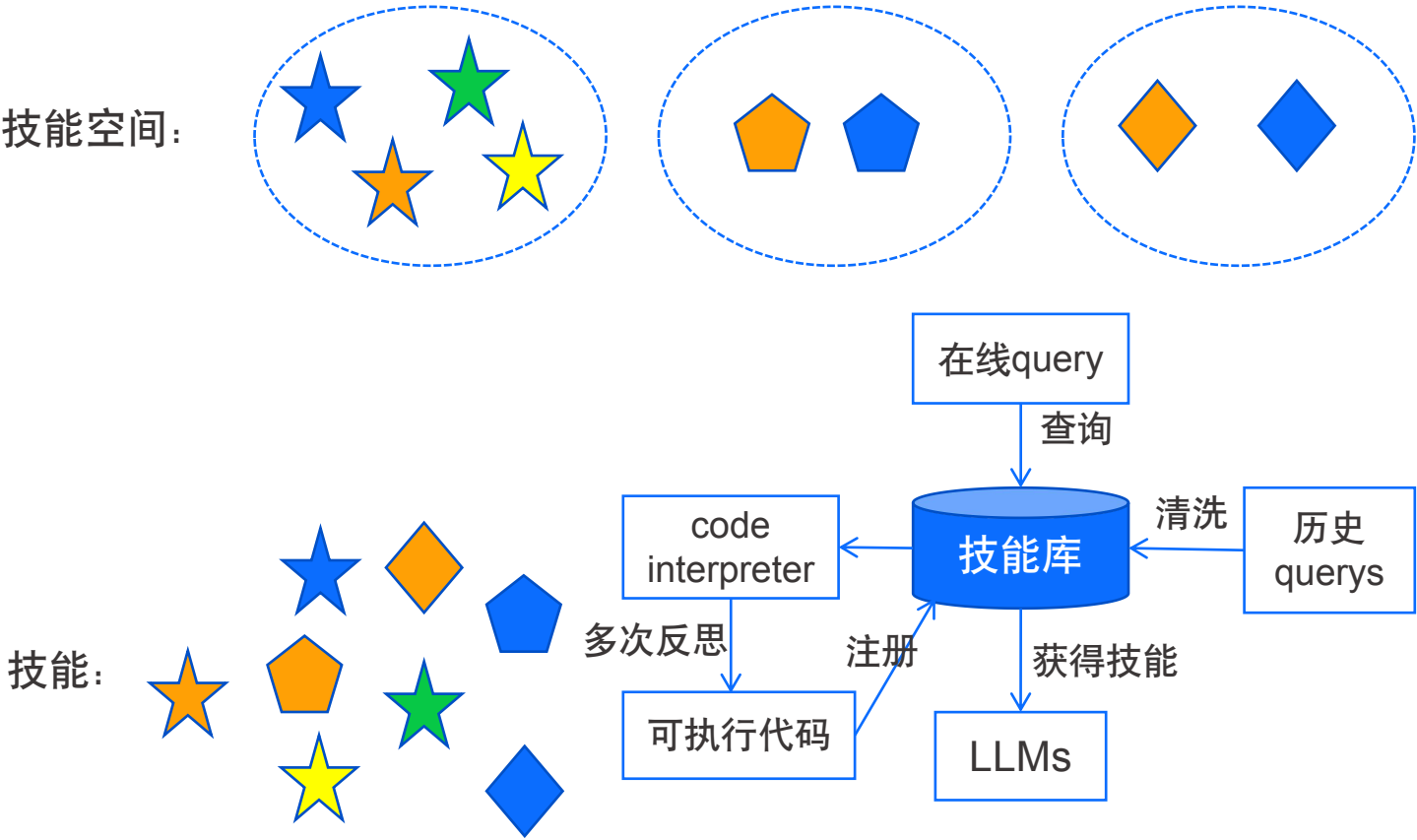
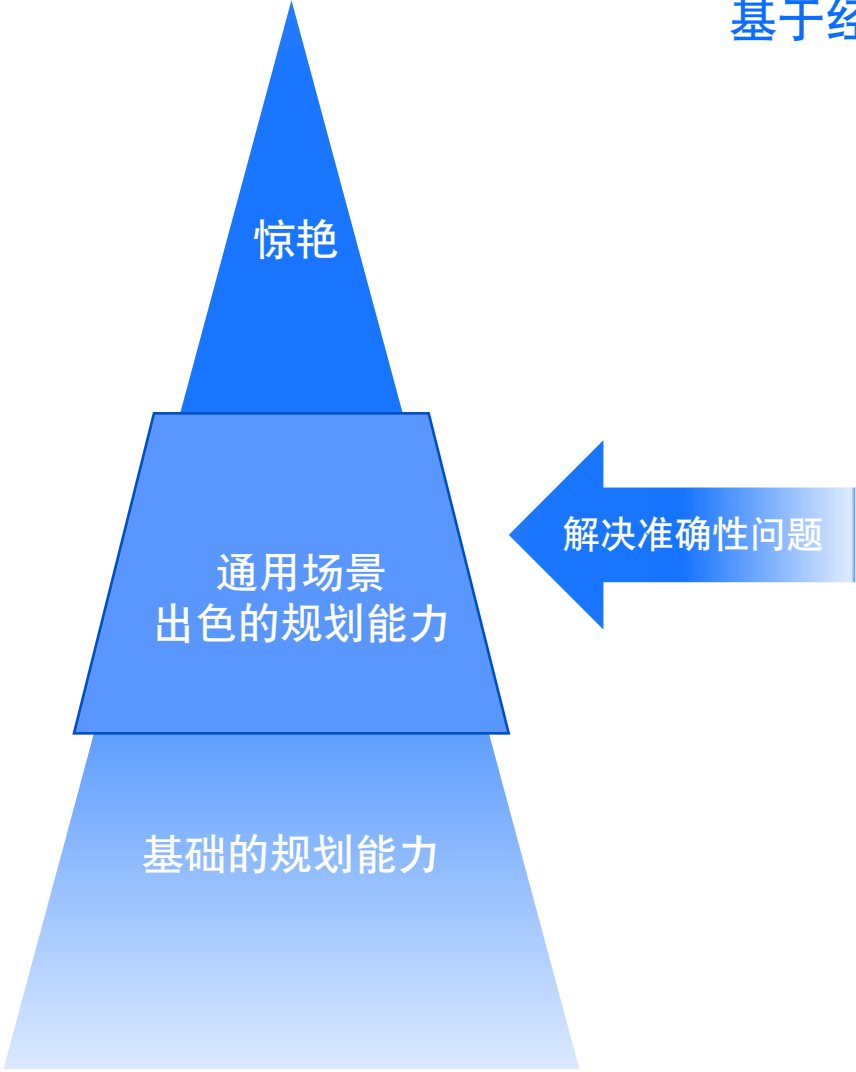
企业内的**私密数据**不会作为基座模型的训练集

大模型是**预训练**模型



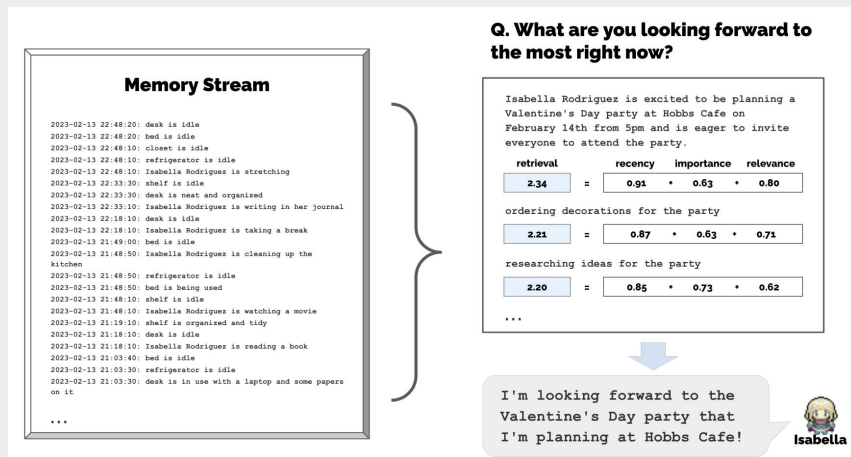
技能库与技能空间

基于经验，枚举通用场景的问题，缩小问题空间



记忆召回

短期记忆召回



$$\gamma_E = \alpha_{recency} \cdot S_{recency} + \alpha_{relevance} \cdot S_{relevance} + \alpha_{importance} \cdot S_{importance}$$

取数和分析重要度高，计算重要度低

长期记忆召回

字典
keywords



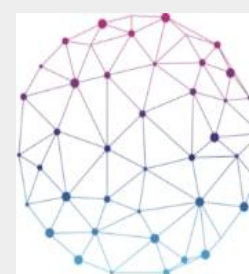
关键词索引
BM25算法

语义记忆
embedding



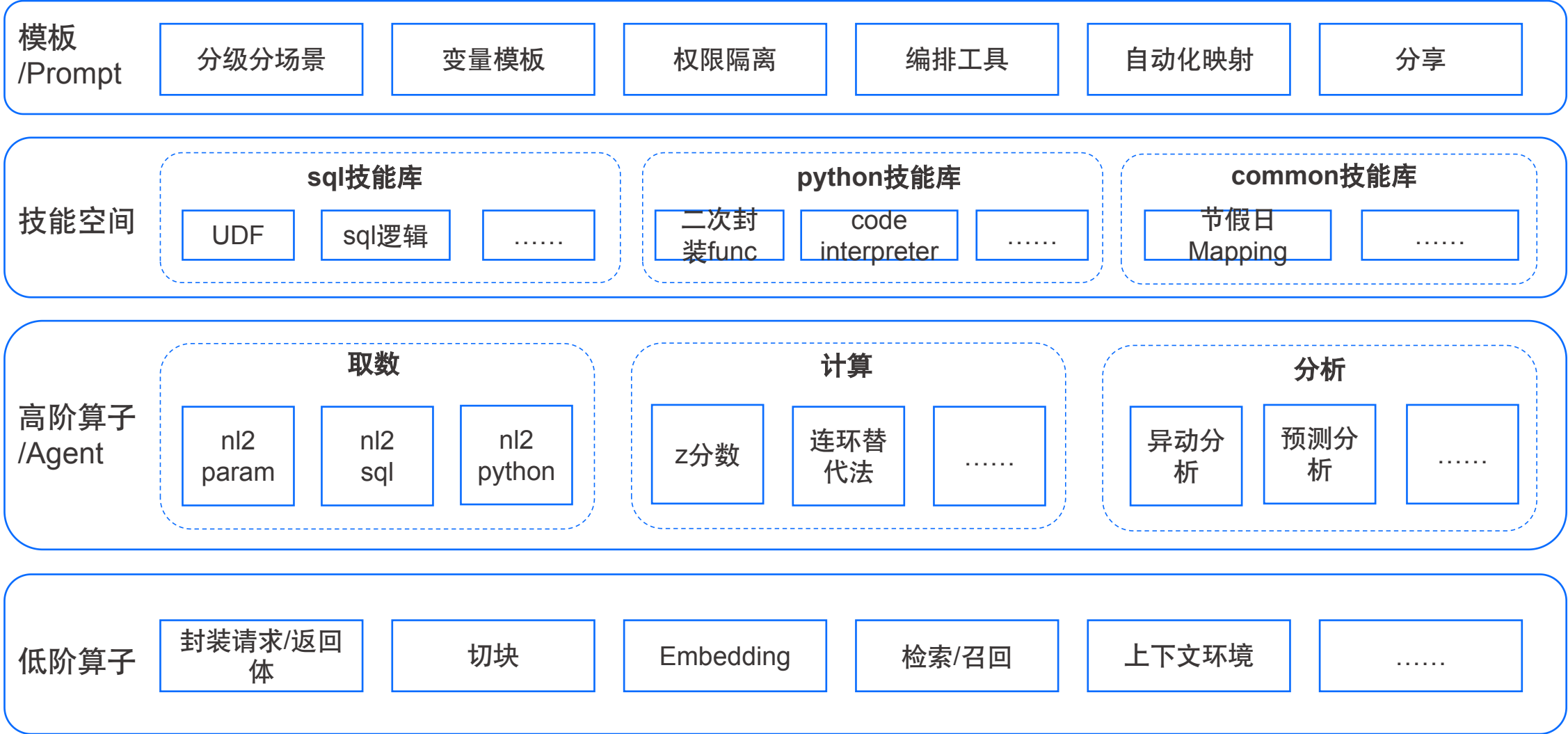
语义相似匹配

知识图谱
graph



关系检索
属性检索

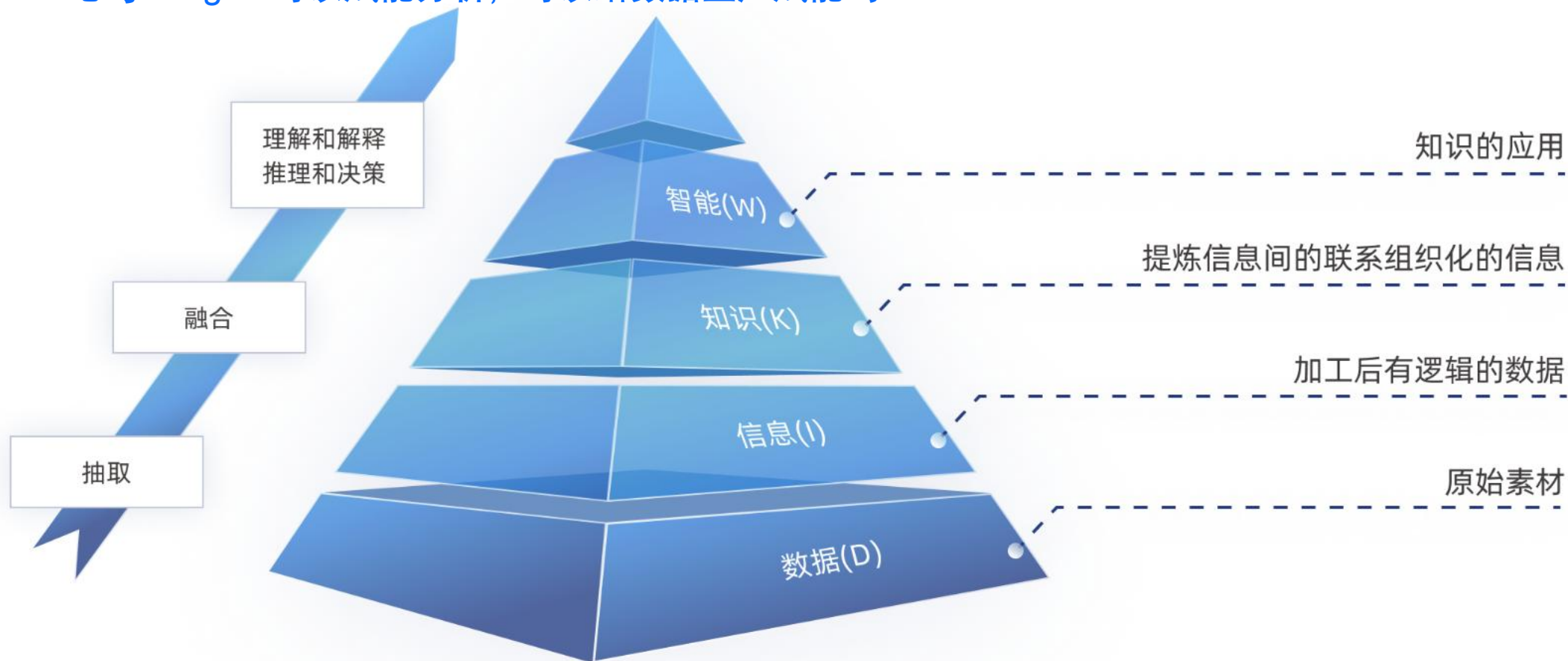
复杂场景解决方案：Agentic WorkFlow



总结

之前**会**分析的人，可以通过Tigi将方法论具象化、模板化。
之前**不会**分析的人，在Tigi的加持下，变得**会**分析。

思考：Agent可以赋能分析，可以给数据生产赋能吗？

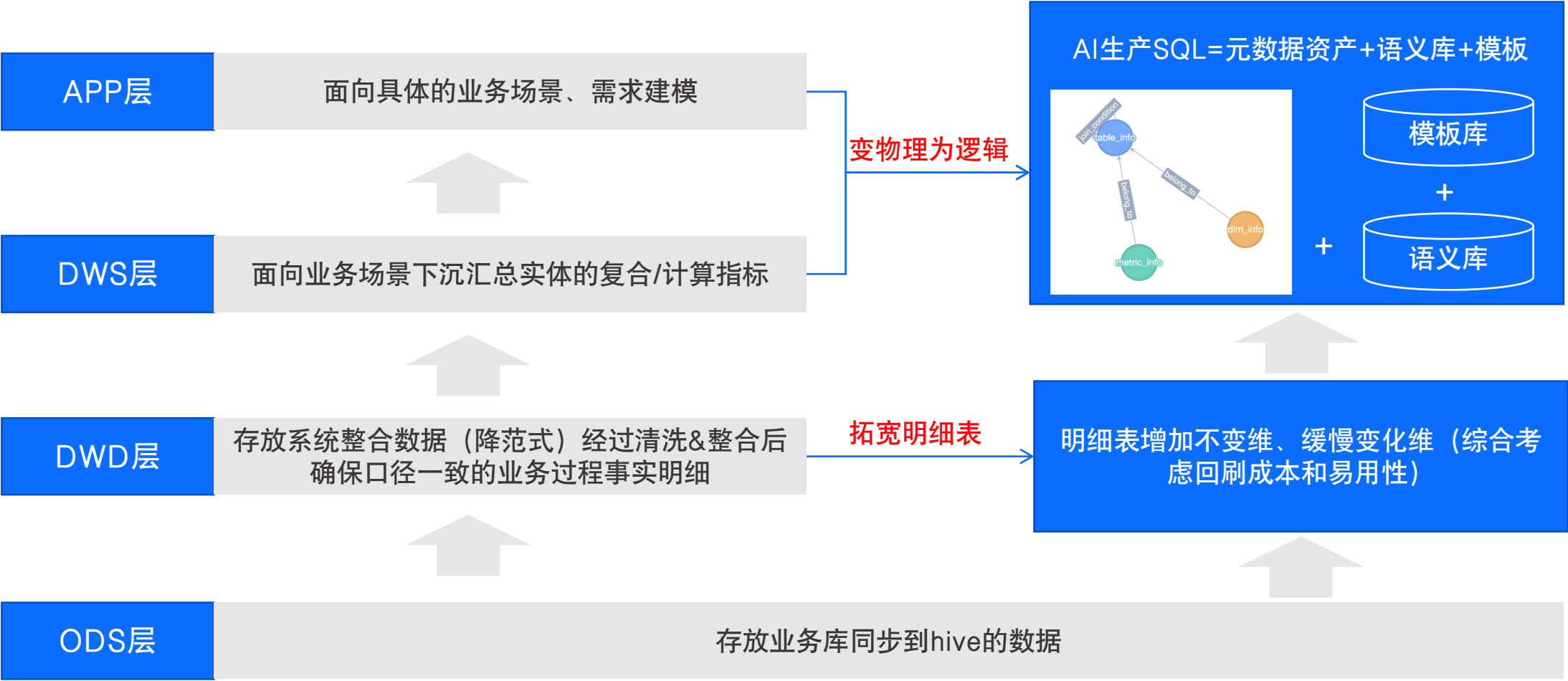




Contents 目录

- 01** 什么是Agent?
- 02** 企业落地面临的挑战
- 03** 数据分析领域的Agent探索-Tigi
- 04 Agent驱动数据生产变革**
- 05** 后续迭代方向与思考

变革一：做厚数据集市 -> 做宽明细数据层



变革二：进一步挖掘元数据的价值

人能看懂

AI能看懂

数据治理

统一metadata，保障数据一致性和准确性，驱动数仓建模规范落地
优化查询逻辑，保障数据查询及时性

应用治理

通过热度元数据，可以识别冷热看板，为后续的看板整合下线提供输入



数据生产

将元数据资产整合成语义库、模板库和知识图谱
直接作用于AI生成SQL的输入

AI加持下，元数据将直接应用于数仓生产

传统元数据的价值体现在数据和应用的治理上

元数据（用户历史的query请求、表schema、指标业务口径、计算公式、维度枚举值、血缘、热度、查询耗时等）



Contents 目录

- 01** 什么是Agent?
- 02** 企业落地面临的挑战
- 03** 数据分析领域的Agent探索-Tigi
- 04** Agent驱动数据生产变革
- 05** 后续迭代方向与思考

大模型幻觉

流程上每个节点的幻觉都会累加

思路

确认查询语句

数据集ID

*日期类型:

日周月

*日期选择: 2024-05-20 至 2024-06-30

*查询指标:

查询维度:

大类

筛选条件:

对比类型:

确认查询

关键环节需要人工二次确认，不断优化人工确认环节的用户体验（减少确认次数）

RAGs建模

元数据采集问题

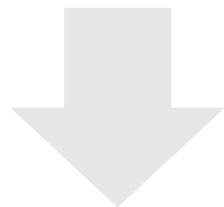
元数据散落在各个平台上，治理平台、指标管理平台，数据挖掘加工的结果

元数据建模不规范问题

元数据组织模式尚处于探索中，多种范式并存，给使用带来复杂度

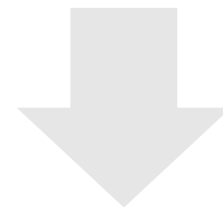
元数据重复存储问题

元数据存储了多份，并且随着存储引擎的增加而增加



思路：

元数据采集同步规范
元数据加工规范
元数据建模规范



探索新型数据库
(兼备Embedding和图构建)

防安全风险

训练中敏感数据清洗

宗旨：用户隐私信息不入训练集，密文信息训练时不降低可用性

推理中敏感数据防外泄

宗旨：用户不感知加密解密过程，大模型基于密文输入进行推理不损失效果

效果测评

ABTest对比加密前后的大模型输出效果

加密组件

加密码表Mapping、数值混淆加密等

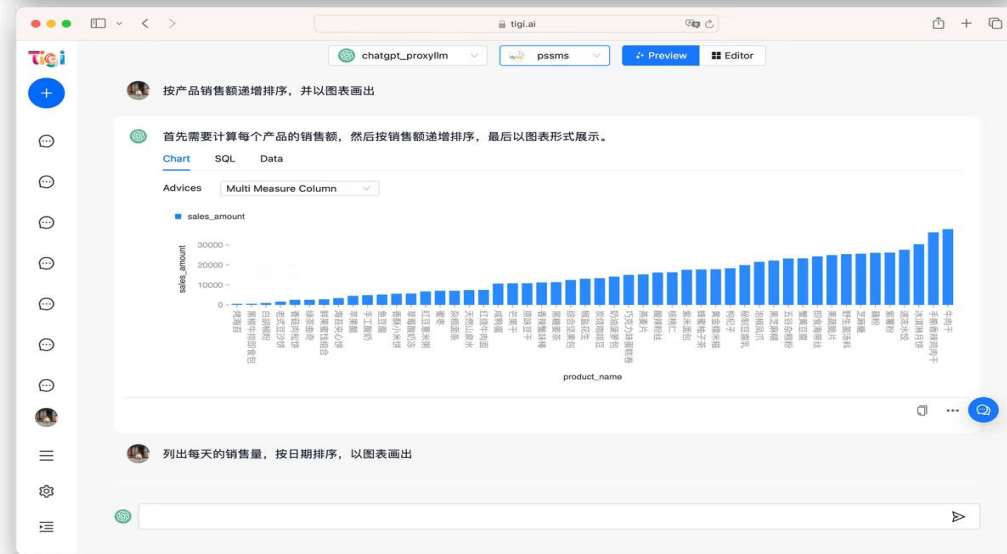
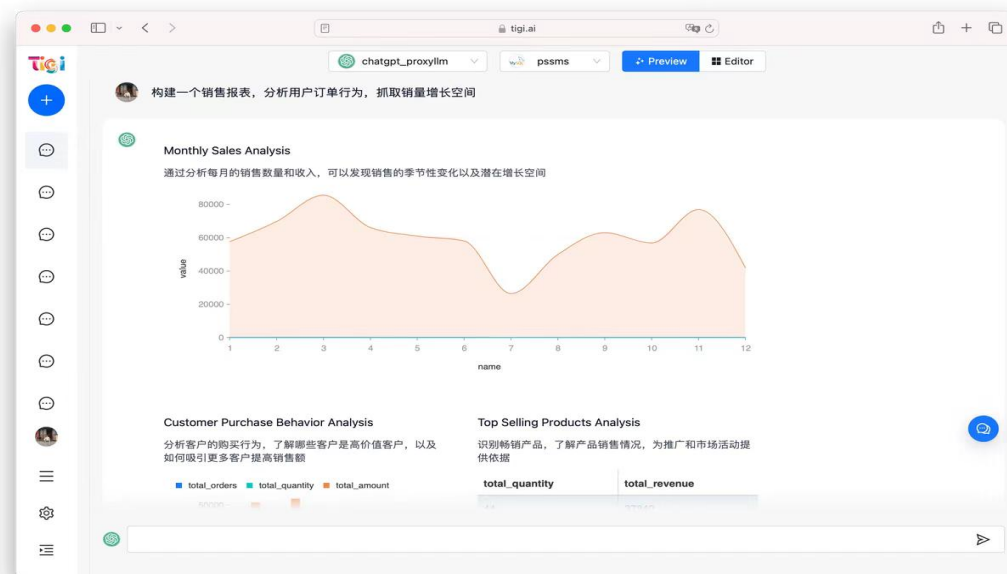
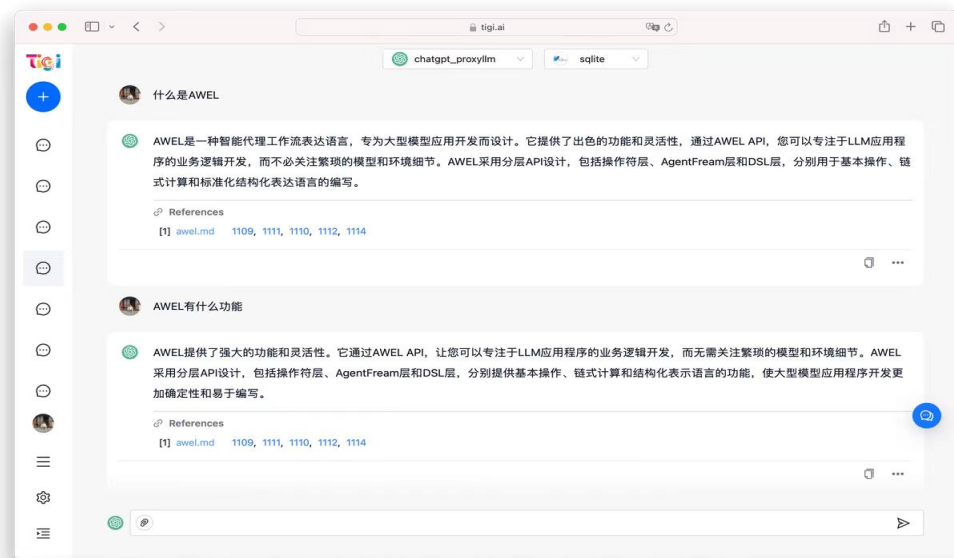
识别隐私数据组件

正则、字典表、模糊匹配、OCR等

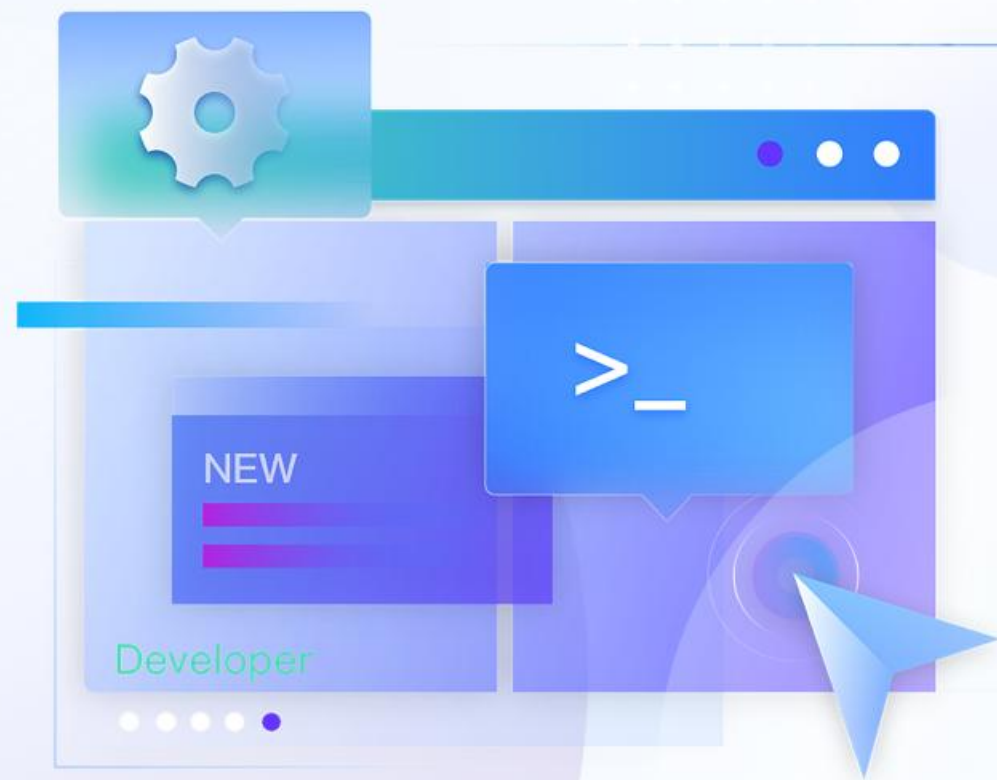
输入数据

个人信息：身份证、手机号、Email、地址等
ID类：订单ID、商品ID、门店ID等
标识类：MAC地址、IP地址、WX_OPEN_ID等

产品演示



期待与大家
一起学习和交流





微信关注Tigi，了解更多数智决策解决方案

Tigi体验链接: <https://tigiai.cn/>



一起探索Agent在数据分析领域的落地