

数据驱动的智能诊断系统：多智能体系统在生产环境中的技术落地与实践

演讲人：赵庆杰

阿里云 Serverless 基础架构负责人 & AgentRun 产研负责人

AiCon
全球人工智能开发与应用大会

目录

01

引言：智能体实践背景与挑战

02

实践路线：运维智能体的探索路径

03

总结与展望

极客邦科技 2026 年会议规划

促进软件开发及相关领域知识与创新的传播



参会咨询



查看会议

北京

1200人

QCon

全球软件开发大会

会议时间：4月16-18日

- Agentic Engineering
- AgentOps
- 下一代模型架构与推理优化
- AI 原生基础设施
- 知识工程实践
- AI 安全

深圳

1000人

AiCon

全球人工智能开发与应用大会

会议时间：8月21-22日

- Agentic AI
- 轻量化与高效推理
- 多模态应用
- AI + IoT 场景实践
- AI 工业化落地

北京

1000人

AiCon

全球人工智能开发与应用大会

会议时间：12月18-19日

- 大模型架构创新
- 多模态 AI 产业融合
- 具身智能
- AI for Science
- 大模型安全

4月

6月

8月

10月

12月

AiCon

全球人工智能开发与应用大会

会议时间：6月26-27日

- AI Infra 系统工程
- 多 Agent 协作与实践
- 多模态融合
- 模型训练与推理创新
- 数据平台与特征服务

上海

1000人

QCon

全球软件开发大会

会议时间：10月22-24日

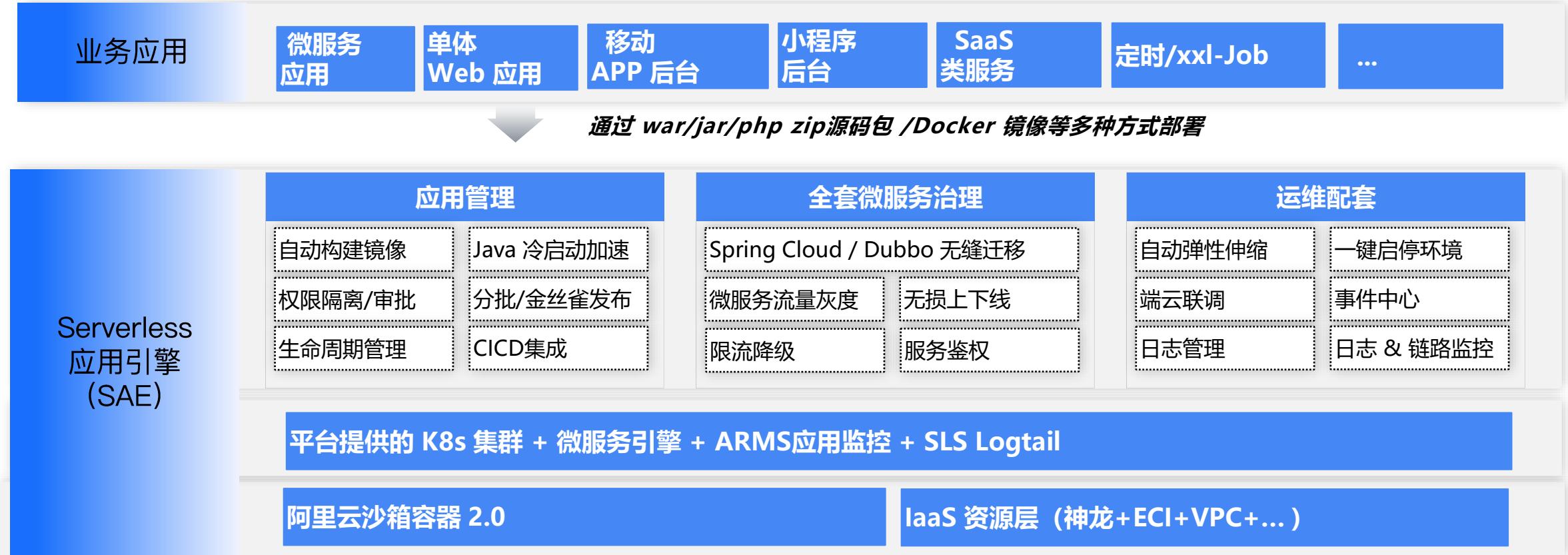
- AI Agent
- Vibe Coding
- 智能可观测
- 推理基建
- 模型攻防
- AI x 创造力

上海

1200人

01 引言：智能体实践背景与挑战

实践背景—承载百万级别的应用的容器 PaaS 系统





重复的疑难问题耗费人力



实例健康检测
OOM
物理机故障
网络故障

实例重启



机器负载过高
过保机型
机器组件故障
机型不匹配

性能问题



P2报警过多
无法发现问题
没有有效的总结和处理

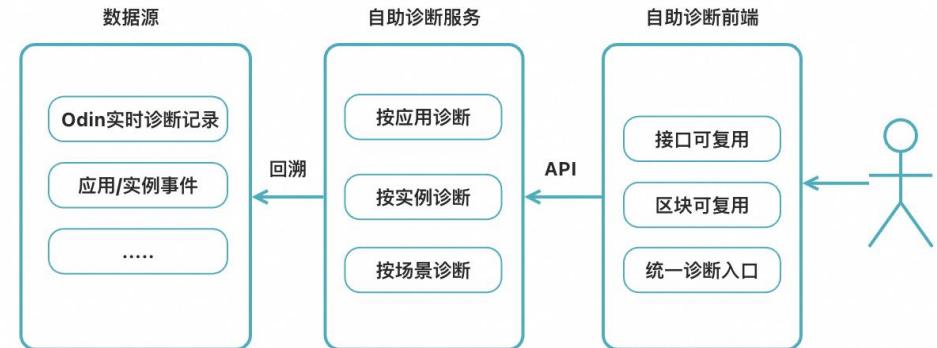
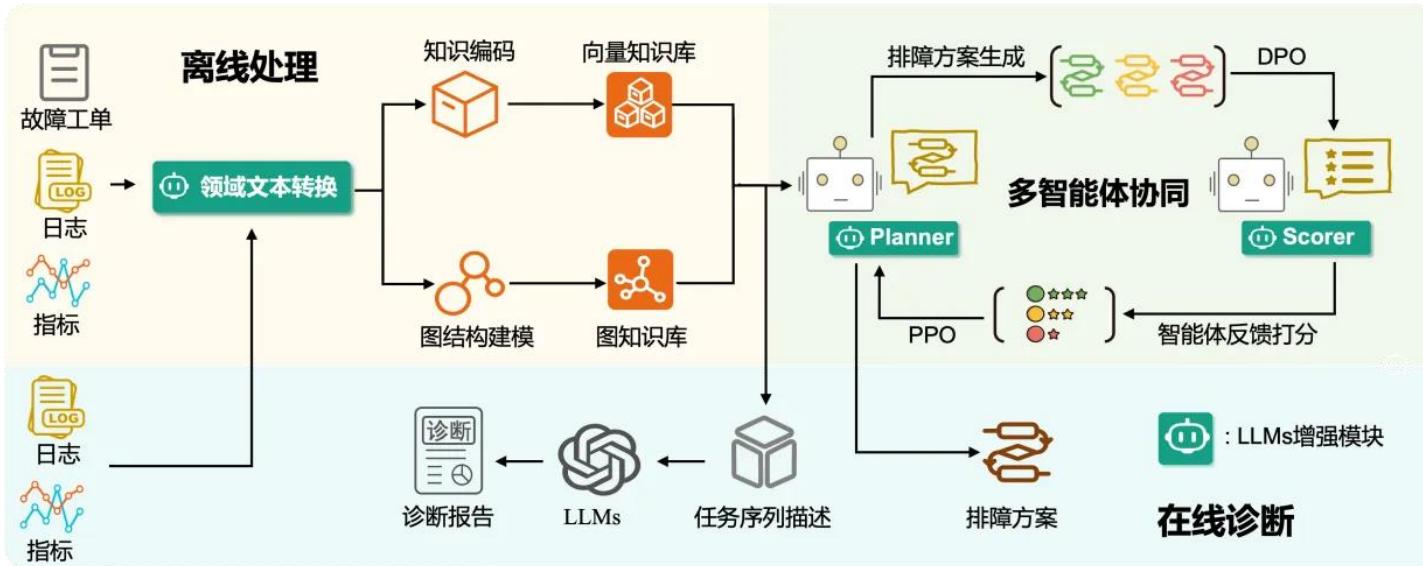
报警处理



传统协议改造成本高
交互模块太多
中间过程步骤多

Chat Ops

■ Serverless：智能诊断智能体



01 运维智能体的探索路径

三个关键阶段



01 阶段一：静态工作流



挑战



数据丰富度不足
通用思考模式不成熟
在 RAG 之上做补充

LLM 未成熟

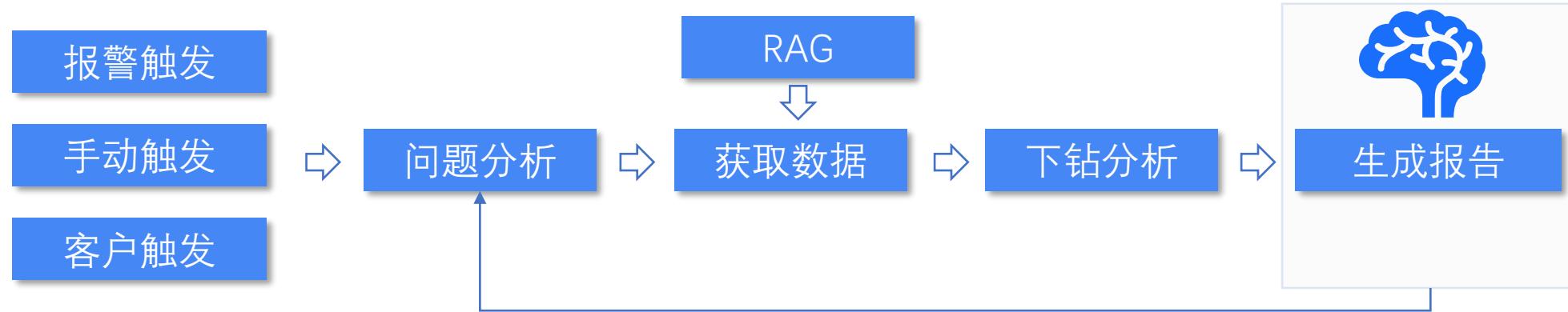


各类数据查询不统一
数据的连接性无法表达
数十种不同数据来源
RAG 的补充

数据查询不统一

原始状态，简单处理日常的答疑

架构方案



This screenshot shows a workflow definition interface with a sidebar listing various services and their deployment status. The main area displays a complex workflow graph with nodes like '开始' (Start), '模型选择' (Model Selection), '查询订单或物流信息' (Query Order or Logistics Information), '提取订单ID' (Extract Order ID), 'ORDERS', '查询物流信息结束' (End of Query Logistics Information), '归纳问题' (Summarize Problem), 'AGENT-查询物流信息' (Agent - Query Logistics Information), '百炼知识库检索' (Bailian Knowledge Base Search), and '退换货流程结束' (End of Return/Refund Process). Red boxes highlight specific nodes: '模型选择' with a note about AI gateway configuration, '提取订单ID' with a note about tool realization, and '百炼知识库检索' with a note about MCP Function Calling.

- orders: Web 服务 已部署
- express: Web 服务 已部署
- competent-gould: 流程编排 待部署
- distracted-torvalds: 模型服务 待部署
- frosty-lalande: 模型服务 待部署
- focused-fermi: 流程编排 待部署

痛点

模型

- LLM 能力偏弱，不具备plan的能力
- 重点在优化提示词
- 依赖 Token数，无法继续提升体验

数据

- 数据之间的联系描述复杂，只能靠流程代码保证
- 数据源很多，如何收敛

流程

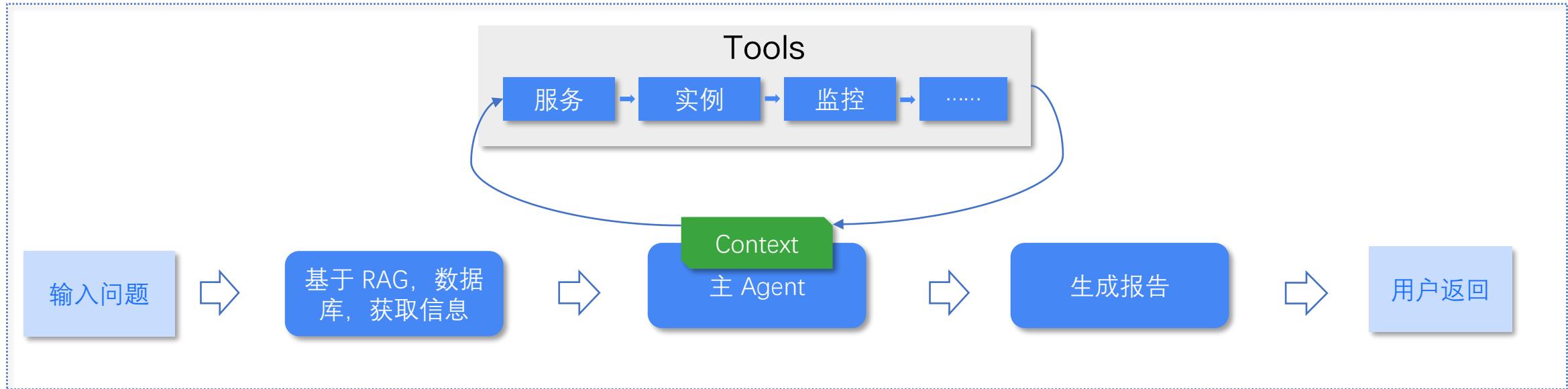
- 流程固定，无法适应后查询诉求
- 无法回溯历史，无法自行进行进化

能覆盖35% 左右的场景，远不够完美

02

阶段二：单 Agent 架构

架构方案



主 Agent

Agent 接收指令后，进行思考推理，并生成步骤：To Do List, 每一个步骤，可以使用已经安装的工具（MCP）调用，以及 sandbox::浏览器等能力去获取所需的信息，当执行完所有步骤，综合评估分析生成指令要求的报告。

工具 Tool

初期使用 FunctionCall 的能力去扩展实时查询的需求，逐步演变到 MCP 工具，通过统一的 MCP 协议，对接了服务，实例，监控，日志等等 MCP 工具，MCP 工具替代了“search”片段，提供了更完整的信息给大模型做参考。

Sandbox

沙箱也是非常好用的工具，特别是在文章中有链接但没有具体信息的场景下，通过浏览器沙箱获取到网页内容，并且可以并发执行，相互不影响，用完后删除。但是要依赖强大的弹性算力基建。



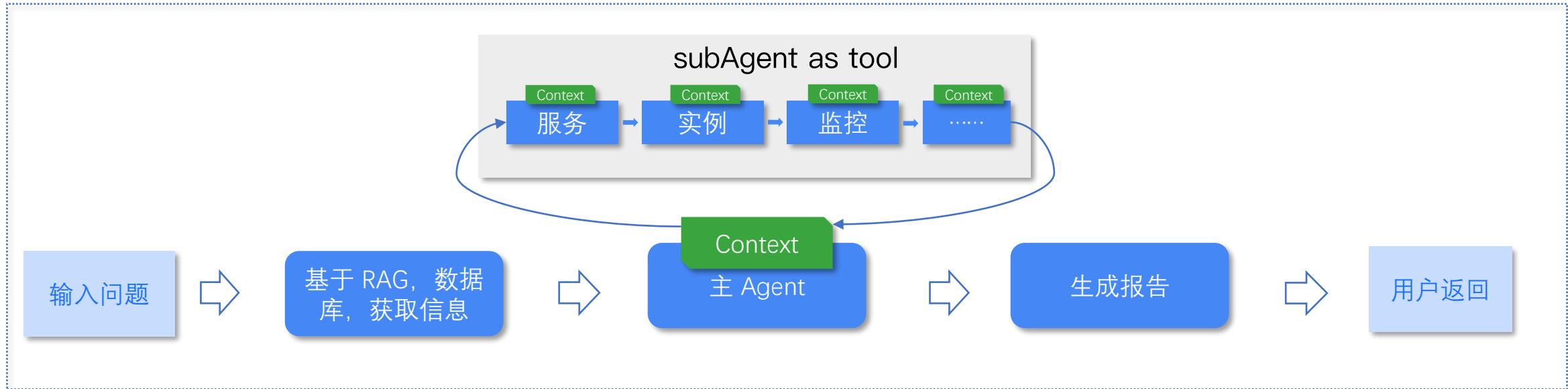
痛点

能覆盖50% 左右的场景，上下文管理是
最大瓶颈

03

阶段三：Multi-Agents尝试

初始方案



主 Agent

第一个Agent最好使用具备思考推理能力的大模型，接收指令后，先进行思考推理，并生成 Plan，根据当前注册的SubAgent进行分配，并依次发送指令给subAgent进行执行返回，主Agent结合返回信息继续推进完成整个任务。

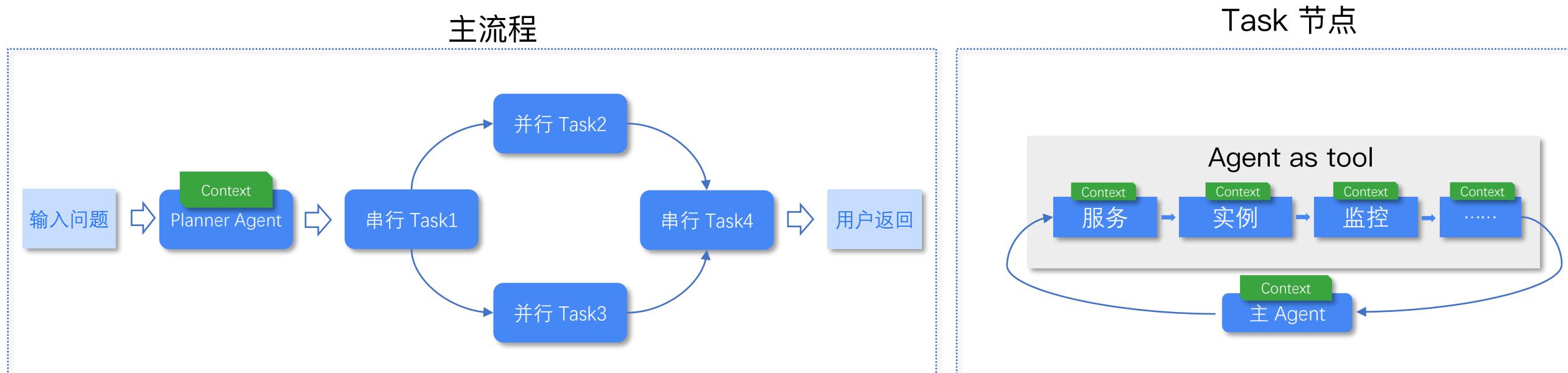
工具 Tool

每一个SubAgent有自己的大脑（LLM），但不都是相同的LLM，根据需求选择符合的LLM，防止大材小用浪费成本，SubAgent与主Agent是贡献上下文的，以此减轻传递负担。

挑战

流程设计比较理想，但是落地存在很大不足，整个流程不可控，常常出现时间超时，或者步骤丢失，导致整体的效果相比静态的工作流效果还差，如何解决Workflow的确定性执行，是我们面临最大的挑战。

迭代方案



Planner Agent

为了解决执行过程中的不确定性，采用有向 DAG 工作流的方式去确定性的执行 Planner 拆分后的 Action，每一个action 由工作流的机制保证容错，确保每一个步骤有一个结果，并且可以某一些工作流可以并行，从而加速 Agent 执行的时间。

Task 节点

基于DAG的每一个节点会独立完成一个任务，这个任务是由 Muti-Agents 完成，采用 Sub-Agents 的模式，Agent 之间共享上下文信息，从而完成独立的一项任务，任务完成后，会输出给 DAG 工作流平台中，由 DAG 决定是否继续。

优势

通过工作流机制解决了时间不可控的问题，通过 Muti-Agents 模式解决了上下文过大问题，整体方案借助了静态工作流与 A2A 的方案优势，从而完成目标工作。



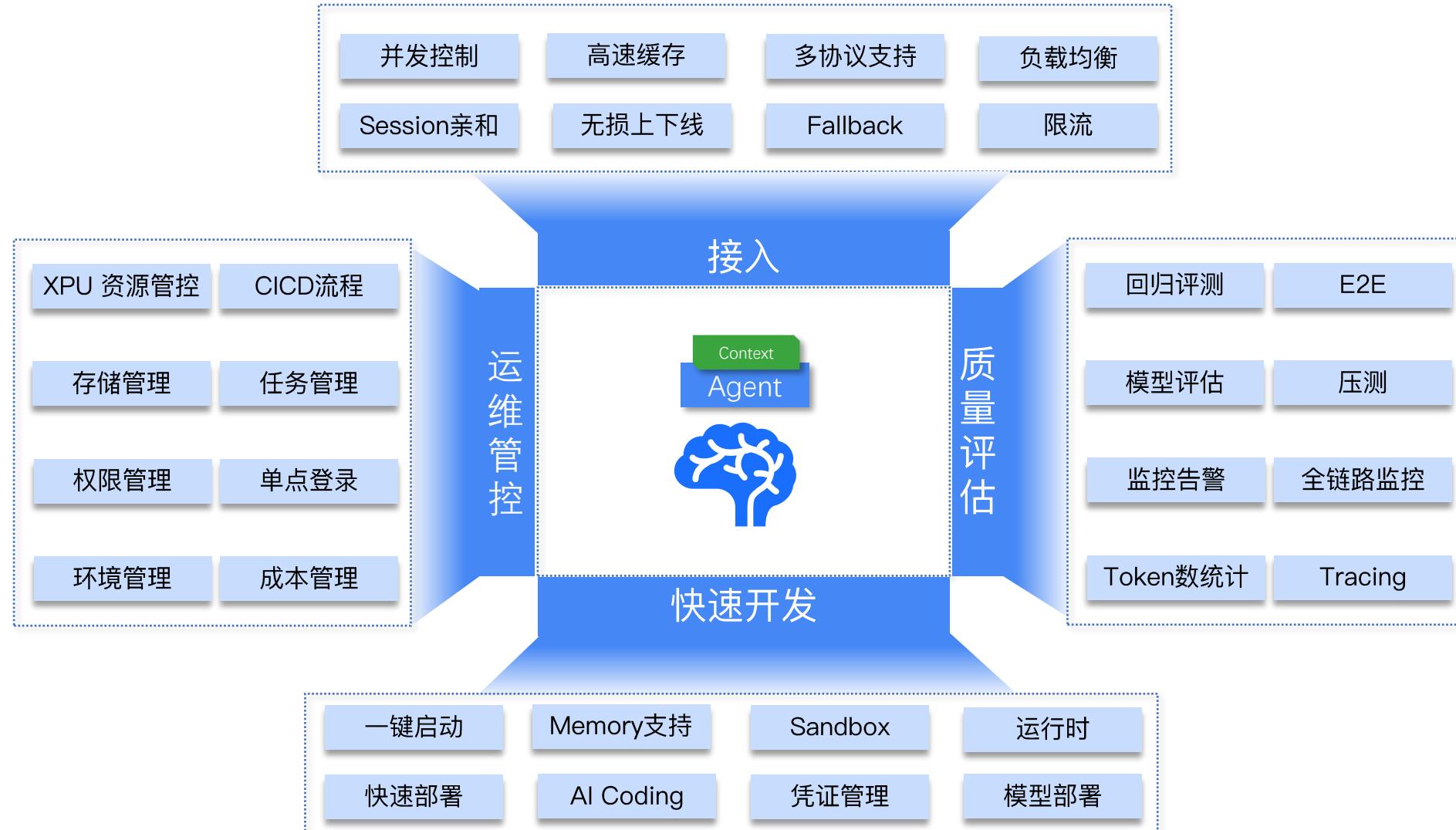
痛点

能覆盖90% 左右的场景，准确率仍需要提升

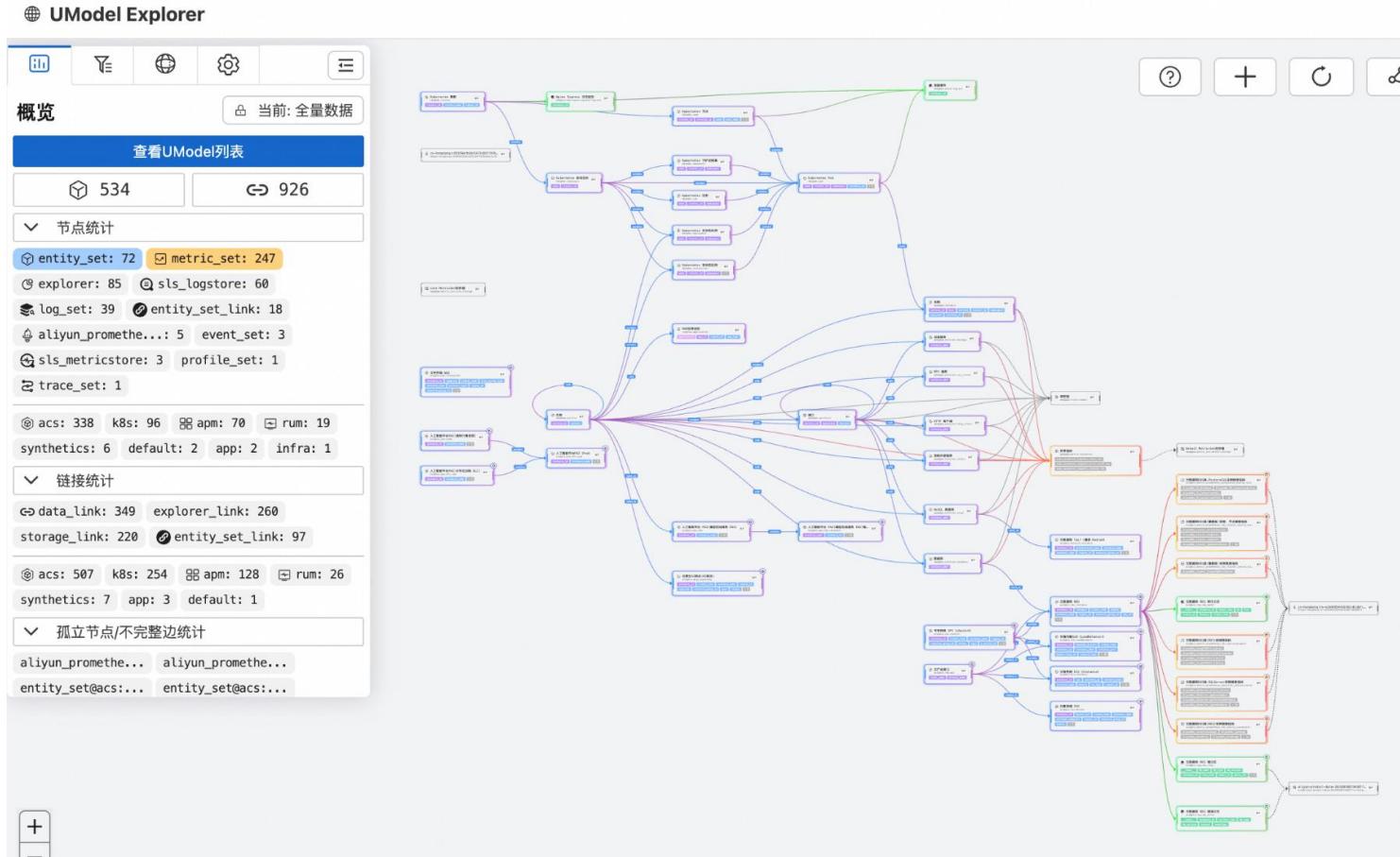
02

项目复盘总结

强大的基础设施支撑高效开发



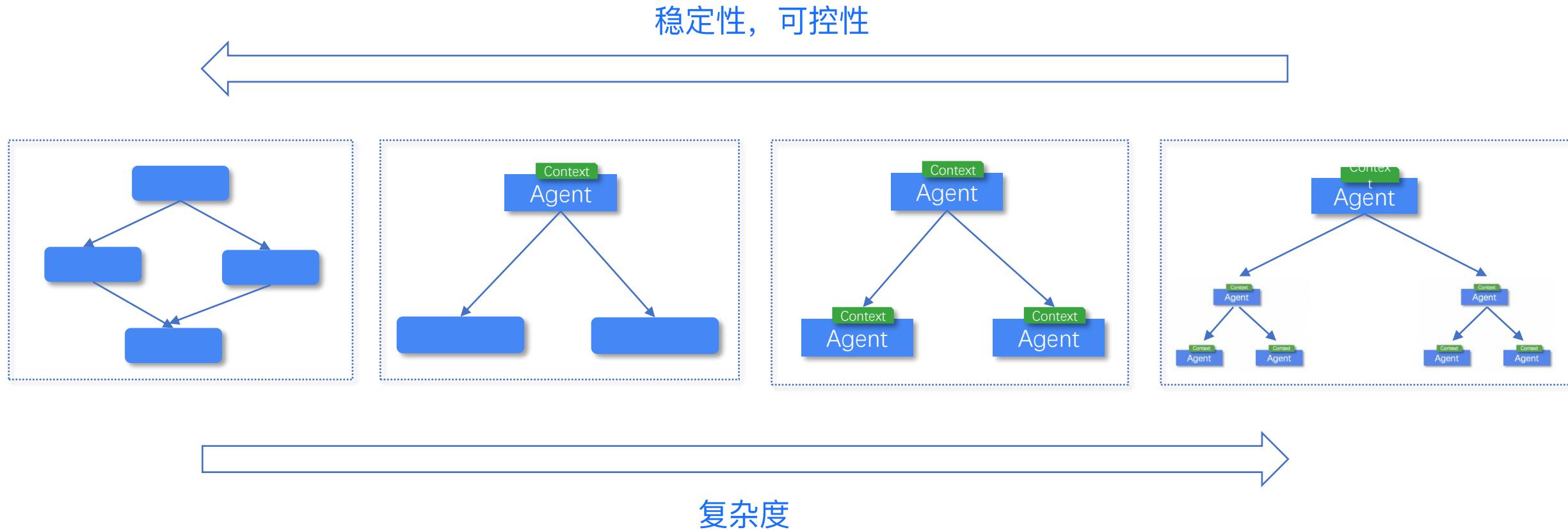
Data For AI



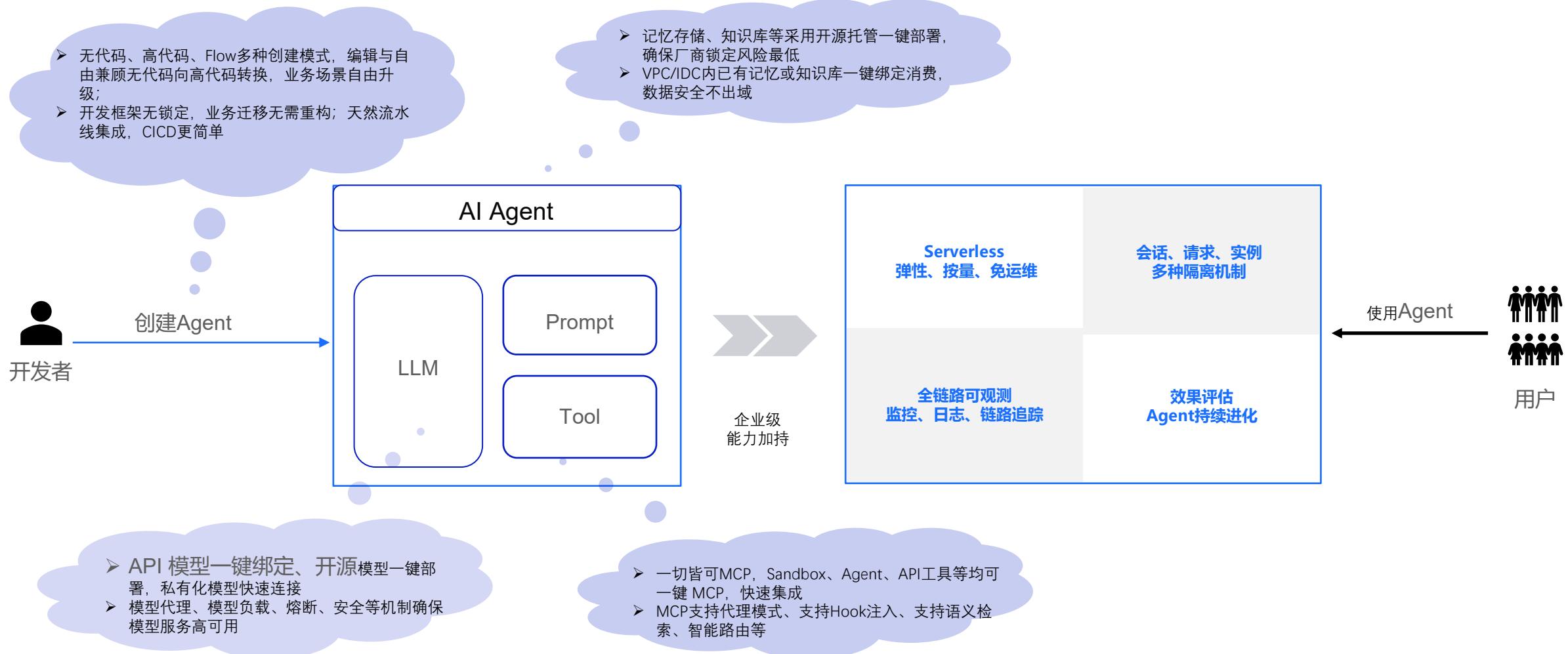
如何让模型理解你的数据

- 数据源管理
- 数据关系的表达

架构选择平衡



快速开发



■ 上下文工程

长任务如何保证执行完整性

拆分的任务数过多，会步骤遗忘，幻觉放大

- DAG 有向工作流
- 记录中间过程
- 上下文共享

如何使用 Rag && Memory

Agent 的核心竞争力：决定性于：Rag & 记忆。

- 向量数据库是否足够？
- Rag 的演进方向？
- 记忆要记录什么？何时读取？何时写入？

工具管理

工具过多导致效果下降

并不是工具越多越好，最重要的是选出适合的工作

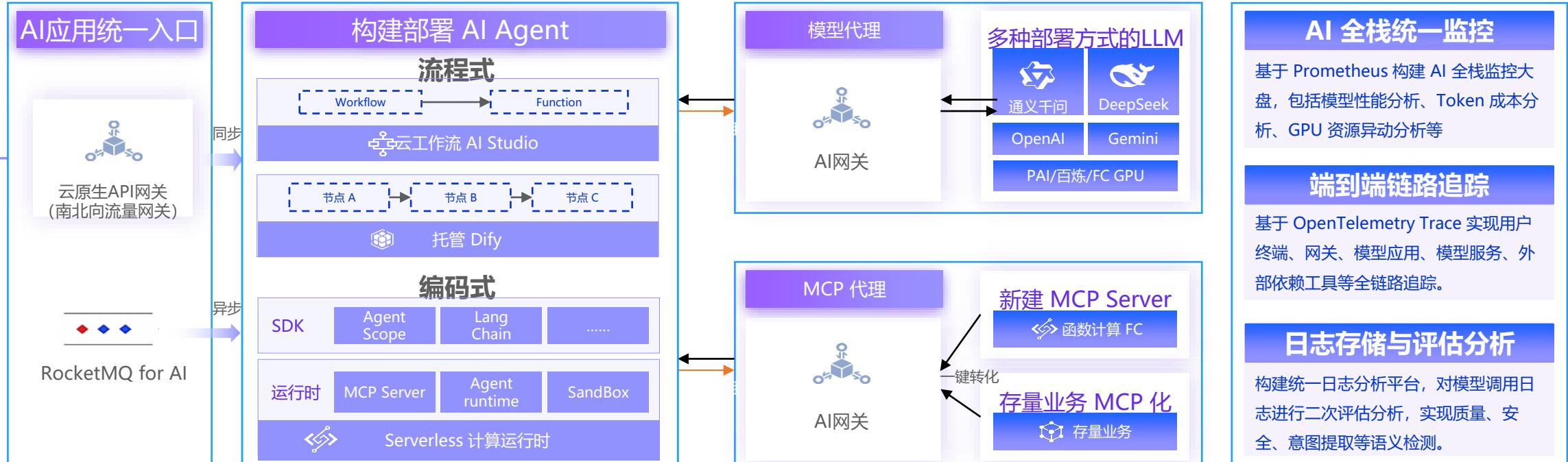
- 如何选出最适合的工具
- 工具的参数的正确性如何保证
- 工具的参数如何界定

工具如何做到复用

我们会使用很多的工具，各自为战？如何最大化的复用工具。

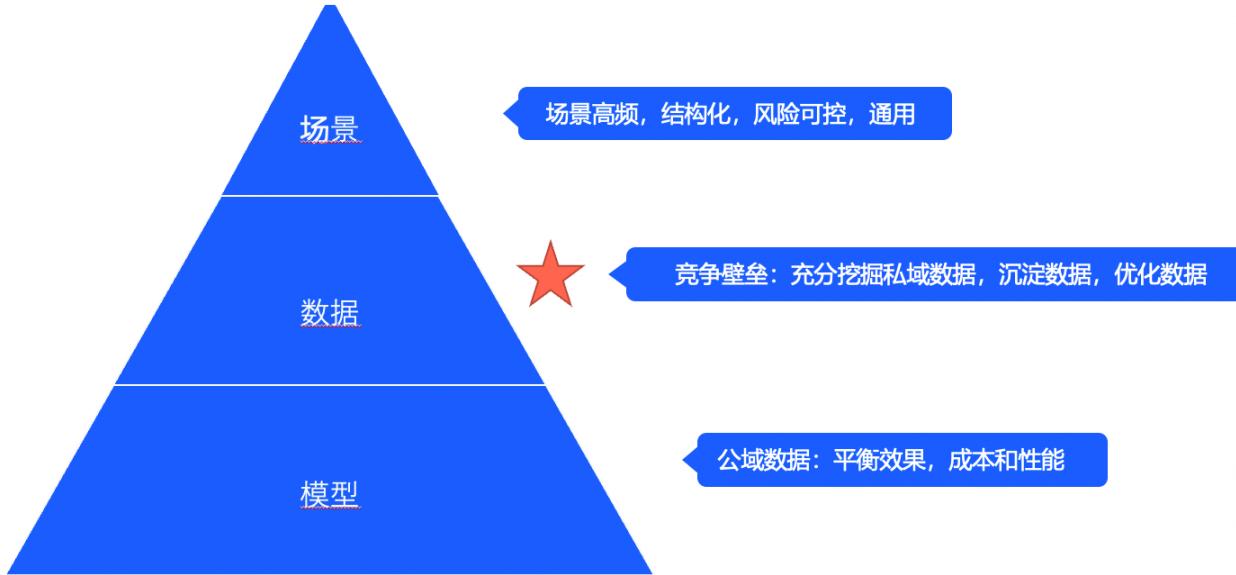
- 工具市场的重要性
- 如果保证工具的稳定性

全链路可观测



数据驱动的迭代

找到核心提效场景，构建高质量数据壁垒，借助大模型大势快速迭代

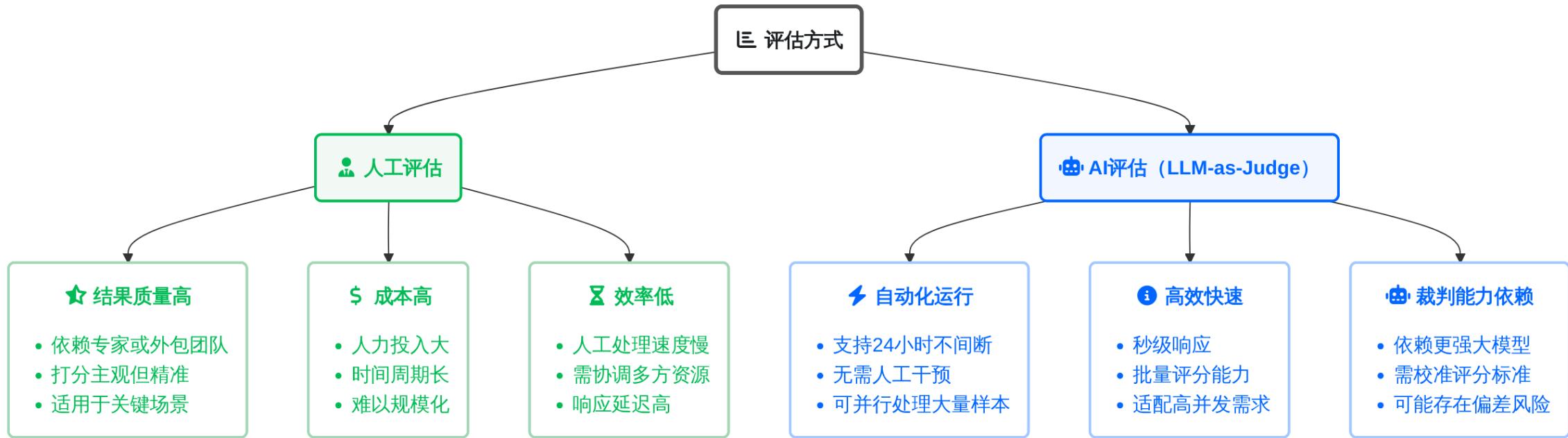


AI 原生应用 数据飞轮

充分挖掘私域数据：客户数据可沉淀，行业数据可演进，评估数据可量化，反馈数据可持续



模型评估的重要性



02

未来规划



数据驱动 Agentic AI 全链路研发闭环



极客邦科技 2026 年会议规划

促进软件开发及相关领域知识与创新的传播



参会咨询



查看会议

北京

1200人

QCon

全球软件开发大会

会议时间：4月16-18日

- Agentic Engineering
- AgentOps
- 下一代模型架构与推理优化
- AI 原生基础设施
- 知识工程实践
- AI 安全

深圳

1000人

AiCon

全球人工智能开发与应用大会

会议时间：8月21-22日

- Agentic AI
- 轻量化与高效推理
- 多模态应用
- AI + IoT 场景实践
- AI 工业化落地

北京

1000人

AiCon

全球人工智能开发与应用大会

会议时间：12月18-19日

- 大模型架构创新
- 多模态 AI 产业融合
- 具身智能
- AI for Science
- 大模型安全

4月

6月

8月

10月

12月

AiCon

全球人工智能开发与应用大会

会议时间：6月26-27日

- AI Infra 系统工程
- 多 Agent 协作与实践
- 多模态融合
- 模型训练与推理创新
- 数据平台与特征服务

上海

1000人

QCon

全球软件开发大会

会议时间：10月22-24日

- AI Agent
- Vibe Coding
- 智能可观测
- 推理基建
- 模型攻防
- AI x 创造力

上海

1200人

THANKS

探索 AI 应用边界

Explore the limits of AI applications

AiCon

全球人工智能开发与应用大会