

51CTO WOT

World Of Tech 2024

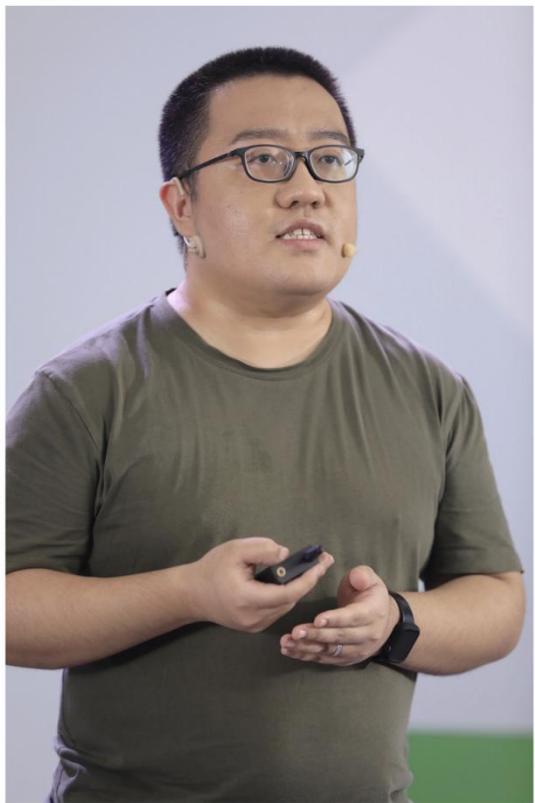
WOT全球技术 创新大会

智启新纪
慧创万物



去哪儿旅行机票售前客服 Agent&RAG实践

李佳奇
去哪儿旅行 技术总监



李佳奇 去哪儿旅行/机票目的地事业群 技术总监

- 10年机票业务研发工作经验
- 技术中心TC委员/机票TC主席/机票AIGC落地负责人
- 在业务架构/DDD落地/高并发高可用系统设计/技术团队建设等方面有丰富的经验
- 多次参与各种公司内外技术分享和行业大会分享

目录

- 背景介绍
- 技术方案
- 成果总结

背景介绍

团队迷茫

技术储备是否足够
是否会被AI代替
该如何顺应趋势



老板焦虑

是否会出现降维打击
如何继续保持领先
如何尽快拿出产品

用户期待

杀手级应用何时出现
体验极大提升
智能效果提升

AIGC 时代来临我们该如何应对

我们的应对方法



打基建

开发基建、测试基建、运行基建

找机会

业务*AI矩阵

建团队

AI周会/培训/项目练兵

目前AIGC在去哪儿机票的落地现状

技术基建

01

Langchain4J Qunar框架

RAG

prompt/智能体平台

团队建设

02

覆盖团队20%成员的AIGC技术储备

AI周会/AI圆桌会

技术分享/技术交流

项目落地

03

业务侧：用户体验提升/收益提升/智能营销/归因分析

技术侧：工单自动化/SQL自动生成/checklist智能生成

售前客服项目介绍

项目背景

售前客服覆盖缺失

基于用户主流程的导购可行性验证

用户问卷调研需求

技术准备

RAG/Embedding基建准备

Agent设计落地技术准备

技术人员培养

达成路径

mvp版本入口小流量验证

封闭式问题/开放式问题

开源大模型/商用大模型



技术方案

【开放式回答】初始版

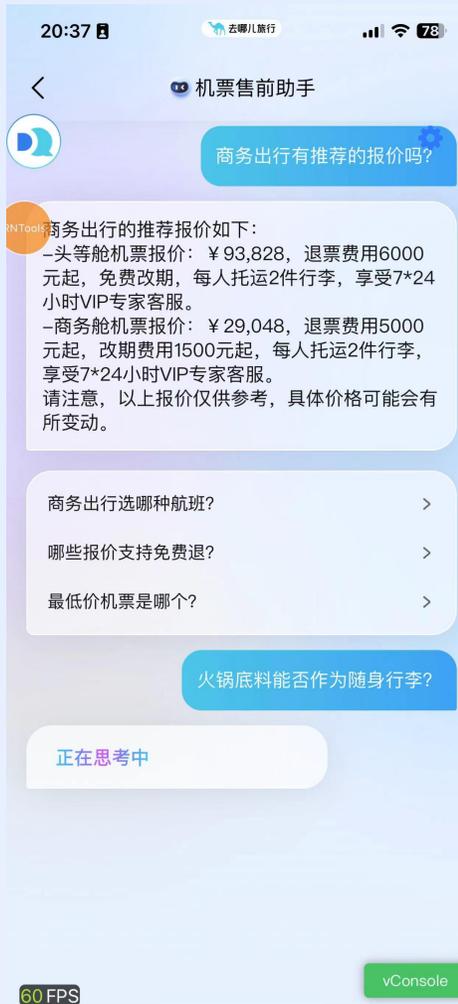
静态/动态问题回答



【封闭式回答】文字版

问题分类

问题推荐



报价推荐

图文消息

【封闭式回答】图文版





Models

GPT3.5/4/4v DALL-E2/3

Gemini

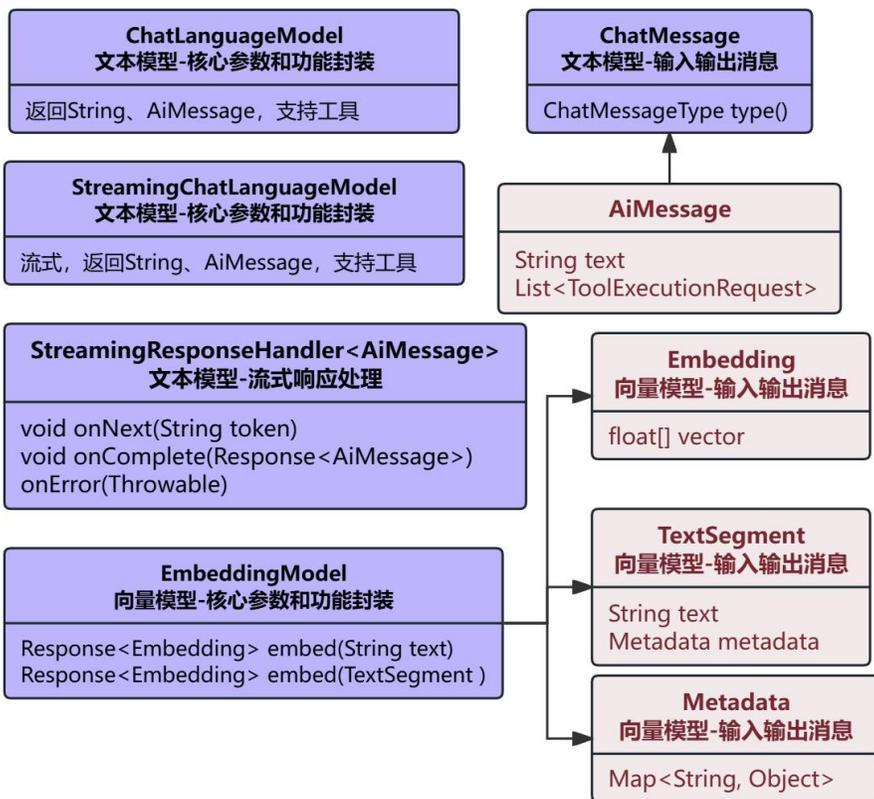
文心一言

通义千问

.....

Model本质封装了和LLM供应商的HTTP API, 文本、向量、图片等不同类型的API不同, 封装用到的对象也不同。

langchain4j-core
LLM Model核心组件设计介绍



售前助手-代码示例

```

/**
 * 开源版本 <br/>
 * 智谱4 <br/>
 * 模型: glm-4<br/>
 */
@Bean
public ZhipuAiStreamingChatModel zhipuAiChatModel4(@Value("${glm4.api.key}") String apiKey) {
    return ZhipuAiStreamingChatModel.builder().model("glm-4").logRequests(true).logResponses(true).apiKey(apiKey)
        .build();
}
    
```

```

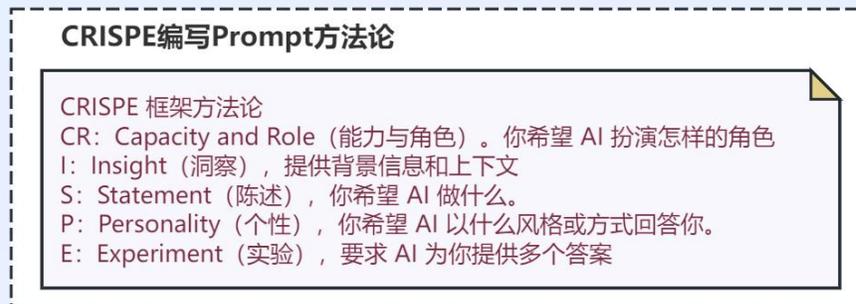
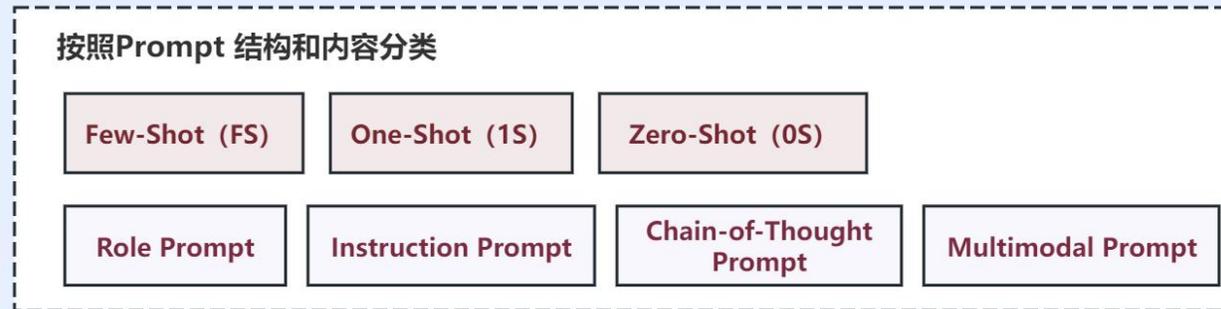
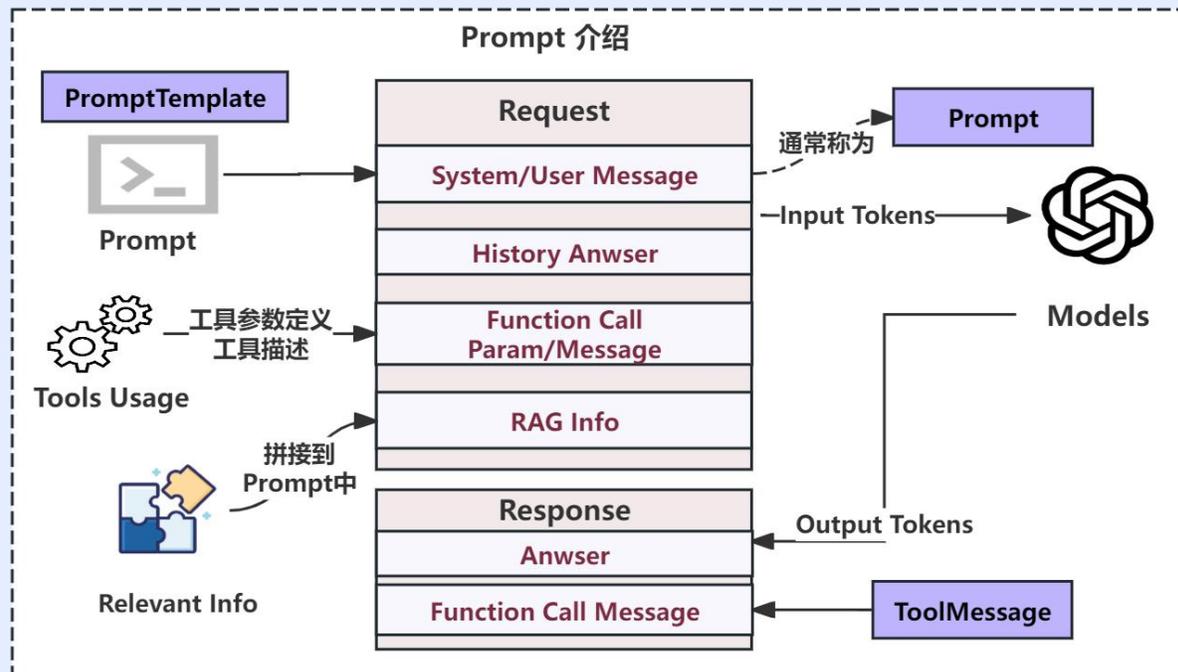
/**
 * 去哪儿算法组版本 <br/>
 * 算法组封装的llmModel, 远程Http调用算法组封装的sdk接口, 算法组实现了llmModel的统一账号管理、风控拦截、llm供应商对接等。<br/>
 * 模型: gpt-4-turbo <br/>
 * 问题推荐Agent使用, 输入和输出比较少, 和回答Agent并行执行
 */
@Bean
public QunarAiChatModel qunarAiChatModelGpt4(@Value("${gpt.api.key}") String key,
    @Value("${gpt.api.password}") String password) {
    DefaultQunarAiClient client = DefaultQunarAiClient.builder(key, password).logLevel(ClientLogLevel.ALL).build();
    ClientRequiredConfig requiredConfig =
        ClientRequiredConfig.builder().project("presale_ai_assist").password(password)
            .userIdentityInfo("fanmao.meng").build();
    ClientOptionalConfig optionalConfig = new ClientOptionalConfig();
    optionalConfig.setTopP(0);
    optionalConfig.setChatApiType(ChatComplexRequest.ApiType.GPT_4_TURBO);
    optionalConfig.setChatApiVersion(ChatComplexRequest.ApiVersion.V2023_12_01_PREVIEW);
    return QunarAiChatModel.builder().client(client).requiredConfig(requiredConfig).optionalConfig(optionalConfig)
        .build();
}
    
```

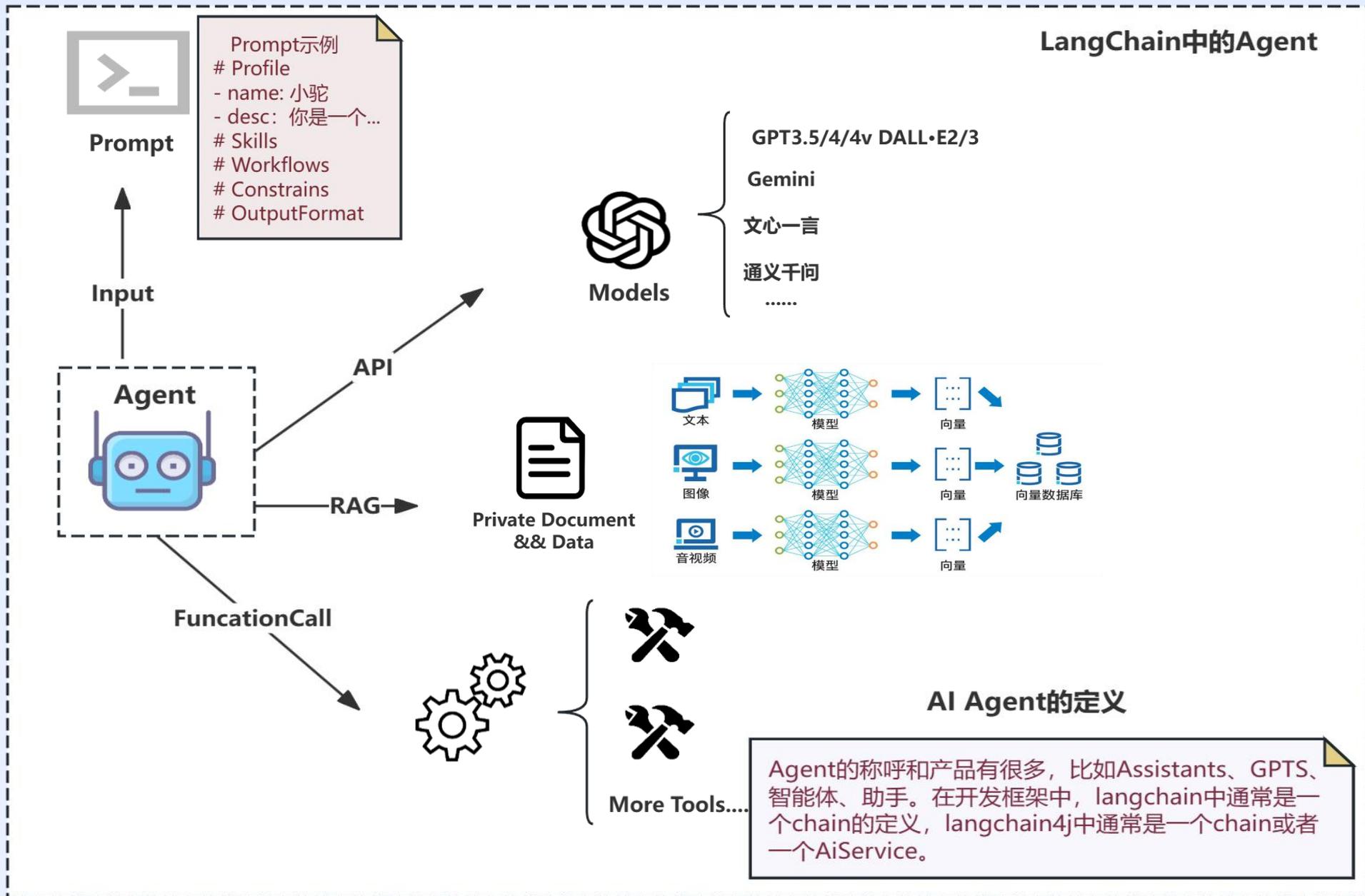


Prompt

Prompt示例

```
# Profile
- name: 小驼
- desc: 你是一个...
# Skills
# Workflows
# Constrains
# OutputFormat
```

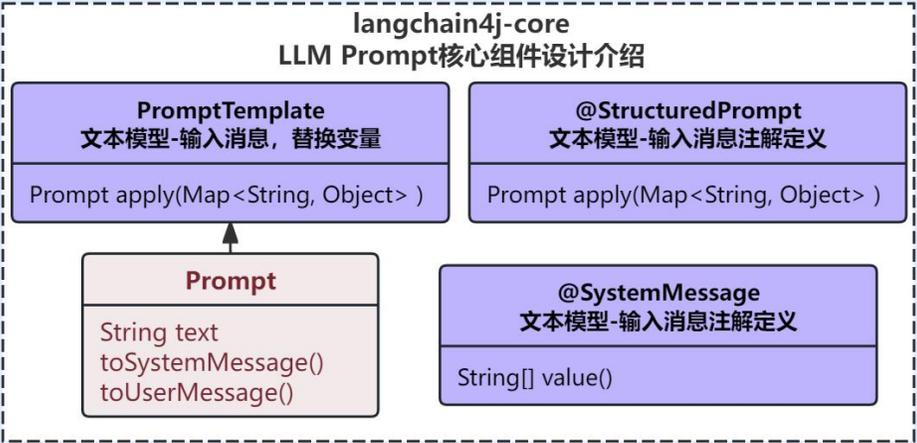






Prompt

```
Prompt示例
# Profile
- name: 小驼
- desc: 你是一个...
# Skills
# Workflows
# Constrains
# OutputFormat
```



售前助手-代码示例

Prompt定义

```
no usages new *
@SystemMessage("""
# Profile
你是忠实的去哪儿网机票用户，你会收到一个问题，你的工作不是回答收到的问题，而是要思考，作为一个机票用户，问了这个问题之后，你还会问哪些问题
# Skills
## 接下来要问的问题
1. 你能够准确的理解给出的问题
2. 你不会直接回答这个问题，而是从给出的**参考问题列表**中选择3个你接下来要问的问题
3. 将你接下来要问的问题整理成列表
# WorkFlows
1. 第一步：仔细理解给出的问题
2. 第二步：接下来要问的问题
3. 第三步：从第二步得到的问题列表中，挑选最多3个你接下来最想了解的问题，并给出
4. 第四步：仔细思考，上一步得到的问题都是你接下来最想了解的问题吗
5. 第五步：请注意：问题列表中的所有问题必须包含在**参考问题列表**中，如果有不存在的问题，请去掉该问题
# Initializ
**参考问题列表**： {{frequentlyAskQuestion}}
# Constrains
- 你不能直接回答给出的问题
- 请严格参照回答格式进行回答，不要添加任何无关的信息
- 最终得出的问题列表允许为空集合，但不允许包含与原始问题相同的问题
- 最终得到的问题列表必须包含在**参考问题列表**中
# OutputFormat
你的回答必须是Json格式的问题列表，参考结构如下：
{"questions":<问题列表>
""")
Response<AiMessage> questionGenerateX(@UserMessage String userMessage, @V("frequentlyAskQuestion") String frequentlyAsk,
@UserName String username);
```

Prompt变量替换、转为请求LLM参数

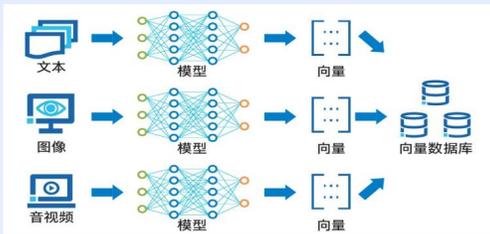
```
protected Optional<dev.langchain4j.data.message.SystemMessage> prepareSystemMessage(Method method, Object[] args) {
    Parameter[] parameters = method.getParameters();
    Map<String, Object> variables = getPromptTemplateVariables(args, parameters);

    dev.langchain4j.service.SystemMessage annotation = method.getAnnotation(dev.langchain4j.service.SystemMessage.class);
    if (annotation != null) {
        String systemMessageTemplate = getPromptText(
            method,
            type: "System",
            annotation.fromResource(),
            annotation.value(),
            annotation.delimiter()
        );

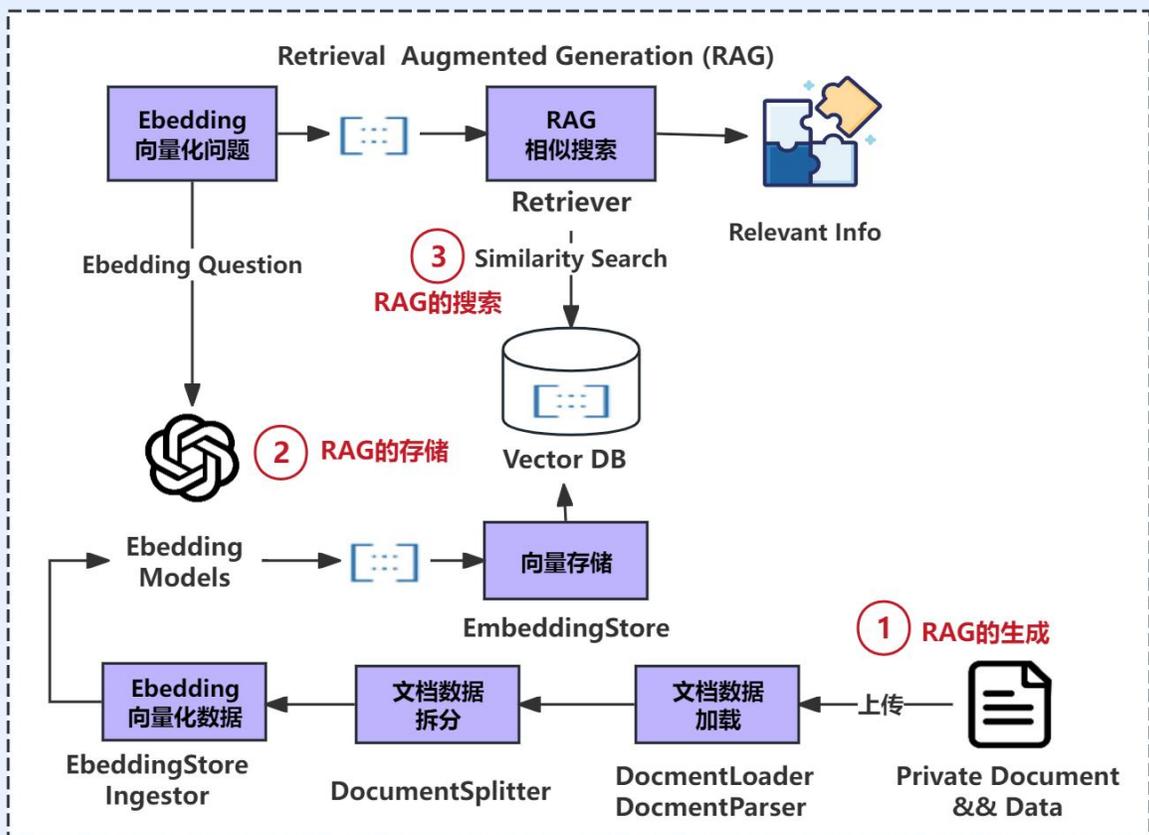
        Prompt prompt = PromptTemplate.from(systemMessageTemplate).apply(variables);
        return Optional.of(prompt.toSystemMessage());
    }

    return Optional.empty();
}
```

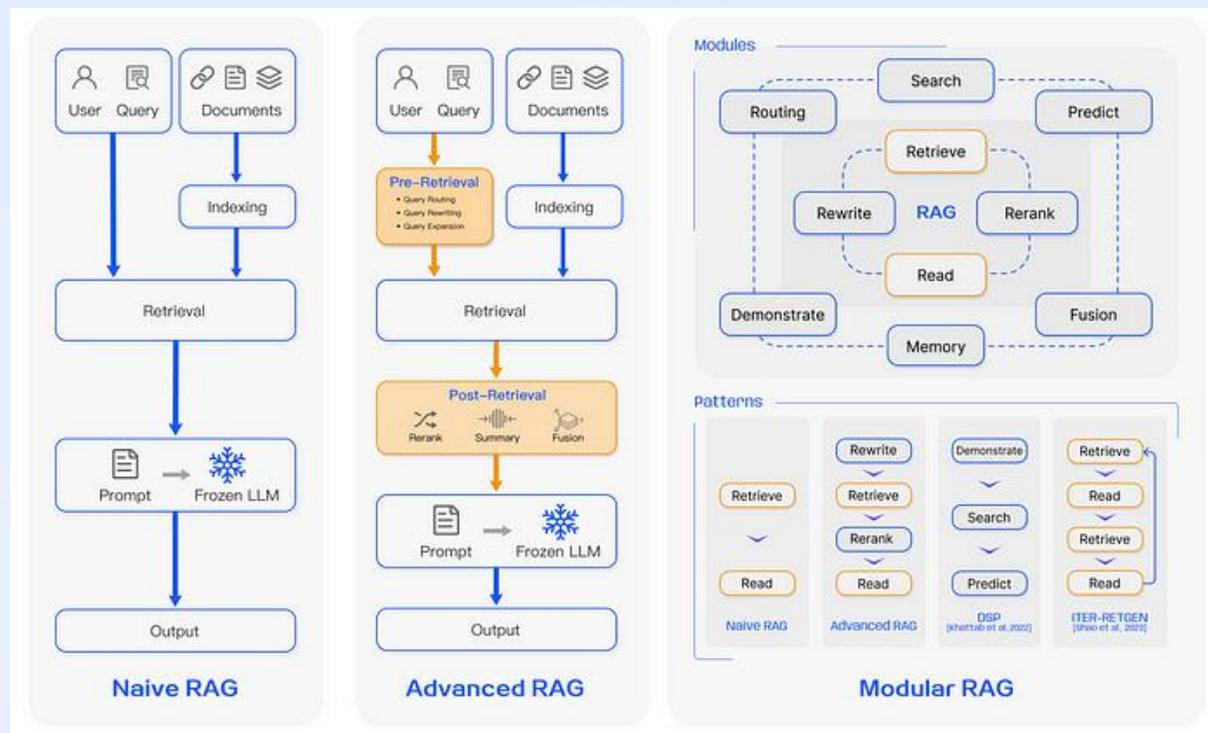
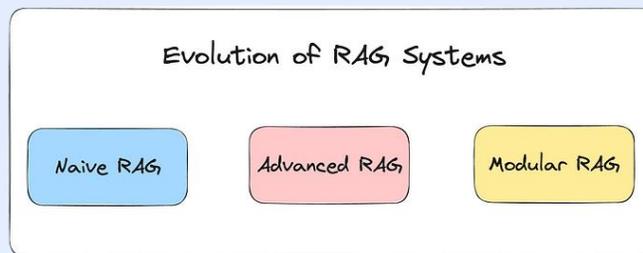
Private Document & Data

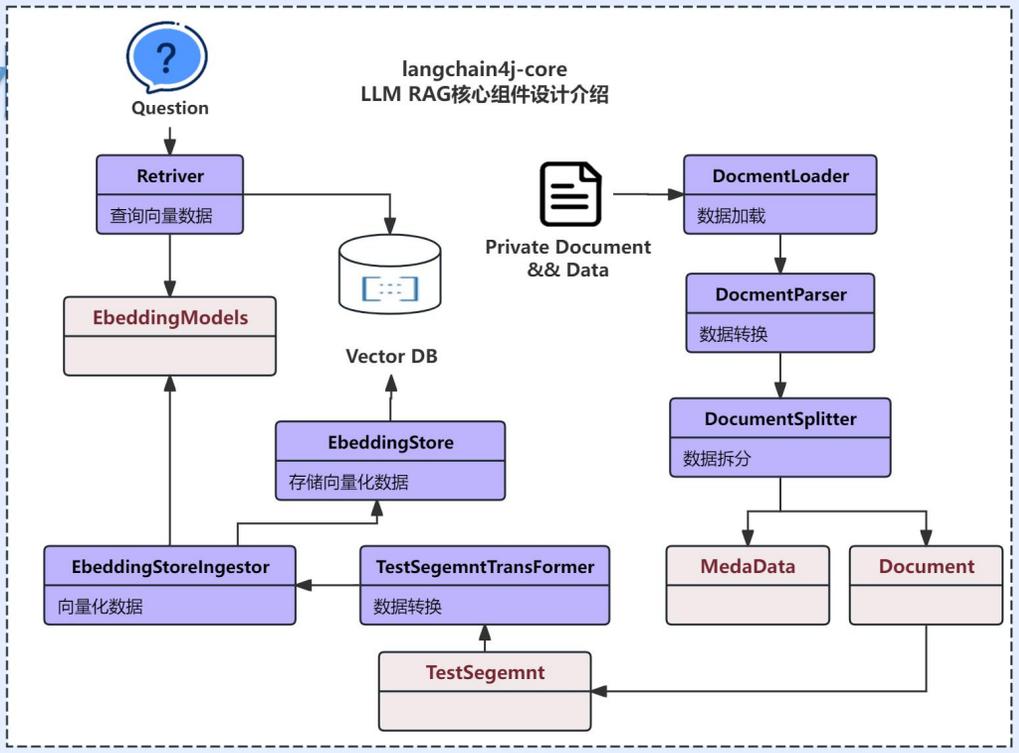


RAG生成和查询流程介绍



三种 RAG 范式介绍





```

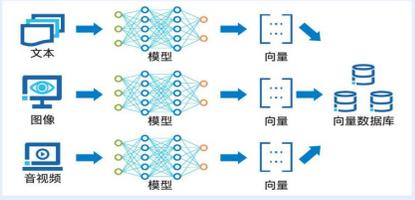
@Slf4j
public class QAPairTransform implements TextSegmentTransformer {

    no usages
    public static final QAPairTransform INSTANCE = new QAPairTransform();

    @Override
    public TextSegment transform(TextSegment segment) {

        String text = segment.text();
        String[] qaPais = text.split(regex: "##");
        if (qaPais.length != 2) {
            log.info("解析问题对失败, text:{}", text);
            return segment;
        } else {
            return TextSegment.from(qaPais[0], segment.metadata());
        }
    }
}
    
```

Private Document && Data



售前助手-代码示例

```

防晒隔离水能否带上飞机? ##meta:{"category": "规则类, 化妆品"}####
折叠式婴儿车能否进客舱? ##meta:{"category": "规则类, 儿童物品"}####
能否带十几条烟上飞机? ##meta:{"category": "规则类, 烟酒"}####
    
```

```

public class QuestionDocumentSplitter extends DocumentByHashTagSplitter {
    protected TextSegment createSegment(String text, Document document, int index) {
        String[] elements = text.split(regex: "##meta:");
        String questionStr = elements[0];
        Map<String, String> metadataMap = Maps.newHashMap();
        metadataMap.putAll(document.metadata().asMap());
        return TextSegment.from(questionStr, Metadata.from(metadataMap));
    }
}
    
```

```

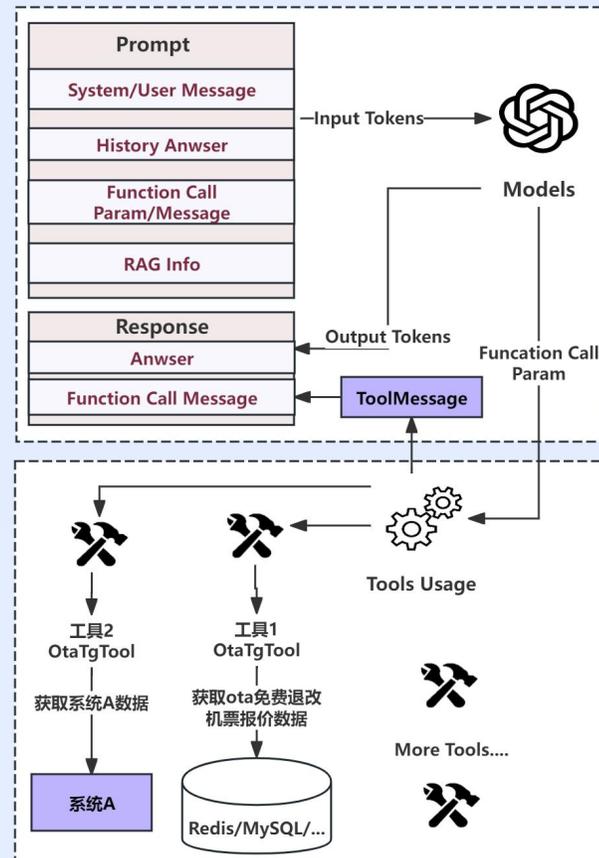
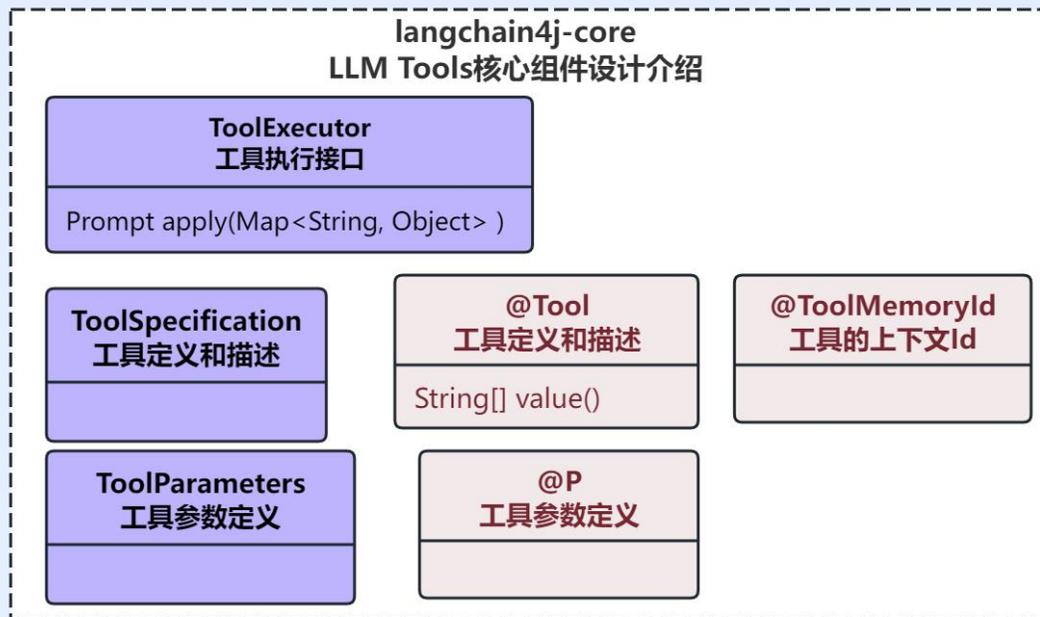
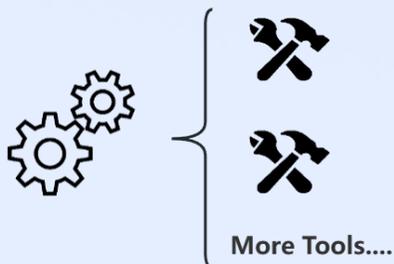
/**
 * 去哪儿封装的EmbeddingModel, 远程Http调用算法组接口生成向量, 算法组底层使用自有的模型
 */
@Bean
public EmbeddingModel embeddingModelQunar() {

    String key = modelAccountConfigProcess.getKey();
    String password = modelAccountConfigProcess.getPassword();
    String userIdentityInfo = modelAccountConfigProcess.getUserIdentityInfo();

    DefaultQunarAiClient client = DefaultQunarAiClient
        .builder(key, password)
        .build();

    ClientRequiredConfig requiredConfig = ClientRequiredConfig.builder()
        .project("presale_ai_assist")
        .password(password)
        .userIdentityInfo(userIdentityInfo).build();

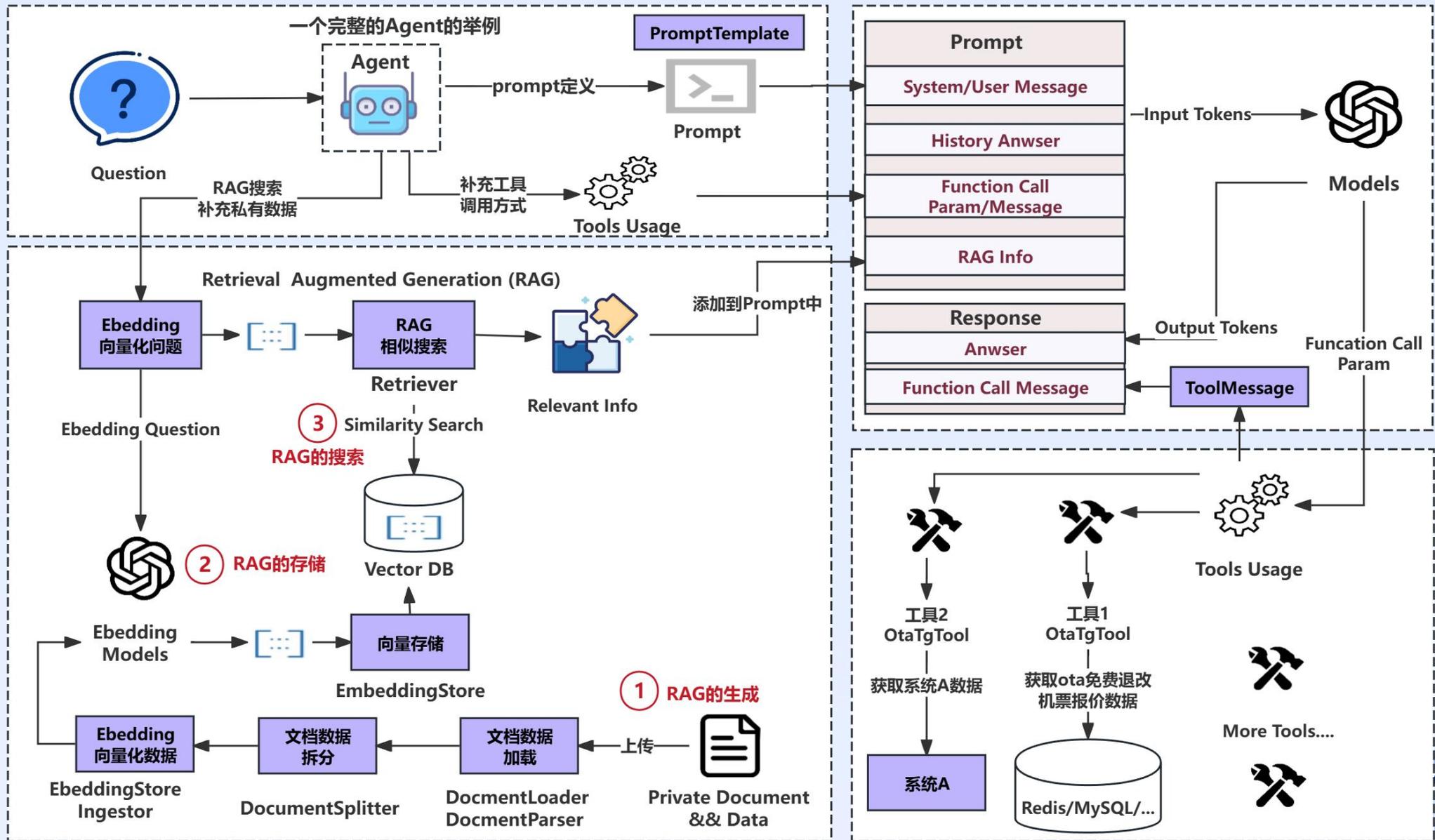
    ClientOptionalConfig optionalConfig = new ClientOptionalConfig();
    optionalConfig.setEmbeddingApiType(modelAccountConfigProcess.getEmbeddingApiType());
    return QunarEmbeddingModel.builder().client(client).requiredConfig(requiredConfig).optionalConfig(optionalConfig).build();
}
    
```



售前助手-代码示例

```

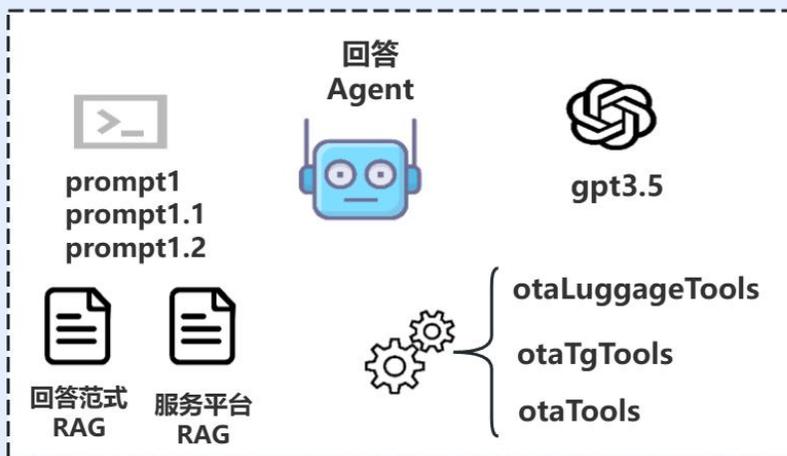
@Tool(name = "getBusinessCabinProduct", value = "可以通过这个工具查询商务舱、公务舱、头等舱的机票产品信息，返回的报价都是商务舱、公务舱、头等舱")
private String getBusinessCabinProduct(@P("ota报价缓存Key") String otaCacheKey) {
    log.info("商务舱、公务舱、头等舱报价查询工具, key:{}", otaCacheKey);
    String result = interRedisService.getKey(otaCacheKey);
    if (StringUtils.isNotEmpty(result)) {
        OtaData otaData = JsonUtil.fromJson(result, OtaData.class);
        return getBusinessCabinProduct(otaData);
    }
    return "查询结果为空，请引导用户购买别的机票产品";
}
    
```



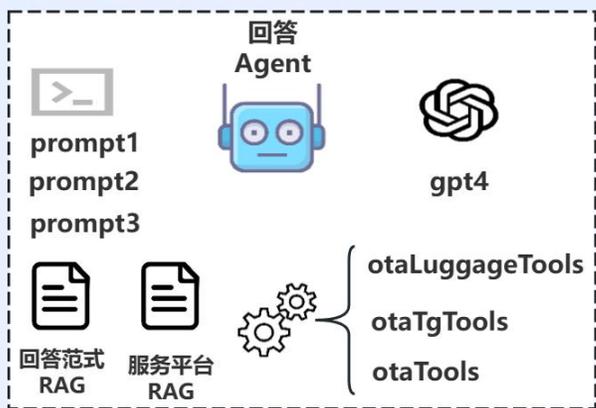
V1.0版 单Agent协作 一个上下文



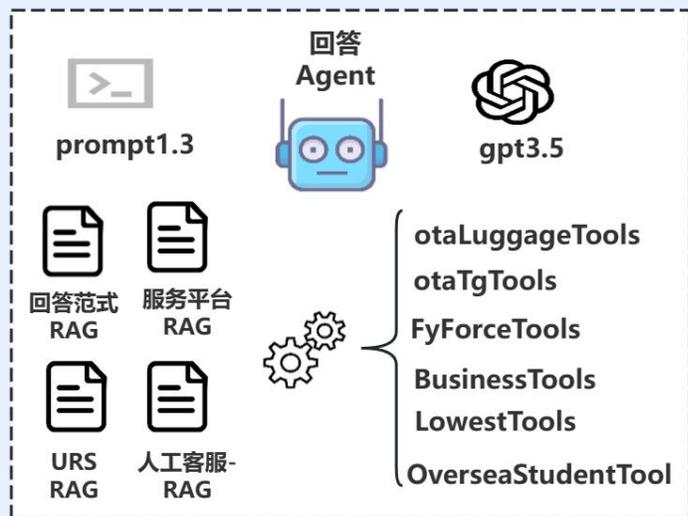
V1.1版 单Agent协作 扩展tools/优化prompt1/扩展RAG



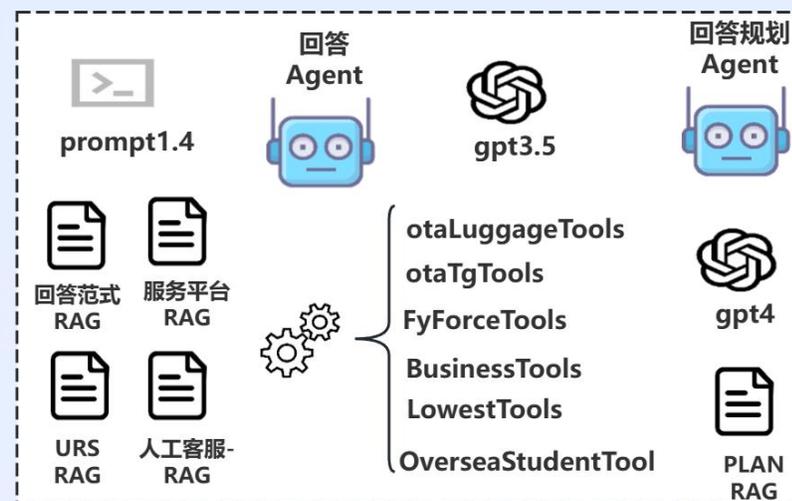
V1.2版 单Agent协作
Prompt合并+更换GPT4模型



V1.3版 单Agent协作
简化Prompt+增加知识覆盖+动静知识补充校准+Tools拆分和接入



V1.4版 单Agent协作
增加辅助自我反思Prompt+基于PLAN-RAG+隐式多轮



V1.0 版本问题

- 1) RAG数据少，答案过时不正确
- 2) Prompt长，gpt35不完全遵循规则
- 3) 开放式回答准确率低
 - 无意义、流水式回答
 - 报价引导弱，无相关报价位置

V1.1 尝试解决

- 1) RAG校准、补充回答范式RAG知识
- 2) 细化Tools 补充报价位置信息，拆分退改、行李等tools
- 3) 优化Prompt，配置化

过时不正确

如何购买儿童票

目前仅部分机票支持购买儿童票。如果您需要购买，请在搜索航班时勾选携带儿童选项，勾选后搜索出的航班信息是可以购买儿童票。如果您已经购买了成人票，请联系人工客服尝试购买。

请帮我介绍下这些报价有什么不同？

您所查的当前机票，都不包含免费退票，但第1个报价提供了更好的退改手续费。<去哪儿旗舰店>

我想买一个最低价的机票，有推荐的吗？

流水回答、引导弱

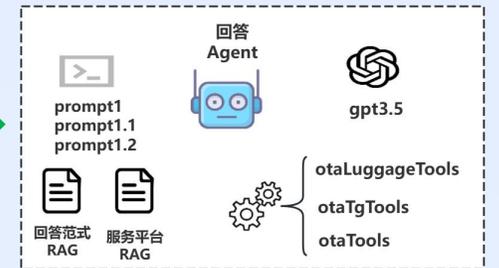
根据最新的机票报价数据，我为您找到了以下推荐的最低价机票：

- 1.经济舱特惠票：价格为¥520，每人托运1件x23KG行李，每人手提1件x11KG行李（三边之和不超过115cm）。退改详情请在购买成功后预计48小时内完成出票。
- 2.去哪儿优选经济舱：价格为¥520，每人托运1件x23KG行李，每人手提1件x11KG行李（三边之和不超过115cm）。退改详情请在购买成功后预计60分钟内完成出票。
- 3.商务舱：价格为¥400，每人托运1件x23KG行李。退改详情请享受尊享7*24小时VIP专家客服。
- 4.头等舱：价格为¥400，每人托运1件x23KG行李。退改详情请享受尊享7*24小时VIP专家客服。请注意，以上价格和退改签政策仅供参考，具体以实际购买页面显示为准。祝您购票顺利！

V1.0版 单Agent协作 一个上下文



V1.1版 单Agent协作 扩展tools/优化prompt1/扩展RAG



V1.2版 单Agent协作 Prompt合并+更换GPT4模型



V1.3版 单Agent协作 简化Prompt+增加知识覆盖+动静知识补充校准+Tools拆分和接入



V1.4版 单Agent协作 增加辅助自我反思Prompt+基于PLAN-RAG+链式多轮



V1.1 版本问题

- 1) 默认推荐问题, 小孩关键字命中风控
- 2) 回答不稳定

V1.2 尝试解决

- 1) Prompt问题同义词默认知识增加
- 2) 升级模型到GPT4

售前助手-问题示例

带儿童适合买哪个报价

小驼开小差了, 请稍后再试

关键字命中风控 带儿童适合买哪个报价

小驼开小差了, 请稍后再试

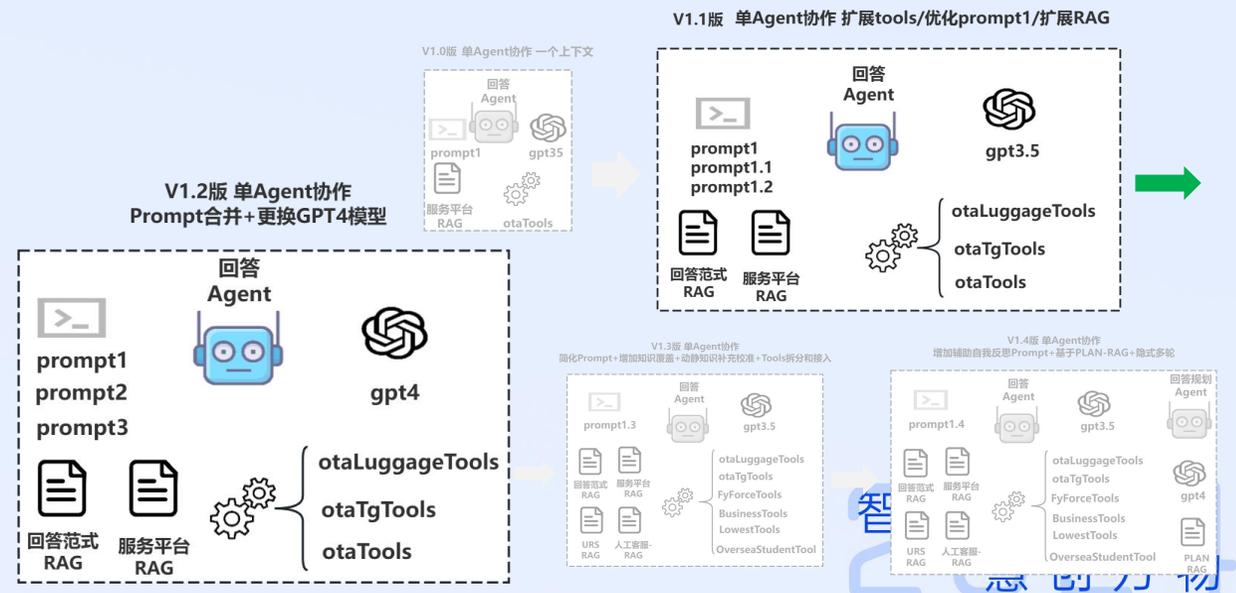
回答不稳定

可以免费退票吗?

根据航班报价规则, 您可以点击您选定的航班报价或在航班信息下方点击【退改签规则】, 以查看该报价的自愿退款/改签手续费规则。如果您已经订票, 可以在APP的【订单】频道找到对应机票订单, 点击进入订单详情页后点击【我要退票/改签】按钮, 在选择相应的退改原因后就可以完成提交操作。祝您顺利办理退票手续!

可以免费退票吗?

抱歉, 我无法获取到关于免费退票的信息。建议您查看您选定的航班报价, 或直接联系航空公司官网、APP、或客服, 以获取详细的退票规则和手续费信息。祝您顺利办理退票手续!



V1.2 版本问题

- 1) GPT4耗时慢 平均35s
- 2) RAG参考回答质量不高，回答不够用户

V1.3尝试解决

- 1) 降级回GPT 3.5，拆分Agent职责，多Agent并行
- 2) RAG拆分和补充，Q-RAG/ QA-RAG/URS-RAG/人工客服 RAG
- 3) RAG分类和校准，【静态答案】校准、【动态参考回答】优化
- 4) Tools继续拆分

售前助手-问题示例

退票保障是什么?

退票保障是一种增值服务，您可以在您需要退票时为您提供一定比例的退款手续费补偿。这项服务通常需要额外支付费用，但可以在您自愿退票时减轻经济损失。具体的补偿比例和规则可能会根据不同的服务包而有所不同。购买退票保障后，如果您的航班发生变动或取消，您可以按照服务条款获得相应的退款或补偿。

这几个报价要怎么选择?

选择机票报价时，您可以考虑以下几个因素：

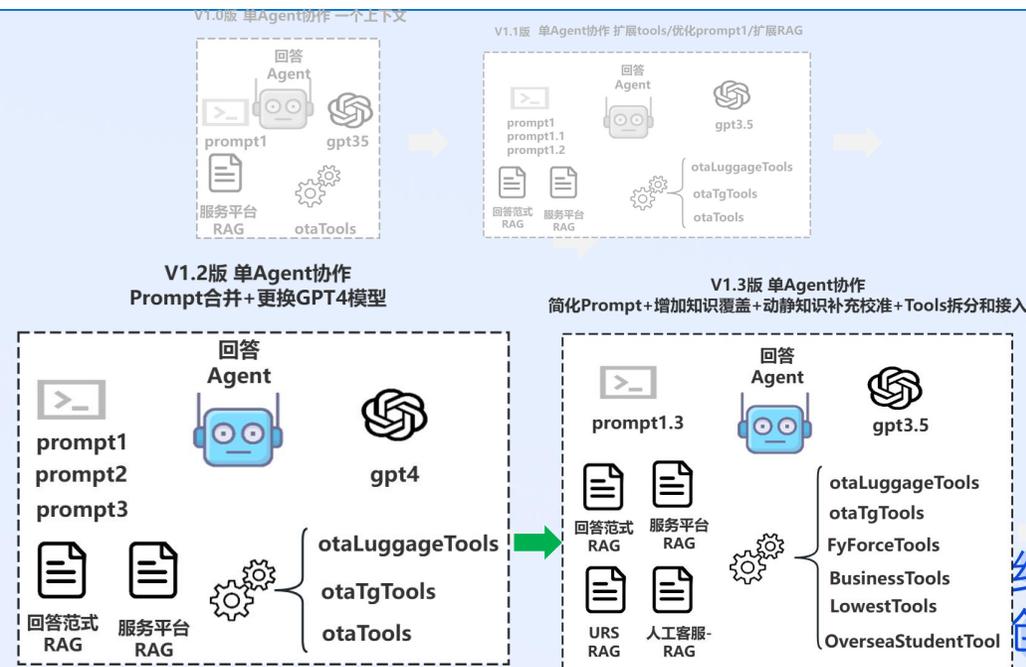
1. 价格：根据您的预算选择适合的价格区间。
2. 舱位等级：经济舱、商务舱等不同舱位等级提供不同的服务和舒适度，根据个人需求选择。
3. 退改签政策：查看退票和改期的政策，选择符合您需求的报价。
4. 行李额度：根据您的行李需求选择包含免费托运行李的报价。
5. 增值服务：一些报价可能包含增值服务，如退票保障、改期保障等，根据个人需求选择是否购买。

综合考虑以上因素，您可以选择最符合您需求的机票报价。

无工具，有RAG
35s返回

有RAG、有工具回调
52s返回

回答不够用户



V1.3 版本问题

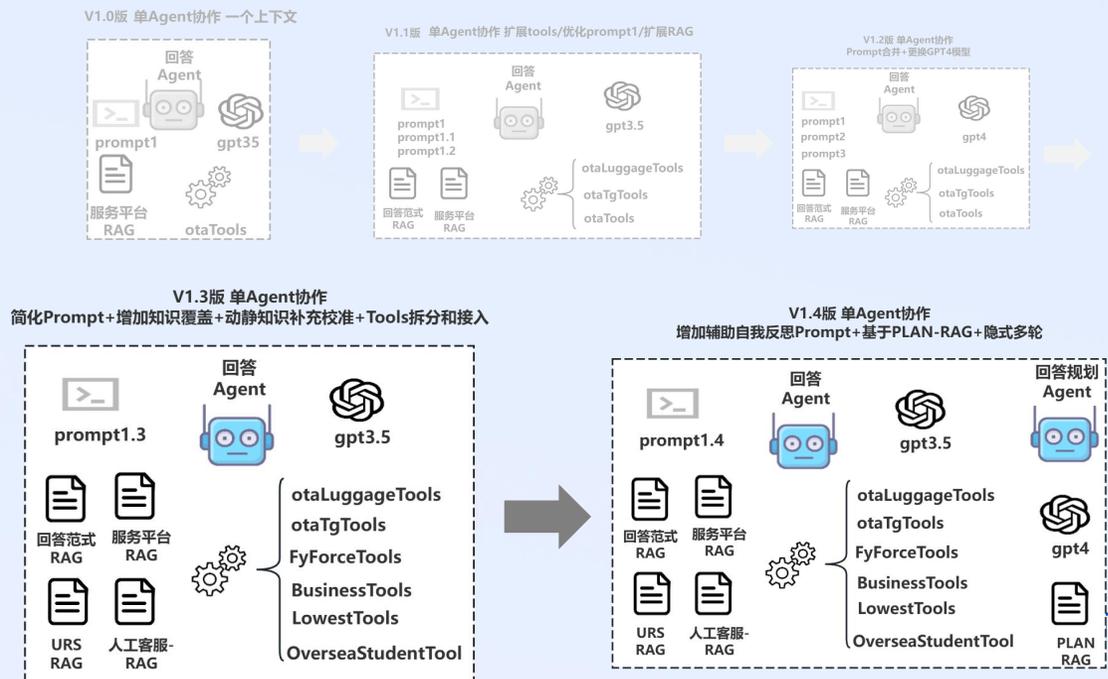
- 1) 【动态参考回答】优化成本高、效果慢
- 2) 开放式问题回答准确率、一致率低

V1.4 尝试解决

- 1) 测试其他模型性能，探索可能性
- 2) 增加回答规划Agent，生成PLAN-RAG
- 3) Prompt增加自我反思流程
- 4) 尝试隐式多轮，假设性问答
- 5) 增加回答复盘Agent

售前助手V1.4-计划中

【开放式回答】初始版--->【开放式回答】优化版



Agent测试方案版本迭代记录

单Agent版本	v1.0	v1.1	v1.2	v1.3
版本功能	<ol style="list-style-type: none"> 1) 基于RAG数据回答 2) 基于OAT实时数据回答 3) 打招呼功能, 默认推荐3个问题 	<ol style="list-style-type: none"> 1) RAG校准、补充回答范式RAG知识 2) 细化Tools 补充报价位置信息, 拆分退改、行李等tools 3) 优化Prompt, 并且QConfig配置化 	<ol style="list-style-type: none"> 1) Prompt问题同义词默认知识增加 2) 升级模型到GPT4 	<ol style="list-style-type: none"> 1) 降级回GPT 3.5, 拆分Agent职责, 多Agent并行 2) RAG拆分和补充, Q-RAG/ QA-RAG/URS-RAG/人工客服RAG 3) RAG分类和校准, 【静态答案】校准、【动态参考回答】优化 4) Tools继续拆分
版本问题	<ol style="list-style-type: none"> 1) RAG数据少, 答案过时不正确 2) Prompt长, gpt35不完全遵循规则 3) 回答准确率低 3.1 无意义、流水式回答 3.2 报价引导弱, 无相关报价位置 4) 默认问题推荐不合理 	<ol style="list-style-type: none"> 1) 默认推荐问题, 小孩关键字命中风控 2) 回答不稳定 	<ol style="list-style-type: none"> 1) 耗时慢 平均35s 	
Agent组件图	<p>V1.0版 单Agent协作 一个上下文</p> 	<p>V1.1版 单Agent协作 扩展tools/优化prompt1/扩展RAG</p> 	<p>V1.2版 单Agent协作 合并多agent的Prompt Tools+RAG+更换GPT4模型</p> 	

成果总结

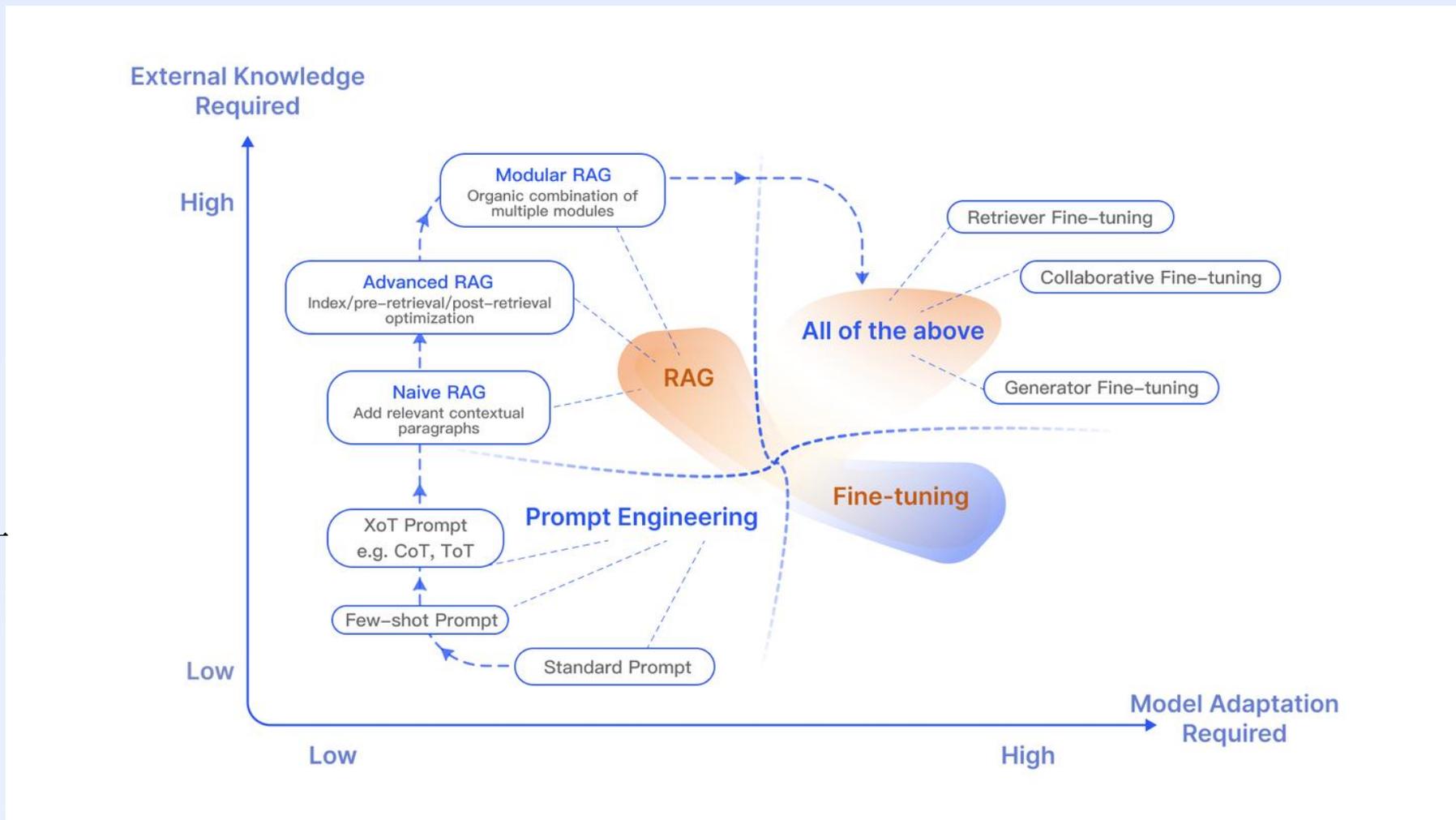
Agent优化的四个方向

一、优化Prompt

二、优化RAG

三、优化Tools

四、Model和Agent切换设计



Prompt 从Zero-Shot到Few-Shot

1) 默认通用知识、同义词知识补充

```
## 默认知识  
- 什么是增值服务：退票保障、改期保障等服务介绍，帮助用户了解增值服务的价格和作用。  
- 小孩的含义：在中文语意中，小孩、孩子、儿童是同等语意，可以互换使用
```

2) RAG回答范式补充Prompt

售前助手-示例

商务出行有推荐的报价吗？

商务出行通常需要考虑航班的准时性、舒适度以及服务质量。建议选择以下几个方面的航班：

- 直飞航班**：直飞可以减少中转的时间和不确定性，提高出行效率。
 - 头等舱或商务舱**：提供更宽敞的座位和更好的服务，有助于商务旅客在飞行途中休息或工作。
 - 有良好口碑的航空公司**：选择服务质量高、准时率高的航空公司。
 - 灵活的退改签政策**：商务行程可能会有变动，选择退改签政策较为灵活的机票可以减少损失。
 - 附加服务**：如优先值机、额外行李额等服务，可以提供更多便利。
- 如果您需要具体的航班推荐，请提供出发地、目的地、出行日期等信息，我可以为您推荐合适的机票报价。

优化前

VS

RAG回答范式补充Prompt

```
## 参考回答格式：  
问题：商务出行有推荐的报价吗？## 回答：您好，为了更好地旅行体验，我们推荐您考虑【头等/公务/商务舱】机票。 <br>  
一个<商务舱/公务舱的机票报价产品信息> <br>  
一个<头等舱的机票报价产品信息>。
```

商务出行有推荐的报价吗？

商务出行的推荐报价如下：

- 头等舱机票报价：¥ 93,828，退票费用6000元起，免费改期，每人托运2件行李，享受7*24小时VIP专家客服。
 - 商务舱机票报价：¥ 29,048，退票费用5000元起，改期费用1500元起，每人托运2件行李，享受7*24小时VIP专家客服。
- 请注意，以上报价仅供参考，具体价格可能会有所变动。

(以下的模块RAG, 包含 6 种模块类型、14 个模块和 40 多个运算符, 即优化方式)

RAG模块和运算符选择

1、附加元数据Metadata Attachment

可以使用元数据信息来丰富块, 例如分类、特征、摘要等, 随后根据该元数据过滤检索, 从而限制搜索范围。

2、尝试使用元数据过滤

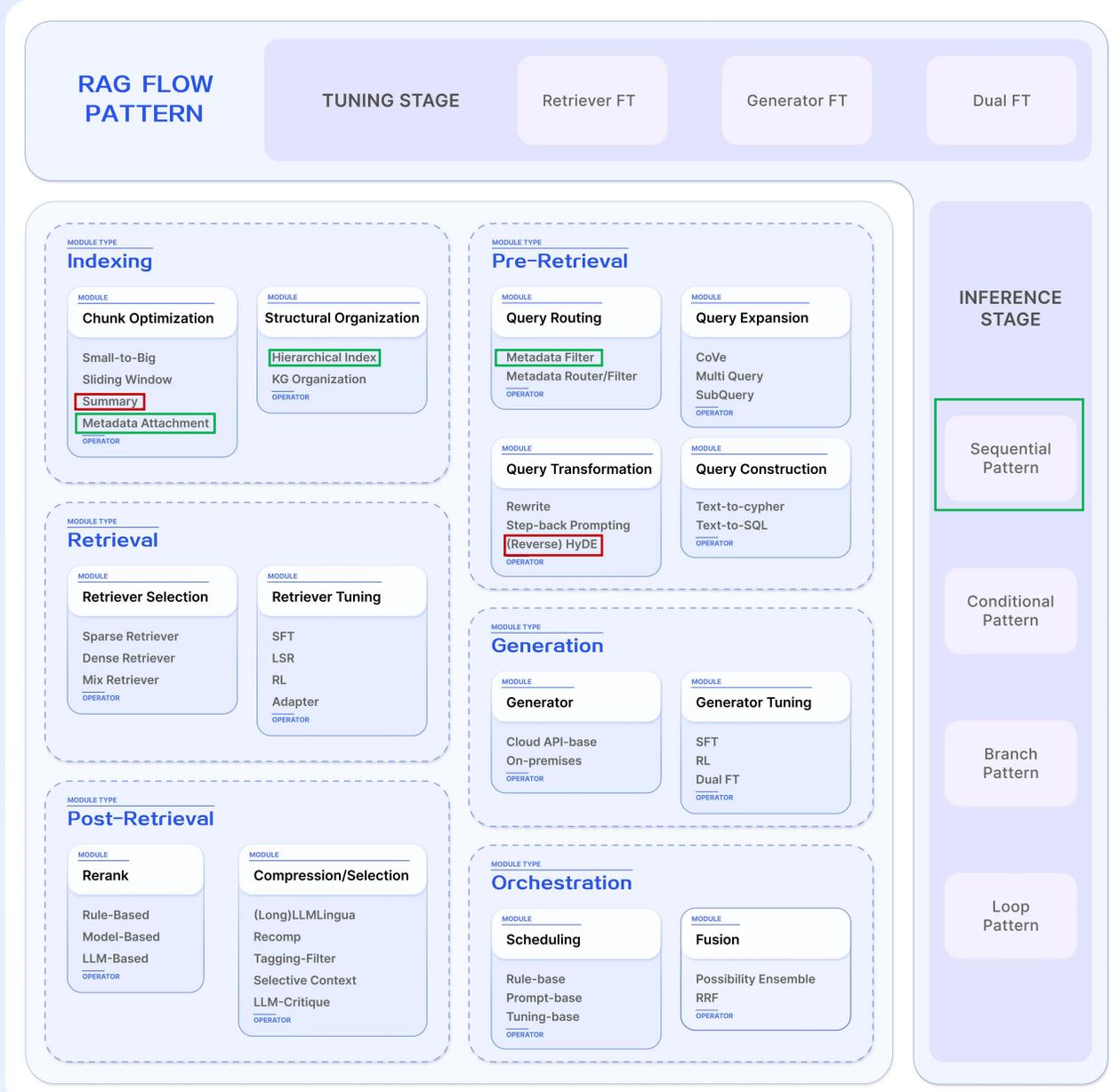
相似 \neq 相关, 根据块中的关键字和元数据进行过滤以缩小搜索范围, 增加知识相关性

3、层次索引Hierarchical Index

在文档的层次结构中, 节点以父子关系排列, 块链接到它们, 增加知识的递进性。

RAG Flow流程选择

微调阶段: 无
推理阶段: 顺序模式。

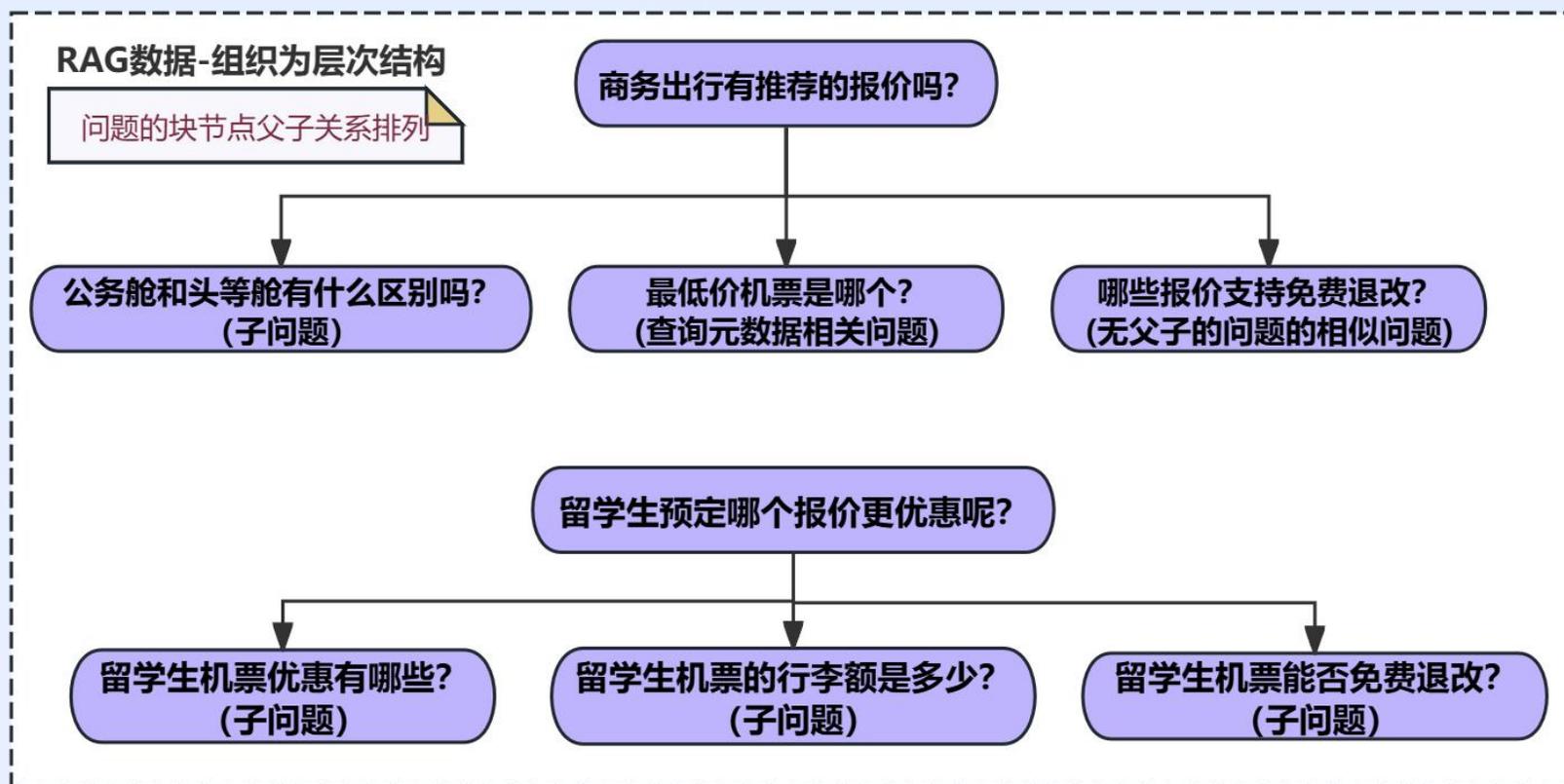


RAG模块和运算符选择

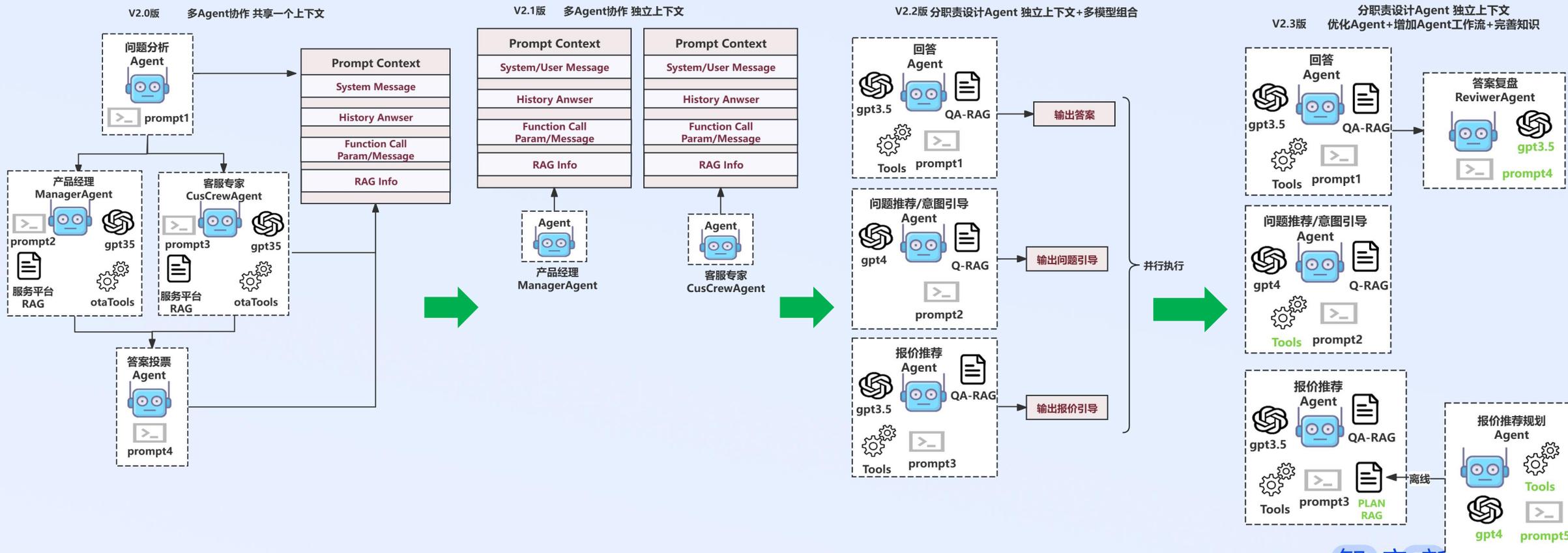
层次索引Hierarchical Index

在文档的层次结构中，节点以父子关系排列，块链接到它们，增加知识的递进性。

question_embedding表的id 和 question_relation表的id进行关联



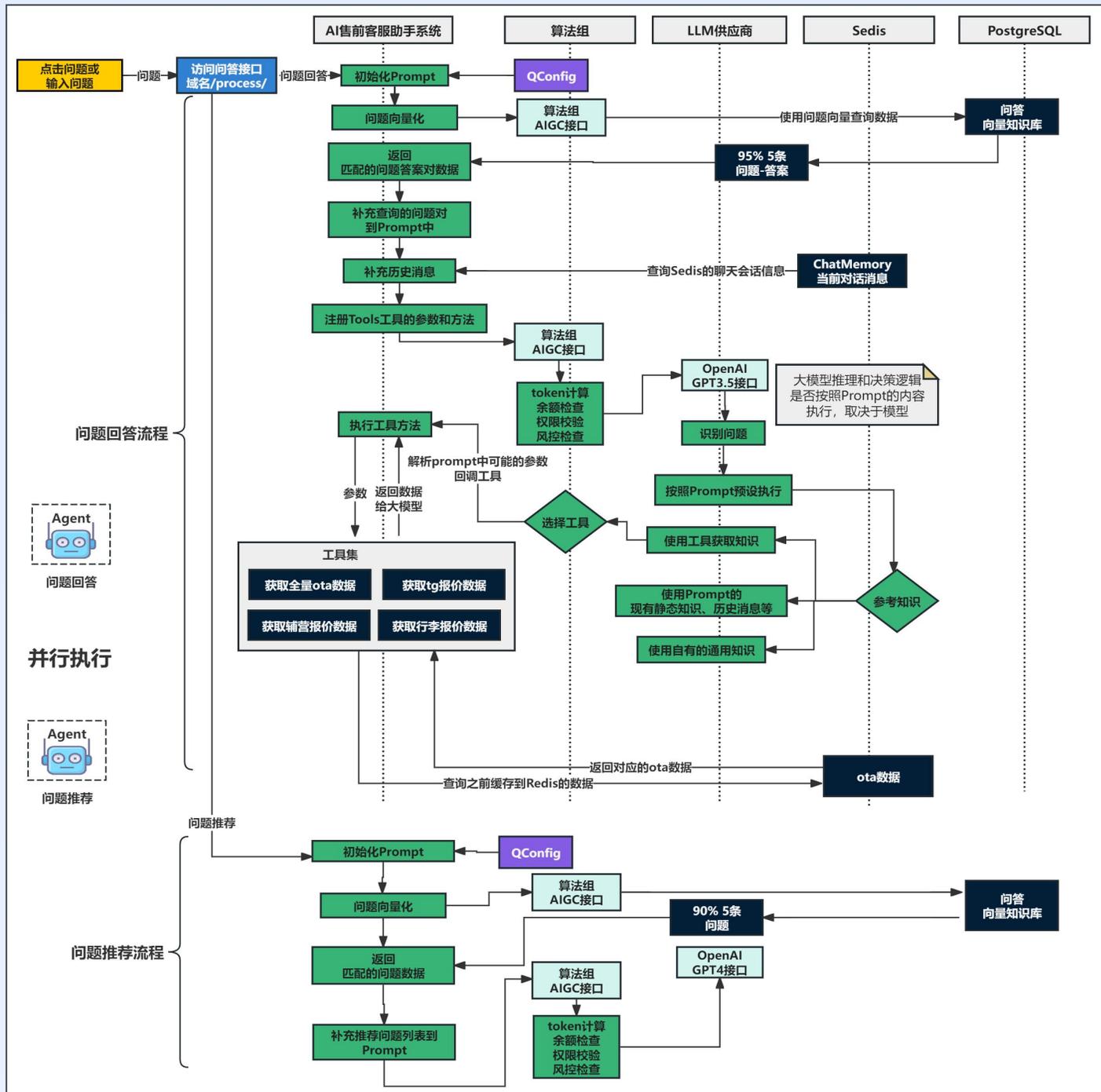
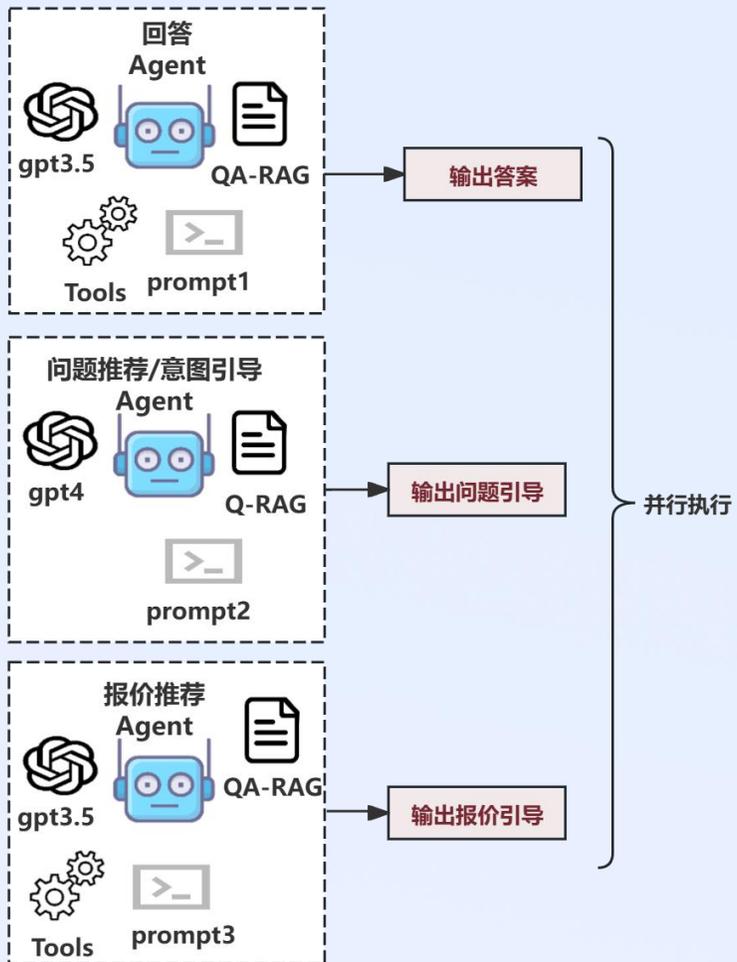
多Agent



多Agent

售前助手-示例

V2.2版 分职责设计Agent 独立上下文+多模型组合





谢谢观看

THANKS