

智能体总体架构

一、核心概念：Agent 的双层记忆架构 (Dual-Memory System)

在流程图开始前，我们先明确 Python 智能基座层需要实现的两层记忆：

1. 短期记忆 (Short-term / Working Memory)：

- **载体**：Redis (`ChatSession` + `MemoryBuffer`)。
- **作用**：维持当前这“一趴”对话的连贯性。通常保留最近的 N 轮对话历史（如最近 10 轮或 4000 个 Token 内的历史）。超出窗口的老旧对话会被丢弃或压缩。

2. 长期记忆 (Long-term / Episodic Memory)：

- **载体**：ES Serverless (Vector 存储) + MySQL (`UserProfile`)。
- **作用**：跨会话的用户画像与事实留存。
- **实现机制（也是一种工具）**：Python 后台实现一个 `update_memory` 工具。当大模型在聊天中敏锐地捕捉到用户的关键事实（例如：“我平时主要用 Java 和 React”、“我有个下个月要交的 PPT”），模型会主动调用这个工具，将事实提取、向量化并持久化到 ES 中。下次用户再来时，系统会先通过向量检索把这些长期事实捞出来，作为系统 Prompt 注入。

四、架构师详细逻辑说明 (针对 Agent 与 Memory)

1. 记忆系统的运作闭环 (Memory Lifecycle)

- **读记忆（前置动作）**：用户一句“帮我把那张图做成汇报”，模型怎么知道是哪张图？系统会在调用 LLM 前，去 Redis 读取最近 10 轮对话（知道是上一轮发的架构图），同时去 ES 长期记忆库拉取用户偏好（知道用户的汇报对象通常是技术高管）。
- **写记忆（工具动作）**：这是一个极其巧妙的设计。大模型的工具清单里有一个名为 `update_long_term_memory(key_fact, confidence_score)` 的函数。在多次对话中，如果模型发现用户反复提及某个核心项目，模型会自主决定调用这个工具，将该事实永久写入 ES。

2. “While 循环”的破局 (Agentic Loop)

传统的接口是：`Request -> LLM -> Response`。

现在的 Agent 接口是：`Request -> While(LLM.wants_to_use_tool) { Execute_Tool ->`

`Feed_back_to_LLM } -> Final Response`。

在这个循环中，大模型可能会自己跟自己“聊”好几个回合。比如：先调用“**Todo 工具**”发现有任务，再调用“**知识检索工具**”查相关的资料，最后再调用“**多媒体工具**”生成一段总结音频。所有的内部思考过程对用户是透明的，前端只会显示一个优雅的“Agent 思考中：正在检索知识... 正在生成图表...”。

3. 展现层的多模态渲染 (Artifacts UI)

由于我们赋予了模型 `ToolCode` (生成交互代码) 和多媒体等能力，前端 (React 19) 收到 SSE 流时，不能仅仅把 Markdown 当文本丢给解析器。前端需要实现**自定义渲染拦截器**：

- 收到纯文本：走正常的打字机效果。
- 收到特定标记的 React/HTML 代码块：在页面右侧或对话流中展开一个沙盒视窗 (Iframe 或 Shadow DOM)，直接将代码编译渲染为可交互的 UI 组件。
- 收到音频流链接：渲染原生的 Audio 播放器。



