

Manus AI Agent 终极技术解析：从爆火到反思的全链路架构深度复盘



AI Agent 研究

2025-12-31 北京

阅读 6 分钟

本文是目前关于 Manus AI Agent 最全面、最深入的技术解析，整合了 Manus 联合创始人季逸超的官方技术复盘、核心开发团队的工程实践总结以及第三方深度技术分析，涵盖从架构设计、核心机制到经验教训的完整链路，总计 **12 大核心维度**，超过 **100 个技术细节点**，为 AI Agent 开发者提供一份兼具技术深度与实践价值的参考范本。

一、Manus 的技术抉择：上下文工程 vs 端到端模型

1.1 路线之争：AI Agent 的两种进化路径

技术路线	核心特点	迭代周期	模型依赖	灵活性	典型代表
端到端 自训练 模型	需大规模标注数据，任务专用模型	数周-数月	强依赖，模型淘汰则产品失效	低，难以适配新任务	早期专用 AI 系统
上下文 工程	基于提示工程 + 架构优化	数小时-数天	弱依赖，随模型进步	高，无缝切换	Manus、Cursor

技术路线	核心特点	迭代周期	模型依赖	灵活性	典型代表
	化，与底层模型解耦		步而提升	模型生态	

1.2 Manus 的战略选择：押注上下文工程的核心逻辑

- 历史教训：**团队吸取前一轮创业中“自研模型被淘汰”的惨痛经历，放弃重资产路线
- 工程效率：**产品迭代周期从数周缩短到数小时，快速响应市场需求
- 模型解耦：**作为“船”随潮水而动，而非“搁浅的柱子”，底层模型进步时自动受益
- 生态适配：**支持自托管模型与 API 调用无缝切换，覆盖私有化部署与云服务场景

1.3 上下文工程的本质：“手工随机梯度下降”

Manus 团队在短时间内**四次重构 Agent 框架**，每次都是对上下文塑造方式的深度反思。上下文工程并非简单的提示堆砌，而是通过架构搜索、提示调整、性能评估的循环，不断逼近局部最优解的实验科学。

二、分布式多智能体架构 (MAS)：分工协作的神经中枢

2.1 核心设计：规划-执行-验证三层闭环

Manus 摒弃传统单体大模型，采用专业化子智能体协同机制，将复杂任务拆解为三大核心模块，形成“分而治之”的高效协作体系。

Agent类型	核心职责	技术实现	关键指标	典型应用场景
Agent)	优先级调度、生成 todo.md	+思维链推理	响应 <500ms	杂任务

Nginx 源码分析

[注册](#) [登录](#)
[主页](#) [关于](#) [RSS](#)

Agent)	优先级调度、生成 todo.md	+思维链推理	响应 <500ms	杂任务
执行代理 (Execution Agent)	调用工具/ 代码/API 完成具体 操作、环 境交互	CodeAct 机制 +Docker 沙盒执行 环境	工具调用 成功率 >90%，执 行延迟<2s	数据爬 取、文 档生 成、网 页自动 化
验证代理 (Verification Agent)	结果交叉 校验、逻 辑验证、 错误检 测、迭代 优化	多轮测试 +反馈闭环 +模型自评 估	错误识别 率>85%， 结果修正 率>80%	财务报 表核对、 代 码审 计、 数 据一 致 性验 证

2.2 协同机制：Validation Loop 相互校验系统

- 状态同步：各 Agent 通过共享内存实时同步任务进度与执行结果
- 错误传递：执行失败时自动触发验证代理介入，分析失败原因并生成修正方案
- 动态调整：规划代理根据验证结果实时调整子任务优先级与执行路径
- 鲁棒性提升：通过该机制，系统整体容错能力提升约 60%，避免单一模块出错导致任务失败

三、KV 缓存：Agent 系统的性能生命线

3.1 核心原理：输入输出比例失衡下的成本优化

Manus 的平均输入输出 token 比例高达 **100:1**（复杂任务可达 1000:1），KV 缓存通过存储模型前缀的中间状态，实现重复前缀的高效复用，将推理成本降低 **10 倍以上** (Claude Sonnet：缓存输入 0.3 美元/百万 token，未缓存 3 美元/百万 token)。

3.2 提高 KV 缓存命中率的五大实践准则

1. **稳定提示前缀**：避免在系统提示中引入时间戳等动态内容，哪怕一个 token 差异也会导致缓存失效
2. **上下文仅追加不修改**：确保序列化过程确定性，避免 JSON 键顺序变化等隐性破坏缓存的因素
3. **显式标记缓存断点**：手动插入断点，合理分配缓存空间以防过期
4. **分布式一致性路由**：自托管模型时，通过会话 ID 确保跨工作器的缓存一致性
5. **工具定义固化**：保持工具描述稳定，避免频繁修改导致缓存失效

3.3 缓存失效的代价与应对策略

- **代价**：缓存失效会导致推理延迟增加 3-5 倍，成本上升 10 倍，用户体验显著下降
- **应对**：实现缓存预热机制，在任务开始前预加载高频工具描述与系统提示；建立缓存失效预警系统，实时监控命中率并自动调整策略

四、MCP 协议：智能体通信的高速公路 (真相澄清)

4.1 关于 MCP 的核心事实

重要澄清：Manus 联合创始人季逸超在技术复盘中明确表示，Manus **并未采用 MCP 协议**，而是借鉴其思想，自研了基于 CodeAct 的通信机制。市场上流传的“MCP 是 Manus 秘密武器”属于误读。

4.2 Manus 实际采用的通信架构

1. **CodeAct统一接口**: 将所有 Agent 行动转化为可执行的 Python 代码，实现跨模块通信标准化
2. **基于 SSE 的实时传输**: 采用 Server-Sent Events 技术支持异步通信，实现双向上下文传输与数据同步
3. **状态管理模块**: 自研 Lifecycle 状态管理系统，处理任务创建、执行、暂停、取消、完成全流程
4. **消息流优化**: 实现 Ping/Cancellation/Progress 等机制，确保多任务场景下的高效消息传递与高容错

五、CodeAct 机制：“代码即行动”的执行引擎

5.1 核心创新：将任务执行转化为代码编写与运行

Manus 的核心执行范式，在云端 Linux 沙盒环境中实时执行 Python/JavaScript 代码，实现精准的复杂操作，区别于传统 Agent 的“工具调用+参数传递”模式。

5.2 技术实现的六大关键细节

5.2.1 沙盒隔离环境

- 每个任务独立运行于 **Docker Ubuntu 22.04 LXC 容器**，避免资源冲突与安全风险
- 支持**异步执行**：用户离线后任务仍可继续运行，完成后通过邮件/通知反馈结果
- 资源限制：CPU≤2核，内存≤4GB，网络带宽≤10Mbps，防止恶意操作占用资源

5.2.2 16 种核心 Action 空间（完整列表）

操作类型	具体功能	适用场景	准确率
Computer Use	屏幕截图、拖拽、点击、打字、文件读写	桌面应用自动化	95%+
Browser Use	网页导航、DOM 操作、表单填写、元素提取	网页数据爬取、自动化测试	92%+
Shell Command	执行 Linux 命令、安装软件、管理文件	系统配置、环境搭建	98%+
Code Execution	Python/JavaScript 代码运行、结果解析	数据分析、模型训练	90%+
API Call	RESTful API 调用、参数生成、响应解析	第三方服务集成	88%+
File Management	创建、读取、修改、删除文件/目录	文档处理、数据存储	99%+
PDF Processing	提取文本、合并/拆分 PDF、添加水印	文档分析、报告生成	93%+
Image Processing	截图裁剪、格式转换、OCR 识别	图像分析、内容提取	91%+

5.2.3 自纠错重试机制

1. 执行失败时自动分析错误类型（语法错误、运行时错误、环境错误）
2. 根据错误类型生成修正方案，重新编写代码并执行（最多重试3次）
3. 失败重试率达**85%以上**，在GAIA基准测试中Level 3任务胜率达**57.7%**，远超OpenAI同类产品

5.2.4 多模态交互增强

- 结合截图进行what/how/when判断，提升复杂UI操作的准确性
- 实现“视觉-行动”闭环：截图→分析界面元素→生成操作代码→执行→验证结果
- 非标准网页处理能力显著优于传统工具调用型Agent

5.2.5 成本优化策略

- 代码片段复用：缓存高频执行代码，减少重复生成成本
- 资源动态调度：根据任务复杂度自动调整容器规格，降低闲置资源消耗
- 单任务成本控制在**2美元**以内，显著低于行业平均水平（5-10美元）

5.2.6 安全防护措施

- 代码审查：执行前检查代码是否包含恶意操作（如rm -rf /、curl恶意链接）
- 网络隔离：限制访问敏感网站与服务，防止数据泄露
- 操作审计：记录每一步执行日志，支持追溯与复盘，确保合规性

六、外部化记忆：文件系统作为“无限上下文”

6.1 长上下文的三大痛点与传统解决方案局限

痛点类型	具体表现	传统解决方案	方案缺陷
观察结果庞大	网页、PDF等 内容超上下文 窗口	上下文截断/压 缩	不可逆信息丢 失，影响决策 准确性
模型性能下降	上下文过长导 致注意力分散	摘要生成	关键细节丢 失，无法恢复 原始信息
传输成本高昂	长输入即使缓 存，传输与预 填充耗时	分段处理	增加系统复杂 度，易出现上 下文断裂

6.2 Manus 的创新：文件系统作为终极外部记忆

1. **核心思想：**将文件系统视为“无限上下文”，模型可按需读写文件，将长期状态外部化
2. **可逆压缩策略：**始终采用可恢复压缩，保留 URL 即可删除网页内容，保留路径即可省略文档内容
3. **操作流程：**
 - 模型读取大型内容时，自动保存到文件系统并生成引用链接
 - 后续需要时通过链接重新读取完整内容，不占用上下文窗口
 - 上下文长度缩短 **80%+**，同时不丢失关键信息

6.3 外部记忆的实践应用

- **todo.md 动态更新：**在复杂任务中不断重写该文件，将全局计划复述到上下文末尾，确保模型始终关注核心目标
- **错误日志持久化：**将失败操作与观察结果保存到文件，为模型提供自我修正的完整依据
- **多 Agent 记忆共享：**通过文件系统实现跨 Agent 状态同步，提升团队任务处理效率

七、工具管理：屏蔽而非移除的动态状态机

7.1 动作空间膨胀的困境

随着 Agent 能力扩展，工具数量爆炸式增长（Manus 支持 **47 种原生工具** + 用户自定义工具），导致：

- 模型选择错误动作的概率增加
- 工具定义位于上下文前部，任何变动都会导致缓存失效
- 历史操作引用已删除工具，引发模型困惑与幻觉

7.2 Manus 的解决方案：上下文感知的状态机

1. **核心机制：**不直接删除工具，而是在解码期间通过**屏蔽 token logits** 约束动作选择
2. **实现细节：**
 - 设计一致前缀的动作（如 browser_、shell_），便于分组屏蔽
 - 响应预填充机制，灵活实现自动/必需/禁止三种函数调用模式
 - 保持工具定义稳定，最大化 KV 缓存命中率

管理方式	优点	缺点	KV 缓存影响	适用场景
动态增删工具	灵活，动作空间最小化	缓存失效，历史引用混乱	极大负面	工具极少且固定的场景
屏蔽 token logits	保持上下文稳定，缓存友好	需设计一致前缀，屏蔽逻辑复杂	极大正面	Manus 等多工具复杂场景

八、注意力操纵：通过“背诵”减少任务漂移

8.1 任务漂移的风险与成因

复杂任务往往需要数十次工具调用，LLM 驱动的 Agent 极易“跑题”，遗忘早期目标，尤其在长上下文下更为突出。成因包括：

- 注意力机制倾向于关注近期内容
- 长上下文导致早期目标被稀释
- 多步骤推理中目标逐步模糊

8.2 Manus 的实践：动态复述目标的三大策略

1. **todo.md 全局计划更新**：每 5-10 步操作后重新生成该文件，将核心目标推送到模型近期注意力范围
2. **上下文末尾追加摘要**：在每次工具调用后，自动生成任务进度摘要并添加到上下文末尾
3. **自然语言强化提示**：通过“现在我需要完成 ...”“我的最终目标是 ...”等句式强化模型对目标的关注

8.3 注意力操纵的效果

- 多步任务完成率提升 **40%+**
- 任务漂移概率降低 **65%**
- 无需架构改动，仅通过自然语言提示即可实现目标一致性

九、错误处理：保留而非隐藏的自我修正机制

9.1 错误是 Agent 的常态

在多步骤任务中，模型幻觉、环境异常与工具故障等错误在所难免：

- 工具调用失败率约 10%
- 执行结果错误率约 15%
- 任务陷入死循环概率约 5%

9.2 保留错误的核心价值

- 1. 自我修正基础：**将失败操作与观察结果保留在上下文中，模型可隐式调整内部信念，降低重复犯错概率
- 2. 错误模式识别：**通过分析历史错误日志，模型自动学习常见错误类型与应对策略
- 3. 用户信任提升：**透明展示错误与修正过程，增强结果可解释性

9.3 容错策略全解

异常类型	处理机制	恢复率	重试次数	降级方案
工具调用失败	重试 2-3 次 → 更换工具 → 求助人类	90%+	3 次	简化任务流程，使用替代工具
执行结果错误	验证代理 校验 → 重新执行 → 参数调整	85%+	2 次	降低精度要求，提供近似结果
任务陷入死循环	超时检测 (默认 30 分钟) → 中断任务 → 重新规划	95%+	1 次	拆分任务为更小模块，分步执行
模型幻觉	交叉验证 → 事实核查 → 重新生成	80%+	2 次	引用权威数据源，标注不确定性

十、分层记忆管理：Working-Hot-Cold 三层体系

10.1 核心架构

Manus 采用 **Working-Hot-Cold Memory Orchestration** 三层记忆体系，实现高效数据协同与实时更新。

记忆层级	存储内容	生命周期	技术选型	访问频率	容量限制
Working Memory	当前任务上下文、执行步骤、临时结果	任务存续期	大模型上下文窗口+本地缓存	极高	8k-32k tokens
Hot Memory	近期任务经验、用户偏好、高频工具调用记录	7-30天	Redis+向量索引	高	1M+ tokens
Cold Memory	领域知识库、历史任务归档、长期经验	永久	Chroma/Milvus 向量库+知识图谱	中低频	无限（文件系统存储）

10.2 关键优化

1. **自编辑记忆 (Self-editing Memory)**: Agent 可自主更新记忆内容，剔除无效信息，提升检索效率
2. **LangGraph Store 优化**: 实现记忆的结构化存储与高效检索，检索响应时间 < 100ms
3. **记忆优先级调度**: 根据任务类型动态调整各层级记忆的访问权重，确保关键信息优先被检索

十一、Human-in-the-Loop: 人机协作的深度融合

11.1 实时用户交互反馈闭环

1. **动态断点 (Breakpoints)**: 在关键任务节点自动暂停，等待用户确认后继续执行
2. **Streaming 与异步技术**: 支持实时进度展示与用户中断，提升交互体验
3. **用户反馈机制**: 用户可直接修改执行结果，模型自动学习并应用到后续步骤

11.2 Time Travel 功能: Agent 的"后悔药"

1. **核心能力**: 支持任务状态回溯与重放，用户可查看任意步骤的执行过程与结果
2. **技术实现**: 通过文件系统记录每一步操作的完整状态，包括上下文、工具调用、执行结果
3. **应用场景**:
 - 调试复杂任务流程
 - 分析失败原因并重新执行
 - 合规审计与操作追溯

11.3 权限分级控制

权限级别	操作范围	适用用户	安全保障
完全自主	无需用户干预完成全流程	信任用户，简单任务	沙盒隔离+操作审计
半自主	关键节点需要用户确认	普通用户，复杂任务	动态断点+用户反馈
人工主导	每步操作需用户批准	新用户，高风险任务	权限最小化+实时监控

十二、Manus的经验教训与AI Agent的未来趋势

12.1 四大核心教训

- 上下文工程的科学化：**需建立系统化的评估指标与优化流程，避免"手工随机梯度下降"的盲目性
- 外部记忆的必要性：**文件系统作为外部记忆将成为Agent突破长期依赖瓶颈的关键
- 错误恢复能力的重要性：**这将成为衡量Agent架构的核心指标，而非仅在理想条件下的"完美表现"
- 人机协作的不可替代性：**无论Agent多么智能，人类在复杂决策、价值判断与创意生成方面仍具优势

12.2 未来趋势预测

- SSM+外部记忆融合：**状态空间模型（SSM）若能掌握基于文件的外部记忆，有望突破Transformer的长期依赖瓶颈
- 上下文工程自动化：**通过大模型自动生成优化提示与架构设计，减少人工干预
- 多Agent生态协同：**从单一Agent到Agent团队，通过标准化协议实现跨平台、跨组织的智能协作

4. 安全与合规强化：Constitutional AI (CAI) 与 RLHF 技术将成为 Agent 标配，确保输出的安全可靠

终极架构总结：Manus的核心竞争力

Manus的技术架构通过“**上下文工程+CodeAct执行+外部记忆+动态容错**”四大核心优势，实现了从“AI工具”到“AI协作者”的质变：

核心维度	Manus 创新点	行业平均水平	提升幅度
任务完成率	复杂任务完成率 85%+	50% 左右	70%+
推理成本	单任务成本 2 美元	10-20 美元	80%+
上下文效率	长度缩短 80%+, 无信息丢失	长度缩短 30-50%, 信息丢失严重	200%+
错误恢复能力	综合恢复率 90%+	60% 左右	50%+
迭代周期	数小时-数天	数周-数月	10 倍+

Manus的技术复盘让Agent研发回归本质：技术创新，工程为王，细节决定成败。未来的AI Agent 竞争，既是模型能力的竞赛，更是工程科学的较量。唯有在实践中不断试错、总结与优化，才能在智能体的浪潮中立于不败之地。

赞

收藏

分享

阅读 689 · 发布于 2025-12-31

举报



AI Agent 研究

7.2k 声望 · 12.8k 粉丝 · 63

一群有AI的人 研究AI-Agent的开发，做优秀的AI应用；

[关注作者](#)[« 上一篇](#)[下一篇 »](#)[AI Agent 完整设计指南（全维度、... 新的一年，如何好好学习 AI Agent...](#)

引用和评论

推荐阅读

2026年第二周学习——规划与工具调用原理

AI Agent 研究

【实测有效】Gemini 3 / Google Antigravity 授权登录没反应、账号无权限？解决办法汇总指南

uiuihaoAICG · 赞 2 · 阅读 53.2k · 评论 3

一步封神：Unity 环境搭建终极全宇宙级攻略（Win/Mac/云）

黄花花 · 赞 3 · 阅读 4.2k · 评论 3

AI 推理硬件选型指南：CPU 与 GPU 的抉择

思否编辑部 · 赞 1 · 阅读 9.2k

Gemini 3 国内使用指南~（支持最新 Gemini 3 Pro、GPT-5 和 Claude 4.5）

爱逃课的篮球 · 赞 1 · 阅读 17.3k



Akamai 推出 Akamai Inference Cloud (AI 推理云)，重新定义人工智能的应用场景与实现方式

思否编辑部 · 赞 1 · 阅读 9.1k



腾讯云 Agent 应用创新大赛收官，智能体迎来加速时刻

思否编辑部 · 赞 1 · 阅读 18.8k

0 条评论

得票

最新



撰写评论 ...



提交评论

评论支持部分 Markdown 语法：**粗体** _斜体_ [链接]
(<http://example.com>) `代码` - 列表 > 引用。你还可以使用 @
来通知其他用户。

©2026 Nginx 源码分析

除特别声明外，作品采用《署名-非商业性使用-禁止演绎 4.0 国际》进行许可

使用 SegmentFault 发布

SegmentFault - 凝聚集体智慧，推动技术进步

服务协议 · 隐私政策 · 浙ICP备15005796号-2 · 浙公网安备33010602002000号