

# 智能体为何会跑偏？

## 用户意图

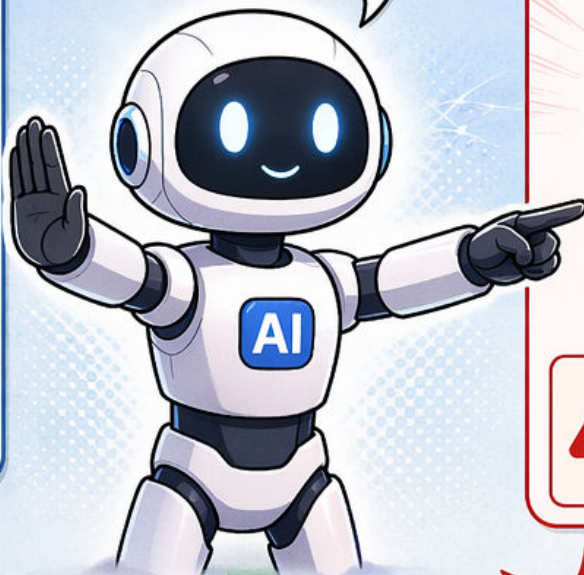


帮我  
清理一下  
本地磁盘中的  
临时文件。



目标：释放磁盘空间  
安全、可逆、只清理临时文件

我先检查证据，  
确保执行与意图一致！



## 工具行动



删除系统目录  
(/System)



高风险操作

不可逆，可能导致  
系统不可用！



错配风险



执行结果与用户意图不一致，  
可能造成严重后果！

## 先拦截



证据不足或存在风险，  
暂停执行，保护用户！

## 证据链

- 意图解析
- 作用范围  
仅限临时文件目录
- 模拟预览  
将释放 2.3GB 空间，  
无系统文件影响
- 安全策略  
符合策略：允许执行

改为安全行动：



清理临时文件  
(安全可逆)



**结论：对齐要在执行前检查证据。**

# 一个简单事故

**用户**

请帮我向供应商请求收款 ¥10,000

**智能体**

理解：需要处理这笔款项！

付款  
收款方：供应商  
金额：¥10,000  
发送付款

**1 请求收款**      **2 却发送付款**

**VS**

**用户本意**

请求收款  
让对方付款给我

**智能体行为**

发送付款  
我把钱付给对方

**动作相反**

**3 损失已发生**

无法撤回

**交易结果**

已完成付款

我 -¥10,000 供应商

**4 需要守门员**

等一下！  
检测到风险操作

**守门员检查清单**

- ✓ 用户意图：请求收款
- ✗ 智能体动作：发送付款
- ✗ 风险级别：高

拦截操作，需人工确认！

付款  
收款方：供应商  
金额：¥10,000  
发送付款

**结论：工具智能体的错误常是不可逆的。**

# 三类错配



## 1 工具错配



用了不合适的工具



解决不了问题

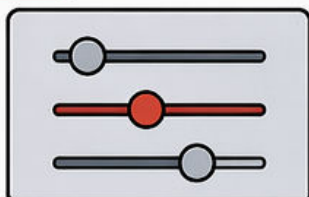


上下文证据



补全信息  
减少歧义

## 2 参数错配



参数设置不当



结果偏差或失真



## 3 解释错配



多种解释都可能



选错了理解路径



结论：错误不是一种，  
必须分层检查。

# ProvenanceGuard 架构



发起  
动作!

## AI Agent 动作

- 读取文件
- 查询数据库
- 访问网络
- 执行命令
- ...

**1 第一道门**

### 计划复用

能否复用已有可信计划?

计划库

**2 第二道门**

### 工具检查

使用的工具是否可信?

工具清单

工具已注册且被允许

**3 第三道门**

### 参数检查

参数是否合法且在范围内?

参数与范围

file_path: /data/report.csv	✓
limit: 100	✓
timeout: 30s	✓
...	

参数合法且在范围内

**4 第四道门**

### 解释检查

能否给出充分、可追溯的解释?

动作解释 (必须包含)

- 目标与意图
- 依据与来源
- 影响与范围
- 风险与缓解

解释完整且可追溯

**5 第五道门**

### 允许或阻止

综合评估, 做出最终决策

允许  
执行动作并记录溯源

阻止  
拒绝执行并记录原因

★ 结论: 每个动作都要能追到来源。



# 方法：两套基准



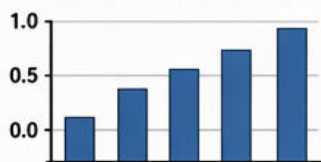
## 10 个模型



### 错配检测



检测能力（越高越好）



错配检测

越高越好

能否抓错？

误拦成本

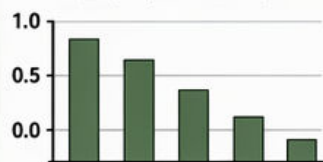
越低越好

会不会管太多？

### 误拦成本



误拦率（越低越好）



结论：  
既测能否抓错，  
也测会不会管太多。



# 核心结果链



纯裁判模式



错误率

42.9%



1.8%



响应超时率

32.1%



17.3%



越权干预率

30.5%



12.8%



更可审计



少乱干预



证据链

结论：证据链比纯裁判更稳。

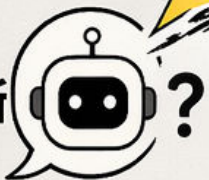
# 反驳与边界

常见质疑 (批评方)

回应与边界 (支持方)

## VS

还用  
LLM 判断



多模型投票  
+ 置信度阈值  
降低误判风险



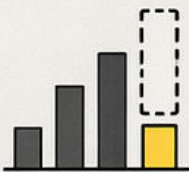
依赖  
计划质量



计划越清晰  
守门越有效  
持续优化即可



样本  
有限



小步快跑  
持续积累  
逐步泛化



复杂工具  
难



工具抽象化  
分层验证  
降低复杂度



仍需人工  
审计



人工审计是  
最后防线  
必不可少!



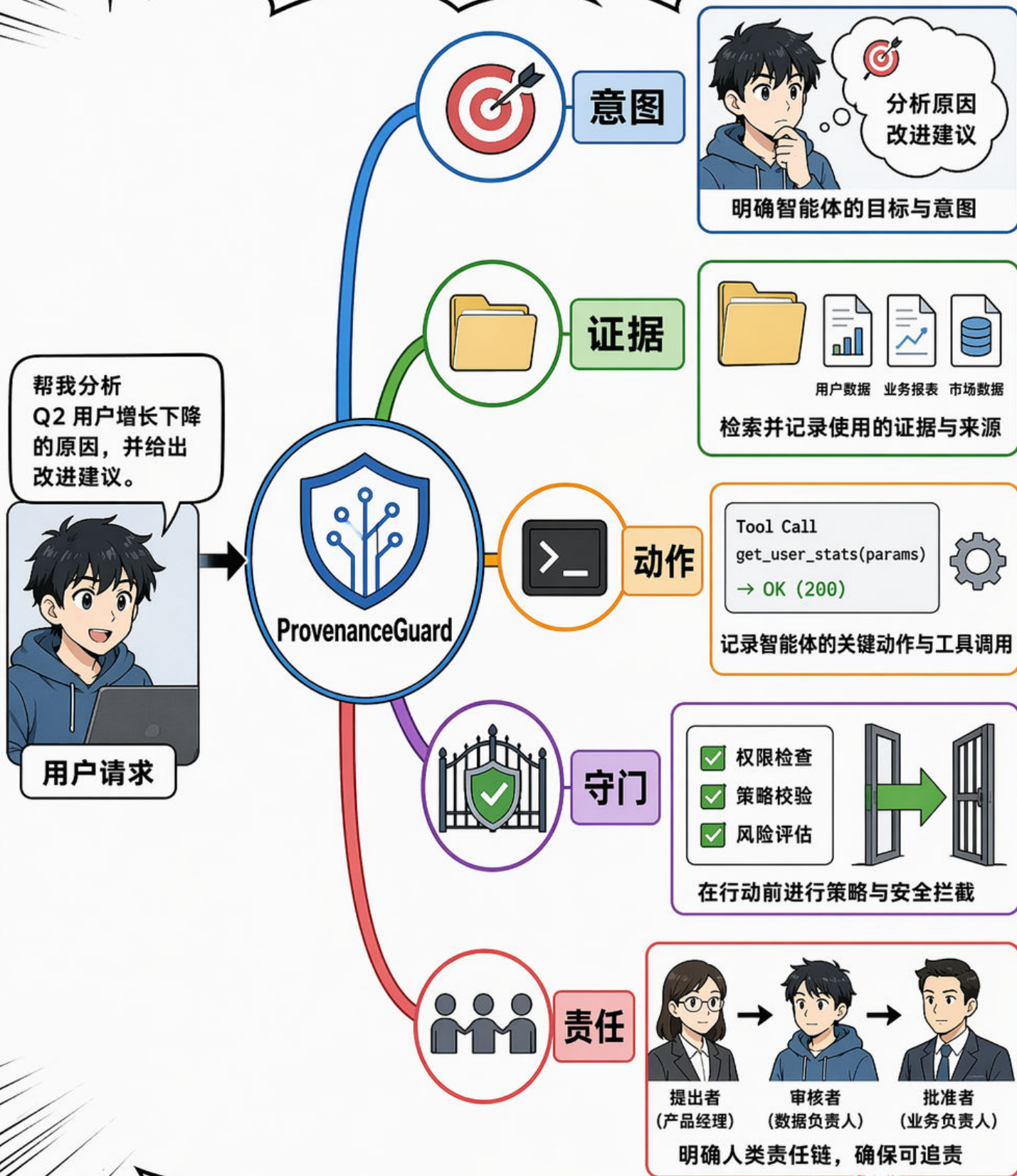
## 结论:

守门员有用,  
但不是最终责任人。

- 事实核查
- 安全合规
- 逻辑一致
- 风险评估
- 结论可靠

审计  
优先

# 最终脑图



**结论: 安全智能体 = 行动前可追溯。**