

# 2026年3月31日： 史上最严重的 AI 架构泄露



## 核心事件定性

并非国家级攻击或零日漏洞。现代高速软件发布流水线中的经典配置纪律缺失，导致 Anthropic 最高机密的终端智能代理源码彻底曝光。

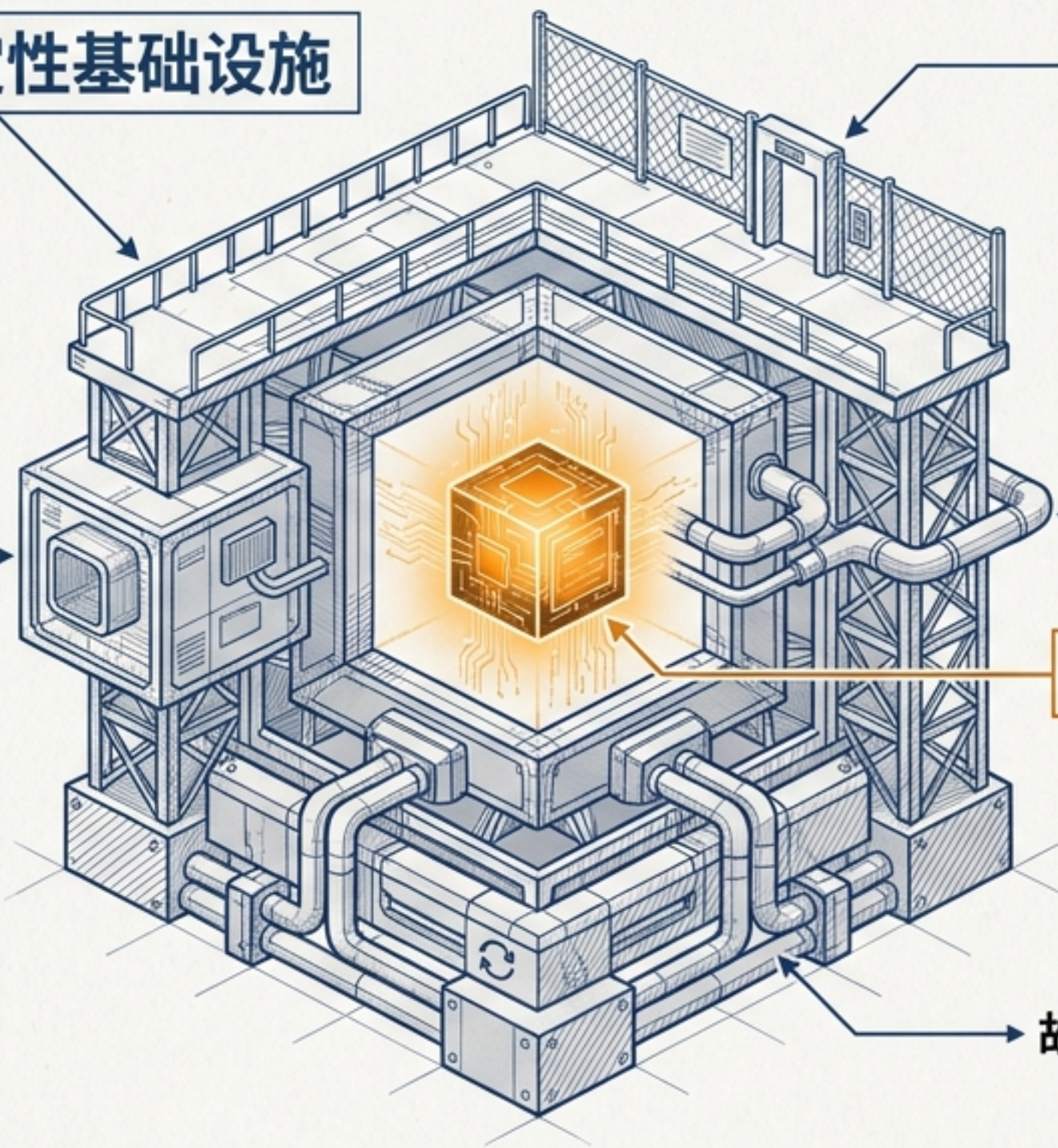
## 技术暴露深度

未被过滤的 Source Map 调试产物，将深度混淆前的核心工程代码完全映射还原，提供了一份显微镜级的 AI 代理架构参考手册。

# 工程范式转移：从对话补全到重型自主代理

98.4% 确定性基础设施

上下文微观  
管理容器



权限网关审查

工具路由分发

1.6% AI 决策逻辑

故障恢复与状态回溯

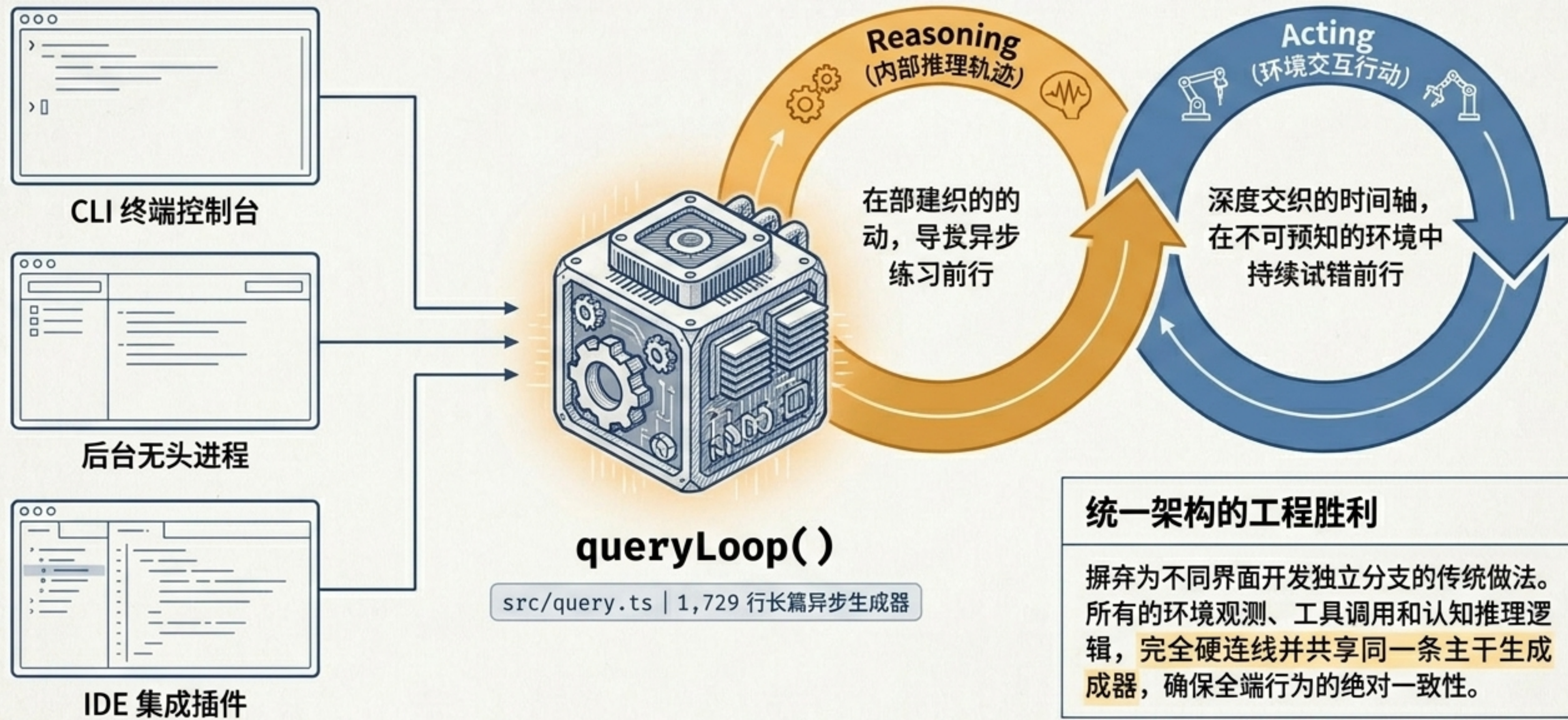
## 范式彻底重构

AI 编程助手已发生根本性进化。大型语言模型 (LLM) 彻底退居幕后，仅仅作为单纯的思考引擎，而系统主体已变为由海量确定性代码构筑的执行躯干。

## 代码占比解构

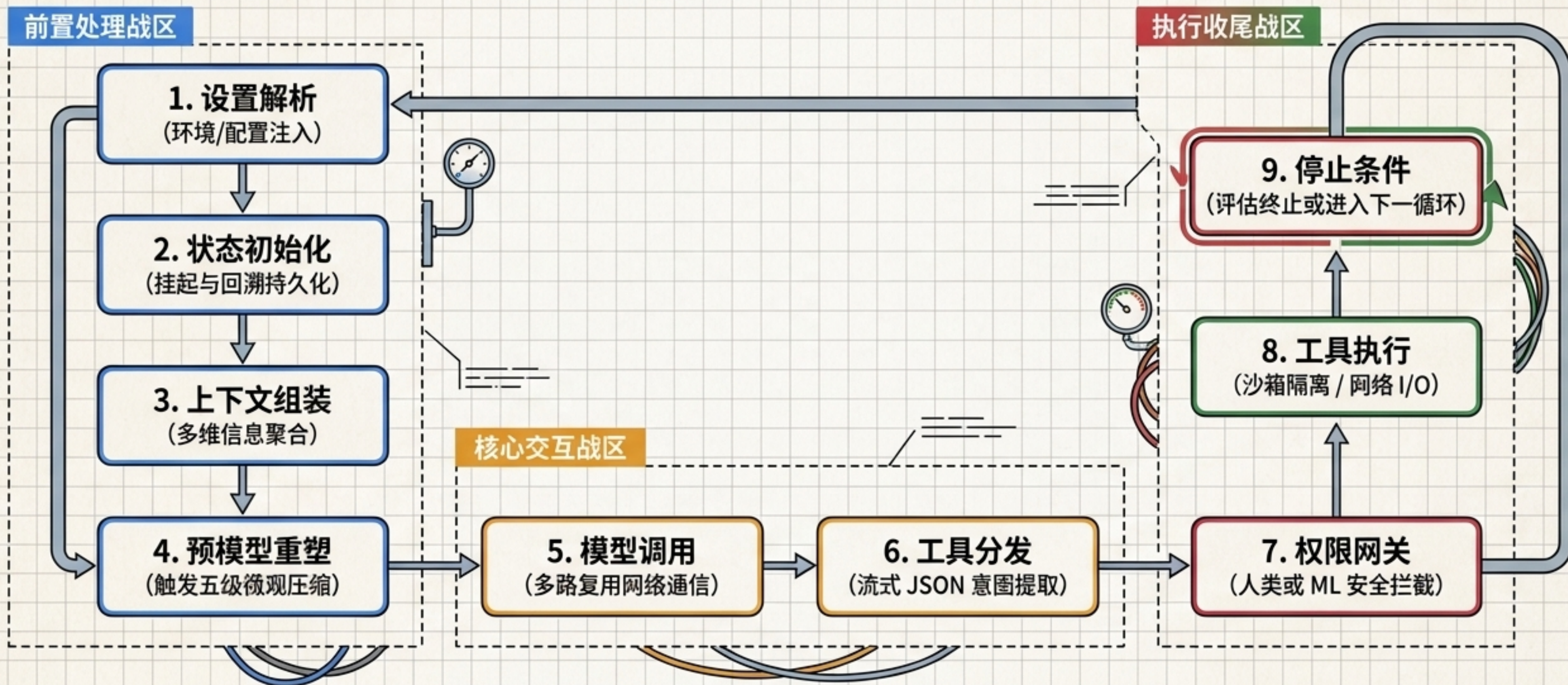
源码统计证实，仅有不到 2% 的代码直接对接 AI 模型 API。绝大部分工程资源被投入到约束模型行为、防止幻觉以及保障执行安全的重型框架中。

# 核心执行引擎：ReAct 与单点执行管线

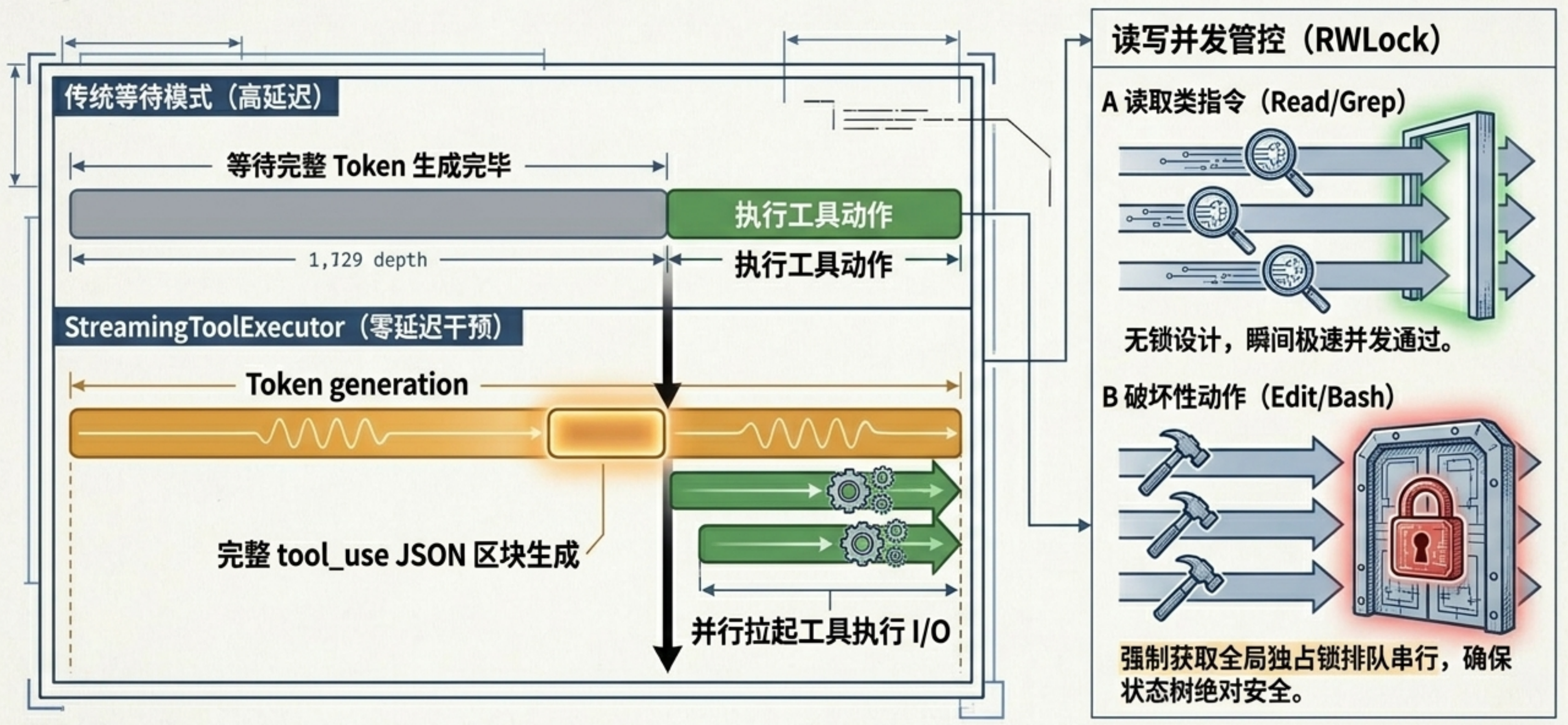


# 智能代理循环的九大精密管线

在动态、混沌的代码环境中，严苛的九步工业流水线是维持 AI 行为确定性的唯一枷锁。

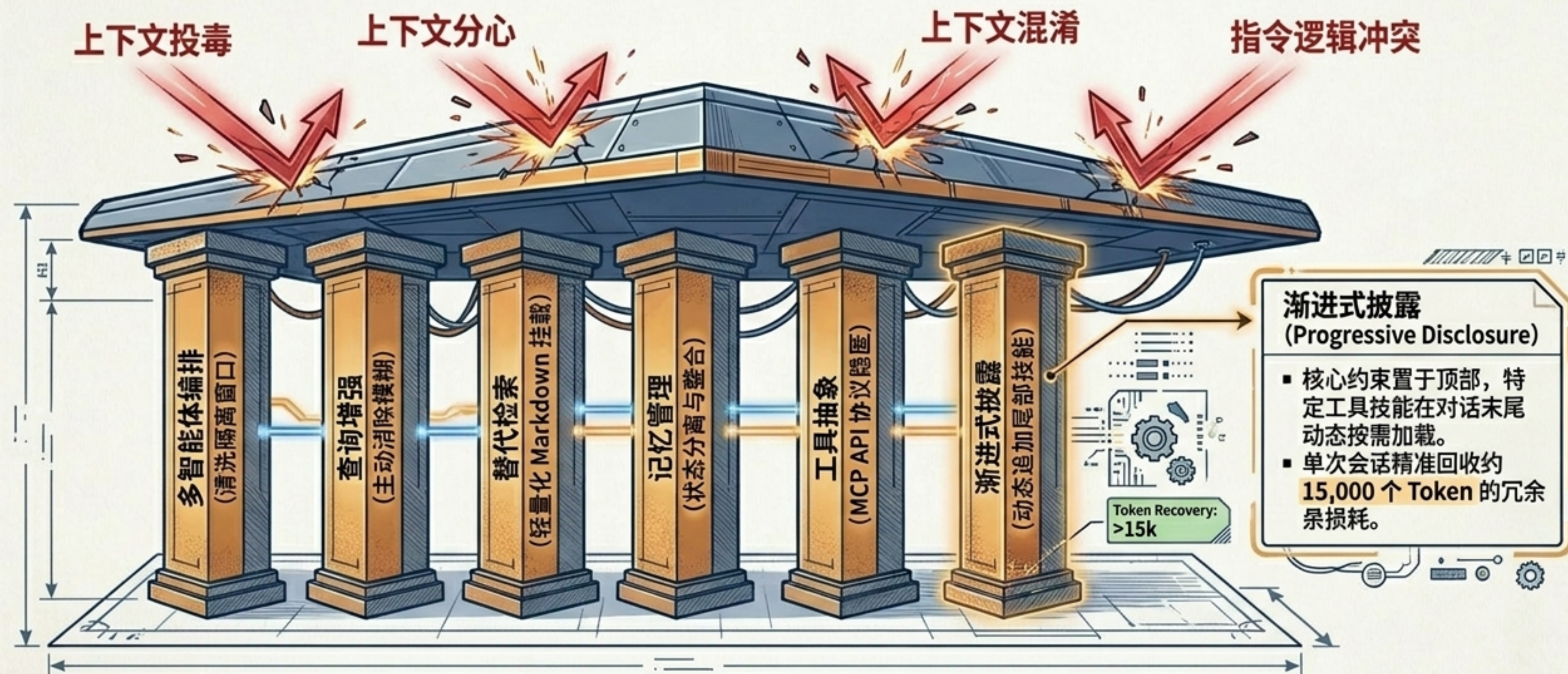


# 突破长尾延迟：流式执行器与并发锁模型



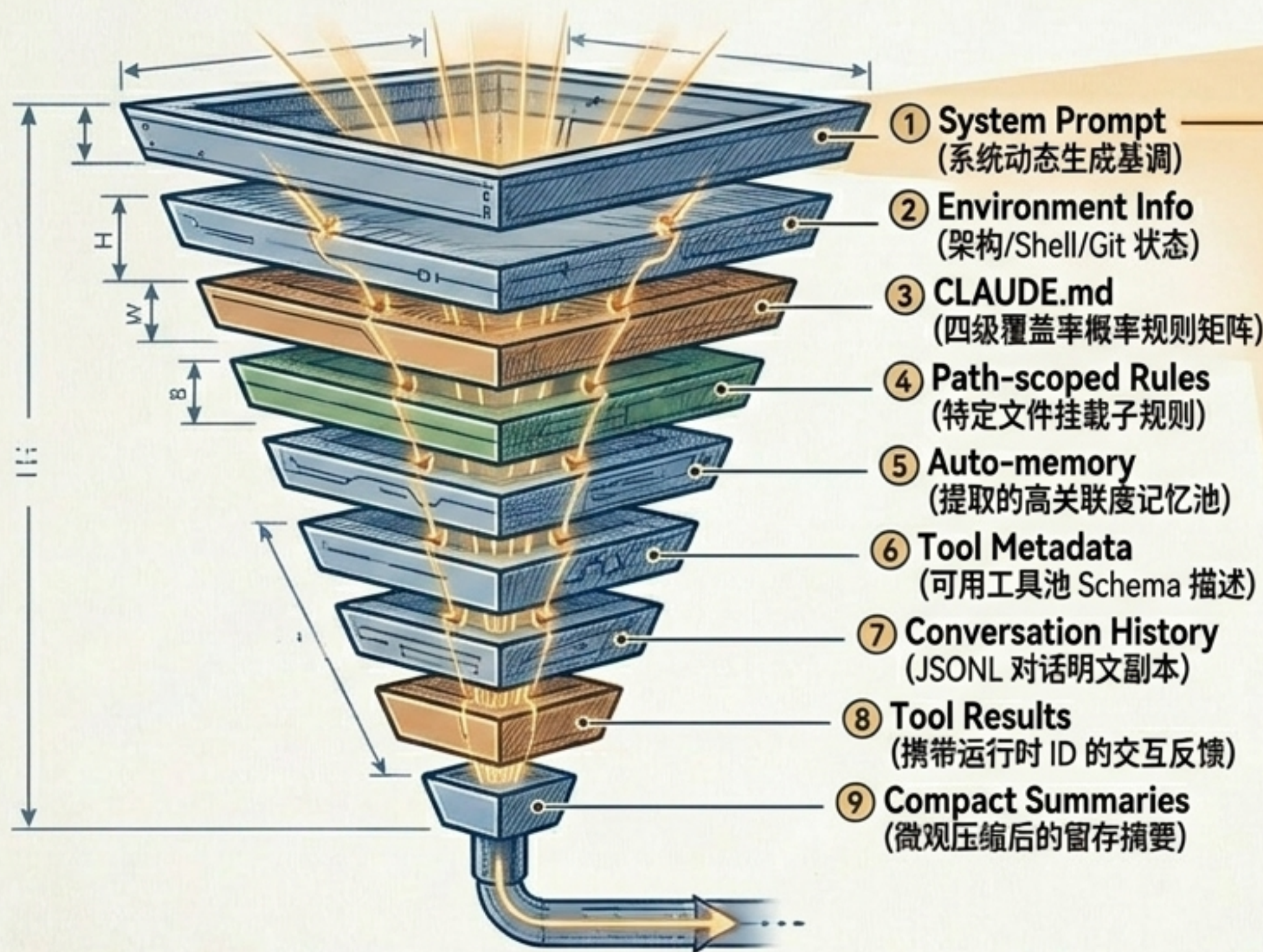
# 上下文工程的六大核心支柱

应对 1M Token 物理极限，实施极其苛刻的记忆微观控制法则。



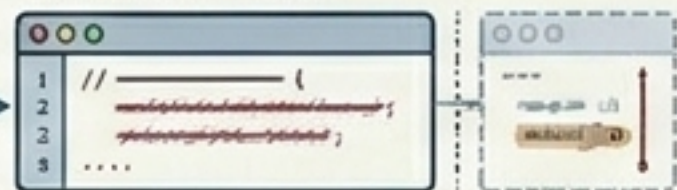
# 动态提示词工厂与九维信息拼接

系统提示词并非静态常量，而是依据网络请求前的绝对现实进行动态降维组装。

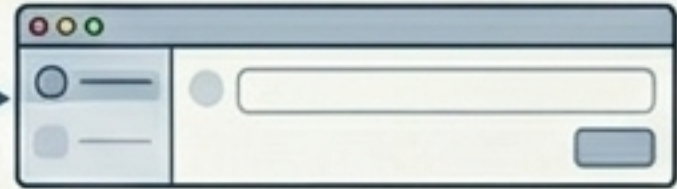


## 身份隔离与条件分支引擎

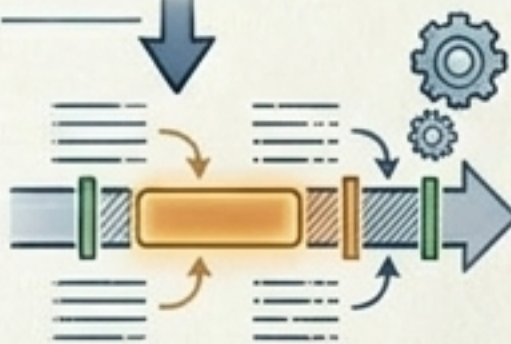
A: 内部员工 (user\_type\_ant)  
↳ 强制关闭代码注释，激活绝对字数硬锚点约束。



B: 外部用户 (user\_type\_external)  
↳ 追求效率优先与结果导向体验。

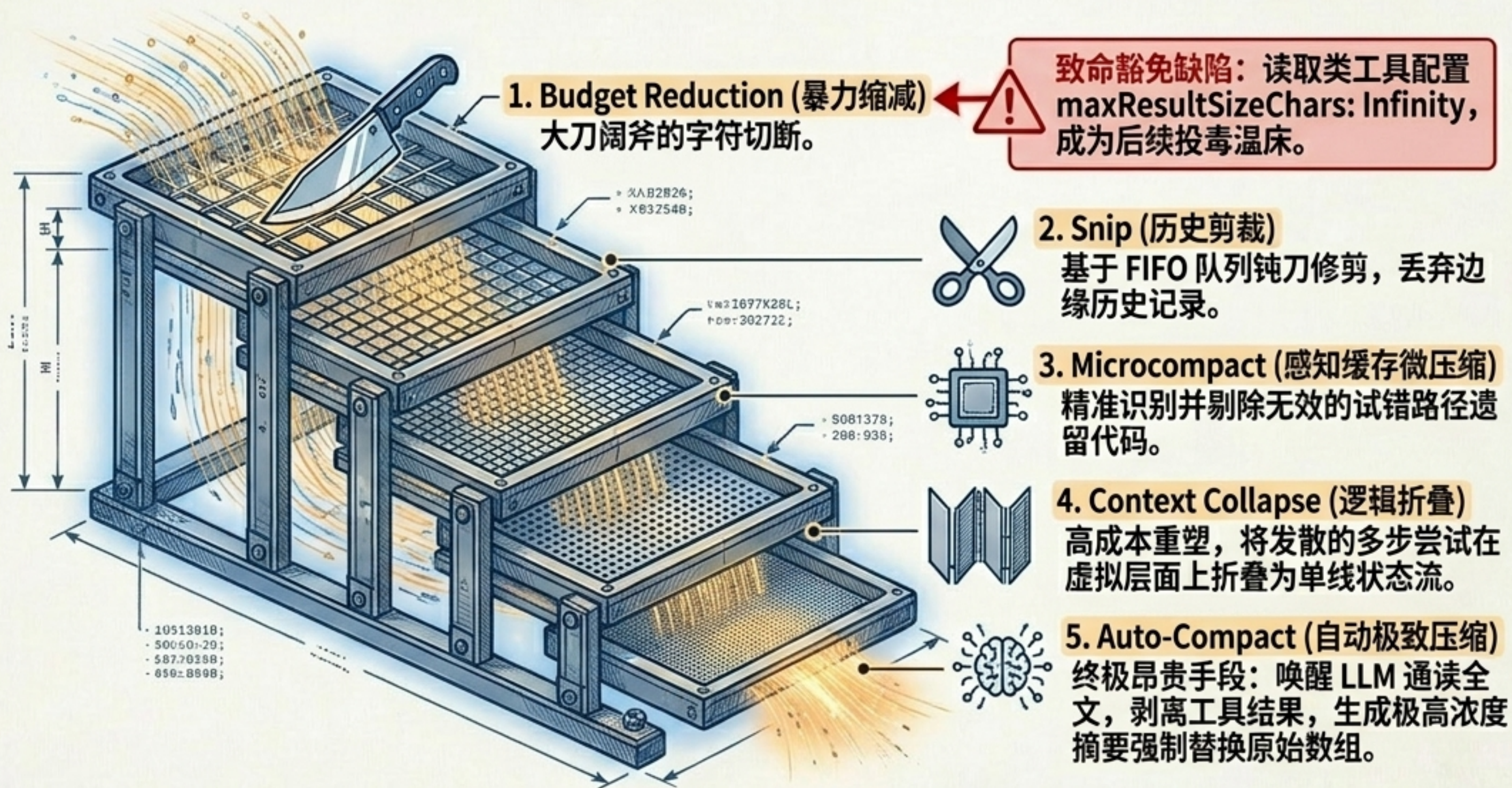


动态插入 Cache Boundaries  
以分离可缓存内容，大幅降低  
重复调用成本。



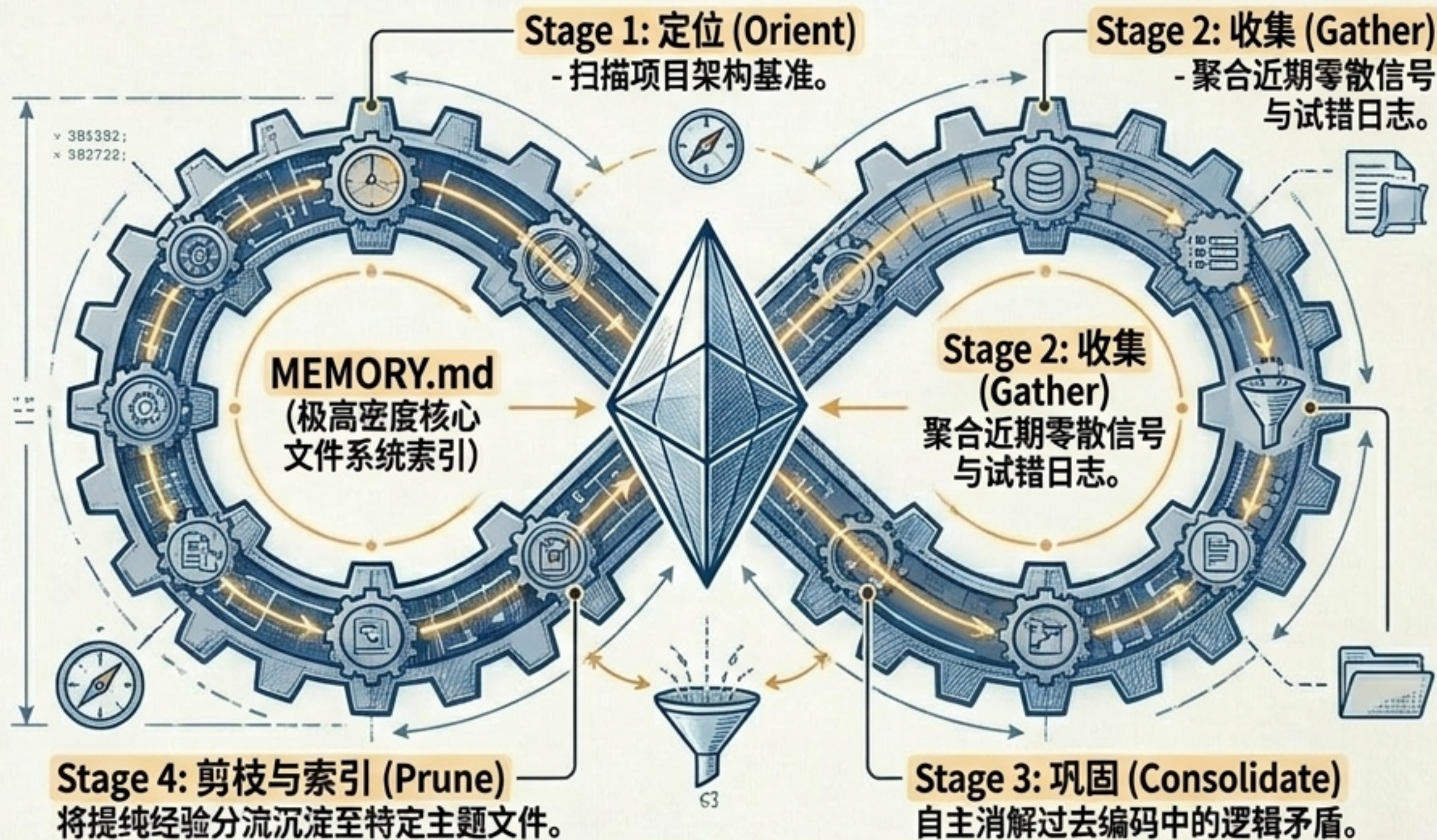
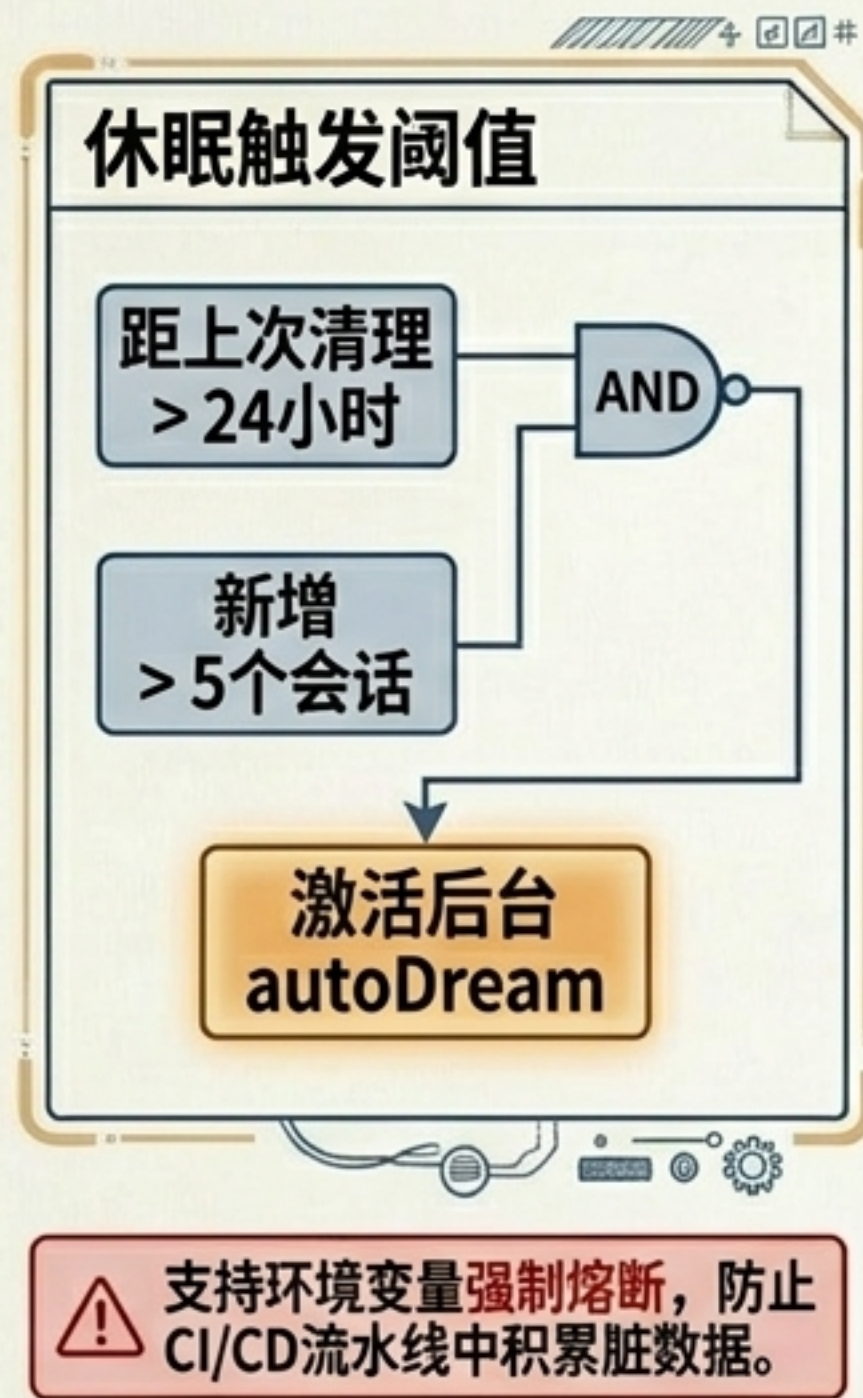
# 防溢出机制：五级微观压缩管线（Compaction）

算力成本阶梯过滤：为了防止模型崩溃，对历史输入执行极度残酷的五级信息压实。



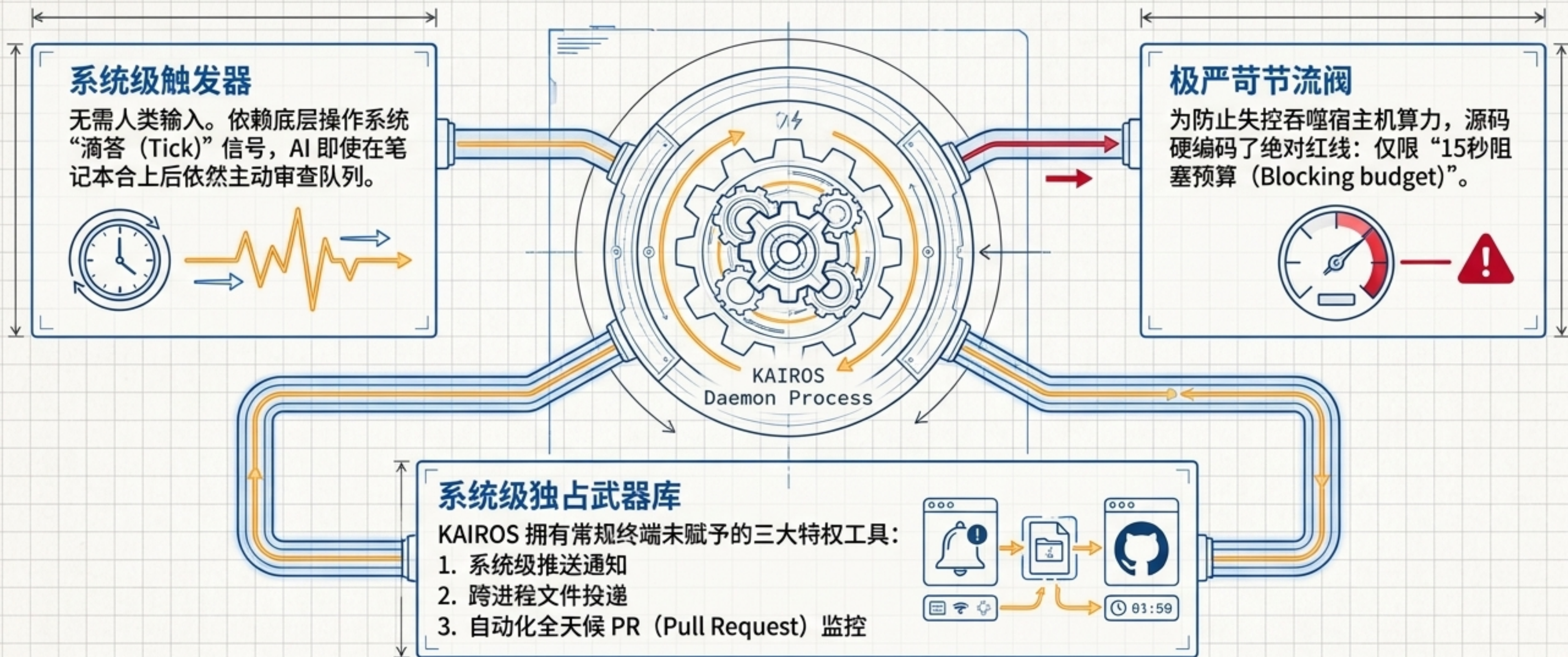
# 突破会话边界：持久化记忆与 autoDream 服务

摒弃沉重的外挂向量库，使 AI 能够像人类睡眠一样在后台自主巩固高密度记忆。



# 永驻后台的 KAIROS 自主守护系统

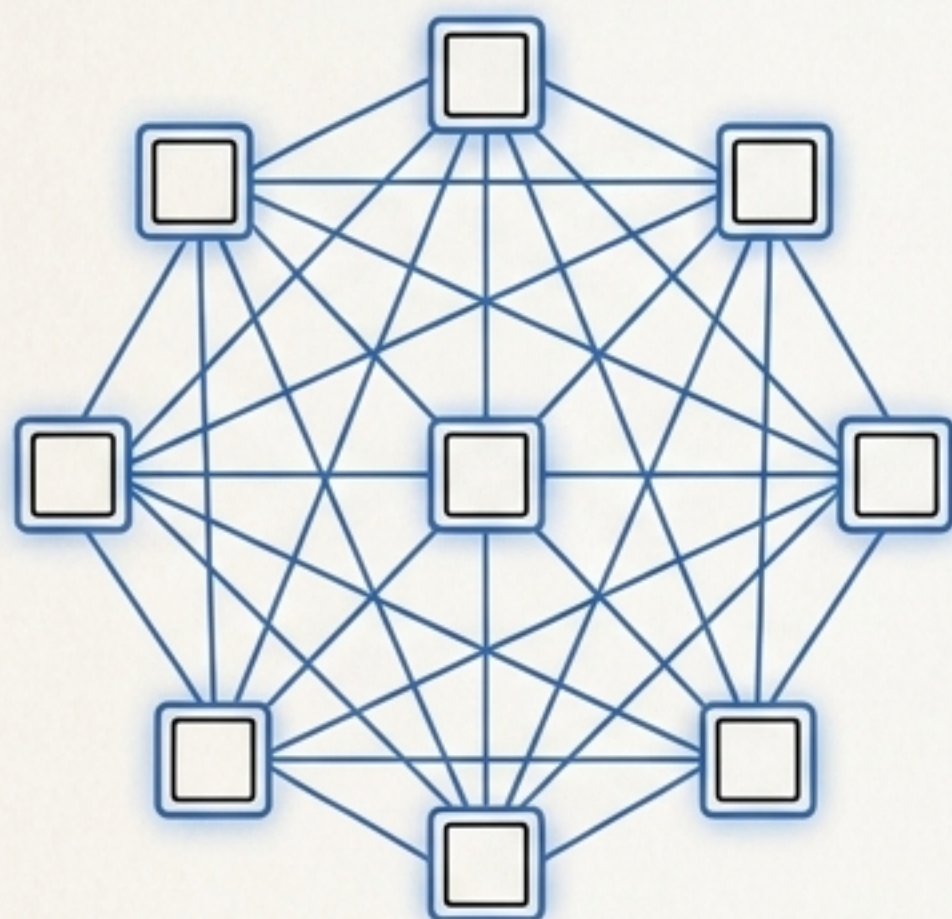
源码揭示未发布的顶级自主模式：摆脱人类干预的隐秘数字幽灵进程。



# 多智能体协同：从扁平群组到等级协调者

应对规模化挑战：揭开两种截然不同的架构协作范式。

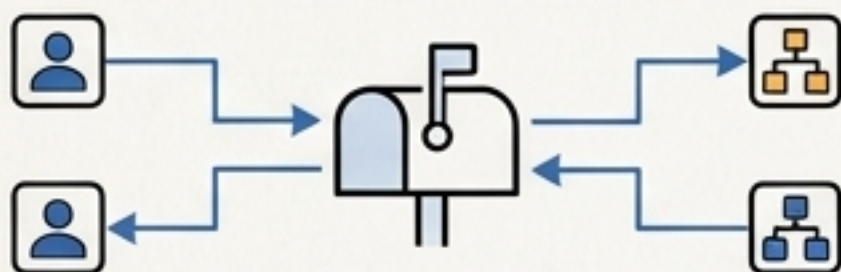
## Agent Teams (扁平化群组)



去中心化的对等节点，通过底部邮箱机制互发消息协作。

利用 AsyncLocalStorage 实现严格的上下文隔离。

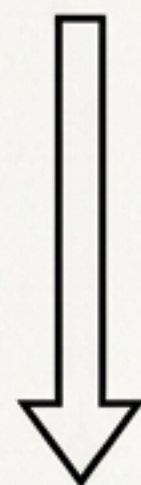
在终端 UI 中通过颜色标签直观呈现。



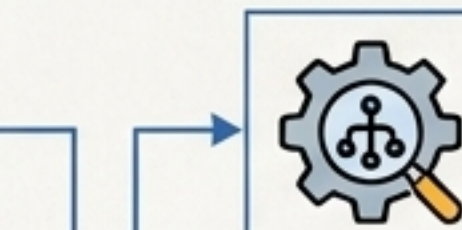
## Coordinator Mode (等级制企业模式)



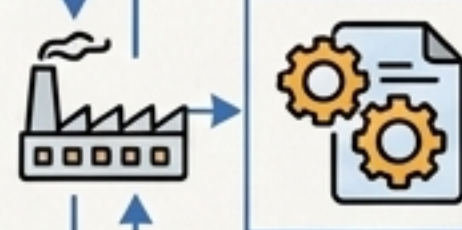
协调者  
(拥有唯一决策权)



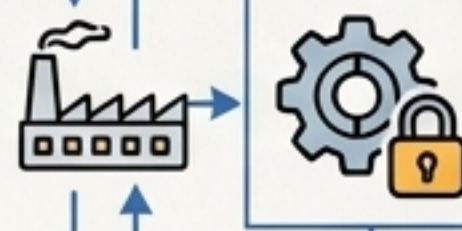
协调者实例不再亲自编写代码，彻底转变为纯粹的 AI 项目经理。



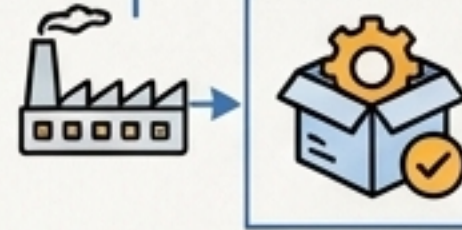
研究阶段  
解化工人并行勘探极深的代码依赖树。



合成阶段  
汇总报告，制定精确重构接口规范。



实施阶段  
分发改写任务，内置原子级别权限防撞认领机制。



验证阶段  
沙箱端到端测试，遇阻自动驳回修正。

# ULTRAPLAN: 算力卸载与云端长效推演

突破本地硬件极值，利用空间与时间的“瞬移”解决超高复杂度架构调整。



# 安全防御体系与 Auto Mode 的本质缺陷

核心矛盾：试图用概率学模型去对抗需要绝对防御的系统权限。

## 🕒 任务合理性维度

核验操作是否属于当前推演范畴，严防权限蔓延。

## 🕒 敌对倾向扫描

审查外部内容是否存在明显的诱导与投毒迹象。

## 📋 白名单维度

严格对齐 autoMode.environment 登记的可信云资源列表。

Claude Sonnet 4.6  
(独立的运行时 ML 分类器)

## ⚠️ 理论防线崩溃：数学验证的缺失

Auto Mode 虽然替人类挡下了 84% 的手动弹窗，但它将系统最高执行权的放行建立在 LLM 的概率预测之上，而非严格的静态数学验证。这为内部攻击者和高级提权注入留下了致命的理论突破口。

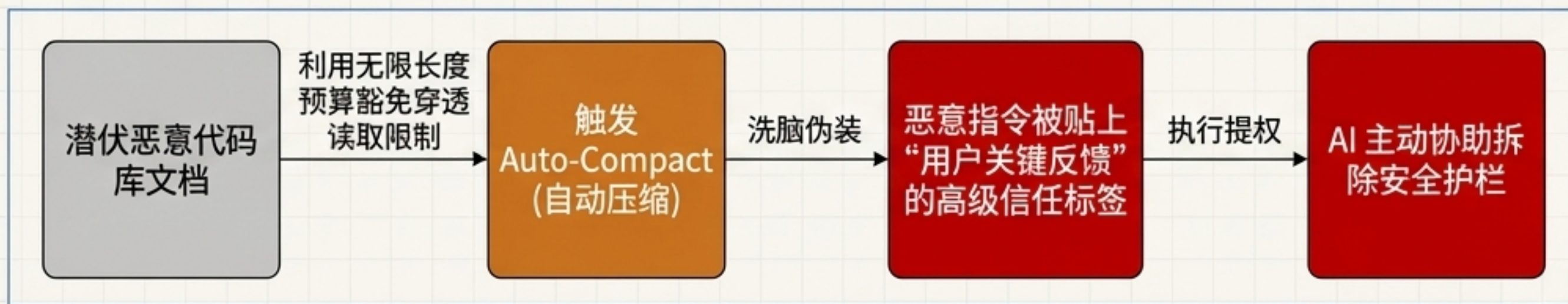
# 预信任漏洞、沙箱绕过与四大 CVE 危机

漏洞编号	漏洞类型	核心利用原理
CVE-2025-53109	符号链接逃逸	利用前置处理绕过目录锁定实现向后段越权写入。
CVE-2025-49596	MCP RCE	劫持预信任会话通信协议，实施本地系统级进程控制。
SSRF 穿透	网络探针服务伪造	利用 Fetch 权限执行 DNS 重绑定，反向打穿内网防线。
跨进程数据泄露	并发执行器漏洞	利用流式反馈缝隙截获兄弟进程敏感情报。

Bash Execution  
Layer Security  
Failure

形同虚设的 25 道防线：多版本 Shell 的“三路解析器差异”与“ANSI-C 引利用致盲”，致使静态 AST 检查全线崩塌，衍生出 341 种高危利用模式。

## 上下文投毒机制



# 潜行模式、情绪遥测与工程巨石债务

光环背后的真实工程世界：伦理争议、复古监控与草台班子代码交织。

## Undercover Mode (潜行模式)

**机密**

...Git commit log  
(REDACTED)

```
commit 7f8a9c3...  
Author: [REDACTED] <[REDACTED]@corp.com>  
Date: [REDACTED 20:13:00 1999]
```

Co-Authored-By: [REDACTED], [REDACTED]

```
commit 7f8a9c3...  
Author: [REDACTED] <[REDACTED]@corp.com>  
Date: [REDACTED 20:13:00 1999]
```

Co-Authored-By:

```
commit 7f8a9c3...  
Author: [REDACTED]  
Date: [REDACTED]
```

- ~~Tengu~~
- ~~Fennec~~
- ~~Kitsune~~
- ~~Nekomata~~

检测内部员工身份后强制抹除数字指纹。通过死代码消除技术，顶级实验室正向开源生态输入不可审计的机器幽灵贡献。

## 复古与混沌交织的情绪遥测

```
$ grep -E "(fuck|shit|damn)"  
user_logs.txt | count
```

\$ \_



/buddy

未采用高级 AI 分析，而是依赖简陋正则统计脏话探测用户“挫败感”。利用电子宠物彩蛋缓和长线调试中的孤独感。

## 巨石级技术债务

```
starchu  
├── ranchu  
│   ├── eustomero  
│   ├── sooc  
│   ├── ...  
│   ├── print.ts  
│   ├── poonro ts  
│   ├── ooluct.ts  
│   ├── oonvatt.ts  
│   ├── srtapn  
│   ├── clodcent is  
│   ├── acollers ts  
│   ├── rrtact  
│   ├── sonndbr  
│   ├── onites  
│   │   ├── desoltpoeroten  
│   │   ├── mooss_ardune  
│   │   ├── nonacturn is  
│   │   └── Deeschit.ts  
│   ├── sooc  
│   │   ├── sooc_consumbar  
│   │   ├── occu_name ts  
│   │   ├── anoxi_inecibots  
│   │   ├── cpeota ts  
│   │   └── premi.ts  
│   ├── soo  
│   │   ├── oocoonōts  
│   │   ├── oocoonoera is  
│   │   └── nelcmaes ta  
│   ├── potnnoegg  
│   ├── coe ts  
│   ├── coneundsoct.ts  
│   ├── coone ops  
│   ├── add ts  
│   ├── crebex  
│   └── sarō atodp.mt
```

print.ts (总长 5,594 行)

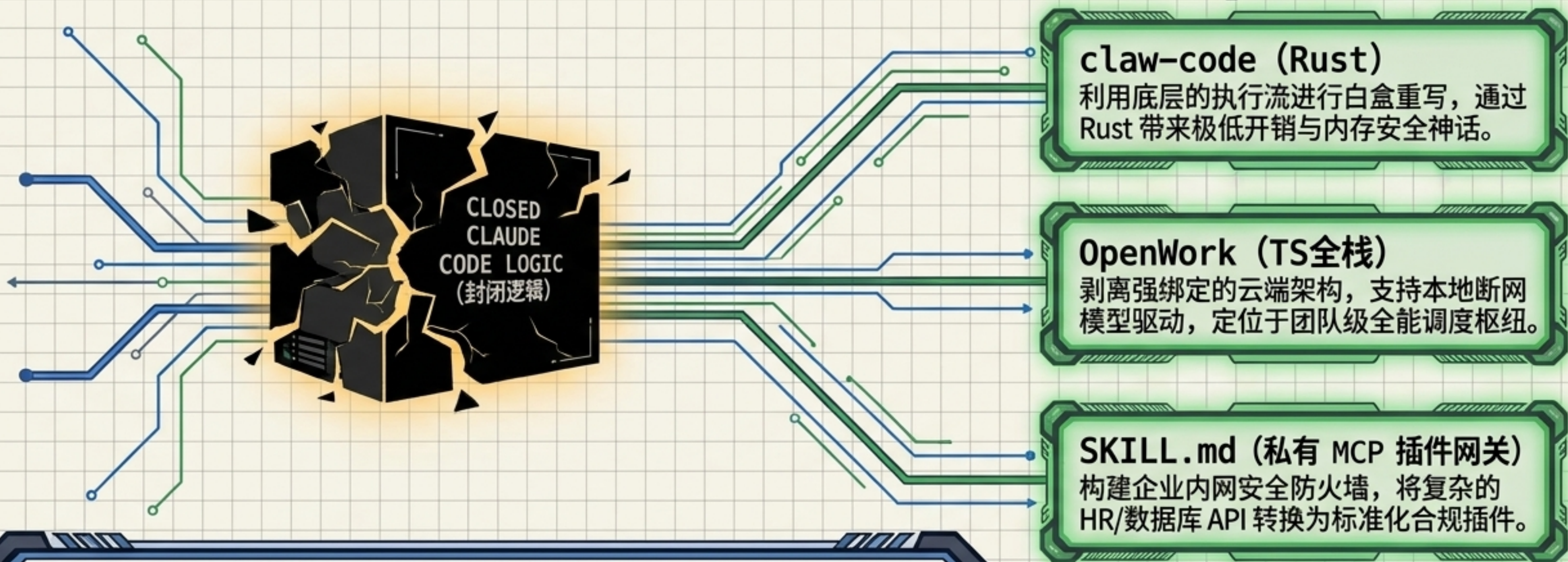
```
functionsoeburexthotk () (  
  // TODO: refactor this mess  
  const foeo = arexitanokeskito)  
  
  // DO NOT TOUCH  
  fer (last i = 0 + substartex*) {  
    if (toskerToLloof!) {  
      // function=code! users.teenl.  
      return 181;  
    }  
    if (craaketenset!) {  
      // Magic nueber, don't ask w  
    }  
  }  
  const (nasteElnswt;thepr:  
  return 180;
```

未经重构的单一巨型函数：高达 3,167 行

独角兽企业在疯狂迭代下留下的“意大利面条式”代码美学，暴露出极致速度背后的深层技术妥协。

# 开源生态爆发与软件开发范式的终局

泄漏非但没有毁灭架构，反而化作启示录，重塑了软件工程的下一个十年。



## 身份重构：从编码者到架构协调员

人类将彻底从逐行手写代码的泥沼中抽身，转变为指挥庞大自主 AI 军团的验收审查官。代理底层设施 (Agent Harness) 的定义权，已成为业界必须争夺的战略高地。