

# 什么是Harness Engineering?

## Harness Engineering 介绍

基于以下11篇原始资料整理：OpenAI、Anthropic、Thoughtworks、LangChain、HumanLayer、Inngest、学术论文（CAR框架）及 walkinglabs 课程资料。整理日期：2026-04-20。

### 一、定义

#### 核心公式

不同来源高度一致地使用同一个核心等式：

AI Agent = Model (模型) + Harness (挽具/环境)

**Harness** (中文常译作"挽具"或"智能体环境") 是围绕 AI 模型构建的一切——工具、指令、状态管理、验证机制、运行时基础设施——它让模型的智能变得"可用"。

"The model contains the intelligence and the harness is the system that makes that intelligence useful."

—— LangChain 《The Anatomy of an Agent Harness》

#### 各大厂对定义的定义对比

机构	定义侧重点	主要比喻/框架	与其他机构的差异
LangChain	最宽泛：Harness = 除模型以外的一切（文件系统、工具、沙箱、编排、运行时基础设施）。"Agent = Model + Harness" 这个等式最早由 LangChain 的 Vivek Trivedy 系统性提出。	工作引擎 (Work Engine)	定义范围最广，包括所有围绕模型的技术层
Anthropic	实用工程视角：Harness 是让 Agent 在多个上下文窗口中持续推进任务的一整套环境脚手架——包括初始化脚	环境脚手架 (Environment Scaffolding)	更聚焦于长期运行任务的连续性与状态管理，强调"clean state"理念

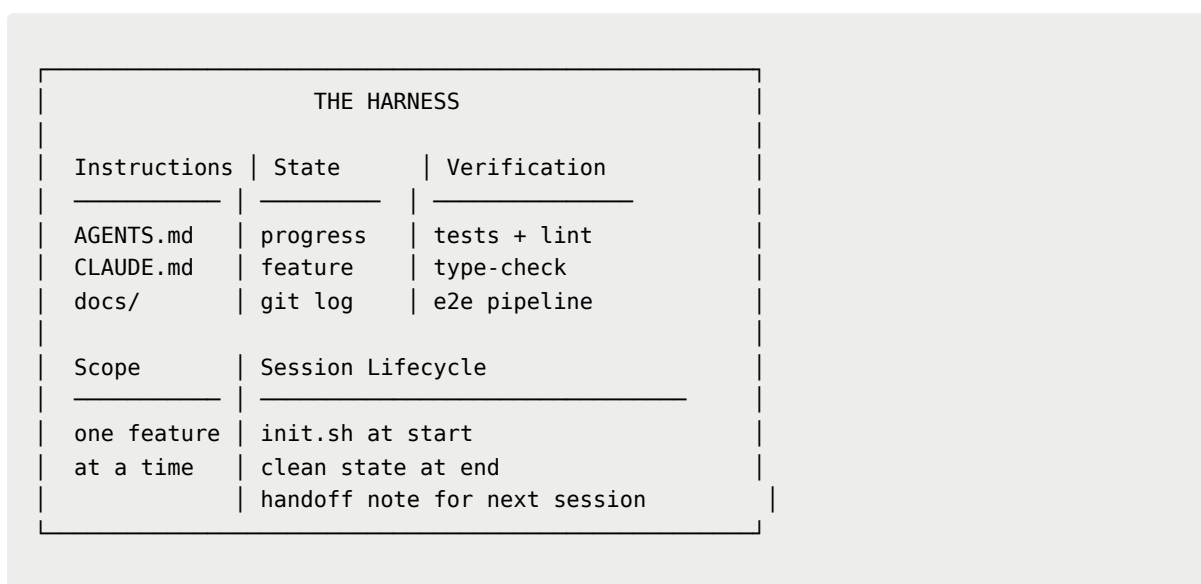
机构	定义侧重点	主要比喻/框架	与其他机构的差异
	本、功能列表、进度文件、会话移交机制。Claude Agent SDK 本身被描述为"通用智能体挽具"。		
<b>OpenAI</b>	<b>注意：</b> OpenAI 官方文章《Harness engineering: leveraging Codex in an agent-first world》中，"harness"这个词本身在正文只出现 <b>一次</b> （且是指 evals，而非 harness engineering 本身）。他们更多使用"环境设计（environment design）"、"反馈回路（feedback loops）"、"控制系统（control systems）"等术语。实质上，其 harness 指：结构化文档库、架构约束 linter、可观测性栈、以及"垃圾回收"机制。	代码仓库即知识系统（Repository as System of Record）	术语使用最不明确，与其他机构存在 <b>显著命名差异</b> ；强调"零人工代码"的极端自动化场景
<b>Thoughtworks</b>	在"限界 上下文"中定义：对于编码智能体用户，Harness 分为 <b>内层</b> （工具内置：系统提示、代码检索）和 <b>外层</b> （用户自建：规则、技能、静态分析、Review Agent）。将 Harness 比作赛博内廷（cybernetic governor）。	前馈导引（Feedforward Guides）+ 反馈传感（Feedback Sensors）	唯一明确区分"构建者挽具"与"用户挽具"两个层次；引入 <b>计算型（Computational）vs 推理型（Inferential）</b> 的控制分类
<b>HumanLayer</b>	将 Harness Engineering 定义为 <b>上下文工程（Context Engineering）的子集</b> ——专门通过编码智能体的配置接入点（AGENTS.md、MCP 服务器、技能、子智能体、钩子）来管理上下文窗口。	配置即杠杆（Configuration as Leverage）	强调 Harness Engineering 不等于 Context Engineering，而是其子集；最关注防止"上下文腐烂"
<b>Inngest</b>	基础设施视角：Harness 是"连接、保护、编排组件	布线挽具/安全带（Wiring	唯一将 Harness 等同于 <b>持久化事件驱动基础</b>

机构	定义侧重点	主要比喻/框架	与其他机构的差异
	——而不亲自执行任务"的层。LLM 是引擎，工具是外设，内存是存储，Harness 是将它们连接起来的基础设施（耐久执行、事件路由、状态持久化、并发控制）。	Harness / Safety Harness)	<b>设施</b> 的视角，来自 DevOps/云原生工程背景
<b>学术界 (CAR框架)</b>	将 Harness 层分解为三个维度： <b>Control (控制)</b> ——哪些指令保持权威； <b>Agency (智能体能力)</b> ——哪些行动可用； <b>Runtime (运行时)</b> ——状态如何延续、故障如何处理。提出"Harness-sensitive"概念：部分 Agent 性能提升可能来自 Harness 改进，而非模型本身。	CAR (Control-Agency- Runtime) 三元框架	唯一来自学术界的系统性框架；提出 HarnessCard 报告格式，类比模型卡 (Model Card)

**无差异之处：**所有来源都认同"Agent = Model + Harness"的核心等式，以及 Harness 的目标是让模型输出更可靠、更可控。

## 二、特点

### 1. 五大子系统 (walkinglabs 综合框架)



- **Instructions (指令)**: 告诉 Agent 做什么、按什么顺序、读什么文件。采用渐进式披露 (Progressive Disclosure), 而非一个巨型文件。
- **State (状态)**: 追踪已完成什么、正在进行什么、接下来是什么。持久化到磁盘, 确保会话间的连续性。
- **Verification (验证)**: 只有通过测试才算完成。Agent 不能在没有可运行证据的情况下宣告任务完成。
- **Scope (范围)**: 将 Agent 约束到每次一个功能。防止过度扩展和半途而废。
- **Session Lifecycle (会话生命周期)**: 开始时初始化, 结束时清理, 为下一次会话留下清晰的重启路径。

## 2. 计算型 vs 推理型控制 (Thoughtworks)

- **计算型 (Computational)**: 确定性、快速, 由 CPU 执行。测试、linter、类型检查、结构分析。可在每次变更时运行, 结果可靠。
- **推理型 (Inferential)**: 语义分析、AI 代码审查、LLM as Judge。由 GPU/NPU 运行, 速度慢、成本高、结果具有不确定性, 但能处理语义判断。

## 3. 前馈导引 + 反馈传感 (Thoughtworks)

- **前馈导引 (Feedforward Guides)**: 在 Agent 工作前注入上下文——AGENTS.md、技能文件、参考文档、如何引导脚本。提高首次生成正确的概率。
- **反馈传感 (Feedback Sensors)**: 在 Agent 工作后检测问题——静态分析、日志、浏览器测试、AI 代码审查。提供自我修正循环。

## 4. 三类调控维度 (Thoughtworks)

- **可维护性挽具 (Maintainability Harness)**: 调控代码内部质量——重复代码、圈复杂度、测试覆盖率、架构漂移。
- **架构适应性挽具 (Architecture Fitness Harness)**: 调控架构特征——性能要求、可观测性标准、依赖方向规则。
- **行为挽具 (Behaviour Harness)**: 调控功能正确性——规格说明、测试套件、端到端验证。(目前最难解决的维度)

# 三、优势

1. **显著提升任务完成质量**: Anthropic 的受控实验表明, 同一模型 (Opus 4.5)、同一提示 (构建2D复古游戏编辑器), 无 Harness 时花费 \$9/20 分钟, 产出无法运

行；有完整 Harness 时花费 \$200/6 小时，产出可以实际游玩的游戏。**模型没有变，Harness 变了。**

2. **突破单次上下文窗口限制**：通过 初始化 Agent + 会话移交文件 (progress.md、feature\_list.json、git log)，使 Agent 能跨多个上下文窗口持续工作数小时乃至数天。
3. **减少人工监督开销**：精心设计的 Harness 能在 Agent 输出到达人眼之前自动捕获并修正大量问题，降低审查负担、减少 token 浪费。
4. **工程化可复制**：将最佳实践编码为 Harness (AGENTS.md、验证脚本、结构测试)，可在团队中共享和复用，而不是每次项目重头来过。
5. **架构约束自动化执行**：通过自定义 linter 和结构测试，将架构规范 (如依赖方向、命名规范、文件大小限制) 自动化执行，"一旦编码，处处生效" (OpenAI)。
6. **防止上下文腐烂 (Context Rot)**：子智能体作为"上下文防火墙"，工具结果截断、进度压缩等机制防止上下文窗口随时间降级，保持 Agent 在"智能区间"内工作。
7. **提高对 Agent 产出的信任**：在代码进入人眼之前完成静态分析、测试运行、架构检查，为开发者提供质量保证。

## 四、劣势

1. **前期构建成本高**：设计有效的 Harness 需要投入大量精力——编写 AGENTS.md、init.sh、feature\_list.json、验证脚本、结构测试等。初期会比直接提示 Agent 更慢。
2. **运行成本显著增加**：多 Agent 架构 (规划者 + 生成者 + 评估者) 可能使成本提升 20 倍以上 (Anthropic 实验：单 Agent \$9 vs 完整 Harness \$200)。
3. **维护与漂移**：Harness 文件 (尤其是 AGENTS.md) 会过时。"单一巨型文件"方式会成为陈旧规则的墓地，需要持续的"文档园艺 (doc-gardening)" 维护 (OpenAI)。
4. **模型过拟合风险**：模型会在其训练所用的 Harness 上过度拟合。Codex 模型高度依赖 `apply_patch` 工具，在不同 Harness 下性能下降。Terminal Bench 2.0 数据显示，Opus 4.6 在 Claude Code 内部排名第33，但在不同 Harness 中排名第5。 **最优 Harness 不一定是模型训练所用的 Harness。**
5. **行为正确性仍难保证**：Harness 善于捕获结构性和维护性问题，但对"功能是否符合用户意图"的保证仍然有限。依赖 AI 生成的测试套件存在误判风险

(Thoughtworks 指出这是目前最大的悬而未决问题)。

6. **复杂度膨胀陷阱**: 过度工程化的 Harness 会让开发者花更多时间调优配置而非交付功能。HumanLayer 建议: "偏向交付; 只在 Agent 实际失败后才添加配置。"
7. **可测性差**: 如何评估 Harness 本身的质量? 目前还缺少类似代码覆盖率或变异测试的 Harness 覆盖率指标 (Thoughtworks 提出的开放问题)。

## 五、适用场景

### 高度适用

场景	原因
长期运行的编码任务 (数小时乃至数天)	解决跨上下文窗口连续性问题; 防止 Agent 过早宣告任务完成
企业级/遗留代码库中的自动化编码	需要架构约束执行、上下文管理、风格一致性保证
从零到一构建产品 (无人工代码)	OpenAI 实验: 3个工程师 × 5个月 × 约100万行代码, 吞吐量随团队规模增长而提升
反复出现相同类型失败的 Agent  workflows	"每当 Agent 犯错, 就将修复编入 Harness, 让这个错误永远不再发生" (Mitchell Hashimoto)
多 Agent 协作系统	规划者 + 生成者 + 评估者架构, 各 Agent 职责隔离
需要高可靠性的生产环境	通过验证回路、钩子、结构测试减少人工监督需求

### 不适用或性价比低

场景	原因
一次性简单任务	Harness 的前期成本超过其带来的价值
高度不确定或创意性任务 (主观设计审美等)	行为 Harness 目前对功能正确性保证有限
模型能力覆盖的任务 (随模型迭代)	更强的模型会让部分 Harness 组件变得冗余——需要定期重新评估哪些组件仍是必要的
团队不具备工程化能力编写验证脚本和结构测试	没有可运行验证的 Harness 只是一堆过时文档

## 六、参考来源

来源	链接
OpenAI	<a href="#"><u>Harness engineering: leveraging Codex in an agent-first world</u></a>
Anthropic	<a href="#"><u>Effective harnesses for long-running agents</u></a>
Anthropic	<a href="#"><u>Harness design for long-running application development</u></a>
Anthropic	<a href="#"><u>Building Effective AI Agents</u></a>
Thoughtworks (martinfowler.com)	<a href="#"><u>Harness engineering for coding agent users</u></a>
LangChain	<a href="#"><u>The Anatomy of an Agent Harness</u></a>
HumanLayer	<a href="#"><u>Skill Issue: Harness Engineering for Coding Agents</u></a>
Inngest	<a href="#"><u>Your Agent Needs a Harness, Not a Framework</u></a>
学术论文 (CAR框架)	<a href="#"><u>Harness Engineering for Language Agents: The Harness Layer as Control, Agency, and Runtime</u></a>
walkinglabs	<a href="#"><u>Learn Harness Engineering</u></a>
walkinglabs	<a href="#"><u>Awesome Harness Engineering</u></a>