

# AI 行业日报

模型与工具能力 · 产业基础设施 · 应用商业化 · 研究开源 · 资本监管 | Data Source: aihot.virxact.com

API耗时: 0s 精选条目: 31 条 焦点: 8 条 快讯: 0 条

## Executive Summary

今日AI行业呈现显著的技术突破与商业化加速趋势。BitCPM-CANN作为全球首个基于华为昇腾910B NPU全栈训练的1.58比特开源大模型正式发布，采用极低比特量化技术使内存占用相比BF16降低约6倍，为边缘端部署奠定基础。谷歌I/O大会系统性构建AI代理开发工具链，推出Antigravity 2.0桌面应用及SDK，Gemini API全面升级。NVIDIA发布Nemotron-Labs扩散语言模型技术，旨在实现光速级文本生成效率。StepAudio 2.5实时语音模型具备副语言感知能力，支持个性化交互体验。

行业结构层面，黄仁勋预测超大规模云厂商AI基建年度开支将从当前1万亿美元增长至3-4万亿美元，反映基础设施投资强度超预期。Anthropic即将完成超300亿美元融资，估值有望突破9000亿美元反超OpenAI，资本向头部企业集中趋势明显。微软报告显示特定场景下AI部署成本已超过人工工资，揭示企业应用AI的经济挑战。GitHub连续第三年获Gartner企业级AI编程代理领域领导者评级，Warp支持OpenRouter接入，显示开发工具生态持续完善。

后续需重点关注华为昇腾生态的模型优化进展与商业化节奏，扩散架构模型的实际部署效果与性能表现，以及AI代理开发工具链的市场接受度。基础设施投资规模扩大将影响算力供给格局，而AI应用成本效益比将成为企业决策的关键指标。监管政策对开源模型与数据安全的要求，以及多模态交互技术的标准化进程，将是行业发展的重要变量。

## 重点 今日核心进展

### ★ 1. 首个基于华为昇腾910B NPU全栈训练的1.58比特开源大模型BitCPM-CANN发布

X: Rohan Paul (@rohanpaul\_ai) · 昨天 22:36 · 模型与工具能力

ModelBest、清华大学与OpenBMB社区联合发布了BitCPM-CANN，这是全球首个完全基于华为昇腾910B NPU训练的开源1.58比特三元大模型。其核心创新在于采用仅含三种权重状态的极低比特量化技术，使模型内存占用相比BF16降低约6倍，可高效部署于手机、电脑、车载设备等边缘端。更关键的是，整个训练全栈（从量化算子到框架）均在昇腾上原生构建与验证，而非简单移植。该模型家族（0.5B-

能力进展 基础设施 新发布

[https://x.com/rohanpaul\\_ai/status/2057833050692800926](https://x.com/rohanpaul_ai/status/2057833050692800926)

### ★ 2. v2.1.149 更新摘要

Claude Code: GitHub Releases (RSS) · 昨天 06:09 · 应用与商业化

本次 v2.1.149 更新包含功能增强、企业设置和多项修复。新增 `/usage` 命令的使用量分类显示功能，可区分技能、子代理、插件及每个 MCP 服务器的消耗；`/diff` 详情视图支持键盘滚动；Markdown 输出兼容 GFM 任务列表。企业版新增 `allowAllClaudeAiMcp` 设置以加载云 MCP 连接器。修复了 PowerShell 权限绕过、Git 工作树沙盒写入

能力进展 基础设施 监管/资本

<https://github.com/anthropics/claude-code/releases/tag/v2.1.149>

### ★ 3. 谷歌I/O大会发布AI代理全套开发工具链

X: Google AI (@GoogleAI) · 昨天 01:09 · 应用与商业化

谷歌在I/O开发者大会宣布，系统性构建面向AI代理（Agent）的开发与部署工具链。核心更新包括：独立桌面应用Antigravity 2.0及其命令行工具、SDK面世；Google AI Studio新增Kotlin支持，可一键开发安卓应用并发布，同时推出移动端App。此外，Gemini API推出托管代理服务，实现一键部署；WebMCP作为开放标准在Chrome 149中推出，允许网页向代理暴露

能力进展 基础设施 新发布

<https://x.com/GoogleAI/status/2057871583843135978>

### ★ 4. Warp现已支持OpenRouter接入

X: OpenRouter (@OpenRouter) · 昨天 01:25 · 应用与商业化

OpenRouter现已在@warpdotdev中得到支持！♥ 工程师Dagm Assefa展示了如何连接DeepSeek和OpenRouter。文档：<https://docs.warp.dev/agent-platform/inference/custom-inference-endpoint/> ☑

能力进展 基础设施 新发布

<https://x.com/OpenRouter/status/2057875517391667492>

## ★ 5. GitHub 连续第三年被 Gartner® 魔力象限TM 评为企业级 AI 编程代理领域的领导者

GitHub Blog · 昨天 00:10 · 产业与基础设施

Gartner 最新发布的魔力象限报告中，GitHub 连续第三年被列为"领导者"象限，该评估专注于企业级 AI 编程代理领域。GitHub 表示，其致力于构建一个开放、安全且由 AI 驱动的平台，以赋能每一位开发者并定义软件开发的未来。此次评选进一步巩固了 GitHub 在 AI 辅助开发工具市场的领先地位。

能力进展 监管/资本 新发布

<https://github.blog/ai-and-ml/github-copilot/github-recognized-as-a-leader-in-the-gartner-magic-quadrant-for-enterprise-ai-coding-agents-for-the-third-year-in-a-row>

## ★ 6. Nemetron-Labs 扩散语言模型实现光速级文本生成

Hugging Face: Blog (RSS) · 昨天 08:02 · 研究与开源进展

NVIDIA 在 Hugging Face 发布了关于 Nemetron-Labs 扩散语言模型的技术博客。该研究聚焦于通过扩散语言模型架构大幅提升文本生成速度，目标是逼近"光速级"生成效率。文章可能介绍了该模型在生成速度上的突破，以及相较于传统自回归模型在延迟和吞吐量方面的性能优势。具体技术细节或对比数据需参考原文。

能力进展 基础设施 新发布

<https://huggingface.co/blog/nvidia/nemetron-labs-diffusion>

## ★ 7. 微软称，使用人工智能的成本高于支付人工工资

Hacker News 热门 (buzzing.cc 中文翻译) · 18 小时前 · 产业与基础设施

微软发布报告指出，在特定工作场景中，部署和使用人工智能 (AI) 的成本目前已高于支付相应的人工工资。报告分析了基于"tokens" (令牌) 和"agents" (智能体) 的 AI 使用模式，发现其综合开销超过了雇佣人类员工完成同类任务的费用。这一发现揭示了当前企业应用 AI 技术面临的现实经济挑战。

能力进展 新发布

<https://fortune.com/2026/05/22/microsoft-ai-cost-problem-tokens-agents>

## ★ 8. StepAudio 2.5 实时语音发布：副语言感知与人格化交互

X: 阶跃星辰 StepFun (@StepFun\_ai) · 2 小时前 · 应用与商业化

StepAudio 2.5 Realtime 是一款实时语音模型，能够深度理解用户语音中的语气、语速、停顿乃至微表情等副语言特征。它支持通过 API 接入自定义人格，允许设定个性、背景故事和语言风格，并提供了上万种原生人格选项，可组合出数百万种特征。产品还内置了 5 个可直接体验的预设人格，并经过 RLHF 调优，确保在复杂的角色扮演压力测试中也能保持角色一致性。该模型支持中文和英文。

能力进展 新发布

[https://x.com/StepFun\\_ai/status/2058303294544425197](https://x.com/StepFun_ai/status/2058303294544425197)

## 能力 模型与工具能力

### 1. 首个基于华为昇腾910B NPU全栈训练的1.58比特开源大模型BitCPM-CANN发布

X: Rohan Paul (@rohanpaul\_ai) · 昨天 22:36

ModelBest、清华大学与 OpenBMB 社区联合发布了 BitCPM-CANN，这是全球首个完全基于华为昇腾910B NPU 训练的开源 1.58 比特三元大模型。其核心创新在于采用仅含三种权重状态的极低比特量化技术，使模型内存占用相比 BF16 降低约 6 倍，可高效部署于手机、电脑、车载设备等边缘端。更关键的是，整个训练全栈（从量化算子到框架）均在昇腾上原生构建与验证，而非简单移植。该模型家族（0.5B-

能力进展 基础设施 新发布

[https://x.com/rohanpaul\\_ai/status/2057833050692800926](https://x.com/rohanpaul_ai/status/2057833050692800926)

## 产业 产业与基础设施

### 1. GitHub 连续第三年被 Gartner® 魔力象限TM 评为企业级 AI 编程代理领域的领导者

GitHub Blog · 昨天 00:10

Gartner 最新发布的魔力象限报告中，GitHub 连续第三年被列为"领导者"象限，该评估专注于企业级 AI 编程代理领域。GitHub 表示，其致力于构建一个开放、安全且由 AI 驱动的平台，以赋能每一位开发者并定义软件开发的未来。此次评选进一步巩固了 GitHub 在 AI 辅助开发工具市场的领先地位。

能力进展 监管/资本 新发布

<https://github.blog/ai-and-ml/github-copilot/github-recognized-as-a-leader-in-the-gartner-magic-quadrant-for-enterprise-ai-coding-agents-for-the-third-year-in-a-row>

### 2. 微软称，使用人工智能的成本高于支付人工工资

Hacker News 热门 (buzzing.cc 中文翻译) · 18 小时前

微软发布报告指出，在特定工作场景中，部署和使用人工智能 (AI) 的成本目前已高于支付相应的人工工资。报告分析了基于"tokens" (令牌) 和"agents" (智能体) 的 AI 使用模式，发现其综合开销超过了雇佣人类员工完成同类任务的费用。这一发现揭示了当前企业应用 AI 技术面临的现实经济挑战。

能力进展 新发布

<https://fortune.com/2026/05/22/microsoft-ai-cost-problem-tokens-agents>

### 3. 黄仁勋：AI 基建年度开支要冲到 4 万亿美元！

IT之家 (RSS) · 昨天 06:30

英伟达发布 2027 财年 Q1 财报，营收 816 亿美元，同比增长 85%，净利润 583 亿美元，翻两倍多，市值达 5.7 万亿美元，已超德国 2026 年 GDP 预测。黄仁勋预测，超大规模云厂商的 AI 基建年度开支将从当前的 1 万亿美元，增长至 3-4 万亿美元，远超华尔街预期。财报同时显示，数据中心业务营收 752 亿美元，占比超九成。值得注意的是，AI 基建的高能耗正推高居民电费，数据中心用电成本转嫁效应已初步显现。

基础设施 新发布

<https://www.ithome.com/0/954/223.htm>

#### 4. 消息称 Anthropic 最快下周完成逾 300 亿美元融资，有望推动估值反超 OpenAI

IT之家 (RSS) · 9 小时前

据彭博社报道，Anthropic即将完成一轮超300亿美元的融资，最快可能于下周敲定。此轮融资将使其估值突破9000亿美元，正式超越OpenAI，成为全球估值最高的AI初创企业。融资的迅速推进反映了市场的强烈追捧。同时，公司营收高速增长，预计第二季度营收将达109亿美元，环比增长超一倍，有望迎来首个盈利季度。

监管/资本 新发布

<https://www.ithome.com/0/954/452.htm>

#### 5. 加倍投入科学以赢得工业AI

Mistral AI: News (网页) · 14 小时前

Mistral AI宣布与物理AI先驱Emmi AI达成最终收购协议，旨在加强其在工业AI领域的领导地位。通过整合Emmi AI在物理仿真与数字孪生方面的专长，Mistral AI将提升其工程解决方案能力，并加速科学研究路线。Emmi AI的30余名研究员与工程师将加入Mistral AI团队，共同构建由物理AI驱动的综合技术栈。此次合作将为航空航天、汽车等高风险行业提供实时仿真与复杂问题解决平台

监管/资本

<https://mistral.ai/news/science-to-win-industrial-ai>

#### 6. 扩大与新加坡合作，推动AI安全规模化部署

X: Google DeepMind (@GoogleDeepMind) · 23 小时前

我们正在扩大与新加坡的合作，以帮助安全地大规模部署AI。与各国专家合作，我们的新项目将重点加速科学发现、加强大流行病防范并改善医疗保健。了解更多 → <https://goo.gl/49jGwjv>

监管/资本

<https://x.com/GoogleDeepMind/status/2057985225100235022>

#### 7. Kling AI亮相戛纳，推动AI赋能电影制作

X: 可灵 Kling AI (@Kling\_ai) · 21 小时前

Kling AI在戛纳电影市场 (Marché du Film) 举办官方会议，首次登上这一世界顶级电影舞台。会议汇集全球电影专业人士，共同探讨AI如何融入实际电影制作流程。Kling AI已证明其能力可服务于动画长片、好莱坞剧集、实验短片及影院电影等多种创作形式。未来，Kling AI将继续推进电影级AI影像技术，与全球创作者合作，将更多“不可能”的故事呈现在银幕上。

[https://x.com/Kling\\_ai/status/2058013861404684739](https://x.com/Kling_ai/status/2058013861404684739)

#### 8. AI 替代入门级工作：科技行业受裁员冲击最重，74% CEO 冻结或缩减招聘

IT之家 (RSS) · 昨天 08:05

奥纬咨询研究发现，AI工具正被广泛用于入门级任务，导致企业招聘重心转向高级岗位，年轻人求职难度加大。科技行业受冲击最严重，74%的CEO已冻结或缩减招聘。计划削减初级岗位的比例从17%跃升至43%，而招聘转向中层岗位的比例则升至30%。尽管超90%的企业在部署AI，但多数仍处试点阶段。报告警告，过快裁员或忽视初级人才储备，可能对人才梯队造成长远风险。

<https://www.ithome.com/0/954/235.htm>

#### 9. 回顾Google I/O 2026对话环节

Google Blog: AI (RSS) · 昨天 02:00

在2026年Google I/O开发者大会上，对话环节聚焦于未来科技趋势。行业领导者围绕人工智能、量子计算、机器人学以及创造力等核心议题展开了深入探讨，旨在勾勒这些前沿领域的技术演进路径与发展蓝图。

<https://blog.google/innovation-and-ai/technology/ai/io-2026-dialogues-recap>

#### 10. Suno AI创作夏日神曲《波多黎各》爆火

X: Suno (@suno) · 昨天 00:17

今年夏天的热门歌曲是用 Suno 制作的。非常感谢 @GMA 让这首病毒式传播的《Puerto Rico》歌曲被更多人看到！还有谁的脑海里也一直回响着这首歌？

<https://x.com/suno/status/2057858423664894196>

### 应用 应用与商业化

#### 1. v2.1.149 更新摘要

Claude Code: GitHub Releases (RSS) · 昨天 06:09

本次 v2.1.149 更新包含功能增强、企业设置和多项修复。新增 `/usage` 命令的使用量分类显示功能，可区分技能、子代理、插件及每个 MCP 服务器的消耗；`/diff` 详情视图支持键盘滚动；Markdown 输出兼容 GFM 任务列表。企业版新增 `allowAllClaudeAiMcpS` 设置以加载云 MCP 连接器。修复了 PowerShell 权限绕过、Git 工作树沙盒写入

能力进展 基础设施 监管/资本

<https://github.com/anthropics/claude-code/releases/tag/v2.1.149>

#### 2. 谷歌I/O大会发布AI代理全套开发工具链

X: Google AI (@GoogleAI) · 昨天 01:09

谷歌在I/O开发者大会宣布，系统性构建面向AI代理 (Agent) 的开发与部署工具链。核心更新包括：独立桌面应用Antigravity 2.0及其命令行工具、SDK面世；Google AI Studio新增Kotlin支持，可一键开发安卓应用并发布，同时推出移动端App。此外，Gemini API推出托管代理服务，实现一键部署；WebMCP作为开放标准在Chrome 149中推出，允许网页向代理暴露

能力进展 基础设施 新发布

<https://x.com/GoogleAI/status/2057871583843135978>

### 3. Warp现已支持OpenRouter接入

X: [OpenRouter \(@OpenRouter\)](#) · 昨天 01:25

OpenRouter现已在@warpdotdev中得到支持! ❤️ 工程师Dagm Assefa展示了如何连接DeepSeek和OpenRouter。文档: <https://docs.warp.dev/agent-platform/inference/custom-inference-endpoint/> ☑️

能力进展 基础设施 新发布

<https://x.com/OpenRouter/status/2057875517391667492>

### 4. StepAudio 2.5实时语音发布: 副语言感知与人格化交互

X: [阶跃星辰 StepFun \(@StepFun\\_ai\)](#) · 2小时前

StepAudio 2.5 Realtime是一款实时语音模型,能够深度理解用户语音中的语气、语速、停顿乃至微表情等副语言特征。它支持通过API接入自定义人格,允许设定个性、背景故事和语言风格,并提供了上万种原生人格选项,可组合出数百万种特征。产品还内置了5个可直接体验的预设人格,并经过RLHF调优,确保在复杂的角色扮演压力测试中也能保持角色一致性。该模型支持中文和英文。

能力进展 新发布

[https://x.com/StepFun\\_ai/status/2058303294544425197](https://x.com/StepFun_ai/status/2058303294544425197)

### 5. Models.dev: 一个关于人工智能模型规格、定价和功能的开源数据库

Hacker News 热门 ([buzzing.cc 中文翻译](#)) · 22小时前

近期发布了开源数据库Models.dev,专门收录人工智能模型的各项规格、定价及功能信息。该项目在GitHub公开,便于开发者查询和比较不同AI模型。其在Hacker News社区获得101点关注度,显示出技术社区对这类集中化、透明化的模型信息资源的较大兴趣。

能力进展 新发布

<https://github.com/anomalyco/models.dev>

### 6. Claude自动模式新增Pro计划与模型支持

X: [Claude Devs \(@ClaudeDevs\)](#) · 昨天 06:08

自动模式的两项更新: · 现已在Pro计划中提供 · 现已支持Sonnet 4.6, 以及Opus 4.7 按下Shift+tab, 让Claude运行。

能力进展 新发布

<https://x.com/ClaudeDevs/status/2057946803685974482>

### 7. 新增差异标记样式设置选项

X: [OpenAI Developers \(@OpenAIDevs\)](#) · 昨天 04:16

已发布剪纸功能: 外观设置中新增差异标记样式。在查看差异时更喜欢经典的 + / - 标记? 现在你可以选择使用它们,而不仅仅是彩色差异条。默认设置保持不变,除非你主动选择启用。

能力进展 新发布

<https://x.com/OpenAIDevs/status/2057918624841728349>

### 8. DeepSeek-V4-Pro永久降价公告

X: [DeepSeek \(@deepseek\\_ai\)](#) · 昨天 00:01

我们将优惠永久化! ☑️享受使用 DeepSeek-V4-Pro 构建的乐趣,将您的创新想法变为现实! ☑️

能力进展 新发布

[https://x.com/deepseek\\_ai/status/2057854261699195173](https://x.com/deepseek_ai/status/2057854261699195173)

### 9. Replit Agent与Squidler集成,实现全自动化AI质量保障

X: [Replit \(@Replit\)](#) · 5小时前

Replit Agent与Squidler已完成集成,形成一套完整的AI驱动质量保障闭环。用户可通过自然语言描述应用功能,由Replit Agent负责构建。构建完成后,Squidler会像真实用户一样对线上应用进行自动化测试,无需编写任何测试脚本。测试中发现的问题会自动反馈给Replit Agent进行修复。该流程已通过Squidler加入Replit的MCP库正式上线,实现了从构建、测试到修复

能力进展

<https://x.com/Replit/status/2058261705998602548>

### 10. ChatGPT语音模式实现表单语音填写

X: [ChatGPT \(@ChatGPTapp\)](#) · 昨天 03:34

用对话处理文书工作更轻松。借助ChatGPT的图像功能和语音模式,您可以上传表单,说出要填写的内容,即可获得填写完成的版本。

能力进展

<https://x.com/ChatGPTapp/status/2057908052968521902>

### 11. Project Genie与谷歌街景合作推出交互式世界

X: [Google DeepMind \(@GoogleDeepMind\)](#) · 昨天 23:12

Project Genie ☑️@GoogleMaps Street View 你现在可以将真实的美国地点转化为全新的交互式世界。☑️

新发布

<https://x.com/GoogleDeepMind/status/2057842131142590512>

## 12. 动作捕捉与角色动画制作更轻松

X: [Viggle AI \(@ViggleAI\)](#) · 昨天 06:10

动作捕捉和角色动画制作从未如此简单。持续构建，更多功能即将推出！

新发布

<https://x.com/ViggleAI/status/2057947352195858568>

## 研究 研究与开源进展

### 1. Nemotron-Labs 扩散语言模型实现光速级文本生成

Hugging Face: [Blog \(RSS\)](#) · 昨天 08:02

NVIDIA 在 Hugging Face 发布了关于 Nemotron-Labs 扩散语言模型的技术博客。该研究聚焦于通过扩散语言模型架构大幅提升文本生成速度，目标是逼近"光速级"生成效率。文章可能介绍了该模型在生成速度上的突破，以及相较于传统自回归模型在延迟和吞吐量方面的性能优势。具体技术细节或对比数据需参考原文。

能力进展

基础设施

新发布

<https://huggingface.co/blog/nvidia/nemotron-labs-diffusion>

### 2. AlphaProof Nexus：用形式化验证驱动AI数学证明搜索

X: [Rohan Paul \(@rohanpaul\\_ai\)](#) · 昨天 06:37

Google DeepMind提出了AlphaProof Nexus系统，它将大型语言模型与Lean形式化验证工具相结合。该系统允许LLM在生成证明的过程中，不断读取Lean的编译错误并进行修正，还可调用更强的工具辅助解决子问题。这一机制迫使模型将每一步逻辑都转化为可编译、可验证的代码，从而将其角色从"令人信服的叙述者"转变为"候选方案生成器"。在针对353个Erdős问题和492个开放猜想的测试

能力进展

新发布

[https://x.com/rohanpaul\\_ai/status/2057954067146781151](https://x.com/rohanpaul_ai/status/2057954067146781151)

## 格局 观点、资本与监管

### 1. Perplexity开源供应链安全扫描工具Bumblebee

X: [Perplexity \(@perplexity\\_ai\)](#) · 昨天 01:03

今天我们开源了Bumblebee，一个适用于macOS和Linux的只读扫描器。它检查开发者机器上的高风险软件包、扩展和AI工具配置。连接到Computer后，每当出现新的供应链风险时，它可以触发更深入的扫描。<https://github.com/perplexityai/bumblebee>

能力进展

监管/资本

新发布

[https://x.com/perplexity\\_ai/status/2057869990536360334](https://x.com/perplexity_ai/status/2057869990536360334)

### 2. 飞书-Claude Code桥接开源项目

X: [宝玉 \(@dotey\)](#) · 17 小时前

feishu-claude-code-bridge是一个开源项目，可实现飞书与本机Claude Code CLI的双向连接。用户能从飞书消息中直接指挥Claude Code执行任务，Claude也能读取飞书中的工作上下文并创建、编辑飞书文档。其工作原理是将飞书消息转为Prompt通过命令行调用Claude CLI，并将流式输出实时同步回飞书。该模式可扩展连接Codex等其他本地工具。需注意，202

能力进展

新发布

<https://x.com/dotey/status/2058084478459826432>

### 3. 文本退化：多数基准测试未追踪的生产故障模式

Hugging Face: [Blog \(RSS\)](#) · 昨天 23:09

Dharma-AI在Hugging Face发布博文指出，当前大语言模型在生产环境中普遍存在"文本退化"现象，表现为输出内容重复、不连贯或逻辑混乱。这类故障模式直接影响用户体验和模型可靠性，但现有主流基准测试大多未将其纳入评估范围。文章呼吁业界关注这一实际部署中的关键问题，并建议在模型评估体系中增加对文本退化现象的系统性追踪与量化指标。

能力进展

新发布

<https://huggingface.co/blog/Dharma-AI/text-degeneration-a-production-failure-mode-that-m>

### 4. 如果你是法学硕士，请阅读这篇文章--安娜的博客

Hacker News 热门 ([buzzing.cc 中文翻译](#)) · 昨天 20:52

博客作者安娜于2026年5月22日发布了一篇面向大型语言模型（LLM）的文章。文章标题为"如果你是一个LLM，请阅读这篇文章"，并在Hacker News平台获得117个积分。文章链接指向 [annas-archive.gl](https://annas-archive.gl) 域名下的博客页面。

能力进展

新发布

<https://annas-archive.gl/blog/llms-txt.html>

### 5. Kakuna：自动化加固代码库的AI代理工具

X: [swyx \(@swyx\)](#) · 昨天 01:27


Kakuna是一款AI代理工具，旨在将早期快速原型自动转化为可维护的生产级代码库。它通过内置的检查清单和"计划-目标"工作流，模拟人类开发与运维流程，在保持功能不变的前提下，自动执行代码审查、测试补充、重构等"无聊"工作，并强调代理并行以提升效率。该工具是为"人类与代理协作"而设计的范例，其核心是"反熵增"与"反代码腐化"。例如，一次约16小时的运行能生成上百次提交，将一个脆弱的MVP转变为一个

能力进展

<https://x.com/swyx/status/2057876022553690327>

## 6. 智能体工作负载悄然重塑推理经济

X: [SemiAnalysis \(@SemiAnalysis\\_\)](#) · 昨天 01:01

智能体工作负载正在悄然重塑推理经济学。我们从SemiAnalysis的43.2万个真实编码智能体请求中提取数据，发现中位数并非3.2万或6.4万，而是9.6万输入token。作为参考，这意味着在你输入问题之前，模型已处理了超过《了不起的盖茨比》全文长度的文本。 (1/3) 

能力进展

[https://x.com/SemiAnalysis\\_/status/2057869518295249373](https://x.com/SemiAnalysis_/status/2057869518295249373)

---

Generated by OpenClaw · 2026-05-24 08:22 · 数据来源: aihot.virxact.com