

人工智能行业专题（15）

从全球模型巨头的发展历程，思考模型企业的 壁垒与空间

行业研究 · 行业专题

互联网 · 互联网 II

投资评级：优于大市（维持）

证券分析师：张伦可

0755-81982651

zhanglunke@guosen.com.cn

S0980521120004

证券分析师：张昊晨

zhanghaochen1@guosen.com.cn

S0980525010001

证券分析师：刘子谭

liuzitan@guosen.com.cn

S0980525060001

- 根据Semi-Analysis报道，26Q1预计Anthropic单季度ARR的净增规模将首次超越OpenAI，成为全球AI收入规模增长最快的大模型公司。我们认为伴随模型跨越Agentic Coding拐点，当前正处于Agent爆发的起点，OpenClaw仅用2个月就成为GitHub最多星标数的软件项目。当前发展的速度和技术变革所带来的商业化影响，我们认为已经可以与2000年互联网变革的高度相比较。
- 本篇报告意在对比Anthropic、OpenAI和谷歌Gemini的发展历程、产品矩阵、商业化策略等，去思考AI时代大模型企业发展的核心壁垒和未来商业化空间。**我们认为今天Anthropic的快速增长，更核心的源自创始人Dario的敏锐的技术品味（Taste）所驱动的发展决策。AI的未来是未知的，无法通过后视镜前瞻去选择商业化空间最大的发展方向，技术驱动产品发展是更成功的路径，因此技术领袖的战略判断和选择非常重要。**比如三年前，OpenAI认为要“大力出奇迹”，先把模型弄聪明，再用人类反馈（RLHF）去修补它；Google更侧重于打造模型全家桶赋能其自身生态；而Anthropic认为模型必须从底层架构上就是可控的、讲逻辑的、和严格遵守规则的，因此Anthropic选择聚焦编程场景训练。而过去两年模型训练中RLVR（可验证强化学习）的方式恰好在编程领域取得了更加明显的能力提升，最终Anthropic率先实现Agentic Coding能力的跨越式拐点，即Opus 4.5模型的推出，自此开启模型自主完成任务、调用工具的时代，推动OpenClaw风靡全球，拉动模型厂商API类型收入快速增长。
- **伴随模型能力快速提升，我们发现模型和应用的边界正在逐渐模糊。比如当Claude Opus 4.5能够实现自任务时长明显提升后，模型自主调用工具完善Agent任务，实际上正在对软件/互联网应用层过去的工作流设计、用户使用习惯产生明显变化，用户可以通过构建各类型Skills（实际更简单）的方式完成过去应用层的功能。**我们观察到海外头部明星AI应用，Cursor、Perplexity 等由于缺乏底层模型壁垒，也开始面临用户数的冲击等问题。
- **投资建议：我们认为AI时代应该重点关注ARR快速增长的前沿大模型厂商，以及已经降本增效或增收明显的公司。**
- **风险提示：宏观经济波动风险、下游需求不及预期风险、核心技术水平升级不及预期的风险、AI快速迭代平权化下竞争加剧等。**

Anthropic、Google、OpenAI对比

图：三家模型厂商多维度对比

	Anthropic	Google	OpenAI
创始人背景	Dario Amodei : 前OpenAI 研究副总裁, 领导了 GPT-2/3 开发 Daniela Amodei : 前 OpenAI 安全与政策副总裁, Dario 的妹妹	Sundar Pichai : 谷歌CEO, 04年加入谷歌, 最初担任产品管理副总裁 Demis Hassabis : Deepmind联合创始人兼CEO, 曾领导团队开发了 AlphaGo	Sam Altman : OpenAI 创始人兼CEO, 创业孵化公司Y Combinator 前总裁
企业员工人数	~4000	190,820 (25Q4)	~4000
AI产品矩阵	Claude、Claude Code、Claude Cowork	Gemini、NotebookLM、Flow Whisk、Antigravity	ChatGPT、Pulse、Codex、Sora
发展策略	坚持2B路线和Coding场景, 在产品布局上相对克制	专注模型多模态能力, 围绕AI原生及AI赋能传统产品两条路线布局应用	通过2C场景打造核心壁垒, 开始发力企业业务, 同样强调多模态路线
用户数		MAU 7.5亿	WAU 9.1亿
算力储备	25年末与微软签订300亿算力资源, 投资500亿建设自有数据中心。同时与谷歌新签订100万张TPU合同。	25年CAPEX914亿美元, 同比+74%, 预计2026年资本开支1750-1850亿元, 同比+91%-102%	星际之门项目目标建设10GW算力, 此外25年还与微软和AWS分别签订2500亿美元和380亿美元的算力服务
估值(亿美元)	3800	37000	8400
25年全年收入体量	45亿美元		131亿美元
26Q1 ARR(预测)	190亿美元		250亿美元
26Q1 净增ARR(预测)	100亿美元		36亿美元

资料来源: 各公司官网、路透社、The Information、国信证券经济研究所整理

- **一、Anthropic：凭专业生产力打造高毛利护城河**
 - 核心团队成員、经营理念、算力储备
 - 模型能力：Coding、Agent场景下的SOTA模型
 - 商业模式：极简产品矩阵，API贡献主要收入
 - 财务表现：最强模型带来token溢价，28年有望迎来现金流转正
- **二、谷歌：多模态能力领先，生态优势明显**
 - 模型能力：围绕多模态能力打造模型矩阵，综合性能领先
 - 商业模式：原生AI应用+Gemini赋能传统产品，云与广告受益增长
- **三、OpenAI：C端产品领导者，开始发力企业市场**
 - 模型能力：模型路线从分化到统一
 - 商业模式：C端产品领导者，发力企业市场
 - 财务表现：收入预测持续上修，预计2030年超过2800亿
 - 算力储备：星际之门项目持续推进
- **四、静态理解模型的商业化市场空间**

Anthropic: 核心团队成员来自OpenAI, 当前估值3800亿美元



- Anthropic成立于2021年5月, 核心团队成员来自OpenAI, 联合创始人Amodei兄妹曾分别担任OpenAI研发副总裁和安全与政策副总裁。20年末由于对认为OpenAI对安全的投入不足, 以及与微软过度绑定会导致公司过度商业化, 对技术的控制减弱, 因此选择离开创建Anthropic。24年开始, 陆续引入了外部来自谷歌、Meta、Salesforce等公司的高管。
- 估值: 26年2月Anthropic完成300亿美元G轮融资, 其中包含25年11月来自微软和英伟达的150亿美元融资, 投后估值已达到3800亿美元。

图: Anthropic核心高管团队

姓名	职位	背景 / 关键贡献	任命/加入时间
Dario Amodei	联合创始人兼 CEO	前 OpenAI 研究副总裁, 领导了 GPT-2/3 开发	2021年1月
Daniela Amodei	联合创始人兼 President	前 OpenAI 安全与政策副总裁, Dario 的妹妹	2021年1月
Ami Vora	产品负责人 (Head of Product)	前 Meta 高管、公司产品副总裁; 2026年接手商业化产品线	2026年1月
Mike Krieger	Labs 团队联合负责人	Instagram 联合创始人; 2026年卸任 CPO 转向前沿实验研发	2024年5月
Rahul Patil	首席技术官 (CTO)	前 Stripe 基础设施负责人, 负责全球规模化系统架构	2025年11月
Jared Kaplan	联合创始人兼首席科学官 (CSO)	约翰·霍普金斯大学物理教授, AI “缩放法则”共同发现者	2021年1月
Jan Leike	对齐科学负责人	前 OpenAI 超级对齐团队 负责人, AI 安全领域顶尖专家	2024年5月
Paul Smith	首席商业官 (CCO)	前 Salesforce 核心高管, 负责 10 亿美元级大客户订单谈判	2025年1月
Chris Ciauri	国际业务董事总经理	前 Google Cloud EMEA (欧非中东) 总裁, 负责全球扩张	2025年9月
Jack Clark	联合创始人, 政策与战略负责人	前 OpenAI 策略总监, 负责与全球政府及监管机构对接	2021年1月
Chris Olah	联合创始人, 可解释性研究负责人	前 OpenAI 研究员, 机械可解释性 (解构神经网络) 的开创者	2021年1月

图: Anthropic 融资及估值

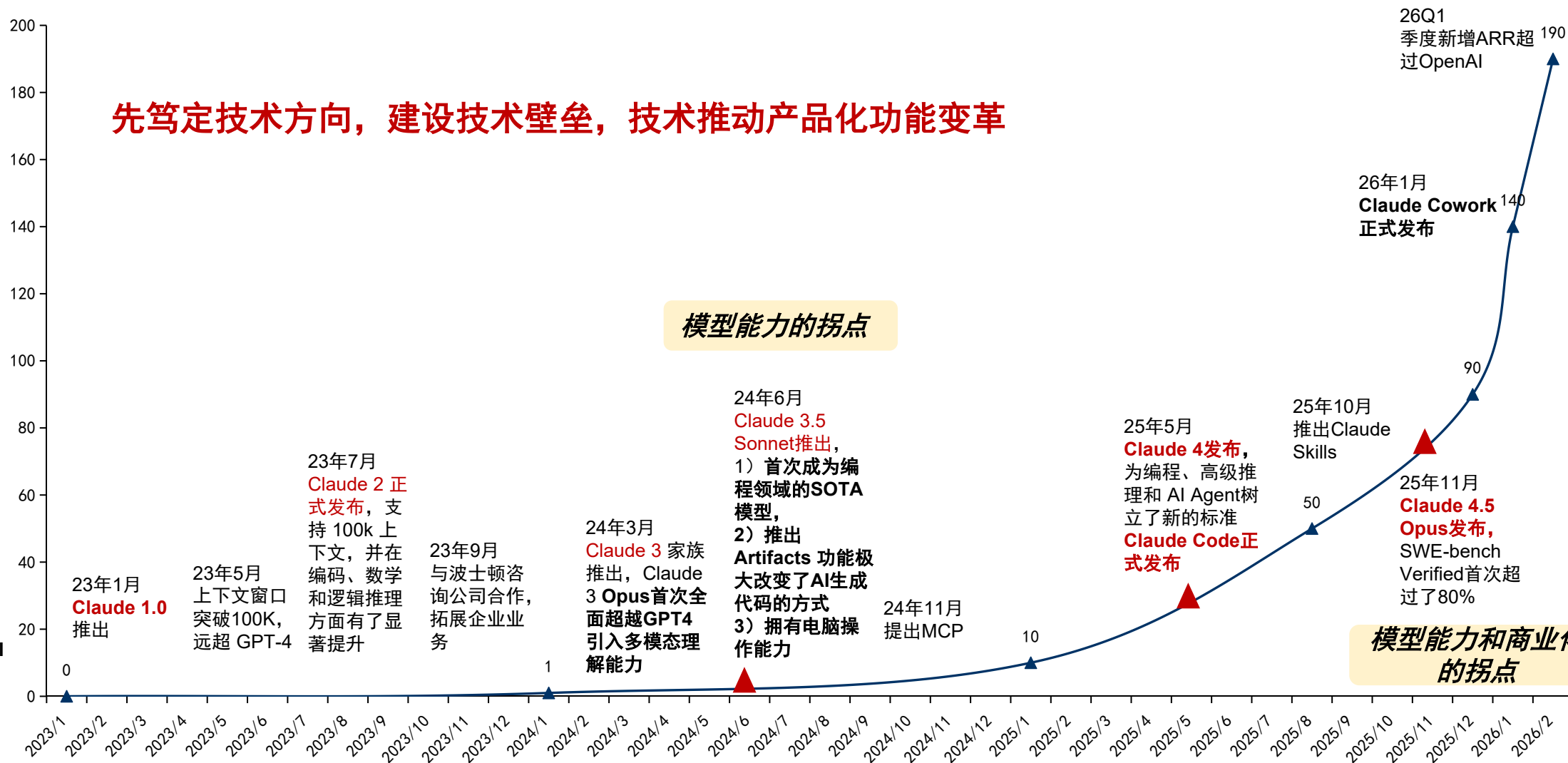
融资轮次	时间	融资金额 (亿美元)	投后估值 (亿美元)	领投方 / 主要投资者
Series A	2021年5月	1.24		Jaan Tallinn (Skype联合创始人), Eric Schmidt等
Series B	2022年4月	5.8		FTX首席执行官Sam Bankman-Fried领投, Caroline Ellison, Jim McClave等
Series C	2023年5月	4.5	50	Spark Capital 领投, 谷歌、Salesforce Ventures、Sound Ventures、Zoom Ventures 等
Series D	2023年9月	12.5	200	Amazon
Series E	2024年3月	27.5	184	Amazon
Series E+	2025年3月	45	615	Lightspeed Venture Partners 领投, Bessemer Venture Partners等
Series F	2025年9月	130	1,830	ICONIQ Capital, Fidelity, Lightspeed等
Series G	2026年2月	300	3,800	NVIDIA (100亿), Microsoft (50亿) 等

通过梳理Anthropic创始人Dario在不同时期访谈传递的内容主旨，我们总结Anthropic在AI发展趋势下有以下几点重要思考和策略：

1. 当OpenAI认为要“大力出奇迹”，先把模型弄聪明，再用人类反馈（RLHF）去修补它；Google认为要融入全家桶生态时，Dario认为模型必须从底层架构上就是可控的、讲逻辑的、严格遵守规则的。
2. Anthropic 专注于提供企业端服务，认为 AI 在企业端的应用（主要是在企业、初创公司、开发者和高效能用户之间的生产力场景里）会超过消费端。因为商业客户的需求（如在生物化学等专业领域）能为提升模型核心智能提供 stronger 的激励，这比面向普通消费者的应用更能推动技术突破。
3. 重视编程：1) 编程是AI构建的基础技能，因此也最快会被颠覆。技能与构建AI的人员距离越远，AI对其造成颠覆性影响所需的时间就越长；2) 模型在编程上变强后，也会帮助训练下一个更强的模型，形成正循环。
4. 定位平台公司，针对客户的核心需求开发垂直产品。对Anthropic的定位是一家平台型公司，在部分领域推出自己的产品（例如Claude code等）主要是因为：1) 直接接触达终端用户，能精准了解用户的使用场景、核心需求；2) 很多传统企业直接基于 API 进行开发，门槛较高，需要为他们提供更易上手的方案，要么是配套的开发工具包，要么是现成的应用程序。
5. 商业模式选择的思考：由于模型的快速变化，任何固定的产品形态，都有可能很快变得过时。API 的价值在于，它始终提供最接近底层能力的接口，让开发者基于最新技术构建。
6. Dario认为，我们正站在指数曲线的终点，而技术曲线与经济曲线之间存在天然的时间差，所以很多人还处于体感不明显的阶段，但是技术已经达到难以想象的水平。1) 技术对生产力的提升需要形成闭环，如果只是在某个中间环节插入AI，而没有重构整个工作方式，收益会很有限，甚至是负的；2) AI完成90%的代码到完成100%的代码是生产力数量级的差异。

Anthropic模型和产品发展历程

图：Anthropic ARR（亿美元）



数据来源：Anthropic、The Information、国信证券经济研究所整理

Anthropic采取多云路线，25年末加大算力建设投入

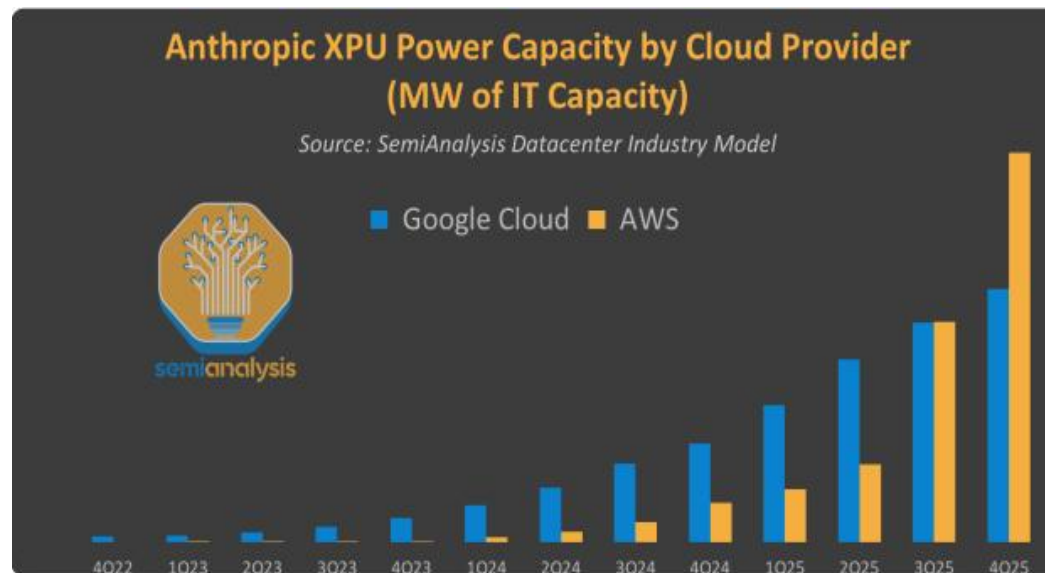
- Anthropic采取多云路线，同时使用多家芯片。根据SemiAnalysis报道，Anthropic算力来源最初主要是作为投资方的谷歌和亚马逊，25H1及以前谷歌占比较多。
- 25年开始Anthropic随着收入的快速增长，25年末连续签订了多笔算力采购/建设合同。分别与谷歌（获得100万TPU/1GW支持）、微软+英伟达（300亿美元/1GW合同）签订了新的合作协议，同时宣布将投资500亿美元与Fluidstack合作进行自有数据中心的建设。

图：Anthropic算力采购及建设情况

合作方	采购/自建	金额	具体细节
谷歌	采购		23年双方开始合作，谷歌一直是Anthropic的主要算力提供方。25年10月双方宣布扩展1GW的计算容量，26年上线，Anthropic将获得100万个TPU芯片
亚马逊	采购		AWS与Anthropic合作开展Project Rainier，25年10月已投入使用，到年底采用超过100万颗Trainium 2芯片进行训练和推理任务。其中仅印第安纳州项目就斥资110亿美元，发电量超过2.2GW。
微软+英伟达	采购	300亿美元	25年11月，Anthropic宣布采购300亿美元Azure计算容量，由英伟达Blackwell和Rubin提供支持，购买算力容量上限达到1GW，以获得英伟达和微软150亿美元的投资
Fluidstack	自建	500亿美元	2025年11月Anthropic宣布将投资500亿美元，与Fluidstack合作在德克萨斯州和纽约州建设数据中心

资料来源：Anthropic、国信证券经济研究所整理

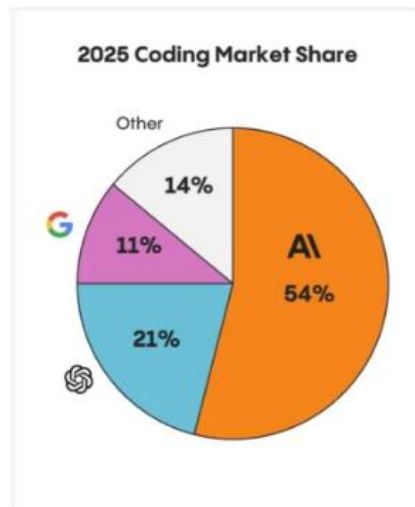
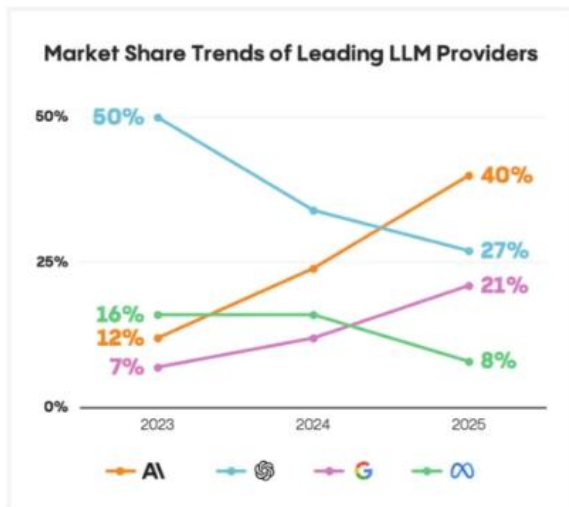
图：Anthropic XPU容量提供方



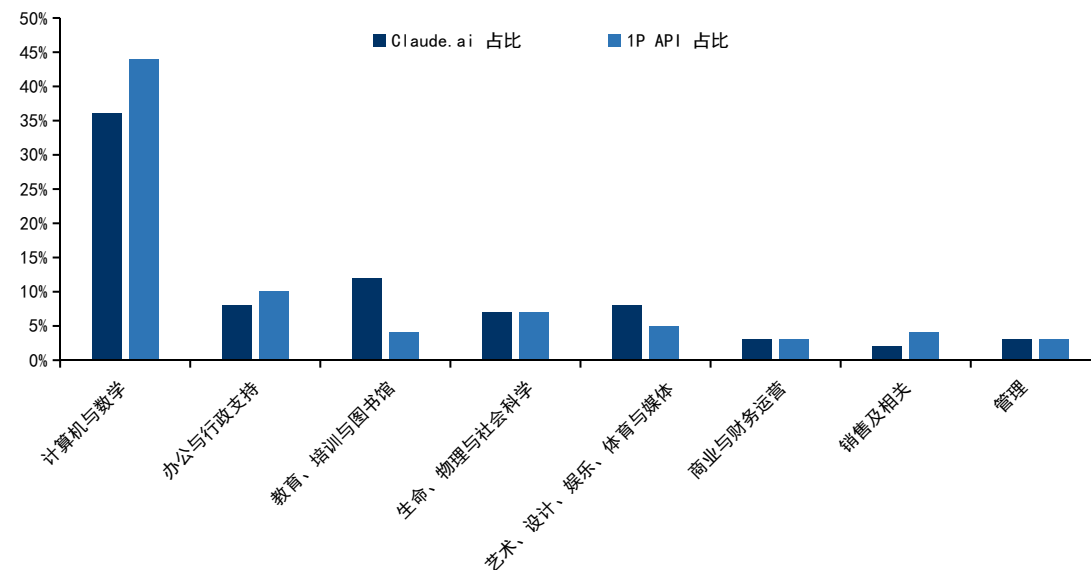
资料来源：SemiAnalysis、国信证券经济研究所整理

- **Anthropic Coding能力突出，是代码开发场景的首选。**Anthropic 的势头起于 2024 年 6 月发布的 Claude Sonnet 3.5，到了 2025 年 2 月的 Claude Sonnet 3.7 更是首次展示出“Agent 优先”的 LLM 雏形。到 2025 年 5 月，随着 Claude Sonnet 4、Opus 4 以及 Claude Code 的推出，其领先优势已被彻底坐实。根据Anthropic数据，44%的流量都是与计算机和数学类职业相关，细分的使用场景中前三名的分别是调试 Web 应用程序、解决技术问题以及构建专业商业软件。此外，得益于在编程场景中的出色表现，根据Menlo Ventures，25年Anthropic在企业大模型API市场份额已经达到40%，Coding 市场份额则达到54%。

图：企业大模型API调用份额



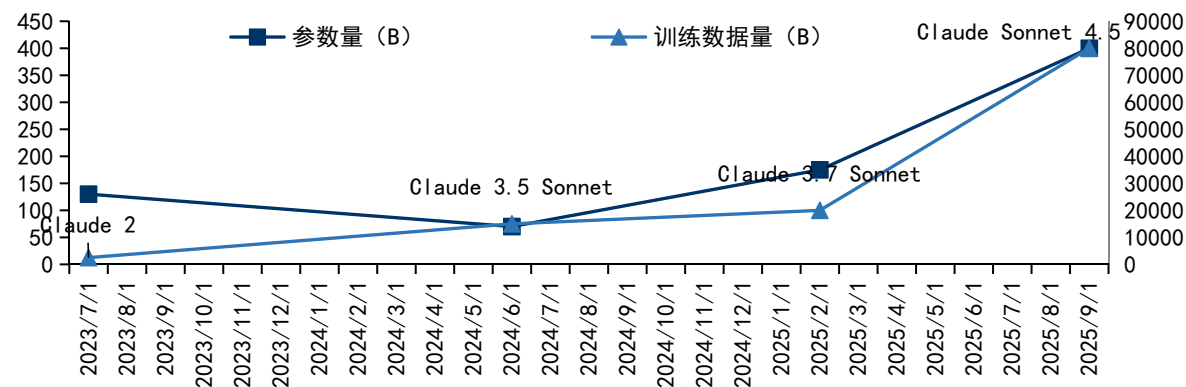
图：Anthropic Claude对话用途结构（通过Claude应用和API调用）



Anthropic: 25年Opus系列的迭代引领了Agent时代的开启

- 首款Claude模型上线于23年3月，并在24年3月Claude3开始分为不同尺寸的版本，Claude 3 Opus首次全面超越同时期的GPT-4。
- **25年Opus系列的迭代引发了Agent的拐点。**25年5月，Opus时隔一年推出了新一代版本，开启Agent时代，并且在年末的Opus 4.5推出后在编程、工具调用等维度均达到了SOTA水平，此后开始基于Opus4.5出现OpenClaw、Claude Cowork等优秀的Agent产品。

图：Claude部分版本参数量和训练数据量预估



资料来源：lifearchitect.ai、国信证券经济研究所整理

图：Claude版本迭代情况

大版本	版本型号	发布时间	特点
Claude 1	Claude 1.0 / Instant	2023年3月	首款商用模型；引入“宪法AI”（Constitutional AI）理念。
Claude 2.0	Claude 2.0	2023年7月	首次开放网页端，上下文达到100K，逻辑与代码大幅增强
Claude 2.1	Claude 2.1	2023年11月	上下文提升至 200k，幻觉率大幅降低，首次正式支持工具调用功能
Claude 3	Opus / Sonnet / Haiku	2024年3月	首次拆分为三个版本，由强到弱依次是Opus / Sonnet / Haiku；首次引入多模态视觉理解；Opus 性能首次全面超越同时期的 GPT-4。
Claude 3.5	Sonnet	2024年6月	Artifacts 功能上线；确立了全球最强代码模型的地位。凭借极速响应和卓越的代码逻辑反超了 3.0 Opus
	Haiku	2024年10月	
Claude 3.7	Sonnet (Reasoning)	2025年2月	首个“混合推理”模型；发布 Claude Code 终端开发工具。
Claude 4	Opus / Sonnet	2025年5月	Agent 时代开启；支持“扩展思考”模式与大规模工具并行调用。
Claude 4.5	Sonnet	2025年9月	进一步提升了任务稳定性、降低了编程错误率、又花了计算机操控能力
	Haiku	2025年10月	性能追平Sonnet4，首次引入扩展推理
	Opus	2025年11月	20万上下文窗口，实现多智能体协作和长程任务自主性
Claude 4.6	Opus	2026年2月	上下文窗口延伸至100万，在agent编程、计算机使用、工具使用、搜索、金融领域达到行业最强

资料来源：Anthropic、国信证券经济研究所整理

Opus 4.5在复杂任务的交付率上实现质的飞跃

- **Opus 4.5在编码上扮演了自主性极高的AI工程师角色。**举例：以前让模型“写一个 Python 函数来抓取网页”，然后“把这个基于 Django 的老项目迁移到 FastAPI，并重构数据库模型”，模型容易混淆格式，导致代码跑不起来。但是Opus 4.5 能“脑补”出整个项目的 50+ 个文件之间的引用关系。修改 A 文件时，它会自动意识到 B、C、D 文件也需要调整。Opus代码的一次性通过率显著优于其他代模型。
- 此外，Opus 4.5在定价上采取了更平衡性的市场策略，通过“effort”参数提供了前所未有的成本与性能调控精度。

图：Opus4.5特点

编码能力

Opus 4.5的编码优势为一种端到端从高效执行到全流程的自主软件工程能力。

- **复杂项目独立开发：**任务分解、系统架构设计、跨文件编写、重构优化代码。
- **代理式智能协作：**自主调用超过十种工具，协调处理多种跨系统复杂 workflows。
- **专业级成果输出：**胜任对代码质量和长期可维护性有高要求的大型商业项目。

性价比

Opus 4.5的定价策略实现性能与成本平衡的突破，以低成本提供旗舰级智能。

- **价格门槛大幅降低：**以显著低价提供同等甚至更优的顶级智能体验。
- **重塑市场竞争力：**企业用户获得行业领先的编码、推理与多模态能力，降低了尖端AI技术投入的总体拥有成本。

“effort”参数

Opus 4.5 引入“effort”参数通过资源调控器，实现精细化成本与性能管理。

- **场景化性能匹配：**根据任务的重要与复杂性，在高、中、低模式间灵活切换。
- **高效的资源利用：**在各模式下，都能在性能超越前代模型的同时，大幅减少资源消耗，极大提升计算资源投入效率。
- **企业级部署优化：**使企业IT基于业务实现大规模、可持续的AI应用部署。

资料来源：Anthropic官网、Openrouter、X, 国信证券经济研究所整理

Opus 4.5与不同工具的交互能力达到生产级别可用

- Opus 4.5在工具与生态上的演进，对内模型能力、对外开发生态、对下部署平台三位一体的协同设计，让AI Agent从概念验证，更近一步走向了规模化落地，标志其角色从单一的模型调用转变为智能体生态系统的核心引擎。
- Claude 3.5 时期推出的 Computer Use（操作电脑）在 4.5 Opus 上达到了生产级可用。** Claude 4.5 Opus 的能力包含：1) 像人一样看屏幕：它能直接看 GUI（图形界面），它能处理“去 SAP 系统里把上个月的财务报表导出来，然后发邮件给张总”；2) 视觉与逻辑的融合：如果网页弹出了一个“从没见过的广告窗”挡住了按钮，以前的 Agent 会卡死或报错。Opus 4.5 能理解弹窗。并模拟人类点 X 把它关掉，再继续操作。意义：这直接打通了所有没有 API 的老旧企业软件（Legacy Enterprise Software）。它就是一个不知疲倦的 RPA 机器人，但不需要写规则。

图：Opus4.5工具与生态

模型内核

Opus 4.5在底层能力上更强大可靠，为智能体注入自主发现与精准执行的能力。

- 动态工具发现：**支持动态工具搜索，按需筛选并加载当前任务工具，不背负冗余资源。
- 精准工具调用：**支持直接嵌入调用示例，显著提升复杂工具调用的准确性和可靠性。
- 程式化工具调用：**支持程式化工具调用，开发者可以在代码中直接结构化地调用工具。

开发生态

Opus 4.5从API延伸到开发者日常工具中，将顶尖能力无缝嵌入专业 workflow。

- 深度集成开发环境：**通过全面升级，将编码与智能体能力嵌入IDE，能获得从代码生成、bug修复到系统重构的端到端AI辅助。
- 赋能浏览器与办公软件：**通过Chrome扩展和Excel升级，直接作用于网页内容分析、浏览器任务自动化以及复杂的电子表格处理。

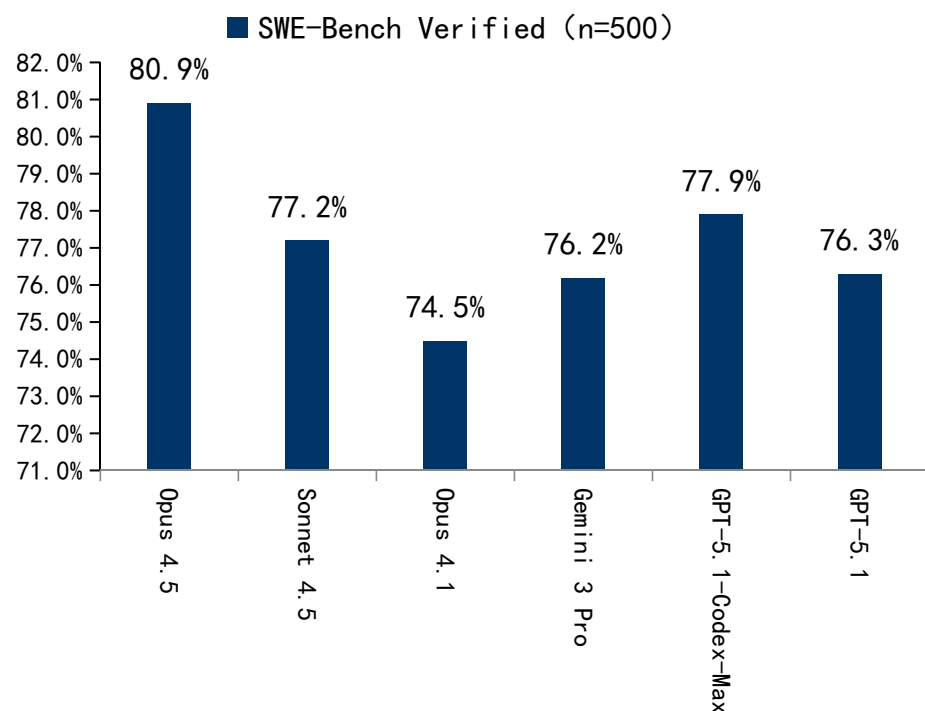
企业平台

Opus 4.5设计与Amazon Bedrock 深度集成，提供生产级智能体的部署与管理底座。

- AgentCore 基础设施：**提供持久化内存、Tool Gateway以及内置的安全与访问控制。
- 生产级可观测性：**通过集成，企业可以实时跟踪智能体 workflow 中的Token信息，实现透明的成本与性能管理。
- 支持长时任务：**提供长时间 workflow 支持，处理数小时的复杂分析、开发或自动化流程。

- 根据官方测试反馈，Opus 4.5对模糊需求的理解力得到了明显提升，复杂Bug自行定位也更稳定。
- 在真实场景的软件工程测试SWE-Bench Verified里，它是第一个拿到80%以上分数的模型；在视觉、推理和数学方面的测试都比前代模型更强，并且在多个重要领域都达到业界领先水平。

图：各模型的软件工程测试准确率(%)



资料来源：Anthropic官网，国信证券经济研究所整理

表：各模型在不同任务下的测试准确率

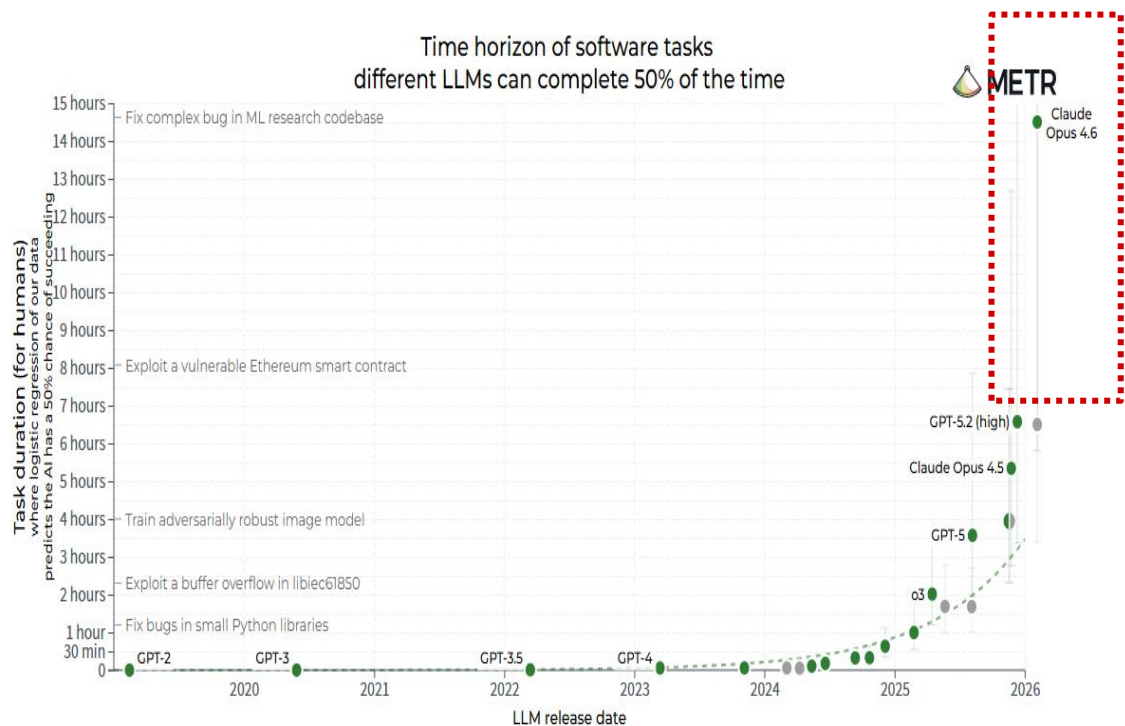
任务类别	测试名称	Opus 4.5	Sonnet 4.5	Opus 4.1	Gemini 3 Pro	GPT-5.1
Agentic coding (代理式编码)	SWE-bench Verified	80.90%	77.20%	74.50%	76.20%	76.30%
Agentic terminal coding (代理式终端编码)	Terminal-bench 2.0	59.30%	50.00%	46.50%	54.20%	47.60%
Agentic tool use (代理式工具调用)	t ² -bench (Retail)	88.90%	86.20%	86.80%	85.30%	—
	t ² -bench (Telecom)	98.20%	98.00%	71.50%	98.00%	—
Scaled tool use (规模化工具调用)	MCP Atlas	62.30%	43.80%	40.90%	—	—
Computer use (计算机使用)	OSWorld	66.30%	61.40%	44.40%	—	—
Novel problem solving (创新问题解决)	ARC-AGI-2 (Verified)	37.60%	13.60%	—	31.10%	17.60%
Graduate-level reasoning (研究生级推理)	GPQA Diamond	87.00%	83.40%	81.00%	91.90%	88.10%
Visual reasoning (视觉推理)	MMMU (validation)	80.70%	77.80%	77.10%	—	85.40%
Multilingual Q&A (多语言问答)	MMMLU	90.80%	89.10%	89.50%	91.80%	91.00%

资料来源：Anthropic官网，国信证券经济研究所整理

Claude系列模型支持Agent独立完成任务长度的时间明显领先

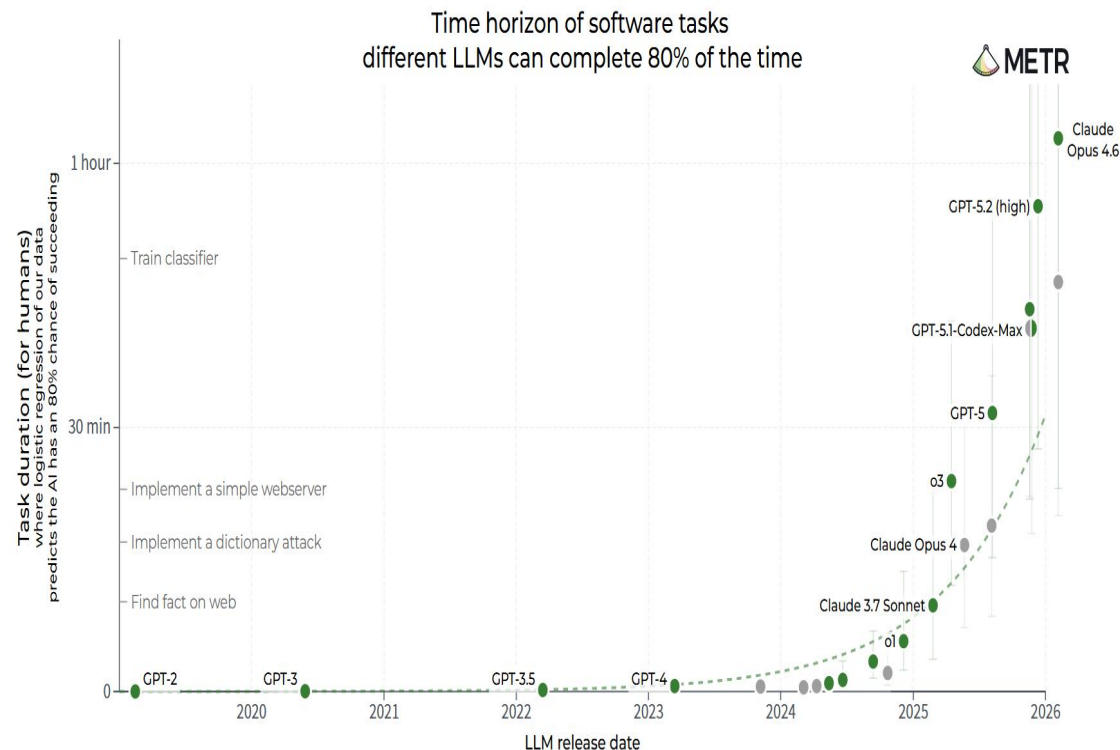
- Claude模型迭代的一个非常重要的变化是：能够完成的任务时长增加。根据METR，过去六年里，这一指标持续呈指数级增长，平均每七个月翻一番，目前Opus 4.6已经突破了1小时，由此推断，不到十年，我们将看到人工智能代理能够独立完成目前人类需要花费数天甚至数周才能完成的大部分软件任务。

图：通用前沿模型智能体能够以50%的可靠性自主完成的任务长度



资料来源：METR、国信证券经济研究所整理

图：通用前沿模型智能体能够以 80% 的可靠性自主完成的任务长度



资料来源：METR、国信证券经济研究所整理

Agent维度来看，Claude领跑模型能力榜单

- 根据Artificial Analysis，在Agent的测评维度中（GDPval-AA, τ^2 -Bench Telecom），Claude Opus 4.6目前保持领先地位。

图：主流模型综合智能 及 Agent能力比较



资料来源：Artificial Analysis、国信证券经济研究所整理测算

商业模式：API调用为Anthropic主要收入来源和产品形态

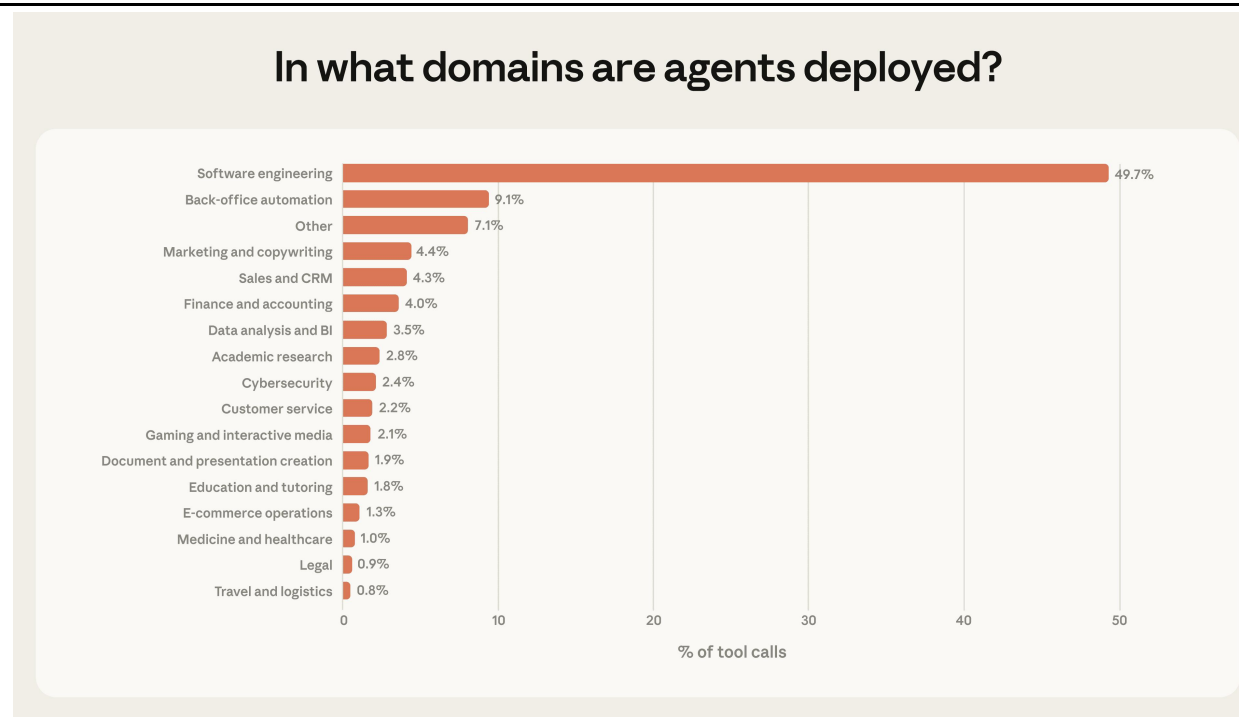
- **API调用**：API调用是当前Anthropic的主要收入来源，主要通过三大云厂的平台为客户提供服务，Claude三个系列中，Sonnet最近多个版本保持价格稳定，轻量版Haiku持续提价，最强模型Opus在4.5发布后降价2/3，在4.6推出的fast模式中，速度提升2.5倍的同时，价格提升6倍，进一步分层以适应不同的用户需求。
- 根据Anthropic的官方统计，软件工程领域占工具调用总数的近 50%，此外还涵盖商业智能、客户服务、销售、金融和电子商务等领域，但占比并不高，表明我们正处于智能体应用的早期阶段。软件工程师率先大规模构建和使用智能体工具

图：Claude不同模型版本输入/输出价格

单位：美元/百万 tokens	Opus	Sonnet	Haiku
Claude 3	15/75	3/15	0.25/1.25
Claude 3.5		3/15	0.8/4
Claude 3.7		3/15	
Claude 4	15/75	3/15	
Claude 4.1	15/75		
Claude 4.5	5/25	3/15	1/5
Claude 4.6	5/25		
Claude 4.6(fast)	30/150		

资料来源：Anthropic、国信证券经济研究所整理

图：Claude按领域划分的API调用分布



资料来源：Anthropic、国信证券经济研究所整理

商业模式：企业订阅计划&第三方合作伙伴

- **企业订阅计划：**针对企业用户Anthropic推出了 Team Plan和Enterprise Plan两个模式，其中Team Plan更适合小规模团队。
- **与第三方合作伙伴共同打造企业服务：**25Q4开始Anthropic在企业客户市场加大了跟第三方的合作，陆续与Salesforce、Cognizant、Snowflake、埃森哲等公司合作，通过将模型嵌入其AI产品或专门打造团队加速企业客户渗透率的提升。

图：Claude企业版的定位及功能

	Enterprise Plan	Team Plan
定位	面向数百人规模以上的大型企业，或对数据隐私、审计、上下文容量有极高要求的金融、法律、生物科研机构。	适合 5 到 75 人左右的小型初创公司、工作室或企业内部的特定职能部门
	Team Plan的全部功能，另增： <ul style="list-style-type: none"> • 增强的上下文窗口 • 高级席位下的Claude Code • Google Docs 目录编目 • 基于角色的访问权限及精细化权限配置 • 跨域身份管理系统（SCIM） • 审计日志 • 用于可观测性与监控的合规性 API • 自定义数据留存控制 • 网络级访问控制 • IP 白名单 • 提供符合 HIPAA 标准的方案 	<ul style="list-style-type: none"> • 包含 Claude Code 和 Cowork • 可连接 Microsoft 365、Slack 等应用 • 支持企业级跨组织搜索 • 提供统一账单与管理功能 • 支持单点登录（SSO）与域名捕获 • 提供针对远程和本地连接器的管理员控制 • 支持企业级部署 Claude 桌面应用
适用场景	金融风控、全公司级的 AI 代码助手部署、需要处理海量合规文件的法律部门。	需要 AI 协助写代码、写文案或整理日常文档，但不需要复杂的安全审计和超大上下文的团队。
定价	定制化报价	标准席位 20美元/月 高级席位 100美元/月（用量限制是标准版的5倍）

资料来源：Anthropic、国信证券经济研究所整理
请务必阅读正文之后的免责声明及其项下所有内容

图：Anthropic企业业务合作的第三方伙伴



1) 将 Claude 直接集成到 Snowflake 的企业数据中。
2) 此外还将建立一项联合市场推广 (GTM) 计划，在全球最大的企业中部署 AI 代理

1) 组建埃森哲 Anthropic 业务集团，埃森哲围绕 Claude 建立专门的业务部门。3 万名埃森哲专业人员将接受 Claude 培训
2) 推出一项新的联合服务，帮助CIO衡量人工智能的价值，并在整个工程组织中采用人工智能。
3) 联合开发行业解决方案

Cognizant 将 Claude 模型、Claude 代码、MCP 和 Agent SDK 与 Cognizant 的软件开发和人工智能平台相结合

印度IT服务公司，将 Anthropic 的 Claude 模型和 Claude Code 与企业级人工智能套件和服务组合 Infosys Topaz 相结合。

将 Claude 打造为 Salesforce Agentforce 平台的首选模型

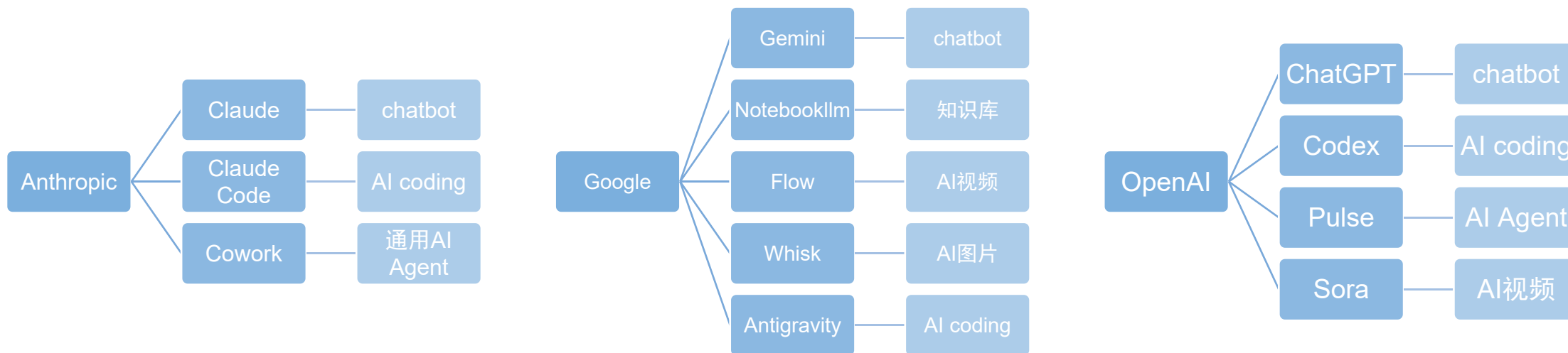
1) 企业将能够使用 Anthropic 的 Claude Agent SDK 在 Intuit 平台上构建和定制安全、精准的人工智能代理
2) Intuit的税务、财务、会计和营销专业知识和工具将直接应用于 Anthropic 的产品中，包括 Cowork、Claude for Enterprise 和 Claude.ai。

资料来源：Anthropic、国信证券经济研究所整理

商业模式：个人订阅制，产品侧重点与OpenAI和Gemini有明显区别

- 个人订阅制下，海外几大模型厂商都采取了打包产品矩阵进行订阅的形式，从而形成了更强的差异化。Anthropic的个人订阅制包含免费版、Pro版（20美元/月）、Max版（5x是100美元/月，20x是200美元/月，差异在于tokens的限额不同）。作为最基础的Pro版，既包含了Claude Code和Cowork的使用权，同时还能够获得嵌入excel等工具的插件能力，以及所有Claude模型的使用权。
- 产品策略差异对比：1) **Anthropic**：并不关心流量入口争夺，推出的编程和agent产品同时还以插件形式内嵌到浏览器、excel、PPT等工具中，并推出了多个垂类场景的Cowork插件（财务、法律、市场营销等）；2) **Google**：与OpenAI都更聚焦C端入口争夺，同时由于多模态能力出色，还推出了专门的知识库、视频、图片工具，此外还嵌入传统产品（搜索、Gmail等）；3) **OpenAI**：核心是C端流量入口ChatGPT，功能丰富度远高于前两者，例如引入健康、智能购物、生图生视频功能也直接嵌入ChatGPT中。

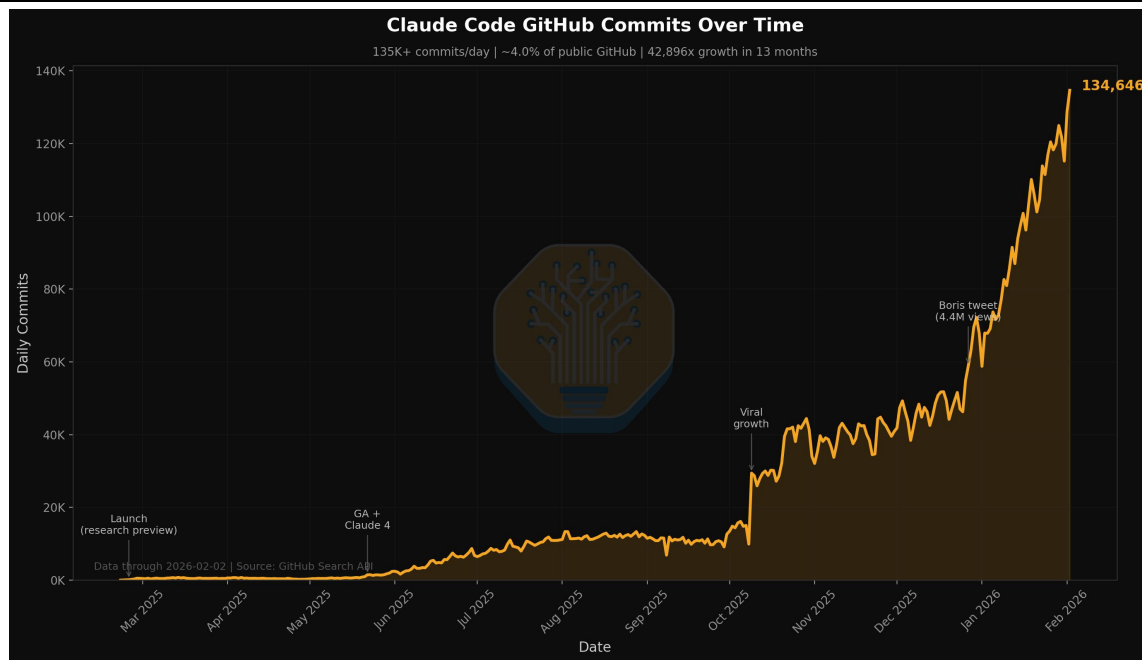
图：Anthropic、Google、OpenAI产品矩阵



Claude Code: 收入和用户26年以来进入加速期

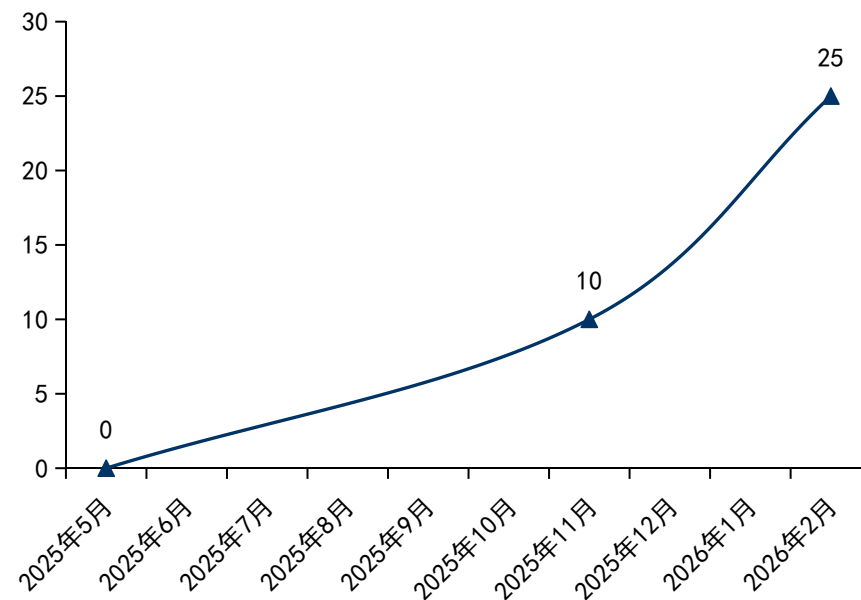
- Claude Code 是由 Anthropic 公司于 2025年2月24日推出的基于命令行的 AI 编程工具，不同于IDE插件，Claude Code直接运行于终端。采用CLI的模式一方面由于产品在设计之初的成本因素（由一个设计者发起，IDE插件的设计相对复杂），另一方面，**设计者认为当前模型处于快速迭代过程中，CLI能够更快响应模型的变化。**
- 26年初以来Claude Code的采用量和收入进入加速期。1) 据SemiAnalysis，截至26年2月，全球所有 GitHub 公开提交中，有 4% 是由 Claude Code 创建的——比一个月前的比例翻了一番。2) 收入端，25年11月ARR达到10亿美元，并且在26年2月达到25亿美元，且自2026 年初以来，Claude Code 的企业订阅用户数量增长了四倍，企业用户收入已占 Claude Code 总收入的一半以上。3) 自26年年初到2月，Claude Code 的每周活跃用户数量也翻了一番。

图：Claude Code 创建的 GitHub 公开提交



资料来源：SemiAnalysis、国信证券经济研究所整理

图：Claude Code ARR (亿美元)

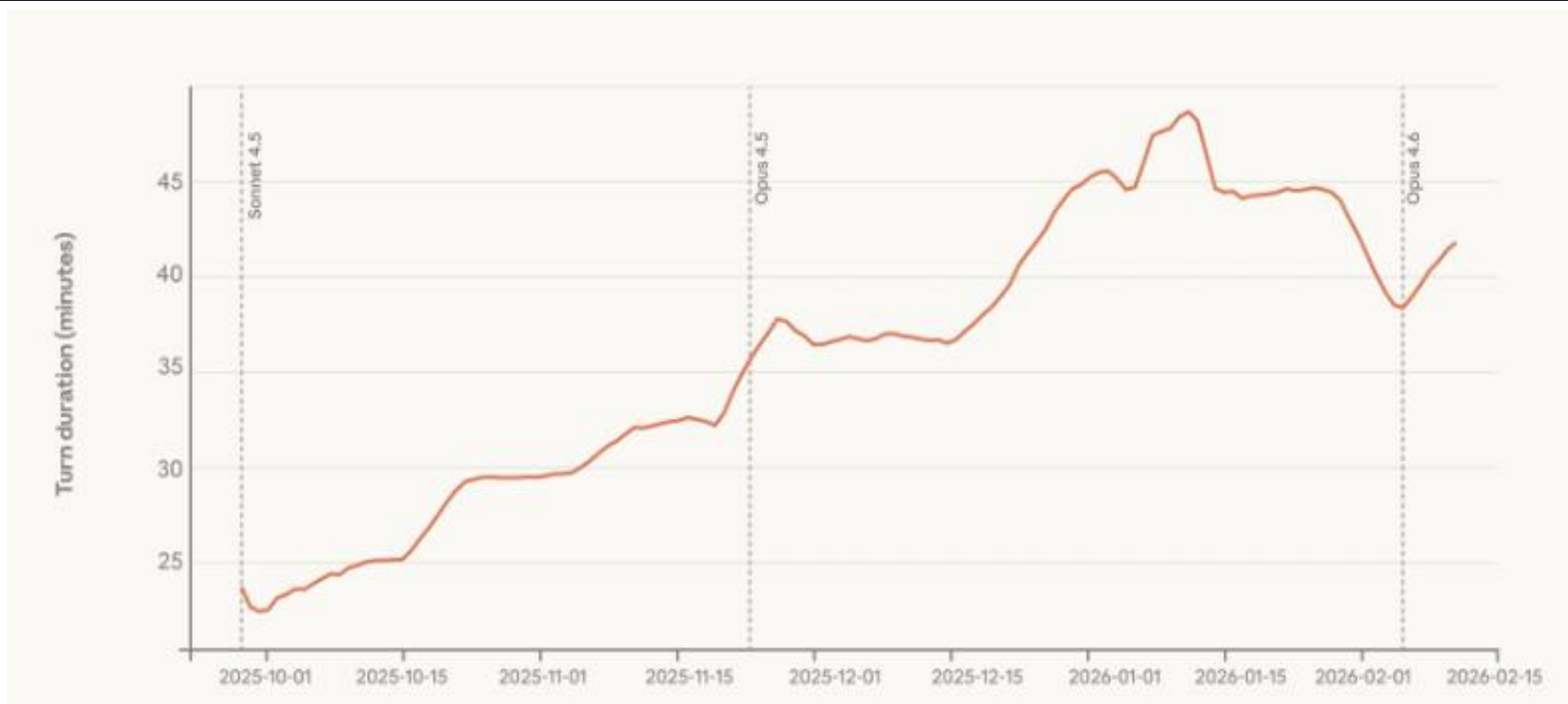


资料来源：Anthropic、国信证券经济研究所整理

Claude Code：持续使用时间不断增长

- 在运行时间最长的会话中，Claude Code 停止运行前的持续时间在三个月内几乎翻了一番，从9月下旬的不到 25 分钟增加到超过 45 分钟。这种增长在各个模型版本中都保持平稳，这表明其增长并非完全源于性能的提升，而是现有模型的自主能力远超其实际应用水平。

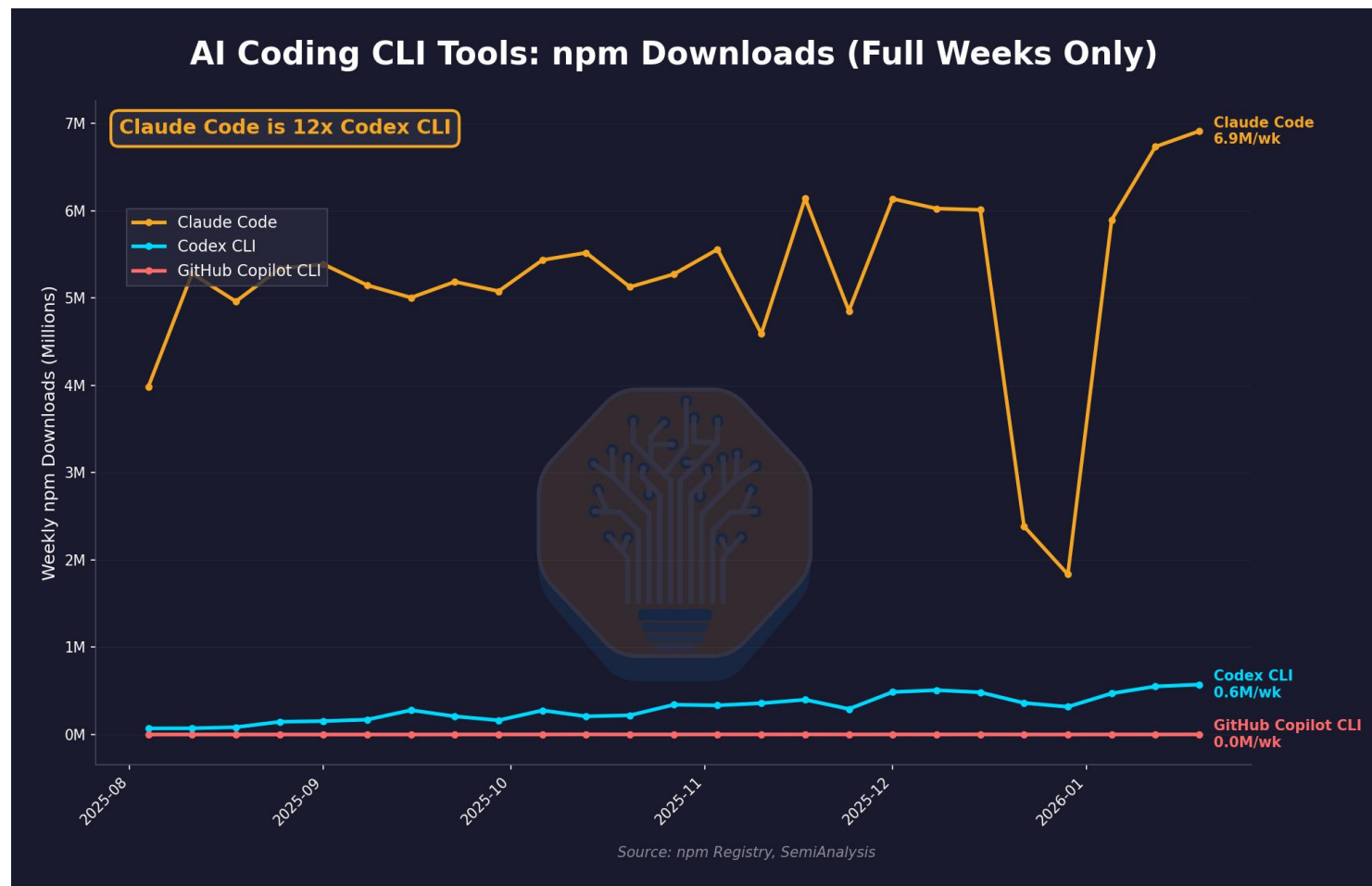
图：Claude Code 每次操作所持续的时间（使用时长）



资料来源：Anthropic、国信证券经济研究所整理

- 根据SemiAnalysis, 在 npm (全称 Node Package Manager, JavaScript 世界的“应用商店”或“插件库”) 上, Claude Code 的下载量大幅领先于Codex Cli。

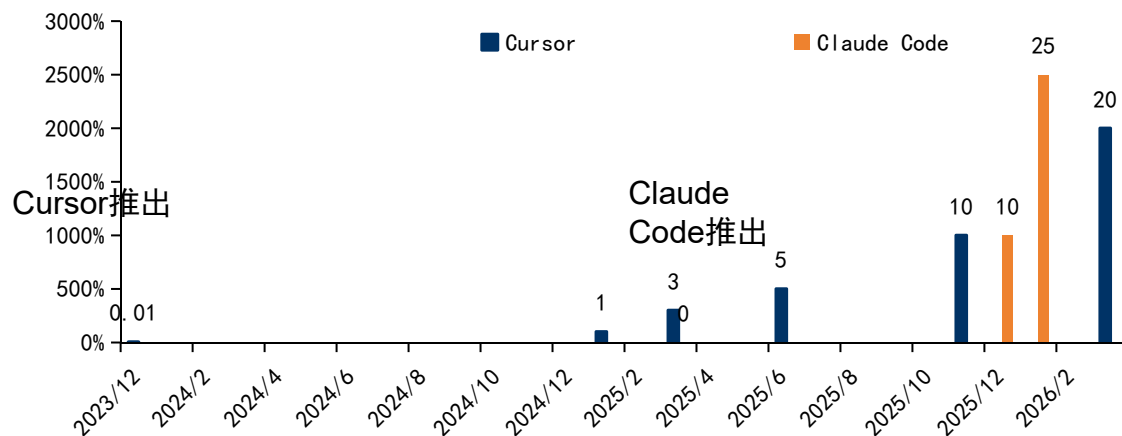
图: NPM AI Coding CLI工具的下下载量



资料来源: SemiAnalysis、国信证券经济研究所整理

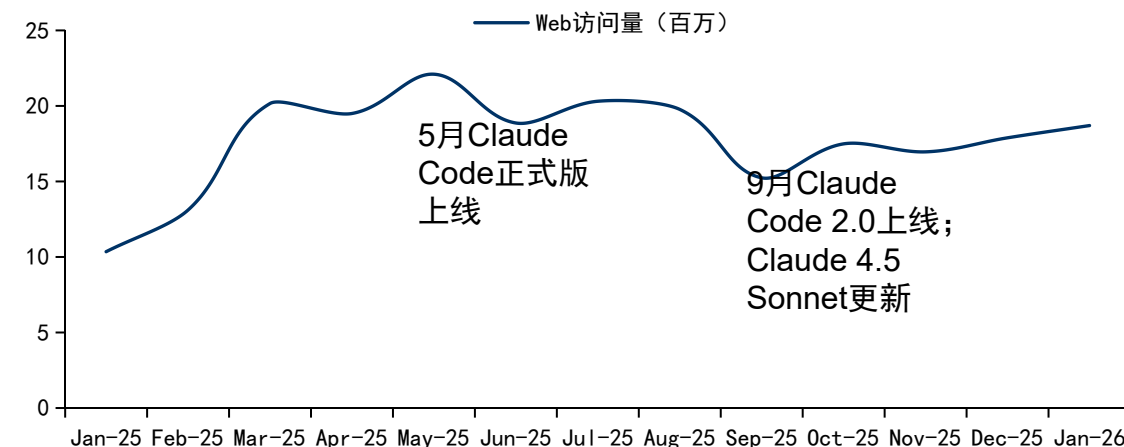
- **产品形态逐步统一。**在刚推出早期，Cursor更多是作为一款IDE，主打智能Tab补全和快速代码生成，而Claude Code 则主要是一个基于终端的AI coding/agent产品。但同时，Claude Code具有集成在VS Code、JetBrains等IDE的能力，而Cursor也在后续推出了Cursor CLI，方便直接终端使用。对于此前的用户而言，IDE界面更适合日常编码、重复性工作和快速原型开发，能够根据上下文智能预测代码，让编程变得流畅高效，而终端界面则便于复杂系统设计、代码重构和技术决策。但现在二者的差异正在快速缩小。
- **25年5月Claude Code正式版全面上线后，对Cursor的web流量就产生了冲击，并且双方的ARR水平快速拉开差距。**为此Cursor积极寻求差异化，1) 一方面在25年10月推出自研模型Composer 1，26年2月的Composer 1.5进一步强化，在Terminal -Bench2.0的测试中超过Sonnet 4.5，虽然仍弱于Opus 4.6，但只要订阅用户就可以免费使用，有良好的性价比优势。2) 另一方面，Cursor CEO也在2月发文表示，编程的Tab时代即将结束，新时代的典型特征是，Agent能在更长的时间跨度内、更少人工干预下，独立完成更大规模的任务，Cursor也在往这个方向迭代。

图：Cursor和Claude Code ARR对比



资料来源：The Information、Anthropic、Cursor、国信证券经济研究所整理

图：Cursor网站访问量



资料来源：AI产品榜、国信证券经济研究所整理

Claude Cowork: Agent时代最重要的产品雏形出现

- **Claude Cowork:** 2026年1月12日, Anthropic推出了Cowork——“适用于通用计算的Claude Code”。四名工程师仅用10天就完成了开发。大部分代码由Claude Code本身编写。其架构与Claude Agent SDK、MCP及子代理相同。该软件能根据收据生成电子表格, 按内容整理文件, 并从零散笔记中起草报告。Cowork更像是终端功能的Claude Code, 再加上桌面端界面。
- **通过开源插件库覆盖各垂类场景,** 客户可以利用这些插件将 Claude 打造为特定角色或团队的专家, 例如, 销售插件可以将 Claude 连接到CRM 和知识库, 让它了解销售流程, 并提供从潜在客户研究到电话跟进等各种操作的命令。
- **与OpenClaw相比,** Cowork的灵活性更低, 但有更强的安全性和更低的使用门槛。

图: Cowork 目前官方推出的插件

插件	主要功能	插件	主要功能
生产力	管理任务、日历、日常工作流程和个人事务	法律事务	审核文件、标记风险并跟踪合规情况
企业搜索	查找公司所有工具和文档中的信息	市场营销	撰写内容、策划营销活动、管理产品发布
插件创建	从零开始创建和自定义新插件	客户支持	问题分类、撰写回复并提出解决方案
销售	研究潜在客户、准备交易并遵循销售流程	产品管理	编写规格说明、确定产品路线图优先级并跟踪进度
财务数据	分析财务数据、构建模型并跟踪关键指标 查询、可视化和解释数据集	生物学研究	检索文献、分析结果、设计实验
		人力资源	支持员工生命周期内的人事运营, 从起草录用通知和制定入职计划到撰写绩效考核和进行薪酬分析。
设计	通过生成评论框架、撰写用户体验文案、运行可访问性审核和构建用户研究计划来加速设计工作流程。	工程	简化日常工程工作流程, 例如编写站会总结、协调事件响应、构建部署清单和起草事后分析报告。
运营	管理核心业务运营, 包括流程文档、供应商评估、变更请求跟踪和运行手册创建。	品牌声音	分析您现有的文档、营销材料和对话, 将您的品牌声音提炼成清晰、可执行的准则。
财务分析	支持每位财务分析师所需的基本工作流程, 从市场和竞争对手研究到财务建模和PowerPoint 模板创建和质量检查。	投资银行	加快交易流程, 包括审查交易文件、构建可比公司分析和准备推介材料。
股票研究	简化研究工作流程, 例如解析盈利报告、根据新的指导意见更新财务模型以及撰写研究报告。	私募股权	通过审查大量文件集、提取标准化财务数据、模拟场景以及根据投资标准对机会进行评分, 为交易搜寻和尽职调查提供支持。
财富管理	帮助顾问分析投资组合, 识别偏差和税务风险, 并大规模生成再平衡建议。		

图: Cowork 和 Openclaw对比

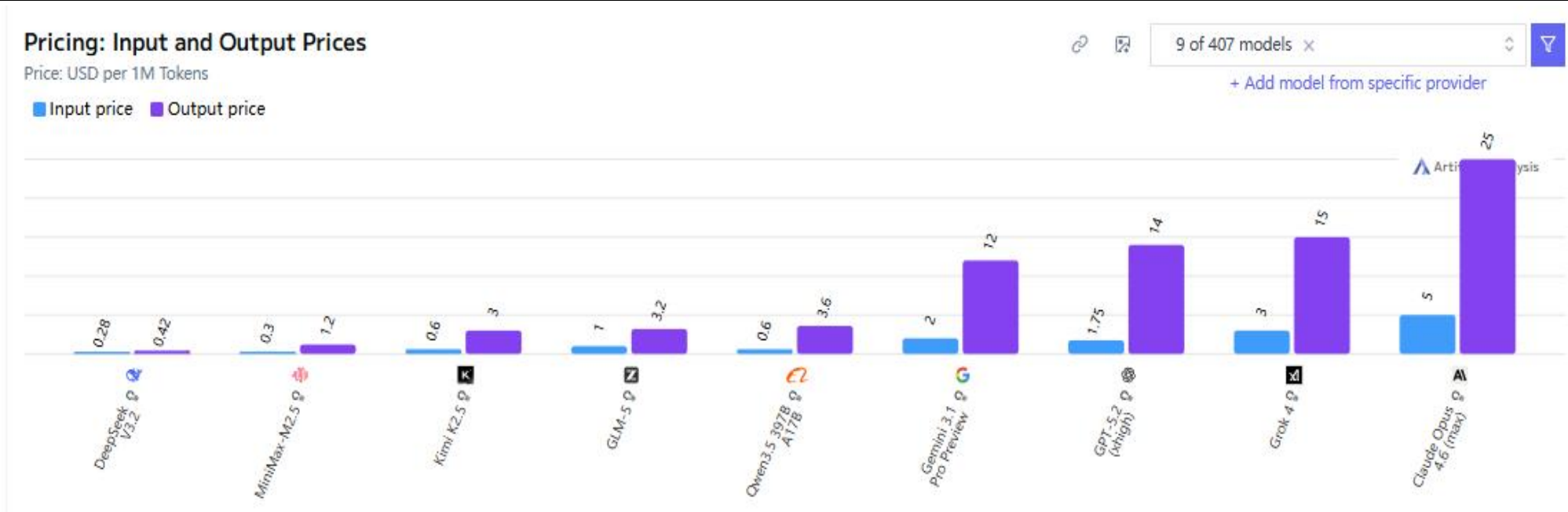
	Claude Cowork (Anthropic)	OpenClaw (开源项目)
产品性质	Anthropic 官方出品, 集成在 Claude Desktop App 中。	由开发者 Peter Steinberger 发起的开源项目 (现已加入 OpenAI 基金会)。
交互入口	桌面原生应用界面, 点击式操作, 适合非技术人员。	主要是 IM 聊天软件 (WhatsApp、Telegram、Slack 等), 随时随地远程控制。
运行环境	在受控的沙盒 (虚拟机) 中运行, 仅访问你授权的文件夹。	需本地或服务器 (如 Mac Mini、Docker) 常驻运行, 支持 24/7 主动执行任务。
技术门槛	极低。安装桌面端后即可使用, 无需配置 API 或终端。	较高。需要配置 API Key、Node.js 环境或部署服务器, 适合开发者。
核心功能	整理文件、写报告、处理表格、浏览器自动化 (侧重办公产出)。	自动退订邮件、监控提醒、多模型编排、远程任务下发 (侧重全能管家)。

资料来源: Anthropic、Openclaw、国信证券经济研究所整理

财务表现：最强模型带来token溢价，28年有望迎来现金流转正

- 以企业客户为主，API贡献主要收入来源，订阅收入占比持续提升。根据Dario Amodei，Anthropic目前大约80%的收入来自企业客户，即使是Claude code的ARR中也有一半来自于企业客户。形式上以API为主，根据The information，在25年8月的Anthropic的ARR结构中，API占比大约60%。但随着Claude code的用户快速增长，订阅收入在Anthropic收入结构中的占比正快速提升。
- 最强模型享受token溢价，因此毛利率水平也强于OpenAI。在最新的模型中，Claude Opus 4.6是定价最贵的模型，百万tokens输出价格为25美元，远高于其他厂商的SOTA模型，因此Anthropic拥有更好的毛利率水平，据The Information，25年Anthropic毛利率已回到40%，并且未来几年仍将持续提升，预计28年达到约75%。考虑到随着收入规模的增加，训练成本（研发费用）占比将持续下降，预计28年Anthropic有望最先迎来现金流的转正。

图：主流模型输入输出定价（美元/百万tokens）

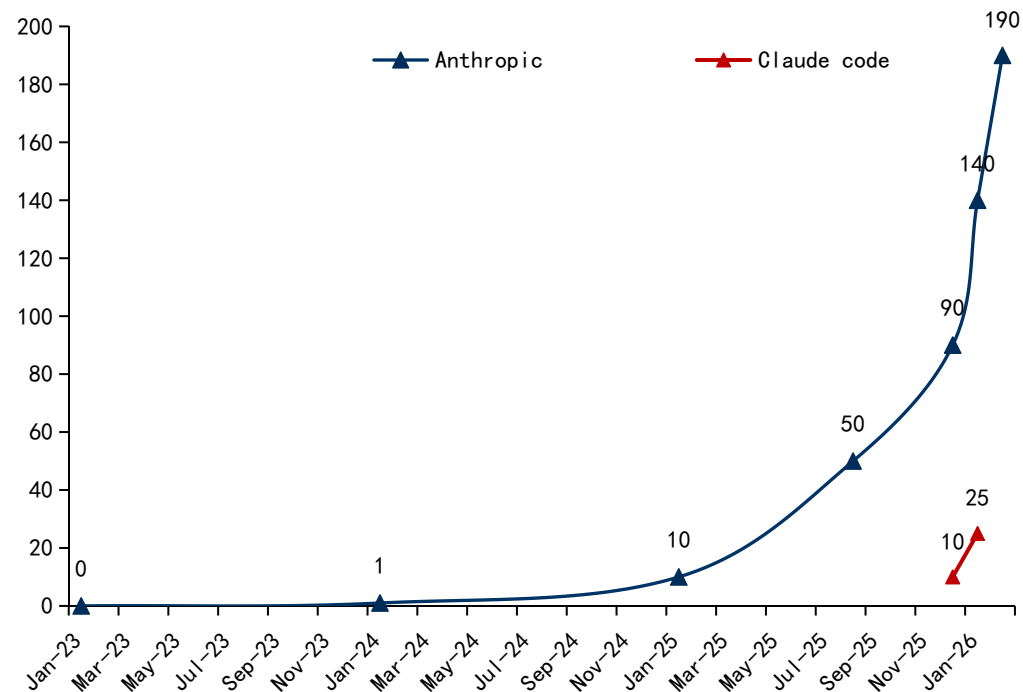


资料来源：artificial analysis、国信证券经济研究所整理

财务表现：26年初以来，Anthropic ARR收入呈现加速增长

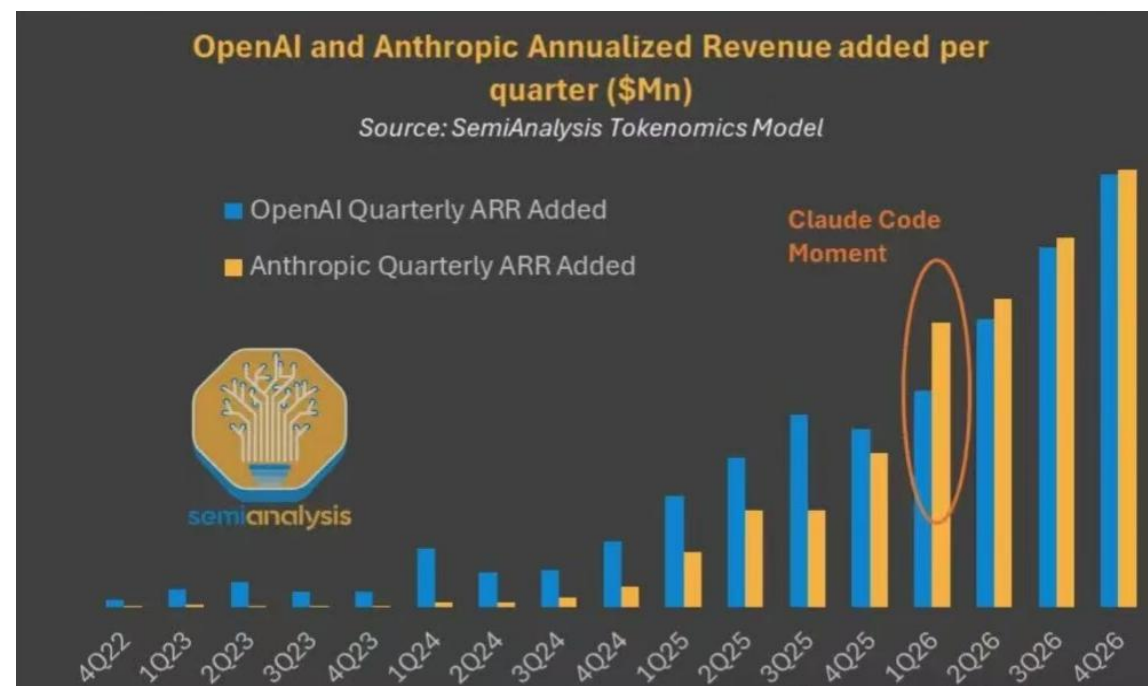
- 26年初以来，伴随Agent产品全球爆发，Anthropic ARR在前两月较25年12月实现翻倍以上增速，较去年下半年重新加速。根据Bloomberg报道，Anthropic在2月末ARR已达到190亿美元，根据Semi Analysis，在Claude Code的驱动下，预计26Q1其增量ARR有望反超OpenAI。

图：Anthropic ARR（亿美元）



资料来源：The information、国信证券经济研究所整理

图：Anthropic ARR（亿美元）



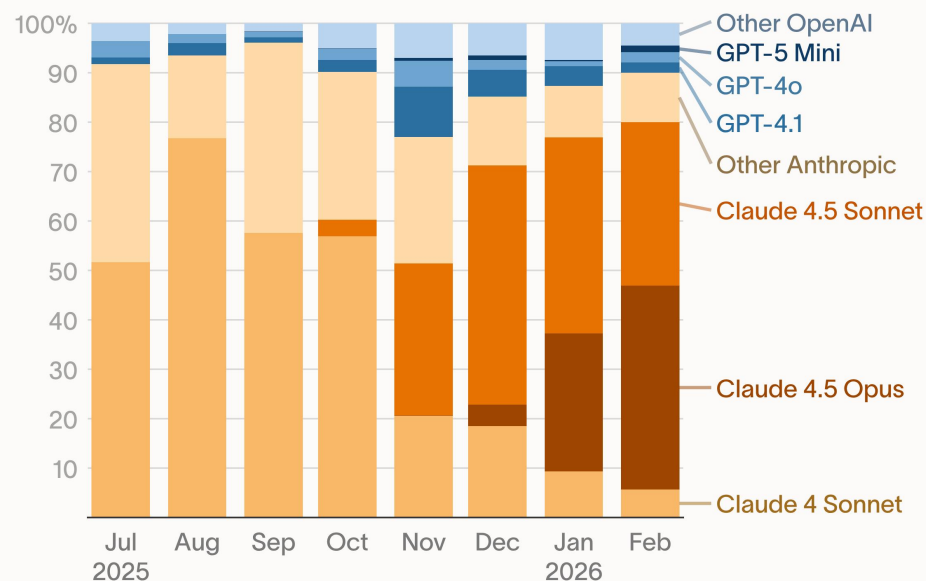
资料来源：Semi Analysis、国信证券经济研究所整理

图：Anthropic与OpenAI API 收入份额

Anthropic leads business API spend

Invoices ending February 2026

Market share of API spend by U.S. businesses



Source: Ramp Economics Lab (ramp.com/data); excludes spend on AI models not from OpenAI and Anthropic. Corporate card and bill pay data from 50,000+ U.S. businesses on Ramp's financial operation platform.



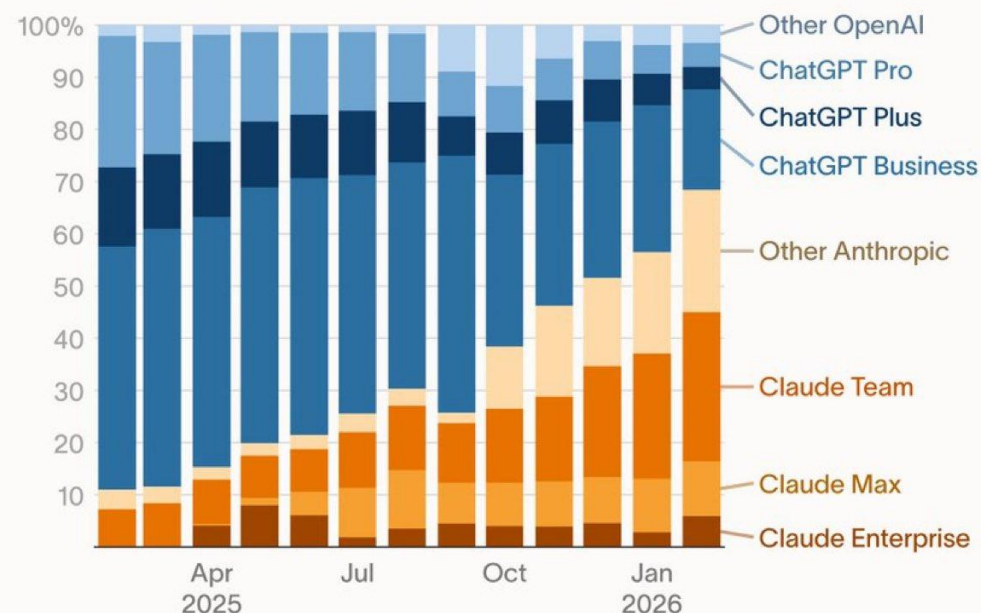
资料来源：Ramp、国信证券经济研究所整理

图：Anthropic与OpenAI 聊天机器人 收入份额

Anthropic leads AI chat for businesses

Invoices ending February 2026

Market share of spend for AI chat subscriptions by U.S. businesses



Source: Ramp Economics Lab (ramp.com/data); excludes spend on AI models not from OpenAI and Anthropic. Corporate card and bill pay data from 50,000+ U.S. businesses on Ramp's financial operation platform.



资料来源：Ramp、国信证券经济研究所整理

财务表现：最强模型带来token溢价，28年有望迎来现金流转正



- **算力总投入持续增加，与多家云厂合作。**据The information，预计26年Anthropic的推理+训练支出将达到约190亿美元，28年将增长至555亿美元，Anthropic与多家云厂合作，亚马逊Project Rainier为公司提供约100万颗Trainium芯片的算力支持，25年10月与谷歌签署了将租用100万颗TPU芯片（价值数百亿美元），同样在25年底还承诺向微软购买300亿美元的算力。
- **财务情况：**据The information，25年公司整体收入约45亿美元，由于需求端的持续超预期增长，公司在最新的预测中上调了对未来收入的展望，预计到28年将实现超过1000亿美元收入，上调幅度约40%，由于需求激增带来的模型相关费用也有所增加，因此预计现金流转正的时间点将推迟一年，28年迎来盈利。

图：Anthropic收入毛利率预测

Anthropic	2024A	2025E	2026E	2027E	2028E	2029E
单位：百万美元						
收入	381	4,500	20,000	55,000	102,000	148,000
YoY			344%	175%	85%	45%
营业成本（推理成本）	739	2,700	7,400	16,500	25,500	34,040
毛利率	-94%	40%	63%	70%	75%	77%
毛利润	-358	1,800	12,600	38,500	76,500	113,960
研发费用（训练成本）	1,500	4,000	12,000	23,000	30,000	42,000
研发费用率	394%	89%	60%	42%	29%	28%
销售&管理费用	114	1,125	10,600	24,000	35,000	50,000
销售&管理费用率	30%	25%	53%	44%	34%	34%
经营利润	-1,972	-3,325	-10,000	-8,500	11,500	21,960
OPM	-518%	-74%	-50%	-15%	11%	15%

资料来源：The information、国信证券经济研究所整理测算

- **一、Anthropic：凭专业生产能力打造高毛利护城河**
 - 核心团队成員、经营理念、算力储备
 - 模型能力：Coding、Agent场景下的SOTA模型
 - 商业模式：极简产品矩阵，API贡献主要收入
 - 财务表现：最强模型带来token溢价，28年有望迎来现金流转正
- **二、谷歌：多模态能力领先，生态优势明显**
 - 模型能力：围绕多模态能力打造模型矩阵，综合性能领先
 - 商业模式：原生AI应用+Gemini赋能传统产品，云与广告受益增长
- **三、OpenAI：C端产品领导者，开始发力企业市场**
 - 模型能力：模型路线从分化到统一
 - 商业模式：C端产品领导者，发力企业市场
 - 财务表现：收入预测持续上修，预计2030年超过2800亿
 - 算力储备：星际之门项目持续推进
- **四、静态理解模型的商业化市场空间**

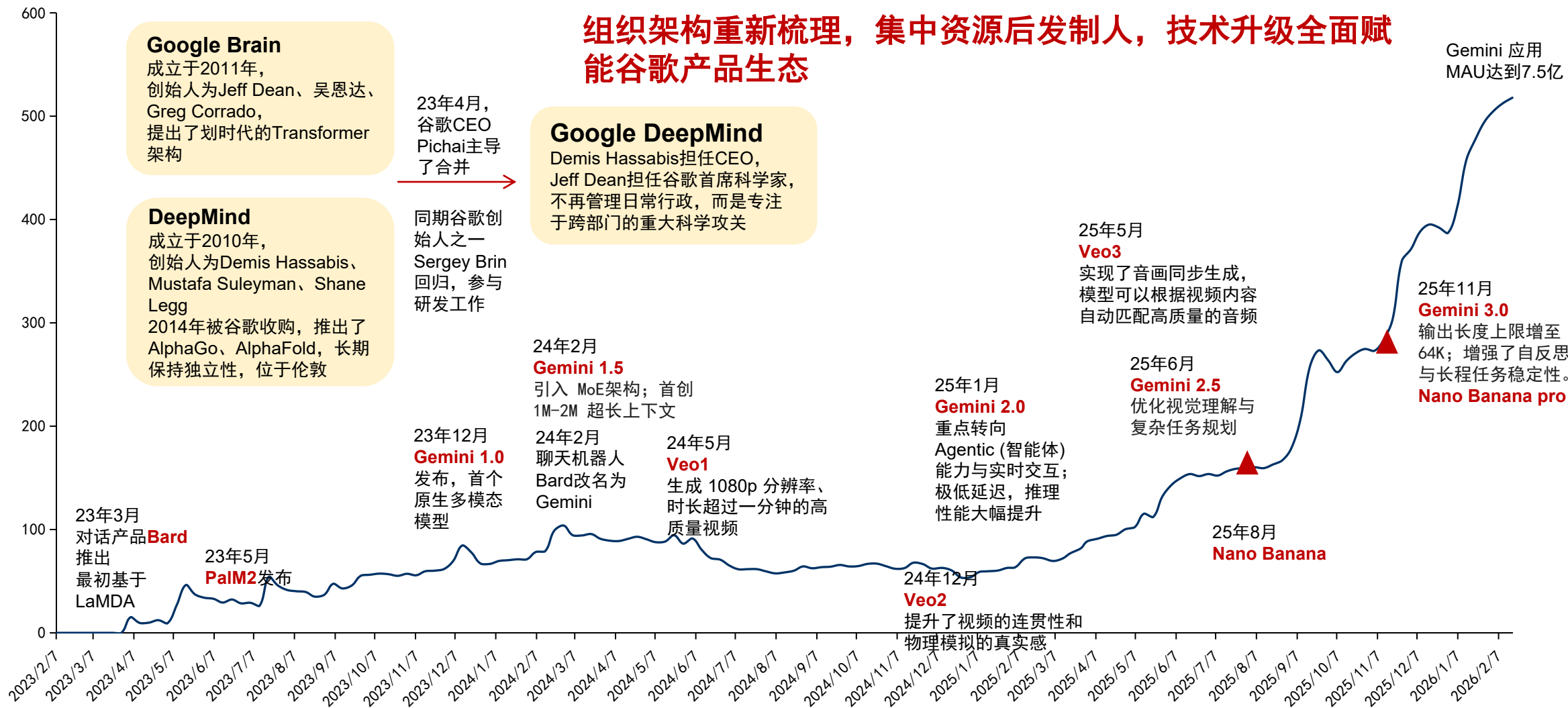
谷歌：围绕多模态能力打造模型矩阵

- 谷歌的模型分为两条路线：1) 原生多模态模型Gemini，围绕Gemini的多模态能力进一步衍生出image生图能力，从token价格来看，Gemini一直处于提价周期中；2) 垂类内容制作模型：生图模型imagen、视频模型Veo、音乐模型Lyria、世界模型Genie，其中imagen系列与nano banana系列在模型架构上天然存在差异（transformer和diffusion）。

图：谷歌Gemini系列模型迭代情况及价格变化

模型版本	发布时间	型号	重要代际变化	输入价格（每百万tokens）	输出价格（每百万tokens）
Gemini 1.0	2023.12	pro	原生多模态架构（文本/图/音/视）	0.5	1.5
		ultra		0.5	1.5
Gemini 1.5	2024.02	flash	引入 MoE（混合专家）架构；首创 1M-2M 超长上下文；推出高性能 Flash 版。	0.075（文本长度≤128k） 0.15（文本长度>128k）	0.3（文本长度≤128k） 0.6（文本长度>128k）
		pro		1.25（文本长度≤128k） 2.5（文本长度>128k）	5（文本长度≤128k） 10（文本长度>128k）
Gemini 2.0	2025.01	flash-lite	重点转向 Agentic（智能体）能力与实时交互；极低延迟，推理性能大幅提升。	0.075	0.3
		flash		0.1	0.4
		pro		1.25	5
Gemini 2.5	2025.06	flash-lite	优化视觉理解与复杂任务规划；推出 Deep Think 模式应对逻辑难题。	0.1	0.4
		flash		0.3	2.5
		pro		1.25（文本长度≤200k） 2.5（文本长度>200k）	10（文本长度≤200k） 15（文本长度>200k）
	2025.08	flash image (nano banana)	0.30（文字 / 图片）	30	
Gemini 3.0	2025.11	flash	新一代基础架构；输出长度上限增至 64K；增强了自反思与长程任务稳定性。	0.50（文字 / 图片 / 视频） 1.00（音频）	3
		pro		2（文本长度≤200k） 4（文本长度>200k）	12（文本长度≤200k） 18（文本长度>200k）
		pro image (nano banana pro)		2.00（文字/图片）	12.00（文字和思考） 120.00（图片）
Gemini 3.1	2026.02	pro	最新版本。改进了自定义工具调用（Custom Tools）的优先级与精准度。	2（文本长度≤200k） 4（文本长度>200k）	12（文本长度≤200k） 18（文本长度>200k）

图：Gemini 网站周访问量（百万）



数据来源：Anthropic、The Information、国信证券经济研究所整理

- 25年开始Gemini从Gemini 2.5到Gemini 3逐渐成为模型性能综合榜单的领跑者，目前最新的模型Gemini 3.1也在Artificial Analysis综合了多维度的十项评分中排名第一。

图：主流模型综合智能 及 Agent能力比较

Artificial Analysis Intelligence Index

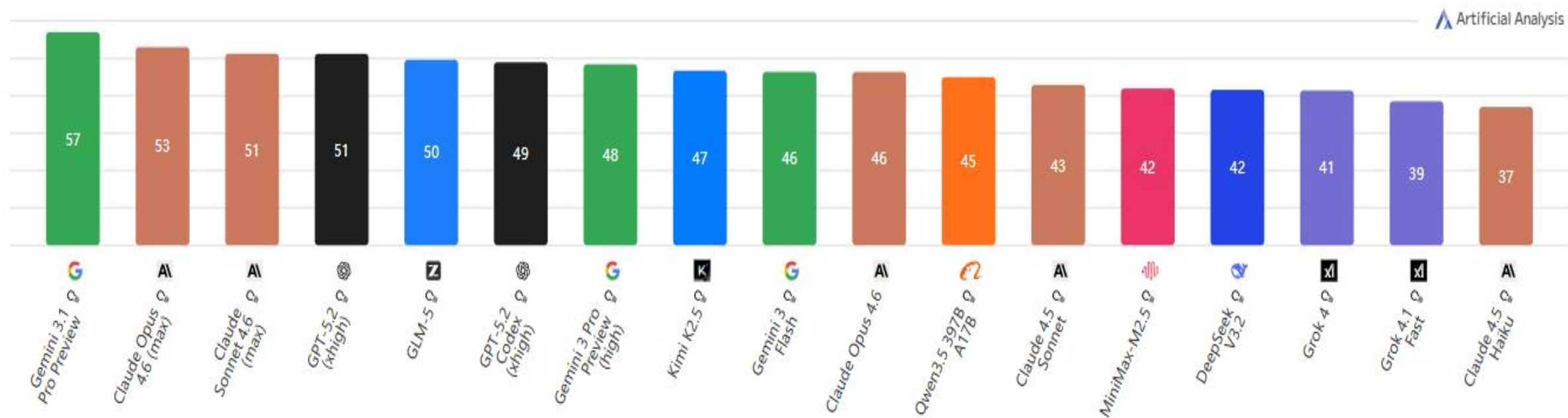
Artificial Analysis Intelligence Index v4.0 incorporates 10 evaluations: GDPval-AA, τ^2 -Bench Telecom, Terminal-Bench Hard, SciCode, AA-LCR, AA-Omniscience, IFBench, Humanity's Last Exam, GPQA Diamond, CritPt



17 of 410 models ×



+ Add model from specific provider



资料来源：Artificial Analysis、国信证券经济研究所整理测算

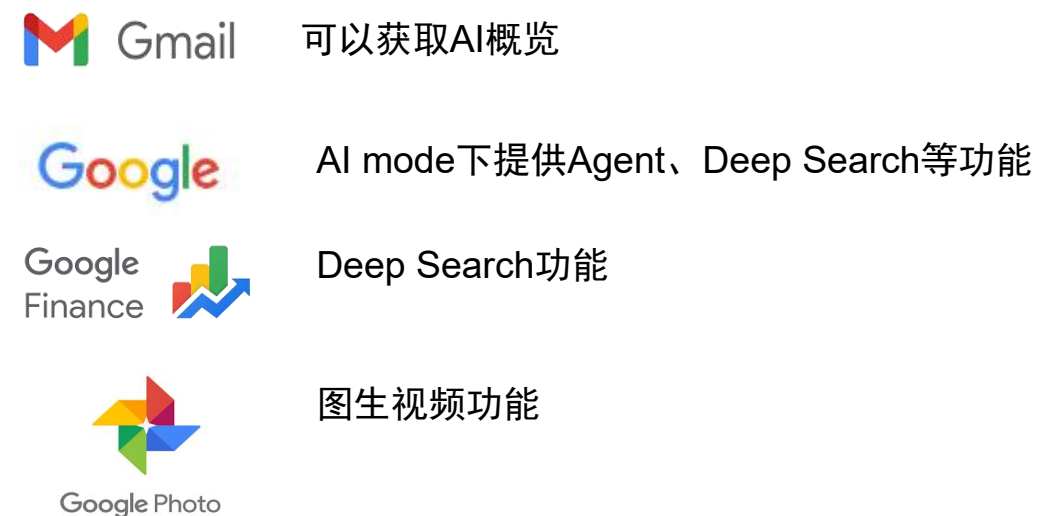
- 谷歌通过Google AI plans的订阅模式对外提供一揽子AI相关的产品服务，其丰富度远高于Anthropic，除了原生的AI产品 Chatbot Gemini、AI视频 Flow、AI图像 Whisk、AI编程 Antigravity外，还有多个功能集成于谷歌体系内的原有产品中，打造原生AI应用+老应用AI化的产品矩阵。相关订阅收入计入Google Subscriptions, Platforms, and Devices中。

图：谷歌C端 Google AI plans的产品矩阵

AI原生产品



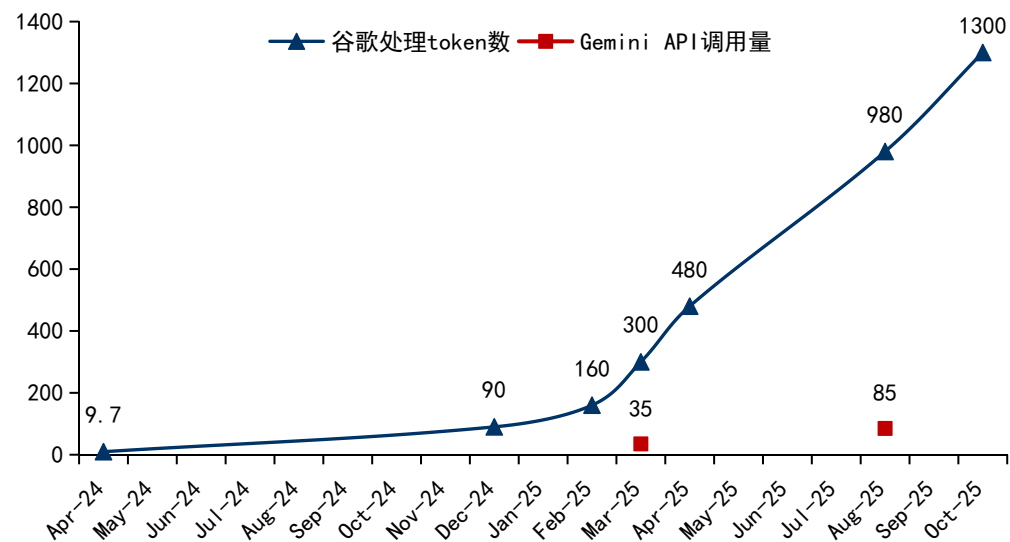
传统产品接入AI功能



商业模式： Gemini系列API调用拉动谷歌云收入快速增长

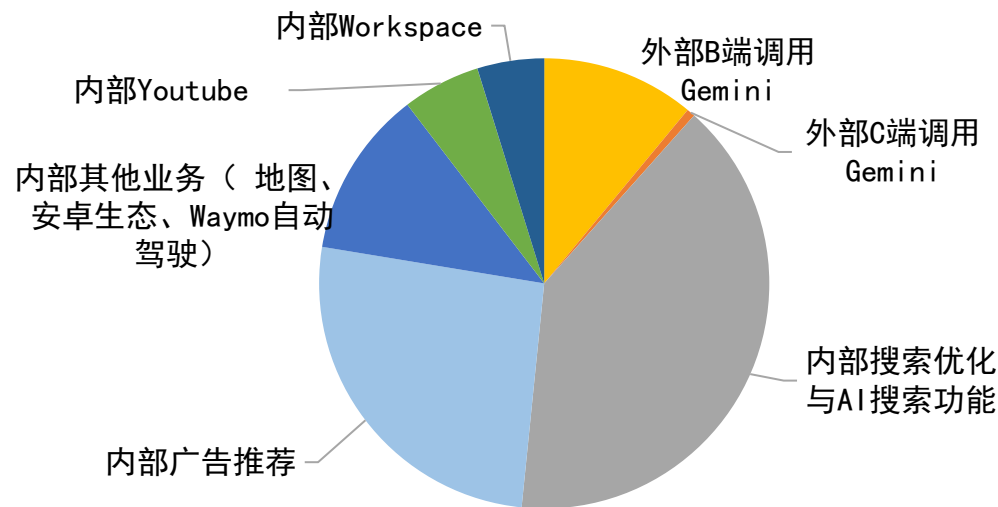
- 谷歌对外提供API服务主要通过两个渠道：Google AI Studio和Google Cloud Vertex AI，相关收入计入谷歌云中，其中Google AI Studio主打快速上手和原型开发，Google Cloud Vertex AI强调稳定、安全和大规模集成，专门为企业级生产环境准备，Vertex AI除了针对使用模型按照token收费，还会根据提供的服务如计算节点、智能搜索、数据存储等进行收费。
- 根据The Information， Gemini的API调用量25年快速增长，从3月的350亿增长至8月的850亿，在谷歌Token消耗中的占比约为10%。

图：谷歌月度tokens消耗及Gemini API调用tokens消耗



资料来源：公司官网，国信证券经济研究所整理测算

图：谷歌月度tokens调用量场景占比

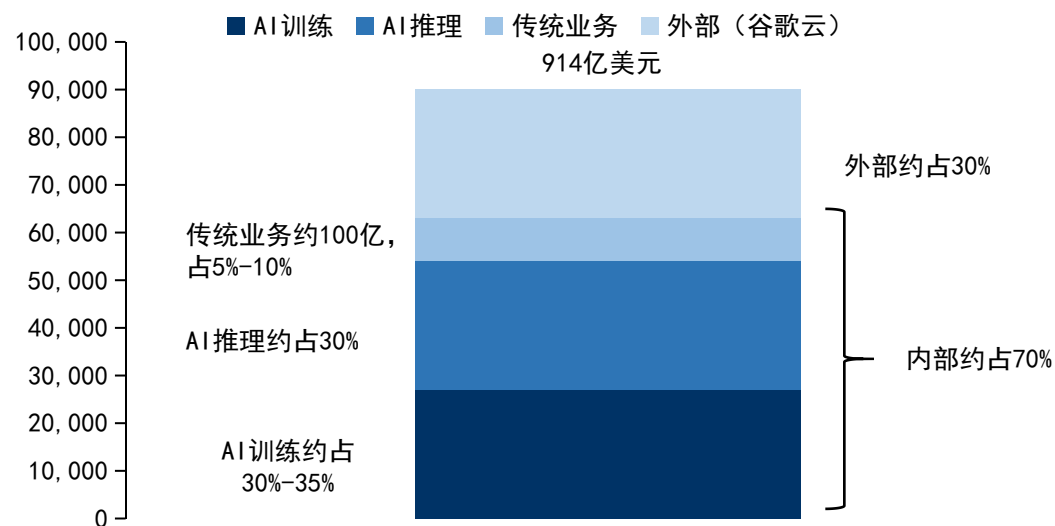


资料来源：公司官网，国信证券经济研究所整理测算

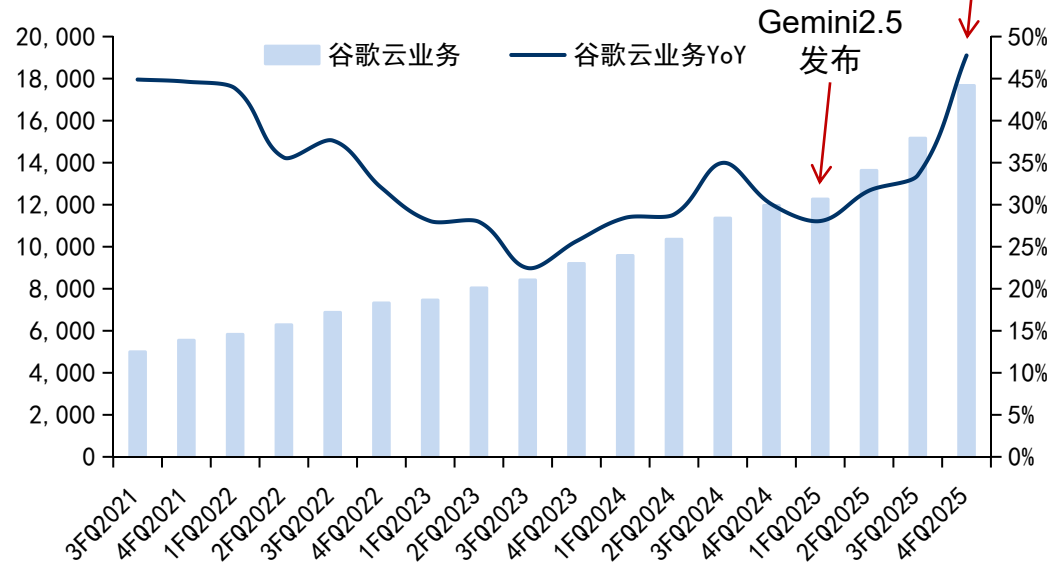
谷歌算力资源分配：内部占比较高，外部通过云提供服务

- 谷歌2025年CAPEX达到914亿美元，我们预计其中外部使用（谷歌云）约占30%，内部使用主要用于模型训练、推理以及传统业务。
- ① 谷歌内部模型训练：结合OpenAI、Anthropic等头部模型公司情况，尽管25年推理需求快速增长，我们预计Gemini训练相关CAPEX占比仍略高于推理。
- ② 谷歌内部应用推理：Gemini模型正在全面赋能谷歌各项业务，包括部署搜索算法优化与AI搜索功能、广告推荐系统改造、Youtube、Workspace、地图、安卓生态等，提升用户体验、交互效率，同时2025Q4 Gemini MAU达到7.5亿，Gemini 3 pro日均token消耗是2.5的3倍，Gemini的token消耗中，API调用每分钟处理量从25Q3的70亿增长至100亿。
- ③ 谷歌传统业务：谷歌传统业务自身同样需要一定的算力投入，比如广告业务的推荐效率改善。
- ④ 外部（谷歌云GCP）：我们测算谷歌CAPEX中约30%用于外部，通过谷歌云给客户提提供算力支持。

图：谷歌CAPEX结构



图：谷歌云收入及增速



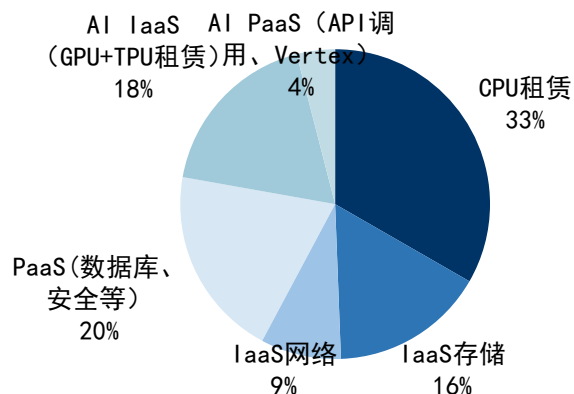
资料来源：彭博，国信证券经济研究所整理测算

资料来源：公司财报，国信证券经济研究所整理

谷歌云：AI云驱动收入加速增长，26预计API收入将实现爆发式增长

- AI云已成为谷歌云增长主要动力。25年谷歌云收入同比+36%，我们测算其中AI云有望实现同比+187%，收入占比达到22%。GPU和TPU租赁相关收入占比将达到18%，同比+135%，26年将延续这一趋势；API调用相关收入占比达到4%，并且预计26年将实现爆发式增长达到25年的4-5倍。
- 当前算力的需求方仍以头部公司为主，1) GPU/TPU租赁业务中，大型互联网公司为主要客户。我们预计TPU的收入主要来自苹果和Anthropic，小部分来自Safe Superintelligence和Physical Intelligence等人工智能实验室，GPU租赁收入则主要来自Meta、字节、腾讯等公司。
- 2) API调用收入客户结构相对分散。据公司25Q3业绩会，过去12个月，近150家Google Cloud客户使用我们的模型处理了约1万亿个tokens，用于各种应用。

图：谷歌云收入拆解



资料来源：公司财报，国信证券经济研究所测算整理测算

图：谷歌云收入拆解

单位：百万美元	2024年	2025年	2026E
云收入	43,229	58,705	87,399
YoY		36%	49%
增长绝对值		15,476	28,694
传统云	38,729	45,790	52,658
YoY		18%	15%
CPU租赁	12,000	14,000	17,000
YoY		17%	21%
AI云	4,500	12,915	34,741
YoY		187%	169%
占收比		22%	40%
GPU和TPU租赁	4,500	10,567	23,000
YoY		135%	118%
TPU %		40%	55%
GPU %		60%	45%
占收比		18%	26%
API调用		2,348	11,741
占收比		4%	13%

资料来源：公司财报，国信证券经济研究所测算整理测算

- Google Workspace 是 Google 专为企业、学校和个人组织提供的一套基于云端的全功能办公协作套件，24年谷歌将AI功能作为 Workspace的插件提供Gemini Business和Gemini Enterprise两个服务给客户，对其进行额外收费，25年开始，谷歌直接将Gemini的核心功能整合进Workspace各个产品中，停止销售单独插件，并对Workspace进行了进行了约 15%-20% 的价格上调。
- 2025年10月谷歌重新推出了一个全新的、独立的 Gemini Enterprise 平台，新平台包含了Gemini模型、无需编写代码的工作台、预构建的 Google Agent、可以安全连接到客户公司的数据，并且拥有超过 10 万个合作伙伴。
- **截至25Q4末， Gemini Enterprise 已向 2800 多家公司售出超 800 万个付费席位，包括纽约梅隆银行（BNY）、维珍邮轮（Virgin Voyages）等， 25Q4 Gemini Enterprise 处理了超 50 亿次客户互动，同比增长 65%。**
- **寻找外部合作伙伴拓展企业客户：**1) 在Gemini Enterprise 中集成第三方Agent产品的功能，比如ServiceNow AI Agent Fabric、基于Agentforce构建的Agent等；2) 利用谷歌云的合作（分销）伙伴资源，依靠埃森哲、德勤、Cognizant等帮助拓展客户

图：谷歌Gemini Enterprise付费权益



The image shows two pricing cards for Gemini Enterprise. The left card is for the 'Business' plan, which is described as '非常适合小型企业和组织内的团队。无需进行 IT 设置。' The price is '\$21 美元/席位/月'. The right card is for the 'Standard/Plus' plan, described as '非常适合需要企业级 IT 控制的大型组织'. The price is '\$30 美元/席位/月'. Both cards include a '开始 30 天试用' button and a '与销售人员联系' link. The Business plan lists several key features like integration with Google Workspace and Microsoft 365, and the Standard/Plus plan lists features like higher quotas and VPC Service Controls.

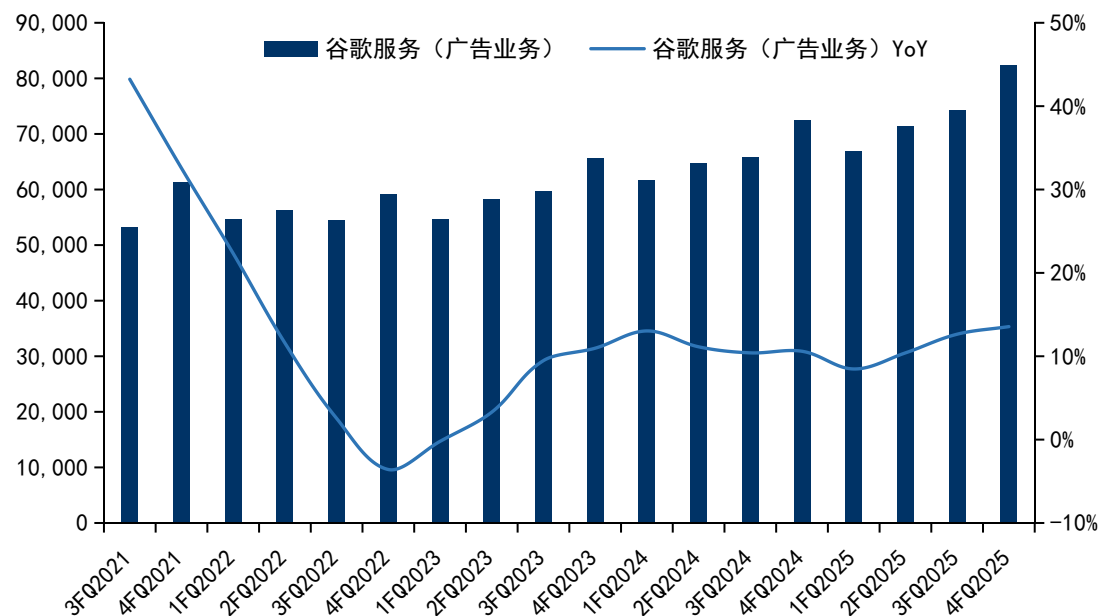
版本	描述	起价
商务版	非常适合小型企业和组织内的团队。无需进行 IT 设置。	\$21 美元/席位/月
标准版/Plus 版	非常适合需要企业级 IT 控制的大型组织	\$30 美元/席位/月

资料来源：公司官网，国信证券经济研究所整理

AI 带动广告转化率改善，广告收入稳健增长

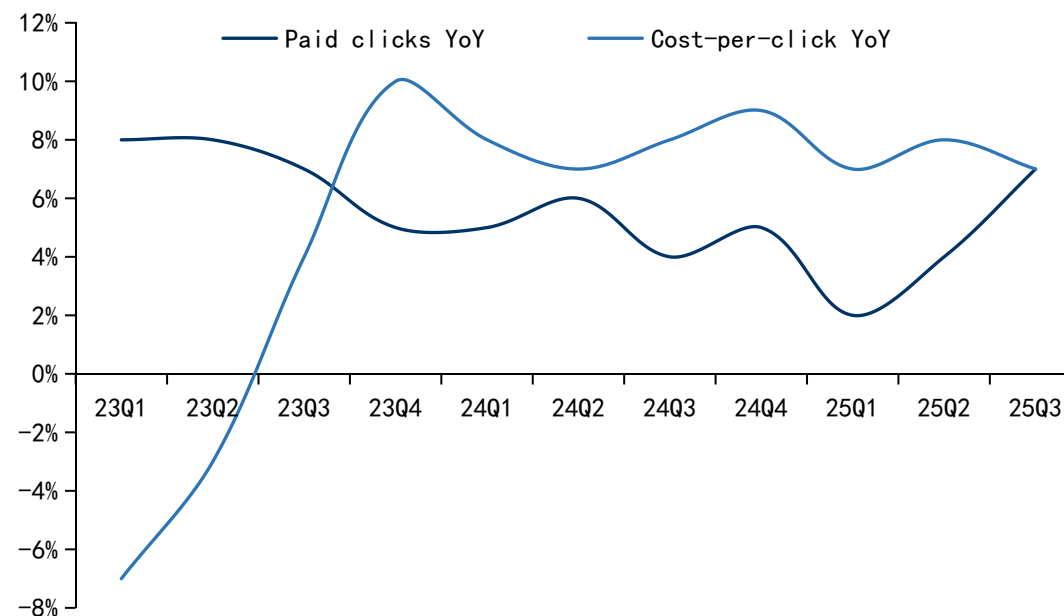
- AI 对谷歌广告产品的持续升级带来广告主转化效率的不断提升（24年以来新功能推出后通常能带来20%左右的效率改善），比如1) 启用 AI Max 的广告客户通常能获得 14% 的额外转化量；2) 使用智能竞价探索（Smart Bidding Exploration）的广告系列平均转化量增长 19%；3) 目前已有超过 200 万广告客户使用谷歌的人工智能驱动资产生成工具投放广告，较去年同期增长 50%。

图：谷歌广告收入情况



资料来源：公司财报，国信证券经济研究所整理

图：谷歌搜索广告单价和付费点击增速情况



资料来源：公司财报，国信证券经济研究所整理

AI对谷歌广告业务的带动主要体现在自动生成广告素材功能、智能出价功能、智能匹配几个方面。推出的产品包括 DemandGen（基于推荐而不是搜索的AI广告产品，旧版叫Discovery ads）、AI MAX（利用AI扩展关键词和潜在用户、素材优化）、Smart Bidding Exploration（智能出价，通过AI触达更多流量）。

图：24年起AI对谷歌广告业务带来变化

时间	变化	效果
24Q1	将Gemini 功能引入PMax 平台，帮助企业生成文本图片素材	使用 P-Max 素材生成功能的广告主发布广告效力良好或卓越的广告的可能性提高了 63%。而那些将 PMax 广告效力提升至卓越的广告主， 平均转化次数可提高 6% 。采用ACA（自动生成广告素材）功能的企业在 Search 和 PMax 广告中，平均每次转化费用相似， 但转化次数提高了 5% 。
24Q2	由人工智能驱动的利润优化工具已扩展到PMax广告和shopping广告	与仅关注收入的竞价相比，使用利润优化和智能出价的广告主平均利润提升了 15%。
	DemandGen	DemandGen 与 Search 或 PMax 结合使用时，平均可提升 14% 的转化率。
		Google 营销团队使用 DemandGen 为 Pixel 8 广告活动创建了近 4,500 个广告变体，该广告活动在 YouTube、Discover 和 Gmail 上展示，点击率提高了一倍，成本却降低了近四分之一。
24Q4	向客户正式开放了我们的营销组合模型 Meridian，帮助更多企业将资金重新投入到他们熟知的创意和媒体购买策略中。	基于谷歌 AI 技术的 YouTube 视频广告系列的广告支出回报率比手动广告系列高出 17%。
25Q1	在 Demand Gen 中，广告主可以更精准地管理全球 YouTube、Gmail、Discover 和 Google 展示广告网络的广告投放，并了解哪些资产在渠道层面效果较优。	使用 Demand Gen 的企业现在在购买和潜在客户等目标上， 每美元支出的转化率平均同比增长 26% 。而当 Demand Gen 与产品 Feed 结合使用时，平均而言，每美元支出的转化率同比增长了一倍以上。
25Q2	推出搜索广告的AI MAX功能	广告主通常会获得 14% 的转化率提升。
	智能出价探索，允许广告主更频繁地对不太显眼但潜在价值更高的查询进行竞价	使用智能出价探索的广告系列平均转化率提升了 19%。
	创意方面，推出了 Asset Studio，帮助企业生成创意素材	超过 200 万广告主使用 Google 的 AI 素材生成工具来投放广告，比去年同期增长了 50%。
25Q3	AI MAX 9月全球上线	已有数十万广告主在使用，仅本季度，该功能就带来了数十亿次新增搜索请求
25Q4		Q4广告主通过 AI Max 和 Pmax 的文本定制功能，借助 Gemini 制作了近 7000 万份营销素材。

- **一、Anthropic：凭专业生产能力打造高毛利护城河**
 - 核心团队成員、经营理念、算力储备
 - 模型能力：Coding、Agent场景下的SOTA模型
 - 商业模式：极简产品矩阵，API贡献主要收入
 - 财务表现：最强模型带来token溢价，28年有望迎来现金流转正
- **二、谷歌：多模态能力领先，生态优势明显**
 - 模型能力：围绕多模态能力打造模型矩阵，综合性能领先
 - 商业模式：原生AI应用+Gemini赋能传统产品，云与广告受益增长
- **三、OpenAI：C端产品领导者，开始发力企业市场**
 - 模型能力：模型路线从分化到统一
 - 商业模式：C端产品领导者，发力企业市场
 - 财务表现：收入预测持续上修，预计2030年超过2800亿
 - 算力储备：星际之门项目持续推进
- **四、静态理解模型的商业化市场空间**

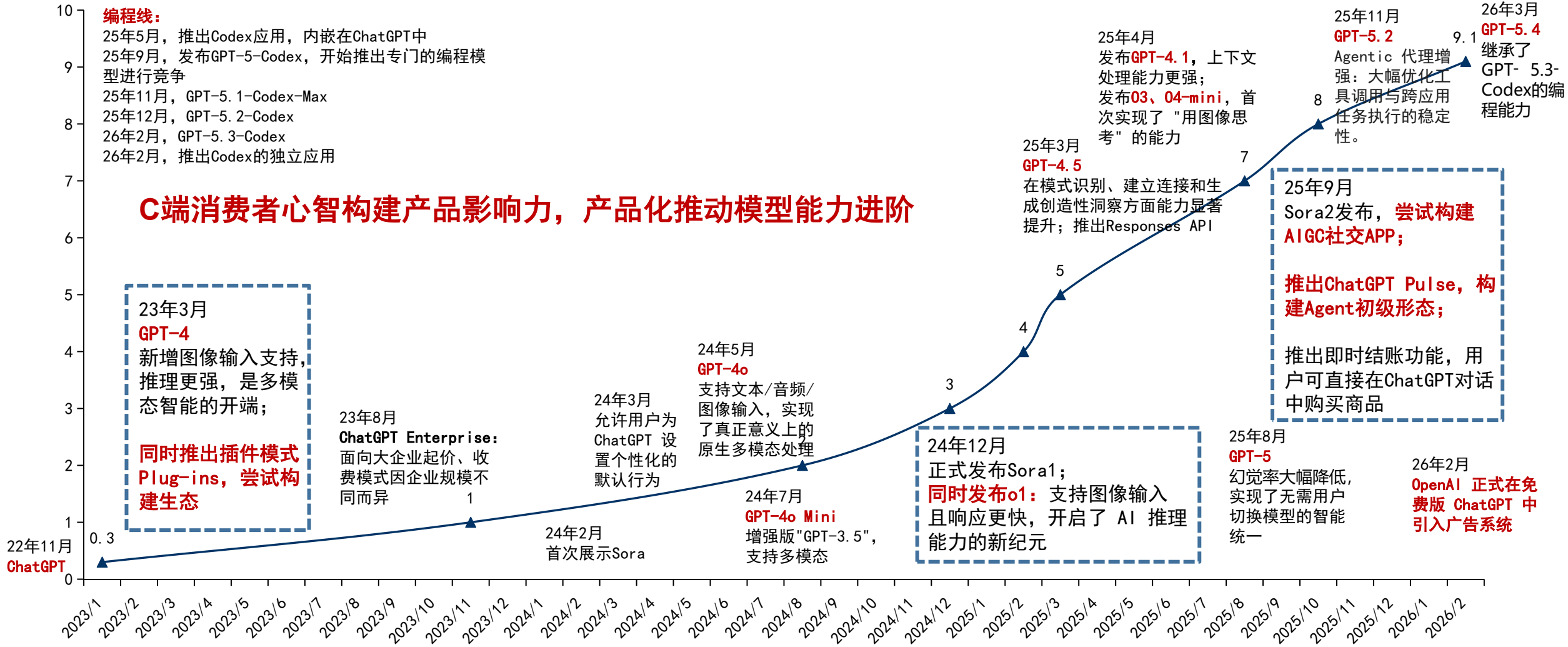
OpenAI：模型路线从分化到统一

- 作为LLM的引领者，OpenAI的路线曾发生过变化，在23年初推出GPT4后，开始走模型分工路线，陆续推出了GPT4o（多模态）、o1/o3（推理模型），但由于过于复杂的模型路线，从GPT5开始，又将各种模型能力集成到一个模型中。

图：谷歌Gemini系列模型迭代情况及价格变化

模型版本	发布时间	型号	重要代际变化	输入价格（每百万tokens）	输出价格（每百万tokens）
GPT-1	2018.06		开创性架构：首次验证 Transformer 的“预训练+微调”范式。		
GPT-2	2019.02		零样本学习：参数量提升10倍；展示了强大的通用文本生成潜力。		
GPT-3	2020.06		规模效益（Scaling）：1750亿参数；确立了 Few-shot（少样本）提示工程模式。	2	2
GPT-3.5	2022.11		RLHF 革命：引入人类反馈强化学习；推出 ChatGPT，开启 AI 平民化时代。	0.5	1.5
GPT-4	2023.03		逻辑质跃：支持图像输入；复杂推理与编程能力大幅超越前代。	30	60
	2023.11	turbo		10	30
GPT-4o	2024.05	mini	原生全模态：实现文字/音/画端到端实时交互；响应速度接近人类。	0.15	0.6
				2.5	10
OpenAI o1	2024.09	mini	强化学习推理：引入“思维链”（CoT）；在数学、奥数与代码领域表现卓越。	1.1	4.4
		pro		15	60
GPT-4.1	2025.04	nano	端侧优化：推出高性能轻量级版本，主要适配移动端及离线场景。	150	600
		mini		0.1	0.4
				0.4	1.6
OpenAI o3	2025.04	mini	深度推理巅峰：解决 AIME 与 Codeforces 等顶尖难题；推理能力进一步通用化。	2	8
				1.1	4.4
		pro		2	8
GPT-5	2025.08	nano	System 3 架构：具备长期记忆（Long-term Memory）与自我纠错；支持 1M 超长上下文。	20	80
		mini		0.05	0.4
		pro		0.25	2
GPT-5.1	2025.11		引入了自适应推理能力	1.25	10
GPT-5.2	2026.02		Agentic 代理增强：大幅优化工具调用（Tool Call）与跨应用任务执行的稳定性。	1.75	14
		pro		21	168
GPT-5.4	2026.03		原生电脑操作能力；支持工具搜索；100万token的上下文窗口	2.5	15
		pro		30	180

图：ChatGPT周活变化



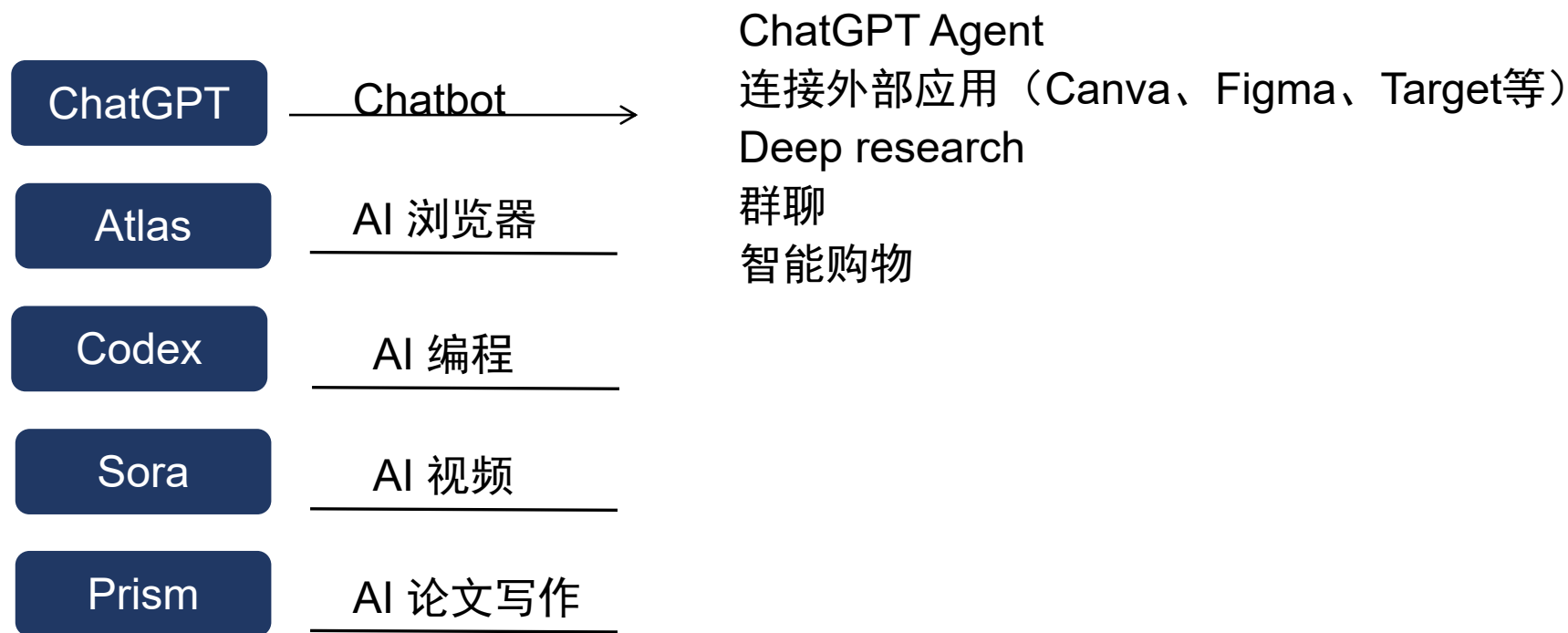
数据来源：ChatGPT官网、CNBC、国信证券经济研究所整理

请务必阅读正文之后的免责声明及其项下所有内容

商业模式：率先开拓“消费者优先”2C战略，ChatGPT最早获得全球最多C端AI用户

- OpenAI从一开始就坚持“消费者优先”的2C战略。OpenAI围绕C端订阅同样提供了矩阵式的产品服务，但生态的丰富度略弱于谷歌，核心仍然依靠ChatGPT吸引C端用户，目前在ChatGPT内部已经打造了丰富的软件功能，覆盖Agent、购物、深度研究等，同时打通了多款外部应用。
- 记忆是OpenAI的护城河。奥特曼在12月17日Big Technology Podcast的采访中表示，相比模型性能更重视粘性，认为chatgpt的用户粘性和护城河来自习惯、一致性和个人体验积累出的信任感。由于ChatGPT最先抢占市场构建起了超过8亿用户的基本盘，用户在ChatGPT有更多的历史交流记录，也意味着回答能够更贴近用户的需求。

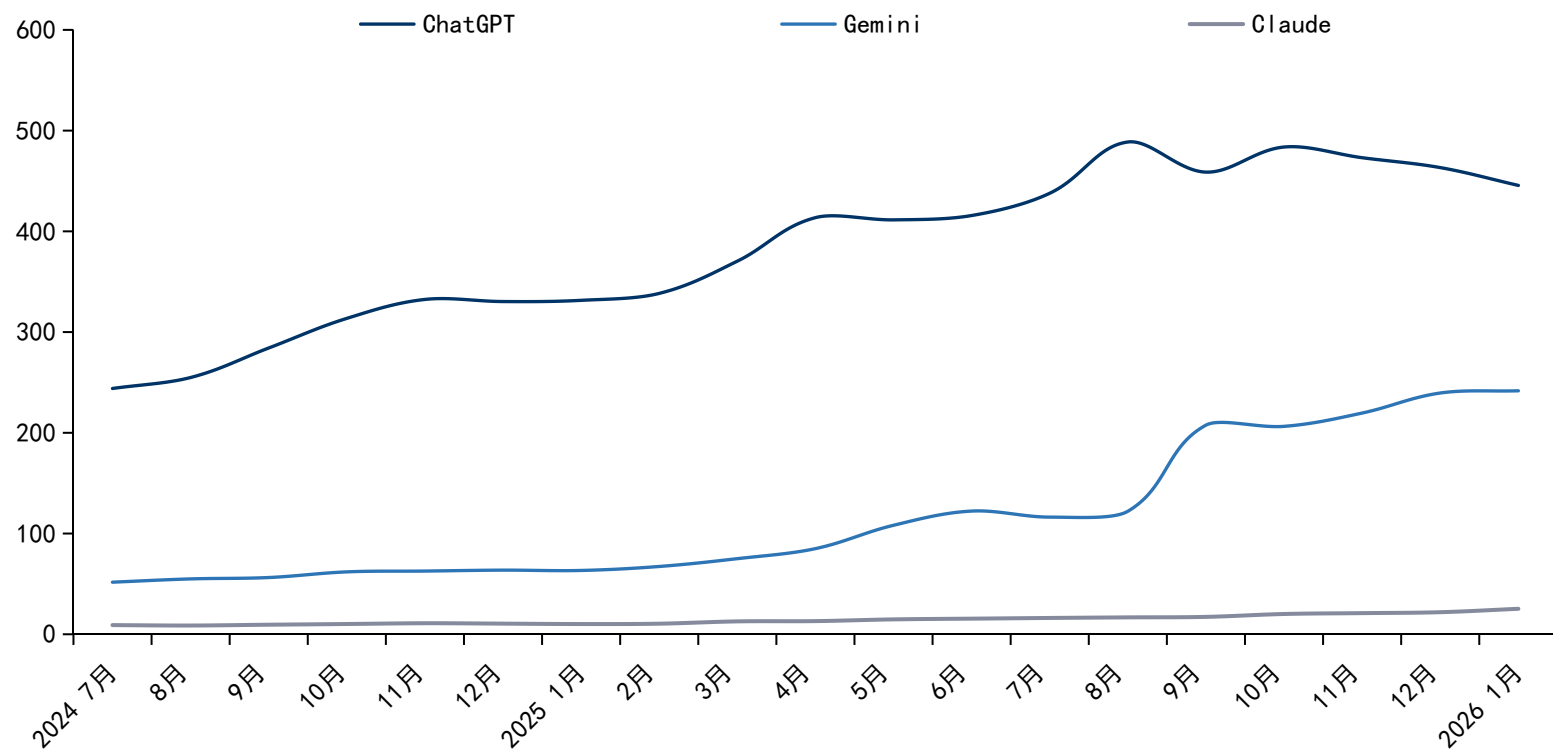
图：OpenAI C端订阅产品矩阵



用户：伴随其他模型厂商快速崛起，OpenAI月活领先优势逐渐缩小

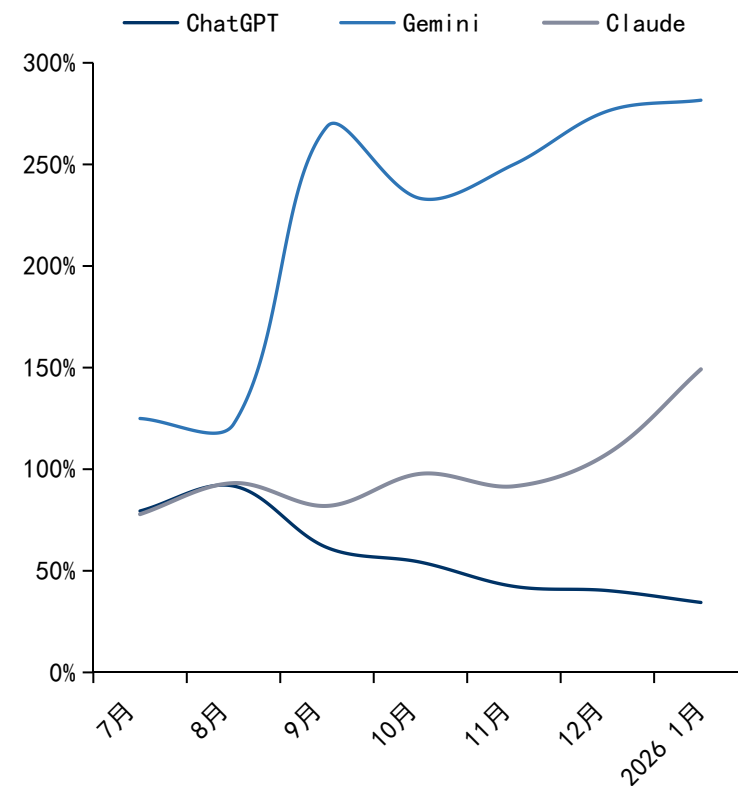
- 作为25年最重要的2款模型，下半年伴随Gemini3、Claude Opus 4.5的发布，Gemini、Claude的用户规模开始加速增长，尤其是Gemini正快速追赶ChatGPT用户，网页端MAU已达到其大约一半的水平，而同期ChatGPT的用户增长则逐步放缓。

图：ChatGPT、Gemini、Claude网站MAU对比（百万）



资料来源：Similar Web、国信证券经济研究所整理测算

图：ChatGPT、Gemini、Claude网站MAU同比增速



资料来源：Similar Web、国信证券经济研究所整理测算

商业模式：同样提供API形式，客户多为传统SAAS企业

- 据量子位，OpenAI调用Tokens超过1万亿的客户名单中多数为SaaS类公司。其中主要包括协同办公、开发者工具（coding）、电商（客服、自动生成商品内容、广告优化）、搜索推荐、数据分析等场景。
- 参考当前GPT5价格，输入价格为1.25美元/百万tokens，输出价格为10美元/百万tokens，通常输入的内容量远低于输出，按照10美元/百万tokens单价计算，30位调用量1万亿的客户贡献收入约3亿美元。
- 模型技术的快速发展为SaaS行业带来了颠覆性的机会，一方面模型厂商本身能够提供平台给企业自己构建Agent产品，另一方面基于AI的SaaS类产品在使用体验上有明显提升。后续云厂/模型厂的API调用收入增长，也依赖于Agent等产品的爆发。

图：OpenAI tokens调用量超过1万亿的客户

公司	业务简介
Duolingo	语言学习 App, 以游戏化课程设计著称
OpenRouter	多模型聚合平台, 统一 API 调用 GPT、Claude、Gemini 等模型
Indeed	全球最大招聘网站之一, 提供求职与招聘服务
Salesforce	企业级 CRM 与云计算服务提供商
CodeRabbit	AI 代码审查与自动化开发平台
iSolutionsAI	企业 AI 解决方案与自动化服务提供商
Outtake	AI 视频与内容生成平台
Tiger Analytics	数据分析与 AI 咨询公司
amp	企业费用与支出管理平台
Abridge	医疗 AI 公司, 自动转录与总结医患对话
Sider AI	AI 代码评审与文档生成工具
Warp.dev	开发者终端软件, 集成 AI 命令补全与自动化
Shopify	全球电商平台, 为中小企业提供建站与支付解决方案
Notion	生产力工具, 集笔记、任务、知识库于一体
WHOOP	可穿戴设备公司, 专注健康与运动数据监测
HubSpot / Dashworks	营销自动化与 CRM 平台, 集成 AI 知识搜索
JetBrains	软件开发工具公司, 包含 IntelliJ IDEA、PyCharm 等
Delphi	企业 AI 助理平台, 专注知识问答与自动化决策
Decagon	AI 企业服务平台, 提供商业分析与客户交互工具
Rox	初创公司, 开发 AI 社交与交互产品
T-Mobile	美国大型电信运营商, 布局 AI 客服与网络优化
Zendesk	客服系统与 AI 客服解决方案提供商
Harvey	AI 法律助手, 服务律师事务所与企业法务
Read AI	会议分析工具, 自动生成会议纪要与洞察
Canva	在线图形设计平台, 整合 AI 图像生成与排版
Cognition	AI 代码生成公司, 打造可编程智能体
Datadog	云监控与数据可视化平台
Perplexity	AI 原生搜索引擎, 基于 LLM 生成答案
Mercado Libre	拉美最大电商与支付平台
Genspark AI	企业级 AI 研发公司, 提供 AI 产品与工具平台

资料来源：OpenAI、国信证券经济研究所整理

商业模式：26年将发力企业业务

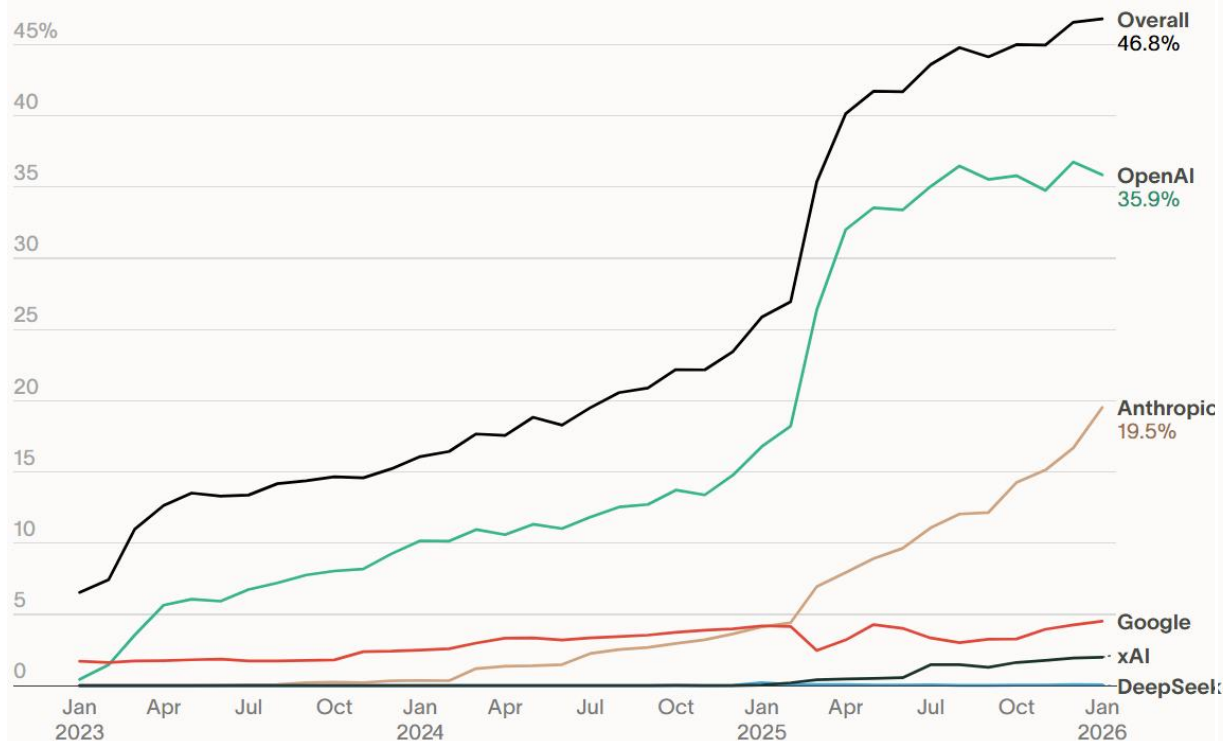
- OpenAI 23年推出ChatGPT Enterprise，面向大型机构，提供企业级安全、隐私和管理功能，24年起为中小型团队推出ChatGPT Team，后改名ChatGPT Business。
- 2025年11月OpenAI宣布全球企业客户突破100万（包括企业订阅服务和API），截至25年11月已有超过700万个ChatGPT for Work席位，仅在两个月内就增长了40%。ChatGPT Enterprise席位的数量同比增长了9倍。
- 25年末OpenAI CEO Sam Altman强调，26年OpenAI会发力企业业务，从25年末开始，OpenAI持续寻找第三方合作伙伴进行战略布局。1) 与Snowflake、ServiceNow等公司合作在其产品中嵌入GPT模型；2) 26年2月公司宣布成立Frontier Alliances，与波士顿咨询集团（BCG）、麦肯锡公司、埃森哲和凯捷将携手帮助客户制定战略、整合系统、重塑工作流程，并在全球范围内扩展部署。
- 根据Ramp AI Index数据，截至26年1月OpenAI在美国企业中模型采用率是36%。（根据Ramp统计的企业支出数据）

图：美国企业模型采用率

Ramp AI Index: Model Adoption Rate

Share of U.S. businesses with paid subscriptions to AI models, platforms, and tools

View by Overall Sector Size Model



Source: Ramp AI Index, business spend data from Ramp. Overall includes businesses subscribed to any AI product or service based on Ramp spend data. • Get the data • Embed • Download image

ramp

资料来源：Ramp、国信证券经济研究所整理

财务表现：25年收入131亿美元，预计2030年超过2800亿

- 收入端：OpenAI当前收入以C端订阅为主，我们预计OpenAI 2025年收入131亿美元，其中C端订阅约占65%，目前整体用户中付费渗透率约5%，付费用户约4000万+。26年起公司将持续发力B端业务，同时在C端ChatGPT中引入广告，2030年在多种变现手段驱动下整体收入达到2800亿美元。
- 利润端：由于API定价低于Anthropic等因素，OpenAI毛利率水平更低，25年仅为33%，预计26年开始持续回升未来有望达到接近70%，2030年有望实现OP转正。

图：OpenAI财务情况预测

单位：百万美元	2024A	2025A	2026E	2027E	2028E	2029E	2030E
OpenAI收入	3,700	13,100	30,000	62,000	113,000	184,000	284,000
YoY		254%	129%	107%	82%	63%	54%
ChatGPT (C端, 含广告)	2,700	8,500	17,000	35,000	58,000	94,000	150,000
YoY		215%	100%	106%	66%	62%	60%
ChatGPT (B端业务)		2,000	8,000	15,000	32,000	48,000	70,000
YoY			300%	88%	113%	50%	46%
API调用	1,000	2,600	4,900	10,500	18,000	30,000	47,500
YoY		160%	88%	114%	71%	67%	58%
其他收入 (硬件等)			100	1,500	5,000	12,000	16,500
YoY				1400%	233%	140%	38%
营业成本 (推理成本)	2,220	8,777	14,400	27,280	45,200	66,240	93,720
毛利	1,480	4,323	15,600	34,720	67,800	117,760	190,280
毛利率	40%	33%	52%	56%	60%	64%	67%
销售&管理费用	1,000	6,157	8,100	11,780	14,690	18,400	28,400
费用率	27%	47%	27%	19%	13%	10%	10%
研发费用 (训练成本)		8,300.00	32,000	65,000	90,400	119,600	127,800
费用率		63%	107%	105%	80%	65%	45%
经营利润		-10,134	-24,500	-42,060	-37,290	-20,240	34,080
OPM		-77%	-82%	-68%	-33%	-11%	12%

资料来源：The information、国信证券经济研究所整理测算

OpenAI 算力建设: 星际之门以及硬件厂商合作计划

- **星际之门(Stargate) 项目:** 2025 年1月21日, 特朗普与OpenAI正式宣布启动, 新成立Stargate公司, 软银、Oracle、OpenAI、阿联酋MGX等计划在未来四年内投资5000亿美元并获取Stargate股权, 目标建设10GW电力容量的AI基础设施(目前已确认锁定甲骨文5.5GW+软银1.5GW)。首期1000亿美元已到账或签署托管协议, 后三期靠绿色债券滚动和战略股东跟投。
- **其他外部采购:** 25年10月与微软签订了多年2500亿美元的Azure服务采购, 11月与 AWS签署了一项7年、约380亿美元的合作协议。
- **硬件厂商合作计划:** 25年9月英伟达承诺向OpenAI的基础设施投资1000亿美元, 锁定未来10GW数据中心的芯片供应权。25年10月6日, OpenAI将部署6GW AMD GPU, AMD向OpenAI授予最多1.6 亿股普通股的认股权证。

表: 星际之门合作方与参与方式

资金来源	金额 (亿美元)	性质	备注
软银愿景基金	1500	股权+可转债	牵头方, 分5期认购普通股
阿联酋MGX+ADQ	750	股权	MGX 600 亿、ADQ 150 亿, 换取 15% 优先股
甲骨文	400	股权 + 设备抵股	200 亿现金+200 亿云基础设施设备折价入股, 获取3000亿美元数据中心承建权, 超5.5GW
英伟达	200	设备抵股	以 GB200 芯片组、NVLink 交换机折价入股。不包括25年9月英伟达承诺向OpenAI投资1000 亿美元, 这是在 "星际之门" 项目框架之外的独立合作。
CoreWeave	224	合作协议	协议由Coreweave提供高性能 AI 基础设施
美国及各州政府	300	补贴 + 税收减免	能源部贷款担保150 亿、各州税收减免 150 亿
项目债 (绿色债券)	1500	债务	由软银、高盛、摩根大通承销, 期限10年
其他战略投资者	250	股权	Arm、富士康、台积电、苹果、AWS 等小额跟投
总计	5000	--	--

数据来源: 公司官网、华尔街新闻、维基百科、国信证券经济研究所整理

请务必阅读正文之后的免责声明及其项下所有内容

表: OpenAI 数据中心建设计划

时间	地点	负责方与容量	建设说明
2025-2026	德克萨斯州阿比林市	Oracle负责, 目前0.7GW, 后续可达1.4GW	两栋建筑完工, 每栋约60,000个GB200 NVL72 服务器
2025-2026	俄亥俄州洛兹敦、德克萨斯州米拉姆县	软银负责, 未来18个月内可出售1.5GW	Vera Rubin 芯片
2026-	阿联酋的Stargate UAE、挪威的 Stargate 等	阿联酋2026年上线200MW, 后续计划规模5GW	星际之门阿联酋
2026-2027	德克萨斯州沙克尔福德县、新墨西哥州多纳阿纳县、美国中西部	Oracle负责, 超4GW	
总计	未来四年内Stargate计划容量达到近10吉瓦 (GW), 投资超5000亿美元; 英伟达投资1000亿美元锁定未来10GW数据中心的芯片供应权。AMD授权最多1.6亿股 (AMD的10%) 绑定未来6GW芯片供应权。		

数据来源: 公司官网、华尔街新闻、彭博新闻、国信证券经济研究所整理

理

星际之门计划



初期投资400亿美元（包括现金与设施折价），获取3000亿美元数据中心承建权。



软银初期投资1500亿美元并负责项目的财务规划和资金筹集。孙正义持股90%的Arm提供芯片技术支持。软银旗下的可再生能源公司SB Energy提供至少1.5GW数据中心能源建设。



与OpenAI 签订总价值 224 亿美元的协议，提供高性能 AI 基础设施。



2025 年 11 月，OpenAI 与 AWS 签署了一项 7 年、约 380 亿美元的合作协议，租用 AWS 上数十万块最新 Nvidia GPU（如 GB200/GB300 系列）来扩展模型训练与推理能力。



25年9月4日，博通宣布第四个定制AI芯片的主要客户（OpenAI）承诺了价值 100 亿美元的订单。通过博通开发 3nm ASIC，计划 2026 年量产，目标降低训练成本 30%。



25 年 9 月 22 日，NVIDIA 承诺投资 OpenAI 1000 亿美元，锁定未来 10GW 数据中心的芯片供应权，通过 GPU 销售回收资金。既锁定 OpenAI 长期采购，分享 OpenAI 未来增长红利，又形成“资本-硬件-生态”闭环，巩固其在 AI 芯片领域的垄断地位。



截止25年10月，微软总投资约为130亿美元，持有约OpenAI 27%股权。微软获得知识产权共享、收入分成（20%以内）、API 独家经营权。OpenAI 享有 Azure 优先供给以及 30% Azure 折扣。10月双方签订了2500亿美元的多年 Azure 服务合同。



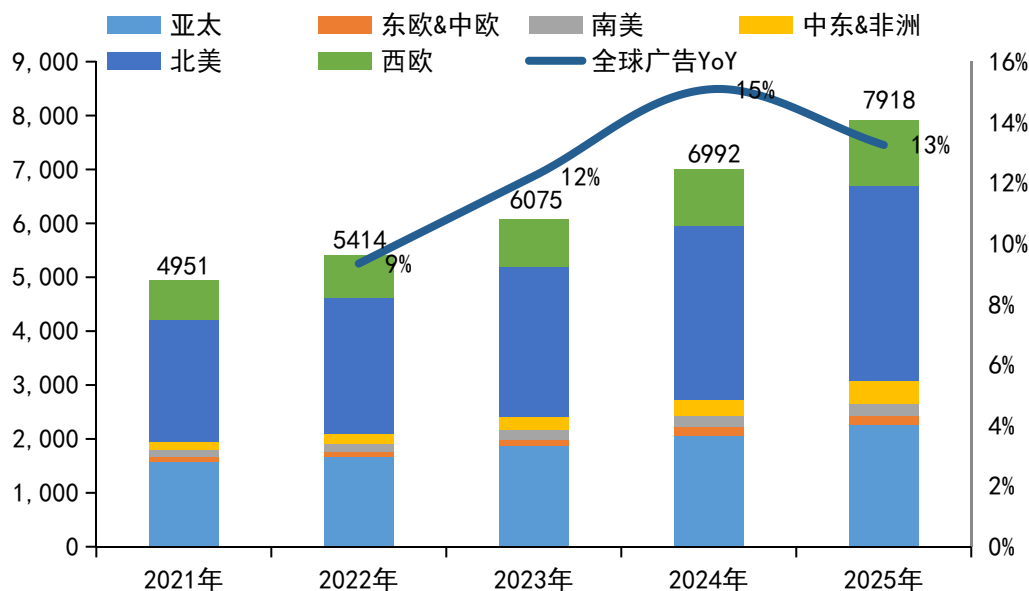
25 年 10 月 6 日，OpenAI 通过承诺部署 6GW AMD GPU 算力换取 1.6 亿 AMD 股权。前提 AMD 股价达标且完成算力部署，形成“算力换股权”的新型合作模式。AMD 通过 OpenAI 旗舰客户背书获取 AI 入场券与 AI 算力市场份额，OpenAI 降低对 NV 的单一依赖且可获得约 10% AMD 股权。

- **一、Anthropic：凭专业生产力打造高毛利护城河**
 - 核心团队成員、经营理念、算力储备
 - 模型能力：Coding、Agent场景下的SOTA模型
 - 商业模式：极简产品矩阵，API贡献主要收入
 - 财务表现：最强模型带来token溢价，28年有望迎来现金流转正
- **二、谷歌：多模态能力领先，生态优势明显**
 - 模型能力：围绕多模态能力打造模型矩阵，综合性能领先
 - 商业模式：原生AI应用+Gemini赋能传统产品，云与广告受益增长
- **三、OpenAI：C端产品领导者，开始发力企业市场**
 - 模型能力：模型路线从分化到统一
 - 商业模式：C端产品领导者，发力企业市场
 - 财务表现：收入预测持续上修，预计2030年超过2800亿
 - 算力储备：星际之门项目持续推进
- **四、静态理解模型的商业化市场空间**

C端流量入口变化影响接近1.5万亿美元市场

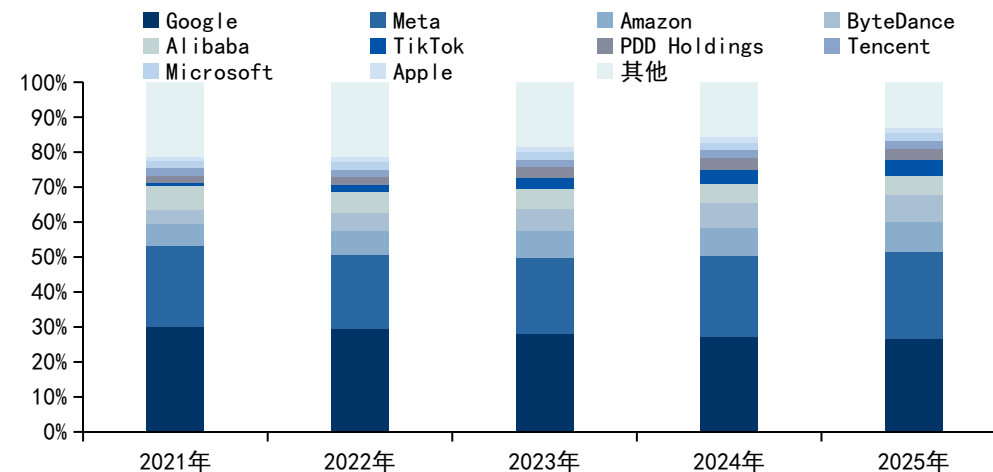
- C端流量入口变化可能带来相关变现收入的转移，流量变现的主要方式包括广告、电商、订阅等，我们测算25年合计接近1.5万亿美元
- **广告：**全球数字广告市场约8000亿美元，中国约1500亿美元，美国3500亿美元。
- Chatbot带来的C端流量变化可能对搜索广告、电商广告最先产生冲击，社交广告影响相对有限，全球市场来看，头部公司中搜索广告（谷歌）+电商广告（亚马逊+阿里+拼多多等）占比接近50%。

图：全球数字广告市场规模（亿美元）



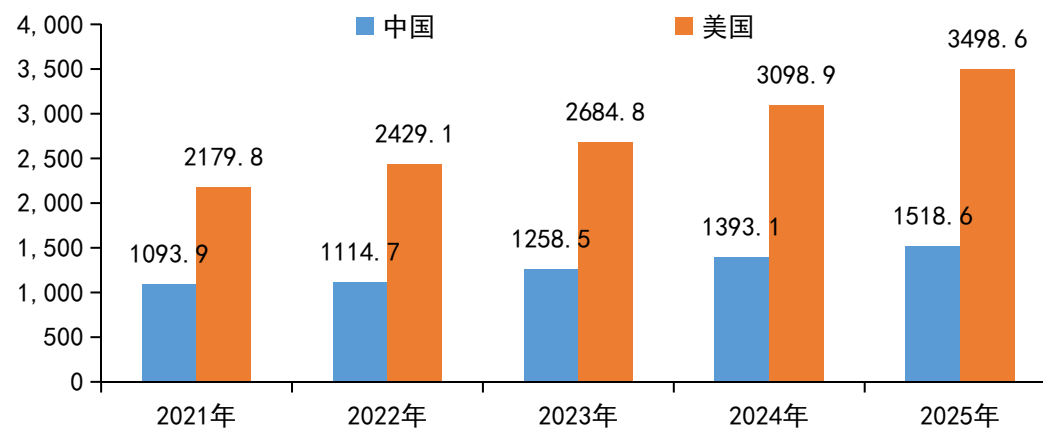
资料来源：emarketer，国信证券经济研究所整理

图：全球数字广告市场份额



资料来源：emarketer，国信证券经济研究所整理

图：中国&美国数字广告市场规模（亿美元）

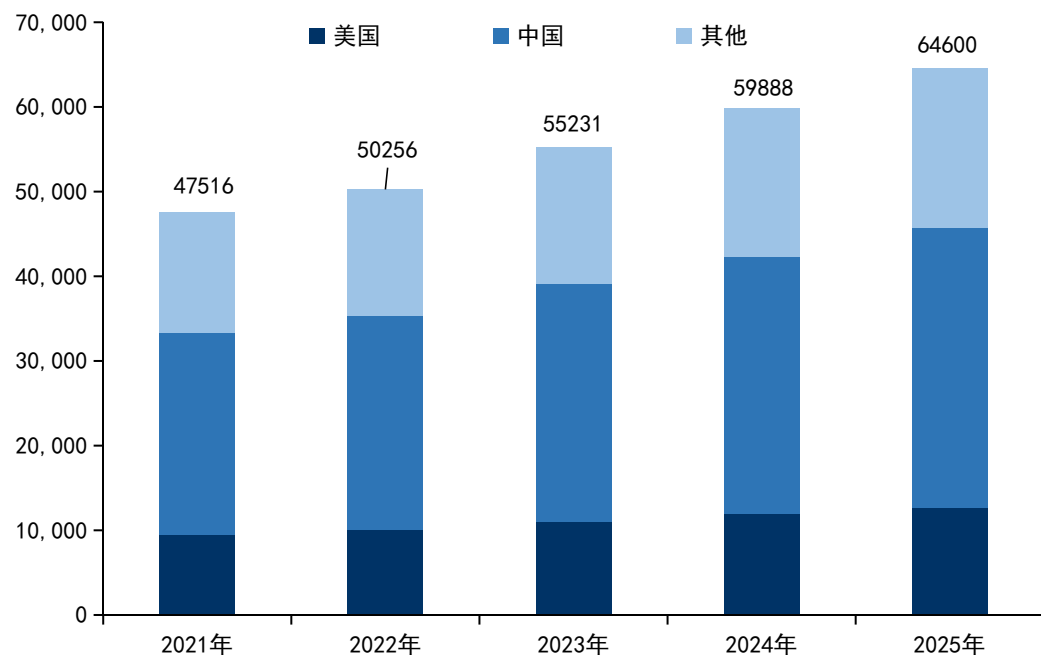


资料来源：emarketer，国信证券经济研究所整理

C端流量入口变化影响接近1.5万亿美元市场

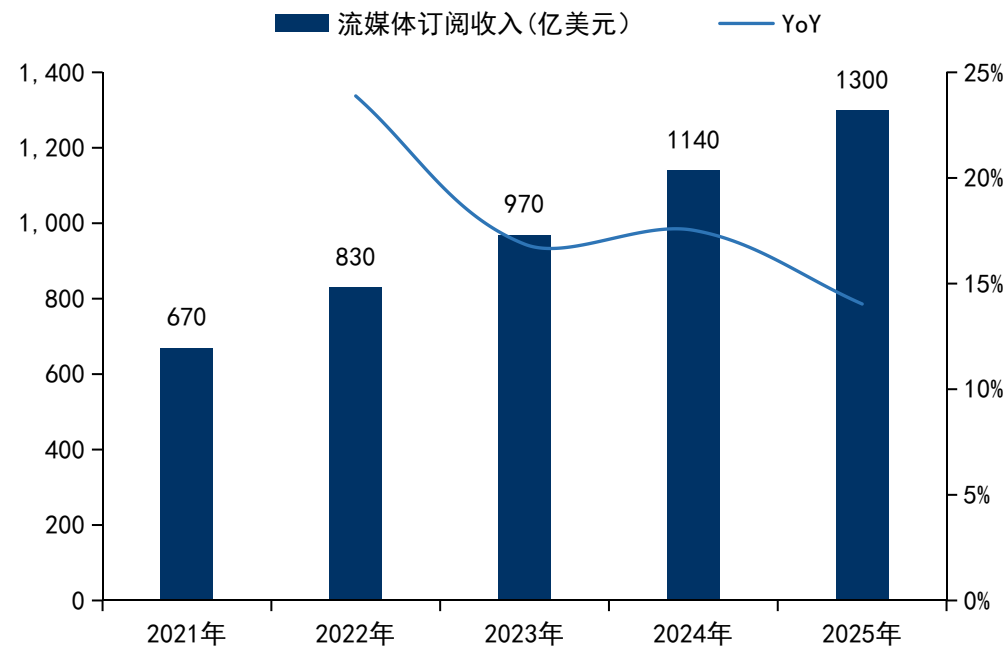
- **电商：**根据Emarketer测算，25年全球电商市场规模约6.5万亿美元；我们假设海外佣金率为10%，国内佣金率为5%，则25年对应电商佣金约5000亿美元。
- **订阅：**订阅服务主要集中于视频、音乐、阅读等内容消费场景，根据PWC数据，全球流媒体订阅收入规模25年约1300亿美元，数字音乐流媒体订阅收入约150亿-200亿美元。

图：全球电商市场规模（亿美元）



资料来源：emarketer，国信证券经济研究所整理

图：全球流媒体订阅市场规模（亿美元）

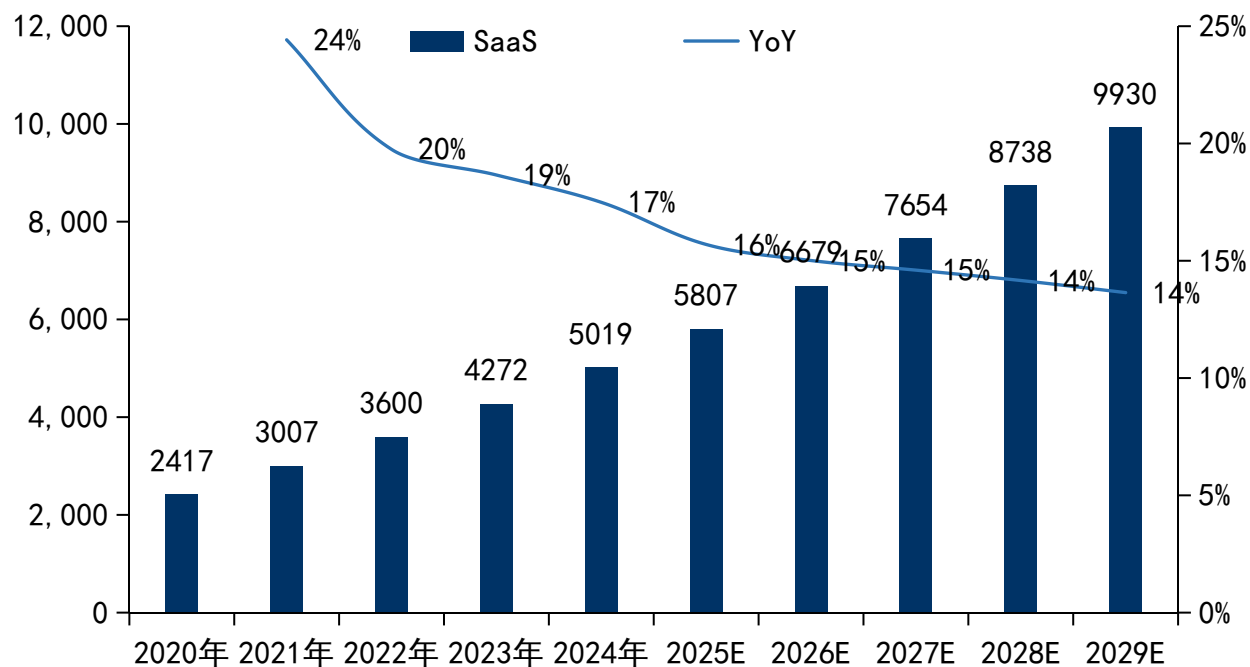


资料来源：PWC，国信证券经济研究所整理

B端5800亿美元SaaS市场将被重构

- B端：SaaS市场将被AI重构，部分公司被大模型替代，部分公司完成AI转型，还有部分原生AI应用贡献增量。
- 据IDC数据，预计2025年全球SaaS市场将达到5800亿美元，同比+16%，预计2029年将达到近1万亿美元规模。预计大模型会颠覆轻量级的工具，以及对准确度要求没有那么高的工具，比如营销工具、翻译工具、邮件回复机器人等。拥有数据壁垒，在垂类细分场景中布局，软件定义工作流程较复杂，或对准确度要求极高的行业，被大模型替代的风险较小，比如医疗、能源、会计、安全等领域。

图：全球SaaS市场规模（亿美元）及增速

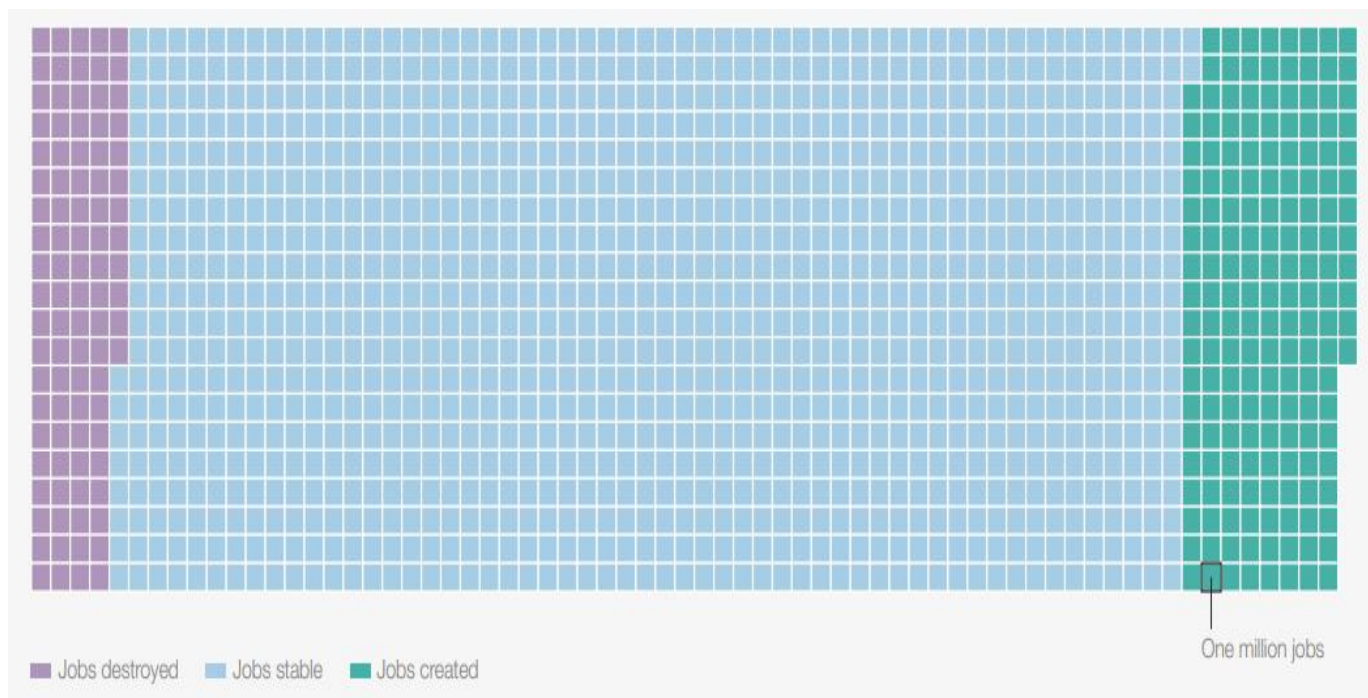


资料来源：emarketer，国信证券经济研究所整理

劳动力替代规模或达到12万亿美元

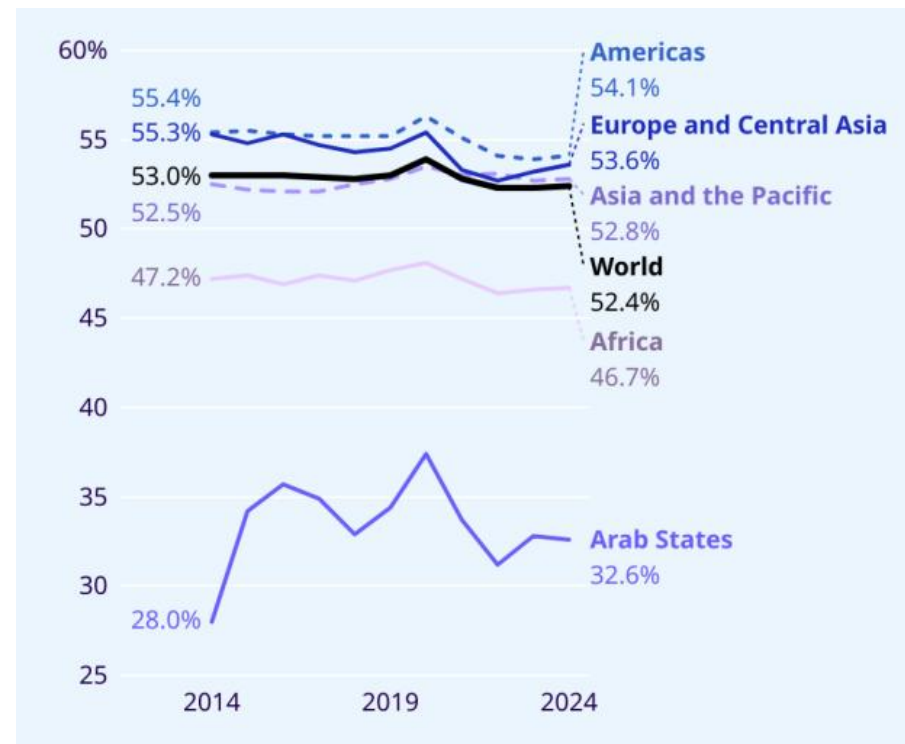
- 根据世界经济论坛（WEF）在《2025年未来就业报告》中指出，到2030年将有 9200 万个岗位被取代，约占现有岗位的 22%，但同时会创造约1.7 亿个新岗位。
- 根据国际劳工组织（ILO）数据，全球劳动收入占GDP比例稳定在52%–53%左右，2025年全球名义GDP超过110万亿美元，则劳动收入约55万亿美元，假设其中22%将被AI替代，则对应的市场规模约12万亿美元。

图：到2030年全球劳动力市场变化



资料来源：WEF，国信证券经济研究所整理

图：全球劳动收入占GDP比例



资料来源：国际劳工组织，国信证券经济研究所整理

第一，宏观经济波动。若宏观经济波动，公司业务、产业变革及新技术的落地节奏或将受到影响。

第二，下游需求不及预期。若下游AI需求不及预期，相关的AI研发投入增长或慢于预期，致使行业增长不及预期。

第三，核心技术水平升级不及预期的风险。AI大模型研发进度落后，AIGC相关产业技术壁垒较高，核心技术难以突破，影响整体进度。

第四，AI快速迭代、平权化下竞争加剧，影响云业务利润率。

国信证券投资评级

投资评级标准	类别	级别	说明
报告中投资建议所涉及的评级（如有）分为股票评级和行业评级（另有说明的除外）。评级标准为报告发布日后6到12个月内的相对市场表现，也即报告发布日后的6到12个月内公司股价（或行业指数）相对同期相关证券市场代表性指数的涨跌幅作为基准。A股市场以沪深300指数（000300.SH）作为基准；新三板市场以三板成指（899001.GSI）为基准；香港市场以恒生指数（HSI.HI）作为基准；美国市场以标普500指数（SPX.GI）或纳斯达克指数（IXIC.GI）为基准。	股票投资评级	优于大市	股价表现优于市场代表性指数10%以上
		中性	股价表现介于市场代表性指数±10%之间
		弱于大市	股价表现弱于市场代表性指数10%以上
		无评级	股价与市场代表性指数相比无明确观点
	行业投资评级	优于大市	行业指数表现优于市场代表性指数10%以上
		中性	行业指数表现介于市场代表性指数±10%之间
		弱于大市	行业指数表现弱于市场代表性指数10%以上

分析师承诺

作者保证报告所采用的数据均来自合规渠道；分析逻辑基于作者的职业理解，通过合理判断并得出结论，力求独立、客观、公正，结论不受任何第三方的授意或影响；作者在过去、现在或未来未就其研究报告所提供的具体建议或所表述的意见直接或间接收取任何报酬，特此声明。

重要声明

本报告由国信证券股份有限公司（已具备中国证监会许可的证券投资咨询业务资格）制作；报告版权归国信证券股份有限公司（以下简称“我公司”）所有。本报告仅供我公司客户使用，本公司不会因接收人收到本报告而视其为客户。未经书面许可，任何机构和个人不得以任何形式使用、复制或传播。任何有关本报告的摘要或节选都不代表本报告正式完整的观点，一切须以我公司向客户发布的本报告完整版本为准。

本报告基于已公开的资料或信息撰写，但我公司不保证该资料及信息的完整性、准确性。本报告所载的信息、资料、建议及推测仅反映我公司于本报告公开发布当日的判断，在不同时期，我公司可能撰写并发布与本报告所载资料、建议及推测不一致的报告。我公司不保证本报告所含信息及资料处于最新状态；我公司可能随时补充、更新和修订有关信息及资料，投资者应当自行关注相关更新和修订内容。我公司或关联机构可能会持有本报告中所提到的公司所发行的证券并进行交易，还可能为这些公司提供或争取提供投资银行、财务顾问或金融产品等相关服务。本公司的资产管理部门、自营部门以及其他投资业务部门可能独立做出与本报告意见或建议不一致的投资决策。

本报告仅供参考之用，不构成出售或购买证券或其他投资标的的要约或邀请。在任何情况下，本报告中的信息和意见均不构成对任何个人的投资建议。任何形式的分享证券投资收益或者分担证券投资损失的书面或口头承诺均为无效。投资者应结合自己的投资目标和财务状况自行判断是否采用本报告所载内容和信息并自行承担风险，我公司及雇员对投资者使用本报告及其内容而造成的一切后果不承担任何法律责任。

证券投资咨询业务的说明

本公司具备中国证监会核准的证券投资咨询业务资格。证券投资咨询，是指从事证券投资咨询业务的机构及其投资咨询人员以下列形式为证券投资人或者客户提供证券投资分析、预测或者建议等直接或者间接有偿咨询服务的活动：接受投资人或者客户委托，提供证券投资咨询服务；举办有关证券投资咨询的讲座、报告会、分析会等；在报刊上发表证券投资咨询的文章、评论、报告，以及通过电台、电视台等公众传播媒体提供证券投资咨询服务；通过电话、传真、电脑网络等电信设备系统，提供证券投资咨询服务；中国证监会认定的其他形式。

发布证券研究报告是证券投资咨询业务的一种基本形式，指证券公司、证券投资咨询机构对证券及证券相关产品的价值、市场走势或者相关影响因素进行分析，形成证券估值、投资评级等投资分析意见，制作证券研究报告，并向客户发布的行为。



国信证券

GUOSEN SECURITIES

国信证券经济研究所

深圳

深圳市福田区福华一路125号国信金融大厦36层

邮编：518046 总机：0755-82130833

上海

上海浦东民生路1199弄证大五道口广场1号楼12楼

邮编：200135

北京

北京西城区金融大街兴盛街6号国信证券9层

邮编：100032