

『翻墙』: DeepSeek-V4的解题思路

□ 记者 郑小芳

工程技术跟科学研究的不同之处或许在于,当正门不通时,我们可以“翻墙”进去。——题记

近日,在经历了爆火后的沉寂,深度求索全新系列模型 DeepSeek-V4 预览版正式上线并同步开源,在全球 AI 圈引起广泛关注。

此次 DeepSeek-V4 之所以成为行业焦点,在于它实现了技术与产业层面的双重突破。技术层面,V4 标配了过去只在海外闭源模型高端版本中才具备的 100 万 Token(词元)超长上下文能力,相当于《三体》三部曲总字数,可一次性承载海量文档并完成全局推理,彻底告别长文本处理的切片割裂痛点。产业层面,V4 实现了与华为昇腾、寒武纪等国产算力芯片的深度耦合,为国产 AI 生态的自主构建按下了加速键。

在 DeepSeek-V4 发布前,国产芯片一直面临一个尴尬的境地:用户少,生态建不起来。这一困境的根源在于,传统 Transformer 架构在处理长上下文时,显存占用随文本长度呈平方级增长。国产芯片受制于制程工艺,虽在算力峰值上奋力追赶,但显存带宽与容量方面仍与英伟达 H100/Rubin 存在代际差距。在传统架构下用国产芯片处理长文本几乎是“不可能的任务”。

但 DeepSeek-V4 给出了一条全新的解题思路:既然硬件条件暂时追不上,那就从架构层面重构“生产关系”。

V4 引入了 CSA(压缩稀疏注意力)与 HCA(高度压缩注意力)相结合的混合架构。通俗地讲,它不再逐字逐句比对百万字文档,而是学会了“抓重点”和“做摘要”。通过对 token 进行智能压缩,V4 大幅降低了对计算和显存的需求,在提升性能的同时实现了资源高效利用,破解了大模型“高性能必高成本”的行业难题。在百万 Token 上下文设定下,V4-Pro 的单 token 推理算力仅为上一代 V3.2 的 27%,KV 缓存只需 10%;更经济的 V4-Flash 版本这两项指标分别仅为 10%和 7%。

如此一来,原本需要 80GB 显存才能胜任的长上下文推理任务,现在可能只需 8GB 显存就能运行。英伟达的高带宽显存稀缺性,在这套新架构面前被大幅削弱,国产芯片的显存短板便不再成为瓶颈。

此次 V4 没有按行业惯例给予英伟达早期测试权限,而是将提前适配的机会独家开放给华为昇腾和寒武纪,目标直指从 CUDA 生态到华为 CANN 框架的整体迁移。可以预见,当一款开源模型的架构创新能够系统性降低显存门槛,跑通国产芯片并实现推理成本大幅下降,英伟达凭借显存优势筑起的护城河便不再是铁板一块。

这正是 DeepSeek“曲线救国”的深层逻辑:不执着于单卡性能上的硬碰硬,而是凭借系统级优化、软硬深度协同和架构创新,重新定义了竞争维度。

□ 牧龙

最近,A 股市场上演了一场惊心动魄的“股王”争夺战。短短两周多的时间里,两家科技公司——源杰科技与寒武纪,先后超越贵州茅台,登顶“股王”宝座。截至 5 月 6 日收盘,寒武纪以 1825 元/股排在第一位,源杰科技以 1531.01 元/股排在第二位,贵州茅台则以 1375 元/股位列第三。而去年同期(2025 年 5 月 7 日收盘),源杰科技的股价不过 132.17 元/股,短短一年涨幅超 1000%。

受益于全球 AI 产业爆发,算力需求暴涨,光模块、GPU、存储芯片、光芯片、云计算等领域备受关注,作为光芯片赛道龙头企业的源杰科技顺势而起,成为了 AI 风口上的“股王”,其创始人张欣刚也因此被推到了聚光灯下。

► 选择最难的路

2013 年的一天,陕西咸阳某酒店的停车场内,一位日本供应商望着眼前一辆老旧的奥拓,眉头紧锁,迟迟不敢上车。

他一度认为车里的人是骗子。毕竟,对方开口声说要造高端光芯片,却开着这样一辆破旧的小车带他参观所谓的“先进工厂”。这巨大的反差,让他实在无法把眼前的场景与芯片研发联系起来。车里的人,就是张欣刚。

那位日本供应商绝不会想到,眼前这个人会在之后十几年的时间里,将企业打造成中国光芯片领域的龙头,在全球光芯片市场挤占一席之地。

1970 年出生的张欣刚,本科毕业于清华大学材料系,后赴美深造,先后取得南加州大学材料科学硕士、博士学位。2001 年,博士毕业的张欣刚进入光通信公司 Luminent 工作,他从基层研发员起步,逐步晋升至研发经理,深度参与光芯片从研发设计到制造封测的全流程工作。2007 年,

人物

13 年的『光速突围』

新晋『股王』源杰科技创始人张欣刚



欣刚就是其中之一。

张欣刚于 2010 年回国创业,在北京成立了一家光电技术公司,但因合伙人挪用资金,导致公司停摆,首次创业失败。但张欣刚没有止步。2013 年,他决定再次创业。当时,陕西咸阳正好推出一些科创扶持政策,张欣刚便将落脚点定在咸阳,成立陕西源杰半导体技术有限公司。

那个时候,国内高端光芯片几乎被海外厂商全盘垄断。为了打破这一局面,张欣刚选择了一条最难的路——IDM 模式,即从芯片设计、外延生长、晶圆制造到封装测试,每一步都自己掌握。这种模式前期投入大、回报周期长,但张欣刚坚信,只有把全流程攥在自己手里,才能真正突破卡脖子限制,做出性能稳定、可靠的高端芯片。

► 布局 AI 算力赛道

创业初期,一切都很难。张欣刚面临的是国内产业基础薄弱、国际技术封锁、公司资金匮乏等一系列问题。他一方面筹措资金,另一方面

Luminent 与另一家光通信设备商 Fiberson 整合,成立索尔思光电,张欣刚出任研发总监。

索尔思光电在全球光模块行业中占据重要地位,不少业内大佬都是从索尔思走出来的,张

豆包收费引发热议

AI 正在加速变现

5 月 4 日,字节跳动旗下的 AI 助手豆包在苹果应用商店悄然挂出一则服务声明——标准版 68 元/月,加强版 200 元/月,专业版 500 元/月,这个坐拥 3.45 亿月活的超级 AI 应用,正计划推开付费订阅的大门。显然,国内大模型即将告别那个“靠免费换规模”的粗放时代,硬生生闯入需要“真金白银”说话的下半场。



□ 沈毅斌

当每个用户都是一笔“算力账单”

豆包收费背后,首先是一笔经济账。

截至今年 3 月,豆包大模型日均 Token 消耗量从 2024 年 5 月发布时的 1200 亿,飙升至 120 万亿,增长约 1000 倍。火山引擎披露,累计 Token 调用量超一万亿的企业客户从 100 家增长到 140 家。

Token 时代,新的经济规律已初现端倪:边际成本效应消失,用户规模越大,算力成本越高。

一份流传于中文 IT 技术社区(CSDN)的成本拆解模型显示,豆包单次推理成本中,硬件折旧占 58%、电力占 29%。简言之,每一次用户提问、每一轮对话生成,都直接对应扎实的算力消耗,成本并没有随着规模扩大而下降。

与此同时,收入端还没有打通。过去,中国互联网的变现逻辑是:用户越多,广告价值和商业变现空间越大。AI 时代,这条路刚刚开始试水。今年 5 月,OpenAI 推出自助式广告平台,目标收入是 2026 年 25 亿美元、2030 年 1000 亿美元。但质疑声甚嚣尘上,毕竟人们因 AI 而放弃搜索引擎,原本是为了得到一个有效结果,而不是带着

“金钱味道”的推荐广告。

豆包也在探索广告模式,但至少在目前,它还很难彻底跨过 GEO(生成式引擎优化)和“大模型投毒”的界限。QuestMobile 数据显示,截至 3 月,豆包月活用户接近 3.45 亿,日活用户约 1.4 亿,月人均使用次数达 54.8 次。

一面是居高不下的成本,一面是还未成熟的商业模式,即便是字节跳动这样公认的“家有余粮公司”,也越来越烧不起 Token 了。有媒体报道,2025 年字节跳动总营收预计接近 2000 亿美元,但净利润同比下滑超 70%,核心原因正是对 AI 业务的投入。

尽管抖音集团副总裁李亮在微博发文澄清,实际经营利润率降幅“远没有那么大”,但他也承认,抖音电商增速放缓和新兴业务相关投入增大,导致下半年经营利润率有小幅下滑。有人戏称,“抖音赚的钱,正在被豆包烧光”。

国内 AI 大模型是否跟进?

豆包收费消息一出,所有人的目光都投向了腾讯元宝、阿里千问等其他面向 C 端的大模型。它们跟还是跟?截至目前,腾讯元宝和阿里千问并未传出要收费的消息,甚至有业内人士猜测,元宝和千问很可能会借此机会,抢夺豆包用户。

3 月以来,腾讯云连续两次涨价:3 月 11 日大模型

API 输入价格单轮暴涨 463%,4 月 9 日又宣布 5 月 9 日起 AI 算力相关产品上调 5%;阿里云也宣布从 4 月 18 日起,平头哥真武 810E 等算力卡上调 5%~34%,CPFS(智算版)上调 30%,5 月 15 日起,阿里云百炼平台部分模型单元服务涨幅从 2%~7%不等。

目前来看,两家友商的节奏是 B 端先动、C 端观望,以生态吸引更多用户。

只是,元宝和千问能承接多少豆包流失的用户?能坚持免费策略多久?同样要看能否找到用户数量和 Token 成本的平衡点。

豆包能“物有所值”吗?

毕竟,全球大模型厂商都玩不起打折竞赛了。

进入 2026 年后,全球 AI 市场涨声一片。随着 AI 逐渐成为有效生产力工具,B 端用户的付费意愿逐渐提升。

Claude Code 的官网显示,在企业级部署中,平均成本约为每位开发人员每天 13 美元、每月 150~250 美元,90%的用户每天成本低于 30 美元,按上限测算,年成本约为 1.1 万美元,而数据显示,一个硅谷软件工程师的平均年薪约为 14 万美元。

智谱 AI 年内三度提价,API 调用定价累计上调 83%,而 Token 消耗量反而增长 400%。智谱 CEO 张鹏表示,智能上限决定定价权,Token 消耗规模决定价值体量。大模型正在从“谁便宜买谁”转向“谁好用谁值钱”。

对于已经积累足够用户黏性和产品口碑的头部玩家来说,率先尝试付费分层,其实是主动探索可持续的商业闭环。一旦付费模式跑通,就能给整个赛道吃下一颗“定心丸”,也让市场对 AI 应用的商业化有着更清晰的预期。

豆包的付费探索,正是这一转向在 C 端的具体化体现。与全球付费体系相比,豆包的定价不算 outlier。

目前 ChatGPT 已建立完整的四层定价体系:Go 版约 8 美元/月,Plus 版 20 美元/月,Pro 版 100 美元/月及 200 美元/月。其中 Pro 档位直接对标 Anthropic Claude Max,面向重度编程用户。豆包标准版的 68 元人民币(约合 9.9 美元)与 ChatGPT Go/Plus(8~20 美元)处于相近区间,专业版 500 元(约 73 美元)则大致位于 Pro 版和中国市场消费水平交叉点上。

凯基证券的分析报告认为,豆包付费标志着中国大模型行业从“野蛮增长”转向“价值兑现”,竞争重心将从“模型参数”转向“商业 ROI”。

但豆包需要解决的问题是,能否提供与其收费水平匹配的能力,让具有付费能力和意愿的用户觉得“物有所值”。

豆包能走通这条“基础功能免费+高算力功能付费”道路吗?它所保留的免费功能,足以让它留住大部分活跃用户吗?分层订阅能满足高阶用户需求吗?

在这个新开启的 AI 时代,人们只能等待答案自己浮现。据《IT 时报》

资讯

6G 技术试验频率获工信部批复

本报讯 为进一步推动我国 6G 技术研发、标准研制与产业化进程,工业和信息化部近日向 IMT-2030(6G)推进组批复 6GHz 频段 6G 试验频率使用许可,支持其在部分地区开展 6G 技术试验,面向国际电信联盟确定的 6G 典型场景与关键性能指标,开展技术研发攻关和测试验证。

上海发布全球科技伙伴计划

本报讯 上海市政府近日印发《上海全球科技伙伴计划》,提出国际联合研究、人才交流访问、国际科技组织引育、开放科学、产业创新国际化、创新生态优化六大行动,支持海外科研机构、企业、组织和科技创新创业人才等成为上海科技伙伴,建设具有全球竞争力的开放创新生态,提升上海科技创新的全球影响力和引领力,助力共建全球科技共同体。

英伟达高端 AI 芯片在华市占归零

本报讯 近日,英伟达首席执行官黄仁勋表示,该公司在中国 AI 加速器市场的份额目前已降至零。据统计,2025 年中国市场的 AI 加速卡出货量约为 410.6 万块,其中英伟达以 55% 的市场份额排名第一,这与该公司 2021 年在华 95% 的市场份额相比,已大幅下滑。