

GenAI

中国GenAI市场洞察

企业级大模型调用全景研究

State of GenAI: Foundational Model
in Chinese Enterprise, 2025

■ 目录

引言	-----	03
章节1 中国企业级大模型发展综述	-----	05
章节2 中国企业级大模型调用现状及趋势	-----	09
章节3 中国企业级大模型调用行为分析	-----	15
章节4 中国企业级大模型调研企业画像与方法论	-----	21

■ 引言

沙利文联合头豹研究院发布《中国GenAI市场洞察：企业级大模型调用全景研究》，在本次研究中，我们对700位企业IT部门负责人、技术总监/经理及AI项目负责人进行了问卷调研，覆盖了金融、制造、互联网、消费电子、汽车等多个重点行业。调研对象所在企业类型多样，覆盖不同营收层级和AI投入规模的企业。

本研究旨在评估企业级大模型市场开源与闭源模型的部署情况，并洞察企业选择偏好背后的动因，形成对中国企业级大模型应用现状与趋势的结构化理解，为中国企业级大模型市场提供具有指导性的洞见。

■ 性能与选择丰富性提升驱动开源生态繁荣

在企业级大模型应用中，开源与闭源模型的性能差异不断收敛。Qwen、DeepSeek等开源体系已在对话、代码生成、数据分析、逻辑推理等场景实现“闭源替代”，兼具高可用性和性价比。

丰富的多类型、多尺寸模型矩阵，让企业可根据业务复杂度与响应速度需求，精准选型，实现算力与成本双优化。

开源模型的核心优势在于后训练灵活性，企业能针对自身数据分布和业务逻辑深度优化通用基础模型，这正是开源生态与闭源服务的本质区别与关键壁垒。

■ 企业从追求“大而全”模型转向业务适配的“灵活集成+技术自主可控”方案演进

企业级大模型应用的决策重心正从“追求单一最强模型”转向“为特定业务场景寻求最优解”，核心在于算力性价比、系统灵活性与安全可控性的最佳平衡，推动用户需求由对单一模型的依赖，演化为“多模态、多尺寸、多场景”模型矩阵的应用模式。

相比闭源服务，开源模型以其透明、高度自主可控和部署灵活的优势，赋能企业构建自有技术体系并沉淀核心技术资产；随着社区不断迭代，其性能在多种场景中展现出强大竞争力，满足日益复杂的业务需求。

对成本与可控性的双重诉求，促使企业采取更务实的部署策略——无论是大规模、多尺寸地应用开源模型，还是开源与闭源混合使用，都通过灵活的资源配置，构建弹性且经济高效的AI基础设施，以应对快速变化的商业挑战。

Key Findings

关键发现

2025H1 中国企业级大模型日均调用量已达10万亿tokens

单位：亿tokens



2025年上半年

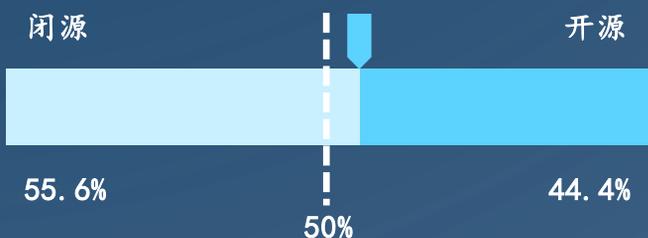
大模型市场日均调用量达到

10万亿 tokens

较2024年下半年的2万亿

tokens实现约**363%**的增长

呈现出爆发式放量态势



开源模型快速获得企业的认可，

调用占比已提升至**44.4%**，较闭源的**55.6%**已非常接近

中国企业级大模型调用市场，

阿里通义系列大模型调用量

占比达**17.7%**

成为目前市场选择最多的模型



INTRODUCTION

01

中国企业级大模型 发展综述

1. 企业级大模型发展综述

1.1 在人工智能快速迭代的背景下，大模型正形成“开源与闭源并行发展”的双轨格局。

在人工智能快速迭代的背景下，大模型正形成“开源与闭源并行发展”的双轨格局。开源大模型凭借算力性价比最优、系统集成灵活以及透明可验、社区驱动等特征，逐渐成为企业实现低成本落地与自主可控的优选路径；其开放的结构与全球开发者协作机制，不仅提升了模型的适配性与长期价值，也推动了技术的民主化和应用的多样化。闭源大模型则以商业闭环、黑盒可控、集中优化与高成本壁垒为特征，通过集中资源构建高性能模型并保障服务稳定性，适合对可靠性、算力效率与服务闭环要求较高的企业客户。

图1：大模型市场演进趋势，2025年

闭源大模型	讯飞星火X1	文心大模型X1 Turbo	doubao-seed-1.6系列			
	Kimi k1.5	SenseNova 6.0	MiniMax-Hailuo-02			
	Baichuan-M1-preview	豆包Seedream3.0	Hunyuan3D-PolyGen			
	文心大模型4.5	Doubao-1.5-thinking-pro	星火医疗 V2.5			
	文心大模型X1	Doubao-1.5-thinking-pro-vision	讯飞星火X1升级版			
	混元Turbo S	豆包·语音播客模型	SenseNova V6.5			
	豆包Seed-Thinking-1.5	SeedEdit 3.0	混元Large-Vision			
文心大模型4.5 Turbo	Seedance 1.0 Pro					
开源大模型	MiniMax-01系列	R1-Omni	Qwen3 全系列	MiniMax-M1	ThinkSound	GLM-4.5系列
	DeepSeek R1	HunyuanVideo 12V	GLM-4-32B/9B	Kimi-Dev	HumanOmniV2	ARC-Hunyuan-Video
	Qwen2.5-VL	DeepSeek-V3-0324	Kimi-Audio	Kimi-VL-Thinking	Kimi K2	混元0.5B/1.8B/4B/7B
	Wan2.1	Qwen2.5-VL-32B-Instruct	DeepSeek-Prover-V2	Hunyuan-A13B	Qwen3-235B-A22B-Instruct-2507	Qwen-Image
	QwQ-32B	Qwen2.5-Omni-7B	Wan2.1-VACE	文心大模型4.5系列	Qwen3-Coder	Baichuan-M2
	Baichuan-M1-14B	DeepSeek-R1-0528	Qwen3-Embedding	混元3D世界模型1.0	Intern-S1	Qwen-Image-Edit
	Kimi-VL系列	InternVL 3.0	混元3D 2.1	GLM-4.1V-Thinking	Wan2.2	DeepSeek-V3.1
						Intern-S1-mini

当前，大模型演进趋势正呈现出“开源为主流、参数趋轻量、应用更普及”的方向。开源大模型的快速涌现正在重塑产业格局，仅2023年全球就有149个基础模型问世，其中开源模型占比达65.7%，显示其正逐步取代闭源方案成为生态主导力量。

从技术演进看，参数规模的轻量化趋势显著降低了算力门槛，使模型在训练与推理环节更具能效优势和部署灵活性，推动大模型向更多垂直行业与中小企业场景加速渗透。与此同时，开源模式依托透明可验的模型结构、全球开发者社区的协作机制以及灵活的微调与定制能力，不仅降低了企业研发与应用的边际成本，还显著增强了技术的自主可控性和长期迭代潜力。相比之下，闭源模型虽在高性能训练、集中优化与商业闭环服务上具备优势，但其高成本和封闭特性限制了大规模普及。

未来，预计超过80%的企业将在智能化建设中采用开源大模型，开源生态的战略地位将进一步凸显，并成为驱动产业普及化、促进技术民主化、支撑数字化转型的重要核心力量。



Model Adoption

02

中国企业级大模型
调用现状及趋势

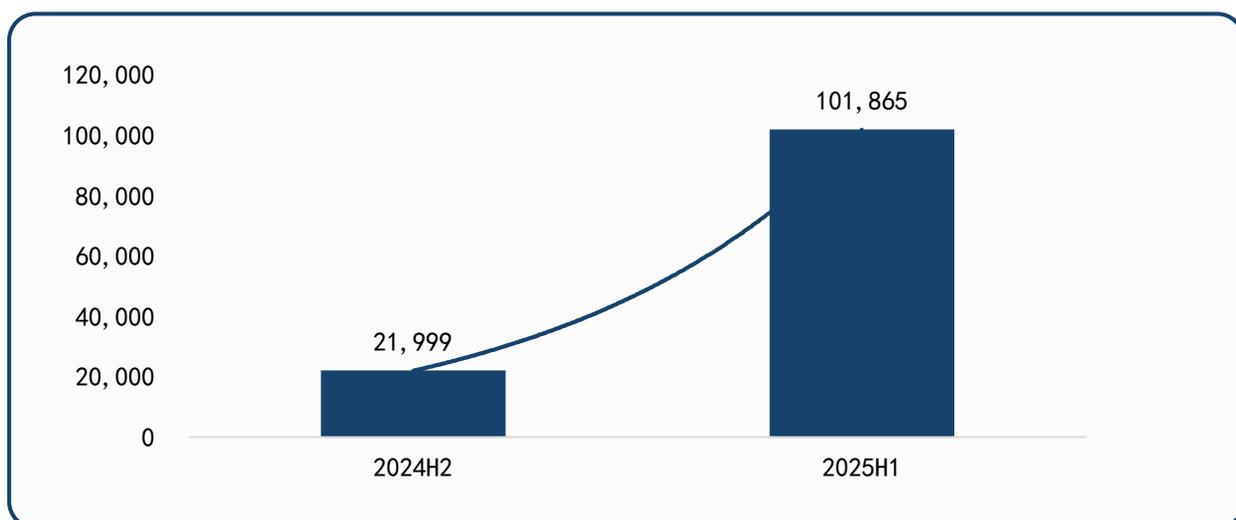
2. 企业级大模型调用现状及趋势

2.1 2025年上半年，大模型市场日均调用量达到101,865亿 tokens，较2024年下半年的21,999亿tokens 实现约363%的增长，呈现出爆发式放量态势。

这一跃升不仅意味着市场需求的全面释放和持续攀升，也伴随着算力与存储等基础资源消耗的显著增加，更清晰地表明大模型正快速走出试点验证期，进入规模化落地的新阶段。随着调用量的高速增长，大模型逐渐成为企业数字化与智能化升级的核心支撑力量，同时也推动产业链上下游在算力供给、数据治理、应用开发与服务交付等环节加速重构，产业格局正呈现出由点状突破向系统化集成演进的趋势。

图1：2025H1中国企业大模型市场日均调用量

单位：亿tokens



调用量激增的背后，驱动因素主要有以下三个方面：

1. 供给端的模型与算力扩容是基础驱动。过去一年，以DeepSeek、通义为代表的开源体系在模型性能和成本效率上实现突破，降低了企业自建与定制化应用的边际成本；与此同时，阿里云、火山引擎、百度智能云等闭源平台加速释放算力供给，通过预置API接口和SaaS服务拓展了中小企业用户基数，形成了调用量爆发的“基础设施推力”。
2. 需求端的场景化渗透显著提速。金融、政务、制造、教育等行业在政策激励与竞争压力下，加快了RAG（检索增强生成）、行业Agent和数字员工的落地部署，推动模型调用由“研发测试”走向“生产级运行”。特别是在金融与制造业，风险控制、质检优化、供应链调度等关键场景已进入常态化调用，直接拉动日均调用总量。
3. 生态外溢效应正在形成规模扩散。开源社区的快速繁荣带动工具链、微调框架与中间件的普及，使企业在多模型、多任务的并行部署上更加高效；闭源厂商则通过算力、数据和生态的绑定效应，推动垂直行业解决方案大规模复制，形成正向循环。两者共同作用下，市场日均调用量呈现出指数型曲线增长。

总体看来，中国大模型产业已进入“规模化应用驱动”的新阶段，市场竞争焦点正在由单纯的模型性能比拼，转向算力保障、数据合规、生态整合与行业深度适配。



Behavior Analysis 03

中国企业级大模型
调用行为分析

3. 企业级大模型调用行为分析

3.1 在企业级应用场景中，“业务价值”是开闭源模型选型的核心出发点，闭源模型胜在“省心可靠”，开源模型则更契合企业对于灵活性和自主可控性的考量。

图8：开源与闭源的选择驱动因素



在进行模型选型时，无论是开源还是闭源，“业务价值”始终是其核心的出发点。企业通常将“性能表现最佳”作为首要考量因素，但这并非单纯追求最高的绝对指标，而是强调模型在特定业务场景下的适配度与满足度。例如，在边缘计算等应用中，开源模型对小尺寸部署的支持及其在性能与延迟间的平衡，可能更符合企业的实际需求。这一原则决定了企业会根据自身情况，在不同技术路径中做出权衡。

在此背景下，闭源模型的核心优势在于其“省心可靠”。凭借稳定的性能、成熟的生态体系、完善的技术支持以及快速的迭代修复能力，闭源模型目前占据市场主流。厂商丰富的落地案例和信誉背书为企业提供了低风险、高效率的解决方案，使其能够快速获得可靠的产品与服务，这是一种以最小化试错成本来保障业务价值的稳健路径。

与此相对，开源模型的最大魅力则体现在其强大的“自主可控性”。其最独特的优势在于支持深度定制化开发，允许企业将模型与自有数据和业务逻辑深度融合，进行精细优化。更重要的是，选择开源意味着企业能够完全掌握模型的知识产权和数据安全，无需将核心数据传输至第三方，并将训练成果固化为企业专有的数字资产。这从根本上规避了供应商锁定的长期风险，是一种通过掌握核心技术来构建长期、自主业务价值的战略选择。

■ 联系我们

沙利文联合头豹研究院发布《中国GenAI市场洞察：企业级大模型调用全景研究》，在本次研究中，我们对700位企业IT部门负责人、技术总监/经理及AI项目负责人进行了问卷调研，覆盖了金融、制造、互联网、消费电子、汽车等多个重点行业。调研对象所在企业类型多样，覆盖不同营收层级和AI投入规模的企业。

本研究旨在评估企业级大模型市场开源与闭源模型的部署情况，并洞察企业选择偏好背后的动因，形成对中国企业级大模型应用现状与趋势的结构化理解，为中国企业级大模型市场提供具有指导性的洞见。

报告完整版登录www.leadleo.com

搜索《中国GenAI市场洞察：企业级大模型调用全景研究》

袁栩聪-首席分析师

弗若斯特沙利文大中华区



联系邮箱：

oliver.yuan@frostchina.com

李庆-研究总监

弗若斯特沙利文大中华区



联系邮箱：

livia.li@frostchina.com

FROST & SULLIVAN

沙利文

