

多重因素引致美股剧烈调整，Meta 发布最新开源大模型 Llama-3

计算机行业

推荐

维持评级

摘要 | 04.12-04.19

- 股指动态** 美股三大指数、中概股及港股科技、A股计算机指数全线下跌。标普500指数-3.05%，纳斯达克综合指数-5.52%，费城半导体指数-9.23%；TAMAMA科技指数-7.02%；纳斯达克中国金龙指数-2.09%；恒生科技指数-5.65%；计算机-3.73%。
- 个股表现** 热门科技股巨幅调整，集体伴随市值大幅缩水；英伟达市值跌破2万亿美元。据统计，相比4月12日收盘价，4月19日盘后，苹果合计-6.54%，英伟达-13.59%，特斯拉-14.03%，谷歌-2.31%，亚马逊-6.18%，META-6.02%，微软-5.40%，ARM-30.98%，英特尔-4.17%，高通-7.97%，AMD-10.19%。美股七大科技公司（苹果、英伟达、特斯拉、谷歌、亚马逊、Meta、微软）单周总市值合计蒸发9661.66亿美元（约合6.9万亿人民币），刷新历史记录。
- 10年期国债及汇率** 美联储发言偏鹰派，10年期美债持续走弱；美债利率再走高或使科技股阶段性再承压。周内，美国10年期国债利率持续走高，累计上浮12bps至4.62%；中国10年期国债利率下降至2.25%，累计下降2.97bps。4月19日，美元兑人民币中间价报7.10；较4月12日价累计调升79个基点。
- 核心观点**

近期海外部分地区冲突导致地缘政治风险提升、避险情绪加剧，致使风险资产遭抛售。美联储发言偏鹰派，美债利率再走高或使科技股阶段性再承压。受手机及PC下游需求预期疲软影响，台积电下调预期，致使全球半导体资金大幅回撤；超微电脑未如期披露业绩预告，引发市场担忧AI科技股业绩不及预期风险，从而引发AI概念股剧烈回调。我们认为，上述多重利空因素致使热门科技股巨幅调整。短期来看，四月下旬美股科技将陆续披露季报，在此前的高市场预期背景下，各公司业绩相对市场预期具较大不确定性。

Meta发布最新开源大语言模型Llama-3，性能较上一代大幅提升、赶超竞品。相较Llama-2，在参数规模、训练数据集、模型架构及安全性方面均有不同程度提升，被认为是“首个开源GPT-4级别的模型”。Llama-3可大幅扩充Meta AI应用功能，实现与Facebook、Instagram、WhatsApp等App的无缝集成；此外，网页版Meta AI与OpenAI ChatGPT-3.5类似，无需注册/登陆，Llama-3预计打开开源大模型新生态。

- 风险提示**

地缘政治风险；技术迭代不及预期风险；科技巨头竞争加剧风险；下游需求不及预期风险。

分析师

吴砚靖

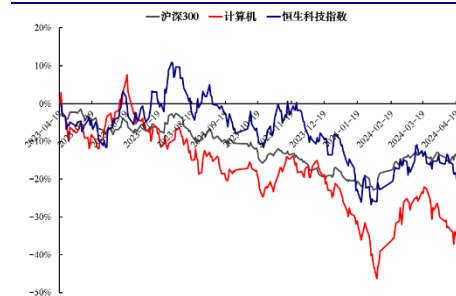
☎：(8610) 66568589

✉：wuyanqing@chinastock.com.cn

分析师登记编码：S0130519070001

国内表现

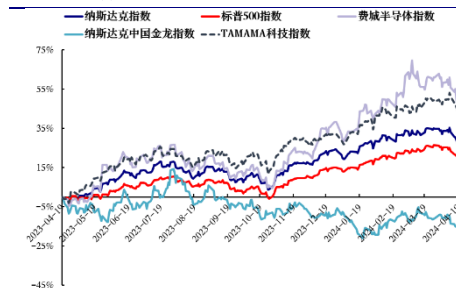
2024-04-19



资料来源：Wind，中国银河证券研究院

全球行情

2024-04-19



资料来源：Wind，中国银河证券研究院

相关研究

【银河计算机】全球科技动态追踪_Gaudi 3、MTIA v2 相继发布，AI产业链加速升级

【银河计算机】全球科技动态追踪_ChatGPT 开放免费注册使用，AI聊天机器人或将颠覆传统搜索引擎

【银河计算机】全球科技动态追踪_OpenAI 公布模型 Voice Engine，AI安全问题再受关注

目 录

一、全球市场表现.....	3
(一) 股市动态.....	3
(二) 债市及汇率情况.....	3
(三) 重点公司表现.....	3
二、行业要闻.....	5
(一) 算力及终端.....	5
(二) 大模型及云应用.....	7
三、风险提示.....	9

一、全球市场表现

(一) 股市动态

近期海外部分地区冲突导致地缘政治风险提升；4月19日，标准普尔全球评级将以色列主权信用评级由AA-下调至A+。受此影响，避险情绪加剧，国际基准布伦特原油价格飙升至90美元/桶，致使风险资产遭抛售，股市承压。

美股三大指数、中概股及港股科技、A股计算机指数全线下跌。标普500指数-3.05%，纳斯达克综合指数-5.52%，费城半导体指数-9.23%；TAMAMA科技指数-7.02%；纳斯达克中国金龙指数-2.09%；恒生科技指数-5.65%；计算机-3.73%。

表1：主要股指周变动

指数代码	指数简称	涨跌幅%					市盈率 PE (TTM)
		本周	上周	本月	本年度	2023	
SPX.GI	标普500指数	-3.05	-1.56	-5.46	4.14	24.23	24.82
IXIC.GI	纳斯达克指数	-5.52	-0.45	-6.70	1.80	43.42	39.60
SOX.GI	费城半导体指数	-9.23	-1.54	-12.20	3.15	64.90	43.98
8884057.WI	TAMAMA科技指数	-7.02	1.03	-5.43	6.89	67.81	34.57
HXC.GI	纳斯达克中国金龙指数	-2.09	-3.22	-4.89	-10.17	-3.39	19.06
HSTECH.HI	恒生科技指数	-5.65	0.68	-5.74	-12.92	-8.83	20.40
CI005027.WI	计算机	-3.73	-4.22	-10.31	-18.20	8.90	78.10

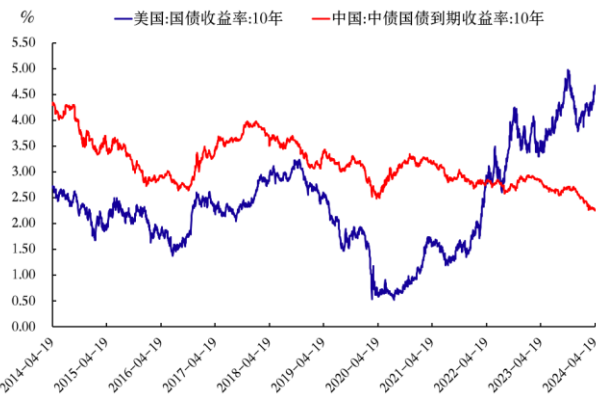
资料来源：Wind，中国银河证券研究院

(二) 债市及汇率情况

美联储发言偏鹰派，市场通胀担忧反复，叠加美国经济数据强于预期，10年期美债持续走弱；考虑利空因素边际影响预期收窄，预计美债利率上涨趋势放缓，但仍在底部区间内震荡。美债利率再走高或使科技股阶段性再承压。周内，美国10年期国债利率持续走高，累计上浮12bps至4.62%；中国10年期国债利率下降至2.25%，累计下降2.97bps。

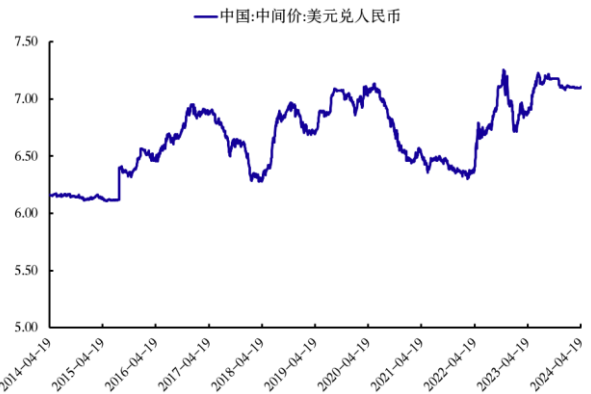
4月19日，美元兑人民币中间价报7.10；较4月12日价累计调升79个基点。

图1：国债收益率（10年期）



资料来源：Wind，中国银河证券研究院

图2：美元兑人民币汇率(中间价)



资料来源：Wind，中国银河证券研究院

(三) 重点公司表现

多重利空因素致使热门科技股巨幅调整。据统计，相比 4 月 12 日收盘价，4 月 19 日盘后，苹果合计-6.54%，英伟达-13.59%，特斯拉-14.03%，谷歌-2.31%，亚马逊-6.18%，Meta-6.02%，微软-5.40%，ARM-30.98%，英特尔-4.17%，高通-7.97%，AMD-10.19%。

美股科技集体伴随市值大幅缩水；英伟达市值跌破 2 万亿美元。美股七大科技公司（苹果、英伟达、特斯拉、谷歌、亚马逊、Meta、微软），单周总市值合计蒸发 9661.66 亿美元（约合 6.9 万亿人民币），刷新历史记录。周内，英伟达市值缩水近 3 千亿美元，市值跌破 2 万亿美元。

短期来看，四月下旬美股科技将陆续披露季报，在此前的高市场预期背景下，各公司业绩相对市场预期具较大不确定性。

表 3：美股七大科技公司总市值变化情况（相对前一周）

公司	总市值(亿美元)		绝对值变化(亿美元)
	2024-04-12	2024-04-19	
苹果	27262.64	25479.10	-1783.54
英伟达	22046.50	19050.00	-2996.50
特斯拉	5447.58	4683.23	-764.35
谷歌	19694.94	19251.99	-442.96
亚马逊	19352.28	18156.61	-1195.68
META	13050.40	12264.42	-785.98
微软	31349.01	29656.36	-1692.65
总计			-9661.66

资料来源：Wind，中国银河证券研究院

二、行业要闻

（一）算力及终端

【SK 海力士与台积电达成战略合作，共同开发 HBM4 芯片】

SK 海力士与 2023 年 4 月 19 日宣布公司就下一代 HBM 产品生产和加强整合 HBM 与逻辑层的先进封装技术，将与台积电公司密切合作。双方近期签署了谅解备忘录。公司计划与台积电合作开发预计在 2026 年投产的 HBM4，即第六代 HBM 产品。

两家公司将首先将致力于针对搭载于 HBM 封装内最底层的基础裸片（Base Die）进行性能改善。HBM 是将多个 DRAM 裸片（Core Die）堆叠在基础裸片上，并通过 TSV1 技术（在 DRAM 芯片打上数千个细微的孔，并通过垂直贯通的电极连接上下芯片的技术）进行垂直连接而成。基础裸片也连接至 GPU，起着对 HBM 进行控制的作用。此外，双方将协力优化 SK 海力士的 HBM 产品和台积电的 CoWoS2（台积电独有的制程工艺，是一种在称为硅中阶层（Interposer）的特殊基板上搭载并连接 GPU、xPU 等逻辑芯片和 HBM 的封装方式。其技术在 2D 封装基板上集成逻辑芯片和垂直堆叠（3D）的 HBM，并整合成一个模组，因此也被称为 2.5D 封装技术）技术融合，共同应对 HBM 相关客户的要求。

【甲骨文宣布在未来 10 年在日本投资超 80 亿美元建设人工智能基础设施】

据路透社 4 月 17 日报道，甲骨文表示将在未来 10 年内投资超过 80 亿美元，以满足日本对云计算和人工智能基础设施的需求。甲骨文在一份声明中表示，最新的投

资将加强该公司在日本的云计算服务—甲骨文云基础设施（OCI）的服务能力和覆盖范围。甲骨文还将扩大其运营范围，并通过雇佣日本人员支持工程团队。

【微软计划到 2024 年底积累 180 万枚人工智能芯片】

据 Business Insider 4 月 18 日消息，微软计划到 2024 年底积累 180 万枚人工智能芯片，将其拥有的 GPU 数量增加两倍。

Business Insider 称，从当前财年到 2027 财年，微软预计将在 GPU 和数据中心上花费约 1000 亿美元。据 Business Insider 援引的 DA Davidson 的分析师估计，微软去年在英伟达芯片上花费了 45 亿美元，一位微软高管表示，这个数字与微软的实际支出大致相符。

微软内部正在努力设计自己的人工智能芯片，以减少对英伟达的依赖，但一些员工持怀疑态度，因为微软落后英伟达多年，而且技术进步如此之快。

【AMD 发布 AI PC 专用处理器：Ryzen Pro 8040 和 Ryzen Pro 8000】

芯片设计商 AMD 于 4 月 16 日推出了一系列用于人工智能商务笔记本电脑和台式机（AI PC）的新处理器：Ryzen Pro 8040 系列（笔记本处理器）、Ryzen Pro 8000 系列系列（台式机处理器）。AMD 称之为商用 PC 迄今为止最强大的芯片，这两个系列的芯片都采用了 4nm 的生产工艺。

据悉，这些新的 AMD 芯片将从 2024 年第二季度开始，为包括惠普和联想等品牌的 PC 机型提供支持。据 CNBC 援引美国咨询公司 Gartner 估计，到 2024 年，AI PC 的出货量将占所有 PC 的 22%。该研究公司预计，到今年年底，AI PC 的出货量将达到 5450 万台。

【Mentee Robotics 展示了旗下首款人形机器人的原型 Menteebot】

据 The Robot Report 2024 年 4 月 17 日报道，初创公司 Mentee Robotics 展示了旗下首款人形机器人的原型 Menteebot，号称在所有操作层都接入了 AI，是“可以被指导的”个性化 AI 机器人。

Menteebot 能够使用人工智能来理解自然语言命令。该机器人的运动基于一种新的机器学习方法，称为模拟到现实（Sim2Real）。Mentee Robotics 声称，该机器人的目标市场有两个：一是家庭应用市场，够执行一系列任务，包括餐桌布置、餐桌清理、衣物处理；二是仓库应用，仓库自动化机器人能够有效地定位、检索和运输物品，并能够处理重达 25 公斤（55 磅）的负载。

Mentee Robotics 表示，它计划 2025 年第一季度之前发布生产完毕的原型机器人。

图3: Menteebot 机器人



资料来源: Mentee Robotics 官网, 中国银河证券研究院

(二) 大模型及云应用

【Meta 公布其最新大语言模型 Llama-3】

Meta 公布其最新大语言模型 Llama 3, 该模型被认为是“首个开源 GPT-4 级别的模型”。Llama-3 目前有两个参数型号, 分别为 8B 和 70B: 1) Llama-3-8B: 拥有 80 亿个参数, 专门为需要快速推理和较少计算资源的应用场景设计, 同时也能保持较高的性能标准; 2) Llama-3-70B: 拥有 700 亿个参数, 能够处理更为复杂的任务, 提供更深入的语义理解, 有更强的生成能力。Meta 声称再未来几个月会推出具有新功能、更长上下文窗口和增强的性能模型, 并且会发布 Llama-3 的技术报告。

Llama-3 与 Llama-2 相比在多个方面有显著提升。1) 在参数规模上, 两种不同规模的参数模型能够提供更丰富更精确的应用场景; 2) Llama-3 的训练数据集比 Llama-2 扩大了 7 倍, 超过 15Ttoken, 其中包含比 Llama-2 多 4 倍的代码数据, 使得 Llama-3 再代码方面表现出色, 且为了即将到来的多语言应用, Llama-3 预训练数据集的 5% 以上由涵盖 30 多种语言的高质量非英语数据组成; 3) 在模型架构上, Llama-3 采用更高效的分词器和分组查询注意力 (Grouped Query Attention, GQA) 技术, 提高了模型的推理效率和处理长文本的能力; 4) 在安全性上, Llama-3 引入了 Llama Guard 2, 以及 Code Shield 和 CyberSec Eval 2 等安全工具, 增强了模型的安全可靠性。

Llama-3 在性能上优于同等级参数模型。经过指令微调后的 Llama-3-8B 模型在 MMLU、GPQA、HumanEval、GSM-K 和 MATH 等测试中优于 Gemma 7B、Mistral 7B。微调后的 Llama-3-70B 在 MMLU、GPQA、HumanEval、GSM-K 和 MATH 等测试中优于 Gemini Pro 1.5 和 Claude 3 Sonnet。

图4: Llama-3 性能测试对比图

Meta Llama 3 Instruct model performance

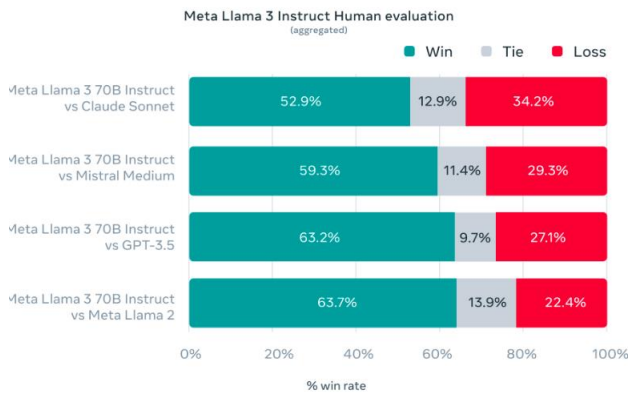
	Meta Llama 3 8B	Gemma 7B - It Measured	Mistral 7B Instruct Measured	Meta Llama 3 70B	Gemini Pro 1.5 Published	Claude 3 Sonnet Published
MMLU 5-shot	68.4	53.3	58.4	82.0	81.9	79.0
GPQA 0-shot	34.2	21.4	26.3	39.5	41.5 CoT	38.5 CoT
HumanEval 0-shot	62.2	30.5	36.6	81.7	71.9	73.0
GSM-BK 8-shot, CoT	79.6	30.6	39.9	93.0	91.7 11-shot	92.3 0-shot
MATH 4-shot, CoT	30.0	12.2	11.0	50.4	58.5 Minerva prompt	40.5

资料来源: Meta 官网, 中国银河证券研究院

请务必阅读正文最后的中国银河证券股份有限公司免责声明。

Meta 开发了一套高质量评估集。该评估集包含 1,800 个提示，涵盖 12 个关键用例：寻求建议、头脑风暴、分类、封闭式问答、编码、创意写作、提取、栖息角色/角色、开放式问答、推理、重写和总结。为了防止模型在这个评估集上意外过拟合，建模团队也无法访问。**通过对 Claude Sonnet、Mistral Medium 和 GPT-3.5 上测试，Llama-3 表现优异，均已超过 50% 的胜率胜出其他模型。**

图5: Llama-3 和同类模型在人类评估集上的测试结果



资料来源: Meta 官网, 中国银河证券研究院

Llama-3 的技术架构基本情况。1) 采用解码器 (decoder-only) 架构, 标准 Transformer 模型架构; 2) Llama-3 采用了具有 128K 个 token 的分词器, 使其拥有更强的性能; 3) 采用 GQA 技术, 该技术将注意力机制中的查询分组, 减少了计算量的同时保持了性能; 4) Llama-3 支持 8192 个 token 的序列, 并且使用掩码技术确保注意力不会跨越文档边界; 5) Llama-3 可以使用指令微调提升模型在特定任务上的表现。

Llama-3 可大幅扩充 Meta AI 应用功能, 实现与 Facebook、Instagram、WhatsApp 等 App 的无缝集成; 此外, 网页版 Meta AI 与 OpenAI ChatGPT-3.5 类似, 无需注册 (登陆), Llama-3 预计打开开源大模型新生态。

【Meta、USC、CMU 和 UCSD 推出无限上下文架构: Megalodon】

本周, Meta、USC (南加州大学)、CMU 和 UCSD 的研究人员突出了具有无限上下文长度的神经网络架构: Megalodon (巨齿鲨)。

Megalodon 有望解决 transformer 有限上下文的难题。现在普遍使用的 transformer 架构在处理上下文时, 会受到二次复杂度, 以及长度外推能力的限制。虽然现已有二次方解决方案, 但它们在预训练效率和下游任务准确性方面的经验表现还不如 Transformer。**Megalodon 基于 Mega 的架构进行了改进, 增加了数个技术组件:** 1) 复杂指数移动 (CEMA) 组件, 可以增强模型处理复杂数据的能力; 2) 时间步长归一化层: 将传统的组归一化技术扩展到自回归序列建模任务中, 允许模型在处理序列数据时, 进行有效的归一化; 3) 为了增强大规模 LLM 预训练的稳定性, 提出了将归一化注意力, 和带有两跳残差的预归一化相结合的配置。

在具体表现上。在与 Llama2 的受控正面比较中, Megalodon 在 70 亿个参数和 2 万个训练令牌的规模上实现了比 Transformer 更好的效率。Megalodon 的训练损失为 1.70, 处于 Llama2-7B (1.75) 和 13B (1.67) 之间。

【OpenAI 开设东京中心，添加针对日本优化的 GPT-4 模型】

OpenAI 在 2024 年 4 月 14 日宣布在东京开设了第一家亚洲办事处，在日本开展业务。

Techcrunch 在报道中强调 OpenAI 此举意义重大。1) 在扩张过程中，OpenAI 可能需要将其技术本地化，适配不同的语言环境；2) 随着人工智能的利弊成为政府、监管机构以及公众讨论的焦点，对 OpenAI 来说，切实了解并影响这些对其有利的趋势变得尤为重要。

针对日本优化后的 GPT-4 模型相比之前性能更佳，运行速度是之前的 3 倍，且增强了对日语中细微差别的理解，包括文化理解等，这对于 GPT-4 的商业应用很重要。目前，OpenAI 正在向一些日本本地企业提供对 GPT-4 自定义模型的早期访问权限，并且“在未来几个月内”通过 OpenAI API 逐步开放访问权限。

三、风险提示

地缘政治风险；技术迭代不及预期风险；科技巨头竞争加剧风险；下游需求不及预期风险。

图表目录

图 1: 国债收益率 (10 年期)	3
图 2: 美元兑人民币汇率(中间价).....	3
图 3: Menteebot 机器人.....	7
图 4: Llama-3 性能测试对比图	7
图 5: Llama-3 和同类模型在人类评估集上的测试结果.....	8

表格目录

表 1: 主要股指周变动.....	3
表 2: 重点公司周数据.....	4
表 3: 美股七大科技公司总市值变化情况 (相对前一周)	5

分析师承诺及简介

本人承诺以勤勉的执业态度，独立、客观地出具本报告，本报告清晰准确地反映本人的研究观点。本人薪酬的任何部分过去不曾与、现在不与、未来也将不会与本报告的具体推荐或观点直接或间接相关。

吴砚靖，TMT/科创板研究负责人。北京大学软件项目管理硕士，10年证券分析从业经验，历任中银国际证券首席分析师，国内大型知名PE机构研究部执行总经理。具备一二级市场经验，长期专注科技公司研究。

免责声明

本报告由中国银河证券股份有限公司（以下简称银河证券）向其客户提供。银河证券无需因接收人收到本报告而视其为客户。若您并非银河证券客户中的专业投资者，为保证服务质量、控制投资风险、应首先联系银河证券机构销售部门或客户经理，完成投资者适当性匹配，并充分了解该项服务的性质、特点、使用的注意事项以及若不当使用可能带来的风险或损失。

本报告所载的全部内容只提供给客户做参考之用，并不构成对客户投资咨询建议，并非作为买卖、认购证券或其它金融工具的邀请或保证。客户不应单纯依靠本报告而取代自我独立判断。银河证券认为本报告资料来源是可靠的，所载内容及观点客观公正，但不担保其准确性或完整性。本报告所载内容反映的是银河证券在最初发表本报告日期当日的判断，银河证券可发出其它与本报告所载内容不一致或有不同结论的报告，但银河证券没有义务和责任去及时更新本报告涉及的内容并通知客户。银河证券不对因客户使用本报告而导致的损失负任何责任。

本报告可能附带其它网站的地址或超级链接，对于可能涉及的银河证券网站以外的地址或超级链接，银河证券不对其内容负责。链接网站的内容不构成本报告的任何部分，客户需自行承担浏览这些网站的费用或风险。

银河证券在法律允许的情况下可参与、投资或持有本报告涉及的证券或进行证券交易，或向本报告涉及的公司提供或争取提供包括投资银行业务在内的服务或业务支持。银河证券可能与本报告涉及的公司之间存在业务关系，并无需事先或在获得业务关系后通知客户。

银河证券已具备中国证监会批复的证券投资咨询业务资格。除非另有说明，所有本报告的版权属于银河证券。未经银河证券书面授权许可，任何机构或个人不得以任何形式转发、转载、翻版或传播本报告。特提醒公众投资者慎重使用未经授权刊载或者转发的本公司证券研究报告。

本报告版权归银河证券所有并保留最终解释权。

评级标准

评级标准	评级	说明
评级标准为报告发布日后的6到12个月行业指数（或公司股价）相对市场表现，其中：A股市场以沪深300指数为基准，新三板市场以三板成指（针对协议转让标的）或三板做市指数（针对做市转让标的）为基准，北交所市场以北证50指数为基准，香港市场以摩根士丹利中国指数为基准。	行业评级	推荐：相对基准指数涨幅10%以上
		中性：相对基准指数涨幅在-5%~10%之间
		回避：相对基准指数跌幅5%以上
公司评级	推荐：相对基准指数涨幅20%以上	
	谨慎推荐：相对基准指数涨幅在5%~20%之间	
	中性：相对基准指数涨幅在-5%~5%之间	
	回避：相对基准指数跌幅5%以上	

联系

中国银河证券股份有限公司研究院

深圳市福田区金田路3088号中洲大厦20层

上海浦东新区富城路99号震旦大厦31层

北京市丰台区西营街8号院1号楼青海金融大厦

公司网址：www.chinastock.com.cn

机构请致电：

深广地区：程曦 0755-83471683chengxi_yj@chinastock.com.cn

苏一耘 0755-83479312suyiyun_yj@chinastock.com.cn

上海地区：陆韵如 021-60387901luyunru_yj@chinastock.com.cn

李洋洋 021-20252671liyongyang_yj@chinastock.com.cn

北京地区：田薇 010-80927721tianwei@chinastock.com.cn

唐嫚羚 010-80927722tangmanling_bj@chinastock.com.cn