

研究所：

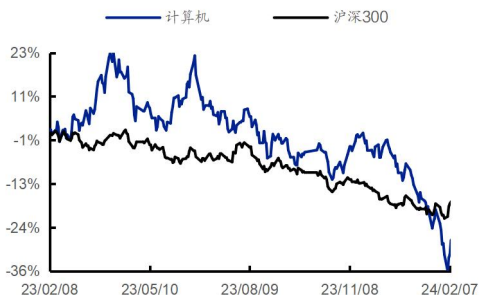
证券分析师：

刘熹 S0350523040001
liux10@ghzq.com.cn

2024 年互联网 AI 资本开支持续提升

——AI 算力“卖水人”系列专题（1）

最近一年走势



行业相对表现

2024/02/08

| 表现 | 1M | 3M | 12M |
|--------|--------|--------|--------|
| 计算机 | -15.1% | -26.1% | -27.9% |
| 沪深 300 | 2.4% | -6.8% | -17.4% |

相关报告

《科技央企国企改革系列研究：央企市值管理考核将全面推开，关注科改新机遇（推荐）*计算机*刘熹》——2024-01-30

《计算机事件点评：超微电脑业绩上修，AI 服务器兑现或超预期（推荐）*计算机*刘熹》——2024-01-21

《AI 算力月度跟踪报告（202401）：英伟达 H20 或将量产，液冷服务器加速渗透（推荐）*计算机*刘熹》——2024-01-21

《计算机行业动态研究：鸿蒙 NEXT 进阶发布，原生生态千帆启航（推荐）*计算机*刘熹》——2024-01-19

《计算机行业专题研究：液冷：算力+双碳提振需求，将迎来规模化推广（推荐）*计算机*刘熹》——2024-01-15

核心逻辑：

北美互联网巨头上调 2024 年 AI 资本开支指引，验证大模型竞赛对 AI 算力的需求提升。我们预计随着大模型的迭代，以及 AI 应用用户规模的增长，AI 算力市场将打开更大的市场空间。

投资要点：

■ 起源：互联网资本开支投向软硬件、云计算等，算力是主要组成

互联网资本开支主要指在互联网企业的 ICT 相关支出，根据 Gartner，2024 年全球 IT 支出预计达 5.1 万亿美元，同比+8%，其中算力是开支的主要组成。大模型的竞赛将驱动互联网资本开支持续增长，根据 Dell’Oro Group，预期 2027 年，AI 基础设施支出将推动数据中心资本支出超 5000 亿美元，全球超过 20% 的服务器部署可能会是加速类型。

■ 转折：ChatGPT 等大模型发布提升算力需求，加大互联网资本开支

1) 海外：①微软：OpenAI（微软）发布预训练大模型 GPT-4V 及 GPTs，微软预期 FY2024Q3 资本支出实现实质性增长；②谷歌：发布对话式 AI 服务 Bard 以及多模态大模型 Gemini，谷歌预计 2024 年资本支出逐渐提升；③Meta：推出开源大语言模型 LLaMA2，Meta 2024 年资本支出指引调整上限至 300-370 亿美元。

2) 国内：截至 2024 年 1 月，国内已经有 42 款大模型通过备案审批。百度文心大模型 4.0 全面升级，参数规模预计超万亿；阿里通义千问 2.0 发布，参数达千亿级；字节跳动“Coze 扣子”AI Bot 开发平台上线。2023Q3，BAT 均表示将持续加大人工智能投入，并表示不仅陆续将 AI 大模型赋能主营业务，还表示将持续加大人工智能模型的投入。

■ 展望：大模型迭代与应用创新将提升算力需求、提升互联网资本开支

1) 训练：“大模型+大数据”成为预训练模型“新范式”，根据 SemiAnalysis，GPT-4 共包含 1.8 万亿参数，AI 算力将从单机走向集群时代。我们认为，随着多模态模型的持续发展，模型训练所消耗的算力有望继续提升。

2) 推理：1 月 11 日，OpenAI 宣布已构建 300 万 GPTs。Meta、微软、谷歌等具备海量用户基础，大模型应用潜在渗透空间高。国内市场，科大讯飞表示，截至 1 月 9 日，星火小助手创建数已突破 51000 款。IDC 预计，到 2027 年，用于推理的工作负载将达到 72.6%。

3) 互联网 AI 资本开支有望持续提升。根据 Counterpoint，预计 2023 年全球多数云服务商的 AI 开支占总资本支出的比例仅在 3-7% 区间内。

■ **行业评级及投资策略：**AI 算力有望持续增长，GPU、AI 服务器、服务器散热、光模块、数据中心等环节有望持续受益。维持对计算机行业“推荐”评级。

■ **相关公司：**

1) **服务器整机：**工业富联、浪潮信息、中科曙光、紫光股份、华勤技术、神州数码、烽火通信、中国长城、高新发展、纬创、广达、英业达。

2) **服务器组件：**①**AI 芯片：**海光信息、寒武纪、龙芯中科、英伟达、AMD；②**散热：**飞荣达、中航光电、曙光数创、英维克、奇鋆科技、双鸿、建准电机；③**主板：**沪电股份、深南电路、胜宏科技、技嘉、华擎。

3) **光模块：**中际旭创、新易盛、光迅科技、华工科技。

4) **数据中心：**奥飞数据、光环新网、宝信软件、数据港、电科数字。

■ **风险提示：**宏观经济影响下游需求、信创政策不及预期、市场竞争加剧、中美博弈加剧、相关公司业绩不及预期等。各公司并不具备完全可比性，对标的相关资料和数据仅供参考。

内容目录

| | |
|--|----|
| 1、 互联网资本开支投向软硬件、云计算等，算力是主要组成 | 5 |
| 2、 大模型驱动算力需求，海内外企业扩大资本开支 | 6 |
| 2.1、 需求提升：AI 芯片性能提升可降 GPT4 训推成本 | 6 |
| 2.2、 海外：微软/谷歌/META 等大模型推进，资本开支指引显著增长 | 8 |
| 2.3、 国内：BAT 积极布局大模型，资本开支有望回升 | 9 |
| 3、 未来：训推算力需求提升，互联网具用户+资金优势，AI 资本开支有望提升 | 11 |
| 3.1、 “大模型+大数据”新范式，训练算力需求提升 | 11 |
| 3.2、 大模型应用创新，驱动推理需求提升 | 12 |
| 3.3、 大厂用户+资金实力深厚，AI 资本支出提升空间大 | 13 |
| 4、 相关公司 | 14 |
| 5、 风险提示 | 15 |

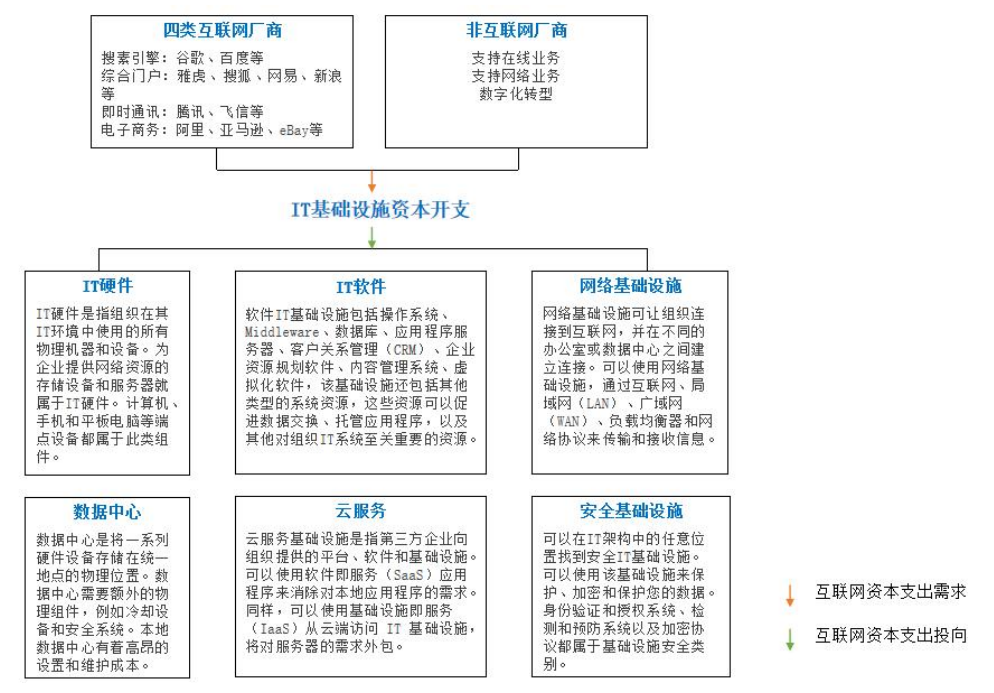
图表目录

| | |
|--|----|
| 图 1: 互联网资本开支介绍 | 5 |
| 图 2: 全球 IT 资本支出预期及分类情况 | 6 |
| 图 3: 2024E 全球 IT 资本支出分类占比情况 | 6 |
| 图 4: 2018-2025E 全球云厂商资本支出情况 | 6 |
| 图 5: 2022 年全球 AI 服务器采购量占比 | 6 |
| 图 6: LLM 所需带宽需求 | 7 |
| 图 7: NVIDIA GPU 性能对比 | 7 |
| 图 8: 2023 年, Meta 和微软预计购买 15 万块 H100 | 8 |
| 图 9: 全球服务器市场规模持续增长 | 8 |
| 图 10: 海外大模型训练和应用进展 | 8 |
| 图 11: 微软、谷歌、META 资本性支出 | 9 |
| 图 12: 阿里、百度、腾讯、字节大模型发展情况 | 10 |
| 图 13: 大模型备案审批的企业名单 | 10 |
| 图 14: 大模型+大数据成为 AI 预训练模型“新范式” | 12 |
| 图 15: 中国智能算力规模及预测(单位: EFLOPS) | 12 |
| 图 16: 大模型训练算力当量, 大模型训练对算力需求将更大 | 12 |
| 图 17: 大模型应用场景 | 13 |
| 图 18: 中国人工智能服务器工作负载预测 | 13 |
| 图 19: 全球部分大型云厂商经营性现金流量情况 | 13 |
| 图 20: 2024 年 1 月最受欢迎社交网络 | 13 |
| 图 21: 2023E 全球云服务厂商资本开支市场份额 | 14 |
| 图 22: 2023E 云服务商 AI 开支占总资本开支比例 | 14 |
| 图 23: 服务器产业链相关公司 | 15 |
| 表 1: 全球部分 AI 芯片厂商基本情况一览 | 7 |
| 表 2: BAT 资本开支情况 | 11 |

1、互联网资本开支投向软硬件、云计算等，算力是主要组成

互联网资本开支（Internet Capital Expenditure）指互联网企业在互联网和信息技术基础设施上的支出，包括公司投资于硬件、软件、网络设备、数据中心、云计算服务等方面的支出，旨在建设、维护和升级与互联网相关的基础设施、提高服务质量与效率。互联网资本开支是企业为了支持其在线业务、网络服务和数字化转型而进行的重要投资。

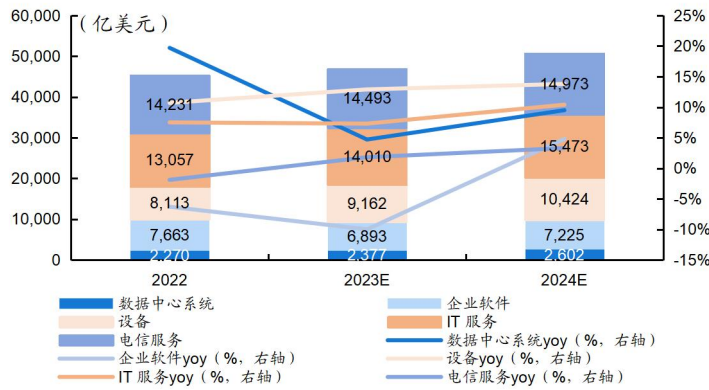
图 1：互联网资本开支介绍



资料来源：AWS 官网，国海证券研究所

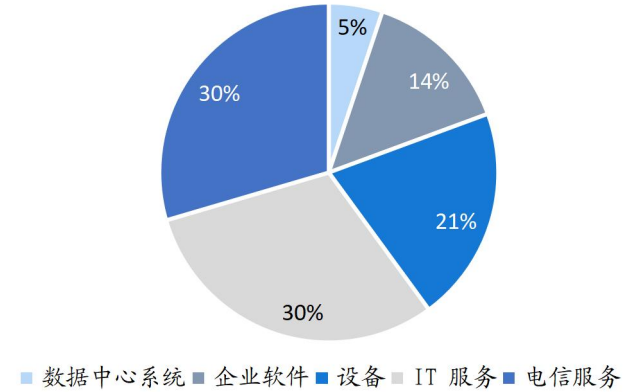
全球 IT 资本支出持续增长，预期 AI 服务器采购量将进一步提升。根据 Gartner 的最新预测，2024 年全球 IT 支出预计将达 5.1 万亿美元，同比+8%。根据 Dell'Oro Group，2022 年，十大云服务商（阿里巴巴/亚马逊/苹果等）在数据中心基础设施共投入 1140 亿美元，占云服务商支出最大份额。Dell'Oro Group 预期到 2027 年，人工智能基础设施支出将推动数据中心资本支出超过 5000 亿美元，全球数据中心资本支出预计将增长 15%，超过 20% 的全球服务器可能会加快部署。

图 2：全球 IT 资本支出预期及分类情况



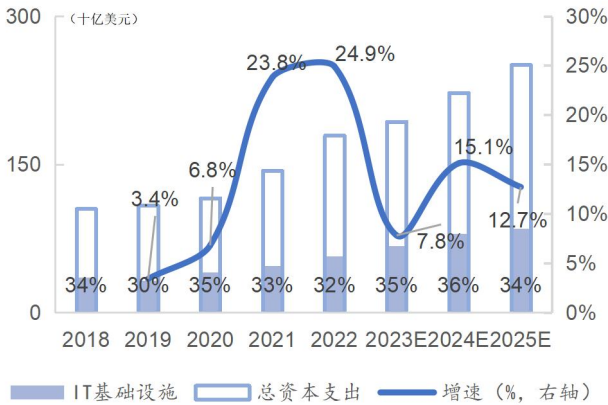
资料来源：Gartner，国海证券研究所

图 3：2024E 全球 IT 资本支出分类占比情况



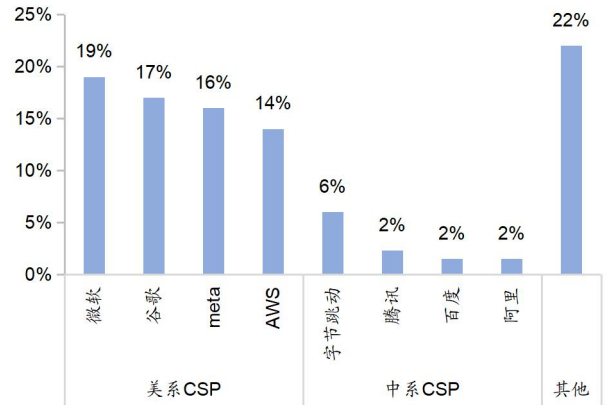
资料来源：Gartner，国海证券研究所

图 4：2018-2025E 全球云厂商资本支出情况



资料来源：Counterpoint，国海证券研究所

图 5：2022 年全球 AI 服务器采购量占比



资料来源：Trendforce，国海证券研究所

注：数据四舍五入原因各份额相加或不为 100%

2、大模型驱动算力需求，海内外企业扩大资本开支

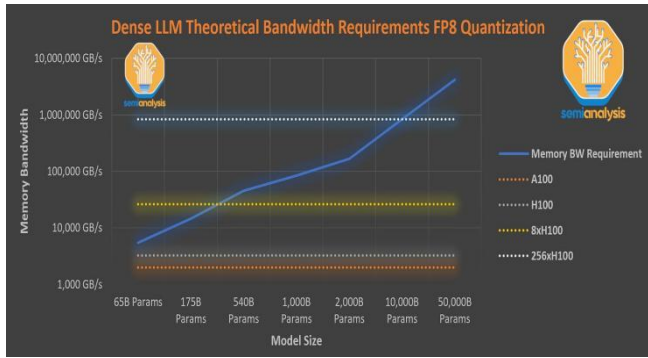
2.1、需求提升：AI 芯片性能提升可降 GPT4 训推成本

GPT4 的参数规模过万亿，带动 AI 芯片需求提升。根据 Semianalysis, ChatGPT4 在 120 层中总共包含了 1.8 万亿参数，而 GPT-3 仅约 1750 亿参数。GPT-4 的参数规模或是 GPT-3 的 10 倍以上，与 175B 参数的 Davinchi 模型相比，GPT-4 的成本是其 3 倍，这表明参数量提升可带动 AI 芯片需求。

1) 训练：根据 Semianalysis, OpenAI 训练 GPT-4 时，在约 25000 个 A100 上训练了 90-100 天，1 美元/每 A100 小时，训练成本约 6300 万美元；但在 2 美元/每 H100 小时，预训练在大约 8192 个 H100 进行，只需 55 天，费用为 2150 万美元。**2) 推理：**128 个 A100 来推理 GPT-4 8k 序列长度，每 1k tokens 成本

0.0049 美分，128 个 H100 推理 GPT-4 8k 序列长度，每 1k tokens 成本 0.0021 美分。我们认为，这表明 AI 芯片的计算性能、内存带宽的提升，有助于提高大模型训推效率与降低成本，也促使大模型厂商投资新一代的 AI 芯片。

图 6: LLM 所需带宽需求



资料来源: semianalysis, 国海证券研究所

图 7: NVIDIA GPU 性能对比

| 性能参数 | V100 PCIe | A100 80GB PCIe | A800 80GB PCIe | H100 80GB PCIe |
|------------------|----------------------------------|------------------------------------|-------------------------------------|--------------------------------------|
| 微架构 | Volta | Ampere | | Hopper |
| FP64 | 7TFLOPS | 9.7TFLOPS | | 26 TFLOPS |
| FP32 | 14TFLOPS | 19.5TFLOPS | | 51 TFLOPS |
| FP16 Tensor Core | | 312TFLOPS | | 756.5 TFLOPS |
| INT8 Tensor Core | 62 TOPS | 624 TOPS | | 1513 TOPS |
| GPU 显存 | 32/16GB HBM2 | 80GB HBM2e | | 80GB |
| GPU 显存带宽 | 900 GB/s | 1935GB/s | | 2TB/s |
| 最大热设计功耗(TDP) | 250W | 300W | | 300-350W |
| 互连技术 | NVLink:300 GB/s PCIe: 32 GB/s | NVLink:600GB/s PCIe 4.0: 64GB/s | NVLink: 400GB/s PCIe 4.0: 64GB/s | NVLink: 600GB/s PCIe 5.0: 128GB/s |

资料来源: NVIDIA, 国海证券研究所

云服务提供商对英伟达 GPU 的需求较强。Omdia 数据显示，2023 年，Meta 和微软两家公司以 15 万块 H100 GPU 的购买量并列位居第一，大多数服务器 GPU 都直接供应给了超大规模云服务供应商。考虑到 B100 的计算力预计可达 H100 的两倍以上，微软或决定 2024 年采购 B100 芯片。

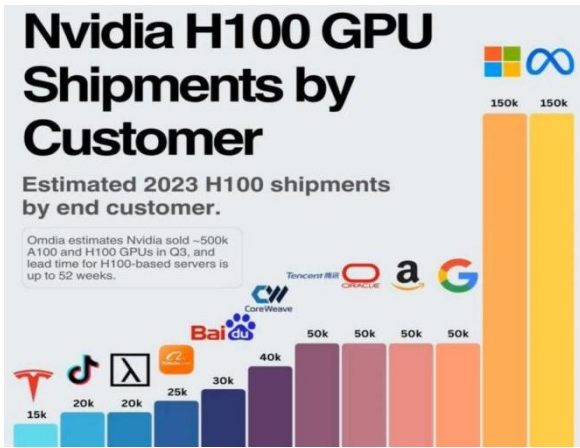
1 月 18 日，Meta 首席执行官扎克伯格宣布，为了搭建能够支持 AGI 愿景的基础设施，Meta 计划在 2024 年年底获得约 35 万块英伟达 H100 GPU，再算上其他 GPU，拥有算力总和接近 60 万块 H100 所能提供的算力。

表 1: 全球部分 AI 芯片厂商基本情况一览

| 厂商 | 主打产品 | 单价预测 | 产能预测 | 2023 年出货量预测 | 应用领域 | 服务器合作代工厂商 | 核心客户 |
|-----|--------------------------|---|---------------------------------------|----------------------------|-------------|------------------------------------|-----------------------------------|
| 英伟达 | H100/A100 H200(2024)等 | H100 3.65 万美元/块 A100 1.5 万美元/块 | 约 100 万块/年 2024 年 H100 将达 200 万块/年 | >150 万块 其中 H100 达 55 万块 | HPC 及 AI 领域 | 浪潮、戴尔、HPE、联想、新华三及工业富联、英业达、广达等 | 微软、亚马逊、谷歌、Meta 及阿里、腾讯、字节跳动、百度、美团等 |
| AMD | MI300A MI300X 等 | MI300X >3 万美元/块 | >15 万块/年 | >18 万块 | HPC 及 AI 领域 | 戴尔、HPE、浪潮、联想、Supermicro、IBM 及工业富联等 | 微软(约 50%)、亚马逊、谷歌、Meta 在测试 |
| 英特尔 | Gaudi2 等 | 约 0.8-1 万美元/块 | >10 万块/年 | >8 万块 | HPC 及 AI 领域 | 浪潮信息、新华三及工业富联、神达等 | 微软、阿里、美团等 |

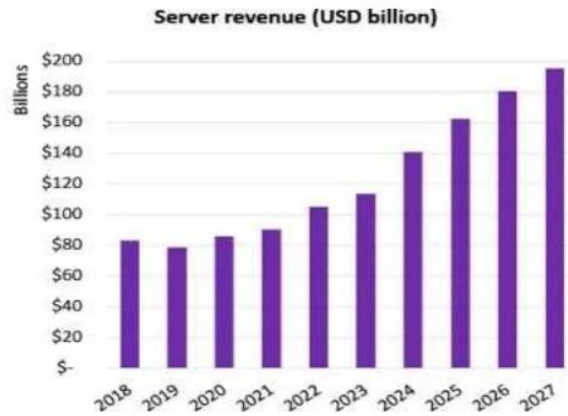
资料来源: IDC, 各公司财报, 芯八哥, 国海证券研究所

图 8: 2023 年, Meta 和微软预计购买 15 万块 H100



资料来源: Omdia Research, IT 之家

图 9: 全球服务器市场规模持续增长



资料来源: Omdia Research, IT 之家

2.2、海外: 微软/谷歌/META 等大模型推进, 资本开支指引显著增长

复盘海外代表厂商发展脉络, 技术形态上大模型围绕多模态和 AI Agent 两大方向升级迭代。以 OpenAI 为例, 陆续发布多模态预训练大模型 GPT-4V 以及 AI Agent 初级形态产品 GPTs。大语言模型 (LLM) 在 2023 年为逼近通用人工智能 (AGI) 提供了一个可能路径。1 月 11 日, OpenAI 宣布, GPT Store 正式上线, 已有超 300 万自定义版本 ChatGPT 应用。2024 年, 随着技术的迭代进步和大模型的深度普及, AGI 时代或即将到来。

图 10: 海外大模型训练和应用进展

| | 2022年11月 | 2023年2月 | 2023年3月 | 2023年4-6月 | 2023年7-9月 | 2023年10-12月 |
|---------------|------------|--------------------------------|--|--|---|---|
| OpenAI | 发布 ChatGPT | | <ul style="list-style-type: none"> 开放ChatGPT和Whisper模型的API 发布多模态预训练大模型 GPT4 | <ul style="list-style-type: none"> 向ChatGPT Plus用户开放联网, 并推出Plugin API新增函数调用, 4倍上下文, 降价25%-75% 推出ChatGPT iOS App | <ul style="list-style-type: none"> 全面开放GPT-4 API 推出安卓版ChatGPT APP 上线企业版ChatGPT GPT-3.5 Turbo支持微调 发布DALL-E 3和GPT-4V | <ul style="list-style-type: none"> 发布GPT-4 Turbo, 价格降低2-4倍, 上下文窗口提升16倍 GPTs(套壳工具)和Assistant API 1亿周活, 200万开发者 |
| 谷歌 | | 宣布推出由LaMDA模型支持的对话式AI服务 Bard | <ul style="list-style-type: none"> 发布视觉语言模型 PaLM-E 开放PaLM API, 并发布 Generative AI App Builder | <ul style="list-style-type: none"> 发布PaLM 2和由PaLM2驱动搜索引擎 发布Duet AI办公全家桶 全面上线基于Vertex AI的生成式人工智能服务 | <ul style="list-style-type: none"> 推出AI芯片TPU v5e 4-7月生成式AI客户项目数量增长150多倍 | <ul style="list-style-type: none"> 官宣其认为规模最大、功能最大的多模态大模型 Gemini Pixel 8 Pro首次搭载Gemini Nano 推出面向云端AI加速的TPU v5p |
| 微软 | | 开放API, 将AI功能嵌入搜索引擎Bing和Edge浏览器 | <ul style="list-style-type: none"> 推出Microsoft 365 Copilot, 由GPT-4驱动 新必应日活破亿 Github推出Copilot X计划, 将ChatGPT引入IDE | <ul style="list-style-type: none"> 推出Azure OpenAI 国际版 全面开放BingChat, 开放第三方插件 | <ul style="list-style-type: none"> 推出语音合成模型Natural Speech2 推出基于微软国际版Azure AI 数字员工方案 | <ul style="list-style-type: none"> 推出AI芯片Azure Maia和云原生CPU Azure Cobalt 发布Copilot Studio Copilot for Microsoft 365面向企业开放商用 |
| Meta | | 推出AI语言模型 LLaMA | | <ul style="list-style-type: none"> 发布SAM, 可准确识别图像中的对象 开源DINOv2, 可直接从图像中学习特征, 而不依赖文本描述 | <ul style="list-style-type: none"> 推出新一代开源大型语言模型 LLaMA2, 并免费开放商用 推出一种高效且可扩展的任意模态增强语言模型AnyMAL | |
| 其他 | | | <ul style="list-style-type: none"> 前OpenAI员工创办的Anthropic推出Claude | <ul style="list-style-type: none"> 亚马逊推出Bedrock服务和大语言模型Titan | <ul style="list-style-type: none"> Anthropic推出Claude2, 支持上传文件, 响应文本更长 苹果正开发推进“Apple GPT”模型 亚马逊推出Agents for Bedrock | <ul style="list-style-type: none"> Anthropic推出Claude2.1, 上下文长度翻番到20万 马斯克旗下xAI团队发布首个AI大模型产品Grok 英伟达推出芯片设计大模型 ChipNeMo |

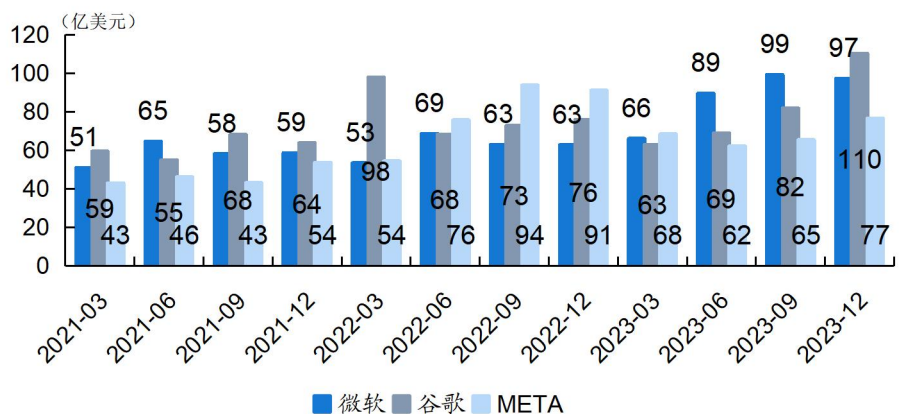
资料来源: OpenAI, 微软, Meta, 36 氪, techweb, IT 之家, 同花顺财经, 华尔街见闻, 澎湃, 芯智讯, 国海证券研究所

微软资本支出持续提升，预期 FY2024Q3 实现实质性增长。FY2024Q2，公司资本支出包括根据融资租赁获得的资产为 115 亿美元，以支持公司的云和 AI 产品的需求，并且微软为财产和设备支付的现金为 97 亿美元。FY2024Q3，公司预计资本开支将出现季环比的“实质性”增长，此次增长反映微软在扩大业务、提升技术能力和优化设施等方面的持续投资。

人工智能投入力度较大，Alphabet(谷歌母公司)资本支出逐渐提升。FY2023Q4，公司的资本支出增长 45%，达到 110.19 亿美元，创多年来最高水平，公司表示其中绝大多数是对技术基础设施的投资，并非是一次性支出，而是对公司未来的投资，包括但不限于 Google DeepMind、谷歌服务、谷歌云、广告商、用户、开发商、企业云、政府等。Alphabet 预期其 2024 年支出将“显著增加”，因为该公司还将继续投资服务器、数据中心等基础设施，为其人工智能产品铺路。

1 月 19 日，据 Meta CEO 扎克伯格表示，META 计划斥资数十亿美元购买英伟达计算机芯片，旨在加强 Meta 在人工智能研究和相关项目上的能力。针对 2024 年资本开支上，Meta 小幅调高了指引上限，从 300-350 亿区间，调到 300-370 亿美元，因为要增加一些 AI 等创新业务等投入。

图 11: 微软、谷歌、META 资本性支出



资料来源：iFind，国海证券研究所。注：此处数据为各公司现金流量表中数据，或与上文字数据有出入，主要系文字部分数据包含根据融资租赁获得的资产等部分

2.3、国内：BAT 积极布局大模型，资本开支有望回升

国内互联网厂商持续推动，国产大模型研发加速。截至 2023 年 11 月，国产大模型有 188 个，其中通用大模型 27 个；截至 2024 年 1 月，已经有 42 款大模型通过备案审批。从 2023 年 8 月第一批企业和机构算起，已公布第四批的企业名单，形成比较稳定的批复节奏。

互联网巨头有望保持领先地位。百度、阿里、腾讯等企业积极布局大模型赛道，加速大模型领域投入。百度文心大模型 4.0 全面升级，据 IT 之家，参数规模预计超万亿，综合水平可比拟 GPT-4；阿里通义千问 2.0 发布，参数达千亿级，综合性能将加速赶超 GPT-4。

图 12：阿里、百度、腾讯、字节大模型发展情况

| | 2019年 | 2021-2022年 | 2023年4-6月 | 2023年8-9月 | 2023年10月-至今 |
|----|--|--|--|---|---|
| 阿里 | <ul style="list-style-type: none"> 第一颗自研AI芯片含光800发布，推理性能达78563IPS | <ul style="list-style-type: none"> 自研CPU芯片倚天710发布采用5nm工艺，单芯片容纳高达600亿晶体管 倚天710已大规模应用，未来两年20%的新增算力将使用自研CPU | <ul style="list-style-type: none"> 通义千问1.0上线，邀请客户进行测试 一站式大模型应用开发平台“百炼”全新亮相 | <ul style="list-style-type: none"> 通义千问70亿参数模型Qwen-7B开源 陆续开源视觉理解模型Qwen-VL、Qwen-14B等 通义千问1.0向公众开放 | <ul style="list-style-type: none"> 通义千问2.0发布，参数规模达千亿级，综合性能超过Llama-2-70B和GPT-3.5，加速追赶GPT-4 |
| 腾讯 | <ul style="list-style-type: none"> 启动编解码芯片“沧海”研发 | <ul style="list-style-type: none"> 发布AI推理芯片“紫霄”、视频转码芯片“沧海”和智能网卡芯片“玄灵” | <ul style="list-style-type: none"> 视频编解码芯片“沧海”已量产并投用数万片 | <ul style="list-style-type: none"> 自研混元大模型发布，并正式通过腾讯云对外开放 参数规模超千亿级，预训练语料超2万亿tokens | <ul style="list-style-type: none"> 面向C端用户陆续开放体验，已有超180个内部业务接入混元 混元大模型开放文生图功能 |
| 百度 | <ul style="list-style-type: none"> 3月，文心大模型1.0发布 7月，升级至2.0 | <ul style="list-style-type: none"> 7月，文心大模型3.0发布，首个知识增强百亿参数大模型 12月，ERNIE-ViLG全球最大中文预训练生成模型发布 | <ul style="list-style-type: none"> 文心大模型升级至3.5，与3.0相比，模型效果累计提升超过50%，训练速度提升2倍，推理速度提升了17倍 | | <ul style="list-style-type: none"> 文心大模型4.0发布，参数规模预计超万亿，综合水平可比拟GPT-4 文心4.0整体效果提升32% |
| 字节 | | <ul style="list-style-type: none"> 发布大规模 Training & Serving 方案 Monolith 自研K8s 存储 KubeBrain发布，突破etcd限制，支撑线上超过20,000节点的超大规模 Kubernetes 集群的稳定运行 | <ul style="list-style-type: none"> 开源 Shmipc，零拷贝，引入的同步机制具有批量收割IO的能力 开源分布式训练调度框架 Primus | <ul style="list-style-type: none"> 云雀大模型发布，参数达1300亿 发布基于云雀开发的AI对话产品“豆包”，预置了英语学习助手和写作助手两个功能 | <ul style="list-style-type: none"> 2024.02.01，“Coze扣子”AI Bot开发平台正式上线，具有无限拓展的能力集、丰富的数据源等优势 |

资料来源：各公司大模型公众号，DoNews，DeepTech 深科技，观察者网，中国新闻周刊网，芯智讯，快科技，财联社，证券时报，智东西，大模型之家，澎湃新闻，腾讯网，福布斯中国，未来智库，锦缎研究院，东方财富网，新京报，深科技，国海证券研究所

图 13：大模型备案审批的企业名单

| 批次 | 公司 | 大模型产品 | 批次 | 公司 | 大模型产品 |
|-----|-----------|-----------------|---------------|------------|--------------|
| 第一批 | 百度 | 文心一言 | 第三批 | 京东 | 言犀大模型 |
| | 抖音 | 云雀大模型 | | 抖音 | 福祿瓜视觉大模型 |
| | 智谱AI | GLM大模型 | | 快手 | 快意大模型 |
| | 百川智能 | 百川大模型 | | 红棉小冰科技 | 小冰大模型 |
| | 商汤 | 日日新大模型 | | 聆心智能 | CharacterGLM |
| | MiniMax | ABAB大模型 | | 澜舟科技 | 孟子GPT |
| | 中科院 | 紫东太初大模型 | | 中科闻歌 | 雅意大模型 |
| | 上海人工智能实验室 | 书生通用大模型 | | 深言科技 | 语鲸大模型 |
| 第二批 | 美团 | 暂无公开信息 | 云知声 | 山海大模型 | |
| | 蚂蚁集团 | 百灵大模型 | 零一万物 | 零一万物大模型 | |
| | 知乎 | 知海图AI | 第四范式 | 式说大模型 | |
| | 出门问问 | 序列猴子大模型 | 衍远科技 | 品商大模型 | |
| | 昆仑万维 | 天工大模型 | 识音智能 | 慕小仙大模型 | |
| | 月之暗面 | Moonshot AI大模型 | 新壹科技 | 一叶轻舟大模型 | |
| | 面壁智能 | 面壁Luca大模型 | 创思远达 | 新壹视频大模型 | |
| | 网易有道 | 子曰大模型 | 步刻科技 | 魔方大模型 | |
| | 好未来 | 九章大模型 (MathGPT) | BOSS直聘 | 微步情报智脑大模型 | |
| | 金山办公 | WPS AI | 智联招聘 | 南北图大模型 | |
| | 360公司 | 奇元大模型 | 脉脉 | “AI改简历”新功能 | |
| | | | 小米 | “智能问答”新功能 | |
| | | 什么值得买 | “小爱同学AI助手”新功能 | | |
| | | 掌阅 | AI问答机器人 | | |
| | | | “阅爱聊”微信小程序 | | |

资料来源：智能涌现，国海证券研究所

大模型规划下，互联网厂商资本开支有望回升。百度、阿里、腾讯等公司密集发布大模型，2023Q3，BAT 均表示将持续加大人工智能投入，并表示不仅陆续将 AI 大模型赋能主营业务，还表示将持续加大人工智能模型的投入。其中，阿里巴巴集团首席执行官吴泳铭表示，阿里明确战略重心和优先级方向，抓住 AI 科技变革带来的全新机会，创造更多的客户价值。

表 2: BAT 资本开支情况

| 公司 | 资本开支情况 |
|----|---|
| 阿里 | ● 2023Q3，资本开支达 51.50 亿元，公司表示，AI 相关的投入和探索仍将继续，并计划将 ROIC（资本回报率）提升至双位数水平。 |
| 腾讯 | ● 2023Q3，资本开支达 80 亿元，同比+237%。3Q23 业绩会上公司表示，2023 年经营性资本支出占比 3%-3.5%，如果 2024 年可以获得更多 GPU，预计会在此基础上增加 1 个百分点。 |
| 百度 | ● 2023Q3，资本开支达 35.29 亿元，环比+30.41%，资本开支力度加大，坚定 AI 方向投入。根据 Donews，公司表示在未来几个季度里，将继续坚定投资人工智能领域，特别是在大语言模型和生成式人工智能方面。 |

资料来源：东方财富网、各公司公告、海豚投研、Donews、Canalys、腾讯网、福布斯中国、未来智库、锦缎研究院、新京报、深科技，国海证券研究所

国内算力供应链逐渐改善。根据路透社等，2 月，英伟达已开始接受经销商对中国特供版 AI 芯片 H20 的预订，每块卡 12000-15000 美元（8.6-10.77 万元人民币），产品形态包括计算卡和搭载 8 张 H20 计算卡的服务器，华为 910B 的售价约为 12 万元人民币。我们认为，后续需关注英伟达 H20 性价比及国内大模型厂商对其需求情况，以及华为昇腾 910B 等版本的供给与适配进展。

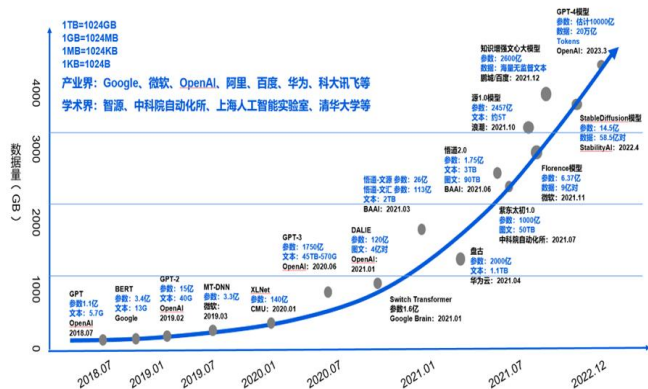
3、未来：训推算力需求提升，互联网具用户+资金优势，AI 资本开支有望提升

3.1、“大模型+大数据”新范式，训练算力需求提升

“大模型+大数据”成为预训练模型的“新范式”。近年新推出的大语言模型所使用的数据量和参数规模呈现“指数级”增长，根据 SemiAnalysis，GPT-4 共包含 1.8 万亿参数，GPT-3 只有约 1750 亿参数，参数量级提升带来算力消耗增加。我们预计未来大模型所消耗的算力将远超过目前已有的大模型。

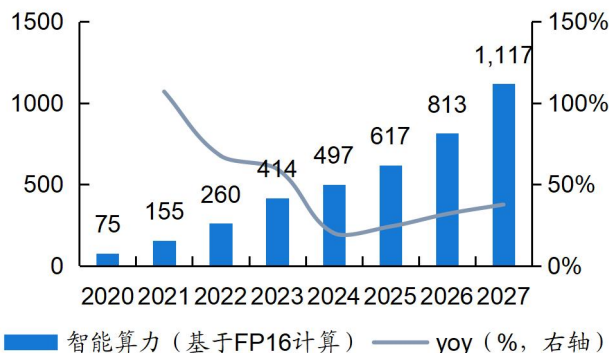
大模型对算力的需求将指数级增长，AI 算力将从单机走向集群时代。大算力集群是业界厂家 AI 模型开发的共同选择，如腾讯数万卡规模的星海 AI 集群，华为数万卡规模昇腾 AI 集群等，AI 芯片、AI 服务器等的需求量提升。

图 14: 大模型+大数据成为 AI 预训练模型“新范式”



资料来源: A Survey of Large Language Models, 新智元, 国海证券研究所

图 15: 中国智能算力规模及预测 (单位: EFLOPS)



资料来源: IDC, 国海证券研究所

图 16: 大模型训练算力当量, 大模型训练对算力需求将更大

| 模型名称 | BERT-Large | GPT-2 | GPT-3 | T-5 | MT-NLG | PaLM | PaLM-2 | Switch-Transformer | Chinchilla | LLaMA | 源1.0 |
|------|------------|-------|--------|------|--------|---------|---------|--------------------|------------|--------|--------|
| 参数量 | 3亿 | 15亿 | 1750亿 | 110亿 | 5300亿 | 5400亿 | 3400亿 | 1.6万亿 | 700亿 | 650亿 | 2450亿 |
| 算力当量 | 2.4PD | 8.7PD | 3640PD | 26PD | 9900PD | 29000PD | 85000PD | 46PD | 6795PD | 6330PD | 4095PD |

资料来源: IDC, 浪潮信息 (注: PD = PetaFlops/s-day)

3.2、大模型应用创新, 驱动推理需求提升

GPT 商店内 AI 应用数量增加, 提升大模型推理需求。1月11日, OpenAI 宣布正式推出 GPT Store 和 ChatGPT Team 服务, 社区成员已经构建了 300 万个 GPTs, 已有 91100 个 GPTs 涌现于公共网络。科大讯飞表示, 以讯飞星火为核心、基于国产自主底座的“GPT 商店”蓬勃发展, 截至 1月9日, 星火小助手创建数已突破 51000 款, 不断助力企业数字新生。

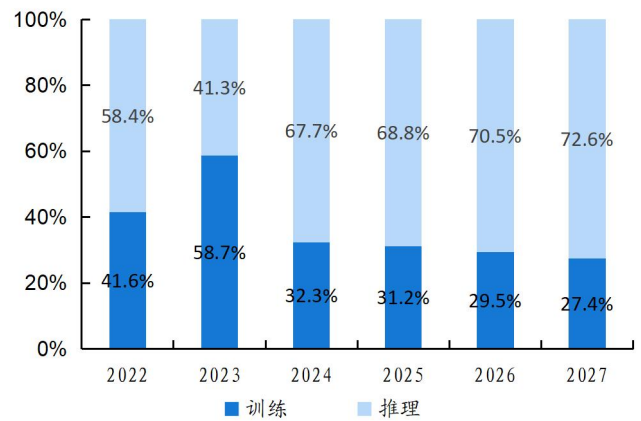
IDC 数据显示, 在中国, 随着训练模型的完善与成熟, 模型和应用产品逐步进入投产模式, 处理推理工作负载的人工智能服务器占比将随之攀升。IDC 预计, 到 2027 年, 用于推理的工作负载将达到 72.6%。

图 17: 大模型应用场景



资料来源: 从大模型到人工智能—机遇与挑战专题论坛, 国海证券研究所

图 18: 中国人工智能服务器 workload 预测

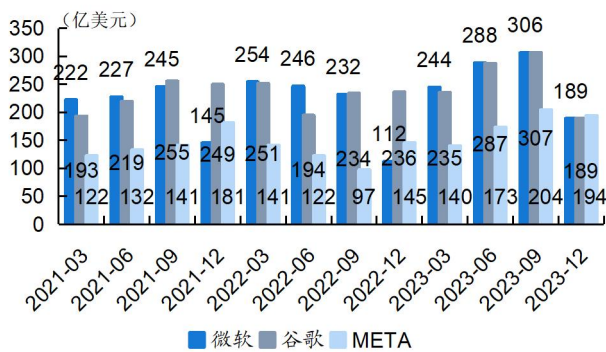


资料来源: IDC, 国海证券研究所

3.3、大厂用户+资金实力深厚, AI 资本支出提升空间大

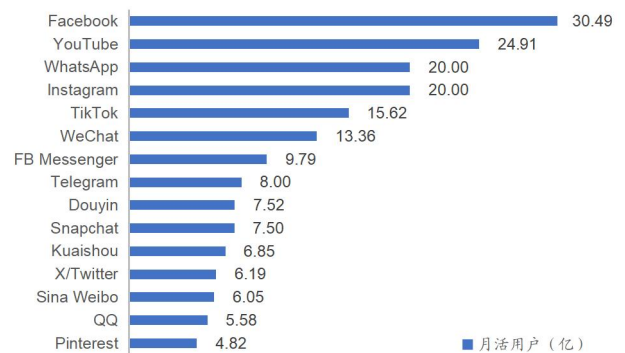
互联网大厂与云厂商具备充足用户规模与现金流量, 具备实施 AI 投资的基础。2023 年 10-12 月, 各大云厂商具有较为充足的经营性现金流量, 大约为 188-194 亿美元, 可以很好地支持生成式人工智能所需的资金投入。根据 statista, 2024 年 1 月, FACEBOOK 等社交媒体月活量较高, 具备较好用户基础, 也促使各大互联网厂商发展大模型, 满足用户需求。

图 19: 全球部分大型云厂商经营性现金流量情况



资料来源: iFind, 国海证券研究所

图 20: 2024 年 1 月最受欢迎社交网络

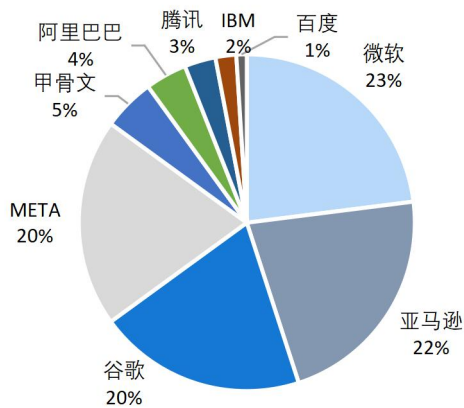


资料来源: statista, 国海证券研究所

大多云服务商的 AI 开支占总资本支出的比例在 3-7% 区间内, 还有较大提升空间。根据 Counterpoint 统计, 由于人工智能和网络设备的投资, 预计 2023 年全球云服务提供商的资本支出中, 35% 专门用于包括服务器和网络设备在内的 IT

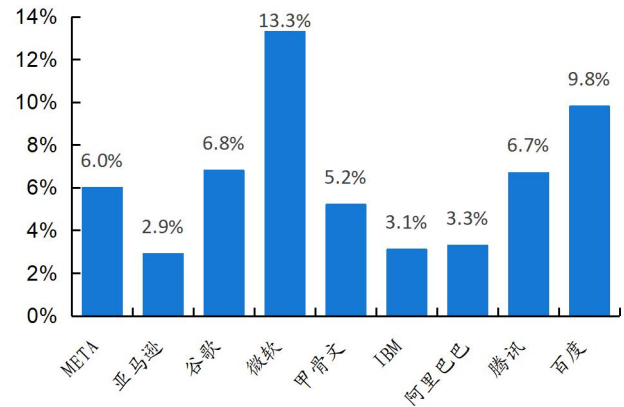
基础设施。微软和亚马逊的数据中心开发投资总额排名最靠前，微软将把超过13%的资本支出用于人工智能基础设施。

图 21：2023E 全球云服务厂商资本开支市场份额



资料来源：Counterpoint，国海证券研究所

图 22：2023E 云服务商 AI 开支占总资本开支比例



资料来源：Counterpoint，国海证券研究所

4、相关公司

行业评级及投资策略：AI 算力有望持续增长，GPU、AI 服务器、服务器散热、光模块、数据中心等环节有望持续受益。维持对计算机行业“推荐”评级。

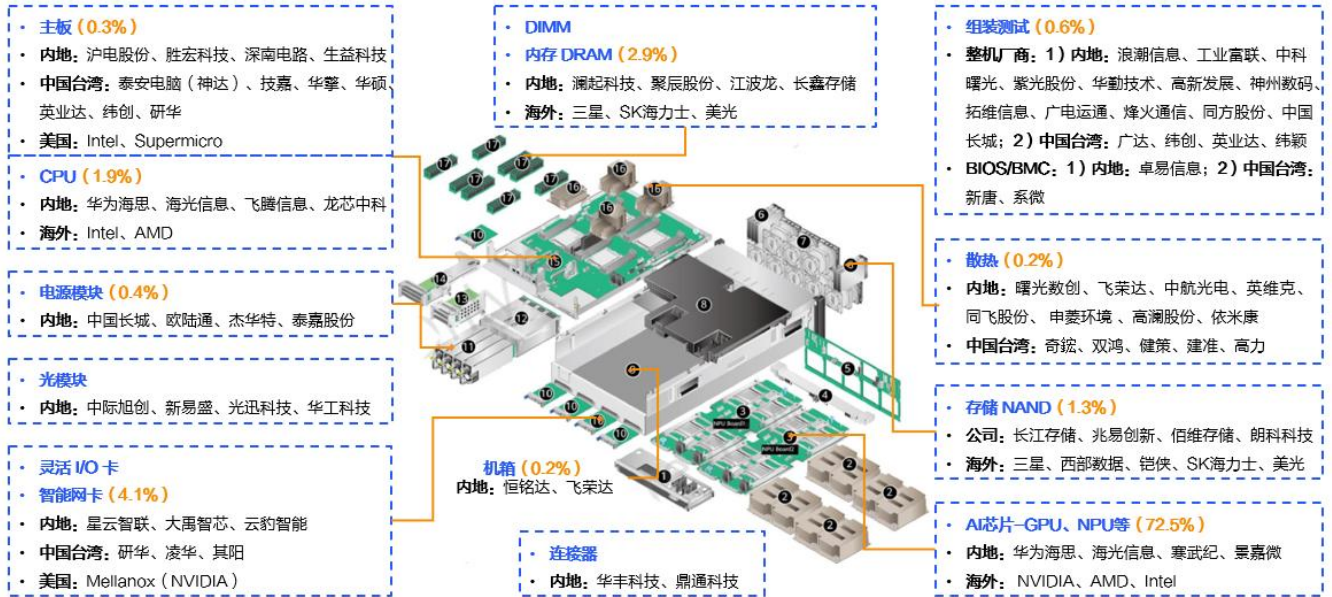
1) 服务器整机：工业富联、浪潮信息、中科曙光、紫光股份、华勤技术、神州数码、烽火通信、中国长城、高新发展、纬创、广达、英业达。

2) 服务器组件：
①AI 处理器：海光信息、寒武纪、龙芯中科、英伟达、AMD；
②散热：飞荣达、中航光电、曙光数创、英维克、奇鋆科技、双鸿、建准电机；
③主板：沪电股份、深南电路、胜宏科技、技嘉、华擎。

3) 光模块：中际旭创、新易盛、光迅科技、华工科技。

4) 数据中心：奥飞数据、光环新网、宝信软件、数据港、电科数字。

图 23: 服务器产业链相关公司



注: 蓝字为零部件, 括号中橙色数字为价值量。

上图展示以华为 Atlas 800 训练服务器为例(鲲鹏 920 * 4, 昇腾 910 * 8), 其中CPU集成在主板上, NPU集成在NPU板上; 具体价值量数字对标Nvidia DGX H100服务器, 具体计算方法见下文

资料来源: 华为昇腾官网, Semianalysis, 各公司公告, 各公司官网, 百度百科, 爱采购, 国海证券研究所

5、风险提示

- 1) 宏观经济影响下游需求:** 宏观经济环境下行, 将影响客户对信息化基础设施的采购需求;
- 2) 信创政策不及预期:** 行业主要驱动因素之一是信创政策持续落地, 若信创产业推进不及预期, 或导致行业内公司订单增速下行;
- 3) 市场竞争加剧:** IT 产品和服务行业是成熟且完全竞争的行业, 新进入者可能加剧整个行业的竞争态势;
- 4) 中美博弈加剧:** 国际形势持续不明朗, 美国不断通过“实体清单”等方式对中国企业实施打压, 若中美紧张形势进一步升级, 将可能导致中国半导体供应链供应受到影响;
- 5) 相关公司业绩不及预期:** 市场环境变化、公司治理情况变化、其他非主营业务经营不及预期等原因或将导致相关公司的整体业绩不及预期。
- 6) 各公司并不具备完全可比性, 对标的相关资料和数据仅供参考。**

【计算机小组介绍】

刘熹，计算机行业首席分析师，上海交通大学硕士，多年计算机行业研究经验，致力于做前瞻性深度研究，挖掘投资机会。新浪金麒麟新锐分析师、Wind 金牌分析师团队核心成员。

【分析师承诺】

刘熹，本报告中的分析师均具有中国证券业协会授予的证券投资咨询执业资格并注册为证券分析师，以勤勉的职业态度，独立，客观的出具本报告。本报告清晰准确的反映了分析师本人的研究观点。分析师本人不曾因，不因，也将不会因本报告中的具体推荐意见或观点而直接或间接收取到任何形式的补偿。

【国海证券投资评级标准】

行业投资评级

推荐：行业基本面向好，行业指数领先沪深 300 指数；
中性：行业基本面稳定，行业指数跟随沪深 300 指数；
回避：行业基本面向淡，行业指数落后沪深 300 指数。

股票投资评级

买入：相对沪深 300 指数涨幅 20%以上；
增持：相对沪深 300 指数涨幅介于 10%~20%之间；
中性：相对沪深 300 指数涨幅介于-10%~10%之间；
卖出：相对沪深 300 指数跌幅 10%以上。

【免责声明】

本报告的风险等级定级为 R3，仅供符合国海证券股份有限公司（简称“本公司”）投资者适当性管理要求的客户（简称“客户”）使用。本公司不会因接收人收到本报告而视其为客户。客户及/或投资者应当认识到有关本报告的短信提示、电话推荐等只是研究观点的简要沟通，需以本公司的完整报告为准，本公司接受客户的后续问询。

本公司具有中国证监会许可的证券投资咨询业务资格。本报告中的信息均来源于公开资料及合法获得的相关内部外部报告资料，本公司对这些信息的准确性及完整性不作任何保证，不保证其中的信息已做最新变更，也不保证相关的建议不会发生任何变更。本报告所载的资料、意见及推测仅反映本公司于发布本报告当日的判断，本报告所指的证券或投资标的的价格、价值及投资收入可能会波动。在不同时期，本公司可发出与本报告所载资料、意见及推测不一致的报告。报告中的内容和意见仅供参考，在任何情况下，本报告中所表达的意见并不构成对所述证券买卖的出价和征价。本公司及其本公司员工对使用本报告及其内容所引发的任何直接或间接损失概不负责。本公司或关联机构可能会持有报告中所提到的公司所发行的证券头寸并进行交易，还可能为这些公司提供或争取提供投资银行、财务顾问或者金融产品等服务。本公司在知晓范围内依法合规地履行披露义务。

【风险提示】

市场有风险，投资需谨慎。投资者不应将本报告视为作出投资决策的唯一参考因素，亦不应认为本报告可以取代自己的判断。在决定投资前，如有需要，投资者务必向本公司或其他专业人士咨询并谨慎决策。在任何情况下，本报告中的信息或所表述的意见均不构成对任何人的投资建议。投资者务必注意，其据此做出的任何投资决策与本公司、本公司员工或者关联机构无关。

若本公司以外的其他机构（以下简称“该机构”）发送本报告，则由该机构独自为此发送行为负责。通过此途径获得本报告的投资者应自行联系该机构以要求获悉更详细信息。本报告不构成本公司向该机构之客户提供的投资建议。

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。本公司、本公司员工或者关联机构亦不为该机构之客户因使用本报告或报告所载内容引起的任何损失承担任何责任。

【郑重声明】

本报告版权归国海证券所有。未经本公司的明确书面特别授权或协议约定，除法律规定的情况外，任何人不得对本报告的任何内容进行发布、复制、编辑、改编、转载、播放、展示或以其他方式非法使用本报告的部分或者全部内容，否则均构成对本公司版权的侵害，本公司有权依法追究其法律责任。