

## AI“大行其道” 英伟达“坐享其成”

本报记者 秦泉 北京报道

ChatGPT 在全球的关注度持续火热，再次将 AI 产业推到聚光灯下，科技巨头争相谋局落子，继微软、谷歌之后，国内企业百度、阿里巴巴等也先后发布大模型，并进行用户测试和企业应用接入。随着 AI 产业迎来“iPhone 时刻”，算力需求正在持续释放，以 AI 服务器核心零部件 GPU（图像处理、加速芯片）为代表的供给端走俏，其价格也在不断上涨，而在 AI 芯片 GPU 市场占据绝对优势的英伟达也赚得盆满钵满。

多位业内人士在接受《中国经营报》记者采访时表示，大型模型通常需要庞大的算力和存储资源来进行训练，GPU 已成为 AI 加速芯片通用性解决方案，越来越多的企业和个人开始使用 GPU 来训练大型深度学习模型。这种需求的增加可能导致 GPU 的价格上涨，从而导致显卡价格的上涨。此外，由于供应链问题、半导体短缺等因素的影响，显卡价格的波动也可能受到一定程度的干扰。

## 英伟达大秀肌肉

目前主流 AI 厂商都进入了“千亿参数时代”，多采用了英伟达的 GPU。

AI 场景需要多核、高并发、高带宽 AI 芯片。AI 芯片，也被称为 AI 加速器或计算卡，即专门用于处理人工智能应用中的大量计算任务的模块。当前，AI 芯片主要分为 GPU、FPGA，及以 TPU、VPU 为代表的 ASIC 芯片，而 GPU 凭借其高性能、高灵活度特点成为 AI 加速方案首选。据 IDC 数据，预计到 2025 年，GPU 仍将占据 AI 芯片 80% 市场份额。

资料显示，2018 年 OpenAI 开发的 GPT-1 的预训练大模型参数为 1.1 亿，2019 年发布的 GPT-2 提高至 15 亿，2020 年 GPT-3 的预训练大模型参数已经提高至 1750 亿。而为了训练 ChatGPT，OpenAI 构建了由近 3 万张英伟达 V100 显卡组成的庞大算力集群，GPT-4 更是达到了 100 万亿的参数规模，其对应的算力需求同比大幅增加。

TrendForce 分析认为，要处理

近 1800 亿参数的 GPT-3.5 大型模型，需要 2 万颗 GPU 芯片，而大模型商业化的 GPT 需要超过 3 万颗。GPT-4 则需要更多。

不仅如此，目前主流 AI 厂商都进入了“千亿参数时代”，多采用了英伟达的 GPU。以科大讯飞星火认知大模型为例，其使用了英伟达的 T4 Tensor Core GPU 进行训练。除了 T4 Tensor Core GPU，科大讯飞还使用了其他类型的 GPU 和其他硬件设备来支持其深度学习平台的开发和应用。这些硬件设备包括英伟达的 PaddlePaddle、NVIDIA Tesla V100、AMD EPYC 等，以及多种 CPU、内存、网络设备。

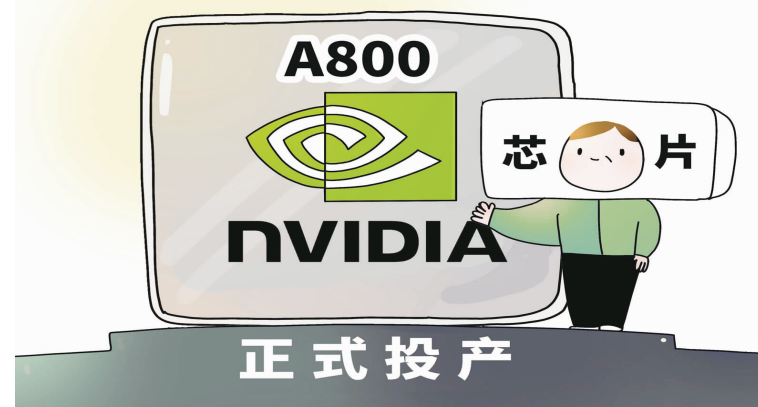
昆仑万维集团 CEO 方汉也表示：“超过千亿级别的大模型，它的训练大概需要 1000—2000 张 A100 的卡，没有 2000 张 A100 的卡，实验都做不了。”

招商证券指出，从通用服务器

到 AI 服务器，一个最显著的变化就是 GPU 取代了 CPU 成为整机最核心的运算单元以及价值量占比最大的部分，传统服务器通常至多配备 4 个 CPU+相应内存和硬盘，在 AI 服务器中，通常是 2 颗 CPU+8 颗 GPU，部分高配 4U 服务器中甚至可以搭配 16 颗 GPU，预计 AI 服务器中 GPU+CPU+存储的价值量占比有望达到 80% 以上的水平。

据统计，英伟达当前在售的用于大模型训练的 GPU 卡至少有 9 款型号，其中高性能的有 4 款，分别是 V100、A800、A100 及 H100。而此轮 AI“军备竞赛”也让用于上述显卡的价格一路高涨。其中，A100 此前售价在 1.5 万美元（约合人民币 10.35 万元），但目前在一些平台上，此款显卡价格上涨到 15 万元左右。

英伟达也借机赚足了“钱包”。TrendForce 数据显示，如果



英伟达也借这场“AI 军备竞赛”之机赚足了“钱包”。

视觉中国/图

以英伟达 A100 显卡的处理能力计算，GPT-3.5 大模型需要 2 万块 GPU 来处理训练数据。目前英伟达 A100 显卡的售价为 10000~15000 美元之间，预估英伟达可以赚 3 亿美元（约 20 多亿元人民币）。

值得注意的是，英伟达还在源源不断地为这场军备竞赛输送弹药。在此前 GTC 开发者大会上，英伟达推出了新的 Hopper

CPU——配有双 GPU NVLink 的计算，GPT-3.5 大模型需要 2 万块 GPU 来处理训练数据。目前英伟达 A100 显卡的售价为 10000~15000 美元之间，预估英伟达可以赚 3 亿美元（约 20 多亿元人民币）。

不过，即使价格上涨，目前市面上几乎“一卡难求”。一位业内人士对记者表示，客户对英伟达 A100/H100 芯片需求强劲，后者订单能见度已至 2024 年，更紧急向代工厂台积电追单。

## 国产厂商的机遇

目前中国 GPU 开发者大多使用国外厂家提供的 IP，自主性不高，不过经过多年沉淀是能够实现国产替代的。

虽然国内外的大模型项目接连落地，但除了百度、阿里巴巴等企业采用自研芯片外，国内大多数企业仍难求高端 GPU。据透露，国内可用于训练 AI 大模型的 A100 大约有 4 万—5 万个。

英伟达在去年收到美国政府的许可，通知称：“若对中国（含中国香港）和俄罗斯的客户出口两款高端 GPU 芯片——A100 和 H100，需要新的出口许可。”不仅如此，该许可证还要求还包括未来所有的英伟达高端集成电路，只要其峰值性能

能和芯片间 I/O 性能均大于或等于 A100 的阈值，以及包括这些高端电路的任何系统，都被纳入许可证限制范围。

不过，英伟达针对中国客户推出了替代型号 A800，与原有的 A100 系列计算卡相比，A800 系列的规格基本相同，比较大的区别在于 NVLink 互连总线的连接速率，前者为 600GB/s，后者限制在了 400GB/s。综合使用效率只有 A100 的 70% 左右。前不久英伟达还发布了特供版的

H800，作为其旗舰芯片 H100 的替代版。

4 月 14 日，腾讯云正式发布新一代 HCC（High-Performance Computing Cluster）高性能计算集群。据悉，该集群采用腾讯云星海自研服务器，搭载英伟达最新代次 H800 GPU，H800 基于 Hopper 架构，对跑深度学习推荐系统、大型 AI 语言模型、基因组学、复杂数字孪生等任务的效率提升非常明显。与 A800 相比，H800 的性能提升了 3 倍，在显存带宽上也

有明显的提高，达到 3TB/s。

伴随着近期宏观经济回暖以及国内互联网企业纷纷加大 AI 算力布局，PC 和服务器的需求上升有望为国内 GPU 市场带来整体拉动效应。

目前，国内已涌现出一批优秀的 GPU 设计和制造厂商。

其中，海光信息目前已经成功掌握高端协处理器微结构设计等核心技术，并以此为基础推出了性能优异的 DCU 产品。其深算一号产品和英伟达 A100 及 AMD 高端

GPU 产品（MI100）进行对比，单芯片产品基本能达到其 70% 的性能水平。

值得注意的是，上述业内人士表示，虽然国内的 GPU 厂商取得了一些成绩，但是由于 GPU 研发难度大、开发周期长，例如 A100，英伟达只用了三个月的时间便研发出替代方案，而国内却并不多见。而且目前中国 GPU 开发者大多使用国外厂家提供的 IP，自主性不高，不过经过多年沉淀是能够实现国产替代的。

## 三星大变局：以 AI 重构商业版图

本报记者 吴清 北京报道

在全球消费电子和半导体产业持续低迷的背景下，老牌巨头三星全线出击，正在酝酿一场全球业务的大变局。

5 月 15 日，据日媒报道，韩国三星电子计划将投资 300 亿日元（约合人民币 15.4 亿元）在日本横滨建设半导体后段封装测试产线，预计将在 2025 年完工量产。而就在前一日，据韩媒报道，三星电子与

Naver 计划联合开发一款用于企业的生成式 AI（人工智能）工具，对标 ChatGPT。三星和 Naver 还计划在下半年推出 AI 芯片。

而更值得关注的是，三星掌舵者、三星电子会长李在镕刚刚结束为期 22 天的访美行程。在此期间，他分别会见了特斯拉 CEO 马斯克、英伟达 CEO 黄仁勋、微软 CEO 纳德拉、谷歌 CEO 皮查伊等 20 多位科技公司的高管，遍及互联网、AI、通信、生物等行业。

《中国经营报》记者注意到，尽管此次访美行程紧凑，但李在镕花了很多时间与当地 AI、半导体产业的“领袖”互动交流，顺应近期 AIGC（AI 生成内容）的热潮，AI 和半导体芯片成为李在镕此行的重点。一位三星内部人士告诉记者：“在全球信息与通信技术不景气的情况下，三星需要培育增长型业务作为公司的新支柱。”

今年 4 月末，三星电子公布的业绩报告显示，公司第一季度营业

利润为 6402 亿韩元（约合人民币 33.2 亿元），同比下降 95.5%，为 14 年来的最低季度营业利润。上述内部人士补充道：“在这个关键时刻，李在镕会长决定亲自激活全球网络，探索新的业务战略，并取得了突破。”

而就在此前，李在镕和苹果 CEO 库克等一道到访中国并低调出席了一系列活动，去年以来三星手机也通过线下开店、线上广告等形式“重返”中国市场。

## 抢滩 AIGC 全线出击

作为近年来 AI 领域最重大的突破，近期全球掀起的 AIGC 和 ChatGPT 浪潮，三星也未错过。

5 月 14 日，也就是李在镕结束访美回国后的第三天，据韩媒报道，三星电子与韩国最大在线搜索引擎运营 Naver 计划联合开发一款用于企业的生成式 AI 工具，对标 ChatGPT。

两家公司计划最早于今年 10 月发布该工具，并将首先用于三星的设备解决方案（DS）部门，该部门包括半导体业务。三星还将在测试后，将其应用范围扩大到三星的其他业务，包括负责智能手机和家电业务的设备体验（DX）部门。

“安全”或是三星下场自研 AIGC 工具的重要考量因素之一。上述韩媒报道指出，通过与 Naver 合作，三星可在提高公司生产力的同时规避使用其他平台带来的商业机密泄露风险。此前出于安全考虑，三星电子已禁止员工使用 ChatGPT、谷歌 Bard 和微软必应等流行的生成式 AI 工具。

同时，三星和 Naver 还计划在今年下半年推出 AI 芯片。两家公司已为 AI 芯片业务筹备数月。自三星和 Naver 于 2022 年 12 月正式宣布合作开发 AI 芯片，短短 5 个月时间内，已完成了相关 AI 芯片的技术验证，计划最快 8 月做成可编程芯片并正式测试。

在强化上游 AI 芯片等领域合作的同时，三星在下游消费电子领域的动作也没落下。从去年开始，中国的消费者在电视和线下门店中看到了更多三星手机的身影。

市场分析机构 Canalys 公布的 2023 年第一季度全球智能手机厂商排名显示，三星以 6030 万部出货量重回市场第一，苹果公司以 5800 万部的出货量位居第二；小米、OPPO 和 vivo 分别以 3050 万部、2660 万部和 2090 万部位列第三、四、五名。2023 年第一季度全球智能手机出货量同比下降 13%，降至 2.7 亿部。

正是在全球手机市场下滑背景下，以及面对来势汹汹的苹果及中国手机品牌，“重返”中国这个全球最大消费电子市场成为三星虽然艰难但必须完成的动作。三星手机曾一度占据中国约 30% 的手机市场份额，但此前受各种因素影响已降至 1% 左右，与其全球第一的地位形成鲜明对比。

据悉，2021 年 12 月，三星电子副会长兼 CEO 韩钟熙挂帅的“中国业务创新团队”开始实施线下流通渠道扩张战略。此前，一位三星电子内部人士也向记者证实，折叠屏手机是三星布局中国市场的重点。

此外，三星近期组织业务调整频频。分析人士认为，李在镕可能调整三星电子的组织架构，使集团文化更为灵活，以适应最新的产业潮流所带来的快速变化。

“三星拥有雄厚的资金、技术、人才积淀，全产业链布局优势以及良好的上下游合作伙伴关系，业绩压力之下，三星多路出击、全线布局。”上述半导体行业人士对记者表示，不能低估一个少壮派领导重振三星的意志和三星作为老牌跨国巨头复苏的潜力。

## 布局半导体封测 应对台积电竞争

据前述日媒报道，此次三星电子投资在日本横滨建立半导体后段封装测试产线，主要是因为当地有三星的研发机构，这样可以相互配合发展，同时将会获得日本政府“芯片法案”的配套补贴。

从日本“芯片法案”条款来看，三星此次投资将可获得约 100 亿日元（约合人民币 5.12 亿元）的税收减

免和资金补助，这是日本政府为在当地投资半导体芯片制造领域企业提供的激励政策。

目前尚不清楚该项目的具体投资计划，不过根据已公布的信息，该项目接下来会招聘数百个工作岗位以进行产线上的工作。

过去多年来，三星一直专注于半导体前端的先进制程技术发展。不过在当前摩尔定律“失效”，先进

制程推进越发困难的背景下，三星近年来也开始积极发展先进封装技术，以更先进的封装技术来弥补制程技术提升的放缓。

有意思的是，三星在半导体领域的长期竞争对手台积电早在 2021 年宣布在日本熊本与 SONY、DENSO 合作建立半导体晶圆厂。因此有业内人士认为，三星在日本的投资，其实也是针对竞争对手台

积电而来。

不过，三星此举可能还有更为现实的考量。“日本一直是半导体设备及材料强国，前两年因为半导体材料被‘卡脖子’曾让三星陷入被动，此次，三星希望借投资深化与上游日本材料和设备供应商的合作，寻求半导体制造更大的突破。”一位半导体行业人士对记者表示。

## 22 天旋风访美强化 AI、芯片领袖合作

而就在此前，李在镕刚刚结束为期 22 天的访美之行，在这并不长的时间内，李在镕“旋风”般地与 20 多位跨国企业的高管进行了会面。

据悉，李在镕于今年 4 月底前往美国，5 月 12 日返回韩国，在此期间，李在镕从美国东海岸一路飞到西海岸，是其自 2014 年以来最久的访美行程。在此行中，李在镕与谷歌、微软、英伟达、特斯拉、辉达、百健、欧嘉隆等公司高层会晤，而这些公司都是全球 AI、车用芯片、通讯芯片及生物领域的领导者。

比如，5 月 10 日李在镕和黄仁勋在硅谷一家餐厅一对一会面，讨论了未来在 AI 芯片技术和晶圆代工合作的计划；辉达则是全球首屈

一指的生成式 AI 绘图处理器（GPU）设计业者，这场会面提高了辉达下单三星高频宽半导体芯片（HBM）的市场预期。

备受各方关注的是，李在镕在硅谷的三星电子半导体研究中心，与马斯克的会晤，这也是两人第一次私人会谈。消息人士说，三星和特斯拉就共同开发自动驾驶汽车芯片、新一代信息通信技术进行了交流，这场会面将为双方扩大车用芯片合作铺路。

车用芯片被认为是半导体行业未来发展的重中之重，前景广阔。市场研究机构 Strategic Analytics 和 Research & Markets 的数据显示，全球汽车芯片市场预计将在

明年升至 4000 亿美元（约合人民币 2.79 万亿元），到 2028 年，这个数字将进一步增至 7000 亿美元（约合人民币 4.87 万亿元）。

今年 4 月，三星刚获得了知名自动驾驶科技公司 Mobileye 的高性能芯片订单。业内认为，三星基于特斯拉全新一代自动驾驶芯片生产经验，进一步增强其在汽车芯片市场的影响力。

作为全球最大的 DRAM 和 NAND Flash 生产商之一，也是全球最大的半导体代工厂商之一，半导体业务一直是三星的营收和利润的主力军。不过受累于低迷的行业形势，今年一季度该业务却亏损 4.58 万亿韩元（约合人民币 236.79

亿元）。三星方面还预计，由于需求大幅下滑，今年全球芯片市场规模将萎缩 6% 至 5630 亿美元（约合人民币 3.92 万亿元），困难状况可能持续全年。“在此背景下，三星急需稳固原有合作关系的同时，激活更多市场需求。”上述半导体行业人士表示。

此次行程中，李在镕也和纳德拉、皮查伊等企业领袖举行一连串会议，并积极与全球知名的 AI 学者建立人脉，讨论 AI 的多元应用和未来的合作，显示出三星想重新连接受新冠疫情影响的国际合作关系，为 AI 技术未来成长引擎的蓝图奠定基础。三星目前已视 AI 为未来公司成长的关键引擎之一。