

# 从基础设施到具身智能 中国科技力量闪耀GTC 2026

◎记者 孙小程

素有“AI 界春晚”之称的英伟达 GTC 2026 已拉开帷幕。

与往年相比，中国科技领域代表企业参与的广度和深度明显提升，在主题演讲、技术展示和生态合作等多个维度中展现创新实力。

联想集团、吉利汽车、宇树科技……此次与会的中国企业，业务范围涵盖 AI 基础设施硬件、垂直行业应用及前沿具身智能等多个领域。在与英伟达的合作中，它们实现了从应用到生态的系统性参与。可以预见，中国科技企业正将产业动能与全球创新技术紧密链接，共同加速智能技术的普惠与落地。

## 助力 AI 基础设施建设

“今年将会是属于你们的一年，我感受到了。”在 GTC 2026 展会上，英伟达首席执行官黄仁勋对联想集团董事长兼 CEO 杨元庆说。杨元庆回应称：“联想集团的业务板块非常强劲，甚至比以前更强。”

联想集团与英伟达合作多年，是每年 GTC 的核心参与者之一。在 GTC 2026 上，联想集团与英伟达联合发布新一代联想 Hybrid AI Advantage（混合式 AI 优势集）解决方案，旨在加速 AI 落地，缩短首 token 时间（TTFT），并在个人、企业和云环境中带来可量化的商业成果。

该解决方案是联想集团与英伟达三十年合作的全新成果，也是混合式 AI 技术走向产业化的重要标志。杨元庆表示，联想集团和英伟达共同具备独特优势，能够帮助各类机构将 AI 投入运营——从实验阶段到企业生产，再到 AI 云级巨型工厂。

除了联想集团，还有多家同样专注于 AI 基础设施的中国企业展示了最新布局。例如，在液冷领域，领益智造旗下子公司立敏达（Readore）进入新一代 Rubin 架构 Manifold（分水器）生态，展示了包括 UOD/MOD 快接头、Inner Manifold 在内的核心液冷产品。

据领益智造介绍，立敏达是一家以热管理产品为核心的服务器综合硬件方案供应商。立敏达核心客户覆盖海外算力行业头部客户及其供应链相关合作伙伴、相关服务器代工厂、电源解决方案公司，并且已建立深厚且长久的客户关系。

## 赋能汽车产业转型

以人工智能为核心的智能化体验，是近年来汽车产业的角力点。在 GTC 2026 上，中国汽车产业链的参与凸显了这一趋势。

吉利汽车集团发布了超级 Eva——首个打通智能座舱、智能辅助驾驶、数字生态的超级智能体，并宣布极氪 8X 将首发搭载。超级 Eva 由吉利、千里科技和阶跃星辰共同研发，其核心驱



3月16日，英伟达公司首席执行官黄仁勋在英伟达GTC大会上发表主题演讲

动大脑来自阶跃星辰自研的 Step 3.5 Flash 基座模型。

同时，吉利还宣布将携手英伟达在智能辅助驾驶、智能座舱、智能制造与研发、云端与 AI 基础设施等领域，持续深化合作与战略协同。除了整车企业，禾赛科技与速腾聚创等激光雷达企业也参与了 GTC 2026，重点呈现与英伟达的技术及平台合作。

禾赛科技宣布，公司已加入 NVIDIA Holos AI 系统检测实验室，这是首个获得美国国家标准学会国家认可委员会认证的全球项目。作为成员单位，禾赛将在这一统一框架下对其激光雷达平台开展功能安全、网络安全以及 AI 合规方面的评估与验证。

禾赛科技联合创始人及 CEO 李一帆表示：“只有在安全性与可靠性达到最高标准的前提下，自主系统才能实现广泛应用。随着自动驾驶汽车和智能机器人逐步迈向规模化部署阶段，安全将成为整个系统最核心的要求。”加入 NVIDIA Holos 生态，将进一步提升自动驾驶系统的安全能力，并推动自动驾驶技术实现规模化应用。

速腾聚创介绍，作为全面加入 NVIDIA Jetson、DRIVE、Omniverse 三大生态系统的合作伙伴，公司凭借数字化激光雷达产品，与英伟达生态形成深度协同，为全球汽车与机器人头部企业共同打造更高性能的感知解决方案。

目前，搭载速腾聚创“主+补盲”数字化激光雷达产品组合与 NVIDIA DRIVE AGX Thor 芯片的系统方案，已成为 L4 级产业伙伴打造下一代自动驾驶系统的优选技术组合。在 GTC 2026 现场，文远知行、小马智行等均展示了搭载了这一自动驾驶系统方案的 Robotaxi

车型，展示无人出行服务的全球规模化、商业化运营格局。

## 共推具身智能落地

黄仁勋从不掩饰对具身智能的看好，他曾言：“具身智能是 AI 的下一浪潮”。围绕具身智能，英伟达构建了从底层算力、核心模型、仿真训练的全栈技术生态。

GTC 2026 上亦出现了多家头部具身智能企业的身影，包括宇树科技、智元机器人、自变量等。

大会披露的英伟达 AI 生态核心伙伴中，“AI for Robotics”（即具身智能机器人板块）共有 Physical Intelligence、Figure、Skild AI、智元机器人、自变量机器人等 9 家具身智能企业。黄仁勋宣布，将与这些企业开展物理 AI 的大规模部署，推动智能机器人在工厂、物流、交通、医疗和基础设施等领域应用，加速机器人在真实场景的落地。

在合作方向上，智元机器人介绍，其与英伟达围绕机载超算大脑开展联合硬件适配，支撑百亿参数 VLA 模型与世界模型在机器人本体本地高效推理，满足电子制造、汽车总装、物流作业等场景下 0.1 毫米级精密操作与长程自主规划需求。作为英伟达 Isaac GR00T N 系列全球首批生态伙伴，智元机器人将 GR00T 模型体系与自研机器人操作系统深度融合，形成可复制的通用技能迁移与灵巧操作部署流程。

宇树科技创始人兼 CEO 王兴兴发表了主题演讲。他表示，具身智能的 ChatGPT 时刻，意味着机器人可通过语音或文字指令，在 80% 的陌生场景中顺利完成约 80% 的任务，而实现这一目标需要更多全球合作。

纤和 CPO 方面持续扩产。”黄仁勋在介绍其技术路线图时说。

黄仁勋还“剧透”了下一代计算架构 Feynman，它将首次实现铜线与 CPO 的共同部署。同时，英伟达正在研发部署在太空的数据中心级 AI 计算能力部署到卫星和轨道数据中心（ODC），并强调其面向在轨推理、实时地理空间智能和自主航天任务。

## OpenClaw 是新时代“操作系统”

谈及最近风靡 AI 圈的“龙虾”，黄仁勋将开源项目 OpenClaw 形容为“人类历史上最受喜爱的开源项目”，称其仅用几周时间就超越了 Linux 在过去 30 年取得的成就。他表示，OpenClaw 本质上就是 Agent 计算机的“操作系统”。

“每一个 SaaS（软件即服务）公司都将变成 AaaS（智能体即服务）公司。”黄仁勋介绍，为了让这种具备访问敏感数据和执行代码能力的智能体安全落地，英伟达推出了企业级的 NemoClaw 参考设计，增加了策略引擎和隐私路由。

物理 AI 则是具身化的智能体。据介绍，本次 GTC 共有 110 款机器人亮相，几乎囊括全球所有机器人研发企业。英伟达提供训练计算机、仿真计算机、机载计算机和完整的软件栈及 AI 模型。

自动驾驶方面，黄仁勋认为，自动驾驶的“ChatGPT 时刻”已经到来。“今天，我们宣布四家新合作伙伴加入英伟达 RoboTaxiReady 平台：比亚迪、现代、日产、吉利，合计年产量达 1800 万辆。加上此前的奔驰、丰田、通用，阵容进一步壮大。我们同时宣布与 Uber 达成重大合作，将在多个城市部署并接入 RoboTaxi Ready 车辆。”他说。

工业机器人方面，ABB、Universal-Robotics、KUKA 等众多机器人企业与英伟达合作，将物理 AI 模型与仿真系统相结合，推动机器人在全球制造产线的落地。

## 算力大变局

（上接1版）

在 GTC 大会上，黄仁勋抛出一个极具冲击力的预测：到 2027 年，市场对英伟达 Blackwell 和 Vera Rubin 系统的订单需求将带来至少 1 万亿美元的营收。这一数字较去年他对 2026 年 5000 亿美元市场需求的预测直接翻倍。他甚至强调，1 万亿美元只是保守预测，实际全球 AI 算力需求会比这一数字高得多。

其次，“上天”（太空算力）和“入地”（端侧算力）成为产业演进主线。2026 年，人类一边仰望星空，突破地球物理边界；一边深耕大地，让智能渗透进每一个角落。

再次，中美走向两条侧重不同的发展道路。美国侧重于前沿突破，驱动力来自资本与顶尖研发的良性循环。目前，英伟达 Vera Rubin 平台七颗芯片全部投产，从 3nm 走向 1.6nm，从地面延伸至太空，从 GPU 扩展至 LPU，黄仁勋的每一步都在探索“下一代算力长什么样”。

中国则更重视向深处扎根，驱动力来自国家战略和产业升级的现实需求。“我们的应用空间更广阔，工业场景和商业链路更完整，所以我们在优先解决算力如何真正‘用起来’。”方海声说。

最后，全球算力加速形成“一个世界，两套系统”的竞争格局。中国主导的生态更开放，有望覆盖新兴市场；美国主导的生态仍牢牢占据高端。

## 两大隐忧

持续飙升的资本开支，是全球算力产业的“总开关”。

2026 年，全球 AI 总支出预计达 252 万亿美元，同比增长 44%。最近两年，中国各地已掀起算力中心扩建的浪潮。在 A 股市场，一些上市公司动辄斥资数十亿甚至上百亿元采购算力，以进军算力中心业务。

如火如荼的算力扩建背后，也有人焦虑与担心，狂热的投资会不会一地鸡毛？

近期，记者调研南方某国家级智算中心时，看到“冰火两重天”的景象。机房相关负责人指着左边几乎满负荷的机架说，这是装了英伟达 GPU 的服务器，出租率 90% 以上。他又指向右边说，这是装了国产 GPU 的服务器，价格便宜很多，但出租率不到 50%。

记者询问缘由，“能效差距很大，尤其在生态上。”对方回答。

飞腾信息技术有限公司副总经理郭御风也有类似感触。他调研发现，不少智算中心算力利用率不足 30%，大量算力资源长期闲置。他认为，行业“重算力、轻应用”“重建设、轻实效”的结构性问题突出。

“当前算力结构供需错配。低端算力过剩而高端智能算力不足；西部通用算力利用率偏低，东部产业急需的智能算力供应紧张。同时，算力孤岛现象严重，跨区域、跨主体算力资源难以高效流通。”北京国际城市发展研究院副院长连玉明说。

还有业内人士透露：不少地方数据中心配的芯片仍以 CPU 为主，适配传统 IT 与云服务场景，难以满足 AI 大模型训练与推理的需求；有些地方芯片搭配的计算架构不合理，导致应用场景区，即便电力成本低，也无法使用。

这些现象背后，是中国本轮数据中心建设热潮，既缘于 AI 真实需求，也与地方政府和资金方追逐热点、盲目上马项目相关。

## 上天入海

尽管科技巨头对 AI 的未来相当乐观，但要花一笔钱“变成”一个个智算中心，正变得越来越困难。能源供给、散热能力与耗水量正逐渐成为地面算力增长的瓶颈。

于是，人类将目光投向星空。2 月，马斯克旗下 SpaceX 完成了对 AI 公司 xAI 的全资并购，标志着航天与 AI 进入深度融合阶段。据马斯克预测，未来 2 至 3 年太空将成为全球 AI 算力成本最低的区域。目前，SpaceX、亚马逊、谷歌、英伟达等巨头纷纷布局太空算力。

同一时间，在浙江杭州之江实验室的“三体计算星”指挥中心，中国科学院院士、之江实验室主任王坚透露，该星已实现空间组网突破，在轨协同完成 10 个人工智能模型与应用的部署验证。

“今年‘三体计算星’还将发射 50 颗卫星，计划 2032 年完成 1000 颗计算卫星组网，形成能互联互通为人工智能服务的太空算力星屋。届时，总算力将达每秒百亿亿次。”在王坚看来，把算力送到太空的价值堪比电的发明，将催生诸多想象不到的新价值。

此外，中国还在悄悄试水算力“下海”。2 月 10 日，东海之滨，全球首个“海风直联”海底数据中心在上海临港投产。它将海上风电与海底数据中心直接联通，绿电占比高达 95%，利用海水自然冷却。

除了算力的“上天下海”，端侧算力（可称之为“入地”）的发展更是一场触手可及的大变革，2026 年将成为驱动消费电子和汽车产业升级的核心引擎。从豆包手机形态到 Openclaw 带火的 Mac Mini，

标志性案例层出不穷。汽车具备高算力芯片、人机交互界面和充足电源，成为端侧 AI 硬件落地的理想场景。

“AI 历经多轮演进，已迈入以推理为核心的全新阶段。”3 月 17 日，在 2026 华为数据存储新春发布会上，华为存储产品线副总裁、闪存领域总裁谢黎明表示。当日，华为正式发布针对 AI 推理场景的全新 AI 数据基础设施：面向中心推理场景的 FusionCube A1000 AI 超融合一体机。

## 国产提速

算力比拼，首当其冲的就是算力芯片。岁末年初，国产芯片的“小龙”们——摩尔线程、沐曦股份、壁仞科技、天数智芯先后登陆资本市场，这既是资本市场助力科创企业的最新实践，也是中国算力芯片行业进入资本化高歌时刻。

“今年芯片供给将从满足‘有没有’，向提供‘好不好、准不准’的差异化、场景化解决方案演进。”天数智芯相关负责人说。

“2026 年，芯片会更聚焦易用性、安全性和高能效比。”瀚博半导体创始人兼 CTO 张磊对记者说，“随着下游 AI 推理、工业质检、数字孪生等需求释放，国产算力的竞争力正在从硬件参数向全栈解决方案能力拓展。”

国产算力芯片虽然发展迅猛，但要真正做到从“跟跑”到“并跑”，仍任重道远。

在设计与制造环节，先进 EDA 工具匮乏、高端工艺产能不足，仍是主要瓶颈。多家国产算力芯片公司反映，国产先进工艺产能“抢手”，即便设计出了先进芯片，大规模量产时却难上加难。“芯片设计行业都是给晶圆厂打工的。”一位半导体行业资深分析师表示，“在‘天花乱坠’的芯片参数外，我们更关心企业能否流片及量产。”

更大的差距在软件生态。一位不愿具名的芯片业内人士坦言：“真正的壁垒不是把芯片做出来，而是让全球开发者愿意用你的软件平台。英伟达之所以能统治市场，不光因为芯片强，更因为它花了 20 年建立起来的 CUDA 软件生态。”

“这非常重要。”快思慢想研究院院长田丰说，“英伟达因为有 CUDA 软件层才使得它的显卡在流体力学、纳米科研算法等细分领域的算法加速表现优异。”

上述芯片业内人士还透露：目前国内大部分的 AI 大模型开发和训练，依然基于英伟达的 CUDA 生态；国产算力芯片的软件适配、算力优化、开发工具链完善度，与英伟达存在代际差距。

中国工程院院士邓中翰近日表示，国产高端算力芯片规模化应用面临三大核心挑战：一是技术适配性不足，现有芯片多对标传统架构，与 AI 大模型、智算集群的多元计算需求不匹配，存在“算力适配难、场景落地贵”的现实痛点；二是生态体系有短板，软硬件协同、标准统一、场景验证未形成闭环，无法满足规模化应用的稳定性要求；三是算力利用率低，行业“暴力计算”模式大幅拉高能耗成本，也让国产芯片在工艺制程受限的情况下，难以发挥架构创新的优势。

## 算电协同

在太平洋的另一侧，美国算力产业也有自己的“烦恼”。

日前，一则消息从华盛顿传出：美国政府官员正要求微软、Alphabet 等科技公司作出承诺，确保其数据中心不会推高电价，不会给消费者带来其他负担。

这一行动旨在应对全美范围内数据中心扩张引发的政治与公共安全问题。在部分地区，科技公司已遭遇来自民众的日益高涨的抵制声浪。一些维权团体正在积极行动，反对建设高耗能数据中心，理由是这些数据中心会挤占当地基础设施、可用水资源和电力供应。包括亚特兰大和新奥尔良在内的多个地区已经对新建数据中心实施限制措施。

美国的现状揭示了一个全球性难题：算力繁荣的社会成本，正在从隐性走向显性。

如何化解上述问题？“算电协同”概念应运而生。这是 2026 年中国新基建战略中的核心概念，2026 年首次被写入政府工作报告。

运营商认为，在 AI 进入“拼电力”时段，中国的“算电协同”模式更有可能在大规模、可持续的绿电算力供给上胜出，因为它解决了能源与算力在地理和时间上的根本性错配。而美国若不解决电网碎片化和扩容缓慢的问题，其 AI 发展的能源瓶颈可能会日益凸显。对于全球而言，中国的实践提供了一种将数字经济与能源转型深度融合的“中国方案”。

人工智能的浪潮已经势不可挡。作为 AI 核心基础设施，智能算力的研发、迭代、扩张，算力与电力的协同发展，同样势不可挡。在算力之争就是国力之争的大背景下，算力大变局的戏码，或许才刚刚开始。

## 黄仁勋：数据中心将成为Token“工厂”

◎记者 赛世平 郑维汉

当地时间 3 月 16 日，英伟达 GTC 2026 大会正式开幕，英伟达创始人兼 CEO 黄仁勋发表了主题演讲。

针对市场高度关注的订单与营收“天花板”，黄仁勋在大会上给出了极为乐观的业绩预期。“到 2027 年，至少有 1 万亿美元的需求。事实上，我们甚至会供不应求，实际的计算需求会比这高得多。”黄仁勋说。

同时，黄仁勋还介绍了 Rubin、Groq、Feynman 等最新数据和进展，并认为 Open-Claw（“龙虾”）是新时代的操作系统，自动驾驶的“ChatGPT 时刻”已经到来。

## 数据中心将成为Token“工厂”

“未来的数据中心不再是存储文件的仓库，而是生产 Token（AI 生成的基本单位）的‘工厂’。”黄仁勋说。

他解释称，每一座数据中心、每一座工厂，从定义上来说都是受电力限制的。一座 1GW（吉瓦）的工厂永远不会变成 2GW，这是物理和原子的定律。在固定的功率下，谁的每瓦 Token 吞吐量最高，谁的生产成本就最低。

Token 经济学正在塑造未来的职场新形态。硅谷最新的招聘筹码是：“你的 offer 里带多少 Token？”黄仁勋预计，在未来，公司的每一位工程师都需要一个年度 Token 预算，他们的基础年薪可能是几十万美元，英伟达会在此基础上再拿出大约一半的金额作为 Token 额度给他们，让他们实现 10x 的效率提升。

随着模型越来越大、上下文越来越长，AI 会变得更聪明，但 Token 的生成速率会降低。黄仁勋表示，在这个 Token 工厂里，吞吐量和 Token 生成速度将直接转化为明年的精确收入。英伟达的架构能够让客户在免费层实现极高的吞吐量，同时在最高价值的推理层级上，将性能提升惊人的 35 倍。

目前，英伟达 60% 的业务来自排名前五的

超大型云服务商，另 40% 的业务则广泛分布于主权云、企业、工业、机器人和边缘计算各个领域。

## LPU 预计三季度发货

“过去提到 Hopper，我会举起一块芯片，那很可爱。但提到 Vera Rubin，大家想到的是整个系统。在这个 100% 液冷、完全消灭了传统线缆的系统中，过去需要两天安装的机架，现在只需两小时。”提及 Vera Rubin 系统时，黄仁勋表示。

据黄仁勋介绍，通过极致的端到端软硬件协同设计，Vera Rubin 在同一座 1GW 数据中心里创造了惊人的数据跨越：在短短两年时间内，我们将 Token 的生成速率从 2200 万提升到 7 亿，实现了 350 倍的增长。

为了解决极速推理（如 1000Tokens/秒）条件下的带宽瓶颈，英伟达推出了单柜容纳 256 个 LPU 的 Groq 3 LPX 机柜，计划与今年晚些时候交付客户的 Vera Rubin 机架系统并排放置。黄仁勋表示，Groq LPX 机柜可使旗下 Rubin GPU 的每瓦能耗性能提升 35 倍。

黄仁勋说，英伟达通过 Dynamo 软件系统，将需要海量计算和显存的“预填充（Pre-fill）”阶段交给 Vera Rubin，将对延迟极度敏感的“解码（Decode）”阶段交给 Groq。黄仁勋还对企业算力配置给出了建议：“如果你的工作主要是高吞吐，100% 使用 Vera Rubin；如果你有大量高价值的编程级别的 Token 生成需求，拿出 25% 的数据中心规模给 Groq。”

据透露，由三星代工的 Groq LP30 芯片已在量产，预计第三季度出货，而首个 Vera Rubin 机架已在微软 Azure 云上运行。

此外，针对光互联技术，黄仁勋展示了全球首款量产的共封装光学（CPO）交换机 Spectrum X。“铜缆扩展、南向光学扩展（Scale-Up）、北向光学扩展（Scale-Out）三条路线并行推进，我们需要所有合作伙伴在铜缆、光