

芯片公司的“上天”“落地”探索

本报记者 张心怡

算力“上天”

应对 YottaFlops 级计算需求

AI 算力似乎永远不够用。表面上看，是 AI 用户数量急剧膨胀和推理任务日益复杂导致算力需求井喷，但更深层次的原因在于：AI 的思考过程越来越详细。许多大模型不再像初代 ChatGPT 那样直接输出答案，而是有了思考路径，通过持续地自我验证乃至自我否定来寻求最优解，每多出一个思考步骤，都意味着计算量的增长。

“OpenAI o1 模型的引入是 AI 的转折点，推理不再是一次性给出答案，而是一个思考过程。为了教会 AI 思考，强化学习和大量计算被引入后训练阶段，让计算机通过自我尝试来学习如何执行任务，导致用于预训练、后训练、测试时缩放的计算量呈爆炸式增长。现在，我们每进行 1 次推理，都可能生成 2 个 token 而不是 1 个，测试时缩放导致模型生成的 token 数量每年增加 5 倍。”英伟达创始人兼首席执行官黄仁勋在 CES 2026 演讲中指出。

而智能体 (Agent) 的盛行，正在推动 AI 从“被动响应”向“主动决策”的根本性转变，并进一步推高算力需求。

“当我们把 AI 扩展到更广泛的智能体时，全球计算基础设施需求激增的趋势将更加深远。我们需要将计算能力再提高 100 倍，在未来 5 年内达到超过 10 Yotta-Flops，也就是我们在 2022 年所拥有算力的 1 万倍。”AMD 首席执行官苏姿丰 (Lisa Su) 在 CES 2026 开幕演讲中表示。

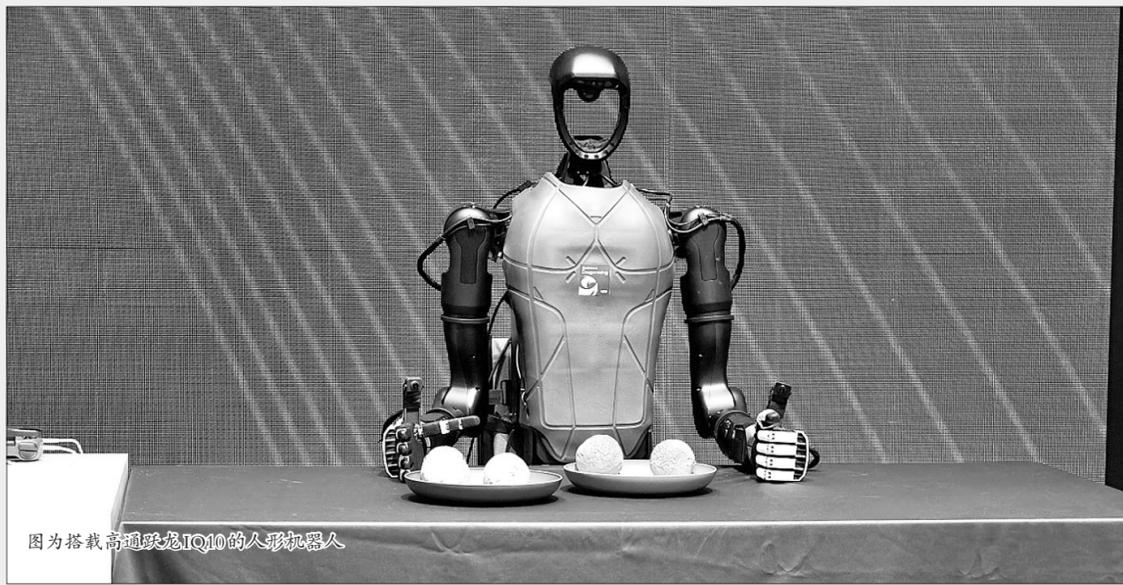
面向 AI“深度思考”及智能体时代的算力需求，AMD、英伟达等企业基于新一代计算单元和不同组件之间的协同设计，以更快的迭代频率、更系统化的计算平台，应对全球 AI 开发者的算力焦虑。

Helios，是 AMD 面向 Yotta 级 AI 算力需求的新一代机架级平台。据介绍，每个 Helios 机架拥有超过 1.8 万个 GPU 计算单元和超过 4600 个 Zen6 CPU 核心，提供 2.9 Exa-flops 性能。每个机架还具备 31TB 的 HBM4 内存，260TB/s 的纵向扩展带宽，以及 43 TB/s 的横向扩展带宽。

Helios 的核心是计算托盘，每个托盘包含 4 个 Instinct MI455X GPU，搭配新一代 EPYC“Venice”CPU 和 Pensando Volcano NIC (网卡)，通过开放的 ROCm 软件栈整合。

其中，AMD 新一代 Instinct MI455X 加速器被苏姿丰形容为“我们制造过的最先进的芯片”，拥有 3200 亿个晶体管 (比 MI355 多 70%)，包含 12 个 2 纳米和 3 纳米制程的计算及 I/O Chiplet，以及 432GB 的 HBM4 显存，所有

一边“拔高”上限，一边务实“落地”——在 CES 2026 上，芯片继续成为最大的看点之一。AMD 和英伟达等企业面向 AI“深度思考”和智能体趋势，向 YottaFlops (10 的 24 次方浮点计算) 级算力需求挺进，推动平台级别和机架级别的算力跃迁。与此同时，高通、英特尔、Arm、AMD、英伟达等企业，都在强调物理 AI，打通 AI 奔赴现实世界的“最后一公里”。



图为搭载高通骁龙 IQ10 的人形机器人

单元通过 AMD 下一代 3D 芯片堆叠技术连接。驱动 MI455X GPU 的是代号“Venice”的下一代 EPYC CPU，采用 2 纳米工艺，拥有 256 个最新高性能 Zen 6 核心；在机架规模下，Venice 能全速为 MI455X 供应数据，实现协同工程。以上组件与 800G 带宽的 Pensando Volcano 网络芯片、Salina DPU 集成，提供超高带宽和超低延迟，使成千上万的 Helios 机架能够在数据中心内扩展。

另据苏姿丰透露，下一代 MI500 系列的开发已经在进行中，该系列基于 AMD 下一代 CDNA 6 架构，采用 2 纳米工艺，并使用更高速的 HBM4E 内存。

“随着 2027 年 MI500 系列的推出，我们有望在 4 年内实现 1000 倍的 AI 性能提升，让更强大的 AI 惠及所有人。”苏姿丰说道。

而英伟达也一改“每一代新平台最多迭代一两颗芯片”的原则，一口气推出集成 Vera CPU、Rubin GPU、NVLink 6 交换机、ConnectX-9 SuperNIC、BlueField-4 DPU、Spectrum-6 以太网交换机 6 款全新芯片的新一代 AI 计算平台 NVIDIA Rubin。据悉，Rubin 平台在 MoE 模型训练中使用的 GPU 数量仅为 Blackwell 平台的 1/4，生成每个 token 的成本低至 1/10。

有意思的是，Rubin GPU 在 NVFP4 (英

伟达提出的 4 位浮点格式) 精度下的推理性能达到 50 PFLOPS，是 Blackwell 的 5 倍，但晶体管数量只有 Blackwell 的 1.6 倍，实现这一目标的关键是英伟达的创新技术：NVFP4 Tensor Core。这种采用新格式的运算引擎不是简单地在数据路径中嵌入某种 4 位浮点数，而是一个完整的处理器单元，懂得如何动态、自适应地调整精度和结构，以应对 Transformer 模型的不同计算阶段，从而在允许损失精度时实现更高的吞吐量，在需要的时候再恢复到最高精度。该技术使 Rubin GPU 能够以相对更少的晶体管增长来实现更大幅度的性能提升，也有助于进一步降低 AI 算力成本。

物理 AI“落地”

打通“奔赴现实世界”最后一关

走进 CES 高通展台，搭载高通骁龙 (Dragonwing) IQ10 系列的人形机器人正在捡拾水果。在展台上，摆放着红色、绿色的盘子各一个，红色、绿色的塑料水果各两个，机器人会将水果抓起来，放到与水果颜色一致的盘子里。无论工作人员或参观者如何变动水果的位置，还是交换两个盘子的位置，机器人依然

会稳稳地拿起水果，放到同色的盘子中。

在本次 CES 上，高通、英特尔、英伟达、AMD、Arm 等芯片公司，都在强调“面向真实世界部署”的物理 AI。高通公司 AI 产品技术中国区负责人万卫星曾在公开演讲中，将 AI 应用的演进分为四个阶段。一是感知 AI，比如传统的自然语言处理、语音降噪、图片识别；二是生成式 AI，基于训练数据创作内容从而响应用户提示；三是智能体 AI，能够自主行动和决策；四是物理 AI，可以理解真实的物理世界，并根据物理定律做出反馈和响应。

要实现“AI 无处不在”的愿景，物理 AI 显然是必不可少的“最后一公里”。

汽车、机器人是高通构建物理 AI 的主要抓手。汽车方面，高通的骁龙数字底盘解决方案已经被全球超过 4 亿辆汽车采用。机器人方面，高通推出下一代完整的机器人技术栈架构，集成硬件、软件和复合 AI，并发布了最新高性能机器人处理器高通骁龙 IQ10 系列。此外，在物联网方面，高通推出全新高通骁龙 Q-8750 和 Q-7790 处理器，聚焦终端侧 AI、多媒体能力、安全特性以及其他先进功能，从而更好地支持无人机、视觉系统、智能摄像头和 AI 电视等广泛的物联网产品形态。

在最新处理器的基础上，高通强调技术

栈构建和完整的产品组合。比如，在机器人领域，高通提供了通用型机器人架构，结合视觉语言动作模型 (VLA) 和视觉语言模型 (VLM) 等端到端 AI 模型，支持先进感知和运动规划，从而赋能泛化操作能力以及人与机器人的交互能力。搭载高通骁龙 IQ10 的通用型机器人架构提供了完整的技术栈，包括异构边缘计算、边缘 AI、混合关键级系统、软件、机器学习运维和 AI 数据飞轮，并依托合作伙伴生态系统与开发者工具套件，使机器人能更高效地进行推理并智能地适应时空环境，经优化后能够在多种形态下实现工业级可靠性的规模化部署。

英特尔在 CES 2026 正式发布第三代英特尔酷睿 (Core) Ultra 处理器，覆盖从 PC 到边缘领域的应用。该处理器是首款基于 Intel 18A 制程打造的计算平台，旗舰型号酷睿 Ultra X9 388H 配备 16 个 CPU 核心 (4 个性能核、8 个能效核、4 个低功耗能效核)、12 个 Xe 核心 (核显) 和 50 TOPS NPU 算力，将赋能超过 200 多款 PC 产品设计。

值得注意的是，英特尔在 3 系列处理器上，首次实现了边缘处理器与 PC 版本同步发布，并首次获得了针对嵌入式和工业边缘场景的测试与认证，包括宽温范围支持、确定性以及 7×24 小时全天候可靠性，加速 AI 在具身智能、智慧城市、自动化与医疗领域的部署。据悉，搭载第三代英特尔酷睿 Ultra 处理器的边缘系统预计将于 2026 年第二季度开始面市。

Arm 将物理 AI 视为 AI 领域发展的核心动能，基于算力支撑，赋能汽车、机器人及各类设备感知、理解现实环境，并在实际场景中安全可靠地运行。特斯拉新一代 AI5 芯片基于 Arm 计算平台打造，其 AI 性能相较上一代提升 40 倍。基于 Arm 架构的 NVIDIA DRIVE Thor 平台为文远知行 L4 级自动驾驶出租车 GXR 所搭载的联想 HPC 3.0 高性能计算平台提供算力支撑。HERE Technologies 借助基于 Arm 架构的 Amazon Graviton 基础设施，更高效地将工作负载从云端迁移到生产环境。

在演讲中，苏姿丰将物理 AI 视为 AI 技术领域最严峻的挑战之一，需要构建能够无缝集成多种类型处理器的机器，以理解环境、做出实时决策，并在无须任何人工输入的情况下采取精确行动，且整个过程对误差零容忍。目前，AMD 技术已经用于训练模拟物理因果关系的大模型系统，支持 World Labs (世界实验室) 构建遵守物理定律和动力学的空间智能，并应用于机器人和太空探索。

“交付物理 AI 需要全栈式的方法，包括用于运动控制和协调的高性能 CPU，用于处理实时视觉和环境数据的专用加速器，以及开放的软件生态系统，使开发者能够快速行动，并在平台和应用程序之间无缝迁移。”苏姿丰表示。

奋力谱写新型工业化发展新篇章

