

数据密集型HPC产业趋势分析

智能计算芯世界 2021-10-18 00:00

按照广义 HPC 市场领域划分的全球本地自建 HPC 收入（百万美元）

	2019	2020	2021	2022	2023	2024	2019-2024 复合 年增长率
服务器	\$13,595	\$13,744	\$13,741	\$16,197	\$17,708	\$18,977	6.9%
存储	\$5,379	\$5,520	\$5,605	\$6,675	\$7,478	\$8,075	8.5%
中间件	\$1,599	\$1,618	\$1,640	\$1,946	\$2,142	\$2,310	7.6%
应用	\$4,647	\$4,682	\$4,643	\$5,380	\$5,783	\$6,092	5.6%
服务	\$2,218	\$2,186	\$2,131	\$2,421	\$2,552	\$2,636	3.5%
总收入	\$27,438	\$27,750	\$27,761	\$32,619	\$35,662	\$38,090	6.8%

来源：Hyperion Research，2021 年 5 月



↑ 点击蓝字，轻松关注

在过去的25年中，本地自建广义HPC市场不断扩大，从1996年的72亿美元增加到2019年的约279亿美元，产业价值几乎翻了两翻，预计2024年将达到377亿美元，成为世界上发展最快的 IT 市场之一。推动这种增长的关键要素包括科研和工程研究对算力日益增长的需求、国家间为拥有最快的超级计算机而进行的长期竞争，以及现有商业技术促使 HPC 走向大众，使 HPC 计算更亲民，甚至许多中小企业也可以用得起。

内容来自Hyperion Research “数据密集型HPC产业趋势白皮书”。

下载地址：

数据密集型HPC产业趋势白皮书
2020年HPC市场总结和预测报告

近年来，市场增长受到新的因素推动，尤其是将 HPC 资源应用于领先的人工智能（AI）和其他高性能数据分析（HPDA）任务，包括这些资源不断被转移到企业数据中心以支持实时业务运营带来的数据分析需求。虽然 HPC 整体市场预计在未来五年（2019-2024 年）将以 6.8%的复合年增长率（CAGR）增长，但HPDA 的市场份额（包括支持 HPC 的人工智能）预计将以 5 年平均 17%的 CAGR 迅猛增长，而 AI 份额的 5 年 CAGR 则达到更高的 33%。

从传统 HPC 建模/仿真应用向新的 HPDA/AI/ML/DL 应用演进的主要特点是从计算密集型负载向数据密集型负载转变。这一转变凸显了存储架构在为研究人员、工程师和业务数据分析师提供最佳性能的 HPC 基础设施中发挥出的关键作用，帮助其最快获得研究和分析结果。从市场角度来看，存储约占整个 HPC 市场的 20%，预计到 2024 年本地 HPC 的存储收入为 80 亿美元。

Hyperion Research 发现，随着数据密集型应用和负载的不断普及，对 HPC 生态系统的需求也不断地发生变化。HPDA/AI 的快速发展同样也推动着传统 HPC 建模/仿真应用的不断转型。HPDA/AI/ML/DL 技术产生越来越多的数据，给现有的 HPC 存储生态系统带来巨大压力，要解决和优化这两种类型的负载就需要高度关注HPC存储基础设施。

市场整体趋势

预计有几个因素将推动 HPC 领域所有细分市场的持续增长，并且这种增长率很可能会超过企业通用IT 领域的预期增长率。传统的 HPC 建模和仿真环境不断扩大并推动市场发展，更多的企业和政府用户正在寻求更快的周转时间，同时增加问题规模、建模保真度和迭代次数。下图总结了对本地自建广义 HPC 的整体市场预测，需要注意的是，存储是广义 HPC 市场中增长最快的领域，约占本地 HPC 市场支出的 20%。

按照广义 HPC 市场领域划分的全球本地自建 HPC 收入（百万美元）

	2019	2020	2021	2022	2023	2024	2019-2024 复合 年增长率
服务器	\$13,595	\$13,744	\$13,741	\$16,197	\$17,708	\$18,977	6.9%
存储	\$5,379	\$5,520	\$5,605	\$6,675	\$7,478	\$8,075	8.5%
中间件	\$1,599	\$1,618	\$1,640	\$1,946	\$2,142	\$2,310	7.6%
应用	\$4,647	\$4,682	\$4,643	\$5,380	\$5,783	\$6,092	5.6%
服务	\$2,218	\$2,186	\$2,131	\$2,421	\$2,552	\$2,636	3.5%
总收入	\$27,438	\$27,750	\$27,761	\$32,619	\$35,662	\$38,090	6.8%

来源：Hyperion Research，2021 年 5 月

在 HPC 市场中，HPDA/AI 细分市场的增长明显大于整个 HPC 市场的增长。具体而言，在 HPC 生态系统的存储领域中，HPDA 存储的复合年增长率为通用 HPC 市场的 2 倍，而 AI 存储的复合年增长率几乎为通用HPC 市场的 4 倍。

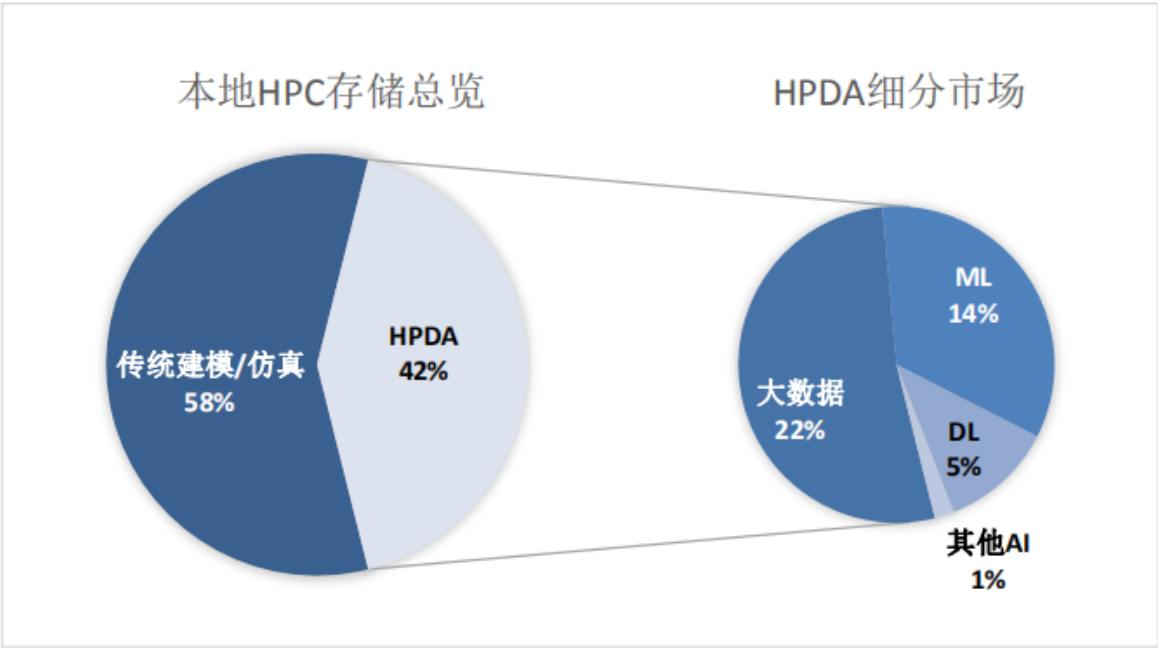
全球本地自建 HPC、HPDA 以及 AI 的存储收入预测（百万美元）

	2019	2020	2021	2022	2023	2024	2019-2024 复合年增 长率
本地 HPC 存储收入	\$5,379	\$5,520	\$5,605	\$6,675	\$7,478	\$8,075	8.5%
HPDA 存储所有收入	\$1,531	\$1,697	\$2,074	\$2,417	\$2,931	\$3,406	17.3%
基于 HPC 的 AI 存储收入	\$391	\$472	\$613	\$800	\$1,242	\$1,619	32.9%

来源：Hyperion Research，2021 年 5 月

HPDA 存储收入是本地自建 HPC 存储收入的子集，基于 HPC 的 AI 存储收入是 HPDA 存储收入的子集。下图显示了 HPC、HPDA 和 AI（机器学习（ML）、深度学习（DL）及其他）等细分场景之间的关系，包括 2024 年相对 HPC 存储市场的分配预测。

基于细分场景的全球本地 HPC 存储份额-2024



来源: Hyperion Research, 2021 年 5 月

影响整个 HPC 市场增长的因素有很多，其中有些因素对数据密集型 HPC 领域的影响尤其深远。HPDA/AI基础设施的应用与日俱增，与之相关的特别值得注意的项目包括：

全球 E 级竞赛将推动多台价值 6 亿美元及以上的 E 级超级计算机发展，其中包括预期内的一台（日本理研的 Fugaku 系统）价值超过 10 亿美元的超级计算机和几台价值超过 1 亿美元的近 E 级超级计算机。全球 E 级竞赛正在为最高端的技术市场注入资金，并为传统建模/仿真和新HPDA/AI/ML/DL 负载的 HPC 基础设施带来许多利益。包括中国、欧盟、日本和美国在内的国家和地区政府越来越关注该行业，这不仅是为了开发领先的系统用于国家关键科学、工程和安全应用，还因为 HPC 对其国内工业和经济具有重要影响，许多国家政府愈发将其视为国家战略资源，要求关键组件不能依赖外国供应来源。因此，众多国家和地区对该项耗资巨大的本土技术开发工作进行了长期投资。在全球疫情大流行期间和随之而来的经济竞争中，这些超大型机器相对按照计划进行安装，并对随后近 E 级超级计算机产生涓滴效应。此外，即便全球疫情大流行，不同国家或地区之间的第一之争也在一直持续。

来自传统 HPC 应用的更高精度建模和仿真正在产生越来越多待分析和存储的数据。因此，用于存储升级和采购的支出带来了总体支出的增长。□ 机器学习和深度学习，新硬件、算法、应用和场景的推出将继续激发人们对 HPC 的兴趣，它们可以为最为严苛的训练、推理和决策支持任务提供快速强大的性能。□ 有充分证据表明，HPC 与企业市场在不断融合，竞争正驱动越来越多的商业组织将其复杂的业务规划和运营要求与其数据和相关数据分析基础相结合，并将更多的业务运营推向实时/近实时环境。此外，为了共同解决传统建模和仿真中的数据密集型 AI/ML/DL 负载以及不断增长的其他大数据分析作业的问题，HPC 和企业计算也在进行融合。

数据密集型 HPC 负载将继续推动新的存储需求，未来的架构将支持同步和空间分布式计算和存储广泛存在于整个 HPC 基础设施中。同样，对于解决如图形分析这类对非模式内存访问有要求的关键负载，物理分布、全局共享的内存技术将变得更加重要。最后，对日益复杂的存储/内存层次的有效管理需要包含嵌入式智能来有效管理数据流，以确保可随时随地使用资源。

HPC 负载使用云资源的情况也在增加。大多数用户将云计算视为对传统 HPC 采购的补充，而不是替代。混合设置通常支持容器开发，这些容器帮助用户负载在本地自建和 HPC 云平台之间编排计算、网络和存储基础设施。企业 HPC 负载也正在向本地自建的私有云扩展。新一代 HPC 架构越来越类似云，HPC 工作流跨越多个容器，每个容器都动态配备适当的硬件和软件资源。最近，Hyperion Research 研究预测，到 2024 年，用户在云中运行 HPC 负载的支出

将达到 88 亿美元，其中大约三分之一（29 亿美元）的支出用于云存储。注意，此项支出完全是指用户将在 HPC 云资源上的支出，不包括云服务提供商(CSP)在支持 HPC 云服务的基础设施上花费的支出。

对数据密集型 HPC 存储要求的展望：定义和挑战

HPDA 泛指利用 HPC 资源的数据密集型负载，包括大数据和 AI 负载。HPDA 问题的特点是数据量大、时效性强以及算法复杂，这对于工资单、电子邮件和一般会计等传统企业业务负载来说影响并不明显。AI 负载是 HPDA 问题的一个重要子集，增长迅速，HPDA 问题寻求从数据本身而不是主要从仿真物理模型中提取价值。为了进一步阐明和定义 AI 负载，适用以下定义：

AI：广泛的通用术语，表示计算机能够做人类想做的事情（但无法以人类思考的方式思考）。AI 包括机器学习、深度学习和其他方法论。

ML：使用示例来训练计算机识别特定模式的过程，例如蓝睛模式或表示欺诈的数字模式。计算机无法学习超过其训练范围的东西，在识别过程中需要人工监督。计算机遵循给定的基本规则。ML 是 AI 的一个子集。

DL：一种先进的机器学习形式，它使用数字神经网络使计算机能够不受其训练内容限制并自行学习，无需额外的显式编程或人工监督。计算机自己制定规则。DL 是 ML 的子集。

HPDA 垂直领域

垂直领域	HPDA（包括 AI） 场景
高级科学与工程研究	减少对建模和仿真的需求（例如，通过消除问题空间的不相关区域）以快速筛选大量文献或源数据；使用代理模型来扩充仿真代码。
自动驾驶	通过在已完成的参数化建模和数据处理中“规模优化”仿真需求，增强 CFD、碰撞、NVH 和实体的建模。
航空航天	加速设计新飞机的仿真，包括为不同的设计挖掘新的解决方案空间。
生物生命科学	通过开发基于 AI 的模型来帮助诊断和治疗疾病，从而加速药物研发。
网络安全	评估来自组织内部和外部成员的威胁；检测欺诈和异常。
国防	监控和信号处理；加密；指挥、控制、通信和情报（C3I）；地理空间图像管理和分析；国防研究；武器设计；面向数据库的模式匹配
能源	地震建模和油藏仿真模型；电网仿真模型；新能源建模
经济学/金融	计量经济学建模、投资组合管理、股票市场和经济预测以及财务分析
政府实验室和研究中心	不受经济限制（例如癌症研究）的广泛科学研究和国防相关项目
人文社科	考古学；文化人类学；历史语言学
物联网(IoT)、边缘计算和智慧城市	城市规划；交通拥堵；空气质量
市场营销与销售	销售分析；收入优化；亲和力营销
天气预报	大气建模、气象学、天气预报和气候建模

来源：Hyperion Research, 2021 年 6 月

此外，AI 工作流通常包括三个阶段：收集、训练和推理。在收集阶段，数据被加载到模型中，通常是加载到大型矩阵中，数据集越大，模型就越准确。加载数据后，通过执行多次计算和比较（矩阵运算），根据参考数据集和预定标准进行数据分配和权重评估来开展训练。执行矩阵运算通常需要大量并行计算能力。训练完成后，可以执行推理（应用预定规则来确定新信息和结果）。

数据密集型 HPC 的挑战

HPDA 和 AI 负载一直在推动 HPC 系统需求突破传统 HPC 系统架构负载的需求，跨越 HPC 系统架构的所有要素。尤其是对存储的需求已经让传统 HPC 存储解决方案达到性能极限，亟需进行多方面创新。

用于传统建模和仿真的传统 HPC 存储通常包括项目文件共享、Scratch 和归档的负载，AI 工作流程则带来一组不同的负载：数据收集和注入、数据准备、训练、推理和归档。有的拥有像传统 HPC 负载那样的存储属性，而有的则推动了新的或更严苛和极端的要求。

HPC 和 AI 负载通常表现出不同的 I/O 模型。传统的 HPC 负载通常基于顺序大 I/O 型，而 AI 负载需要顺序大和随机小 I/O 型的混合，用于 AI 数据集标记的元数据管理需要快速的随机小 I/O 型。

应用场景还催生了各种耐用性和弹性解决方案需求。归档需要极具高性价比的解决方案，没有苛刻的性能要求。传统的临时应用需要高性能，能够将临时结果转移到持久存储以防止出现故障。AI 和 HPDA 解决方案需要混合存储需求满足高性能、瞬态存储和持久弹性存储的要求，包括大块顺序和小块随机 I/O模型的平衡混合。

传统 HPC 和 HPDA 负载

负载	场景	说明
传统 HPC	项目文件共享	<ul style="list-style-type: none">- 通常称为主目录或用户文件- 用于捕获和共享建模和仿真的最终结果- 混合带宽和吞吐量需求，利用混合闪存和 HDD 存储解决方案
	Scratch	<ul style="list-style-type: none">- 用于执行建模和仿真的工作空间容量- 包括元数据容量（高吞吐量 [I/Os/秒] 和闪存化）和原始数据容量以及检查点写入，用于在长时间仿真运行（高带宽 [GB/s]）期间防止系统组件故障，传统上基于 HDD，但现在主要是混合闪存和 HDD。
	归档	<ul style="list-style-type: none">- 长期数据保留
		<ul style="list-style-type: none">- 无关键延迟要求的可扩展存储- 主要是基于云要素持续增加的近线机械化硬盘系统- 通常是文件或对象数据类型

HPDA	数据收集和注入	<ul style="list-style-type: none">- 快速加载各种不同来源的大量数据，以便对数据进行标记、规范化、存储和快速检索以供后续分析- 需要大规模的高带宽 (GB/s)性能来维持检索数据速率，通常是基于大容量 HDD 存储的对象并且逐渐云化
	数据准备	<ul style="list-style-type: none">- 通常称为数据分类或数据标记，需要吞吐量和带宽的平衡组合（混合闪存和 HDD 存储系统）
	训练	<ul style="list-style-type: none">- 利用 ML 和/或 DL 为研究人员、工程师和业务分析师构建准确的模型，用于研究、设计和业务需求- 需要高吞吐量 (IOs/秒)和低延迟才能对数据进行连续和重复的计算分析，通常是基于闪存的存储
	推理	<ul style="list-style-type: none">- 利用该模型进行实验和分析，得出并提供有针对性的科学或商业见解- 还需要高带宽和低延迟，通常基于闪存，带有缓存层
	归档	<ul style="list-style-type: none">- 无关键延迟要求的可扩展存储- 主要是基于近线 HDD 的系统，云要素增加- 通常是文件或对象数据类型

来源：Hyperion Research，2020 年 10 月

最后，数据类型和访问方法推动了对不同类型存储系统的需求发展。结构化和非结构化数据采用不同的访问方式，如文件、块和对象协议。每种访问方式都需要独特的协议支持，通常，这些协议由多个独立的专用系统或一个系统内的不同单元提供，数据通常需要保存多个副本。

数据密集型 HPC 存储方案的机遇

与生活的方方面面一样，挑战往往伴随着巨大的机遇，HDPA/AI 负载的存储也不例外。通过适当关注每个 HPDA/AI 负载和场景提出的要求，存储系统架构师和供应商可以围绕整个 HPC 系统的性能和易用性优化来开展创新。针对 HPDA 负载和数据湖解决方案中经常出现的非结构化数据，我们可以在如下方面进行创新：

支持不同 I/O 模型。提供具有单一文件系统的单一存储架构，可以同时支持大块顺序访问 (TB/s) 和小块随机访问 (IOPS) 的应用性能需求，将消除独立系统的费用并简化存储管理员所需的管理和支持。

支持多协议访问。AI 工作流的不同阶段通常需要使用不同的协议为数据提供服务。可以使用 S3收集对象数据，而训练是通过 NFS 文件访问实现的，在整个 HPDA/AI 工作流中使用的其他协议还包括 MPI-IO、SMB 和 HDFS。与以前需要多份数据副本场景不同，单个系统支持多个接入协议和一份数据能够服务于多个应用，可以节省多系统和额外容量的费用。

支持各种数据访问频率。热数据需要最高的带宽和 IOPS 以支持频繁且及时的访问，而冷数据的访问频率相对较低，对性能的要求也不高。SSD 盘和 HDD 盘具备按需扩展容量和新技术的能力，可以满足热冷数据分级的需求。一个能够随时随地提供合适类型数据的经济高效存储平台是广受HPC 社区欢迎的。

高密度高效设计。HPC 存储解决方案可保留大量数据，这些数据存储占地面积大，通常需要多个设备机架以及关联功率和散热。使用合适的材料促进散热并注意可维护性，也将有助于提供经济高效的解决方案，从而优化 TCO。

下载地址：
数据密集型HPC产业趋势白皮书
中国AI平台市场报告（汇总）
《2021年中国AI开发平台市场报告》

《2021中国AI商业落地市场研究报告》

《中国AI开放平台精品报告》

2020年HPC市场总结和预测报告

ARM架构参考手册及文档

ARM的体系结构与编程.pdf

ARM架构参考手册.pdf

ARM架构参考手册ARM V9.pdf

CPU之战：ARM vs Intel.pdf

ARM系列处理器应用技术完全手册。

CPU和GPU研究框架合集

本号资料全部上传至知识星球，更多内容请登录 **智能计算芯知识（知识星球）** 星球下载全部资料。

END

免责声明：本号聚焦相关技术分享，内容观点不代表本号立场，**可追溯内容均注明来源**，发布文章若存在版权等问题，请留言联系删除，谢谢。

电子书<**服务器基础知识全解(终极版)**>更新完毕，知识点深度讲解，提供182页完整版下载。

获取方式：点击“**阅读原文**”即可查看**PPT**可编辑版本和**PDF**阅读版本详情。

温馨提示：

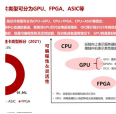
请搜索“**AI_Architect**”或“**扫码**”关注公众号实时掌握深度技术分享，点击“**阅读原文**”获取更多**原创技术**干货。



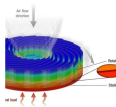
Read more

People who liked this content also liked

大模型算力：AI服务器行业（2023）
智能计算芯世界



CPU处理器散热技术
智能计算芯世界



UCIe封装与异构算力集成
智能计算芯世界



