RISC-V 突破不断 正强势崛起为芯片架构第三极

本报记者 李玉洋 上海报道

全球首款RISC-V大小核处理 器面市、全球首款RISC-V笔记本 正式交付、全球首款开源万兆 RISC-V网络交换机亮相、RISC-V融合服务器全球首发、平头哥推 出首个RISC-VAI平台……近段 时间,RISC-V产业链不断取得新 突破。作为×86、ARM之外的芯片 架构第三极,RISC-V正在全球尤 其是在中国强势崛起。

RISC-V是一个开发、免费的

事件。

关注和应用。

RISC-V生态高速发展

渐进入复杂生态系统,全栈能力不断提高。

在李春强看来,随着RISC-V

比如RISC-V服务器芯片的发

的技术、生态、产业的高速发展,近

年来还涌现出一些里程碑式的积极

布。基于玄铁C910的算能SG2042

服务器芯片和融合服务器,展现了

RISC-V在服务器市场广阔的应用

前景。RISC-V开源开放的特性,

在高性能方向上得到了越来越多的

特尔、英伟达、高通以及平头哥等

13家企业发起的全球RISC-V软

件生态计划"RISE",其主要的目的

是加速RISC-V的软件生态建设及

应用商业化进程,将进一步加速

再比如,今年6月,由谷歌、英

指令集架构,是由加州大学伯克利 分校图灵奖得主 David Patterson 教 授及其课题组基于RISC的CPU指 令集架构,历经30多年研发、迭代 五次而成,2015年加州伯克利大学 将RISC-V指令集架构开源。从此 之后,芯片架构在×86和ARM之 外又多了一个新的选择。

尽管略显年轻,但RISC-V架 构在2022年年底就实现了100亿颗 芯片的出货量。"ARM架构用了17 年完成了这一里程碑,而RISC-V 只用了12年。"电子创新网CEO张

RISC-V在移动通信、数据中心、边

缘计算及自动驾驶等领域的技术和

展历程及特征,李春强表示可大

致分为三个阶段,"在RISC-V发

展早期,大家把RISC-V处理器多

应用于专用芯片,如RF通信、电

源管理芯片等;随着RISC-V指令

集的逐渐完善,越来越多的 IoT

(物联网)、MCU(微控制单元)类

的 SoC 芯片(系统级芯片)采用

RISC-V,包括蓝牙、Wi-Fi、智能

AI、AP类的芯片也得到很好的推

进。"李春强指出,随着RISC-V从

"最近两年,RISC-V在高性能

对于当前 RISC-V 架构的发

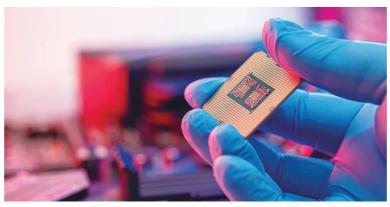
国斌告诉《中国经营报》记者。

"作为开源开放的指令架构, RISC-V对CPU、芯片行业有着深 远的意义。"平头哥玄铁RISC-V软 件研发负责人李春强表示,在短短10 余年间,RISC-V处理器的产业应用 达到100亿颗出货量,显示出这个新 出现的处理器架构旺盛生命力。

平头哥是国内RISC-V产业链 的重要玩家之一。在近日推出首个自 研RISC-VAI平台时,平头哥生态副 总裁杨静甚至说"随着软硬件生态的 逐步成熟,创新的形态不断涌现,所有

芯片都值得用RISC-V做一次"。

多位采访对象向记者表示,在 经历多年的快速成长后,RISC-V 架构向高性能 AI、AP 类等芯片(比 如以上提及的服务器、网络交换机) 推进的趋势愈发明显。"RISC-V的 兴起,在很多人的概念里,应该是自 嵌入市场而起。但实际上,在高性 能计算市场,加速计算的重要性增 强,CPU与其配合的灵活性也变得 重要。"半导体行业资深产业分析师 黄烨锋表示,RISC-V指令集通往 高性能市场是必然会发生的。



RISC-V应用范围正在不断拓展,产业生态不断完善。

本报资料室/图

低成本到高性能的不断拓展, RISC-V能覆盖到对算力要求更高 的领域,比如AI,而参与到RISC-V 生态中来的企业也越来越多。

"首先是操作系统厂商,包括 Google 开展 Android 系统拥抱 RISC-V架构等。其次是AI、车载 和通用计算等领域的芯片公司也逐 渐开始关注或使用RISC-V指令架

构。最后是更多高性能的应用终端 厂家往RISC-V转。"李春强说。

可以看到,在过去的2~3年里, RISC-V从支持小生态系统的芯片 架构,已逐渐进入复杂生态系统,全 栈能力不断提高。"随着软硬件生态 的逐步成熟,创新的形态不断涌现, 所有芯片都值得用RISC-V做一 次。"杨静表示。

高性能应用落地加速

目前,RISC-V高性能全栈技术在多领域展开规模化落地。

在过去的2~3年里,RISC-V从支持小生态系统的芯片架构,已逐

商业化进程。

语音芯片等"。

需要注意的是,在数据中心、服 务器等高性能应用领域将会产生新 的RISC-V需求爆发点。

根据投资机构 ARK Invest 的 预测,到2030年,ARM和RISC-V 可能成为新的处理器标准,在云业 务领域取代英特尔×86架构, ARM + RISC-V 的组合所占据的 市场份额,将从2020年的零,增加 至2030年的71%。

黄烨锋表示,虽然71%这个数 字有待商榷,但ARM和RISC-V

在数据中心拿下更多市场是板上 钉钉的。"数据中心有不同类别的 应用,如果再泛化到网络设备,包 括 networking、存储、HPC(高性能 计算机群)、AI等,虽然软件栈依赖 仍然不见得变少,但对初创企业而 言,解决好其中一两个应用市场与 生态系统,就有机会将CPU产品 做起来,那么企业就有机会存活。"

这是亚马逊自研 CPU 能够快 速部署到自家服务器上的重要原 因,同样RISC-V也有这样的机 会。"目前在数据中心高性能领域, RISC-V架构方面有代表性的企业 典型如平头哥、Ventana、Tenstorrent 等。"黄烨锋指出,即便RISC-V架 构进驻HPC领域的企业大多还很 年轻,但这股力量中的许多新锐绝 对不容小觑。

平头哥在 2023 RISC-V 中国 峰会上发布了首个玄铁 RISC-V 高性能全栈技术,从处理器 IP 到 芯片平台、编译器、工具链等软硬 件技术全面升级,并实现RISC-V 与 Debian、Ubuntu、安卓、 OpenKylin、OpenHarmony、龙蜥、 酷开 WebOS 等主流操作系统的

据悉,目前RISC-V高性能全 栈技术在多领域展开规模化落 地。平头哥携手合作伙伴,实现首 个基于玄铁高性能芯片的安卓商 业化项目落地,量产多款视频视觉 类产品,推出云计算、智能电视等 多场景应用。

中国企业和机构是主要贡献者

"很多本土IC公司、工具公司都在支持RISC-V产业发展,生态日益完善,做了很多好的探索。"

聚焦国内RISC-V产业,业界 形成了哪些业务发展路线? 哪些场 景的产品已经落地?

张国斌表示:"很多本土IC公 司、工具公司都在支持RISC-V产 业发展,生态日益完善,做了很多好 的探索,比如做RISC-V架构的 MCU、蓝牙芯片等。"

深度科技研究院院长张孝荣则 表示相较于×86和ARM,RISC-V

的应用领域仍较少,但在国内企业 不断探索中,逐渐形成了一些发展 路线和落地的产品场景。

"例如,国内一些芯片设计公 司已经开始设计和生产基于 RISC-V架构的处理器芯片,用 于嵌入式设备、物联网设备等领 域;一些大型互联网企业也开始 在自己的服务器和数据中心中采 用RISC-V架构的处理器,以提

高性能和降低成本。"他还表示, 一些初创公司专注于RISC-V的 创新应用,如人工智能芯片、边缘 计算等。

平头哥也认为, AI 正成为 RISC-V的新机遇。"越来越多的AI 引擎采用RISC-V,有直接采用 RISC-V Vector、Matrix指令实现弹 性算力的,也有采用RISC-V作为 主控,实现 NPU(网络处理器)加 速引擎的。"李春强表示,RISC-V 在AI方向的技术将是众多芯片公 司发力的方向。

此外,中国企业和机构已是 RISC-V国际社区的主要贡献者 之一。在RISC-V国际基金会中, 平头哥投入很大力量,参与了30 余个技术方向的标准制定,主导负 责了安卓、数据中心等12个关键 技术小组。

倪光南:应建立合理的"算存比"

越重要。

本报记者 秦枭 哈尔滨报道

ChatGPT 发布至今, AI 大模型 正在进入全新的生态模式,大模型 时代,数据决定AI智能的高度。作 为数据的载体,数据存储成为AI大 模型的关键基础设施。

算力成本约占整个成本的25%,而

数据清洗、预处理等工作,在不算

数据存储硬件的情况下,占到成

本的22%。从这个角度看,数据机

器存储过程,在大模型时代越来

的片面性。"倪光南认为,真正的人

工智能不仅需要算力,还需要存力、

运力,三者缺一不可,只有三者平衡

"大家对算力的理解存在一定

中国工程院院士倪光南表示:

"数据存储产业正成为国家的战略 性、基础性产业与新的国际竞争高 地,我们必须高度重视中国数据存 储产业发展,抓住中国数据存储产

配置、均衡发展,才能充分发挥算力

算,我国存力相对不足,存在重算

力、轻存力的倾向。"在其看来,以数

据存储能力、信息计算能力、网络运

载能力为代表的存力、算力、运力都

是我国信息产业发展的核心和基

不仅如此,国内的存力水平与

础,是建设科技强国的战略支撑。

倪光南表示:"经过存算比的测

的作用。

业面临的重大机遇和挑战,实现科 技自立自强,高质量发展,为科技强 国建设和掌握数字经济竞争主动权 提供坚实支撑。"

海外相比也有一定差距,IDC 公布

的《2023年第一季度中国企业级存

储市场跟踪报告》显示,中国企业级

数据存储市场销售额同比增长

3.45% 至70.14亿元,全闪存储销售

额15亿元,市场占比25%,混闪存

储销售额38亿元,市场占比54%,

相比全球全闪存储市场份额41.3%

的局面,中国全闪存储市场还有很

重算力,轻存力

随着大模型产业的快速发展, AIGC模型预训练数据量呈现指数 级增长,带动算力需求爆发。《中国 经营报》记者了解到,模型参数量从 GPT-1的1.17亿增加到GPT-3的 1750亿;训练数据量也由 GPT-1的 5GB, 增加到 GPT-3 的 45TB。这 也就导致面向AI大模型的数据准 备时间长、数据来源分散、归集慢。

数据统计显示,目前大模型

我国算力总规模已位居全球第 二。工信部发布的数据显示,截至 2022年年底,我国数据中心标准机 架总规模超过7000万架,服务器规 模约2400万台,算力总规模达到 180EFLOPS(每秒 18000 京次浮点 运算)。

算力中心时出现某些倾向,均衡部 署、均衡发展,抢占新一轮科技革命

而目前我国算力中心中采用 SSD 的占比较低,在我国存储市 场中,闪存所占比例还较低,倪 光南认为,要加强政策引导,力推 SSD取代HDD。存储设备应优先 支持SSD,不得以单一价格指标 作为评标依据;重要基础设施的 存储设备,应加强政策引导。

大的发展空间。

不仅如此, 倪光南指出, 国产 数据存储产业正迎来发展窗口 期。存储领域企业应以此为契机, 重视对存储产品的安全审查,包括 供应链安全、信息安全、数据安全 等。尤其是对于存储整机、主控芯 片、存储文件系统等关键核心技 术,建议参照信创工作经验,实施 "自主可控测评"。

快手人局大模型: 站内检索先落地 AIGC提升中

本报记者 李昆昆 李正豪 北京报道

近日,快手公布基于自研 大语言模型应用的最新进展: "快手AI对话"功能已经在安 卓版本开放内测。用户在快手 搜索首页右上角AI图标即可进 人产品首页,选择输入内容对 话或点击引导问题便可开启对 话,或者触发"搜索智能问答" 卡片,通过底部人口可以启用 该功能。

《中国经营报》记者了解 到,"快手AI对话"依托于快 手站内社区内容生态,通过 互动对话的形式,可以帮助 用户快速查找短视频、达人、 百科等内容,体验全新的信 息获取形态。同时,利用大 模型背后海量的数据知识, "快手AI对话"为快手搜索 打破了站内内容的生态壁 垒,有望提供全网检索服务, 为用户提供全新的信息获取 形态。

快手AI对话

记者在快手App里搜索AI 对话,并问到快手为什么要在 App 里做 AI 对话,而不是像文 心一言那样做独立的App。快 手AI对话回复称,"快手选择 在快手App里做AI对话,是因 为它更符合用户的使用习惯 和需求。快手App拥有庞大的 用户群体和丰富的内容生态, 而AI对话能够与这些内容相 结合,提供更加智能、便捷、个 性化的服务。此外,在快手 App 里做 AI 对话也有利于快 手更好地了解用户需求,收集 用户反馈,进一步优化和改进 产品功能。"

盘古智库高级研究员江瀚 告诉本报记者,快手选择将AI 对话功能集成到其主应用程序 中,这可能意味着他们希望为 用户提供更综合和无缝的体 验,使用户在同一个应用程序 中能够完成更多的任务。如果 像文心一言那样做独立的 App 可能需要更多的投入和更高的 成本。这两种方法各有优缺 点,具体取决于快手的产品策 略和用户需求。

"现在有做独立App的,也 有不做独立App的,对于公司 来讲,他可能考虑的一个核心 原因是成本问题。另外,新的 App 出来以后,还要去吸引用 户,做一些推广,这也要花一些 成本。"财经专栏作者马继鹏向 本报记者说。

快手自研的大语言模型叫

"快意"(KwaiYii), 为业务团队 提供标准API和定制化项目合 作方案。快手内部人士称,"快 意"背后的AI团队自2022年至 今已经在国际顶级会议和期刊 上发表了130多篇论文,覆盖机 器学习、自然语言处理、图像识 别、语音识别、搜索推广等多项 研究领域。

据了解,大语言模型需要四 个层:芯片层、框架层、模型层、 应用层。快意数据大语言模型 属于模型层,快手AI对话是在 模型层基础上延伸出来的一个 应用,属于应用层。

江瀚表示,"快手自研的大 语言模型快意与快手AI对话之 间是密不可分的。大语言模型 是实现 AI 对话的关键技术之 一,而快手自研的快意大语言模 型可能为快手AI对话提供了强 大的技术支持。因此,我认为快 手自研的大语言模型快意与快 手AI对话之间是一种支持和互 补的关系。"

和传统大模型相比,快手 AI对话是对于搜索新场景的 探索。一方面将快手站内大 量的内容资源作为索引,解决 大语言模型AI 幻觉的问题,提 升回答准确性;另一方面也用 更加有效的资源组织形式满 足用户多元化需求,不仅覆盖 生活常识、服务查询等内容, 用户还可以进行追问,在个性 化的场景中寻找到更适合自

持续投入AIGC

抖音也在做AI。近日,字 节跳动旗下LLM人工智能机器 人"豆包"现已开始小范围邀请 测试,用户可通过手机号、抖音 或者 Apple ID 登录。据称,"豆 包"是字节跳动公司基于云雀模 型开发的AI工具,提供聊天机 器人、写作助手以及英语学习助 手等功能,它可以回答各种问题 并进行对话,帮助人们获取信 息,支持网页Web平台、iOS以 及安卓平台,但iOS需要使用 TestFlight安装。

除AI对话以外,在短视频 平台中,当前整个内容领域也或 多或少地受到了AI的冲击,AI 写作、AI绘画的技术已经相当 普及,AI短视频、AI直播的技术 离成熟也不远了。比如在抖音 平台的视频中,也有不少AI制 作的虚拟场景,不经提示很难分 辨出来。

"大模型做得好不好,需要 几个硬实力。首先是芯片层, 芯片层在国内做的其实不多, 但大部分公司没有能力做芯 片。"马继鹏说,在芯片层之外, 框架层其实没什么门槛,都是 一些开源框架,包括国内大模 型,都是在开源框架上训练出 来的大模型。

上述人士称,"所以大模型 到底好不好,其实除了技术实力 之外,最重要的就是训练次数有 多少,多少人能参与。ChatGPT 之所以那么火,是因为参与的人 数特别多,全球那么多人参与, 一上线就有几百万人参与,所以 会越来越聪明。但是国内其实 现在不太一样,监管方面也出台 了相关政策。对每一个公司的 大模型,根据监管要求,这些大 模型在公测时也要符合相关监 管规定。"

快手科技创始人兼首席执 行官程一笑介绍,基于目前在 大模型的技术积累,快手已经 实现了多个应用场景的落地。 首先是在搜索方面,于7月8日 启动了智能问答产品的内测,8 月8日启动"AI对话"内测,并 于8月18日在快手App安卓版 本开放内测"快手AI对话"功 能,这是短视频和直播行业首 个基于大语言模型落地的智能 问答产品,在搜索场景为用户 带来智能问答和文本创作等新 功能。

其次是在AIGC方面,快手 已经打造"全模态大模型 AIGC 解决方案"。基于自研的基座大 模型,为用户提供包括文本生 成、图像生成、3D素材生成、音 频生成、视频生成等在内的多种 技术能力。

江瀚称,他对快手在AI方 面的发展持乐观态度。一方面 快手在短视频领域拥有庞大的 用户基础和强大的社交属性,这 为快手提供了丰富的数据资源, 可以为AI技术的发展提供有力 的支持;另一方面,随着人工智 能技术的不断发展,AI在各个 领域的应用越来越广泛,这也为 快手在AI方面的发展提供了广 阔的市场前景。

IDC 预计全球数据量到 2025

掌握先进数据存储主动权

年将达到175ZB,其中我国的数据 量也将由 2018 年的 7.6ZB 增至 48.6ZB,跃居全球第一,而拥有强 大、先进的数据存储产业作为支撑, 才能有发展的主动权。

对此,倪光南建议,产业发 展,标准先行。他提出,为促进产 业更好发展,以"行标"或"团标"的 方式,发布《算力中心建设指南》,

提出"算力"与"存力"的适当比率

倪光南指出,要避免大力发展

和产业变革的制高点。