

## 高性能计算和数据中心融合网络研究综述\*

陆平静,董德尊,赖明澈,齐星云,熊泽宇,曹继军,肖立权  
(国防科技大学 计算机学院,湖南长沙 410073)

**摘要:**随着高性能计算、大数据处理、云计算和人工智能计算呈融合发展趋势,高性能计算网络和数据中心网络的融合网络成为互连网络发展的重要趋势。分析当前融合网络研究现状;针对当前最具代表性的融合网络进行详细阐述,全面展示该领域的最新技术和动态;提出融合网络面临的技术挑战;基于技术挑战对融合网络的发展趋势进行展望,包括融合网络协议栈设计中融合与分化并存、基于在网计算实现融合网络性能加速、面向新兴应用需求优化融合网络性能。

**关键词:**高性能计算网络;数据中心网络;融合网络;在网计算

**中图分类号:**TP301 **文献标志码:**A **开放科学(资源服务)标识码(OSID):**

**文章编号:**1001-2486(2023)04-001-10



与作者互动  
听语音  
聊科研

## Survey on converged networks of high-performance computing network and data center network

LU Pingjing, DONG Dezun, LAI Mingche, QI Xingyun, XIONG Zeyu, CAO Jijun, XIAO Liqun

(College of Computer Science and Technology, National University of Defense Technology, Changsha 410073, China)

**Abstract:** With the convergence trend of high performance computing, big data processing, cloud computing and artificial intelligence computing, the converged network of high-performance computing network and data center network becomes an important trend. The current research status of converged network was analyzed, and the representative converged networks were described in detail to comprehensively show the latest technologies and trends. The challenges faced by the converged network were put forward, and the trends of converged network were proposed, including the convergence and differentiation coexistence of the converged network protocols, the performance acceleration of the converged network based on in-network computing, and the performance optimization of the converged network for emerging applications.

**Keywords:** high performance computing network; data center network; converged network; in-network computing

传统上认为高性能计算(high performance computing, HPC)服务于能力型应用,数据中心(data center, DC)服务于容量型应用。因此,高性能计算网络(high performance computing network, HPCN)和数据中心网络(data center network, DCN)在部署方式,运营理念,网络可用性、可靠性与安全性要求<sup>[1-2]</sup>,多租户虚拟化要求,以及应用程序和编程模型需求<sup>[3-4]</sup>等方面均有很大区别<sup>[5]</sup>,如:HPCN以东西向流量为主,通常采用低直径拓扑<sup>[6-7]</sup>,而DCN以南北向流量为主,通常采用Clos拓扑;DC一般采用增量式部署并向后兼容,而HPC通常高度集中化部署,升级计划通常在初始安装之前制定<sup>[5]</sup>;在HPC应用程序部署中,经常考虑局部性,而DC在部署时一般不需考虑物理相关性,因为容量是按

时间顺序部署的,按物理相关性部署会使虚拟机(virtual machine, VM)分配策略复杂化<sup>[5]</sup>;DC运营商普遍使用虚拟化和多租户来提高管理和资源利用率,而HPCN一般没有虚拟化或多租户需求。

近年来,HPC技术飞速发展,已经进入E级(百亿亿次级)时代<sup>[8-9]</sup>。随着高性能计算机的发展,尤其是使用成本的不断下降,其应用领域从具有国家战略意义的核武器研制、信息安全、石油勘探和高端科学计算领域向更广泛的国民经济主战场快速扩张。国际高性能计算排行榜TOP500<sup>[10]</sup>中,大部分机器并非应用在传统的科学计算领域,而是应用在新兴的互联网云计算和大数据领域。HPC应用从过去的高精尖向更广更宽的方向发展。HPC与云计算、大数据、AI

\* 收稿日期:2023-02-21

**基金项目:**国家重点研发计划资助项目(2021YFB0300101);湖南省自然科学杰出青年基金资助项目(2021JJ10050)

**作者简介:**陆平静(1984—),女,安徽淮北人,副研究员,博士,硕士生导师,E-mail:pingsinglu@nudt.edu.cn;

熊泽宇(通信作者),男,湖南岳阳人,副研究员,博士,硕士生导师,E-mail:xiongzeyu08@nudt.edu.cn

不断融合创新, HPCN 正与互联网技术进行融合<sup>[11]</sup>, 拓展传统 HPCN 支持 DCN 协议栈已成为当前国际高速互连领域的重要发展趋势。

近年来, DC 也迅速发展。2021 年, 全球已经有 326 个超大型 DC<sup>[12]</sup>。亚马逊、谷歌或微软的 DC 规模比最大的单一 HPC 系统还要大<sup>[5]</sup>。如图 1(a)<sup>[5]</sup>所示, 传统 DCN 的角色主要是面向外部客户端提供数据, 并支持在 DC 中运行简单分布式应用程序, 驱动从服务器到 Internet 的南北通信。然而, 随着分布式数据分析和 AI 广泛应用, DC 逐步成为 HPC、大数据、AI、云计算的社会基础设施, DC 的负载模式已经从数据集中、松耦合转化为紧耦合, 大型 DC 流量类似于传统的 HPC 应用程序, DCN 需求迅速融入传统 HPC, 互连网络的吞吐量和延迟需求稳步增长, 与服务器之间通信相关的东西向流量在数量级上占主导地位(如图 1(b)<sup>[5]</sup>所示), 超低延迟通信对于 DC 至关重要。随着大型 DC 采用具有更高带宽需求的高性能加速器, DCN 通过支持远程直接内存访问(remote direct memory access, RDMA)、RDMA 融合以太网<sup>[13-15]</sup>(RDMA over converged Ethernet, RoCE)、互联网广域 RDMA 协议<sup>[16]</sup>(Internet wide area RDMA protocol, iWARP)等新技术不断向

HPCN 融合。

综上所述, 随着 HPC、大数据和 AI 计算呈融合发展趋势, 高性能计算机和数据中心之间的界限越来越模糊, HPCN 和 DCN 融合网络(下文简称融合网络)成为互连网络发展的重要趋势, 从而支撑同一套基础设施高带宽、低延迟运行 HPC、云计算、大数据处理和 AI 计算多领域应用, 降低网络成本<sup>[12]</sup>。

### 1 研究现状分析

当前高性能计算机系统, 既是实现高性能计算的主体, 又是承载云计算、AI 和大数据处理的主体, HPC、云计算、大数据处理和 AI 计算呈融合发展趋势, HPCN 需要支持 DCN 协议, 以满足多领域应用需求。另外, 在当今分布式数据分析和机器学习(machine learning, ML)的时代, DCN 需求迅速融入了传统 HPC 的各个方面。苏金树教授提出: 大规模高效网络计算中的网络技术发展趋势主要包括 3 个方面: 融合、分化、优化<sup>[12]</sup>。融合体现在不同领域的网络技术没有明显分界线; 分化体现在不同领域的独特解决方案或者新应用需求下的创新方案; 优化体现在针对特定场景的技术优化实现。邬江兴院士提出<sup>[17-18]</sup>, 现有的自进化的网络技术发展范式已经不能满足新的发展需求, 必须从思维方式、方法论和实践规范等方面进行全面改革, 并提出了多态智能网络环境, 为网络技术的研究与开发提供了新的范式和思路。如图 2 所示, 融合网络从 HPCN 和 DCN 两个方向不断融合发展。

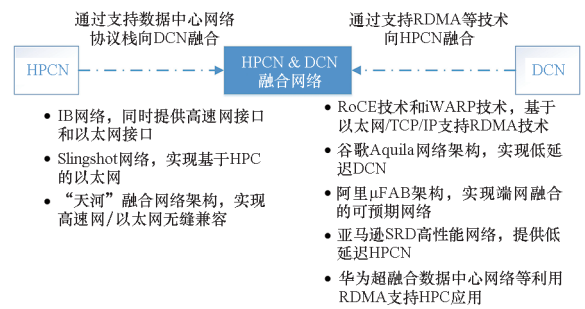
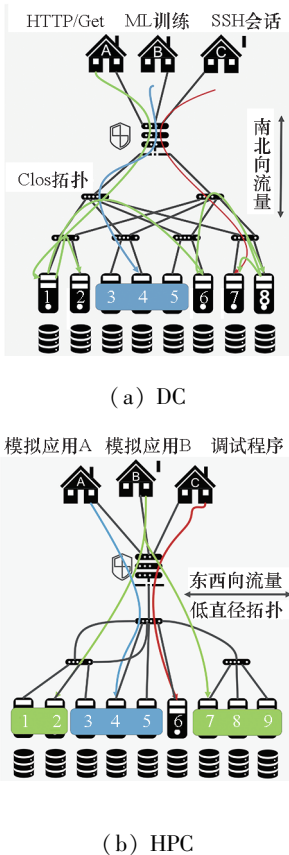


图 2 融合网络研究现状

Fig. 2 Research status of converged network

一方面, 高性能计算机通过拓展传统 HPCN 以支持 DCN 协议栈, 从而不断向 DCN 融合, 满足云计算、大数据处理和 AI 计算多领域应用需求。在 HPC 场景中, 当前有两种主流方案来承载 RDMA: 专用无限带宽(InfiniBand, IB)网络和以太网。IB 高性能互连网络通过开发多模芯片、

图 1 数据中心和高性能计算的使用场景<sup>[5]</sup>

Fig. 1 Usage scenarios of DC and HPC<sup>[5]</sup>

设计基于 IB 的以太网(Ethernet over InfiniBand, EoIB)协议向以太网融合,已经推出多款多网络融合的芯片产品<sup>[19-21]</sup>,具有低延迟和高带宽等高性能,可以极大地提高高性能计算系统和数据中心的性能。Cray 的 Slingshot 技术以 HPC 为中心增加了以太网兼容性,其交换机兼容传统以太网并对 RoCE 的一些不足进行了改进,同时支持高性能计算和数据中心<sup>[22-24]</sup>,交换机平均延迟是 350 ns<sup>[12]</sup>。国防科技大学在自主定制高速互连网络<sup>[25-27]</sup>的基础上提出一种融合网络架构,实现高速网/以太网无缝兼容,灵活支持科学计算和云计算等多领域应用<sup>[28-29]</sup>。但是,由于 IB 网络以及其他私有定制网络采用私有协议,架构封闭,难以与现有大规模的 IP 网络实现很好的兼容互通,同时,IB 网络运维复杂,开销居高不下。因此用以以太网承载 RDMA 数据流,已应用在越来越多的 HPC 场景。

另一方面,DCN 通过支持 RDMA、RoCE、iWARP 等技术实现低延迟网络,不断向 HPCN 融合,数据中心通过支持 RDMA 协议来支撑 HPC 应用场景。许多大型数据中心运营商使用或计划使用 RDMA 机制实现低延迟网络,RDMA 用于数据中心内部通信,数据中心之间仍使用 TCP 通信<sup>[5]</sup>。例如,谷歌推出全新低延迟数据中心网络架构 Aquila,通过 1RMA 协议支持高性能计算、机器学习训练和网络分解,同时与数据中心规模的 TCP/IP/以太网流进行互操作<sup>[30-31]</sup>;阿里提出了下一代数据中心网络架构“可预期网络”,通过“阿里云全栈自研+端网融合技术”实现的高性能可预期网络将 DCN 从“低时延大带宽”演进到“确定性可预期”,不仅支持新兴的大算力和高性能场景,也适用于通用计算场景,成为融合传统网络和未来网络的产业趋势<sup>[32-34]</sup>,端网融合的可预期网络成为达摩院 2023 十大科技趋势之一<sup>[35]</sup>;微软云平台 Azure 采用 IB 互连连接 GPU 服务器<sup>[15,36]</sup>;微软数据中心、谷歌数据中心、微软 Azure 云计算、百度机器学习和腾讯云等优化现有的以太网络,利用 RDMA 来满足在线服务、大规模数据中心和云计算对网络延迟、吞吐量和 CPU 计算性能的严格要求,形成了“一切都在 RDMA 之上”的局面<sup>[12,37]</sup>;美国伊利诺伊大学提出 RDMA 驱动的加速框架 RAMBDA,可以在数据中心实现微秒级通信延迟<sup>[38]</sup>。

下面分别从上述两方面对代表性融合网络技术进行介绍和分析。

## 1.1 HPCN 融合 DCN,支持云计算、大数据、AI 应用

当前, TOP500 中采用的典型 HPCN 技术包括: Mellanox 公司(2019 年被 NVIDIA 公司收购)的 IB 网络、Cray 公司(2019 年被 HPE 公司收购)的 Slingshot 互连、Intel 公司的 Omni-path 网络、以太网、富士通的 Tofu 网络、Bull 公司的 BXI 以及国内的 Sunway、TH Express 等定制/专属网络。随着 HPC 系统不断承载云计算、大数据处理和 AI 计算, Mellanox、Cray、国防科技大学分别推出了融合网络,使得 HPC 系统在满足传统 HPC 应用的同时,支持云计算、大数据、AI 应用。

### 1.1.1 Mellanox IB 网络

Mellanox IB 网络是 HPC 领域活跃着的第一大互连网络,随着 HPC、AI 和模拟数据对低延迟和加速的需求日益增加,IB 已成为 TOP500 的首选网络。根据 2022 年 11 月国际高性能计算排行榜<sup>[10]</sup>, TOP10 系统中有 4 台采用了 IB 互连, TOP500 的系统中有 194 台采用了 IB 网络。

IB 网络通过开发多模芯片、设计 EoIB 协议等向以太网融合,已经推出多款多网络融合的芯片产品<sup>[19]</sup>。网络接口芯片方面, ConnectX - 5、ConnectX - 6、ConnectX - 7 采用虚拟协议互连(virtual protocol interconnect, VPI),支持 IB 和以太网接口,以太网接口支持 RoCE 技术,兼容性强,并通过智能加速引擎提高性能,可为数据中心提供灵活的高性能解决方案<sup>[19-20,24]</sup>。网络交换芯片方面, QM8700、QM9700 系列产品中 Quantum 交换机,通过通信加速器、可伸缩的分层聚合和归约协议(scalable hierarchical aggregation and reduction protocol, SHARP)实现在网计算,满足从 HPC 到 ML 的网络带宽和时延需求<sup>[19,21,24]</sup>。其最新的网卡 ConnectX - 7 提供下一代数据速率(next data rate, NDR) 400 Gbit/s 网络接口; ConnectX - 6、ConnectX - 7 通过在网计算加速引擎进行性能加速,满足 HPC、AI 和超大规模云数据中心的需求<sup>[20]</sup>。

最新研究中, NVLink 设计实现了用于高通信带宽 SuperPODs 的 NVIDIA 交换芯片 NVLINK4<sup>[39]</sup>,基于 NVLINK4 的 NVSWITCH 除了是历史上规模最大、带宽最高的交换机以外,还增加了计算能力为 400 GFlops 的 FP32 SHARP 加速器,实现归约操作的卸载计算,对于 AI 应用中通信密集型操作,数据吞吐率几乎翻倍。物理层兼容 400 Gbit/s 以太网和 IB 专用高速网, DGX H100 服务器提供 8 个 400 Gbit/s 的 ConnectX - 7

以太网/IB 端口, 以及 2 个双端口 BlueField-3 DPU; 支持多轨 IB/以太网; 有效满足 HPC、AI 和超大规模云数据中心的需求。

### 1.1.2 Cray Slingshot 网络

Slingshot 网络是 Cray 公司于 2019 年推出的新一代高性能网络互连技术, 是 Cray 继 SeaStar<sup>[40]</sup>、Gemini<sup>[41]</sup>、Aries<sup>[6,42]</sup> 之后新一代高性能网络互连技术。根据 2022 年 11 月国际高性能计算 TOP500 排行榜<sup>[10]</sup>, TOP10 系统中有 3 台采用了 Slingshot 网络, 包括排名第一的 Frontier 系统, 美国即将推出的另外两台 E 级系统 (EI Capitan 和 Aurora) 也将采用 Slingshot 互连。TOP500 中, 有 30 台系统采用 Slingshot 互连。

与 Cray 以前的互连不同, Slingshot 将高速互连协议建立在标准以太网之上, Cray 称之为“HPC 以太网”。Slingshot 将专有 HPCN 的优势带入了高度可互操作的以太网标准, 即 Slingshot 交换机首先使用标准以太网协议进行操作, 但是当连接的设备支持高级“HPC 以太网”功能时, 将尝试协商高级功能, 同时为高性能计算和数据中心提供支持<sup>[22]</sup>。基于 Slingshot 交换机, Cray 使用 dragonfly 拓扑构建大型 HPC 系统, 但 Slingshot 互连可支持任何数量的拓扑。由于 Slingshot 网络的底层是标准以太网协议, 其交换延迟与 IB 网络相比仍有差距, 在 300 ~ 400 ns 之间近正态分布, 平均交换延迟 350 ns, 这其中依然包含了约 150 ns 的纠错附加延迟, 但相比标准以太网 450 ns 的交换延迟, 其交换延迟性能还是有较大提升<sup>[22-24]</sup>。

### 1.1.3 “天河”融合网络

“天河”高性能计算机系统采用的是由国防科技大学自主研发的高速互连网络<sup>[25-27]</sup>。

针对现有高速互连网络无法同时支持高速网和以太网、无法有效支持计算密集和数据密集型应用的问题, 提出一种融合网络架构<sup>[28-29]</sup>, 如图 3 所示<sup>[28]</sup>, 该架构包含 PCIE 主机接口处理模块、高速网网卡核心逻辑、交叉开关 XBAR、以太网网卡核心逻辑、以太网报文拆分/拼装模块、物理层逻辑、高速网/以太网报文转换模块 (Ethernet over high performance express, EoH) 以及高速网/以太网可配的网络端口。EoH 将高速网虚拟为以太网, 使得连接在高速网中的节点直接与连接在以太网中的节点通信, 通过高速网传输以太网报文, 实现高速网/以太网无缝兼容, 在一套物理硬件上灵活支持科学计算和云计算应用。

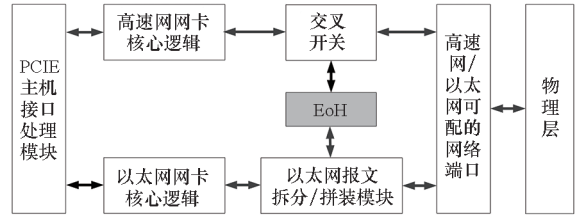


图 3 高速网/以太网融合网络接口卡结构示意图<sup>[28]</sup>  
Fig. 3 Structure diagram of the interconnect network/Ethernet converged network interface card<sup>[28]</sup>

## 1.2 DCN 融合 HPCN, 支持高性能计算应用

DCN 为了追求高吞吐量、低延迟、低成本、易于管理, 一般采用以太网。随着数据中心的负载模式从传统数据集中、松耦合转化为紧耦合, 数据中心对超低延迟通信的需求越来越强烈。谷歌、亚马逊、阿里等数据中心通过支持 RDMA、1RMA 协议等实现低延迟通信。

### 1.2.1 RoCE 技术和 iWARP 技术

RDMA 技术最早在 IB 专用传输网络上实现, 技术先进, 性能最优, 但价格高昂, 应用局限在 HPC 领域。随着以太网性能的大幅提升, 越来越多的人想要选择能兼容传统以太网的高性能网络解决方案, 而传统 TCP/IP 堆栈应用无法支撑 HPC 网络通信。业界厂家把 RDMA 技术移植到传统以太网上, 降低了 RDMA 的使用成本, 推动了 RDMA 技术普及。如图 4 所示<sup>[16]</sup>, 根据协议栈融合度的差异, 分为 RoCE 和 iWARP 两种技术, 而 RoCE 又包括 RoCEv1 和 RoCEv2 两个版本。

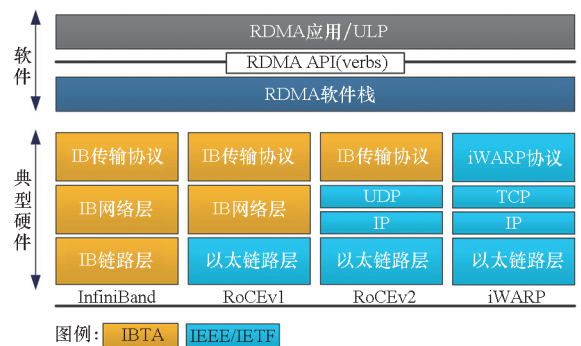


图 4 RoCEv1、RoCEv2 和 iWARP 技术<sup>[16]</sup>

Fig. 4 RoCEv1, RoCEv2 and iWARP technology<sup>[16]</sup>

RoCEv1 协议基于以太网承载 RDMA, 保留了 IB 与应用程序的接口、传输层和网络层, 将 IB 的链路层和物理层替换为以太网的链路层和网络层。由于 RoCEv1 协议没有继承以太网的网络层, 其报文结构是在原有 IB 报文上增加二层以太网报文头, 并没有 IP 字段, 因此 RoCEv1 数据包不能被三层路由, 数据包的传输被局限在二层网

络中路由,只能部署于二层网络<sup>[16]</sup>。RoCEv2 协议对 RoCEv1 协议进行了改进,将 RoCE 协议保留的 IB 网络层替换为以太网网络层和使用 UDP 协议的传输层,并利用以太网网络层 IP 数据包中的区分服务代码点(differentiated service code point, DSCP)和显示拥塞通知(explicit congestion notification, ECN)字段实现了拥塞控制,基于 UDP/IP 协议承载 RDMA,可部署于三层网络。RoCEv2 支持基于源端口号 hash,采用等价多路径(equal cost multi-path, ECMP)实现负载均衡,提高了网络的利用率<sup>[16]</sup>。

iWARP 基于标准 TCP/IP 协议栈支持 RDMA<sup>[16]</sup>。iWARP 在标准的网络层和传输层上运行,TCP 协议栈提供流量控制和拥塞管理,并且不需要无损以太网。iWARP 实现高度路由可扩展的 RDMA 技术,但需要支持 iWARP 的特殊网卡来支持在标准以太网交换机上使用 RDMA。

基于传统以太网承载 RDMA 是 RDMA 大规模应用的必然。比较 RoCE 技术和 iWARP 技术:iWARP 技术由于 TCP 协议的限制失去了绝大部分 RDMA 的性能优势,已经逐渐被业界所抛弃;RoCEv2 架构生态和应用不断成熟,使用 RoCEv2 承载高性能分布式应用已经成为一种趋势。

### 1.2.2 谷歌 Aquila 架构

2022 年 4 月谷歌在顶级会议 NSDI 上发表了 Aquila 架构方案<sup>[30]</sup>,Aquila 是一种实验性的数据中心网络架构,将超低延迟作为核心设计目标,同时也支持传统的数据中心业务。

Aquila 芯片架构<sup>[30-31]</sup>基于 GNet 协议设计了融合交换和网卡的定制芯片,具有低延迟远程存储访问(remote memory access, RMA),如图 5 所示<sup>[30]</sup>,Aquila 芯片架构包含 100 Gbit/s 的 IP 网卡、1RMA 网卡、基于信元的 GNet 交换芯片以及 IP 协议引擎。当流量进入交换机时,一部分通过 IP 网卡走传统的基于数据包的以太网交换,一部分通过 1RMA 网卡走基于信元的 GNet 交换。芯片中间的 IP 协议引擎负责两种交换单位的转换,将 IP 数据包切割处理为多个信元或者将信元重新组装为 IP 数据包。Aquila 架构能够实现 40  $\mu$ s 以下的 IP 流量拖尾结构往返时间(round-trip time, RTT)和低于 10  $\mu$ s 的跨数百台主机的 RMA 执行时间,尾部延迟大幅减少。Aquila 架构采用全连接的 dragonfly 拓扑,实现了 Aquila 控制平面和以太网控制平面的融合统一,使数据中心兼容以太网且具有超低传输时延成为可能<sup>[31]</sup>。

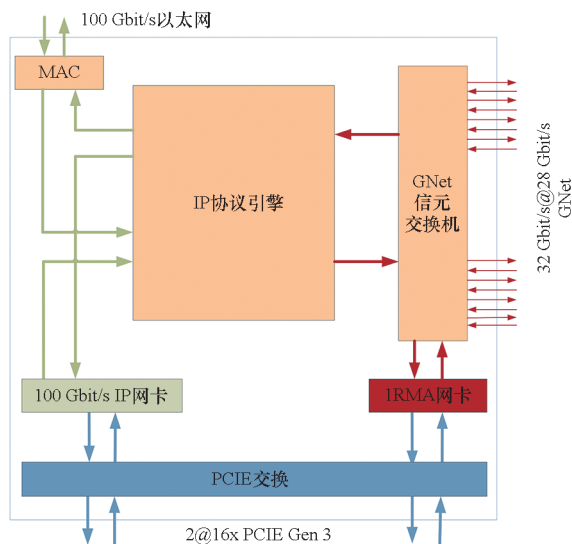


图 5 Aquila 芯片架构<sup>[30]</sup>

Fig. 5 Aquila chip architecture<sup>[30]</sup>

### 1.2.3 阿里高性能网络

随着数据中心的资源池化成为主流,以及 ML/HPC 等大型分布式系统和应用对网络极低延迟要求,网络延迟成为数据中心网络的性能瓶颈,高性能网络成为数据中心性能发挥的关键。针对此,阿里自主研发了高性能网络,先后经历了三个阶段的演进<sup>[43-44]</sup>:第一阶段(2017—2018年),RDMA 大规模落地,阿里云在多个产品中实现 RDMA,通过端到端的流控实现应用性能调优,通过建立端网协同的运营体系,去优先级流控(priority flow control, PFC)实现无损 RDMA,消除 RDMA 稳定性风险;第二阶段(2019—2020年),自研高性能网络协议、高性能网卡和高性能通信库<sup>[45-46]</sup>,实现自主高性能网络;第三阶段(2021年至今),以应用为中心,通过“阿里云全栈自研+端网融合技术”,实现高性能可预期网络。2022年8月阿里在顶级会议 SIGCOMM 上发表端网融合架构<sup>[32]</sup>,提出了要将数据中心网络从“低时延大带宽”演进到“确定性可预期”的目标,开启了确定性数据中心网络研究的新纪元<sup>[31]</sup>。

阿里云自研高性能网络通过 Solar-RDMA RD 传输模式、高精度拥塞控制(high precision congestion control, HPCC)算法和多路径乱序传输协议,成功解决了大规模网络的扩展性难题,同时可以更好容忍网络异常,并消除了慢输入输出<sup>[43-44]</sup>(slow IO)。进一步,Nimitz 容器网络通过单根 IO 虚拟化(single root IO virtualization, SR-IOV)实现了设备的 IO 虚拟化,并通过开放虚拟网卡(open virtual switch, OVS)流表的硬件加速实现了容器间 RDMA 通信的极致性能。最后,通

过  $\mu$ FAB 技术<sup>[31-32]</sup> 和网络统一服务架构 (network unified service architecture, NUSA) 实现了网络性能以及在极端场景下的行为可预期高性能网络<sup>[43-44]</sup>, 成为阿里云高性能集群的坚实技术底座。

#### 1.2.4 亚马逊高性能网络技术

考虑到 TCP 协议的延迟问题, 以及 RoCEv2 需要 PFC 优先级流量控制, 在大规模网络上会造成队头阻塞、拥塞扩散和偶尔的死锁, 不适合亚马逊云服务 (Amazon web service, AWS) 可扩展性要求, AWS 针对超大规模数据中心对现有 AWS 网络进行了优化, 构建了新的网络传输协议——可扩展的可靠数据报 (scalable reliable datagram, SRD), 不保留数据包顺序, 而是通过尽可能多的网络路径发送数据包, 同时避免路径过载。为了最大限度地减少抖动并确保对网络拥塞波动的最快响应, 在 AWS 自定义 Nitro 网卡中实施了 SRD, 并添加了新的 HPC 优化网络接口, 作为 Nitro 网络功能的扩展<sup>[47-48]</sup>。SRD 由 EC2 主机上的 HPC/ML 框架通过 AWS 弹性结构适配器 (elastic fabric adapter, EFA) 内核旁路接口使用, 允许 HPC/ML 应用程序在 AWS 公有云上大规模运行, 使客户能够在 AWS 上大规模运行紧密耦合的应用程序。特别是, EFA 支持运行 HPC 应用程序和 ML 分布式训练, 目前支持多种消息传递接口 (message passing interface, MPI) 实现: OpenMPI、Intel MPI 和 MVAPICH, 以及 NVIDIA 聚合通信库<sup>[48]</sup>。AWS 高性能网络技术 SRD 既能提供 HPC 应用程序所需的持续低延迟, 同时又能保持公有云的优势: 可扩展性强、按需弹性容量、成本效益以及快速采用更新的 CPU 和 GPU<sup>[48-49]</sup>。

## 2 融合网络技术发展挑战

### 2.1 HPC 大规模部署 RoCE 面临延迟挑战

RoCE 技术在 HPC 应用上的延迟与 HPC 专用网络相比相对较高。HPC 对网络的诉求主要集中在高吞吐率和低时延, RoCEv2 基于无连接协议的 UDP 协议, 相比面向连接的 TCP 协议更加快速、占用 CPU 资源更少, 但其不像 TCP 协议那样有滑动窗口、确认应答等机制来实现可靠传输, 一旦出现丢包, RoCEv2 需要依靠上层应用检查到了再做重传, 会大大降低 RDMA 的传输效率。Mellanox 以太网交换机 SN3000 系列的 200 Gbit/s 以太网交换机官方延迟数据是 425 ns, 而 Mellanox IB 交换机 QM8700 系列 HDR 交换机的

官方延迟数据是 130 ns<sup>[19]</sup>。Slingshot 网络平均交换延迟 350 ns<sup>[22]</sup>, 相比标准以太网网络 450 ns 的交换延迟有较大提升, 但与 IB 网络相比仍有差距, 而且 Cray 上一代定制网络 Aries 交换机延迟大约是 100 ns<sup>[6]</sup>, “天河 2A” 的交换机延迟是 81 ns。基于 RoCE 技术的以太网交换机与 IB 交换机以及一些定制的 HPCN 的延迟性能有一定差距。TOP500 中高端 HPC 依然倾向于采用 HPC 专用网络, 如 2022 年 11 月 TOP500 中, 其 TOP10 系统中, “富岳” “神威·太湖之光” “天河 2A” 均采用定制网络, 4 台系统采用 IB 网络, 其他 3 台采用的是 Slingshot 网络。RoCE 技术在 HPC 应用上实现大规模部署, 延迟仍有一定优化空间。

### 2.2 数据中心大规模部署 RoCE 面临安全和性能挑战

RDMA 如何同时实现低延迟和高吞吐量仍然是一个悬而未决的问题。数据中心为了让上层业务能充分利用 RDMA 带来的高性能能力, 需要正确配置、使用一系列 RDMA 关键网络技术, 并根据业务场景的需求, 做相应的调整和优化。

首先, RoCEv2 中使用 PFC 实现无损网络并不能保证低延迟而且可能会带来安全挑战。PFC 是一种粗粒度机制, 它以端口 (或端口加优先级) 级别运行, 不区分流, 当网络发生拥塞时, 可能会产生队列并导致拥堵蔓延, 进而有不公平现象、受害者流、PFC 死锁、暂停帧风暴、慢收现象<sup>[15]</sup> 等一系列性能问题。队列和 PFC 暂停帧都会增加网络延迟, 在规模更大、跳数较多的网络中, 可能会给网络运营带来较大的稳定性风险, 对 RDMA 传输会带来一些安全挑战。要想让应用真正利用 RDMA 的高性能优势, 设计、使用和配置符合业务场景需求的流量控制、拥塞控制机制, 缓解 PFC 缺陷, 使得在迅速变化的流量模式下仍然能保持低延迟, 是数据中心大规模部署 RoCE 面临的重大挑战。

其次, RDMA 无损网络会牺牲部分流<sup>[5]</sup>。尽管许多数据中心运营商使用或计划使用 RDMA 机制, 但在数据中心的 RDMA 和 TCP/IP 流之间的缓冲区和带宽共享可能会牺牲一些流量。在当今的 RDMA 实现中, 简单的基于硬件的重传机制依赖于无损传输层。然而, 大多数 DCN 传统上像互联网一样, 使用有损路由器, 当队列满时丢弃数据包。RDMA 对无损网络的要求为在保守的数据中心环境中采用 RDMA 提出了技术挑战。

### 3 融合网络技术趋势分析

#### 3.1 融合网络协议栈设计中融合与分化并存

融合网络协议栈设计中融合与分化并存:

1) 融合网络融合趋势主要表现在网络层、链路层和物理层。如 RoCEv2 保留了 IB 与应用程序的接口、传输层,使用以太网的链路层、网络层和基于 UDP 协议的传输层,并且利用以太网网络层 IP 数据报中的 DSCP 和 ECN 字段实现了拥塞控制的功能,基于 UDP/IP 协议承载 RDMA。网络层、链路层和物理层融合,可以共用底层硬件和部分设计逻辑,节省面积、存储等资源,HPC 和数据中心都可以使用相同的硬件基础设施,从而显著降低成本。

2) 融合网络分化趋势主要表现在传输协议层、用户接口和软件协议栈针对性设计与优化。传输层协议针对特殊需求进行了针对性的设计,例如针对高端 HPC 开发了先进的定制化互连系统,为数据中心集成了内部路由器和网卡,实现了基于快速 UDP 互联网连接(quick UDP Internet connections, QUIC)和传输层安全的 HTTP/3,优化了数据安全传输效率<sup>[12]</sup>。此外,不同的网络协议可能要求不同的软件栈设计,针对多网络融合,需要完成软件栈移植。如基于高速互连网络驱动语义,借鉴 Open Fabrics 网络接口标准,进行消息数据传输服务、RDMA 数据传输服务、完成事件服务等通信库兼容性封装,采取标准化接口完美支持 MPI、分块全局寻址空间等编程模型;为了支持 TCP/IP 服务功能,利用地址主动注册/地址单播查询机制实现地址解析协议(address resolution protocol, ARP)功能,并且基于 Open Fabrics 内核封装接口实现虚拟层叠网络通信机制,提供基于 TCP/IP 高带宽通信机制,从而支持基于对象的分布式文件系统 Lustre 等;基于高速互连网络驱动语义设计网络管理编程接口,采取专用管理报文流控制机制及报文处理流程,简化网络管理部署并增强管理功能鲁棒性。

#### 3.2 基于在网计算实现融合网络性能加速

虽然 HPC 和数据中心的网络需求在很大程度上是相似的,但在设计和部署理念、运营理念、服务多样性、协议栈、网络利用率、应用程序编程接口等方面特性需求上仍然存在一定差异<sup>[5]</sup>,主要区别在于协议堆栈的上层。基于智能网卡和交换机的在网计算使用特定于应用程序的协议来融合 HPC 网络和数据中心网络,从而实现工作负载

定制化,弥补二者之间的差异。

数据中心的在网计算<sup>[38,50-52]</sup>一般使用智能网卡和交换机的网络多功能流处理用于网络加速,如 DPU、IPU 或 NPU。微软的 Catapult 和 AWS 的 Nitro 网卡都是作为基础设施支持,用于提高多租户主机的安全性、效率和成本。AMD 将 FPGA 和 DPU 的优势用在了网络接口控制器上,实现了 AMD 400 Gbit/s 自适应智能网卡,通过公共的存储和片上网络,融合 ASIC 逻辑、可编程逻辑和嵌入式处理器等不同类型的加速部件以卸载 CPU 中不同类型的计算<sup>[52]</sup>。美国伊利诺伊大学提出一种 RDMA 驱动的加速框架 RAMBDA,可以在数据中心实现微秒级通信延迟<sup>[38]</sup>。

HPC 的在网计算一般通过将聚合通信卸载到网卡和交换机上实现。协议卸载的目标是相对复杂的计算功能,进一步实现计算和网络的融合。如 Mellanox 从 ConnectX-6 开始通过在网计算加速引擎进行性能加速,提出了聚合通信协处理器的思想,在数据传输过程中处理数据,开发了 SHARP 协议,通过将聚合通信操作卸载到网卡芯片,有效地减少了 MPI 聚合通信时间,网卡数据吞吐量提高了 2 倍以上<sup>[39]</sup>;释放了大量的 CPU 资源,实现了计算和通信的高度重叠,有效地加快了机器学习应用的速度,满足 HPC、AI 和超大规模云数据中心的需求。NVLink 通过 SHARP 加速器,实现对 AllReduce 中归约操作的卸载计算,对于 AI 应用中通信密集型操作,数据吞吐率几乎翻倍。高性能计算系统将大规模部署智能网卡,并将应用程序定制的协议卸载到专门的硬件<sup>[5,53]</sup>。

#### 3.3 面向新兴应用需求优化融合网络性能

首先,新兴应用需求将改变互联网基本的传输机制。随着云边端网络融合,互联网端到端的设计原则重边缘轻核心,尽力转发模型也无法适配异构终端和未来多样化应用的需求,迫切需要针对大延迟间歇性连通链路、强实时高带宽连接、内容寻址等异构网络环境设计新的应用驱动的传输机制。

其次,面向 HPC 和数据中心日趋增长的系统规模,内嵌智能计算将增强网络自治自愈能力。网络路由协议收敛时间长、网络配置错误引发振荡、网络故障诊断费时费力,这些问题严重降低了网络可用性和健壮性。内嵌智能计算的网路系统可以挖掘时空多维度规律,从而为网络自动配置、智能诊断、自免疫和自愈合等提供技术支撑。

再次,安全大数据分析、云数据中心部署将驱动网络安全能力提升。安全大数据分析的方法将

向基于机器学习和人工智能的“认知”方向演进,实时数据和历史数据、机器学习方法和人的经验将得到有效结合,从而支撑更有效的未知攻击检测,更全面、更深度的安全态势感知,自动化程度更高的实时检测和响应等能力。面向云数据中心多租户环境,采用虚拟化、大数据、机器学习等使能技术以及 IPSec 和 TLS 进行安全隔离,信任根、安全根、安全固件升级、授权容器实现全方位安全防护。进一步,通过 RDMA 虚拟化可将 RDMA 用到云上<sup>[36,54]</sup>。例如,面向 HPC、大数据分析、AI 应用多租户、高安全、高性能需求,Mellanox 结合 NVIDIA “BlueField DPU”架构和 NVIDIA Quantum IB 交换机构建云原生超算平台,面向云计算环境提供多租户安全隔离,以及最优的裸金属性能。

最后,RoCEv2 中使用 PFC 保证无损传输,但在大规模数据中心大规模部署 RoCE 技术,对 RDMA 传输会带来一些安全挑战<sup>[36]</sup>。优化无损网络性能成为业界研究热点与技术趋势<sup>[55-59]</sup>。华为的超融合数据中心网络智能无损技术首创了基于 AI 的网络智能无损技术。另外,随着新一代 RDMA 网卡片上资源的增加,不依赖于 PFC 的 RoCEv2 网络成为可能,新一代 RDMA 网卡片上实现了更为高效的丢包恢复机制和更好的端到端流控来约束传输中的数据包,从而让有损网络的性能有潜力不输于无损网络,这也将有潜力把 RDMA 推广到规模更大、跳数更多的网络中<sup>[39]</sup>。

## 4 结论

高性能计算网络和数据中心网络多网络融合,从而支撑同一套基础设施高带宽、低延迟运行 HPC、云计算、大数据处理和 AI 计算多领域应用,降低网络成本,是当前互连网络发展的重要趋势。本文对当前融合网络研究现状进行分析。针对 HPC 大规模部署 RoCE 面临延迟挑战,以及数据中心大规模部署 RoCE 面临安全和性能挑战,提出融合网络的技术发展趋势,包括:融合网络协议栈设计中融合与分化并存——融合网络协议在网络层、链路层和物理层日趋融合,传输协议层、用户接口和软件协议栈针对性设计与优化;基于在网计算实现融合网络性能加速;面向多领域新兴应用需求优化融合网络性能。希望我们所做的工作能为未来在该领域的探索提供有用的指导,供相关系统设计者和研究人员参考。

## 参考文献 (References)

[1] ROTHBERGER B, TARANOV K, PERRIG A, et al.

- ReDMark: bypassing RDMA security mechanisms [C]// Proceedings of the 30th USENIX Security Symposium, 2021.
- [2] TARANOV K, ROTHBERGER B, PERRIG A, et al. sRDMA—efficient NIC-based authentication and encryption for remote direct memory access [C]// Proceedings of the 2020 USENIX Annual Technical Conference, 2020.
- [3] HOEFLER T, DINAN J, THAKUR R, et al. Remote memory access programming in MPI-3 [J]. ACM Transactions on Parallel Computing, 2015, 2(2): 1-26.
- [4] SOUMAGNE J, CARNIS P, ROSS R. Advancing RPC for data services at exascale [J]. Bulletin of the Technical Committee on Data Engineering, 2020, 43(1): 23-34.
- [5] HOEFLER T, HENDEL A, ROWETH D. The convergence of hyperscale data center and high-performance computing networks[J]. Computer, 2022, 55(7): 29-37.
- [6] KIM J, DALLY W J, SCOTT S, et al. Technology-driven, highly-scalable dragonfly topology [C]//Proceedings of International Symposium on Computer Architecture, 2008.
- [7] KATHAREIOS G, MINKENBERG C, PRISACARI B, et al. Cost-effective diameter-two topologies: analysis and evaluation[C]//Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, 2015.
- [8] ZHENG W M. Research trend of large-scale supercomputers and applications from the TOP500 and Gordon Bell Prize[J]. Science China Information Sciences, 2020, 63(7): 1-14.
- [9] DONGARRA J. 高性能计算及其未来需求[J]. 中国计算机学会通讯, 2023, 19(1): 1-5.
- DONGARRA J. An overview of high performance computing and future requirements [J]. Communications of CCF, 2023, 19(1): 1-5. (in Chinese)
- [10] MEUER H, DONGARRA J. TOP 500 supercomputer sites [EB/OL]. [2023-02-21]. <http://www.top500.org>.
- [11] 张云泉, 袁良, 袁国兴, 等. 2022 年中国高性能计算机发展现状分析与展望[J]. 数据与计算发展前沿, 2022(6): 3-12.
- ZHANG Y Q, YUAN L, YUAN G X, et al. State-of-the-art analysis and perspectives of China HPC development in 2022[J]. Frontiers of Data & Computing, 2022(6): 3-12. (in Chinese)
- [12] SU J S, ZHAO B K, DAI Y, et al. Technology trends in large-scale high-efficiency network computing [J]. Frontiers of Information Technology & Electronic Engineering, 2022, 23(12): 1733-1746.
- [13] InfiniBand™ Trade Association. Supplement to InfiniBand™ architecture specification volume 1 release 1.2.1 annex A16: RDMA over converged ethernet (ROCE) [R]. InfiniBand™ Trade Association, 2010.
- [14] InfiniBand™ Trade Association. Supplement to InfiniBand™ architecture specification volume 1 release 1.2.1 annex A17: ROCEv2 (IP Routable ROCE) [R]. InfiniBand™ Trade Association, 2014.
- [15] GUO C X, WU H T, DENG Z, et al. RDMA over commodity Ethernet at scale [C]//Proceedings of the 2016 ACM SIGCOMM Conference, 2016.
- [16] Cavium, Chelsio, Intel. iWARP\* RDMA here and now [EB/OL]. [2023-02-20]. <https://www.intel.com/content/dam/www/public/us/en/documents/technology-briefs/iwarp-rdma-here-and-now-technology-brief.pdf>.
- [17] WU J X. Revolution of the development paradigm of network



- technology system—network of networks [ J ]. *Telecommunication Science*, 2022, 38(6): 3–12.
- [18] JI X S, WU J X, JIN J, et al. Discussion on a new paradigm of endogenous security towards 6 G networks[J]. *Frontiers of Information Technology & Electronic Engineering*, 2022, 23(10): 1421–1450.
- [19] NVIDIA. NVIDIA InfiniBand switches [EB/OL]. [2023–02–21]. <https://www.nvidia.com/en-us/networking/infiniband-switching/>.
- [20] NVIDIA. ConnectX–7 400 G adapters [EB/OL]. [2023–02–20]. <https://nvdam.widen.net/s/cs8rnmqwl/infiniband-ethernet-datasheet-connectx-7-ds-nv-us-2544471>.
- [21] NVIDIA. NVIDIA quantum-2 QM9700 series-scaling out data centers with 400 G IB smart switches [EB/OL]. [2023–02–21]. <https://nvdam.widen.net/s/k8sqcr6gzb/infiniband-quantum-2-qm9700-series-datasheet-us-nvidia-1751454-r8-web>.
- [22] DE SENSI D, DI GIROLAMO S, MCMAHON K H, et al. An in-depth analysis of the Slingshot interconnect [ C ]// *Proceedings of SC20: International Conference for High Performance Computing, Networking, Storage and Analysis*, 2020.
- [23] HPE. HPE Slingshot: the interconnect for the exascale era [EB/OL]. [2023–02–21]. <https://assets.ext.hp.com/is/content/hpedam/documents/a50002000-2999/a50002368/a50002368enw.pdf>.
- [24] LU P J, LAI M C, CHANG J S. A survey of high-performance interconnection networks in high-performance computer systems [ J ]. *Electronics*, 2022, 11(9): 1369.
- [25] WANG R B, LU K, CHEN J, et al. Brief introduction of TianHe exascale prototype system [ J ]. *Tsinghua Science and Technology*, 2020, 26(3): 361–369.
- [26] LIAO X K, PANG Z B, WANG K F, et al. High performance interconnect network for Tianhe system [ J ]. *Journal of Computer Science and Technology*, 2015, 30(2): 259–272.
- [27] PANG Z B, XIE M, ZHANG J, et al. The TH Express high performance interconnect networks [ J ]. *Frontiers of Computer Science*, 2014, 8(3): 357–366.
- [28] 肖立权, 常俊胜, 赖明澈, 等. 一种融合网络接口卡、报文编码方法及其报文传输方法: CN111641622B [ P ]. 2022–01–07.  
XIAO L Q, CHANG J S, LAI M C, et al. Converged network interface card, message coding method and message transmission method thereof; CN111641622B [ P ]. 2022–01–07. (in Chinese)
- [29] XIAO L Q, CHANG J S, LAI M C, et al. Converged network interface card, message coding method and message transmission method thereof; US20210367906 [ P ]. 2021–11–25.
- [30] GIBSON D, HARIHARAN H, LANCE E, et al. Aquila: a unified, low-latency fabric for datacenter networks [ C ]// *Proceedings of 19th USENIX Symposium on Networked Systems Design and Implementation*, 2022.
- [31] 黄玉栋. 谷歌阿里竞速: 开启确定性数据中心新纪元 [EB/OL]. (2022–09–20) [2023–02–21]. <https://cloud.tencent.com/developer/article/2114063>.  
HUANG Y D. Google Alibaba competition: opening a new era of deterministic data centers [EB/OL]. (2022–09–20) [2023–02–21]. <https://cloud.tencent.com/developer/article/2114063>. (in Chinese)
- [32] WANG S, GAO K H, QIAN K, et al. Predictable vFabric on informative data plane [ C ]// *Proceedings of the ACM SIGCOMM 2022 Conference*, 2022.
- [33] MIAO R, ZHU L J, MA S, et al. From luna to solar: the evolutions of the compute-to-storage networks in Alibaba cloud [ C ]// *Proceedings of the ACM SIGCOMM 2022 Conference*, 2022.
- [34] Alibaba Cloud. 块存储性能 [EB/OL]. [2023–02–21]. <https://www.alibabacloud.com/help/doc-detail/25382.htm>.  
Alibaba Cloud. EBS performance [EB/OL]. [2023–02–21]. <https://www.alibabacloud.com/help/doc-detail/25382.htm>. (in Chinese)
- [35] 达摩院. 达摩院十大科技趋势: 端网融合的可预期网络 [EB/OL]. (2023–01–11) [2023–02–21]. [https://www.xdyanbao.com/doc/jc1p8m4ps4?bd\\_vid=11055259518058075307](https://www.xdyanbao.com/doc/jc1p8m4ps4?bd_vid=11055259518058075307).  
DAMO Academy. Top ten technology trends of DAMO academy: predictable fabric [EB/OL]. (2023–01–11) [2023–02–21]. [https://www.xdyanbao.com/doc/jc1p8m4ps4?bd\\_vid=11055259518058075307](https://www.xdyanbao.com/doc/jc1p8m4ps4?bd_vid=11055259518058075307). (in Chinese)
- [36] ZHU Y B, ERAN H, FIRESTONE D, et al. Congestion control for large-scale RDMA deployments [ J ]. *Computer Communication Review*, 2015, 45(4): 523–536.
- [37] GAO Y X, LI Q, TANG L B, et al. When cloud storage meets RDMA [ C ]// *Proceedings of 18th USENIX Symposium on Networked Systems Design and Implementation*, 2021.
- [38] YUAN Y F, HUANG J H, SUN Y, et al. RAMBDA: RDMA-driven acceleration framework for memory-intensive  $\mu$ s-scale datacenter applications [ C ]// *Proceedings of IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, 2023.
- [39] ISHII A, WELLS R. The NVLink-network switch: NVIDIA's switch chip for high communication-bandwidth SuperPODs [ C ]// *Proceedings of 2022 IEEE Hot Chips 34 Symposium (HCS)*, 2022.
- [40] BRIGHTWELL R, PEDRETTI K T, UNDERWOOD K D, et al. SeaStar interconnect: balanced bandwidth for scalable performance [ J ]. *IEEE Micro*, 2006, 26(3): 41–57.
- [41] ALVERSON R, ROWETH D, KAPLAN L. The Gemini system interconnect [ C ]// *Proceedings of 18th IEEE Symposium on High Performance Interconnects*, 2010.
- [42] ALVERSON B, FROESE E, INC C, et al. Cray XC series network [ R ]. Cray Inc, 2012.
- [43] 张彭城. 阿里高性能网络探索与实践 [EB/OL]. (2022–06–20) [2023–02–21]. <https://www.modb.pro/doc/65923>.  
ZHANG P C. Exploration and practice of Alibaba high performance network [EB/OL]. (2022–06–20) [2023–02–21]. <https://www.modb.pro/doc/65923>. (in Chinese)
- [44] StrokMitream. 2021 中国智能网卡研讨会回顾 [EB/OL]. (2021–10–23) [2023–02–21]. <https://zhuanlan.zhihu.com/p/424957555>.  
StrokMitream. Review of 2021 China smart network card seminar [EB/OL]. (2021–10–23) [2023–02–21]. <https://zhuanlan.zhihu.com/p/424957555>. (in Chinese)
- [45] LI Y L, MIAO R, LIU H H Q, et al. HPCC: high precision congestion control [ C ]// *Proceedings of the ACM Special Interest Group on Data Communication*, 2019.
- [46] DONG J B, CAO Z, ZHANG T, et al. EFLOPS: algorithm and system co-design for a high performance distributed

- training platform [ C ]//Proceedings of IEEE International Symposium on High Performance Computer Architecture (HPCA), 2020.
- [47] SHALEV L, AYOUB H, BSHARA N, et al. A cloud-optimized transport protocol for elastic and scalable HPC[J]. IEEE Micro, 2020, 40(6): 67–73.
- [48] 软硬件融合. 用于弹性可扩展的 HPC 云优化传输协议 [ EB/OL ]. ( 2021 – 11 – 22 ) [ 2023 – 02 – 21 ]. <https://aijishu.com/a/1060000000255945>. Software and Hardware Integration. Cloud optimized transmission protocol for elastic and scalable HPC [ EB/OL ]. ( 2021 – 11 – 22 ) [ 2023 – 02 – 21 ]. <https://aijishu.com/a/1060000000255945>. ( in Chinese )
- [49] AWS. Amazon EBS 功能 [ EB/OL ]. [ 2023 – 02 – 21 ]. <https://aws.amazon.com/ebs/features/>. AWS. Amazon EBS features [ EB/OL ]. [ 2023 – 02 – 21 ]. <https://aws.amazon.com/ebs/features/>. ( in Chinese )
- [50] SEYEDROUDBARI H, VANAVASAM S, DAGLIS A. Turbo: SmartNIC-enabled dynamic load balancing of  $\mu$ s-scale RPCs [ C ]//Proceedings of IEEE International Symposium on High-Performance Computer Architecture (HPCA), 2023.
- [51] WANG Z L, LUO L Y, NING Q S, et al. SRNIC: a scalable architecture for RDMA NICs [ C ]// Proceedings of the 20th USENIX Symposium on Networked Systems Design and Implementation, 2023.
- [52] DASTIDAR J, RIDDOCH D, MOORE J, et al. AMD 400 G adaptive SmartNIC SoC: technology preview [ C ]// Proceedings of 34th Symposium of Hot Chips, 2022.
- [53] DE SENSI D, DI GIROLAMO S, ASHKBOOS S, et al. Flare: flexible in-network allreduce [ C ]//Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis, 2021.
- [54] KIM D, YU T L, LIU H H. FreeFlow: software-based virtual RDMA networking for containerized clouds [ C ]// Proceedings of 16th USENIX Symposium on Networked Systems Design and Implementation, 2019.
- [55] MITTAL R, LAM V T, DUKKIPATI N, et al. TIMELY: RTT-based congestion control for the datacenter [ C ]// Proceedings of the ACM Conference on Special Interest Group on Data Communication, 2015.
- [56] BABLANI G. Introducing new product innovations for SAP HANA, expanded AI collaboration with SAP and more [ EB/OL ]. [ 2023 – 02 – 21 ]. <https://azure.microsoft.com/en-us/blog/introducing-new-product-innovations-for-sap-hana-expanded-ai-collaboration-with-sap-and-more/>.
- [57] KUMAR G, DUKKIPATI N, JANG K, et al. Swift: delay is simple and effective for congestion control in the datacenter [ C ]// Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, 2020.
- [58] AWS. Elastic fabric adapter: 大规模运行 HPC 和 ML 应用程序 [ EB/OL ]. [ 2023 – 02 – 21 ]. <https://aws.amazon.com/hpc/efa/>. AWS. Elastic fabric adapter: run HPC and ML applications at scale [ EB/OL ]. [ 2023 – 02 – 21 ]. <https://aws.amazon.com/hpc/efa/>. ( in Chinese )
- [59] 梦豪. 工业级大规模 RDMA 技术杂谈 [ EB/OL ]. ( 2022 – 05 – 31 ) [ 2023 – 02 – 21 ]. <https://zhuanlan.zhihu.com/p/510323418>. MENG H. Discussion on industrial scale RDMA technology [ EB/OL ]. ( 2022 – 05 – 31 ) [ 2023 – 02 – 21 ]. <https://zhuanlan.zhihu.com/p/510323418>. ( in Chinese )