

# 面向海量数据迁移的高吞吐量

Research on High Throughput Transmission  
Technology for Massive Data Migration

## 传输技术研究

张 锴,晏家豪,洪 伟(中讯邮电咨询设计院有限公司,北京 100048)

Zhang Kai,Yan Jiahao,Hong Wei(China Information Technology Designing & Consulting Institute Co.,Ltd.,Beijing 100048,China)

### 摘要:

在应对海量数据广域传输任务时,IP网络的效率显著不足,难以满足如东数西算等远距离、大规模的数据迁移场景的需求。深入剖析了导致传输效率低下的原因,并在此基础上,提出了一系列提升传输效率的技术方案建议。

### 关键词:

传输控制协议;带宽时延积;高吞吐量;远程直接数据访问

doi:10.12045/j.issn.1007-3043.2026.04.011

文章编号:1007-3043(2026)04-0062-07

中图分类号:TN913

文献标识码:A

开放科学(资源服务)标识码(OSID):



### Abstract:

The efficiency of IP networks in dealing with massive data wide area transmission tasks is significantly insufficient, making it difficult to support long-distance and large-scale data migration scenarios such as Eastern data and Western computing. It deeply analyzes the reasons for low transmission efficiency. Based on this, it proposes a series of technical solutions aimed at improving transmission efficiency.

### Keywords:

TCP; Bandwidth-delay product; High throughput; RDMA

引用格式:张锴,晏家豪,洪伟. 面向海量数据迁移的高吞吐量传输技术研究[J]. 邮电设计技术,2026(4):62-68.

## 1 概述

### 1.1 研究背景

随着数字中国战略与东数西算战略的深入实施,国家正积极促进算力资源、业务及数据的高效流通,催生了一系列跨地域的海量数据传输需求。若采用低速率专线电路,海量数据的传输周期过长,难以体现时间效率;而若选择高带宽专线,则受限于IP网络中传输控制协议(TCP)的流量与拥塞控制技术瓶颈,难以充分发挥高带宽专线的传输潜能。国内某运营

商的实地网络测试结果显示,在5 000 km的远距离传输场景下,即便采用100G带宽的电路进行海量数据传输,其实际有效吞吐量也仅约能达到电路带宽的14%。因此,需要采用更加高效、可靠的数据传输技术来提高数据传输效率及传输质量。

### 1.2 海量数据迁移场景

#### 1.2.1 东数西算场景

在国家提出的东数西算工程中,通过“东数西存”“东数西渲”“东数西算”优化数据中心布局,实现东西部存力、算力资源共享关系的合理匹配,这必将需要频繁地进行海量数据的远距离传输。

#### 1.2.2 智算中心场景

收稿日期:2026-02-16

在人工智能领域,基础模型和行业模型的训练均依赖海量的模型参数与样本数据。为了提升训练的精度,需要的训练样本数据量也日益增大,通常可达TB至PB量级。在大模型训练阶段,这些数据需被导入至智算中心;训练结束后,训练数据及结果还需回传给用户;在应用阶段,训练完成的模型及其参数则需部署至推理服务器。

### 1.2.3 科学计算场景

在科研领域,特别是在超算快速发展的背景下,超算中心面临着大规模数据的导入与导出需求。以FAST天文数据计算为例,FAST每年需处理约200多个观测项目,每个项目产生的观测数据量可达TB至PB量级,其年产数据量约为15PB。

### 1.2.4 其他行业场景

在其他行业同样存在海量数据远距离传输的需求。以影视素材的传送为例,影视节目的拍摄素材需经后期制作公司进行剪辑、渲染。根据拍摄和制作周期,这些拍摄素材需批量传输至后期制作公司所在地。一部大型综艺或影视节目的原始素材数据量可达PB量级,而单次传输的数据量通常在10TB至100TB量级。

值得注意的是,这里所提及的海量数据传输特指以文件或数据块形式进行传输,对数据完整性和无差错传输有较高的要求,不包括流媒体、在线交易等类型的数据。当这些数据通过专线电路传输时,会展现出数据量巨大、数据流少、传输持续时间长等特点,这类数据流通常被称为大象流。

## 2 IP网数据传输技术

### 2.1 IP网络数据传输协议

IP网络传输层协议主要包括用户数据报协议(UDP)和传输控制协议(TCP)。

UDP本身不提供重传机制,利用UDP协议传输数据时,需通过应用层系统来进行流控及差错检查。当UDP数据包丢失或出现差错时,需要应用层来进行数据重传,因而采用UDP进行数据可靠传输时效率较低。

TCP是一种面向连接的、可靠的、基于字节流的传输层通信协议。它在数据发送前需要进行三次握手来建立连接,并具备重传和流量控制功能,以确保数据包能够准确无误地传输到接收端。TCP因此成为数据可靠传输的首选。

### 2.2 TCP数据传输技术

#### 2.2.1 TCP数据传输方式

在数据传输时,TCP协议需要接收方对收到的分组发送确认消息。为确保发送的数据分组不超出接收方的处理能力,TCP协议采用了称为滑动窗口协议的方法,允许发送方在停止并等待接收确认消息前可以连续发送不超出窗口大小的多个分组。由于发送方不必每发一个分组就停下来等待确认,该协议可以加速数据的传输<sup>[1]</sup>。

#### 2.2.2 TCP流量控制

接收方维护接收窗口,而发送方维护发送窗口。接收端只允许发送端发送接收端缓冲区所能接纳的数据,这样可以防止发送较快的主机导致发送较慢的主机的缓冲区溢出。由于窗口的大小限制了收、发端可接收或发送数据量的大小,因此滑动窗口协议可以实现流量控制<sup>[2]</sup>。

#### 2.2.3 TCP拥塞控制

TCP协议的拥塞控制与流量控制非常类似。流量控制侧重于主机角度,而拥塞控制则从中间网络的角度出发。拥塞控制通过拥塞窗口cwnd来实现,该窗口的大小是发送方维护的一个状态变量,会根据网络的拥塞程度进行动态调整。发送窗口swnd的大小取决于拥塞窗口cwnd和接收窗口rwnd中的较小值。为了实现拥塞控制,RFC2581<sup>[3]</sup>定义了4种算法:慢开始、拥塞避免、快重传和快恢复。其中,慢开始和拥塞避免算法构成了1988年提出的TCP Tahoe版本的拥塞控制策略。为了进一步提升TCP性能,1990年的TCP Reno版本在此基础上增加了快重传和快恢复2个新的拥塞控制算法。

TCP协议拥塞控制过程示意如图1所示。

慢开始机制指的是,在数据发送初期,发送端会将拥塞窗口的大小设定为单个最大报文段(MSS)的容

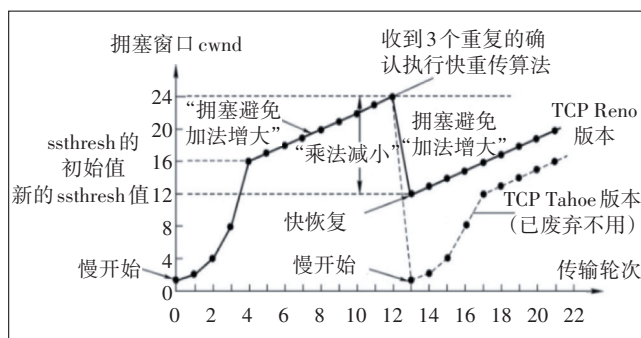


图1 TCP协议拥塞控制过程示意

量。随后,每当接收到一个新的报文段确认,发送端便会将拥塞窗口的数值扩大一倍。通过这种方式,发送端的数据发送速率得以按照指数规律稳步增长。

拥塞避免是在慢开始阶段结束之后使用的机制。当拥塞窗口值达到慢开始门限(ssthresh)初始值时,发送端停止使用慢开始算法,转而使用拥塞避免算法。在拥塞避免阶段,发送端的拥塞窗口每经过一个轮次就增加一个MSS的大小,这样发送速率按线性规律缓慢增加。

当拥塞窗口增加到一定值并检测到网络拥塞时,马上将拥塞窗口降为1,并开始执行慢开始算法,进入一个新的循环过程。

有时个别报文段会在网络中丢失,但实际上网络并未发生拥塞,这将导致发送方超时重传,并误认为网络发生了拥塞,于是错误地启动慢开始算法,因而降低了传输效率。快重传算法对此进行了改进:当TCP发送方连续收到3个重复的ACK时,会立即重传丢失的数据包。这不需要等待重传计时器超时,从而减少了重传延迟。

快恢复是在快重传之后,将拥塞窗口的阈值减半,并进入拥塞避免阶段,然后每经过一个轮次就增加一个MSS的大小。

### 2.3 TCP拥塞控制协议对网络通量的影响

IP网络原本是为了支持传统的互联网业务而设计的,其IP转发机制采用了统计复用和尽力而为的策略,并且相关的流量控制和拥塞控制协议也是基于这些传统业务需求制定的。因此,在处理如网页浏览、电子邮件等传统业务时,IP网络通常能够提供较高的吞吐率和良好的用户体验。

然而,当IP网络用于进行远距离的海量数据传输,特别是传输大象流时,其性能就会受到显著影响,网络吞吐率会大幅下降。这是因为在专线链路传输大象流数据时,传统业务场景中的流量复用现象会大大减弱甚至消失。在这种情况下,提高大象流的传输效率成为提升链路吞吐率的关键。远距离传输还带来了另一个问题,即往返时延(RTT)的显著增加。这会导致发送端难以及时感知到网络拥塞、丢包等网络事件,从而对这些事件的处理产生滞后。此外,TCP早期版本中采用的加性增加、乘性减少(AIMD)拥塞窗口调整算法在长距离传输中显得过于保守。一旦发生拥塞,拥塞窗口会大幅减小,并且由于时延较大,窗口需要很长时间才能恢复,这进一步限制了网络链路的

通量。

### 2.4 窗口大小对网络通量的影响

在数据通信中,链路的通量不仅取决于传输速率本身,还取决于带宽时延积(BDP)。带宽时延积指的是一个数据链路的带宽 $B(\text{bit/s})$ 与来回通信延迟 $\text{RTT}(\text{s})$ 的乘积。

$$\text{BDP} = B \times \text{RTT}$$

带宽时延积代表了在某一时刻,网络线路上已传输但尚未得到确认的最大数据量,同时也是确保发送方与接收方在链路上实现最大吞吐量所必需的缓冲空间容量。TCP流控窗口的大小同样体现了通信过程中可发送但尚未被确认的数据量的上限。若该窗口大小小于链路的带宽时延积,则意味着链路的利用率将无法达到最优<sup>[4]</sup>。

在系统内核层面,TCP发送端与接收端的滑动窗口分别对应 send\_buffer 和 receive\_buffer 这2个缓冲区,且窗口大小是动态变化的参数,其最大值受限于相应缓冲区的大小。为分析链路的通量,采用缓冲区作为研究的基础。

在已知网络链路带宽及路径基本确定的情况下,缓冲区大小 $C(B)$ 应不小于带宽时延积 $\text{BDP}/8$ ,即:

$$C \geq \text{BDP} \div 8 \geq B \times \text{RTT} \div 8$$

时延RTT主要由3个部分构成:传输时延、网络节点的转发时延以及末端系统的处理时延。其中,传输时延与传输链路的距离密切相关,而转发时延和处理时延则呈现出较大的不确定性。在进行量化分析时,暂将时延RTT近似取定为双向传输时延的2倍。

对于100G带宽电路,在不同的传输距离或时延的情况下,电路带宽如果要得到充分利用,系统所需的缓冲区大小(GB)需不小于表1中的数值。

表1 电路带宽充分利用所需缓冲区大小

指标	距离(km)/RTT(s)						
	100/2	500/10	1 000/20	2 500/50	4 000/80	5 000/100	10 000/200
所需缓冲区/GB	0.025	0.125	0.25	0.625	1	1.25	2.5

在RFC 793中,TCP报文的Windows Size字段占2B,窗口最大为65 535 B,即64 KB。对于常被称为长肥网络的高速长距离网络,其带宽时延积远大于64 KB,因而过小的窗口大小将导致链路管道容量的利用率低。

在已知缓冲区容量及RTT的情况下,收发两端之

间传输链路可以达到的最大通量为:

$$B_{\max} = \frac{C \times 8}{RTT}$$

链路的最大吞吐量将随着RTT的增加而下降,随着缓存容量或滑动窗口大小的增加而提升。在不同的系统缓冲区容量及RTT值的条件下,链路可达到的最大吞吐量(Gbit/s)如表2所示。

表2 不同距离及缓冲区大小对应的电路最大吞吐量

指标		距离(km)/RTT(s)						
		100/2	500/10	1 000/ 20	2 500/ 50	4 000/ 80	5 000/ 100	10 000/ 200
最大 吞吐量/ (Gbit /s)	缓冲区为 0.01 GB	40	8	4	1.6	1	0.8	0.4
	缓冲区为 0.1 GB	400	80	40	16	10	8	4
	缓冲区为 0.2 GB	800	160	80	32	20	16	8
	缓冲区为 0.5 GB	2 000	400	200	80	50	40	20
	缓冲区为 1 GB	4 000	800	400	160	100	80	40

以采用100G带宽的电路为例,在某个取定的缓冲区数值和传输距离条件下,若链路的最大吞吐量小于100G,则无法充分利用100G的带宽,造成带宽资源的浪费。

### 3 提升链路通量的技术措施

为了实现IP网海量数据的远距离高效传输,可以采用优化传输协议、新型传输协议、降低网络时延等措施,并可将这些措施组合使用以进一步提升链路的吞吐量。

#### 3.1 优化TCP传输协议

优化TCP传输协议包括扩展窗口大小和优化拥塞控制算法2个方面。

##### 3.1.1 扩展TCP窗口大小

TCP窗口扩展在RFC 1072中被引入,并在RFC 1323中进行了改进。TCP窗口扩展将16位窗口字段扩展为32位长度。解决方案是定义TCP的“选项”字段以指定16位窗口大小位移数,TCP标头字段应按位移数移位以产生更大的窗口值。

在RFC 1323中,TCP报文段的“选项”字段关于“窗口扩大选项TCP Window Scale Option (WSopt)”的格式及取值定义如表3所示<sup>[5]</sup>。

表3 TCP WSopt的格式及取值

格式	种类 Kind=3 (二进制为0000011)	长度 Length=3 (二进制为0000011)	扩大因子 shift.cnt
取值/B	1	1	1

要启用窗口扩大选项,通信双方必须在各自的SYN报文中发送这个选项。

在启用窗口扩大选项的情况下,要将TCP首部中的16 bit窗口值移shift.cnt位,以获得实际的通告窗口大小。新的窗口值=首部中定义的窗口值×2<sup>窗口扩大因子</sup>。

shift.cnt的取值范围为0~14,即最大TCP序号限定为2<sup>16</sup> × 2<sup>14</sup> = 2<sup>30</sup>。通过使用窗口扩大因子,窗口大小最大可以扩展到2<sup>30</sup> B,即1 GB。

窗口扩大选项对于提升高带宽时延积链路的数据传输效率非常有效。对于一个RTT为200 ms的广域网络,若窗口大小为64 KB,此链路最大吞吐量的理论值仅为2.56 Mbit/s,若窗口大小扩展为1 GB,链路最大吞吐量的理论值将可以达到40 Gbit/s。

##### 3.1.2 优化TCP拥塞协议

针对TCP在高带宽时延积网络中存在的问题,目前国内外已经提出了一些拥塞控制优化算法,这些算法通过优化拥塞控制机制,缓解了TCP在高带宽时延积网络中链路利用率不高的问题,在一定程度上提高了TCP的性能。

针对TCP基础拥塞控制技术的优化协议有很多,一些典型的高速TCP改进协议如下。

- a) 基于分组丢失反馈的改进协议:HSTCP、STCP、BIC、CUBIC。
- b) 基于时延反馈的改进协议:Hybrid Slow Start TCP、Vegas。
- c) 基于分组丢失和时延反馈的混合反馈改进协议:CTCP、Africa、YeAH、Illinois。
- d) 基于可用带宽测量的改进协议:Westwood、Fusion、ARENO。
- e) 基于显式反馈的改进协议:XCP、VCP、EVLFTCP、JetMax。

拥塞控制主要是靠发送端维护拥塞窗口cwnd和慢开始门限ssthresh这2个变量实现,各种拥塞控制协议优化算法的本质都是在这2个变量的初始值和如何调整上进行研究。

TCP New Reno是对TCP Reno中快速恢复阶段的重传进行改善的一种改进算法,New Reno在网络低错

误率时和选择确认 SACK 相当,在高错误率时运行效率优于 Reno。

TCP BIC 旨在优化高速高延迟网络的拥塞控制,其拥塞窗口算法使用二分搜索算法,以尝试找到能长时间保持的拥塞窗口最大值。

CUBIC 则是比 BIC 更温和和系统化的分支版本,它使用三次函数代替二分算法作为其拥塞窗口算法,并且使用函数拐点作为拥塞窗口的设置值。

TCP PRR 旨在恢复期间提高发送数据的准确性,该算法确保恢复后的拥塞窗口大小尽可能接近慢开始阈值。

TCP BBR 基于模型主动探测,而以往大部分拥塞算法是基于丢包来作为降低传输速率的信号。

TCP Vegas 算法和其他拥塞控制算法的不同之处在于 Vegas 算法并不急于以丢包来判断是否发生了拥塞,而是通过数据包延迟来判断。Vegas 通过 RTT 来决定增加或者减小拥塞窗口,它能够在拥塞将要发生时就避免拥塞,而不是等到拥塞已经发生后再减小发送速度,因此能够减小重传和超时的概率。

TCP Westwood 改良自 New Reno,不同于以往其他拥塞控制算法使用丢失来测量,其通过测量确认包来确定一个“合适的发送速度”,并以此调整拥塞窗口和慢开始阈值。Westwood 将慢开始阶段算法改良为“敏捷探测(Agile Probing)”,并且设计了一种持续探测拥塞窗口的方法来控制进入“敏捷探测”,使连接尽可能地使用更多的带宽。Westwood+算法使用更长的带宽估计间隔和优化的滤波器来修正 Westwood 对 ACK 压缩场景的带宽估计过高的问题。

CTCP(复合 TCP)是微软实现的一种 TCP 拥塞控制算法,该算法通过同时维护 2 个拥塞窗口,来实现在长肥网络中有较好的性能而又不损失公平性。CTCP 维护 2 个拥塞窗口分别为常规的 AIMD 窗口以及基于延迟的窗口,最终实际使用的滑动窗口大小是这 2 个窗口的和。AIMD 窗口与 Reno 的增加方式相同;如果延迟小,基于延迟的窗口将迅速增加以提高网络的利用率。一旦经历了排队,延迟窗口将逐渐减小以补偿增加的 AIMD 窗口。这样的目的是保持两者的总和大致恒定<sup>[6]</sup>。

STCP 算法是在传统 TCP 基础上进行的改进算法。与传统 TCP 算法不同,STCP 采用的是乘性增加、乘性减少(MIMD)策略,相比加性增加、乘性减少(AIMD)策略,该策略的窗口增加更快减少更慢。

采用拥塞协议优化算法可以更加合理地调整高带宽时延积链路中的拥塞窗口的幅度,窗口恢复得更快,使网络通量的曲线波动不再剧烈,从而提高链路的通量。

## 3.2 其他措施

### 3.2.1 采用 RDMA 等技术替代 TCP

TCP 协议在设计之初并未考虑大象流传输及长肥网络的特殊需求,因此面临一系列局限性。TCP 协议的复杂性导致其延迟相对较高,流控及拥塞控制机制较为保守,带宽利用率低下。此外,TCP 协议包头开销较大,这进一步降低了传输效率。尽管存在各种 TCP 优化协议,但它们仅能在一定程度上改善其性能,无法从根本上解决这些固有的问题。

针对具有高带宽需求、高并发特性、海量数据处理能力及严格要求延迟的应用场景,采用 RDMA 等技术替代 TCP 成为了一种优选方案。相较于 TCP,RDMA 凭借其高吞吐量和超低延迟特性,能显著提升网络通量。

由于 RDMA 技术对丢包等差错敏感,通信节点间传输电路需要有较高的传输质量,否则将影响数据传输的效果。

### 3.2.2 优选传输路径

收发节点通过广域网进行互联,会存在多种可能的传输路径,选择时延最小的路径可得到更高的网络通量。

在电信运营商的广域网络中,由运营商通过 SRv6、SDN、切片技术为远程数据传输规划最佳候选路径。

SRv6 Policy 利用 Segment Routing 的源路由机制,可以实现业务的端到端需求,是实现 SRv6 网络编程的主要机制。通过在 SRH 中封装一系列的 SRv6 Segment ID,可以显式地指导报文按照规划的路径进行转发,实现对转发路径端到端的细粒度控制,满足业务的高可靠、大带宽、低时延等 SLA 需求。

结合网络切片技术,可以为大象流传输提供专享的通道与带宽,避免与其他业务流相互影响,减少网络拥塞,提高传输效率。

为了识别出特定的用户数据流,可采用 APN 来进行标识。APN ID 由运营商网络侧边缘设备生成,并为用户的流量添加应用标记<sup>[7]</sup>。

SRv6/SDN 技术规划最佳候选路径方案如图 2 所示。

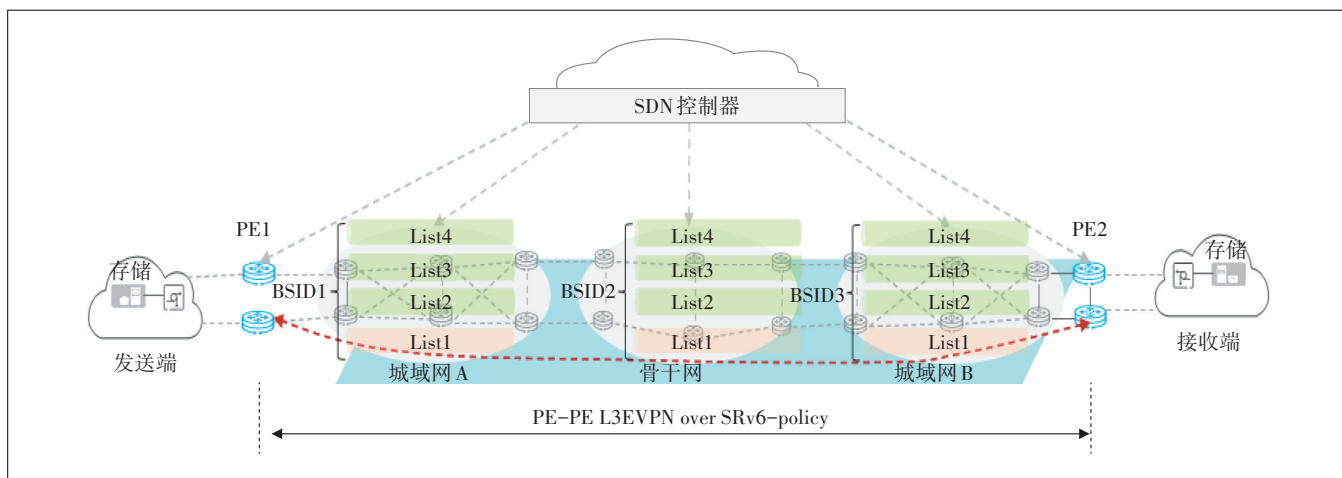


图2 SRv6/SDN 技术规划最佳候选路径方案

SDN 控制器为用户在收发两端的 PE 设备间开通 L3 EVPN OVER SRv6-Policy, 在每个网络域内根据每条链路上可使用的带宽情况, 针对一个 BSID 生成多条 List。发送端根据五元组识别出数据流, 打 APN6 标签, 并完成 APN6-ID 与隧道的映射关系。在每个域内, 根据流量情况, 每个用户的流量在多条 List 之间进行负载分担, 实现所有 List 的流量吞吐量最大。

### 3.2.3 使用确定性网络

当网络时延抖动范围显著时, TCP 协议会频繁地调整窗口大小及门限值, 然后再经历一个逐步趋于稳定的恢复过程, 这直接导致网络吞吐量持续波动。同时, 丢包率的上升也会引发数据的反复重传, 进而降低数据传输效率。因此, 海量数据的传输要求网络具有稳定和可靠的状态。

确定性网络是相对于传统的尽力而为网络而言的。尽力而为网络的问题根源是数据传输的稳定性不够, 比如带宽时高时低, 时延时大时小。互联网是最典型的尽力而为型的网络, 它可以满足很多生活类的应用, 这类应用对传输的确定性要求不高。

确定性网络可以提供确定性的服务质量 (Quality of Service, QoS), 其中 5 种典型的确定性 QoS 包括低时延 (上限确定)、低抖动 (上限确定)、低丢包率 (上限确定)、高带宽 (上下限确定) 和高可靠 (下限确定)<sup>[8]</sup>。典型的确定性网络技术较多, 广域 IP 网络相关的确定性网络包括灵活以太网 (FlexE)、确定网 (DetNet) 以及确定性 IP (DIP) 网络, 目前国内已有这些技术的实验网络或商用网络。

在确定性网络中进行海量数据的远距离传输, 可

以减小网络通量的波动, 由高确定性的网络质量来维持网络的高吞吐量。

### 3.2.4 采用高性能的存储设备及存储网络

在数据迁移过程中, 收发端对存储设备的读写操作也占据着重要地位。即便数据传输网络具备高吞吐量, 若搭配低速存储设备, 整体的传输效率仍会大打折扣。鉴于此, 采用高性能的存储设备及存储网络尤为重要。

随着 NAND Flash 技术的进步、NVMe 协议的迭代, SSD 及采用 SSD 的存储设备已逐渐成为企业数据中心应用的主流。

早期高性能存储网络多采用 FC 协议, 主要应用于存储局域网 (SAN), 但 FC 网络最大带宽只有 32G, 已跟不上 NVMe 设备带宽的迭代提升速度, 满足不了业务发展需求。2016 年标准化组织推出了 NVMe-oF (NVMe over Fabric, NoF), NoF 存储网络应运而生。

NVMe-oF 集成了现有的 NVMe 和高速低延迟传输网络的技术, 可极大地释放数据中心端到端 NVMe 性能。NoF 是在不同种类网络中传输存储协议的技术路线总称, NVMe over Fabric 中的“Fabric”是 NVMe 的承载网络, NoF 可以分为在 FC 网络上传输的 NVMe over FC、在 IP 网络上基于 TCP 协议传输的 NVMe over TCP 以及基于 RDMA 技术的 NVMe over RDMA。

对于 NVMe over RDMA, RDMA 是承载 NoF 的原生网络协议, 是一种无损网络传输技术。当前 RDMA 技术的实现方式主要有 InfiniBand、RoCE、iWARP 这 3 种, NVMe over RDMA 可利用 RDMA 的技术优势提供高效的节点间网络通信。

NVMe-oF 不仅可以应用于 DC 内部的数据存储,还可以实现多 DC 间存储网络互通。图 3 为一个采用

NVMe over RoCE 进行 DC 互联的存储网络架构<sup>[9]</sup>。

在图 3 所示的存储网络中,服务器不仅能从 DC 内

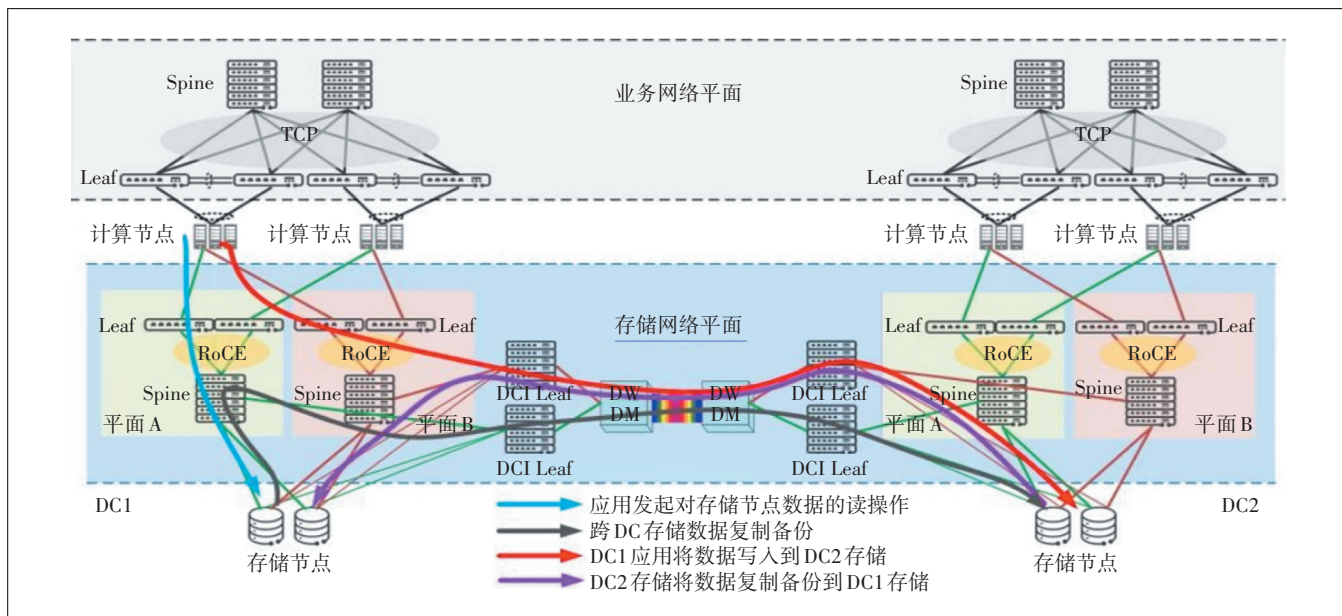


图 3 NVMe over RoCE 多 DC 互联存储网络架构

的本地存储节点读取数据并写入至远端 DC 的存储节点,还能直接将 DC 存储的数据复制到远端的 DC 存储设备,无需服务器的数据传输介入,从而简化了处理流程。得益于 DC 节点间采用的 RDMA 技术,该网络能够确保高效的数据传输效率。

#### 4 结束语

TCP 协议在应对海量数据的远距离高效传输方面存在局限,亟需对流量控制与拥塞控制算法进行优化升级,或者探索全新的网络协议作为替代方案。当前,业界正积极研发更高效的高通量网络传输协议,并已取得一系列技术突破,诸如特斯拉 TTPoE 协议以及国内自主研发的高通量以太网 ETH+ 协议等。这些高性能网络协议旨在构建高通量数据网络,以契合高性能计算、海量数据传输等特殊应用场景的需求。

#### 参考文献:

- [1] IETF. Transmission Control Protocol: RFC 793 [S/OL]. [2025-08-26]. <https://datatracker.ietf.org/doc/rfc793/>.
- [2] FALL K R, STEVENS W R. TCP/IP 详解 卷 1: 协议 [M]. 吴英, 张玉, 许昱玮, 译. 2 版. 北京: 机械工业出版社, 2016.
- [3] ALLMAN M, PAXSON V, STEVENS W. TCP congestion control: RFC 2581 [S/OL]. [2025-08-26]. <https://www.rfc-editor.org/rfc/rfc2581>.

rfe2581.

- [4] FOROUZAN B A, FEGAN S C. TCP/IP 协议族 [M]. 谢希仁, 译. 3 版. 北京: 清华大学出版社, 2006.
- [5] JACOBSON V, BRADEN R, BORMAN D. TCP extensions for high performance; RFC 1323 [S/OL]. [2025-08-26]. <https://www.rfc-editor.org/rfc/rfc1323.html>.
- [6] 佚名. 各种 TCP 拥塞控制算法 [EB/OL]. [2025-08-26]. <https://zhuanlan.zhihu.com/p/544139753>.
- [7] APN6. 什么是 APN6? [EB/OL]. [2025-08-26]. <https://info.support.huawei.com/info-finder/encyclopedia/zh/APN6.html>.
- [8] 第五届未来网络发展大会组委会. 未来网络白皮书——确定性网络技术体系白皮书 (2021 版) [EB/OL]. [2025-08-26]. <http://www-file.huawei.com/-/media/corporate/pdf/news/future-network-whitepaper.pdf?la=zh>.
- [9] 华为技术有限公司. (eBook) NoF+ 存储网络解决方案 [EB/OL]. [2025-08-26]. <https://support.huawei.com/enterprise/zh/doc/EDOC1100211900>.

#### 作者简介:

张锴, 高级工程师, 硕士, 主要从事数据通信及数据中心网络相关咨询设计工作; 晏家豪, 高级工程师, 学士, 主要从事 IP 城域网和 IDC 网络相关咨询设计工作; 洪伟, 高级工程师, 学士, 主要从事 IP 城域网和 IDC 网络相关咨询设计工作。