



ZOMI

高速互联：RDMA 基本概述



Content

1. AI时代为什么需要RDMA

2. RDMA的核心原理

3. RDMA与传统网络对比

4. RDMA支持的三种主流协议：

- InfiniBand
- RoCE (RDMA over Converged Ethernet)
- iWARP



01

为什么需要RDMA



AI时代为什么需要RDMA

- GPU算力增长远超网络带宽提升
 - AI训练中万亿参数模型（如GPT-4）需数百GPU协同，单卡通信带宽需求超200Gbps。
 - 计算/通信比失衡：在8卡集群中，通信耗时占比高达90%（单卡仅10%），**网络成为系统瓶颈。**
- 参数同步的实时性要求
 - 梯度同步需在毫秒级完成，传统TCP/IP延迟（50—100 μ s）导致GPU长时间闲置。**RDMA将延迟降至1—5 μ s，提升集群利用率30%以上。**
 - Meta实测：AI工作负载33%时间浪费在网络等待，RDMA显著缩短作业完成时间（JCT）
- 传统TCP/IP通信在AI时代的瓶颈：
 - 多次数据拷贝：数据需在应用内存与内核缓冲区间复制，消耗CPU资源；
 - 高延迟：内核协议栈处理增加微秒级延迟；
 - CPU利用率高：大规模数据传输时CPU成为瓶颈



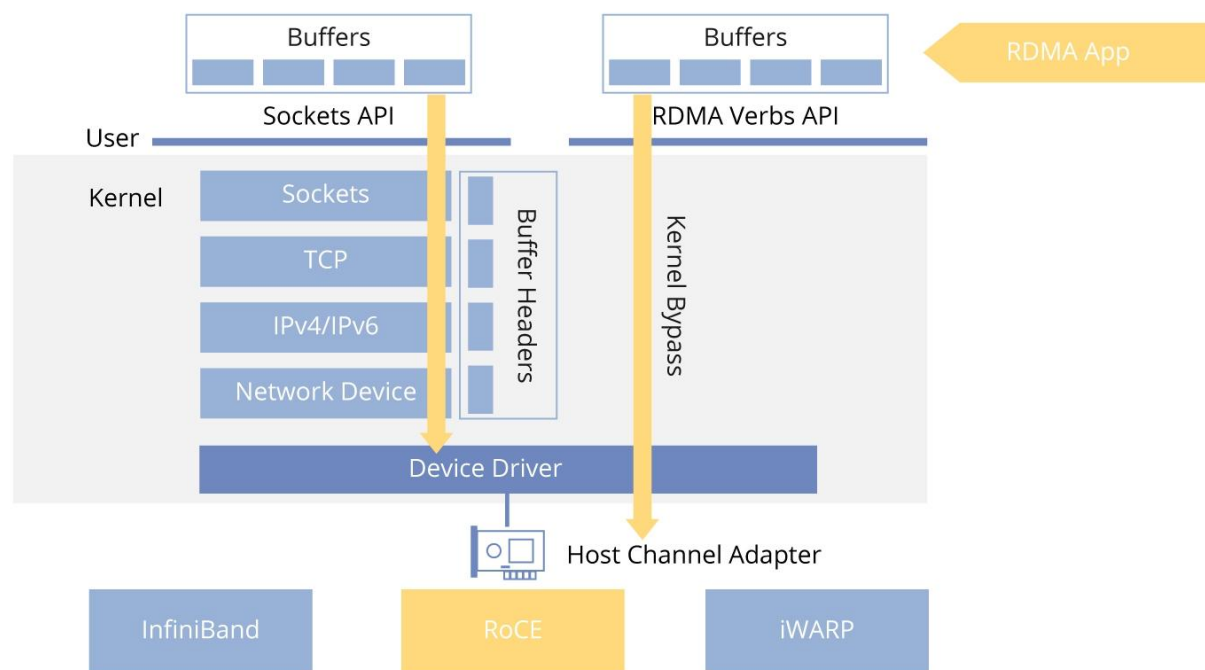
02

RDMA的核心原理



什么是RDMA

- RDMA (Remote Direct Memory Access, 远程直接内存访问) 是一种高性能网络通信技术, 允许计算机直接通过网络访问另一台计算机的内存, **无需操作系统内核和CPU介入**。
- 其目标是实现**零拷贝、零中断、低延迟**的数据传输, 它解决了传统网络通信中的性能瓶颈问题, 尤其在数据中心、高性能计算 (HPC) 和AI等场景中至关重要。



RDMA核心原理

RDMA的核心是通过网卡硬件（如支持RDMA的RNIC）实现跨节点的**直接内存访问**，其核心特性包括

- **零拷贝（Zero-Copy）**

- 数据直接从发送端应用内存写入接收端应用内存，无需经过内核缓冲区、套接字缓存等中间复制环节。
- 传统通信需多次内存拷贝（应用层→内核→网卡），而RDMA通过网卡硬件直接读写用户内存，减少数据搬运开销。

- **内核旁路（Kernel Bypass）**

- 应用程序通过用户态驱动直接操作RDMA网卡（RNIC），无需内核协议栈处理数据。
- 优势：避免上下文切换、系统调用等CPU开销，降低延迟至微秒级（传统TCP/IP为毫秒级）。

- **CPU卸载（CPU Offload）**

- 数据传输协议（如分段、校验和、流量控制）由智能网卡硬件处理，CPU仅负责发起传输指令，不参与数据搬运。
- 效果：释放CPU资源，使其专注于计算任务，提升系统整体效率。



03

RDMA与 传统网络对比



RDMA与传统网络对比

对比维度	RDMA	传统网络（TCP/IP）
传输机制	网卡直接读写远程内存（零拷贝）	数据需经内核协议栈封装/解析
性能表现	微秒级（通常1—5μs）	毫秒级（10—100μs）
CPU负载	极低（<5%）	高（30%—100%，随带宽提升）
内存拷贝次数	零拷贝	4—6次（用户态↔内核态↔网卡）
典型应用场景	HPC、分布式存储（Ceph）、AI训练、实时交易	Web服务、文件传输、普通计算
部署成本与复杂度	较高（需RDMA网卡+无损网络）	低（通用网卡+标准以太网）



04

RDMA

支持的三种主流协议



RDMA支持的三种主流协议

RDMA（远程直接内存访问）技术支持的三种主要协议为 InfiniBand（IB）、RoCE（RDMA over Converged Ethernet）和 iWARP（Internet Wide Area RDMA Protocol）

特性	InfiniBand (IB)	RoCE	iWARP
协议类型	原生RDMA专用协议	基于以太网的RDMA	基于TCP/IP的RDMA
网络要求	专用IB交换机、网卡和线缆	支持RoCE的网卡+无损以太网交换机	支持iWARP的网卡+标准以太网交换机
适用场景	超算中心、AI万卡集群（如NVIDIA DGX）	云数据中心、分布式存储、AI训练	广域网跨数据中心通信



InfiniBand (IB)

- 技术架构：
 - 专为RDMA设计的端到端协议，涵盖物理层到传输层。
 - 通过专用网卡（HCA）和交换机实现内存直接访问，完全绕过操作系统内核，支持零拷贝和内核旁路。
- 性能优势：
 - 超低延迟（可低于 $1\mu\text{s}$ ），高带宽（HDR InfiniBand达400Gbps）
 - 原生拥塞控制机制（基于信用算法），保证无丢包传输。
- 局限性：
 - 生态封闭：需全套专用设备（如NVIDIA Mellanox方案），成本高昂；
 - 跨数据中心扩展困难，主要限于本地集群。



RoCE (RDMA over Converged Ethernet)

- 技术演进

- RoCEv1: 基于以太网链路层（二层），仅支持同子网通信，通过Ethertype 0x8915标识报文；
- RoCEv2（主流）：基于UDP/IP协议栈，支持三层路由（UDP端口4791），可通过ECMP实现负载均衡。

- 关键依赖：无损以太网，需启用以下技术避免丢包

- PFC（基于优先级的流量控制）：为RDMA流量预留专用队列，触发反压机制；
- ECN（显式拥塞通知）：通过IP头标记拥塞，动态降速。

- 性能表现

- 延迟可低至2 μ s（25GbE下比TCP带宽提升30%）；
- 成本显著低于IB，兼容现有以太网基础设施。



- **技术原理**

- 将RDMA语义封装在TCP/IP协议栈中，通过网卡硬件卸载TCP处理（TOE技术），减少CPU开销。

- **优势与局限**

- 支持标准以太网交换机和广域网路由，扩展性最佳；
- TCP协议栈引入额外延迟（序列化/重传机制），性能弱于IB/RoCE；
- 大量TCP连接占用内存资源，大规模集群效率低。

- **应用场景：**

- 跨数据中心通信、金融行业广域网低延迟交易（如Chelsio方案）





Thank you

把 Allinfra 带入每个开发者、每个家庭、
每个组织，构建万物互联的智能世界

Bring AI Infra to every person, home and
organization for a fully connected,
intelligent world.

Copyright © 2025 [Infrasys-AI](#) org. All Rights Reserved.

The information in this document may contain predictive statements including, without limitation, statements regarding the future financial and operating results, future product portfolio, new technology, etc. There are a number of factors that could cause actual results and developments to differ materially from those expressed or implied in the predictive statements. Therefore, such information is provided for reference purpose only and constitutes neither an offer nor an acceptance. [Infrasys-AI org.](#) may change the information at any time without notice.



GitHub github.com/Infrasys-AI/Allinfra

Book infrasys-ai.github.io



引用与参考

1. <https://community.fs.com/blog/remote-direct-memory-access-rdma.html>
2. <https://cloud.tencent.com/developer/article/2508336>
3. <https://doc.mbalib.com/view/903fb2f7b2b20781e8423af9dfa7829e.html>
4. https://blog.51cto.com/u_16213715/13786410
5. <https://blog.csdn.net/gs80140/article/details/145033932>
6. <https://blog.csdn.net/zuopiezia/article/details/144340106>
7. https://mp.weixin.qq.com/s?__biz=MzA3NTY1NjAyMw==&mid=2247486165&idx=1&sn=594c9eba7088a396134e02017d3f61eb&chksm=9e3422c20f13945922c45958caf838dab92f230e78da1d29f6f1d7e4f3624e207173395bd60d#rd
8. https://mp.weixin.qq.com/s?__biz=MzkwOTYwMDYxMQ==&mid=2247504048&idx=1&sn=a44986e6a3570a76a73d773585a1a2dd&chksm=c04e027f60a11f22cbcacaca65d95b7163e439056b77d4e9e163afb69c352992e3517af75cff#rd
9. https://mp.weixin.qq.com/s?__biz=MzA3NTY1NjAyMw==&mid=2247486165&idx=1&sn=594c9eba7088a396134e02017d3f61eb&chksm=9e3422c20f13945922c45958caf838dab92f230e78da1d29f6f1d7e4f3624e207173395bd60d#rd

