# RDMA基础和UDISK实践

# Agenda

- RDMA概念
- RDMA编程接口
- RDMA在UDisk的实践

# RDMA – what is it?

❖**R**emote
- – data transfers between nodes in a network

❖**D**irect
- – no Operating System Kernel involvement in transfers
- – everything about a transfer offloaded onto Interface Card

❖**M**emory
- – transfers between user space application virtual memory
- – no extra copying or buffering

❖**A**ccess
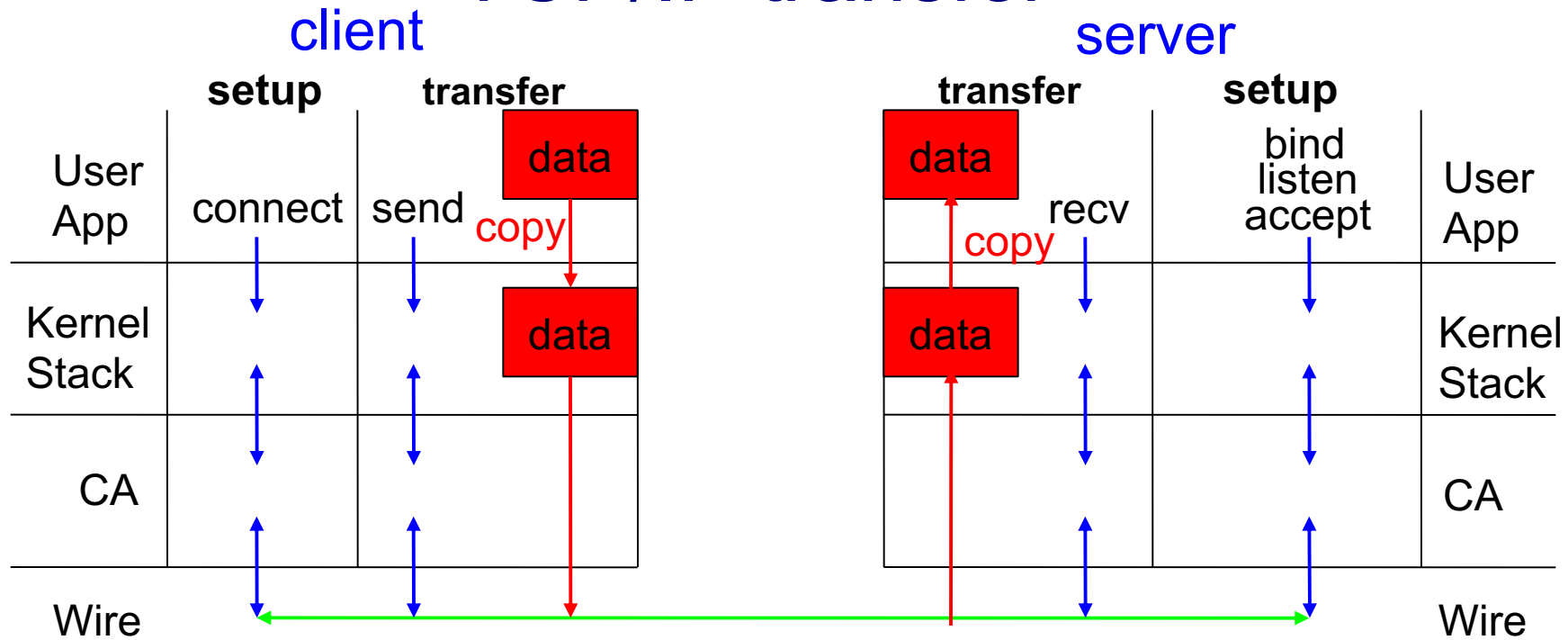- – send, receive, read, write, atomic operations

# RDMA Benefits

❖High throughput

❖Low latency

❖High messaging rate

❖Low CPU utilization

❖Message boundaries preserved

❖Asynchronous operation

# How RDMA differs from TCP/IP

❖ "zero copy" – data transferred directly from virtual memory on one node to virtual memory on another node

❖ "kernel bypass" – no operating system involvement during data transfers

❖ asynchronous operation – threads not blocked during I/O transfers
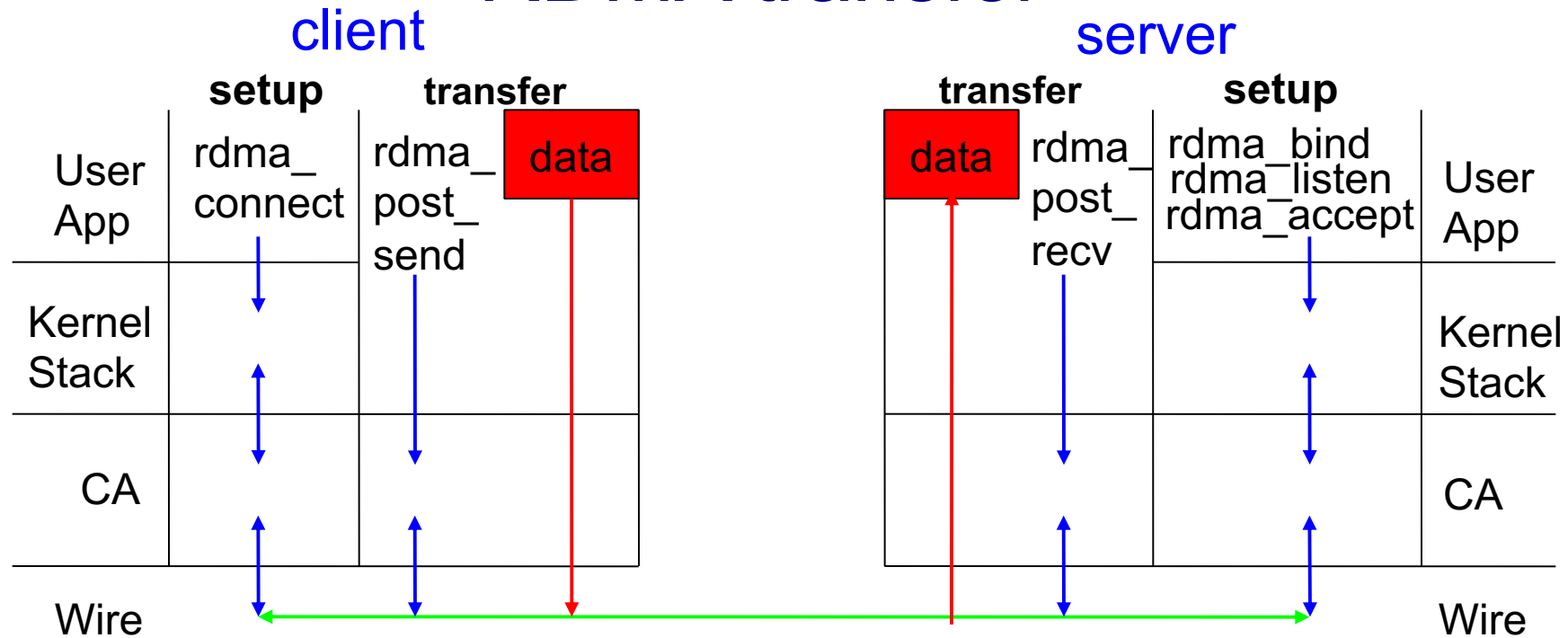
# TCP/IP transfer



blue lines: control information

red lines: user data
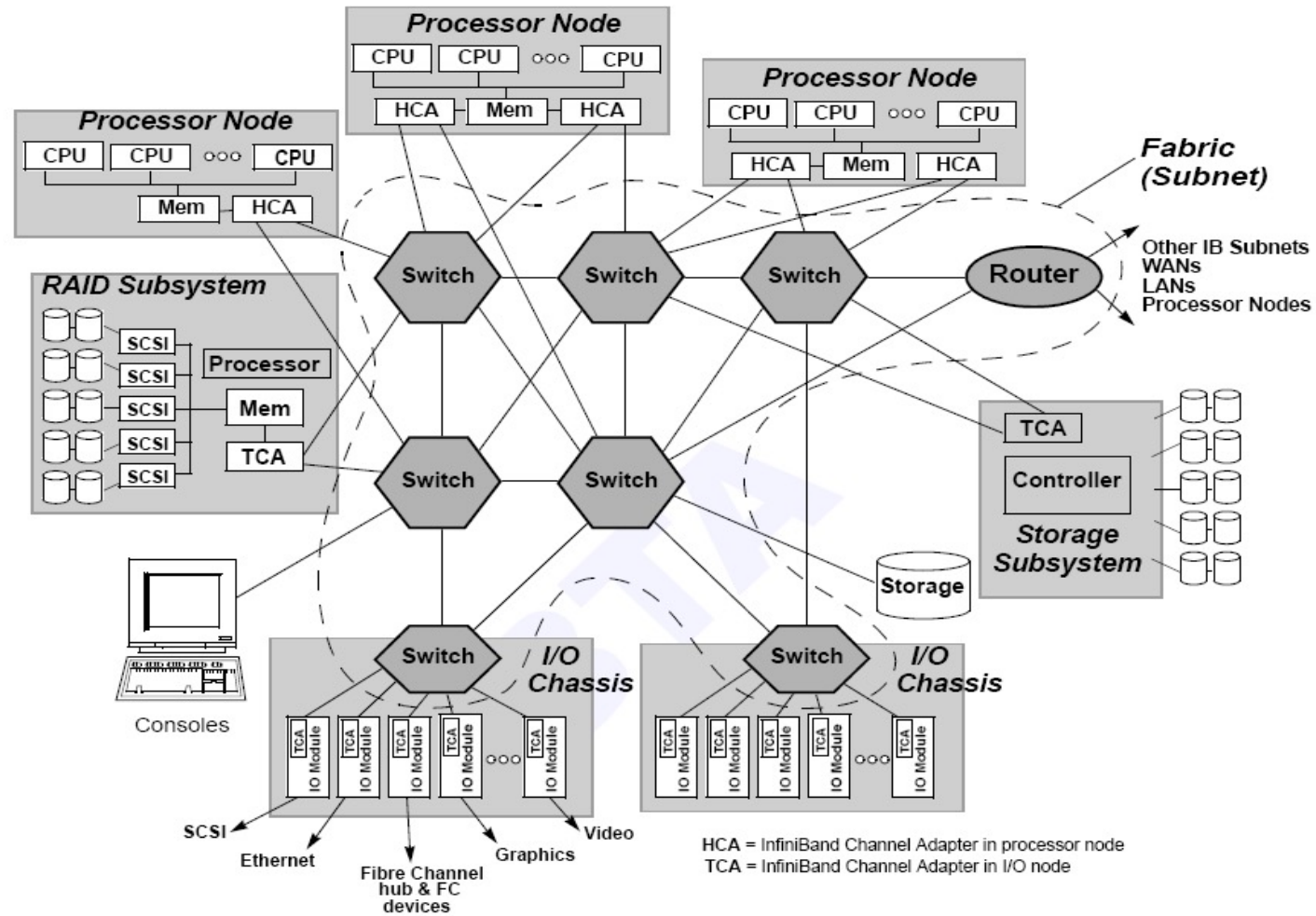
green lines: control and data

# RDMA transfer



blue lines: control information

red lines: user data

green lines: control and data

# Infiniband architecture overview

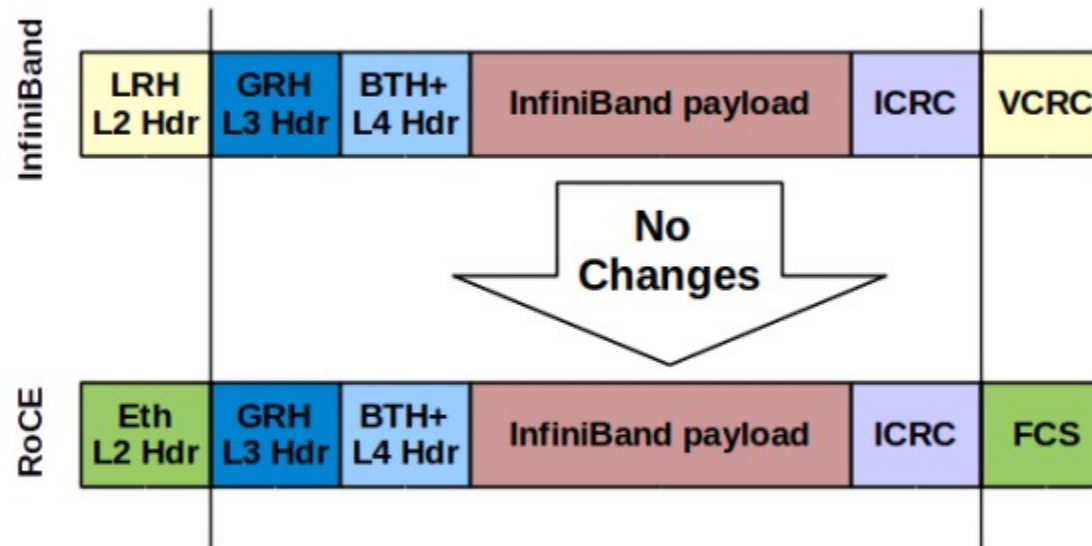# RoCE – RDMA over Converged Ethernet



Figure 13 – Differences between IB and RoCE frames [20]
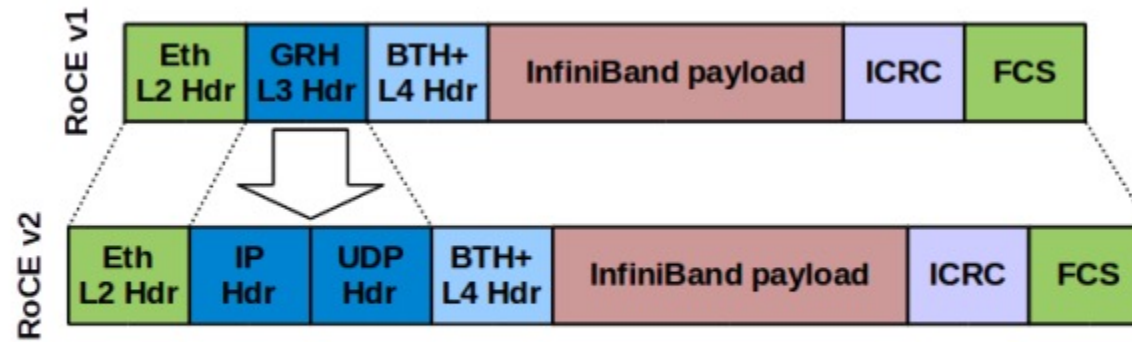
# RoCEv2 – Based on UDP



Figure 14 – Differences between RoCE v1 and v2 frames [21]

# RoCE packet

```
> Frame 46: 1094 bytes on wire (8752 bits), 1094 bytes captured (8752 bits)
> Ethernet II, Src: Mellanox_ce:ec:7b (ec:0d:9a:ce:ec:7b), Dst: Mellanox_34:0a:be (ec:0d:9a:34:0a:be)
    > Destination: Mellanox_34:0a:be (ec:0d:9a:34:0a:be)
    > Source: Mellanox_ce:ec:7b (ec:0d:9a:ce:ec:7b)
      Type: RDMA over Converged Ethernet (0x8915)
> InfiniBand
  > Global Route Header
        0110 .... = IP Version: 6
        .... 0000 0010 .... = Traffic Class: 2
        .... .... .... 0000 0000 0000 0000 0000 = Flow Label: 0
        Payload Length: 1040
        Next Header: 27
        Hop Limit: 255
        Source GID: fe80::ee0d:9aff:fece:ec7b
        Destination GID: fe80::ee0d:9aff:fe34:abe
  > Base Transport Header
        Opcode: Reliable Connection (RC) - SEND Middle (1)
        0... .... = Solicited Event: False
        .1.. .... = MigReq: True
        ..00 .... = Pad Count: 0
        .... 0000 = Header Version: 0
        Partition Key: 65535
        Reserved: 00
        Destination Queue Pair: 0x00020f
        0... .... = Acknowledge Request: False
        .000 0000 = Reserved (7 bits): 0
        Packet Sequence Number: 2
      Invariant CRC: 0x389f2670
```

# RoCEv2 packet

```
> Frame 47: 174 bytes on wire (1392 bits), 174 bytes captured (1392 bits)
∨ Ethernet II, Src: Mellanox_ce:ec:7b (ec:0d:9a:ce:ec:7b), Dst: Mellanox_34:0a:be (ec:0d:9a:34:0a:be)
  > Destination: Mellanox_34:0a:be (ec:0d:9a:34:0a:be)
  > Source: Mellanox_ce:ec:7b (ec:0d:9a:ce:ec:7b)
    Type: IPv4 (0x0800)
∨ Internet Protocol Version 4, Src: 2.2.2.1, Dst: 2.2.2.2
    0100 .... = Version: 4
    .... 0101 = Header Length: 20 bytes (5)
  > Differentiated Services Field: 0x5a (DSCP: AF23, ECN: ECT(0))
    Total Length: 160
    Identification: 0x103e (4158)
  > Flags: 0x4000, Don't fragment
    Time to live: 64
    Protocol: UDP (17)
    Header checksum: 0x21af [validation disabled]
    [Header checksum status: Unverified]
    Source: 2.2.2.1
    Destination: 2.2.2.2
∨ User Datagram Protocol, Src Port: 52067, Dst Port: 4791
    Source Port: 52067
    Destination Port: 4791
    Length: 140
    [Checksum: [missing]]
    [Checksum Status: Not present]
    [Stream index: 3]
∨ InfiniBand
  > Base Transport Header
    Invariant CRC: 0x9ba57e2f
∨ Data (116 bytes)
    Data: 58494f50000000001400000000000000640000000000001000...
    [Length: 116]
```
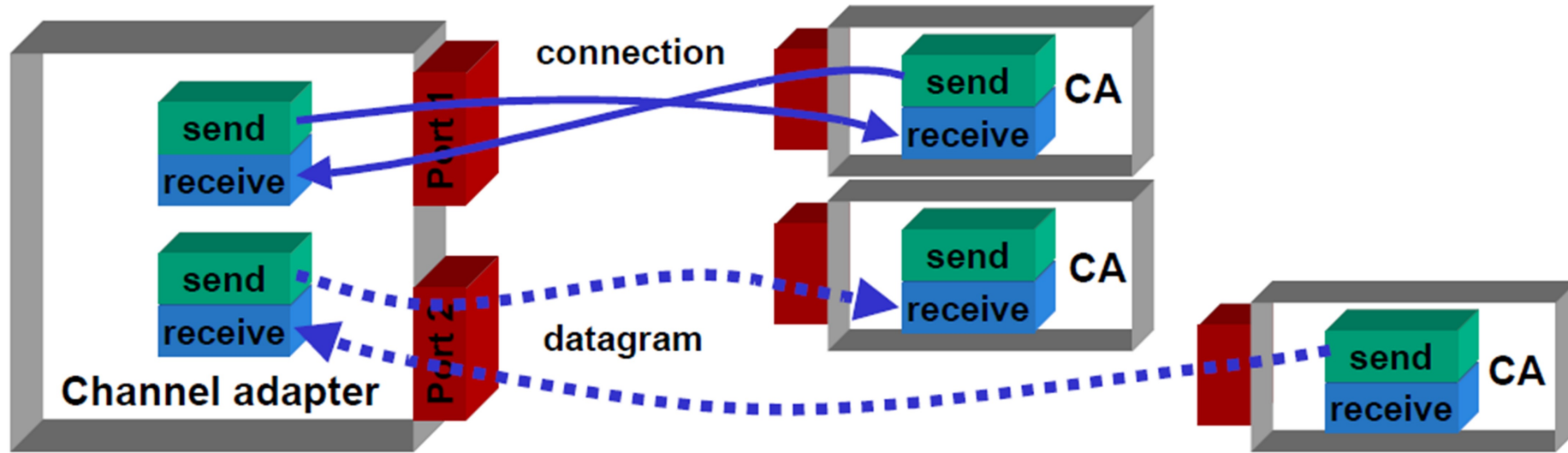
# 4K ping-pong test with 25Gbps Card

- Latency 8us –depth=1
- IOPS 700K –depth=128

# 相关名词

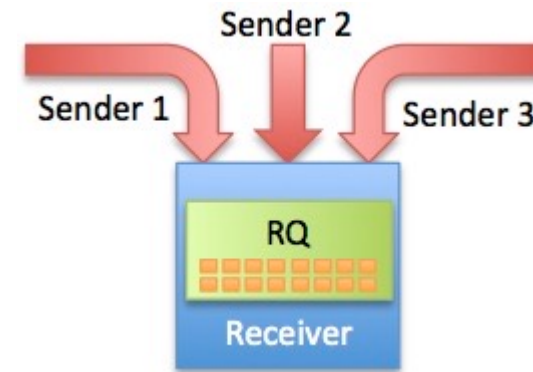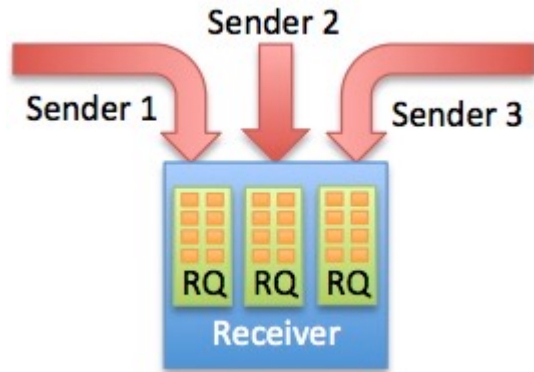| | | |
|---|---|---|
| PD | Protection Domain | Glues queue pairs and memory regions |
| MR | Memory Region | Registered memory region that HCA can read from or write to. Contains R_Key and L_Key |
| QP | Queue Pair | Send / Receive work queue. Send or receive work requests are placed onto a queue pair |
| CQ | Completion Queue | Completion Queue. Completed work requests, so called work completions are placed onto a completion queue. Is associated with queue pair. |
| WR | Work Request | Either send or receive work request. Specifies action to be processed and will be put onto send or receive queue (QP). References scatter/gather element |
| SGE | Scatter/Gather Element | Defines address(es) in memory to read from or to write to. Must be given L_Key or R_Key to authenticate access to memory region |
| WC | Work Completion | After a work request has been completed the work completion delivers result |

# QueuePair

# Send/Recv Queue/Completion Queue

# Steps

- Create QueuePair
- Exchange QueuePair Info between Nodes
- Modify QueuePair state to RTS （via TCP）
- Poll Complete Queue

# RDMA in UDisk

- RC QueuePair
- RoCEv2
- Share CQ – per thread/core
- SRQ – per thread/core

# RDMA in UDisk

- Keep Alive
  - ➢ Do RDMA WRITE with 0 size data

- Large Packet
  - ➢ Split into small packets
  - ➢ Do RDMA READ – difficult when there is a failure

- Memory recycling on failure
  - ➢ Handle different kinds of WC error
  - ➢ Move qp to error state and wait all WC is done before destroy it

- Integration with Spdk
  - ➢ Create Memory Pool From spdk_dma_zmalloc

# Udisk Performance –Single Chunk/Core

| RandRead性能数据 | | | | | |
|---|---|---|---|---|---|
| 队列 | IOPS | 平均延迟 | 95% | 99.95% | 99.99% |
| 1 | 13.0k | 71.16 | 110 | 163 | 165 |
| 8 | 108k | 73.71 | 100 | 126 | 149 |
| 16 | 204k | 78.09 | 97 | 118 | 131 |
| 32 | 364k | 87.72 | 102 | 119 | 139 |
| 64 | 566k | 112.64 | 126 | 178 | 198 |
| 128 | 639k | 199.97 | 239 | 343 | 371 |

# Udisk Performance –Cluster

16*24 RandWrite
Lat 358us
IOPS 1252k