

# RDMA Programming - Mellanox

Saturday, December 23, 2017 12:11 PM

## Infiniband - IBTA

Read Chapter - IBTA Specs ( 1 & 4)

MPI Programming

Socket Programming

## Infiniband Advantage

1. Easy Scaling
2. Low latency
3. High Throughput

Using RDMA has the following major advantages:

- Zero-copy - applications can perform data transfer without the network software stack involvement and data is being send received directly to the buffers without being copied between the network layers.
- Kernel bypass - applications can perform data transfer directly from userspace without the need to perform context switches.
- No CPU involvement - applications can access remote memory without consuming any CPU in the remote machine. The remote memory machine will be read without any intervention of remote process (or processor). The caches in the remote CPU(s) won't be filled with the accessed memory content.
- Message based transactions - the data is handled as discrete messages and not as a stream, which eliminates the need of the application to separate the stream into different messages/transactions.
- Scatter/gather entries support - RDMA supports natively working with multiple scatter/gather entries i.e. reading multiple memory buffers and sending them as one stream or getting one stream and writing it to multiple memory buffers

From <<http://www.rdmamojo.com/2014/03/31/remote-direct-memory-access-rdma/>>

A single lane is a serial link operating at one of five data rates:

- Single Data Rate (SDR) - **2.5Gb/s**
- Double Data Rate (DDR) - **5Gb/s**
- Quad Data Rate (QDR) - **10Gb/s**
- Fourteen Data Rate (FDR) - **14Gb/s**
- Enhanced Data Rate (EDR) - **25Gb/s**

IB Widths (combined lanes): **1X, 4X**

RDMA is completely message-oriented, so all application messages are sent and received as units, unlike TCP/IP, which treats network communication as a stream of bytes.

## RDMA Technologies

1. IB
2. iWarp
3. ROCE



**Quality of Service (QoS)** is the ability to provide different priority to different applications, users, or data flows, or to guarantee a certain level of performance to a data flow.

QoS can be achieved by:

- Defining I/O channels at the adapter level
- Defining Virtual Lanes at the link level

It allows the control of the congestion on the network.

IB enables a scaling of up to

# 48,000

nodes in a single subnet



## CPU Offloads



The IB architecture supports packets transportation in the fabric with minimum CPU intervention.

That is achieved thanks to:

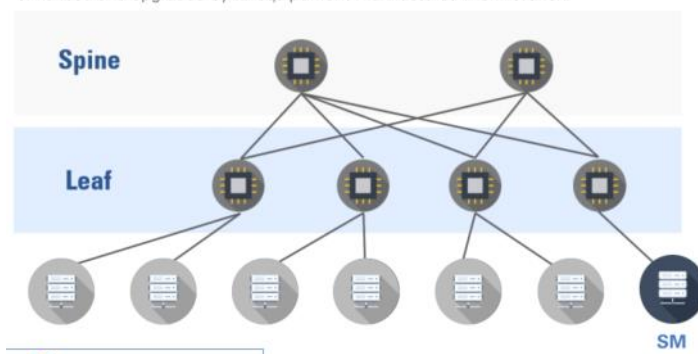
- Hardware based transport protocol
- Kernel Bypass
- Reliable transport - no need for additional CPU cycles for lost packets
- RDMA support - the ability to share and utilize data that is loaded in the memory directly between nodes in the fabric

## Subnet Manager

The Subnet Manager makes the fabric management very simple:

- Plug & Play end nodes environment
- Centralized route manager

As part of the IB spec it offers a variety of in-band diagnostics & management tools that can be further enhanced and upgraded by IB equipment manufactures like Mellanox.



$$\text{EDR} = 25 * 4 = 100 \text{ Gbps}$$

## Transport Layer - Responsibilities



The transport portion of the packet delivers the packet to the proper QP (Queue Pair) and instructs the QP how to process the packet's data.

The transport layer is responsible for segmentation/reassembly:

- **Segmenting** an operation into multiple packets when the message's data payload is greater than the maximum transfer unit (MTU) of the path
- The QP on the receiving end **reassembles** the data into the specified data buffer in its memory.

Responsible for configuring each QP with a certain class of **service type**:

- Connection-oriented vs. Datagram
- Acknowledged vs. Unacknowledged



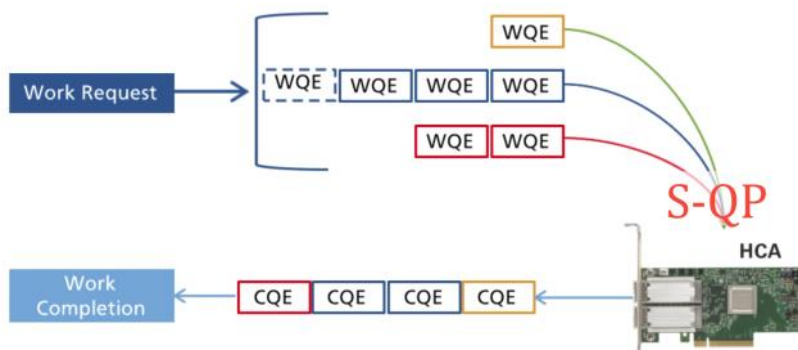
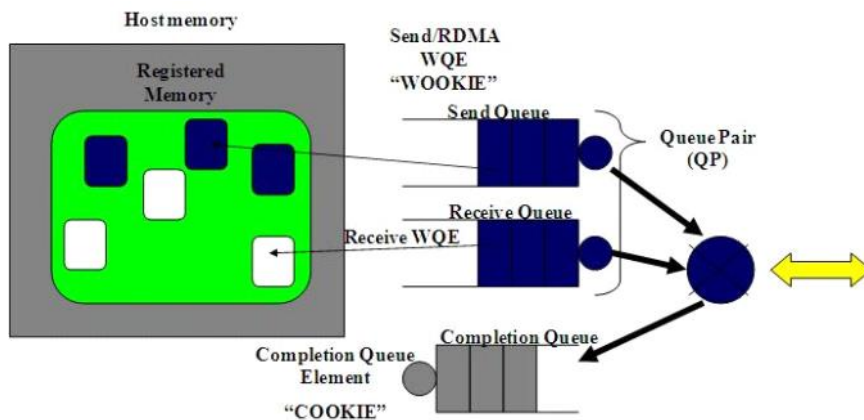
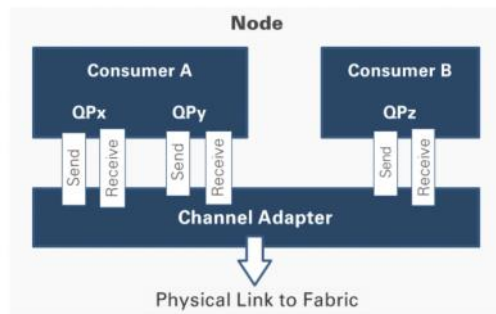


## Queue Pairs (QP)



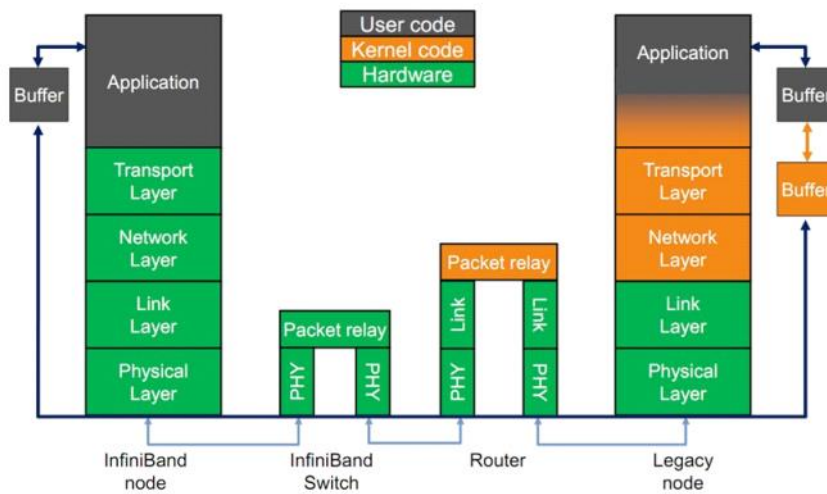
A QP consists of a **Send Queue** and a **Receive Queue**

- Each QP represents one end of a channel
- Send and receive queues are created as a pair
- A Queue Pair is identified by its **QP Number**
- In order to avoid OS involvement, applications have direct access to these QPs
- This is accomplished by mapping each application's virtual address space directly into the QPs



### Zero Copy

- No copy of data from APP to kernel
- use DMA and MMU (memory mapping)
- Sendfile/sendfile64 API



- IB HCA - Performance Everything in Hardware , So if SM is present it get the SUBNET ID without OS. IB HCA Understand Verbs and executes it.
- The caches in the remote CPU(s) won't be filled with the accessed memory content.

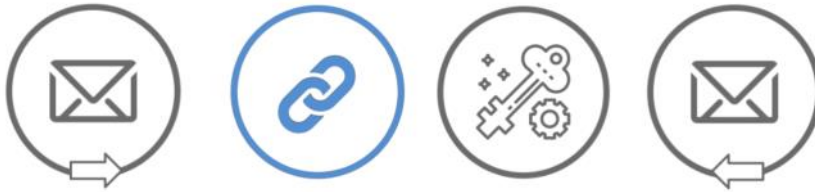
B mainly consists of 5 layers.

No	Layers	Functions
1st	Software Transport Verbs and Upper Layer Protocols	Interface between applications and hardwares Define methodology for management functions
2nd	Transport	Delivers packets to the appropriate QP node Message Assembly/De-assembly Access Right
3rd	Network	Route packets between different partitions/subnets
4th	Data Link(Symbols and framing)	Route packets on the same partition/subnet
5th	Physical	Signal levels/Media/Connections

From <<http://icf94.com/2016/06/27/2016-06-27-rdma/>>



For a SEND operation, the WQE specifies a block of data in the consumer's memory space for the hardware to send to the destination, letting a receive WQE already queued at the destination specify where to place that data.



### Memory Binding

Memory Binding instructs the hardware to alter memory registration relationships. It associates (binds) a memory Window to a specified range within an existing memory Region.



### Remote Direct Memory Access (RDMA)

RDMA supports zero-copy networking by enabling the network adapter to transfer data directly to or from application memory, eliminating the need to copy data between application memory and the data buffers in the operating system.

For an RDMA operation, the WQE specifies a block of data in the consumer's memory space for the hardware to send to the destination, letting a receive WQE already queued at the destination specify where to place that data.

In addition, it also specifies the address in the remote consumer's memory.

Thus an RDMA operation does not need to involve the receive work queue of the destination.

There are 3 types of RDMA operations: RDMA-WRITE, RDMA-READ, and ATOMIC.



### Receive Queue Operation

A RECEIVE WQE specifies where the hardware should place the data received from another consumer when that consumer executes a SEND operation.

Each time the remote consumer successfully executes a SEND operation, the hardware takes the next entry from the receive queue, places the received data in the memory location specified in that receive WQE, and places a CQE on the completion queue,

indicating to the consumer that the receive operation has completed. Thus the execution of a SEND operation causes a receive queue operation at the remote consumer.

### RDMA Opcodes

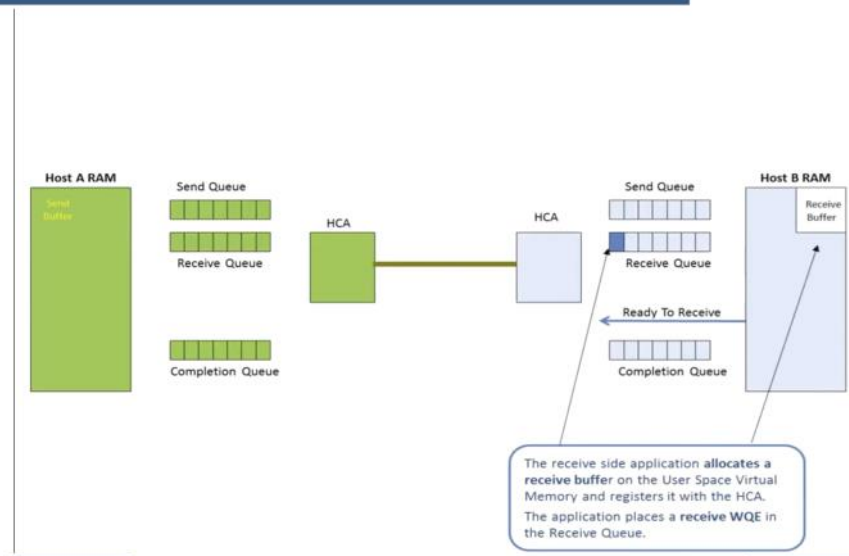
1. RDMA Send
2. RDMA Recv
3. RDMA Write
4. RDMA Read

## RDMA Send/Receive ( 2 Sided Transfer)

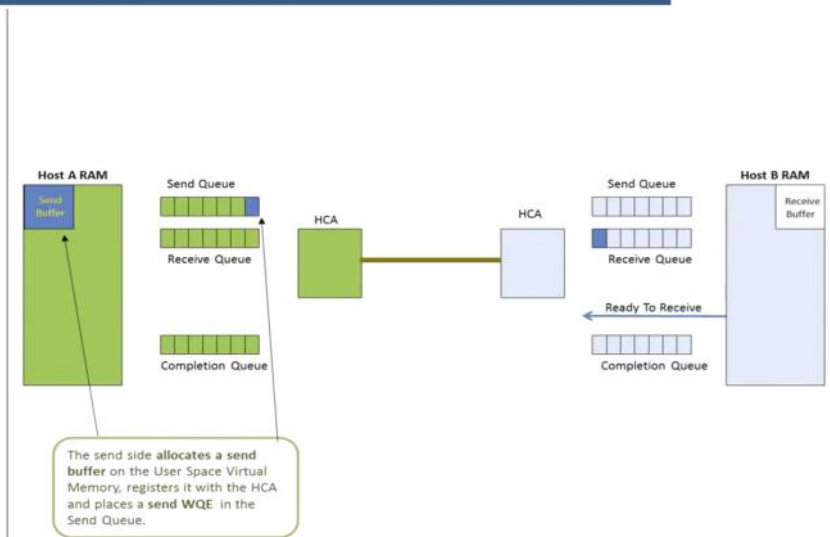
- A receiver first posts a RECV work request that describes a virtual memory area into which the adapter should place a single message.
- The sender then posts a SEND work request describing a virtual memory area containing the message to be sent.
- The network adapters transfer data directly from the sender's virtual memory area to the receiver's virtual memory area without any intermediate copies.
- Since both sides of the transfer are required to post work requests, this is called a "two-sided"

### Send Operation

#### Transport Layer - Send Operation Example

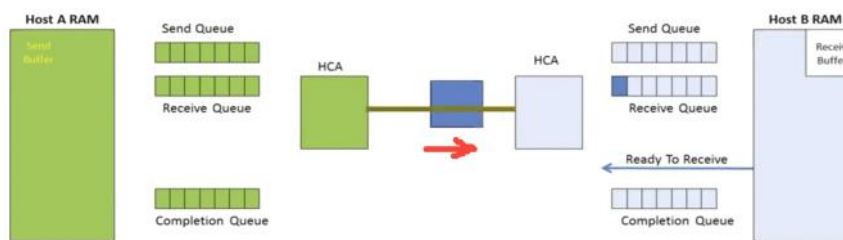
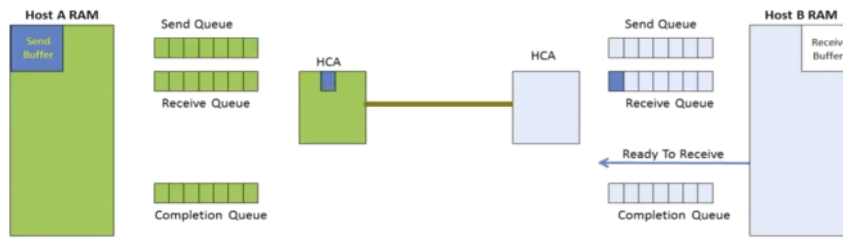


#### Transport Layer - Send Operation Example

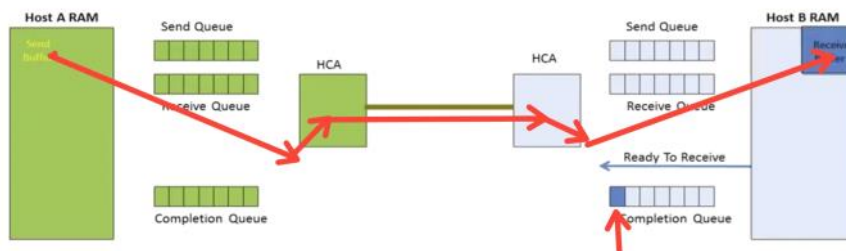
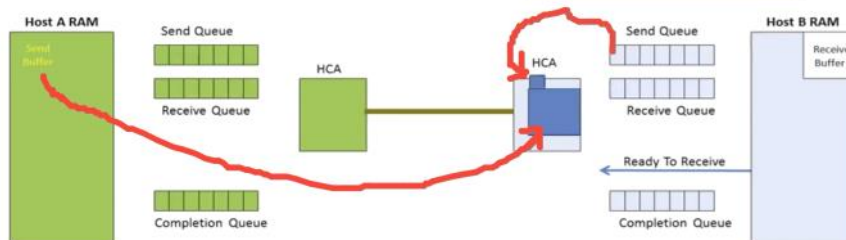




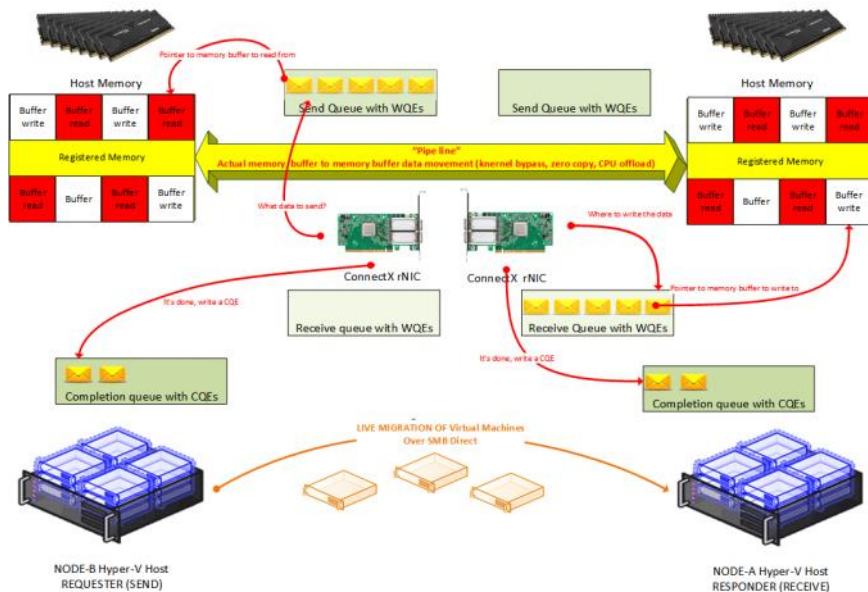
The HCA executes the send request, reads the buffer of the Host Ram and sends to the remote side (HCA).



When the packet arrives at the HCA, it executes the receive WQE Commands, places the buffer CONTENT in the appropriate location and generates a Completion Queue.



In RDMA, actions are specified by verbs which convey requests to the network adapter.



To initiate a transfer, **ibv\_post\_send** places a work request data structure describing the transfer onto a network adapter queue.

Data transfers are all asynchronous: once a work request has been posted, control returns to the user-space application which must later use the **ibv\_poll\_cq** function to remove a work completion data structure from a network adapter's completion queue.

This completion contains the status for the finished transfer and tells the application it can again safely access the virtual memory used in the transfer

### RDMA Write ( 1 -Side Transfer)

A sender posts a RDMA WRITE request that "pushes" a message directly into a virtual memory area that the receiving side previously described to the sender.

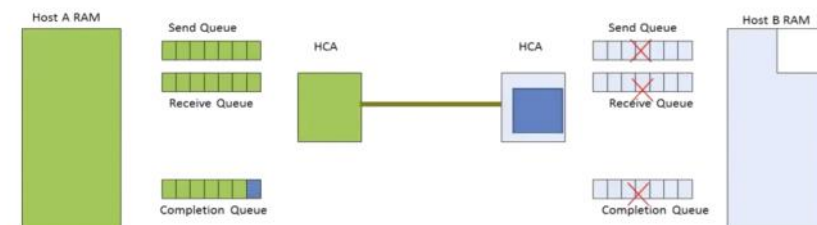
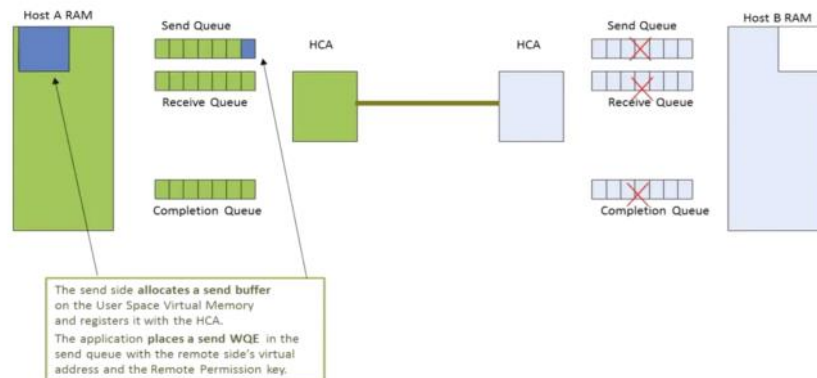
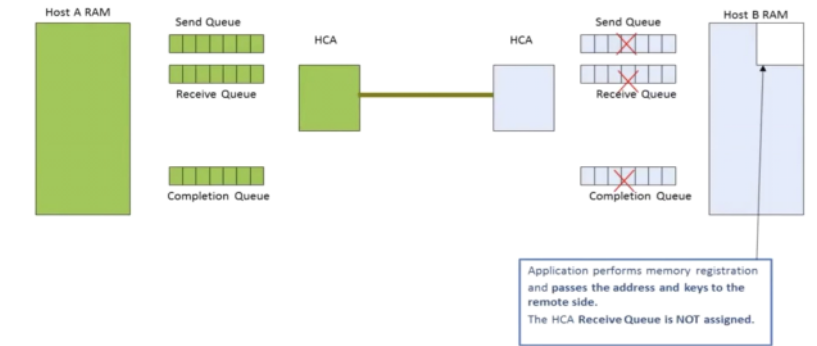


The receiving side's CPU is completely "passive" during the transfer, which is why this is called "one-sided."

Menu

Resources Glossary

## Transport Layer - RDMA Write Example



### RDMA Read ( 1- Side Transfer)

Receiver posts a RDMA READ request that "pulls" a message directly from the sending side's virtual memory, and the sending side's CPU is completely passive.

### RDMA Write with IMM

in RDMA write operation, since Passive do not when does data has been Completed,

Sender instead of RDMA write posts a RDMA WRITE WITH IMM , which is same as RDMA as write but but the send work request also includes 4 bytes of immediate (out-of-band) data that is delivered to the receiver on completion of the transfer.

The receiving side posts a RECV work request to catch these 4 bytes, and the work completion for the RECV indicates the status and amount of data transferred in the message.

- **Inline Data:** The API provides an optional “inline” feature that allows an interface adapter to copy the data from small messages into its own memory as part of a posted work request. This immediately frees the buffer for application reuse, and makes the transfer more efficient since the adapter has the data ready to send and does not need to retrieve it over the memory bus during the transfer

## Processing Completion Queue

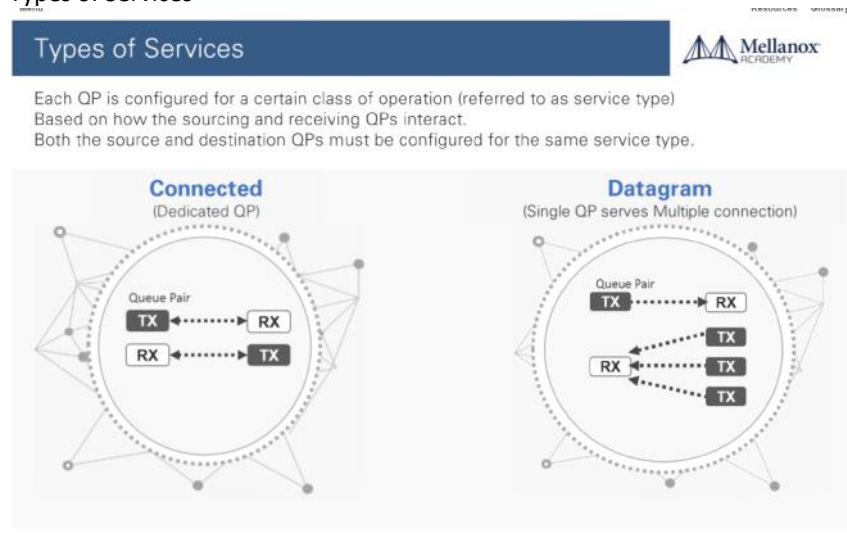
### 1. Busy Polling

The first completion detection strategy, called “busy polling”, is to repeatedly poll the completion queue until a completion becomes available. It allows immediate reaction to completions at the cost of very high CPU utilization, but requires no operating system intervention.

### 2. Event Notification

The second strategy, called “event notification”, is to set up a completion channel that allows an application to wait until the interface adapter signals a notification on this channel, at which time the application obtains the work completion by polling. It requires the application to wait for the notification by transferring to the operating system, but reduces CPU utilization significantly

## Types of Services



Connection Oriented

Reliable Connection (RC)	Unreliable Connection (UC)
<ul style="list-style-type: none"> <li>Each connection will always have a dedicated QP for all connection life</li> <li>If packet is dropped there is a way to track and resend it</li> <li>The most common service in use               <ul style="list-style-type: none"> <li>Implemented mainly with <b>RDMA</b> operations</li> <li>Good for <b>MPI applications, Storage</b></li> </ul> </li> <li>Each QP stores information about itself therefore many QPs consume more quota</li> </ul>	<ul style="list-style-type: none"> <li>Fixed dedicated QPs</li> <li>Good for <b>streaming data</b></li> <li>Packets cannot be tracked and resent in case of error</li> </ul>

## Datagram Oriented

Reliable Datagram (RD)	Unreliable Datagram (UD)
<ul style="list-style-type: none"> <li>Enables one-to-many paradigm</li> <li>Different connections might use the same Queue Pair</li> <li>Free Memory Resources because there is only 1 QP that stores information about itself (QP context)</li> <li>Management might be more complicated</li> <li>Not implemented in Mellanox products</li> </ul>	<ul style="list-style-type: none"> <li>A session does not require ACKs &amp; NACKs from the other side</li> <li>Call completion will be triggered by the HCA at the moment the data was sent to the other side</li> <li>Message size = MTU size because there is only 1 QP message</li> <li>Cannot be divided into packets</li> </ul>

- Most common is Reliable connection & unreliable for real time traffic.
- ACK are only send in reliable type of service ( RC and RD)

## Transport Layer - Retransmission



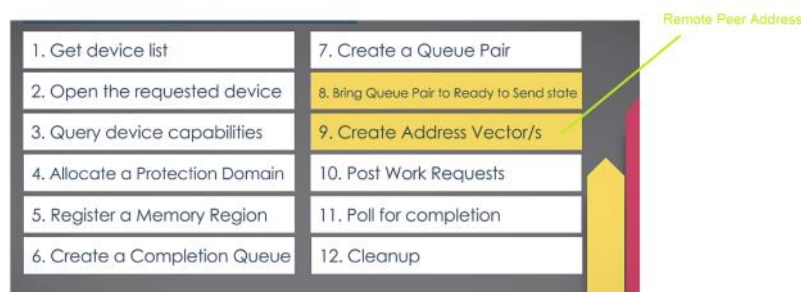
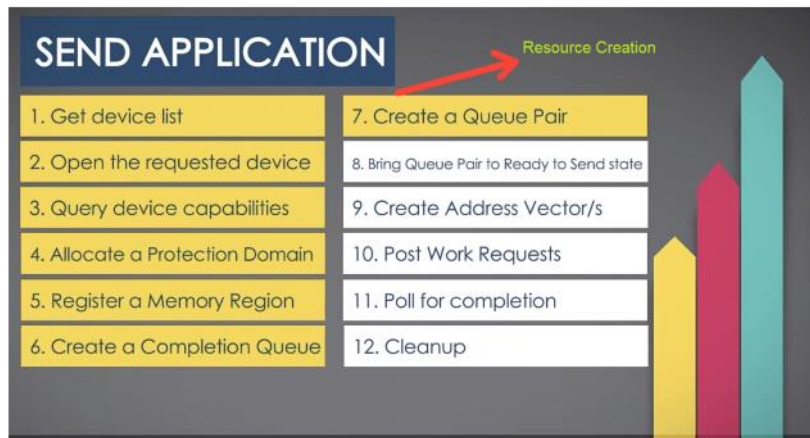
For **reliable transport services**, QPs maintain the flow of packets and **retransmit** in case a packet has dropped:

Each packet has a Packet Sequence Number (PSN) that is used by the receiver to identify lost packets. The PSN is a mechanism to maintain packet flow.

The receiver will send ACKs if packets arrive in order and negative ACKs otherwise. The send QP maintains a timer to catch packets that did not arrive at the receive QP or when an ACK was lost.

Retransmission is considered a "bad flow" which reduces performance or may break a connection (bad links, bad cables, and so on may cause a lot of retransmission)





- Connect two devices
- Install latest MLNX\_OFED
- Run the command "ibstat"
- Change port type: IB/ETH
- Check that the ports are enabled and up
- IB: make sure OpenSM is running
- Check for connectivity - "ping"

## FAQs

### Why do we always see Responder CQE Errors with RoCE RDMA?

*RDMA receivers need to post receive WQEs to handle incoming messages if the application does not know how many messages are expected to be received (e.g. by maintaining high-level message credits) they may post more receive WQEs than will actually be used. On application tear-down, if the application did not use up all of its receive WQEs the device will issue completion with the error for these WQEs to indicate HW does not plan to use them, this is done with a clear syndrome indication of "Flushed with error".*

## 1 .Get IB devices

MLX4\_0

Returns Array of dev\_list ( which is structure to ibv\_device)

lbv\_get\_device\_name -> string name

lbn\_free\_device\_list -> to free the dev\_list.

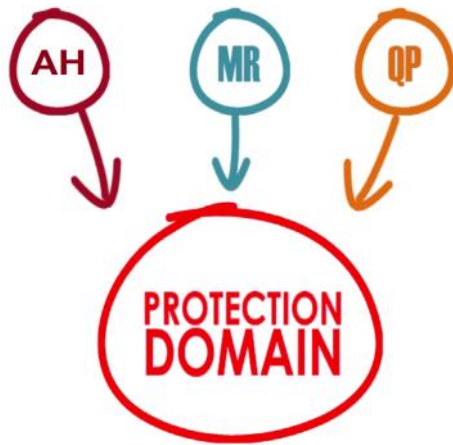
```
device
-----
mlx4 0
```

```
l-fw01 > ibstat
CA 'mlx4_0'
    CA type: MT4103
    Number of ports: 2
    Firmware version: 2.30.8006
    Hardware version: 0
    Node GUID: 0x0002c903000e95840
    System image GUID: 0x0002c903000e95843
    Port 1:
        State: Active
        Physical state: LinkUp
        Rate: 10
        Base lid: 0
        LMC: 0
        SM lid: 0
        Capability mask: 0x00010000
        Port GUID: 0x0202c9fffee95841
        Link layer: Ethernet
    Port 2:
        State: Active
        Physical state: LinkUp
        Rate: 40 (FDR10)
        Base lid: 3
        LMC: 0
        SM lid: 3
        Capability mask: 0x0251486a
        Port GUID: 0x0002c903000e95842
        Link layer: InfiniBand
```

## 2 .Open Device and Get the Attribute







AH- Address Vector or Address Handler ( remote Information)

MR - Local Memory Information

QP - SQ + RQ

`struct ibv_context *context = ibv_open_device(device);`--->returns context

`struct ibv_pd *pd = ibv_alloc_pd(context);`--> for device context return pd pointer

#### 4. Allocation of MR

Pinning of Memory for RMDA transfer



```
char buf[0x1800];

struct ibv_mr *mr = ibv_reg_mr(pd, buf, sizeof(buf), 0);
if (!mr) {
    fprintf(stderr, "Couldn't register MR\n");
    goto close_pd;
}

close_pd:
    ibv_dealloc_pd(pd);

close_device:
    ibv_close_device(context);

free_dev_list:
    ibv_free_device_list(dev_list);
```

Local read
  Local write
  Remote read
  Remote write
  Atomic operations

Permissions by default is Local Read.

## 5. Setting up Completion Queue ( CQ)

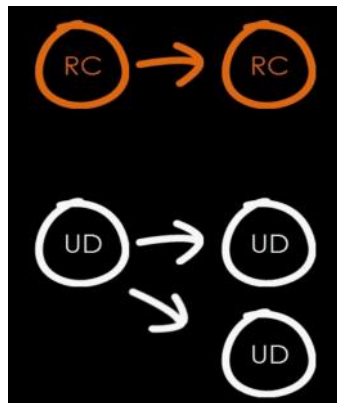
Used by Adapter to Notify application of status of transfer or completed work request, so that MR can be freed.  
Each entry in CQ is called CQE ( CQ element)  
Status of one or more completed WR.  
There can be single CQ for All the QP.

```
#define CQ_SIZE 0x100  
  
struct ibv_cq *cq = ibv_create_cq(context, CQ_SIZE, NULL,  
                                0, 0);  
  
if (cq) {  
    fprintf(stderr, "Couldn't create CQ\n");  
    goto free_mr;  
}
```

CQ\_SIZE is minimum size request, It can be more than that .

```
#define MAX_NUM_SENDS 0x10  
#define MAX_GATHER_ENTRIES 2  
#define MAX_SCATTER_ENTRIES 2
```

```
struct ibv_qp_init_attr attr = {  
    .send_cq = cq,  
    .recv_cq = cq,  
    .cap = {  
        .max_send_wr = MAX_NUM_SENDS,  
        .max_recv_wr = 0,  
        .max_send_sge = MAX_GATHER_ENTRIES,  
        .max_recv_sge = MAX_SCATTER_ENTRIES,  
    },  
    .qp_type = IBV_QPT_UD,---> Unreliable Datagram  
};
```



## 6. Creating a QP

QP = SQ + RQ

Application post WR to SQ contains memory Information to be Transferred, if it wants to Send Data  
Application post WR to RQ containing Memory information where Data has to Be written , once  
Data is received from Sender.

### A ) Queue Pair Attributes & QP Creation

```
define MAX_NUM_SENDS 0x10  
#define MAX_GATHER_ENTRIES 2  
#define MAX_SCATTER_ENTRIES 2
```

```

struct ibv_qp_init_attr qp_attr = {
    .send_cq = cq,                      -----> defines CQ for receive and send QP/
    .recv_cq = cq,
    .cap = {
        .max_send_wr = MAX_NUM_SENDS, ----->Max no of WR for TX
        .max_recv_wr = 0,
        .max_send_sge = MAX_GATHER_ENTRIES,
        .max_recv_sge = MAX_SCATTER_ENTRIES,
    },
    .qp_type = IBV_QPT_UD,----->UD type of QP
};

```

```

struct ibv_qp *qp = ibv_create_qp(pd, &qp_attr);
if (!qp) {
    fprintf(stderr, "Couldn't create QP\n");
    goto free_cq;
}

```

B ) QP Change state

```
#define WELL_KNOWN_QKEY 0x11111111
```

```

qp_modify_attr.qp_state    = IBV_QPS_INIT;---> QP State it should later RTS ( for TX) or RTR ( for RX)
qp_modify_attr.pkey_index  = 0; --->
qp_modify_attr.port_num    = dev_port; --->( Port no in HCA)
qp_modify_attr.qkey        = WELL_KNOWN_QKEY;---> only Used in UD and is same both side (RX and TX)
if (ibv_modify_qp(qp, &qp_modify_attr,
    IBV_QP_STATE |
    IBV_QP_PKEY_INDEX |
    IBV_QP_PORT |
    IBV_QP_QKEY)) {
    fprintf(stderr, "Failed to modify QP to INIT\n");
    goto free_qp;
}

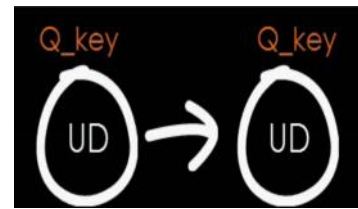
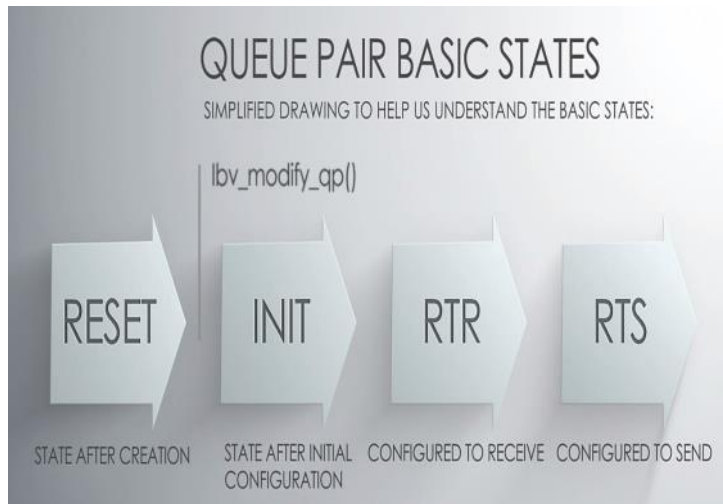
```

```

qp_modify_attr.qp_state    = IBV_QPS_RTR;

if (ibv_modify_qp(qp, &qp_modify_attr, IBV_QP_STATE)) {
    fprintf(stderr, "Failed to modify QP to RTR\n");
    goto free_qp;
}

```



### RESET

During creation

### RTR state

In this state, the QP handles incoming packets. If the QP don't be used as a requester (i.e. Work Requests won't be posted to its send queue), the QP may stay in the RTR state.

### RTS state

In most of the applications, the QPs will be transitioned to the RTS state. In this state, the QP can send packets as a requester and handle incoming packets.

Even if your QP won't act as a requester, it can be transitioned to this state.