



Ethernet Services over IPoIB

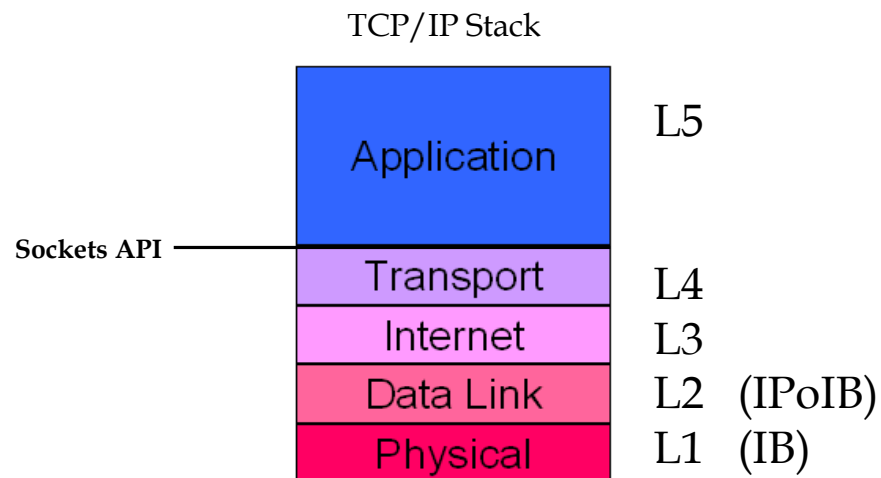
Ali Ayoub, Mellanox Technologies
March, 2012

Background

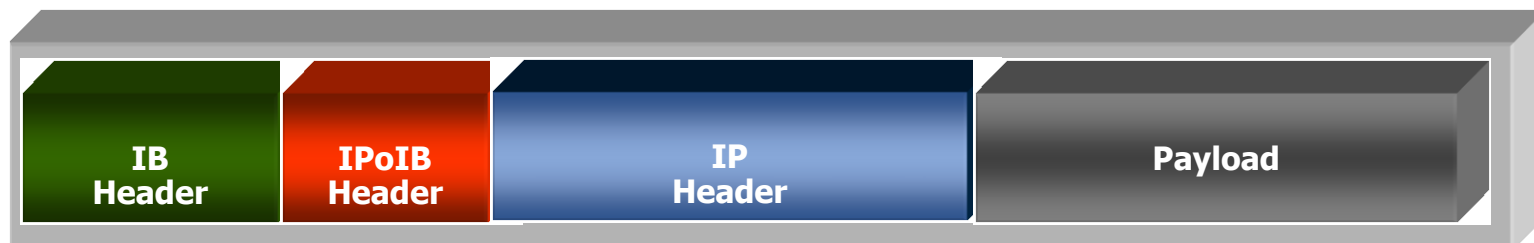
- What is IPoIB?
 - IP encapsulation over InfiniBand
 - RFC 4391/4392
 - Provides IP services over InfiniBand fabric
- Benefits:
 - Acts like an data-link within the TCP/IP Stack
 - Socket-based apps run transparently
 - Allow users to run IP-based unmodified applications on InfiniBand
 - Supports NIC offloads: CSUM, TSO, etc.. And other performance features such as: NAPI, TSS, RSS, etc..

IPoIB Packet Format

- IPoIB integrated as Data Link layer
- Sockets API is not affected



- Packet frame:

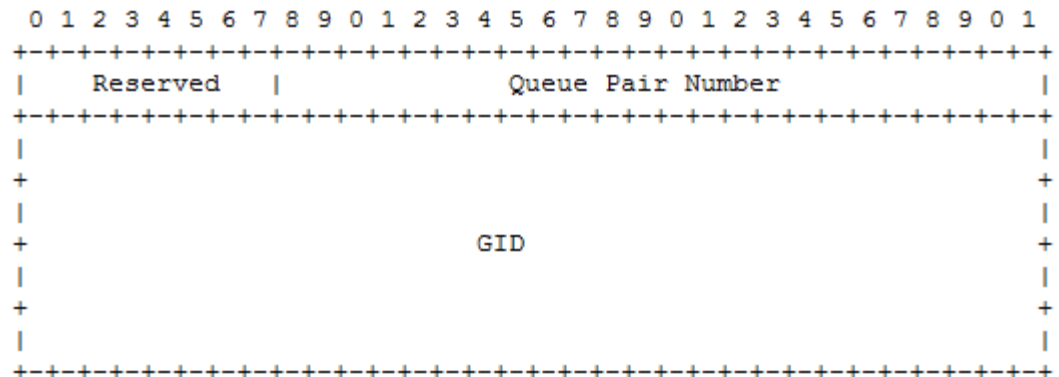


Limitations

- IPoIB is limited for IP applications only
 - No Ethernet Header encapsulation
- IPoIB network interface doesn't act like a standard Ethernet Network Device
 - Non-standard Ethernet interface utilities & characteristics
 - 20 Bytes link-layer address
 - Ethernet link-layer (MAC) address is 6 Bytes
 - DHCP requires using dhcp-client identifier
 - Host administrator must be aware of PKEYs
 - While Ethernet interfaces support VLANs
 - vconfig command cannot be used

Limitations cont.

- Link Layer setting, modification, migration
 - IPoIB Link Layer address is based on QPN/GID
 - Cannot be controlled by the user
 - Host administrator cannot set or re-configure the link layer address (MAC)

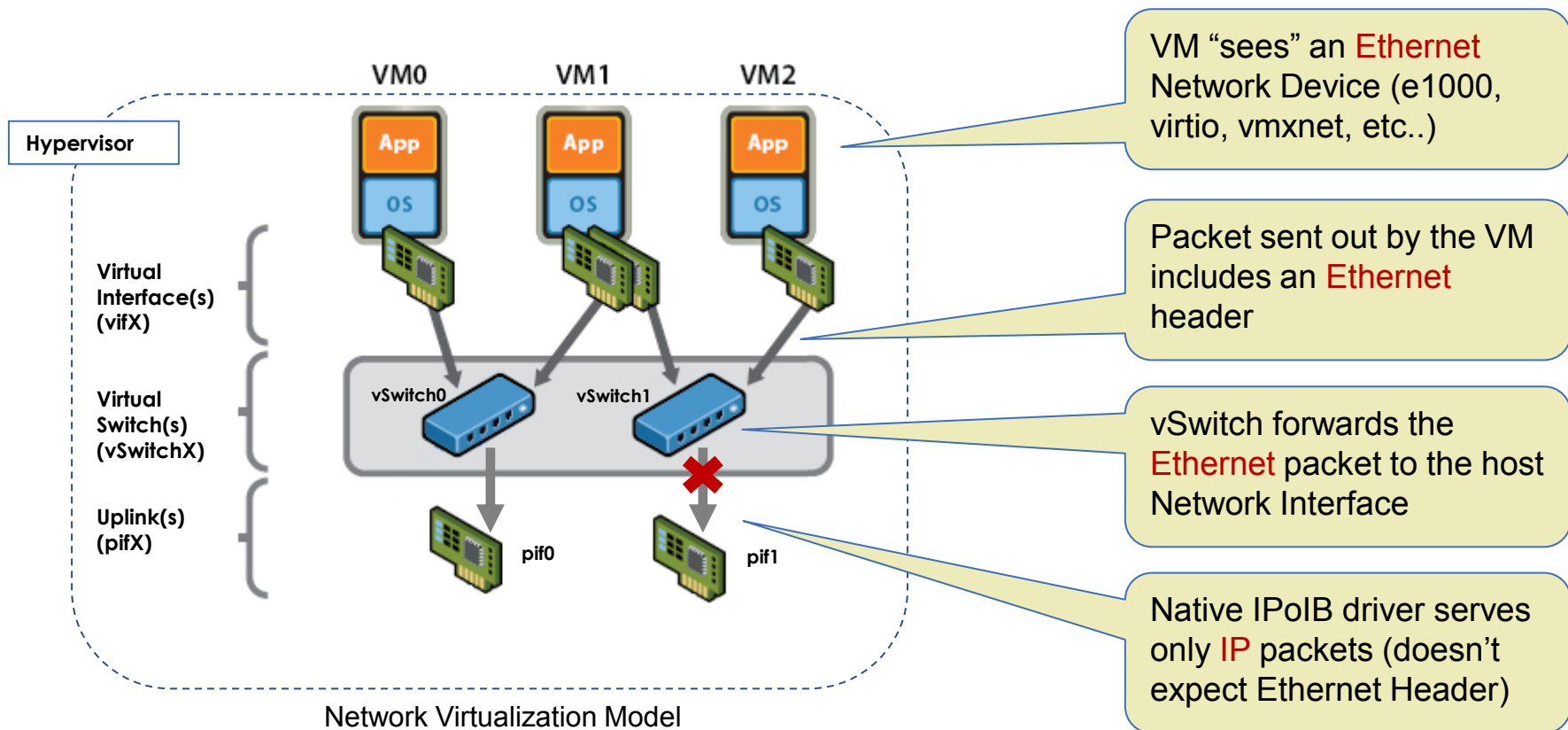


IPoIB Link Layer (RFC 4391)

Limitations cont.

- IPoIB cannot be used in fully-virtualized or para-virtualized environments
 - Hypervisor networking model normally use “bridged mode”
 - Virtual Switch (vSwitch) is Ethernet L2 switch
 - IPoIB NIC cannot be enslaved to a vSwitch
- Promiscuous mode is unsupported
 - Simplifies vSwitch functionality

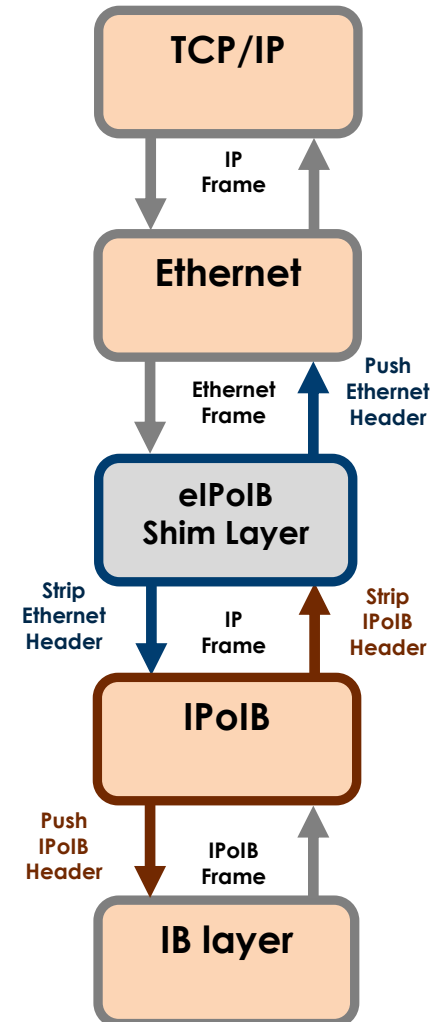
IPoIB & Network Virtualization



Reference <http://www.vmware.com/resources/techresources/997>

Solution

- Create a new Ethernet Network Device over IPoIB interface
- Managed by “eIPoIB” kernel module
- Registers a standard Ethernet Interface into the Operating System
- Same “Look & Feel” as an Ethernet NIC
 - ifconfig, vconfig, ethtool, promiscuous, etc..
- eIPoIB is a **shim layer** driver between the TCP/IP stack and the native IPoIB



eIPoIB vs. IPoIB Interface

- `ifconfig ib0`

```
root@dev-l-vrt-034:~  
[root@dev-l-vrt-034 ~]# ifconfig ib0  
Ifconfig uses the ioctl access method to get the full address information, which limits hardware addresses to 8 bytes.  
Because Infiniband address has 20 bytes, only the first 8 bytes are displayed correctly.  
Ifconfig is obsolete! For replacement check ip.  
ib0      Link encap:InfiniBand  HWaddr A0:00:00:50:FE:80:00:00:00:00:00:00:00:00:00:00:00:00:00:00:00:00  
  
inet addr:11.134.34.1  Bcast:11.134.255.255  Mask:255.255.0.0  
inet6 addr: fe80::202:c903:4b:da63/64 Scope:Link  
UP BROADCAST RUNNING MULTICAST  MTU:65520  Metric:1  
RX packets:883 errors:0 dropped:0 overruns:0 frame:0  
TX packets:26 errors:0 dropped:0 overruns:0 carrier:0  
collisions:0 txqueuelen:1024  
RX bytes:52544 (51.3 KiB)  TX bytes:1908 (1.8 KiB)  
  
[root@dev-l-vrt-034 ~]# █
```

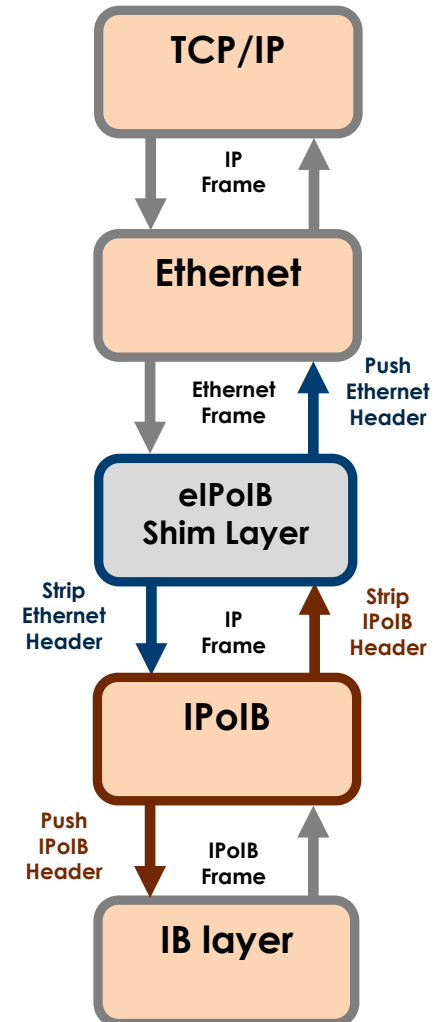
eIPoIB vs. IPoIB Interface

- ifconfig eth3

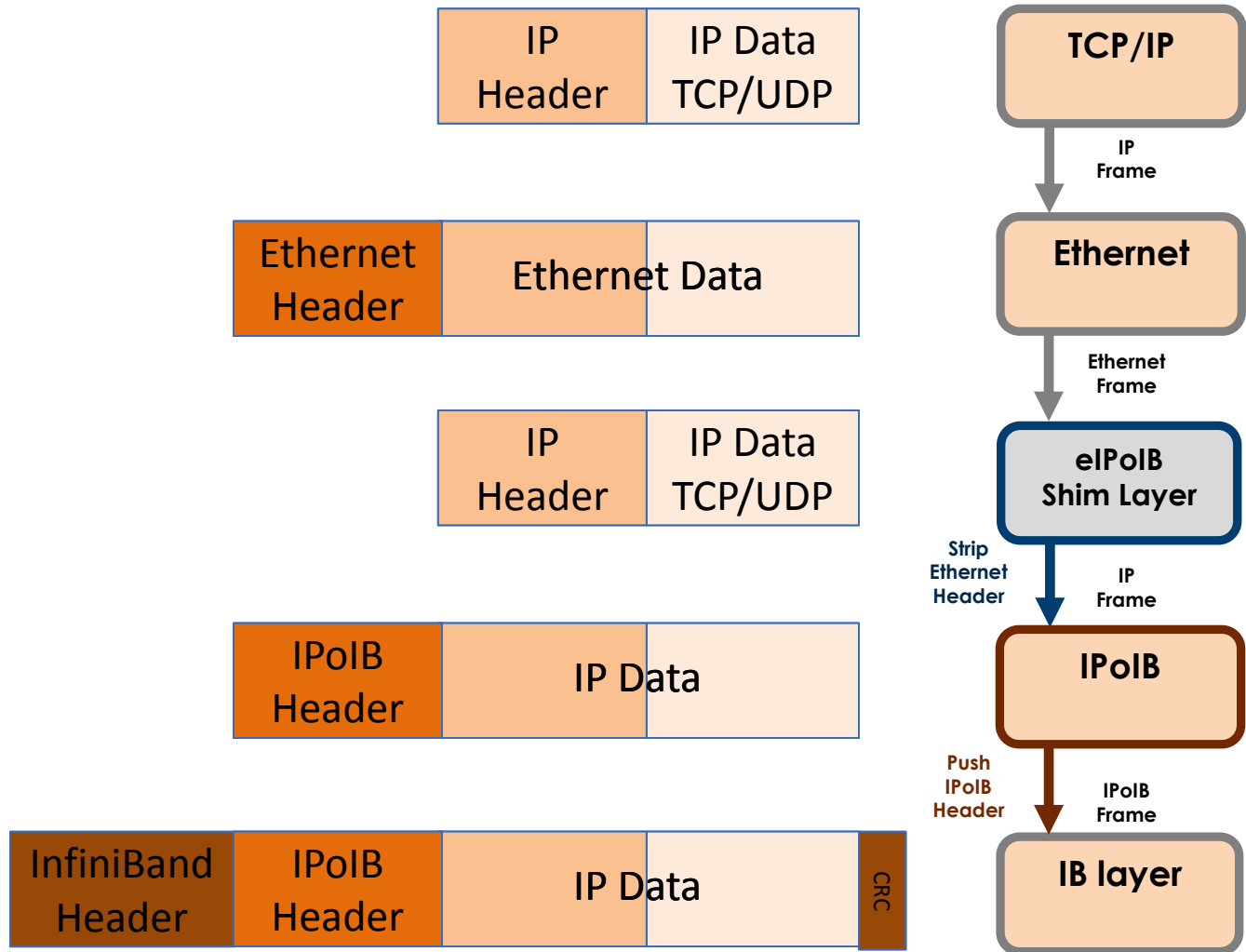
```
root@dev-l-vrt-034:~  
[root@dev-l-vrt-034 ~]# ifconfig eth3  
eth3      Link encap:Ethernet  HWaddr 02:02:C9:4B:DA:63  
          inet addr:12.134.34.1  Bcast:12.134.255.255  Mask:255.255.0.0  
          inet6 addr: fe80::2:c9ff:fe4b:da63/64  Scope:Link  
          UP BROADCAST RUNNING PROMISC MASTER MULTICAST  MTU:65520  Metric:1  
          RX packets:4606 errors:0 dropped:0 overruns:0 frame:0  
          TX packets:898 errors:0 dropped:0 overruns:0 carrier:0  
          collisions:0 txqueuelen:0  
          RX bytes:285056 (278.3 KiB)  TX bytes:78156 (76.3 KiB)  
  
[root@dev-l-vrt-034 ~]# ethtool -i eth3 | head -n1  
driver: eth_ipoib:ib0  
[root@dev-l-vrt-034 ~]#
```

eIPoIB Operations in a Nutshell

- Initialization:
 - Register Ethernet network interface (eth0)
 - Map it to native IPoIB interface (ib0)
- Transmit:
 - Packet to be sent has Ethernet Header
 - eIPoIB strips the Ethernet Header
 - Pass the IP packet to the underlying IPoIB interface.
 - Non IP/ARP/RARP packets are dropped
 - Native IPoIB interface sends the packet out
- Receive:
 - Packet received is an IP packet
 - Native IPoIB hands over the IP packet to eIPoIB layer
 - eIPoIB pushes the Ethernet header
 - Packet forwarded to upper layers as regular Ethernet frame

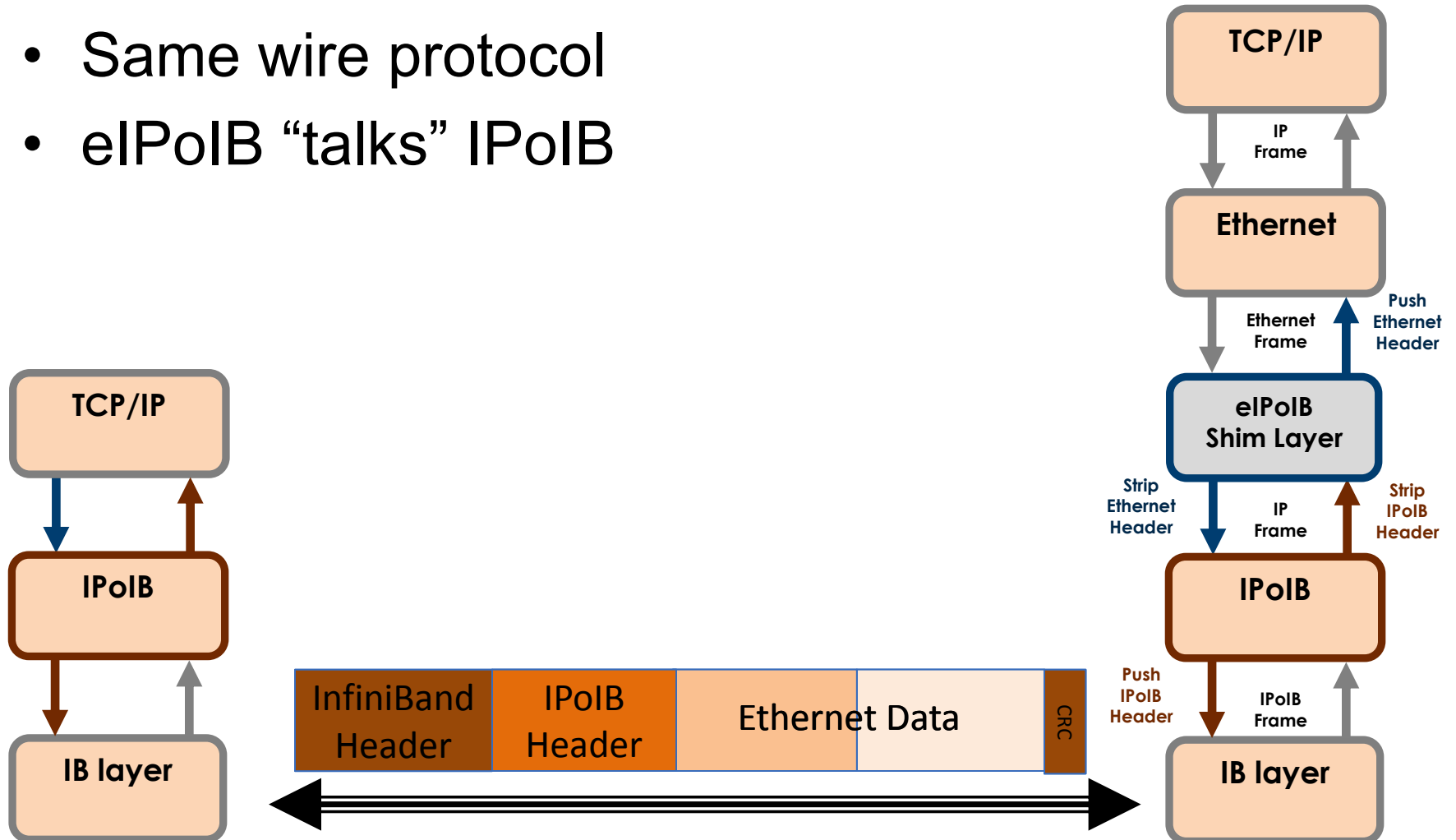


eIPoIB Frame Flow

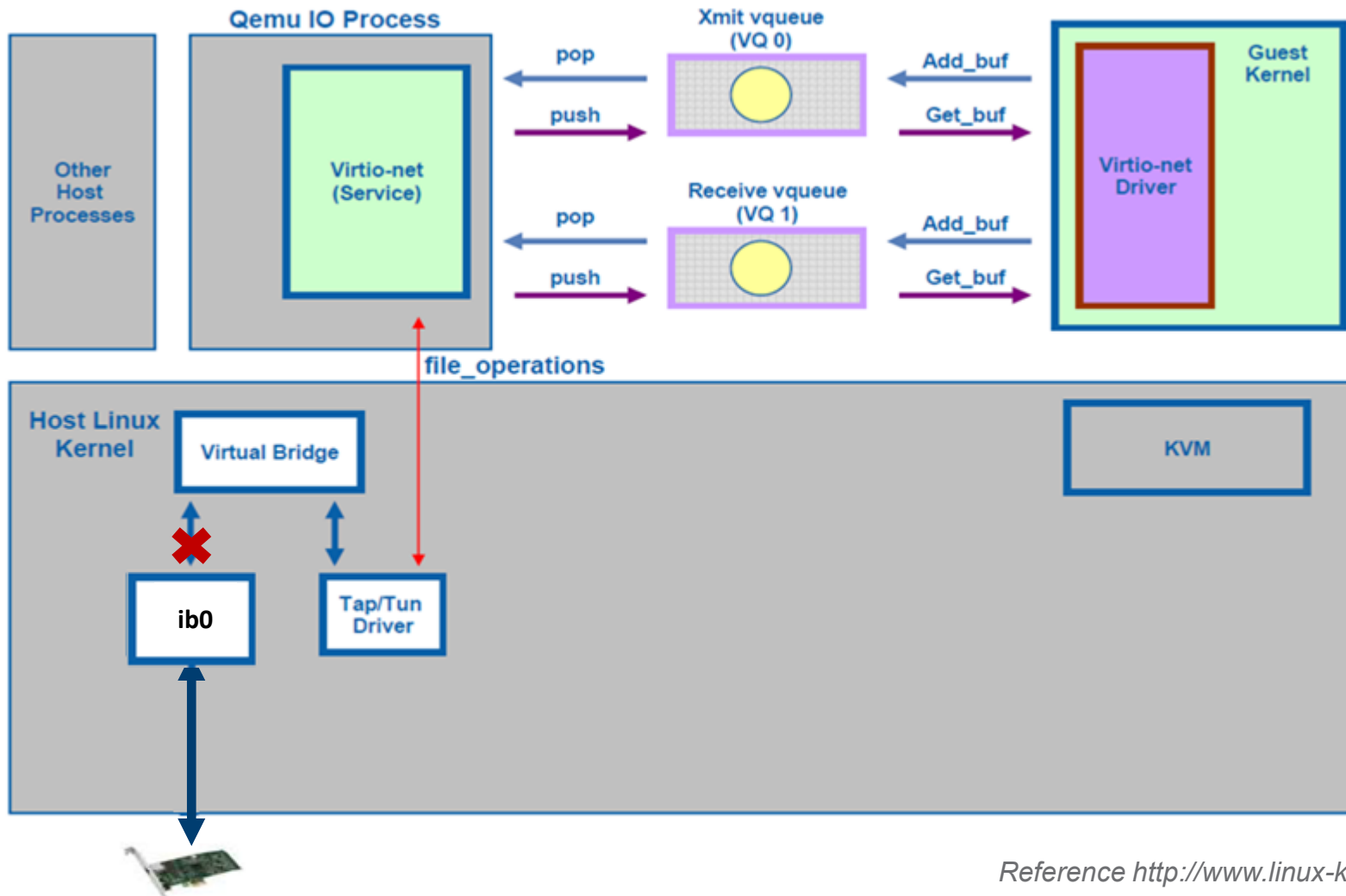


eIPoIB/IPoIB Interoperability

- Same wire protocol
- eIPoIB “talks” IPoIB

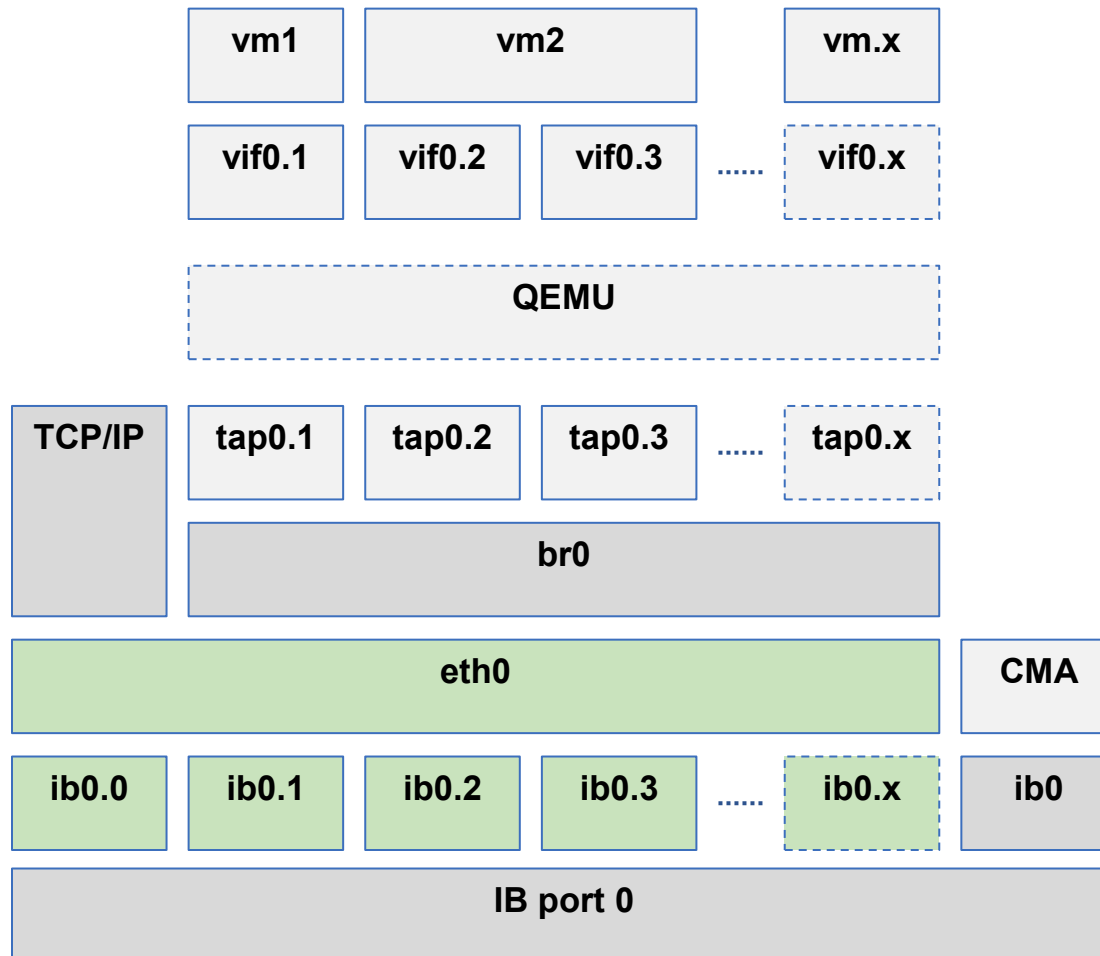


eIPoIB Networking Model (KVM)



Reference <http://www.linux-kvm.org/wiki>

eIPoIB Design (KVM)



eIPoIB – Closer Look

- Promiscuous mode
- MAC Translation

Promiscuous Mode

- Normally vSwitch uplink is put in promiscuous mode
 - InfiniBand doesn't support promiscuous mode
- Promiscuous mode is simulated by:
 - a. Snooping the src.mac and vlan of outgoing packets
 - b. OS notifies the driver when a new MAC/VLAN need to be “served”. For example: the driver can get a notification from libvirt library (available for KVM/XEN) when a new VM virtual NIC is created
- Multicast promiscuous support requires ***IGMP/MLD Snooping*** in eIPoIB level

MAC Translation

- Requirements:
 - eIPoIB exposes an Ethernet MAC
 - Local Ethernet MAC (LEMAC) 6 bytes length
 - eIPoIB neighbors are seen as Ethernet neighbors
 - Remote Ethernet MAC (REMAC) 6 bytes length
 - IPoIB Local MAC used on the wire
 - Local IPoIB MAC (LIMAC) 20 bytes
 - IPoIB neighbors' MAC used on the wire
 - Remote IPoIB MAC (RIMAC) 20 bytes

MAC Translation cont.

Receive Flow

- Receiver QP => dst.mac
- SQPN/SLID => src.mac
 - Remember QPN/LID to IPoIB-MAC (QPN/GID) mapping
- IPoIB header => Ethernet ethertype
- Replace IPoIB header by Ethernet header

Transmit Flow

- src.mac => QP (child interface)
 - Source MAC normally controlled by the host/hypervisor admin
- dst.mac => IPoIB-MAC
- Ethernet header ethertype => IPoIB Header
- Strip Ethernet header, and handover packet to IPoIB

ARP/NDP packets are modified in TX/RX flow, so the MAC addresses in the packet payload are updated accordingly



Questions?