# Thesis Background Presentation

Emerson Ford

October 5, 2020

University of Utah School of Computing

# Containers

## Container Overview

- Increasingly popular framework to distribute and deploy applications.
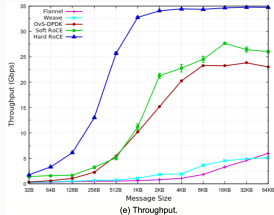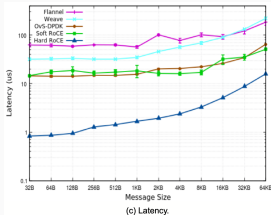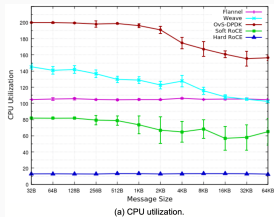- Tools like **Kubernetes** have become popular for container orchestration.

# Container Requirements

- Isolation
  - namespaces
  - cgroups
  - network policy
- Portability
  - migration
- Performance
  - low isolation overhead

# Container Networking Requirements

- Control Plane Policies
  - firewall
  - routing
  - vlans
- Data Plane Policies
  - QoS
  - metering
  - fairness

(a) CPU utilization.
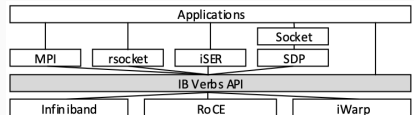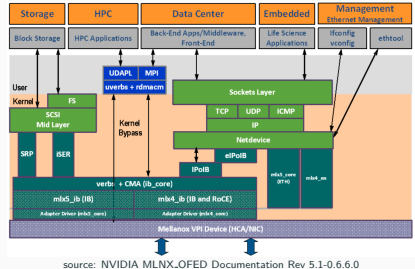
(c) Latency.

(e) Throughput.

- Current networking isolation requires pretty significant performance sacrifices.
- Less than ideal for HPC applications.

source: A Performance Comparison of Container Networking Alternatives by Ubaid Abbasi, El Houssine Bourhim, Mouhamad Dieye, and Halima Elbiaze
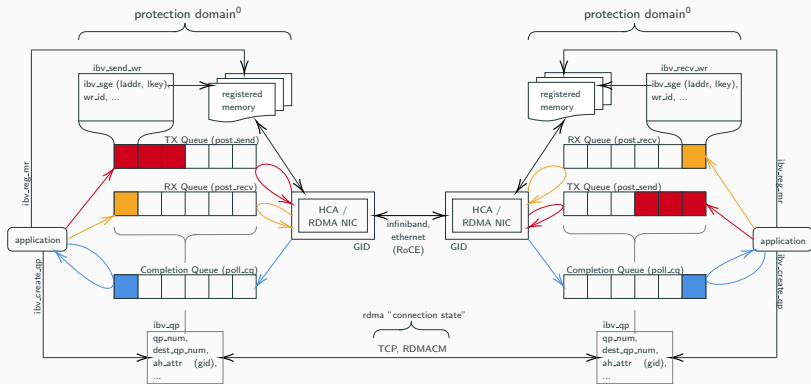
# RDMA Overview

- Form of kernel bypass networking
- `libibverbs` is the "narrow waist" of RDMA operations
- Extremely low latency, high throughput



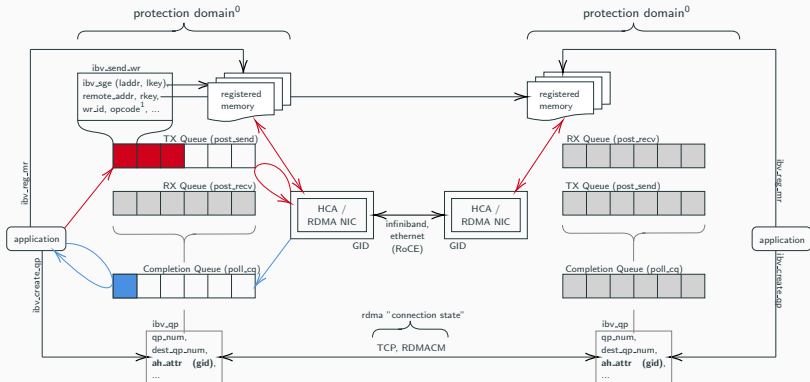source: NVIDIA MLNX_OFED Documentation Rev 5.1-0.6.6.0



source: FreeFlow Paper Figure 3

# RDMA Two-Sided Ops



[0] Every object in the protection domain is mapped in the application's virtual address space. The HCA can access every object in the protection domain.

protection domain[0]

protection domain[0]

ibv_send_wr
ibv_sge (laddr, lkey),
remote_addr, rkey,
wr_id, opcode[1], ...

registered memory

registered memory

ibv_reg_mr

TX Queue (post_send)

RX Queue (post_recv)

RX Queue (post_recv)

TX Queue (post_send)

application

HCA /
RDMA NIC

infiniband,
ethernet
(RoCE)

HCA /
RDMA NIC

application

ibv_reg_mr

GID

GID

ibv_create_qp

Completion Queue (poll_cq)

Completion Queue (poll_cq)

ibv_create_qp

ibv_qp
qp_num,
dest_qp_num,
**ah_attr   (gid)**,
...

rdma "connection state"

TCP, RDMACM

ibv_qp
qp_num,
dest_qp_num,
**ah_attr   (gid)**,
...

---

[0] Every object in the protection domain is mapped in the application's virtual address space. The HCA can access every object in the protection domain.

[1] opcode is one of `IBV_WR_RDMA_WRITE`, `IBV_WR_RDMA_READ`, `IBV_WR_ATOMIC_CMP_AND_SWP`, `IBV_WR_SEND`

- RDMA significantly improves HPC application performance.
- Containers are quickly becoming a common framework for application distribution and deployment, but container networking isolation is slow.
- **Note**: similar research is being done for RDMA use in VMs in the cloud

**Problem Statement**

How can we enable the use RDMA in containers while preserving container requirements and performance?
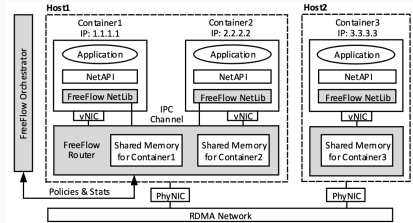
# Software Approach

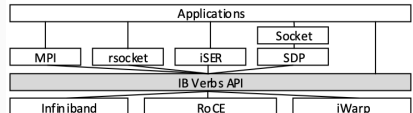**Microkernel / Paravirtualized Approach:**

- FreeFlow
- MasQ

**Virtualized RDMA:**

- SoftRoCE

- RDMA client (FreeFlow Library / FFL)
- RDMA server (FreeFlow Router / FFR)
- Communicate with IPC and shared memory
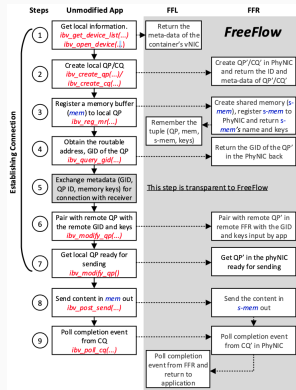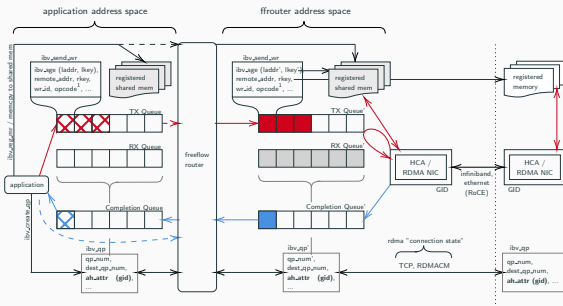- Only need `LD_PRELOAD` to make a FreeFlow compatible application



source: FreeFlow Paper Figure 4



source: FreeFlow Paper Figure 3

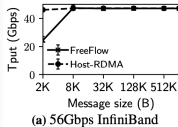source: FreeFlow Paper Figure 5

- IPC communication
  - Latency can be $\geq 5\mu s$
- Fastpath
  - Move `TX Queue` and `RX Queue` to shared memory with FreeFlow router
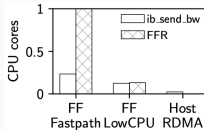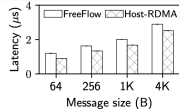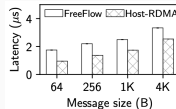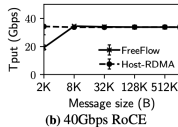  - FreeFlow router spin reads these queue pairs (with cache flushes)

- Use shared memory to support and speed up one-sided operations.
  - Requires `malloc` highjack to page align memory.
  - New functions `ibv_malloc` and `ibv_free` to avoid this.
  - ffrouter must replace `laddr` with `laddr'`
- Utilize `libibverb`'s built in struct flattening to avoid deep copies in RPC.
- Multiple Unix sockets for parallel RDMA queue pairs to avoid head of line blocking.

## FreeFlow Benefits

- Control plane policy enforcement on queue pairs.
  - QoS and network overlay enforcements
- RDMA vNIC can be assigned a private IP, allowing for non-live container migration.
  - ffrouter can query network overlays to get private IP -> public IP translation.
- FreeFlow library can run TCP over RDMA using `rsocket` in `libibverbs`.

# FreeFlow Performance



(a) 56Gbps InfiniBand
(b) 40Gbps RoCE
source: FreeFlow Paper Figure 9
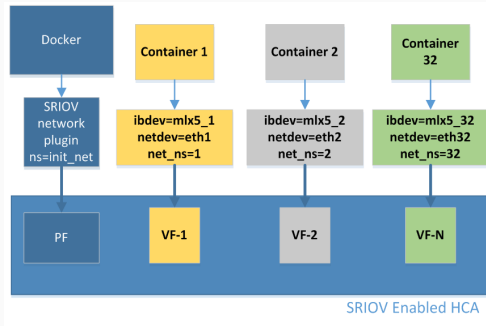source: FreeFlow Paper Figure 10
source: FreeFlow Paper Figure 12

- Approx 33% increase in latency for small messages
- Small message sizes bound in tput due to Fastpath single thread bottleneck
- Non-Fastpath CPU util overhead scales with actual load
  - Fastpath requires at least a single CPU core
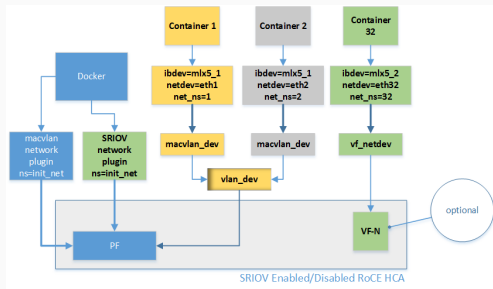- Applications finish at close to host RDMA speeds

# Hardware Approach

- Not portable.
- Control plane policy enforcement relies on switch reconfiguration.

- Control plane policy enforcement relies on switch reconfiguration?
- "GID table entries are created whenever an IP address is configured on one of the Ethernet devices of the NIC's ports."

## RDMA Shared Device

- RDMA namespaces
- RDMA cgroups
- ConnectX6 NICs allow for hardware rules?