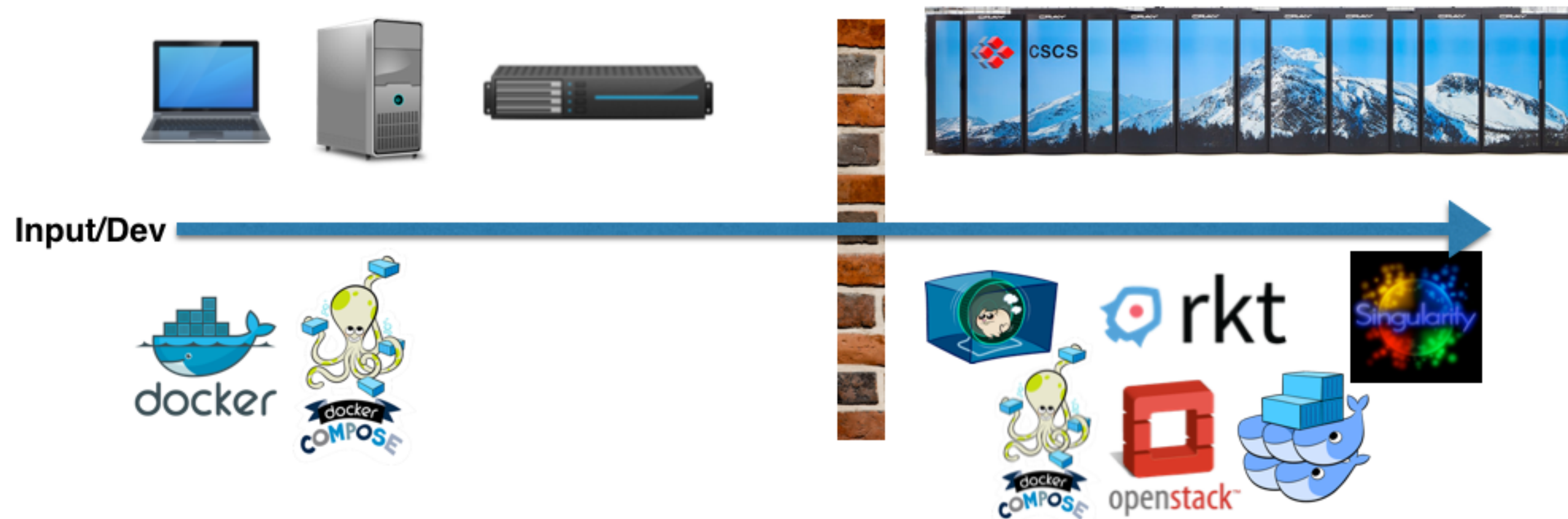# Missing Pieces for HPC

# Container Mobility

Spinning up production-like environment is…

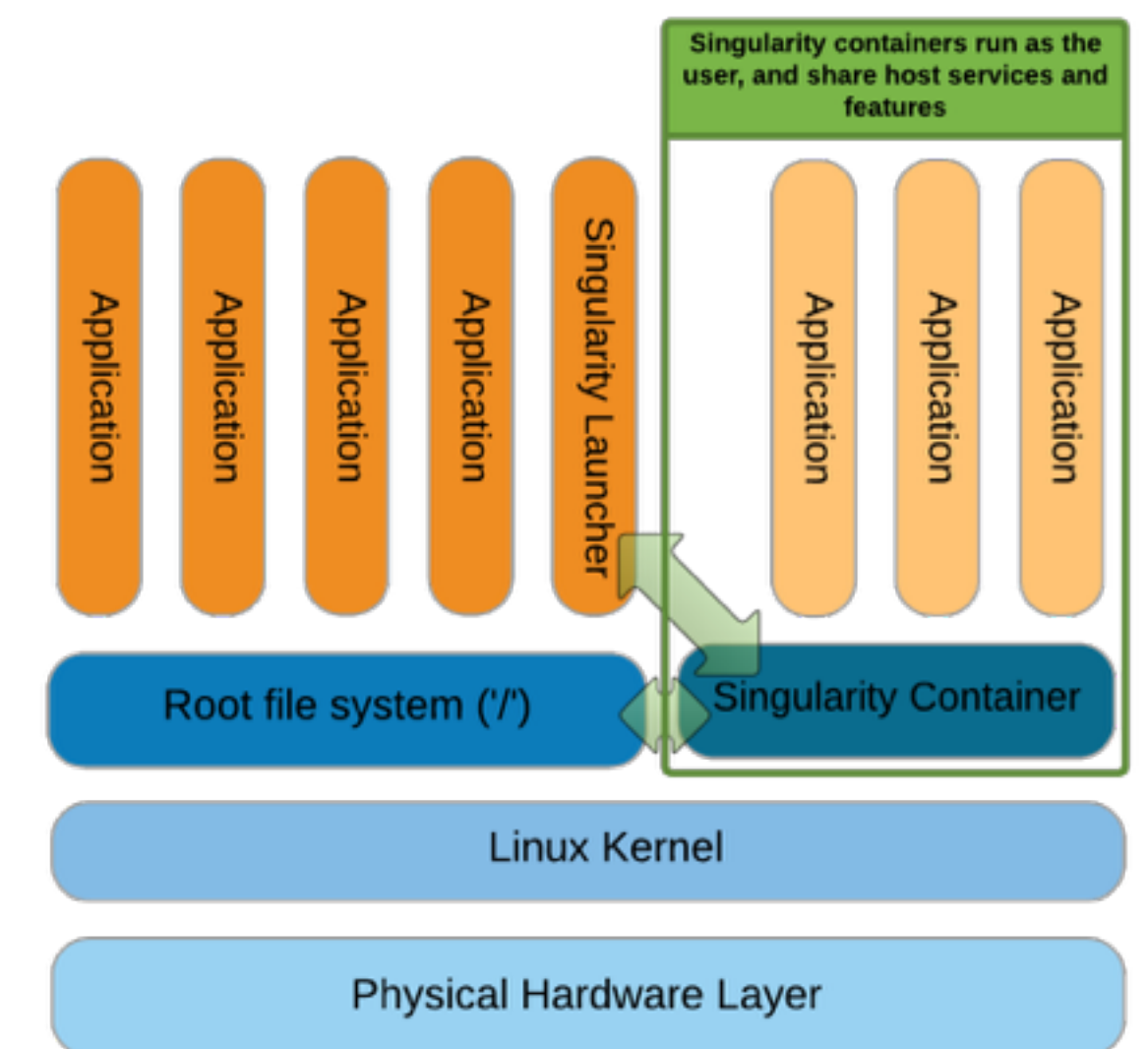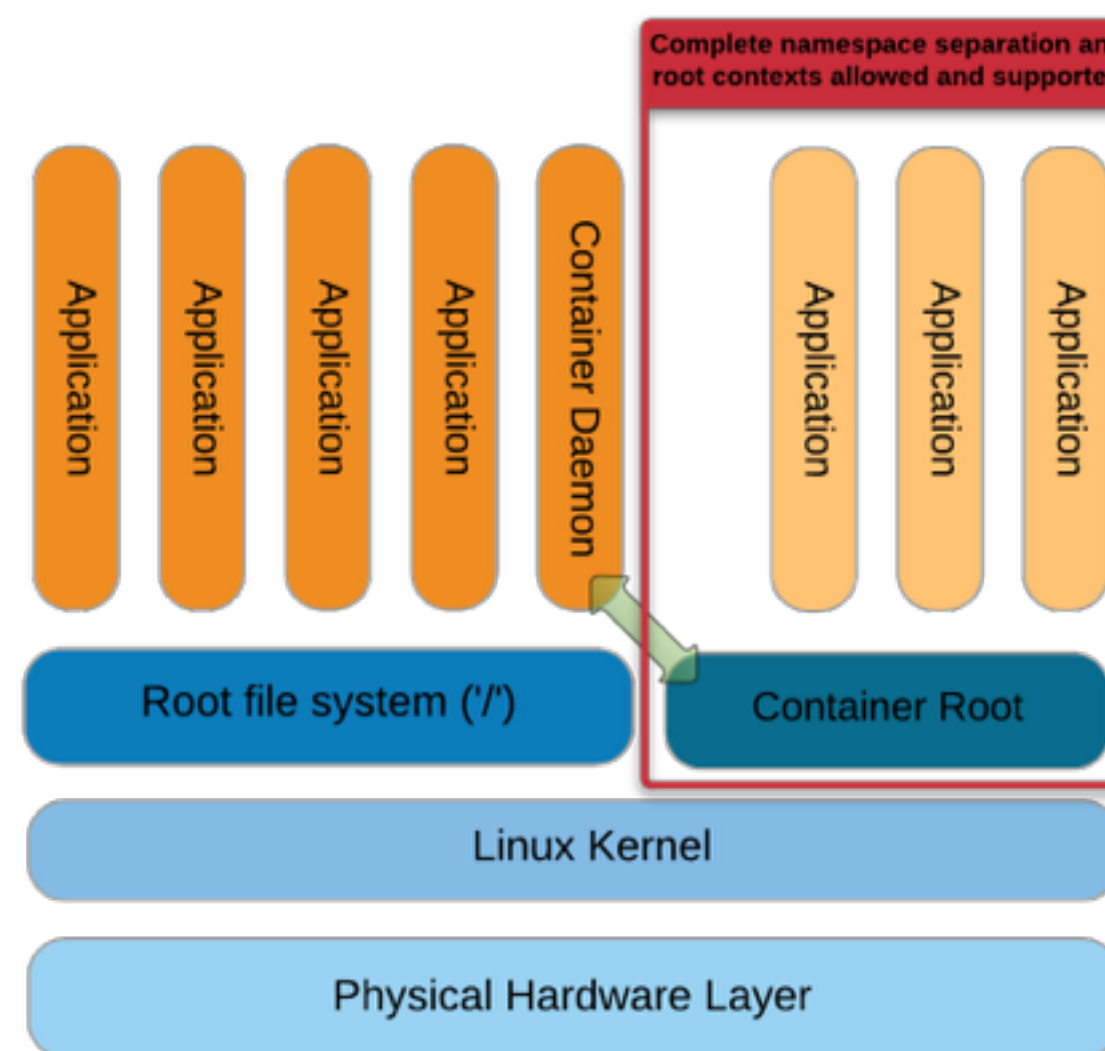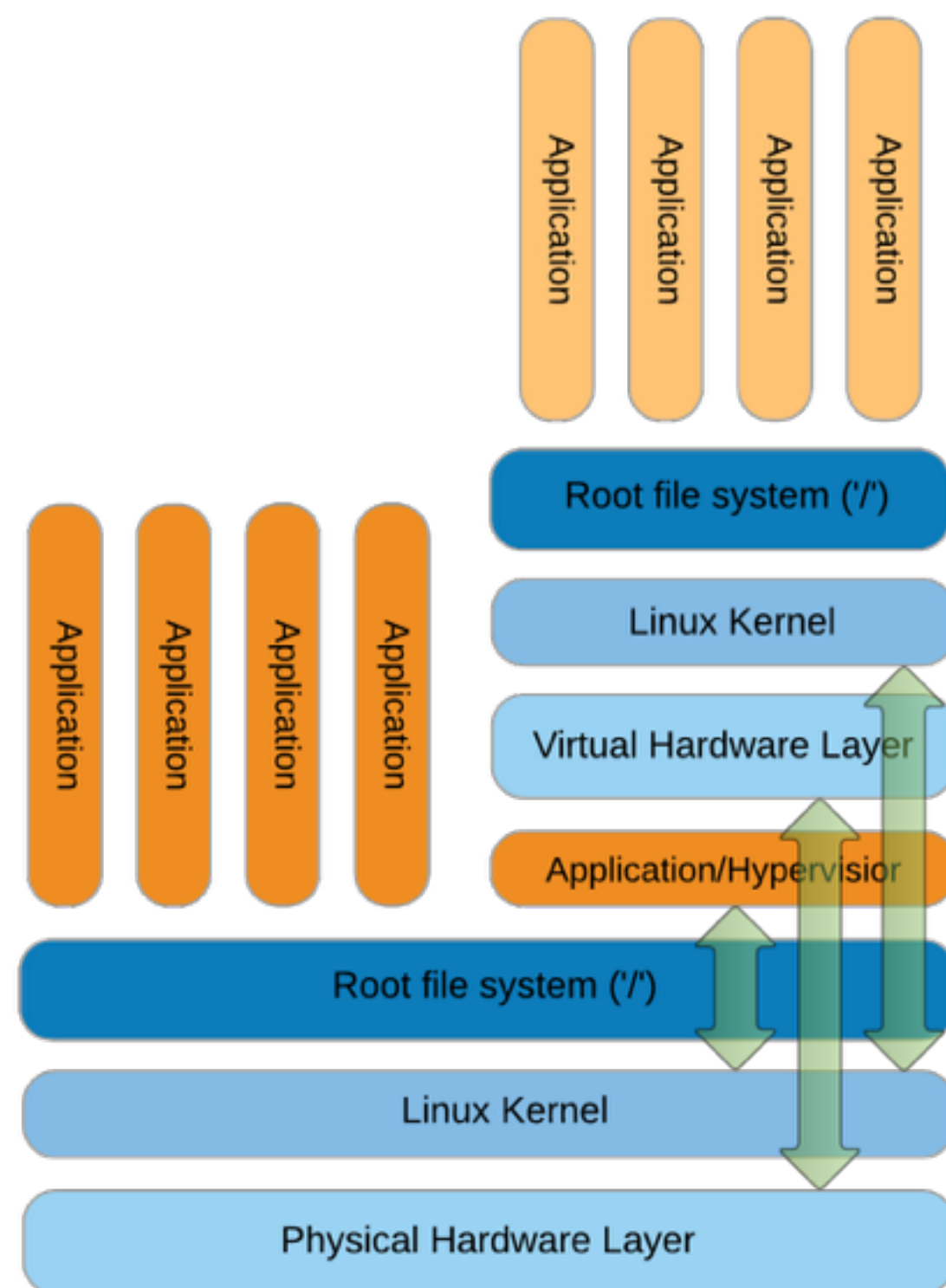▷   …not that easy

## Singularity

▷ User-land container, leveraging some Namespaces

▷ creates an executable
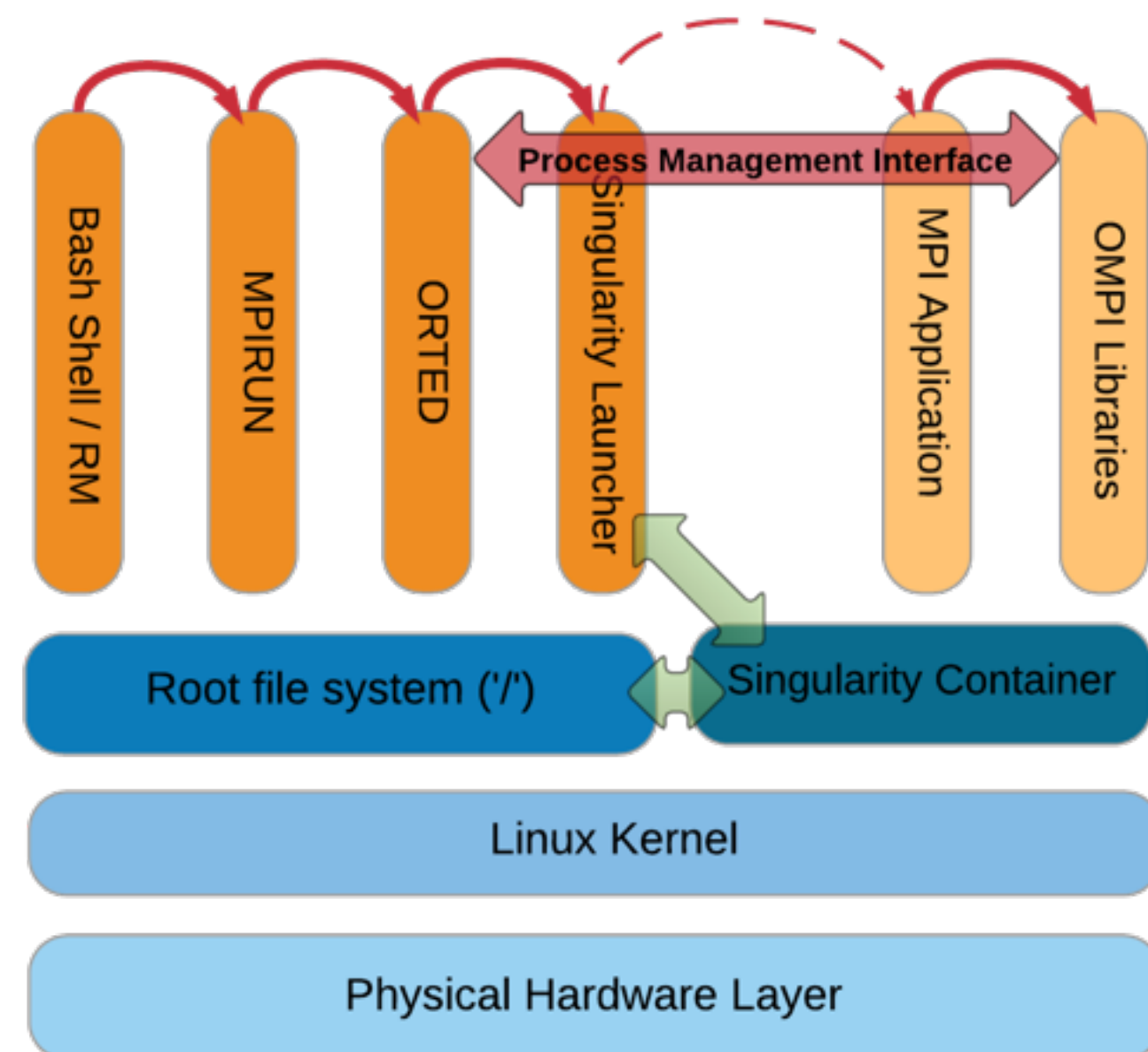
# MPI Singularity

Container Tech creates new NS

- ▷ singularity start from were the user is and trims down from there

- ▷ by doing so it put a barrier on-top of what a user can do

# MPI Singularity

Singularity is not bound to a daemon w/ API calls

▷ thus, it integrates (w/ latest) OpenMPI

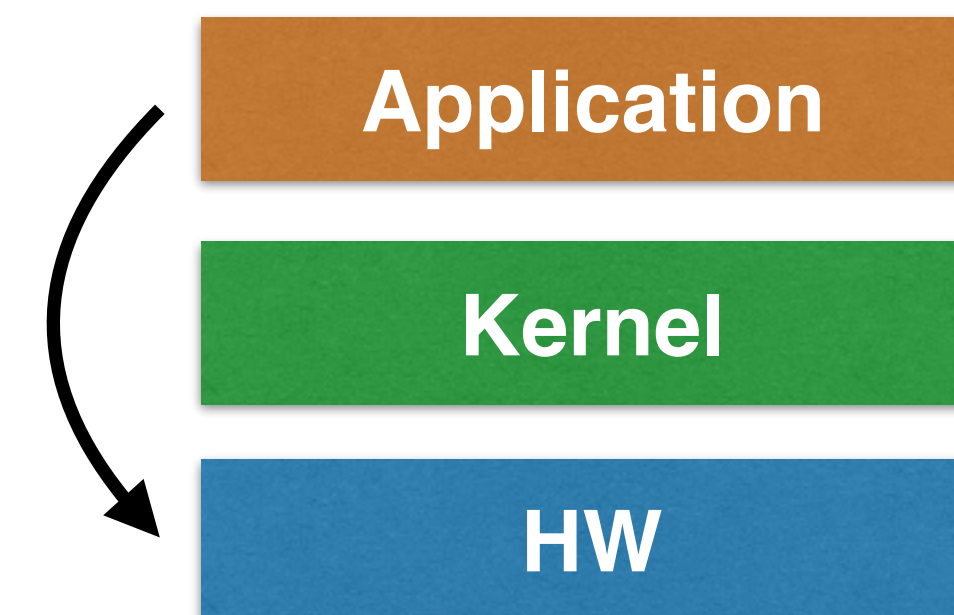▷ comply with the current workflow in HPC w/o changes

**Demo**

# RDMA Namespace

# Kernel By-passing

Perfect for single tenancy

- application talks directly to hardware

- nothing in it's way
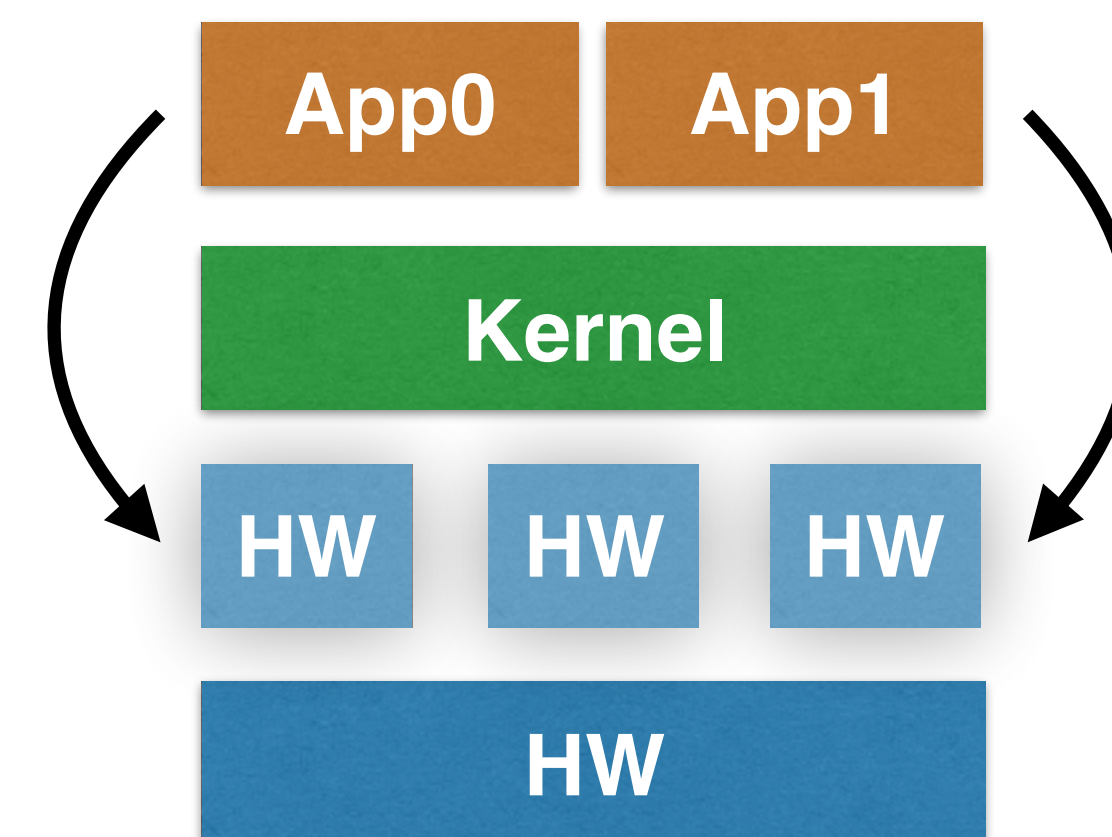
1. In VM environment we introduced SR-IOV

- hardware is exposed as multiple devices within system

- multiplexing is done in HW

2. Linux containers are extremely volatile

- if only we could leverage the kernel again,
  as it's aware of containers

# RDMA Namespace

Mellanox pushes code for RDMA

- Namespace

- CGroups

1. By doing so, the Kernel is in charge again

- knows about the resource needs of the all processes

# RDMA Namespace #2

## STATUS

- **InfiniBand RDMA CM support in v4.4**
- **RDMA cgroup patches submitted**
- **RDMA cgroup Docker patches are ready, will be submitted once kernel patches are accepted**
- **Working on RoCE net namespace support**
- **Future work**
  - InfiniBand: limit P_Key usage in verbs applications
    - Perhaps extend the RDMA cgroup
  - QoS: limit container's bandwidth usage, SL, or VLAN priority
  - Raw Ethernet support

# SLURM / MPI

Fake a ssh-client

```
[bob@a7b1e6e98cb1 ~]$ cat /opt/qnib/src/dssh
#!/bin/bash

REMOTE_HOST=$1
shift
set -x
docker -H unix:///var/run/docker.sock exec -i -u ${USER} ${REMOTE_HOST} $@
[bob@a7b1e6e98cb1 ~]$
```

▷ connects via **docker exec** instead of **ssh**

# SLURM needs patching

SLURM uses ***slurmstepd*** to spawn children

# If only

1. Slurmctld should use docker exec instead of fork the process directly (when in docker-mode).

2. MPI could also use docker exec to introduce process on remote system.