August 22, 2023

**Matthew Sgambati**

# Software Quality Assurance for High Performance Computing Containers utilizing MVAPICH2 and the MOOSE framework
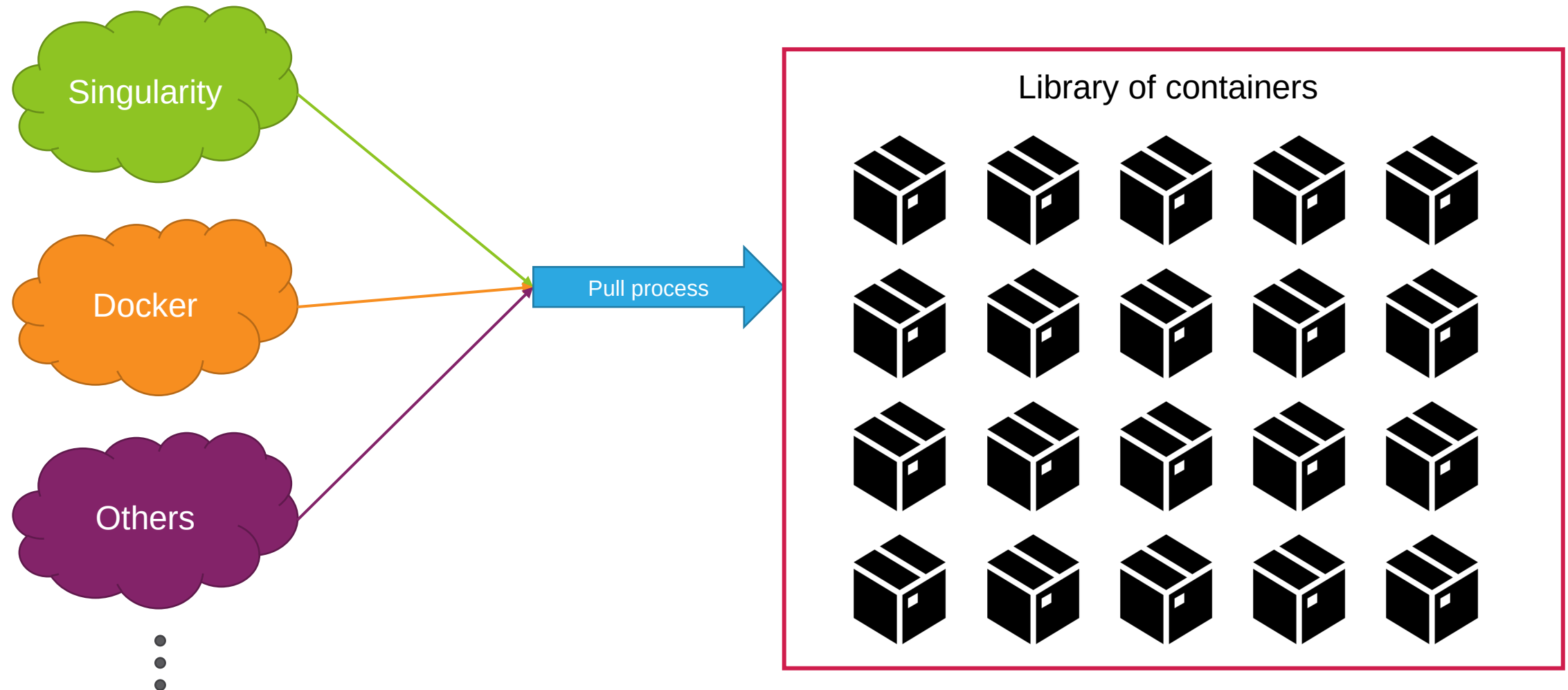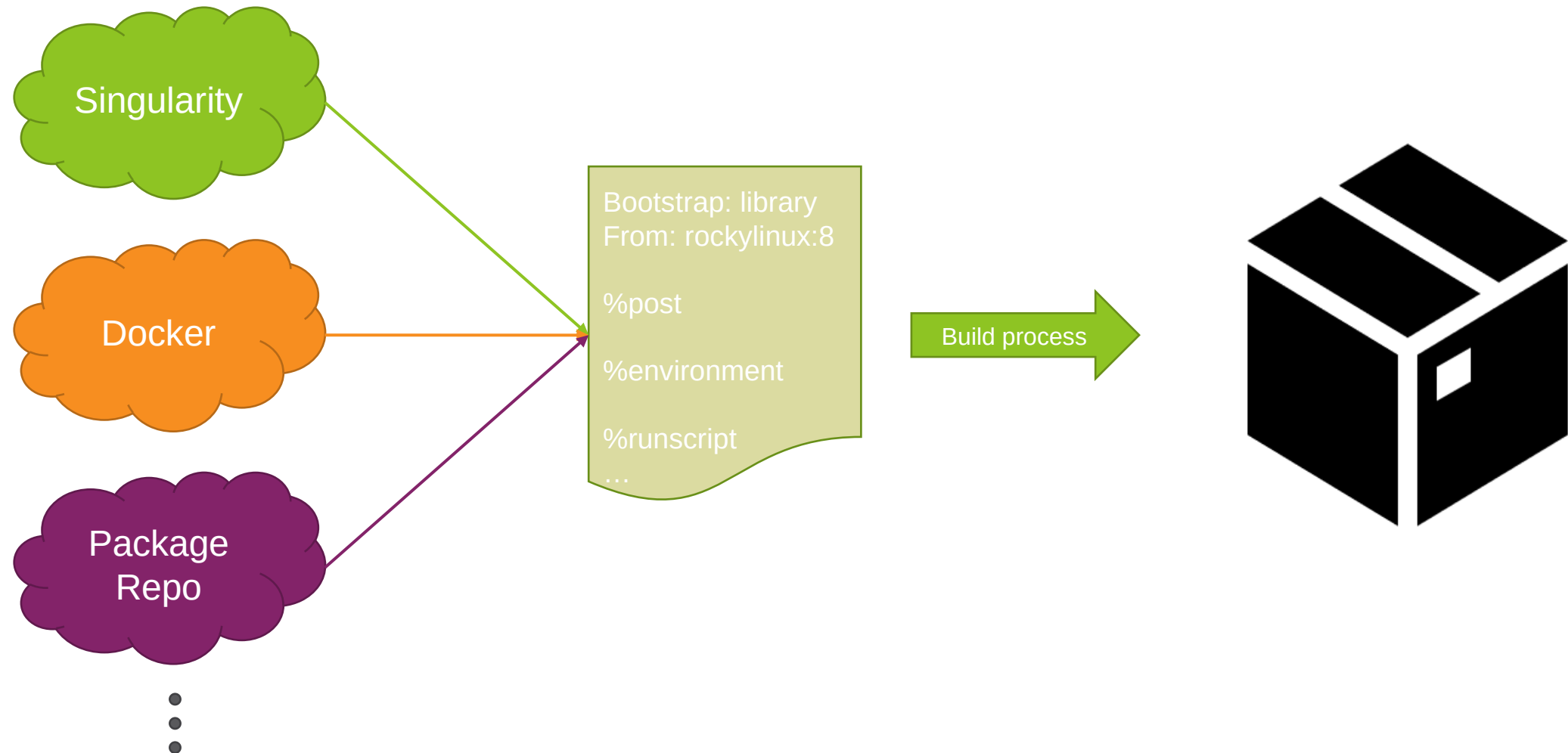
MUG 2023

INL/CON-23-74293

# Why containers?

- Verification and Validation:  Is my code running correctly still?  Is it passing all unit tests and giving the same results as my last few papers?

- Reproducibility:  Can I run my same application 10 years later?  Can I rebuild my container 10 years later?

- Portability across multiple supercomputer architectures:  Can I run the same container on multiple supercomputers with different networks and drivers *without rebuilding and without sacrificing performance*?

- Simplification: one environment, customized packages (independence from system administrators); make it easier on the user:  Can I work without needing intervention from the system administrator?

- Security (insulation against package versioning issues):  Does a security update on a supercomputer break my code?
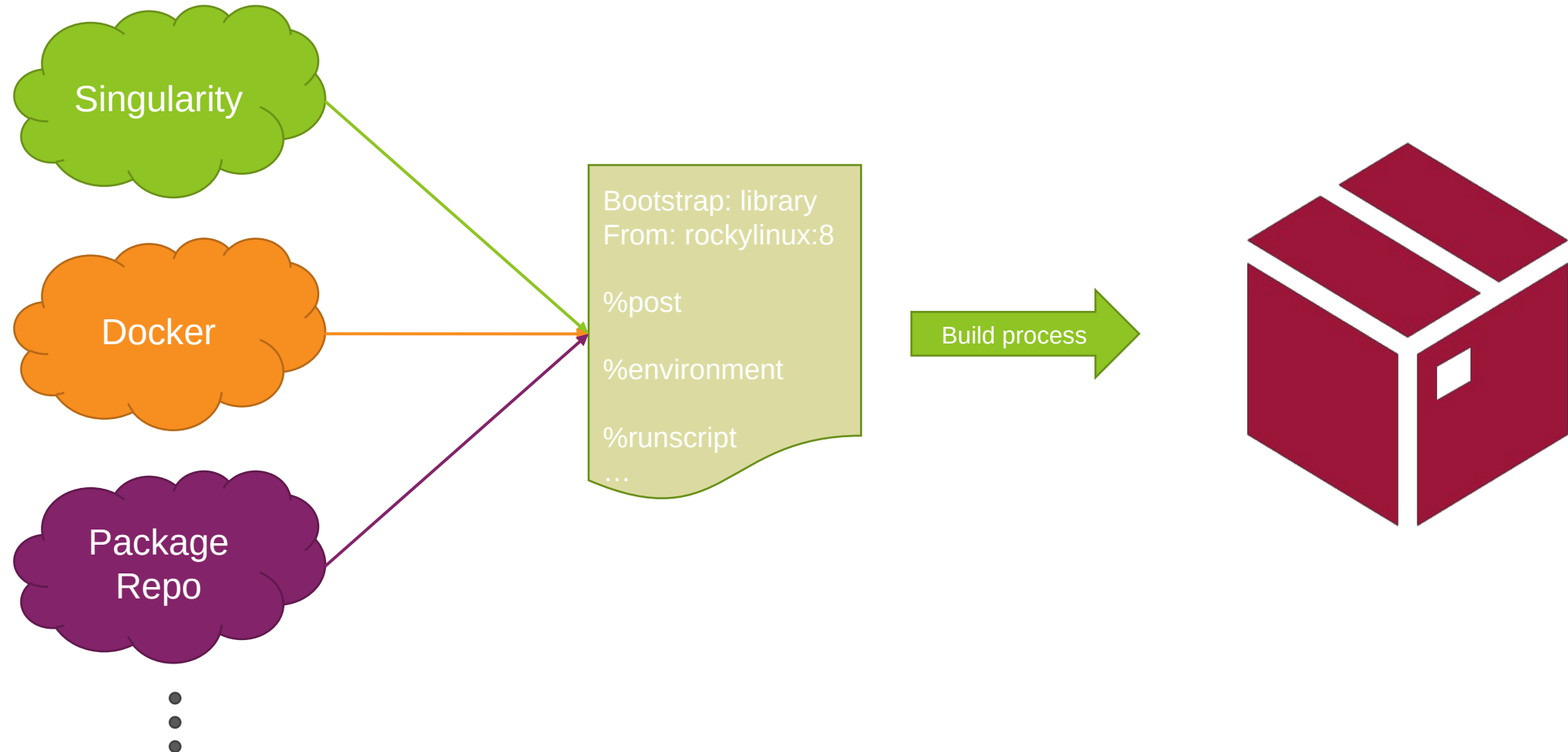
# Containers – Creation

Singularity

Docker

Others

Pull process

Library of containers

# Create a container with def file

Singularity

Docker

Package
Repo

Bootstrap: library
From: rockylinux:8
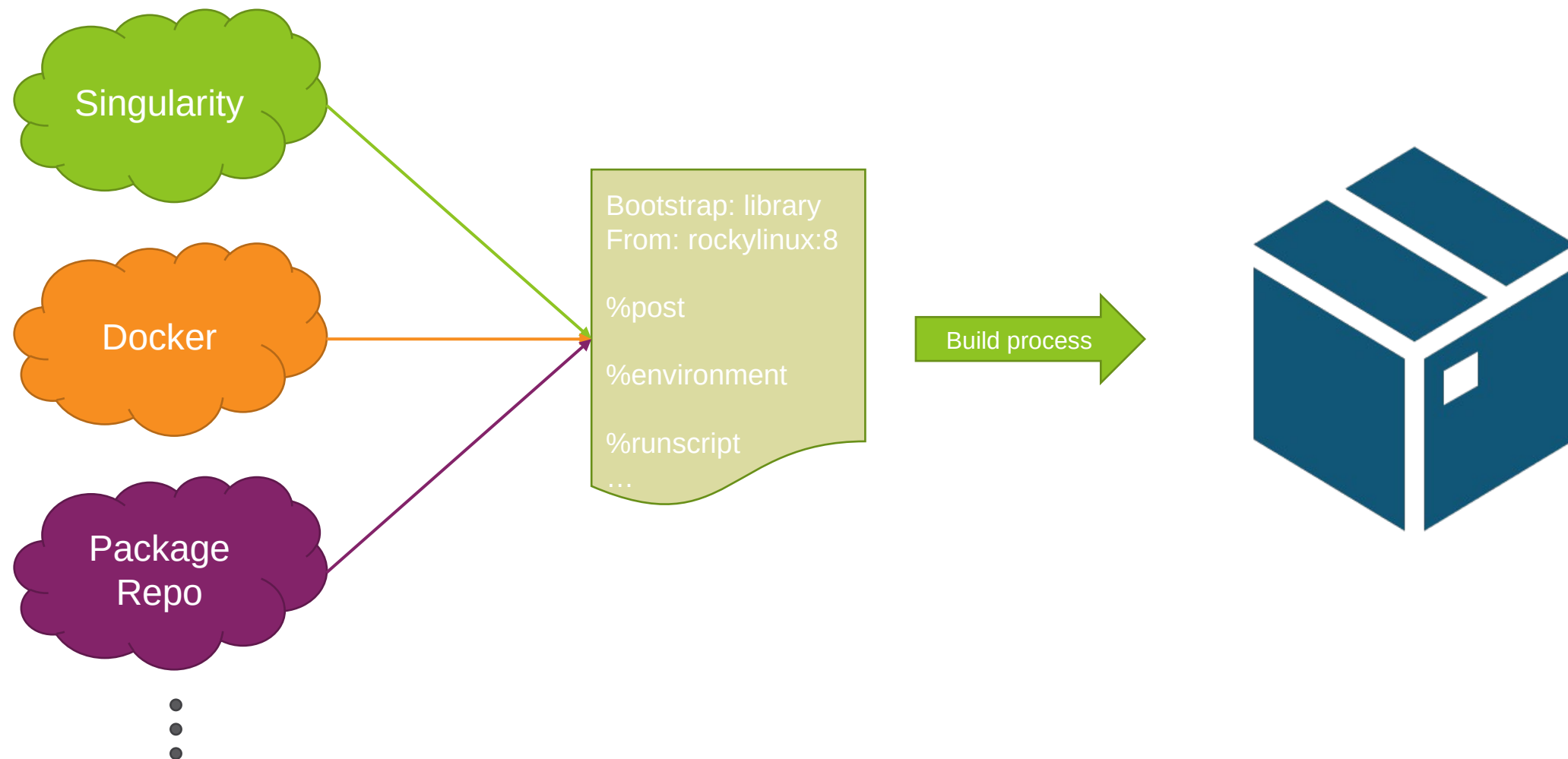
%post

%environment

%runscript
…

Build process

# Create a container – Reproduceable?

# Create a container – Reproduceable?
# Layer system

# Create a container – Reproduceable?
# Layer system



Bootstrap: localimage
From: <path_to_image>

%post

%environment

%runscript
…

Build process

# Create a container – Reproduceable?
# Layer system



Bootstrap: localimage
From: <path_to_image>

%post

%environment

%runscript
…

Build process

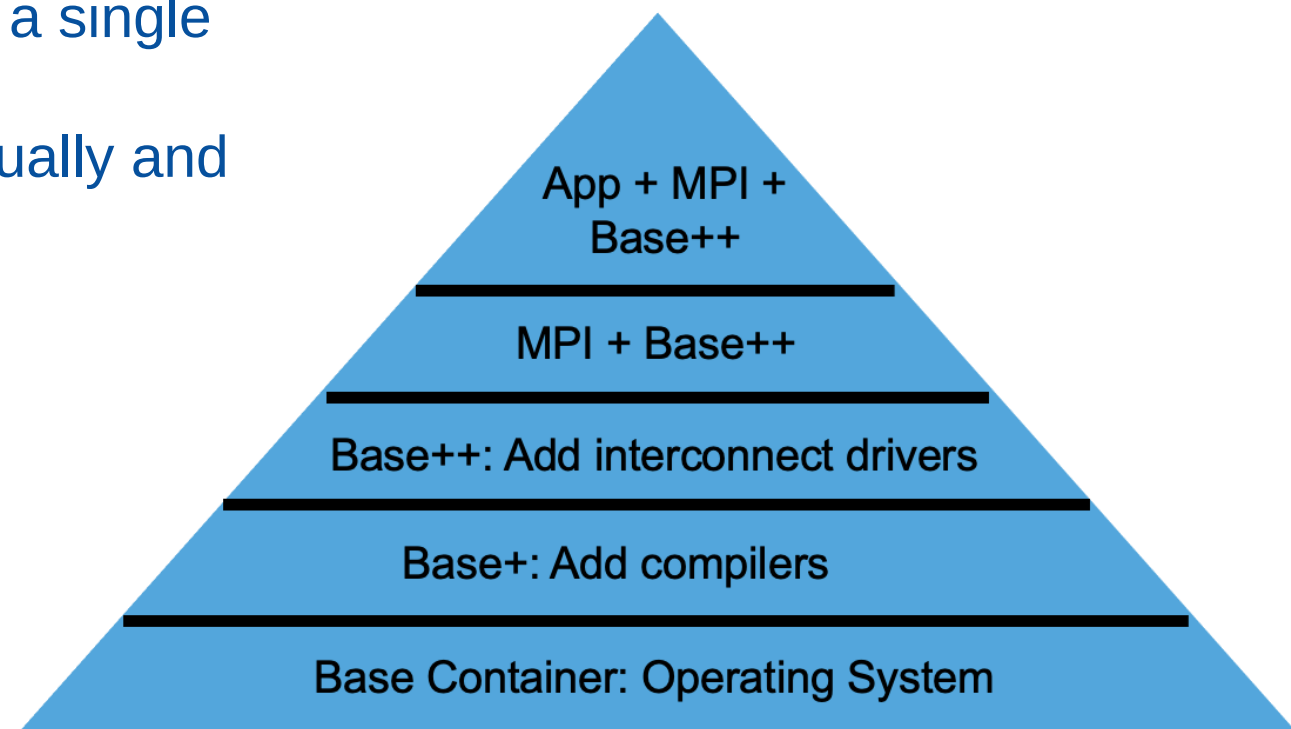# Containers with MPI and IB – Two Approaches

- Hybrid
  - System MPI
  - Container MPI
  - System IB
  - Container IB

- Bind
  - System MPI
  - System IB

NOTE:

**This requires the Host OS and Container OS to be compatible**

# Container Strategy

- There are three key components to the strategy:
  - HPC application containers built via layers
  - Key components are grouped into a single layer
  - Each layer can be updated individually and reproducibly via local mirrors

App + MPI + Base++

MPI + Base++

Base++: Add interconnect drivers

Base+: Add compilers

Base Container: Operating System

# Container Strategy – Assumptions

- Host systems will use a long-term support (LTS) version of drivers and/or software stacks when possible

- Host system administrators will install an ABI-compatible version of MPI if one does not exist on the target system

# Base Container Definition File

```
1    Bootstrap: docker
2    From: rockylinux:8.6
3
4    %post
5        # Capture useful system information
6        echo "Architecture" $(lscpu | grep "^Architecture" | cut -d ':' -f 2) >> "${SINGULARITY_LABELS}"
7        echo "CPU" $(lscpu | grep "^Model name" | cut -d ':' -f 2) >> "${SINGULARITY_LABELS}"
8        echo "uname" $(uname -srvpio) >> "${SINGULARITY_LABELS}"
9
10   %test
11       if grep -q 'NAME="Rocky Linux"' /etc/os-release; then
12           echo "SUCCESS: Container base is Rocky Linux as expected."
13       else
14           echo "ERROR: Container base is not Rocky Linux."
15           exit 1
16       fi
17
18   %labels
19       Authors Matthew.Sgambati@inl.gov Matthew.Anderson2@inl.gov
20       Version 1.0.0
21
22   %help
23       Rocky Linux 8.6 Base Container
24
```

# Base+ Container Definition File


Base+: Add compilers
Base Container: Operating System

```
1   Bootstrap: oras
2   From: <container_registry>/hpcbase/base_00:1.0.0
3
4   %post
5       # Change repos to point to local static mirror
6       sed -i 's/^mirrorlist/#mirrorlist/' /etc/yum.repos.d/Rocky-*.repo
7       sed -i 's#.*baseurl=http://dl.rockylinux.org/\$contentdir#baseurl=http://<local_static_mirror>/repos/rocky-linux/20221208#'
        ↪  /etc/yum.repos.d/Rocky-*.repo
8
9       dnf clean all
10      dnf makecache
11
12      # Install commonly used packages for building code and modifiying files
13      dnf install -y bzip2 gcc gcc-gfortran gcc-c++ gdb git make python39 python39-pip python39-setuptools tar vim
14      dnf clean all
15
16      # Capture useful system information
17      echo "Architecture" $(lscpu | grep "^Architecture" | cut -d ':' -f 2) >> "${SINGULARITY_LABELS}"
18      echo "CPU" $(lscpu | grep "^Model name" | cut -d ':' -f 2) >> "${SINGULARITY_LABELS}"
19      echo "uname" $(uname -srvpio) >> "${SINGULARITY_LABELS}"
20
21  %test
22      GCC=$(which gcc)
23      if [ $? -eq 0 ]; then
24          echo "SUCCESS: gcc is available at ${GCC}"
25      else
26          echo "ERROR: gcc is not installed"
27          exit 1
28      fi
29
30  %labels
31      Authors Matthew.Sgambati@inl.gov Matthew.Anderson2@inl.gov
32      Version 1.0.0
33
34  %help
35      Rocky Linux 8.6 Base Container
36      Extra dependencies for building code and modifiying files are installed in this layer.
37
```
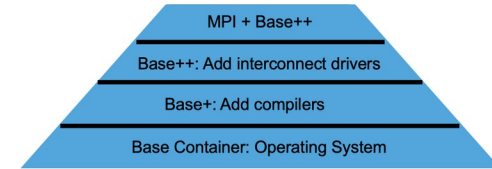
# Base++ Container Definition File



```
1   Bootstrap: oras
2   From: <container_registry>/hpcbase/extra_01:1.0.0
3
4   %post
5       # Create environment variables of software versions
6       MLNX_OFED=MLNX_OFED_LINUX-5.4-3.6.8.1-rhel8.6-x86_64
7       OPX=CornelisOPX-Basic.RHEL86-x86_64.10.12.1.0.7
8
9       # Make src in /opt to hold source code
10      mkdir -p /opt/src
11
12      # Download the MLNX_OFED and OPX files to /opt
13      cd /opt
14      curl -O http://<local_static_mirror>/interconnects/MLNX/20221208/${MLNX_OFED}.tgz
15      curl -O http://<local_static_mirror>/interconnects/OPX/20230212/${OPX}.tgz
16
17      # Install required packages for InfiniBand
18      dnf install -y python39 pciutils lsof ethtool tcsh libnl3 tk numactl-libs tcl
19
20      # Extract MLNX_OFED and install it
21      tar xf ${MLNX_OFED}.tgz -C /opt/src
22      /opt/src/${MLNX_OFED}/mlnxofedinstall --without-fw-update --skip-unsupported-devices-check --basic --user-space-only --distro RHEL8.6
        ↪  --without-depcheck -q
23
24      # Install required packages for OmniPath
25      dnf install -y irqbalance kernel-modules-extra kmod libgcc perl perl-Getopt-Long perl-Socket opensm-libs python2 libatomic
26      dnf download ibacm*x86_64
27      rpm -ivh --nodeps ibacm*.rpm
28
29      # Extract OPX and install it
30      tar xf ${OPX}.tgz -C /opt/src
31      cd /opt/src/${OPX}
32      ./INSTALL -i intel_hfi -i opa_stack --user-space
33
34      # Clean up tarballs/downloads and source directories
35      cd /opt
36      rm -rf /opt/src
37      rm -f /opt/*.tgz
38      rm -f /opt/*.rpm
39
40      # Capture useful system information
41      echo "Architecture" $(lscpu | grep "^Architecture" | cut -d ':' -f 2) >> "${SINGULARITY_LABELS}"
42      echo "CPU" $(lscpu | grep "^Model name" | cut -d ':' -f 2) >> "${SINGULARITY_LABELS}"
43      echo "uname" $(uname -srvpio) >> "${SINGULARITY_LABELS}"
44
```

# MPI + Base++ Container Definition File



```
1   Bootstrap: oras
2   From: <container_registry>/hpcbase/interconnect_02:1.0.0
3
4   %post
5       # Install required packages
6       dnf install -y hwloc findutils
7
8       # Create directories to store files and source code
9       mkdir -p /opt/mpi/examples
10      mkdir -p /opt/src
11      mkdir -p /opt/tars
12
13      # Create mpitest.c program from here doc
14      cat <<- EOF > /opt/mpi/examples/mpitest.c
15          #include <mpi.h>
16          #include <stdio.h>
17          #include <stdlib.h>
18
19          int main (int argc, char **argv) {
20                  int rc;
21                  int size;
22                  int myrank;
23
24                  rc = MPI_Init (&argc, &argv);
25                  if (rc != MPI_SUCCESS) {
26                      fprintf (stderr, "MPI_Init() failed");
27                      return EXIT_FAILURE;
28                  }
29
30                  rc = MPI_Comm_size (MPI_COMM_WORLD, &size);
31                  if (rc != MPI_SUCCESS) {
32                      fprintf (stderr, "MPI_Comm_size() failed");
33                      goto exit_with_error;
34                  }
35
36                  rc = MPI_Comm_rank (MPI_COMM_WORLD, &myrank);
37                  if (rc != MPI_SUCCESS) {
38                      fprintf (stderr, "MPI_Comm_rank() failed");
39                      goto exit_with_error;
40                  }
41
42                  fprintf (stdout, "Hello, I am rank %d/%d\n", myrank, size);
43
44                  MPI_Finalize();
45
46                  return EXIT_SUCCESS;
```

```
76      # Install MPICH
77      MPICH_VERSION=3.4.3
78      MPICH_NAME="mpich-${MPICH_VERSION}"
79      MPICH_URL="http://<local_static_mirror>/mpi/mpich/20221208/${MPICH_NAME}.tar.gz"
80      MPICH_DIR="/opt/mpi/${MPICH_NAME}"
81
82      echo "Installing MPICH-${MPICH_VERSION}..."
83      ## Download
84      cd /opt/tars
85      curl -O ${MPICH_URL}
86      tar xf ${MPICH_NAME}.tar.gz -C /opt/src
87
88      ## Compile and install
89      cd /opt/src/${MPICH_NAME}
90      mkdir build
91      cd build
92      ../configure --prefix=${MPICH_DIR} --with-ucx=${UCX_DIR}
93      make -j 16 |& tee log.make
94      make check |& tee log.make_check
95      make install |& tee log.make_install
96
97      # Install OpenMPI
98      OPENMPI_VERSION=4.1.4
99      OPENMPI_NAME="openmpi-${OPENMPI_VERSION}"
100     OPENMPI_DIR="/opt/mpi/${OPENMPI_NAME}"
101     OPENMPI_URL="http://<local_static_mirror>/mpi/openmpi/20221208/${OPENMPI_NAME}.tar.gz"
102
103     echo "Installing OPENMPI-${OPENMPI_VERSION}..."
104     ## Download
105     cd /opt/tars
106     curl -O ${OPENMPI_URL}
107     tar xf ${OPENMPI_NAME}.tar.gz -C /opt/src
108
109     ## Compile and install
110     cd /opt/src/${OPENMPI_NAME}
111     mkdir build
112     cd build
113     ../configure --prefix=${OPENMPI_DIR} --with-ucx=${UCX_DIR}
114     make -j 16 |& tee log.make
115     make check |& tee log.make_check
116     make install |& tee log.make_install
```
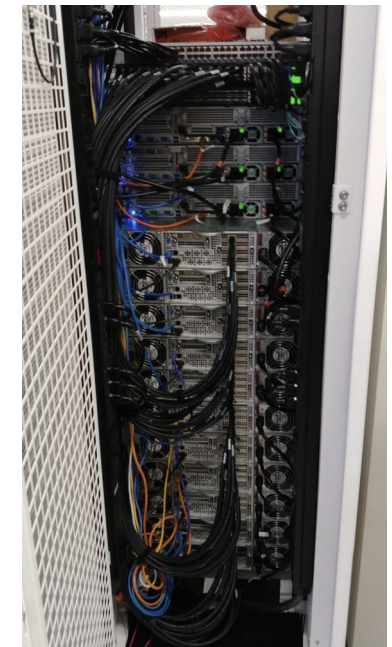
# Container Systems for Testing

| System Name | Core Count | Chipset | Interconnect / Version | OS |
|---|---|---|---|---|
| Sawtooth[1,2] | 99,792 | Intel Xeon 8268 | InfiniBand EDR / 4.9-4.1.7 | CentOS Linux release 7.9.2009 (Core) |
| Lemhi[1,2] | 20,160 | Intel Xeon 6148 | OmniPath / 10.11.0.2-1 | Rocky Linux release 8.7 (Green Obsidian) |
| Hoodoo[1,2] | 352 | AMD EPYC 7302 | InfiniBand HDR / 5.5-1.0.3 | Rocky Linux release 8.5 (Green Obsidian) |
| Galena[1] | 40 | Intel Xeon E5-2698 | InfiniBand EDR / 5.4-3.5.8 | Ubuntu 20.04.5 LTS (Focal Fossa) |



Sawtooth

Lemhi

Hoodoo

Galena

# LULESH – Sawtooth



LULESH Weak Scaling: 1000 iterations

OpenMPI 4.1.4

MVAPICH2 2.3.5

# LULESH – Lemhi and Hoodoo



Lemhi
OpenMPI 4.1.1

Hoodoo
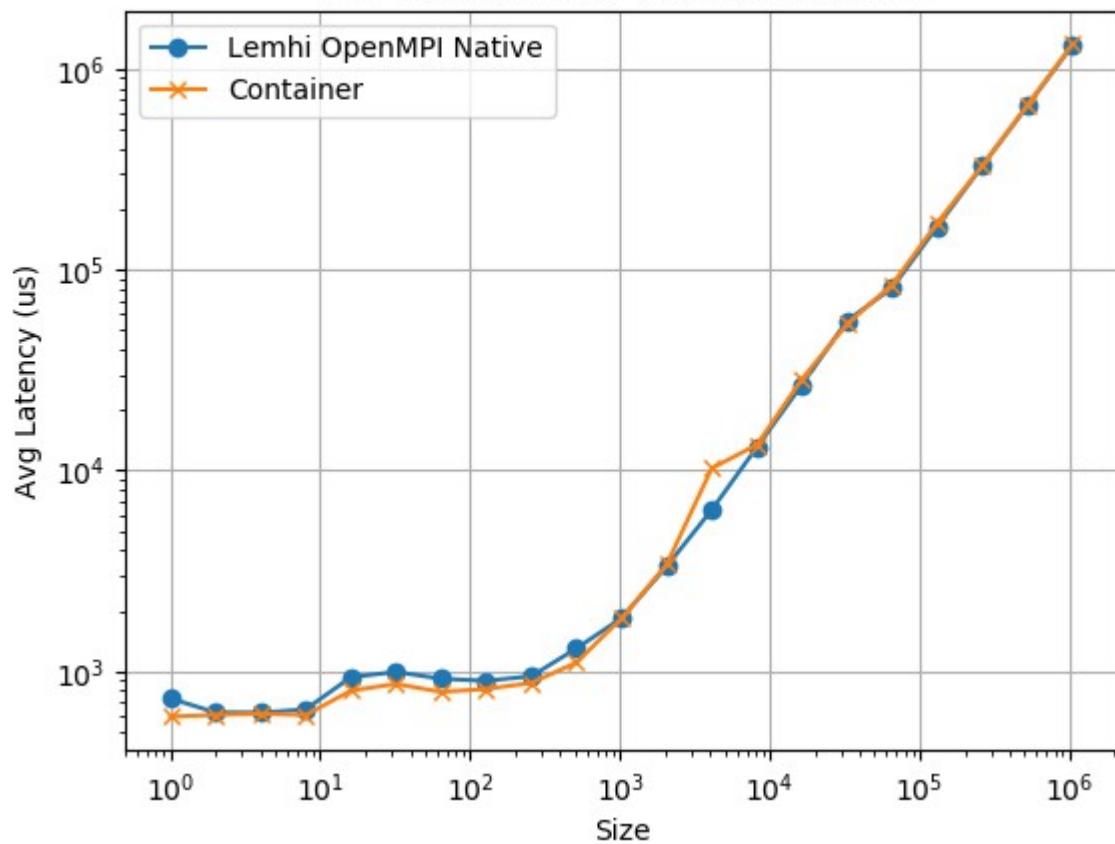OpenMPI 4.0.5

# OSU MPI All-to-Allv 10 Nodes – Sawtooth



OpenMPI 4.1.4

MVAPICH2 2.3.5

# OSU MPI All-to-Allv 10 Nodes – Lemhi and Hoodoo
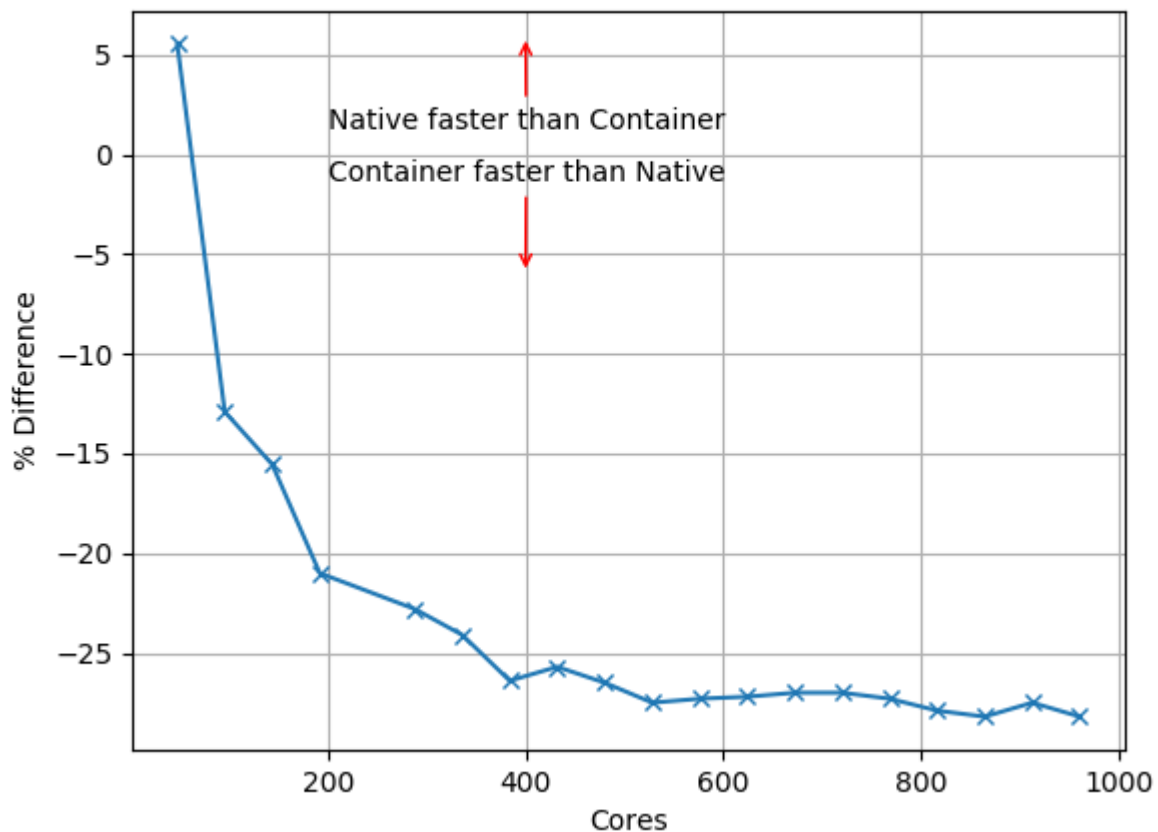


Lemhi
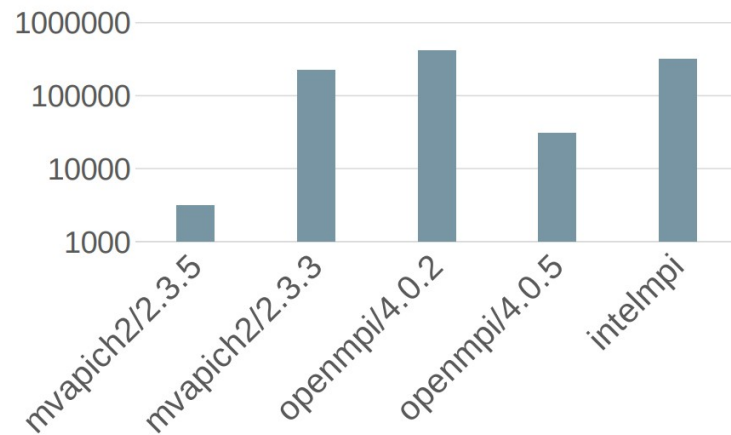OpenMPI 4.1.1

Hoodoo
OpenMPI 4.0.5

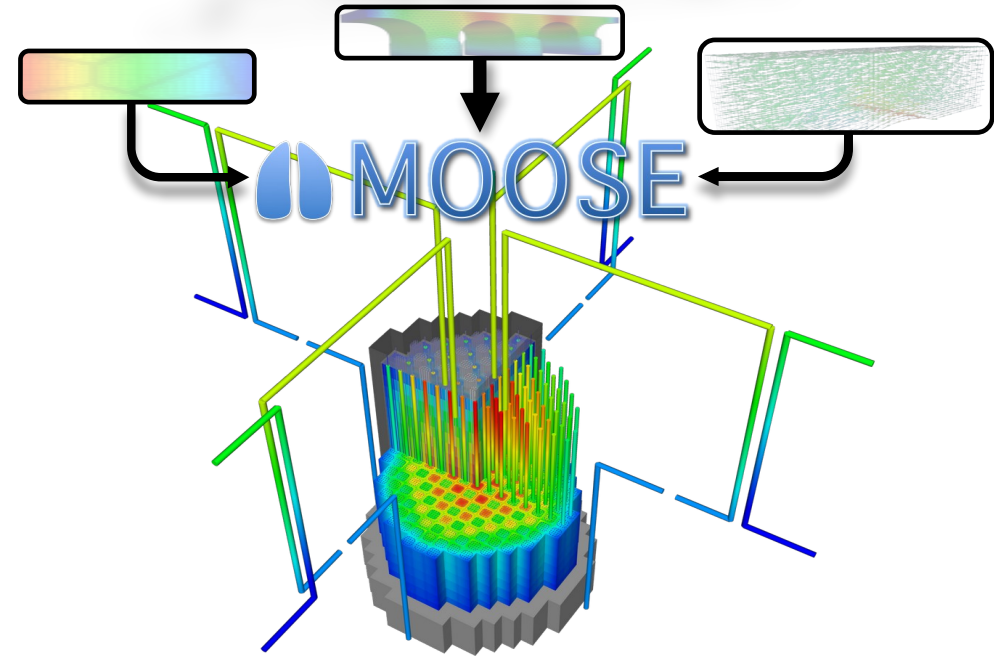# MCNP – Container vs Native Performance
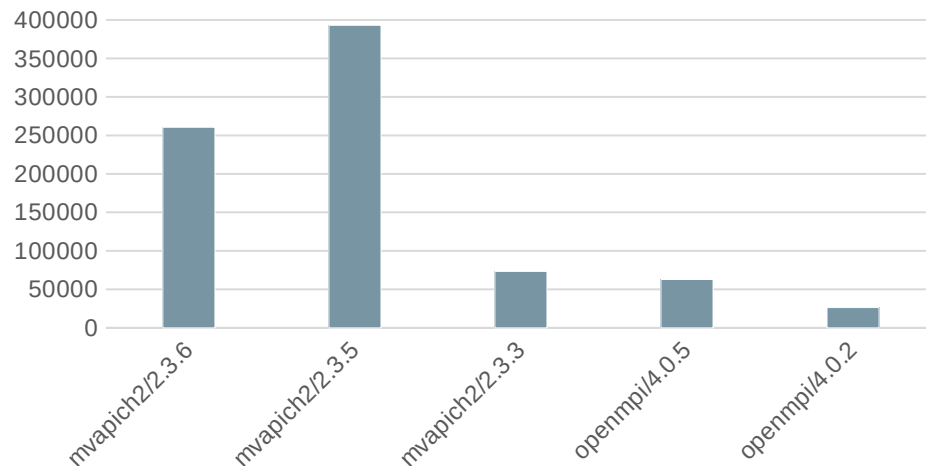


Sawtooth
OpenMPI 4.1.4

Lemhi
OpenMPI 4.1.1

# MVAPICH2 and MOOSE

MPI library usage 1 Jan 2021—1 Sep 2021



MPI Library usage since 1 Jan 2022





Multiphysics Object-Oriented Simulation Environment

- Open-source finite element framework
  - Basis for most nuclear modeling and simulation tools developed at INL

# MOOSE Container layers – outcome

```
28  BootStrap: {{ APPTAINER_BOOTSTRAP }}
29  From: {{ APPTAINER_FROM }}
30
31  %setup
32      # Load jinja vars
33      ROOT_BUILD_DIR={{ ROOT_BUILD_DIR }}
34      APPLICATION_DIR={{ APPLICATION_DIR }}
35      MOOSE_DIR={{ MOOSE_DIR }}
36
37      # Build directory in the container
38      BUILD_DIR=${APPTAINER_ROOTFS}${ROOT_BUILD_DIR}
39      mkdir ${BUILD_DIR}
40
41      # Copy application into the container
42      cp -r ${APPLICATION_DIR} ${BUILD_DIR}
43
44      # Where the application ends up; needed for MOOSE logic below
45      APPLICATION_NAME=$(basename ${APPLICATION_DIR})
46      APPLICATION_BUILD_DIR=${BUILD_DIR}/${APPLICATION_NAME}
47
48      # Figure out where moose is; regardless %post will expect
49      # it to be in {{ ROOT_BUILD_DIR }}/moose
50      MOOSE_BUILD_DIR=${BUILD_DIR}/moose
51      MOOSE_RELATIVE_PATH=$(realpath --relative-to ${APPLICATION_DIR} ${MOOSE_DIR})
52      # MOOSE_DIR is not in the application; we need to copy it
53      if [[ $MOOSE_RELATIVE_PATH = ..* ]]; then
54          mkdir ${MOOSE_BUILD_DIR}
55          cp -r ${MOOSE_DIR}/. ${MOOSE_BUILD_DIR}
56      # MOOSE_DIR is the application (combined-opt)
57      elif [[ '{{ BINARY_NAME }}' == 'moose-combined' ]]; then
58          # do nothing
59          :
60      # MOOSE_DIR is in the application, setup a symlink
61      else
62          ln -s ./${APPLICATION_NAME}/${MOOSE_RELATIVE_PATH} ${MOOSE_BUILD_DIR}
63      fi
64
65  {%- if MOOSE_SKIP_DOCS is not defined %}
66      # Need large_media for documentation
67      cd ${MOOSE_BUILD_DIR}
68      git submodule update --init large_media
69  {%- endif %}
```

- Utilizes the "MPI + Base++" container as a base for continuous integration (non-HPC) and most HPC execution
- Improved reproducibility and portability
- Integrated into CI/CD for building workflows
  - Reduced job builds to 20 per week, which are all automated
- Simplified build process

# Conclusion

- This strategy has shown the following:
  - Portability
    - Across multiple supercomputer architectures
  - Reproducibility
    - Due to layers and local mirrors
  - Traceability
    - Def files only have most recent changes
  - Simplification
    - HPC staff only needs to focus on MPI ABI and interconnect series compatibility
  - Security
    - Insulation against host system updates

# Questions?

*Battelle Energy Alliance manages INL for the U.S. Department of Energy's Office of Nuclear Energy. INL is the nation's center for nuclear energy research and development, and also performs research in each of DOE's strategic goal areas: energy, national security, science and the environment.*

WWW.INL