



Updates and Innovations with the Apptainer Platform

Forrest Burt - Solutions Architect, CIQ

February 3, 2024

FOSDEM 2024

Brussels, Belgium

Developments in Apptainer over the past few years

- Leveraging the User Namespace
- New recommendations for containerized MPI
- Increased adoption of ORAS (for Software Supply Chain)

Leveraging the User Namespace

Namespace types

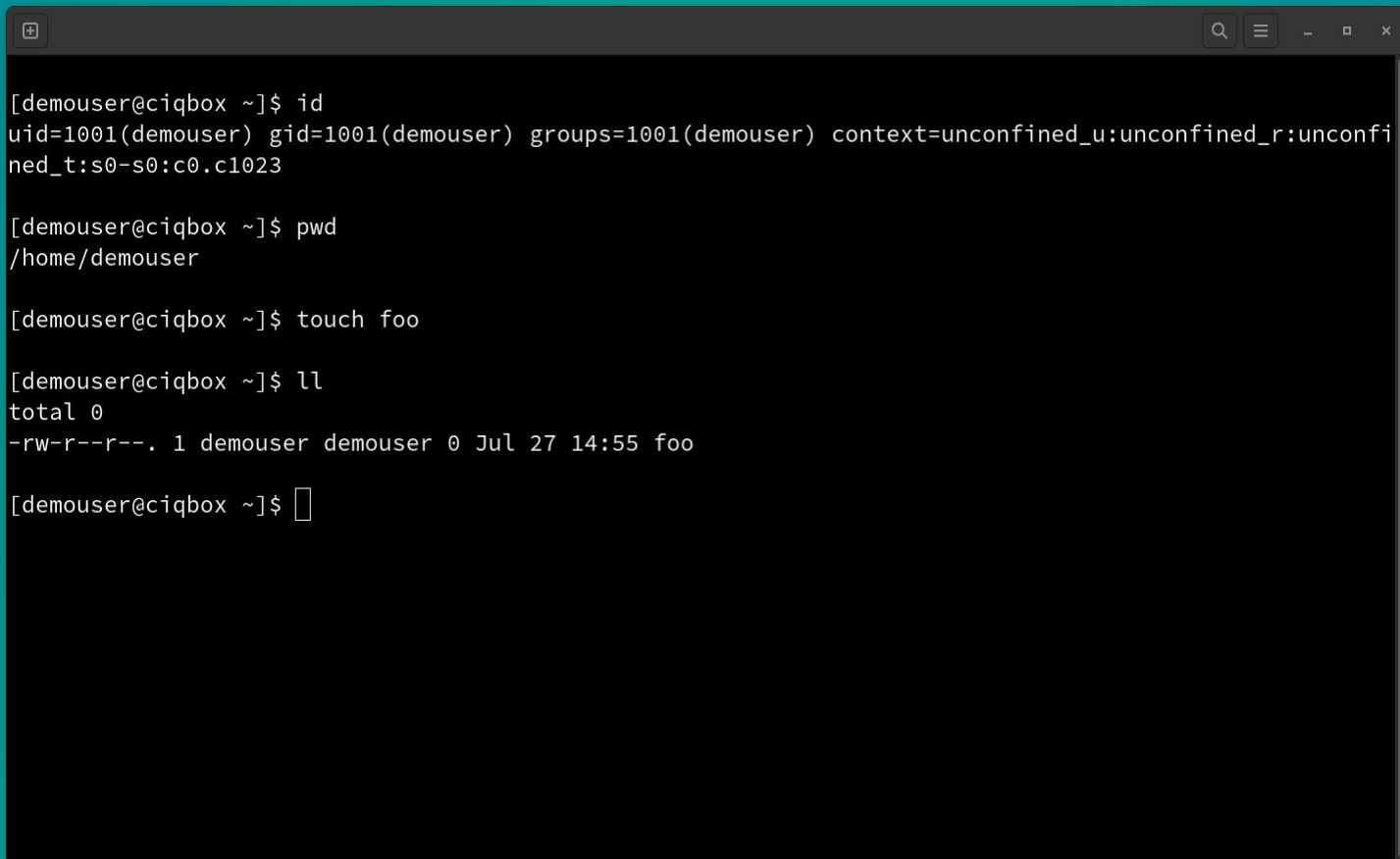
The following table shows the namespace types available on Linux. The second column of the table shows the various APIs. The third column identifies the manual page that provides details on the namespace type. The fourth column identifies the namespace type.

| Namespace | Flag | Page | Isolates |
|-----------|------------------------------|------------------------------------|--------------------------------------|
| Cgroup | <code>CLONE_NEWCGROUP</code> | <code>cgroup_namespaces(7)</code> | Cgroup root directory |
| IPC | <code>CLONE_NEWIPC</code> | <code>ipc_namespaces(7)</code> | System V IPC, POSIX message queues |
| Network | <code>CLONE_NEWNET</code> | <code>network_namespaces(7)</code> | Network devices, stacks, ports, etc. |
| Mount | <code>CLONE_NEWNS</code> | <code>mount_namespaces(7)</code> | Mount points |
| PID | <code>CLONE_NEWPID</code> | <code>pid_namespaces(7)</code> | Process IDs |
| Time | <code>CLONE_NEWTIME</code> | <code>time_namespaces(7)</code> | Boot and monotonic clocks |
| User | <code>CLONE_NEWUSER</code> | <code>user_namespaces(7)</code> | User and group IDs |
| UTS | <code>CLONE_NEWUTS</code> | <code>uts_namespaces(7)</code> | Hostname and NIS domain name |

The namespaces API

As well as various `/proc` files described below, the namespaces API includes the following system calls

Leveraging the User Namespace (for `--fakeroot`)

A terminal window with a dark background and light text. The window title bar shows standard Linux window controls (maximize, search, menu, close). The terminal output shows the user 'demouser' running several commands: 'id' to show their identity, 'pwd' to show their current directory, 'touch foo' to create a file, and 'll' to list the file's permissions and details.

```
[demouser@ciqbox ~]$ id
uid=1001(demouser) gid=1001(demouser) groups=1001(demouser) context=unconfined_u:unconfined_r:unconfi
ned_t:s0-s0:c0.c1023

[demouser@ciqbox ~]$ pwd
/home/demouser

[demouser@ciqbox ~]$ touch foo

[demouser@ciqbox ~]$ ll
total 0
-rw-r--r--. 1 demouser demouser 0 Jul 27 14:55 foo

[demouser@ciqbox ~]$
```

Leveraging the User Namespace (for `--fakeroot`)

```
[demouser@ciqbox ~]$ aptainer shell --fakeroot docker://alpine
Apptainer> id
uid=0(root) gid=0(root) groups=0(root)
Apptainer> pwd
/root
Apptainer> ls -l
total 0
-rw-r--r--  1 root    root      0 Jul 27 14:55 foo
Apptainer> touch bar
Apptainer> ls -l
total 0
-rw-r--r--  1 root    root      0 Jul 27 14:58 bar
-rw-r--r--  1 root    root      0 Jul 27 14:55 foo
Apptainer> exit

[demouser@ciqbox ~]$ ll
total 0
-rw-r--r--. 1 demouser demouser 0 Jul 27 14:58 bar
-rw-r--r--. 1 demouser demouser 0 Jul 27 14:55 foo

[demouser@ciqbox ~]$
```

Leveraging the User Namespace (for build)

```
[demouser@ciqbox ~]$ cat example.def
Bootstrap: docker
From: rockylinux:8

%post
    whoami
    dnf -y update
    dnf -y install vim
    # ^ note that these commands require privs

[demouser@ciqbox ~]$ aptainer build example.sif example.def
+ whoami
root
+ dnf -y update
Rocky Linux 8 - AppStream          594 kB/s | 10 MB      00:17
Rocky Linux 8 - BaseOS            460 kB/s | 4.9 MB    00:10
Rocky Linux 8 - Extras            40 kB/s | 13 kB      00:00
Dependencies resolved.
=====
Package                        Architecture  Version                                Repository  Size
=====
Upgrading:
krb5-libs                       x86_64       1.18.2-25.el8_8                       baseos      841 k
platform-python                 x86_64       3.6.8-51.el8_8.1.rocky.0             baseos      86 k
python3-libs                     x86_64       3.6.8-51.el8_8.1.rocky.0             baseos     7.8 M
```

Leveraging the User Namespace (for `build`)

```
Running transaction test
Transaction test succeeded.
Running transaction
  Preparing      :                                1/1
  Installing     : which-2.21-18.el8.x86_64      1/5
  Installing     : vim-filesystem-2:8.0.1763-19.el8_6.4.noarch 2/5
  Installing     : vim-common-2:8.0.1763-19.el8_6.4.x86_64    3/5
  Installing     : gpm-libs-1.20.7-17.el8.x86_64             4/5
  Running scriptlet: gpm-libs-1.20.7-17.el8.x86_64           4/5
  Installing     : vim-enhanced-2:8.0.1763-19.el8_6.4.x86_64 5/5
  Running scriptlet: vim-enhanced-2:8.0.1763-19.el8_6.4.x86_64 5/5
  Running scriptlet: vim-common-2:8.0.1763-19.el8_6.4.x86_64 5/5
  Verifying      : gpm-libs-1.20.7-17.el8.x86_64             1/5
  Verifying      : vim-common-2:8.0.1763-19.el8_6.4.x86_64 2/5
  Verifying      : vim-enhanced-2:8.0.1763-19.el8_6.4.x86_64 3/5
  Verifying      : vim-filesystem-2:8.0.1763-19.el8_6.4.noarch 4/5
  Verifying      : which-2.21-18.el8.x86_64                 5/5

Installed:
  gpm-libs-1.20.7-17.el8.x86_64          vim-common-2:8.0.1763-19.el8_6.4.x86_64
  vim-enhanced-2:8.0.1763-19.el8_6.4.x86_64  vim-filesystem-2:8.0.1763-19.el8_6.4.noarch
  which-2.21-18.el8.x86_64

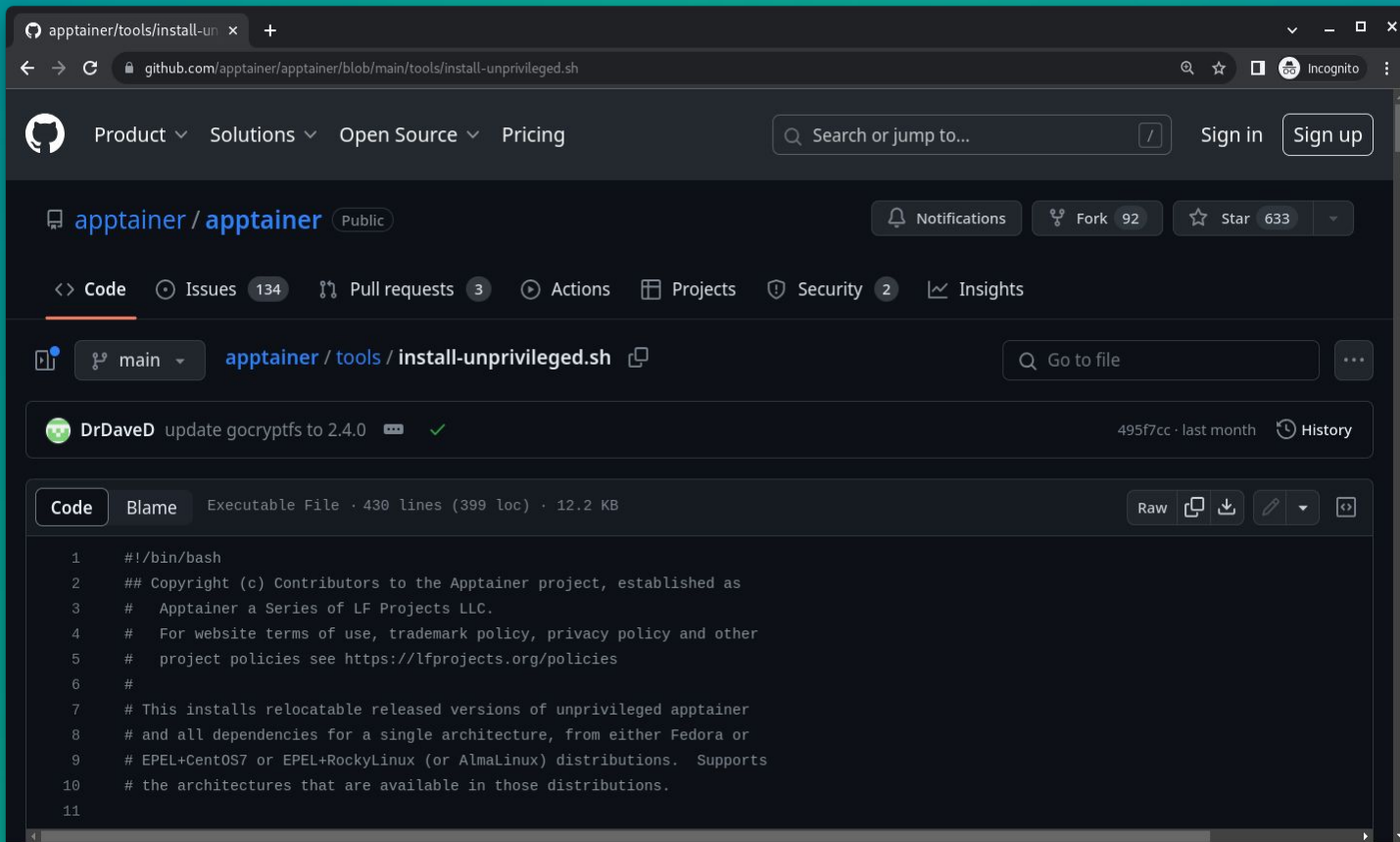
Complete!

[demouser@ciqbox ~]$
```

Leveraging the User Namespace (for installation)

- Default Apptainer installation is now unprivileged
- Squashfuse is used to mount the squashfs file system
- Apptainer enters a new User Namespace and then creates a new Mount Namespace to present the new root filesystem to processes

Leveraging the User Namespace (for installation)



The screenshot shows a web browser displaying the GitHub repository page for the file `install-unprivileged.sh` in the `apptainer/tools` directory. The repository is public and has 92 forks and 633 stars. The file is an executable script with 430 lines of code. The code content is as follows:

```
1  #!/bin/bash
2  ## Copyright (c) Contributors to the Apptainer project, established as
3  #   Apptainer a Series of LF Projects LLC.
4  #   For website terms of use, trademark policy, privacy policy and other
5  #   project policies see https://lfprojects.org/policies
6  #
7  # This installs relocatable released versions of unprivileged apptainer
8  # and all dependencies for a single architecture, from either Fedora or
9  # EPEL+CentOS7 or EPEL+RockyLinux (or AlmaLinux) distributions. Supports
10 # the architectures that are available in those distributions.
11
```

Leveraging the User Namespace (for installation)

Important Note:

Because users can install in their own space they can manage their own configuration. This means that **things like ECL may be rendered invalid!**

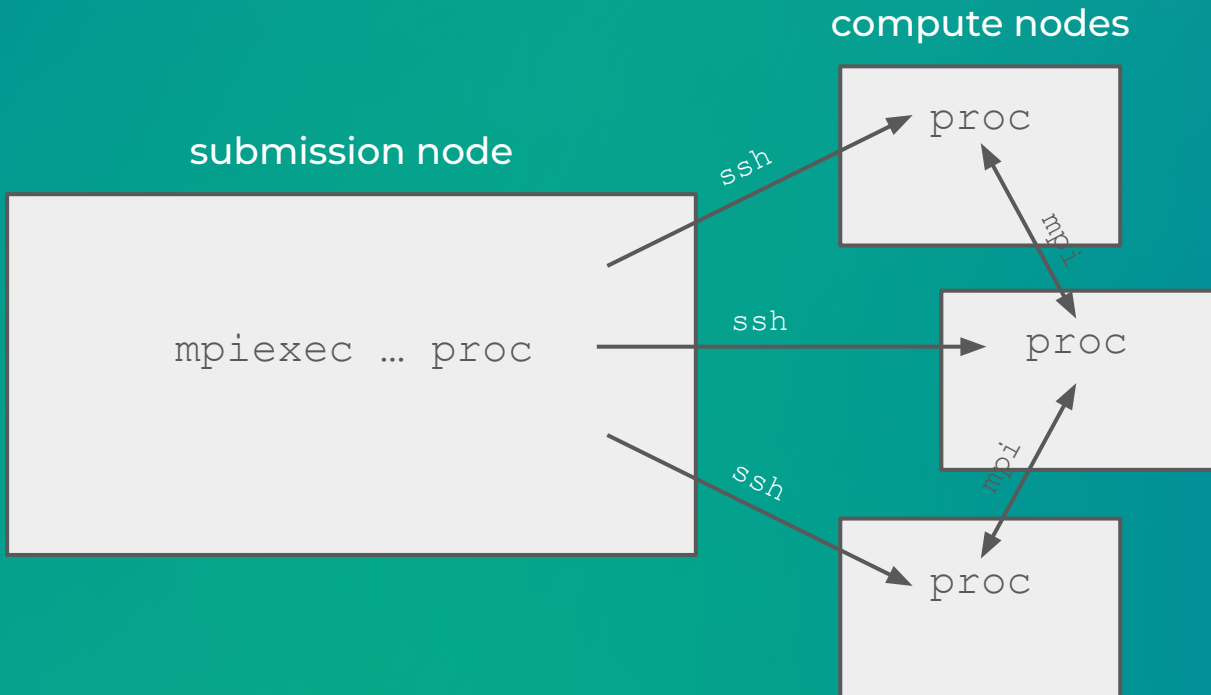
Must disable User Namespace in OS if this is a problem.

New recommendations for MPI jobs

- Containerized MPI jobs can present difficulties
 - wire-up
 - fabric adapter
- Historical approaches are difficult and lack portability
- PMI and libfabric are recommended to increase container portability
- Container-friendly package managers (like Spack) can reduce difficulty

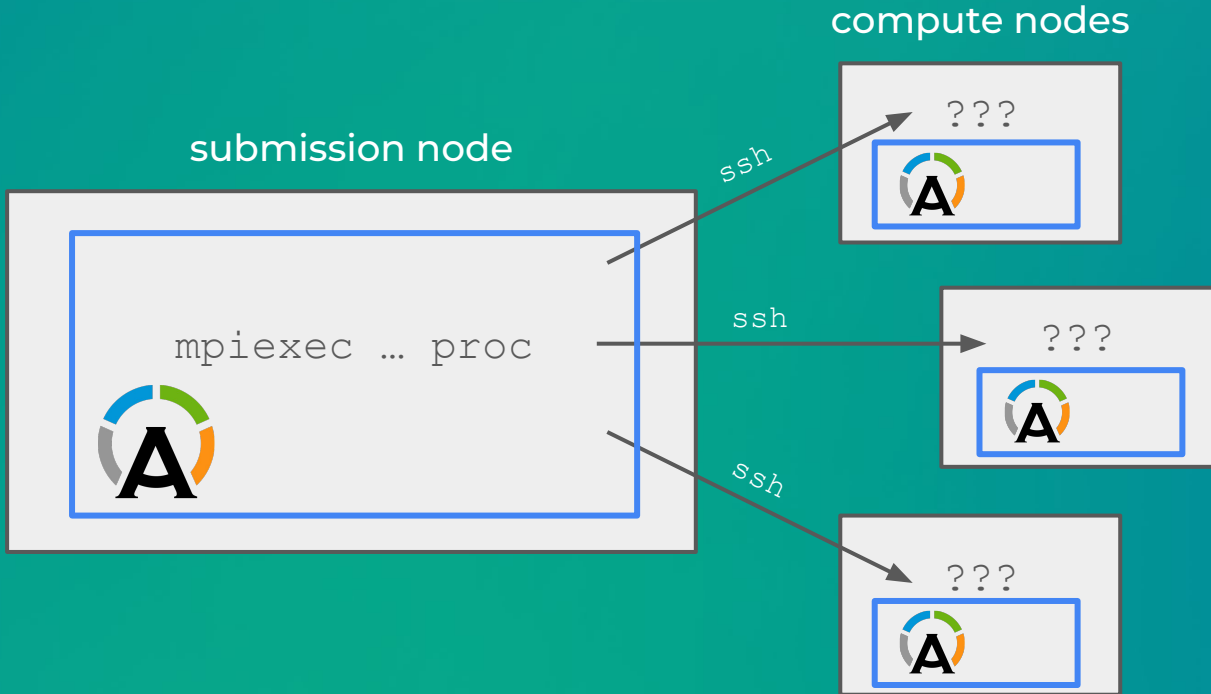
New recommendations for MPI jobs

Review of the wire-up problem



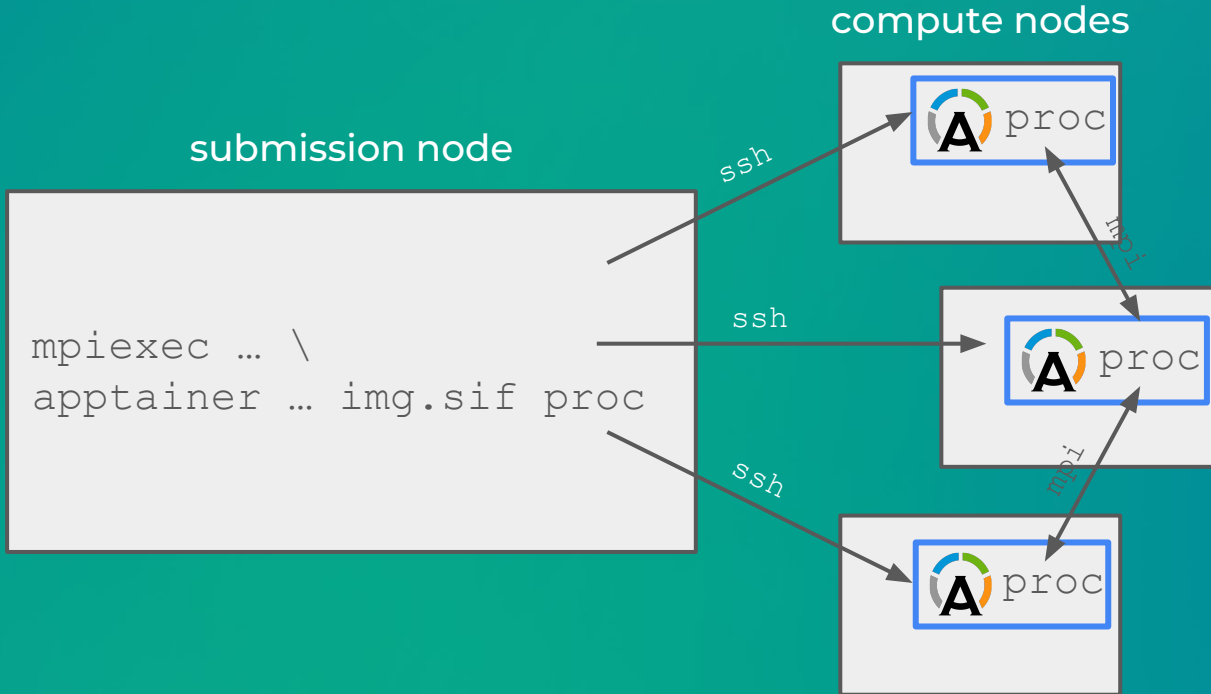
New recommendations for MPI jobs

Review of the wire-up problem



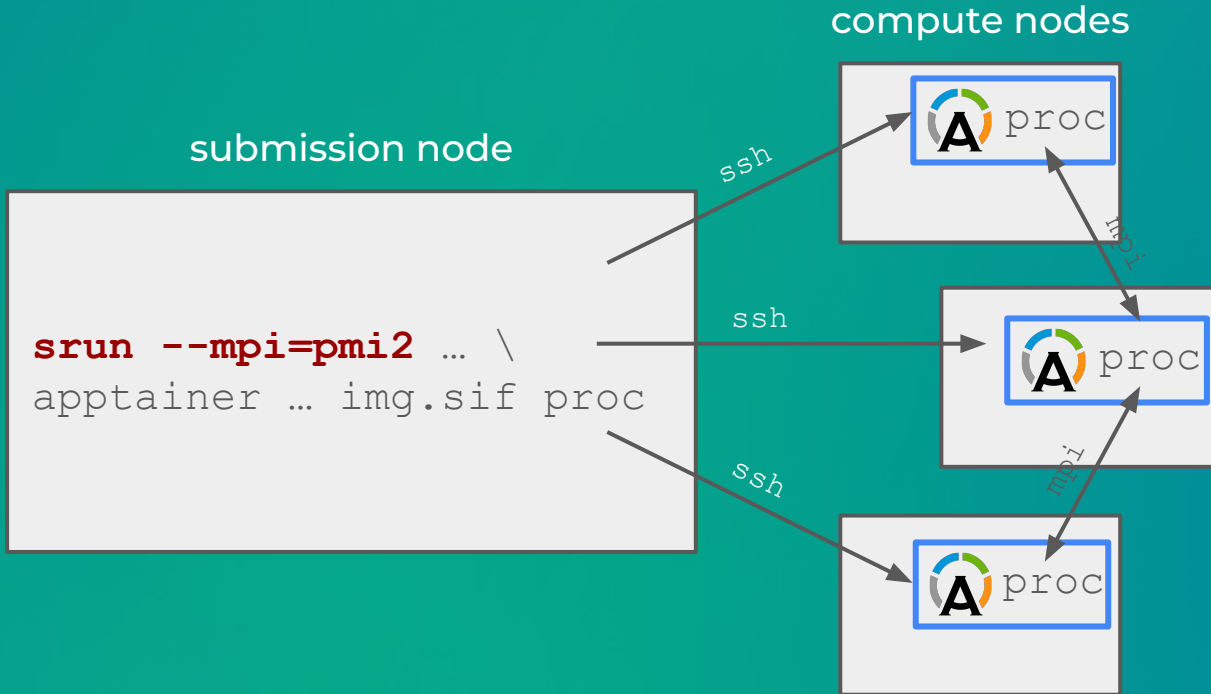
New recommendations for MPI jobs

Review of the wire-up problem



New recommendations for MPI jobs

Solving the wire-up problem with PMI



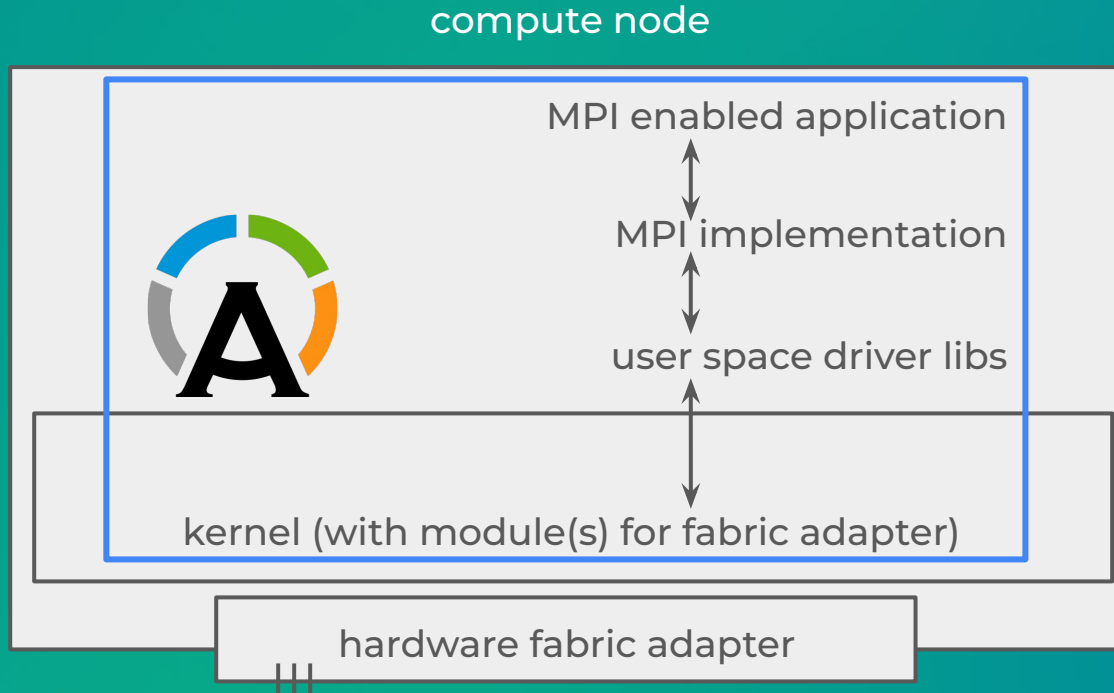
New recommendations for MPI jobs

Links to detailed resources

- [Blog post with detailed \(copy/paste-able\) scripts](#)
- Jonathon Anderson presenting this work at the [2023 HPC-AI Advisory council meeting at Stanford](#)
- Dave Godlove and Jonathon Anderson discussing this work in a [CIQ Webinar](#)

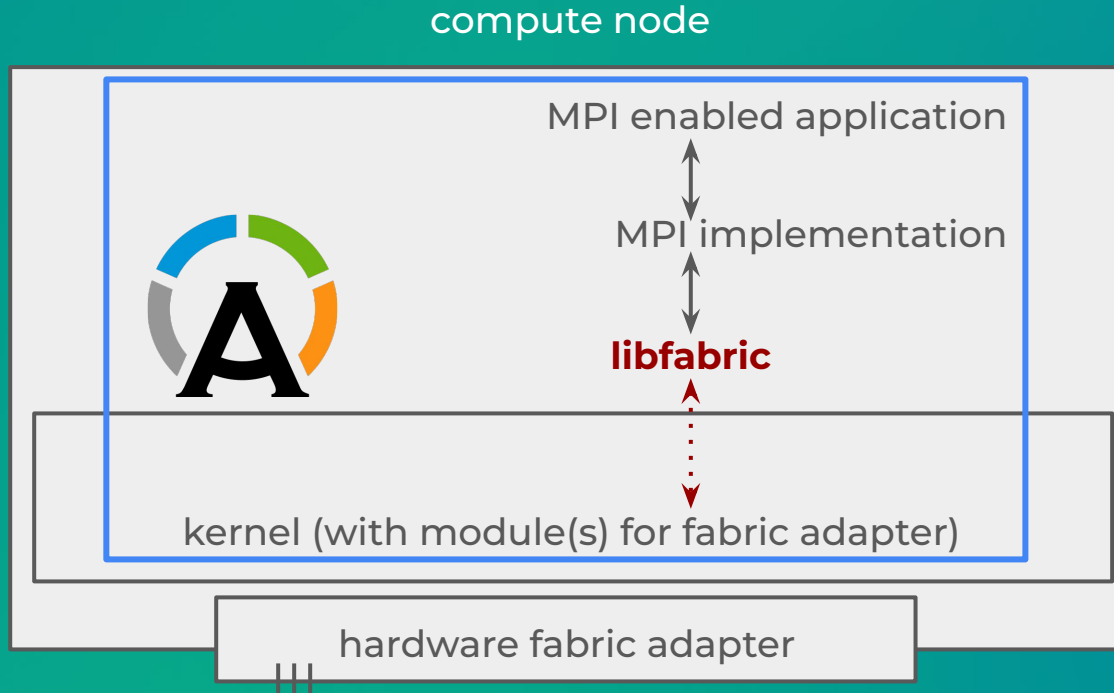
New recommendations for MPI jobs

Review of the fabric adapter problem



New recommendations for MPI jobs

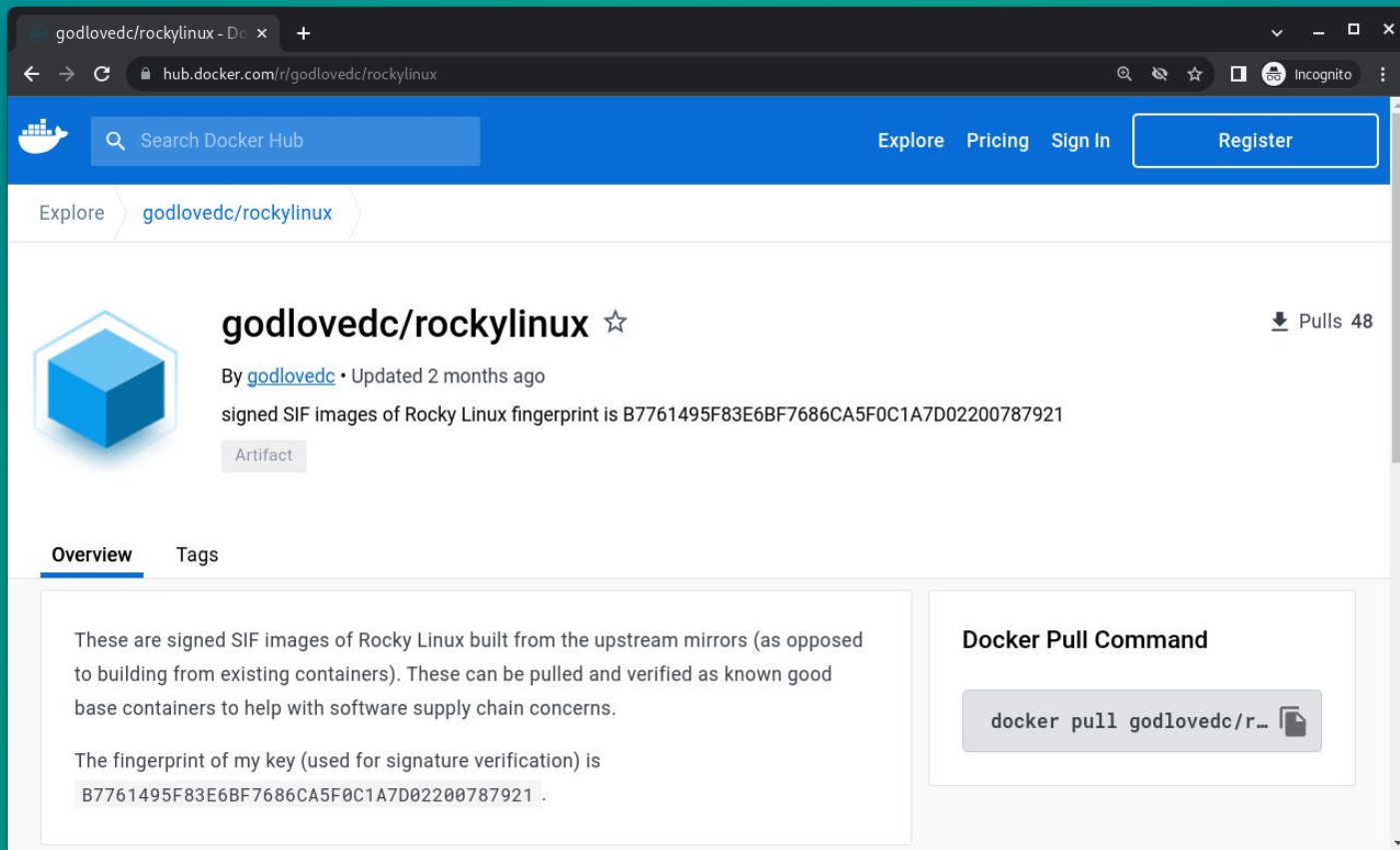
Solving the fabric adapter problem with libfabric



ORAS for Software Supply Chain

- The OCI Registry As Storage (ORAS) protocol allows native SIF files to be stored on OCI registries (like Docker Hub)
- You can use advanced features like signing, verifying, and encryption without giving up the convenience of OCI registries!

ORAS for Software Supply Chain



The screenshot shows a web browser window displaying the Docker Hub page for the repository `godlovedc/rockylinux`. The browser's address bar shows `hub.docker.com/r/godlovedc/rockylinux`. The page header includes a search bar, navigation links for "Explore", "Pricing", "Sign In", and a "Register" button. The repository page features a blue cube icon, the repository name `godlovedc/rockylinux` with a star icon, and a "Pulls 48" indicator. Below the repository name, it states "By [godlovedc](#) · Updated 2 months ago" and "signed SIF images of Rocky Linux fingerprint is B7761495F83E6BF7686CA5F0C1A7D02200787921". A grey "Artifact" label is positioned below the fingerprint. The "Overview" tab is selected, showing a description: "These are signed SIF images of Rocky Linux built from the upstream mirrors (as opposed to building from existing containers). These can be pulled and verified as known good base containers to help with software supply chain concerns." Below this, it says "The fingerprint of my key (used for signature verification) is B7761495F83E6BF7686CA5F0C1A7D02200787921 .". To the right, the "Docker Pull Command" section displays a code block with the command `docker pull godlovedc/r...` and a copy icon.


godlovedc/rockylinux - Docker Hub

hub.docker.com/r/godlovedc/rockylinux

Search Docker Hub

Explore Pricing Sign In Register

Explore godlovedc/rockylinux

 **godlovedc/rockylinux** ☆ ↓ Pulls 48

By [godlovedc](#) · Updated 2 months ago

signed SIF images of Rocky Linux fingerprint is B7761495F83E6BF7686CA5F0C1A7D02200787921

Artifact

Overview Tags

These are signed SIF images of Rocky Linux built from the upstream mirrors (as opposed to building from existing containers). These can be pulled and verified as known good base containers to help with software supply chain concerns.

The fingerprint of my key (used for signature verification) is B7761495F83E6BF7686CA5F0C1A7D02200787921 .

Docker Pull Command

```
docker pull godlovedc/r...
```

ORAS for Software Supply Chain

```
[demouser@ciqbox ~]$ aptainer pull oras://docker.io/godlovedc/rockylinux:8
INFO:   Downloading oras image

[demouser@ciqbox ~]$ aptainer verify rockylinux_8.sif
INFO:   Verifying image with PGP key material
[REMOTE] Signing entity: David Godlove (production key) <davidgodlove@gmail.com>
[REMOTE] Fingerprint: B7761495F83E6BF7686CA5F0C1A7D02200787921
Objects verified:
ID |GROUP |LINK |TYPE
-----
1  |1     |NONE |Def.FILE
2  |1     |NONE |JSON.Generic
3  |1     |NONE |FS
INFO:   Verified signature(s) from image 'rockylinux_8.sif'

[demouser@ciqbox ~]$
```

ORAS for Software Supply Chain

```
[demouser@ciqbox ~]$ aptainer inspect --deffile rockylinux_8.sif
BootStrap: yum
OSVersion: 8
MirrorURL: http://dl.rockylinux.org/pub/rocky/{OSVERSION}/BaseOS/x86_64/os/
Include: dnf

%labels
  Author: davidgodlove@gmail.com

%post
  dnf -y update
  dnf install -y epel-release file

%environment
  LC_ALL=C

[demouser@ciqbox ~]$
```


ORAS for Software Supply Chain

```
[demouser@ciqbox ~]$ apttainer build sigexample.sif sigexample.def
WARNING: 'nodev' mount option set on /tmp, it could be a source of failure during build process
INFO: Starting build...
INFO: Downloading oras image
INFO: Checking bootstrap image verifies with fingerprint(s): [B7761495F83E6BF7686CA5F0C1A7D02200787921]
INFO: Running post scriptlet
+ echo 'only install if signed and verified!'
only install if signed and verified!
INFO: Creating SIF file...
INFO: Build complete: sigexample.sif

[demouser@ciqbox ~]$
```


ORAS for Software Supply Chain

```
[demouser@ciqbox ~]$ aptainer build sigexample.sif sigexample.def
WARNING: 'nodev' mount option set on /tmp, it could be a source of failure during build process
INFO: Starting build...
INFO: Downloading oras image
INFO: Checking bootstrap image verifies with fingerprint(s): [B7761495F83E6BF7686CA5F0C1A7D02200787922]
FATAL: While performing build: conveyor failed to get: while checking fingerprint: image not signed by required entities

[demouser@ciqbox ~]$
```

Developments in Apptainer over the past few years

- Leveraging the User Namespace
- New recommendations for containerized MPI
- Increased adoption of ORAS (for Software Supply Chain)

Get involved!

- Website: <https://apptainer.org/>
- Get Started: <https://apptainer.org/get-started/>
- GitHub: <https://github.com/apptainer/apptainer>
- Community Slack, mailing list, etc links:
<https://apptainer.org/support/>
- User docs: <https://apptainer.org/docs/user/main/>
- Admin docs: <https://apptainer.org/docs/admin/main/>

THANK YOU