

---

# Singularity

The Inner Workings of Securely Running User  
Containers on HPC Systems



Michael Bauer

---

# About Me


Michael Bauer

**@bauerm97** on GitHub

bauerm@umich.edu

Singularity

NewsDocsQuick LinksPeople



# Singularity

Singularity enables users to have full control of their environment. This means that a non-privileged user can “swap out” the operating system on the host for one they control. So if the host system is running RHEL6 but your application runs in Ubuntu, you can create an Ubuntu image, install your applications into that image, copy the image to another host, and run your application on that host in it’s native Ubuntu environment!

Register your Cluster

Add a Publication

Singularity also allows you to leverage the resources of whatever host you are on. This includes HPC interconnects, resource managers, file systems, GPUs and/or accelerators, etc. Singularity does this by enabling several key facets:

- Encapsulation of the environment
- Containers are image based
- No user contextual changes or root escalation allowed
- No root owned daemon processes

## Getting started

Jump in and [get started](#).


Singularity

InformationDownload / InstallationContributingGetting HelpDocumentation

A

A

GitHub, Inc. github.com/singularityware/singularity

 Personal Open source Business Explore Pricing Blog Support

This repository

Sign in 

Sign up

 singularityware / singularity

Watch40

Star143

Fork64

<> Code

Issues68

Pull requests9

Projects1

Pulse

Graphs

Singularity: Application containers for Linux <http://singularity.lbl.gov/>

2,042 commits

8 branches

6 releases

33 contributors

Branch: master 

New pull request

Find file

Clone or download

 gmkurtzer committed on GitHub Merge pull request #469 from satra/patch-1 ... Latest commit 1c72426 19 hours ago

 .github	fixing 'i' with a bad head cold	3 months ago
 alpinelinux	initial port for alpinelinux	3 months ago
 bin	Added support for shell debugging specific output	4 months ago
 debian	QL: d/rules - New version needs sexec-suid set suid	3 months ago
 docs	replace doc files with README pointing to docs repo	a month ago
 etc	Decided not to set TZ in init at all	8 days ago
 examples	fixing container creation	a day ago
 libexec	updating shub tests to support image ids in new database	10 days ago
 man	fixing all pointers to old repo, and updating badge in README	3 months ago
 src	Merge pull request #448 from bbockelm/reentrant_hashmap	11 days ago

## Graphs

Copy path

1c45cbb on Nov 30, 2016

4 contributors    

22 lines (19 sloc) | 691 Bytes

## History

```

1 Project Lead:
2     Gregory M. Kurtzer <gmkturtzer@lbl.gov>
3
4 Developers:
5     Brian Bockelman <bbockelm@cse.unl.edu>
6     Krishna Muriki <kmuriki@lbl.gov>
7     Michael Bauer <bauerm@umich.edu>
8     Vanessa Sochat <vsochat@stanford.edu>
9
10 Contributors:
11     Amanda Duffy <aduffy@lenovo.com>
12     Ángel Bejarano <abejarano@ntropos.com>
13     Bernard Li <bernardli@lbl.gov>
14     Dave Love <d.love@liverpool.ac.uk>
15     Felix Abecassis <fabecassis@nvidia.com>
16     Jarrod Johnson <jjohnson2@lenovo.com>
17     Jason Stover <jason.stover@gmail.com>
18     Maciej Sieczka <msieczka@sieczka.org>
19     Nathan Lin <nathan.lin@yale.edu>
20     Ralph Castain <rhc@open-mpi.org>
21     Yaroslav Halchenko <debian@onerussian.com>

```

# What are Containers?

# What is a Virtual Machine?

“In computing, a virtual machine (VM) is an emulation of a computer system. Virtual machines are based on computer architectures and provide functionality of a physical computer.”



## Examples:



VirtualBox

# Pros

- Run different OS on one set of hardware
- Save money (e.g. buy one laptop, have Windows, OSX, and Linux)
- Easy maintenance

# Cons

- Slower performance
- Memory/storage reqs.

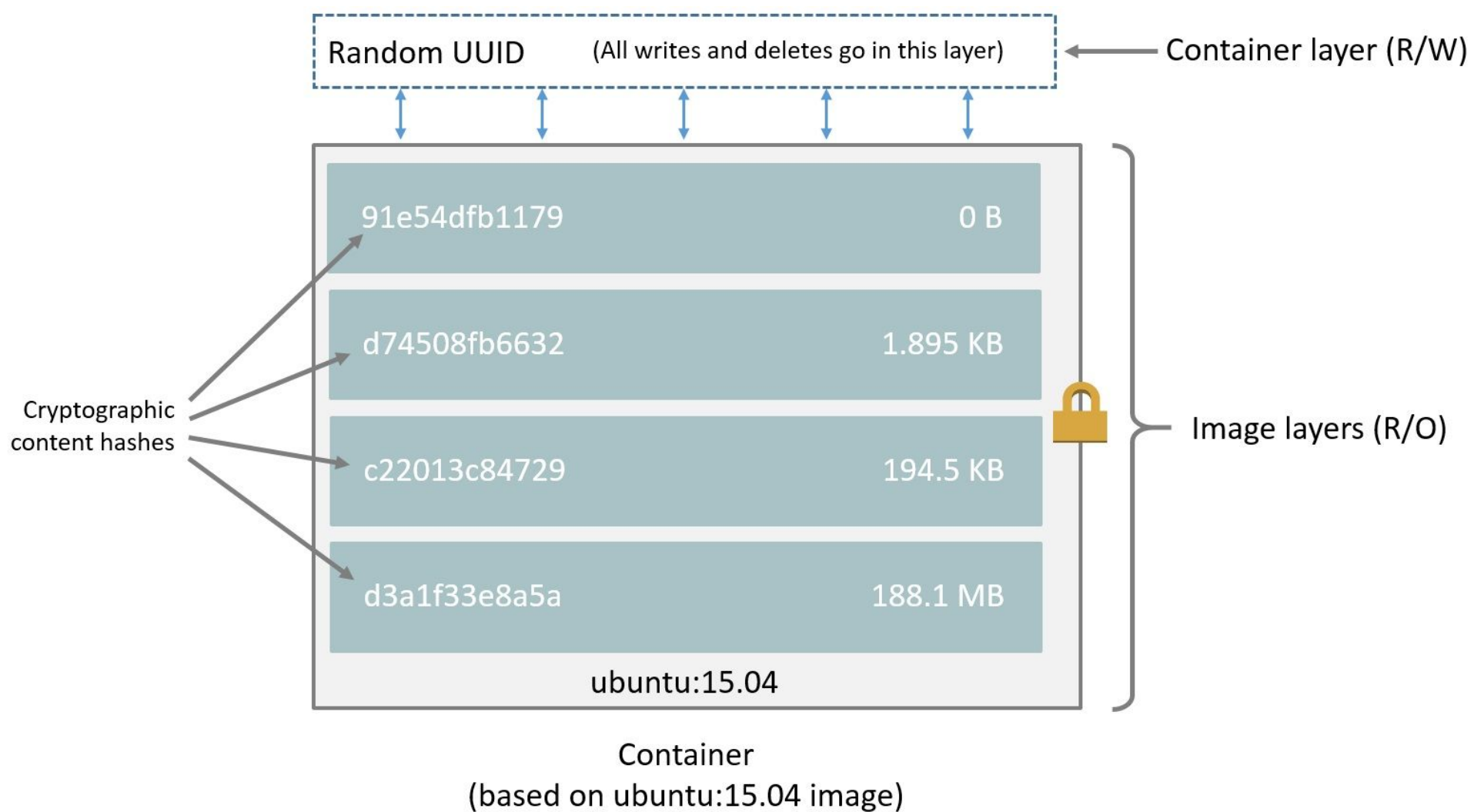


# What are containers?

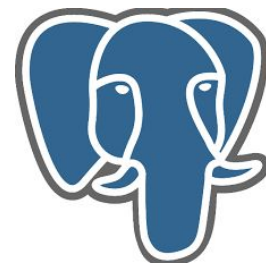
- Similar goal as VMs
- No kernel emulation
- Not architecture level virtualization, but rather software level

# What does that mean?

- Don't waste extra ~5% performance doing emulation
- Smaller footprint (~500 MB vs ~20 GB VM)
- Very small startup time interval (~1 s vs ~1 min VM)
- Multiple instances can share one “container image”



# Who uses containers?

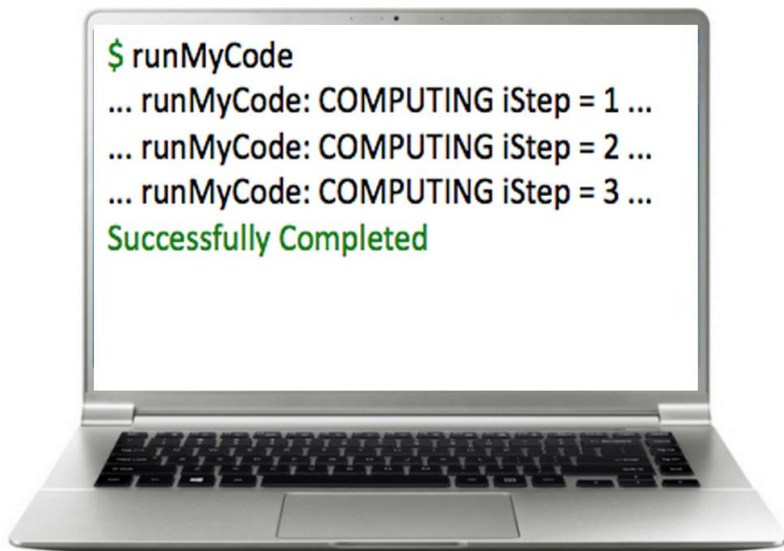


# Containers for Scientific Computing

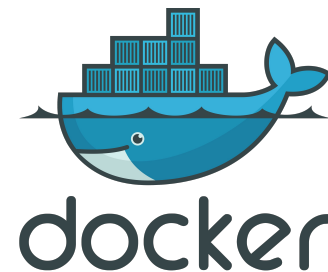
# Why do we want containers in HPC?

- Escape “dependency hell”
- Local and remote code works identically every time
- One file contains everything and can be moved anywhere

# Environment Matters



# Needs for HPC containers



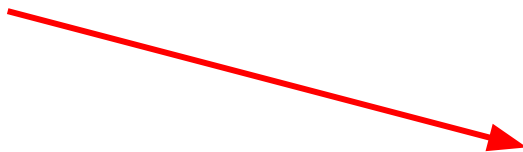
- Any user can run containers without special privileges (root) ✓
- Integrate seamlessly into existing infrastructure ✗
- Portability between many systems ✓
- Users created and provided containers (no administrative oversight) ✓



**HPC container software can never touch root**

# What about Docker?

- Root level process spawns containers as child processes
- Integration into infrastructure difficult (e.g. must consider other scheduler systems)



# Singularity



# Needs for HPC containers



- Any user can run containers without special privileges (root)
- Integrate seamlessly into existing infrastructure
- Portability between many systems
- Users created and provided containers (no administrative oversight)



# Singularity



- Any container can be run by any user - same user inside container and on host
- No workflow changes necessary to use
- Single .img file contains everything necessary
- Safe to run any container without screening its contents

	Shifter	Charlie Cloud	Docker	Singularity
Privilege model	SUID	UserNS	Root Daemon	SUID/UserNS
Support current production Linux distros	Yes	No	No	Yes
Internal image build/boostrap	No*	No*	No**	Yes
No privileged or trusted daemons	Yes	Yes	No	Yes
No additional network configurations	Yes	Yes	No	Yes
No additional hardware	Maybe	Yes	Maybe	Yes
Access to host filesystem	Yes	Yes	Yes***	Yes
Native support for GPU ★	No	No	No	Yes
Native support for InfiniBand ★	Yes	Yes	No	Yes
Native support for MPI ★	Yes	Yes	No	Yes
Works with all schedulers	No	Yes	No	Yes
Designed for general scientific use cases	Yes	No	No	Yes
Contained environment has coorrect perms	Yes	No	Yes	Yes
Containers are portable, unmodified by use	No	No	No	Yes
Trivial HPC install (one package, zero conf)	No	Yes	Yes	Yes
Admins can control and limit capabilities	Yes	No	No	Yes

\* Relies on Docker

\*\* Depends on upstream

\*\*\* With security implications



Site or Organization	System Name	Size (cores)	Purpose of the System
CSIRO	bragg-gpu	2048	broad base scientific
GSI Helmholtz Center	Greencube	300,000	Heavy Ion Physics
Holland Computing Center at UNL	Crane and Tusker	14,000	General purpose campus cluster
HPC-UGent	golett	2500	research across all scientific domains
Lunarc	Aurora	360	Research
Microway	Microway Research Cluster	192	Scientific benchmarking
MIT	openmind	1,176	Neuroscience
National Institute of Health HPC	Biowulf	54,000	General purpose biomedical research
Purdue University	Rice	11520	Campus HPC resource
Purdue University	Conte	78880	Campus HPC resource
Purdue University	Snyder	2220	Campus HPC resource
Purdue University	Hammer	3960	Campus HPC resource
Purdue University	Carter	10560	Campus HPC resource
R Systems NA, Inc.	Oak1	1024	Shared commercial/academic resource
R Systems NA, Inc.	Oak2	2048	Shared commercial/academic resource
R Systems NA, Inc.	HOU1	5376	Shared commercial/academic resource
San Diego Supercomputer Center	Gordon	16384	HPC cluster for XSEDE users
San Diego Supercomputer Center (SDSC)	Comet	47776	HPC Cluster for XSEDE users
Texas Advanced Computing Center	Stampede	102400	NSF key resource, all fields
UFIT Research Computing at the UF	HiPerGator	51,000	research computing cluster
Ulm University, Germany	JUSTUS	550	Computational Chemistry
University of Chicago	midway.rcc.uchicago.edu	24196	University cluster
University of Manitoba	GreX	3840	General purpose HPC cluster
Georgia State University	Orion	362	research
UNF	Stark	64	Functional MRI analysis of the Brain
Genentech, Inc.			Research
Rutgers University	sirius	32	scientific SMP machine
Stanford University	sherlock	12764	Compute for Stanford researchers
Stanford University	scg4	3920	Genomics at Stanford
The University of Leeds	MARC1	1236	Bioinformatics, data analytics
McGill HPC Centre/Calcul Québec	guillimin	22300	Compute Canada cluster
University of Arizona	Ocelote	10000	General Research
University of Arizona	ElGato	2300	GPU cluster
Washington University in St. Louis		2000	General purpose cluster



## Singularity: Scientific Containers for Mobility of Compute

Gregory M. Kurtzer<sup>1</sup>, Vanessa Sochat<sup>2\*</sup>, Michael W. Bauer<sup>3,4</sup>

**1** High Performance Computing Services, Lawrence Berkeley National Lab, Berkeley, CA, USA

**2** Stanford Research Computing Center and School of Medicine, Stanford University, Stanford, CA, USA

**3** Department of Electrical Engineering and Computer Science, University of Michigan, Ann Arbor, MI, USA

**4** Experimental Systems, GSI Helmholtzzentrum für Schwerionenforschung, Darmstadt, Germany



# Basic Usage of Singularity

# Singularity Workflow

## 1. Create image file

```
$ sudo singularity create [image]
```

## 2. Bootstrap image

```
$ sudo singularity bootstrap [image] [definition.def]
```

## 3. Run image

```
$ singularity shell [image]
```

```
$ singularity exec [image] [/path/to/executable]
```

```
$ singularity run [image]
```

```
$ ./image
```

# Singularity Workflow

<https://asciinema.org/a/100297>

<b>Format</b>	<b>Description</b>
<i>directory</i>	Standard Unix directories containing a root container image
<i>tar.gz</i>	Zlib compressed tar archives
<i>tar.bz2</i>	Bzip2 compressed tar archives
<i>tar</i>	Uncompressed tar archives
<i>cpio.gz</i>	Zlib compressed CPIO archives
<i>cpio</i>	Uncompressed CPIO archives

# Docker Integration

<https://asciinema.org/a/101984>

# SLURM Integration

```
#!/bin/bash -l
```

```
#SBATCH --image=~ /centos7/latest
```

```
#SBATCH -p debug
```

```
#SBATCH -N 64
```

```
#SBATCH -t 00:20:00
```

```
#SBATCH -J my_job
```

```
#SBATCH -L SCRATCH
```

```
#SBATCH -C haswell
```

```
srun -n 4096 ./mycode.exe    # an extra -c 1 flag is optional for fully packed pure MPI with  
hyperthreading
```

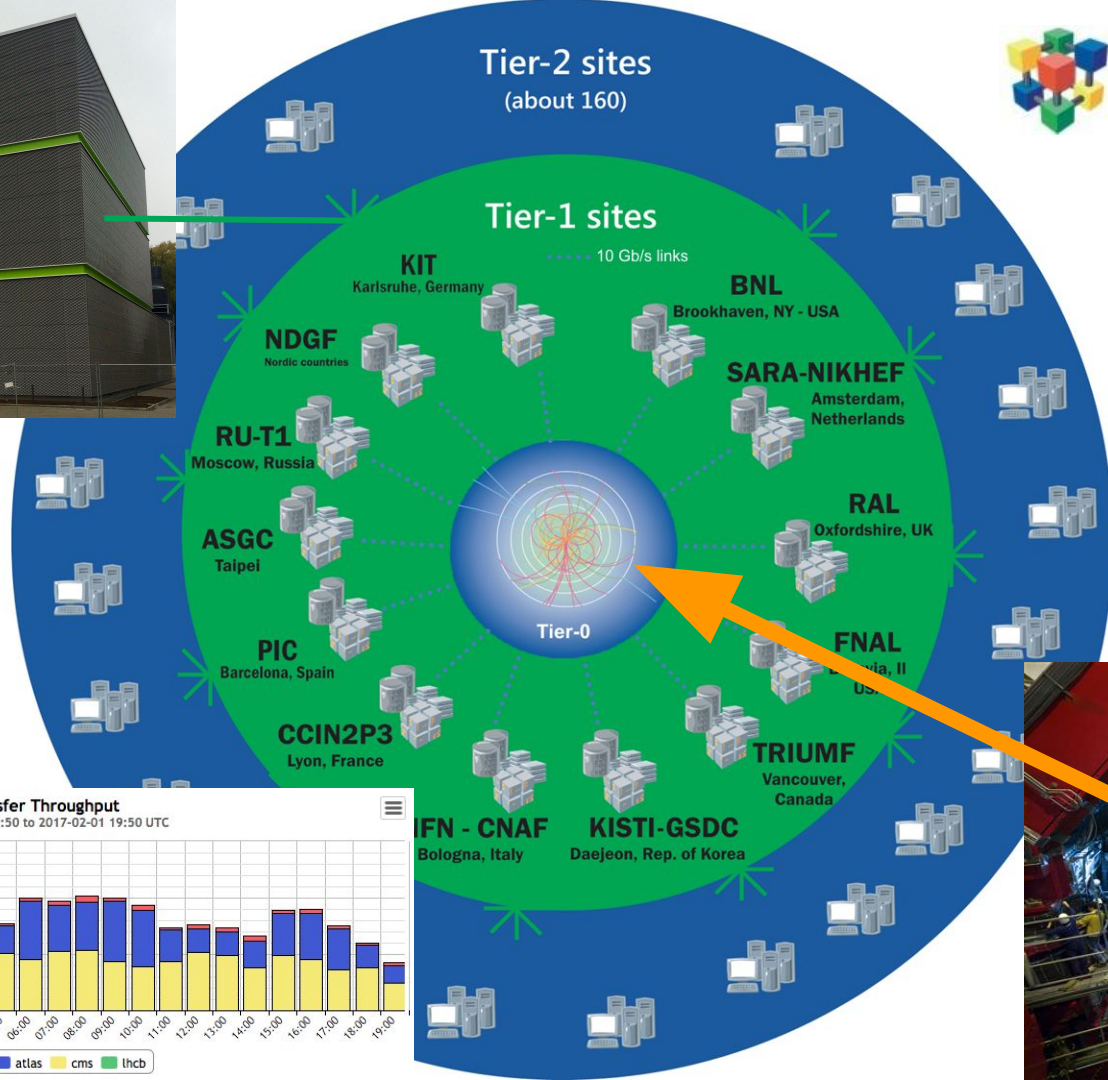
# ALICE Tier 2 Use Case



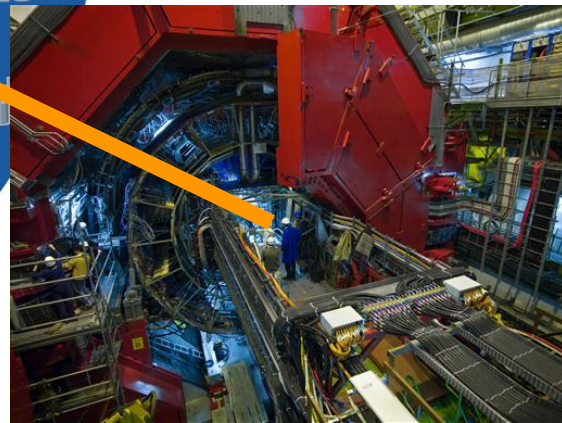
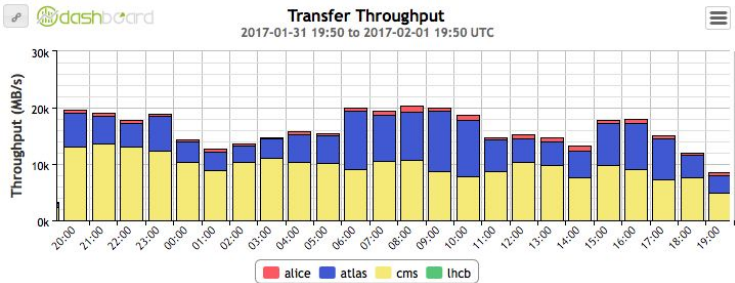
GSI Green Cube  
Darmstadt  
Germany



**WLCG**  
Worldwide LHC Computing Grid



ALICE Detector LHC  
Geneva  
Switzerland





# ALICE Tier 2: Problem

- Run ALICE jobs on ~2k jobs at any time
- Host machines run Debian 7.x kernel 3.16
- ALICE expects Scientific Linux 6 (SL6)
- Library incompatibilities cause frequent errors (much higher than expected)

# ALICE Tier 2: Pre-Singularity Solution

- Correct library versions mounted in Lustre
- SLURM job submission script alters `$LD_LIBRARY_PATH` to point to Lustre
- And maybe more?

Big Ugly Hack

## ALICE Tier 2: Singularity Solution

- Package Scientific Linux 6 into container
- Modify SLURM submission script to run container
- No need to mount Lustre for access to library files
- Can test container locally before deploying to HPC

ALICE GitHub  
Repository



**How does it work**

# How

- Installed as SUID binary owned by root
- All necessary files are mounted
- Hide inside namespaces when possible and requested

# SUID Binary Security

- Make code possible for anybody to audit
- Function call to change eUID to 0
- Function call to change eUID to calling user UID
- Wrap actions that require root in these two calls

# User Namespaces?

- Map UID on host : UID inside user namespace
- Looks and feels like root (mostly)
- Potentially breaks portability



If any user code gets executed as root user, the system should be considered compromised.



# Security and Singularity

# Core Principles

1. Never run user code as real user root
2. Only use eUID 0 when necessary (escalate and drop permissions accordingly)
3. Drop permissions and capabilities when forking into new thread



<http://jdfinley.com/important-papers/safe/>

# Isolation



- Bind mount image file into host's filesystem
- Use chroot to move into the mounted image's root filesystem
- Mount /dev/, /etc/hosts, etc... into container filesystem
- Use namespaces when possible and requested (e.g User, PID, etc...)
- No isolation of network means no extra network configuration
- Can use host's physical devices inside container (e.g. IPoIB)

# Security Cont.



- User invokes non-SUID Singularity binary, which in turn calls SUID binary.  
Only non-SUID binary can run SUID binary
- singularity\_priv\_escalate() called to escalate privileges to eUID 0  
singularity\_priv\_drop() called to drop privileges to eUID of calling user
- Only escalate privileges when necessary
- MS\_NOSUID and O\_CLOEXEC flags set when necessary

<b>Global Options</b>	
<i>-d - --debug</i>	Print debugging information
<i>-h - --help</i>	Display usage summary
<i>-q - --quiet</i>	Only print errors
<i>- - version</i>	Show application version
<i>-v - --verbose</i>	Increase verbosity +1
<i>-x - --sh - debug</i>	Print shell wrapper debugging information
<b>General Commands</b>	
<i>help</i>	Show additional help for a command
<b>Container Usage Commands</b>	
<i>exec</i>	Execute a command within container
<i>run</i>	Launch a runscript within container
<i>shell</i>	Run a Bourne shell within container
<i>test</i>	Execute any test code defined within container
<b>Container Management Commands (requires root)</b>	
<i>bootstrap</i>	Bootstrap a new Singularity image
<i>copy</i>	Copy files from your host into the container
<i>create</i>	Create a new container image
<i>export</i>	Export the contents of a container via a tar pipe
<i>import</i>	Import/add container contents via a tar pipe
<i>mount</i>	Mount a Singularity container image

## Graphs

Copy path

1c45cbb on Nov 30, 2016

4 contributors    

22 lines (19 sloc) | 691 Bytes

## History

```

1 Project Lead:
2     Gregory M. Kurtzer <gmurtzer@lbl.gov>
3
4 Developers:
5     Brian Bockelman <bbockelm@cse.unl.edu>
6     Krishna Muriki <kmuriki@lbl.gov>
7     Michael Bauer <bauerm@umich.edu>
8     Vanessa Sochat <vsochat@stanford.edu>
9
10 Contributors:
11     Amanda Duffy <aduffy@lenovo.com>
12     Ángel Bejarano <abejarano@ontopos.com>
13     Bernard Li <bernardli@lbl.gov>
14     Dave Love <d.love@liverpool.ac.uk>
15     Felix Abecassis <fabecassis@nvidia.com>
16     Jarrod Johnson <jjohnson2@lenovo.com>
17     Jason Stover <jason.stover@gmail.com>
18     Maciej Siczka <msiczka@siczka.org>
19     Nathan Lin <nathan.lin@yale.edu>
20     Ralph Castain <rhc@open-mpi.org>
21     Yaroslav Halchenko <debian@onerussian.com>

```

We are always looking for more Collaborators!

