**Sandia National Laboratories**

*Exceptional service in the national interest*

# THE CASE FOR CONTAINERIZATION WITH HPC SCIENTIFIC SOFTWARE

## Samuel E. Browne

ASC S3C at NLIT 24, Seattle, WA
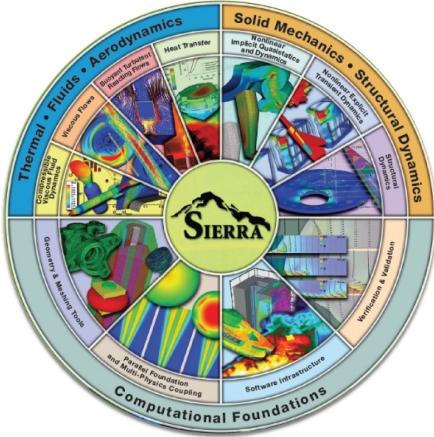
April 8-11, 2024

# ACKNOWLEDGEMENTS

- SIERRA DevOps Team

    - Mark Hamilton, Jake Healy, Andrew Kimler, Justin Lamb, Matt Mosby, Tony Nguyen

- Trilinos Framework Team

    - Anderson Chauphan, Joe Frye, Justin LaPre

# CODES BACKGROUND



SIERRA is Sandia's engineering mechanics code suite. It is a large C++ code under development for more than 30 years, and is widely used across many government HPC customers. It leverages the Trilinos software stack.



Trilinos is an object-oriented software framework for the solution of large-scale, complex multi-physics engineering and scientific problems. It is primarily developed at Sandia, but is open-source and hosted on GitHub, with many developers from the open-source community.

# BACKGROUND

- I have been a member of the SIERRA DevOps team since 2016, primarily working on areas surrounding the build system

- I wrote a custom TPL builder for the SPARC team in 2018 to support the ATS2 'Sierra' platform, which was then extended to all other HPCs supported by SPARC (ATS1, CTS1, Vanguard 'Astra', etc.)

- Since 2022, I have led the Trilinos Framework team which is responsible for the CI build/test infrastructure (including environment management)

# INDUSTRY'S PREFERRED ENCAPSULATION METHOD

Containers are **packages of software that contain all of the necessary elements to run in any environment.** In this way, containers virtualize the operating system and run anywhere, from a private data center to the public cloud or even on a developer's personal laptop.
– Google Cloud

**>53% of software development companies report using containers***

* Survey responses were all over the map, so the lowest found percentage was used

So why don't more HPC development teams use containers for environment management?
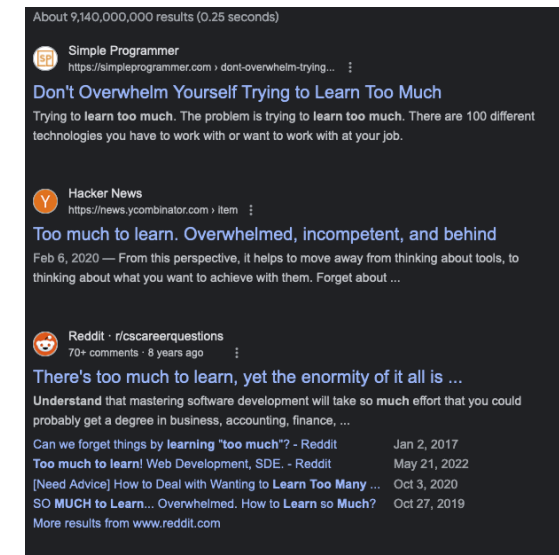
# THE BIGGEST PROBLEM

- Codes are often developed by physicists/hard-science engineers rather than career software engineers, so workflows must be *as simple as is reasonably possible*, freeing the developers to focus on the complex physics

- The "too much to learn" issue:

Amusing note: of the top twelve Google results for this string, six are about software engineering

# CHANGING WORKFLOWS

Current developer workflow

```
module load gcc
module load openmpi
module load hdf5
…
```

Proposed new workflow

```
podman run <containerid>
…
```

These are only the setup commands, not the entire workflow.

The setup may be simpler, but **what about other aspects of the workflow**?

# HPC DEVELOPMENT REQUIREMENTS

- Scientific software development workflows have different requirements than "normal" software

  - Large code suites (e.g. SIERRA base code repository ~25GB)

  - LARGE test suites (e.g. SIERRA test repository ~155GB)

  - Deep and complex third-party library (TPL) dependency trees

    - No universal C++ package manager ala npm for JavaScript or pip for Python

  - MPI parallelism

  - GPU/other device accelerator support

  - "Real" runtime environment is not a container (yet), so concerns about reproducibility

**Are containers incapable of meeting any of the above requirements?**

# LARGE CODE/TESTS

# LOCAL DISKS INTRODUCE COMPLEXITY

- The large code/tests issue from earlier appears!

- Cloning ~200GB of code each time a container is started is not feasible



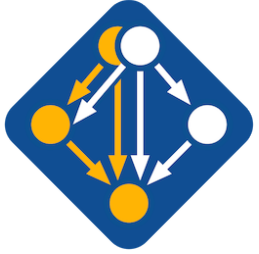Overweight dev container accident

- Enter mounting complexity

```
podman run --mount type=bind,src=/path/to/my/code,dst=/code <containerid>
```

# TPL COMPLEXITY

# CONTAINERS REDUCE NEEDED CONTEXT

Spack – a flexible package manager for HPCs

- Using Spack is demonstrably cheaper than maintaining builds of TPL stacks manually
- Using Spack is complicated, because Spack solves a **very complicated problem**
  - Example: Multiple parallel installs of different compiler-based software stacks

> If we want to use Spack to solve our problem, but as simply as possible, we must reduce the complexity of the problem to be as simple as possible

- Containers make using Spack at least an order of magnitude easier
- *Programming environment teams can now provide Spack for downstream developers*
  - They can also pre-build TPLs for developers who don't want/need to build them
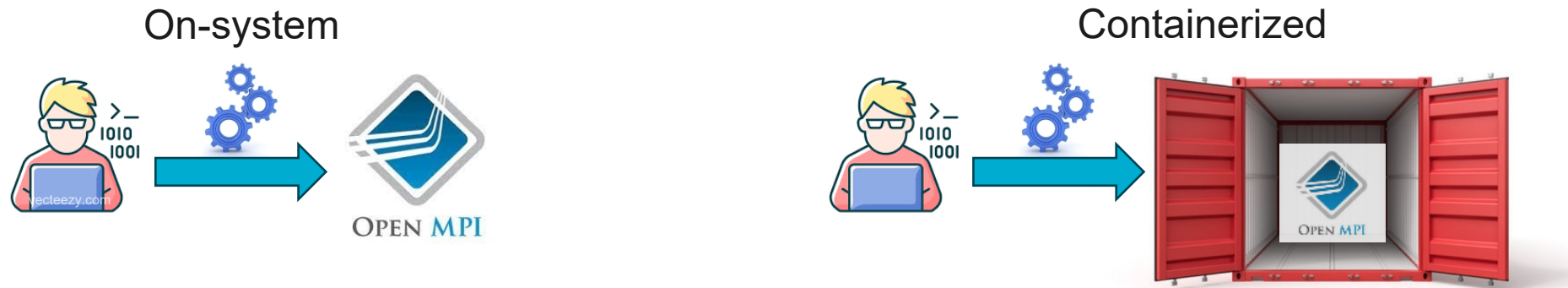
# MPI

# IS IT REALLY ANY DIFFERENT?

- Most HPC scientific software is NOT developed on HPC-style machines

- Developers have large-ish (~16-32 core) workstations that can run MPI software (albeit single-node)

- How does the MPI get there?

On-system

Containerized

This is the same from the user's perspective, and as with other TPL dependencies, much more robust than on-hardware builds

Note: In other scenarios with vendor-specialized MPIs, this may be increased complexity, but still has solutions

# GPUS

# IS IT REALLY ANY DIFFERENT?

- Key Concept: Most code development does not currently happen on machines with GPUs

  - However, this may change with the evolving hardware landscape

- GPU hardware/drivers **can** be accessed from within containers

  - See work with ECP/E4S and ParaTools testing large stacks of software on a continuous basis using containers and many different GPU architectures

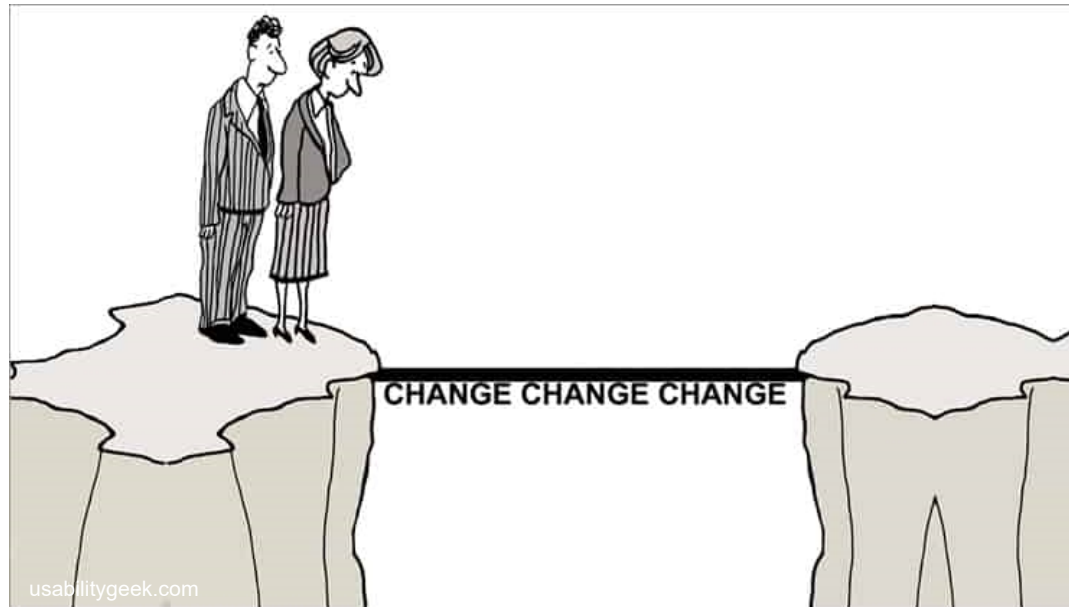Differences in CPU-only vs. GPU development workflows are likely, but technical paths forwards exist

# RUNTIME ENVIRONMENT REPRODUCIBILITY

# MY CUSTOMERS DON'T WANT TO USE A CONTAINER

- Nobody wants to change their workflows at all, ever, and this include shifting to using environments that are managed via containers instead of on-system environment modules



Taking a step back, is it a problem for development to use containerized environments if the end-user does not want to use one?

# CASE STUDY: SIERRA DEVOPS CODE TESTING

- SIERRA has a set of Python-based tools that wrap building the code, as well as a specialized distributed test harness that can interface with HPC queues

- The team desired to move to GitLab CI to manage their testing for integration and coverage visibility reasons

  - Restriction: Can not run as GitLab user, must run as user who pushed

  - Jacamar or containers?

  - Manufactured a RHEL UBI container that contained what we knew to be our runtime dependencies and added it to testing ~02/2023

**Concern: We are no longer testing in the environment that our users are using**

# CASE STUDY: SIERRA DEVOPS CODE TESTING

- Bugs in our code allowed through testing due to environment differences to-date: 0*

- It has repeatedly allowed us to distinguish issues with the runtime environment vs. our code
  - This allowed us to pass relevant debugging information back to the developers of those environments



- * Once, when attempting to upgrade the version of Git in the container, we exposed that our code would not work with the new Git, *prior to deploying that new Git to our customers!*

# A NEW WAY TO MANAGE CI ENVIRONMENTS

- Contributors to the Trilinos project have historically had problems reproducing the environments represented by CI testing

- Containers will allow us to share testing environments

  - Externally AND internally!

  - Developers can help drive TPL selection/upgrade/configuration process

- Results already matching (and with increased stability) on-hardware testing

# A NEW WAY TO MANAGE CI ENVIRONMENTS – BETTER WORKFLOW

- As with SIERRA DevOps, Trilinos wants better integration with its CI processes and its repository manager (GitHub)

- Using containers will allow for secure use of self-hosted runners and GitHub Actions

# CONCLUSION: CONTAINERS SUIT THE TECHNICAL NEED

But remember, the real issue is that workflows need to be as simple as reasonably possible

- Can a containerized development workflow be as simple as what is done today?
    - MAYBE – there is always complexity involved with solving a complex problem
    - Can the complexity be shifted, or hidden where appropriate (without making it impossible to expose!) to simplify the development process for HPC scientific software?