





Containers-First User Environments on HPE Cray EX

Felipe A. Cruz, Alberto Madonna cruz@cscs.ch May 2024

Overview

- Alps Research Infrastructure at CSCS
- On Managing User Environments in HPC systems
- Containers-First
 - The Approach
 - The Architecture
- Pathfinder Implementation
 - Environment Definition File
 - Container Ecosystem Integration for HPC
 - Running example
 - Enhancing functionality and usability
- Use-cases: Al









Image generated by ChatGPT

Alps Research Infrastructure





- General-purpose compute and data Research Infrastructure
 - One Infrastructure
 - Many Uses
- Enabled by vCluster Tech
 - Versatile Clusters
 - Tailored solutions
 - Infrastructure as Code









On Managing User Environments for HPC vClusters

On Managing User Environments in HPC systems

Difficult and getting harder

- Integration Complexity
- Flexibility vs. Stability
- Efficiency in Deployment









On Managing User Environments in HPC systems

Difficult and getting harder

- Integration Complexity
- Flexibility vs. Stability
- Efficiency in Deployment









On Managing User Environments in HPC systems

Difficult and getting harder

- Integration Complexity
- Flexibility vs. Stability
- Efficiency in Deployment















- Decoupling the User Environment
- Enhancing Flexibility and Efficiency
- Increasing Consistency







- Decoupling the User Environment
- Enhancing Flexibility and Efficiency
- Increasing Consistency







- Decoupling the User Environment
- Enhancing Flexibility and Efficiency
- Increasing Consistency





Containers-First Main Features

Leverage Containers

- Isolation of Environments
- Reproducibility and Consistency
- Efficient Deployability
- Extended Compatibility
- Flexibility for Experimentation
- Integration with HPC Systems
- Frictionless User Environments





Containers-First Main Features

- Leverage Containers
 - Isolation of Environments
 - Reproducibility and Consistency
 - Efficient Deployability
 - Extended Compatibility
 - Flexibility for Experimentation
- Integration with HPC Systems
- Frictionless User Environments





Containers-First Main Features

- Leverage Containers
 - Isolation of Environments
 - Reproducibility and Consistency
 - Efficient Deployability
 - Extended Compatibility
 - Flexibility for Experimentation
- Integration with HPC Systems
- Frictionless User Environments









Containers-First Architecture

Containers-First Architecture

- Container Ecosystem Tools
 - Runtime
 - Workload manager integration
 - Image Build
 - Image Conversion
- Environment Definition File
- Container Hooks













Pathfinder Implementation

Environment Definition File – Highlights

- Key Concept for Containers-First
- Standardization and Simplification
- Reproducible Environments
- Customization and Extensibility
- Integration with HPC Tools



Example EDF

```
base environment = [
 2
       "parent_env0",
 3
       "research_group_env"
 4
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
6 containerfile = "./Containerfile.devtools"
 7
  workdir = "/workdir"
  entrypoint = false
 9 writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
      "${SCRATCH}/data:/workdir/data"
14
15
16 [env]
17 MELLANOX VISIBLE DEVICES = "none"
18 CONFDIR = "/workdir/configs"
19 UNSET_BEFORE_RUNNING = ""
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23 com.hooks.ofi_nccl.version = "cuda12"
```



```
base environment = |
       "parent_env0",
 3
       "research group env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
       "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
  MELLANOX VISIBLE DEVICES = "none"
17
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
  com.hooks.cxi.enable = true
23
  com.hooks.ofi nccl.version = "cuda12"
```

- Line 1 4: Defining a base environment from other Env.
- Inherit settings from parents
- Build upon existing environments



```
base environment = [
       "parent_env0",
       "research_group_env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
      "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
  MELLANOX VISIBLE DEVICES = "none"
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
  com.hooks.cxi.enable = true
23
  com.hooks.ofi nccl.version = "cuda12"
```

- Line 5: Specifies the container image to use.
- Example:
 - NVIDIA HPC container
 - from NVIDIA's container registry
 - development tools and CUDA support, based on Ubuntu 22.04



```
base environment = [
       "parent_env0",
       "research_group_env"
 3
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
       "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
  MELLANOX VISIBLE DEVICES = "none"
17
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23
  com.hooks.ofi nccl.version = "cuda12"
```

- Line 6: Points to a Containerfile that defines how to build the container image.
- Basically specify a Dockerfile
- Build dockerfile rather than pulled



```
base environment = [
       "parent_env0",
      "research_group_env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
      "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
17
  MELLANOX VISIBLE DEVICES = "none"
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23
  com.hooks.ofi nccl.version = "cuda12"
```

- Line 7: Sets the working directory inside the container.
- When the container starts, this is the directory from which it will operate.



```
base environment = [
       "parent_env0",
 3
       "research_group_env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
       "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
  MELLANOX VISIBLE DEVICES = "none"
17
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23 com.hooks.ofi nccl.version = "cuda12"
```

- Line 8: Disables the default entrypoint script provided by the container image.
- Override the start-up behavior of the container.



```
base environment = [
       "parent_env0",
 3
       "research_group_env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
       "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
  MELLANOX VISIBLE DEVICES = "none"
17
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23 com.hooks.ofi nccl.version = "cuda12"
```

- Line 9: Allows the file system within the container to be writable.
- By default containers are immutable
- Enable modifications at runtime



```
base environment = [
       "parent_env0",
 3
       "research_group_env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
      "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
  MELLANOX VISIBLE DEVICES = "none"
17
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23 com.hooks.ofi nccl.version = "cuda12"
```

 Lines 10-14: Specifies bind mounts, mapping directories from the host system into the container.



```
base environment = [
       "parent_env0",
 3
       "research_group_env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda_multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
       "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
  MELLANOX VISIBLE DEVICES = "none"
17
  CONFDIR = "/workdir/configs"
18
  UNSET BEFORE RUNNING = ""
19
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23 com.hooks.ofi nccl.version = "cuda12"
```

- Lines 16-19: Sets environment variables within the container.
- To configure how applications behave inside the container.



```
base_environment = [
       "parent_env0",
 3
       "research_group_env"
  image = "nvcr.io/nvidia/nvhpc:23.11-devel-
      cuda multi-ubuntu22.04"
  containerfile = "./Containerfile.devtools"
 6
  workdir = "/workdir"
  entrypoint = false
 9
  writable = true
10 \text{ mounts} = [
11
      "/user/experiment:/workdir",
12
      "/configs:/workdir/configs",
13
       "${SCRATCH}/data:/workdir/data"
14
15
16
  [env]
17 MELLANOX VISIBLE DEVICES = "none"
  CONFDIR = "/workdir/configs"
18
  UNSET_BEFORE_RUNNING = ""
19
20
21
  [annotations]
22
  com.hooks.cxi.enable = true
23
  com.hooks.ofi nccl.version = "cuda12"
```

CSCS

• Lines 21-23: Annotations that can be used to trigger hooks.

EHzürich

- CXI
- NCCL

Container Ecosystem Integration for HPC

- WLM integration and runtime based on Pyxis and Enroot project
- WLM enhancements
 - TOML parsing for EDF support
 - --environment option
 - EDF_PATH variable
 - Advanced feature support
- Runtime enhancements
 - Direct image deployment
 - Layer caching
- User workflow simplification
 - Extended tooling integration
 - Extended container capabilities via hooks





Running example

\$ srun --environment=pytorch.toml --pty bash

- On instantiation, statically linked SSH server starts as daemon in the container
- VS Code editor and debugger available
 - Edit, run, break and inspect the program state on CN











TH zürich

Enhancing functionality and usability

Type of Hook	Description
Slurm	Integrates Slurm workload management capabilities within containers, allowing for the use of Slurm commands directly in containerized jobs.
Libfabric for Slingshot 11	Ensures compatibility with the HPE Slingshot 11 high-speed network, enabling high-performance networking capabilities.
NVIDIA GPU	Provides access to NVIDIA GPUs, crucial for simulations and deep learning models.
OFI NCCL Plugin	Integrates the NVIDIA Collective Communications Library (NCCL) with libfabric, for inter-GPU communications for distributed computing tasks.
SSH	Enables SSH access within containers, facilitating remote debugging and management of containerized applications.
MPI Replacement	Replaces the container's MPI installation with an ABI-compatible version from the host.







Use-cases

Use-cases

Weather and Climate

```
image = "${SCRATCH}/fcn-collaboration/ml+
      fourcastnet+22.12.sqsh"
2
3
  mounts = [
       "${SCRATCH}/fcn-collaboration:/workdir",
4
       "${SCRATCH}/fcn-results/:/output",
5
       "${SCRATCH}/init_dir_fcn:/init_dir",
6
7
       "${SCRATCH}/fcn-collaboration/
      profiler_output:/profiler_output",
8
       "${SCRATCH}/m01/p levels raw:/input"
9
10
11
  workdir = "/workdir"
12
13
  [annotations]
14 com.hooks.cxi.enable = true
15 com.hooks.ofi_nccl.version = "cuda11"
```

Listing 3: Example of an EDF used to train FourCastNet on a CSCS cluster.





Use-cases

Machine Learning

```
1
  mounts =
23
       "${HOME}/.bash_history:/.bash_history",
       "${SCRATCH}:${SCRATCH}",
4
       "/tmp:/tmp"
5
6
7
  [env]
8
  HISTFILE = "/.bash_history"
9
10
  [annotations]
11
  com.hooks.cxi.enable = true
12
  com.hooks.ofi_nccl.version = "cuda12"
```

Listing 6: Example of an EDF providing system settings to be reused as defaults.

1 base_environment = "cscs-mlp"
2 image = "nvcr.io/nvidia/pytorch:24.01-py3"

Listing 7: Example of an use-case specific EDF using a base environment to inherit system defaults on a CSCS cluster for Machine Learning.







Conclusion

Conclusion

- Addresses the current needs of User Environments
- Isolates user environments from the system software stack
- Reproducibility and streamlining deployment processes
- Current pathfinder implementation
 - o Accessible,
 - Modular,
 - o Immediate,
 - and extensible
- For details, examples, and discussion see the paper!











Thank you! Felipe A. Cruz (cruz@cscs.ch)