

Containerizing All the Things: Using Containers for Research

Robert Kalescky

HPC Research Scientist

Adjunct Professor of Data Science

Research Technology Services

Office of Information Technology

Southern Methodist University



Research Support

Containers

Docker

Apptainer/Singularity

Spack

Kubernetes

Research Support



| Domain | Name | Email |
|--|---------------------|--|
| Data Science | Dr. Eric Godat | egodat@smu.edu |
| High-Performance Computing | Dr. Robert Kalescky | rkalescky@smu.edu |
| | Dr. John LaGrone | jlagrone@smu.edu |
| Machine Learning & Artificial Intelligence | Dr. Tue Vu | tuev@smu.edu |
| Custom Devices (IOT, wearables, etc.) | Guillermo Vasquez | guillermov@smu.edu |

Table 1: The OIT Research Technology Services team provides research computing support, consultations, and collaborations.



- Provides research computing tools, support, and training to all faculty, staff, and students using research computing resources
- help@smu.edu with [HPC] as part of the subject line
- Documentation and account requests at https://southernmethodistuniversity.github.io/hpc_docs/

Containers



Before Containerization



- Goods had to be loaded and unloaded individually
- Inefficient - it was not uncommon to spend more time loading and unloading goods than transporting them
- Insecure - goods had to be handled by many people, increasing the chance for loss and theft
- Inaccessible - Long distance shipping only available to the wealthy





- Standardized - containers are all the same size and weight allowances
- Efficient - containers are easy to load and unload and transfer to other modes of transportation
- Secure - goods may be secured in containers from source to final destination
- Available - cost effective to ship goods across the world





- My software doesn't build on this system...
- I'm missing dependencies...
- I need version 1.3.2 but this system has version 1.0.2..
- I need to re-run the exact same thing 12 months from now...
- I want to run this exact same thing somewhere else...
- I want my collaborators to have the same exact software as me...
- I've heard about these Containers, can I just run that?
- Can I run docker on this HPC system?



- It's common to run on multiple systems with different requirements
- We would like to avoid installing the same sets of software again and again
- We would like other people to run our software without our help
- We would like to preserve a known configuration that our software works in



- What are Containers?
- Uses a combination of Kernel “cgroups” and “namespaces” to create isolated environments
- Long history of containers Solaris Zones (2005), LXC(2008), LMCTFY/Google and then Docker(2013).
- Entire ecosystem has grown around containers including open standards and governance.



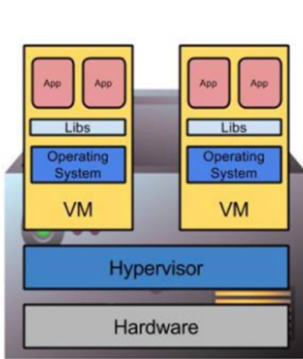
- A lightweight collection of executable software that encapsulates everything needed to run an application
 - Minus the OS kernel
 - Based on Linux only
- Processes and all user-level software is isolated
- Creates a portable* software ecosystem
- Think **chroot** on steroids
- Docker is the most common tool today
 - Available on all major platforms
 - Widely used in industry
 - Integrated container registry via Dockerhub



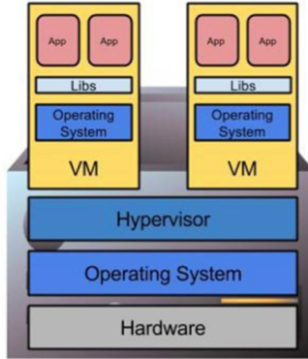
- Containers offer the ability to run fully customized software stacks, *e.g.* based on different Linux distributions and versions
- Containers are not virtual machines, where an entire hardware platform is virtualized, rather containers share a common kernel and access to physical hardware resources



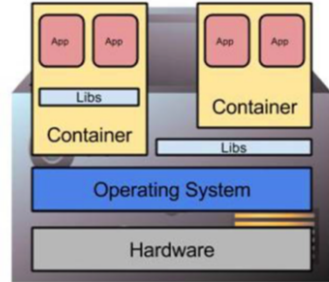
- Type 1 hypervisors insert layer below host OS
- Type 2 hypervisors work as or within the host OS
- Containers do not abstract hardware, instead provide “enhanced chroot” to create isolated environment using a common kernel
- Location of abstraction can have impact on performance
- All enable custom software stacks on existing hardware



Type 1 Hypervisor



Type 2 Hypervisor



Containers



Performant Containers can perform at near native performance.

Flexible Install (almost) any software you need.

Reproducible Define complex software environments that are verifiable.

Compatible Built on open standards that works on all major Linux distributions.

Portable Build once and run (almost) anywhere.



- Hardware** Containers are limited to the same CPU architecture (x86_64, ARM, Power, etc.) and binary formats
- Software** Requires glibc and kernel compatibility between host and container. Other kernel level APIs may also need to be compatible (e.g. CUDA/GPU drivers, network drivers, etc.)
- Filesystem** Paths can be different when viewed from inside or outside of a container



Image A read-only template that defines how to create a container

Container An instantiation of an image, a running instance

Container Runtime Tool or service to execute and manage containers

Registry A service that is used to store and distribute images



- We don't allow direct Docker use on SMU HPC systems
- Docker's security model is designed to support users "trusted" users running "trusted" containers (e.g. users who can escalate to root access)
- Docker is not designed to support scripted / batch based workflows
- Docker is not designed to support parallel applications



- Containers are a single image file
- No root owned daemon processes
- User inside containers are the same as users outside the container (no contextual changes)
- Supports shared, multi-user environments
- Supports HPC hardware such as GPUs and Infiniband networks
- Supports HPC applications like MPI



- Converting Docker containers to Apptainer/Singularity
- Building and running software that require newer systems and libraries
- Running commercial software binaries that have specific requirements

Docker





- Is a “recipe” for how to construct an image.
- Starts **FROM** a defined base image.
- Several basic commands (**ADD**, **COPY**, **RUN**, etc.) can be applied to mutate the image to the desired state.
- Metadata labels can also be added to provide information about the image.
- In addition to the file system changes, the Dockerfile can also control settings like the environment, the starting directory, and default commands.



```
1 FROM ubuntu:20.04
2
3 ENV DEBIAN_FRONTEND noninteractive
4
5 RUN apt-get update &&\
6     apt-get -y install\
7     python3-pip\
8     python3-numpy\
9     python3-pandas
10
11 RUN pip3 install\
12     jupyterlab
13
14 ENTRYPOINT ["python3"]
15
```



- There are several public and private sources for Docker images.
- Images can be used as the base image for custom images.
- Already optimized images can help with reproducible and efficient development workflows.



Docker <https://hub.docker.com>

GitHub Packages <https://github.com/features/packages>

Quay.io <https://quay.io>

NVIDIA <https://catalog.ngc.nvidia.com>

Intel <https://hub.docker.com/u/intel>

AMD <https://hub.docker.com/u/amdih>



```
1  #!/usr/bin/env sh
2
3  pull_and_check() {
4      name=${1}
5      tag=`echo ${name} | cut -d':' -f2`
6      docker pull ${name}
7      docker image ls | egrep "REPOSITORY|${tag}"
8      docker run --rm -it ${name} bash
9  }
10
11  images[0]="ubuntu:jammy"
12  images[1]="nvcr.io/nvidia/nvhpc:22.3-devel-cuda_multi-ubuntu20.04"
13  images[2]="nvcr.io/nvidia/nvhpc:20.7-runtime-cuda10.1-centos7"
14
15  for image in ${images[@]}; do
16      pull_and_check ${image}
17  done
18
```



- Images are CPU-architecture specific
- Docker supports multi-architecture builds
 - Platforms: amd64, arm32v5, arm32v6, arm32v7, arm64v8, i386, ppc64le, and s390x
 - `docker build --platform` with single platform
 - `docker buildx --platform` with list of platforms
- Builds on non-native platforms will be slower as it is running via emulation



```
1  #!/usr/bin/env sh
2
3  # Create builder to build images
4  docker buildx create --name builder --use
5
6  # Build images for x86_64 and ARM64
7  docker buildx build --no-cache --platform\
8    linux/amd64,linux/arm64 -t rkalescky/python3:latest\
9    -f python3.dockerfile --push .
10
11 # Inspect the built images
12 docker buildx imagetools inspect rkalescky/python3:latest
13
```



- Images with build tools can be very large.
- Use the needed image for building.
- Use the smallest image for running.
- Define both the build and execution in a single Dockerfile.



```
1 FROM nvcr.io/nvidia/nvhpc:22.3-devel-cuda_multi-ubuntu20.04 as builder
2 WORKDIR /build
3 COPY hello_world.cpp ./
4 RUN nvc++ -Bstatic -o hello_world hello_world.cpp
5
6 FROM alpine:3.15.4
7 WORKDIR /opt/hello/bin
8 COPY --from=builder /build/hello_world ./
9 ENTRYPOINT ["/opt/hello/bin/hello_world"]
10
```




```
1  #!/usr/bin/env sh
2
3  # Build image using multi-stage Dockerfile
4  docker build -t hello:20.04 -f hello_world.dockerfile .
5
6  # Run the image
7  docker run hello:20.04
8
9  # Note the size difference
10 docker image ls | egrep "hello|22.3-devel"
11
```

Apptainer/Singularity



- Singularity has it's own image **definition language**.
 - Requires (re)writing the definition file.
 - Requires root or “fakeroot”, which is not widely available on HPC systems.
 - Can be done on a Linux system with Singularity installed and then copying the image.
 - Not generally recommended as there would be two definition files to maintain, presumably Docker and also Singularity.
- Pull from Docker registries.
 - Requires pushing and pulling of Docker images.
- Build from Docker archives.
 - Requires exporting, copying, and conversion of Docker images.



```
1  #!/usr/bin/env sh
2
3  # Build images for both x86_64 and ARM64
4  . ./docker_buildx.sh
5
6  # Consume the image on M3 via Singularity
7  ssh m3 'bash -l -c "module load apptainer\'
8  && apptainer run docker://rkalescky/python3 -c \'import numpy as np;
   ↪  print(np.pi)\'"'
```



```
10  # Pull Docker image to a Singularity image
11  ssh m3 'bash -l -c "module load apptainer\
12      && apptainer pull -F python3_3.9.13-slim.sif docker://python:3.9.13-slim\
13      && ls -lh ./python3_3.9.13-slim.sif\
14      && apptainer exec ./python3_3.9.13-slim.sif python3 -c \"import sys;
    ↪   print(sys.version)\""'
15
16  # Singularity mount points
17  ssh m3 'bash -l -c "module load apptainer\
18      && echo \"$APPTAINER_BIND"'
19
```



- Build your Singularity containers on a local system you have root or sudo access. Alternatively build a Docker container
- Transfer your container to the HPC system. If you used Docker, you will need to convert the image
- Run your Singularity containers



```
15  # Export, upload, convert, and run on M3 via Singularity
16  docker save python3:20.04 | ssh m3 'bash -l -c "n=python3_20_04\
17    && cat > ~/$n.tar\
18    && module load apptainer\
19    && apptainer build -F $n.sif docker-archive:$HOME/$n.tar\
20    && apptainer shell $n.sif'
```



- Build application, separating build and deployment containers
- Script container image build and make Lmod module file for ease of use
- Convert Docker image to Apptainer/Singularity image using the docker2singularity tool



```
1  # docker build -t smuresearch/molden:latest .
2  # docker run -it smuresearch/molden:latest
3
4  FROM ubuntu:22.04 AS build
5
6  # Tarball available at https://www.theochem.ru.nl/molden/howtoget.html
7  ENV MOLDEN_VERSION=6.9
8  ENV DEBIAN_FRONTEND noninteractive
9
10 RUN apt-get update &&\
11     apt-get install -y\
12     build-essential\
13     gfortran\
14     libx11-dev\
15     libglu1-mesa-dev\
16     freeglut3-dev\
17     xutils-dev\
18     vim
```



```
20 COPY molder${MOLDEN_VERSION}.tar.gz /
21
22 # `makefile` Edits
23 # 1. -fallow-argument-mismatch to allow for type mismatches
24 # 2. Add quotes due to above flag.
25 # 3. Don't do desktop integration steps
26 RUN tar -xf molder${MOLDEN_VERSION}.tar.gz &&\
27   cd molder${MOLDEN_VERSION} &&\
28   sed -i '/FC = gfortran/{s/$/ -fallow-argument-mismatch/}' makefile &&\
29   sed -i 's/FC=${FC}/FC="\${FC}"/g' makefile &&\
30   sed -i 's/\$(EXTENZ\{0,1\})//g' makefile &&\
31   cat makefile &&\
32   make -j all
```



```
34 FROM ubuntu:22.04
35
36 ENV MOLDEN_VERSION=6.9
37 ENV DEBIAN_FRONTEND noninteractive
38
39 RUN apt-get update &&\
40     apt-get install -y\
41     libgfortran5\
42     libx11-6\
43     libglu1-mesa\
44     freeglut3 &&\
45     apt-get clean &&\
46     rm -rf /var/lib/apt/lists/*
47
48 COPY --from=build /molden${MOLDEN_VERSION}/bin /usr/local/bin/
```



```
1  #!/usr/bin/env zsh
2
3  # Project-specific variables
4  name="molden"
5  version="6.9"
6
7  # Stop on failure
8  set -e -o pipefail
9
10 # Set log file
11 exec > >(tee build.log) 2>&1
```



```
13  # Set platform
14  case $1 in
15  native)
16      platform=$(uname -m)
17      ;;
18  *)
19      platform="amd64"
20      ;;
21  esac
```



```
23  # Singularity image name
24  img="${name}_${version}_${platform}_${date "+%Y_%m_%d_%H_%M_%S"}.sif"
25
26  # Build container with Docker
27  docker build --no-cache --progress=plain\
28    --platform linux/${platform} -t ${name}:${version} .
```



```
30  # Convert Docker image to Singularity image
31  docker run -v /var/run/docker.sock:/var/run/docker.sock\
32    -v $PWD:/output --privileged -t --rm singularityware/docker2singularity\
33    -n ${img} ${name}:${version}
34  mv ${img%.sif}.simg ${img}
```



```
36  # Change Singularity image permissions
37  if [[ $(uname -s) == "Linux" ]]; then
38      sudo chown $USER:$USER $img
39  fi
40
41  # Update module file with new Singularity image name
42  sed -i'' -e "s/^local img_name.*/local img_name      = '${img}'/g"\
43      ${name}.lua
```




```
1  load("apptainer")
2
3  local img_name      = 'molden_6.9_amd64_2024_04_02_14_31_08.sif'
4  local img_directory = '/hpc/m3/containers/'
5  local img_path      = pathJoin(img_directory, img_name)
```



```
7  function build_command(cmd)
8      local cmd_beginning = 'apptainer exec '
9      local cmd_ending    = img_path .. ' ' .. cmd
10     local sh_ending      = ' "$@" '
11     local csh_ending      = ' $*'
12     local sh_cmd          = cmd_beginning .. cmd_ending .. sh_ending
13     local csh_cmd         = cmd_beginning .. cmd_ending .. csh_ending
14     set_shell_function(cmd, sh_cmd, csh_cmd)
15 end
```



```
17 local executables = {
18     'ambfor',
19     'ambmd',
20     'gmolden',
21     'gmolden.exe',
22     'molden',
23     'surf'
24 }
25
26 for _, executable in ipairs(executables) do
27     build_command(executable)
28 end
```

Spack



- Build images defined by Spack environments.
- Spack-based build optimizations are preserved.
- Intermediate Dockerfile uses multi-stage builds
- Currently does not work for multi-architecture builds.



```
1  spack:
2    specs:
3      - gromacs+mpi
4      - mpich
5    container:
6      format: docker
7      images:
8        os: "ubuntu:20.04"
9        spack: develop
10
```



```
1  #!/usr/bin/env sh
2
3  # Define the Spack environment
4  cat spack.yaml
5
6  # Build container definition file
7  spack containerize > Dockerfile
8
9  # Build the container image
10 docker build -t gromacs:latest .
11
```

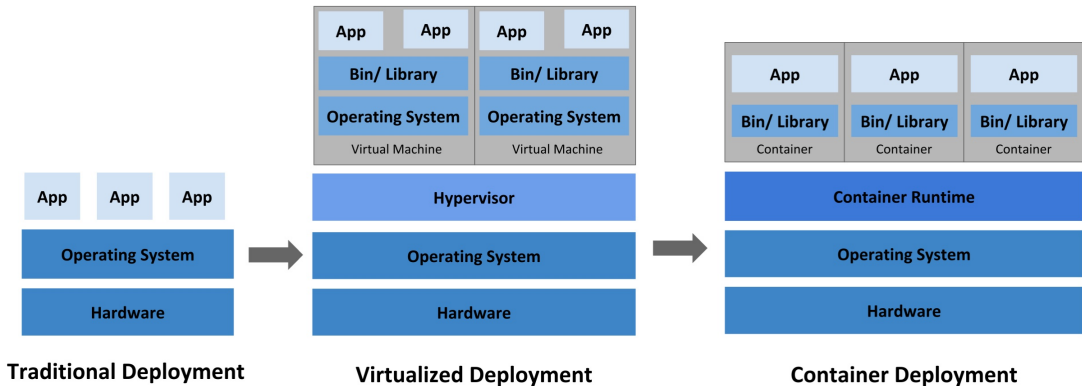
Kubernetes



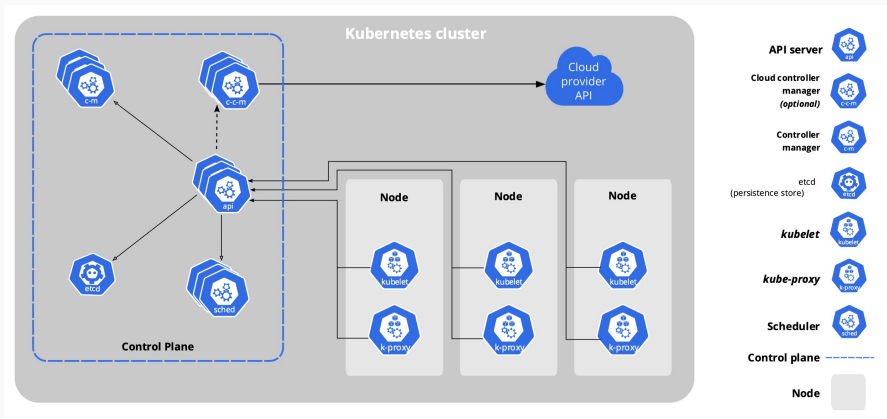
- Platform for managing containerized workloads and services
 - Portable
 - Extensible
 - Open Source
 - Originally developed by Google for managing service deployment
- Kubernetes (K8s)
 - Kubernetes, from Greek, means helmsman or pilot
 - “K8s” comes from the 8 characters between the “k” and the “s”



- Persistence
- Load balancing
- Self-healing
- Automated rollouts and rollbacks
- Resource optimization
- Infrastructure as code



K8s Cluster Components





- kubeadm
- Minikube
 - Local K8s clusters on macOS, Linux, and Windows
 - Commonly used for developing applications
- K3s
- GKE (Google Kubernetes Engine)
- AKS (Microsoft Azure Kubernetes Services)
- EKS (Amazon Elastic Kubernetes Service)
- RKE (SUSE Rancher Kubernetes Engine)
- Ubuntu
 - Charmed Kubernetes
 - MicoK8s





- Job definition
 - One or more containers
 - Specification on how to run them
- Shared resources
 - Compute
 - Storage
 - Networking
 - Context
 - Location and scheduling



Deployments Manage changes to running Pods as specified rate

ReplicaSet Specify the number of available and identical Pods for high-availability

StatefulSets Deployment of Pods with a specific ordering, e.g. Pods with unique network IDs

DaemonSets Every node gets a Pod, which dynamically changes as the cluster size changes

Jobs Batch scheduling of workloads (similar to Slurm Arrays, but with redundancy)

Cronjobs Scheduled workloads



- Specify the Pod level of isolation
- Pod Security Standards
 - Privileged (allows for privilege escalation)
 - Baseline (prevents privilege escalation)
 - Restricted (hardened security)
 - Security-critical applications
 - Low-trust users
- Levels of enforcement
 - Enforce
 - Audit (log annotation)
 - Warn (user-facing warning)



Need help or have questions?

rkalescky@smu.edu

jlagrone@smu.edu

help@smu.edu (include HPC in the subject line)