# Container Technology

FH Sustainable Computing Workshop

Christoph Wissing (DESY) March, 2025





### Containers: Game changer in global economics and modern IT

### **Treatment of Goods – Hamburg Harbour ~1950**

#### Individual packages



- Individual sacks, bags, vessels
- Many works needed in the harbours



#### Photo references: Deutsche Fotothek, Germany - In Copyright - Educational Use Permitted. https://www.europeana.eu/de/item/437/item\_FRTGT7KZTCNQBKLYEYK5UZDSCK36N7E6 Deutsche Fotothek, Germany - In Copyright - Educational Use Permitted. https://www.europeana.eu/de/item/440/item\_SD7XJGL3OJOLQTLDK5GTCUTXOSUCTOGT Deutsche Fotothek, Germany - In Copyright - Educational Use Permitted. https://www.europeana.eu/de/item/440/item\_RWV43GBGVPDMWLJOJ4NLL4IJACI7UL5F

https://pixabay.com/de/users/dendoktoor-14802912/

### **Treatment of Goods Today**

#### **Containers, Containers, Containers**



- Backbone of global shipping: ISO Containers
- Introduced end of the 1950s
- ISO Containers
  - 20" (TEU 20 foot Equivalent Unit): 5,898 m x 2,352 m x 2,393 m
  - 40" (FEU 40 foot Equivalent Unit): 12,032 m x 2,352 m x 2,393 m
- Modern container ships transport over 20,000 TEUs
- Highly automised habour logistics
- Layup times several hours compared to days before



"https://pixabay.com/de/users/stocksnap-894430/

### **Containers: Game Changer in IT**

#### **Again Containers everywhere**



### **kubernetes** Web page https://kubernetes.io



#### **Planet Scale**

Designed on the same principles that allows Google to run billions of containers a week, Kubernetes can scale without increasing your ops team.

- Kubernetes is a container orchestration frame work
  - Start (and stop) containers in an elastic mode depending on demand
  - Designed for scalability
- Let's get started from the beginning...

### Container Lecture

- Compute Architectures
- Operating Systems
- Virtual Machines
- Techniques for Container Virtualisation
- (Brief) Intro to two Container Ecosystems:
  - Docker
  - Singularity
- Some Hand On



#### Container Lecture

- Compute Architectures
- Operating Systems
- Virtual Machines
- Techniques for Container Virtualisation
- (Brief) Intro to two Container Ecosystems:
  - Docker
  - Singularity
- Some Hand On



#### Most topics fill a full semester of in computer science!

#### How to get around with this in ~2h???

We learn it the 'physicist way'

- Skip academic precision at many places
- Understand things to level it required for our science
- That's the way we also learn maths, statistics, chemistry, engineering, electronics ...

Today we will be extra sloppy!

### **Computer Architectures**

### **Common CPU Families**

#### **x86**

- Most common CPU type for desktops and laptops since decades
- Very common also for servers and High Performance Computers
   In the more recent past only, perhaps 20 years
- Main (competing) manufactures: Intel and AMD



https://pixabay.com/de/users/elasticcomputefarm-1865639/

### ARM

- Used in almost any smartphone or tablet
- Increasing use also in servers, desktops even HPCs
  - Recent Apple Macintosh 'M1,M2 ...' is basically ARM
- Many manufactures

#### **Classical Workstation architectures**

- In the 1980-90ies each major vendor had its own
  - SUN: UltraSPARC
  - HP: PA-RISC
  - SGI: MIPS
  - Motorola: 68k Used in early Macs, Amiga, Atari ST
- Of relevance still today: POWER by IBM

### **From Code to Executable Programs**

Compiler



#### Source code

- Written in programming language (e.g. C, FORTRAN, Python, Go ...)
- Simple enough code can be compiled for different architectures

### **Compiler & Assembler**

- In practice one step
- Assembler code architecture dependent (could in principal optimized by a developer)
- (Some) Compiler can generate assembler for a different architecture (*cross-compiler*)

#### **Machine Code**

• Can be executed only on the architecture it was compiled for

### **Operating System**

#### Usually an Operating System (OS) controls a modern computing device

- Not strictly required a specially prepared application could run directly (very inconveniently)
- OS controls all resources
- Applications run on top of OS
  - Common functionality comes from (core) libraries, e.g. glibc in Linux
- Modern OSes provide support for multiple users





### **Examples of popular Operating Systems**

### **Microsoft Windows**

- Most popular OS for desktop PCs
- Market share over 90%
- Available (mainly) for x86
- · Windows 11

### Android

- Most popular OS for mobile devices (phones...)
- Relatively young, appeared in 2008
- Based on Linux
- Mostly used on ARM architecture (but also x86 and others)



### MacOS

• OS for Apple Macintosh (Usage on other hardware technically possible – but violated license)



- Runs on various architecture
  - Motorola 68k (until early 1990s)
  - PowerPC (until early 2000s)
  - X86 (until ~early 2020s)
  - M1,M2 ... ARM64 based (recent)

### Linux

Strictly speaking: Linux is a Kernel no OS!



- OSes are the distributions
   Popular: Ubuntu, Debian, CentOS, Fedora, Arch
- Leading OS for Server, HPCs but also embedded devices
- Originally for x86, ported to many architectures

### **Virtual Machines**

### "A software computer"

### Virtualisation is a key technology in IT

- Flexible abstraction layer
- Resources can be separated
- Virtualisation goes nowadays way beyond Virtual Machine: There is virtualisation for e.g. storage or network

### Virtual Machine Manager (VMM) or Hypervisor

- Controls (and separates) the individual Virtual Machines (VMs)
- Offers the guest OS inside the VM the "virtualised hardware"
  - Can pass-through to real Hardware
  - Can provide e.g. a software version of a processor different from the real HW
  - VMM Type 1: Runs directly on the hardware, no OS
  - VMM Type 2: Runs as an application inside an OS





### Virtualisation is (conceptually) not new

### **Historical Remark**

- Concepts invented ~50 years ago
- Early implementations in operation systems for main frames
- Implementations for PCs came much later
  - Limitations in early x86 processors
  - Not need in typical PC environment?



TYPE I VCS

Page 25

### **Examples for popular Virtualisation Solutions**

### VMWare vSphere

- Vendor has many other solutions
- VSphere ESXi hypervisor:



- Type 1VMM running directly on bare metal
- No OS or any other software running in parallel
- Certainly a data center solution
- VMs are controlled from dedicated control application vCenter

### KVM/QEMU



- Kernel Virtual Machine (KVM) part of Linux kernel
- QEMU: Quick EMUlation
  - Virtualises full computer
  - Can emulate a processor (e.g. POWER on x86)
  - Typically makes use of KVM for x86 VMs on x86
- Type 2 VMM

### **Oracle VirtualBox**

- Came to Oracle when they bought SUN, who bought VirtualBox before
- Type 2 VMM for x86\_64
- Supports many OSs as Hypervisor (VMM)
  - Linux, MacOS, Windows, Solaris
- Popular on Desktops/Laptops

### **Parallels Desktop**

- Popular commercial solution to run VMs on MacOS
  - e.g. run Windows as VM on MacOS
- Type 2 VMM



### **Containers vs Virtual Machines**

#### Virtual machine provides a "full computer" in software

Container virtualsation layer is much thinner

- Encapsulation of environment
- Employs kernel of the host operating system
- Container has to *match the host operating system and the host architecture* A Linux container only runs on Linux

There are products offering Linux containers on Windows. How about that?

- Behind the scenes works "Windows Subsystem for Linux"
- It's actually a light weight Virtual Machine



Infrastructure



### **Comparison Virtual Machines vs Container?**

Criteria	Virtual Machine	Container
Resources	large	light to medium
Boot-up Time	long (many seconds)	short (less than a second)
Performance	lowish, if more than 1 VM	(almost) native speed
Scaling	difficult	easy
Efficiency	low	high
Portability	medium	high



nttps://twitter.com/b0rk/status/1237744128450072578

### Why are containers so popular?

#### "The usual Problem"



Environments of developer and tester are different: Linux distribution, versions of (other) software

### Why are containers so popular?

### "Or a typical Python analysis chain"



#### Hard to port ...

### Why are containers so popular?

"A Solution"



Container preserves the proper environment and all dependencies - looks the same "everywhere"

### **Container Frameworks**

### Combine low level features into a (easier to use) framework

- Components: chroot, cgroups, kernel namespaces, security frameworks (SELinux or AppAmor), ....
- User tools to build, run and manage containers

### 2 prominent examples:

- Docker
- Likely the most common container framework on Linux (even adopted for Windows and MacOS)
- Was the "break through" for containers (on Linux initially)
- Frequently used by developers
- Apptainer (formerly known as Singularity)
- Very common solution for HPC environments
- Targeted at security
- Most adopted container solution at DESY, particularly the NAF





dotCloud, Inc. https://commons.wikimedia.org/w/index.php?curid=52332268

### **Docker Architecture**



From: https://docs.docker.com/get-started/overview/

### **Docker Architecture**



Modified from: https://docs.docker.com/get-started/overview/

### **Docker: hello-word**

#### Docker runs as a daemon that needs root privileges

• Commands: systemctl status docker, systemctl start docker, systemctl stop docker

#### Very instructive hello-world



Which <u>images</u> are on the host: docker image ls (We see 'hello-world' with a few details listed)

What containers are running: docker ps (Likely no container listed right now)

If you have a Linux(!) laptop or desktop with root access you might play around with Docker.

Using <u>Docker Desktop</u> on Windows or Macs is forbidden at DESY due to license restrictions!!!

### **Apptainer**

### Apptainer is a program – no system daemon (like Docker)

- Images live in directories, typically of the user creating them
- An image can be a single "SIF file" or unpacked (as sandbox) in a sub-directory
- Apptainer can build images from Docker registries or from its own definition files



Figure inspired (and updated) from https://doi.org/10.1371/journal.pone.0177459

### **Apptainer: lolcow saying hello world**

### Apptainer is a program – no system daemon (like Docker)

- Images live in directories, typically of the user creating them
- Hint regarding disk space, when \$HOME is of limited capacity (e.g. AFS home dir at DESY)

export APPTAINER\_CACHEDIR=<Some-path-with-more-space>/apptainer/cache (e.g. DUST on the NAF)





### **Apptainer: Self built container**

### Usually building starts from a thin base-OS container

• Modify container "interactively"

Note: There are limitations, Apptainer's fakeroot feature not always works in all images, so installations via apt-get (Debian or Ubuntu based) or dnf (Alma, Redhat or Fedora based) commands may not work

• Build container from a recipe file

Pull the lolcow container from the Github container registry: apptainer **pull** docker://ghcr.io/apptainer/lolcow

Start a shell in the container:
apptainer shell lolcow latest.sif

Start a command in the container:
apptainer exec lolcow\_latest.sif cowsay "Playing with containers"

Convert the SIF (single file image) to an unpacked 'sandbox": apptainer **build** --sandbox my\_lolcow lolcow\_latest.sif

Start the unpacked container in writable mode to make modifications (--no-home is need, because /afs is missing in the container): apptainer shell --no-home --writable my\_lolcow Inside the container add /afs: mkdir /afs (afterwards the -no-home is not needed when starting with -writable)

Convert back unpacked 'sandbox" to a SIF imalge: apptainer build my\_own\_lowcow.sif my\_lolcow

### **Apptainer: Bindmounts**

#### Essentially mounts a directory of the into the container



Bind mount /data/dust of the host to /mnt in the container:

apptainer exec -B /data/dust:/mnt my\_own\_lolcow.sif ls /mnt
shows
group user

'-B' X:Y: Bindmount X outside container Y inside container

#### Automatic bindmounts in the NAF

- /cvmfs
- \$HOME
- \$PWD
- For other mounts one needs to care

### **Apptainer: Tetris to go for the linux console**

#### Derive from a thin Linux distro "Alpine"

- Alpine is only a few MB
- A few lines of definition file
  - Bootstrap from Alpine Docker image
  - Add micro-tetris package
  - Add a start script

tetris.def:

Bootstrap: docker From: alpine Stage: build

%post
 apk add micro-tetris
 touch /etc/localtime

```
%runscript
   echo "Starting Tetris"
   for s in 5 4 3 2 1; do
      echo -n "$s "
      sleep 1
   done
   /usr/bin/tetris
```

Build image from definition file:

apptainer build --ignore-fakeroot-command my\_tetris.sif tetris.def

```
Run container image:
apptainer run my_tetris.sif
```

It's only a few MB (copy to any Linux box with apptainer and play...)

### Wrap Up

#### Containers are an efficient and lightweight virtualisation approach

- Realized employing recent isolation features of the Linux kernel
- Widely used in DevOps and scientific computing

#### Looked a bit into two commonly used Container Engines

- Docker & Apptainer
- Many things can be achieved using either of them

## **Additional Material**

### **Evolution of Technology for OS-level Virtualisation**



### Very Simple "Container": chroot

### **CHange ROOT**

- Introduce new root entry in existing directory tree
- mkdir ~/mnt
   cd ~/mnt
   tar -xzvf /tmp/minimal\_Debian.tar.gz

# Usual directory tree cat /etc/redhat-release

CentOS Linux release 7.9.2009 (Core)

 # Switch into the chroot-environment sudo chroot /mnt/ cat /etc/debian\_version 10.11



chroot isolates only on filesystem level! In the chroot environment you still see and and even kill other processes launched outsite of the chroot-environment

You might want to play with it using the Hands-On VM

### Linux – namespaces

- Isolation feature in modern Linux Kernels
  - Mount (mnt)
  - Process ID (pid)
  - Network (net)
  - Interprocess Communication (ipc)
  - UTS (UNiX time sharing)
  - User ID (user)
  - And more
- Example

tail -f /tmp/mytestfile.txt

# Get PID of this process e.g. ps aux | grep tail # Plain chroot via unshare (no isolation) sudo chroot ~/fake\_root /bin/bash # killing process from above works # Now with isolation unshare --mount --uts --ipc --net --pid --fork --user --map-root-user chroot ~/fake\_root/ /bin/bash # The above PID cannot be killed!



https://8gwifi.org/docs/linux-namespace.jsp

### Linux – cgroups

### **Control groups**

- Linux kernel feature to limit resource access or consumption:
  - Cpu: Limits or priority for processes to use CPU
  - Cpuset: Which core(s) can be utilized
  - Blkio: Limits block I/O usage of e.g. disks
  - Memory: Memory usage
  - net\_prio: Controls network throughput
  - Find the full list: Is /sys/fs/cgroups
- Process (or groups of processes) can get assignment of certain resources



https://www.redhat.com/ja/blog/world-domination-cgroups-rhel-8-welcome-cgroups-v2

### **Docker: (Simple) self built container - Example**

#### Build own container from a "Dockerfile" (description of an image)

- Simple example: Webserver based on Alpine Linux
  - Caddy: Lightweight but rather featured webserver
  - Alpine Linux: Very small Linux "distribution" very good base for containers size is MBs not GBs



### **Docker: Example continued**

#### Work a bit with the container

- Webserver should be reachable (even from outside of the VM): http://192.168.56.104/ (IP might vary)
- Some more commands

Show running container:

Inspect logging: docker logs mycaddy01

Connect a shell to your container: docker exec -it mycaddy01 In that shell: ps shows only caddy whoami shows that you are root (by default)

#### Stop the container:

docker kill mycaddy01

#### Fully remove the container:

docker container rm mycaddy01

'-it': <u>interactive and terminal</u> attached. Could be used with run instead of '-d'

#### Contact

**DESY.** Deutsches Elektronen-Synchrotron

Christoph Wissing FH – CMS group christoph.wissing@desy.de

www.desy.de