

22/11

SLURM e KUBERNETES practically same things

# ADVANCED CLOUD COMPUTING

COURSE'S STRUCT. → BASIC INSTRUMENTS  
containers & features

- FOCUS ON NETWORKING  
communication btw containers
- STORAGE INTERFACE
- COMPLEX WORKFLOW
- MESSENGING SERVICES 4 communic.

HTTPS never at reverse proxy  
is DECRYPTED

↳ obento al container non cifrata

# CONTAINERS

## DEEP DIVE

# DEFINITIONS

## CONTAINER IMAGE

A file which is pulled down from a Registry Server and used locally as a mount point when starting Containers.

**ATTENTION:** often people say container images talking about repositories (bundle of multiple container Image Layers as well as metadata)

### Formats

- Docker
- Appc
- LXD → ~ VIRTUAL MACHINE
- Open Container Initiative (OCI)
  - most common
  - ↳ DOCKER PULL  
pull a repo with this format

layman terms: container at rest

# DEFINITIONS

## CONTAINER ENGINE

A container engine is a **piece of software that:**

- **accepts user requests** (eg **Command line** options)
- **pulls images**, (repd)
- ~ **runs the container**  
(the **container runtime runs the container**)

Is the software you will mostly use

Container engines:

- **docker**
- **podman**
- RKT

cloud providers, often have their own container engines.

- CRI-O
- LXD.

*thin client  
NO DAEMON  
no need  
to be  
root*

*IMPORTANT  
lightweight  
MINIMUM API TO  
MANAGE KUBERNETES*

*[ CEF : managing storage in HPC  
in some context we can use  
to orchestrate KUBERNETES ]*

# DEFINITIONS

## CONTAINER

A container is the runtime instantiation of a Container

**Image** (Registry in exact terms). UNIX PROCESS + ISOLATE RESOURCES

A container is a **standard Linux process** typically created through a clone() system call instead of fork() or execvp(). Also, containers are often isolated further through the use of cgroups, SELinux or AppArmor.

SUB  
SOFTWARE  
to run

to ISOLATE  
PROCESSES

for KERNEL  
RESOURCES

# DEFINITIONS

## CONTAINER RUNTIME

*run the docker image (clone/Form ...)*

lower level component used in a Container Engine

The OCI Runtime Standard reference implementation is runc.

- **crun**
- railcar
- katacontainers → AMAZON

runtime operations:

- use container mount point
- use ~~container~~<sup>REPOSITORY</sup> metadata
- **start containerized** processes (clone system call)
- **Setting up cgroups**
- Setting up SELinux Policy/App Armor rules

*NOT RESPONSIBLE for NETWORK  
or STORAGE*

# DEFINITIONS

## SYSTEMS CALL

- **clone()**

This system call provide precise control over what pieces of execution context are shared between the calling process and the child process. eg. the caller can control whether or not

- the two processes share the virtual address space,
- the table of file descriptors, and the table of signal handlers.
- child process to be placed in separate namespaces.

→ IMPORTANT

- **fork()**

creates a new process by duplicating the calling process. The new process is referred to as the child process.

- **execvp()** → OLD

replaces the current process image with a new process image.

# DEFINITIONS

STACK of CONTAINER IMAGES

## REPOSITORY

repositories can be approximated with container images, but it's important to realize that these repositories are actually made up of layers and include metadata (manifest.json)



# DEFINITIONS

## REPOSITORY

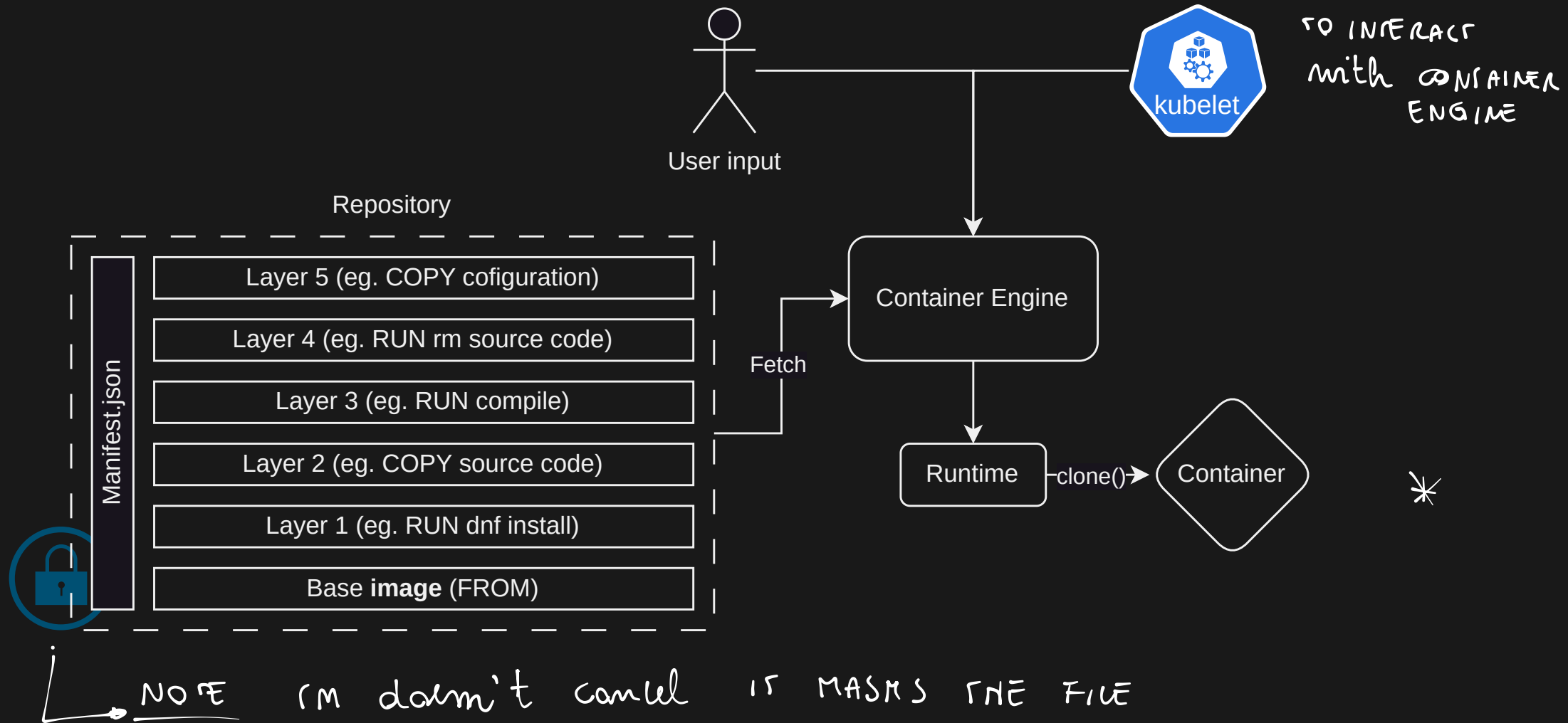
repositories can be approximated with container images, but it's important to realize that these repositories are actually made up of layers and include metadata (manifest.json)

what does it means: `docker pull rhel7`?

- `docker pull registry.access.redhat.com/rhel7:latest`
- `REGISTRY/..optinal groups../REPOSITORY[:TAG]`

*e.g. localhost*

# DEFINITIONS

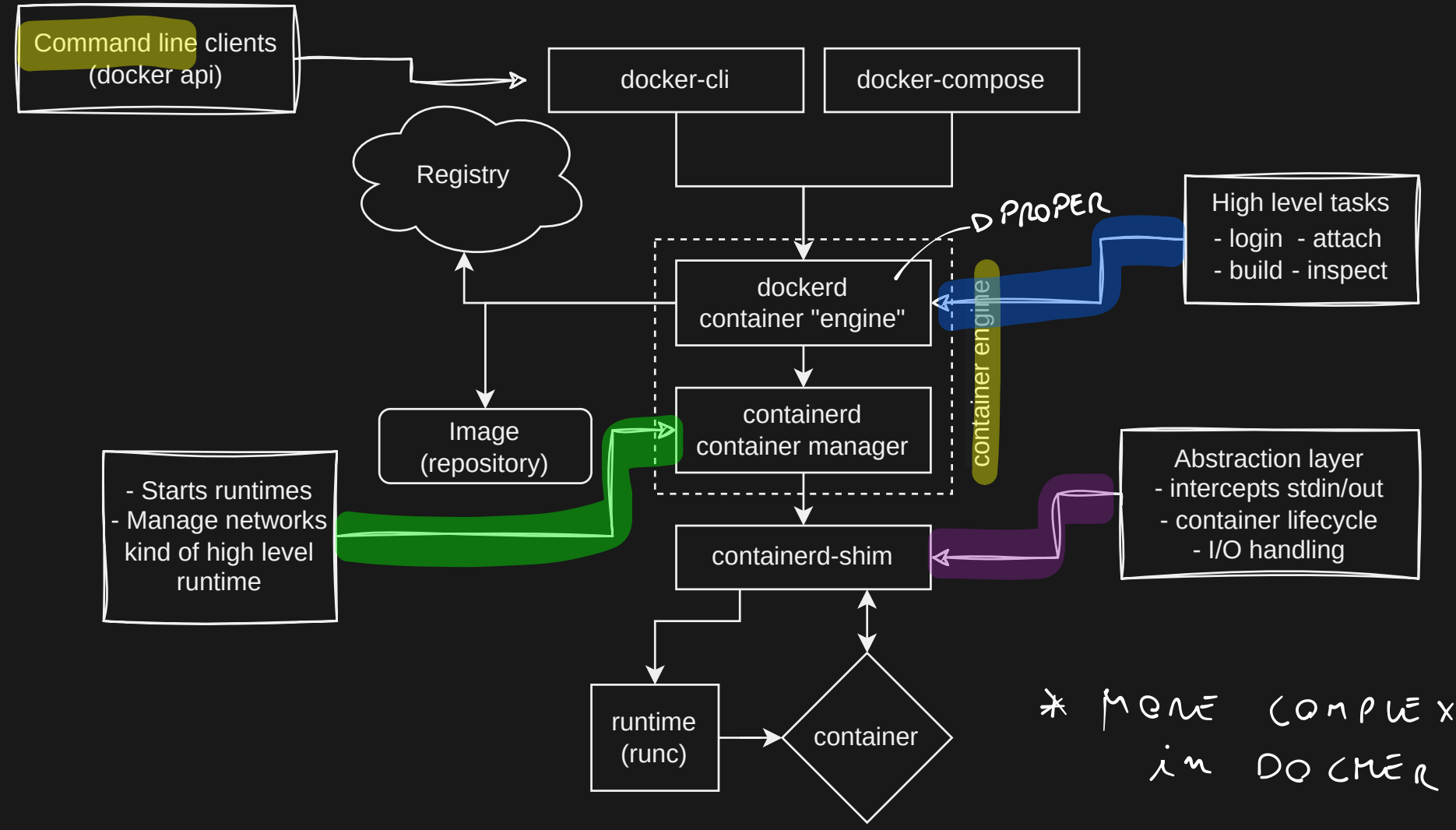


ALL TOGETHER

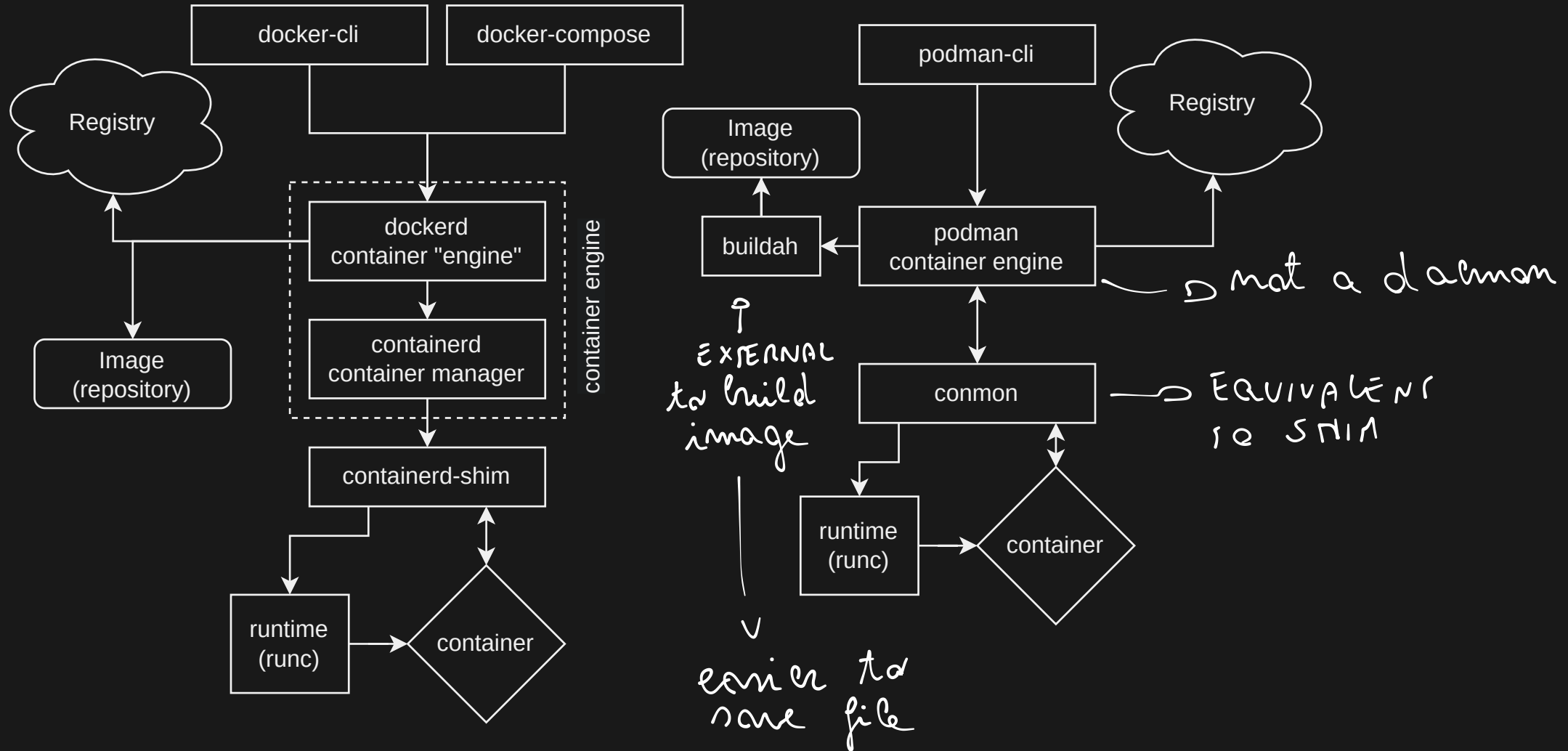
**IMPLEMENTATION**<sub>S</sub>

DOCKER COMPOSE to have in a YAML FILE  
all the containers I  
want for my service

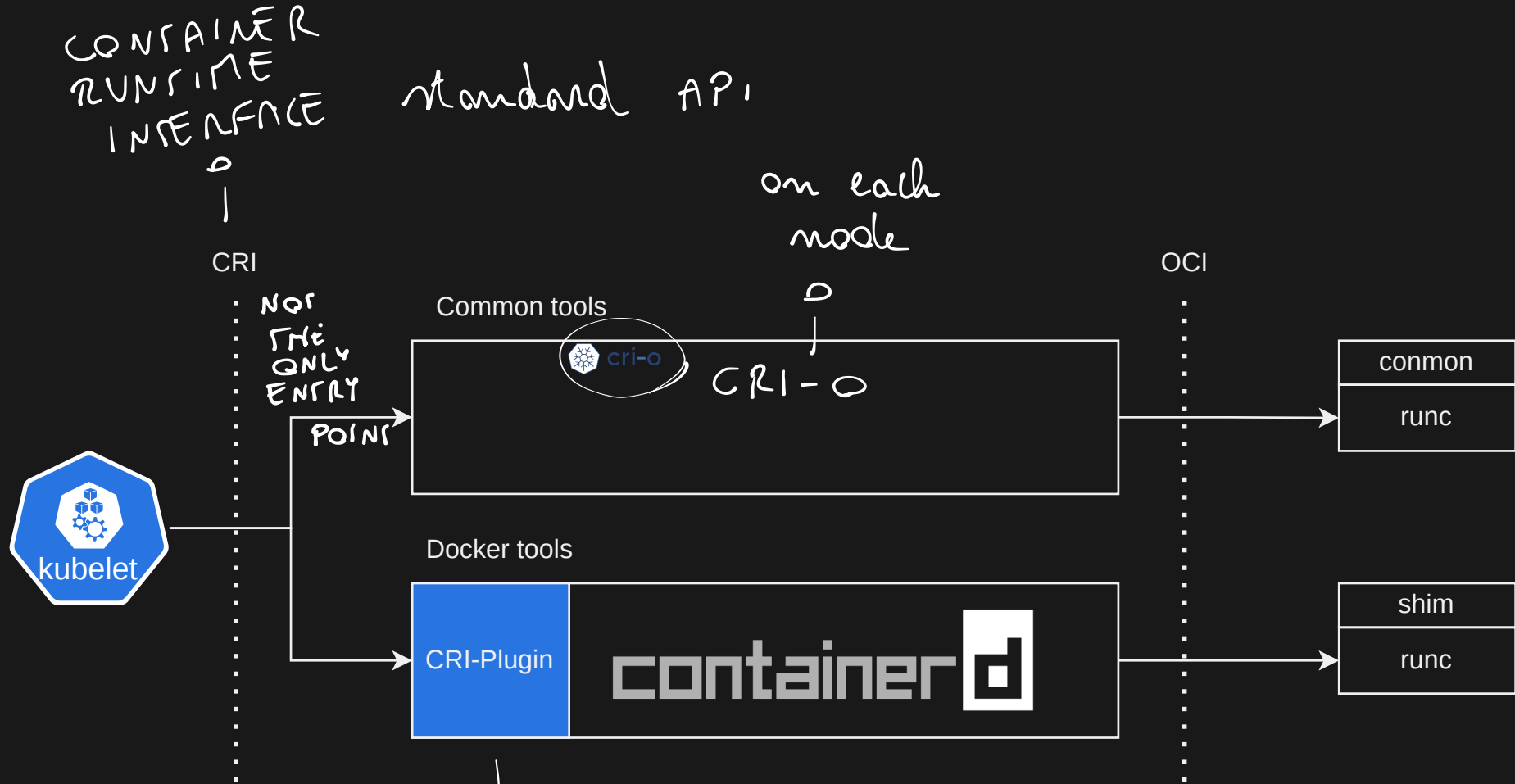
# DOCKER STACK



# DOCKER VS PODMAN



# KUBERNETES



not dockerd  
~ DOCKER BUT MANAGED  
BY KUBERNETES

We can use different  
managers with different  
container runtimes

# LINUX NAMESPACES

- Feature of the linux kernel
- partition kernel resources

a set of processes that shares a namespace set has access to the same resources

even COMBINATION  
of NAMESPACES

NOTE we're building  
a container

PID number identifier

## COMMON NAMESPACES

① INF SYSTEM

A PROCESS IN  
A SET OF  
NAMESPACES

IN A  
CONTAINER

→ NAMESPACES  
ENSURE  
ISOLATION

- Process ID (PID)

- net
- uts
- user
- mnt
- IPC
- cgroups

Isolate the PID number space.

Processes in different PID ns can have the same PID.

From an host point of view the process has two PIDs, one inside and one outside the process ns.

- PID 1 inside each ns
- PID ns can be migrated between hosts
- a process can only view processes inside the ns and sub ns.

if I create a new namespace  
I can put inside a different PID 1



WHEN I INITIALIZE A POD  
network namespace is empty!

## COMMON NAMESPACES

- PID
- network
- uts
- user
- mnt
- IPC
- cgroups

Each net ns has its own resources

- network devices
- IP addresses
- IP routing tables — ROUTING DATA PATH
- nftables
- /proc/net

we'll do it by hand

# COMMON NAMESPACES

- PID
- net
- Unix Timesharing System (uts)
- user
- mnt
- IPC
- cgroups

old name, now is only used for  
hostname segregation

TIME

from a security  
perspective  
IMPORTANT

UID user  
GID group

0 } root  
0 }

mapping user outside  
to root inside the  
container

## COMMON NAMESPACES

- PID
- net
- uts
- user
- mnt
- IPC
- cgroups

Most important. Isolate the user  
and group ID number spaces. A  
process user and group can be  
different inside and outside the ns.

Interest case: having a user ID 0  
inside the namespace but being  
unprivileged outside

# COMMON NAMESPACES

- PID
- net
- uts
- user
- mount
- IPC
- cgroups

Isolate mountpoints to hide files,  
and show different filesystem  
hierarchy. Can be interpreted a bit  
as chroot.

# COMMON NAMESPACES

- PID
- net
- uts
- user
- mnt
- interprocess communication (IPC)
- cgroups

IPCs handle the communication between process using shared memory areas and POSIX message queues and semaphores.

PROCESS COMMUNICATION  
SHARED MEMORY

# COMMON NAMESPACES

- PID
- net
- uts
- user
- mnt
- IPC
- cgroups

repoints the /sys/fs/cgroup folder

- allow for ns migration between hosts
- Avoid leaking sensitive information about the host's resource.

kernel exposes features through files ACCESSIBLE THROUGH VIRTUAL FILESYSTEMS (not real files)  
→ MOUNT POINTS

# CAPABILITIES

categories of processes:

- privileged processes (effective user ID is 0), bypass all kernel permission checks
- unprivileged processes (effective UID is nonzero), subject to full permission checking

privileges associated with superuser are now divided into distinct units, known as capabilities.

PASWD on linux

**Capabilities are a per-thread attribute.**

# WHAT IF YOU WANT TO EXECUTE PRIVILEGED OPERATIONS?

- legacy approach: setuid bit, this allow the user to run the program as the program owner.

```
→ ~ ls -la /usr/bin/passwd  
-rwsr-xr-x. 1 root root 91624 Oct 15 02:00 /usr/bin/passwd
```

READABLE WRITABLE SET UID

- novel approach: leverage capabilities



# CAPABILITY SETS

- Permitted

limiting superset for the effective capabilities that the thread may assume.

- Effective

capabilities used by the kernel to perform permission checks for the thread

- inheritable
- bounding
- ambient

# CAPABILITIES

38 capabilities are responsible for controlling syscall

- CAP\_SYS\_BOOT allow **rebooting** without sudo
- CAP\_IPC\_LOCK allow memory allocation with huge pages via mmap
- CAP\_WAKE\_ALARM Trigger something that will wake up the system
- CAP\_NET\_ADMIN perform network operations  
OPERATION INSIDE NAMESPACE

# WHY SHOULD I CARE?

capability set are tied to a user namespace

this means that each namespace will have its own set of capabilities, however a child namespace can never be granted more capabilities than the creating process. you will have to fix them sometimes on kubernetes

# IN PRACTICE

```
1 → ~ getpcaps $$
2 1021810: cap_wake_alarm=i
3 → ~ getpcaps ①
4 1: =ep
5 → ~ unshare -UinpmrC -f /bin/bash
6 bash-5.2→ whoami
7 root
8 bash-5.2→ getpcaps $$
9 1: =ep
10 bash-5.2→ exit
```

*what capabilities i have?*

*PPID*

# NAMESPACES

# USER NS

CONF INSIDE POD  
all 7 namespaces  
for both  
CONTAINER INSIDE  
6 namespaces (different)  
1 same network

User ns isolate security-related identifiers and attributes

- user IDs/group IDs,
- keyrings  $\rightarrow$  where the password (s) is saved
- capabilities.

| A process's user and group IDs can be different inside and outside a user ns.

MAPPING      INSIDE  $\leftrightarrow$  OUTSIDE  
0                      1000  
you can change

# FFECT OF CAPABILITIES

Having a capability inside a user ns permits a process to perform operations only on resources governed by that ns.

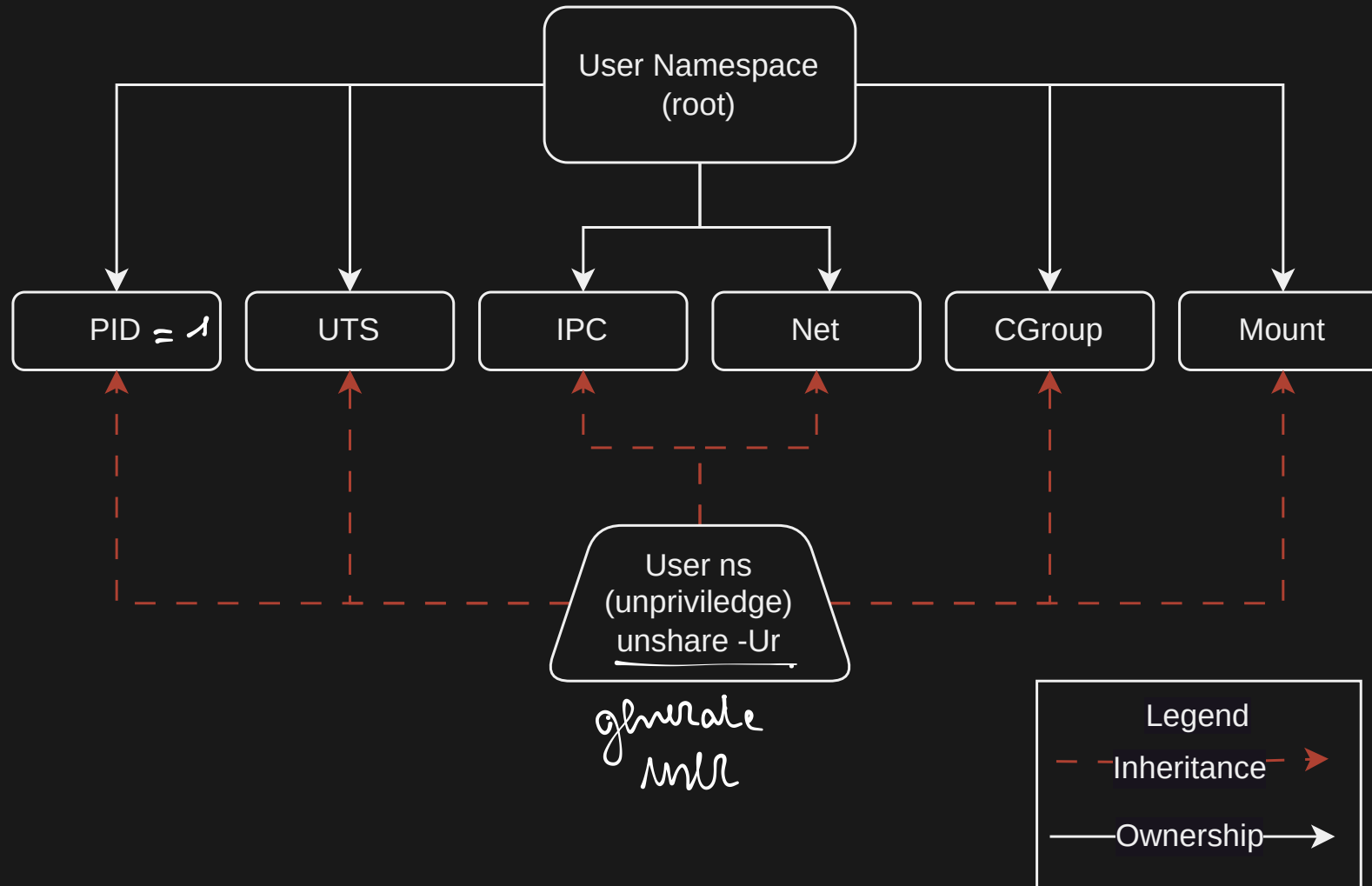
## INTERACTION OF USER NS AND OTHER TYPES OF NS

- unprivileged processes can create user ns
- when nonuser ns is created, it is owned by the user ns in which the creating process is a member.

Privileged operations on resources governed by the nonuser ns require that the process has capabilities in the user ns that owns the nonuser ns.



# HANDSON



# HANDSON



[Stream online](#)



# MOUNT NS

Mount ns provide isolation of the list of mounts seen by the processes in each ns instance. Thus, the processes in each of the mount ns instances will see distinct single-directory hierarchies.

*if I want to mount a filesystem  
I NEED THE CAPABILITY TO MOUNT*

# MOUNT NS

Unexpected behaviour: **a new mount ns is not empty**

actions taken on a poorly configured mount ns will impact the host.

# MOUTPOINT PROPAGATION

ADVANCED

Mountpoints propagates between mount ns because of the **shared subtree** feature.

modify mount point propagation:

- **shared**: A mount that belongs to a peer group. Any changes will propagate to all members.
- **slave**: One-way propagation. master propagate events to a slave, but not viceversa.
- **shared and slave**: Indicates that the mount point has a master, but it also has its own peer group. The master will not be notified of changes to a mount point.
- **private**: Does not receive or forward any propagation events.
- **unbindable**: cannot be bind mounted.

# HANDSON



[Stream online](#)





# HANDSON

