

Containers & Kubernetes

In the [second lecture](#) we have talked about how applications interact with the kernel and the hardware. By now you know that applications running in ring 3 do not have direct access to things like the disk. To access those details they need to execute a [system call](#) to the kernel. The kernel will then give the application the details required, for example a file from the disk.

Let's play a hypothetical game: process A wants to access a file called `test.txt`. It calls the kernel, the file is opened, and process A is now free to read the contents. When process B comes along and does the same, it will receive the same file.

Or does it? In the Linux kernel the filesystem root directory (starting with `/`) is a virtual construct that can be changed on a per-process basis. It is possible to make process A see a completely different folder structure than process B. This mechanism is called [chroot](#). It is, in fact, so old that it predates Linux by more than a decade and was first present in Version 7 Unix in 1979.

Try it out yourself!

Run the following command on a 64 bit x86 Linux in an empty folder:

```
mkdir -p ./bin &&\ncurl "https://www.busybox.net/downloads/binaries/1.21.1/busybox-x86_64" -o ./bin/busybox && \nchmod +x ./bin/busybox && \n./bin/busybox --install ./bin && \nsudo chroot . /bin/sh
```

You should now be able to run `ls -la` or similar commands and see that you are in an almost empty filesystem with no ability to access the files outside.

This alone only isolates the two process on a filesystem level, there are still plenty of opportunities for two processes to interact, for example:

- On the network
- [Interprocess Communication \(IPC\)](#)
- [Sending signals to each other by process ID](#)
- etc.

In the early 2000's there were two projects that attempted to implement process isolation in the Linux kernel: [OpenVZ](#) and [Linux-VServer](#). Running these projects required patching and recompiling the Linux kernel from source.

Starting around 2006 several companies, including Google, IBM, SGI, Bull, and the OpenVZ project itself has put significant effort into taking OpenVZ apart and piece by piece and submitting it to the Linux kernel. The [User Beancounters](#) from OpenVZ, for example, became [cgroups](#) and allow resource allocation to a group of processes.

Do you want to know more about history?

The history of containerization is surprisingly hard to piece together despite the relative young age of the technology. If you want a bit of a deeper dive take a look at [“A Brief History of Containers \(by Jeff Victor & Kir Kolyshkin\)”](#).

All these isolation technologies, together, form what's today known as Linux containers. Windows has undergone a similar development in recent years, and has *virtualized* the typical Windows interaction points between processes: job objects, process table, registry, filesystem, etc.

In essence, *containers don't exist*. They are a collection of technologies that provide process-level isolation for an application. These isolation techniques can, among others, include:

- A separate, virtual network card
- Process separation (no ability to send signals to other containers)
- A separate view of the filesystem
- Restricted IPC between containers
- User separation / user ID mapping

Do you want to build your own container engine?

You can! Take a look at [Liz Rice's Containers From Scratch](#) talk!

Container images

[Docker](#) was not the first technology to create containers using the modern Linux API's, but it was the one that made containers popular. The edge Docker had over its competitors was the way it dealt with images. Other container solutions at the time, like LXC, tried to emulate virtual machines where a full operating system would be installed, and then updated. Docker chose a different route. A `Dockerfile` would contain the instructions needed to build a container image containing the application that should be containerized. These instructions would be executed and resulted in a container image. This container image could then be used to launch one or more copies of it.

Since the container image contains everything the application needs to run these images are ideal operating system independent packages. This solves the age-old problem of having to install the correct version of the runtime environment (PHP, Python, Java, etc) as it is contained within the container image.

When a container is run this container image is the basis for all the data that is contained within. The data is, however, not copied. If the application modifies or creates files only the difference is stored. This makes containers very lightweight on the storage front as well.

Docker went one step further, they introduced a way to share images between machines using a *registry*, a server that stores container images. The public [Docker hub](#) contains a vast array of images built by the community. These technologies were later all standardized under the [Open Container Initiative](#), creating a company-independent format for containers and technologies around them.

Warning

Third party container images should be treated with the same due diligence like installing third party software on your computer! They can contain malicious code and can cause harm despite the containerization.

The container lifecycle

You can, of course, update a running container just as you would a traditional virtual machine. However, that is not the intended, or indeed, optimal way to use them. The intention of containers is *immutable infrastructure*.

Immutable infrastructure means that containers are *never updated*. Instead, they are replaced. When a new version of a software needs to be installed the old container is shut down and a new one is launched from a new image.

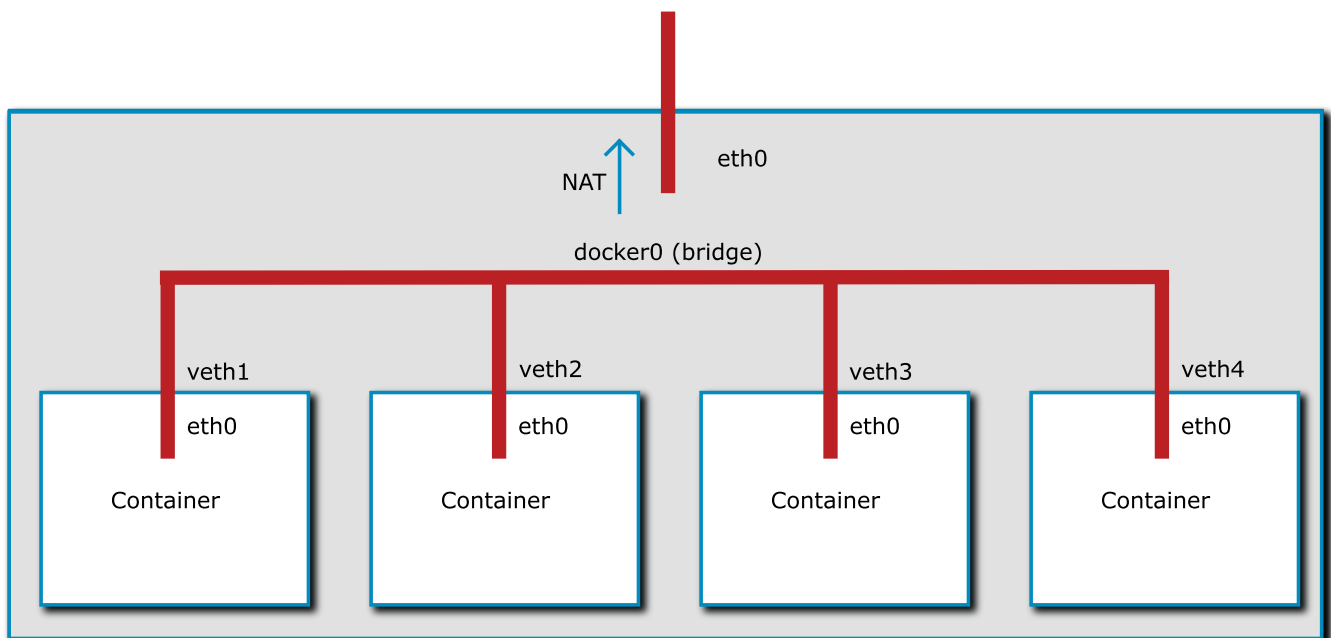
Immutable infrastructure presents a massive benefit: instead of having to deal with a downtime when the upgrade is ran, the updated version can be tested before it is launched.

However, this concept reaches its limits when it comes to running software that needs to store data in a persistent fashion, for example, databases. For this purpose containers can designate special folders to be mounted as *volumes*. Volumes can be persisted by either mounting them as a folder from the host machine, or by mounting a network-based storage system.

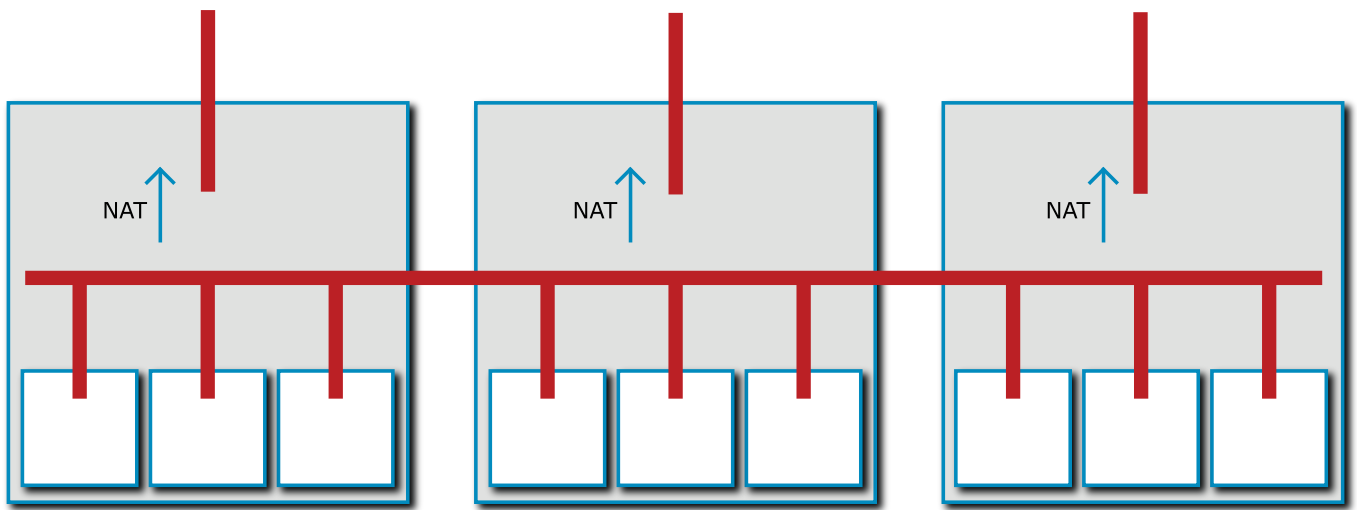
Container networking

As mentioned previously, containers regularly have their own, virtual network interfaces. This virtual network interface can be connected in a number of ways.

If we take a look at Docker's default networking model the virtual network interface is connected to an internal [network bridge](#). This enables containers to connect to each other. For connections to the Internet a [NAT](#) is performed.



More advanced options for container networking create a single virtual network across several hosts and connect the containers to this virtual network. This is an option with [Docker Swarm](#) and is always the case with [Kubernetes](#).



When you look at these graphics, one question may come to mind: the cloud is quite dynamic in its nature, so how does a Kubernetes cluster receive incoming connections? After all, each of the nodes may be replaced with a new one, with a new IP address at any time? This is where the immense feature set of Kubernetes begins to show itself. Kubernetes includes several *controllers* that can interact with cloud providers in a variety of ways. One of these integrations can be, for example, the load balancer integrations. The Kubernetes network can be extended to a load balancer and that load balancer can send traffic directly over the virtual network used by the containers.

Further helping the dynamic nature of the cloud are features like the ability to create internal firewalls. These firewalls can help mitigate possible security vulnerabilities like [SSRF](#). Since these firewalls can be created using the same API as the software deployment, the firewall rules for an application can be versioned in Git, for example.

Since everything required to run an application can be stored in a Dockerfile and a Kubernetes manifest the application source code can be self-contained. This is, of course, only true if the application is written in such a way that works well with the cloud. This is described [in the next lecture](#).

Orchestration (Swarm, Kubernetes, etc)