

Virtualisation: Jails and Unikernels

Advanced Operating Systems
Lecture 18

Lecture Outline

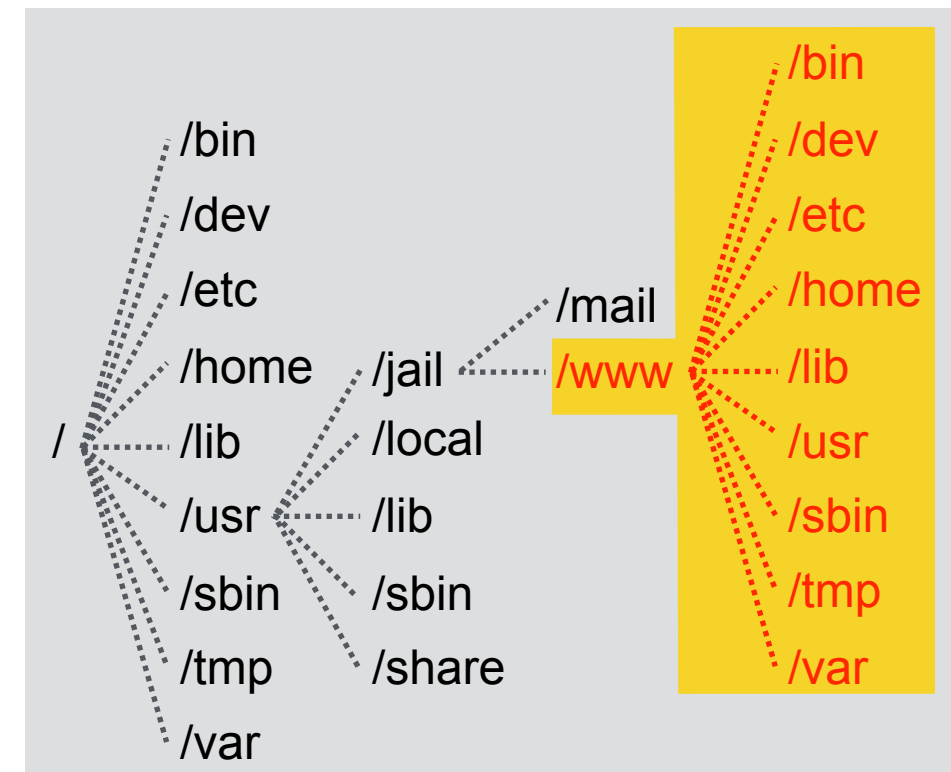
- Alternatives to full virtualisation
- Sandboxing processes:
 - FreeBSD jails and Linux containers
 - Container management
- Unikernels and library operating systems

Alternatives to Full Virtualisation

- Full operating system virtualisation is expensive
 - Need to run a hypervisor
 - Need to run a complete guest operating system instance for each VM
 - High memory and storage overhead, high CPU load from running multiple OS instances on a single machine – but excellent isolation between VMs
- Running multiple services on a single OS instance can offer insufficient isolation
 - All must share the same OS
 - A bug in a privileged service can easily compromise entire system – all services running on the host, and any unrelated data
- A *sandbox* is desirable – to isolate multiple services running on a single operating system

Sandboxing: Jails and Containers

- Concept – lightweight virtualisation
 - A single operating system kernel
 - Multiple services
 - Each service runs within a jail a service specific container, running just what's needed for that application
- Jail restricted to see only a subset of the filesystem and processes of the host
 - A sandbox within the operating system
 - Partially virtualised



Example: FreeBSD Jail subsystem

Jails: Confining the omnipotent root.

Poul-Henning Kamp <phk@FreeBSD.org>

Robert N. M. Watson <rwatson@FreeBSD.org>

The FreeBSD Project

ABSTRACT

The traditional UNIX security model is simple but inexpressive. Adding fine-grained access control improves the expressiveness, but often dramatically increases both the cost of system management and implementation complexity. In environments with a more complex management model, with delegation of some management functions to parties under varying degrees of trust, the base UNIX model and most natural extensions are inappropriate at best. Where multiple mutually untrusting parties are introduced, “inappropriate” rapidly transitions to “nightmarish”, especially with regards to data integrity and privacy protection.

The FreeBSD “Jail” facility provides the ability to partition the operating system environment, while maintaining the simplicity of the UNIX “root” model. In Jail, users with privilege find that the scope of their requests is limited to the jail, allowing system administrators to delegate management capabilities for each virtual machine environment. Creating virtual machines in this manner has many potential uses; the most popular thus far has been for providing virtual machine services in Internet Service Provider environments.

1. Introduction

The UNIX access control mechanism is designed for an environment with two types of users: those with, and without administrative privilege. Within this framework, every attempt is made to provide an open system, allowing easy sharing of files and inter-process communication. As a member of the UNIX family, FreeBSD inherits these security properties. Users of FreeBSD in non-traditional UNIX environments must balance their need for strong application support, high network performance and functionality,

This work was sponsored by <http://www.servetheweb.com/> and donated to the FreeBSD Project for inclusion in the FreeBSD OS. FreeBSD 4.0-RELEASE was the first release including this code. Follow-on work was sponsored by Safepoint Network Services, <http://www.safepoint.com/>

- Unix has long had `chroot ()`
 - Limits process to see a subset of the filesystem
 - Used in, e.g., web servers, to prevent buggy code from serving files outside chosen directory
- Jails extend this to also restrict access to other resources for jailed processes:
 - Limits process to see subset of the processes (the children of the initially jailed process)
 - Limits process to see specific network interface
 - Prevents a process from mounting/un-mounting file systems, creating device nodes, loading kernel modules, changing kernel parameters, configuring network interfaces, etc.
 - Even if running with `root` privilege
- Processes outside the jail can see into the jail – from outside, jailed processes look like regular processes

Benefits of Jails

- Extremely lightweight implementation
 - Processes in the jail are regular Unix processes, with some additional privilege checks
 - FreeBSD implementation added/modified ~1,000 lines of code in total
 - Contents of the jail are limited: it's not a full operating system – just the set of binaries/libraries needed for the application
- Straightforward system administration
 - Few new concepts: no need to learn how to operate the hypervisor
 - Visibility into the jail eases debugging and administration; no new tools
- Portable
 - No separate hypervisor, with its own device drivers, needed
 - If the base system runs, so do the jails

Disadvantages of Jails

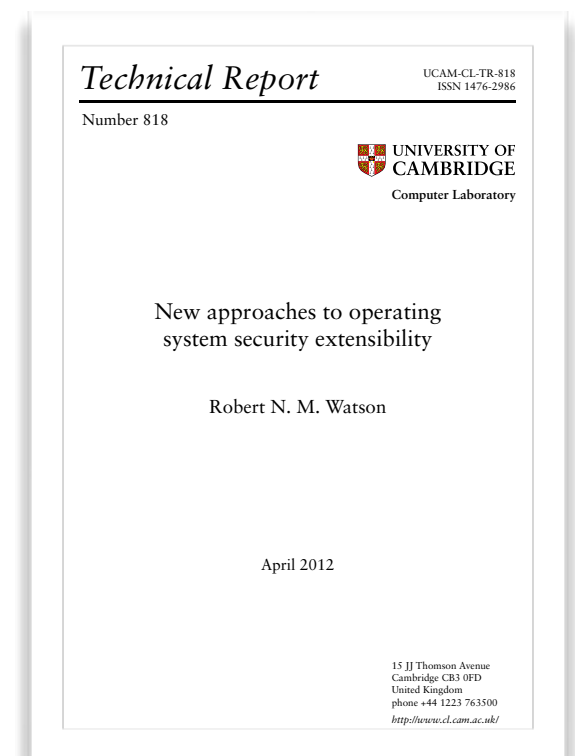
- Virtualisation is imperfect
 - Full virtual machines can provide stronger isolation and performance guarantees – processes in jail can compete for resources with other processes in the system, and observe the effects of that competition
 - A process can discover that it's running in a jail – a process in a VM cannot determine that it's in a VM
- Jails are tied to the underlying operating system
 - A virtual machine can be paused, migrated to a different physical system, and resumed, without processes in the VM being aware – a jail cannot

Example: Linux Containers

- Two inter-related kernel features:
 - Namespace isolation – separation of processes
 - Control groups “cgroups” – accounting for process CPU, memory, disk and network I/O usage
 - Both part of the mainline Linux kernel
- Combined to provide container-based virtualisation
 - Implemented via tools such as “lxc”, “OpenVZ”, “Linux-VServer”
 - Very similar functionality to FreeBSD Jails, but using the Linux kernel as the underlying OS
 - Run a subset of processes within a containers – can be a full OS-like environment, or something much for limited

Example: iOS/macOS Sandbox

- All apps on iOS, and all apps from the Mac App Store on macOS, are sandboxed
 - Sandbox functionality varies with macOS/iOS version – gradually getting more sophisticated and more secure
 - Kernel heavily based on FreeBSD, but sandboxing doesn't use Jails
 - Mandatory access control – <http://www.trustedbsd.org/mac.html>
 - Apps are constrained to see a particular directory hierarchy – but given flexibility to see other directories outside this based on user input
 - Policy language allows flexible control over what OS functions are restricted – restricts access to files, processes, and other namespaces like Jails, but more flexibility in what parts of the system are exposed
 - At most restrictive, looks like a FreeBSD Jail containing a single-process and it's data directories
 - Can allow access to other OS services/directories as needed – more flexible, but more complex and hence harder to sandbox correctly
 - Heavily based on Robert Watson's PhD and work on FreeBSD MAC and Capsicum security framework
 - <http://www.cl.cam.ac.uk/techreports/UCAM-CL-TR-818.pdf>
- Trade-off: flexibility, complexity, security

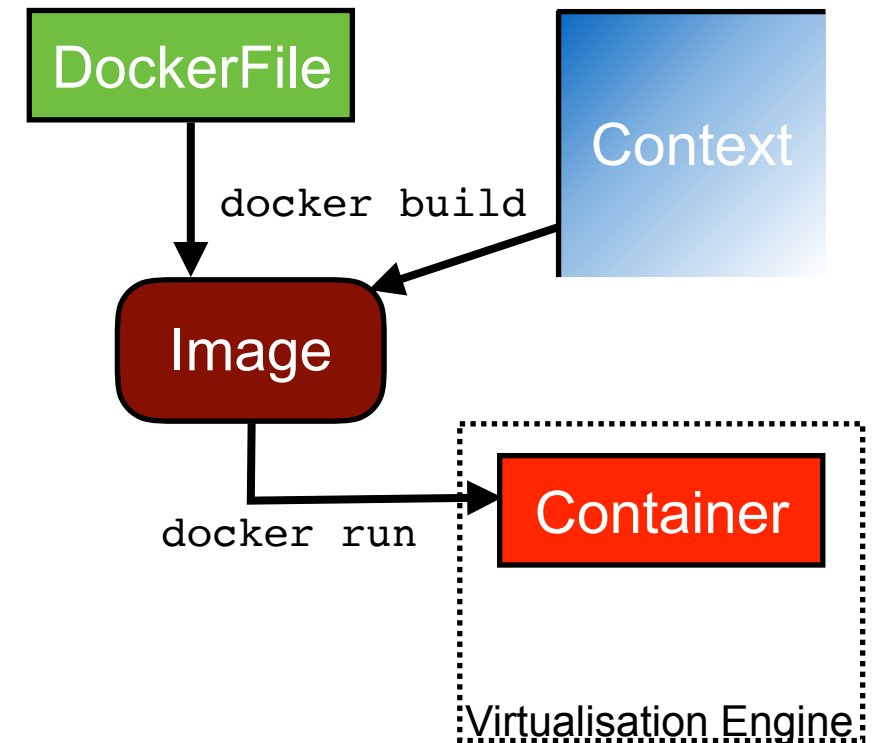


Container Management

- Full virtualisation – VMs run complete OS, and can use standard OS update/management procedures
- Jails run a container with a subset of the full OS
 - Want minimal stack needed to support services running in the container
 - Less software installed in container → security vulnerability less likely
- How to ensure all images are up-to-date?
 - No longer sufficient to patch the OS, must also patch each container
 - Need to minimally subset the OS for the services in the container – can this be automated?
 - How to share service/container configuration?

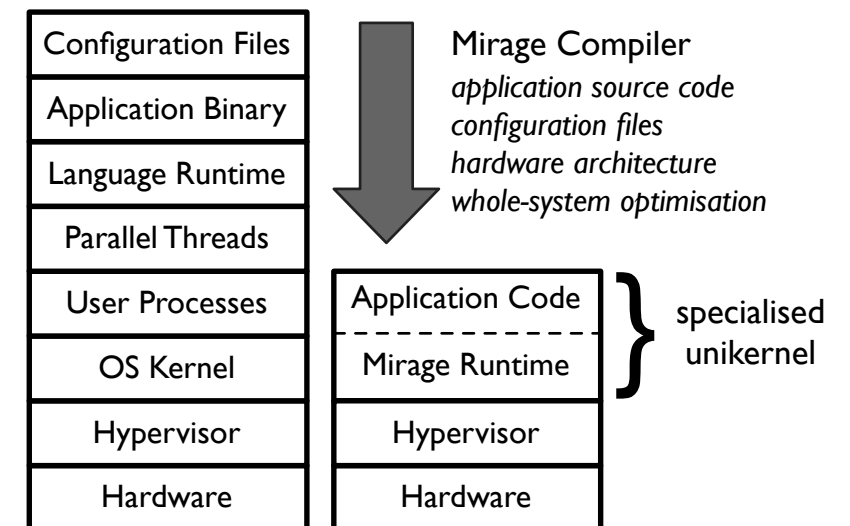
Example: Docker

- DockerFile and context specify an image
 - Base system image
 - Docker maintains official, managed and security updated, images for popular operating systems
 - Images can build on/extend other images
 - Install additional binaries and libraries
 - Install additional data
 - Commands to execute when image is instantiated
 - Metadata about resulting image
- Images are immutable → instantiated in container
 - Virtualisation uses Linux containers, the macOS sandbox, or FreeBSD jails to instantiate images, depending on underlying operating system running the container
 - macOS and Windows run Linux in a hypervisor, and use that to execute containers
 - FreeBSD relies on native Linux system call emulation, and jails
 - Image can specify external filesystems to mount, to hold persistent data accessible by running image
 - A standardised way of packaging a software image to run in a container – plus easy-to-manage tools for instantiating a container



Unikernels: Library Operating Systems

- Taking the idea of service-specific containers to its logical extreme – a library operating system
- Well-implemented kernels have clear interfaces between components
 - Device drivers, filesystems, network protocols, processes, ...
- Implement as a set of libraries, that can link with an application
 - Only compile in the subset of functions needed to run the desired application (e.g., you don't use UDP, it isn't linked in to your kernel, etc...)
 - Link entire operating system kernel and application as a single binary
 - Run on bare hardware, or within a hypervisor
- Typically written in high-level, type-safe, languages, with expressive module systems – not backwards compatible with Unix
 - E.g., MirageOS largely written in OCaml
 - Type-safe languages give clear interface boundaries, to enable automatic minimisation of images



Source: A. Madhavapeddy *et al.*, "Unikernels: Library Operating Systems for the Cloud", Proc. ACM ASPLOS, Houston, TX, USA, March 2013. DOI:10.1145/2451116.2451167

Further Reading

- P.-H. Kamp and R. Watson, “Jails: Confining the omnipotent root”, System Administration and Network Engineering Conference, May 2000.

<http://www.sane.nl/events/sane2000/papers/kamp.pdf>

- Trade-offs vs. complete system virtualisation?
- Overheads vs. flexibility vs. ease of management?
- Benefits of Docker-style configuration of images

- A. Madhavapeddy et al., “Unikernels: Library Operating Systems for the Cloud”, ACM ASPLOS, Houston, TX, USA, March 2013. DOI:10.1145/2451116.2451167

- Is optimising operating system for single application going too far?
- Are unikernels maintainable?
- Relation to containers and hypervisors?

