



# CS 423

## Operating System Design: OS Support for Containers

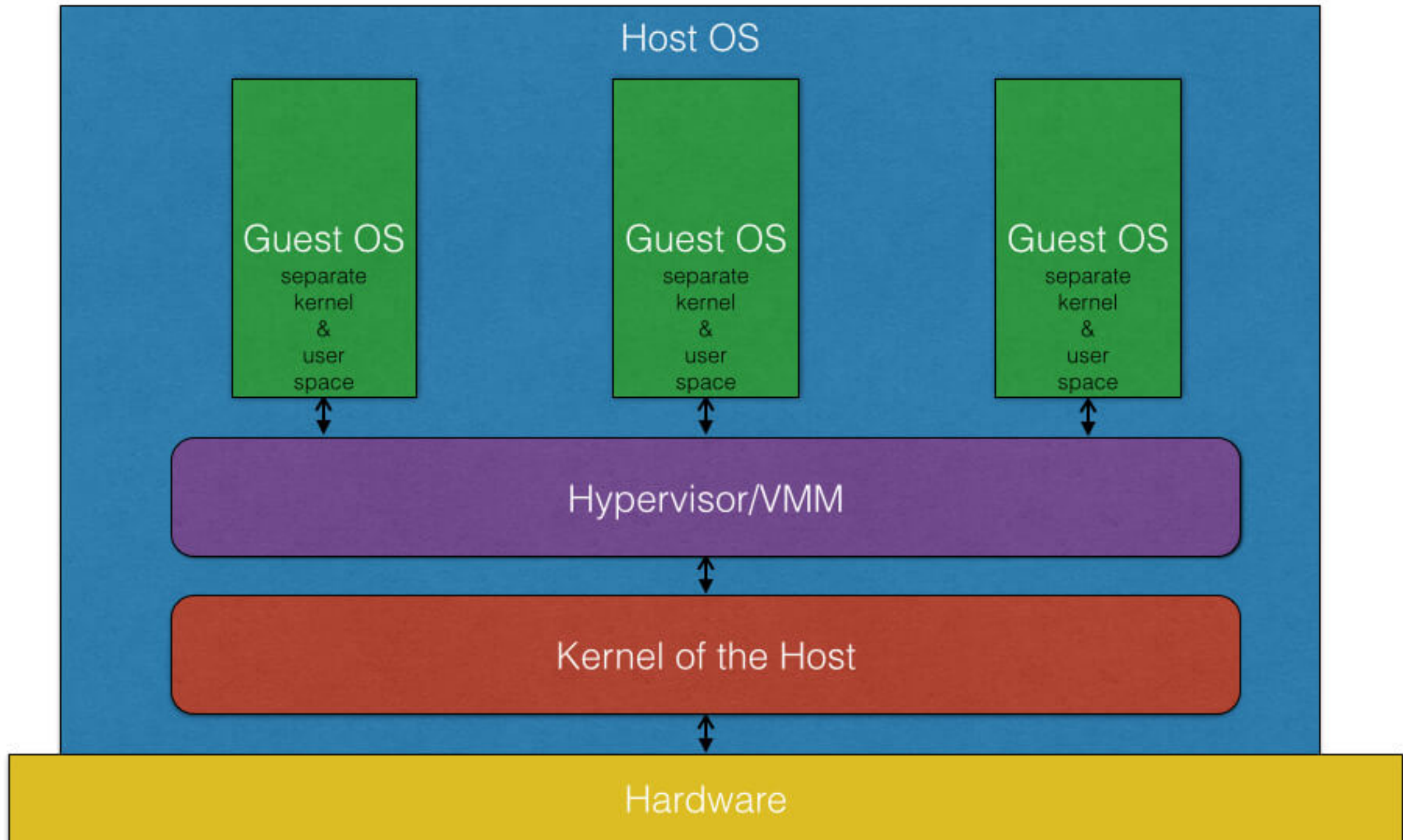
Tianyin Xu

Thanks for Adam Bates and Julie Evans



- Dominated by Infrastructure-as-a-Service clouds (and storage services)
- Big winner was Amazon EC2
- Hypervisors that virtualized the hardware-software interface
- Customers were responsible for provisioning the software stack from the kernel up

# Hypervisors



# Hypervisors



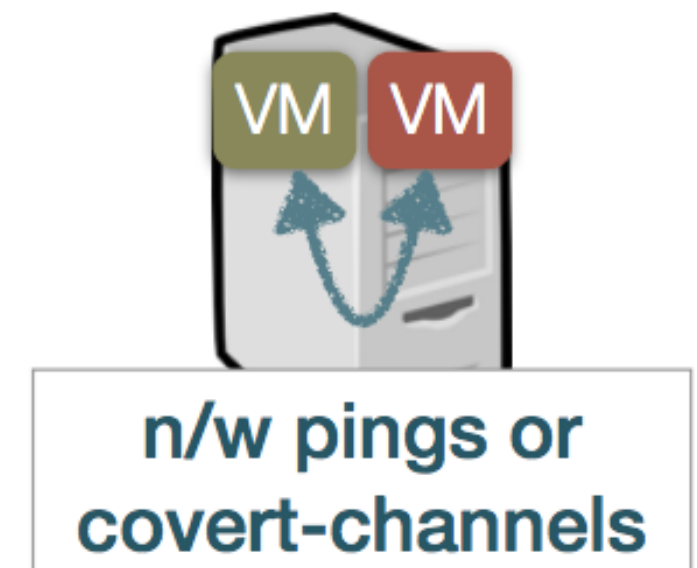
- Strong isolation between different customer's virtual machines
- VMM is 'small' compared to the kernel... less LoC means less bugs means (~)more security.



# Hypervisors



- ‘Practical’ attacks on IaaS clouds relied on side channels to detect co-location between attacker and victim VM
- E.g., we could correlate the performance of a shared resource
  - network RTT’s, cache performance
- After co-resident, make inferences about victim’s activities



# Hypervisors

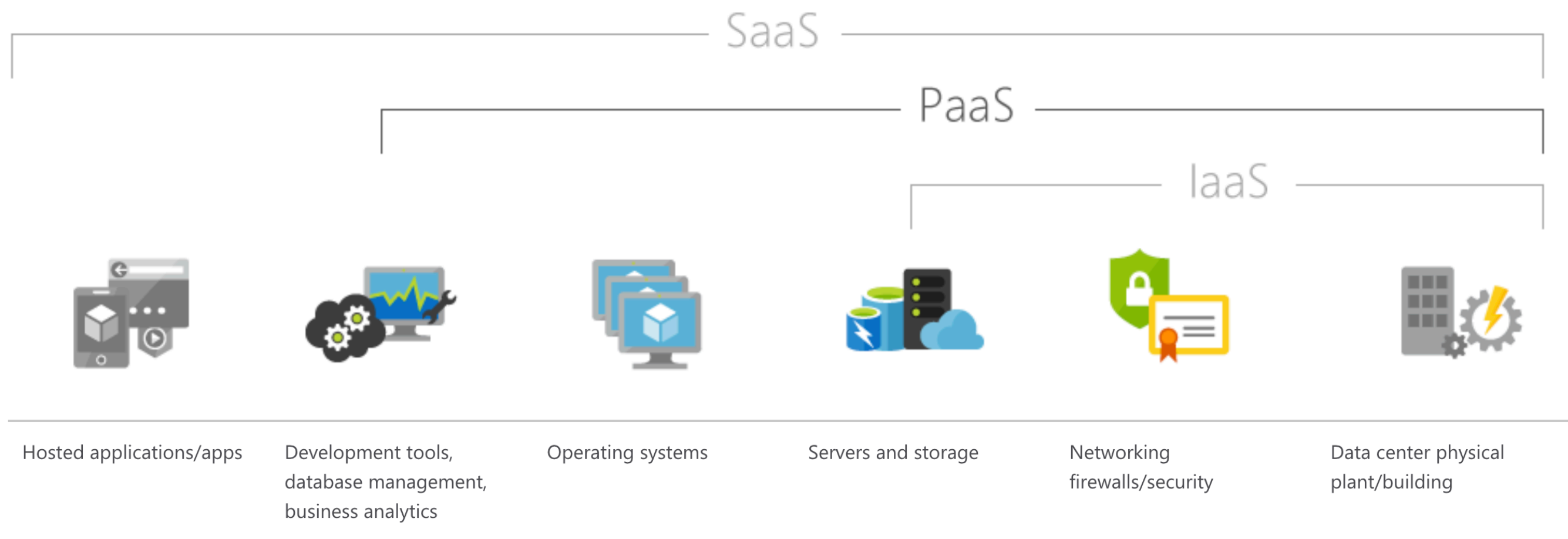


- Strong isolation between different customer's virtual machines
- VMM is 'small' compared to the kernel... less LoC means less bugs means (~)more security.
- High degree of flexibility... but did most customers really need it?

# Cloud Computing (Gen 2)

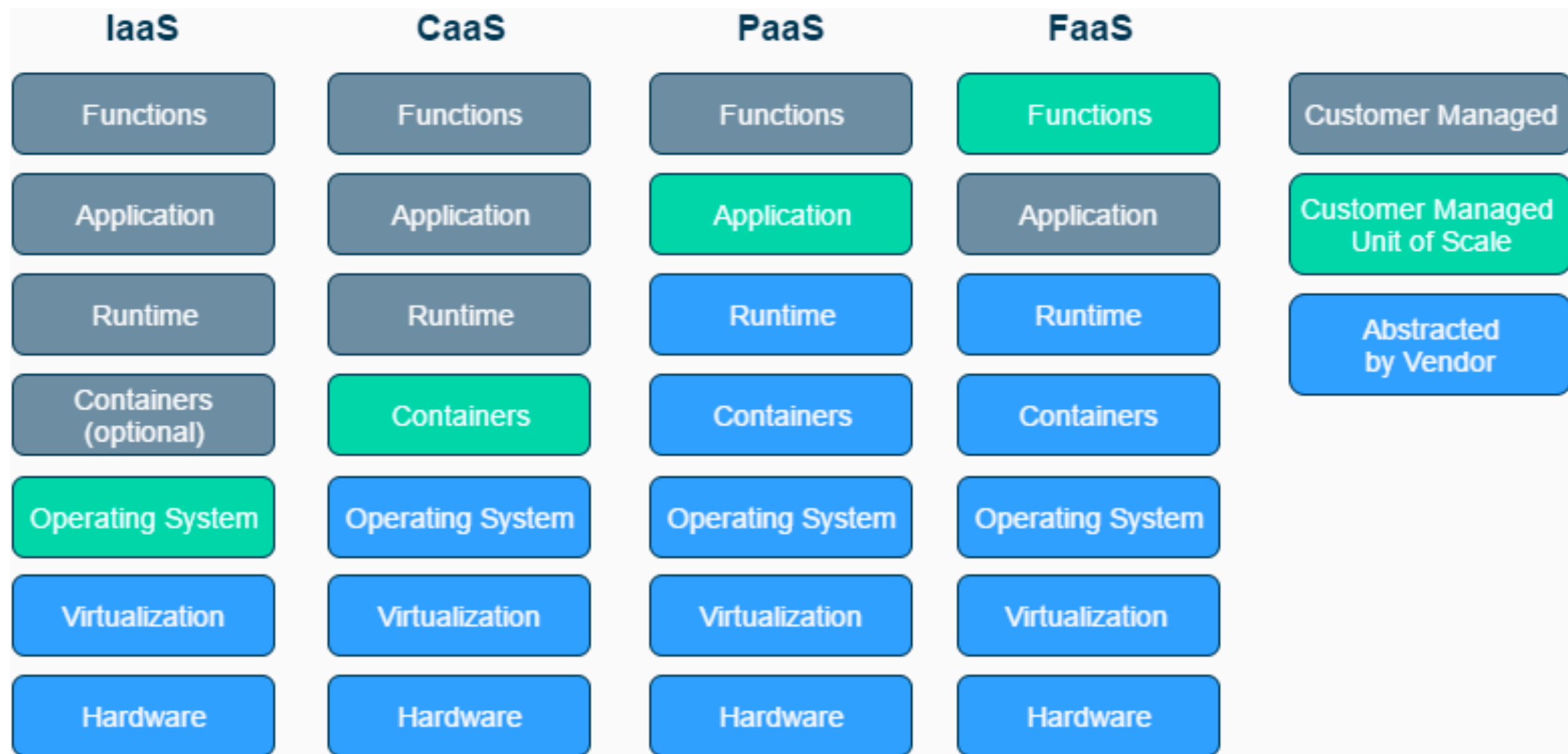


- PaaS: Platform as a Service
- SaaS: Software as a Service





- FaaS: Function as a Service

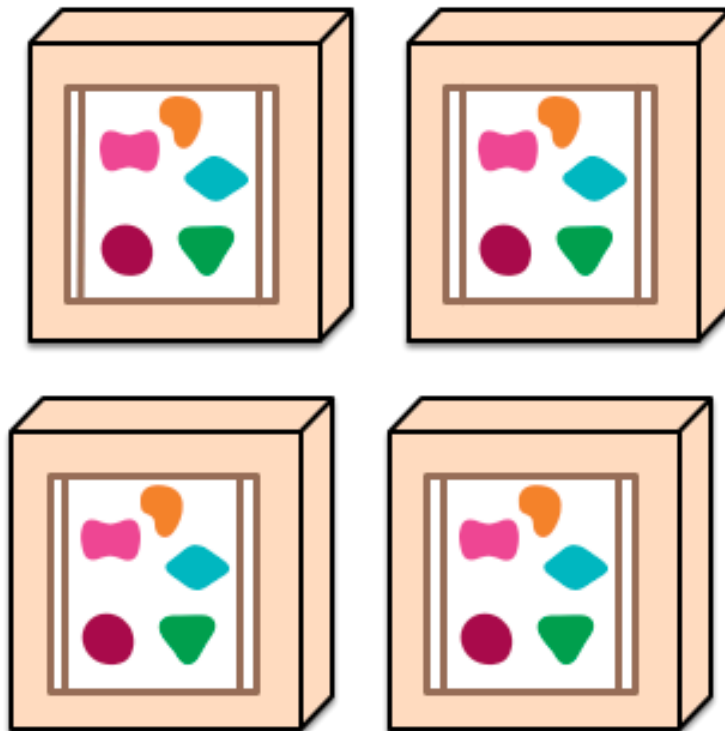


- Microservices

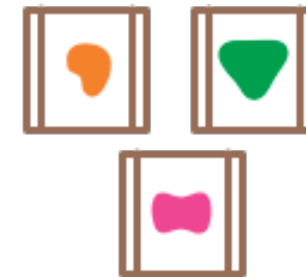
*A monolithic application puts all its functionality into a single process...*



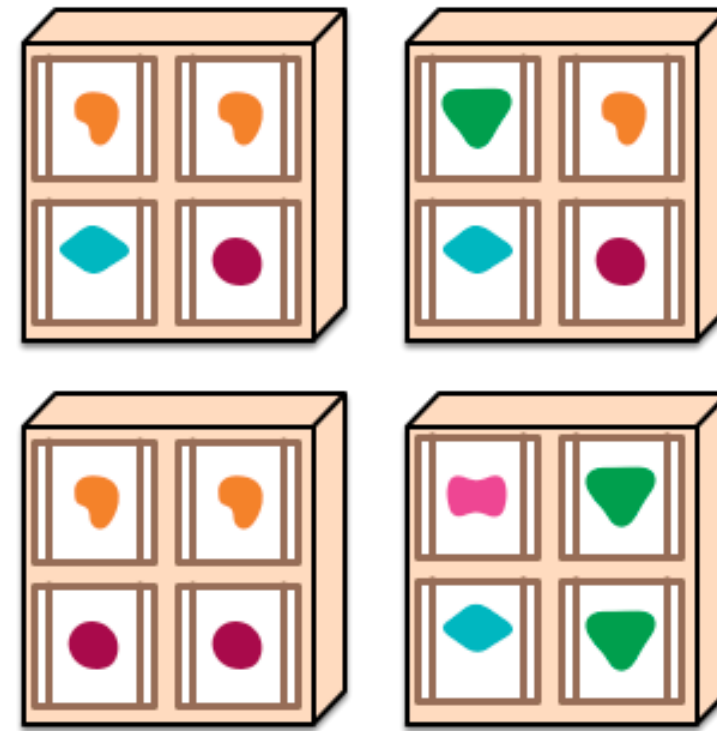
*... and scales by replicating the monolith on multiple servers*



*A microservices architecture puts each element of functionality into a separate service...*



*... and scales by distributing these services across servers, replicating as needed.*

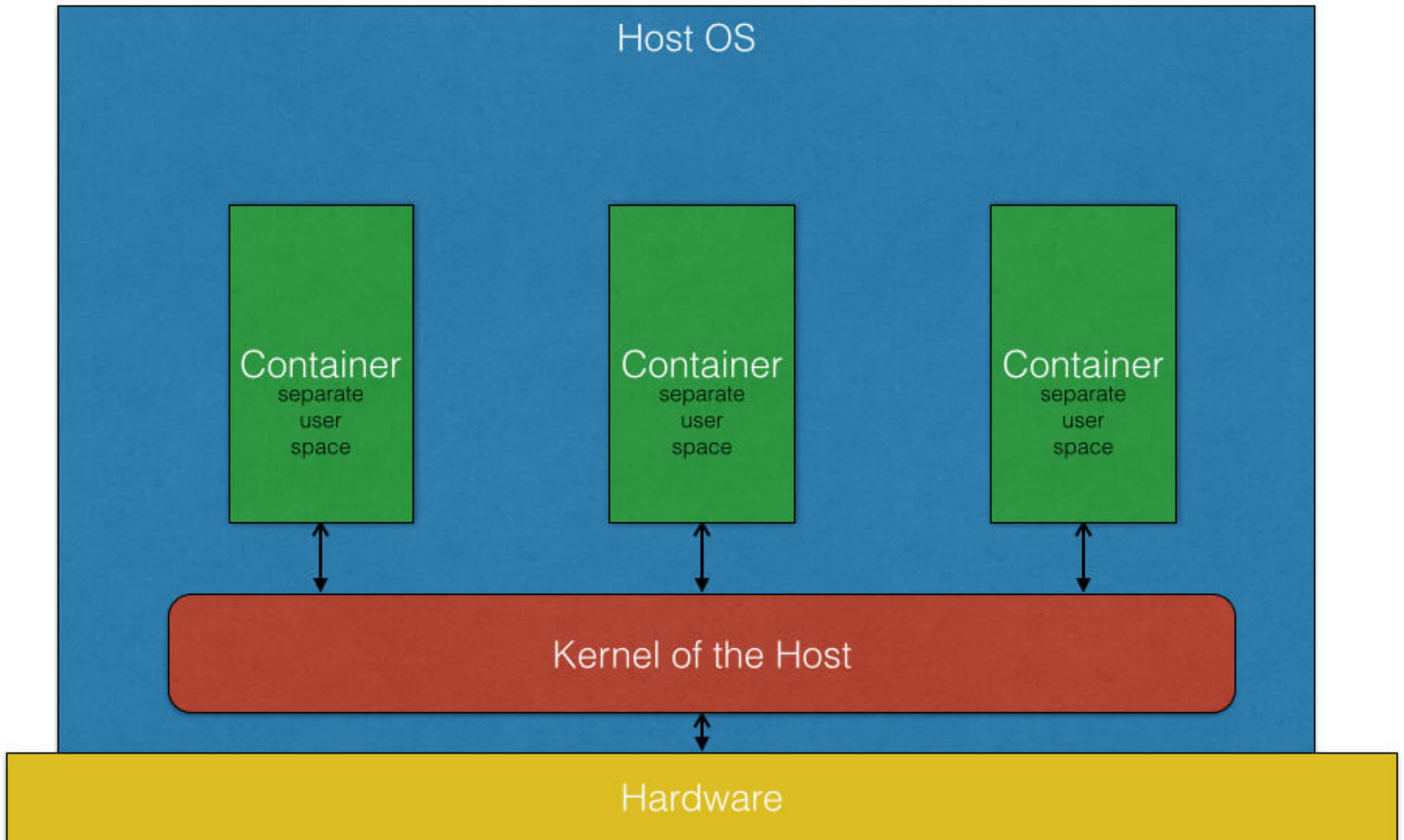


# Enter Containers



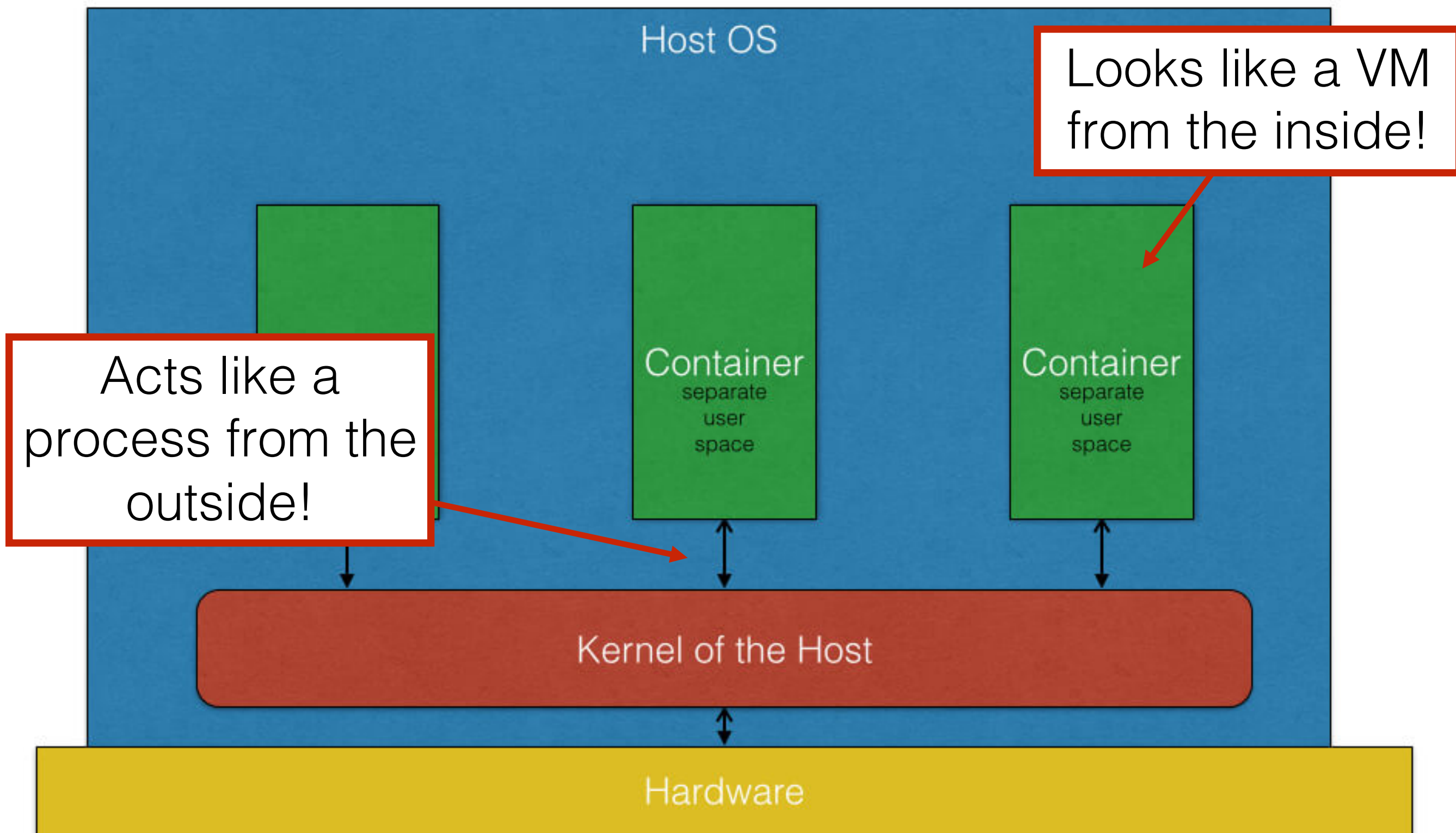
- Rather than virtualize both user space and kernel space... why not just ‘virtualize’ user space?
- Meets the needs of most customers, who don’t require significant customization of the OS.
- Sometimes called ‘OS virtualization,’ which is highly misleading given our existing taxonomy of virtualization techniques
- Running natively on host, containers enjoy bare metal performance without reliance on advanced virtualization support from hardware.

# Enter Containers





# Enter Containers





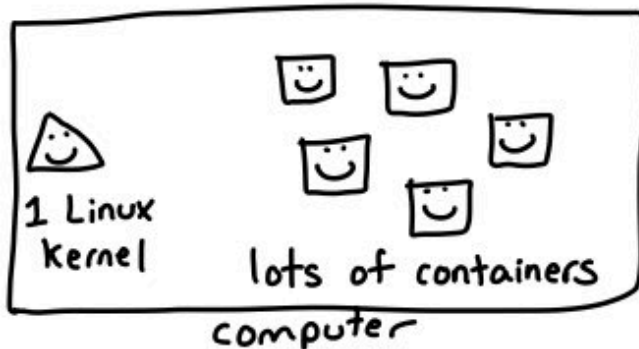
# Containers are processes



JULIA EVANS  
@b0rk

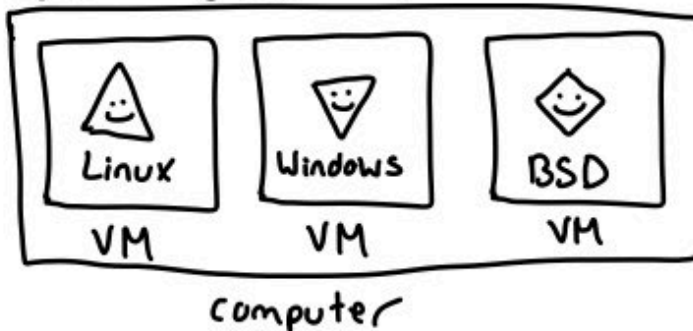
## containers vs VMs

a container is a group of processes



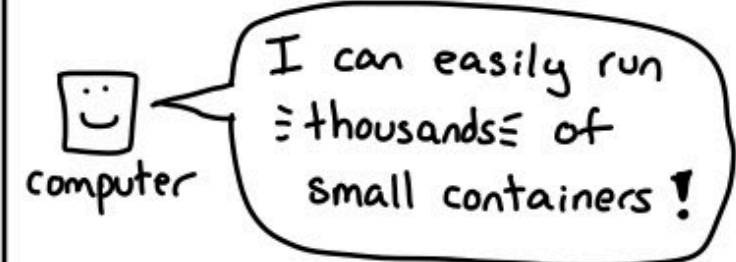
a virtual machine is a fake computer

each one has its own operating system!



containers use less RAM

This is because they share a single Linux kernel.

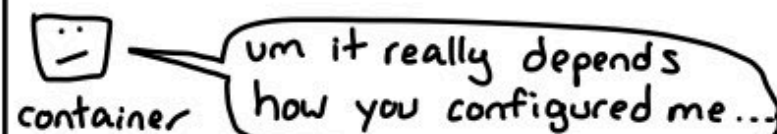
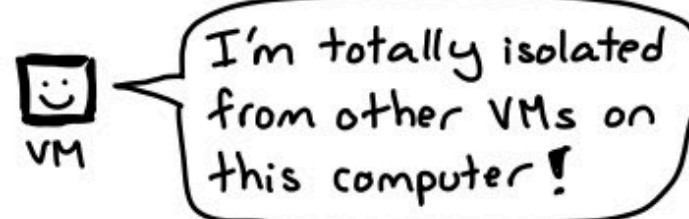


containers start faster

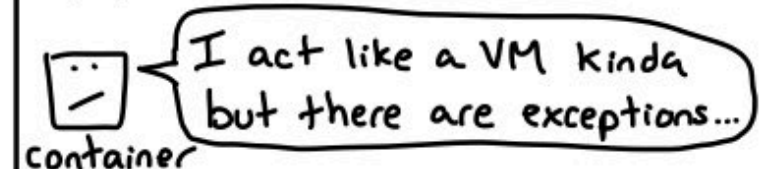
because they're processes and process start fast ♥



containers are more complicated to secure



it's harder to figure out what you can do in a container



# “Container is an old idea”

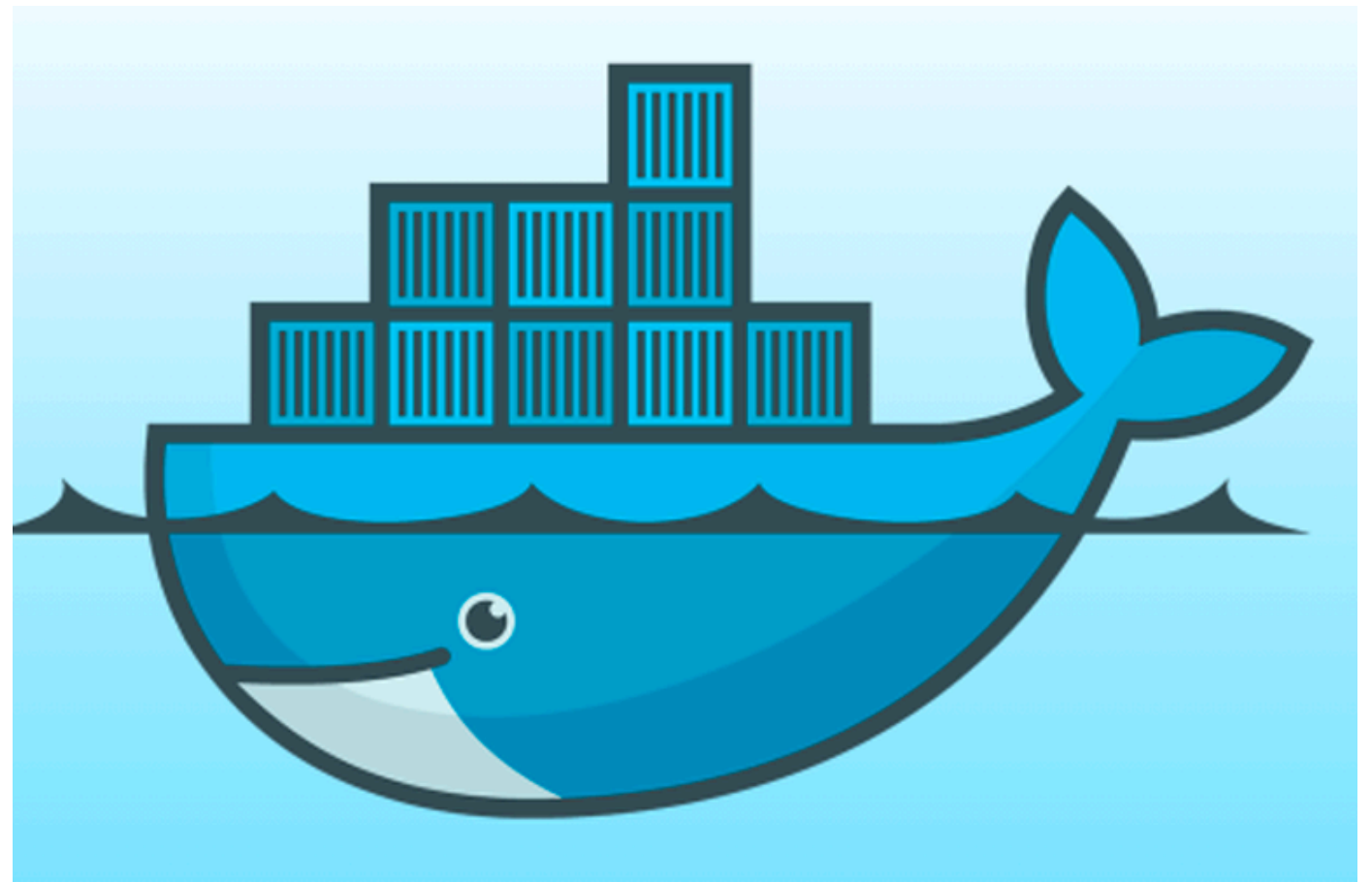


- You didn't heard of it cause it was not called “containers.”
- Linux containers
- BSD Jails
- Solaris Zones

# Docker's Big Idea



- Build, Ship, and Run App, Anywhere
- Debug your app, not your environment --  
Securely build and share any application,  
anywhere





# Docker's Big Idea



JULIA EVANS  
@b0rk

the big idea: include *EVERY* dependency

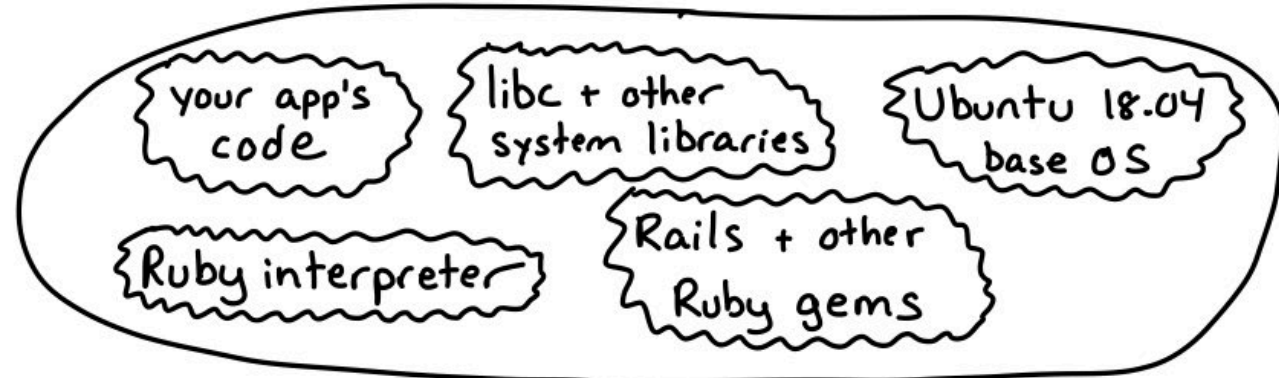
containers package  
*EVERY* dependency  
together



to make sure this  
program will run on  
your laptop, I'm going  
to send you every single  
file on my computer

exaggeration but  
it's the basic idea

a container image is a tarball of a filesystem  
Here's what's in a typical Rails app's container:



## how images are built

0. start with a base OS
  1. install program + dependencies
  2. configure it how you want
  3. make a tarball of the **WHOLE FILESYSTEM**
- (this is what 'docker build' does)

## running an image

1. download the tar ball
  2. unpack it into a directory
  3. Run a program and pretend that directory is its whole filesystem
- (this is what 'docker run' does)

images let you "install"  
programs really easily



Wow, I can get a  
Postgres test database  
running in 45 seconds!



- Linux Containers (LXC):
  - chroot
  - namespace
    - PID, Network, User, IPC, uts, mount
  - cgroups for HW isolation
  - Security profiles and policies
    - Apparmor, SELinux, Seccomp

# containers = chroot on steroids



- `chroot` changes the apparent root directory for a given process and all of its children
- An old idea! POSIX call dating back to 1979
- Not intended to defend against privileged attackers... they still have root access and can do all sorts of things to break out (like `chroot`'ing again)
- Hiding the true root FS isolates a lot; in \*nix, file abstraction used extensively
- Does not completely hide processes, network, etc., though!



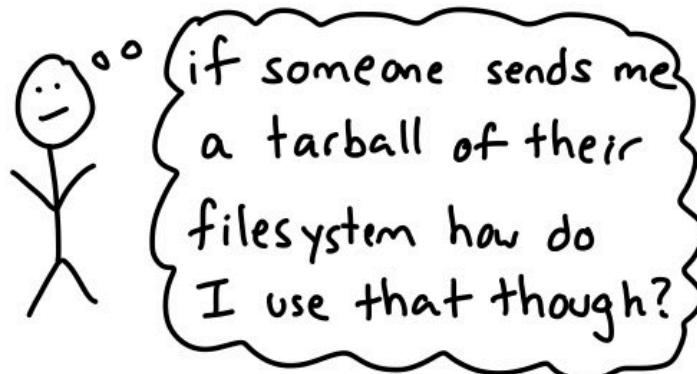


JULIA EVANS  
@b0rk

## chroot

a container image is a  
tarball of a filesystem

(or several tarballs: 1 per layer)



chroot: trick a program into thinking  
it has a different root directory

```
$ ls /path/to/container_filesystem
```

```
bin/ etc/ usr/ var/
```

```
$ sudo chroot /path/to/container_filesystem /bin/bash  
(inside chroot now)
```

```
$ ls /
```

```
bin/ etc/ usr/ var/
```

that's our new fake root directory!  
we tricked ls!

very basic way to  
"run" a Redis container

tarball of filesystem  
with Redis installed

```
$ mkdir redis; cd /redis
```

```
$ tar -xzf redis.tar
```

```
$ chroot $PWD /usr/bin/redis
```

```
# done! redis is running!
```

problems with just  
using chroot

- no CPU/memory limits
- other running processes are still visible
- can't use the same network port as another process
- LOTS of security issues

Docker uses pivot-root  
+ extra isolation features  
to run containers

pivot-root is like chroot  
but harder to escape from

# Namespaces



- The key feature enabling containerization!
- Partition practically all OS functionalities so that different process domains see different things
- Mount (mnt): Controls mount points
- Process ID (pid): Exposes a new set of process IDs distinct from other namespaces (i.e., the hosts)
- Network (net): Dedicated network stack per container; each interface present in exactly one namespace at a time.
- ....



# Namespaces



- The key feature enabling containerization!
- Partition practically all OS functionalities so that different process domains see different things
- Interprocess Comm. (IPC): Isolate processes from various methods of POSIX IPC.
  - e.g., no shared memory between containers!
- UTS: Allows the host to present different host/domain names to different containers.
- There's also a User ID (user) and cgroup namespace

# User Namespace



- Like others, can provide a unique UID space to the container.
- More nuanced though — we can map UID 0 inside the container to UID 1000 outside; allows processes inside of container to think they're root.
- Enables containers to perform administration actions, e.g., adding more users, while remaining confined to their namespace.

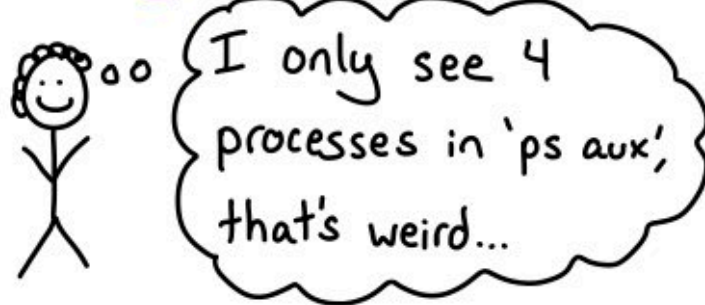
# Namespace



JULIA EVANS  
@b0rk

## namespaces

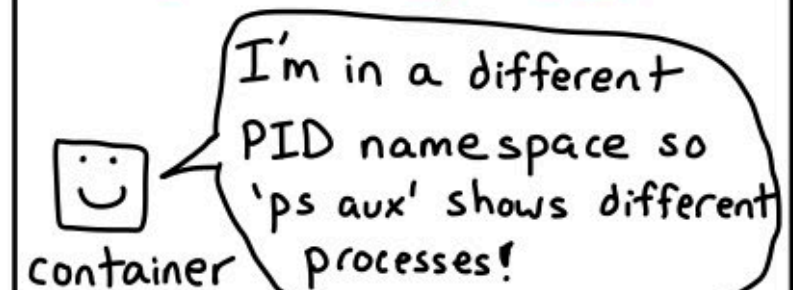
inside a container,  
things look different



Commands that  
will look different

- ps aux (less processes!)
- mount & df
- netstat -tulpn (different open ports!)
- hostname
- ... and LOTS more

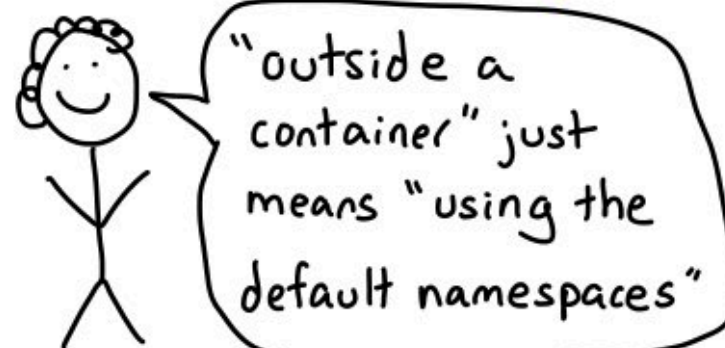
Why those commands  
look different:  
: namespaces :



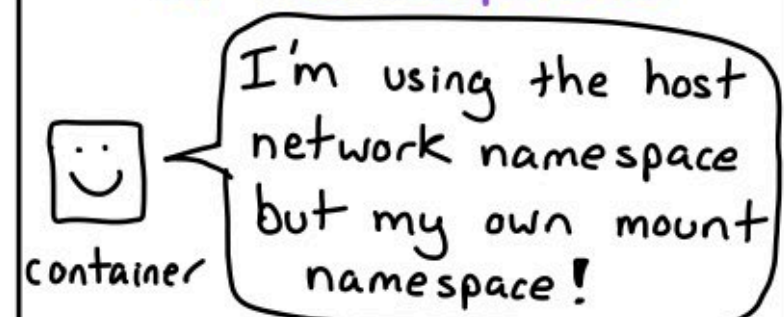
every process has 7  
kinds of namespaces



there's a default  
("host") namespace



processes can have  
any combination  
of namespaces



♥ this? more at [wizardzines.com](http://wizardzines.com)



- Limit, track, and isolate utilization of hardware resources including CPU, memory, and disk.
- Important for ensuring QoS between customers! Protects against bad neighbors
- Features:
  - Resource limitation
  - Prioritization
  - Accounting (for billing customers!)
  - Control, e.g., freezing groups
- The cgroup namespace prevents containers from viewing or modifying their own group assignment



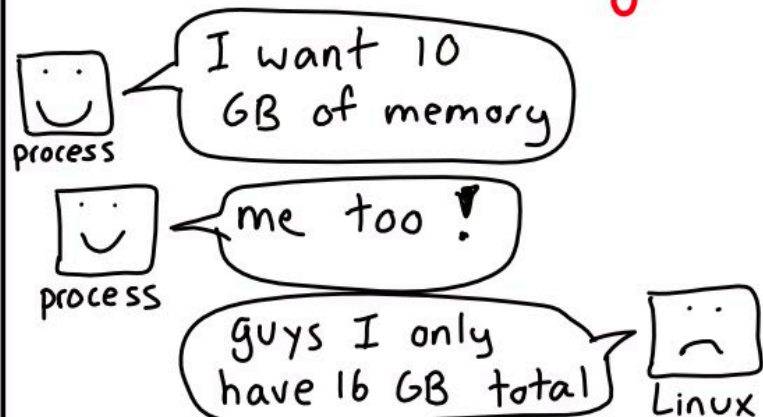
# cgroups



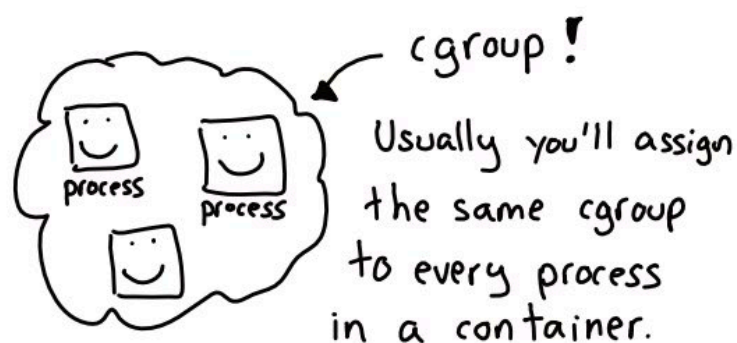
JULIA EVANS  
@b0rk

## cgroups

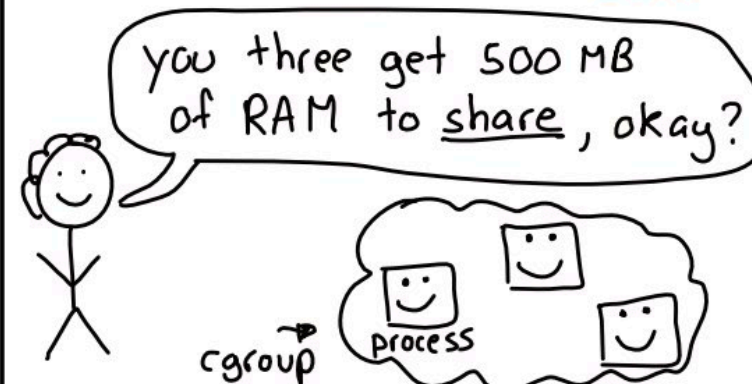
processes can use  
a lot of memory



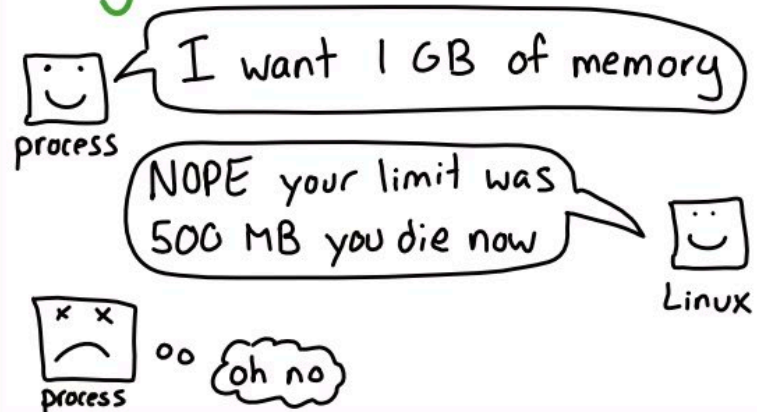
a cgroup is a  
group of processes



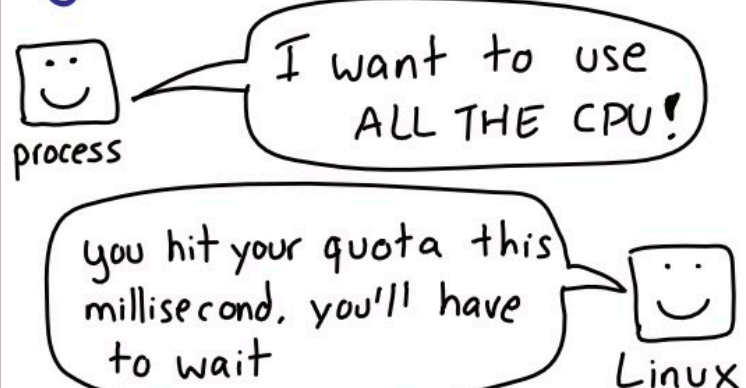
cgroups have  
memory / CPU limits



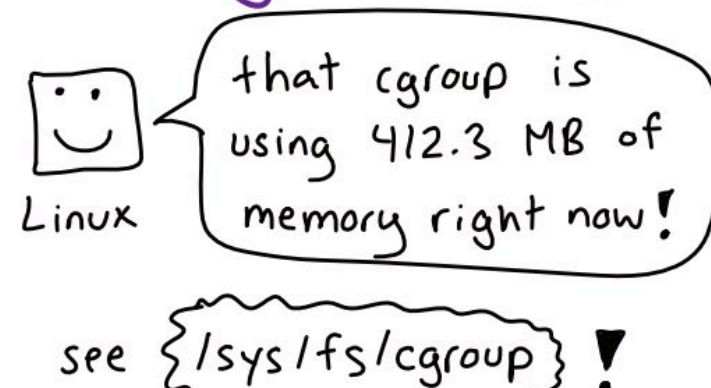
use too much memory:  
get OOM killed



use too much CPU:  
get slowed down



cgroups track  
memory & CPU usage



# Container Security?



*“Containers do not contain.” - Dan Walsh (SELinux contributor)*

- In a nutshell, it's real hard to prove that every feature of the operating system is namespaced.
  - /sys? /proc? /dev? LKMs? kernel keyrings?
- Root access to any of these enables pwning the host
- Solution? Just don't forget about MAC; at this point SELinux pretty good support for namespace labeling.
- SELinux and Namespaces actually synergize nicely; much easier to express a correct isolation policy over a coarse-grained namespace than, say, individual processes



# Capabilities



JULIA EVANS  
@b0rk

## capabilities

the root user can  
do **★anything★**

edit any  
file

change  
network  
config

spy on any  
program's memory

sometimes containers  
need privileged access

I need to update the  
route for 10.1.93.4



container  
process

that container  
shouldn't have  
root access though!



♥ capabilities ♥  
let you grant  
specific permissions

here's CAP\_NET\_ADMIN  
so you can manage  
the network



yay!



container  
process

**CAP\_SYS\_ADMIN**

basically root access. Try to  
use a more specific capability!

**CAP\_NET\_ADMIN**

for changing network  
settings

**CAP\_SYS\_PTRACE**

strace needs this

**CAP\_NET\_RAW**

ping needs this to send  
raw ICMP packets

`$ capsh --print`

run this in a container  
to print its capabilities

`$ getcap /usr/bin/ping`

shows which capabilities  
ping is allowed to use

# Seccomp-bpf



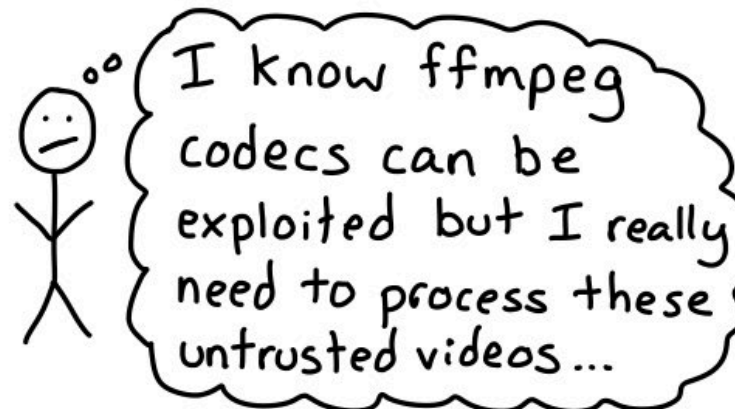
JULIA EVANS  
@b0rk

## seccomp-bpf

all programs use  
system calls



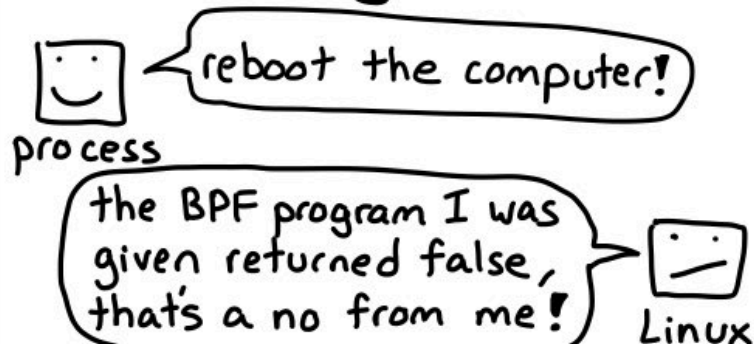
some programs have  
security vulnerabilities



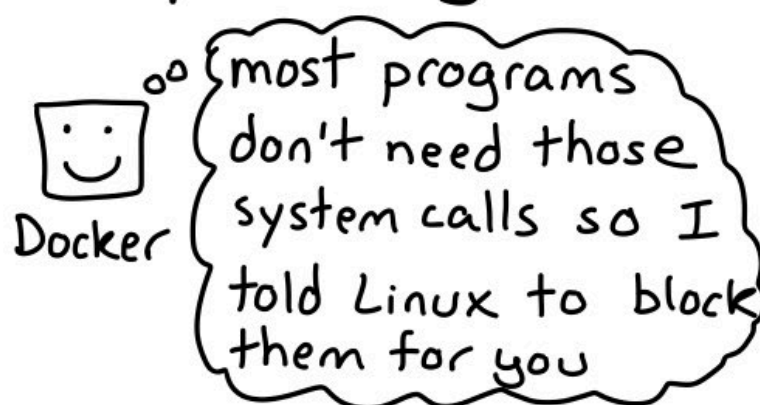
rarely used syscalls  
can help an attacker



**seccomp-BPF**: make  
Linux run a tiny program  
before every system call



Docker **blocks** dozens  
of syscalls by default



2 ways to block  
scary system calls

1. Limit a container's **capabilities**
2. Use a **seccomp-BPF** whitelist

Usually people do both!



# Linux Security Modules

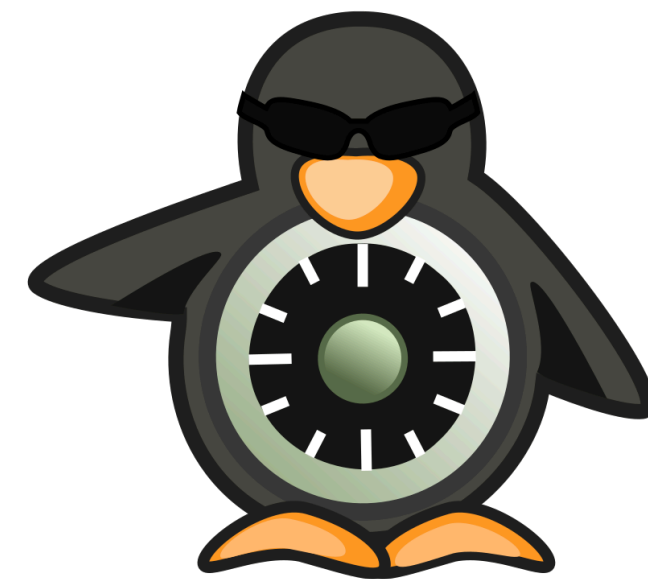
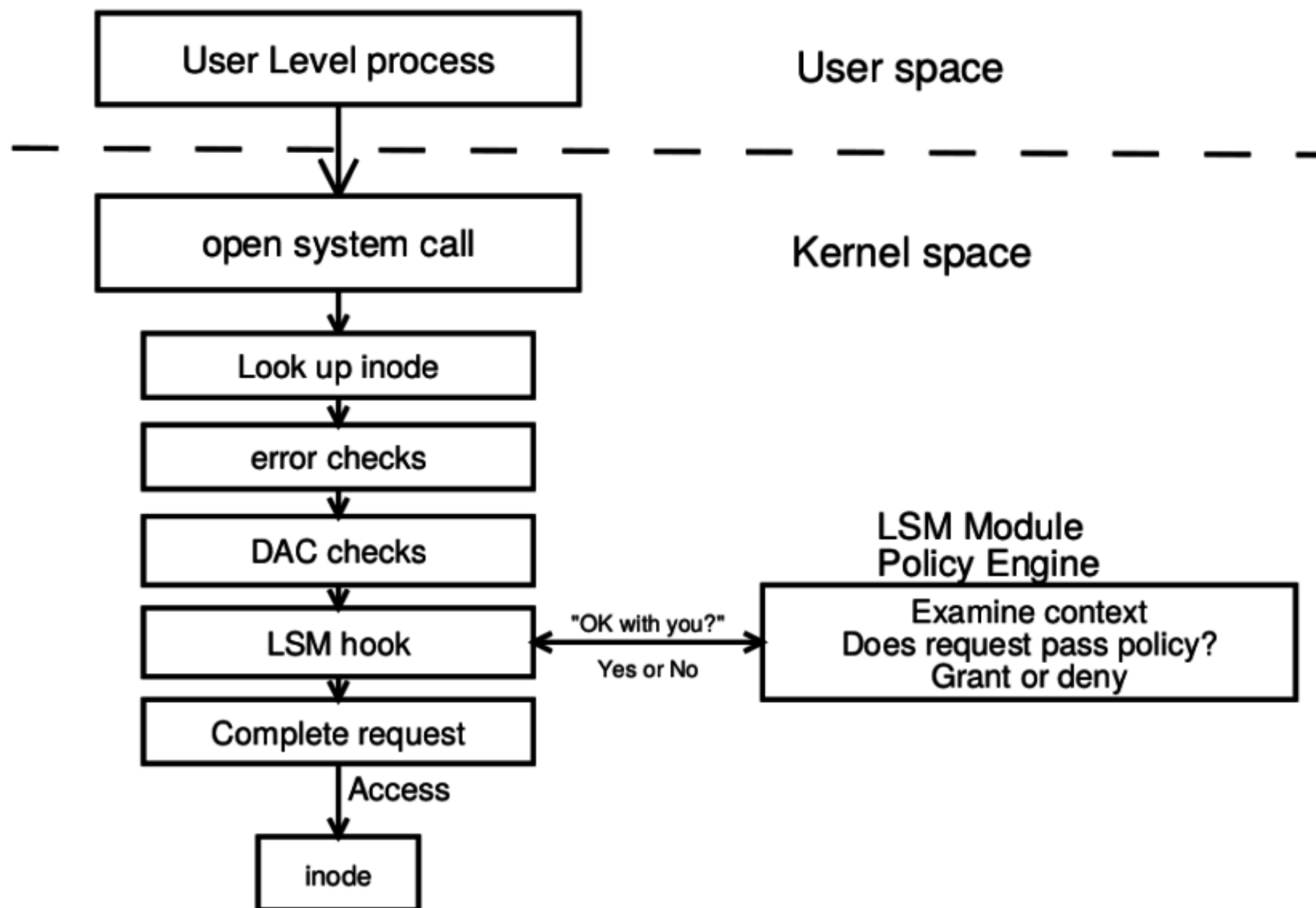


Figure 1: LSM Hook Architecture

# DIY container



JULIA EVANS  
@b0rk

## containers aren't magic

These 15 lines of bash will start a container running the fish shell. Try it!  
(download this script at [bit.ly/containers-arent-magic](http://bit.ly/containers-arent-magic))

```
wget bit.ly/fish-container -O fish.tar           # 1. download the image
mkdir container-root; cd container-root          #
tar -xf ../fish.tar                             # 2. unpack image into a directory
cgroup_id="cgroup_$(shuf -i 1000-2000 -n 1)"    # 3. generate random cgroup name
cgcreate -g "cpu,cpuacct,memory:$cgroup_id"     # 4. make a cgroup &
cgset -r cpu.shares=512 "$cgroup_id"             #    set CPU/memory limits
cgset -r memory.limit_in_bytes=1000000000 \      #
    "$cgroup_id"                                #
cgexec -g "cpu,cpuacct,memory:$cgroup_id" \     # 5. use the cgroup
    unshare -fmuipn --mount-proc \              # 6. make + use some namespaces
    chroot "$PWD" \                             # 7. change root directory
    /bin/sh -c "
        /bin/mount -t proc proc /proc &&
        hostname container-fun-times &&
        /usr/bin/fish"
```

# Summary



JULIA EVANS  
@b0rk

## container kernel features

containers are implemented using these Linux kernel features

You can use any of these on their own. When we use them all we call it a "container".

### ♥ pivot\_root ♥

set a process's root directory to a directory with the contents of the container image

### ★ cgroups ★

limit memory / CPU usage for a group of processes



### ♥ namespaces ♥

allow processes to have their own:

- network
- PIDs
- users
- hostname
- mounts
- + more

### ♥ seccomp-bpf ♥

security: prevent dangerous system calls

### ★ capabilities ★

security: avoid giving root access

### ★ overlay filesystems ★

optimization to reduce disk space used by containers which are using the same image

# Takeaways



- Container support has existing in Linux for many years
- Foundations of containerization has been around for decades!
- Automating LXC for portability (i.e., Docker) has revolutionized cloud computing
- Lasting legacy of containers may be enabling the Function-as-a-Service revolution... cloud customers can now pay by the method invocation without any idle costs.