

# Docker Internals

Rohit Jnagal

jnagal@

Containers @ Google  
Containers at scale.  
Resource Isolation.

[lmcftfy](#)

[libcontainer](#)

[cAdvisor](#)

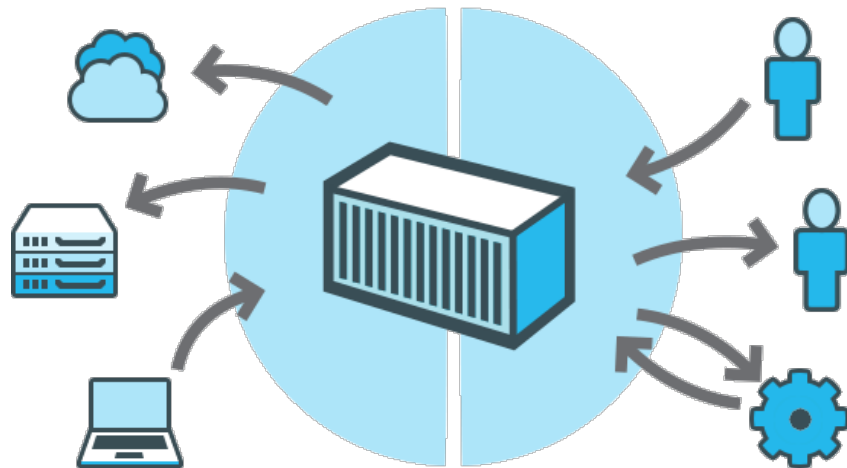
[Kubernetes](#)



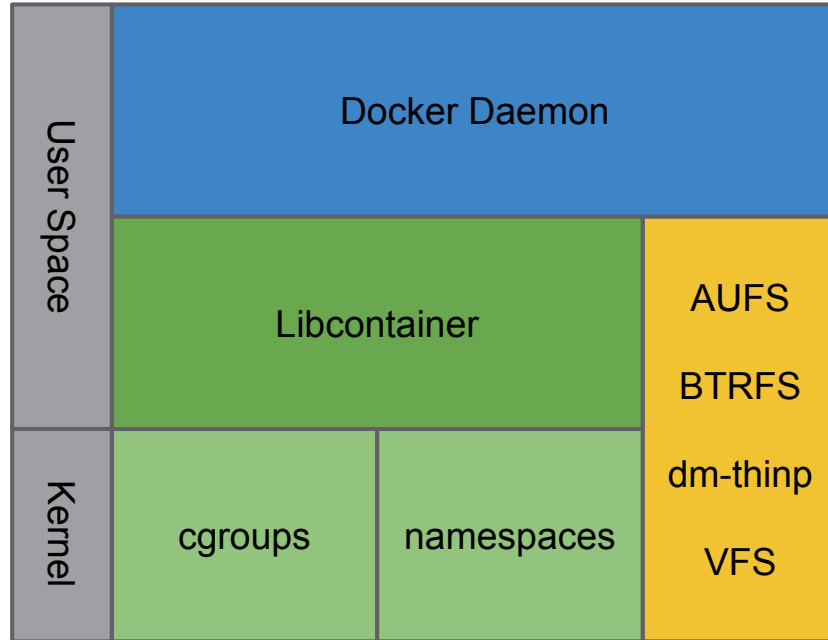
# Docker 101

Build Once, Configure Once.

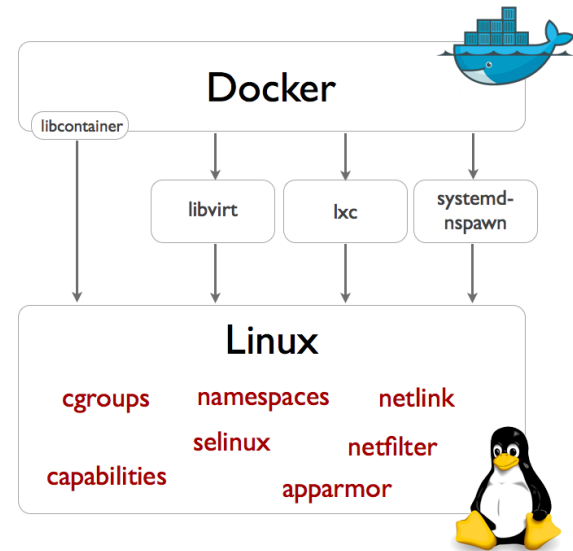
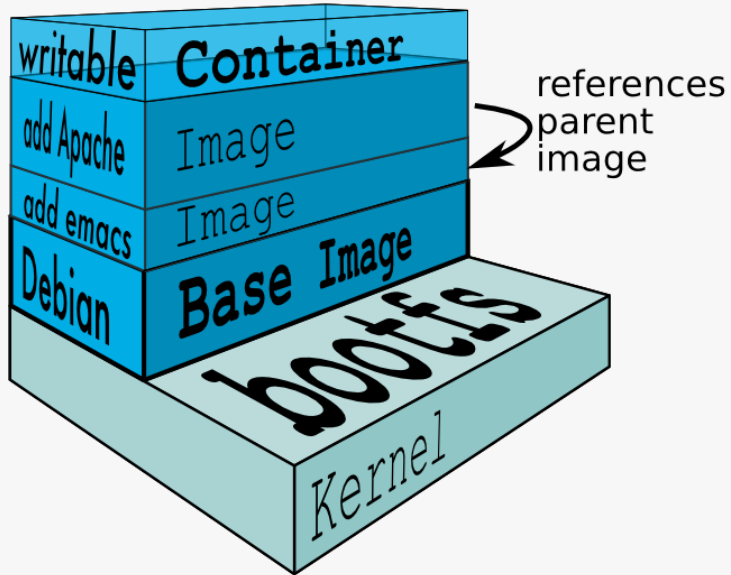
Deploy Everything\*  
Everywhere\*  
Reliably & Consistently  
Efficiently  
Cheaply



# Docker Components



# Docker Components



# Docker Grounds up: Filesystem

## File-system Isolation:

Building a rootfs dir and chroot into it.

With mount namespace, use pivot-root.

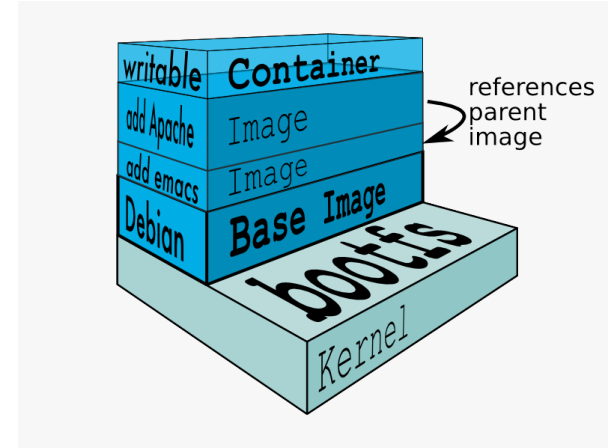
## Features:

Layering, CoW, Caching, Diffing

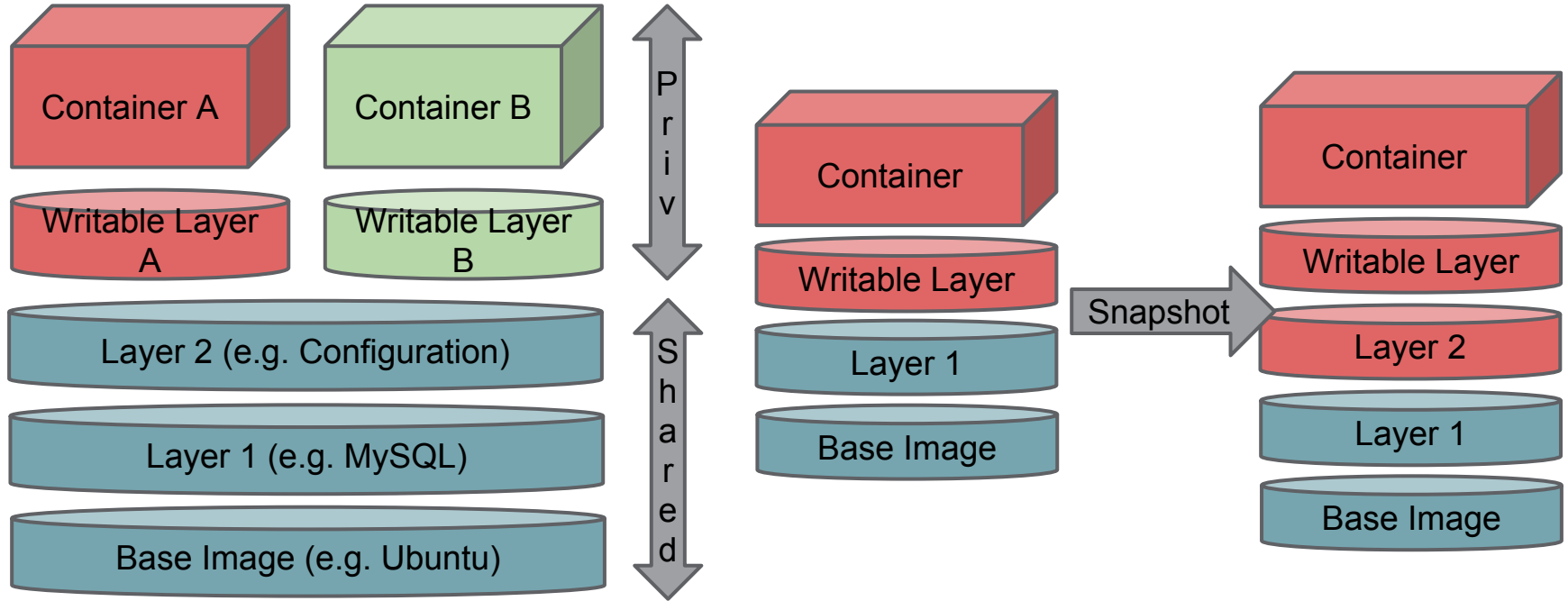
## Solutions:

UnionFS, Snapshotting FS, VFS

***AUFS in action***



# Filesystem



# Docker Grounds up: Filesystem

	Union Filesystems	Snapshotting Filesystems	Copy-on-write block devices
Provisioning	Superfast Supercheap	Fast Cheap	Fast Cheap
Changing small files	Superfast Supercheap	Fast Cheap	Fast Costly
Changing large files	Slow (first time) Inefficient (copy-up!)	Fast Cheap	Fast Cheap
Diffing	Superfast	Superfast	Slow
Memory usage	Efficient	Efficient	Inefficient (at high densities)
Drawbacks	Random quirks AUFS not mainline !AUFS more quirks	ZFS not mainline BTRFS not as nice	Higher disk usage Great performance (except diffing)
Bottom line	<b>Ideal for PAAS and high density things</b>	<b>This is the Future (probably)</b>	<b>Dodge Ram 3500</b>

From: Jérôme Petazzoni

# Docker Grounds up: Namespaces

- Process trees.
- Mounts.
- Network.
- User accounts.
- Hostnames.
- Inter-process communication.

```
pid_t pid = clone(..., flags, ...)
```

CLONE_NEWUTS	hostname, domainname
CLONE_NEWIPC	IPC objects
CLONE_NEWPID	Process IDs
CLONE_NEWNET	Network configuration
CLONE_NEWNS	File system mounts
CLONE_NEWUSER	User and Group IDs

```
setns(int fd, int nstype)
```

```
CLONE_NEWIPC  
CLONE_NEWNET  
CLONE_NEWUTS
```

```
Also: unshare(flags)
```

# Docker Grounds up: Resource Isolation

## Cgroups : Isolation and accounting

- cpu
- memory
- block i/o
- devices
- network
- numa
- freezer

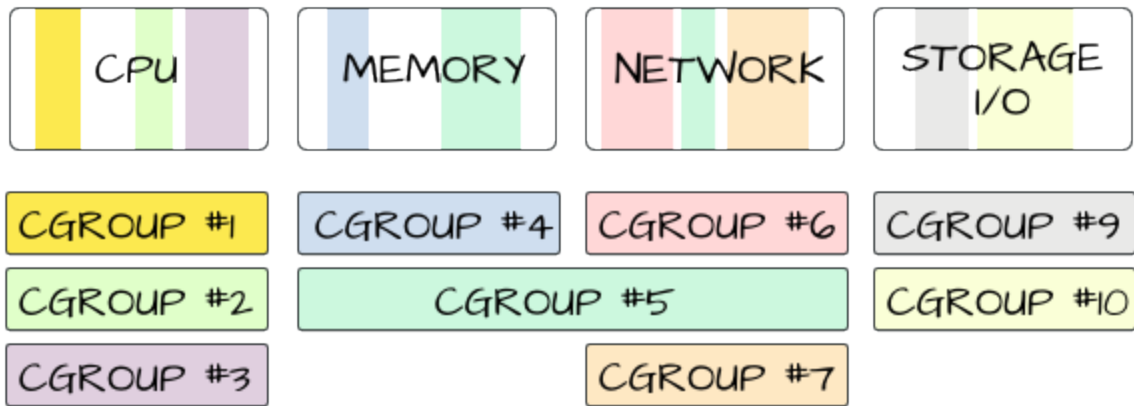


image credit: [mairin](#)

# Docker Grounds up: Add Security

## Security Layers

- Linux Capabilities.
- User namespaces: Unprivileged users.
- nosuid & ro mounts.
- Seccomp-bpf
- GRSEC and PAX
- Device cgroups
- Access Control: SELinux & AppArmor
- Future: Namespace aware sys/proc



image credit: [Leo Reynolds](#)

Questions

Thanks,

jnagal@google  
@jnagal