# Container Storage Best Practices in 2017

Myth-busting and taking state of the present

Keith Resar
Red Hat Solution Architect
January 24th, 2017

# Agenda

- Container Storage Myths
- Container Storage Primer
- Review 6 Storage Drivers
- Chooser a Storage Driver

redhat.
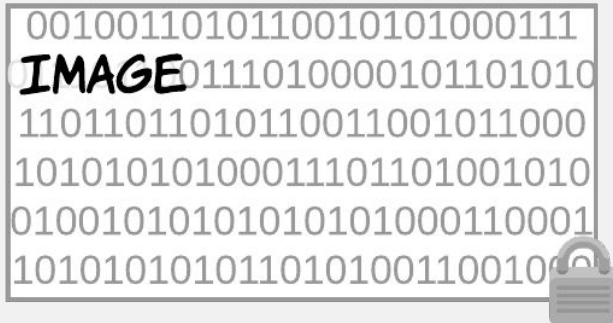
# Container Storage History / Myths

@KeithResar

# Container Storage Level Set

# Image : Container :: Class : Object

For humans, read this to say:
An image is to a container, as a class is to an object.
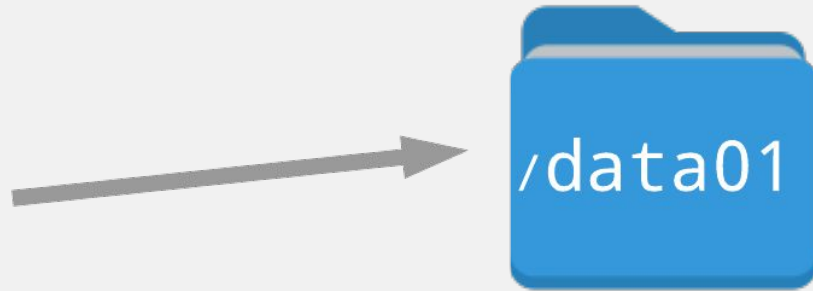
```
00100111010110010101000111
IMAGE     011101000010110101 0
110110110101100110001011000
101010101000111011010010 10
010010101010101010001100 01
1010101010110101001100100
```

```
> ls -l /image; echo $?
0
```

```
> pgrep image; echo $?
1
```

redhat.

# Container Storage

CONTAINER

```
00100110101100101010001 1
READ-WRITE LAYER 101
11011011010110011001011000
```

```
00100110101100101010001 11
IMAGE 011101000010110101 0
11011011010110011001011000
10101010100011101101001010
01001010101010101000110001
101010101011010100110010010
```

# Data Volume Storage

/data01

# Data Volume Storage

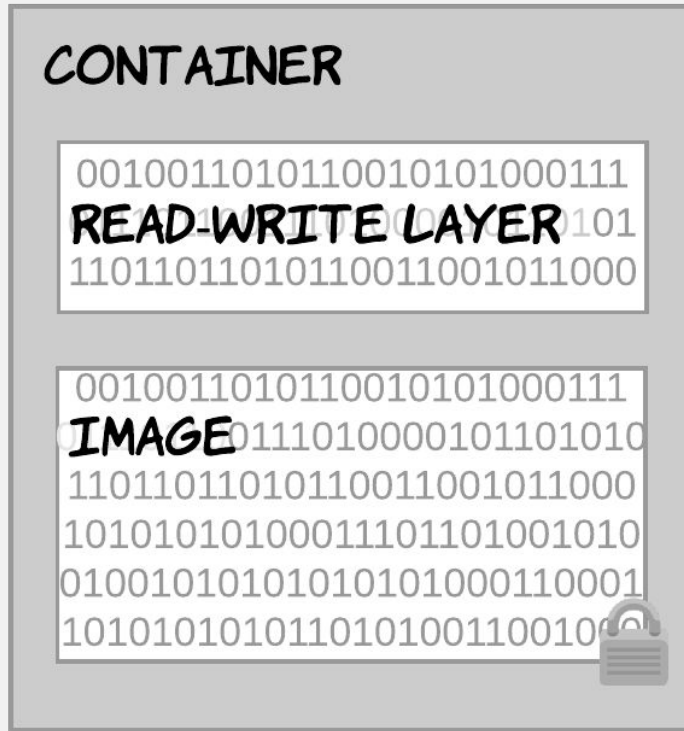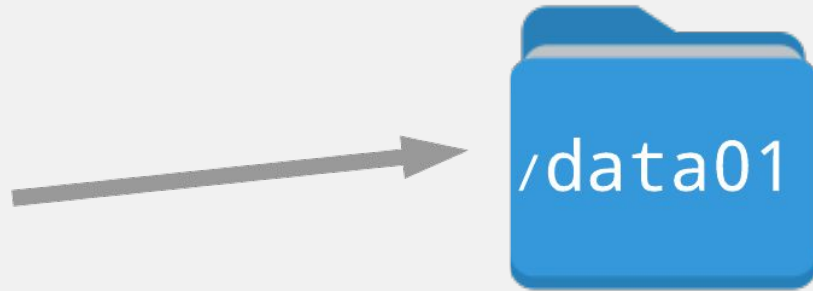**/data01**

Can be extended to support many endpoints and protocols using installable docker plugins.

- Local
- LVM

- GlusterFS
- Ceph
- NFS
- iSCSI

redhat.

# Container Storage

## CONTAINER

```
0010011010110010101000111
READ-WRITE LAYER
11011011010110011001011000
```

```
0010011010110010101000111
IMAGE
11011011010110011001011000
10101010100011101101001010
01001010101010101000110001
10101010101101010011001001
```

# Data Volume Storage

/data01

IMAGE

| | |
|---|---|
| E8E3AAF82AF5 | 55 MB |

httpd

| | |
|---|---|
| AAD5D4C7BBA9 | 192 MB |

rhel7:latest

@KeithResar

@KeithResar

RW LAYER · · · · RW LAYER · · · · RW LAYER

IMAGE

E8E3AAF82AF5  55 MB  httpd

AAD5D4C7BBA9  192 MB  rhel7:latest

redhat.

# Copy-on-write Strategy

# Container Storage Drivers

# Available Storage Drivers

| Technology | Driver | Introduction | File vs. Block |
|---|---|---|---|
| VFS | vfs | origin | * File |
| AUFS | aufs | origin | File |
| OverlayFS | overlay/overlay2 | Aug 2014 (1.11)<br>June 2016 (1.12) | File |
| Device Mapper | devicemapper | Sept 2013 (0.7) | Block |
| Btrfs | btrfs | Nov 2013 | File |

redhat.

# vfs Driver   (1 of 6)

Naive implementation lacking union filesystem and copy-on-write



@KeithResar

# vfs Driver   (1 of 6)

Naive implementation lacking union filesystem and copy-on-write

| The Good | The Bad | Summary |
|---|---|---|
| Reference compatibility model<br><br>Useful for docker-in-docker scenarios to avoid nesting storage drivers | No shared memory, union filesystem, or copy-on-write | Not for production use<br><br>Important support role for storage driver development |

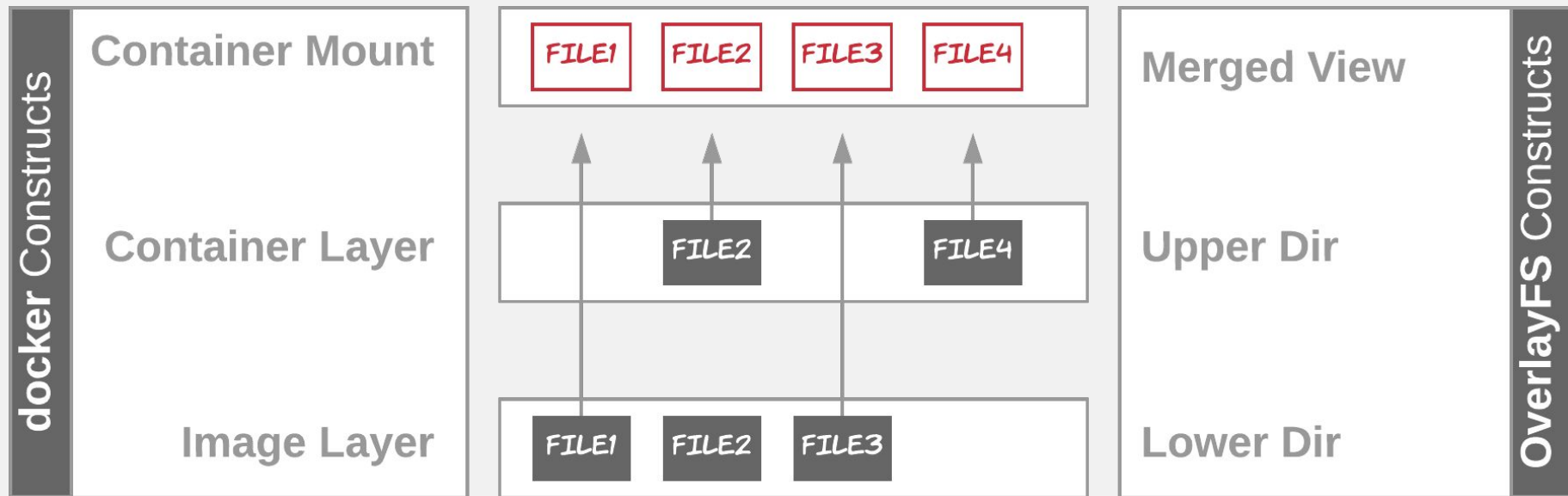redhat.

# AUFS Driver   (2 of 6)

The original docker storage driver



| docker Constructs | | aufs Constructs |
|---|---|---|
| Container Top Layer | FILE1  FILE2  .WH.F3  FILE4 | Union Mount |
| Image Layer 2 | FILE1   X   FILE4 | Branch 3 |
| Image Layer 1 | FILE2 | Branch 2 |
| Image Base Layer | FILE3  FILE4 | Branch 1 |

https://docs.docker.com/engine/userguide/storagedriver/aufs-driver/

@KeithResar

redhat.

# AUFS Driver   (2 of 6)

The original docker storage driver

| The Good | The Bad | Summary |
|---|---|---|
| Battle hardened driver<br><br>Performant and stable for wide range of use cases<br><br>Supports shared memory | Carried patch to mainline Linux kernel limits distro support<br><br>File level implementation impacts copy-on-write | Default for non-RH, will meet majority of needs<br><br>Expectation that it will be supplanted by an Overlay implementation |

redhat.

# Overlay Driver   (3 of 6)
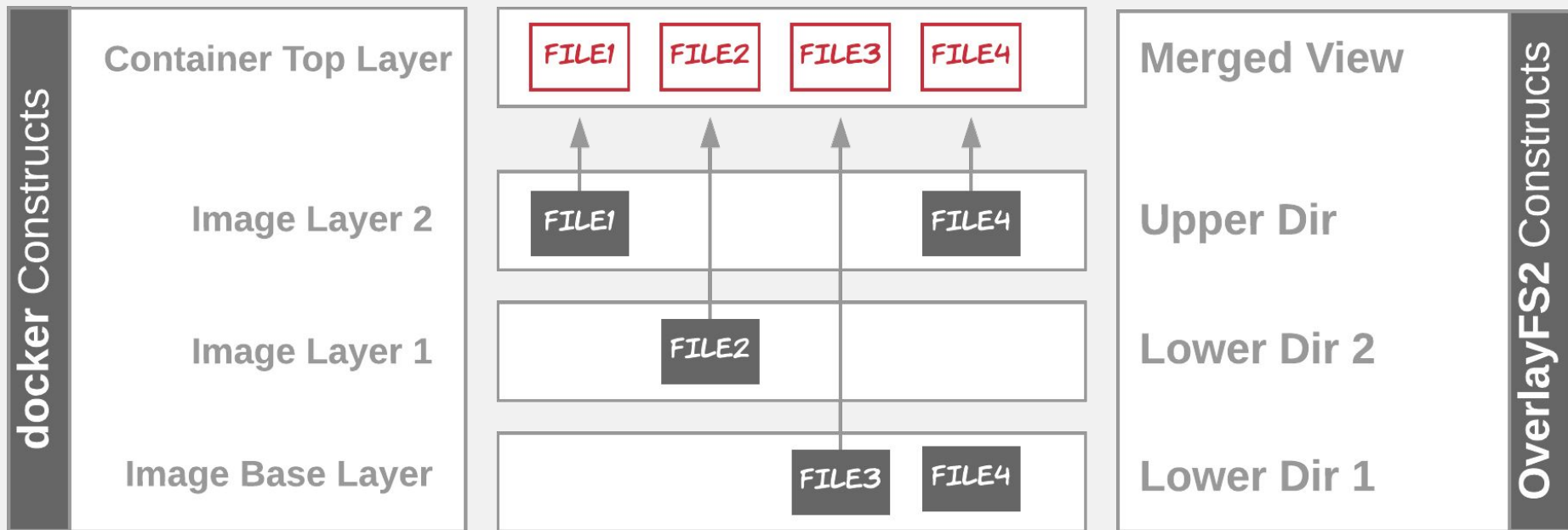
Legacy union filesystem driver, superseded by overlay2

https://docs.docker.com/engine/userguide/storagedriver/overlayfs-driver/

# Overlay Driver   (3 of 6)

Legacy union filesystem driver, superseded by overlay2

| The Good | The Bad | Summary |
|---|---|---|
| Complete union filesystem merged into the mainline kernel<br><br>Shared memory | Architecture drove explosive inode usage, often to the point of exhaustion<br><br>Slow commit performance | Used for backward compatibility in pre-4.0 kernels<br><br>Broad distro support beyond aufs |

# Overlay2 Driver   (4 of 6)

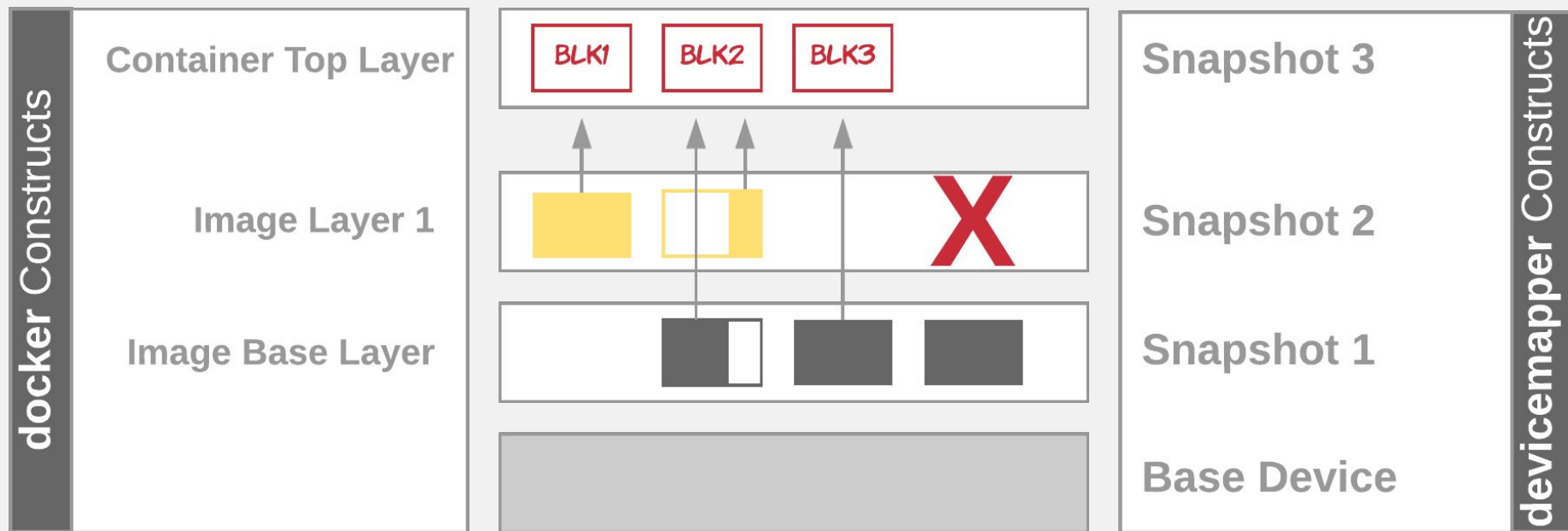Lessons learned from original overlay, and looking forward to continued maturity



@KeithResar

# Overlay2 Driver   (4 of 6)

Lessons learned from original overlay, and looking forward to continued maturity

| The Good | The Bad | Summary |
|---|---|---|
| Retains all benefits of overlay (shared memory, broad distro support)<br><br>Resolves inode exhaustion problems | Relatively young codebase (initial release with  Docker 1.12 in June 2016)<br><br>File-based so copy-on-write operations may be expensive | With maturity may be the best route forward for consistent defaults across many Linux distributions |

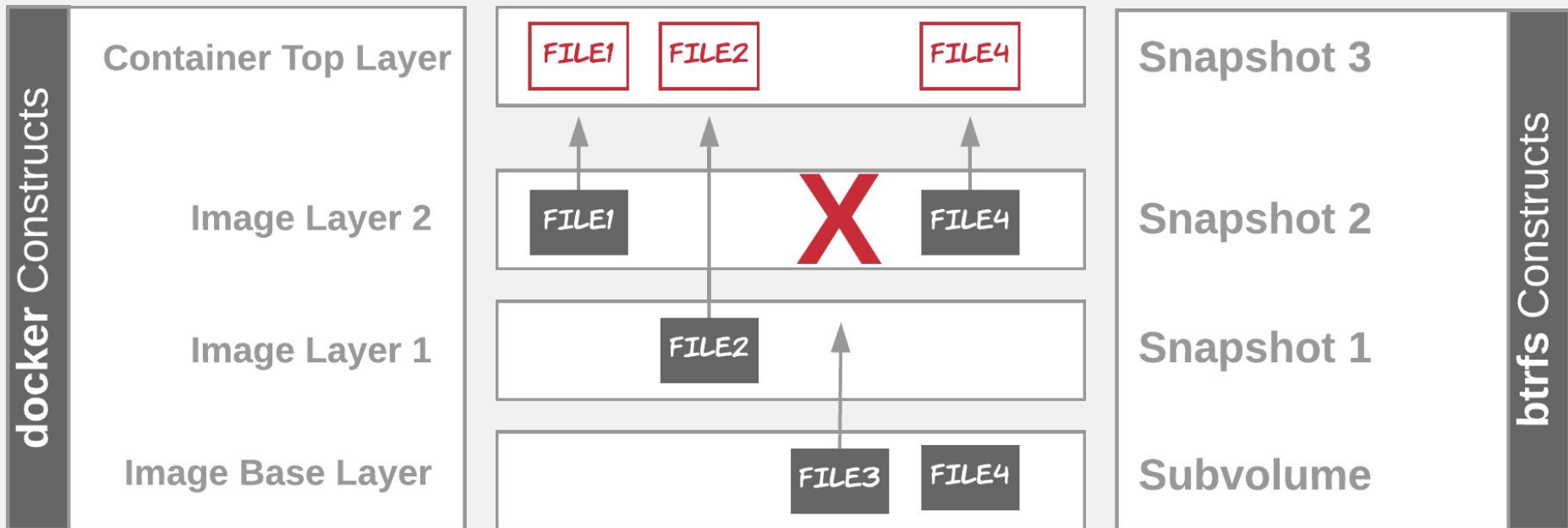# Devicemapper Driver   (5 of 6)

Lvm integrated block-based storage driver

**docker Constructs**

Container Top Layer

Image Layer 1

Image Base Layer

BLK1    BLK2    BLK3

X

Snapshot 3

Snapshot 2

Snapshot 1

Base Device

**devicemapper Constructs**

@KeithResar

redhat.

# Devicemapper Driver   (5 of 6)

Lvm integrated block-based storage driver, default on RHEL

| The Good | The Bad | Summary |
|---|---|---|
| Block-based solution offers efficient copy-on-write<br><br>Quota support<br><br>Available direct and loop modes | Manual setup is intimidating<br><br>No shared memory support | Red Hat go-to graphdriver with mature codebase |

redhat.

# Btrfs Driver   (6 of 6)

Another next generation filesystem, with a continued heavy development requirement

Another next generation filesystem, with a continued heavy development requirement

| The Good | The Bad | Summary |
|----------|---------|---------|
| Now offers SELinux support and quota | No page-cache sharing between containers<br><br>Small writes can lead to out-of-space conditions<br><br>Requires btrfs specific tools rather than Linux native | Btrfs hasn't been a mainstream choice for Linux distros, driving less attention and less testing |

redhat.

# Choosing a Storage Driver

# Benchmark Approach

Benchmarking is treacherous and confusing, and often done poorly - which means that you need to take any benchmark results with a large grain of salt.
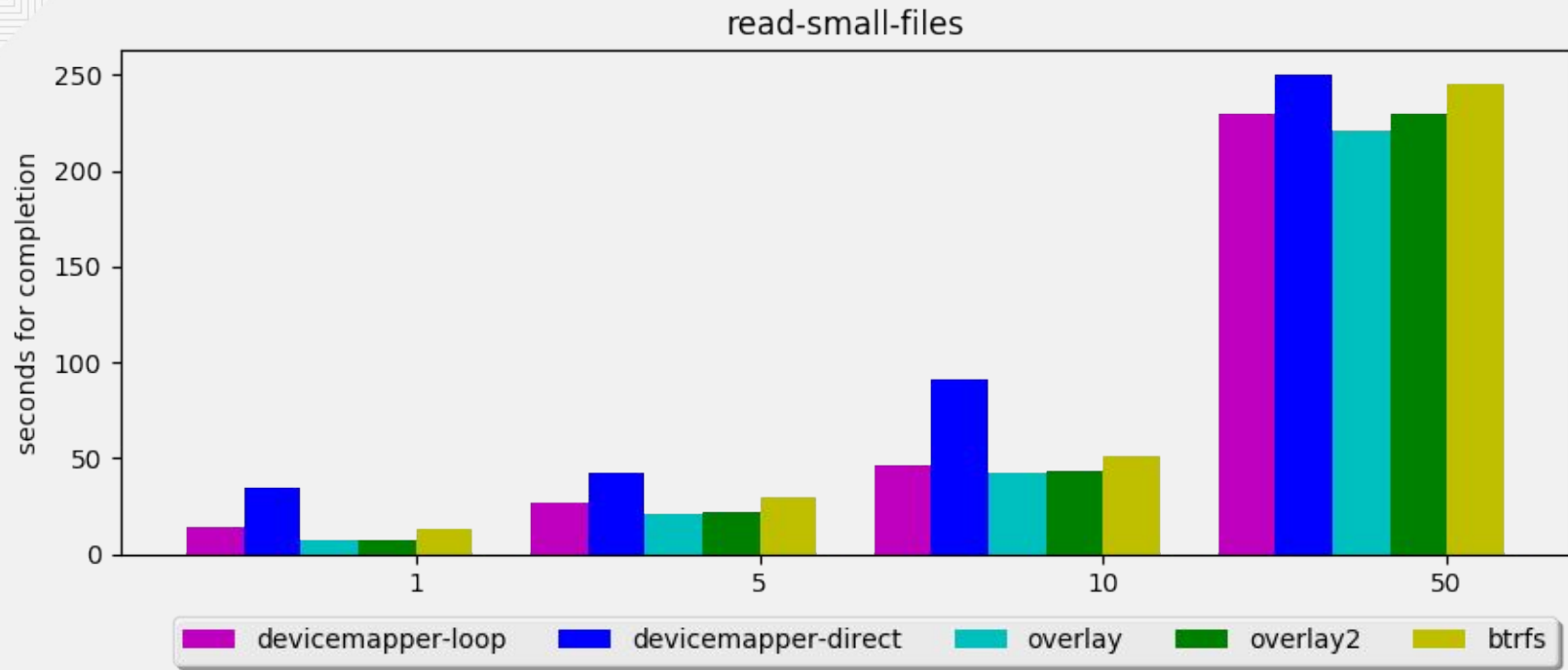
If you've spent less than a week studying a benchmark result, it's probably wrong.

(Running a benchmark is the easy part. Understanding a benchmark can take much longer.)
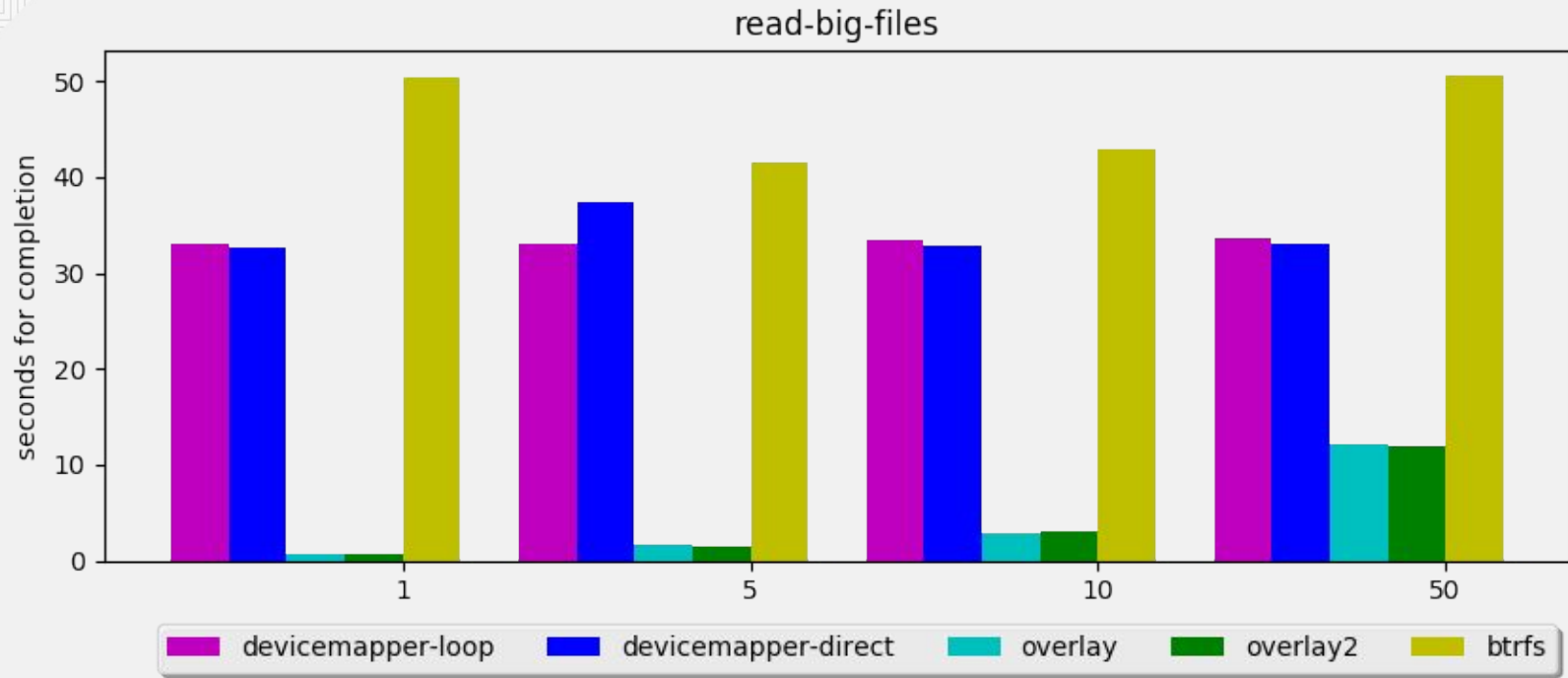
@KeithResar

# Benchmark 1: Reading Files
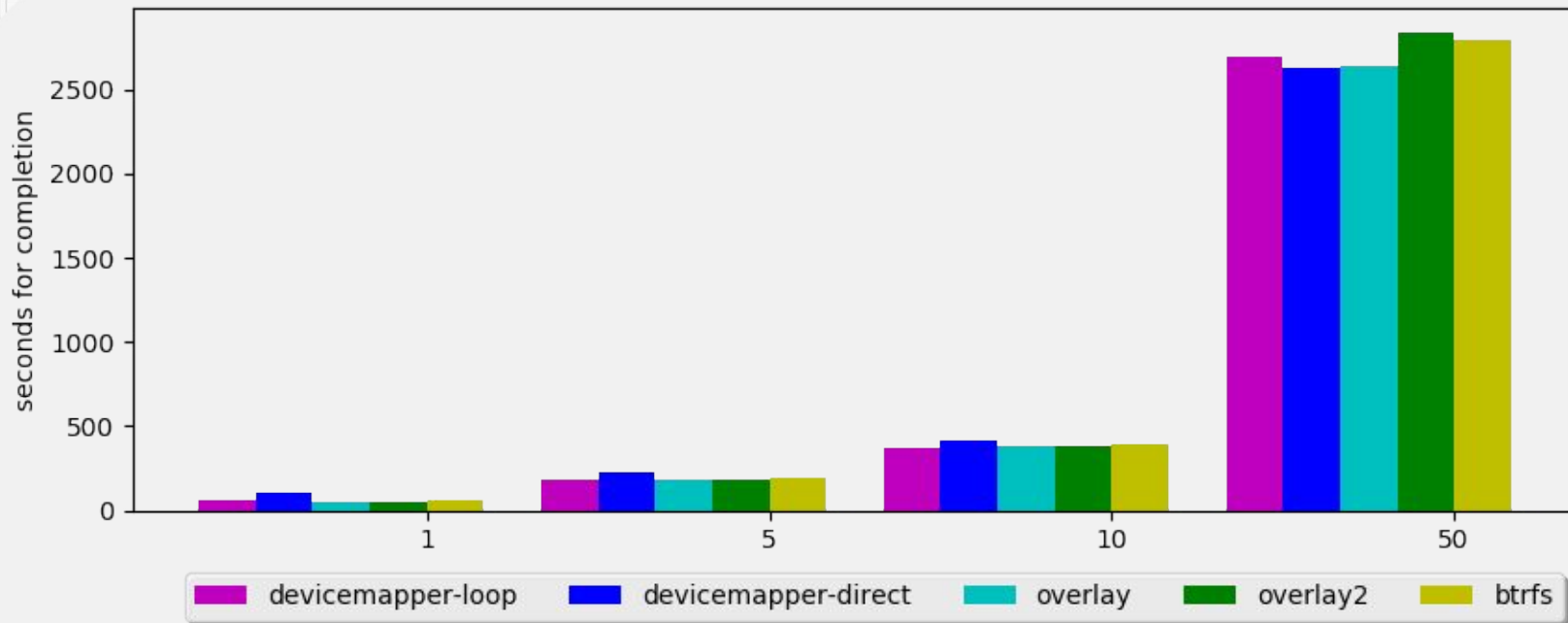
- Reading Small Files
- Reading Large Files
- Reading File Tree

redhat.

read-small-files

Naive benchmarking, for discussion purposes only.  Don't trust this!
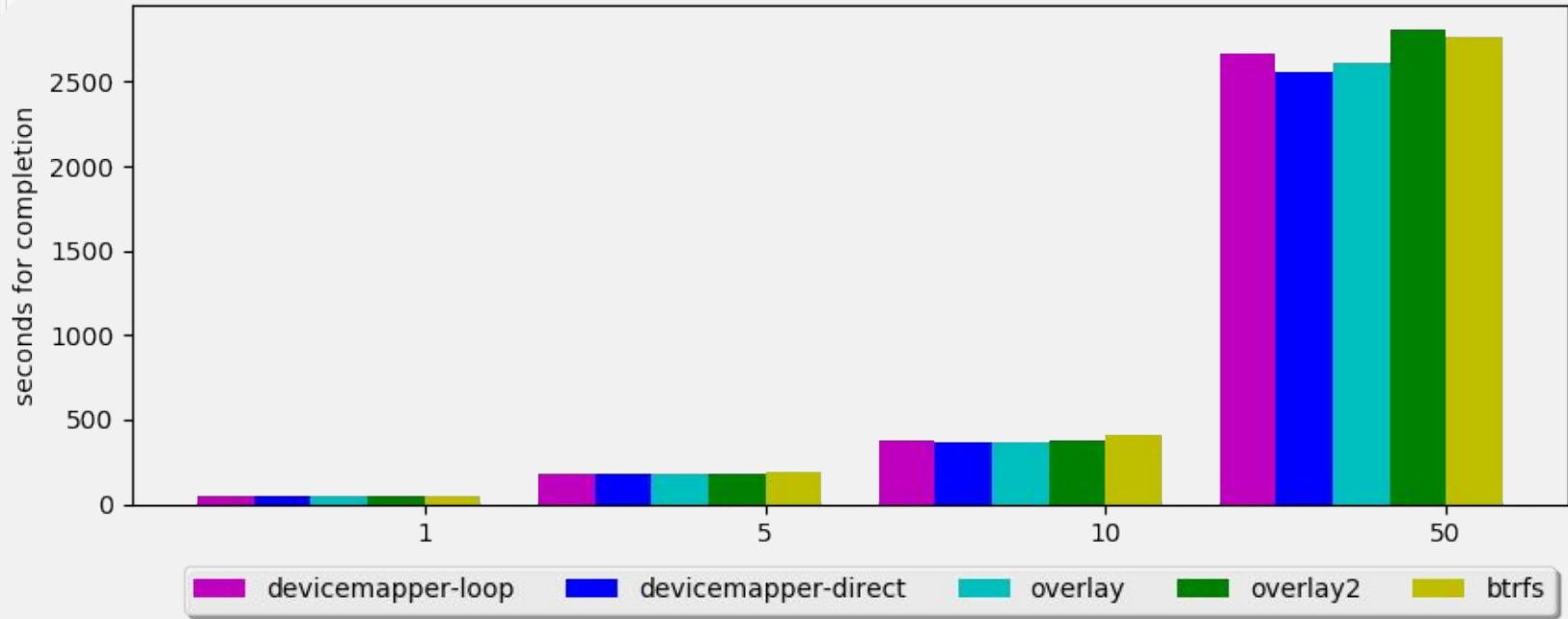
@KeithResar

read-big-files

Naive benchmarking, for discussion purposes only.  Don't trust this!

@KeithResar

read-file-tree

Naive benchmarking, for discussion purposes only.  Don't trust this!
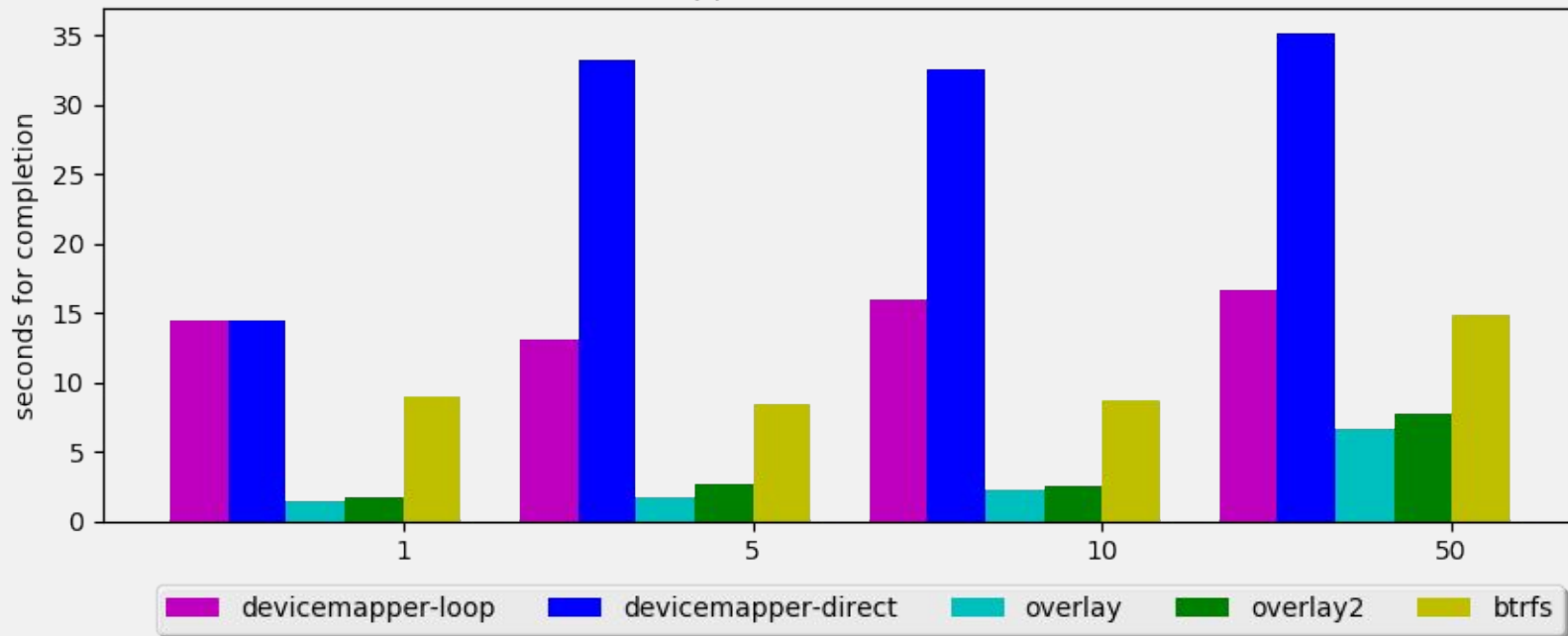
read-file-tree-mounted

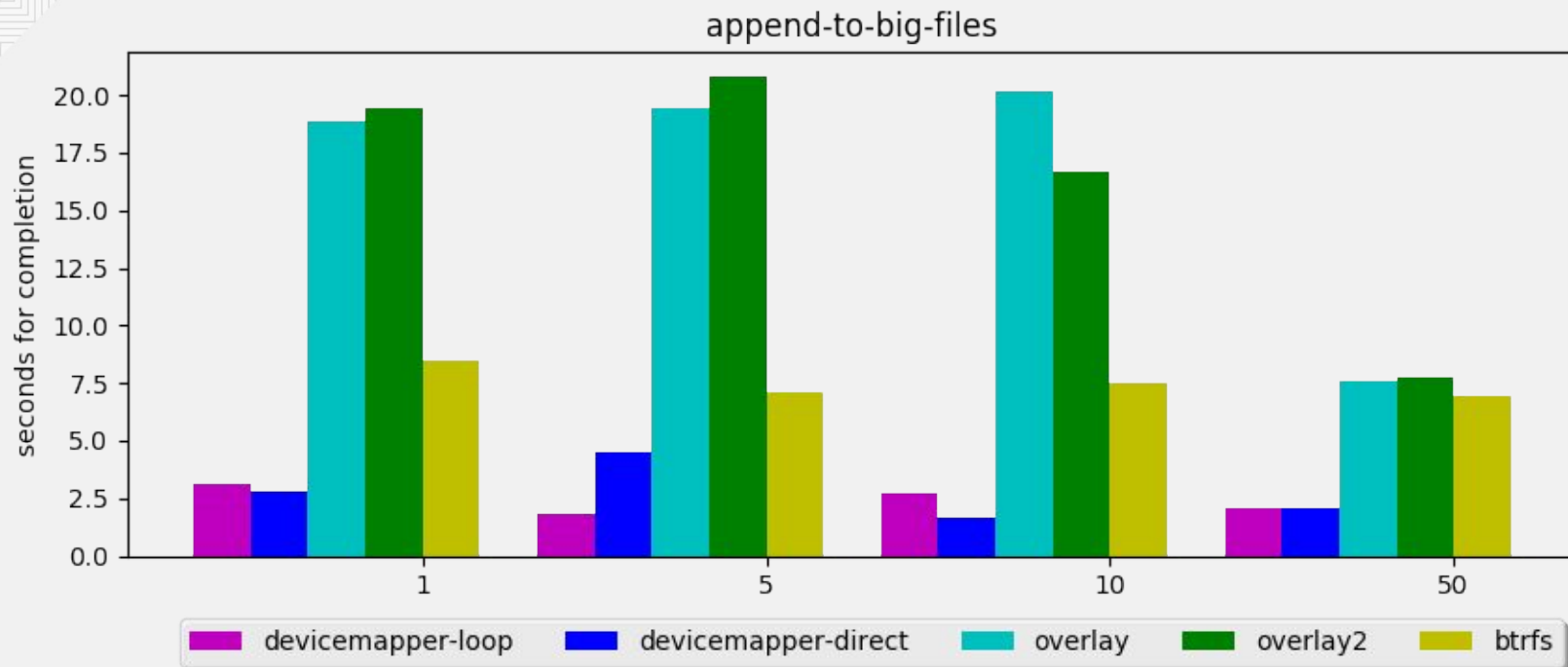Naive benchmarking, for discussion purposes only.  Don't trust this!

# Benchmark 2: Appending to Files

- Appending to Small Files
- Appending to Large Files
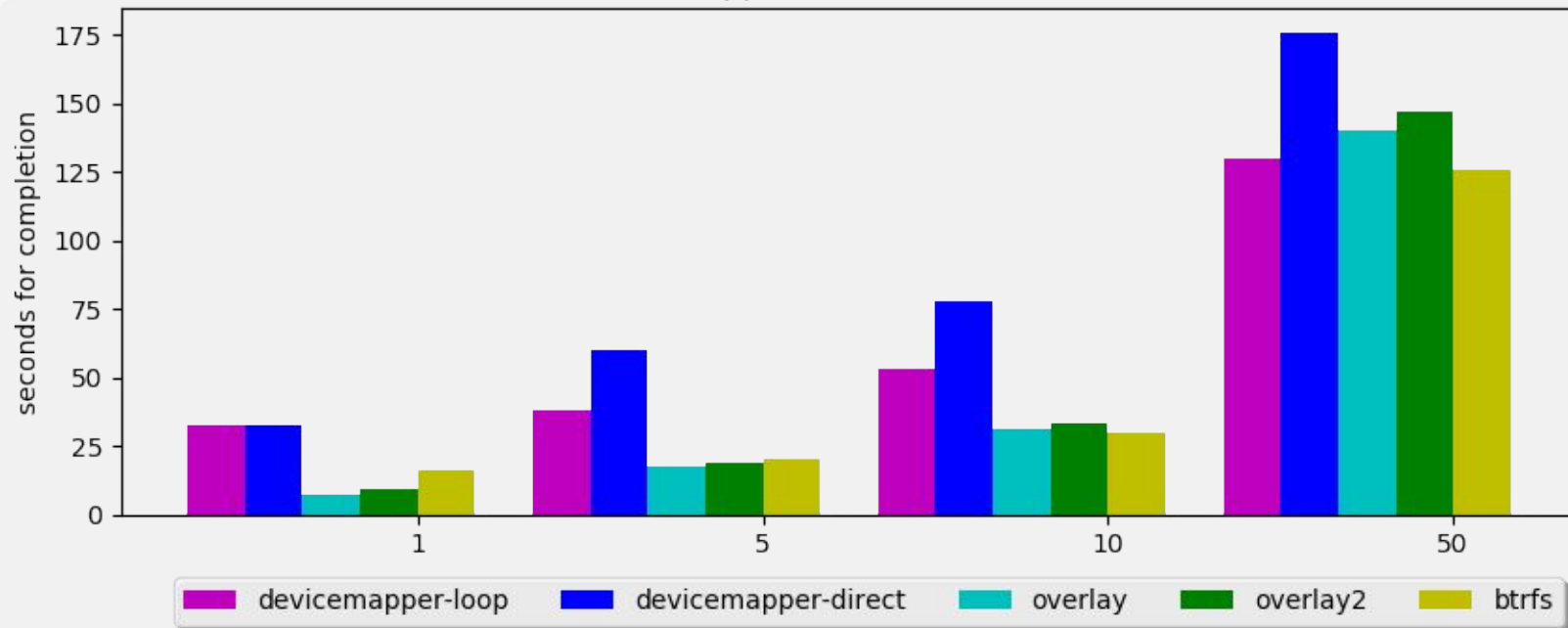- Appending to File Tree

redhat.

append-to-small-files

🛑 Naive benchmarking, for discussion purposes only. Don't trust this!

@KeithResar

append-to-big-files

Naive benchmarking, for discussion purposes only. Don't trust this!

append-to-file-tree

Naive benchmarking, for discussion purposes only.  Don't trust this!

# Storage use cases

| Technology | Attributes | Good Use Case | Bad Use Case |
| --- | --- | --- | --- |
| AUFS | Stable, Production Ready, Good Memory Use | | High Write Activity |
| Btrfs | Mainline Kernel | | High Write Activity |
| Overlay | Stable, Good Memory Use, Mainline Kernel | | Container Churn |
| Devicemapper (loop) | Stable, Mainline Kernel | | Production, Performance |
| Devicemapper (direct-lvm) | Stable, Production Ready, Mainline Kernel | | |

redhat.

# Resources

Storage Drivers in Docker: A Deep Dive
https://integratedcode.us/2016/08/30/storage-drivers-in-docker-a-deep-dive/

The Docker community has documented a good bit of this detail in the official storage driver documentation
https://docs.docker.com/engine/userguide/storagedriver/selectadriver/

Docker Issues and Tips (aufs/overlay/btrfs..)
https://github.com/AkihiroSuda/issues-docker#docker-issues-and-tips-aufsoverlaybtrfs

Comprehensive Overview of Storage Scalability in Docker (2014)
https://developers.redhat.com/blog/2014/09/30/overview-storage-scalability-docker/