

The last word in file systems

Cracow, pkgsrcCon, 2016

ZFS

The Last Word In File Systems

Jeff Bonwick

Bill Moore

www.opensolaris.org/os/community/zfs





WHEEL
S Y S T E M S



FreeBSD

Mariusz Zaborski

<m.zaborski@wheelsystems.com>

<oshogbo@FreeBSD.org>



Outline

1. FreeBSD
2. ZFS
3. Lunch break



Open**ZFS**



freeBSD[®]





FreeBSD

Do you use FreeBSD?





FreeBSD

Do you use ZFS?





FreeBSD

FreeBSD... only “a few” words





FreeBSD - who are we, what do we do?

- modern, open-source operating systems
 - 22 years old
 - BSD license
 - can't use the Internet without using FreeBSD
 - ~20,000 third party software ports
 - A democratically run open source project
-

FreeBSD - who are we, what do we do?



- Complete integrated UNIX system
 - multi-processing, multi-threaded kernel
 - Intel/AMD 32/64-bit, ARM, MIPS, PPC sparc64, ia64
 - UNIX, POSIX, BSD programming interfaces
 - Multi-protocol network stack
 - IPv4, IPv6, IPsec, ATM, SCTP, 802.11, Bluetooth...
-

FreeBSD - who are we, what do we do?



- Complete integrated UNIX system
 - World
 - Kernel
 - Toolchain
 - Extensive documentation





FreeBSD

Who uses FreeBSD? - **You!**

- WhatsApp
- Sony
- NetApp
- Apple
- Yahoo
- Netflix
- Juniper
- Wheel Systems
- Swisscom
- Microsoft Azure
- iXsystems
- Cisco
- Netgate
- NYI
- Verisign
- Dell KACE
- Dell/Compellen
t



FreeBSD

FreeBSD some features:

- UFS2
 - ZFS
 - DTrace
 - Jails
 - Bhyve
 - LLVM/Clang/lldb
 - Capsicum
 - MAC Framework
 - netmap
 - Linuxulator
-



Open**ZFS**

ZFS



Are our hard disks invincible?



Open**ZFS**



Are our hard disks invincible?



Open**ZFS**



Please write to offsets A, B, C, D



Please write to offsets A, B, D



Done, are you happy CPU?

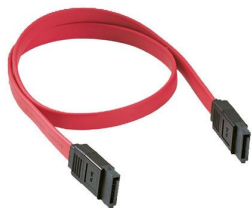
Are our hard disks invincible?



OpenZFS



Please write 0xFF 0x90 0xF9



Please write 0xFF 0x91 0xF9



Done, are you happy CPU?

Are our hard disks invincible?



Open**ZFS**





Thats joke, right?

OpenSSH CVE-2002-0083 - privilege escalation to root **OpenZFS**

```
- if (id < 0 || id > channels_alloc) {  
+ if (id < 0 || id >= channels_alloc) {
```



OpenZFS

Thats joke, right?

Assembly

```
    cmp    $0x0, 0x8(%ebp)
    js     16
    mov    0x4,%eax
    cmp    %eax, 0x8(%ebp)
    j1 30
    mov    0x8(%ebp), %eax
    mov    %eax, 0x4(%esp)
    movl   $0x4c, (%esp)
    call  25
```

```
    cmp    $0x0, 0x8(%ebp)
    js     16
    mov    0x4,%eax
    cmp    %eax, 0x8(%ebp)
    jle 30
    mov    0x8(%ebp), %eax
    mov    %eax, 0x4(%esp)
    movl   $0x4c, (%esp)
    call  25
```

Thats joke, right?

Hex



Open**ZFS**

39 45 08 **7c** 1a 8b 45

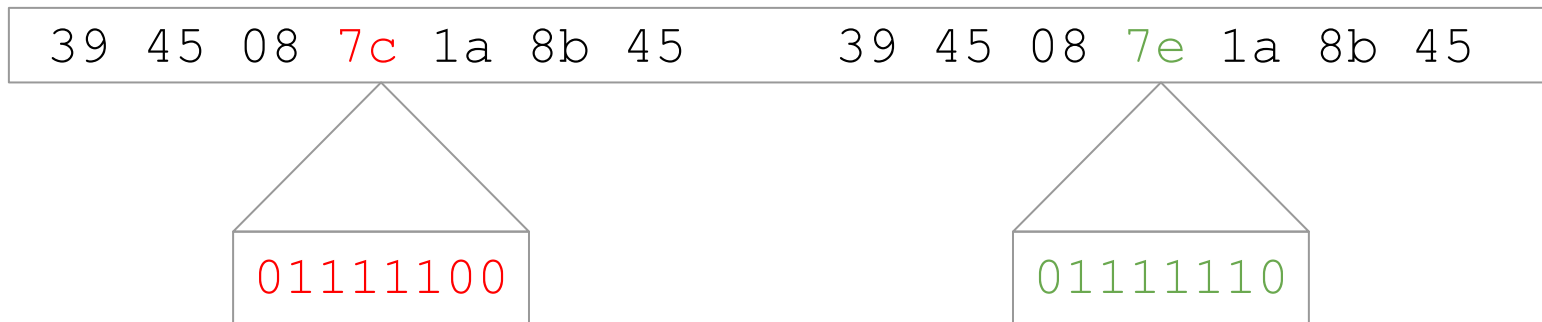
39 45 08 **7e** 1a 8b 45

Thats joke, right?

Binary



OpenZFS

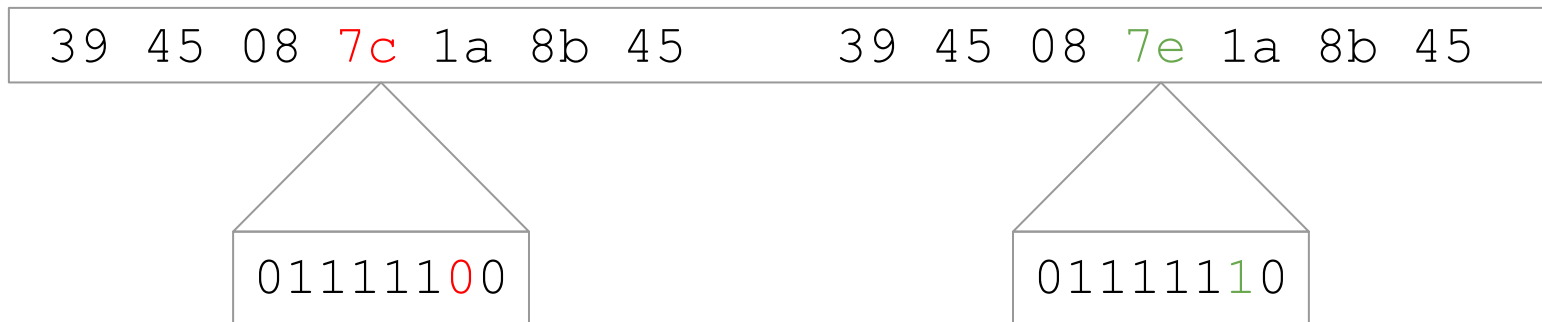


Thats joke, right?

Binary



OpenZFS





ZFS - history

2001: Closed-source development of ZFS started with two engineers at Sun Microsystems.

2005: Source code was released as part of OpenSolaris.

2007: Apple started porting of ZFS to Mac OS X.

2008: A port to FreeBSD was released as part of FreeBSD 7.0.

2008: Development of a native ZFS Linux port started, known as ZFS on Linux.

ZFS - history



2010: OpenSolaris was discontinued, the last release was forked. Further development of ZFS on Solaris was no longer open source.

2010: illumos was founded as the truly open source successor to OpenSolaris.

2013: Official announcement of the OpenZFS project.





Open**ZFS**

Why ZFS is ZFS?

1 kilobyte	1,000,000,000,000,000,000,000,000
1 megabyte	1,000,000,000,000,000,000,000,000
1 gigabyte	1,000,000,000,000,000,000,000,000
1 terabyte	1,000,000,000,000,000,000,000,000
1 petabyte	1,000,000,000,000,000,000,000,000
1 exabyte	1,000,000,000,000,000,000,000,000
1 zettabyte	1,000,000,000,000,000,000,000,000

Why ZFS is ZFS?

1 kilobyte	1,000,000,000,000,000,000,000,000
1 megabyte	1,000,000,000,000,000,000,000,000
1 gigabyte	1,000,000,000,000,000,000,000,000
1 terabyte	1,000,000,000,000,000,000,000,000
1 petabyte	1,000,000,000,000,000,000,000,000
1 exabyte	1,000,000,000,000,000,000,000,000
1 zettabyte	1,000,000,000,000,000,000,000,000



Open**ZFS**

160Gb ~ 275\$

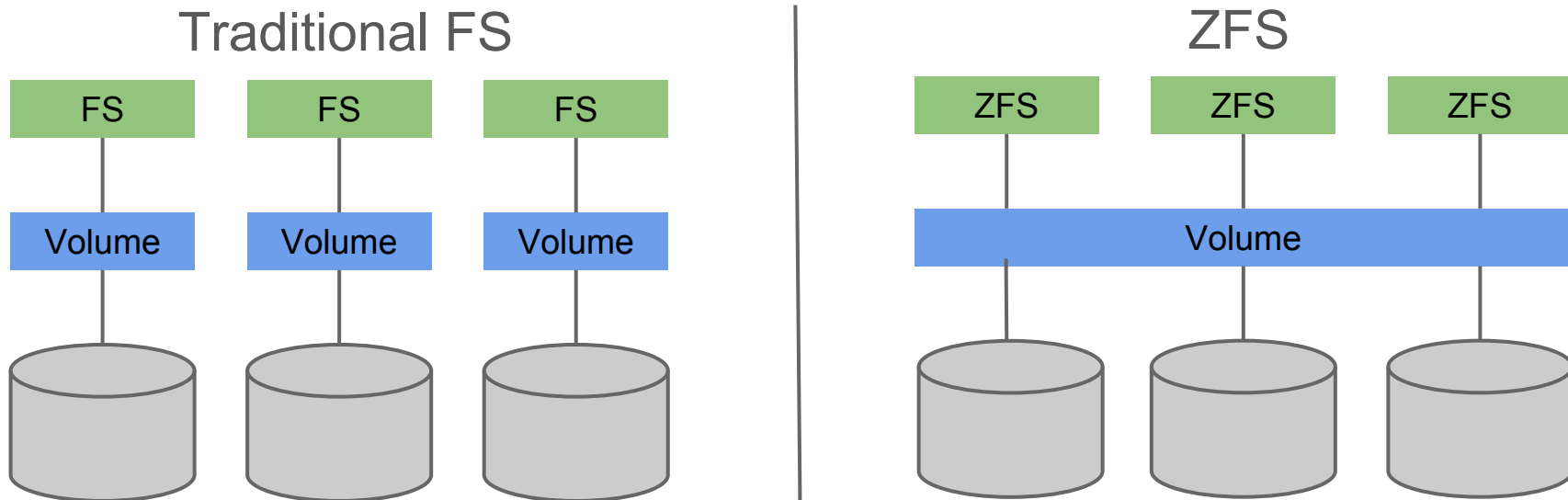


ZFS - overview

- Storage pool, integrated volume manager.
- Dynamic striping.



Open**ZFS**



ZFS - overview



OpenZFS

- Storage pool, integrated volume manager.

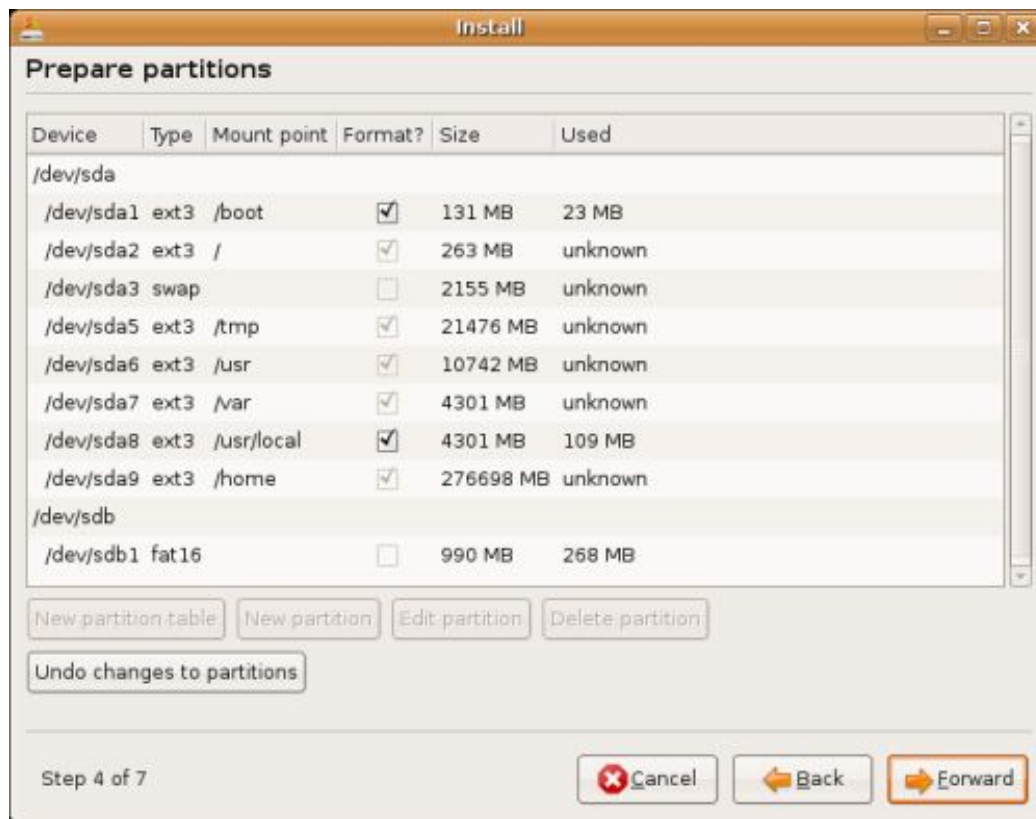
```
$ df -h
```

Filesystem	Size	Used	Avail	Capacity	Mounted on
zroot	155G	526M	155G	0%	/
devfs	1.0K	1.0K	0B	100%	/dev
zroot/tmp	155G	261M	155G	0%	/tmp
zroot/usr	172G	17G	155G	10%	/usr
zroot/usr/home	186G	31G	155G	17%	/usr/home
zroot/usr/ports	156G	1.2G	155G	1%	/usr/ports
zroot/var	159G	4.0G	155G	3%	/var
zroot/var/crash	155G	96K	155G	0%	/var/crash
zroot/var/db	155G	216M	155G	0%	/var/db
zroot/var/log	155G	512K	155G	0%	/var/log
zroot/var/mail	155G	188K	155G	0%	/var/mail
zroot/var/tmp	155G	1.2M	155G	0%	/var/tmp

ZFS



OpenZFS





Open**ZFS**

ZFS - overview

- End-to-end checksumming
 - SHA256, fletcher2, fletcher4
- Copy-on-write
- Transactional operations
 - Always consistent on disk
 - No fsck, no journaling



What is journaling?

- Stores metadata
- Recovery mechanism
- Used in:
 - Ext3
 - UFS with SU+J

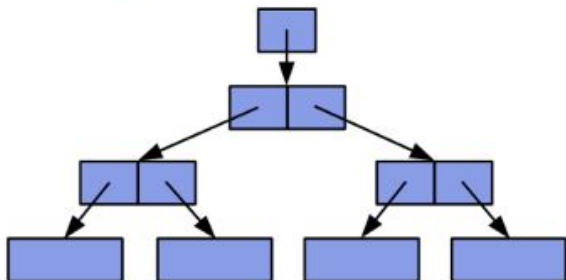


ZFS - CoW

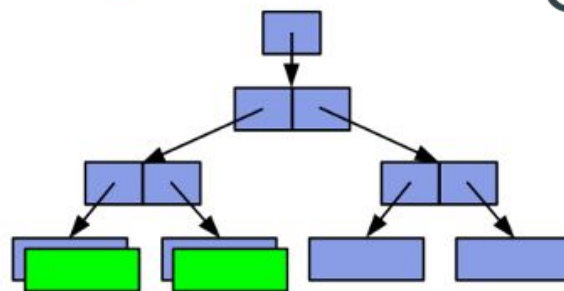


OpenZFS

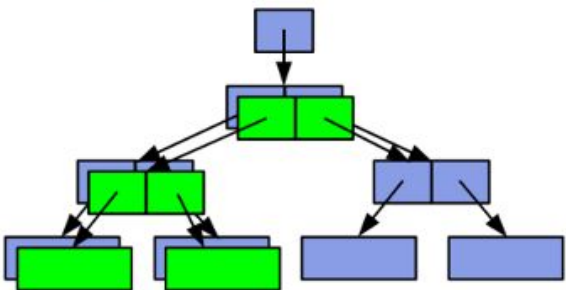
1. Initial block tree



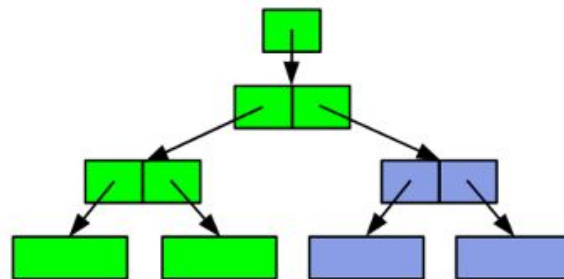
2. COW some blocks



3. COW indirect blocks

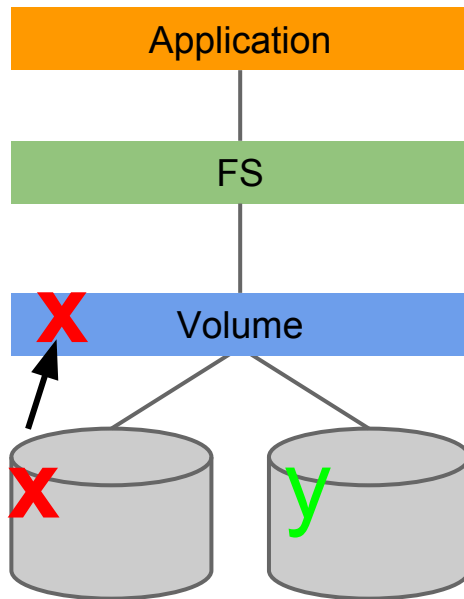


4. Rewrite uberblock (atomic)



ZFS - overview

Traditional FS

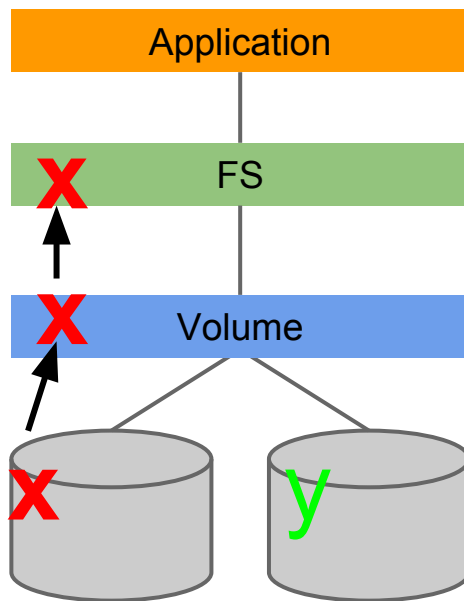


Open**ZFS**



ZFS - overview

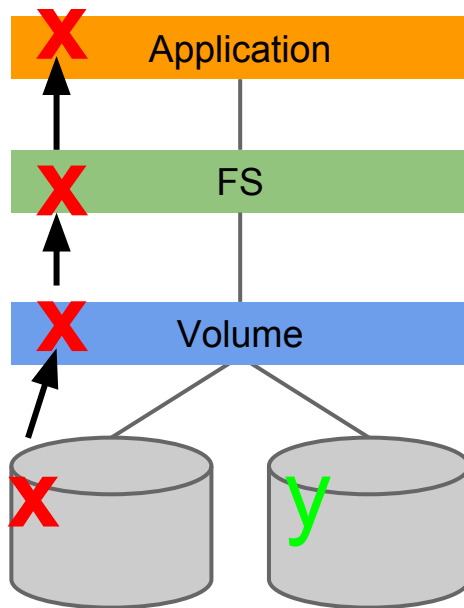
Traditional FS



Open**ZFS**

ZFS - overview

Traditional FS



OpenZFS

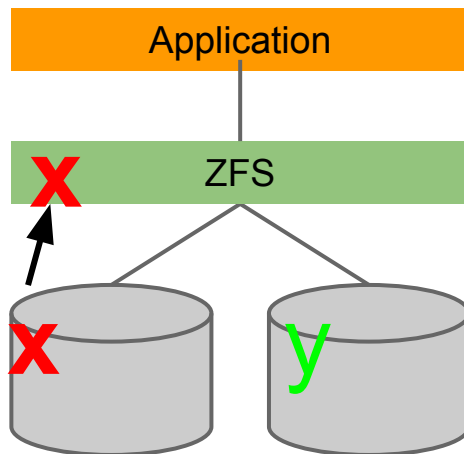


ZFS - overview

ZFS



OpenZFS

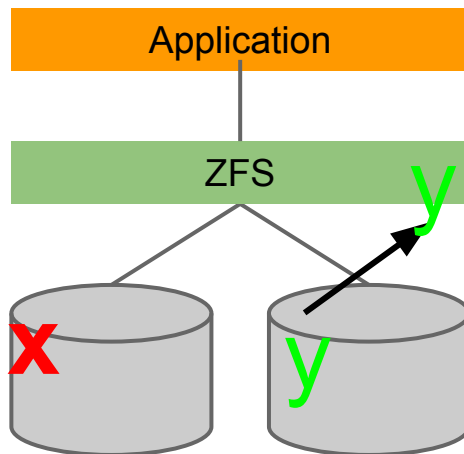


ZFS - overview

ZFS self-healing



OpenZFS

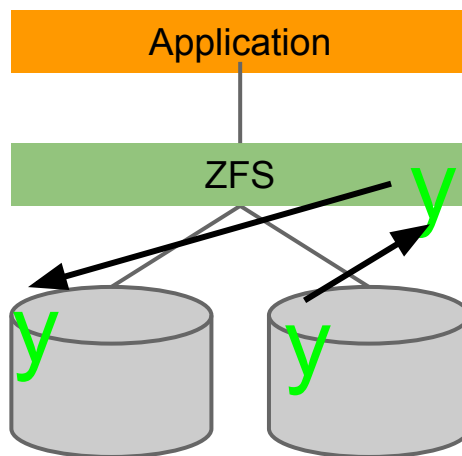


ZFS - overview

ZFS self-healing



Open**ZFS**

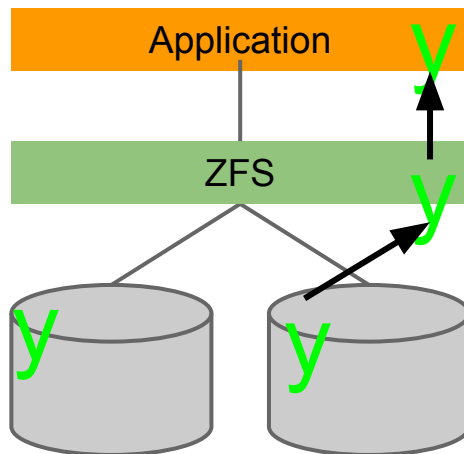


ZFS - overview

ZFS self-healing



Open**ZFS**

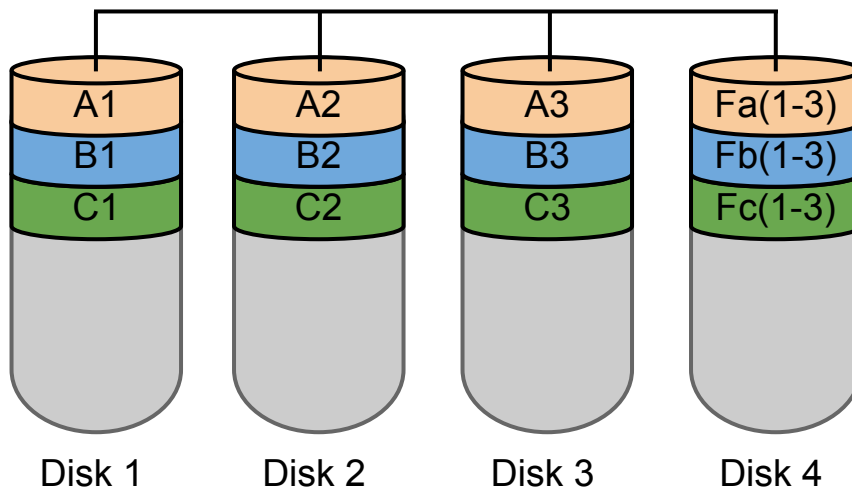


ZFS redundancy

- Stripe
- Mirror (RAID-1)
- Raid-Z{1,2}



Open**ZFS**



ZFS Raidz-Z{1,2,3}



Open**ZFS**

Raid	Minimal hard disk count	Recommended amount of hard disks	Parity hard disk	Data hard disk
RAIDZ1	2	3	1	n - 1
RAIDZ2	3	4	2	n - 2
RAIDZ3	4	6	3	n - 3



ZFS - other features

- Endian-independent
- Compression
 - LZ4
 - LZJB
 - GZIP
 - ZLE



Open**ZFS**

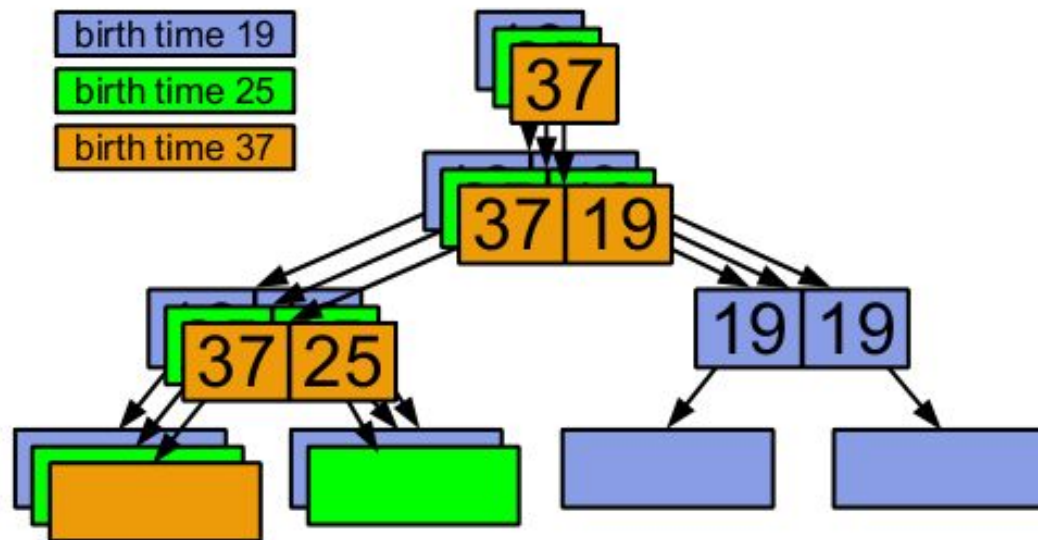


ZFS - other features



Open**ZFS**

- Snapshots
 - Very very cheap - $O(1)$
 - Birth time



ZFS - other features



Open**ZFS**

- Snapshots

```
zfs snapshot tank/home/oshogbo@201606
```

```
zfs rollback tank/home/oshogbo@201512
```

```
cd ~/oshogbo/.zfs/snapshot/201512/
```

ZFS - other features

- Incremental snapshots
- Clones
- ZVOL
- Deduplication
 - verify
 - sha256



Open**ZFS**



ZFS - other features

- Quotas
- Reservations
- NFS
- Resilvering



Open**ZFS**



ZFS - when to use it?



Open**ZFS**

- Backups

```
zfs send -i 201605 tank@201606 | ssh mbackup zfs recv -F obackup
```

- Bhyve/jails

```
# zfs get -o property,value all zroot/vm0001 bhyve:cpus 4
```

```
bhyve:memory 2048
```

```
bhyve:console /dev/nmdm0001A
```

```
bhyve:networks [{type:virtio-net,name:tap7}]
```

```
bhyve:disks [ { type:ahci-hd, path:/dev/zvol/zroot/disk0001 },  
              { type:ahci-hd, path:/dev/zvol/zroot/disk0002 } ]
```

ZFS - when to use it?



Open**ZFS**

- Cluster multi-master

NAME	USED	AVAIL	REFER	MOUNTPOINT
data/data/10000013	40M	6.48G	33K	/data/10000013
data/data/10000013/data	29M	6.48G	29.3M	/data/10000013/data
data/data/10000013/pdfs	11M	6.48G	11.6M	/data/10000013/pdfs
data/data/local	4.74G	6.48G	4.57G	/data/local
data/data/local/data	178M	6.48G	178M	/data/local/data
data/data/local/pdfs	436K	6.48G	402K	/data/local/pdfs

ZFS - ongoing work

- Resumable send/recv (FreeBSD 11)
- FS encryption
- Sending compressed blocks



Open**ZFS**





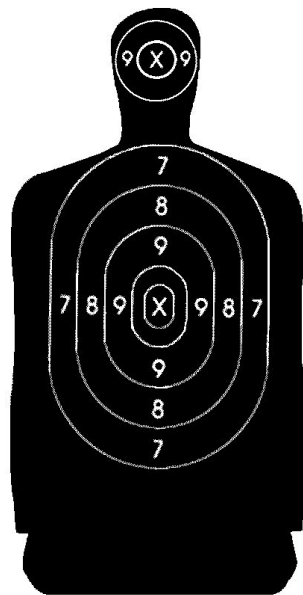
Open**ZFS**

ZFS downsides

- Requires a lot of RAM.
 - An insane amount of RAM (1GB RAM per 1TB of storage).
 - Not so useful for embedded systems?
 - Defragmentation
 - Really slow on small capacity
-



Thank you for your attention!



<m.zaborski@wheelsystems.com>

<oshogbo@FreeBSD.org>

