# An Introduction to the Implementation of ZFS

Brought to you by

Dr. Marshall Kirk McKusick

BSD Canada Conference 2015
June 13, 2015
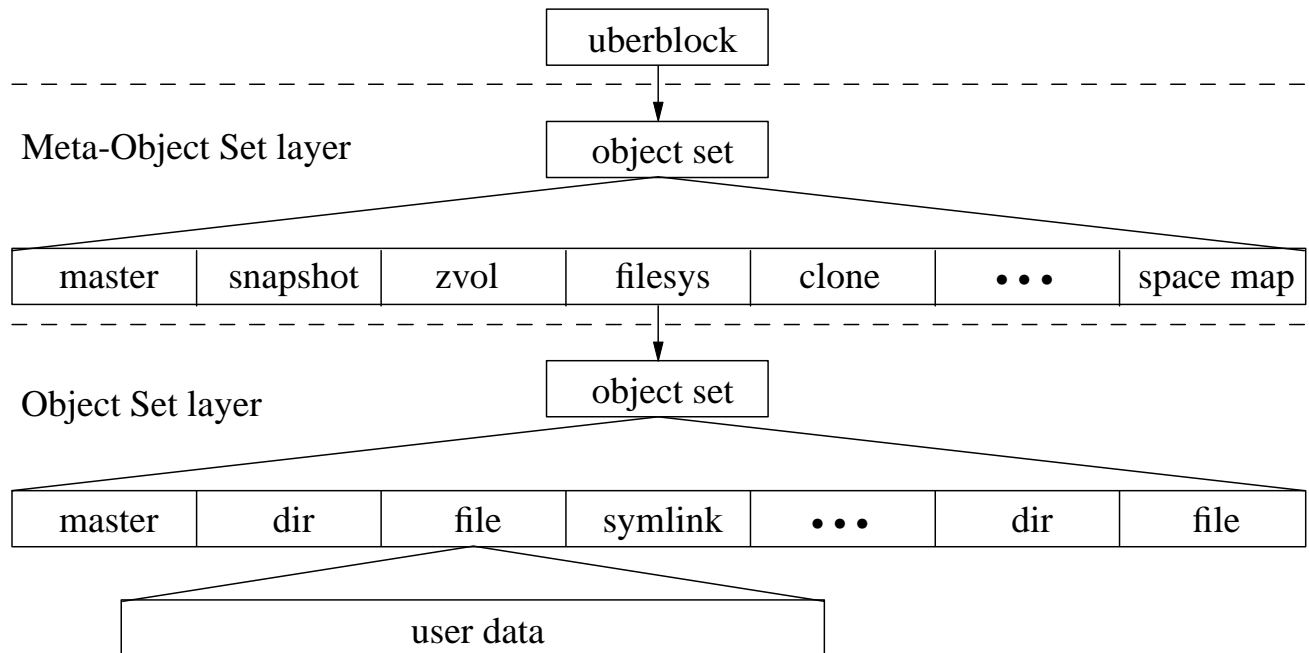
University of Ottawa
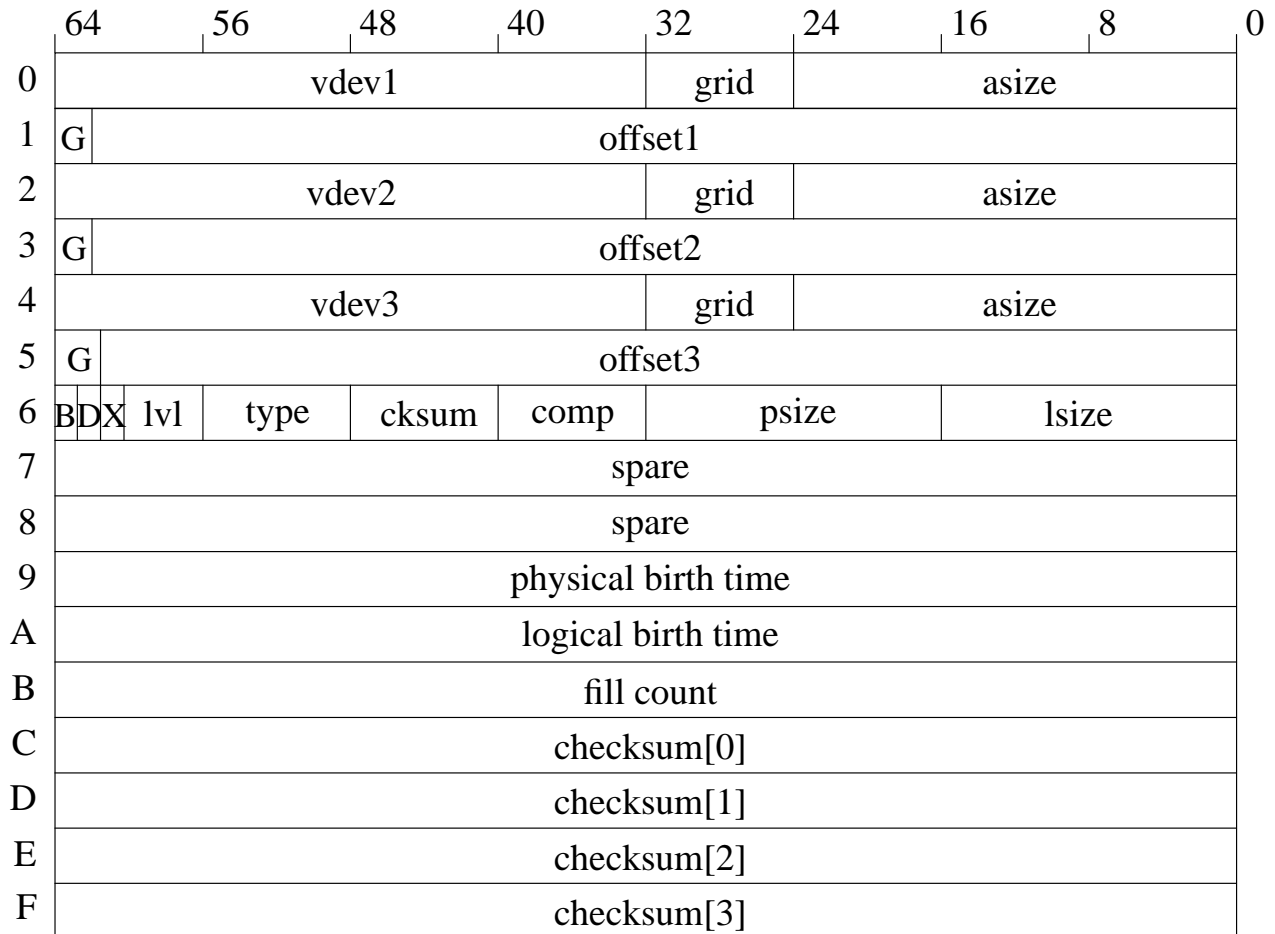Ottawa, Canada

# Zettabyte Filesystem Overview

- Never over-write an existing block

- Filesystem is always consistent

- State atomically advances at checkpoints

- Snapshots (read-only) and clones (read-write) are cheap and plentiful

- Metadata redundancy and data checksums

- Selective data compression and deduplication

- Pooled storage shared among filesystems

- Mirroring and single, double, and triple parity RAIDZ

- Space management with quotas and reservations

- Fast remote replication and backups

# Structural Organization

```
                          ┌──────────────┐
                          │  uberblock   │
                          └──────┬───────┘
- - - - - - - - - - - - - - - - -│- - - - - - - - - - - - - - - - -
                                 ▼
Meta-Object Set layer     ┌──────────────┐
                          │  object set  │
                          └──────────────┘

┌──────────┬──────────┬────────┬─────────┬─────────┬────────┬───────────┐
│  master  │ snapshot │  zvol  │ filesys │  clone  │ • • •  │ space map │
└──────────┴──────────┴────────┴─────────┴─────────┴────────┴───────────┘
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -
                                 ▼
Object Set layer          ┌──────────────┐
                          │  object set  │
                          └──────────────┘

┌──────────┬──────────┬────────┬─────────┬─────────┬────────┬───────────┐
│  master  │   dir    │  file  │ symlink │ • • •   │  dir   │   file    │
└──────────┴──────────┴────────┴─────────┴─────────┴────────┴───────────┘
                     ┌─────────────────────────┐
                     │       user data         │
                     └─────────────────────────┘
```

- Uberblock anchors the pool

- Meta-object-set (MOS) describes array of filesystems, clones, snapshots, and ZVOLs

- Each MOS object references an object-set that describes its objects

- Filesystem object sets describe an array of files, directories, etc.

- Each filesystem object describes an array of bytes
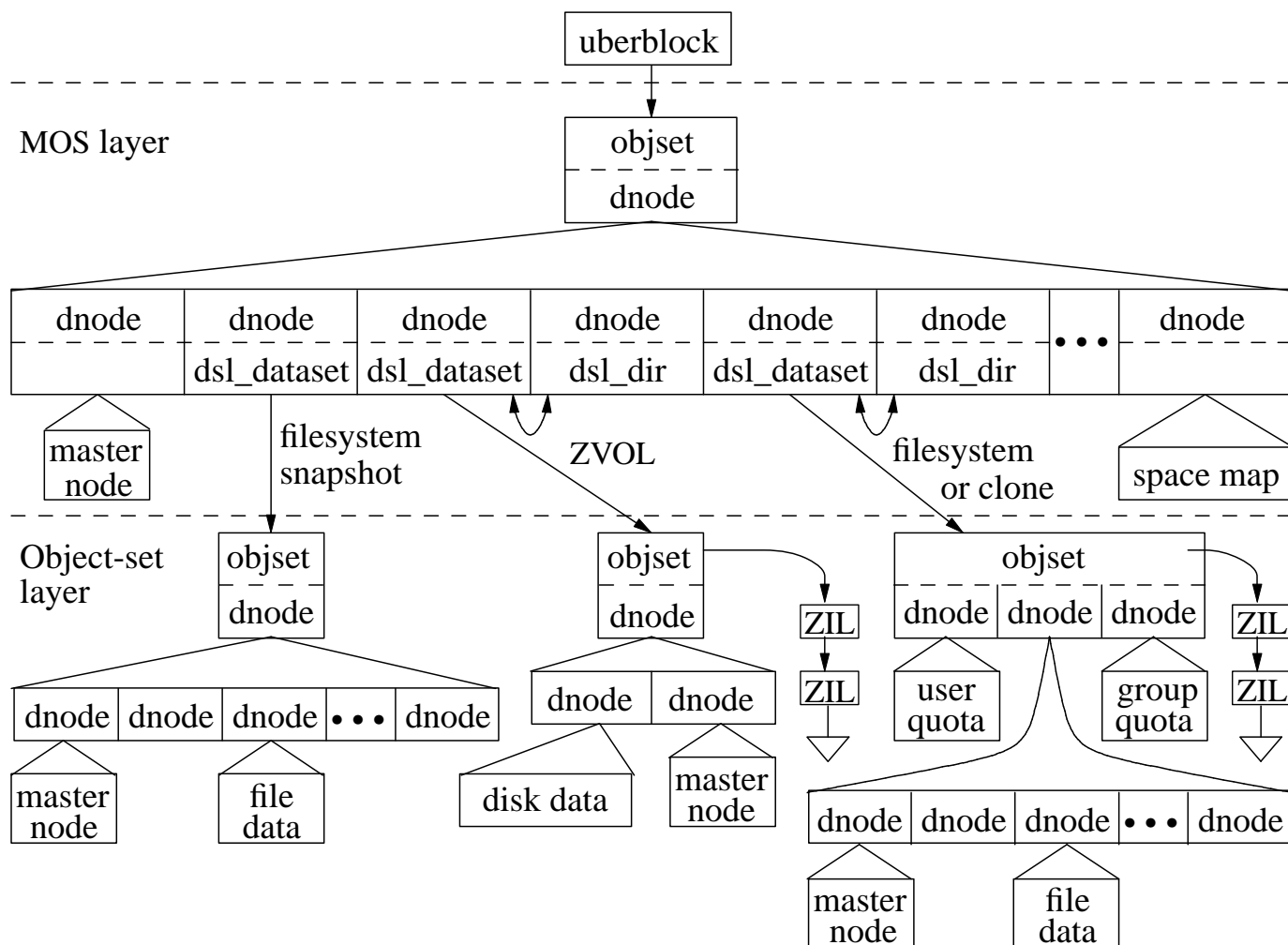
# ZFS Block Pointer

| | 64 | 56 | 48 | 40 | 32 | 24 | 16 | 8 | 0 |
|---|---|---|---|---|---|---|---|---|---|
| 0 | | vdev1 | | | | grid | asize | | |
| 1 | G | | | offset1 | | | | | |
| 2 | | vdev2 | | | | grid | asize | | |
| 3 | G | | | offset2 | | | | | |
| 4 | | vdev3 | | | | grid | asize | | |
| 5 | G | | | offset3 | | | | | |
| 6 | B D X | lvl | type | cksum | comp | | psize | lsize | |
| 7 | | | | spare | | | | | |
| 8 | | | | spare | | | | | |
| 9 | | | | physical birth time | | | | | |
| A | | | | logical birth time | | | | | |
| B | | | | fill count | | | | | |
| C | | | | checksum[0] | | | | | |
| D | | | | checksum[1] | | | | | |
| E | | | | checksum[2] | | | | | |
| F | | | | checksum[3] | | | | | |

- Up to three levels of redundancy

- Checksum separate from data

- Birth time is the transaction-group number in which it was allocated

- Maintains allocated, physical (compressed), and logical sizes

# ZFS Block Management

- Disk blocks are kept in a pool

- Multiple filesystems and their snapshots are held in the pool

- Blocks from the pool are given to filesystems on demand and reclaimed to the pool when freed

- Space may be reserved to ensure future availability

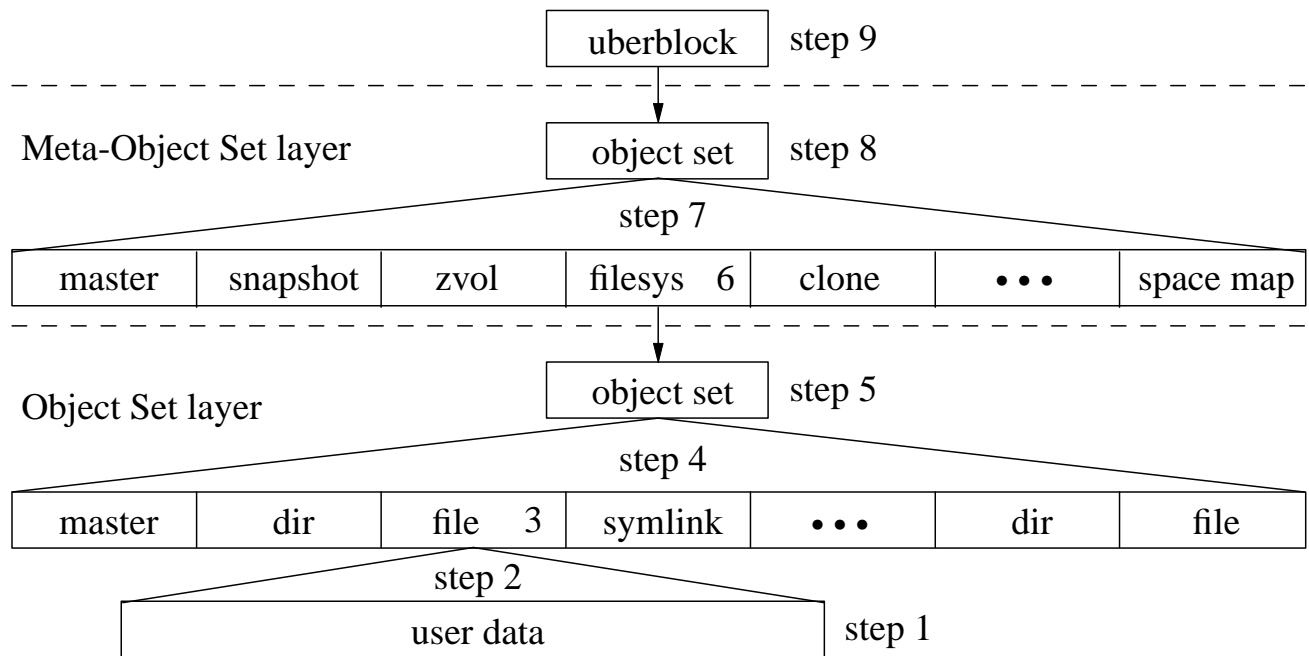- Quotas may be imposed to limit the space that may be used

# ZFS Structure



- MOS layer manages space and objects using that space

- Object-set layer manages filesystems, snapshots, clones, and ZVOLs

# ZFS Checkpoint

- Collect all updates in memory

- Periodically write all changes to an unused location to create a checkpoint

- Last step in checkpoint writes a new uberblock

- Entire pool is always consistent

- Checkpoint affects all filesystems, clones, snapshots, and ZVOL in the pool

- Need to log any changes between checkpoints that need to be persistent

- The **fsync** system call is implemented by forcing a log write not by doing a checkpoint

- Recovery starts from last checkpoint, rolls forward through log, then creates new checkpoint

# Flushing Dirty Data

```
                          ┌─────────────┐
                          │  uberblock  │  step 9
                          └─────────────┘
─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─

Meta-Object Set layer     ┌─────────────┐
                          │  object set │  step 8
                          └─────────────┘
                               step 7
┌─────────┬──────────┬─────────┬──────────┬─────────┬────────┬────────────┐
│ master  │ snapshot │  zvol   │ filesys 6│  clone  │ • • •  │ space map  │
└─────────┴──────────┴─────────┴──────────┴─────────┴────────┴────────────┘
─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─ ─

Object Set layer          ┌─────────────┐
                          │  object set │  step 5
                          └─────────────┘
                               step 4
┌─────────┬──────────┬─────────┬──────────┬─────────┬────────┬────────────┐
│ master  │   dir    │ file  3 │ symlink  │ • • •   │  dir   │   file     │
└─────────┴──────────┴─────────┴──────────┴─────────┴────────┴────────────┘
                       step 2
              ┌──────────────────────────┐
              │        user data         │  step 1
              └──────────────────────────┘
```

Write modified data in this order:

1) new or modified user data

2) indirect blocks to new user data

3) new dnode block

4) indirect blocks to modified dnodes

5) object-set dnode for filesystem dnodes

6) filesystem dnode to reference objset dnode

7) indirect blocks to modified meta-objects

8) MOS object-set for meta-object dnode

9) new uberblock (plus its copies)

text ref: pp. 536-538

# RAIDZ

|  | disk 1 | disk2 | disk 3 | disk 4 | disk 5 |
|---|---|---|---|---|---|
| stripe 1 | P0 | D0 | D2 | D4 | D6 |
| stripe 2 | P1 | D1 | D3 | D5 | D7 |
| stripe 3 | P0 | D0 | D1 | D2 | P0 |
| stripe 4 | D0 | D1 | D2 | P0 | D0 |
| stripe 5 | P0 | D0 | D4 | D8 | D11 |
| stripe 6 | P1 | D1 | D5 | D9 | D12 |
| stripe 7 | P2 | D2 | D6 | D10 | D13 |
| stripe 8 | P3 | D3 | D7 | P0 | D0 |
| stripe 9 | D1 | D2 | D3 | X | P0 |
| stripe 10 | D0 | D1 | X | P0 | D0 |
| stripe 11 | D3 | D6 | D9 | P1 | D1 |
| stripe 12 | D4 | D7 | D10 | P2 | D2 |
| stripe 13 | D5 | D8 |  |  |  |

- Variable size stripes since each block knows its size

- Parity sectors on disk that starts block

- N == 1, up to four blocks per parity sector.

- Blocks must be multiple of N + 1 sectors

# RAIDZ Recovery

- Rebuild by traversing all MOS objects and rebuild their blocks

    Fast when pool not fully allocated as only used blocks are rewritten

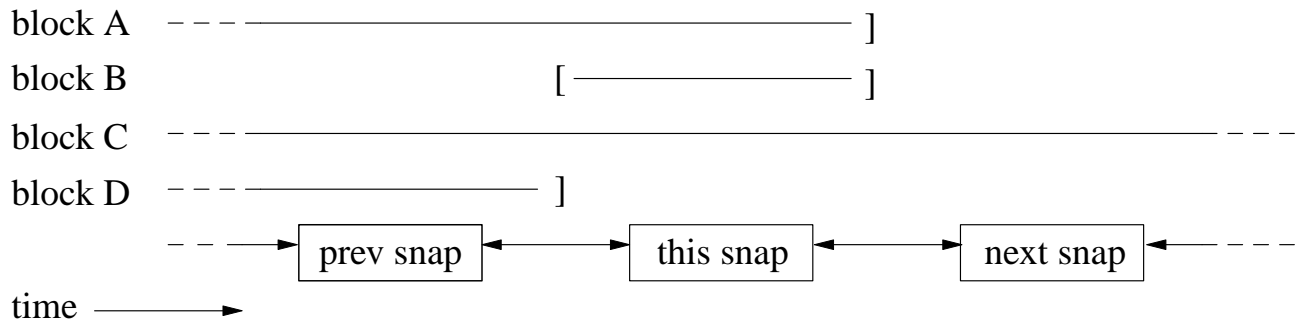    Slow when pool is full as many random reads needed

    use physical birth time to determine blocks that need to be rebuilt

- Never need to recalculate and write parity

# Freeing Filesystem Blocks

- Blocks are tracked using space maps, birth time, and deadlists

- When a block is allocated, its birth time is set to the current transaction number

- Over time, snapshots are taken which also reference the block

- When a file is overwritten, truncated, or deleted, its blocks are released

- For each freed block, the kernel must determine if a snapshot still references it

    if born after most recent snapshot, it can be freed

    otherwise it is added to the filesystem's deadlist

# Freeing Snapshot Blocks

```
block A  – – – ――――――――――――――――――  ]
block B                         [ ―――――――― ]
block C  – – – ―――――――――――――――――――――――――――― – – –
block D  – – – ――――――――――― ]
         – – →[ prev snap ]←――→[ this snap ]←――→[ next snap ]←– – –

time  ――――→
```

- Freeing "this snap"

    Iterate over "next snap" deadlist (blocks A and B)

    Block A born before "prev snap" so added to "this snap" deadlist

    Block B born after "prev snap" so must be freed

    Move deadlist of "this snap" to become deadlist of "next snap"

    Remove "this snap" from list of snapshots and directory of snapshot names

- Never need to make a pass over entire block-allocation map

# ZFS Strengths

- High write throughput

- Fast RAIDZ reconstruction on pools with less than 40% utilization

- Avoids RAID "write hole"

- Blocks move between filesystems as needed

- Integration eases administration (mount points, exports, etc)

# ZFS Weaknesses

- Slowly written files scattered on disk

- Slow RAIDZ reconstruction on pools with greater than 40% utilization

- Block cache must fit in the kernel's address space, thus works well only on 64-bit systems

- Needs under 75% utilization for good write performance

- RAIDZ has high overhead for small block sizes such as 4 Kbyte blocks typically used by databases and ZVOLs.

- Blocks cached in memory are not part of the unified memory cache so inefficient for files and executables using **mmap** or when using **sendfile**

# Questions

More on ZFS:

- "The Design and Implementation of the FreeBSD Operating System, 2nd Edition", Chapter 10

- Manual pages: zfs(8), zpool(8), zdb(8)

Marshall Kirk McKusick

<mckusick@mckusick.com>

http://www.mckusick.com