



ORACLE[®]



Filesystem Features and Performance

Chris Mason

Filesystems

- XFS
 - Well established and stable
 - Highly scalable under many workloads
 - Can be slower in metadata intensive workloads
 - Often the best option for production use today
- Ext4
 - Can read ext3 filesystems
 - Includes delayed allocation and extents
 - Lifts many limits from ext3
- Btrfs
 - Many new features

Filesystem Performance Metrics

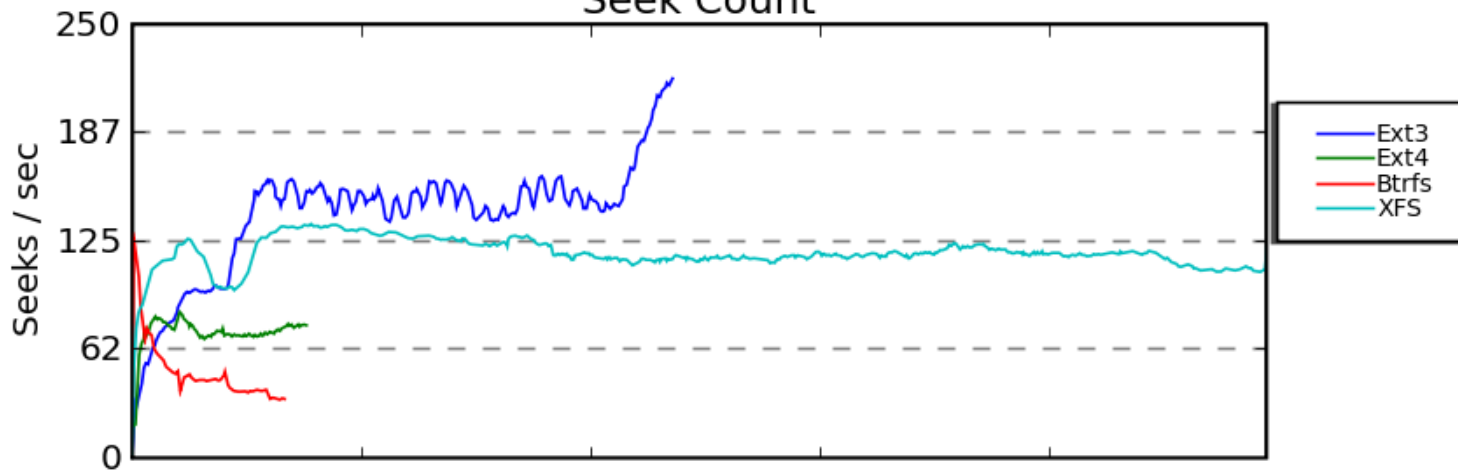
- File allocation and writeback
 - Extents are much more efficient
 - Crash recovery log can be a major performance factor
- Directory indexing
 - Every FS does this differently
 - Major impact on read and write performance
- Synchronous Writes
 - Expose bottlenecks in many different parts of the FS

File allocation and Writeback

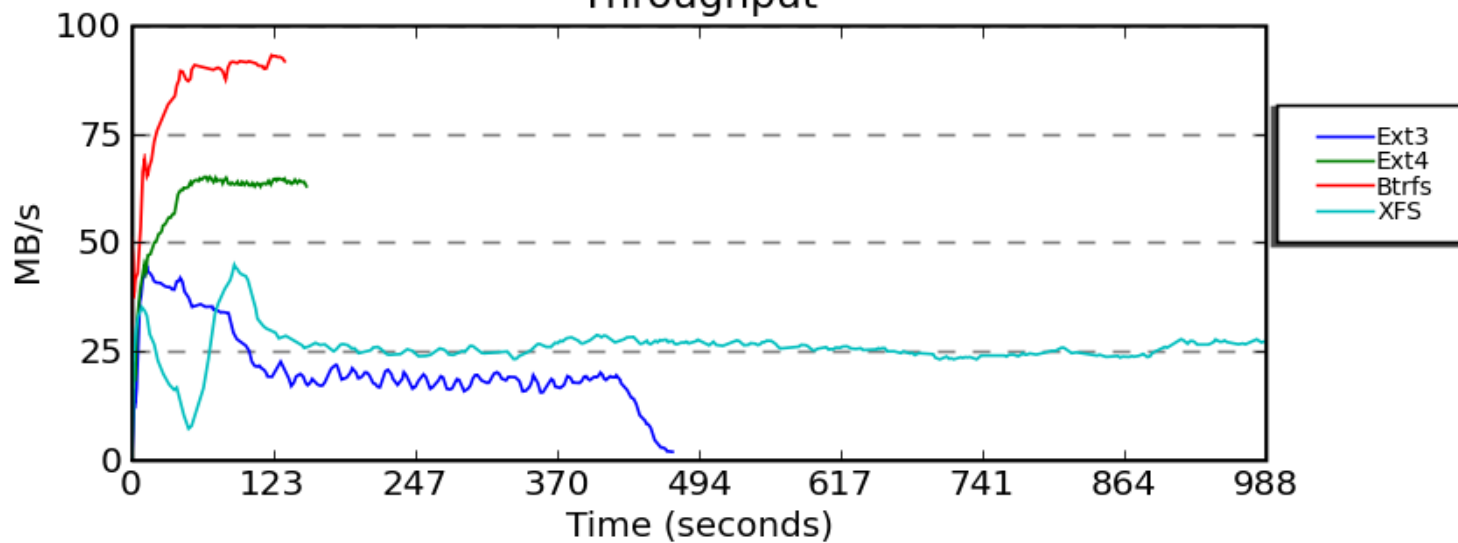
- Delayed allocation
 - Don't allocate blocks on disk until just before dirty pages are written. All 3 filesystems do this
- Extent based storage
 - Allocates a range of bytes instead of a fixed unit
 - Much more efficient for large files
- Transaction log
 - XFS uses a logical log of operations
 - Ext3 and Ext4 use a write ahead block log
 - Btrfs uses copy on write (COW)

Kernel tree copies

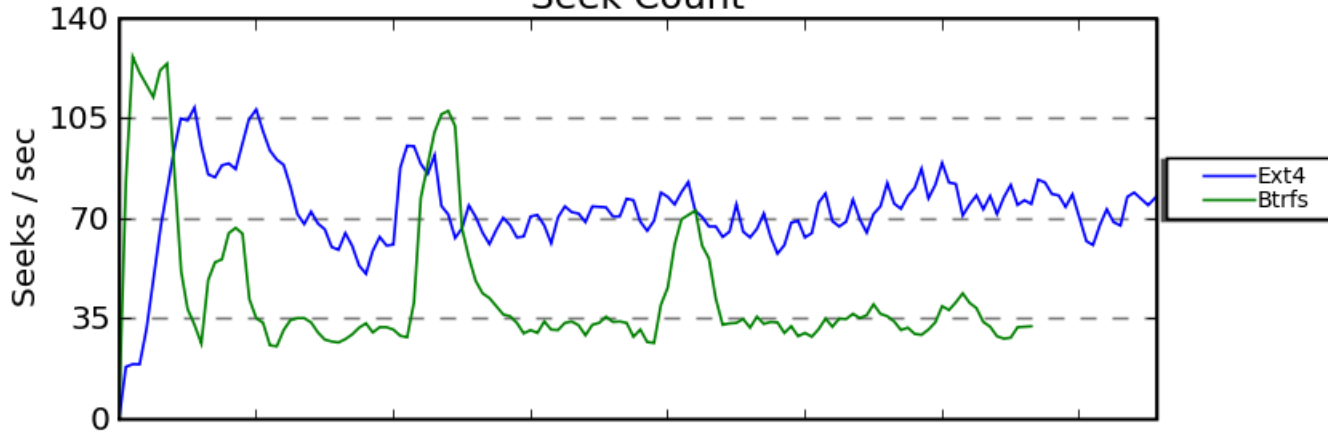
Seek Count



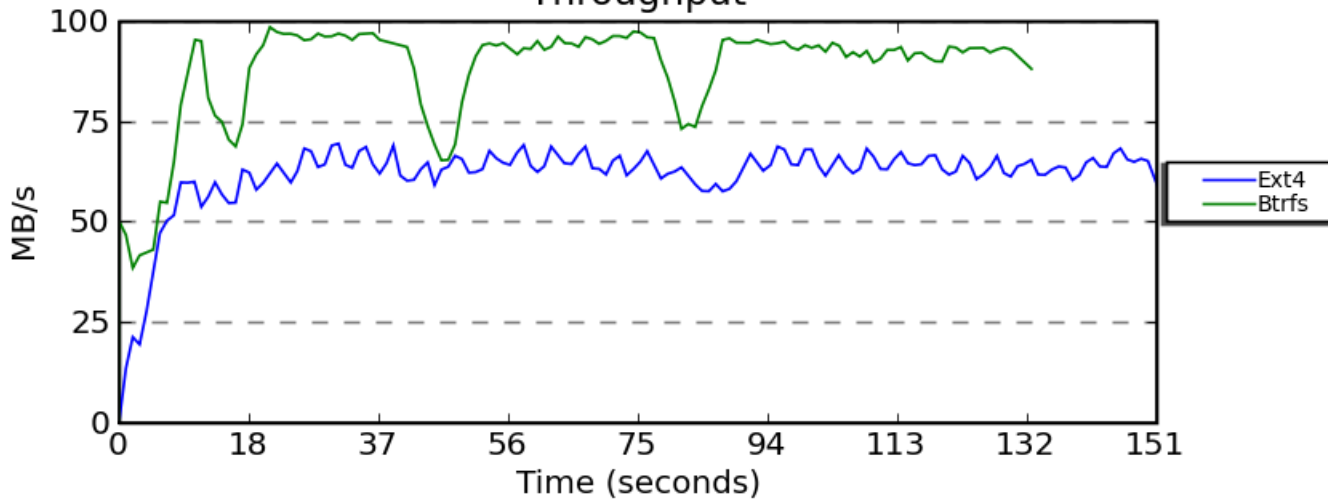
Throughput



Kernel tree copies Seek Count



Throughput



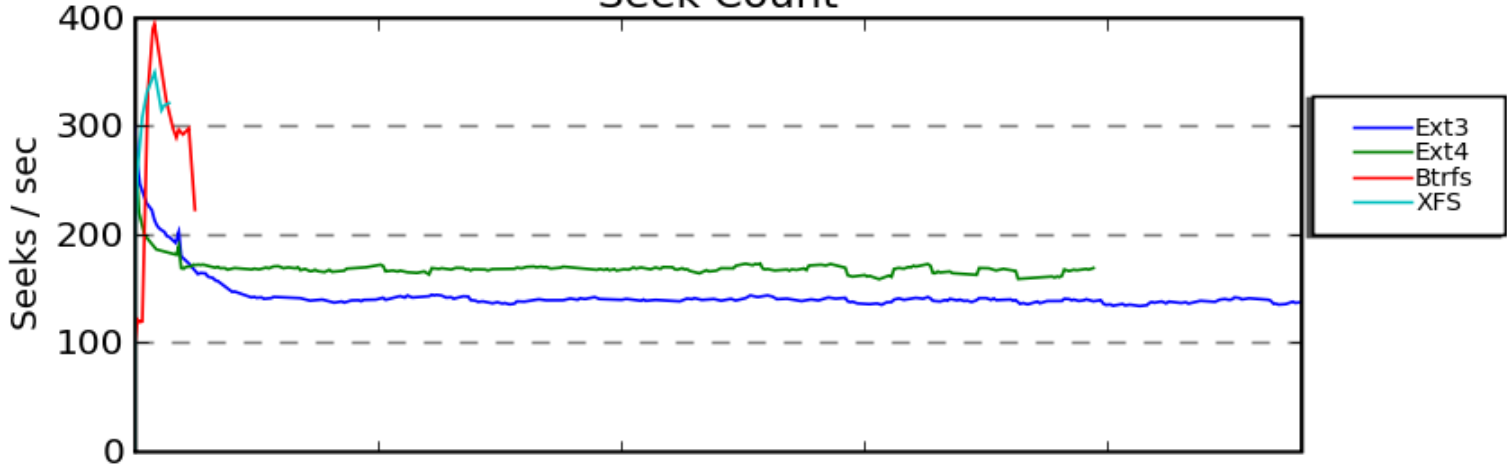


Kernel Tree Copy Movies

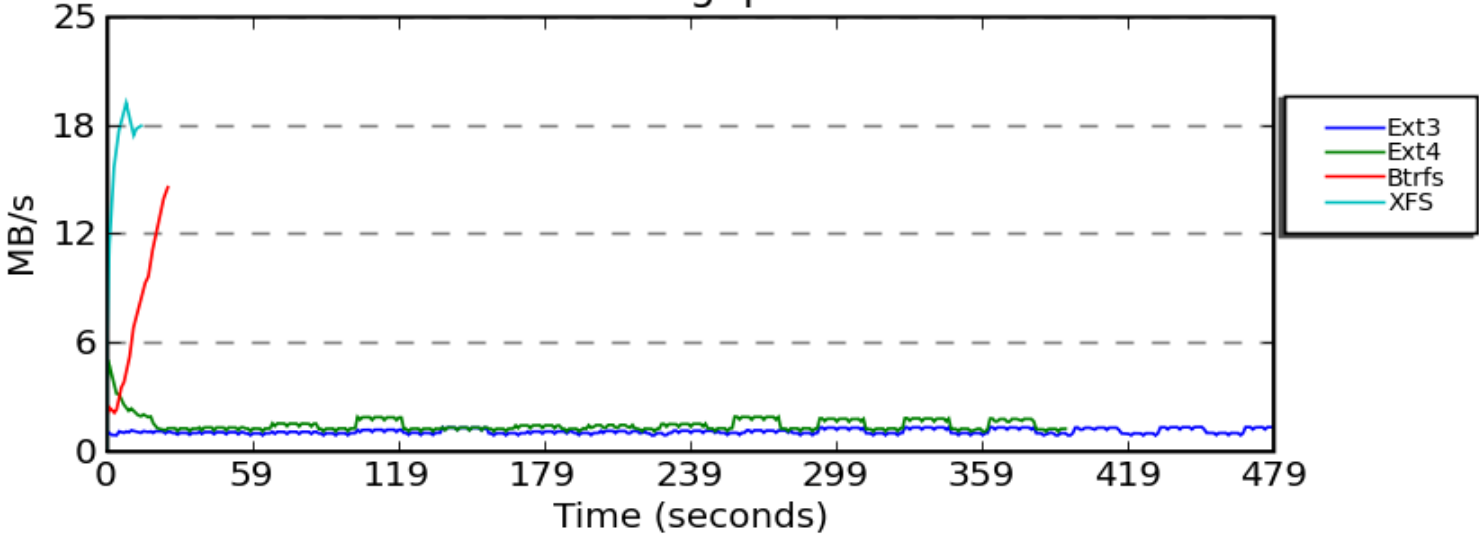
Directory Indexing

- Directory indexes provide better performance in large directories.
- They also determine the order the filesystem returns things to readdir
 - Used by backup programs to read all the files in the directory
- XFS indexes with a specialized btree
- Btrfs indexes once by name and once by sequence number
- Ext3 and Ext4 use htree
 - Very bad results for readdir on large directories

Read large dir Seek Count



Throughput

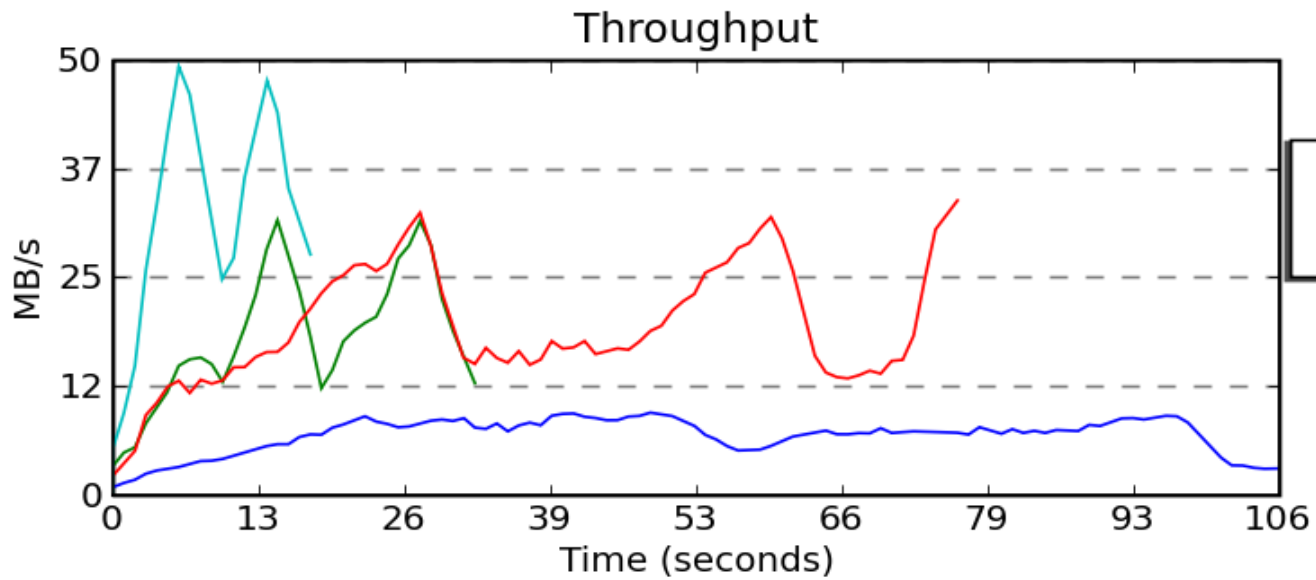
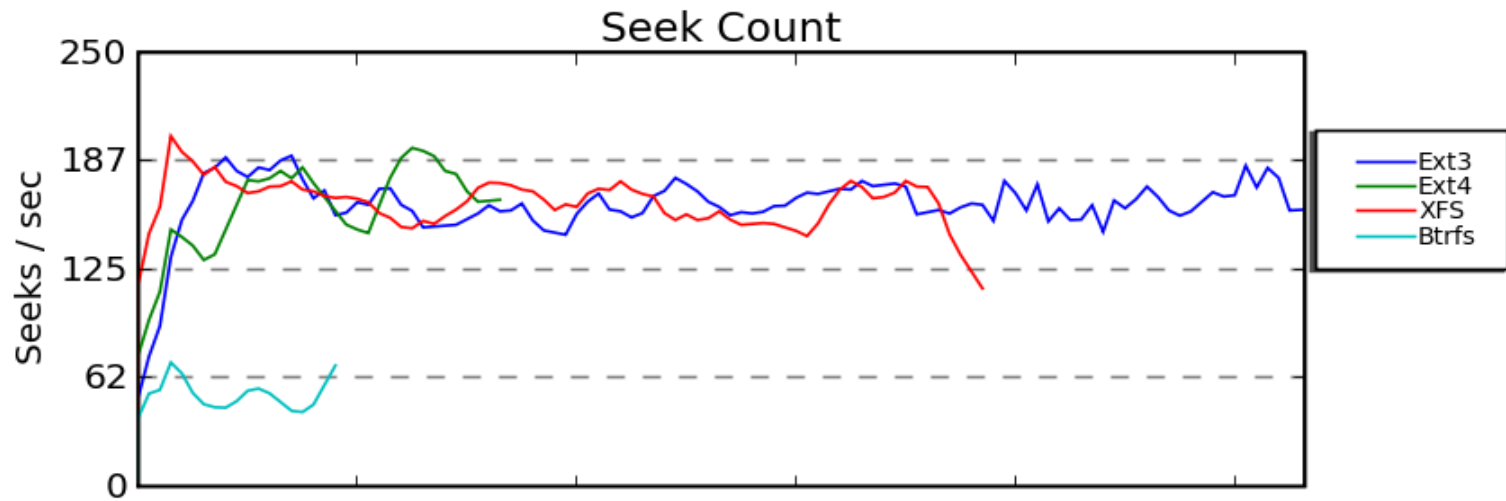




Directory Read Movies

Synchronous Writes

- Performance depends on logging subsystem, data writeback and data allocation policies
- Synchronous writes can have a big impact on other work being done by the filesystem
- Barriers are used by default on Ext4, Btrfs, and XFS to protect against corruption on power failures
- Btrfs uses a specialized log to improve synchronous performance





Synchronous write movies

Btrfs Goals

- General purpose filesystem that scales to very large storage
- Fast development cycle to attract developers and users early
- Feature focused, providing features other Linux filesystems cannot
- Administration focused, easy to run and very fault tolerant
- Perform well in a variety of workloads

Btrfs Features

- Extent based file storage
- Copy on write metadata and data
- Space efficient packing of small files
- Optional transparent compression (zlib)
- Integrity checksumming for data and metadata
- Writable snapshots
- Efficient incremental backups
- Multiple device support
- Offline conversion from Ext3

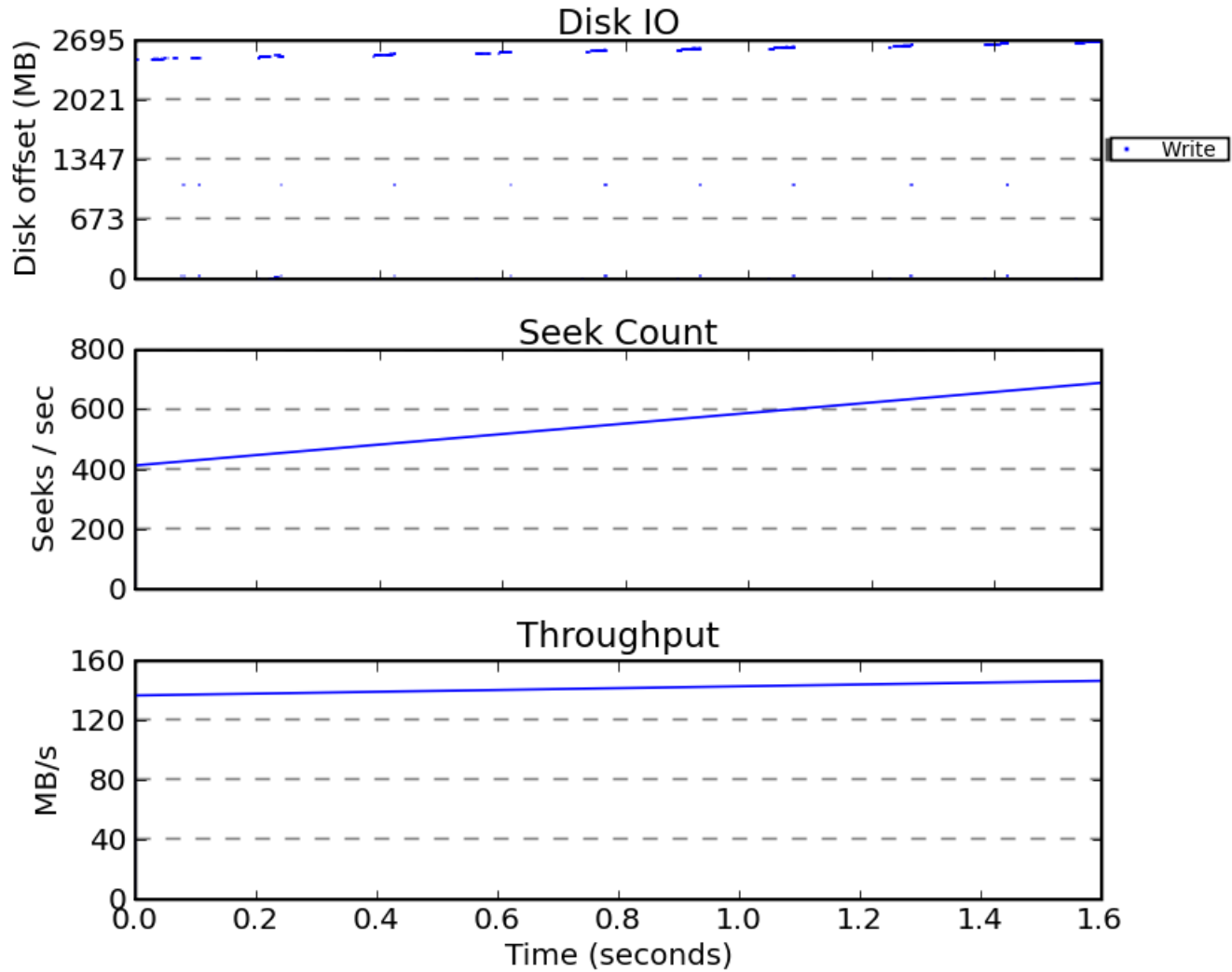
Btrfs Status

- Included in 2.6.29-rc1
- Generally usable in many workloads
- Generally stable

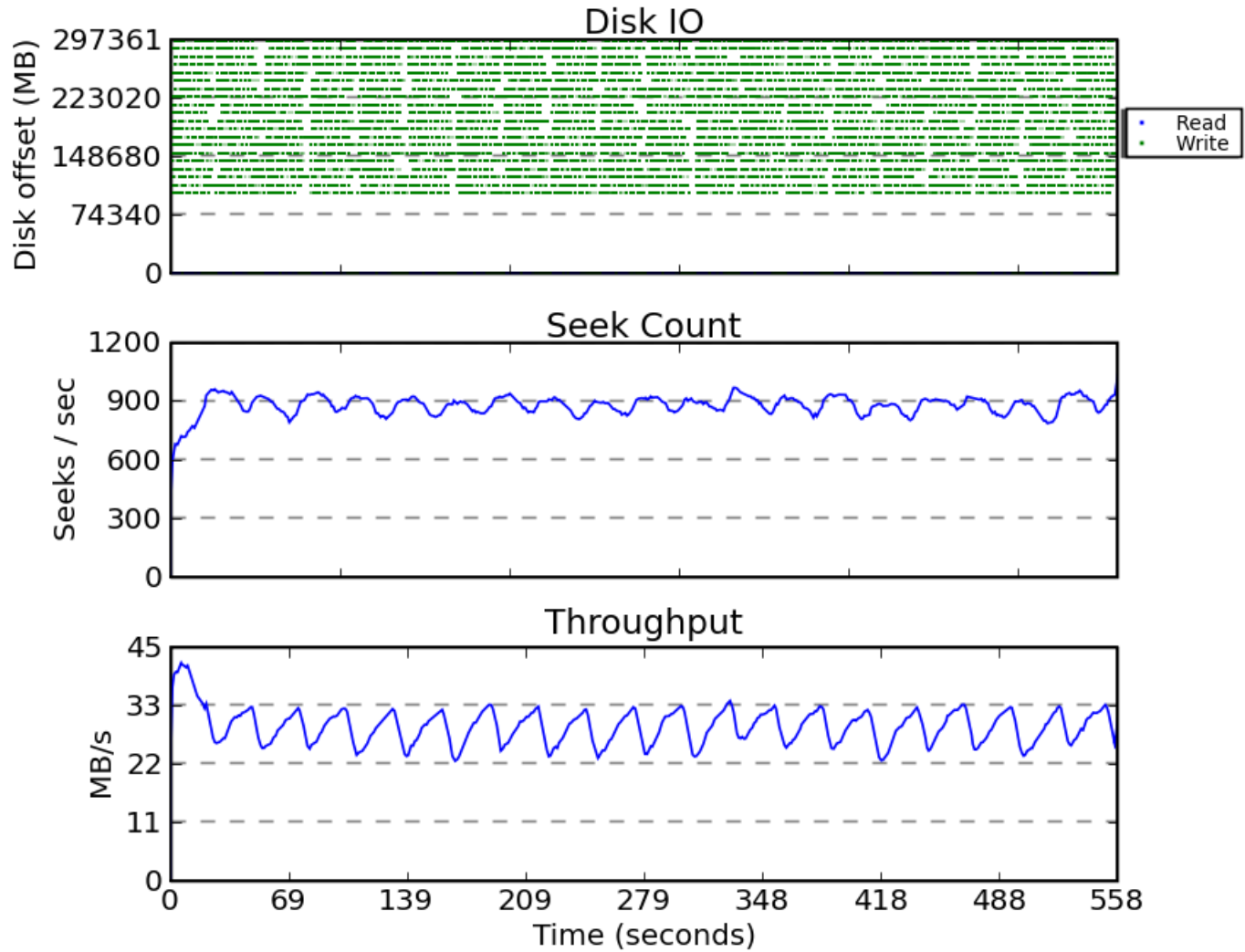
Snapshots and Subvolumes

- Subvolume is the unit of snapshotting
- As efficient as directories on disk
- Every transaction is a new unnamed snapshot
- Every commit schedules the old snapshot for deletion
- Snapshots use a reference counted COW algorithm to track blocks
- Snapshots can be written and snapshotted again
- COW leaves old data in place and writes the new data to a new location

Btrfs COW (20 snapshots 400M)



LVM Snapshot COW (20 snaps 400M)



Multi-device Support

- Devices are added into a pool of available storage
- New logical address space is allocated with a specific RAID configuration and data storage flags
 - System (used by the volume management code)
 - Metadata
 - Data
 - Raid0, raid1, raid10, single-spindle-dup
- Space is allocated from the storage pool in large chunks (1GB or more)
- Devices can be mixed in size and speed

Multiple Device Management

- Create the filesystem
 - `mkfs.btrfs /dev/sdb`
 - `Mount /dev/sdb /mnt`
- Add some devices
 - `Btrfs-vol -a /dev/sdc /mnt`
 - `Btrfs-vol -a /dev/sdd /mnt`
- Redistribute the data to the new devices
 - `Btrfs-vol -b /mnt`
 - This also restripes as required
- Remove a device
 - `Btrfs-vol -r /dev/sdc /mnt`

Online Resizing

- Shrink online by 1GB
 - `Btrfsctl -r -1g /mnt`
- Shrink a specific device
 - `Btrfs-show # find device ids`
 - `Btrfsctl -r 2:-1g /mnt`
 - 2: indicates to shrink devid 2
- Grow a device online
 - `Btrfsctl -r 2:+1g /mnt`
- Interfaces will be expanded to allow raid reconfiguration and advanced device management online

Synchronous Operations

- COW transaction subsystem is slow for frequent commits
 - Forces reCow of many blocks
 - Forces significant amounts of IO writing out extent allocation metadata
- Simple write ahead log added for synchronous operations on files or directories
- File or directory Items are copied into a dedicated tree
 - File back refs allow us to log file names without the directory
 - One log btree per subvolume

Synchronous Operations

- The log tree uses the same COW btree code as the rest of the FS
- The log tree uses the same writeback code as the rest of the FS, and uses the metadata raid policy.
- Commits of the log tree are separate from commits in the main transaction code.
 - fsync(file) only writes metadata for that one file
 - fsync(file) does not trigger writeback of any other data blocks

Mixed Fsync Workload

- 8 Processes creating, writing and fsyncing 20k files
- 1 Process creating 30 copies of the linux kernel sources in sequence

	Kernel tree I/O Rate	Fsync Creates
Btrfs	83MB/s	15,680
Ext4	122MB/s	1280
Ext3	66.85MB/s	3200

data=ordered

- Implemented to make sure that after a crash, files don't point to extents that have not been written
 - Prevents exposing stale data from deleted files
- Ext3 and Ext4
 - Write all dirty data before transaction commits
 - A single fsync will sync all dirty file non-delalloc data in the FS
- XFS and Btrfs
 - Update file metadata only after data is on disk

SSD

- Very fast for reads, no seek penalties
- Becoming very fast for writes, much smaller seek penalties
- Changing the design of storage subsystems, databases and filesystems