#### BetrFS: A Right-Optimized Write-Optimized File System

Amogh Akshintala, Michael Bender, Kanchan Chandnani, Pooja Deo, Martin Farach-Colton, William Jannen, **Rob Johnson**, Zardosht Kasheff, Bradley C. Kuszmaul, Prashant Pandey, Donald E. Porter, Leif Walsh, Jun Yuan, Yang Zhan

Facebook, Farmingdale College, MIT & Oracle, Rutgers, Stony Brook, Two Sigma, UNC, Williams College General-purpose file-systems strive to perform well on a wide variety of applications

- Sequential reads
- Sequential writes
- Random writes
- File/directory renames
- File deletes
- Recursive scans
- Metadata updates

Achieving good performance on all these operations is a long-standing challenge



- Sequential reads
- Sequential writes
- Random writes
  - File/directory renames
  - File deletes
- Recursive scans
- Output A construction of the second secon

#### Example: ext4



Achieving good performance on all these operations is a long-standing challenge

Sequential reads Sequential writes Random writes File/directory renames File deletes Recursive scans Metadata updates

Example: log-based file systems



Logging updates is fast, but logged data can have little locality

#### Some operations seem to require a trade-off



# Main idea of this talk Write optimited data structures

### BetrFS



#### Write-Optimized Data Structures (WODS)

- New class of data structure
  - LSM trees [O'Neil, Cheng, Gawlick, & O'Neil '96]
  - Bε-trees [Brodal & Fagerberg '03]
  - COLAS [Bender, Farach-Colton, Fineman, Fogel, Kuszna,
  - **xDicts** [Brodal, Demaine, Fineman, Iacono, Langerman, & Munro '10]
- WODS perform inserts/updates/deletes orders-ofmagnitude-faster than in a B-tree
  - WODS queries are asymptotically no slower than in a B-tree

ed in 90's

BetrFS uses B<sup>ε</sup>-trees

#### The Disk-Access Machine (DAM) model [Aggarwal & Vitter '88]

#### How computation works:

- Data is transferred in blocks between RAM and disk.
- The number of block transfers dominates the running time.
- Goal: Minimize # of block transfers
  - Performance bounds are parameterized by block size *B*, memory size *M*, data size *N*.



#### Example: B-trees



### B<sup> $\epsilon$ </sup>-trees ( $\epsilon = 1/2$ )



#### The search-insert asymmetry

• Inserts are orders-of-magnitude faster than point queries

Point query:  $O(\log_B N)$ Insert:  $O\left(\frac{\log_B N}{\sqrt{B}}\right)$ 

- But many updates require querying the old value first
- e.g. "Add \$10 to rob's account balance"
  - OldBalance = query(rob)
  - NewBalance = oldBalance + \$10
  - insert(newBalance)

#### Upserts: read-modify-write as fast as an insert



#### B<sup>ε</sup>-tree performance summary

Point query
$$O(\log_B N)$$
As fast as a B-treeInsert/delete/upsert $O\left(\frac{\log_B N}{\sqrt{B}}\right)$ Very fast (10K-100K per second)Range query $O\left(\log_B N + k/B\right)$ Near disk  
bandwidth

To get the best possible performance, we want to do Inserts, deletes, upserts, and range queries, and avoid point queries.

#### The BetrFS schema (version 0.1)

• Maintain two separate B<sup>ε</sup>-tree indexes:

metadata index: full path -> struct stat
data index: (full path, blk#) -> data[4096]

- Implications:
  - Fast directory scans
  - Data blocks are laid out sequentially

### Mapping file-system operations to key-value operations

**Operation** read write metadata update readdir mkdir/rmdir unlink rename



Small, random, unaligned writes are an order-of-magnitude faster



\*lower is better

•1 GiB file, random data

• 1,000 random 4-byte writes

•fsync() at end

BetrFS random writes benefit from B<sup>ɛ</sup>-tree insertion performance

#### Small file creates are an order-of-magnitude faster



• Create 3 million files and write

200-bytes to each

Balanced directory tree with

fanout 128

BetrFS file creates benefit from B<sup>E</sup>-tree insertion performance

#### Sequential I/O



- Write random data to file, 10
   4K-blocks at a time
- Sequentially read data back

#### Recursive directory traversals



 Recursive scans from root of Linux 3.11.10 source

- GNU find scans file metadata
- grep -r scans file contents

BetrFS directory traversals benefit from B<sup>E</sup>-tree range-query performance

Lower is better

#### File deletion



• Write random data to file,

fsync() it

• Delete file

BetrFS

BetrFS deletes require O(n) key-value operations

#### **Directory rename**



BetrFS renames require O(n) key-value operations

Lower is better

#### BetrFS (version 0.1) performance summary

- Sequential reads
- Sequential writes
  - Random writes
- File/directory renames
- Sile deletes
  - Recursive scans
- Metadata updates

Let's fix these problems

## Accelerating rename without slowing down directory traversals

#### Full-path indexing yields fast directory scans



#### Directory Tree (logical)

#### Example: grep -r "key" /home/rob/doc/

| disk<br>head |                           |  |  |  |  |  |
|--------------|---------------------------|--|--|--|--|--|
|              | /home/rob/doc             |  |  |  |  |  |
|              | /home/rob/doc/latex       |  |  |  |  |  |
|              | /home/rob/doc/latex/a.tex |  |  |  |  |  |
|              | /home/rob/doc/latex/b.tex |  |  |  |  |  |
|              | /home/rob/doc/bar.c       |  |  |  |  |  |
|              | /home/rob/local           |  |  |  |  |  |
|              |                           |  |  |  |  |  |
|              |                           |  |  |  |  |  |
|              |                           |  |  |  |  |  |
|              |                           |  |  |  |  |  |
|              |                           |  |  |  |  |  |

Disk (physical)

#### Rename is expensive when using full-path indexing



Directory Tree (logical)

Disk (physical)

#### Zoning: balancing indirection and locality



#### Zone: a subtree of the directory hierarchy



#### Moving the root of a zone is cheap



#### Renaming a subtree of a zone requires copying



#### Managing zone sizes



Large zones  $\rightarrow$  fast directory scans Small zones  $\rightarrow$  fast renames

We can keep zone sizes in a "sweet spot" by splitting large zones and merging small zones

#### How big should zones be?



BetrFS-0.2 uses 512KB zones to balance rename and scan performance

#### BetrFS 0.2 rename performance



Performing sequential writes at disk bandwidth and with full data journaling semantics

#### BetrFS version 0.1 writes everything twice



| 1 -        |            |    |          |              |    |          |   | > |
|------------|------------|----|----------|--------------|----|----------|---|---|
| I          | incort(k1  |    | <b>\</b> | incort (k2   |    | <b>`</b> |   | I |
| I          | Insert(KI, | VI | )        | LIISEI L(KZ, | VZ | )        |   | I |
| <u>د</u> _ |            |    |          |              |    |          | · | / |

#### BetrFS version 0.2: late-binding journal



#### Why don't we use late-binding for small writes?

- Reason 1:
  - Late-binding requires writing out a large (e.g. 4MB) node
  - For small writes, this is huge write-amplification
  - It's more efficient to make small writes durable by simply logging them
- Reason 2:
  - Small, random writes get written to disk several times as they get flushed down the B<sup>ε</sup>-tree
  - So writing them to the log is not a big extra cost

#### Late-binding journal: performance evaluation

#### Sequential write 120 -Fast sequential writes with full data journaling TOO (MB/sec 80 Throughput 60 40 20 0 BetrFS-0.1 BetrFS-0.2 ext4



#### Rangecast delete performance evaluation



# Is BetrFS still fast at other operations?





# What about real application performance?

#### Macrobenchmark: git



#### Macrobenchmark: dovecot imap maildir workload



#### BetrFS (version 0.2) performance summary

Sequential reads Sequential writes Random writes File/directory renames File deletes **Recursive scans** Metadata updates

#### Conclusion

- Write-optimized data structures can enable us to overcome long-standing file-system design trade-offs
- Write-optimized file systems can offer across-theboard top-of-the-line performance
- Write-optimization creates a need/opportunity to revisit many file system design issues

betrfs.org

•Code available at

#### SSD performance preview

