



File Systems Fated for Senescence? Nonsense, Says Science!

Ainesh Bakshi Martin Farach-Colton <u>Rutgers University</u> Yizheng Jiao Yang Zhan Donald E. Porter <u>UNC Chapel Hill</u>

Alex Conway Rutgers University Michael A. Bender William Jannen Rob Johnson <u>Stony Brook University</u> Jun Yuan <u>SUNY Farmingdale</u> Bradley C. Kuszmaul MIT and Oracle Corporation





File Systems Fated for Senescence? Nonsense, Says Science; The Essence of Semperjuvenescense is Coalescence!



Aging is fragmentation over time



Aging is fragmentation over time







Aging is fragmentation over time



SD @



SD @

Aging is fragmentation over time







- Do file systems age?
 - □ How can we reproduce aging?
 - How can we measure aging?
- □ How can we stop aging?











I'm Feeling Lucky

Chris Hoffman at <u>howtogeek.com</u> says:

"Linux's ext2, ext3, and ext4 file systems... [are] designed to avoid fragmentation in normal use."

"If you do have problems with fragmentation on Linux, you probably need a larger hard disk."







I'm Feeling Lucky

Chris Hoffman at <u>howtogeek.com</u> says:

"Linux's ext2, ext3, and ext4 file systems... [are] designed to avoid fragmentation in normal use."

"If you do have problems with fragmentation on Linux, you probably need a larger hard disk." "Modern Linux filesystems keep fragmentation at a minimum...Therefore it is not necessary to worry about fragmentation in a Linux system."





I'm Feeling Lucky

Chris Hoffman at <u>howtogeek.com</u> says:

2017 Storage Developer Conference. © Rutgers University

"Linux's ext2, ext3, and ext4 file systems... [are] designed to avoid fragmentation in normal use."

"If you do have problems with fragmentation on Linux, you probably need a larger hard disk." "Modern Linux filesystems keep fragmentation at a minimum...Therefore it is not necessary to worry about fragmentation in a Linux system." LINUX System Administrators' Guide Lars Wirzenius and Joanna Oja

Nope

Goog	e file system aging Q 	
Scholar	About 2,340,000 results (0.07 sec)	
Articles	File system aging—increasing the relevance of file system benchmarks KA Smith, MI Seltzer - ACM SIGMETRICS Performance Evaluation, 1997 - dl.acm.org	
Case law	Abstract Benchmarks are important because they provide a means for users and researchers to characterize how their workloads will perform on different systems and	
My library	different system architectures. The field of file system design is no different from other areas Cited by 131 Related articles All 15 versions Cite Save	



- Aging happens in real file systems
 - Smith and Seltzer ('97)
- Benchmarks should incorporate aging
 - □ Zhu, Chen and Chiueh ('05)
 - Agrawal, A. Arpaci-Dusseau and R. Arpaci-Dusseau ('09)

- Aging happens in real file systems
 - Smith and Seltzer ('97)

Yep

- Benchmarks should incorporate aging
 - □ Zhu, Chen and Chiueh ('05)
 - Agrawal, A. Arpaci-Dusseau and R. Arpaci-Dusseau ('09)









Nope



Let's do some science!





Inducing Aging

□ We use three different workloads: Developer workload Server workload

- Synthetic workloads



Inducing Aging

We use three different workloads:
 Developer workload
 Server workload
 Synthetic workloads







get coffee



22



get coffee git pull



23



get coffee git pull make







get coffee git pull make get coffee







get coffee git pull make get coffee git pull







SD @

get coffee git pull make get coffee git pull add awesome features





SD @

get coffee git pull make get coffee git pull add awesome features get coffee





SD @

get coffee git pull make get coffee git pull add awesome features get coffee git pull







get coffee git pull make get coffee git pull add awesome features get coffee git pull fix bugs

SD[©]





get coffee git pull make get coffee git pull add awesome features get coffee git pull fix bugs

. . .









SD (7



Use the Linux kernel repo from github.com







Measuring aging





Measuring aging

time grep -r random_string /path/to/filesystem








time grep -r random_string /path/to/filesystem









time grep -r random_string /path/to/filesystem









time grep -r random_string /path/to/filesystem









time grep -r random_string /path/to/filesystem









time grep -r random_string /path/to/filesystem







time grep -r random_string /path/to/filesystem







time grep -r random_string /path/to/filesystem







time grep -r random_string /path/to/filesystem







time grep -r random_string /path/to/filesystem







time grep -r random_string /path/to/filesystem



Then normalize per gigabyte read



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.

46

Do modern file systems age?













How can we be sure this is aging? Let's rule out some alternatives...





2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.

52



Rejuvenate the file system

Copy the logical state to newly formatted file system

- □ After every 100 git pulls
- Measure grep cost







Alternative: This only happens to ext4.







Git workload – HDD

Lower is better

Alternative: Aging only happens on hard drive.



2017 S

Jers University et al. All Rights Reserved.



Git workload – XFS SSD



Git workload – SSD

Lower is better



Git workload – SSD

Lower is better



On hard drive, sequential LBAs are written to nearby physical locations







On hard drive, sequential LBAs are written to nearby physical locations

Non-sequential LBAs may be written to physically separate locations



On SSDs, physical locality is not an issue.

But the hardware can serve small requests faster than the OS can make them!

"Improving File System Performance of Mobile Storage Systems Using a Decoupled Defragmenter" Hahn et al ATC 2017



(17)

On SSD, requests for logically sequential data are efficient

read(512, 16384)





On SSD, requests for logically sequential data are efficient

read(512, 16384)



Alternative: "You probably need a larger hard drive."





Partition size and aging

Experiments performed on a 20GB partition.
Repository size grows from 290MB to 1.2GB.
The partition was thus at most 6% full.

Aging is not a full disk issue!



Partition size and aging



Alternative: This only happens with a cold cache.





What about caching?

Write-caching, or deferred allocation, is generally limited to a few seconds in order to maintain reliability. Aging happens on the scale of days or more, so write-caching cannot help.





What about caching?

Read-caching, or disk cache, can help, but is limited by available memory. When the data set exceeds available memory, aging again becomes a problem.




Git workload – ext4 warm cache

Lower is better





Git workload – ext4 warm cache

SD @

Lower is better



Note that once the data set doesn't fit in memory, it quickly becomes slower than the unaged version with cold cache

- Btrfs, ext4, F2FS, XFS and ZFS all age
 Up to 22x on HDD
 Up to 2x on SSD
- Git lets us replay a real development history
 Induce aging by simulating years of use
 Takes between 5 hours and 2 days



Quick review

- Versions of the aging scripts used are available at betrfs.org
 - Test your own file system for aging
 Use them to perform benchmarks on aged file systems





How to prevent aging





Interfile fragmentation: blocks of the same file are not collocated





- Interfile fragmentation: blocks of the same file are not collocated
 - Avoid breaking large files into small fragments



- Interfile fragmentation: blocks of the same file are not collocated
 - Avoid breaking large files into small fragments
- Intrafile fragmentation: blocks of related files are not collocated





- Interfile fragmentation: blocks of the same file are not collocated
 - Avoid breaking large files into small fragments
- Intrafile fragmentation: blocks of related files are not collocated
 - Cluster logically related small files





- Interfile fragmentation: blocks of the same file are not collocated
 - Avoid breaking large files into small fragments
- Intrafile fragmentation: blocks of related files are not collocated
 - Cluster logically related small files

But what do we mean by small?







Natural transfer size

Higher is better



Natural transfer size

Higher is better



- Interfile fragmentation: blocks of the same file are not collocated
 - Avoid breaking large files into small fragments
- Intrafile fragmentation: blocks of related files are not collocated
 - Cluster logically related small files

Prediction: keeping related data in 4MiB chunks will substantially reduce aging





How to prevent aging (with Btrfs)





reduce How to prevent aging (with Btrfs)



88

How to reduce aging (with Btrfs)







Btrfs: Larger nodes => less aging?

Lower is better



Btrfs: Larger nodes => less aging?

Lower is better



Btrfs: Larger nodes => more writing?



Write Amplification = $\frac{|I/Os|}{|Workload|}$

For a B-tree, bigger leaves will mean more write-amplification





Btrfs: Larger nodes => more writing?



Btrfs: Larger nodes => more writing?



B-tree performance tradeoff







B-tree performance tradeoff



This tradeoff is inherent to B-trees



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.

How to prevent aging (with Betrfs)







In a B-tree, when changes are made to a leaf, the whole leaf must be read, modified and written out



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.



In a B-tree, when changes are made to a leaf, the whole leaf must be read, modified and written out



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.



In a B-tree, when changes are made to a leaf, the whole leaf must be read, modified and written out



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.



In a B-tree, when changes are made to a leaf, the whole leaf must be read, modified and written out



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.



In a B-tree, when changes are made to a leaf, the whole leaf must be read, modified and written out







In a B^ɛ-trees, intermediate nodes have buffers, aggregating changes en route to the leaves.



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.



In a B^ɛ-tree, intermediate nodes have buffers, aggregating changes en route to the leaves.



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.

There is no tradeoff!

- With a B^ε-tree, we only write a node when we've aggregated enough data.
- So we can have big leaves without writeamplification!





Git workload – Betrfs HDD





Git workload – Betrfs HDD



Our Setup: Cold Cache, 3.4 GHz Quad Core, 4GiB RAM, 20 GiB HDD partition - SATA 7200 RPM





Git workload – Betrfs HDD




Git workload – Betrfs SSD



110





Thank you!

Alex Conway betrfs.org Don't be afraid to rewrite your data Don't be afraid to rewrite your data



alexander.conway@rutgers.edu



2017 Storage Developer Conference. © Rutgers University et al. All Rights Reserved.

113