

MASARYK UNIVERSITY
FACULTY OF INFORMATICS



The effects of age on file system performance

BACHELOR'S THESIS

Samuel Petrovic

Brno, Spring 2017

Replace this page with a copy of the official signed thesis assignment and a copy of the Statement of an Author.

Declaration

Hereby I declare that this paper is my original authorial work, which I have worked out on my own. All sources, references, and literature used or excerpted during elaboration of this work are properly cited and listed in complete reference to the due source.

Samuel Petrovic

Advisor: Adam Rambousek

Acknowledgement

This is the acknowledgement for my thesis, which can span multiple paragraphs.

Abstract

This is the abstract of my thesis, which can span multiple paragraphs.

Keywords

filesystem, xfs, IO operation, aging, fragmentation ...

Contents

1	Introduction	1
2	Related work	3
2.1	<i>Aging file system using real-life data</i>	3
2.2	<i>Synthetic aging simulation</i>	4
3	File systems	5
3.1	<i>XFS</i>	6
3.2	<i>Ext4</i>	6
4	Used tools	9
4.1	<i>Beaker</i>	9
4.2	<i>FIO</i>	10
4.3	<i>Fs-drift</i>	10
4.4	<i>Storage generator</i>	11
5	Creating aged file system	15
5.1	<i>Aging process</i>	15
5.2	<i>File system images</i>	15
5.3	<i>Implementation</i>	16
6	Performance testing of aged file system	19
6.1	<i>Benchmark settings</i>	19
6.2	<i>Test structure</i>	19
7	Testing environment	21
7.1	<i>Used machines</i>	21
7.2	<i>Storage</i>	21
7.2.1	<i>HDD</i>	22
7.2.2	<i>SSD</i>	22
8	Results	23
8.1	<i>Performance of aged file system</i>	23
8.2	<i>Differences between XFS and EXT4</i>	23
8.3	<i>Differences across different storage</i>	23

List of Tables

List of Figures

- 4.1 Uniform distribution of file access 12
- 4.2 Normal distribution of file access 12
- 4.3 Moving random distribution 13

1 Introduction

File systems remain an important part of modern storage solutions. Large, growing databases, multi-media and other storage based applications need to be supported by high-performing infrastructure layer of storing and retrieving information. Such infrastructure have to be provided by operating systems (OS) in form of file system.

Originally, file system was a simple tool developed to handle communication between OS and physical device, but today, it is a very complex piece of software with large set of tools and features to go with. Performance testing is an integral part of development cycle of most of produced software. Because of growing complexity of file systems, performance testing took of as an important part of file system evaluation.

The standard workflow of performance testing of file systems is to run benchmark (e.g. testing tool) on a clean instance of OS and on a clean instance of tested file system [1]. Generally, this workflow present stable and meaningful results, yet, it only gives overall idea of file system behavior in early stage of its life cycle.

File systems, as well as other complex software is subjected to progressive degradation, referred to as software aging [2]. Causes of file system aging are many, but mostly fragmentation of free space, unclustered blocks of data and unreleased memory. This degradation cause problems in performance and functionality over time. Understanding of performance changes of aged file system can help developers to implement various preventions of aging related problems.

Testing of file system aging fundamentally consists of two steps. First is to bring fresh file system to an aged state and second is the actual performance test of the aged instance.

To achieve consistency of results and to shorten testing time, file system images are used in this thesis. Once the image of an aged file system is created, it can be stored for later use and reloaded back on device. To save space, only metadata of created file system is used, since content of created files is random and therefore irrelevant. Replayed metadata point at various blocks on device, recreating fragmentation while seldom taking significantly less space.

1. INTRODUCTION

Foremost, this thesis describe implementation of flexible tests which represent the two aforementioned steps. The first test is able to age fresh file system and store the result as an image for later use. The performance statistics collected in the process, as well as resulting layout can be used to evaluate ability of file system to respond to aging. Second test can evaluate resulting image even further by releasing some space and conducting performance test on resulting layout.

Furthermore, using developed tests to test different configurations of file systems and storage is demonstrated. The subject of research are differences between popular Linux file systems (XFS, Ext4) and storage technology (HDD, SSD) in context of aging. Because of nature of collected data, a processing tool was implemented to parse large amount of information into human readable reports. All the generated reports are part of this thesis in form of Appendix A.

The aging tests implemented in this thesis are established as part of testing cycle in Red Hat Kernel Performance team.

In the second chapter, the text present already conducted research of effects of age on file system performance. Third chapter describes used file systems and theirs main features. Chapter 4 introduce tools used in implementation of tests while describing their relevant features. Chapters 5 and 6 describe actual implementation of mentioned tests and remaining chapters present the results obtained by using created tests. In the conclusion chapter, I discuss the effects of age on file system performance as well as further options for file system aging performance testing.

2 Related work

In this chapter I present different approaches of file system aging described and implemented in the past. The first section discuss usage of collected data to create aging workload. The second section discuss possibilities of aging the file system artificially, without pre-collected data.

2.1 Aging file system using real-life data

This approach is based on modeling the aged file system using data collected from file system used in real-life environment.

Such data can be in form of snapshots (e.g. images) of file systems, as was thoroughly described by Smith and Seltzer [3].

The snapshots were collected nightly over a long period of time (one to three years) on more than fifty file systems. By comparing pairs of snapshots in the sequence, performed operations were estimated, resulting in a very realistic aging workload. However, as some studies suggest, most of files have life span shorter than 24 hours [4]. Therefore, as Smith and Seltzer admit, by snapshotting every night, this process does not account for most of the created files, resulting in loss of important part of data.

Furthermore, to age a file system sized 1024 megabytes, 87.3 gigabytes of data had to be written, taking 39 hours to complete, rendering the workflow impractical for in-production testing needs.

Smith and Seltzer also defined a layout score as a method to evaluate fragmentation of a file system. Layout score is defined as a fraction of blocks of file, which are contiguously allocated. Files of one block size are ignored, since they can't be fragmented, and for every file, first block is ignored too. Evaluation of the whole file system is then computed as an aggregated layout score of all files.

The problem resulting from not tracking shortly lived files can be solved by another approach called collecting traces. Traces are sequences of I/O operations performed by OS, captured at various levels (system call, drivers, network, etc.). The sequence of operations can be replayed back to the file system, aging it in a realistic manner.

2. RELATED WORK

Overall, using real-life data to age file systems brings realistic results, but at a cost of higher expenses, such as storing the collected data (meta data snapshot of 500GB Ext4 file system can have up to 500MB). Additionally, to cover cases of different types of file system usage, data from several such file systems have to be collected, expanding the amount of needed data even further. Such data is not always available, rendering this type of approach useful only in cases the researcher is already in possession of said data.

2.2 Synthetic aging simulation

Synthetic aging is a type of aging that does not require real-life data for its running. It relies on purely artificial workload performed on a file system, invoking aging factors, such as fragmentation.

Fast file system aging was described as a part of a trace replay benchmark called TBBT [5]. This type of aging consists of sequence of interleaving append operations on a set of files. By controlling the amount of files involved in the process, researchers had great control over fragmentation. Such workflow, while creating desired fragmentation, is however quite unrealistic, making the results of testing on such file system questionable [1].

The way to conduct synthetic aging, while keeping some amount of realism is to try to mimic the real-life usage by creating requests of random nature.

Aging workload generator such as fs-drift [6], used in this thesis, can be used (while carefully configured) to mimic long term real-life usage. Fs-drift simply creates a sequence of randomly chosen requests to be handled by the file system. The probability of the request type to be chosen is controlled by the workload table. In addition, whole process is highly configurable, making it possible to simulate various types of file system usage. Furthermore, a random distribution of file access can be controlled to mimic real-life user. All mentioned qualities, if used correctly, could result in fast, real-life mimicking file system aging.

3 File systems

In this chapter, I present basic information about file systems and describe main features of chosen file system in regard of fragmentation and scalability, which are important topics when discussing file system aging.

File system is a set of tools, methods, logic and structure to control how to store and retrieve data on and from device.

The system stores files either continuously or scattered across device. The basic accessed data unit is called a block, which capacity can be set to various sizes. Blocks are labeled as either free or used.

Files which are non-continuous are stored in form of extents, which is one or more blocks associated with the file, but stored elsewhere.

Information about how many blocks does a file occupy, as well as other information like date of creation, date of last access or access permissions is known as metadata, e.g. data about stored data. This information is stored separately from the content of files. On modern file systems, metadata are stored in objects called index nodes (e.g. inodes). Each file a file system manages is associated with an inode and every inode has its number in an inode table. On top of that the file system stores metadata about itself (unrelated to any specific file), such as information about bad sectors, free space or block availability in a structure called superblock.

In this thesis, targeted file systems are two most popular Linux file systems, XFS and Ext4 [Lu:2014:SLF:2578042.2560012], which are also main Red Hat supported file systems. These file systems belong to the group of file systems called journaling file systems.

Journaling file system keeps a structure called journal, which is a buffer of changes not yet committed to the file system. After system failure, these planned changes can be easily read from the journal, thus making the file system easily fully operational, and in correct and consistent state again.

3.1 XFS

XFS [**xfs:page**] is a 64-bit journaling file system known for its high scalability (up to 9 exabytes) and great performance. Such performance is reached by architecture based on allocation groups.

Allocation groups are equally sized linear regions within file system. Each allocation group manages its own inodes and free space, therefore increasing parallelism. Architecture of this design enables for significant scalability of bandwidth, threads, and size of file system, as well as files, simply because multiple processes and threads can access the file system simultaneously.

XFS allocates space as extents stored in pairs of B+ trees, each pair for each allocation group (improving performance especially when handling large files). One of the B+ trees is indexed by the length of the free extents, while the other is indexed by the starting block of the free extents. This dual indexing scheme allows efficient location of free extents for I/O operations.

Prevention of file system fragmentation consist mainly of a features called delayed allocation and online defragmentation.

Delayed allocation, also called allocate-on-flush is a feature that, when a file is written to the buffer cache, subtracts space from the free-space counter, but won't allocate the free-space bitmap. The data is held in memory until it have to be stored because of system call. This approach improves the chance, that the file will be written in a contiguous group of blocks, avoiding fragmentation and reducing CPU usage as well.

3.2 Ext4

Ext4, also called fourth extended filesystem is a 48-bit journaling file system developed as successor of Ext3 for Linux kernel, improving reliability and performance features. Ext4 is scalable up to 1 exbibyte (approx. 1.15 exabyte). Traditional Ext2 and Ext3 block mapping scheme was replaced by extent based approach similar to XFS, which positively affects performance.

Similarly to XFS, Ext4 use delayed allocation to increase performance and reduce fragmentation. For cases of fragmentation that still

occur, Ext4 provide support for online defragmentation and e4defrag tool to defragment either single file, or whole file system.

4 Used tools

In this chapter, I present tools which were used to implement automated tests for creating and storing aged file systems and measuring their performance. I will describe the main features and means of their usage. All the presented tools are open source projects.

4.1 Beaker

Beaker is a project aimed at automating testing workflow. The software can reserve machine from a pool of labs, install OS and packages, configure environment and perform tasks. The whole process is guided by sequence of instructions in an XML format. Examples of usage below.

Example 4.1: Specifying OS to be installed

```
1 <distroRequires>
2   <and>
3     <distro_family op="=" value="RedHatEnterpriseLinux7"/>
4     <distro_variant op="=" value="Server"/>
5     <distro_name op="=" value="RHEL-7.2"/>
6     <distro_arch op="=" value="x86_64"/>
7   </and>
8 </distroRequires>
```

Example 4.2: Configuring environment using kickstart

```
9 <kickstart>
10   <![CDATA[
11     install
12     lang en_US.UTF-8
13     skipx
14     keyboard us
15     rootpw redhat
16     firewall --disabled
17     authconfig --enableshadow --enablemd5
18     selinux --enforcing
19     timezone --utc Europe/Prague
20
21     bootloader --location=mbr --driveorder=sda
22     zerombr
23     clearpart --all --initlabel --drives=sda
24     part /boot --fstype=ext2 --size=200 --asprimary --label=BOOT --
       ondisk=sda
25     part /mnt/tests --fstype=ext4 --size=40960 --asprimary --label=MNT
       --ondisk=sda
26     part / --fstype=ext4 --size=1 --grow --asprimary --label=ROOT --
       ondisk=sda
```

4. USED TOOLS

```
27     reboot
28     %packages --excludedocs --ignoremissing --nobase
29     @core
30     wget
31     python
32     dhcpv6-client
33     dhclient
34     yum
35 ]]>
36 </kickstart>
```

Example 4.3: Executing task and passing arguments

```
37     <task name="/kernel_fsperf/storage_generator" role="STANDALONE">
38         <params>
39             <param name="TEST_PARAM_STORAGE_GENERATOR" value="-s create
              -f ext4 -t single -m /RHTSspareLUN1 -d /dev/sdc -T 1
              SASHDD_ext4"/>
40         </params>
41     </task>
```

4.2 FIO

Flexible Input/Output tool is a IO workload generator written by Jens Axboe. It is a tool well known for it's flexibility as well as large group of users and contributors. The flexibility is integral for conducting less artificial and more natural performance tests. However, approaching more natural test behavior, stability of results drop, so ideal equilibrium between these two parameters has to be found.

FIO accept

4.3 Fs-drift

fs-drift is a very flexible aging test, that can be used to simulate lots of different workloads. The test is based on random file access and randomly generated mix of requests. These requests can be writes, reads, creates, appends, truncates or deletes.

At the beginning of run time, the top directory is empty, and therefore *create* requests success the most, other requests, such as *read* or *delete*, will fail because not many files has yet been created. Over time, as the file system grows, *create* requests began to fail and other requests will more likely succede. File system will eventually reach a

state of equilibrium, when requests are equally likely to execute. From this point, the file system would not grow anymore, and the test runs unless one of the *STOP* conditions are met (specified with parameters).

The file to perform a request on is randomly chosen from the list of indexes. If the type of random distribution is set to *uniform*, all indexes have the same probability to be chosen, see 4.1. However, if the type of random distribution is set to *gaussian*, the probability will behave according to normal distribution with the center at index 0 and width controlled by parameter *gaussian-stddev*. This is useful for performing cache-tiering tests. Please note, that file index is computed as modulo maximal number of files, therefore instead of accessing negative index values, the test access indexes from the other side of spectrum, see Figure 4.2

Furthermore, fs-drift offers one more option to influence random distribution. After setting parameter *mean-velocity*, fs-drift will choose files by means of moving random distribution. The principle relies on a simulated time, which runs inside the test. For every tick of the simulated time, the center of bell curve will move on the file index array by the value specified using *mean-velocity* parameter. By enabling this feature, the process of testing moves closer to reality by simulating more natural patterns of file system access (the user won't access file system randomly, but rather works with some set of data at a time). On figure Figure 4.3, you can see bell curve moving by 5 units two times.

4.4 Storage generator

Storage generator is a beaker task developed by Jozef Mikovič. It is capable of automated configuration of storage on a machine. In a single-device mode, storage generator simply creates new partition on a given device and creates and mounts file system. In a recipe mode, storage generator follows set of bash instruction to create more advanced configurations such as merging multiple devices using LVM, creating LVM cache or encrypted volume.

Example 4.4: Configuring storage using storage generator in beaker environment

```
42 | <task name="/kernel_fsperf/storage_generator" role="STANDALONE">
```

4. USED TOOLS

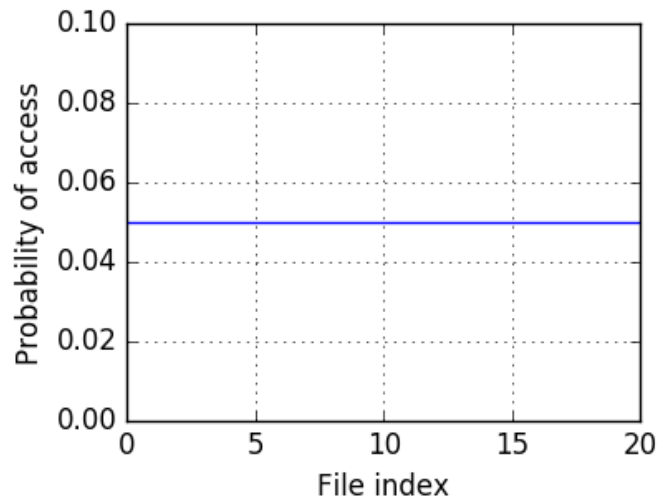


Figure 4.1: Uniform distribution of file access

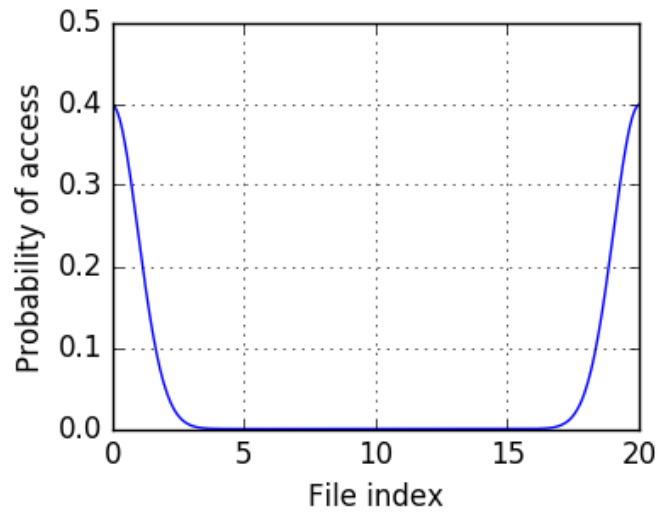


Figure 4.2: Normal distribution of file access

```
43     <params>
44         <param name="TEST_PARAM_STORAGE_GENERATOR" value="-s create
          -f xfs -t lvm -m /RHTSspareLUN1 -r jokerlvm -T 2
          SATASSDLVM_xfs"/>
45     </params>
46 </task>
```

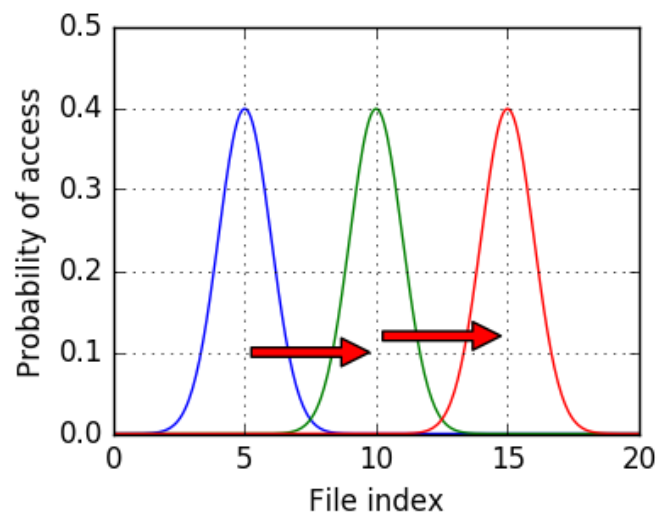


Figure 4.3: Moving random distribution

5 Creating aged file system

In this chapter, I describe possibilities of aging new file system, creating file system images and usage of implemented testing task.

5.1 Aging process

As the creator states in README, to fill up a file system, maximum number of files and mean size of file should be defined such that the product is greater than the available space. So if the workload is supposed to fill 500GB of space, while having maximum file size of 1GB (therefore mean size is 500MB), maximum number of files should be much higher than 1000. Optimal approach is to define seemingly no upper limit to let the fs-drift fill the volume, therefore numbers as high as 10^8 .

toto prepišat podľa 9year old bench study There is an option to specify user-defined file to use as a workload table, which is a desired percentual representation of operations in a workload. Since the goal of this workload is to create fragmented file system in a short time, read and rename operations are irrelevant. Therefore only create, append and delete have representation in this workload. The optimal results were reached when every operation had equal representation, e.g. 33%

Duration of the test is set to 5 hours so the test is usable for testing campaign without oversaturating of the servers.

5.2 File system images

File system images can be created by using tools developed to inspect file systems in case of emergency. For ext based file systems, there is e2image tool and for xfs, there is xfs_metadump. Both tools create images as sparse files, so compression is needed.

E2image tool can save whole ext based file system or just its meta-data and offers compression of image as well. Created images can be further compressed by tools such bzip2 or tar.

Example 5.1: Creating compressed image using e2image

```
$ e2image -Q $DEVICE $NAME.qcow2
```

5. CREATING AGED FILE SYSTEM

Such images can be later replayed back on a device. From that point, file system can be mounted and revised.

Example 5.2: Reloading compressed image

```
$ e2image -r $NAME.qcow2 $DEVICE
```

Xfs_metadump saves XFS file system metadata to a file. Due to privacy reasons file names are obfuscated (can be disabled by -o parameter). As well as e2image tool, the image file is sparse, but xfs_metadump doesn't offer a way to compress the output. However, output can be redirected to stdout from where it can be passed to a compression tool.

Example 5.3: Creating compressed image using xfs_metadump

```
$ xfs_metadump -o $DEVICE -l bzip2 > $NAME
```

Such images, when uncompressed can be replayed back on device by tool xfs_mdrestore. File system can be then mounded and inspected as needed:

Example 5.4: Reloading image using xfs_mdrestore

```
$ xfs_mdrestore $NAME $DEVICE
```

5.3 Implementation

Workflow of image creating is contained in the beaker task drift_job. After extracting fs-drift, the main script starts python script, which handles the process of running fs-drift. Settings of fs-drift are passed as a parameter and are parsed inside the script. Before running the fs-drift, python daemon thread is triggered to log free space fragmentation periodically while fs-drift is running. After the aging process is done, overall fragmentation is logged.

After the aging process, the script use presented system tools to create and compress the image. Information about system is gathered as well and all the logs are archived and sent to data collecting server. Parameters available for drift_job:

1. -s | --sync, flag to signalise wheather or not to send data to server (usefull for developing purposes)
2. -m | --mountpoint

3. -d | -device
4. -r | -recipe, parameters to pass to fs-drift
5. -t | -tag, string to distinguish different tests

6 Performance testing of aged file system

In this chapter, I describe structure of performance test which use images created by previous workflow. In the first section, I present settings of FIO benchmark for the optimal results and in section two, I describe implementation of this test as beaker task.

6.1 Benchmark settings

To ensure stability of test results, I decided to use simple form of standard performance test.

6.2 Test structure

Performance testing of created images is done by a package `recipe_fio_aging`. Upon installation of necessary tools (`libs`, `fio`), the package finds and downloads corresponding file system image according to obtained parameters. As shown, images are stored compressed, therefore decompression is needed after download. Once these steps are successfully completed, the image is restored on the device by using presented tools (`e2image`, `xfs_mdrestore`). If the image restoring completes successfully, file system can be mounted and worked with exactly like it would be just after the aging process.

After image restoration, some amount of the files needs to be deleted to create space for the performance test to take place. The files to be removed are chosen randomly until desired amount of volume has been freed. By using this workflow, e.g. freeing some amount of space, we can simulate aged file system in various phases of aging by using just one image of a very fragmented file system.

When free space is reclaimed, FIO test will take place using parameters given to `recipe_fio_aging`. The overall space occupied by the test should not be larger than available space on the file system, otherwise the test will either fail completely or report incorrect results.

For statistical correctness, the FIO test can run several times in a row. After last iteration, the results are archived and sent to data-collecting server.

6. PERFORMANCE TESTING OF AGED FILE SYSTEM

Parameters available for recipe_fio_aging:

1. -s | -sync, flag to signalise wheather or not to send data to server (usefull for developing purposes)
2. -n | -numjobs, number of test repetitions. For statistical stability
3. -m | -mountpoint
4. -d | -device
5. -r | -recipe, parameters to pass to FIO test
6. -t | -tag, string to distinguish different tests

7 Testing environment

In this chapter, I describe testing environment and storage used for testing with created tests.

7.1 Used machines

The aging process took place on these Machines:

1. Model: LENOVO System x3250 M6
2. CPU: Intel(R) Xeon(R) CPU E3-1230 v5 @ 3.40GHz (4 cores), arch i386 x86_64
3. Memory: 16384 MB
4. Storage:
 - (a) EG0600FBVFP HP Proliant HardDrive
 - (b) Interface: Serial Attached SCSI
 - (c) Capacity: 600 GB
1. Model: IBM x3650 M4
2. CPU: Intel(R) Xeon(R) CPU E5-2620 v2 @ 2.10GHz
3. Memory: 65536 MB
4. Storage:
 - (a) 2xSSDSC2BB480G4i IBM Solid State Drive
 - (b) Interface: Serial ATA
 - (c) Capacity: 480 GB

The system installed on machines is RHEL-7.2 with kernel 3.10.0-514.el7.x86_64

7.2 Storage

HDD is a rotational disk, which requires specific approach from kernel, to ensure the lowest possible seek time. Seek time is a time for moving parts of the device to find next relevant block of data. This affect overall performance greatly, because with large fragmentation, seek time becomes quite high.

As for SSD, this type of device does not have any moving parts, which make perform really well. One of the problems, however, is limited lifecycle of memory cells. SSD manufacturers deal with this

7. TESTING ENVIRONMENT

problem by adding controller with its own scheduler, which make sure, no parts of the device are used significantly more than other parts.

When aging the filesystems, I expect for those grown on HDD to perform significantly slower after aging process, and I expect SSD filesystems not to be affected at all, or maybe significantly less.

7.2.1 HDD

7.2.2 SSD

8 Results

The output of result generator is a html report summarising all information about system, links to raw data and charts of measured values.

8.1 Performance of aged file system

8.2 Differences between XFS and EXT4

8.3 Differences accross different storage

9 Conclusion

Here I will admit, that these results were not really surprising and ABSOLUTELY no breakthrough, however, as noone really research this branch of QE, the results are definitely a step further in this field.

Bibliography

- [1] Avishay Traeger et al. “A Nine Year Study of File System and Storage Benchmarking”. In: *Trans. Storage* 4.2 (May 2008), 5:1–5:56. ISSN: 1553-3077. DOI: 10.1145/1367829.1367831. URL: <http://doi.acm.org/10.1145/1367829.1367831>.
- [2] Domenico Cotroneo et al. “A Survey of Software Aging and Rejuvenation Studies”. In: *J. Emerg. Technol. Comput. Syst.* 10.1 (Jan. 2014), 8:1–8:34. ISSN: 1550-4832. DOI: 10.1145/2539117. URL: <http://doi.acm.org/10.1145/2539117>.
- [3] Keith A. Smith and Margo I. Seltzer. “File System Aging&Mdash;Increasing the Relevance of File System Benchmarks”. In: *SIGMETRICS Perform. Eval. Rev.* 25.1 (June 1997), pp. 203–213. ISSN: 0163-5999. DOI: 10.1145/258623.258689. URL: <http://doi.acm.org/10.1145/258623.258689>.
- [4] John K. Ousterhout et al. “A Trace-driven Analysis of the UNIX 4.2 BSD File System”. In: *Proceedings of the Tenth ACM Symposium on Operating Systems Principles*. SOSP ’85. Orcas Island, Washington, USA: ACM, 1985, pp. 15–24. ISBN: 0-89791-174-1. DOI: 10.1145/323647.323631. URL: <http://doi.acm.org/10.1145/323647.323631>.
- [5] Ningning Zhu, Jiawu Chen, and Tzi-Cker Chiueh. “TBBT: Scalable and Accurate Trace Replay for File Server Evaluation”. In: *Proceedings of the 4th Conference on USENIX Conference on File and Storage Technologies - Volume 4*. FAST’05. San Francisco, CA: USENIX Association, 2005, pp. 24–24. URL: <http://dl.acm.org/citation.cfm?id=1251028.1251052>.
- [6] B. England. *fs-drift*. <https://github.com/parallel-fs-utils/fs-drift>. 2017.