

MASARYK UNIVERSITY
FACULTY OF INFORMATICS



The effect of aging on filesystems performance

BACHELOR'S THESIS

Samuel Petrovic

Brno, Spring 2016

Replace this page with a copy of the official signed thesis assignment and the 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	Filesystems and used tools	3
2.1	<i>Filesystems</i>	3
2.2	<i>XFS</i>	3
2.3	<i>EXT4</i>	4
2.4	<i>FIO</i>	4
2.5	<i>Fs-drift</i>	4
2.6	<i>Snapshots</i>	5
3	Implementation	7
3.1	<i>filesystem_ager</i>	7
3.2	<i>recipe_fio</i>	8
3.3	<i>Inspecting filesystem</i>	10
3.4	<i>Aging recipes</i>	10
3.5	<i>Data processing</i>	11
4	Testing environment	13
5	Results	15
5.1	<i>Comparing against fresh filesystems</i>	15
5.2	<i>Comparing performance of Ext4 and XFS</i>	15
5.3	<i>Comparing overall state of aged EXT4 and XFS</i>	15
6	Conclusion	17

List of Tables

List of Figures

1 Introduction

This thesis will debate the effect of aging and fragmentation on filesystem performance, specifically ext4 and xfs.

Performance of Red Hat supported filesystems is closely watched by our team, but the effect of aging is not included in the testing matrix and is quite lightly understood in general.

Therefore, this thesis will try to shed some light as what the effect could possibly be as well as to design somewhat conceptual way of exploring this field further.

Possible other chapters to explain terms: journaling filesystem, allocation groups, B+ trees

2 Filesystems and used tools

2.1 Filesystems

Filesystem is a system to do mojo with files. Note to self: don't forget to delete this line and fill it with some actual text later.

2.2 XFS

XFS is a 64-bit journaling file system created by Silicon Graphics, Inc(SGI) in 1993. It is known for excellence in execution of paralel I/O operations, because of its architecture based on allocation groups.

Allocation groups are euqally 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 filesystem, as well as its 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 (imrpoving 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 for the highly efficient location of free extents for file system operations.

Extent describes one or more contiguous blocks, which can considerably shorten list of blocks.

Metadata journaling ensures consistency of data in case of emergency situations as f.e. system crash.

Prevention of file system fragmentation consist mainly of *delayed allocation* feature as well as online defragmentation(*xfs_fsr*), that can turururu

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, instead, until it must be flushed (to storage) because of memory pressure when calling Unix *sync*, or when flushing dirty buffers. 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.

Pridat nieco na zaver o XFS?

2.3 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.

Traditionally, ext* systems use an indirect block mapping scheme. Such an approach is generally inefficient for large files, on operations like deleting or truncating. Ext4 use modern approach of *extents*, which positively affect performance and encourage continuous layouts.

When allocating, ext4 use multiblock allocation, which is more efficient way than one block allocation at time, which is present in earlier ext* file systems. Multiblock allocation has far better performance, particularly when in use with delayed allocation and extents.

Similar as xfs, ext4 use delayed allocation design, to increase performance, especially when in use with multiblock allocation and extent-based approach, also reducing fragmentation on the device. For cases of fragmentation that still occur, ext4 provide support for online defragmentation and *e4defrag* tool to defragment either single file, or whole filesystem.

2.4 FIO

2.5 Fs-drift

fs-drift is a workload aging test written by Ben England. It relies on randomly mixed requests generated according to options. These requests can be writes, reads, creates, appends 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 of lack of files and small probability of randomly choosing existing one.

Over time, as the filesystem grows, *create* requests began to fail and other requests succede more. Finally, filesystem will reach a state

of equilibrium, when requests are equally likely to execute. From this point, the filesystem will not grow anymore, and the test will run until one of the *STOP* conditions are met (specified with parameters).

2.6 Snapshots

3 Implementation

3.1 filesystem_ager

This aging tool is a simple approach to write and remove many files of random size.

The tool consist of three scripts and one common library called *functions*. The scripts are named *filesystem_ager.py*, *fio_config_generator.py* and *random_deletor.py*.

The workflow consist of calling *filesystem_ager*, with desired parameters. Script manages triggering *fio_config_generator*, calling *fio* tool on generated config and triggering random deletor. These three actions are repeated given number of times. Parameters of *filesystem_ager* are:

1. Total desired size do be written in one cycle
2. Denominator of total desired size (Total desired size will be divided by this number)
3. Range of size of written files
4. Number of cycles

Although FIO tool has some parameters to randomize the size of files which are written, the management of file sizes and randomisation, as well as naming of files is handled by *fio_config_generator* instead, to provide more control over those qualities. Parameters of this script are:

1. Total deisred size to be written
2. Range of size of written files

The script will generate global settings of a workload, then proceeds to generate jobs for every file that will be written. File size is always the name of that file, and these are gathered to a list, then list of generated files is returned and script ends. Including file size in its name, as well as indexation of files will help effectively search and delete files in the random deletion process, without need to search for files on the disk and examine them for size. Simplistic approach in *fio* config will hopefully result in compatibility and reliability in use with any *fio* version.

After config file is generated, *filesystem_ager* will run *fio* tool on generated config and therefore, files are written on the device.

3. IMPLEMENTATION

The removing of files is handled by `random_deletor` script. Its parameters are:

1. Total written size
2. Denominator of total size
3. Range of size of written files
4. Number of existent files

If denominator equals zero, `random_deletor` won't remove any files and will return empty list. Otherwise, desired range of deletion is estimated. `Random_deletor` then proceeds to remove files while desired volume is not deleted. Files are randomly selected through choosing random integer from zero to number of existent files. This step may seem inefficient, but with large amounts of generated files, the time to perform successful selection will not change dramatically. Selected file name is then parsed for size information, and if it fits into desired volume to remove, it is deleted, through `subprocess` command. Names of removed files are gathered in a list and returned.

Number of deleted files is subtracted from number of existent files. `filesystem_ager` then sums up deleted volume, log it as well as other information and triggers the cycle again.

Here is a sequence diagram to show the structure of `filesystem_ager`.

3.2 `recipe_fio`

Measuring a performance is done by a tool I developed, `recipe_fio`. Similar to `filesystem_ager`, `recipe_fio` uses `fio` tool to handle desired IO operations, but instead of focusing on filling the filesystem, the script uses measurement features of `fio`, which consist of performing IO operations and reporting results.

*talk about used recipe

For purposes of this thesis, I let `recipe_fio` to report bandwidth and operations per second (IOPS).

The main script receives slightly enhanced `fio` configuration file, enriched of some non-`fio` parameters, which are used by test only. These parameters are:

1. used filesystem
2. number of test repetitions. For statistical stability, I decided to run the test several times under same conditions.

3. specifying a snapshot to be tested
4. flag which represents whether or not to rsync data on a result-storing server

After compiling, tool parses the parameters and gathers information about system, which consist of: version of kernel, time and date, hostname, RHEL compose, memory, kernel, mount, system info, system variables and fio version.

Then it proceeds to set environment for testing by:

1. Installing fio tool
2. Creating directory for results
3. Loading given snapshot on a device
4. Randomly removing volume from the device to make room for test

Volume is randomly removed using `random_delete_volume.py` script. This scriptt globs all files in the filesystem, retrieves information about used volume as well as it's overall size. Then proceeds to randomly choosing files to delete and stops when desired volume is freed. The approach of recursively globbing all files may be inefficient, but this way, we can be sure, that volume is deleted from whole device evenly.

Python script `run_tests.py`, manages to parse recipe parameter, resulting in creating one or several fio configuration files in the directory, further adding logging parameters to them. Then for every created file, directory on the desired medium is created and fio tool is triggered. If the script succesfully ends, file OK is created in the results directory.

When the measurement is over, bash script generate the name of results, which consist of:

1. time
2. date
3. used snapshot
4. version of kernel
5. version of RHEL compose

Then proceeds to compress results into tar archive with generated name, and according to `-g` flag will, or will not, rsync the result onto the data server.

3.3 Inspecting filesystem

For determining some overall idea about an extent to which is the filesystem aged or dirty, I wrote scripts that generate histograms representing fragmentation of used space as well as fragmentation of free space. Both scripts use common linux tools and pyplot to generate the graphics. Both scripts can display linear or logarithmic Y scale.

Script `extent_distribution.py` makes use of `xfs_io` `fiemap` tool, which is a tool to display extent distribution of a given file and works correctly even for `ext*` filesystems.

The script will first recursively crawls the whole filesystem from given top folder and makes a list of all files. `Fiemap` is then run over every file separately.

The only data, that are then parsed from the output, is how many non-contiguous extents does the file have. These integers are aggregated to a single list, from which are then counted, and final histogram is made.

Script `free_space_fragmentation.py` use the tool `e2freefrag`, which runs over a device, and outputs the histogram of free space fragmentation in textual form. Script will store this output and then easily parse the histogram and aggregate the data into a graphic form.

3.4 Aging recipes

To determine which `fs-drift` settings will be the most fitting for purposes of this thesis, I wrote a small python script `fs-drift_matrix`.

It is capable of taking matrix of possible `fs-drift` parameters from `json` file and then run it on a device.

After every run, histograms-generating scripts are triggered to store histograms of *free space* and *used space* fragmentation. Also outputs of the `fs-drift` script and `df` command are logged.

As the creator states in README, to fill up a filesystem, maximum number of files and mean size of file should be defined such that the product is greater than the available space.

For the purpose of this thesis, desired usage is 60%-100% with enough fragmentation to consider the device aged.

3.5 Data processing

4 Testing environment

The aging process took place on an ibm3250 machine with following parameters:

1. Intel(R) Xeon(R) CPU, X2460 2.80Ghz, 8 cores
2. RAM 10GB
3. 4x300GB SAS disks + 1x50GB Sas disk

The system installed on machine was — with kernel 3.10.0-229.el7.x86_64

I created a two volume groups, G1 and G2. Each group have a pair of 300GB disks with striping of 2. This setting allows to double the speed of IO operations.

On each volume group i created logical volume of 100GB.

5 Results

The output of result generator is a html report summarising all information about system, links to raw data and charts of measured values.
*talk about highcharts and how do you represent data

5.1 Comparing against fresh filesystems

5.2 Comparing performance of Ext4 and XFS

5.3 Comparing overall state of aged EXT4 and XFS

6 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.