

RobinHood v4 Progress Report

LAD'23 - 6th of October, 2023

Yoann VALERI, yoann.valeri@cea.fr

Commissariat à l'énergie atomique et aux énergies alternatives - www.cea.fr

Focus on RobinHood



Providing efficient and easy to use means to replicate
and query any filesystem's metadata

Reminder of last year's news



- Tools in development:
 - **rbh-lustre-find**: query and filter data in a backend based on Lustre specific attributes
 - **rbh-fsevents**: update a backend based on events provided by a source
 - **CREAT**, **TIME** events and **CLOSE**
 - **rbh-find**: new filters (size, permissions, extended attributes) and actions (listing, sorting, dumping to files)

What's new since last year ?

- **rbh-lustre-find** is near completion:
 - General tool and usage is done
 - Works like **rbh-find** -> supports its filters and additional ones for Lustre
 - Filters currently implemented: FID, HSM state and OST used
 - Usable by invoking **rbh-lfind**

Example of rbh-lfind usage



```
1 # content of `test_dir`
2 one [0x1:0x2:0x3] (archived) [0]
3 two [0x4:0x5:0x6] (released) [1,2]
4 three [0x7:0x8:0x9] (none) [2,3,4]
5
6 # rbh-sync rbh:lustre:test_dir rbh:mongo:test
7 # rbh-lfind rbh:mongo:test -hsm-state archived
8 /one
9
10 # rbh-lfind rbh:mongo:test -ost 1
11 /two
```

```
12 # rbh-lfind rbh:mongo:test -fid [0x7:0x8:0x9]
13 /three
14
15 # rbh-lfind rbh:mongo:test -hsm-state none -or -ost 0
16 /one
17 /three
18
19 # rbh-lfind rbh:mongo:test -ost 2 -and hsm-state released
20 /two
21
22 # rbh-lfind rbh:mongo:test -ost 4 -or \
23 ( -hsm-state archived -and -fid [0x4:0x5:0x6] )
24 /three
```

What's new since last year ? (2)

- **rbh-fsevents** is also near completion:
 - Lustre changelog reader/enricher done
 - All events are managed
 - Update to the database done
 - Deduplication work on-going, soon to be finished
- **rbh-fsevents [--enrich URI] [--lustre] <source> <destination>**
 - The **--lustre** is currently necessary if the source is a changelog reader, but the source will be changed in the future into a URI
 - Destination can either be a URI (for instance to update a Mongo database) or the standard output (with "-")

Example of rbh-fsevents usage



```
1      # touch test_archive
2      # lfs hsm_archive test_archive
3      # setfattr -n user.test -v 42 test_archive
4      # touch test_pfl
5      # lfs migrate -E 1k -c 2 -E -1 -c 1 test_pfl
6
7      # lfs changelog lustre-MDT0000
8      01 01CREAT ... t=[0x200000404:0x8e:0x0] j=touch.0 ... p=[0x200000404:0x8d:0x0] test_archive
9      02 11CLOSE ... t=[0x200000404:0x8e:0x0] j=touch.0 ...
10     03 16HSM    ... 0x280 t=[0x200000404:0x8e:0x0] ...
11     04 16HSM    ... 0x280 t=[0x200000404:0x8e:0x0] j=lhsmttool_posix.0 ... p=[0x200000404:0x8d:0x0]
12     05 16HSM    ... 0x0 t=[0x200000404:0x8e:0x0] j=lhsmttool_posix.0 ... p=[0x200000404:0x8d:0x0]
13     06 15XATTR  ... t=[0x200000404:0x8e:0x0] j=setfattr.0 ... x=user.test
14     07 01CREAT ... t=[0x200000404:0x94:0x0] j=touch.0 ... p=[0x200000404:0x8d:0x0] test_pfl
15     08 11CLOSE ... t=[0x200000404:0x94:0x0] j=touch.0 ...
16     09 12LYOUT  ... t=[0x200000404:0x94:0x0] ... p=[0x200000404:0x8f:0x0]
17     10 11CLOSE ... t=[0x200000404:0x94:0x0] j=lfs.0 ...
```


Example of rbh-fsevents usage (2)



```
18 # rbh-fsevents --enrich rbh:lustre:/mnt/lustre --lustre lustre-MDT0000 rbh:mongo:test
19 # mongo 'test' --eval 'db.entries.find({"ns.name":"test_archive"})'
20 '{ "ns" : [ {"name" : "test_archive", "xattrs" : { "path" : "/test_archive" } " ] },
21   "xattrs" : {
22     "fid" : BinData(0,"BAQAAAIAAACOAAAAAAAAAAAA=="),
23     "hsm_archive_id" : 1, "hsm_state" : 9,
24     "mdt_index" : 0, "ost" : [ 3 ],
25     "user" : { "test" : BinData(0,"NDI=") } } }'
26
27 # mongo 'test' --eval 'db.entries.find({"ns.name":"test_pfl"})'
28 '{ "ns" : [ {"name" : "test_pfl", "xattrs" : { "path" : "/test_pfl" } } ],
29   "xattrs" : {
30     "begin" : [ 0, 1048576 ], "comp_flags" : [ 16, 0 ],
31     "end" : [ 1048576, -1 ], "fid" : BinData(0,"BAQAAAIAAACUAAAAAAAAAAAA=="),
32     "flags" : 0, "gen" : 3, "mdt_index" : 0, "ost" : [ 2, 3, -1 ],
33     "pattern" : [ 0, 0 ], "pool" : [ "", "" ],
34     "stripe_count" : [ 2, 1 ], "stripe_size" : [ 1048576, 1048576 ] } }'
```


What's new since last year ? (3)



- Developments done for a new backend for the European project 
- Based on a component developed in it, the Hestia object store
- Includes:
 - A new backend for **librobinhood**
 - A new tool for query and filtering, **rbh-iosea-find**

What's new since last year ? (4)

- A new retention feature
- Allows one to set an expiration date on files by using an extended attribute
- Support:
 - An absolute expiration date
 - A relative expiration date, based on the maximum between the access and modify time of the file
 - An infinite retention date
- Tied to the Lustre backend currently, but will be changed in the future
- Can be queried by using the **expired** and **expired-at** filters of **rbh-lfind**
- The **expired-at** can also check for entries that expire before or after the given date

Example of the retention feature

```
1 # ls -color test_dir
2 fileA
3 fileB
4 fileC
5
6 # now=$(date +%s)
7 # setfattr -n user.expires -v $((now - 30)) test_dir/fileA
8 # setfattr -n user.expires -v inf test_dir/fileB
9 # setfattr -n user.expires -v +30 test_dir/fileC
10
11 # rbh-sync rbh:lustre:test_dir rbh:mongo:test
12 # rbh-lfind rbh:mongo:test -expired
13 /fileA
14
15 # rbh-lfind rbh:mongo:test -expired-at inf
16 /fileB
17
18 # rbh-lfind rbh:mongo:test -expired-at -$(date +%s --date='1 hour')
19 /fileA
20 /fileC
21
22 # rbh-lfind rbh:mongo:test -expired-at +$(date +%s --date='10 seconds')
23 /fileC
24
25 # sleep 120
26 # rbh-lfind rbh:mongo:test -expired
27 /fileA
28 /fileC
```

Some early benchmarks

- 2 million inodes synchronised
- CentOS Linux release 8.4.2105
- Lustre 2.12.9, Robinhood-lustre 3.1.7
- Same material

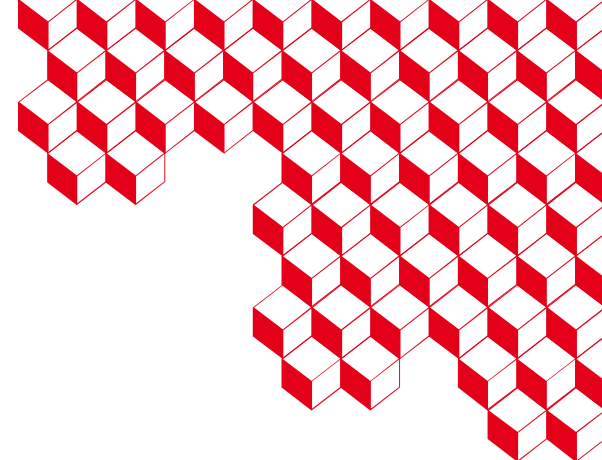
	Time (minutes)	Database size (GB)	Characteristics
Robinhood V3	24	1.58	Mono-process, 16 threads for entry process, 2 for FS scan
RobinHood V3	6	1.58	Parallelized on largest directory, fully optimized
Robinhood V4	60	0.395	Mono-process Mono-thread
RobinHood V4	27	0.395	parallelized on largest directory, mono-thread

What's next ?

- Finish the deduplication in **rbh-fsevents**
- Acknowledgement of events in **rbh-fsevents**
- Pre-production tests of the tools on our systems during 2023/2024
- The **rbh-report** tool to retrieve general information about a filesystem
- Improved synchronization by using **MPIFileUtils**
- More filters for the different **rbh-find** tools

Want to help with the development ?

- Projects have been moved on Github:
 - From one project per tool to one project for the suite
 - The address is: <https://github.com/robinhood-suite/robinhood4>
 - Reviews done on Gerrithub
- Some patches in review:
 - <https://review.gerrithub.io/c/robinhood-suite/robinhood4/+/557610>
 - First patch of the deduplication stack, documentation of the process
 - <https://review.gerrithub.io/c/robinhood-suite/robinhood4/+/1167917>
 - Fix a bug with **rbh-sync** branching function
- Feel free to install the suite and test it out for yourself, any feedback in appreciated!



Thanks for your attention

RobinHood v4 progress report

Yoann Valeri

Software stack developer at CEA

yoann.valeri@cea.fr

Commissariat à l'énergie atomique et aux énergies alternatives – www.cea.fr

Why a version 4 for RobinHood?



	Version 3	Version 4
Scale up	SQL paradigm <i>MariaDB</i>	NoSQL paradigm <i>MongoDB</i>
Code genericity	Software specialised for Lustre filesystems	Generic tools calling specific backends
Inclusion to Linux repositories	Expert system	Library of features and applications
Code refactoring	Heavy code caused by Lustre behaviour evolution	Clean design to better correspond to current filesystems