



DE LA RECHERCHE À L'INDUSTRIE

RobinHood v4 progress report

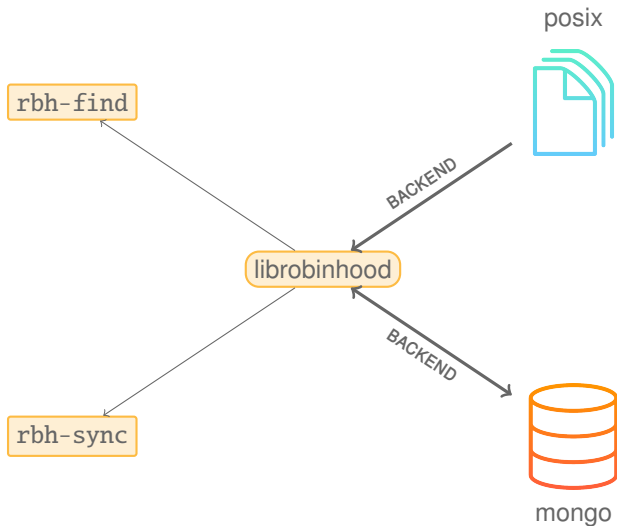
September 27th 2022

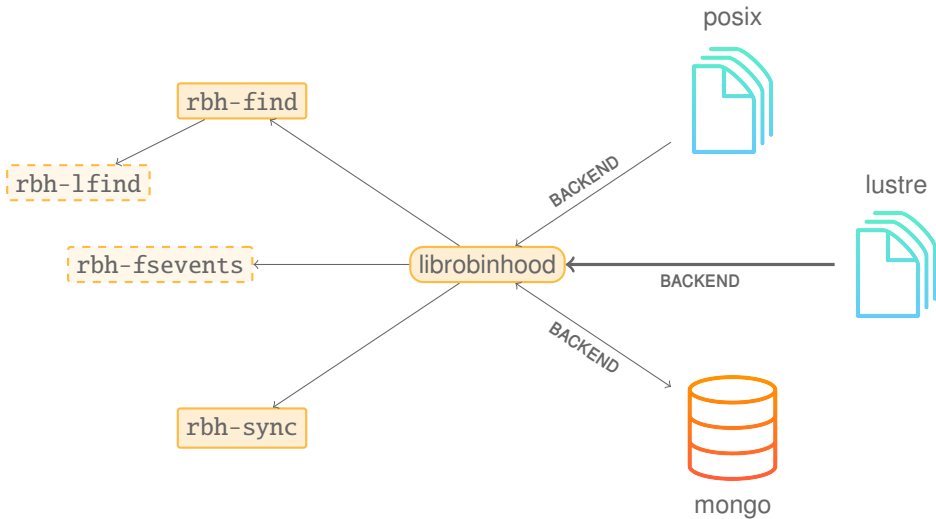
Sebastien Gougeaud

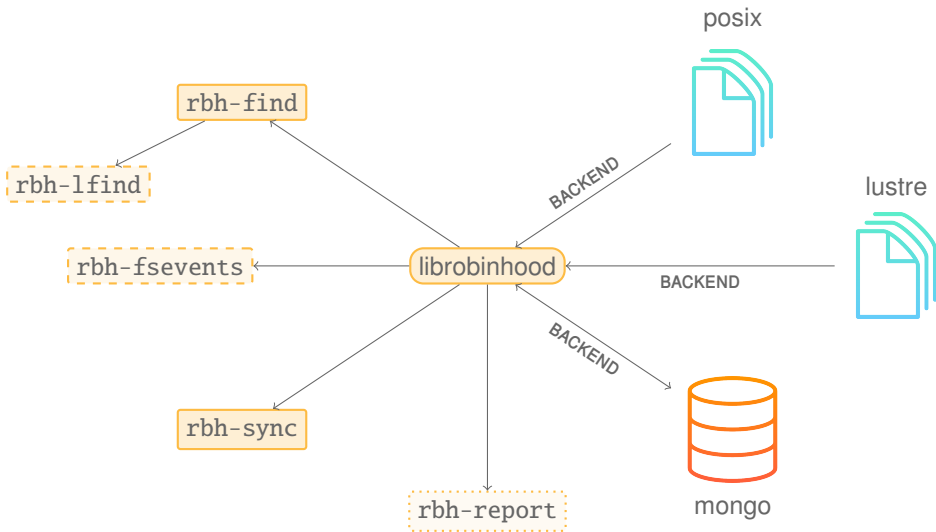
- ▶ The last update at LAD was two years ago
- ▶ RobinHood v4 is still in development by the CEA and former members of the team
- ▶ Developments can be followed on gerrithub and github (<https://github.com/cea-hpc/librobinhood>)
- ▶ Need for having it soon in production on our clusters

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

	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/applications
code refactoring	heavy code caused by Lustre behavior evolution	clean design to better correspond to current filesystems







- ▶ Support of custom extended attributes
 - Attach key-value attributes to entries
 - For instance: the path, the name, Lustre attributes

- ▶ Conversion between RBH FSEvents and YAML
 - FSEvent stands for File System Event: file creation, metadata change, etc.
 - Handling a new backend only need to convert specific log events into FSEvents

- ▶ Filter options are added as long as they may be needed to fetch metadata or define data placement policies:
 - Size --size
 - Permissions --perm
 - User-defined extended attributes --xattr
- ▶ Entries may be sorted using -sort/-rsort options
- ▶ Actions on entries are added:
 - Listing --ls
 - Printing to files --fprint/-fprint0/-fls

```
1 # rbh-sync rbh:posix:/mnt/lustre rbh:mongo:test
2 # rbh-find rbh:mongo:test -type f -sort name
3 /big
4 /huge
5 /small
6
7 # rbh-find rbh:mongo:test -type f -sort size
8 /small
9 /big
10 /huge
11
12 # rbh-find rbh:mongo:test -type f -sort size -ls
13 144...292      4 -rw-r--r-- 1 root root      17 Sep 13 14:00 /small
14 144...293     1024 -rw-r--r-- 1 root root     1048576 Sep 13 14:10 /big
15 144...294 284368 -rw-r--r-- 1 root root   291188736 Sep 13 14:11 /huge
16
17 # rbh-find rbh:mongo:test -type f -rsort size -ls
18 144...294 284368 -rw-r--r-- 1 root root   291188736 Sep 13 14:11 /huge
19 144...293     1024 -rw-r--r-- 1 root root     1048576 Sep 13 14:10 /big
20 144...292      4 -rw-r--r-- 1 root root      17 Sep 13 14:00 /small
21
22 # rbh-find rbh:mongo:test -type f -sort size -ls
23 144...292      4 -rw-r--r-- 1 root root      17 Sep 13 14:00 /small
24 144...293     1024 -rw-r--r-- 1 root root     1048576 Sep 13 14:10 /big
25 144...294 284368 -rw-r--r-- 1 root root   291188736 Sep 13 14:11 /huge
26
27 # rbh-find rbh:mongo:test -type f -sort name -ls
28 144...293     1024 -rw-r--r-- 1 root root     1048576 Sep 13 14:10 /big
29 144...294 284368 -rw-r--r-- 1 root root   291188736 Sep 13 14:11 /huge
30 144...292      4 -rw-r--r-- 1 root root      17 Sep 13 14:00 /small
```

- ▶ Lustre is the file system used in our clusters
- ▶ Storing lustre relevant information in the mirror backend, such as:
 - FID
 - File layout
 - OST
 - HSM state
- ▶ Those information are retrieved during the synchronization and converted to namespace attributes of the entry stored in the database backend
- ▶ Filtering entries using those information
- ▶ Overload of `rbh-find` tool to add lustre-related options

```

1 # mkdir /mnt/lustre/dir-{0..15}
2 # touch /mnt/lustre/dir-{0..15}/file-{0..15}
3 # for i in {0..5}; do
4 > lfs hsm_set --archived /mnt/lustre/dir- $\{i\}$ /file-12
5 > lfs hsm_set --norelease /mnt/lustre/dir- $\{i\}$ /file-11
6 > done
7 # rbh-sync rbh:lustre:/mnt/lustre rbh:mongo:test
8 # rbh-find rbh:mongo:test | head
9 /dir-14/file-13
10 /dir-14/file-14
11 /dir-14/file-15
12 /dir-0
13 /dir-15/file-0
14 /dir-1
15 /dir-15/file-1
16 /dir-2
17 /dir-15/file-2
18 /dir-3
19 # mongo --quiet --eval "db.entries.find({"statx.type": 32768})" | head -n 1
20 { "_id" : BinData(0,"lwA...AAA"), "statx" : { "blksize" : 4194304, "nlink" : 1, "uid" : 0, "gid" : 0, "type" : 32768,
21 "mode" : 420, "ino" : NumberLong("144115205272502451"), "size" : NumberLong(0), "blocks" : NumberLong(0),
22 "attributes" : { "immutable" : false, "append" : false }, "atime" : { "sec" : NumberLong(1663145146), "nsec" : 0 },
23 "btime" : { "sec" : NumberLong(1663145146), "nsec" : 0 }, "ctime" : { "sec" : NumberLong(1663145146), "nsec" : 0 },
24 "mtime" : { "sec" : NumberLong(1663145146), "nsec" : 0 }, "rdev" : { "major" : 0, "minor" : 0 }, "dev" : { "major" : 1273,
25 "minor" : 181606 } }, "xattrs" : { "lustre" : { "lov" : BinData(0,"0Av...AAA="), "trusted" : { "link" : BinData(0,"3/H...S0w"),
26 "lma" : BinData(0,"AAA...AAA"), "lov" : BinData(0,"0Av...AA="), "som" : BinData(0,"BAA...AAA") } },
27 "ns" : [ { "parent" : BinData(0,"lwA...AAA"), "name" : "file-0", "xattrs" : { "path" : "/dir-10/file-0",
28 "fid" : BinData(0,"AQ...A="), "hsm_state" : 0, "hsm_archive_id" : 0, "flags" : 0, "magic" : "LOV_USER_MAGIC_V1", "gen" : 0,
29 "stripe_count" : [ ], "stripe_size" : [ ], "pattern" : [ ], "comp_flags" : [ ], "pool" : [ ],
30 "ost" : [ NumberLong(0) ], "mdt_index" : 0 } } ] } ] }

```


```
1 # rbh-lfind rbh:mongo:test -sort name -ost 0 | head
2 /dir-15/file-0
3 /dir-0/file-0
4 /dir-1/file-0
5 /dir-2/file-0
6 /dir-3/file-0
7 /dir-4/file-0
8 /dir-5/file-0
9 /dir-6/file-0
10 /dir-7/file-0
11 /dir-8/file-0
12
13 # rbh-lfind rbh:mongo:test -sort name -ost 1 | head
14 /dir-15/file-1
15 /dir-0/file-1
16 /dir-1/file-1
17 /dir-2/file-1
18 /dir-3/file-1
19 /dir-4/file-1
20 /dir-5/file-1
21 /dir-6/file-1
22 /dir-7/file-1
23 /dir-8/file-1
24 # lfs_path2fid /mnt/lustre/dir-0/file-0
25 [0x200000401:0x13:0x0]
26 # rbh-lfind rbh:mongo:test -fid [0x200000401:0x13:0x0]
27 /dir-0/file-0
28
29 # rbh-lfind rbh:mongo:test -hsm-state archived
30 /dir-0/file-12
31 /dir-1/file-12
32 /dir-2/file-12
33 /dir-3/file-12
34 /dir-4/file-12
35 /dir-5/file-12
36 # rbh-lfind rbh:mongo:test -hsm-state norelease
37 /dir-0/file-11
38 /dir-1/file-11
39 /dir-2/file-11
40 /dir-3/file-11
41 /dir-4/file-11
42 /dir-5/file-11
```

- ▶ First events are processed:
 - CREAT
 - ATIME/CTIME/MTIME
 - CLOSE
- ▶ Others are following
- ▶ Still thinking about the right moment to acknowledge the events:
 - When the changelog events are converted to FSEvents?
 - When the events are recorded in the backend?

```

1 # touch /mnt/lustre/old
2 # lctl --device lustre-MDT0000 changelog_register
3 # touch -a /mnt/lustre/old
4 # touch /mnt/lustre/new
5 # lfs changelog lustre-MDT0000
6 1 18CTIME 10:45:21.574804672 2022.09.13 0x5
7 t=[0x200000401:0x12:0x0] j=touch.0 ef=0xf u=0:0 nid=0@lo
8 2 11CLOSE 10:45:21.588901555 2022.09.13 0x42
9 t=[0x200000401:0x12:0x0] j=touch.0 ef=0xf u=0:0 nid=0@lo
10 3 01CREAT 10:45:25.915058526 2022.09.13 0x0
11 t=[0x200000401:0x13:0x0] j=touch.0 ef=0xf u=0:0 nid=0@lo
12 p=[0x200000007:0x1:0x0] new
13 4 11CLOSE 10:45:25.922082667 2022.09.13 0x42
14 t=[0x200000401:0x13:0x0] j=touch.0 ef=0xf u=0:0 nid=0@lo
15 # rbh-fsevents --enrich /mnt/lustre lustre-MDT0000 -
16 --- lupsert
17 "id": !!binary lwAAAAEEAAACAAAEEgAAAAAAAAAAAAAAAAAAAA
18 "xattrs": {}
19 "statx":
20   "atime":
21     "sec": 1663065921
22     "nsec": 0
23   "ctime":
24     "sec": 1663065921
25     "nsec": 0
26 ---
27 --- lupsert
28 "id": !!binary lwAAAAEEAAACAAAEEgAAAAAAAAAAAAAAAAAAAA
29 "xattrs": {}
30 "statx":
31   "atime":
32     "sec": 1663065921
33     "nsec": 0
34 ---
35 --- lupsert
36 "id": !!binary lwAAAAEEAAACAAAEEwAAAAAAAAAAAAAAAAAAAA
37 "xattrs":
38   "ns":
39     - "parent": !!binary lwAAAAcAAAACAAAQAQAAAAAAAAAAAAAAAA
40     "name": "new"
41     "xattrs":
42       "path": "new"
43     "fid": !!binary AQQAAAIAAAAAAATAAAAAAAAAAAAA=
44 "statx":
45   "type": "file"
46   "uid": 0
47   "gid": 0
48   "atime":
49     "sec": 1663065925
50     "nsec": 0
51   [...]
52   "blksize": 4194304
53   "attributes":
54     "immutable": n
55     "append": n
56     "rdev":
57       "major": 0
58       "minor": 0
59     "dev":
60       "major": 1273
61       "minor": 181606
62 ---
63 --- lupsert
64 "id": !!binary lwAAAAEEAAACAAAEEwAAAAAAAAAAAAAAAAAAAA
65 "xattrs": {}
66 "statx":
67   "atime":
68     "sec": 1663065925
69     "nsec": 0
70 ---

```

- ▶ Development of a non-lustre related backend, for the  IO-SEA european project, using CORTX-Motr and Phobos as object store
- ▶ FSEvent deduplication
- ▶ Pre-production tests on our systems for a production during 2023

project	repository	
librobinhood	https://github.com/cea-hpc/ {	librobinhood
rbh-find		rbh-find
rbh-find-lustre		rbh-find-lustre
rbh-sync		rbh-sync
rbh-fsevents		rbh-fsevents



Thanks for your attention

RobinHood v4 progress report
Sebastien Gougeaud

Commissariat à l'énergie atomique et aux énergies alternatives
Centre de Bruyères-le-Châtel | 91297 Arpajon Cedex
T. +33 (0)1 69 26 40 00 | F. +33 (0)1 69 26 40 00
Établissement public à caractère industriel et commercial – RCS Paris B 775 685 019

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