

A Ceph plugin for XrootD

Sébastien Ponce
sebastien.ponce@cern.ch

CERN

January 28th 2015

- 1 Introducing Ceph
- 2 The Ceph plugin of XrootD
- 3 Practical usage and status

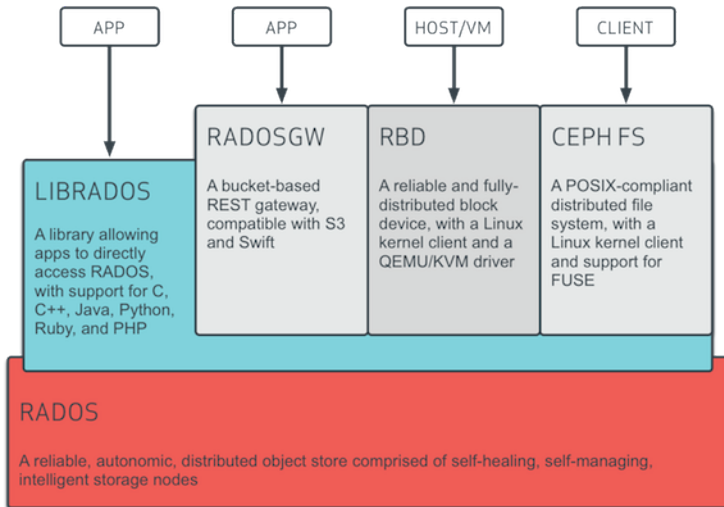
Introducing Ceph

- An Object is
 - a piece of data
 - metadata (extended attributes)
 - a unique identifier
- Different from a file system
 - Get/Put semantic
 - No file descriptor
 - Flat namespace
aka no namespace
 - “small” objects (few MBs)



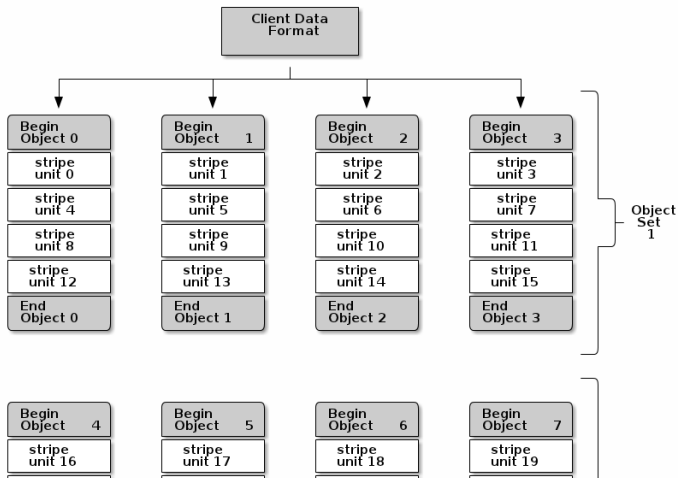
Controlled Replication Under Scalable Hashing
<http://ceph.com/papers/weil-crush-sc06.pdf>

- a clever data placement algorithm ...
 - supporting cluster maps (e.g. rows/cabinet/shells/devices)
 - approximating a uniform probability distribution
 - with completely deterministic mapping
 - but pseudo-random distribution
 - minimizing data movements on cluster evolution
- ... for a very scalable object store
 - no central catalog of object placement
 - placement computation done by clients
 - configurable replication (file by file)
 - erasure coding available



- librados
 - the natural object store interface
 - strong atomicity/snapshotting features
 - has object store limitations :
 - no namespace
 - objects should be small (logical I/O unit, order of MBs)
 - no parallel I/O within an object
- libradosstriper
 - adds striping on top of librados
 - reusing the striping of Ceph FS and RBD
 - available starting with giant release

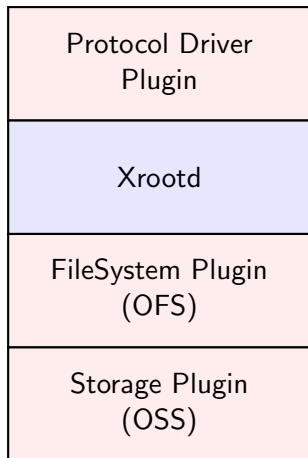
Defaults : 1 stripe, stripe unit = object size = 4 MB



The Ceph plugin of XrootD

OSS = Open Storage System

- XrdOSS / XrdOssDf interfaces
 - implementing entry points to the underlying storage
 - e.g. open, read, write, close
- compatible with any OFS or Protocol plugin
 - e.g. Castor, EOS, HTTP



PROs

- Taking full benefit of the distributed object store
- No file size limit
- Embedded striping and parallel I/O
 - parallel I/O is not activated by default (nb stripes = 1)

A small drawback

- No directory listing

Why to deal with external attributes ?

- Xrootd has no client interface for them
- OSS plugins do not support external attributes
- ... but an OFS plugin may need them !

How they are supported

- via the XrdSysXattr interfaces introduced in Xrootd 4.1
- implemented via an extra plugin

Practical usage and status

Config file syntax

```
ofs.osslib      libXrdCeph.so
ofs.xattrlib    libXrdCephXattr.so
all.export      *?
```

Usage

```
# xrdcp myfile root://myserver/myfile
[1000MB/1000MB][100\%][=====][17.24MB/s]
# rados ls | grep myfile
myfile.000000000000000000
myfile.000000000000000001
# xrdcp root://myserver/myfile myfile2
[1000MB/1000MB][100\%][=====][16.38MB/s]
```

Config file syntax

```
ofs.osslib libXrdCeph.so [[user@]pool]
```

Extended file syntax

```
[[user@]pool:]path
```

Examples

```
xrdcp root://myserver/mypool:myfile ...  
xrdcp root://myserver/myuser@mypool:myfile ...  
xrdcp root://myserver/:file_with_a_in_it ...
```

Complete config file syntax

```
ofs . osslib <lib> [[ user@ ] pool [, <layout > ]]  
  layout : nbStripes [ , stripeUnit [ , objSize ] ]
```

Complete file syntax

```
[[ user@ ] pool [, <layout > ] : ] path
```

Examples

```
myuser@mypool , 4 : myfile  
myuser@mypool , 4 , 65536 : myfile  
mypool , 1 , 33554432 , 33554432 : myfile
```


- code is ready and under testing
- available on github
(<http://github.com/sponce/xrootd>)
- to be integrated soon into the main trunk
- will be part of release 4.2