



The future of storage

Ceph-FS与大数据的恋爱史

目 · 录

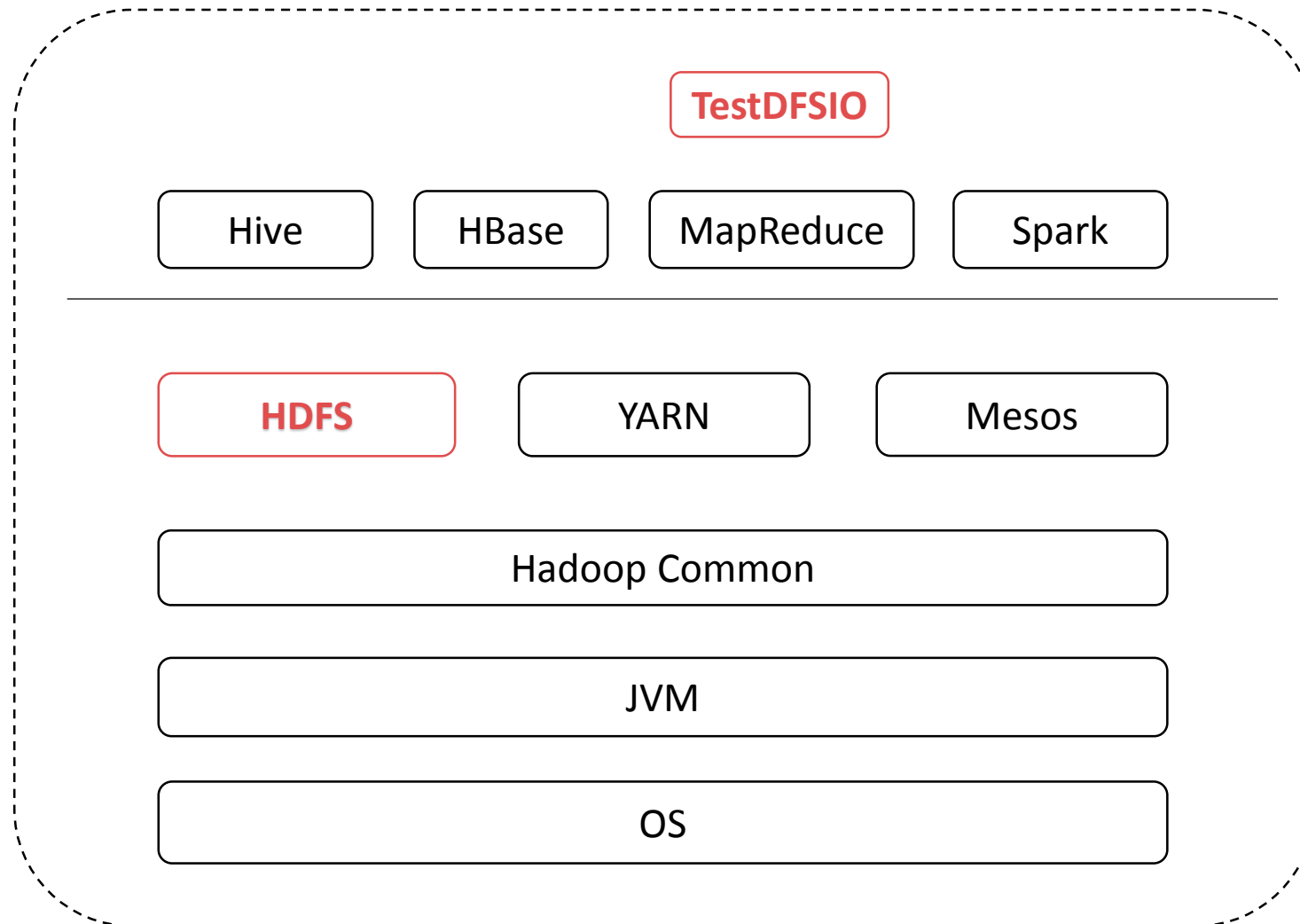
01 故事背景

02 第一次亲密接触

03 热恋期

故事·背景

Hadoop



典型大数据

文件
大

GB-TB

数量
较多

千万级

读多
写少

一次写多次读

典型大数据 & HDFS

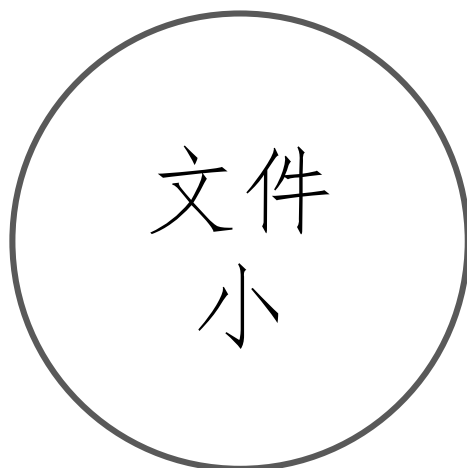
擅长处理大文件

支持千万量文件

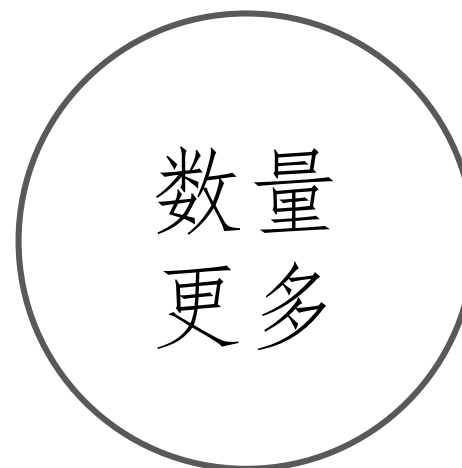
读性能 > 写性能

天生一对

非典型大数据



KB-MB



亿级

非典型大数据 & HDFS

小文件性能差

亿级文件搞不定

元数据内存消耗大

不合拍

HDFS Hold 不住，Ceph-FS 能 Hold 住吗？

Ceph-FS

小
文件

读写性能好

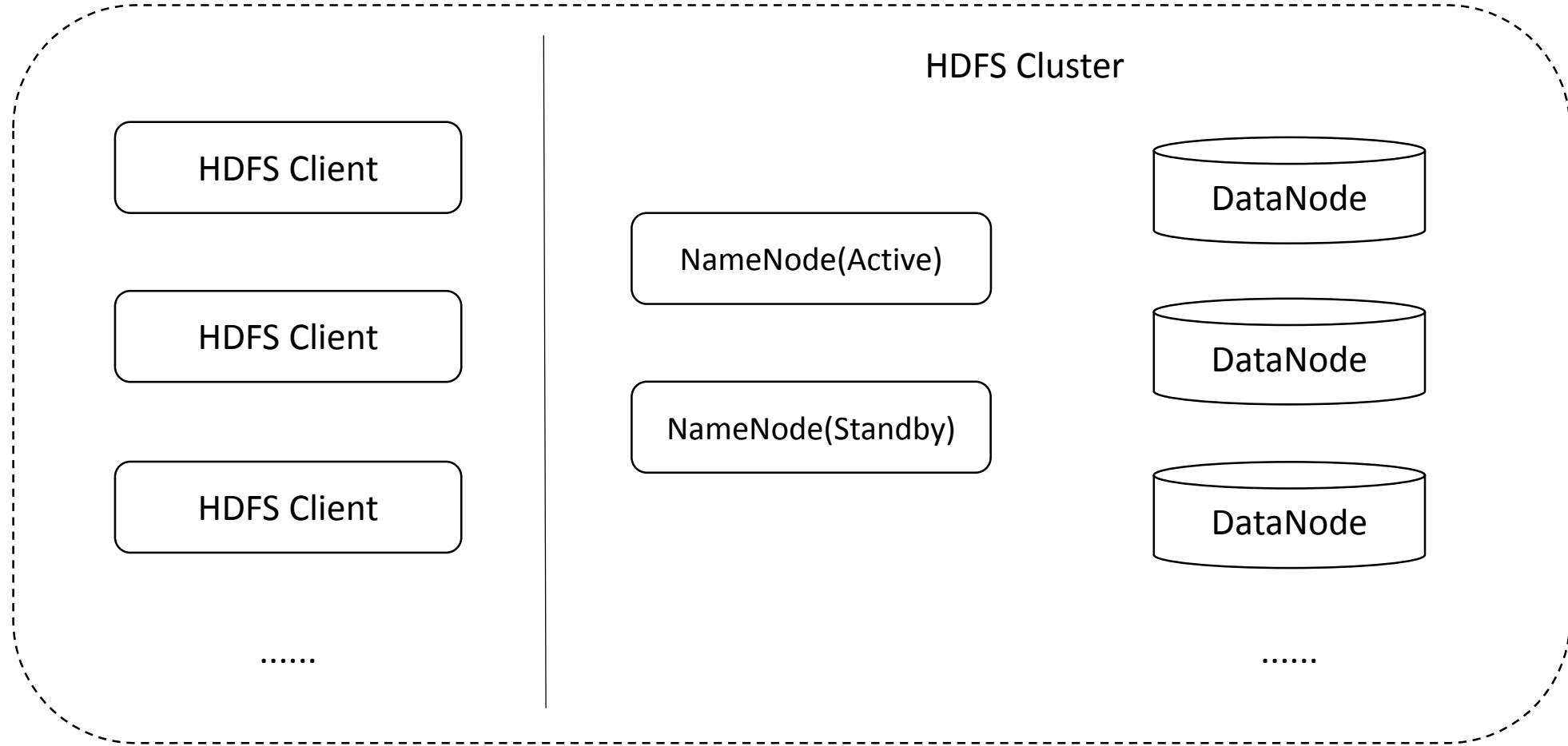
海量
文件

无限制

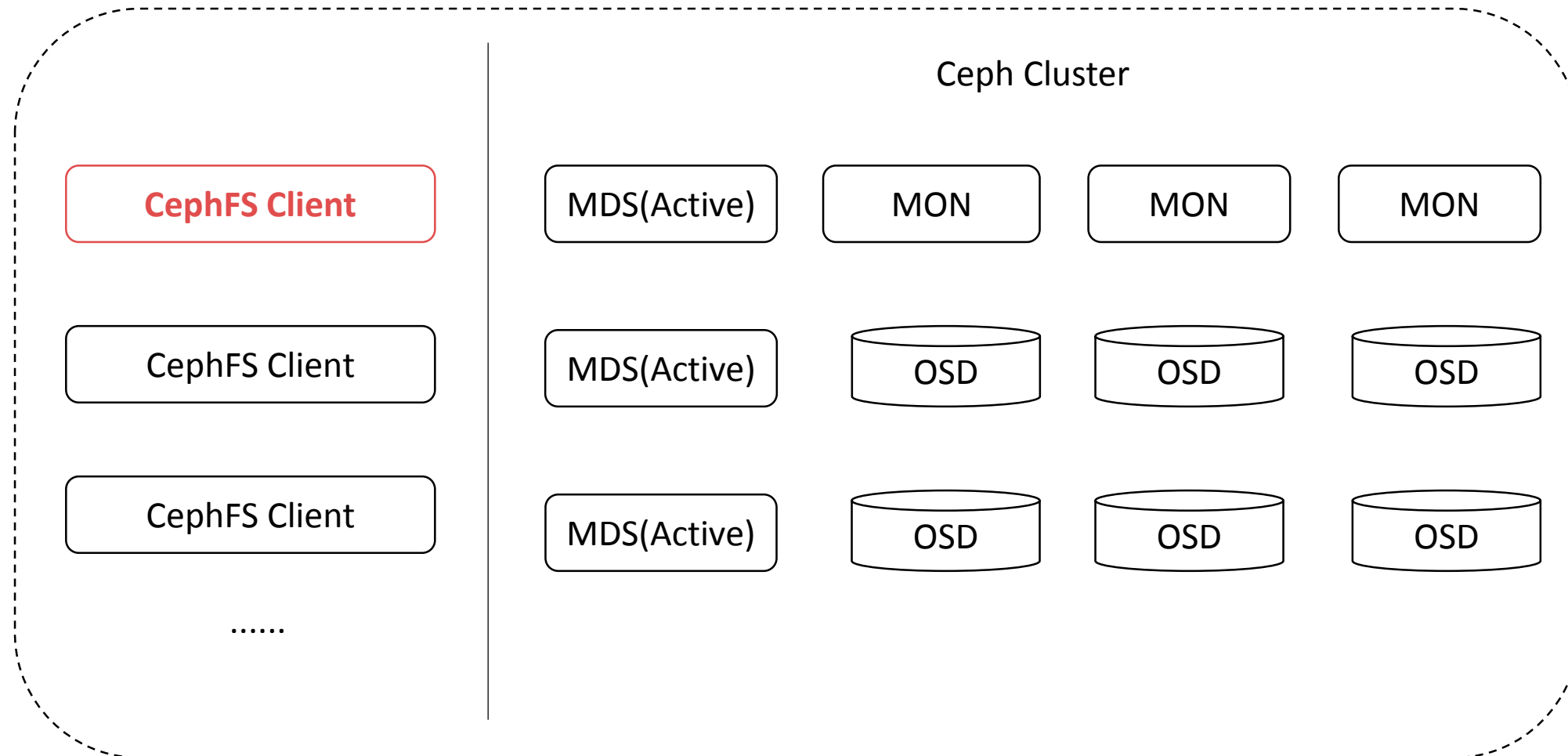
内存
可控

热点元数据

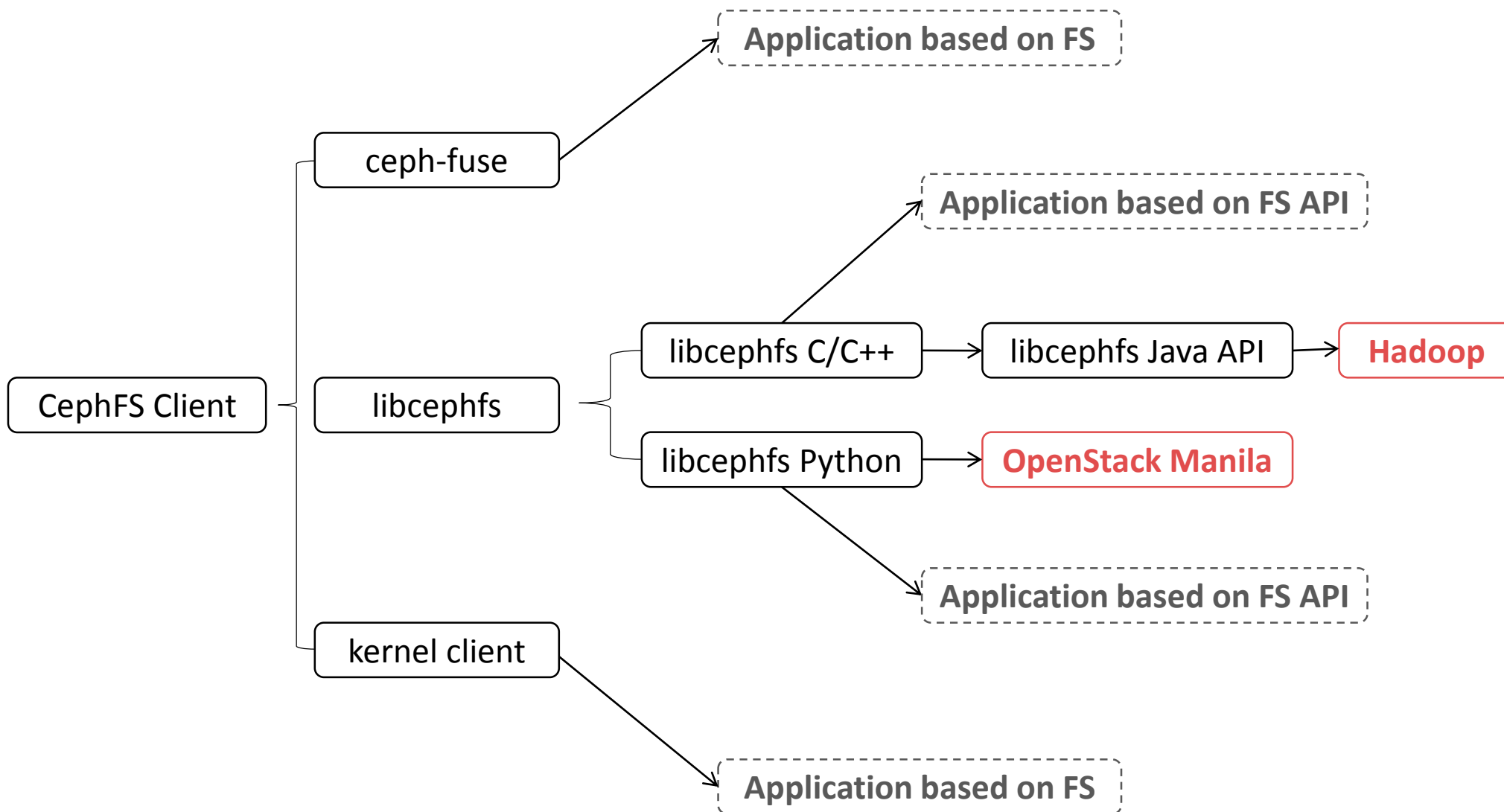
HDFS



Ceph-FS



Ceph-FS



Ceph-FS

Hadoop Cluster

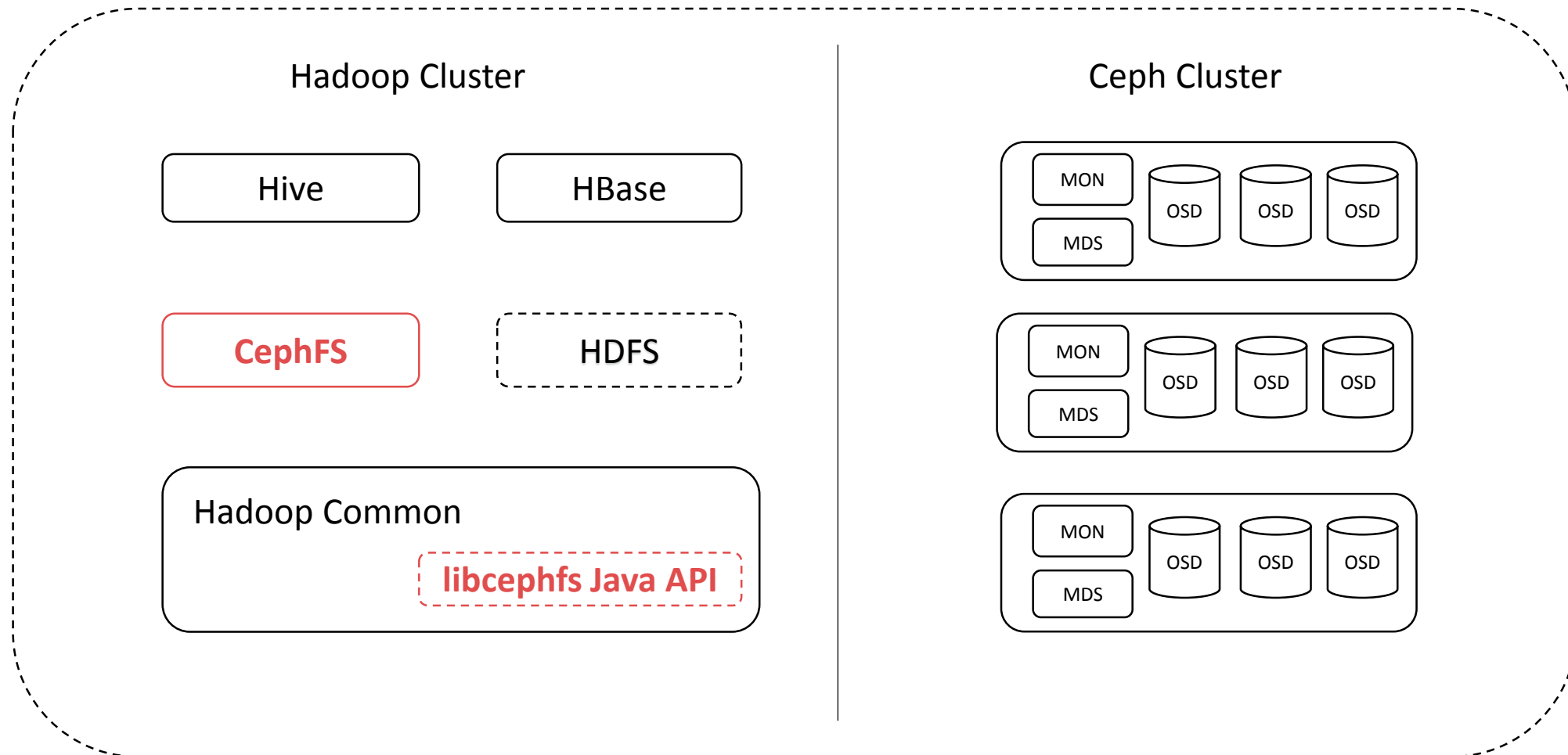
libcephfs Java API

libcephfs JNI API

libcephfs C/C++ API

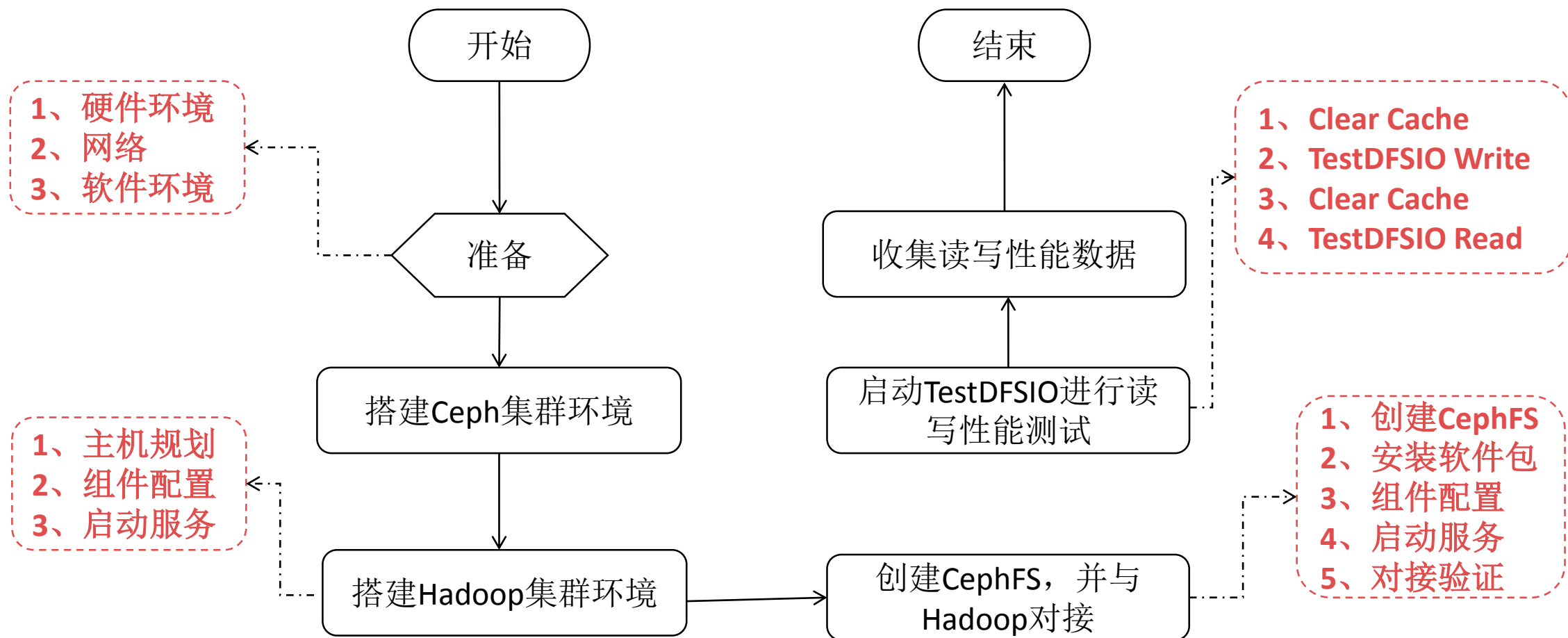
Client(与Ceph集群建立会话)

Hadoop & Ceph-FS



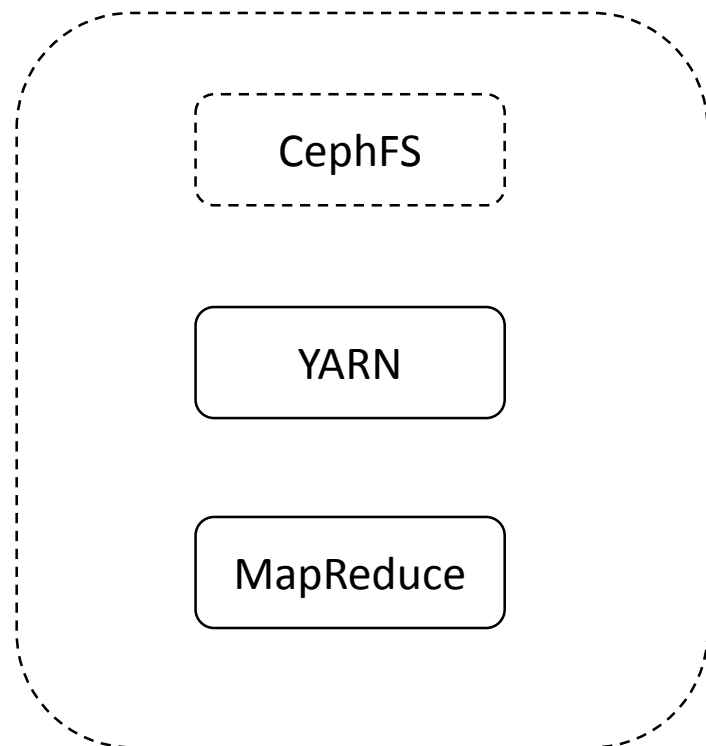
第一次 · 亲密接触

验证流程

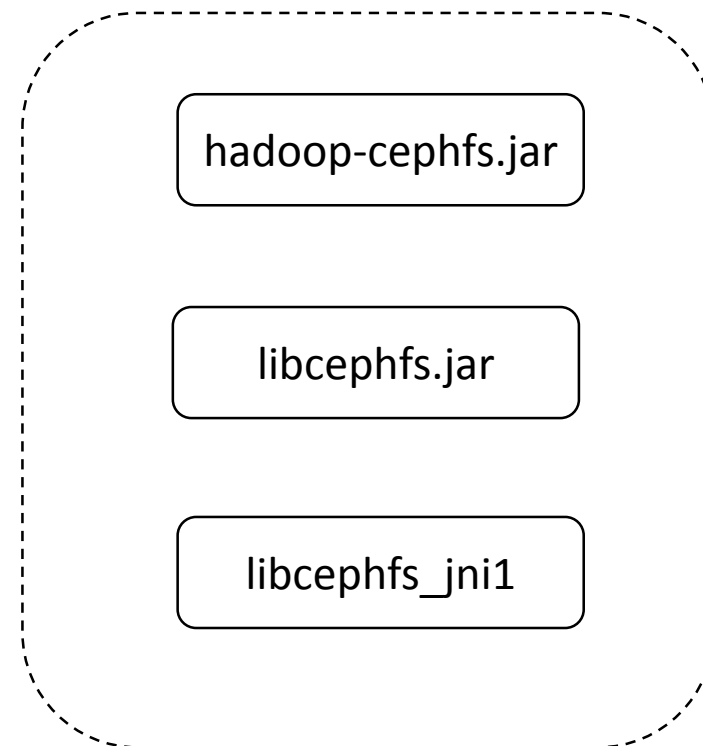


验证准备

Hadoop 集群组件



CephFS 组件



验证障碍

遇到的问题:

- ❑ **FileSystem implementation error - default port -1 is not valid**

hadoop-libcephfs.jar中未定义函数getDefaultPort, 此函数用于获取默认端口。

如果未定义, 则使用0作为默认端口, 在执行过程中进行合法性检查时, 则抛出此异常。

- ❑ **写文件到Hadoop时, 报访问被拒绝错误**

检查ceph权限配置: `ceph auth list` 查看client.hadoop是否有对osd、mon、mds执行相应操作的权限, 在/etc/ceph/ceph.conf中设置:

```
[client]
```

```
client_mount_uid = $(id hadoop)
```

```
client_mount_gid = $(id hadoop)
```

- ❑ **有些节点可以读写cephfs的文件, 有些节点报Operation not permitted错误**

不同节点上hadoop用户的uid和gid不同引起, 修改hadoop用户的uid和gid使其一致

```
id hadoop
```

```
usermod -u $uid hadoop
```

```
groupmod -g $gid hadoop
```

功能验证

Hadoop FS Shell	CephFS	HDFS	Comments
cat	supported	supported	
chgrp	no effect	supported	<i>may support</i>
chmod	supported	supported	
chown	no effect	supported	<i>may support</i>
copyFromLocal	supported	supported	
copyToLocal	supported	supported	
count	supported	supported	
cp	supported	supported	
createSnapshot	unsupported	supported	<i>may support</i>
deleteSnapshot	unsupported	supported	<i>may support</i>
df	unsupported	supported	default values are normally displayed

功能验证

Hadoop FS Shell	CephFS	HDFS	Comments
find	supported	supported	
get	supported	supported	
getfacl	supported	supported	
getfattr	unsupported	supported	<i>may support</i>
ls	supported	supported	
mkdir	supported	supported	
moveFromLocal	supported	supported	
moveToLocal	unsupported	unsupported	
mv	supported	supported	
put	supported	supported	
renameSnapshot	unsupported	supported	<i>may support</i>

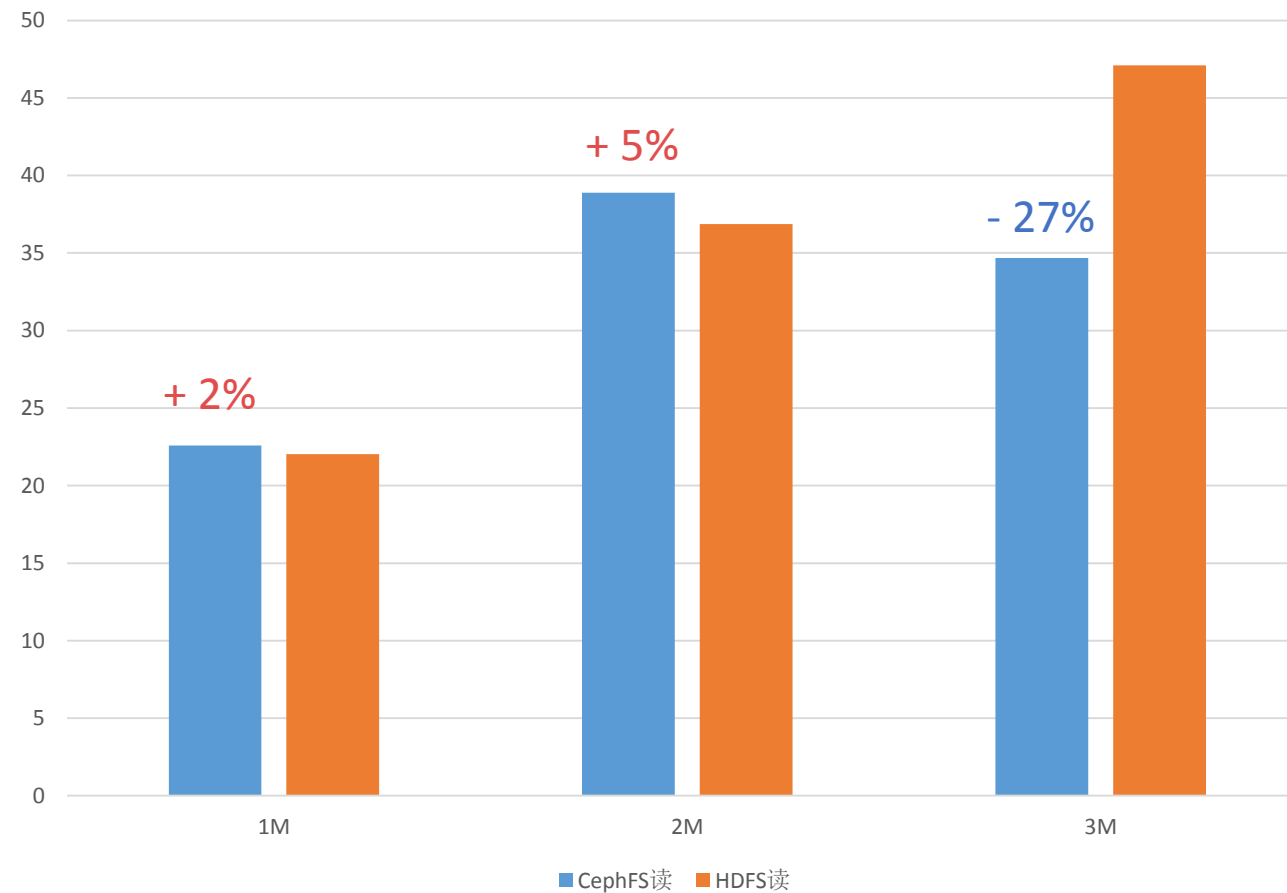
功能验证

Hadoop FS Shell	CephFS	HDFS	Comments
rm	supported	supported	CephFileSystem.delete Method
rmdir	supported	supported	
setfacl	unsupported	supported	<i>may support</i>
setfattr	unsupported	supported	<i>may support</i>
stat	supported	supported	
tail	supported	supported	

<http://hadoop.apache.org/docs/current/hadoop-project-dist/hadoop-common/FileSystemShell.html>
<https://github.com/ceph/cephfs-hadoop>

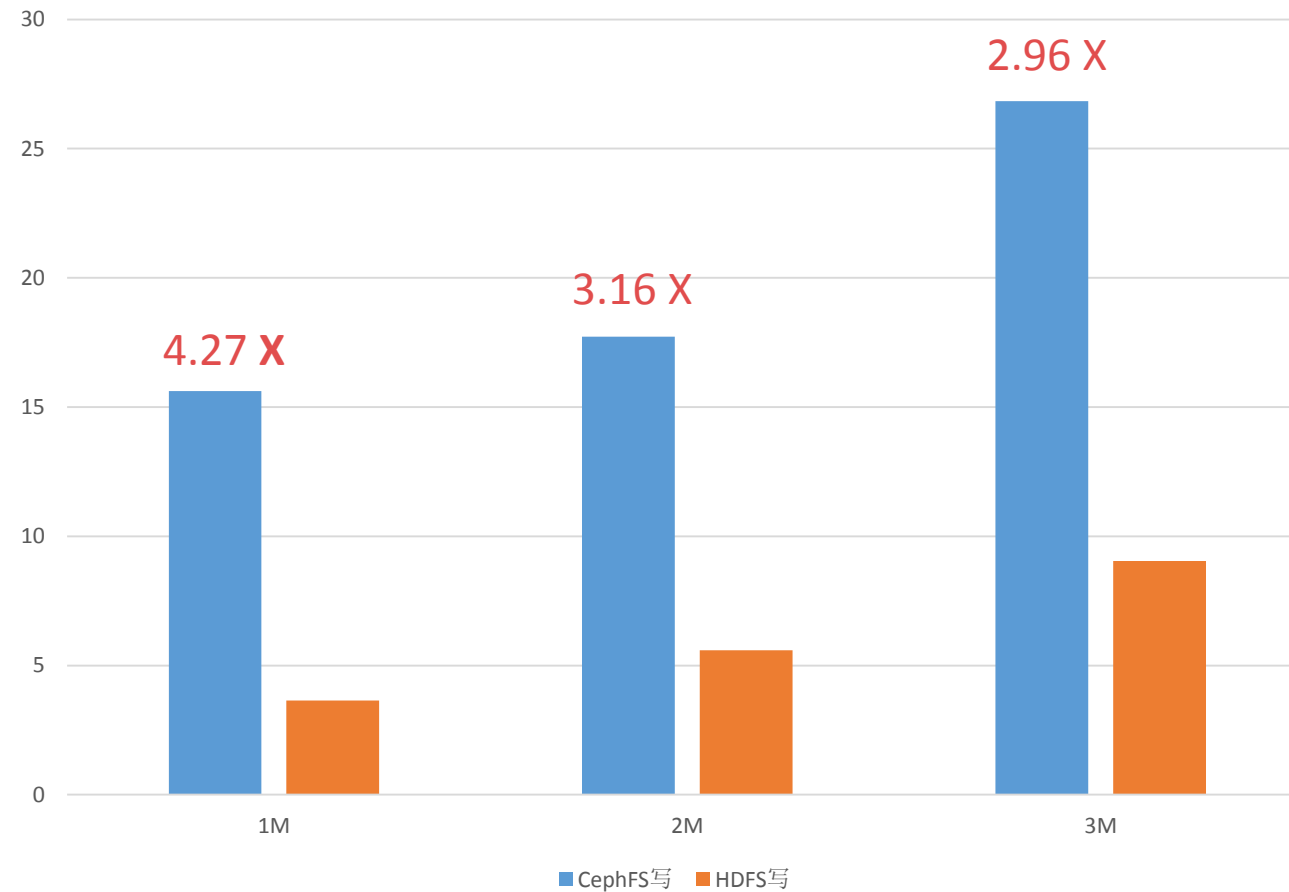
性能验证

CephFS与HDFS读性能比较



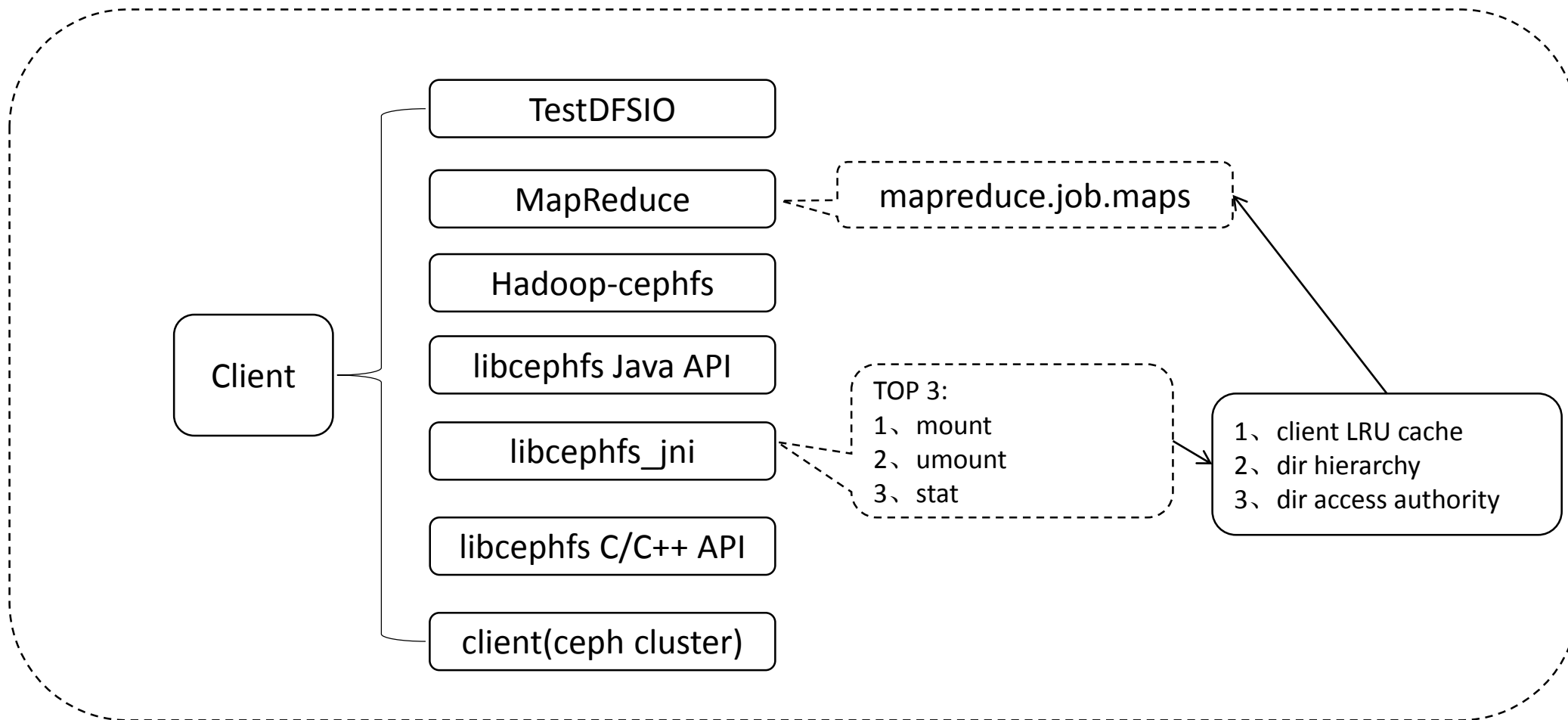
性能验证

CephFS与HDFS写性能比较

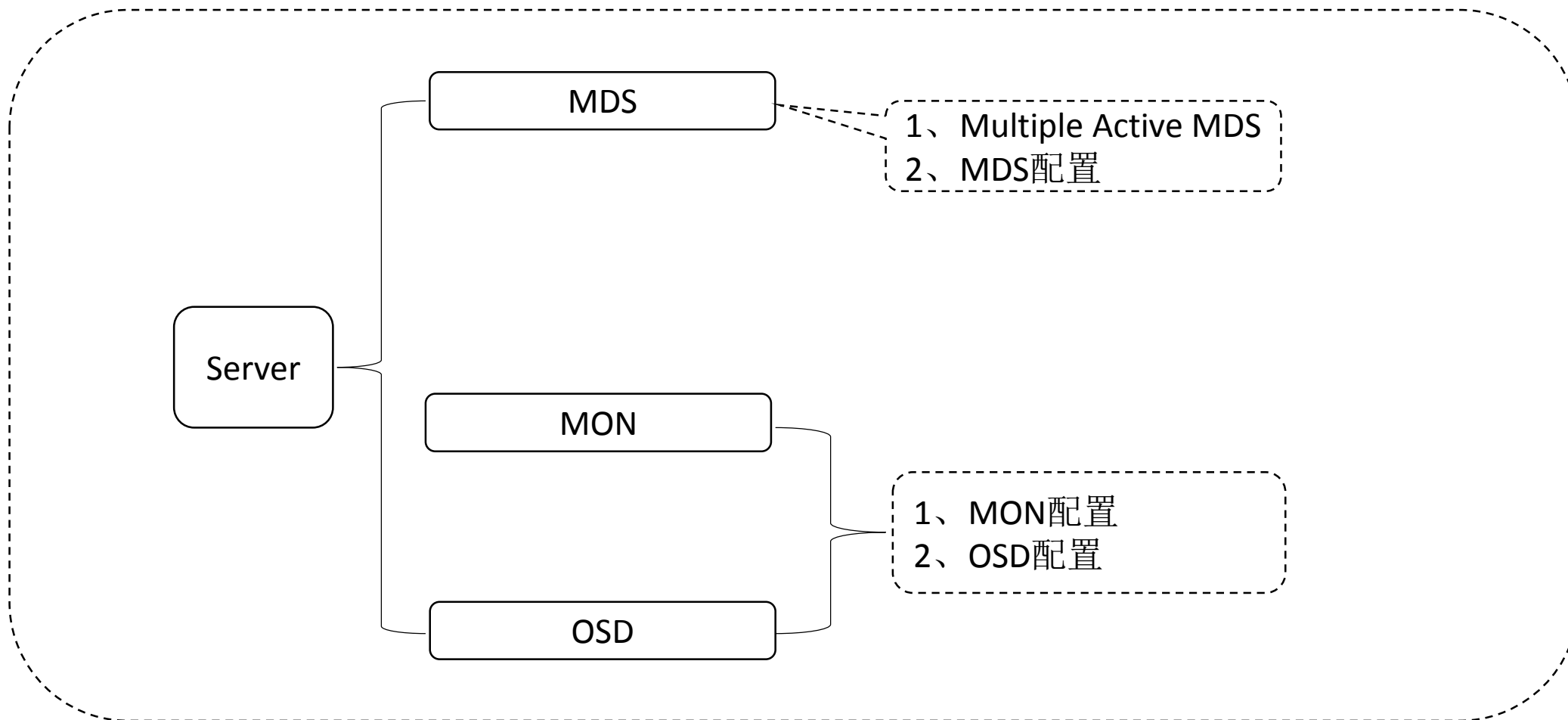


热 · 恋期

性能优化——瓶颈分析(Client)

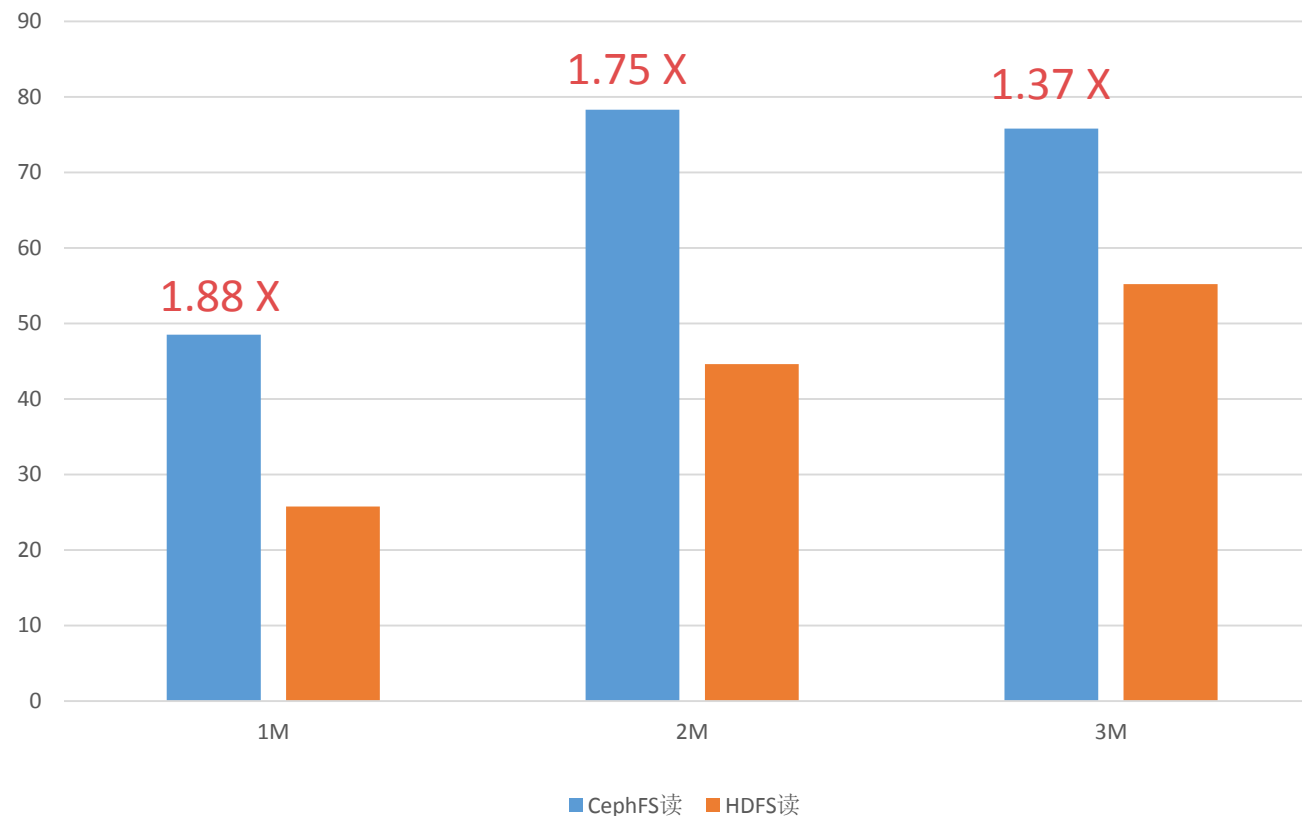


性能优化——瓶颈分析(Server)



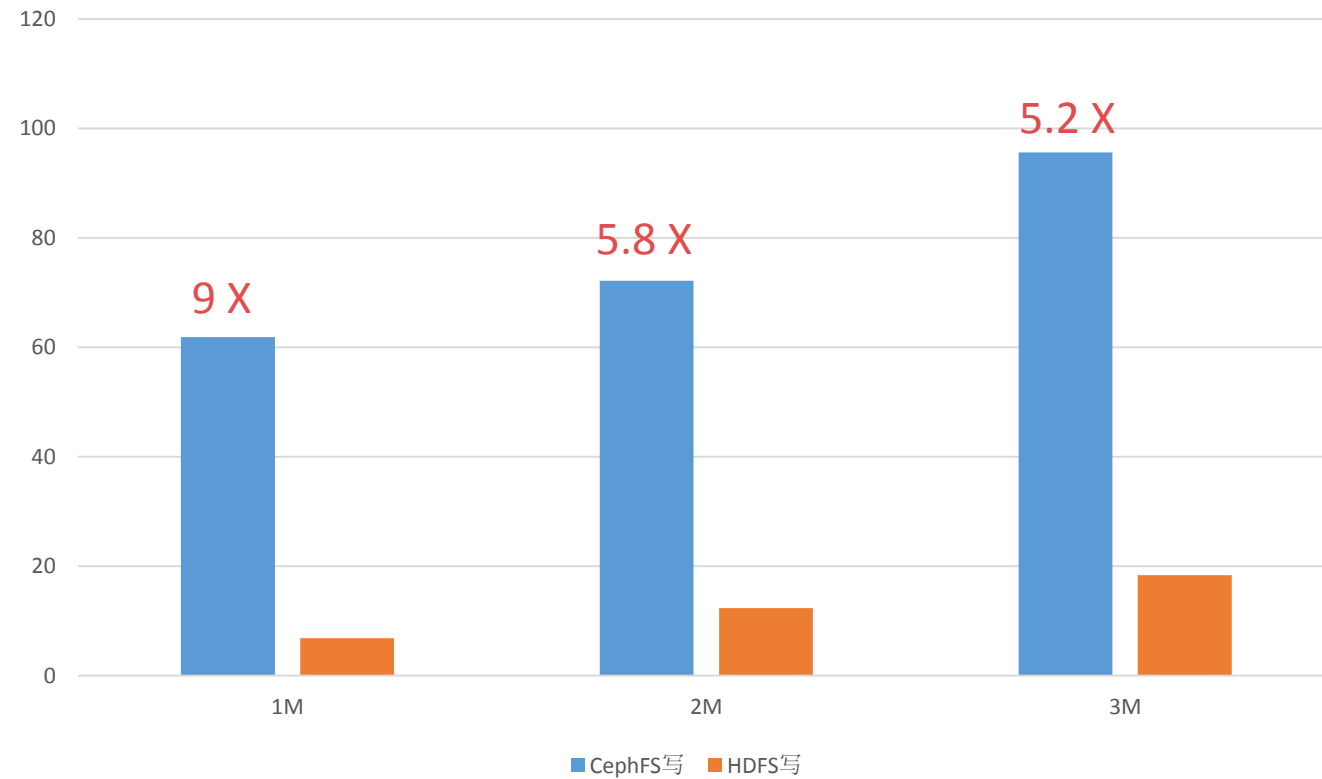
优化后性能验证

CephFS与HDFS读性能比较



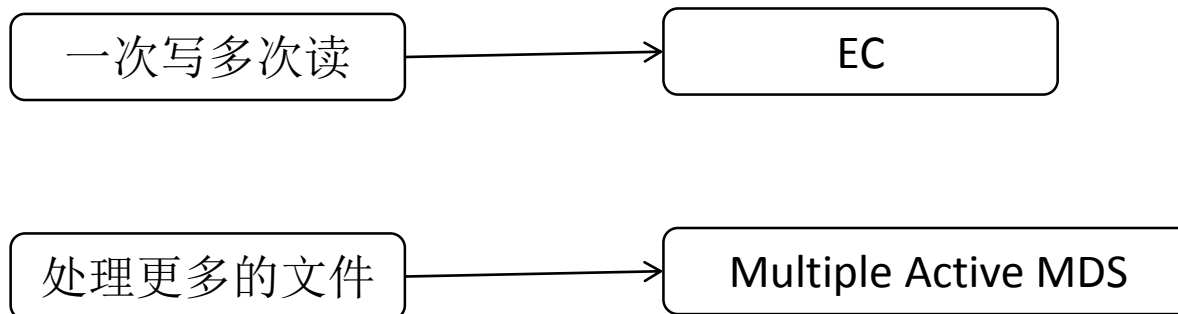
优化后性能验证

CephFS与HDFS写性能比较



未完·待续

后续计划



ZTE中兴