

SHANGHAI JIAOTONG UNIVERSITY

BIG DATA PROCESSING TECHNOLOGY

Project 1: Erasure Code in Hadoop

Author:

FEI Yixiao

118260910031

LI Yanhao

118260910036

LUO Dihao

118260910039

SHEN Shengyang

118260910042

YAN Shuhan

118260910050

Supervisor:

Chentao WU

Xin XIE

October 29, 2019

1 Introduction

Erasur code is a method of data protection in which data is broken into fragments, expanded and encoded with redundant data pieces and stored across a set of different locations or storage media.

In hadoop system, an erasure code called RS code has already been implemented. For every n data units, RS code will have k parity units, which k is set by user, and generate a matrix which size is $(n + k) \times n$, then store the data that has been transformed by this matrix. Although it will cost more disk space, when any data unit becomes unavailable, we could pick n available units for reconstruction. We can pick these n row in the matrix and calculate its inverse matrix, then multiply the vector to get the origin n data unit.

But RS code has its disadvantage, each time we lost a unit, which is the most probable situation, we have to get all n units for reconstruction, it's a huge cost of time, so maybe we could sacrifice some uncommonly occurred situation, for example lost k units at the same time, to decrease the cost of most situation, for example lost only 1 unit.

LRC code is the solution for such situation, similar to RS code, LRC code also has n data units and k parity units, instead of RS code's k parity units are related to every data unit, which we called global parity unit, LRC code has l local parity units only related to certain data units, and r global parity units where $l+r=k$. In ours project, we choose $n = 6, l = 2, r = 2$. The structure shows in Figure 1.

In the two local parity units, each is related to 3 data units, so that if any data unit is lost, we could use the other two data units and the related local parity unit for reconstruction. That only need half the data reading. Such method can'y handle the situation of four certain units lost, which the units are one local parity unit and the three data units it related to. But this situation is almost impossible to happen so the cost is nearly zero.

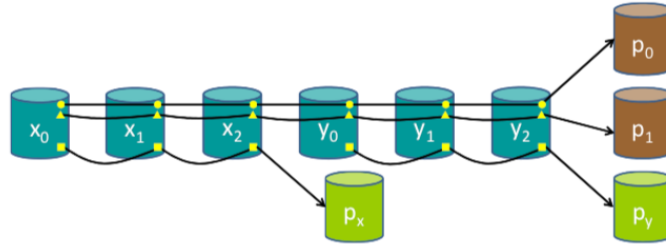


Figure 1: LRC structure

2 Code

The implementation of LRC code is based on JAVA and composed with following files:

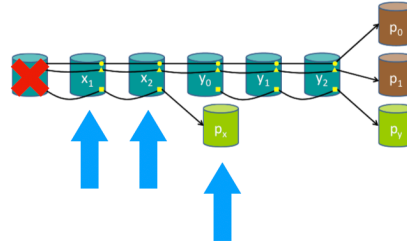
- LRCUtil.java: In this file we define the function to generate the encode matrix and the function to encode data. The encode matrix is composed by an 6×6 identity matrix. Two local parity vectors, one with the first three elements are one and others are zero, another with the last three elements are one and others are zero. Two global parity vectors. All under $GF(2^8)$. The encode matrix shows as below:

$$\begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 1 & 2 & 3 & 4 & 5 & 6 \\ 1 & 4 & 9 & 16 & 25 & 36 \end{bmatrix}$$

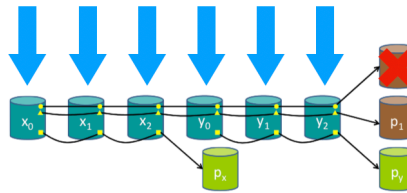
- LRCRawErasureCoderFactory.java: In this file we declare the functions we need and merge the options into LRC code.
- LRCRawEncoder.java: In this file we call the functions in LRCUtil.java to generate encode matrix and do the matrix multiplication to encode data.
- LRCRawDecoder.java: In this file we recover the lost data. We separate the matrix into three parts:
 - First local part: this part includes the first three data units and the related local parity unit.
 - Second local part: this part includes the other three data units and the related local parity unit.
 - Global part: this part includes the two global parity units.

Depend on different situation we will decode under following logic:

- One unit lost: If this unit belongs to the two local part, we will use the remaining three units to recover the lost unit. If it belongs to the global part, we will use the corresponding row of the encode matrix to re-calculate this unit.

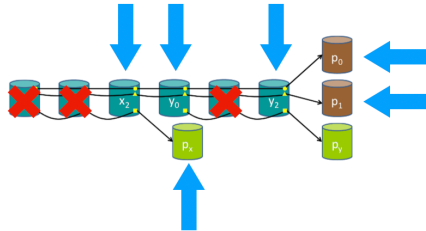


(a) local unit lost



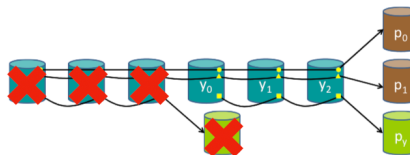
(b) global unit lost

- Two to four units lost: First we get which part does the first unit belong to and use the remaining units in this part, then use the global units if available, if the number of units still lower than 6, we randomly choose units in the last part.



(a) more than one unit lost

- Four units in the same part lost or more than four units lost: This is the situation this LRC code can't recover.



(a) cannot recover

Using the above logic we get the decode matrix either 3×3 or 6×6 , then we calculate the invert matrix and recover the lost unit(s).

3 Example

The feasibility of this LRC code is tested in the hadoop file system in the following steps:

- Compile the whole hadoop project, package it into a native distribution and deploy the hadoop system using docker container.
- Create an empty directory in hadoop file system. Set the erasure code policy to LRC-6-4-1024k.
- Put an 800 MB file in the directory. The file will be encoded by LRC code into 6 data chunks and 4 parity chunks.
- Find the physical locations of file chunks. Choose several chunks and manually delete them.
- Fetch the file from the hadoop file system.
- Verify the identity of the fetched file.

First we compile the whole hadoop project and deploy the hadoop system using docker container. Here our LRC code configuration is composed with 6 data units and 4 parity units, thus we need to deploy 10 datanodes.

Type "hadoop namenode -format" and "\$HADOOP_HOME/sbin/start-all.sh" to launch hadoop file system

[illegible]

Figure 5: start namenode and datanodes

Type "hdfs ec -listPolicies" to check if LRC policy is registered in hadoop.

```
root@h01:~# hdfs ec -listPolicies
Erasure Coding Policies:
ErasureCodingPolicy=[Name=LRC-6-4-1024k, Schema=[ECSSchema=[Codec=lrc, numDataUnits=6, numParityUnits=4]], CellSize=1048576, Id=6], State=ENABLED
ErasureCodingPolicy=[Name=R5-10-4-1024k, Schema=[ECSSchema=[Codec=r5, numDataUnits=10, numParityUnits=4]], CellSize=1048576, Id=5], State=DISABLED
ErasureCodingPolicy=[Name=R5-3-2-1024k, Schema=[ECSSchema=[Codec=r5, numDataUnits=3, numParityUnits=2]], CellSize=1048576, Id=2], State=DISABLED
ErasureCodingPolicy=[Name=R5-6-3-1024k, Schema=[ECSSchema=[Codec=r5, numDataUnits=6, numParityUnits=3]], CellSize=1048576, Id=1], State=ENABLED
ErasureCodingPolicy=[Name=R5-LEGACY-6-3-1024k, Schema=[ECSSchema=[Codec=r5-legacy, numDataUnits=6, numParityUnits=3]], CellSize=1048576, Id=3], State=DISABLED
ErasureCodingPolicy=[Name=XOR-2-1-1024k, Schema=[ECSSchema=[Codec=xor, numDataUnits=2, numParityUnits=1]], CellSize=1048576, Id=4], State=DISABLED
root@h01:~#
```

Figure 6: list erasure code policies

Create an empty directory in hadoop file system. Enable the LRC policy and set the erasure code policy to **LRC-6-4-1024k**.

```
root@h01:~# hadoop fs -mkdir /input
mkdir: '/input': File exists
root@h01:~# hdfs ec -enablePolicy -policy LRC-6-4-1024k
Erasure coding policy LRC-6-4-1024k is enabled
root@h01:~# hdfs ec -setPolicy -path /input -policy LRC-6-4-1024k
```

Figure 7: set LRC policy

Type "hadoop fs -put b.dat /input " to put an 800 MB file b.dat in the directory /input.

```
root@h01:~# hadoop fs -put b.dat /input
2019-10-28 07:03:28,850 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:28,859 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:28,863 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:28,867 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:28,874 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:28,907 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:29,033 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:29,044 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:29,054 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-28 07:03:29,057 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
root@h01:~#
```

Figure 8: put the file b.dat in /input

Use "hdfs fsck /input -files -blocks -locations -racks" to find the physical locations of the file in all nodes.

```

root@h01:~# hdfs fsck /input -files -blocks -locations -racks
Connecting to namenode via http://h01:9876/fsckugiroot&files=1&blocks=1&locations=1&racks=1&path=/input
FSCK started by root (auth:SIMPLE) from /172.18.0.2 for path /input at Mon Oct 28 07:19:32 UTC 2019
/input -dir
/input/b.dat 800004800 bytes, erasure-coded: policyLRC 6:4:1024k, 1 block(s): OK
0 BP-1827639026-172.18.0.2:1572244616297:blk_-9223372036854775721:1004 len=800004800 live_repl=10 (blk_-9223372036854775720:/default-rack/172.18.0.11:9866, blk_-9223372036854775727:/default-rack/172.18.0.19866, blk_-9223372036854775720:/default-rack/172.18.0.2:9866, blk_-9223372036854775721:/default-rack/172.18.0.8:9866, blk_-9223372036854775724:/default-rack/172.18.0.12:9866, blk_-9223372036854775723:/default-rack/172.18.0.7:9866, blk_-9223372036854775722:/default-rack/172.18.0.3:9866, blk_-9223372036854775721:/default-rack/172.18.0.10:9866, blk_-9223372036854775720:/default-rack/172.18.0.5:9866, blk_-9223372036854775719:/default-rack/172.18.0.4:9866)

Status: HEALTHY
Number of data-nodes: 11
Number of racks: 1
Total dufs: 1
Total symlinks: 0

Replicated blocks:
Total size: 0 B
Total files: 0
Total blocks (validated): 0
Minimally replicated blocks: 0
Over-replicated blocks: 0
Under-replicated blocks: 0
Mis-replicated blocks: 0
Default replication factor: 1
Average block replication: 0.0
Missing blocks: 0
Corrupt blocks: 0
Missing replicas: 0

Erasure Coded Block Groups:
Total size: 800004800 B
Total files: 1
Total block groups (validated): 1 (avg. block group size 800004800 B)
Minimally erasure-coded block groups: 1 (100.0 %)
Over-erasure-coded block groups: 0 (0.0 %)
Under-erasure-coded block groups: 0 (0.0 %)
Unsatisfactory placement block groups: 0 (0.0 %)
Average block group size: 10.0
Missing block groups: 0
Corrupt block groups: 0
Missing internat blocks: 0 (0.0 %)
FSCK ended at Mon Oct 28 07:19:32 UTC 2019 in 1 milliseconds

The filesystem under path '/input' is HEALTHY
root@h01:~#

```

Figure 9: find physical locations

From the above figure the 10 chunks are stored in 172.18.0.11, 172.18.0.6, ..., 172.18.0.4 in order, which respectively correspond to $x_0, x_1, x_2, y_0, y_1, y_2, p_x, p_y, p_0, p_1$. We delete the data unit x_0 .

```

root@h01:/home/hadoop3/hadoop/tmp/dfs/data/current/BP-1827639026-172.18.0.2-1572244616297/current/finalized/subdir0/subdir0# ls
blk_-9223372036854775744 blk_-9223372036854775744_1004.meta
root@h01:/home/hadoop3/hadoop/tmp/dfs/data/current/BP-1827639026-172.18.0.2-1572244616297/current/finalized/subdir0/subdir0# rm *
root@h01:/home/hadoop3/hadoop/tmp/dfs/data/current/BP-1827639026-172.18.0.2-1572244616297/current/finalized/subdir0/subdir0#

```

Figure 10: physically delete x_0

Use "hadoop fs -cat /input/b.dat > c.dat" command to fetch the file from hadoop file system.

```

root@h01:~# hadoop fs -cat /input/b.dat > c.dat
2019-10-20 08:33:49,470 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-20 08:33:49,584 WARN ImplBlockReaderFactory: I/O error constructing remote block reader:
Java.io.IOException: Got error, status=0000, status message openReadBlock BP-1076424496-172.18.0.2-157237839587-b1k_-9223372036854775792_1001
537036854775792 is not valid, for OP_READ_BLOCK, self=/172.18.0.2:35462, remote=/172.18.0.2:38866, for file /input/b.dat, for pool BP-1076424496-172.18.0.2-157237839587 block_-9223372036854775792_1001
at org.apache.hadoop.hdfs.protocol.datatransfer.DataTransferProtoUtil.checkLockStatus(DataTransferProtoUtil.java:134)
at org.apache.hadoop.hdfs.protocol.datatransfer.DataTransferProtoUtil.checkLockStatus(DataTransferProtoUtil.java:116)
at org.apache.hadoop.hdfs.client.impl.BlockReaderRemote.checkSuccess(BlockReaderRemote.java:446)
at org.apache.hadoop.hdfs.client.impl.BlockReaderRemote.newBlockReader(BlockReaderRemote.java:408)
at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.getRemoteBlockReader(BlockReaderFactory.java:854)
at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.getRemoteBlockReaderFromTip(BlockReaderFactory.java:758)
at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.build(BlockReaderFactory.java:380)
at org.apache.hadoop.hdfs.DFSInputStream.getLockHeader(DFSInputStream.java:644)
at org.apache.hadoop.hdfs.DFSInputStream.createLockHeader(DFSInputStream.java:264)
at org.apache.hadoop.hdfs.StripeHeader.readStripe(StripeHeader.java:290)
at org.apache.hadoop.hdfs.StripeHeader.readStripe(StripeHeader.java:136)
at org.apache.hadoop.hdfs.DFSInputStream.readDownStripe(DFSInputStream.java:326)
at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:419)
at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:829)
at java.io.DataInputStream.read(DataInputStream.java:180)
at org.apache.hadoop.io.IOUtils.copyBytes(IOUtils.java:94)
at org.apache.hadoop.io.IOUtils.copyBytes(IOUtils.java:68)
at org.apache.hadoop.io.IOUtils.copyBytes(IOUtils.java:129)
at org.apache.hadoop.fs.shell.DisplayCat.printToStdout(Display.java:103)
at org.apache.hadoop.fs.shell.DisplayCat.processPath(Display.java:96)
at org.apache.hadoop.fs.shell.Command.processPathInternal(Command.java:367)
at org.apache.hadoop.fs.shell.Command.processPath(Command.java:331)
at org.apache.hadoop.fs.shell.Command.processPathArgument(Command.java:304)
at org.apache.hadoop.fs.shell.Command.processArguments(Command.java:276)
at org.apache.hadoop.fs.shell.FsCommand.processArguments(FsCommand.java:128)
at org.apache.hadoop.fs.shell.Command.run(Command.java:177)
at org.apache.hadoop.fs.FsShell.run(FsShell.java:327)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:70)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:90)
at org.apache.hadoop.fs.FsShell.main(FsShell.java:386)
2019-10-20 08:33:49,587 WARN Hdfs.DFSClient: failed to connect to /172.18.0.2:38866 for blockBP-1076424496-172.18.0.2-157237839587-b1k_-9223372036854775792_1001
Java.io.IOException: Got error, status=0000, status message openReadBlock BP-1076424496-172.18.0.2-157237839587-b1k_-9223372036854775792_1001
537036854775792 is not valid, for OP_READ_BLOCK, self=/172.18.0.2:35462, remote=/172.18.0.2:38866, for file /input/b.dat, for pool BP-1076424496-172.18.0.2-157237839587 block_-9223372036854775792_1001
at org.apache.hadoop.hdfs.protocol.datatransfer.DataTransferProtoUtil.checkLockStatus(DataTransferProtoUtil.java:134)
at org.apache.hadoop.hdfs.protocol.datatransfer.DataTransferProtoUtil.checkLockStatus(DataTransferProtoUtil.java:116)
at org.apache.hadoop.hdfs.client.impl.BlockReaderRemote.checkSuccess(BlockReaderRemote.java:446)
at org.apache.hadoop.hdfs.client.impl.BlockReaderRemote.newBlockReader(BlockReaderRemote.java:408)
at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.getRemoteBlockReader(BlockReaderFactory.java:854)
at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.getRemoteBlockReaderFromTip(BlockReaderFactory.java:758)
at org.apache.hadoop.hdfs.client.impl.BlockReaderFactory.build(BlockReaderFactory.java:380)
at org.apache.hadoop.hdfs.DFSInputStream.getLockHeader(DFSInputStream.java:644)
at org.apache.hadoop.hdfs.DFSInputStream.createLockHeader(DFSInputStream.java:264)
at org.apache.hadoop.hdfs.StripeHeader.readChunk(StripeHeader.java:290)
at org.apache.hadoop.hdfs.StripeHeader.readStripe(StripeHeader.java:136)
at org.apache.hadoop.hdfs.DFSInputStream.readDownStripe(DFSInputStream.java:326)
at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:419)
at org.apache.hadoop.hdfs.DFSInputStream.read(DFSInputStream.java:829)
at java.io.DataInputStream.read(DataInputStream.java:180)
at org.apache.hadoop.io.IOUtils.copyBytes(IOUtils.java:94)
at org.apache.hadoop.io.IOUtils.copyBytes(IOUtils.java:68)
at org.apache.hadoop.io.IOUtils.copyBytes(IOUtils.java:129)
at org.apache.hadoop.fs.shell.DisplayCat.printToStdout(Display.java:103)
at org.apache.hadoop.fs.shell.DisplayCat.processPath(Display.java:96)
at org.apache.hadoop.fs.shell.Command.processPathInternal(Command.java:367)
at org.apache.hadoop.fs.shell.Command.processPathArgument(Command.java:304)
at org.apache.hadoop.fs.shell.Command.processArguments(Command.java:276)
at org.apache.hadoop.fs.shell.FsCommand.processArguments(FsCommand.java:128)
at org.apache.hadoop.fs.shell.Command.run(Command.java:177)
at org.apache.hadoop.fs.FsShell.run(FsShell.java:327)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:70)
at org.apache.hadoop.util.ToolRunner.run(ToolRunner.java:90)
at org.apache.hadoop.fs.FsShell.main(FsShell.java:386)
2019-10-20 08:33:49,593 WARN Hdfs.DFSClient: [datanode1@out10storage[172.18.0.2:8066,05-ff2dc09-9109-4016-9ff5-e811ab9a882,01K]] are unavailable and all striping blocks on them are lost. Ignored nodes
= null
2019-10-20 08:33:49,594 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-20 08:33:49,636 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-20 08:33:49,718 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-20 08:33:49,791 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-20 08:33:49,866 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-20 08:33:49,954 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
root@h01:~# [

```

Figure 11: fetch b.dat from hadoop and copy it to c.dat

We can see from the above figure that a block of the target file b.dat is first missing, but then is reconstructed from other blocks. The command outputs correct content to c.dat.

Finally use "cmp" command to verify the identity of c.dat

```

root@h01:/home/hadoop3/hadoop/tmp/dfs/data/current/BP-1827639026-172.18.0.2-1572244616297/current/finalized/subdir0/subdir0# cmp c.dat /b.dat
root@h01:/home/hadoop3/hadoop/tmp/dfs/data/current/BP-1827639026-172.18.0.2-1572244616297/current/finalized/subdir0/subdir0#

```

Figure 12: verify the identity of c.dat

We test different missing blocks configurations with successful reconstructions:

- x_0 unit is missing
- x_0 and y_0 units are missing
- x_0 , y_0 and y_1 units are missing
- x_0 , y_1 , p_x and p_0 units are missing

and also an expected failure configuration in which all missing blocks belong to one local part X :

- x_0, x_1, x_2 and p_x units are missing

```

2019-10-29 10:10:10,831 WARN hdfs.DFSClient: [datanodeInfoWithStorage[172.18.0.2:19860_05-c9d7c97f-47fc-4878-9c09-6d9456373caf_015K]] are unavailable and all striping blocks on them are lost. IgnoredNodes = null
2019-10-29 10:10:10,832 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-29 10:10:10,859 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
2019-10-29 10:10:10,866 INFO sasl.SaslDataTransferClient: SASL encryption trust check: localHostTrusted = false, remoteHostTrusted = false
-cats: Too many erased in a local part, data not recoverable

```

Figure 13: error: too many missing blocks in a local part