

Data Integrity (HDFS)

- Hadoop user expect no data will be lost / corrupted during storage / processing
- Every input / output operation on disk / r/w carries small chance of introducing errors - Read / write
- When lot of data flowing (flowing) - chances are high
- checksum - way of detecting corrupted data - first enter the sys, transmit across a channel.
- eg CRC-32 (Cyclic Redundancy Check) - Hadoop Checksum file system
- HDFS - uses more efficient variant CRC-32 C (4 byte long), Storage overhead less than 1%.
- Verify data - from Client / other Data Node
- If DN detects error - Client receive I/o Ex.
- Addition, DN runs Data Block Scanner (daemon) - periodically verifies all blocks stored on DN, to guard against corruption ('bit rot').
- HDFS stores replica - it can heal corrupted data, copy from good replica.
- Disable verification - false, setVerifyChecksum() on filesystem before open().
- Shell → -ignorecrc with -set / -copyToLocal

Local file System

- Perform Client Side checksum
- When write file 'filename', the file system create 'filename.crc' in same directory for each chunk of file.
- Chunk is controlled by -file.bytes-per-checksum (512 bytes - Default).

- CheckSum verified when file read - exception.
- ↳ fairly cheap to compute (Java-native code) adding few percentage overhead - time to R/W.
- Acceptable price to pay for Data Integrity.
- Disable - use RawLogicalFilesystem instead of LFS.

CheekSum File System

- wrapper around FileSystem
- LFS uses this to do its work
- makes it easy to add CheckSums to other FS.
- method - ① getRawFilesystem() - underlying FS is called RawFS.

② getCheckSumfile() - Path of checksum file.

- If error detected - while reading file - it will call its reportCheckSumFailure() method.
- Default implementation does nothing | But LFS moves offending file + CheckSum to a side directory called bad_file.

2)

Compression

- brings 2 benefits
- ① reduce Space - to store files
- ② Speed up data transfer - accessing N/W or disk.

DEFLATE - is a compression algo whose standard implementation is ZLIB.

- gzip is DEFLATE with extra header and footer.

<u>Compression format</u>	<u>Tool</u>	<u>Algo</u>	<u>filename extension</u>	<u>Spittable</u>
DEFLATE	N/A	DEFLATE	.deflate	NO
gzip	gzip	-do-	.gz	NO
bzip2	bzip2	bzip2	.bz2	YES
LZO	lzip	LZO	.LZO	NO
LZ4	N/A	LZ4	.LZ4	NO
Snappy	N/A	Snappy	.snappy	NO

- Space/time trade off \rightarrow fast compression / decom vs smaller space saving
- Nine deft option \rightarrow -1 optimize for Speed - 9 means optimize for Space.
- % gzip -1 file (fastest comp method)
- gzip - middle of space / time trade off
- bzip compresses more effective than gzip but slower
- LZO, LZ4 faster - but compress less
- Spittable compression formats - suitable for MR.

Codec

- Implementation Comp / decom algo.
- represented by impl of Compression Codec interface es GZipCodec - Compression / decom of gzip.
- Compress data written to output Stream use CreateOutputStream method to create a CompressionOutput Stream.
- Call CreateInputStream to obtain CompressionInputStream to read uncompress data from Stream

Compression format

DEFLATE

gzip

bz2

LZO

LZ4

Snappy

Hadoop Compressor/Codecorg.apache.hadoop.io.CompressionDefaultCodec

• GzipCodec

• BZip2Codec

• LzoCodec

• LZ4Codec

• SnappyCodec

Native Library

- Fair performance - it is preferable
- Some Cases - reduced deComp - 50%, Comp - 10%
- Hadoop runs fine while using Java classes

Compression formatJava ImplementationNative Impl.

DEFLATE

Yes

Yes

gzip

Yes

Yes

bz2

Yes

Yes

LZO

No

Yes

LZ4

No

Yes

Snappy

No

Yes

- Hadoop has native impl of certain components for
 - Performance reasons
 - Non availability of Java impl.
 - Some Codec only support in Native
- Available in libhadoop.so (Native Hadoop lib)
- Don't have to change any config setting - use NL extramarks

- Disable - when debugging a computation related prob
↳ setting property io.mnative.lib available to false
- CodecPool - if using NL and doing lots of Com/Decom
Consider this
 - allow to reuse Com/Decom
 - amortizing the cost of creating objects

Compression and Input Split

- eg Uncompress file 1GB - HDFS (128 block size) - 8 blocks
- MR Jobs - 8 InputSplit - Separate Map Tasks
- gzip ⇒ Splitting won't work - impossible to read at arbitrary point in gzip stream
- impossible to read split independently
- gzip use DEFLATE to store Comp/Decom
- problem - start of each block is not distinguished
- gzip don't support splitting
- In this MR not try to split gzip file.
- will work at expense of locality
- 1 map will process 8 blocks
- less granular | takes longer time.

LZO - same problem

- possible to preprocess LZO file using indexer tool
- tool build an index of split points - making them splittable

bzip2 - provide syn maker b/w blocks - support splitting

Which Compression format

- Use Container file format - Seq file, Avro, Parquet
- Support both Comp / Splitters
- Use Comp format that support Splitters - bz2, bzip2, LZO (Compressed)
- Split file into chunk, compress each chunk sep

Compression in MapReduce

- MR decompress automatically as they read using filename extension
- To Compress MR o/p - Set job Configs mapreduce.output.fileoutputformat.compress = true
- Codec - class name

Seq file set

- mapreduce.output.fileoutputformat.compress.type - Control type of compression to use.
- def - RECORD
- Block (Recommended)
- Use Compression for intermediate o/p of map phase
- Use Compression LZO, LZ4, Snappy - fine performance gain
- Volume of data reduced

Map reduce compression properties

- 1 - mapreduce.output.fileoutputformat.compress - boolean / false
- 2 - Default Codec Class → Codec Class name
- 3 - type - String / RECORD