

Map output Compression Properties

- mapreduce.map.output.compress - boolean | false
- mapreduce.map.output.compress.codec - class | DefaultCodec

Serialization

- turn objects into byte stream
- used for transmission over n/w, writing to persistent store
- Deserialization - reverse process
- used in 2 distinct area - IPC (RPC), Persistent storage
- RPC uses serialization to render msg to binary Stream
 - ↳ compact | fast | extensible | Interoperable
- RPC lifespan less than second | Persistent data - 4 year
- RPC 4 properties - persistent also consider
- Hadoop uses its own serialization formats.
- Writable - compact, fast | not easy to extend or use from other language than Java

Writable interface

- 2 methods (write(), readFields())
 - ↳ 1) writing its state to DataOutput binary stream
 - 2) reading its state from DataInput binary stream
- IntWritable - wrapper for Java int.
- IntWritable w = new IntWritable();
 - w.write(100);

or IntWritable w = new IntWritable(100);

- Need help to serialized IntWritable

Serialization
 ByteArrayOutputStream out = new ByteArrayOutputStream();
 DataOutputStream do = new DataOutputStream(out);
 extramarks Writable.write(do);
 do.close()

`ByteInputStream in = new ByteInputStream(bytes)`
`DataInputStream datain = new DataInputStream(in)`
`writable.readFields(datain)`
`datain.close()`

- Inheritable in Element's WritableComparable interface
- It permit implementors to compare record w/o deserializing them object - ^{avoid} overhead of object creation

Writable classes

<u>Java Primitive</u>	<u>writableImpl</u>	<u>Serializable</u> <small>(byte)</small>
-----------------------	---------------------	---

boolean	BooleanWritable	
---------	-----------------	--

byte	ByteWritable	
------	--------------	--

short	ShortWritable	
-------	---------------	--

int	IntWritable	
-----	-------------	--

float	FloatWritable	
-------	---------------	--

long	LongWritable	
------	--------------	--

double	DoubleWritable	
--------	----------------	--

Text - writable for UTF-8 Sep.

- writable equivalent to Strip

Sequence file

- Log file - each key record is a new line of text
- ~~For Log binary type, Plain text is not suitable format~~
- Sequencefile fits - providing persistent data structure for binary Key value pair.

- Key-Value and (key-value), value (mutable) data
- Also containers for smaller files
- Packup files into SequenceFile makes Storing and Processing smaller file - efficient
- work naturally with Hadoop Data type
- includes metadata which identifies the data types of key-value
- Actually 3 file types in one
 - ↳ Uncompressed / Record compressed / Block Compressed
- Used in Map Reduce
 - ↳ Sequencefile InputFormat
 - ↳ Sequencefile OutputFormat

Problems

- Useful but Potential problems
- Only typical accessible via Java API
- If def of key-value object changes, file becomes unreadable.
- Alternate is AVRO
- Append only

Advantages

- 1- More compact than Text files
- 2- Provide support for compression at diff level - Block/Record
- 3- Can be split and processed in parallel
- 4- Solve large number of small file problem.
- 5- Temp output of mapper can be stored in Seq file

Writing a sequence file

```

Configuration conf = new Configuration();
Filesystem fs = filesystem.get(URI.create(uri), conf);
Path path = new Path(uri);
InetSocketAddress key = new InetSocketAddress();
Text value = new Text();
SequenceFile.Writer writer = null;
try {
    writer = SequenceFile.createWriter(fs, conf, Bytes,
        Key.serializer(), value.serializer());
}

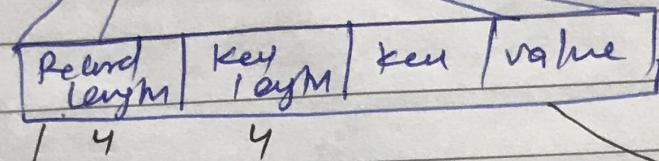
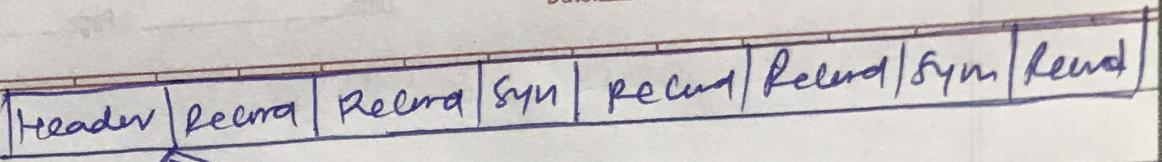
```

Reading a Sequence file

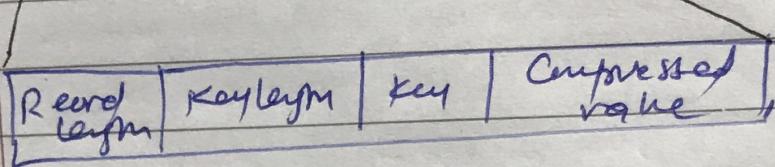
```

Configuration conf = new Configuration();
Filesystem fs = filesystem.get(URI.create(uri), conf);
Path path = new Path(uri);
SequenceFile.Reader reader = null;
try {
    reader = new SequenceFile.Reader(fs, path, conf);
    Writable key = (Writable) ReflectionUtils.newInstance(
        (reader.getKeyClass()), conf);
    Writable value = (Writable) ReflectionUtils.newInstance(
        (reader.getValueClass()), conf);
}

```

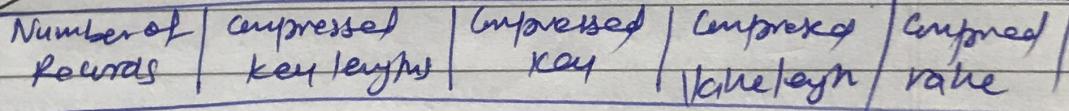
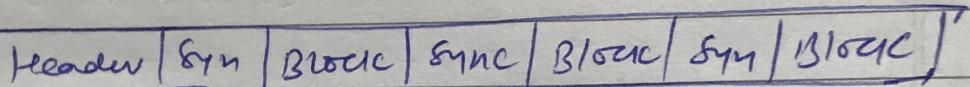


NO compression



Record compression

Sequence file with no compression and with Record compression



Block of compression

Seq file with block compression

Sequence file Programs