

# Hadoop & MapReduce Cheat Sheet

G.B. - v 0.0.1

---

## 1. Key Definitions

- **Big Data (3Vs):** Characterized by **Volume** (scale of data), **Variety** (different forms of data), and **Velocity**.
- **Apache Hadoop:** A scalable, fault-tolerant distributed system for Big Data. It moves the code/computation to the data (data locality) rather than moving data across the network.
- **HDFS (Hadoop Distributed File System):** Distributed storage system that provides global file namespace and stores data redundantly (typically 3 replicas) to ensure persistence and availability.
- **MapReduce Paradigm:** A programming model that separates the *what* (Map and Reduce functions defined by the developer) from the *how* (scheduling, synchronization, and fault tolerance handled by the framework).
- **Shuffle & Sort:** Phase between Map and Reduce where:
  - Data is **partitioned** across reducers
  - Keys are **sorted**
  - Values are **grouped by key**
- **Combiner:** A "mini-reducer" executed **LOCALLY** on the Mapper's output to pre-aggregate data and reduce network traffic during the shuffle and sort phase. It is not guaranteed to execute and it must be **associative** and **commutative**.

## 2. Architecture: Blocks, Mappers, and Reducers

- **Block Size & Splits:** Files in HDFS are split into fixed-size **blocks/chunks** (typically 64-128MB). An **Input Split** is a logical representation of data to be processed. Usually, one split corresponds to one HDFS block.
- **Number of Mappers:** Automatically determined by the Hadoop framework. **One Mapper task is instantiated for each Input Split.** You cannot arbitrarily set the exact number of Mappers. You have **at least one mapper** per input file. Map is called once per line in input file.
- **Number of Reducers:** This is **settable by the user** in the Driver class (e.g., `job.setNumReduceTasks(n)`).
- **Reducer Input:** The framework automatically groups all Mapper outputs by key (Shuffle and Sort phase). The Reducer receives a single key and an **Iterable list of all values** associated with that key: `key : [<value1>, <value2>, ...]`.

## 3. File Input Formats: Standard vs KeyValue

- **TextInputFormat (Standard):** Used for plain text files. Files are broken into lines.
  - **Key:** Position (offset) of the line in the file (`LongWritable`).
  - **Value:** The content of the entire line (`Text`).
- **KeyValueTextInputFormat:** Used for plain text files where each line follows a `key<separator>value` structure (default separator is the TAB character `\t`).
  - **Key:** The text preceding the separator (`Text`).
  - **Value:** The text following the separator (`Text`).

## 4. Counters

Counters are a useful channel for gathering statistics about the job (e.g., counting invalid records, specific events, etc.).

- **How they work:** They are global variables managed by the Hadoop framework. Every Mapper/Reducer can increment a counter, and the framework safely aggregates these increments.
- **Usage:**
  1. Define them using a Java `enum`.
  2. Increment them inside the Mapper or Reducer using `context.getCounter(EnumName).increment(1)`.
  3. Retrieve the final value in the Driver (you cannot see its content otherwise) **after** the job has completed successfully.

## 5. Complete Java Implementation

### DRIVER CLASS (Universal)

```
1 import org.apache.hadoop.conf.*;
2 import org.apache.hadoop.fs.Path;
3 import org.apache.hadoop.io.*;
4 import org.apache.hadoop.mapreduce.*;
5 import org.apache.hadoop.mapreduce.lib.input.*;
6 import org.apache.hadoop.mapreduce.lib.output.*;
7 import org.apache.hadoop.util.*;
8
9
10
11 public class MyDriver extends Configured implements Tool {
12     // [MODIFIABLE] Define Global Counters
13     public static enum COUNTERS {
14         INVALID_RECORDS, SPECIAL_EVENTS }
15     public int run(String[] args) throws Exception {
16         Configuration conf = this.getConf();
17
18         // --- PASSING PARAMETERS TO MAP/REDUCE ---
19         // [MODIFIABLE] Set custom parameters here BEFORE Job.getInstance()
20         conf.set("my.custom.threshold", args[3]);
21
22         Job job = Job.getInstance(conf, "My Job Name");
23         // JOIN MAP-SIDE Add hdfs file businessrules.txt in the distributed cache
24         //job.addCacheFile(new Path("businessrules.txt").toUri());
25         job.setJarByClass(MyDriver.class);
26
27         // [MODIFIABLE] Input / Output paths
28         FileInputFormat.addInputPath(job, new Path(args[0]));
29         FileOutputFormat.setOutputPath(job, new Path(args[1]));
30
31         // [MODIFIABLE] Input / Output formats
32         job.setInputFormatClass(TextInputFormat.class);
33         job.setOutputFormatClass(TextOutputFormat.class);
34
35         // --- MAPPER SETUP ---
36         job.setMapperClass(MyMapper.class);
37         job.setMapOutputKeyClass(Text.class);
38         job.setMapOutputValueClass(IntWritable.class);
39
40         // --- REDUCER SETUP ---
41         // [MODIFIABLE] Comment this line out if you want a MAP-ONLY job
42         job.setReducerClass(MyReducer.class);
43
44         // --- COMBINER SETUP ---
45         // [MODIFIABLE] Uncomment to USE a Combiner
46         // job.setCombinerClass(MyReducer.class);
47
48         // --- NUM REDUCE TASKS (MAP-ONLY SWITCH) ---
49         // [MODIFIABLE] Set to 0 for MAP-ONLY
50         job.setNumReduceTasks(Integer.parseInt(args[2]));
51
52         // --- FINAL OUTPUT TYPES ---
53         // NOTE: If Map-only, these MUST match the MAPPER'S output classes!
54         job.setOutputKeyClass(Text.class);
55         job.setOutputValueClass(IntWritable.class);
56
57         // [FIXED] Run job and wait for completion
58         boolean success = job.waitForCompletion(true);
59
60         // --- READ COUNTERS AFTER JOB COMPLETION ---
61         if (success) {
62             long invalidCount = job.getCounters()
63                 .findCounter(COUNTERS.INVALID_RECORDS)
64                 .getValue();
65             System.out.println("Total invalid records: " + invalidCount);
66         }
67
68         return success ? 0 : 1;
69     }
70
71     public static void main(String[] args) throws Exception {
72         int res = ToolRunner.run(new Configuration(), new MyDriver(), args);
73         System.exit(res);
74     }
75 }
```

## MAPPER CLASS

```
1 import java.io.IOException;
2 import org.apache.hadoop.io.*;
3 import org.apache.hadoop.mapreduce.Mapper;
4
5 // [MODIFIABLE] Generics: <InputKey, InputVal, OutputKey, OutputVal>
6 class MyMapper extends Mapper<LongWritable, Text, Text, IntWritable> {
7     private int threshold; // Variable to store the shared parameter
8     private ArrayList<String> rules;
9     // [OPTIONAL] setup() is executed ONLY ONCE per Mapper -instance-, not per line (HashMap!)
10    protected void setup(Context context) {
11        // [MODIFIABLE] Retrieve the shared parameter. Default value is 0.
12        threshold = context.getConfiguration().getInt("my.custom.threshold", 0);
13        String nextLine;
14        rules = new ArrayList<String>();
15        // Open the business rules file (that is shared by means of the distributed cache mechanism)
16        URI[] CachedFiles = context.getCacheFiles();
17        // This application has one single cached file.
18        // Its path is URIsCachedFiles[0]
19        BufferedReader rulesFile = new BufferedReader(new FileReader(new File(CachedFiles[0].getPath())));
20        // Each line of the file contains one rule
21        while ((nextLine = rulesFile.readLine()) != null) {
22            rules.add(nextLine);
23        }
24        rulesFile.close();
25    }
26
27    protected void map(LongWritable key, Text value, Context context) throws IOException,
28        InterruptedException {
29        String line = value.toString();
30        // [MODIFIABLE] Counter usage: check for invalid/empty records
31        if (line.trim().isEmpty()) {
32            context.getCounter(COUNTERS.INVALID_RECORDS).increment(1);
33            return; // Skip further processing for this record
34        }
35        // [MODIFIABLE] Standard mapping logic
36        String[] words = line.split("\\s+");
37        for (String word : words) {
38            context.write(new Text(word), new IntWritable(1));
39        }
40    }
41    // [OPTIONAL] cleanup() is executed ONLY ONCE per Mapper after all records are processed
42    protected void cleanup(Context context) throws IOException, InterruptedException {
43        // [MODIFIABLE] Emit results based on the in-mapper local variables
44        // for (Map.Entry<String, Integer> entry : localMap.entrySet()) {
45        //     context.write(new Text(entry.getKey()), new IntWritable(entry.getValue()));
46        // }
47    }
48 }
```

## REDUCER CLASS (Can act as Combiner if associative and commutative)

```
1 import java.io.IOException;
2 import org.apache.hadoop.io.*;
3 import org.apache.hadoop.mapreduce.Reducer;
4 // [MODIFIABLE] Generics: <InputKey, InputVal, OutputKey, OutputVal>
5 class MyReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
6     // [MODIFIABLE] Local variable to maintain state across MULTIPLE keys processed by this Reducer
7     // private int totalSumAcrossAllKeys = 0;
8     // [OPTIONAL] setup() is executed ONLY ONCE per Reducer instance
9     protected void setup(Context context) {
10        // You can read parameters here just like in the Mapper
11        // int someParam = context.getConfiguration().getInt("my.param", 0);
12    }
13    // [FIXED] reduce() signature. Notice the Iterable for the values!
14    protected void reduce(Text key, Iterable<IntWritable> values, Context context) throws IOException,
15        InterruptedException {
16        int sum = 0;
17
18        // [MODIFIABLE] Custom logic: iterate over the list of values
19        for (IntWritable val : values) {
20            sum += val.get();
21        }
22
23        // [MODIFIABLE] Emit the final aggregated (Key, Value) pair
24        context.write(key, new IntWritable(sum));
25    }
26
27    // [OPTIONAL] cleanup() is executed ONLY ONCE per Reducer after ALL keys have been processed
28    protected void cleanup(Context context) throws IOException, InterruptedException {
29        // [MODIFIABLE] Emit final global results based on Reducer's local variables
30        // context.write(new Text("TOTAL_PROCESSED_BY_THIS_REDUCER"), new IntWritable(totalSumAcrossAllKeys));
31    }
32 }
```

## CUSTOM DATA CLASS (Personalized Data Type)

```
1 import java.io.*;
2 import org.apache.hadoop.io.Writable;
3
4 // [FIXED] Complex Values must implement Writable.
5 // If used as a Key, it must implement WritableComparable<MyCustomData>
6 public class MyCustomData implements Writable {
7
8     // [MODIFIABLE] Custom internal state
9     private float sum = 0;
10    private int count = 0;
11
12    // [FIXED] A default empty constructor is strictly required by Hadoop
13    public MyCustomData() {}
14
15    public MyCustomData(float sum, int count) {
16        this.sum = sum;
17        this.count = count;
18    }
19
20    // [FIXED] Serialization method (order matters)
21    @Override
22    public void write(DataOutput out) throws IOException {
23        out.writeFloat(sum);
24        out.writeInt(count);
25    }
26
27    // [FIXED] Deserialization method (order matters)
28    @Override
29    public void readFields(DataInput in) throws IOException {
30        sum = in.readFloat();
31        count = in.readInt();
32    }
33
34    public float getSum() { return sum; }
35    public int getCount() { return count; }
36
37    @Override
38    public String toString() { return "Sum: " + sum + ", Count: " + count; }
39 }
```

## 6. MapReduce Design Patterns

| Pattern                         | Examples                                     | Tasks                            | Implementation & Management  |
|---------------------------------|--|----------------------------------|--|
| <b>Numerical Summarizations</b> | Word count, Min/Max/Avg                      | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> Emits ( <code>key</code> , <code>summary_field</code> ).<br><b>Reduce:</b> Computes the final aggregate statistics.<br><b>Comb.:</b> usable if the task is commutative and associative.  |
| <b>Inverted Index</b>           | Web search engines (Word $\rightarrow$ URLs) | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> Emits ( <code>keyword</code> , <code>record_ID</code> ).<br><b>Reduce:</b> Receives IDs and concatenates them into a list.<br><b>Combiner:</b> Usually not useful (no numerical aggregation).  |
| <b>Counting with Counters</b>   | Count records, track specific events         | $N$ Mappers<br><b>0</b> Reducers | <b>Map-Only Job:</b> Mappers evaluate records and simply increment global <b>Counters</b> (no output emitted). Only driver can read and prints results after completion.   |
| <b>Filtering</b>                | Distributed grep, Data cleaning              | $N$ Mappers<br><b>0</b> Reducers | <b>Map-Only Job:</b> Mappers evaluate a condition/rule. Emits the record ( <code>key</code> , <code>record</code> ) only if the condition is met   |
| <b>Top K</b>                    | Outliers analysis,s                          | $N$ Mappers<br><b>1</b> Reducer  | <b>Map:</b> Uses <code>setup()</code> to init a local list, updates it in <code>map()</code> , and emits local Top $K$ in <code>cleanup()</code> with a shared "null" key.<br><b>Reduce:</b> Merges all local lists to find the global Top $K$ .         |
| <b>Distinct</b>                 | Duplicate removal, unique values             | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> Emits the record as key ( <code>record</code> , <code>null</code> ).<br><b>Reduce:</b> Framework groups identical keys; reducer emits ( <code>key</code> , <code>null</code> ) exactly once for each group.                                  |
| <b>Shuffling</b>                | Data anonymization, shuffling                | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> Emits ( <code>random_key</code> , <code>input_record</code> ) .<br><b>Reduce:</b> Emits ( <code>input_record</code> , <code>null</code> ) to randomize order   |
| <b>Job Chaining</b>             | Complex workflows                            | Multiple Jobs                    | <b>Driver:</b> Manages workflow by executing jobs in order . Output of phase/Job $i$ becomes the input for phase/Job $i + 1$ .   |
| <b>Reduce-Side Join</b>         | Natural/Theta/Join on 2 large tables         | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> One Mapper class per table. Emits ( <code>common_attr</code> , <code>table_name + record</code> ) .<br><b>Reduce:</b> Gathers records for a key, computes the "local join". Careful, <b>you must use a list usually (not the iterable)</b> . |
| <b>Map-Side Join</b>            | Join a LARGE table with a SMALL table        | $N$ Mappers<br><b>0</b> Reducers | <b>Map-Only Job:</b> Small table is loaded into main memory via <b>Distributed Cache</b> in the <code>setup()</code> method . Mapper performs local join on the fly for each record of the large table   |

## 7. MapReduce SQL Operators

| Operator                           | Description   | Tasks                            | Implementation & Management  |
|------------------------------------|---|----------------------------------|--|
| <b>Selection</b> ( $\sigma$ )      | Applies a predicate $\forall$ record.                             | $N$ Mappers<br><b>0</b> Reducers | <b>Map-Only Job:</b> Uses the <i>Filtering pattern</i> . If record $r$ satisfies condition $C$ , emit ( <code>r</code> , <code>null</code> ). Otherwise, discard .   |
| <b>Projection</b> ( $\Pi$ )        | Keeps only a subset of attributes and <b>removes duplicates</b> . | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> Selects subset of attributes to form a new record $r'$ , emits ( <code>r'</code> , <code>null</code> ) .<br><b>Reduce:</b> Emits ( <code>r'</code> , <code>null</code> ) exactly once for each group (removes duplicates) .                    |
| <b>Union</b> ( $R \cup S$ )        | Combines two relations with the same schema.                      | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> For each input record $t$ in $R$ or $S$ , emits ( <code>t</code> , <code>null</code> ).<br><b>Reduce:</b> Emits ( <code>t</code> , <code>null</code> ) exactly once for each group to eliminate duplicated records .                           |
| <b>Intersection</b> ( $R \cap S$ ) | Returns records appearing in BOTH relations.                      | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> For $t \in R$ , emits ( <code>t</code> , "R"). For $t \in S$ , emits ( <code>t</code> , "S").<br><b>Reduce:</b> Emits ( <code>t</code> , <code>null</code> ) ONLY if the list of values contains BOTH "R" and "S" (size is 2)                  |
| <b>Difference</b> ( $R - S$ )      | Returns records appearing in R but NOT in S.                      | $N$ Mappers<br>$M$ Reducers      | <b>Map:</b> Needs 2 Mapper classes. For $t \in R$ , emits ( <code>t</code> , "R"). For $t \in S$ , emits ( <code>t</code> , "S") .<br><b>Reduce:</b> Emits ( <code>t</code> , <code>null</code> ) ONLY if the list of values contains ONLY the value "R" . |
| <b>Join &amp; Group By</b>         | Table joins and aggregations.                                     | Varies                           | Implemented using <b>Join Patterns</b> (Map/Reduce side) or <b>Summarization Pattern</b> respectively .  |