# MapReduce Tutorial.

## I. Introduction.

### ☑ Why MapReduce?

Traditional Enterprise System

$$ \boxed{\text{User}} \longleftrightarrow \boxed{\begin{array}{c}\text{Centralised}\\\text{System}\end{array}} \longleftrightarrow \boxed{\begin{array}{c}\text{Relational}\\\text{Database}\end{array}} $$
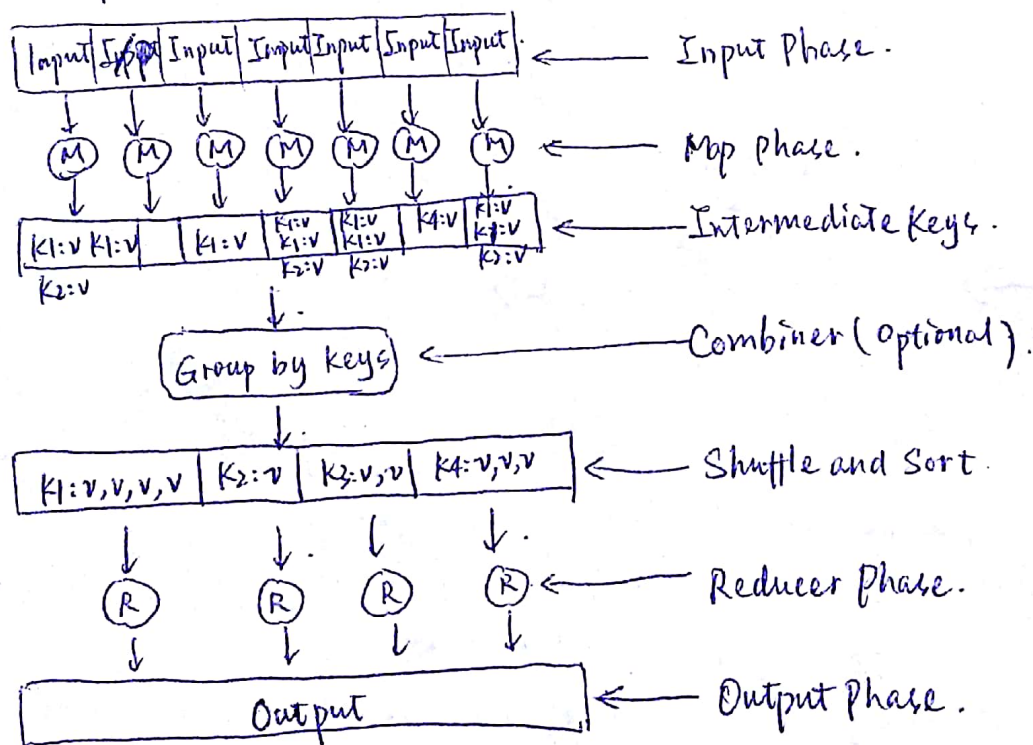
Limitation : Centralised system creates too much of a bottleneck while processing multiple files simultaneously.

Solution: MapReduce Algorithm.

MR. divides a task into small parts and assign them to many computers. Results are collected at one place and integrated to form the result dataset.

### ☑ How MapReduce works?



Input Phase: Record Reader. parses data to the mapper in <key, value>

Map : a user-defined function, <key, value> ⟶ new <key, value>

Combiner: In a mapper, combiner is a local reducer.
intermediate <key, value> ⟶ combiner does aggregation.

**Shuffle and Sort**: Reducer task starts here.
Grouped <k, v> pairs, ⟶ local machine (Reducer here).
<k, v> sorted by key ⟶ data list.

Reducer : grouped. key-value paired data ~~as input~~.
↓
Reducer function on each of them.
aggregated, filtered, combined.

new $<k,v>$ pairs.

Output phase: write $<k,v>$ pairs onto a file.

Example:



Map Phase. Shuffle & Sort — Reduce phase.

Split:
- ABR → A,1 / B,1 / R,1 → A,1 / A,1
- CCR → C,1 / C,1 / R,1 → B,1 / B,1
- ACB → A,1 / C,1 / B,1 → C,1 / C,1 / C,1 → R,1 / R,1

Input:
ABR
CCR
ACB

Reduce:
A,2
B,2
C,3
R,2

## MapReduce Example (twitter).



Input — Twitter Data (tweets)

MapReduce: TOKENIZE | FILTER | COUNTER | AGGREGATE COUNTERS

Data Source Adapter — Hadoop Related Databases

Tokenize : tokenizes tweets into maps of tokens. $<k,v>$ pairs
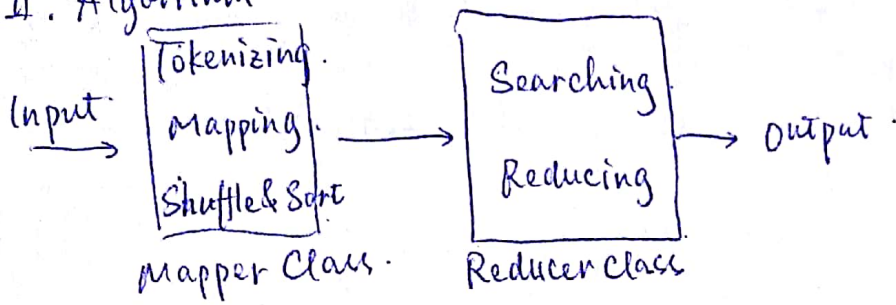Filter : filters unwanted words, writes filtered words → $<k,v>$.
Count : generates a token counter per word.
Aggregate counters: prepare an aggregate of similar counter values.

# MapReduce Tutorial

## II. Algorithm

Input → | Tokenizing / Mapping / Shuffle & Sort | → | Searching / Reducing | → Output

Mapper Class.  Reducer Class

1. Sorting = output of mapper class by their keys.
   $\langle k, v \rangle$.

2. Searching : combiner phase & reducer phase.

3. Indexing : inverted index, TF.-IDF. (Term freq. - inverted doc. freq).

$$TF_{(the)} = \frac{\# \text{ of 'the' terms in the doc.}}{\# \text{ of terms in the doc.}}$$

$$IDF (the) = \log_e \frac{\text{total \# of documents}}{\# \text{ of documents with term 'the'}}$$

## III Installation (skip).

## IV. API.

1. JobContext Interface : • super interface for all classes.
   • defines different jobs in mapreduce.

   Sub interfaces :
   • MapContext $\langle$ KEYIN, VALUEIN, KEYOUT, VALUEOUT $\rangle$.
   • ReduceContext $\langle \cdots \rangle$.

2. Job Class :
   • most important API. in MapReduce.
   • configure the job, submit it, control its execution, query the state.

~~. Constructors of Job Class.~~

3. Mapper Class.
   • defines map job.  input $\langle key, value \rangle \longrightarrow$ intermediate $\langle key, value \rangle$

4. Reducer Class :
   • defines reduce job.
   • 3 primary phases
     (I) Shuffle : copies the sorted output from each mapper using HTTP across the network.
     (II) Sort : merge-sorts reducer inputs by keys.
           Shuffle and sort phases occur simultaneously.

(3). Reduce : In this phase, the reduce (object, iterable, context).
method is called for each <key, (collection of values)> in the
sorted inputs.
- reduce method is called once for each key value.

## V. Hadoop Implementation.

### 1. MapReduce Algorithm.
- Hadoop sends Map and Reduce tasks to appropriate servers in the cluster
- most computing couple takes place on nodes with data locally.
- after computation completes, collect to form appropriate result and sends it
  back to the Hadoop server.

### 2. Input and Outputs.
- <k,v> pairs as input and output.
- Serializable → Writable interface.
- key classes implement the <u>Writable Comparable interface</u> to facilitate
  sorting.

Input     intermediate

$<k_1, v_1> \rightarrow map \rightarrow <k_2, v_2> \rightarrow reduce$

output $<k_3, v_3>$

### 3. MapReduce Implementation. (skip).

## VI. Partitioner.
Properties :
- takes place ~~before~~ after ~~Mapper class~~ Map phase and before Reduce phase.
- # partitioner = # reducer, determine # partitioner according to # reducer.
- partitions intermediate <k,v> from mapper.
- works like a hash fn.

## VII. Combiners.
Properties :
(1) Semi-reducer.
(2) Summarize map output records with the same key.
(3)* reduce the volume of data-transfer between map and reduce.
    Note: usually map output volume is high.

How Combiner works?
(1) no predefined interface. implement the reducer interface's reduce()
    method.
(2) produce summary info. from a large dataset because
    it replaces the original Map out.

# MapReduce Tutorial

## VII. Combiner (con't).

How combiner works?

(3) usually the code and operation for a combiner is similar to that of a reducer.

*(4) The combiner phase takes each key-value pair from the Map phase, processes it, and produces the output as, <u>key-value collection pairs</u>.

e.g. wordcount.java.