# Map Reduce.
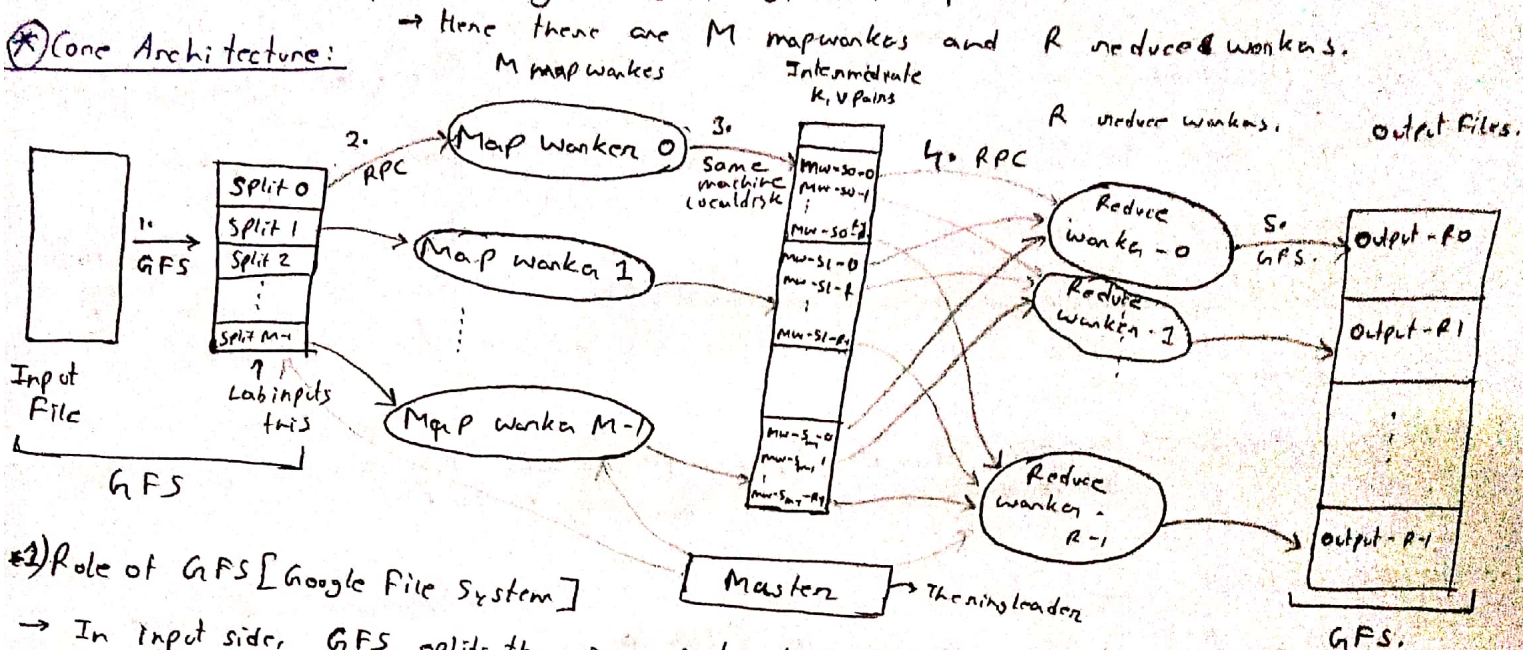
→ It is a batch processing model. [ All maps must be done before reduce starts]

→ Used to compute large amounts of data parallely.

→ Here there are M mapworkers and R reduced workers.

* Core Architecture:



*1) Role of GFS [Google File System]

→ In input side, GFS splits the given into M splits each at size 64 MB. It also replicates a split into 3 more machines for fault tolerance.

→ In output side, GFS the reduce workers writes its output to R separate files mansed by GFS [ 1 output for each Reduce worker]

2) The master knows the location of the splits. When a map worker is idle it gives the map worker a map task. In a map task, the worker take a single input split split passes it and passes it to the Ke user defined Map function where $k = \{$ SplitNumber $\}$ not imp $\quad v = \{$ The contents of the file $\}$.

3. Each map task produces R intermediate files. The logic is to use a partition key for Ex: hash $\left(\begin{smallmatrix} intermediate \\ key \end{smallmatrix}\right)$ % R $\subset$ ⓪→ goes to mw-so-⓪  ②→ goes to mw-so-② . The files The key is not written to file every time as it is inefficient insted it is written to buffers and buffers are periodically flushed to files. It then informs the master the location of R intermediate output files. { The file is on the same computer where the worker runs.

4. The master informs the reduce workers about the locations so they can start working on it. Ex: Since it is already partitoned, Reduce worker O take the intermediate files mw-*-0.txt and process it. Then it sorts the data using given map and user given reduce fn is run where for each key k.1 value is the cumulative sum of all intermediate values. [ for word counting ]

5. The output files from Reduce workers are initially temporary files. Example name output-0-temp.txt and when processing is completed it is renamed to output-0.txt.

## * Master Data Structure:

→ For Each map & reduce task, Stores the

→ we need to detect which task failed & reschedule it.

1. State (idle, in progress, completed)

2. Identity of worker [for non idle tasks]

3. For completed map tasks stores the location of R intermediate files.

## * Fault Tolerance:

→ Master failure : → we can do checkpoints periodically but since this is only one task process we can manually restart it.

→ worker failure:

  ↳ Master Pings worker periodically [Lab 10s]
    ↳ If no response is received from a worker, master worker fails.

Response → All Good.

→ The map tasks completed or inprogress by the worker are reset to 'initial' state and are elisible for rescheduling on other worker.

→ The reduce tasks that are in progress by the worker is reset to initial state and is eligible for rescheduling.

why: Completed map tasks are stored in same machine' as the worker [localdisk] so we need to redo that since disk is inaccessible.

Completed reduce tasks do not need rescheduling since the output is or file system.

→ when a map task is executed first by worker A and later by worker B [cuz A failed] all workers executing the reduce tasks are notified of re execution. Any reduce task that has not already read the data from worker A will now read from worker B.

→ Our system should produce same output as would have been produced by a sequential execution.

→ we rely on atomic commits to achieve this property.

→ Each in progress task writes to its temp files. A map task produces R such files and Reduce tasks one such file.
  ↳ per reduce task.

→ When map task completes it sends its R output message file location & name. to the master. If master receives completion message for a already complete task, it ignores the message else it stores the location in its DS.
  ↳ master

→ When a reduce task is complete, the worker renames its temporary output to final output file. [ Rename is atomic ]

* How Network Bandwidth problem is solved?

→ In 1 we said that GFS stores 3 copies for each split so we run our map task on the machines where the copies are present on if it fails, on same network. This minimizes the network Bandwidth. [ Read is local ]

* Task Granularity :

→ # of Map task & # of Reduce task should be larger than # of worker machines.

why :  i) Faster worker perform more tasks, preventing idle state.
      ii) If a worker fails, its tasks can be spread across all other workers, speeding up re execution.

* StragglerProblem:

→ when the job is almost done, schedule "backups" copies of remaining in-progress tasks. Mark the task done when either the primary or backup finishes.

→  Stragglers ⌒(machines with bad disk) can delay completion so backup
              on high CPU load                        reduces completion time.

* Implementation Note:
We have written master assigns task, which is push model → Hard to do.
Instead we will do pull model, workers ask for task periodically → By RPC.