

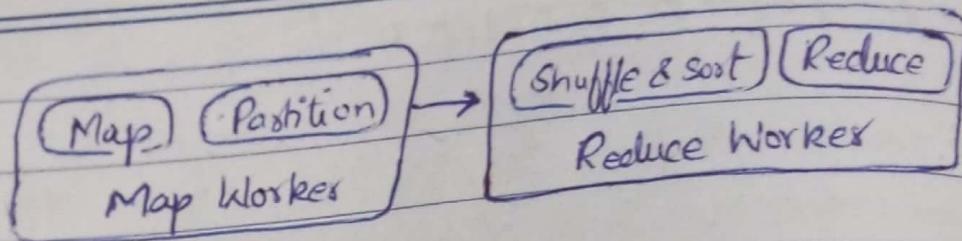
"MapReduce"

- Map reduce performs the processing and generation of large data sets in a distributed and parallel manner.
- MapReduce consist of two distinct task - Map and Reduce.
- Two essential daemons of MapReduce: Job Tracker & Task Tracker
 - it is like SlaveNode
- Map and Reduce work on divide & conquer approach.
 - divide
 - combines

divide the data into key & value pairs

"MapReduce FRAMEWORK"

- It consist of two key functions
 - 1. Map: Filters and sorts data into intermediate key-val pairs
 - 2. Reduce: Aggregates the intermediate data into smaller key-val pairs or final output.



- Map:

↳ Grab the relevant data from the source

↳ Parse into key-val

↳ Write it into an intermediate file

- Partition:

↳ Partitioning: identify which of R reducers will handle which keys

↳ Map partitions data to target it to one of R

Reduce workers based on a partitioning function

(both R & Partitioning function user defined)

Shuffle & Sort:

↳ Shuffle & Sort: Fetch the relevant partition of the output from all mappers

↳ Sort by keys (different mappers may have sent data with the same key)

- Reduce:

↳ Input is the sorted output of mappers

↳ Call the user Reduce function per key with the list of values for that key to aggregate the results.

Day

Date Full Process of the Framework:

- Step #01:
 - Split Input files into chunks (shards)
 - Shard 0 | shard 1 | shard 2 | | shard M-1

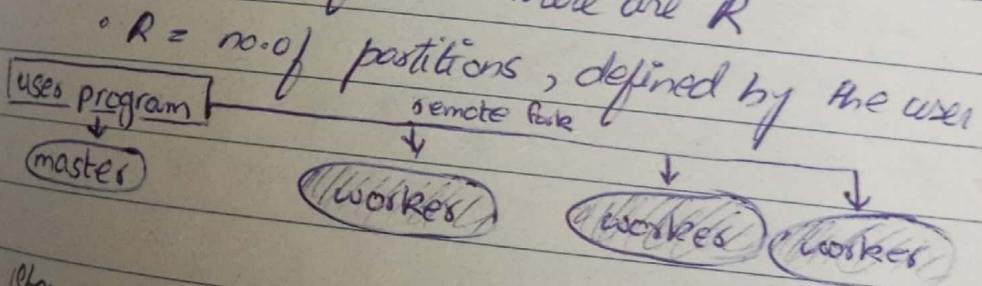
- Step #02: Fork Processes

- Start up many copies of the program on a cluster of machines

- One Master : scheduler & coordinator
 - Lots Workers

- Idle workers are assigned either:

- map tasks (each works on shard) - there are M map tasks
 - reduce tasks (each works on intermediate files) - there are R



- Step #03: Run Map Tasks

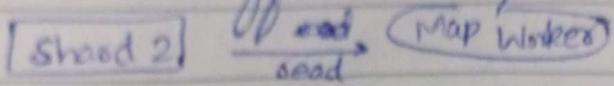
- Reads contents of the input shard assigned to it

- Parses key-val pairs out of input data.

- Parses each pair to the user defined map function

- Produces intermediate key-val pairs

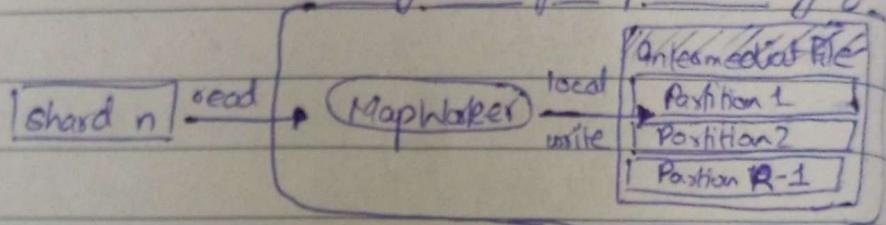
- These are buffer in memory



◆ Step # 04: Create Intermediate Files

- Intermediate key-val pairs produced by the user's map function buffered in memory and are periodically written to the local disk.

- Positioned into R regions by a partitioning fn.



◆ Step # 04(a): Partitioning :

- Map data will be processed by Reduce workers.
 - User's reduce function will be called once per unique key generated by Map
- We need to 1st sort all the key-val data by keys and decide which Reduce worker processes which keys
 - The reduce worker will do sorting

• Partition Function

Decide which of R reduce workers will work on which key

- Default function : $\text{hash}(\text{key}) \bmod R$
- Map workers partitions the data by keys.

Day

- Each Reduce worker will later read their partition from every map worker.

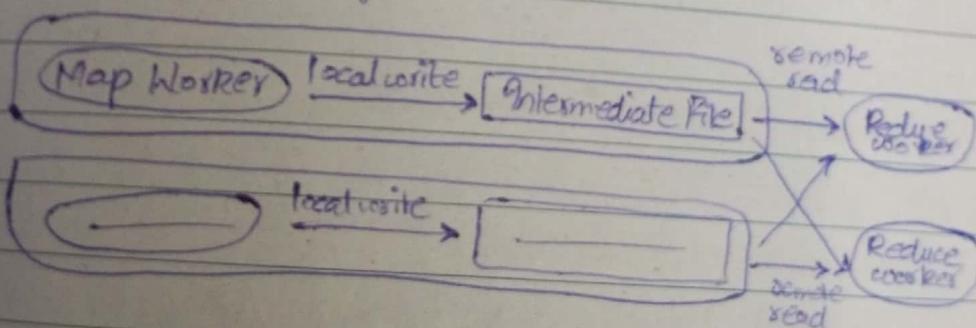
o Step #05 : Reduce Task Sorting

- Reduce workers gets notified by the Master about the location of intermediate files for its partition

Shuffle : Uses RPC's to read the data from the local disks of the map workers.

Sort : When the reduce worker reads intermediate data for its partition

- It sorts the data by the Intermediate keys
- All occurrences of the same key are grouped together.



o Step #06 : Reduce Task - Reduce

- The sort phase grouped data with a unique intermediate key

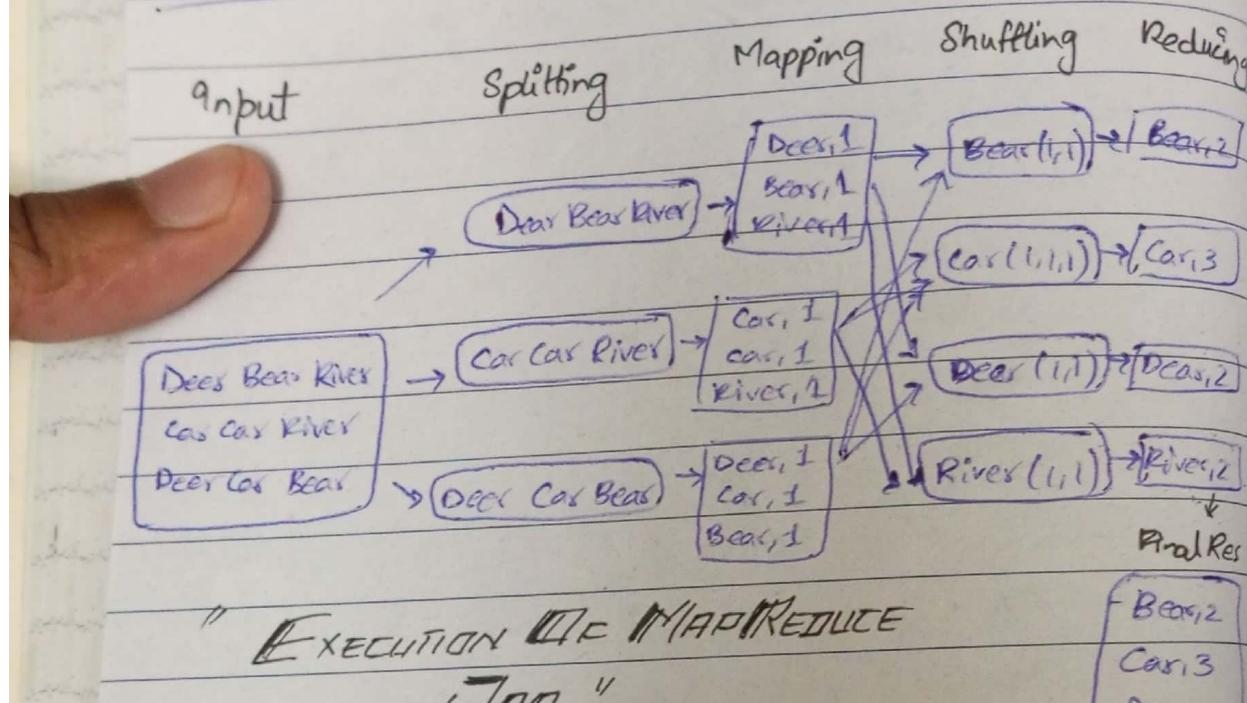
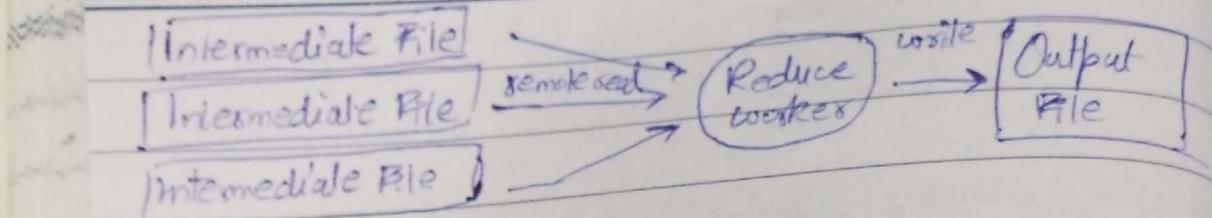
- User's Reduce func is given the key and preset of intermediate values for that key.

$\langle \text{key}, (\text{value}_1, \text{value}_2, \dots) \rangle$

- The output is appended to output file.

Date _____

Day _____



"EXECUTION OF MAPREDUCE JOB"

• MapReduce Cluster Components:

List (K₃, V₃)

a) Job Tracker (Master Node)

• Serves as the master node in a MapReduce cluster.

• Responsibilities:

- Initiates and assigns Map and Reduce tasks to worker nodes

- Monitors task progress and handles task failures.

b) Task Tracker (Worker Nodes)

- Each worker node runs a Task Tracker process
- Responsibilities :
 - Communicates with Job tracker to receive tasks and send status updates.
 - Executes tasks in an isolated environment using Java Virtual Machine (JVM) to ensure process isolation.

• Execution of a MapReduce Job

a) Job Submission :

1. The job tracker receives the job and reads input data from HDFS
2. The data is split into smaller chunks called splits (or shards)

b) Data Localization :

- The Job Tracker assigns shards to worker nodes, prioritizing data locality - preferably selecting nodes that already hold data.
- Benefits of Data localization :

- conserves network bandwidth
- Reduces data transfer time, improving job efficiency.

c) Task Execution:

- Worker Nodes execute tasks assigned by JT and periodically send progress updates via the task tracker
- Tasks are executed indep by each worker in isolated JVM's to prevent conflicts.

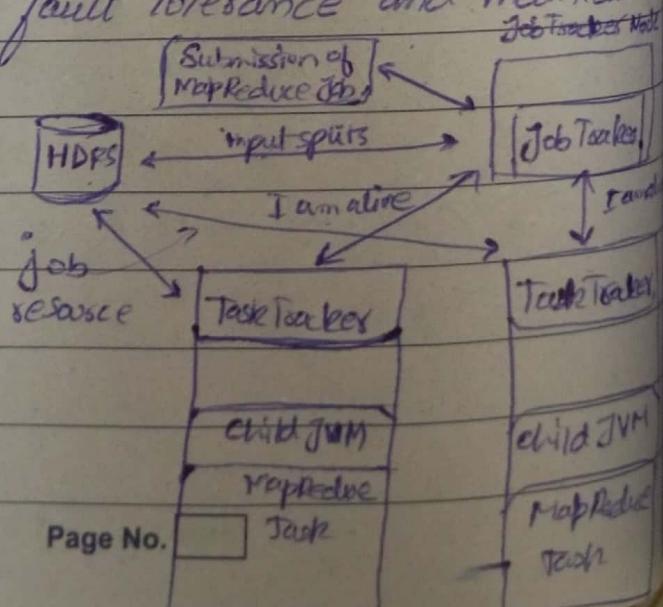
• Staggler Detection and Task Reassignment

a) Staggler Detection

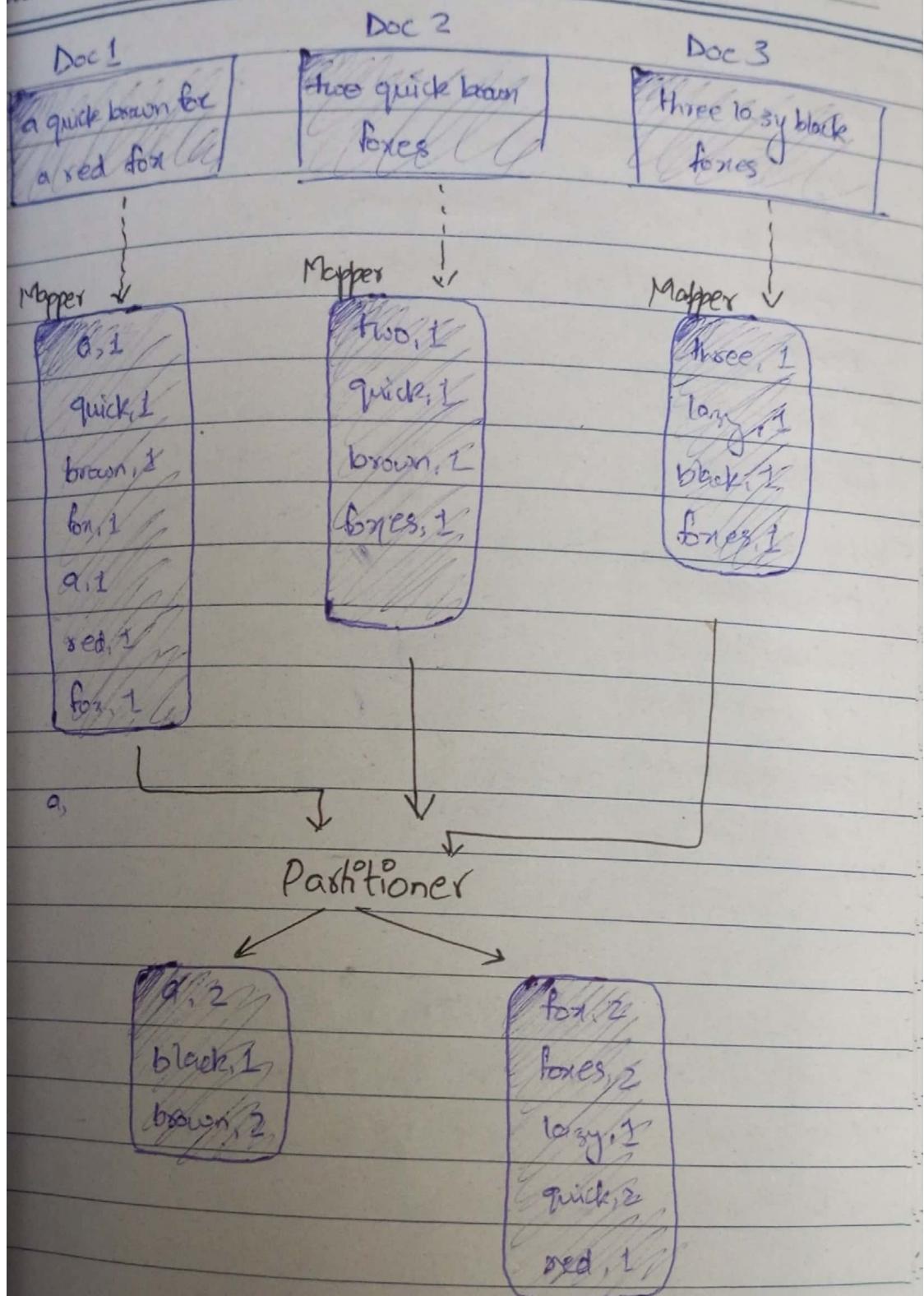
- Worker nodes regularly send updates to the JT
- If no update is received from a Task Tracker within a specified period, the JT assumes the worker has failed.

b) Task Reassignment:

- The job tracker reassigns the task to another available node, ensuring the task is completed without delays
- This process ensures fault tolerance and maintains job progress.



Day

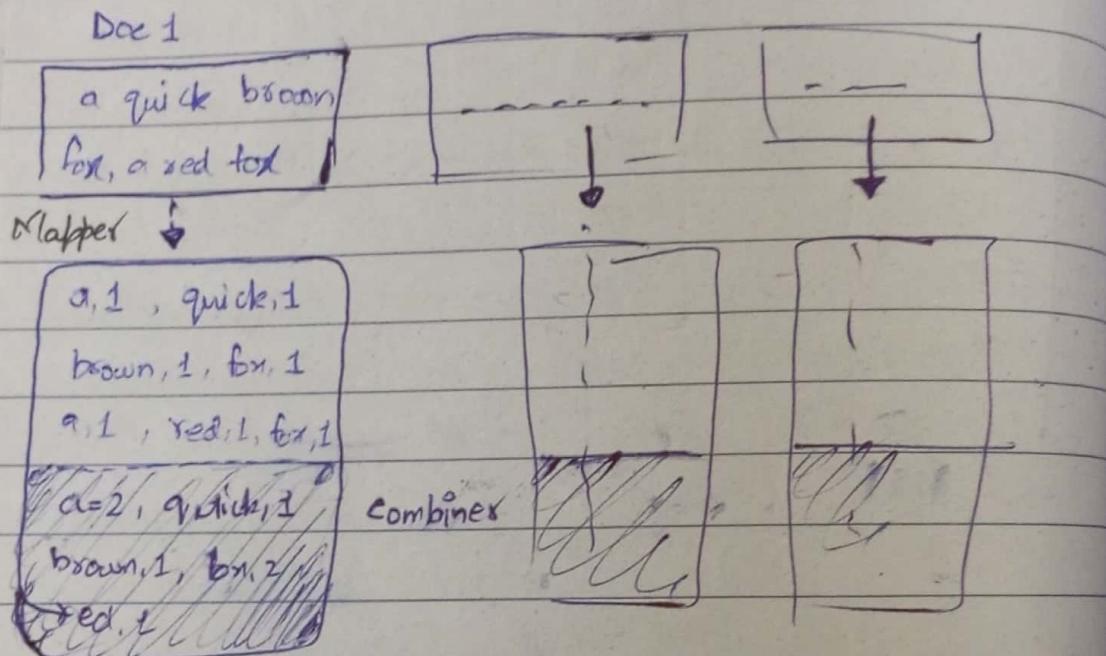


► Partitioner function for MapReduce

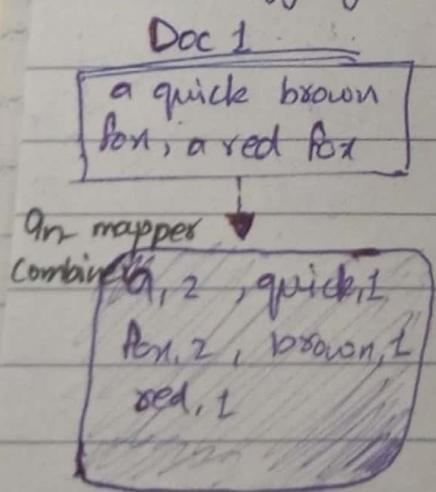
▷ Combined Function in MapReduce Job :

"The task of combiner func is to perform local reduction. That is reduce values with similar keys, locally at each mapper."

- Executed locally at each mapper

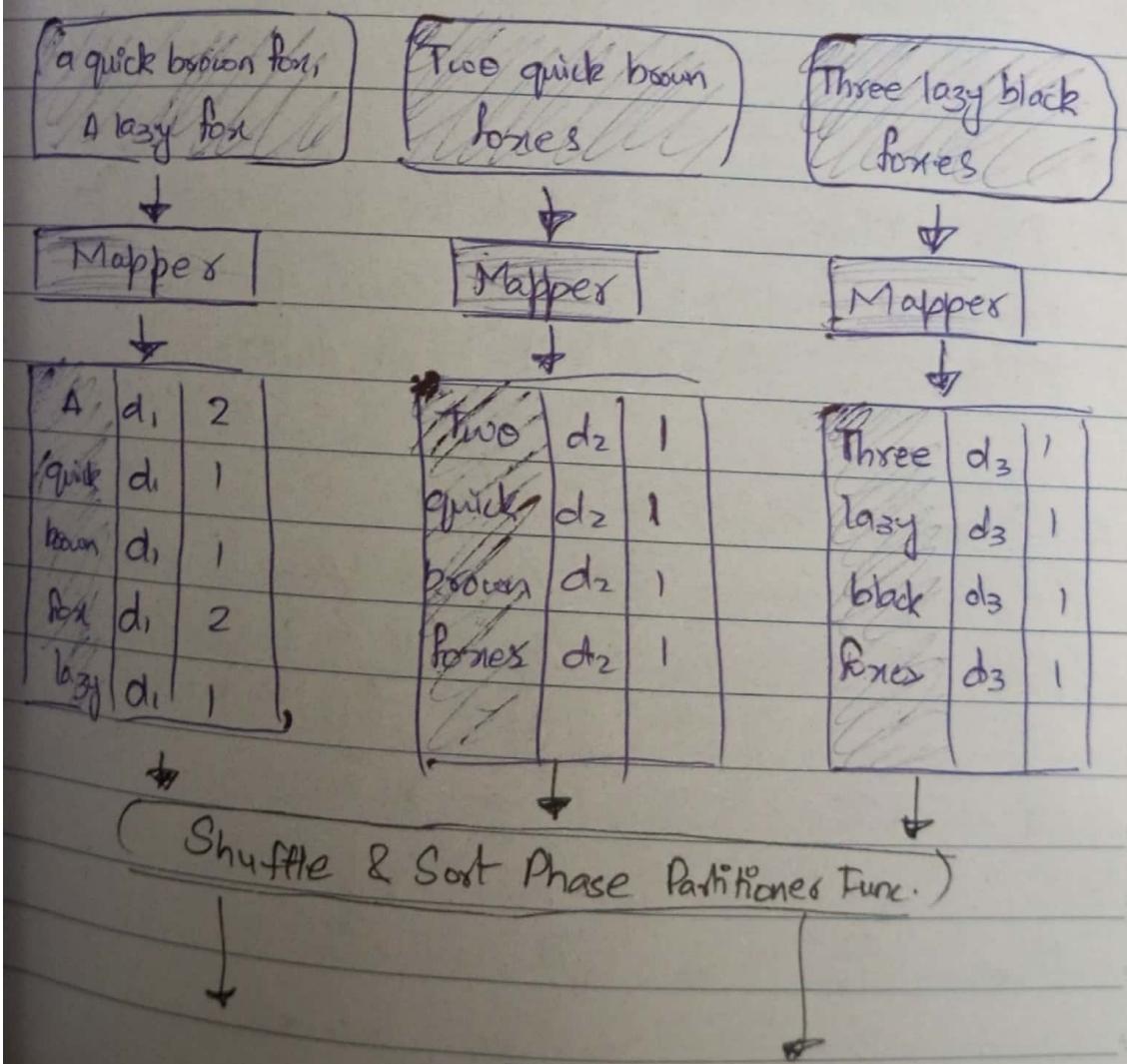


In-Mapper Combiner: It is different than the conventional combiner as the mapper function performs local aggregation before spilling the data.



ale
INVERTED INDEX :

- An inverted index is commonly used for MapReduce Analytics. Using this we can make a map of some terms to a list of items.
- It creates an index from the dataset to search for items in less time.
- It maps terms (words) to the locations (e.g. documents or files) where they appear, enabling efficient full text searches.



Date _____

Day _____

Reduces			Reduces		
↓			↓		
A	d ₁	2	dog	d ₁	1
black	d ₃	1	quick	d ₁	1
brown	d ₁	1	three	d ₃	1
fox	d ₁	2	tow	d ₂	1
foxes	d ₂	1	d ₃	2	

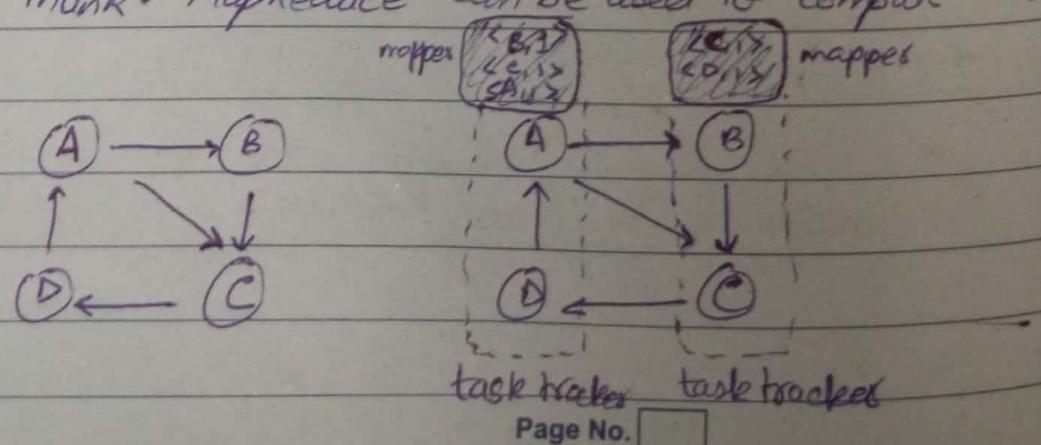
Each term is associated with a list of documents
Identifies where it appears.

Applications:

- ① Google, Bing, Elastic Search use inverted indexes to retrieve web pages matching search queries.
- ② For indexing field in text-based datasets.

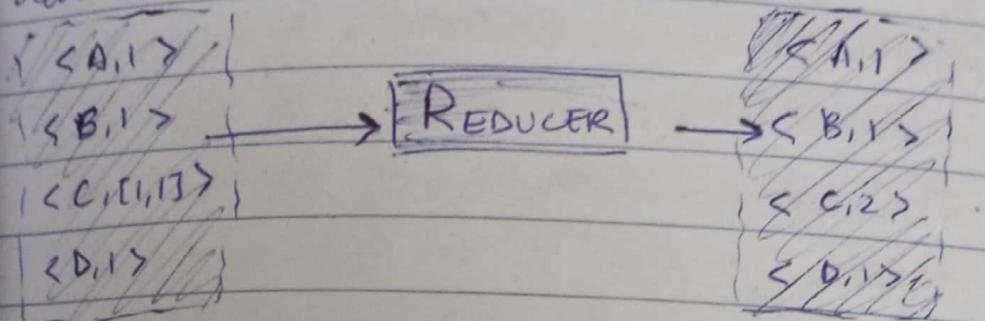
○ Computing Inlinks & Outlink:

The no. of incoming links to a webpage is called inlink. MapReduce can be used to compute inlinks.



Two mappers are created. Each mapper computes initial outlines for two pages.

After that the $\langle \text{key}, \text{value} \rangle$ pairs are sent to the reducer.



Map Reduce For Iterative Jobs :

1. Iterative Computation in MapReduce:

- Iterative algo (eg, PageRank, Clustering) involve repeating a computation process until a convergence point is reached.
- These algo req the output of a Reducer to be sent back as the input to the mapper in the next iteration.

2. Workflow for Iterative Jobs:

- Output of the MapReduce phase is written to HDFS.
- The Mapper in the next iteration reads the Reducer output from HDFS.

Date _____

Day _____

3. Inefficiency of MapReduce for Iterative Tasks.

- **HDFS Write Overhead:**
 - During the Reduce phase, output is written to HDFS, which involves
 - ↳ Transferring data over the network
 - ↳ Writing replicas for fault tolerance.
 - These operations significantly slow down the iterative process.

• **HDFS Read Overhead:**

- Each iteration requires reading back data from HDFS, introducing further delays.

4. Result:

- ① HDFS read/write ops are expensive, as they involve heavy I/O and network usage.
- ② This repeated read/write make MR unsuitable for iterative computations.

MapReduce - Fault Tolerance

→ WORKER FAILURE

a) Worker Failure Detection:

- The master node pings workers periodically.
- If a worker fails to respond within a certain timeout, it is marked as failed by the Master.

b) Task Rescheduling:

• Map Tasks:

- Completed map tasks on the failed worker are re-executed, as their output is stored on the local disk of the failed node & is no longer accessible
- In-progress map tasks are reset to the idle state and rescheduled on other workers.

• Reduce Tasks:

- In progress reduce tasks are also reset to the idle state and rescheduled
- "Completed reduce tasks don't re-execute" as their output is stored in the "global file system"
eg: HDFS.

→ MASTER FAILURE :

- if the master tasks fail:
 - The system can restart the master from the checkpoint state "if such a mechanism is available"
 - In the absence of checkpointing, the current "MapReduce computation is aborted" if on a single node
- Client responsibility:
 - The client can detect master failure and choose to retry the entire MapReduce oper.

→ SEMANTICS IN PRESENCE OF FAILURE :

a) Atomic Output Commits :

- To ensure fault tolerance, all outputs are committed atomically during map and reduce task completion

b) Output Handling by Map Tasks :

- Each map task writes its output to a private temporary file (one file per reducer).
- On completion the worker sends a message to the master with the names of these temporary files.
- If the master receives duplicate completion messages for the same task, it ignores them.

c) Output Handling by Reduce Task :

- Each reduce task writes to a private temporary file during execution.

Page No.

- Once completed, the task atomically renames the temporary file to the final output file.
- If multiple workers execute the same reduce tasks, multiple non-rename calls for the same file occur. However the atomic op. ensures that only one final output file is created.

▷ MapReduce - Data Locality

- Network bandwidth is a limited resource in distributed computing environments.
- Conserving bandwidth is critical to ensure efficient data processing and scalability.
- Data Locality:

The idea is to move the computation to the data rather than transferring large amounts of data across the network.

- How it works?

1. GFS (Google File System)

- Files are divided into 64MB blocks and replicated 3 times across diff machines.

- This replication ensures redundancy and availability of data.

2. Scheduling MapTasks:

- The masterNode uses the file location info from GFS.

- It tries to schedule map tasks:

↳ On a machine with a local copy of the data.
↳ Or near the machine (e.g. on the same rack or network switch)

3. Result:

- A significant portion of input data is read locally avoiding unnecessary network transfers
- Local reads consume "no network bandwidth."

♦ REFINEMENTS :

o PARTITIONING FUNCTION:

• Purpose:

→ Divides intermediate data (key:value) into partitions corresponding to reduce tasks.

→ Ensures all data for a specific key or logical group ends up in the same reducer.

• Default Partitioning:

→ Typically uses $\text{hash}(\text{key}) \bmod R$, where R is the no. of reduce tasks.

• Custom Partitioning:

→ Users can provide a special partitioning function for specific needs.

→ Example:

↳ For URLs, using hash(hostname(urlKey))
mod R ensures all URLs from the same host
go to the same partition / output file.

o ORDERING GUARANTEES :

- Within - Partition Ordering:

→ MapReduce ensures intermediate key-val pairs within a partition are processed in increasing key order.

- Benefits:

→ Enables the generation of sorted output files.

→ Useful for:

↳ Formats that require efficient random access lookups by key

↳ Users who find f sorted data more convenient to process.

o COMBINER FUNCTION :

- Performs local aggregation of intermediate data on the map worker before transferring it to reducer.

- Act as a mini-Reduce function on each map worker.

- Combines data like key counts <word, 1> to reduce the volume sent over the network.

- Benefits:

↳ Reduces network traffic

- ↳ Speeds up operations by minimizing data transfer.
- Implementation
 - ↳ The same code is used for the Reduce & Combiner function
- ↳ Difference
 - ↳ Combiner output is temporary and sent to reducer
 - ↳ Reducer output is final & written to HDFS.

▷ SIDE EFFECTS :

- In MapReduce side effects refer to additional outputs like:

- ↳ logs for debugging
- ↳ Metadata, etc

generated by a task beyond the main key-val output.

• Ensuring Atomicity :

- Atomic writer prevent partially written or corrupted files from being visible to the system.

→ How it works

- ↳ Write side effects to a temporary file
- ↳ Rename the temp file to its final name only after the operation is complete

- This ensures that only fully generated outputs are visible to other parts Page No. [] of system.

→ Limitations :

1. No-Two Phase Committ:

Map Reduce does not support atomic two phase commits, means it cannot guarantee consistency between multiple files written by the same task.

2. Deterministic Behavior Needed:

Task producing multiple outputs must ensure that their execution is deterministic means, given the same input, the task must produce the same output each time.

This is essential when tasks are retried due to failures.

▷ SKIPPING BAD RECORDS:

- Some records/data may cause deterministic crashes in user-defined MapReduce Functions
- When a task crashes on a specific record:
 - ↳ A signal handler catches the crash and logs the record's sequence no.
 - ↳ The MapReduce master tracks records that repeatedly cause crashes
 - ↳ After multiple failures, the master instructs workers to skip the problematic record.
- Ensures progress even if certain records cannot be processed.

⇒ STATUS INFORMATION :

- Master Node's HTTP Server
 - Provides live status pages for monitoring job progress and diagnosing issues, debugging user code.
 - Displayed Information
 - ↳ No. of completed & in-progress tasks
 - ↳ Data processed
 - ↳ bytes of input
 - ↳ bytes of intermediate data
 - ↳ bytes of output
 - ↳ processing rates
 - ↳ Workers failures and tasks they were processing
 - ↳ Links to logs
 - ↳ Standard error files
 - ↳ Standard output files

⇒ COUNTERS :

- Provides facility to count occurrences of various events (e.g.: total words processed, no. of failed tasks).
- User define counters are updated in the Map or Reduce func
- Counter values are periodically sent to the master for aggregation.

- Current counter values are shown on the status page for live monitoring.
- Final aggregated values are returned to the user after job completion.
- Duplicate task executions are handled, ensuring no double counting of counters.

▷ PERFORMANCE

1. CLUSTER CONFIGURATION:

- 1800 machines, each with:
 - Two 2GHz Intel Xeon processor
 - 4 GB memory (1-1.5GB reserved for other tasks)
 - Two 160GB IDE disks
 - Gigabit Ethernet link
- Machines arranged in a two-level tree shaped network with 100-200 Gbps aggregate bandwidth at the root.
- Round-trip time b/w machines < 1 millisecond.
- Data splits into 64 MB chunks for input.

2. GREP :

- Scan 1 billion 100 byte records to find a rare three character pattern.

• Input Data :

↳ split into 64 MB chunks

↳ 15,000 map tasks & 1 reduce task.

• Output Data : All result stored in 1 file

• EXECUTION OVERVIEW

1. Start up Phase :

↳ Takes 60 sec due to

↳ Propagating the program to all worker machines

↳ Interacting with GFS (Google File System)

↳ Opening the set of 1,000 input files

↳ Retrieving information for data locality optimization.

2. Map Phase :

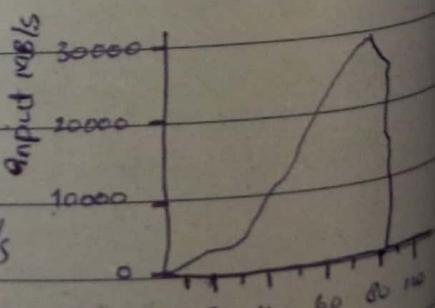
↳ Rate of data scanning

increases as more machine

(workers) are assigned to the job

↳ Peak Input Rate : over 30 Gb/s

with 1,764 workers



3. Completion Phase :

↳ As map task finishes, the rate of data scanning declines

↳ Hit zero around 80 seconds.

Page No. seconds
Data transfer rate over time

- Entire computation takes 150 seconds.
- ↳ 60 sec for startup
 - ↳ 90 sec for actual map & reduce task execution.

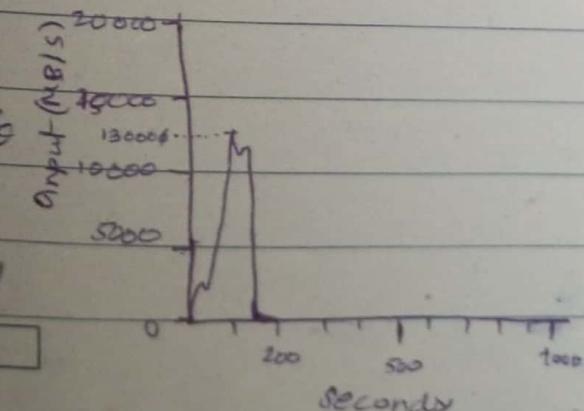
3. SORT :

- Sort 1 billion 100-byte records
- Input Data is divided into 15,000 map tasks
 - each processing 64MB piece
- Output is partitioned into 4,000 files
- Map Function:
 - ↳ Extracts a 10-byte key from each text line
 - ↳ Emits the key as the intermediate key and entire line as the value
- Reduce Function:
 - ↳ Uses an Identity Func., which passes the intermediate key-val pair unchanged as the final output.
- A custom positioning function segregates data into 4000 partitions based on the initial bytes of the sorting key.

1. Input Read Rate

Peaks at 13 GB/s due to data locality optimization

Drops off quickly once all map tasks complete



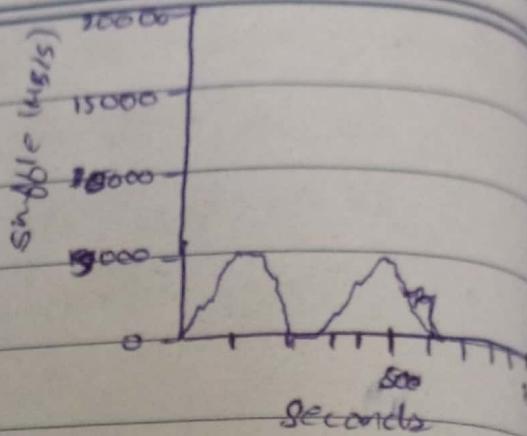
Page No. []

Date

Day

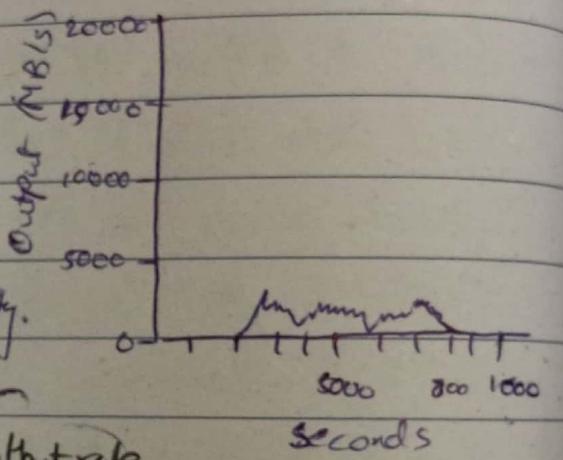
2. Shuffle Rate:

- Data transfer b/w map & reduce begins as soon as the 1st map task completes



3. Output Write Rate

- Lower than Shuffle rate due to replication.
- Two copies of output are written for reliability & availability.



~~Input Rate > Shuffle Rate > Output rate~~

↳ due to
data locality

↳ It is slow bcz
of replication.

Erasure coding instead of replication
can solve this problem.

→ Backup Task 8 Machine Failure Room Slides /