

Wk	M	T	W	T	F	S	S
10	1	2	3	4	5	6	
11	7	8	9	10	11	12	13
12	14	15	16	17	18	19	20
13	21	22	23	24	25	26	27
14	28	29	30	31			

2016

Saturday

January

16

016-350 • WK 03

Video 1

Big data processing pipelines : A dataflow approach

→ summarize what dataflow means and its role in data science

10:00 →

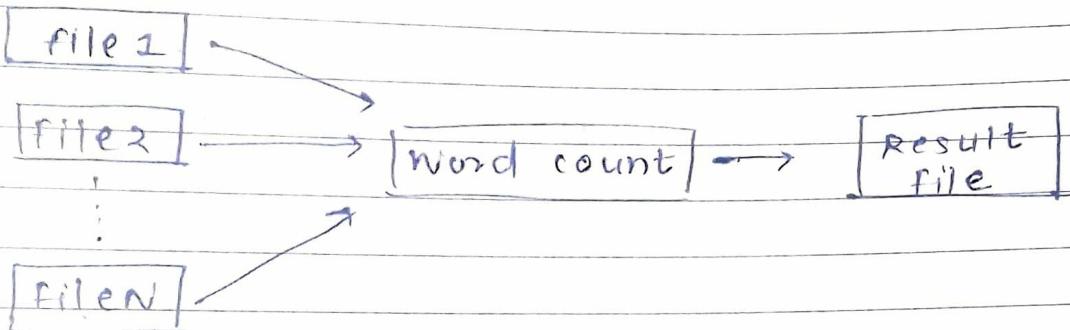
Explain 'split->do->merge' big data pipeline with eg,

→ Define the terms data parallel.

11:00

Example MapReduce App :- wordcount

12:00



1:00

2:00

3:00

steps : split

4:00

Step 2 : Map

5:00

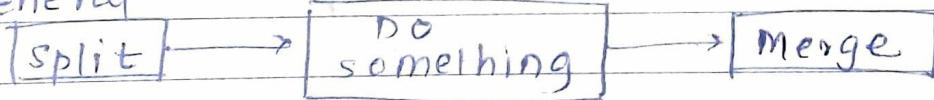
Step 3 : shuffle and sort

5:00

Step 4 : Reduce

In general

6:00



Represents a large no. of applications

7:00

Big data pipelines

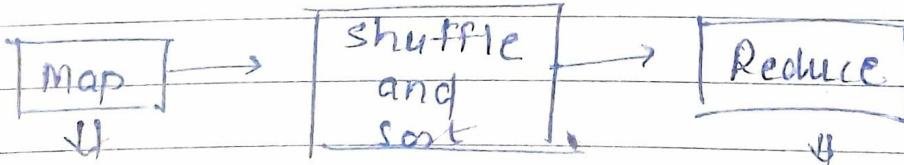
pipe = Unix operation

Word count Example :-

Sunday

17

NOTES



parallelization over the

TP

parallelization

data sorting

parallelization

over data groups

video 2

- * Some high-level processing operations in Big data pipelines

storage → HDFS
→ concan.
→ Many other

Data Transformation

- List common data transformations within big data pipelines
- Design a conceptual data processing pipeline using the basic data transformations.

① Map :- Apply some operation to each member of a collection. → color each member of set

② Reduce :- 'collecting' things that have same key

③ Cross/ cartesian :- Multiplication

④ Match/ join :- selective multiplication

(Do some process to each pair from two sets - which have same key)

⑤ Co-group :- Group common items

→ collect similar things.

→ Apply a process to each collection.

Do some process to each pair from two sets.

⑥ Filter :- Select Elements that match a criteria
Filter even no.

$$x \cdot 2 = 0$$

NOTES

December 2016						
wk	M	T	W	T	F	S
49				1	2	3
50	5	6	7	8	9	10
51	12	13	14	15	16	17
52	19	20	21	22	23	24
53	26	27	28	29	30	31

2016
Friday
October

28
302-064 • WK 44

Aggregation in Big data Pipelines.

- Compare and select the Aggregation operation that you require to solve your problem.
- Explain how you can use Aggregations to compact your dataset and reduce volⁿ (in many cases)
- Design complex operations in your pipelines using a series of Aggregations.

What is Aggregation :- $f(\text{all elements})$
Symbol for any transformation

Aggregation $\rightarrow f(\text{all elements})$

Σ → symbol for summation

- ① GroupBy
- ② Average
- ③ MAX
- ④ MIN
- ⑤ standard deviation

Connecting Aggregations

sum \rightarrow Max $\Rightarrow \text{MAX}(\text{sum})$

sum \rightarrow min $\Rightarrow \text{MIN}(\text{sum})$

Boolean Aggregation :- ① AND
② OR

Sets :- Union, intersection, difference

NOTES strings :- concatenation

Aggregations \rightarrow Organized & compact data

Variety \rightarrow Actionable insights
volⁿ

A
U
G
2 16
22 23
28 29 30

July 2015

Week 28
Day 189 • 176
Date 08 • 07 • 2015

8

Wednesday

- List common analytical operations pipeline.
- Analytical operations within big data
- Describe sample applications for these analytical operations.

Analytical Operations :-



- Purpose :-
 - ① Discover meaningful trends and patterns in data.
 - ② Gain insights into patterns.
 - ③ Make data-driven decisions.

Sample Analytical Operations :-

- ① classification
- ② Path Analysis
- ③ clustering
- ④ Connectivity Analysis

K-means in spark

```
• spark python code for performing k-means on data
data = sc.textFile("data/mllib/kmeans-data.txt")
parsedData = data.map(lambda line:
    array([float(x) for x in line.split(' ')]) )
# cluster the data
clusters = KMeans.train(parsedData, 2, maxIterations=10,
    runs=10, initializationMode="random")
```

8.00 Path Analysis:-

- Path analysis using Cypher on neo4j
- II Finding shortest path b/w specific nodes:
match p = shortestPath((a)-[:TO*]-(c))
where a.Name = 'A' and c.Name = 'P'
return P, length(P) limit 1

II Find all shortest paths:

```
match p = allShortestPaths((source)-[r:TO*]-  
(destination))
```

where source.Name = 'A' and destination.Name = 'P'
return extract(n in nodes(p) | n.Name) as Paths

9.00 Connectivity analysis using Cypher on neo4j

II Find the degree of all nodes

```
match (n:myNode)-[r]-()
```

with n as nodes, count(distinct r) as degree

return degree, count(nodes) order by degree

Graph Analytics Techniques

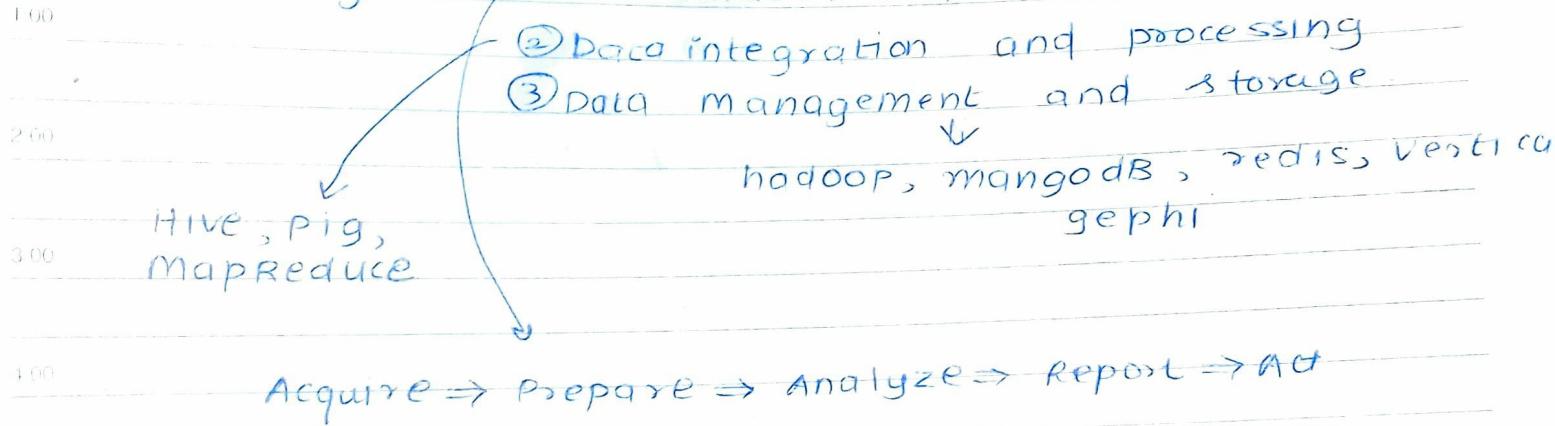
- ① Path analytics
- ② Connectivity Analytics
- ③ community
- ④ centrality

Videos

Overview of Big data processing system

- 9:00 → Hadoop Ecosystem
- Layer Diagram with three layers for data storage, data processing and workflow management.
- 10:00 → Summarize an evaluation criteria for big data processing systems.
- 11:00 → Hadoop, spark, Flink, Beam and storm
- 12:00

Hadoop Ecosystem :-



Categorize of Big Data processing

- Execution model → Batch
- streaming

- 6:00 → Latency
- 7:00 → Scalability
- programming language
- Fault Tolerance

Big data processing system

- Apache storm
- spark
- hadoop
- flink
- beam

NOTES

Hadoop Map Reduce :-

- EM :- Batch processing using disk storage
- Latency :- High
- Programming lang :- Java
- Fault tolerance :- Replication

Spark :- EM :- Batch and stream processing using disk or memory storage

Latency → low-latency for small micro batch size
program language is scala, python, java, R

Flink :- EM :- Batch and stream processing using disk or memory storage

Latency → low

PPL :- Java and scala

Beam :- Batch and stream processing

low-latency

java and scala

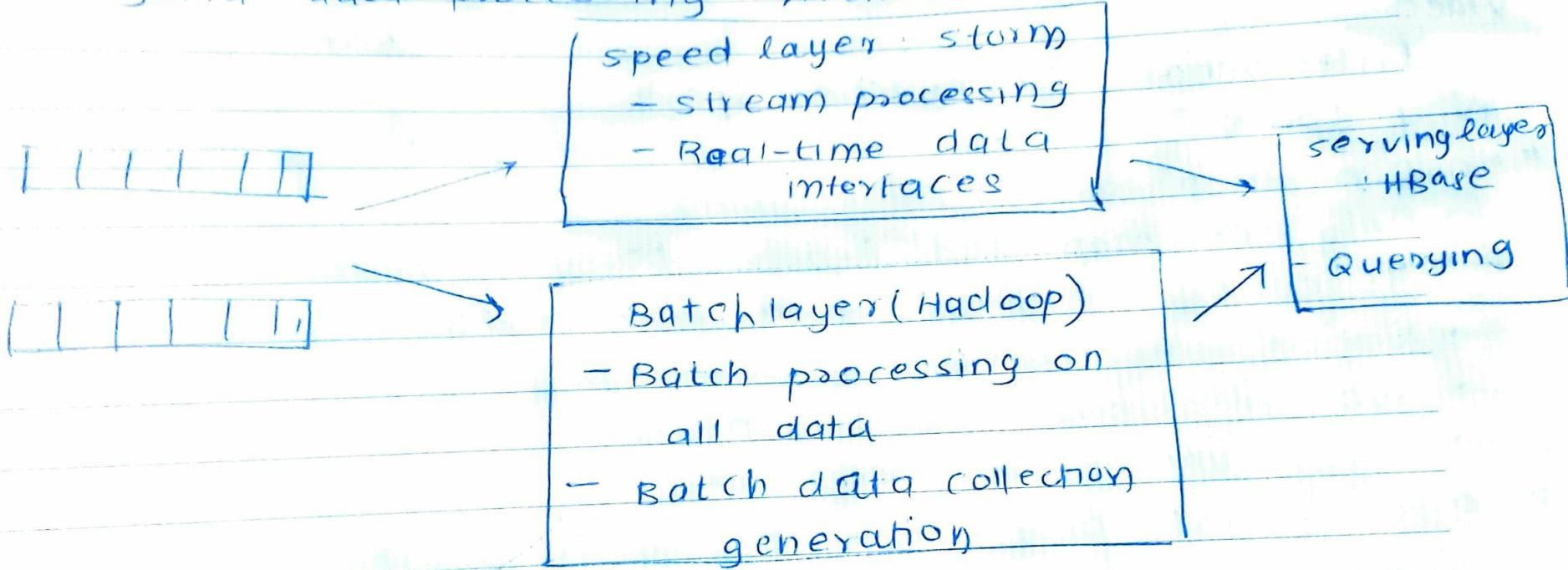
Storm :- stream processing, very low-latency,

many programming langs

NOTES

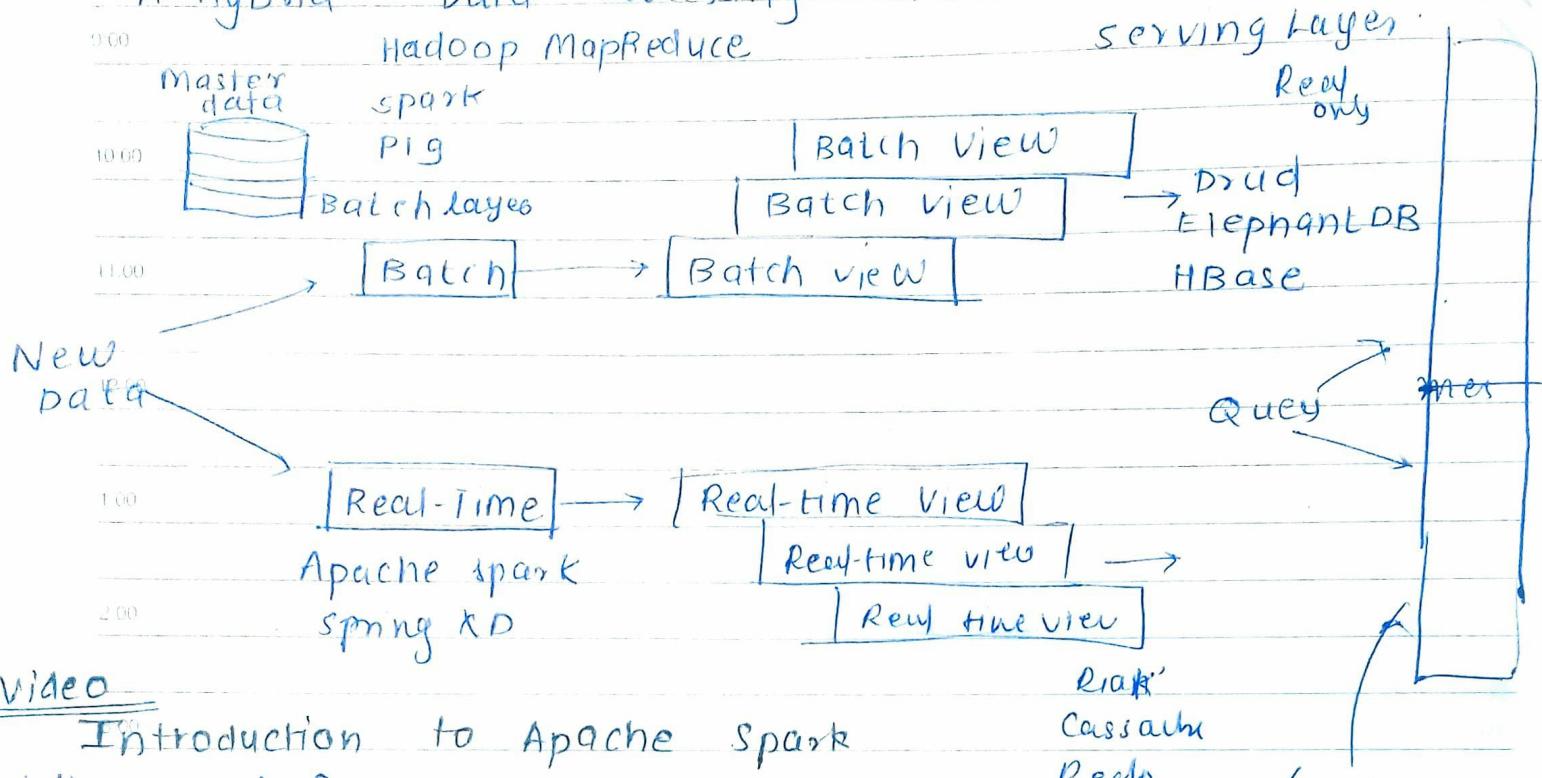
Lambda Architecture :-

A Hybrid data processing Architecture.



TES

A Hybrid Data Processing Architecture



Video

Introduction to Apache Spark

Why spark?

* Hadoop MapReduce Shortcomings :-

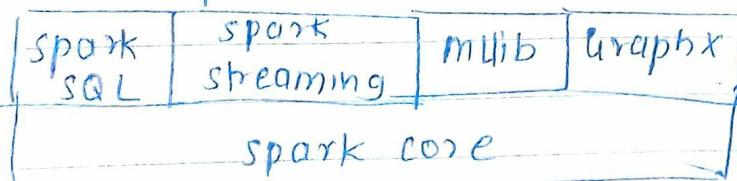
- Only for map and Reduce based computations - (window of data)
- Relies on reading data from HDFS.
- Native support for Java only.
- No interactive shell support.
- No support for streaming.

* Basics of Data Analysis with spark :-

- Expressive programming model
- In-memory processing
- Support for diverse workloads

NOTES → Interactive shell

The spark stack



Explore

→ Build

→ scale

ing started with Spark:-

The Architecture and Basic concepts

Date _____

Describe how Spark does in-memory processing using

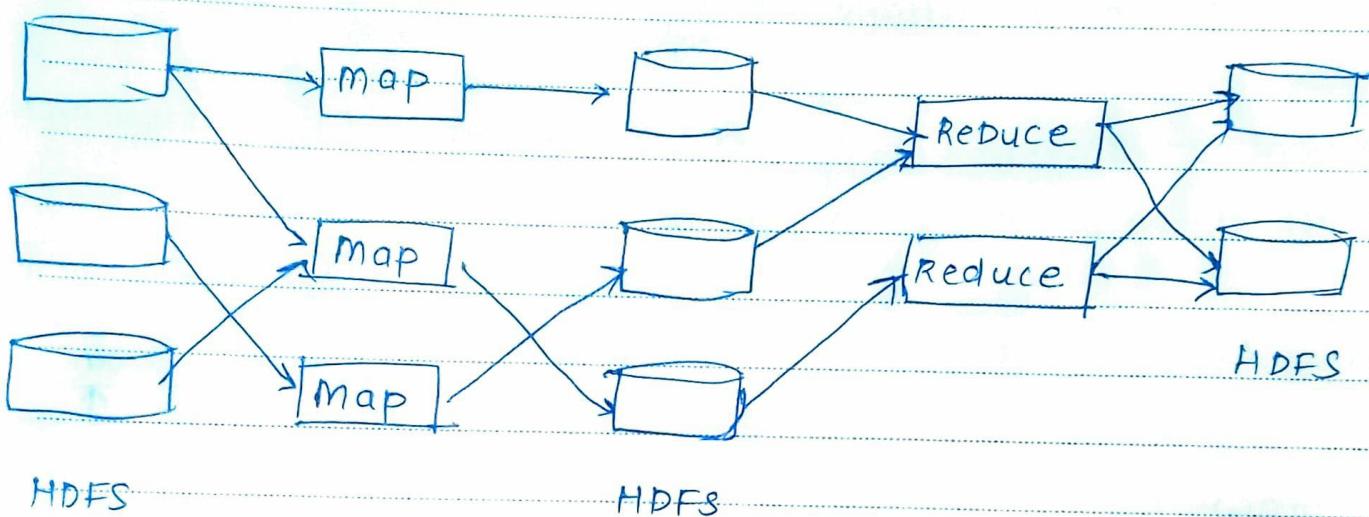
RDD abstraction

→ Explain the inner workings of the Spark architecture

→ Summarize how Spark manages and executes code on clusters

Q. What does in memory processing mean?

MapReduce



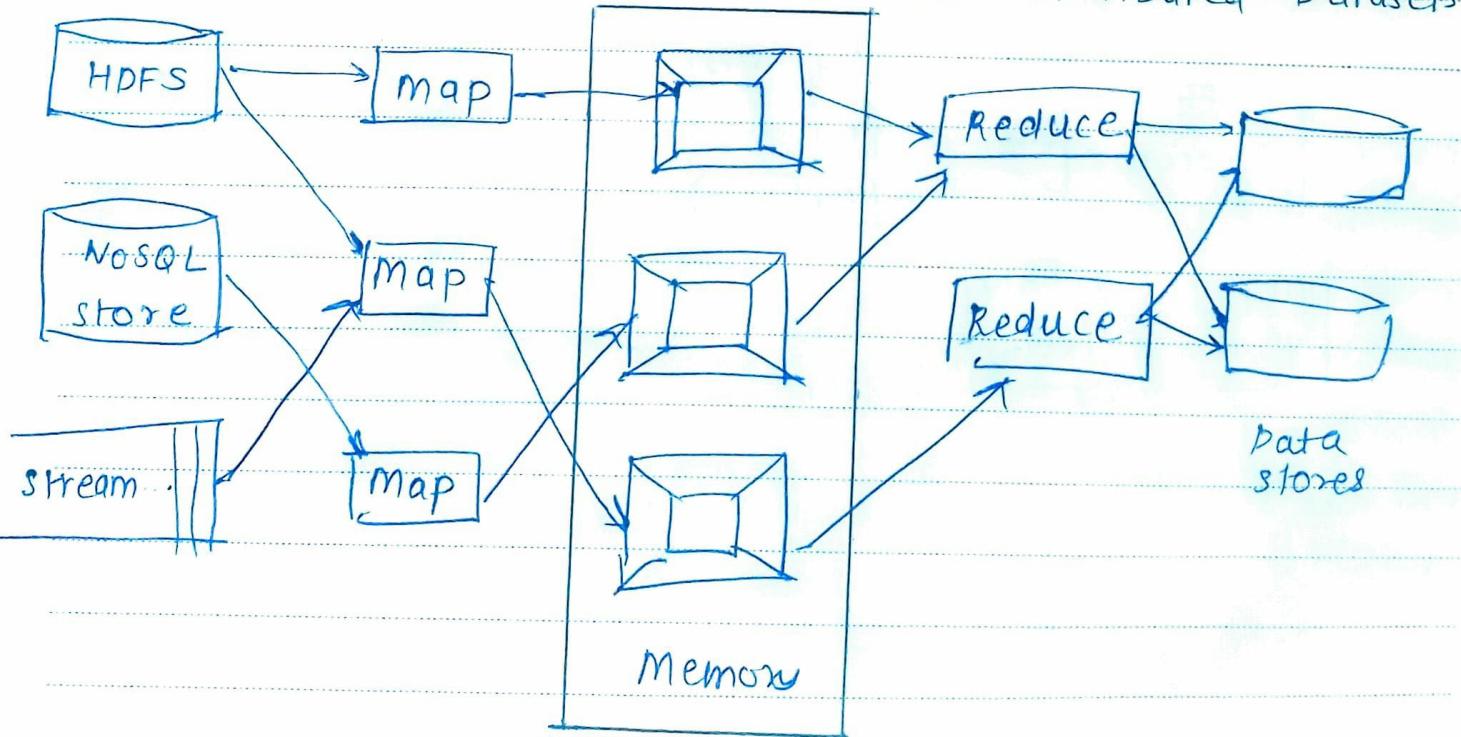
HDFS

HDFS

HDFS

Spark

Resilient distributed datasets



Priority

M/T/W/T/F/S/SU/

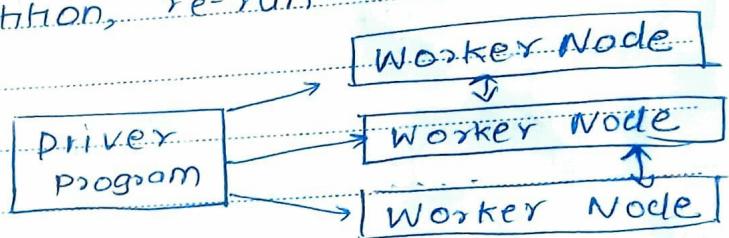
Date.....

5/19 16 23
6/2 15 22
3 28 29
14 21 28
27 30

Resilient distributed Datasets :-

- ① Dataset → Data storage created from:
HDFS, S3, HBase, JSON, text, Local hierarchy
→ or created transforming another RDD.
- ② Distributed → Distributed across the cluster of machines.
→ Divided in partitions, atomic chunks of data
- ③ Resilient → Recover from errors eg. node failure, slow processor
→ Track history of each partition, re-run.

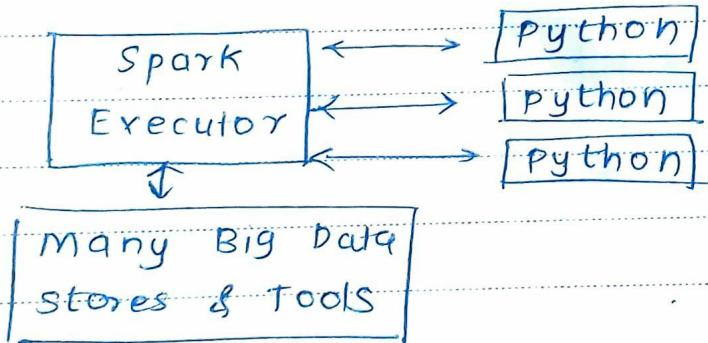
Spark Architecture :-



Driver program :-

```
lines = sc.textFile("txt")
```

Worker Node



(Priority)

Q.

What are the layers of spark?

- GraphX
- Spark core
- MLlib
- spark streaming
- spark SQL

Q. Why is Hadoop is not good platform for machine learning?

- Bottleneck using HDFS
- Map and Reduce Based computation
- No iterative shell and streaming
- Java support only

Q. What are three layers of Hadoop Ecosystem?

- coordination and workflow management
- Data Management and storage
- Data integration and processing

Q. What is Data-parallelism as defined

- Running the same function simultaneously for the partitions of a data set on multiple cores.

Q. Which procedure best generalizes big data

- procedures, such as the map reduce process?
- split → do → merge

NOTES

Q. 5 keypoints in order to categorize big data system

- Execution Model, Latency, Scalability,
- Programming language, fault Tolerance

Q. What is Lambda Architecture?

- A type of hybrid data processing architecture.

Q. Which of the following scenarios not an aggregation operation?

- Removing undefined values.

Q. What happens to data when it aggregated

- Data becomes smaller.

Q. K-means clustering

- Group samples into k clusters.

5.00

Q. Writing data to memory bth pipeline steps

In-memory processing.

Sunday 27

NOTES