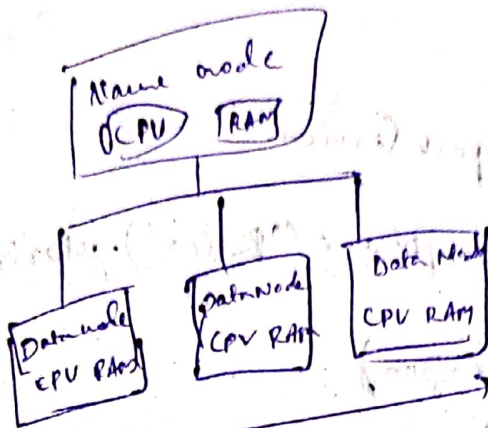


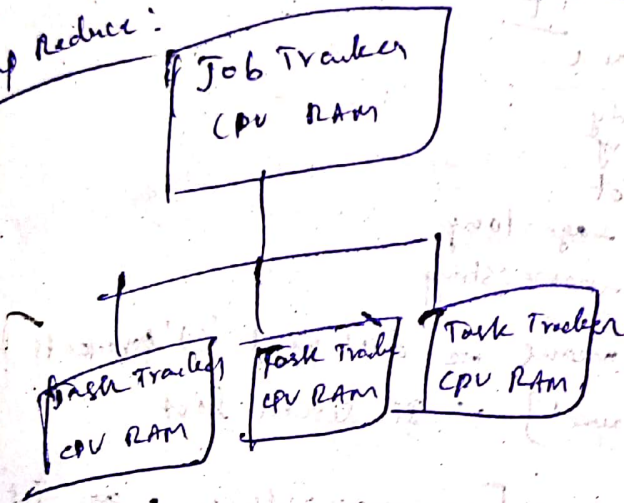
Spark: (It is ^{use to} handle ~~to~~ Big Data)
 → Hadoop Distributed File System (HDFS) (Distributed large dataset)
 ↓
 uses because of fault tolerance



It we use only local machine, if data is lost then we lose all data.
 Coming to HDFS, it distributes data to all processes and maintains availability of data.

(Distribute a Computational task to a distributed dataset).

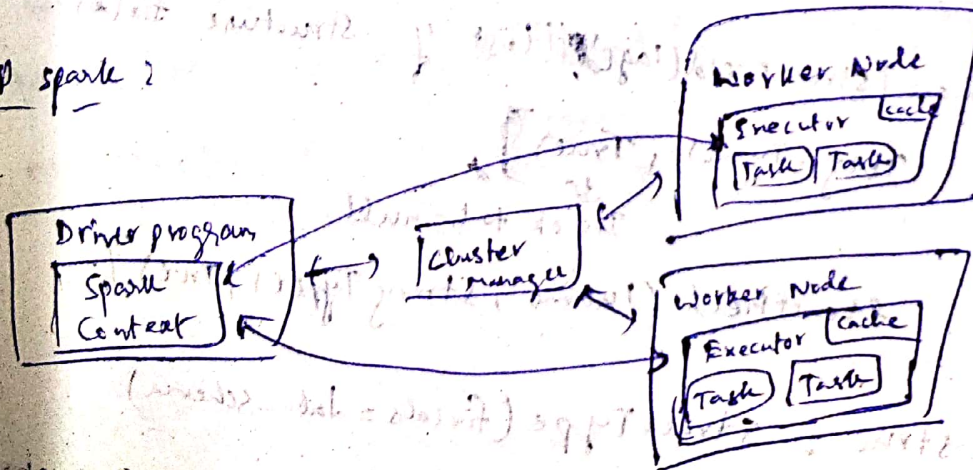
Map Reduce:



Spark vs MapReduce?
 ↓
 100% faster than mapreduce.
 It writes most of data to disk after each map and reduce operations.

→ Spark keeps most of the data in memory

RAP spark:



Python + Spark

Column → Some feature/variable
row → Individual datapoint

Spark dataframe Basics:

```
- from pyspark.sql import SparkSession  
- spark = SparkSession.builder.appName('Basics').getOrCreate()  
df = spark.read.json('people.json')
```

```
df.show()
```

age	name
30	andy
20	ring

```
df.printSchema()
```

```
// root  
|-- age: long  
|-- name: string
```

df.columns // attribute means we don't need '()' brackets
['age', 'name'] → it return list

```
df.describe() // Dataframe [summary=string, age=int]
```

```
df.describe().show()
```

```
from pyspark.sql.types import StructField, StringType,  
IntegerType, StructType
```

```
data-schema = [StructField('age', IntegerType(), True)]
```

// list of structure fields

It's ok to be null.

```
StructField('name', StringType(), True)]
```

```
final-struct = StructType(fields=data-schema)
```

```
df = spark.read.json('people.json', schema=final-struct)
```

```
df.printSchema()
```

root
-- age: Integer
-- name: string

`df['age']` // pyspark.sql.column.Column
`df.select('age')` // DataFrame[age: int] → It returns object
`df.select('age').show()` //

age
null
30
19

 → It returns DataFrame.
`df.head(2)[0]` // Row(age=19, name=Tony)
`df.select(['age', 'name'])` // multiple columns

`df.withColumn('newage', df['age'])` // Creating new Column.
`.show()`

age	name	newage
null	Tony	null
20	Ruby	20

`df.withColumn('double-age', df['age'] * 2).show()`

→ It double the age

`df.withColumnRenamed('age', 'my-new-age').show()`

PLSQL?

`df.createOrReplaceTempView('people')` // Registered as SQL Temp view

`results = spark.sql("select * from people")`

`results.show()` // Table

`new_results = spark.sql("select * from people where age > 30")`

`new_results.show()` //

age	name
30	Andy

df = spark.read.csv('appl-stock.csv', inferSchema=True, header=True)

df.head(3) // first three rows

df.filter("close < 500").show()

df.filter("close < 500").select(['open', 'close']).show()

// it gives the columns of open close where close < 500.

df.filter(df['close'] < 500).show()

df.filter((df['close'] < 200) & (df['open'] > 200)).show()

result = df.filter(df['Low'] == 197).show().collect()

// returns a list

row = result[0]

row.asDict() // converts list to dict

row.asDict()['volume']
// one of the column name

groupBy & Aggregate Operations

df.groupBy("Company").mean().show()

sum()

max()

min()

count()

df.agg({'Sales': 'sum'}).show()

column name

function name

(max, min)

etc

sum(Sales)
870

group_data : df.groupBy("Company")

group_data.agg({'Sales': 'max'}).show()

Company	max
Apple	300
google	850

from pyspark.sql.functions import countDistinct, avg, stddev

df.select (avg('sales')).show() // $\frac{\text{avg(sales)}}{33-33}$ (Average - sales)
 alias('Average_sales')

df.orderBy (df['sales'].desc(), show()) → Descending order

df.orderBy ('sales').show() → Ascending order

Missing Data

thresh = 2 (means, if row has 2 null values, then drop that row)

~~df.dropna()~~ df.na.drop().show()

df.na.drop (how='all').show() → (By default it set as 'any' means if row has any null value then drop it)

df.na.drop (subset=['sales']).show()

df.na.fill ('FILL VALUE').show()

df.na.fill ('name', subset=['name']).show()

⇒ (if name has null value, then 'name' is filled)

Dates & Time Stamps

from pyspark.sql.functions import dayofmonth, hour, dayofyear, month, year, weekofyear, format_number, date_format

df.select (dayofmonth (df['Date'])).show()

df.select (hour (df['Date']), month (df['Date'])).show()

df.withColumn ('year', year (df['Date'])).show()

newdf = groupby ('year').mean().show()

df.select ('year', avg('sales')).show()