# DATA LAKES with SPARK
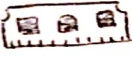
**BIG DATA** : You need a distributed system of several servers to be able to work with the data.

Numbers everyone should know

CPU "Brain of computer" Operation: 0.4ns

Memory / RAM . "Ephemeral storage". Memory Reference : 100ns

Storage : SSD/Magnetic Disk . Random Read : 16μs

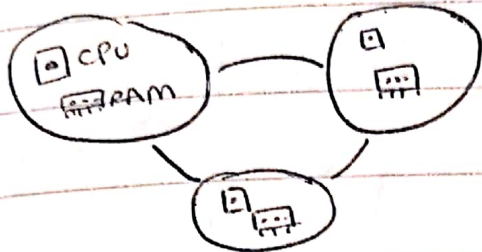Network : Access to outside . Round trip for data from EU to US : 150ms

FASTEST ↑

SLOWEST ↓

2.5 GHz CPU : 2.5 Billion Operations / second.

if 1 operation = 8 bytes → CPU can process 20 Billion Bytes/sec

Memory is 2x faster than storage
But it's more expensive.

# Distributed Computing vs. Parallel Computing

several Nodes accross Network.

Each cpu has itsown memory

1 node

multiple cPUs share a single memory.

→ 200x Faster than RAM
→ 15x " than SSD
→ 20x " than network

## Hadoop Framework

HDFS - Data Storage

MapReduce - Data Processing

YARN - Resources Manager

Hadoop Commons - libraries + Utilities

(+) Other Tools on Top.

Pig - SQL for MR          Hive - SQL for MR

SPARK          Storm - StreamingData

Flink - " "

**MR**

HDFS     MAP     SHUFFLE     REDUCE

Data
split in
chunks

| A | B | B |
| A | C | B | B |
| C | D | A |
| D | E |

key, values     keys, value.

| | keys, value | |
|---|---|---|
| (A, 1) | (A, 1) | (A, 3) |
| (B, 1) | (A, 1) | |
| | (A, 1) | |
| (B, 1) | (B, 1) | (B, 4) |
| (A, 1) | (B, 1) | |
| (C, 1) | (B, 1) | (C, 2) |
| ... | (B, 1) | (D, 2) |
| | ... | (E, 1) |

# DAG: directed Acyclical Graph

Spark builds a step by step directions
of what Functions + data it will need
→ Once it builds the DAG From the code
it checks if there is anything it
can procestinate
→ LAZY EVAWATION

## IMMUTABLE

Spark does not change of
mutate data

## MAPS

Makes a copy of the original input data
and applies a function to it.

*example*
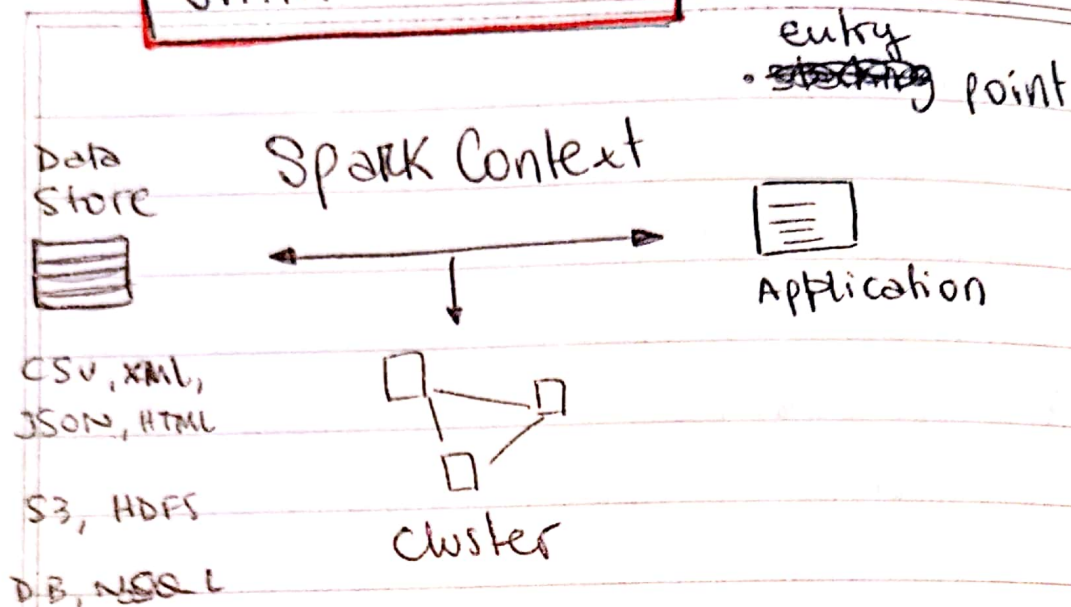
```
import py spark
sc = pyspark . Spark Context (app Name = "---")
distrib -songs = sc . paolellize ( songs )
distrib-songs . map (lambda x: x.lower())

distrib - songs . collect ()
```

# SPARK PROGRAMS

- ~~starting~~ entry point

Data Store

Spark Context

Application

CSV, XML, JSON, HTML

S3, HDFS

DB, NoSQL

cluster

```
from pyspark import SparkContext, SConf
config = SparkConfig().setAppName(" ")
        .setMaster(" IP | local")
sc = SparkContext(conf = conf)
```

To read DF:                            → sc SQL equiv.

```
from pyspark sql import SparkSession
Spark = SparkSession.builder /
        .appName(" ").config(" ")./
        .getOrCreate()
```

# Data wrangling with DF

Read JSON into DF:

```
user-log = spark.read.json(path)
user-log.printschema()
    "        .describe()
            .show(n=1)
            .take(5)
```

Save into CSV:

```
"     .write.save(out-path, format="csv",
                            header=True)
user-2 = spark.read.csv(path, header=True)
        .describe('column').show()
        .count()              .dropDuplicates()
        .select()            .sort("column")
        .where(column == value)
        .groupby(   )
```

draw a Plot of your df

```
pd = df.toPandas()
plt.scatter(pd[x], pd[y])
plt.xlim(-1, 24);
    func = udf(lambda x: ..., IntegerType())
df.withColumn("colname", func(col))
```

# DATA LAKES

Y
#
3

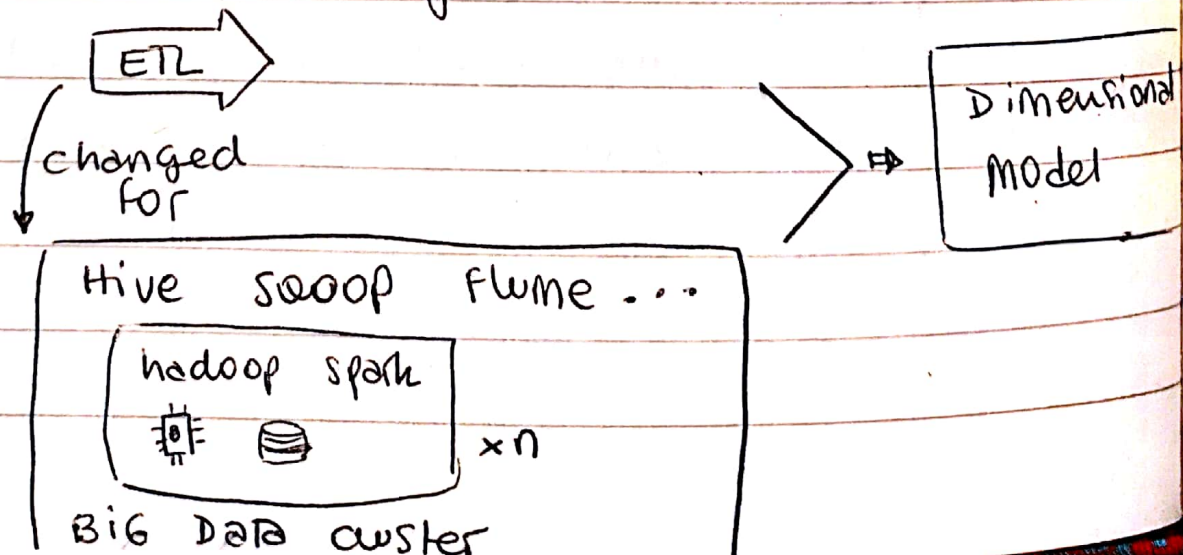- Abundance of unstruct. data
- Unprecedented data volumes
- Rise of Big Data techs.
- Data analysis - new types (ex. ML, NLP,)
- Emergence of new roles (ex. Data Scientist)

*Evolved from Data warehouse to cope with:
- Variety of data formats and structuring
- Support to agile + ad-hoc data exploration.
- Wider data transformation needed

Big Data tech:

(i) ETL Offloading : instead of ETL grid or staging area → same Hw for storage and processing.

(11) **schema - on - read** → As easy to work with files as with a database, without having to → create a DB
→ insert data into DB.

```
df = spark . read . csv ( path,
        inferSchema = "True", header=True,
        sep = ';').
```

TO specify a schema:

```
schema = Struct Type ([
        Struct Field ("col-name", IntegerType ())s
        ... ])
df = spark . read. csv (path,
        schema = schema, sep=';',
        mode = 'DROPMALFORMED')
```
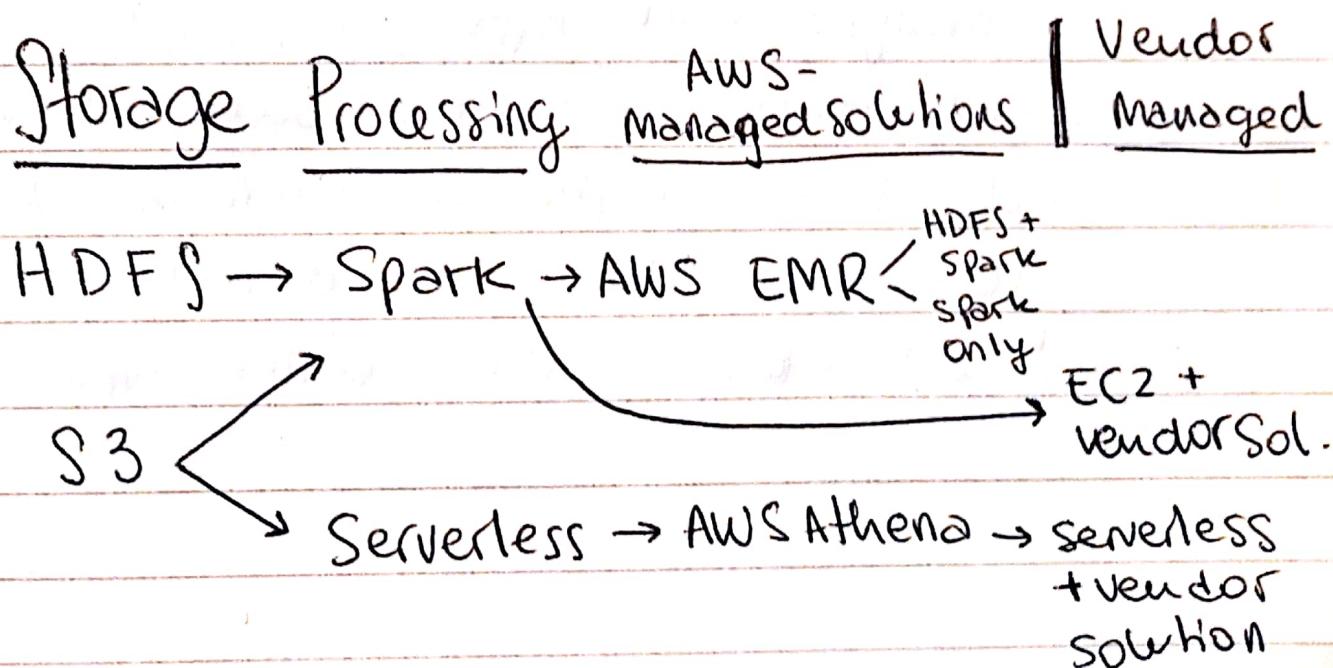
```
df .create or Replace TempView ("Table")
Spark . Sol (" select * from table"),show(s)
```
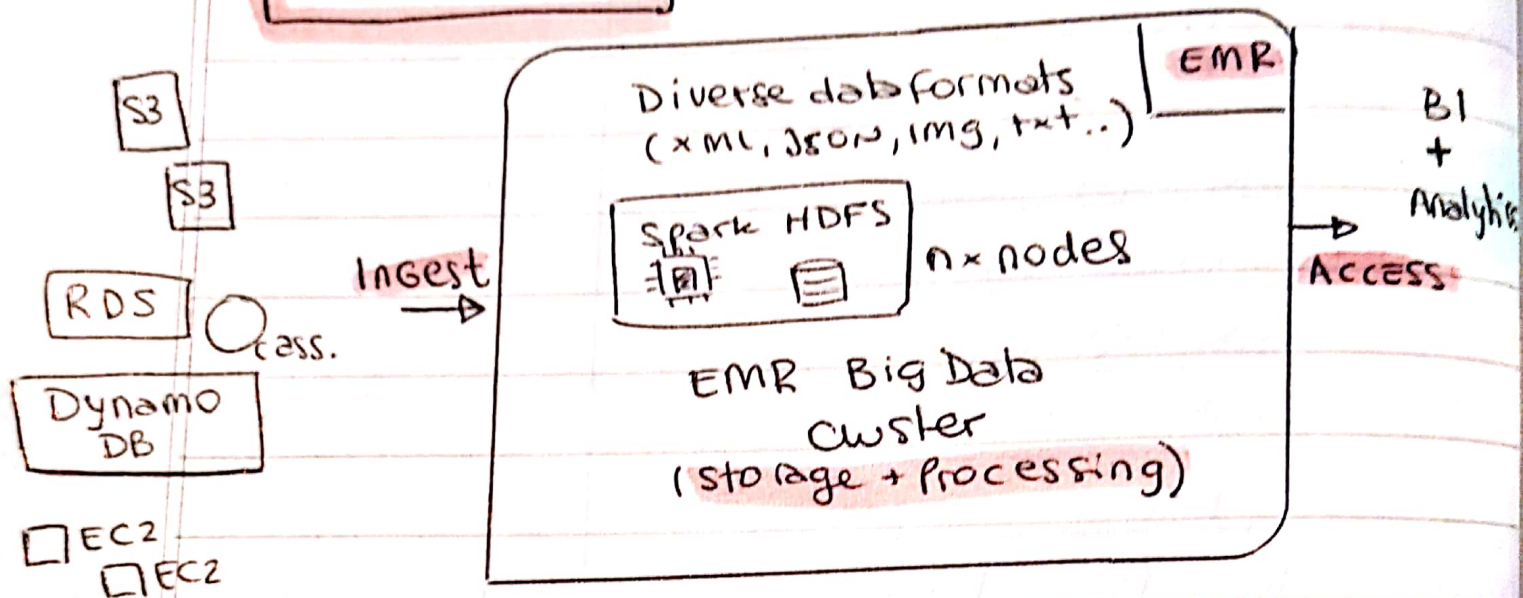
## (iii) Un Structured support

- Spark can read/write files in
  - Text based formats → Text, JSON, CSV..
  - Binary Formats → Avro, Parquet ..
  - Compressed formats → Gzip, snappy

- Read/write from <> FS:
  - Local file system
  - HDFS
  - S3

         & Databases:
  - SQL through JDBC
  - NoSQL → MongoDB, Cassandra, Neo4j
    Spark.read.format ("jdbc"). ...

| | Data Warehouse | Data Lake |
|---|---|---|
| Data Format | Tabular | All formats |
| Ingeshion | ETL | ELT |
| Data model | Star snowflake w/ conformed dimensions or DM or OLAP cubes | Star, snow, OLAP + ad-hoc reps. |
| schema | schema-on-write | schema-on-read |
| Tech | Expensive MPP, disks, netws. | Commodity Hw. |
| Data Q! | High, consistent, chean | mixed, raw, some transfo. |
| users | BA | DS, BA, ML engineers |
| Analysis | Reports, BI viz. | ML, Graphs, data exploration |

# DATA LAKE ON AWS

| Storage | Processing | AWS-managed solutions | Vendor Managed |
|---|---|---|---|

HDFS → Spark → AWS EMR < HDFS + spark / spark only

S3 → EC2 + vendor sol.

S3 → Serverless → AWS Athena → serverless + vendor solution

## AWS EMR  (HDFS + Spark)



**S3**

**S3**

**RDS** Cass.

**Dynamo DB**

☐ EC2
  ☐ EC2

**Sources**

Diverse data formats
(xml, json, img, txt..)   | EMR

Spark  HDFS      n × nodes

EMR  Big Data
        Cluster
(storage + processing)

Ingest →

Access →

BI
+
Analytics
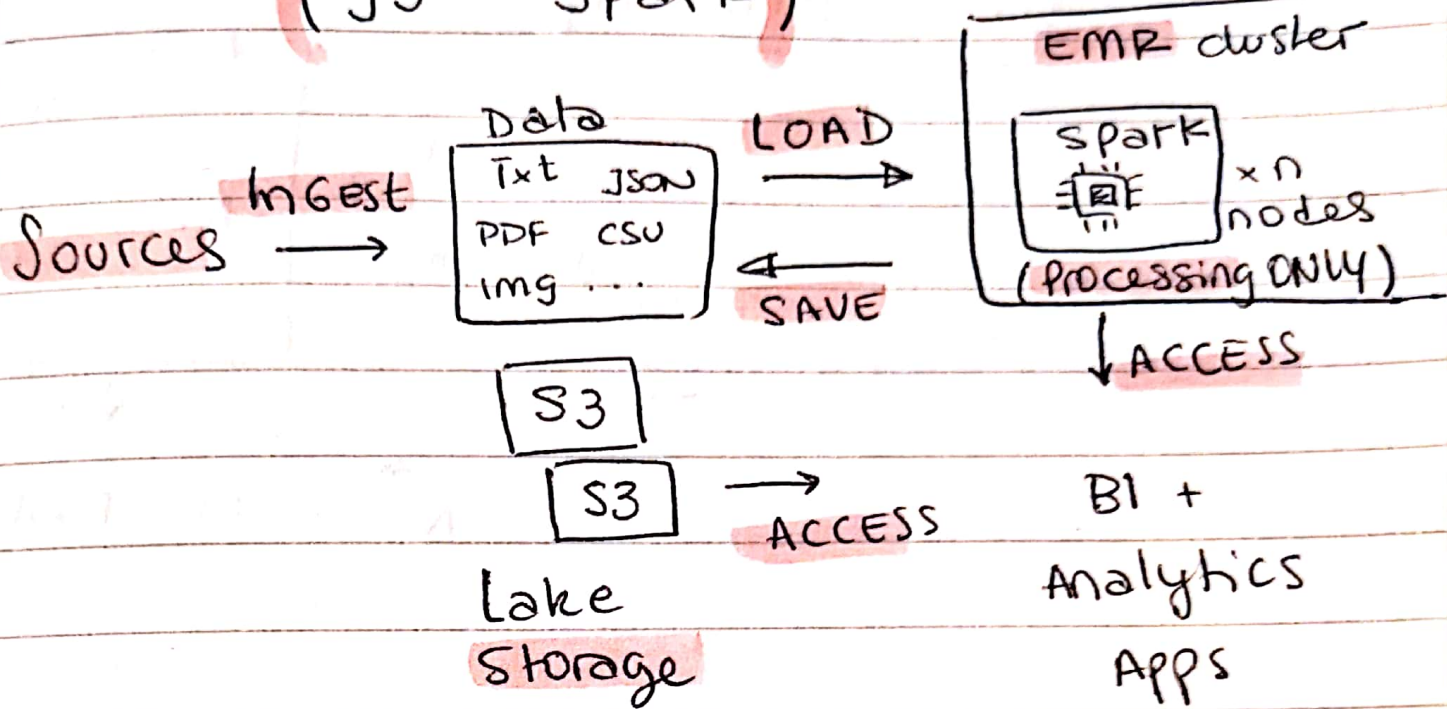
- Once ingested all data is stored on HDFS, processed on the cluster + accessed from Analytics + BI Apps.

- EMR cluster is created once, it can grow but it cannot be shutdown or else it loses all the stored data.
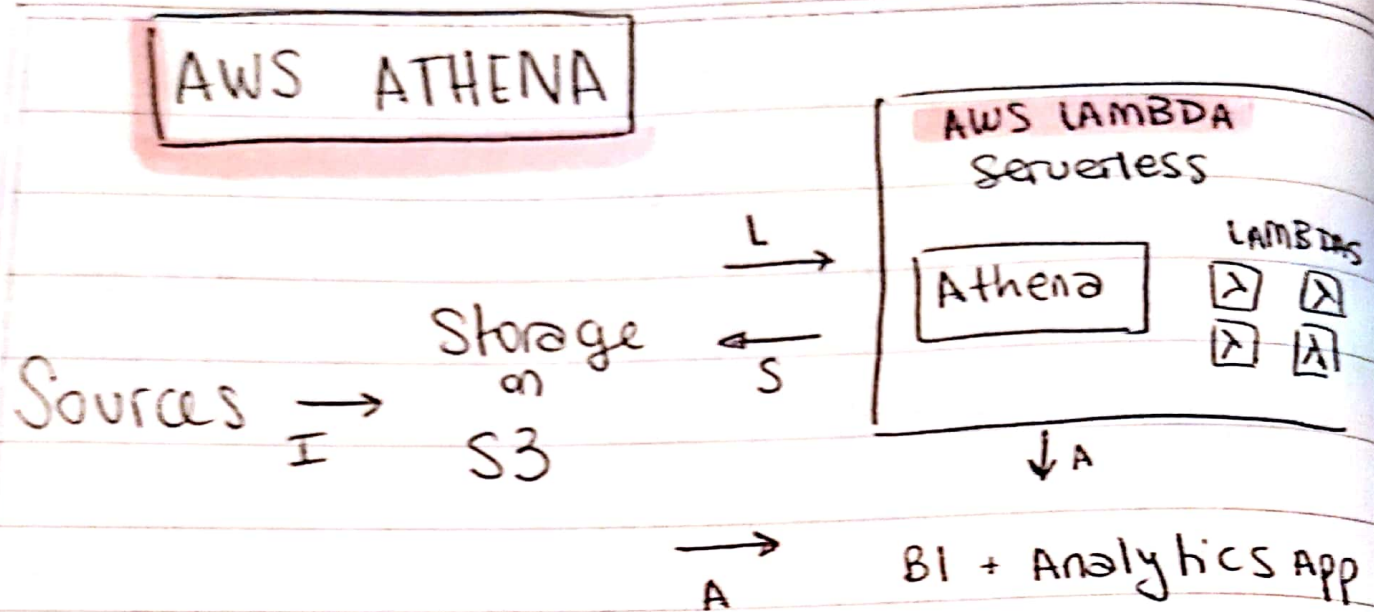
! COST → keeping the cluster running –

# ( S3 + Spark )



- NO HDFS in EMR
- EMR cluster is spun on-demand, shutdown
  otherwise.

<u>Potentially</u> → less costly, easier to manage,
more performant.

## AWS ATHENA

AWS LAMBDA
Serverless

LAMBDAS

Storage on S3

Sources →
I

L →
← S

Athena

A ↓

→
A

BI + Analytics App

LAMBDA → Function as a service
(Pay by execution time)

ATHENA → loads + processes data on serverless lambda resources

• Transparent management of resources