# Big Data

## Lecture 1. Introduction.

4th V does not only apply to Big Data (opinion)   Same with 5th   → added by business people to maintain hype.

big data ≠ volume.    velocity (streaming) ⇒ big data.

- Shared nothing: communication, consistency issues. (Spark)

Map/Reduce strategy → difficult!

- Lambda architecture: (λ):   IoT example with trucks.

  ↳ divide data in batch (persistent; immutable) → master data   it precomputes queries
               speed (real time)
                              ↳ serving layer.

## Lecture 2. Intro to CS/Cloud.

## Lecture 3. Dockerize model.

- Docker (SAAP): container management.
  ↳ has everything on app needs to run.

                                           vs.        Virtual Machines  (good isolation // higher overload)

                   slim
                   easy ship, pack, run
                   easy CI/CD
                   scability

- CLI: intermediario en las cominaciones al servidor.

- docker:   HOST

            containers  )≡(  images
              ⋮             ⋮

docker  run  → (w) 8080:8080   -v  volume-path #: mount  jupyter/datascience-notebook
                                   ‿                ‿
                                your disk      inside container.

- Kubernetes: manages clusters of hosts (against docker images not having enough server)
  ↳ usable by anyone if app is dockerized.

# Lecture 4. No SQL : Not only SQL.

↳ scalable, fast, flexible, easy...
schemas (vs. SQL: rigid)

- Disadvantage: no standarization rules (each DB has its own query language)
  bad with relational data
  ...

- Types: 1) Key value: big hash map (distributed)   (Key, value) pairs.   (Amazon shopping basket)   → Wikipedia
  2) document-based: stores data as JSON objects. MongoDB.   Elastic Search → indexes to categorize similar data.
  3) column-based: BigTable, Kassandra. SQL-likewise, difference in how it is stored.   → we store rows → each one is HashMap
  4) graph-based: model relationships between entities   (each can have its own columns)
  DATA-MODELING. generally not main DB.   Hash-Map of Hash-Map.
  ↳ fraud detection.

SQL queries can crash with really large joins → MongoDB.
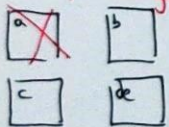Check equivalencies Mongo - SQL.

# Lecture 5. Hadoop & S3

- Hadoop Distributed File System (HDFS) + MapReduce (Computation)
  _____
  storage

do not your data to your local (huge!) ⇒ take your code to the server.

- blocks: storage units (files are divided in blocks).

data distribution
  → if one goes down = corrupted data.



REPLICATE EACH BLOCK THRICE
  ↳ name node shows where each block is. (up to million of files)

  ↳ use this for processing.

- S3: object storage (id, data, metadata)   only to dump data there
  _____      unlimited.
  not splitted.

Simple Storage Service

durable

two types of data: frequent or not accesses. → storage class.
  Standard | Standard-1A, One Zone-1A (for back ups)
  Intelligent Tiering (unknown, adaptative)
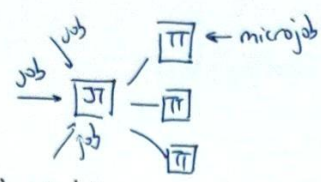  Glacier: logs and archives.

- Hadoop Processing:

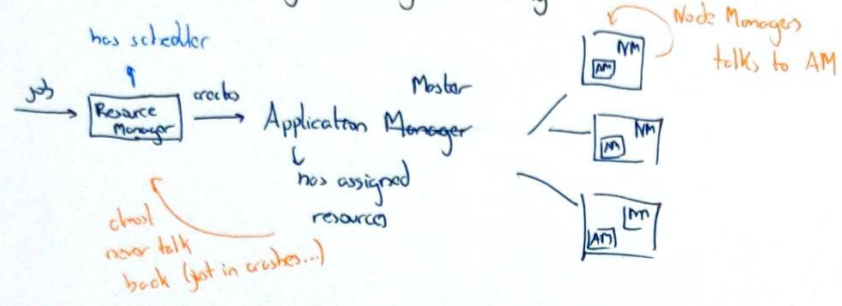  - MapReduce: distribute by key after dividing into machines.
    
    reduce      map

    Input splits → map → shuffle → reduce

    Job/Task-tracker for status managements.

    
    ← microjob

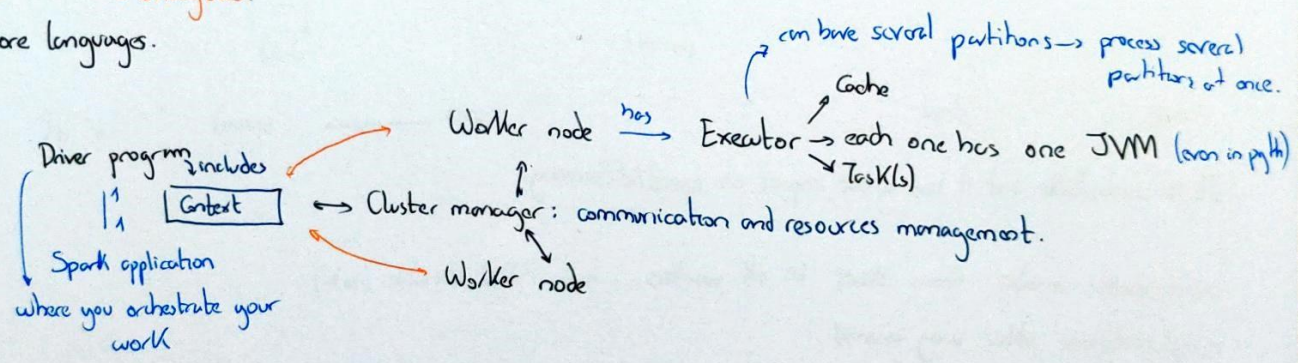  - Hadoop 2: Yarn for resource management + job scheduling.

    has scheduler

    

    Node Managers talks to AM

    NM ≈ TT
    AM ≈ JT (for one job only Hug?)

    chcall near talk back (just in crashes...)

    Spark uses the same architecture.

# L6. Spark Fundamentals

- Not a replacement of Hadoop or BD storage: big data processing and analysis (instead of MapReduce)

  - Automatically distributes the work (fast)
  - Interactive exploration and analytics.
  - Easier code in more languages.

- Architecture:

  

  can have several partitions → process several partitions at once.

  Cache

  Worker node → has → Executor → each one has one JVM (even in python)
  → Task(s)

  Driver program includes → Cluster manager: communication and resources management.
  Context ← → Worker node

  Spark application where you orchestrate your work

  dds: lazily ≡ created only when needed.

- RDD: immutable collection that is partitioned and distributed among workers. → many types.

- Lineage: tree of transformations. → you can recreate a piece if lost and not all. ⟩ DAGs. no cycle graph.

  transformations act on RDDs (data) / action: you want other kind of output count() (if you call twice on action, it will do everything again)

  lazy until you call an

  filter, coalesce (change number of packages)

  (cost) ← save the RDD for reuse!
  (cache ())

  - pair RDD: "dictionaries"

  no loops in RDD (items are distributed)

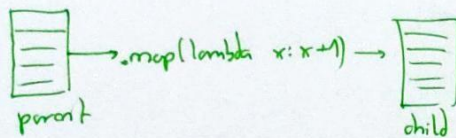  reduce for operations with two items.

# L8. Spark SQL

- Datasets and dataframes
  - Typed: string? int?    untyped.

- Spark Session are the new entrypoint (instead of Spark Context)
  - ↳ better for SQL context.
  - DataFrames can be created from RDDs, hive, data sources
    - ↳ RDDs in the background.
- UDF in SQL: User defined functions.

# L9. Spark Advanced

- Optimization: spark pipeline runs as much code as possible in a some stage.
  
  new stages have to be created when data needs to move across stages (shuffling)

Spark is immutable ⇒ each transformation ≡ a new RDD

- Narrow dependency (fast) each partition of parent RDD is used at most by one partition of child RDD

parent → .map(lambda x: x+1) → child

Wide     "     slow :  "    "    "    "    " used ____ several " s of ...

It is controllable and it has great impact on speed/efficiency.

- broadcast variables: from driver to all workers. (non RDD variables (data))
- accumulators: other way around.

# L10.

- Window SQL Function: group in a window a set of rows.
  - Spark SQL: window functions:
    1) Ranking Functions
    2) Analytic Functions
    3)

- Streaming processing: spark streaming, flink, apache storm.
  - ↳ microbatching: divide stream in batches. with DStreams. (Data Stream)

ou con access to specific RDDs :  [RDD] [RDD] [RDD]

we need