# Big Data Technologies

Trainer: Mr. Nilesh Ghule.

# Big Data Technologies

Contents

① Introduction

② Hive (DWH) — SQL

③ Spark — Python

④ Kafka — Python

⑤ Airflow — Python

⑥ Hadoop — Java → generics, array list & hashmap
exception handling
jdbc
Stream programming
Oops basics,

⑦ HBase —

Prerequisites

Lecture:
8:00 am to 1:00 pm

Lab:
2:00 pm to 7:00 pm

## Evolution - Data engg.

File based data handling

1970+ : RDBMS

1990 : internet + www
Java

1998+ : NoSQL

2000 : MPP

2001- : Google Big Data → 2003 : GFS
                          ↳ 2004 : MR
2004+

2006 : Cloud Computing

## hadoop cluster



## Distributed Systems

* Cluster : set of Computers in a close network doing dedicated task.
* Horizontal scalability
* Distributed Storage + Distributed Computing
* HA, Reliable, Fault tolerance,

## *Distributed Storage challenges :
① Block size
② Data node failure
③ Metadata node fail.

## *Distributed Computing
① Synchronization
② node failure

## * Big Data Characteristics

① Volume
② Velocity
③ Variety
④ Veracity
⑤ Value

data quality →

Structured data
  - fixed schema → RDBMS

Semi-Structured data
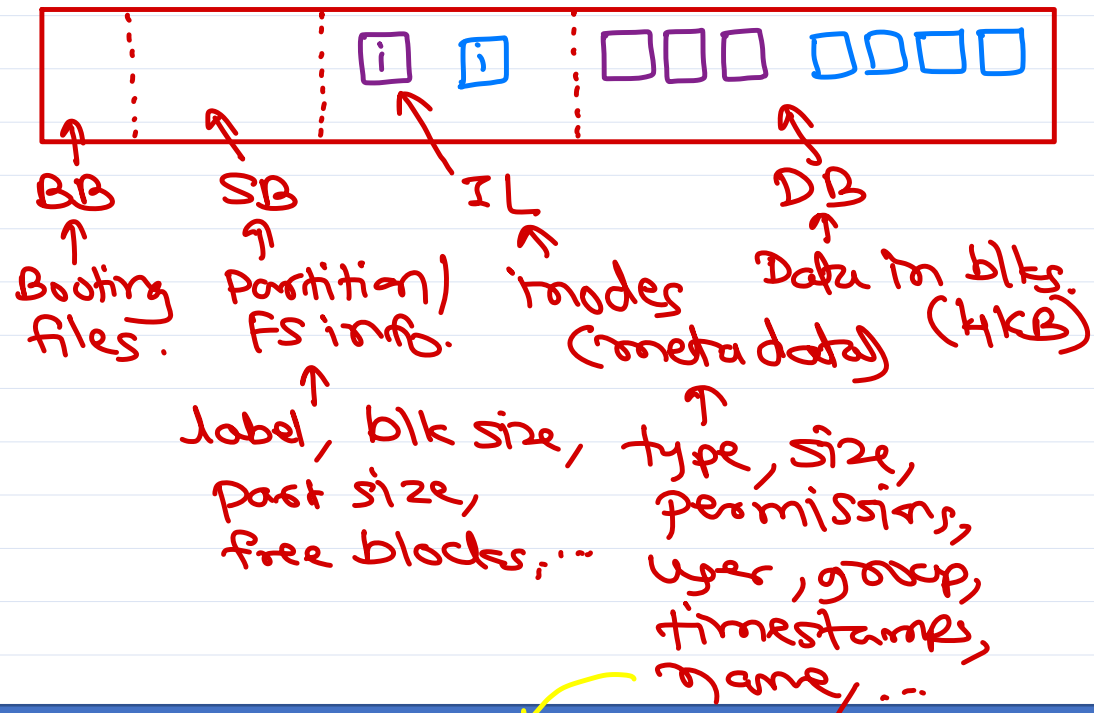  - flexible schema → NoSQL

unstructured data
  - no schema ↳ Big Data

# Local storage vs Distributed storage (HDFS)

File = collection of data/info on a storage device.

File = Data + Info
(contents) (metadata)

File System = Organizing files on disk e.g. FAT, NTFS, EXT3/4,...



BB      SB      IL              DB

↑ Booting files.

↑ Partition/ FS info.
↑ label, blk size, part size, free blocks,...

↑ inodes (metadata)
↑ type, size, permissions, user, group, timestamps, name,...

↑ Data in blks. (4KB)

info about data blocks.

① Data divided into data blocks → stored on multiple nodes in a cluster.

② Data Block Size:
- Bigger block sizes to reduce overheads.
- HDFS block size
  Hadoop 1.x → 64 MB
  Hadoop 2.x+ → 128 MB
- Different /Custom block size can be given to each file.

③ Data Node failure
- Any node failure is handled by redaudent storage.
- HDFS Replication: Each data block copied on multiple Data Nodes. Default: 3
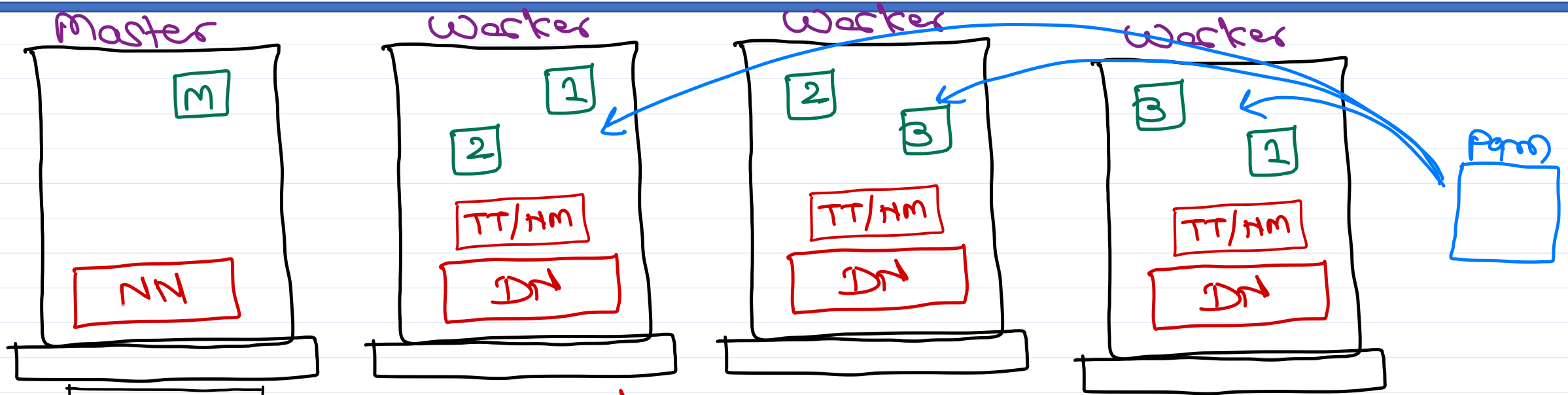- Hadoop 3.x: one more mechanism. i.e. Erasure Coding.
- Replication: Overheads 200%.
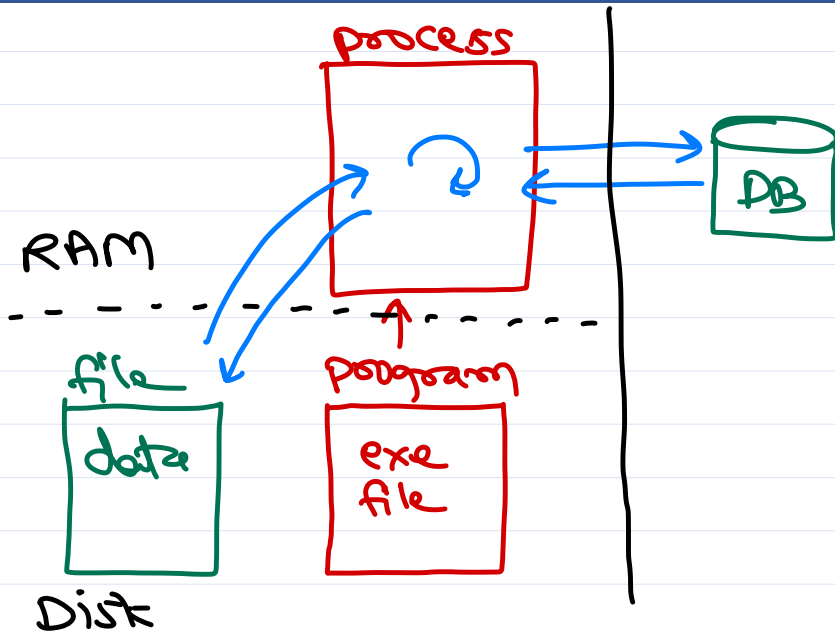- Erasure Coding: overheads 50%. + But need to re compute the lost Block.

# Distributed Storage (HDFS)



**Master**

M

NN

JT/RM

M

SNN

**Worker**

1

2

TT/NM

DN

**Worker**

2

3

TT/NM

DN

**Worker**

3

1

TT/NM

DN

Pgm

④ NameNode failure.
- If metadata is not avail, whole System is down
- Keep metadata backup(s).
- In Hadoop 1.x, NN is SPOF.
  - SNN takes periodic backup
  - On failure, Admin switch SNN to NN (manual)
- In Hadoop 2.x, Standby NN takes active backup of meta.
  - on NN failure, Standby NN can auto become NN.
- Hadoop 3.x can config multiple Standby NNs.

# Local computing vs Distributed computing (Map Reduc
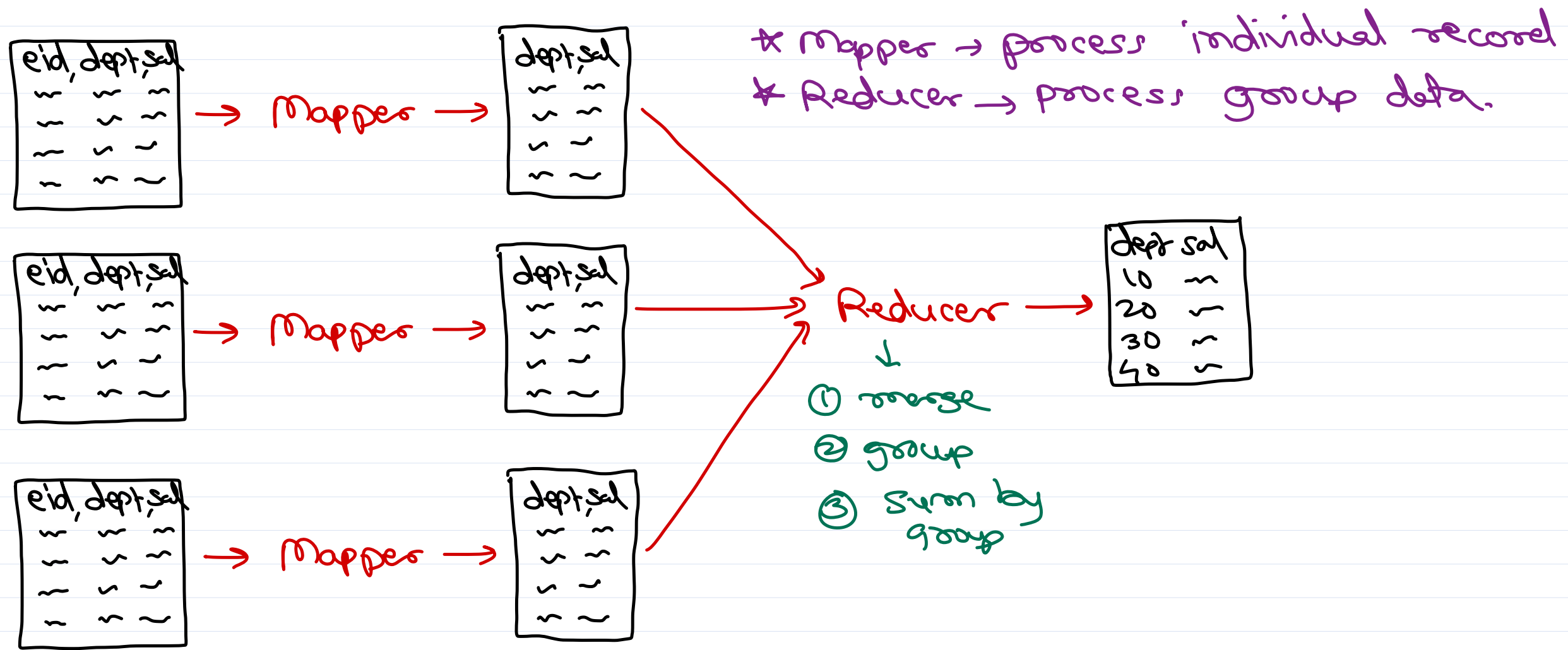


process

RAM

file
data

Disk

program
exe
file

DB

Problems:
① limited memory (RAM)
② limited computing (CPU)
③ limited storage (Disk)
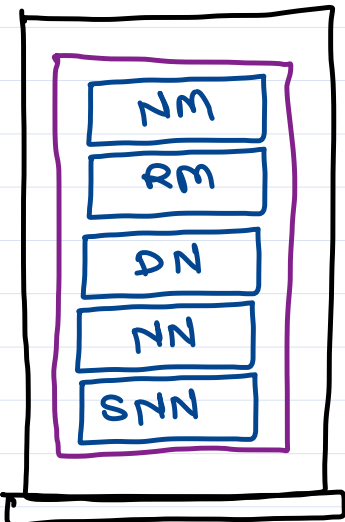④ limited network (Band width)
⑤ Disk/Network
        speed.

✓ Data is distributed in multiple nodes (in blocks).
✓ Program (much smaller in size w.r.t. data size) It will be copied on all data nodes.
✓ Data processing will be planned/tracked/synchronized by special programs.
✓ Each node will have partial data processing and final result will be accumulated later
✓ Any node failure will lead to reassign that task on other node.
✓ Distributed Computing → Map Reduce.

✓ Data to be processed is fetched in program memory and then processing is carried out.
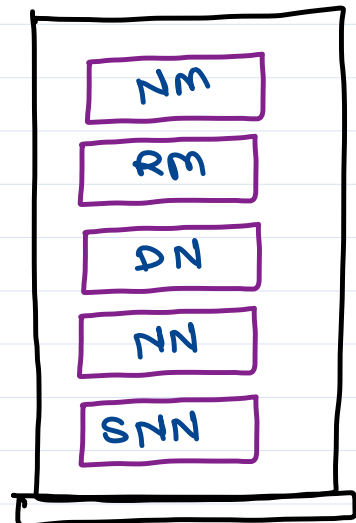✓ Suitable for smaller data volumes.

# Map Reduce - Design pattern



* Mapper → process individual record
* Reducer → process group data.

Reducer
↓
① merge
② group
③ sum by group

| dept | sal |
|------|-----|
| 10   | ~   |
| 20   | ~   |
| 30   | ~   |
| 40   | ~   |

# Hadoop Installation Modes



**Top-left node:** NM, RM, DN, NN, SNN
- ✓ fast (no ipc)
- ✗ small data process only
- ✓ only for testing

local mode

**Bottom-left node:** NM, RM, DN, NN, SNN
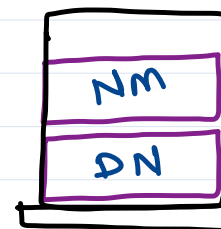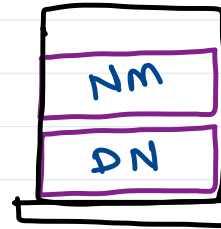- ✓ Dev machine setup
- ✗ slower (ipc)
- ✓ only for dev & testing.
- ✗ Small data process.

single node cluster or psuedo dist mode.

**Top-right (master):** NN, SNN, RM
**Workers:** NM/DN, NM/DN, NM/DN, NM/DN ...

multinode cluster or full dist cluster.

**Bottom-right Workers:** NM/DN, NM/DN, NM/DN and RM, NN, SNN

lab assign cluster

# Thank you!

*Nilesh Ghule <nilesh@sunbeaminfo.com>*