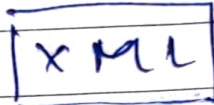


## poly plot Resistance. 2

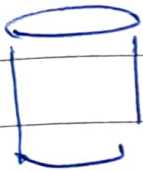
ETL -> Extract Transform Load



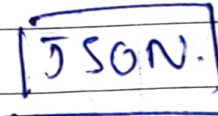
MySQL  
(CRM)



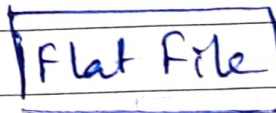
Partner Companies



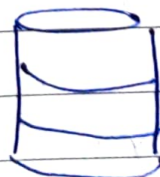
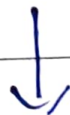
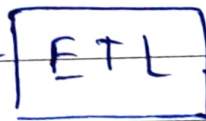
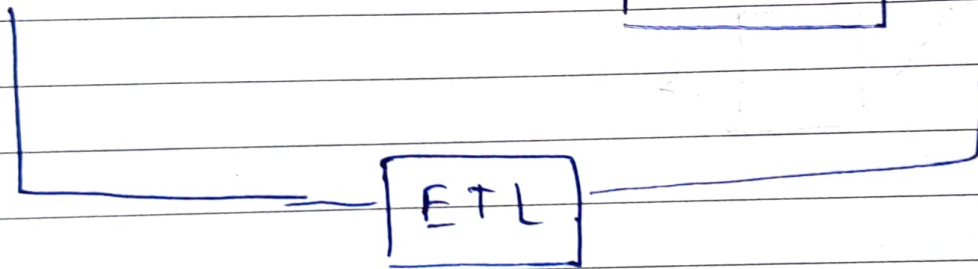
Oracle  
(transactional)



Social Media



Log. files



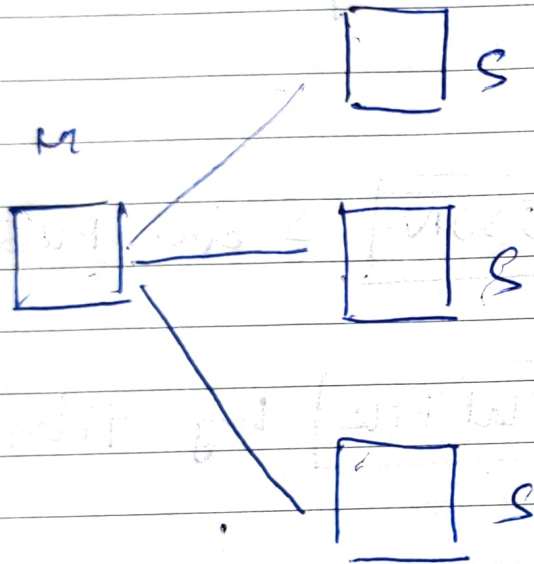
D/W. [Data Warehouse]  
OLAP

BI

MAP → Machine Parallel processing.

Cloudera → Commercial dist of hadoop.

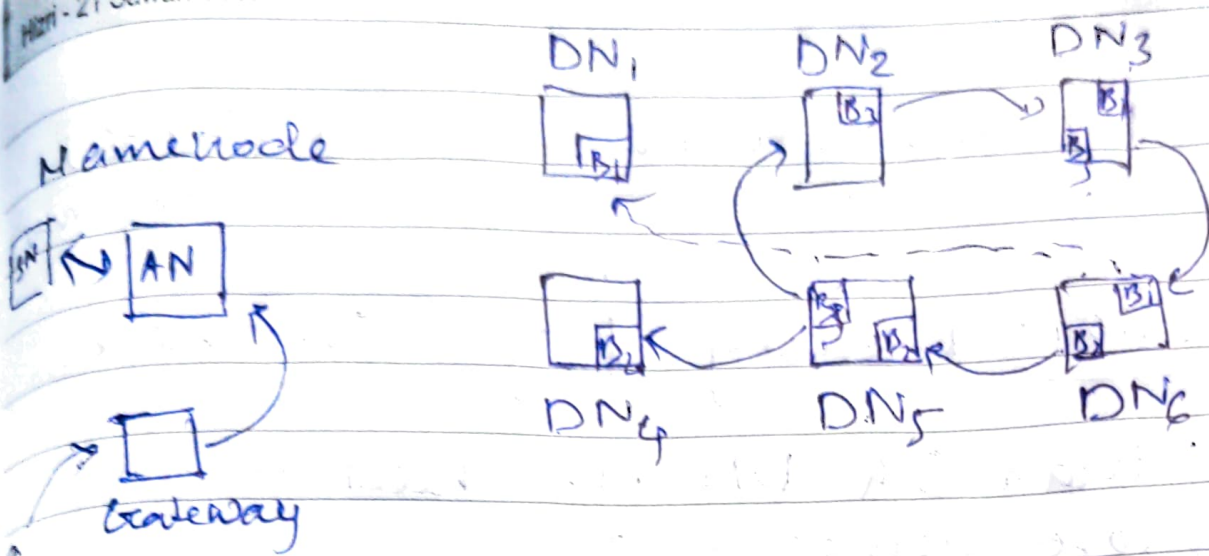
hortonworks, MapR.



HDFS  
↓  
Storage.

YARN  
↓  
Resource Manager.

Map Reduce  
↓  
Processing.

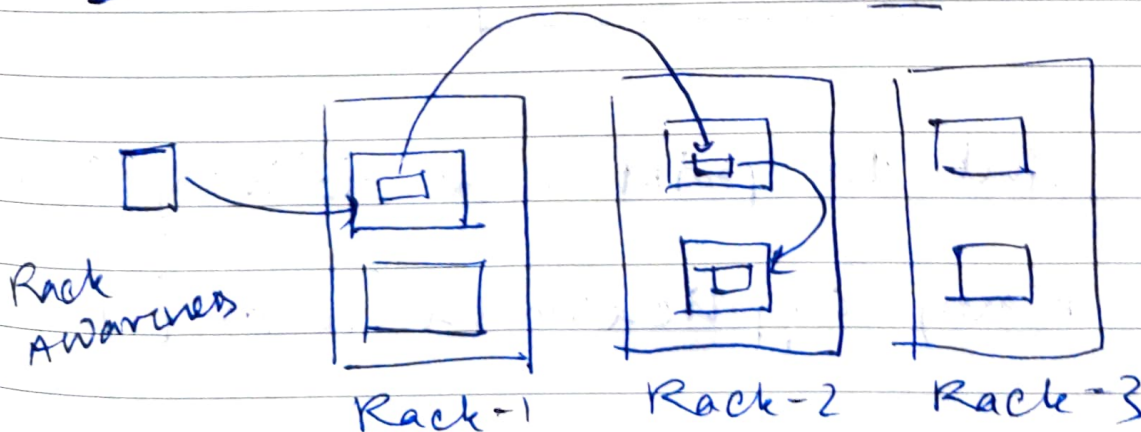


- i) Any format of data can be stored.
- ii) Modification of file is not possible.

Block Size → Max size of data that can be stored.

192 MB →  $\begin{matrix} B_1 & B_2 & B_3 \\ 64 & 64 & 64 \end{matrix}$

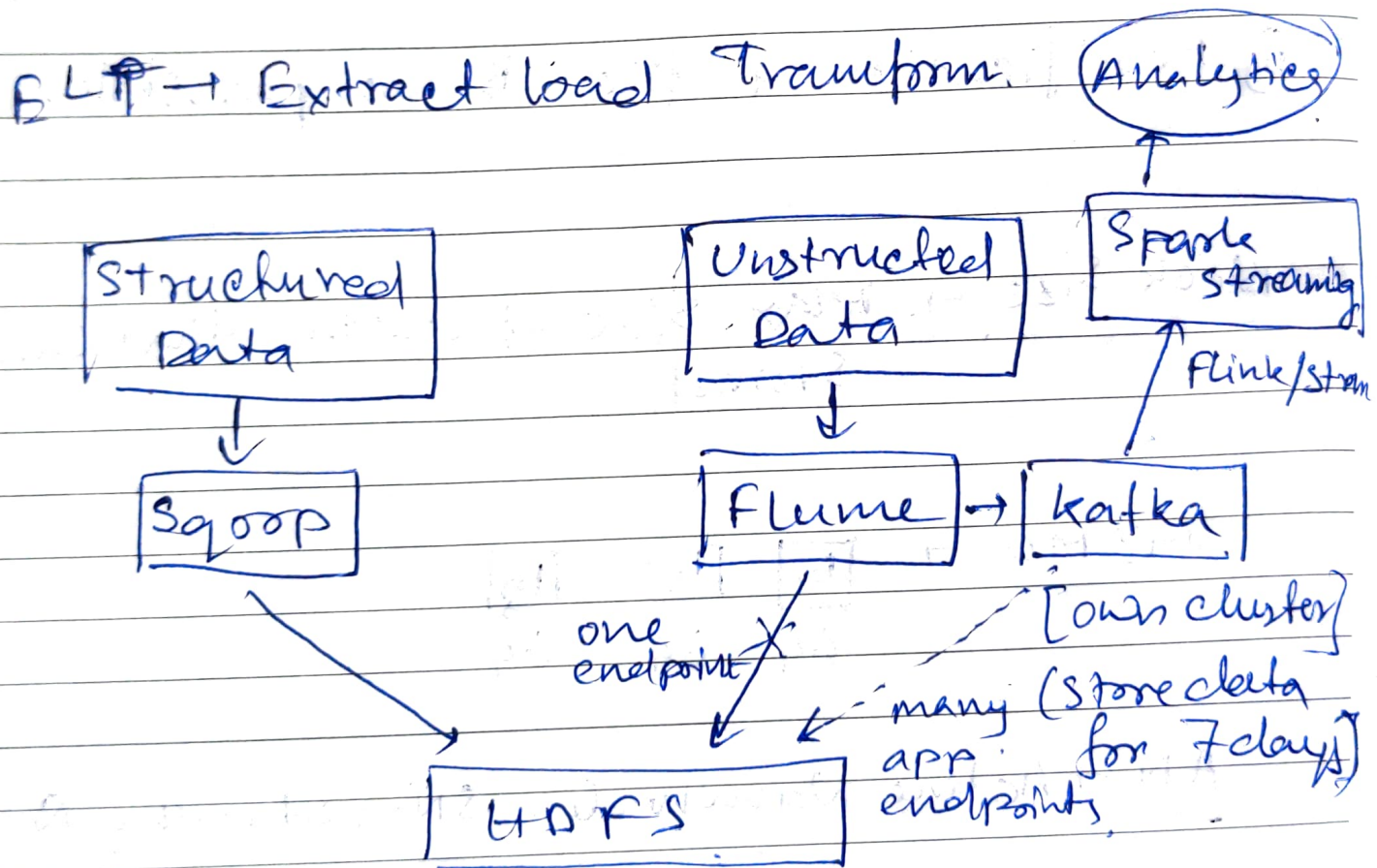
AN → Active NameNode. SN → Standby NameNode.





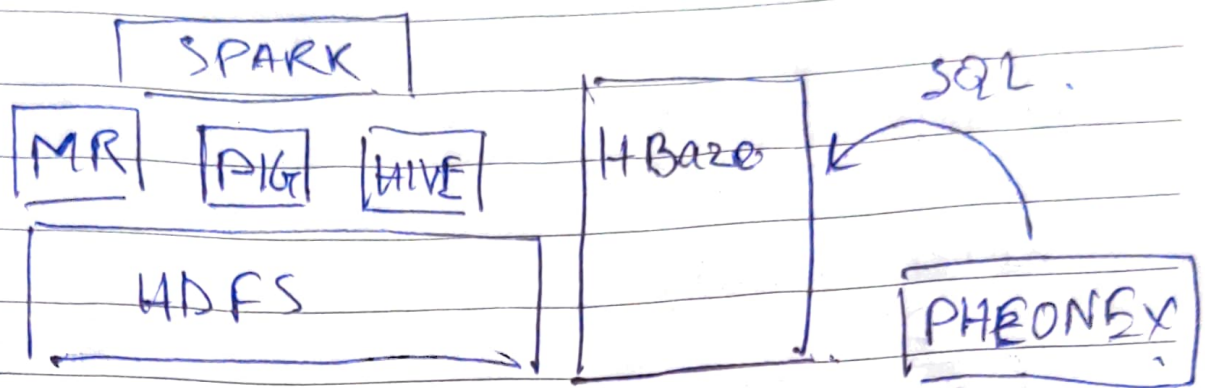
Largest Hadoop Cluster → Yahoo (42k) machines

Federation: Ideally if you ~~can~~ have more than 5k DN you need to add one more name node.



Flume: point to point data delivery system

Kafka: Many application can read from.



HIVE - SQL on top of HDFS

SQL  $\rightarrow$  HIVE (Translator)  $\rightarrow$  MAP RED

SPARK - In memory execution system. (Batch processing)

Big Query  $\rightarrow$  HIVE

Small Query  $\rightarrow$  Impala | HAWQ, LLAP | DRILL  
~~PHEONEX~~ etc.  
 PHEONEX