

HADOOP VS. TRADITIONAL DATA STORAGE AND
PROCESSING

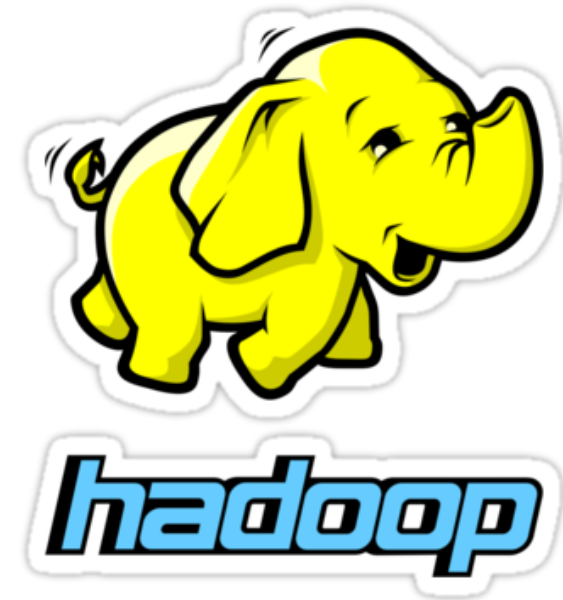
WHAT IS HADOOP?

HADOOP IS FUNDAMENTALLY AN OPEN-SOURCE INFRASTRUCTURE SOFTWARE FRAMEWORK THAT ALLOWS DISTRIBUTED STORAGE AND PROCESSING A HUGE AMOUNT OF DATA I.E. BIG DATA.

IT'S A CLUSTER SYSTEM WHICH WORKS AS A MASTER-SLAVE ARCHITECTURE.

HENCE, WITH SUCH ARCHITECTURE, LARGE DATA CAN BE STORED AND PROCESSED IN PARALLEL.

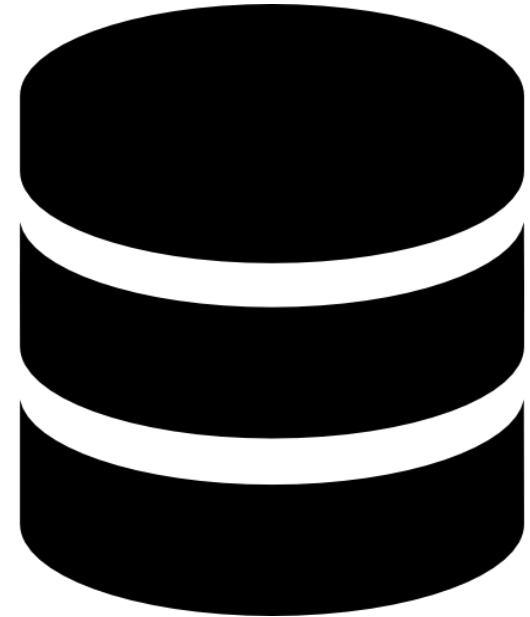
DIFFERENT TYPES OF DATA CAN BE ANALYZED, STRUCTURED (TABLES), UNSTRUCTURED (LOGS, EMAIL BODY, BLOG TEXT) AND SEMI-STRUCTURED (MEDIA FILE METADATA, XML, HTML).



WHAT IS RDBMS?

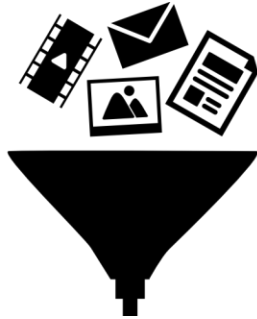
RDBMS STANDS FOR THE RELATIONAL DATABASE MANAGEMENT SYSTEM. IT IS A DATABASE SYSTEM BASED ON THE RELATIONAL MODEL SPECIFIED BY EDGAR F. CODD IN 1970. THE DATABASE MANAGEMENT SOFTWARE LIKE ORACLE SERVER, MY SQL, AND IBM DB2 ARE BASED ON THE RELATIONAL DATABASE MANAGEMENT SYSTEM.

THE DATA REPRESENTED IN THE RDBMS IS IN THE FORM OF THE ROWS OR THE TUPLES. THIS TABLE IS BASICALLY A COLLECTION OF RELATED DATA OBJECTS AND IT CONSISTS OF COLUMNS AND ROWS. NORMALIZATION PLAYS A CRUCIAL ROLE IN RDBMS. IT CONTAINS THE GROUP OF THE TABLES, EACH TABLE CONTAINS THE PRIMARY KEY.



#1 DATA VARIETY

HDOOP



RDBMS



Used for Structured,
Semi-Structured and
Unstructured data

Mainly for Structured
data

#2 DATA STORAGE

HDOOP



RDBMS



Use for large data set
(Tbs and Pbs)

Average size data
(GBS)

#3 QUERYING

HDOOP



RDBMS

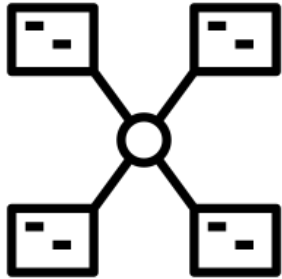


HQL (Hive Query
Language)

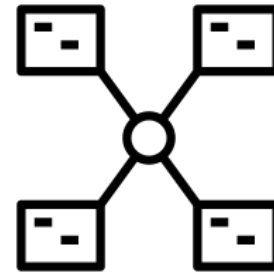
SQL Language

#4 SCHEMA

HDOOP



RDBMS



Required on reading
(dynamic schema)

Required on write
(static schema)

#5 COST

HDOOP



Free



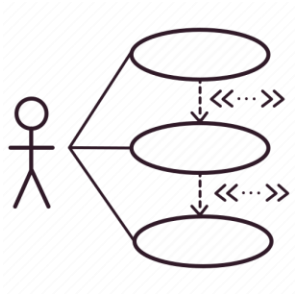
RDBMS



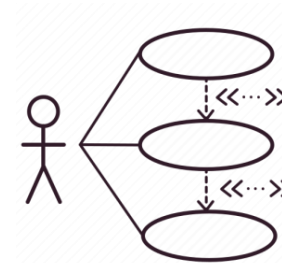
License

#6 USE CASE

HDOOP



RDBMS



Analytics (Audio, video, logs etc), Data Discovery

OLTP (Online transaction processing)

#7 SPEED

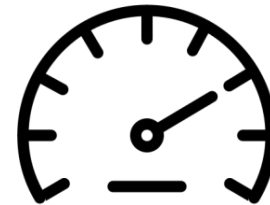
HDOOP



Both reads and writes
are fast



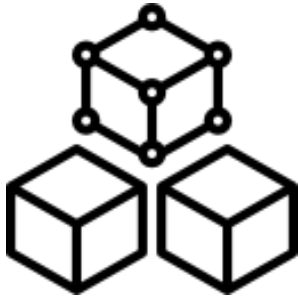
RDBMS



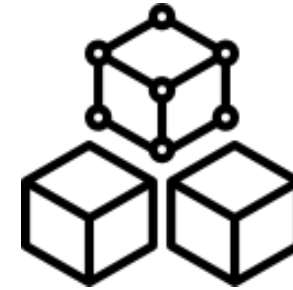
Reads are fast

#8 DATA OBJECTS

HDOOP



RDBMS



Works on Key/Value
Pair

Works on Relational
Tables

#9 THROUGHPUT

HDOOP



RDBMS

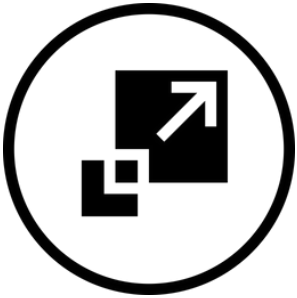


High

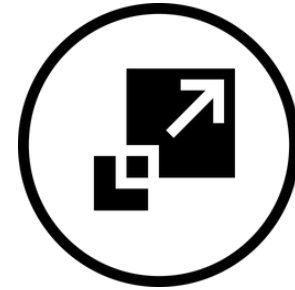
Low

#10 SCALABILITY

HDOOP



RDBMS



Horizontal

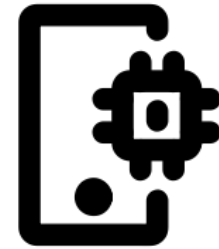
Vertical

#11 HARDWARE PROFILE

HDOOP



RDBMS



Commodity/Utility
Hardware

High-End Servers

#12 INTEGRITY

HDOOP



Low



RDBMS



High (ACID)



SO...

BY THE ABOVE COMPARISON, WE HAVE COME TO KNOW THAT HADOOP IS THE BEST TECHNIQUE FOR HANDLING BIG DATA COMPARED TO THAT OF RDBMS.

AS DAY BY DAY, THE DATA USED INCREASES AND THEREFORE A BETTER WAY OF HANDLING SUCH A HUGE AMOUNT OF DATA IS BECOMING A HECTIC TASK.

ANALYSIS AND STORAGE OF BIG DATA ARE CONVENIENT ONLY WITH THE HELP OF THE HADOOP ECO-SYSTEM THAN THE TRADITIONAL RDBMS.

HADOOP IS A LARGE-SCALE, OPEN-SOURCE SOFTWARE FRAMEWORK DEDICATED TO SCALABLE, DISTRIBUTED, DATA-INTENSIVE COMPUTING.

THIS FRAMEWORK BREAKDOWNS LARGE DATA INTO SMALLER PARALLELIZABLE DATA SETS AND HANDLES SCHEDULING, MAPS EACH PART TO AN INTERMEDIATE VALUE, FAULT-TOLERANT, RELIABLE, AND SUPPORTS THOUSANDS OF NODES AND PETABYTES OF DATA, CURRENTLY USED IN THE DEVELOPMENT, PRODUCTION AND TESTING ENVIRONMENT AND IMPLEMENTATION OPTIONS.