

Big data Modeling and Management

Week 1 :- 1<sup>st</sup> part :- Summary of Introduction to Big data

2<sup>nd</sup> part :- Big data management :-

- Describe what "data management" means
- Identify the primary issues involved in the management of big data.

\* What is Data Management ?

- How do we ingest the data ?
- Where and how do we store it ?
- How can we ensure data quality
- What operations do we perform on the data ?
- How can these operations be efficient
- How to scale up data vol<sup>n</sup>, variety, velocity and access ?
- How to keep the data secure ?

Ingestion Infrastructure :-

- ① How many data source ?
  - ② How large are data items
  - ③ Will the no. of data sources grow
  - ④ Rate of data ingestion ?
  - ⑤ What to do with bad data ?
  - ⑥ What to do when data is too little or too much ?
- Storage Infrastructure :- where hardware meets Data management
- How much data to store
  - How fast do we need to read/write

\* Memory Hierarchy



Data Quality :- • Better quality means better analysis  
decision making means needed for

M/T/W/T/F/S/Su • Quality assurance means needed for  
regulatory compliance and interch

• Quality leads to better engagement and interch  
with external entities

\* Operations of Data :-

- Operation on single data items that produces a sub-item
  - Operation of collections of data items
  - Operation that select a part of a collection
  - Operations that combine two collections
  - Operations that compute a function on a collection
- \* Efficiency of data operations
- Measured by time and space
  - should use parallelism

Achieving Scalability

Scaling up and scaling out

① Vertical scaling (scale up)

① Horizontal scaling (scale out)

Adding more processors and Adding more, possibly less  
RAM, buying a more expensive powerful machine that  
and robust server. interconnected over a network

② many operations perform better ② Parallel operations will  
with more memory, more possible be slower  
cores

Maintenance can be difficult, ③ Easier in practice to add  
expensive more machine

\* The server industry has many sol's for scale up/scale out  
keeping data secure :-

Data security :- a must for sensitive data  
Increasing the no of machines leads to more security

Data in transit must be secure

Encryption and decryption ↑ security but

more data operations expensive

• Clever men are good, but they are not the best •