

Twitter → #tags

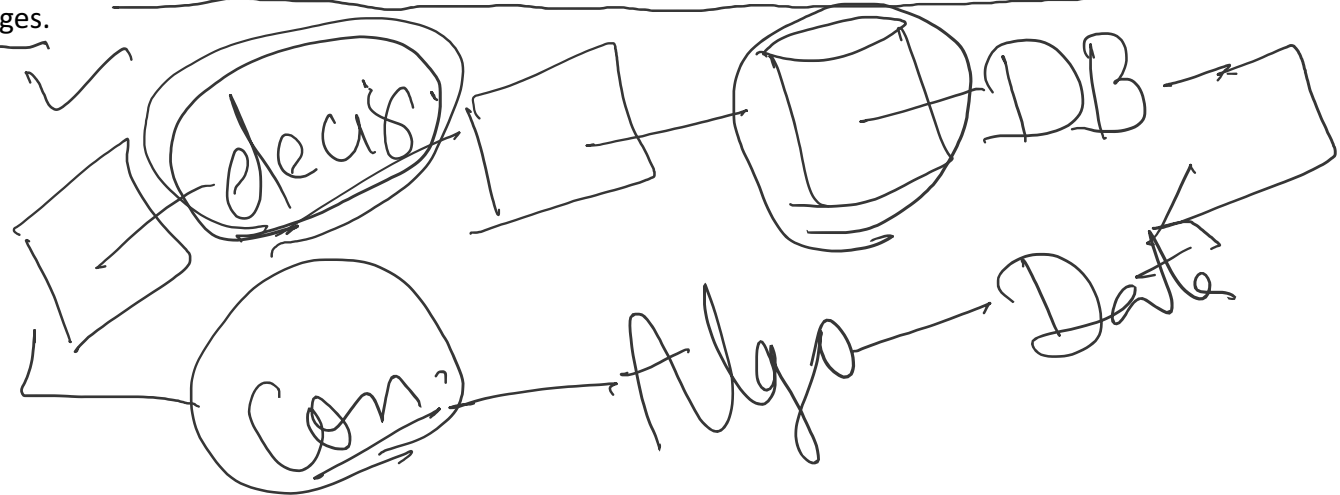
Big Data Analytics 8CS421

How does
big data analytics
works?
lab

GitHub : https://github.com/technoindianjr/Big-Data-Analytics-Lab_8cs4-21

With the advance of IT storage, processing, computation, and sensing technologies, Big Data has become a novel norm of life. Only until recently, computers are able to capture and analysis all sorts of large-scale data from all kinds of fields -- people, behaviour, information, devices, sensors, biological signals, finance, vehicles, astrology, neurology, etc. Almost all industries are bracing into the challenge of Big Data and want to dig out valuable information to get insight to solve their challenges.

Raw data
Pattern
DS. Data



Definition and Characteristics of Big Data

"Big data is **high-volume**, **high-velocity** and **high-variety** information assets that demand **cost-effective**, **innovative** forms of information processing for **enhanced insight and decision making**." -- Gartner

which was derived from:

Meta Data Res (Prod)
→ Date / Time

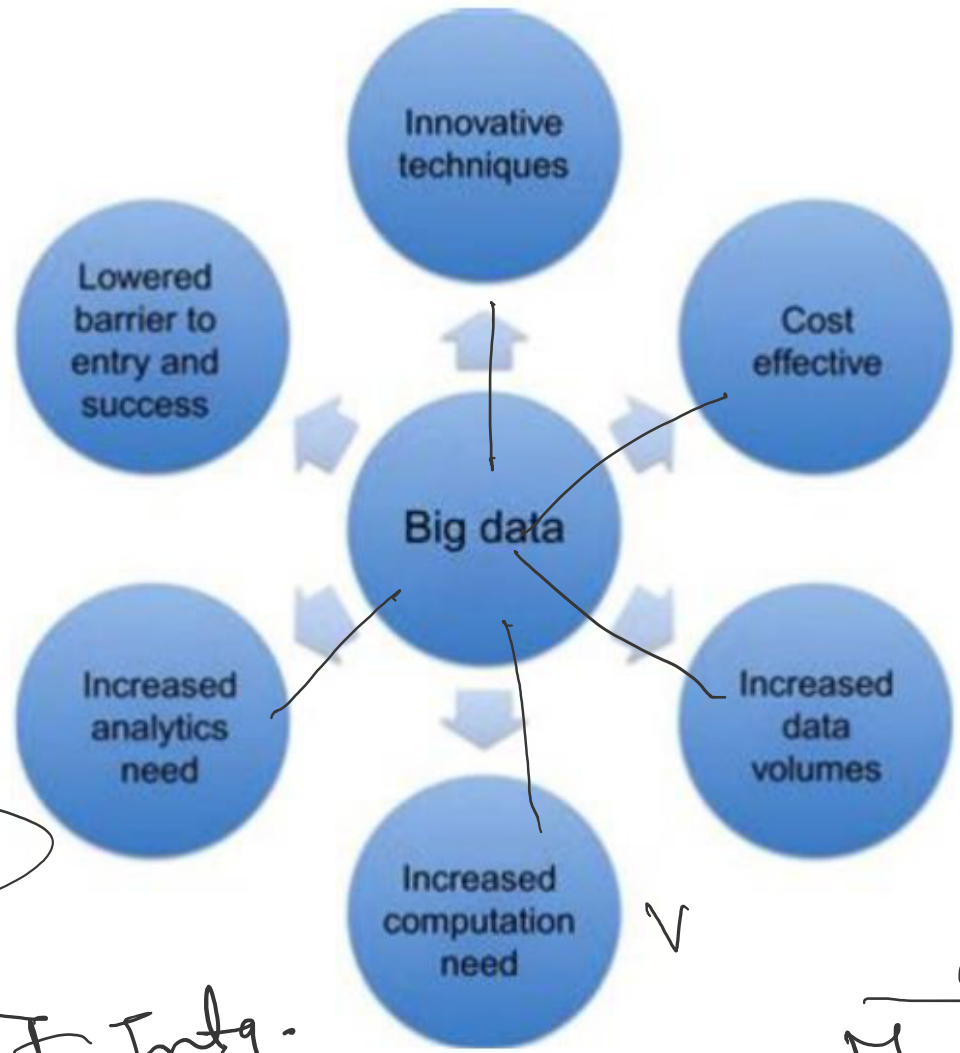
#tag → Video URL href

24x7x365
← Text / Img

"While enterprises struggle to consolidate systems and collapse redundant databases to enable greater operational, analytical, and collaborative consistencies, changing economic conditions have made this job more difficult. E-commerce, in particular, has exploded data management challenges along three dimensions: **volumes**, **velocity** and **variety**. In 2001/02, IT organizations much compile a variety of approaches to have at their disposal for dealing each." – Doug Laney

Sensors Health Loc → V TV Streaming

What made Big Data needed?



SPARK
↑
AWS EMR
MongoDB
Hadoop

V3

PA

SA

DA

NoSQL

Data Lake

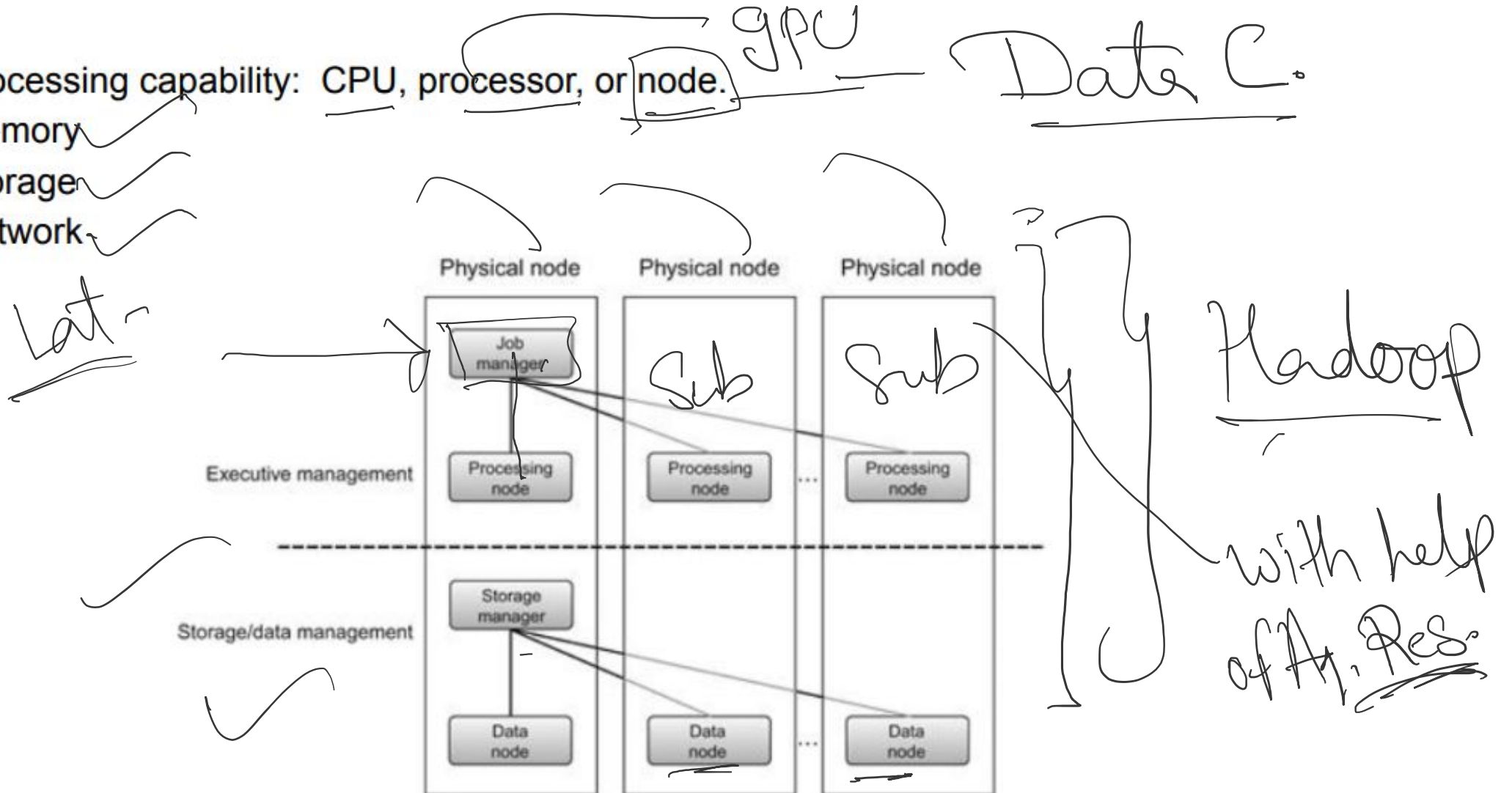
Data Integ.
Software

Data Visual.

BI
Mining
tools

Key Computing Resources for Big Data

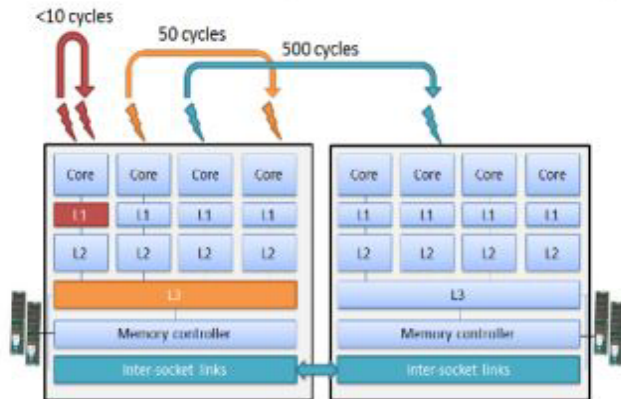
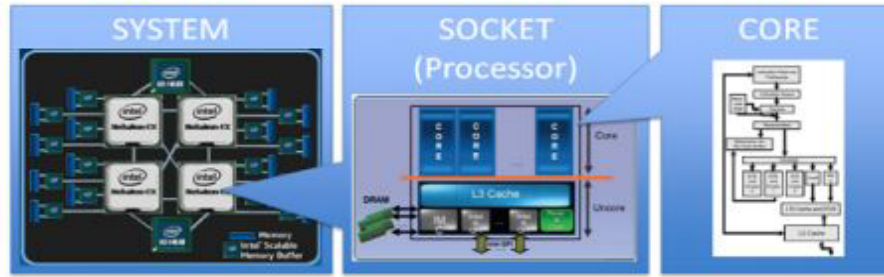
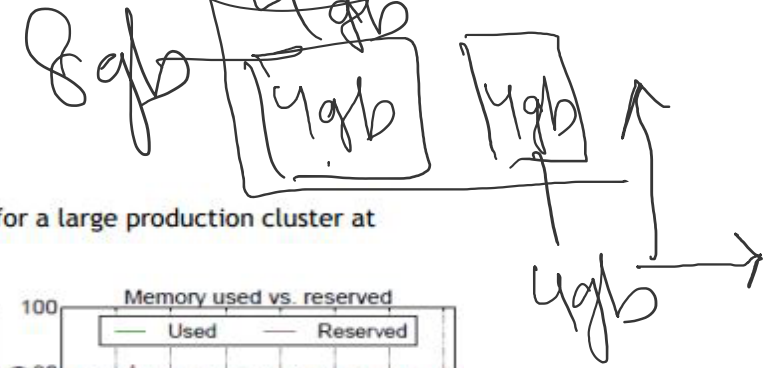
- Processing capability: CPU, processor, or node.
- Memory
- Storage
- Network



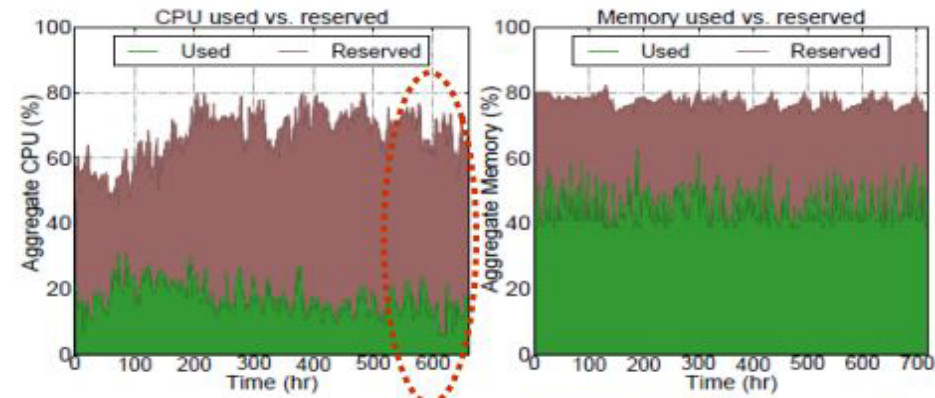
Scalability — Scale Up & Scale Out

- Scale out
 - Use more resources to distribute workload in parallel
 - Higher data access latency is typically incurred
- Scale up
 - Efficiently use the resources
 - Architecture-aware algorithm design

Cap. - 4gb ~ 8gb → 2gb



Example: Resource utilization for a large production cluster at Twitter data center



www.stanford.edu/~cde/2014.asplos.quasar.pdf

- For independent data ==> scale up may not have obvious advantage than scale out
- For linked data ==> utilizing scale up as much as possible before scale out