## Small Data

1970 - Sci EF Codd
Traditional db was designed
RDBMS:- was evolved
→ Structured data → Tables
/Excel files → format
not used these days. (only
10% is structured these days)

unstructured - pics, videos, text

- Store in MB, GB, TB
→ Speed/velocity
(Speed of data increase in
gradually)
- If a amount of data
- is there. 90% of $x$ is
generated recently in last
4-5 years.
- Stored in centralized form
(Structured form)
attendances, library data
- data is locally stored

- Sql server, Oracle,
- Single Node

## (size) Big Data ①

mostly unstructured (90%)
→ 1000TB → PB → EB
petaByte ErabyE
fB, google db
4B data is generated daily
More than 2 Billion users
storage is very high

- Speed of data generated
daily is very high
- Exponentially.
→350 million pics are
uploaded
→4 million likes daily on
fB

- data centres across the
world. by distributing
globaley. Entire traffic
will come at single
place otherwise.
- Hadoop Spack,
Mapreduce, Bigquery
(- Multimode cluster

(These are under
apache now and are
open to use.

↓↓
for DS, ML, AI
↓
train - expert s/w

eg Amazon: many orders → transactions — how to
retain → by big data.

- Big data cant be processed by existing traditional techniques
- Sources :- Stock Exchange (many companies - their own price - their own status) — Big data
  - Social Media (fb, WA, Instagram)
  Billions of uses - post — their own data.
  (Generated daily)
  - Video Sharing Portal: You tube → food, travel, etc.
  many videos - high data
  - Search Engine : Yahoo / google
  eg most searched football player - comes from db
  - Transport : no, owner, distance
  - Banking : we have our account — Transaction
  -

Semi structured data — neither structured nor unstr
that it can be stored in a
Relation
→ JSON, XML is used in this case

Applications :
① call centre (requirements) phn — 2 slots of sim
  2 sim
  high data w.r.t. every no → Data is analysed

| R1 | R2 | R3 |
|----|----|----|
| D | D | D |

Acc. to this calls are made.
Quality Evaluation — recording is also done.
② Social Media :- like - opinion ( companies take their decision for promotional campaigns — sales ↑ — profit ↑

- Shopping :- Amazon, flipkart - Playstore
  variety of items- recommendation
- IoT :- Smart devices - connect - functionality
  ↓
  achieve a
  goal

Data - sent on cloud - decision
eg Healthcare eg ~~Human~~ (Pi) Disease.
  Medication is given to all patients
  ↓
  Sensors are used to capture improvement
  ↓
  decision to continue vaccine or not

- Education :- online exam.
  Particular ~~time~~ - behaviour is recorded.
  - how many qs u take to solve easily
  - how many do u leave difficult (-ve marking)
  - easy option - behavioral analysis
  1 hr lecture - how many pay attention
  facial behaviour - how much % attended

- Entertainment : Netflix.
  earlier only used in foreign - no indian content
  ↓
  demand of hindi movies - data is generated.

challenges of Conventional Sys :-
  → Store     → Manage     → Analyse
  → Specific Time Interval

① large quantities → impossible to process every
② Must be meaningful & shd be collected in
   real time
③ Data is collected from multiple sources

-eg Text, Img, Video, Audio, Video ⇒ Bring it together. Data shd be categorised properly. Requires lot of manpower + time.
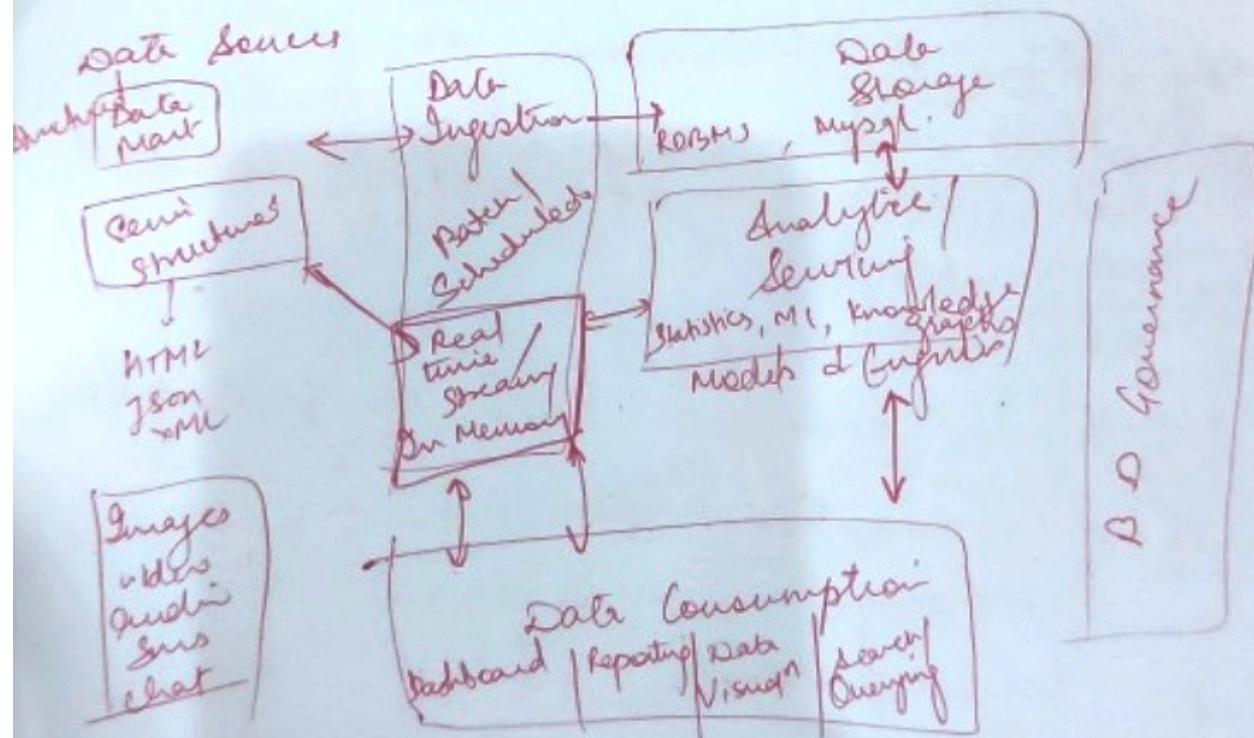- Time can be long.

④ Collect correct data (Remove incompleteness) wrong data else it creates problem.

⑤ Comparison of data using statistics /graphs etc (Complicated job)

⑥ Managers must be able to access data

⑦ Pressure from top (Pressure on risk managers)

⑧ Lack of knowledgeable professionals

To overcome ⇒ Intelligent Data Analysis is used Reduces time, cost & error.

## BIG DATA ARCHITECTURE



Data Sources

Archives Data Mart

Semi structures

HTML JSON XML

Images Video Audio sms chat

Data Ingestion

Batch/ Scheduled

Real time Streaming

On Memory

Data Storage

RDBMS, mysql.

Analytic Serving

Statistics, ML, knowledge graphs

Models & Insights

Data Consumption

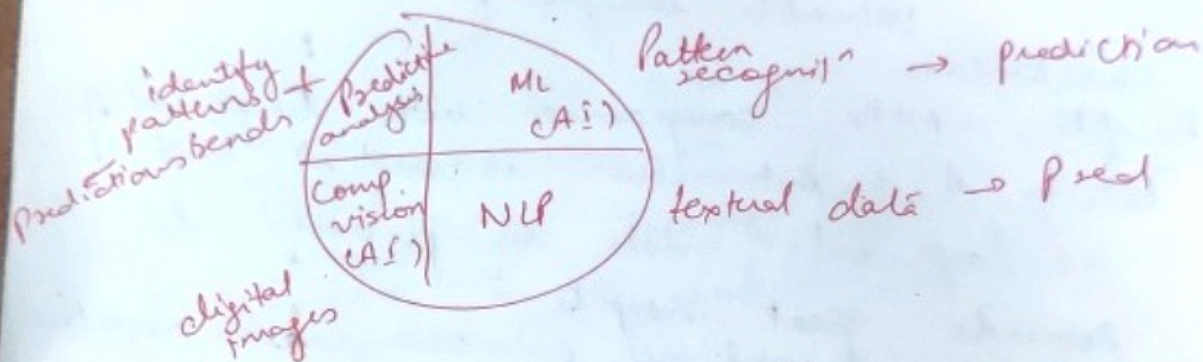Dashboard | Reporting | Data Visual^n | search/ Querying

D Governance

## Technology components

4 components:-
- data capture (social media, sensor readings)
- data storage
- data processing
- data visualization

## fields

identify patterns + trends
Predictions

Predictive analysis

ML (AI) → Pattern recognition → prediction

Comp vision (AI)

NLP → textual data → Pred

digital images

## Technologies

① Apache Hadoop (OS framework) enables distributed processing of large ds on servers. Scalability, CE, flexibility

② Apache Spark OS processing engine (batch + RT analytics on large ds) used with hadoop

③ Apache flink OS stream processing framework high speed analytics on online data streams. User friendly API and scalable arch → growing community

④ Presto : OS SQL engine that supports interactive analysis on huge ds stored in multiple sources. Distributed query processing arch → low latency & strong performance

⑤ Druid : OS Analytical data storage designed for OLAP queries on event based data (eg log files, clickstreams)
fast aggregations & explorations on large ds due to columnar storage format

Others⇒
Apache Hadoop, NOSQL db, MapReduce.
Cloudera, Hortonworks, IBM Biginsyhts,
MapR, Oracle Big Data Appliance.

## Auditing & Analytics :-

↑ use of big data → red$^n$ in oper$^n$ costs, Improved decision making, high customer retention
↑ satisf$^n$ rate.

Banking → Manuf.

RT big data → growth & sustainability

Auditing :- past events → future outcomes. valuable insights for auditors & stakeholders

Benefits :- make comparison with large vol of data. AI + Automation is used in b.d.
→ large vol. of data are processed to provide great insights
Decisions are based on previous cases regarding non compliance & fraud.
latest policy changes — upto date.
financial auditors — can adjust their reporting
process & spot fraudulent transaction
Perform accurate audits
before Analytics → use eff. data aggreg$^n$ & mangmt Sys.
high quality data & Authentic info$^n$.

Human error is overcame by automating.

Analytics :- strong Cybersecurity policy to protect data.
→ change Configurations (Considers x)
→ N/w security (Anti malware) Appls
→ user access
→ Priviledged access mgt.

① set a goal
② Recruit Compete team &
③

① **Features**

**Data wrangling & prep[n]**

before using iterative model.

During model const[n]

② **Data exploration**

visualize data to find insights by using interactive dashboards.

③ **Scalability**

uses less n/w gea & uses less energy → use fast processor & memory

④ **Support for various types of Analytics**

firms make better business decision

⑤ **Version Control :**

Keep track & control changes to s/w code

SD Team keep track of changes

⑥ **Data Management**

Obtain

Store

use data

in cost effective

effective & secure way

Optimizes use of data

⑦ Data Integr[n] :- combine data → make it Simpl

enhance data quality, free resources, lower

IT cost ⑧ Data Governanc:- Accurate, available

usabl

Encryption 6
authorized key (9) Security :- Digital data unauthorized users +
intrusions with firewalls, strong user auth^n, IDS
(10) Visualiz^n :- data driven env.

Appl^ns :

(1) monitoring customer spending + shopping
behaviour

(2) Rec^n

(3) Smart Traffic Sys^m :-
GPS in Ola, Uber
Siri on Apple, Cortana on windows
google asst on Android.

(4) IoT : Sensors- healthcare
(5) energy sector : in every 15 mins
power consump^n to server
Industrial facility → at nyt household
can run heavy machinery.

Compliance :- Collect^n → store → processing
↳ requires more resources.
– Prev^n of fraudulent activities
– managing 3^rd parties risk
– helps in managing cost.
– Disclaimer

Thank You