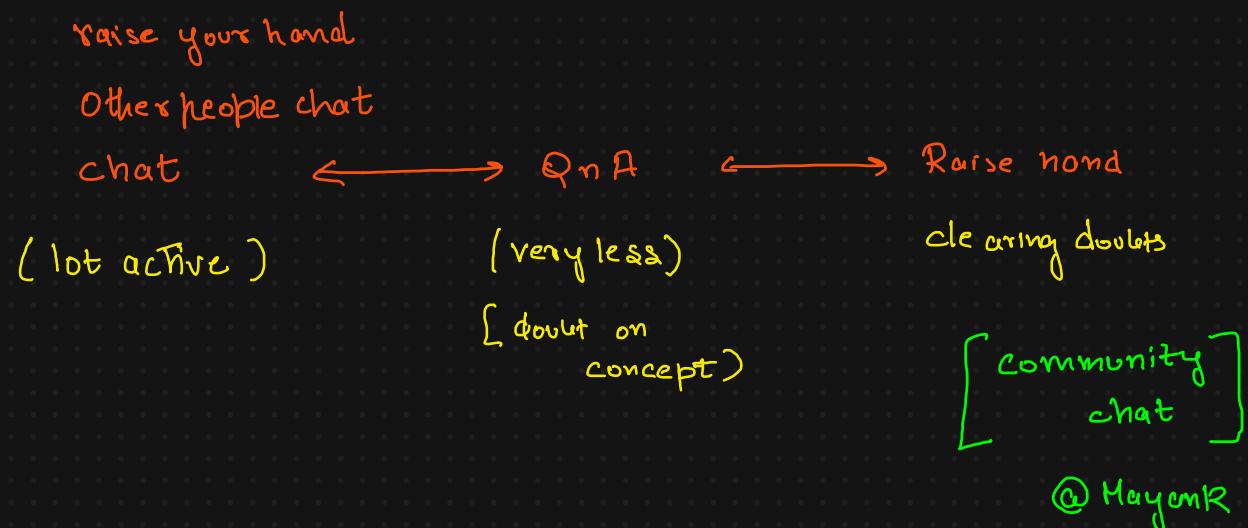


Class-1 - Big Data

✓
✓
✓

Zoom



Agenda

- Who am I?
- Introduction & Getting into Big Data

Mayank Aggarwal

→ 2012 → Coding
love maths → Coders

2013

NSIT, Delhi → 1st year coded
2nd year Data Science → IIIT-D
3rd → ML, DL, AI, DSA
4th → PPO + → DYO

Professional

DYO → 5 months

Goldman Sachs → Big Data or Data engineering

Mindhive → Spark + AI/ML

Ineuron

Scaler + Coding Ninjas

Course

1. Notes + Code + hints + Recording + files + pdf

↳ Critique
↳ Resource section

2. Class 8 - 11 Doubt session till 12

Few Instruction / Good practice

1. Be a leanR paper
2. Make your own notes. [use pens]
3. Don't jump to future
4. Theory + Practical → Projects
5. Great Teacher → everything + how to learn +
vs
Good Teacher → everything
6. Dumb with doubts cleared than smart with having doubts

Feedback → -ve skewed

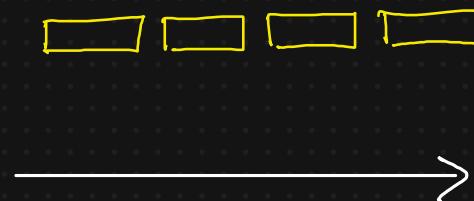
1 we have syllables which is shared

Big Data

Data Science

Kaggle, mail, pendrive, sqe

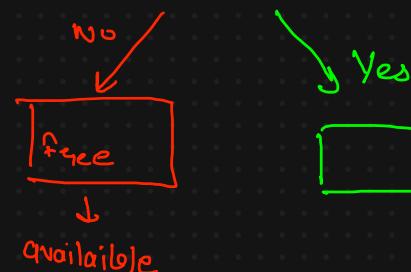
WhatsApp Chat Automation



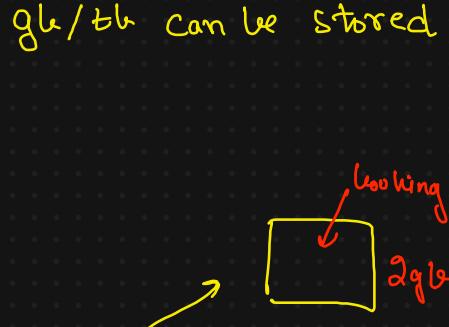
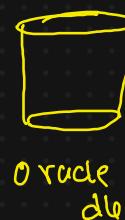
Hi Mayank, we are looking forward to host you for your booking ID SCTP8012 at OYO 12191 Hotel Airlift, Delhi. To confirm your arrival, please reply 'confirm' to this message. If your plans have changed and you need to cancel this booking, please reply 'cancel'.

Problems

- transaction
- small volume

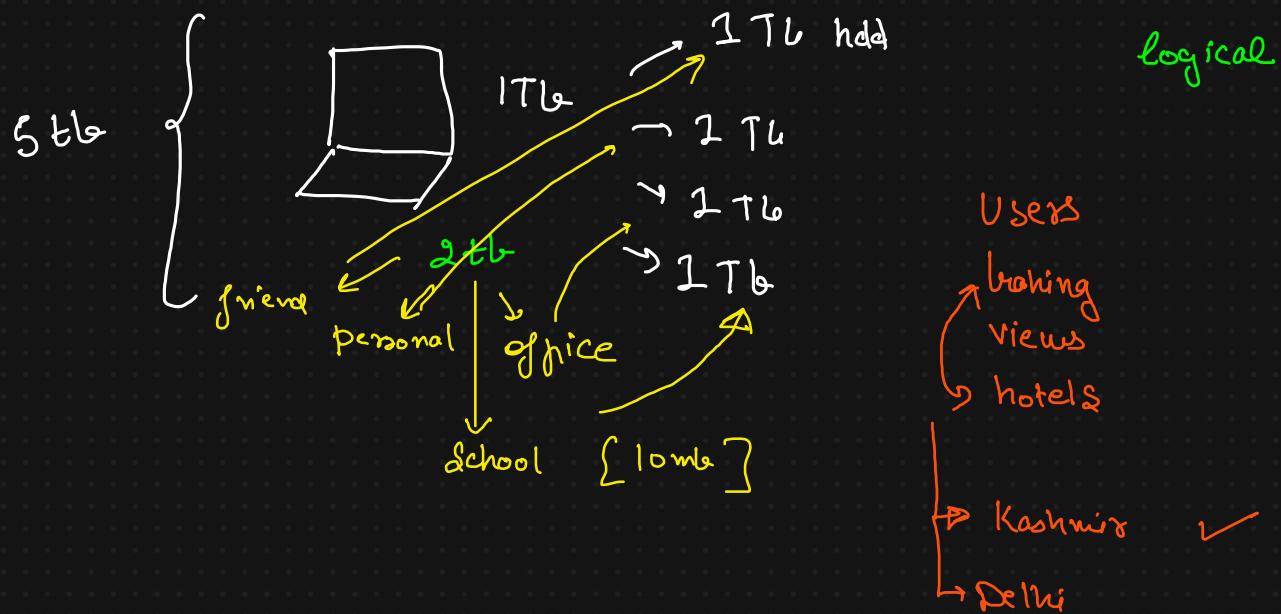


2nd project → Delhi - 2x



:
Storage

- ① 1. There is still a limit which we can expand
- 2. as data increases we have to position our data
- 3. Unstructured data ✓
- 4. Costly.



Big data : huge amount of data that cannot be handled by traditional systems.

These systems deal with data that is too large & complex for traditional system to handle

Big Data is problem
& solution

Why Big data emerge?

1. Internet
2. Social Media
3. IoT

5 V's of Big Data

5 Vs of Big Data are crucial / provide a framework for understanding challenges & characteristic of Big Data.

1. Volume :- Size of data getting generated

OYO vs GS

there is no set volume
when your problem
becomes a big data

2. Velocity: Speed at which data is getting generated & needs to be processed.

Real Time

Bank Transaction

Fraud

live feed

Near Real Time

data is continuously
generated but we are
taking 1min, 2min
rather

Amazon

→ Confirm in 2min

live tracking

Batch

Credit card

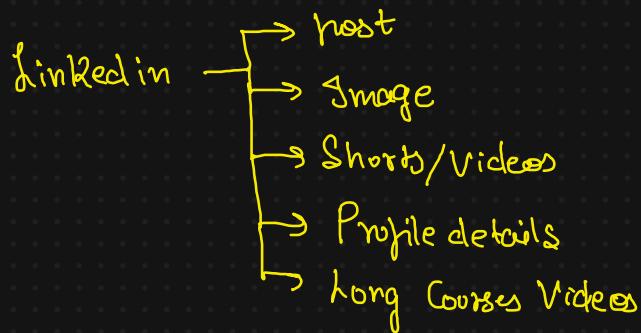
[data engineer]

Once a month



specific time period

3. Variety data can be in different formats & we have to deal with all of them.



Structured

Row & columns

Schema is enforced

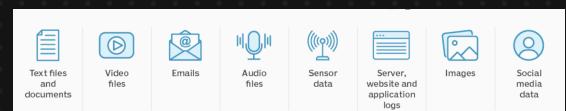
Semistructured

JSON

XML

CSV

Unstructured



Structured data

ID	Name	Age	Degree
1	John	18	B.Sc.
2	David	31	Ph.D.
3	Robert	51	Ph.D.
4	Rick	26	M.Sc.
5	Michael	19	B.Sc.

```
<University>
<Student ID="1">
  <Name>John</Name>
  <Age>18</Age>
  <Degree>B.Sc.</Degree>
</Student>
<Student ID="2">
  <Name>David</Name>
  <Age>31</Age>
  <Degree>Ph.D. </Degree>
</Student>
...
</University>
```

Parquet

```
{
  "customer": "John Doe",
  "gender": "Male",
  "items": [
    {
      "SSN": "123-45-6789",
      "exp": "10/01/2025",
      "moresubnesting": {
        "SSN": "123-45-6789",
        "newfield2": "123 Main Street"
      }
    }
  ]
}
```

4. Veracity: Trustworthiness or Quality of data

→ messy data needs to be corrected.

→ -ve age, Quality issues

5. Value: data should be useful & provide some value

→ help in business decision

→ provides insights

Volume \neq Big data

For a problem to be a big data problem, not all V's need to be satisfied. Generally a combination of 2~3 is enough

Big Data & Distributed Systems.

gta 3 / → gta 5
San Andreas

We need more resources

1. Storage
2. Memory (RAM)
3. Processors

Single System

Monolithic Systems
↓
1



Single, self contained unit

$X \rightarrow 2X \rightarrow 10X$

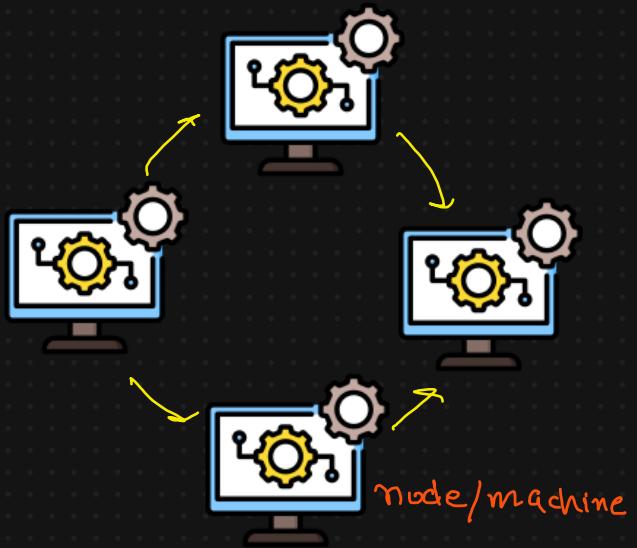
$P \rightarrow 2P \rightarrow 6P$

Scalability problem

cost & performance

Multiple Systems

Distributed architectures



True Scaling

$X \rightarrow 2X \rightarrow 10X$

$P \rightarrow 2P \approx 10P$

they don't have problem
of scalability

Vertical Scaling

Horizontal Scaling

All good big data systems are based on distributed architecture

Designing a Good Big Data System

Cloud vs On-Premise

Database vs Data Warehouse vs

Data Lake

ETL vs ELT

Data Engineer vs Big Data

Engineer

Hadoop Introduction