# Big Data Challenges

**Key Concept:** Resource requirements for Big data

As data complexity, size & variety increases more resources are required

## Resources

1) Storage : HDD, SPD, CD
2) Memory : RAM for processing tasks
3) Performance : CPU cores for computation

## System Architectures for Big Data
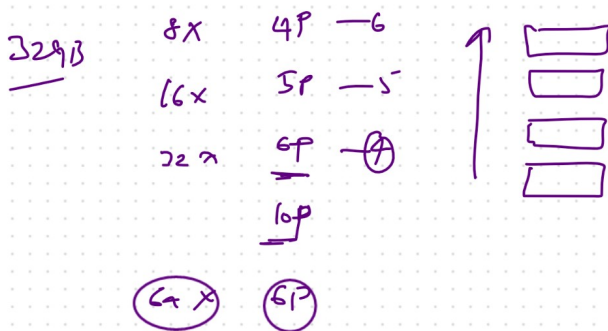
### Monolithic Systems

Distributed System

→ multiple System

* Single System

RAM    Storage
↑ 64GB , 2TB

4 × 32 , 8P
5 × 32 , 10P

32/1TB

32 1T

82/1TB

32 1T

#### Limitations

* Hardware constraint (max Ram or Storage capacity)

efficiently

### Vertical Scale

32GB

| | | |
|---|---|---|
| 8X | 4P | —6 |
| 16× | 5P | —5 |
| 32× | 6P | ④ |
| | 10P | |

(6× ) (6P)

* performance does not scale proportionally with resource upgrades

* Not scale for big data problems

* resources are distributed across nodes

* easily scalable by adding more machines

* achieves true scaling

horizontal scaling

Scaling & massive performance

Example

Buono : 4 Chef 8 chefs
resource ↓ ↓
databases

① **monolithic :** A single kitchen with many chefs, adding more chefs eventually reduces efficiency due to space & workflow constraints

② **Distributed System:** Multiple kitchens, each chef whose specialized in a cuisine. Adding more kitchens to increase efficiency & scalability

# Key Takeaways

## Monolithic Systems

→ limited scalability

→ Not suitable for big data problems

## Distributed Systems:

★ Enabled scalability & efficient resource utilization

★ foundation for all good big data systems

## Big Data System Deployment : On-premise vs Cloud Solution

# On-premise Infrastructure:

## Characteristics

★ Similar to buying an office / house

★ high initial Capital expenditure (CapEx)

★ Organizations must procure all hardware

(eg 20-node cluster ──→ purchasing 20 machines

→ Organizations is responsible for all the setup & maintenance

# Key Factors

**Deployment :** Hardware & s/w hosted within org facilities
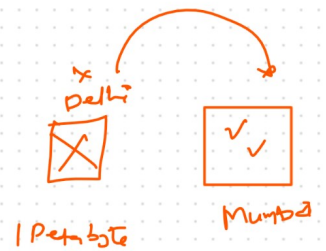
**Cost :** high upfront cost for

→ hardware

→ maintenance

→ IT staff

<u>Scalability :</u> Limitations by org hardware capacity

<u>Maintenance :</u> Org is responsible too managing

* hardware (failure)
* Software updates

Delhi

1 Petabyte

Mumbai

9301.

Flexibility : due to hardware resources

Security : Data remains within org : greater Control

Disaster Recovery : requires internal backup & disaster recovery systems

<u>Cloud Solutions</u>  characteristics :

* Similar to renting a house/co-working space
* pay-as-you-go model
* No ownership of hardware
* Access to machines via internet

Major Cloud Providers

* AWS
* Azure
* GCP
* IBm Cloud

<u>Deployement :</u> resources are hosted on provider's servers, we just need to access via internet

<u>Cost :</u>  pay-use model ⟶ Operational expences

<u>Scalability :</u> highly scalable, add resources on demand instantly

<u>Maintenance :</u> managed by cloud providers

<u>Flexibility :</u> high flexible allows resources to scale up & down as needed

<u>Security :</u> managed by provider, concerns on sensitive data

<u>Disaster Recovery :</u> builtin-recovery & redundancy across multiple locations

# Types of Cloud:

1) <u>Public cloud</u> :

     1, AWS, Azure, GCP, IBM

     2, ideal for startups due to no shifted infrastructure cost

2) <u>Private Cloud</u> :    VMware , OpenStack

      em: banks → ensure security &
                   regulatory compliance

    * different from on-premise by company

                             cloudlike UI & management

3) <u>Hybrid Cloud</u>:

    → combines public & private cloud

    → private cloud for certain usage / sensitive data

    → public data, for computations

    → Balances security & scalability need

4) <u>Community Cloud</u>:

     * Infrastructure shared by multiple organizations with common concern

     * Universities & hospitals share data & resources

     * limited to specific community member