

~~Sharding~~ - type of database partitioning that
separates large database into
smaller & faster to manage parts.

Big table -

Distributed, highly scalable & no SQL database system developed by google.

Designed to handle massive amount of structured data across large no. of servers.

Big table provides → sorted map, multidimensional distributed and sorted map

Data is organized based on (row key, column key & timestamp).

It is optimized for handling large scale workload

& widely used for web indexing, data analysis & time series data storage.

Key features →

① Scalability → petabytes. distributed over 1000+ machine allows for horizontal scaling.

② High Performance

③ Automatic Sharding → BT automatically shards data across servers for efficient storage & retrieval.

④ Replication

⑤ Integrated with GCP -

allows seamless processing & analysis across BigQuery, DataFlow, DataProc.

Overall, BT is flexible & scalable DB for large-scale

col-based data base - stores data by columns rather than rows.

Camlin	Page
Date	/ /

HBase-

open source, distributed, column-oriented DB built on top of Apache hadoop ecosystem

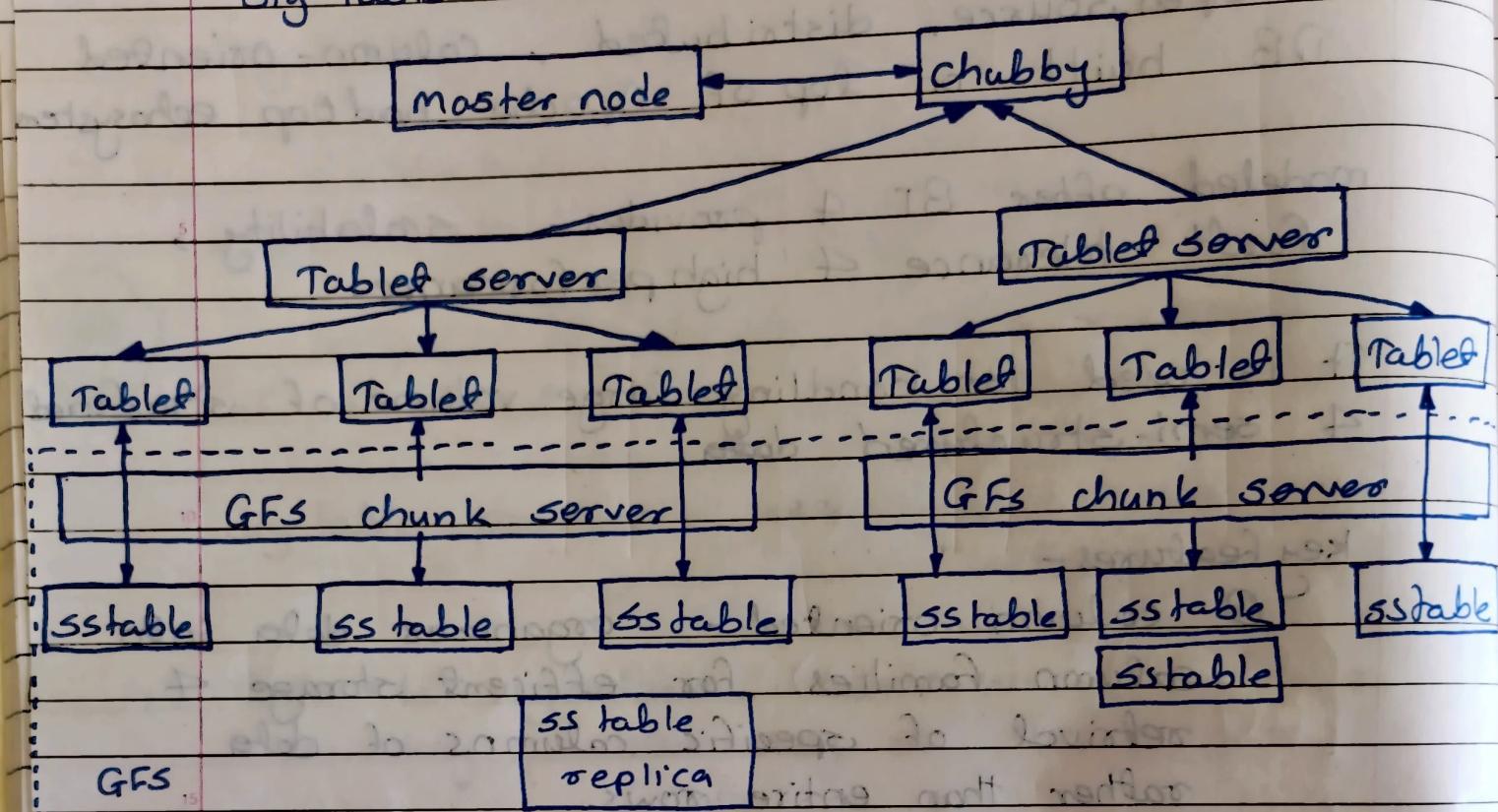
modeled after BT 4 provides - scalability, fault tolerance & high performance.

It is used for handling large volume of structured & semi-structured data.

key features -

- ① column oriented - organizes data into column families for efficient storage & retrieval of specific columns of data rather than entire rows.
- ② Horizontal scaling - (can add more servers)
- ③ High Availability & Fault-Tolerance - (Because of replication)
- ④ Strong Consistency - (for read/write ops)
- ⑤ Integration with Hadoop Ecosystem - HDFS, mapreduce, Hive, for eff processing & analysis.

Big Table arch-



- SS table - file format stores ordered immutable map from key to value. (binary search can be applied)
- Master node - assigns tablets to tablet servers, handles expiration detection of tablets, load balancing of tablet servers, garbage collection in GFS.
- Tablet server - manages set of tablets
 - handles read & write req to tablets.
 - splits the tablets that have grown too large.



HTF this arch works.
working -

- Chubby is used to keep track of tablet servers.

- ① When a tablet server starts, it creates a uniquely named file in specific chubby directory & acquires a lock on it.
- ② To detect if a tablet server is no longer serving its tablets master periodically asks each tablet for status of its locks.
- ③ If tablet server loses its lock, or if master was unable to reach server then it tries to acquire lock externally: exclusive lock on server.
- ④ If master is able to acquire lock ∴ Chubby is live & server is dead.
∴ assigns new server to those tablets & deletes the dead server.

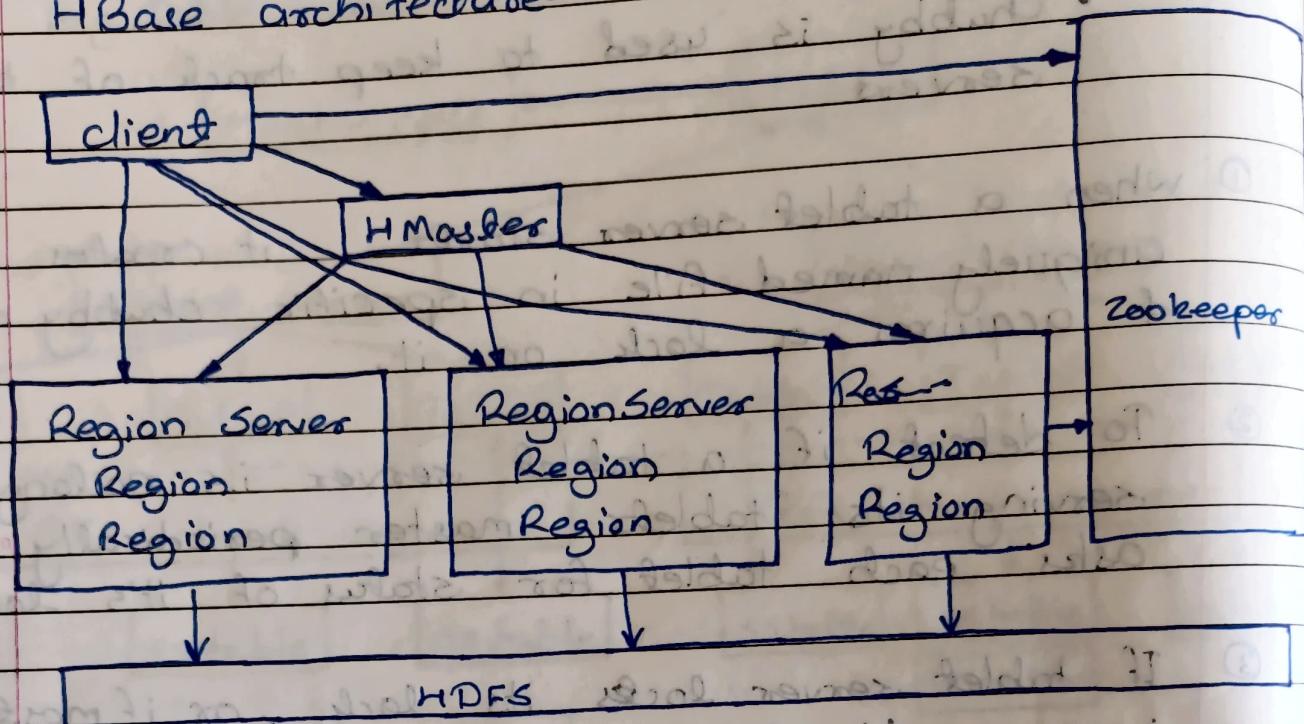
HDFS

- File sys for storing large files
- ② no fast individual record lookup
- ③ only sequential access to data
- ④ for high latency ops
- ⑤ data accessed through Mapreduce Jobs
- ⑥ No random read & write

HBase

- ① database built on HDFS
- ② it has ✓
- ③ Hash table for random access
- ④ for low latency ops.
- ⑤ provides access to single col from billion of recs
- ⑥ can perform

HBase architecture-



- **Zookeeper** - centralized monitoring server
maintains config & info
provides distributed synchronization

If client want to communicate with regions server,
client has to communicate approach Zookeeper.

- **HMaster** - for administrative operations.
of the cluster.

- **HRegions Servers** -
 1. Hosting 4 managing regions
 2. Splitting regions automatically.
 3. Handling read & write request.
 4. Communicating with client directly.

- **HRegions** - maintain a store for each col family
They are designed to accomodate semi-structured
data that could vary in field size,
data types & columns.

Dynamo - is key value pair structured storage system. It can act as database & distributed hash table.

It was created to help scalability issues that Amazon faced during holiday seasons. & then it was used in AWS.

DynamoDB is built on the principles of Dynamo, but Dynamo is based on leaderless replication, DynamoDB is based on single-leader replication.

Principles -

- Incremental scalability - should be able to scale out one node at a time
- Symmetry - every node should have same set of responsibilities
- Decentralization - design should favor peer-to-peer techniques.
- Heterogeneity - system should be able to exploit heterogeneity in the infrastructure.
work distribution should be proportional to capacity of the node.
∴ new nodes with higher capacity can be added

The index layer of ~~is~~ implements core features of Dynamo.

It also inspired many NoSQL databases implementations such as Apache Cassandra, Project Voldemort & Riak.

simpleDB has been departed by AWS and users are encouraged to use DynamoDB as replacement.
(in June 30, 2019) just after our syllabus updated.

Google Cloud Datastore

Amazon SimpleDB

provided by GCP

provided by AWS

Nosql document database service

Nosql Database Service
(now deprecated)

Highly scalable & designed for large amount of structured & semistructured data.

Designed for small scale apps with lower throughput requirements.

Offers strong consistency for read & write operations.

Provides eventual consistency.

Supports schema-less data model with entities and properties.

Offers only schema-less data models.

enable complex query with filtering, sorting & projection.

supports basic queries.

ACID-compliant transaction for data integrity

No built-in support for transactions.

Automatic distribution of data across multiple servers for scalability.

Limited scalability compared to Database.

Integration with GCP.

Integration with Aws.

geo-redundancy - distribution of file in servers across multiple data centres that reside in diff geographical regions.

Camlin Page

Date / /

cloud storage - data accessed over the internet

it is built on virtualization techniques.

advantages - reduce cost, simplify IT management, improve user experience, remote collaboration. Can be accessed via web service API, cloud storage gateway or web-based UI.

any device anywhere access, back ups, easy file sharing.

Cloud Storage Providers -

- Amazon S3 - object storage service, scalable & durable storage for various type of data. Offers features like encryption, versioning & lifecycle management.
- Google Cloud Storage - scalable, highly available object storage service by GCP. integration with Gcloud services.
- Microsoft Azure Blob Storage - optimized for storing large amount of unstructured data. integration with Azure services.
- Dropbox - offers personal & business plans provides file synchronization, sharing & collaboration across multiple devices.
- IBM Cloud Object Storage - for storing & managing large amount of unstructured data.

Features - geo-redundancy, integration with IBM MCS.

Securing the Cloud.

General security advantages of cloud based solutions.

Cloud based solutions offer several security advantages compared to traditional on-premises systems.

① Data Protection & Redundancy -

Cloud providers implement robust data protection, data encryption, access control, regular data backups.

They also have redundant storage sys & geographically dispersed data centres.

② Scalable Secure Infrastructure -

Teams of security experts to monitor, regularly update & patch their sys.

③ Access Control & Identity management -

multi-factor authentication,

role based access control,

OAuth,

reduce risk of unauthorized access.

④ Physical security measures -

surveillance system,

environmental safe guards.

⑤ Compliance & auditing -

Cloud providers adhere to industry standards ISO 27001, SOC 2, HIPAA, & GDPR.

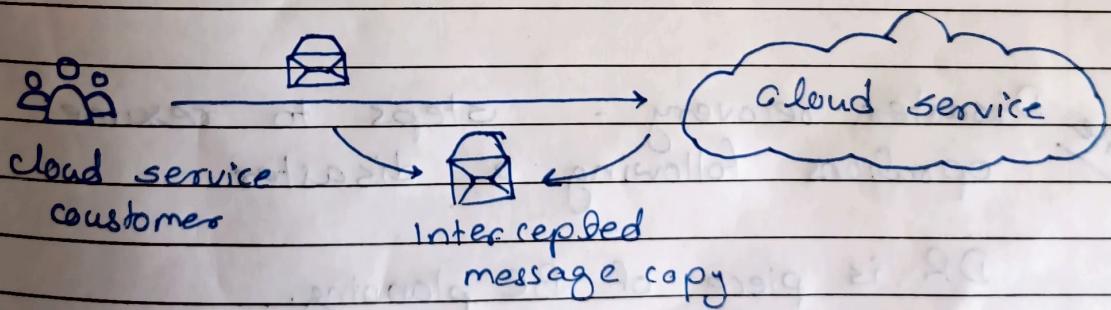
- ⑥ Automatic update & patching-
 reduces the burden of organizations.
 & helps to ensure systems are protected
 against known vulnerabilities.

⑦ Security Monitoring & incident Response-

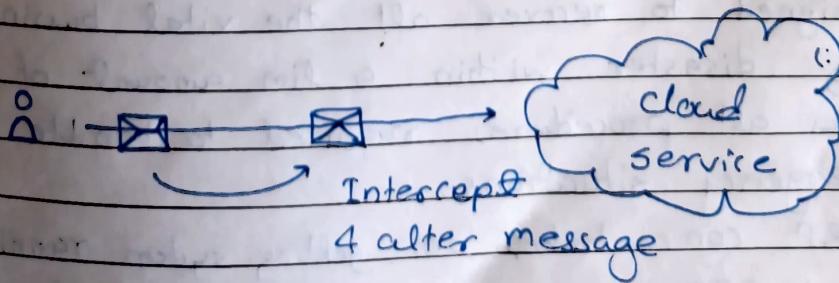
cloud providers employ advanced threat detection systems, perform security incident response and provide visibility into security events, allowing organizations to respond to potential threats more effectively.

• cloud security Threats-

① Traffic eavesdropping

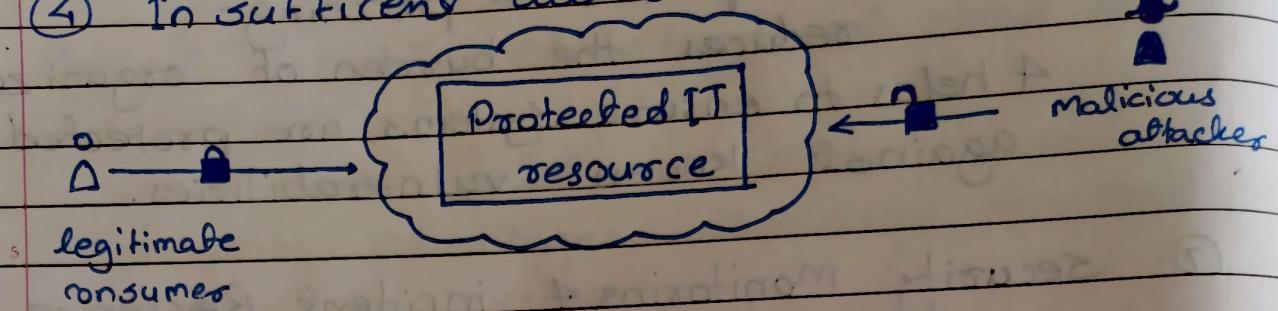


② Malicious intermediary



③ Denial of service - overloading resources so that resources can hardly be allocated to others.

④ Insufficient authorization.



⑤ Virtualization attack -

- Business Continuity & disaster Recovery -

~~Proactive~~ Business Continuity refers to process to ensure that critical function can work during & after a disaster.

Involves planning towards long term challenges.

~~reactive~~ Disaster recovery - steps to resume operations following a disaster.

DR is piece of BC planning.

- Disaster recovery plan -

- designed to recover all the vital business process during disaster within a lim amount of time.
- It has all procedures required to handle the emergency situations.
- DRP concentrates on getting system running asap
- RTO & RPD recovery point time objective and recovery. Point objective are targets of DRP.

- most successful disaster recovery plan is the one which will never be implemented (risk avoidance)

Business continuity plan -

- activities required to keep org running during period of displacement or interruption.
- BCP helps in continuing business even after disaster.
- Business has to stay active otherwise orgs will experience losses.
- legal issues can arise if critical services are not provided to clients.
- hence efficient BCP is required to actively run and maintain business activities.

• Understanding the Threats -

- First step for preparing business for disaster.

Threats : -

- ① Natural & localized disasters, earthquakes, floods, tornadoes & fires.
- ② Failure of IT system - network, file servers, software appliance.
- ③ Power outages, such as utility failures.

A BCP takes such things in account. It also takes in account tangible & non-tangible cost your business may need.