

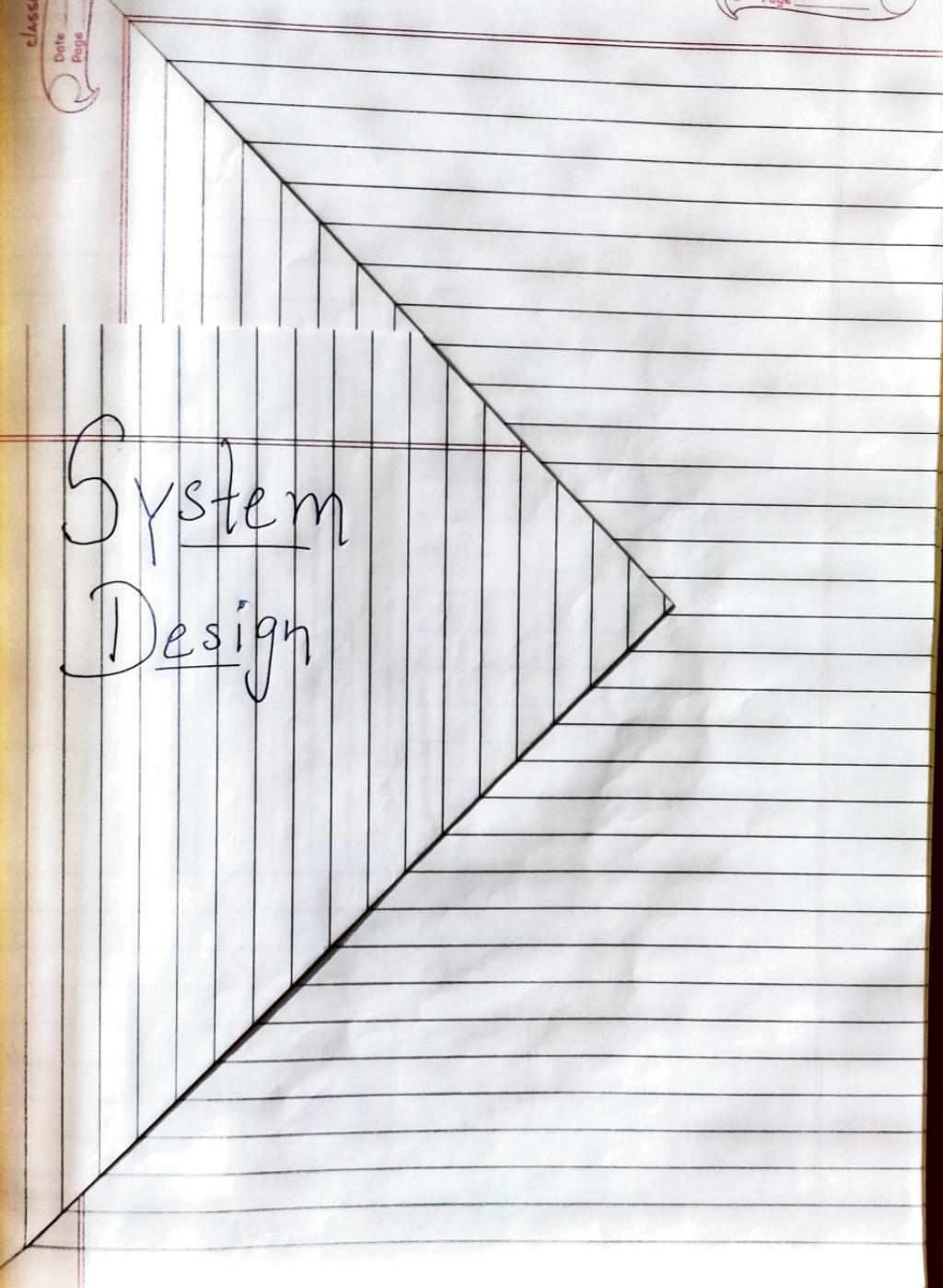
CLASSMATE
Date _____
Page _____

CLASSMATE

Date _____

Page _____

System Design



REST API:

→ Representational State Transfer (REST) REST != HTTP

→ guidelines for between client & server to exchange representational states of data

* Guidelines:

i) Client Server

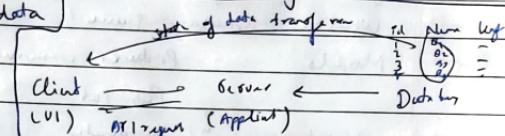
ii) Cacheable

iii) Layered (component knows about the immediate component)

iv) Stateless (server does not know who is requesting)

v) Uniform Interface

vi) Code on Demand



eg REST API implemented over HTTP:

HTTP method: Uniform Resource Identifier (URI) (Response body)

i) GET `https://api.books.com/mybooks/books` {
 ↗ memo: --
 ↗ pmt: --
 ↗ name: --
 ↗ lmt: -- }
 ↗ Representations of data (JSON)
 ↗ server domain
 ↗ (JSON)

ii) POST {
 ↗ id: 85
 ↗ memo: --
 ↗ auth: --
 ↗ pmt: -- }
 ↗ (status: 200
 ↗ id: 85)
 ↗ Payload Response

iii) PUT {
 ↗ id: 85
 ↗ memo: --
 ↗ auth: --
 ↗ pmt: -- }
 ↗ (status: 200
 ↗ id: 85)
 ↗ (payload)
 ↗ Response

iv) DELETE {
 ↗ id: 85
 ↗ memo: --
 ↗ auth: --
 ↗ pmt: -- }

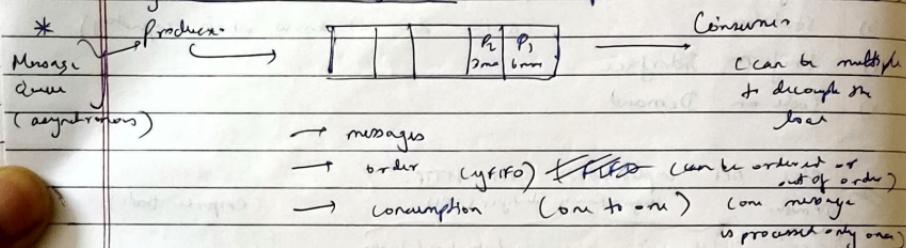
* Path Parameters : `<server-domain>/mybooks/books/1` {
 ↗ id = 98;
 ↗ auth = 1234567890;
 ↗ pmt = 100 }
 ↗ (status: 200
 ↗ id: 98)
 ↗ Every Parameter : `<server-domain>/mybooks/books/1?lmt=20&offst=0`
 ↗ (2 books displayed at a time)

* Security, Authorization & Error Handling :

Message Queues

- Sync vs Async Communication
- What are message queues?
- Model : Producer consumer
- Pub-Sub
- Advantage & Disadvantage
- Use Cases

Asynchronous : Message Queue : decouple the components (Scalability)

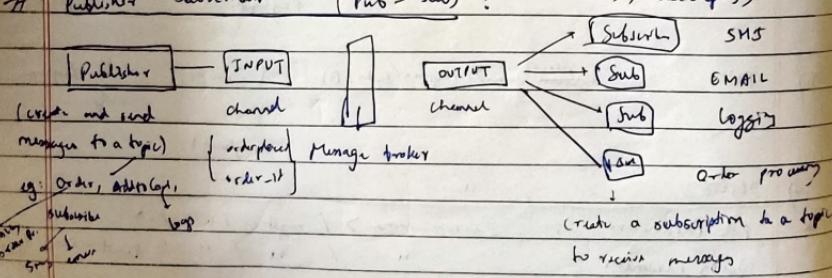


e.g. invoice generation → order doesn't matter

* Asynchronous Communication : Info. is shared independent of time
(Response not necessary)

Synchronous → instant response (real time) ~~change~~

Publisher - Subscriber (Pub - Sub) : (Scalability, decoupling)



Use Cases of pubsub:

- i) Asynchronous Work flow
- ii) Decoupling
- iii) Load balancing (Increase the no. of subscribers)
- iv) Deferred Processing (processing message at a later point of time (off peak time))
- v) Data streaming

Not to be used when: small no. of requests per day
requirement of real time processing

Performance metrics:

- i) Throughput : Work done / Time taken (No. of API calls served per unit time)
- ii) Bandwidth : High resource AND High Bandwidth → High throughput
- iii) Response time : Time taken to serve the request
- iv) Latency

Performance Metrics of Components :

- i) Application: API response time, throughput of API, error occurrence, busy/ugly/delay in code
- ii) Database: time to execute queries, No. of queries exec/time (throughput)
- iii) Cache: latency of writing to cache, No. of cache eviction & invalidation, memory of cache instance
- iv) Message Queue: production/consuming rate, fraction of stale/unprocessed messages
- v) Workers: time to complete job, resources used in processing
- vi) Server instances: Memory / RAM, CPU

Performance metric tools: dashboard of metrics (GCP, AWS, Azure)

fault and failure :

Cause Effect

- The oven ~~for~~ crashed so cake couldn't be baked
fault failure
 - Faults :
 - Transient fault Permanent fault
 - In distributed systems
 - (short time, hard to detect) (system won't work until fix)
 - Solution : Replication (easy to detect)
 - Faults eg : software bug (write tests to detect), network error, memory leakage, high CPU utilization, etc.
 - Even in case of faults, system must fail gracefully.

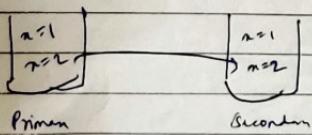
Scaling : eg: 1K request / minute → 10K requests per minute

- Vertical scaling: Increasing capacity of resource (eg: adding more drives to disk)
 - Horizontal scaling: " no. of resources (eg: adding multiple tables)

Properties	Vertical	App. ↑
→ handle increased load	Verbal	1K request / min → look up /
→ not complex	Horizonte	App. ↑
→ performance should not take a hit		1K request / min.

Database Replication :

- handle faults / failures
 - reduces latency (multiple geographical locations)
 - application performance ↑ (by multiple nodes)



At $t = t_1$, $x = 2$ in Primer

Replication lag: e

At $t = t_2$, read is performed

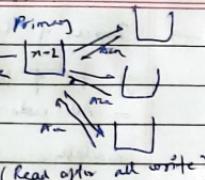
if $e < (t_2 - t_1)$ ↘ else increasing

classmate

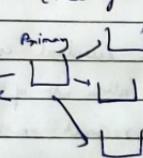
Date _____

Page _____

* Synchronous replication :

- All replicas updated before host is acknowledged
- consistent results but takes more time for writes.
- Problem: one of the replica fails \Rightarrow no acknowledgement


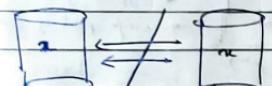
* Asynchronous replication :

- Host is acknowledged after primary DB is updated.
- better performance but may have inconsistency.


(Semi-synchronous \rightarrow wait for one acknowledgement for my backups)
 (or some K)

- * Snapshot \rightarrow store the state of data and roll back in case of failure

CAP (Consistency, Availability, Partitioning) :



A communicating link B
 (if breaks = partition)

Consistency \rightarrow both databases agree on same value
 When the communication link fails, and the system still runs:

Either we can achieve consistency (stop receiving requests on B)
 or availability (receive request on both A and B)
 (not consistent)

- If we tolerate a partition (system can still run and ensure consistency / availability) \rightarrow then it is called partition tolerant system

CAP Theorem :

(Brewer's theorem)

- Having all 3 properties is impossible
- one has to be sacrificed.
- Partition is inevitable hence we must ensure partition tolerance
- Consistency / Availability needs to be chosen according to requirement

Consistency: Any read after a latest write, all nodes should return the latest value

Availability: Every available node should respond in a non-error format to any read request without guarantee of returning latest write

partition Tolerance: System will respond to all read & write even if communication channel between them broken

→ Degrees of Consistency & Availability:

e.g. one of the nodes stop serving reads, but keeps serving writes.
 ⇒ All nodes will give consistent data and also the system is more available.

e.g.: eventual consistency:

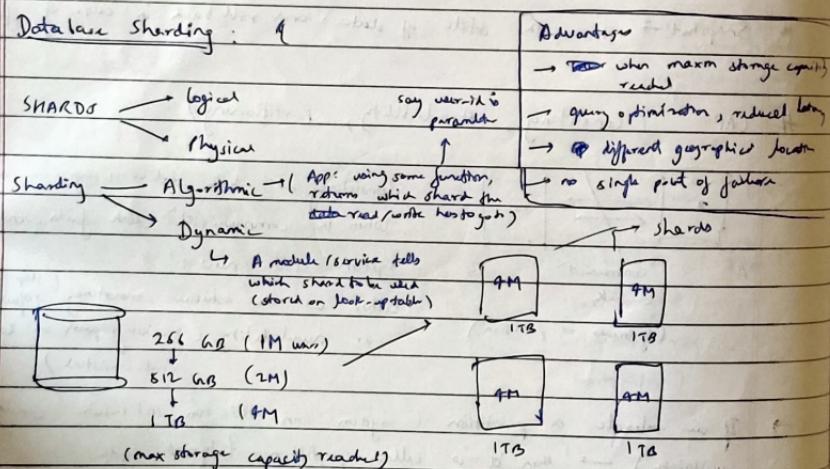
allows both of them to read, but after sometime when link comes back, records on the latest value.

→ Use cases: Consistency → ATM / flight booking / banking

Availability → youtube likes / analytics data

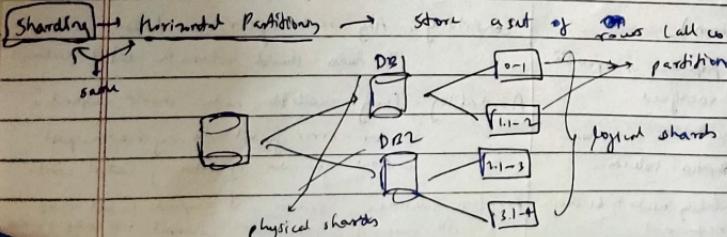
→ Backing network connections to avoid partition (expensive)

Database Sharding:



Vertical Partitioning → store a set of columns on one DB (partition on row of columns)

Horizontal Partitioning → store a set of rows (all columns)



* Key Based Sharding :

- Based on some shard key (may or may not be the primary key)
Hash (shard key) → which shard to use
(Algorithmic sharding)
- Advantages → Equal distribution of data
Disadvantage → when data increases / decreases, no. of shards need to change → hash function has to be changed
migration of data b/w shards is difficult
- Shard key should be static, (do not change frequently)
(multiple columns can also be chosen)

* Range Based Sharding :

- Sharding based on some range
eg: Jan-Feb, Mar-Apr (6 shards)
- Algorithmic sharding (in the application code.)
- Advantage : no hash function so no much migration of data when data grows.
- Disadvantage : more data in one shard (eg Jan-Feb)
Hotspot → uneven distribution of data leading to more data in one shard.
- Use cases → range based queries ($<=$, $>=$, etc)

Directory Based Sharding : (dynamic)

Advantage :

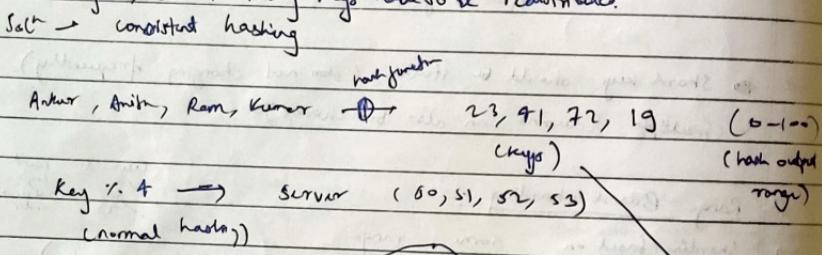
User Id	Zone	Look Up Table		→ easy to add/remove more shards
		Zone	Shard	
101	A	B	1	
102	B	C	2	Disadvantage :
103	C	D	3	→ extra lookup
104	D		4	→ fails if look-up table crashes

App → DB/Shard (replica)

(extra look up required) for every K/V

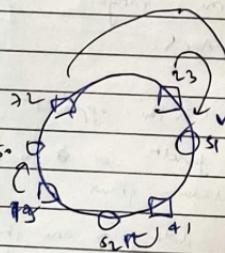
A Hashing:

- String / Objects / Nat \rightarrow Integers (fixed by 32 bits)
- Hash function eg: $f(x) = x \% n$
- Marked value can be used as array index for quick access
- eg use case \rightarrow mapping users to servers ($u_i \rightarrow s_j$)
- Disadvantage of Normal hashing: (eg $n \% n$)
If no. of servers increase / decrease, the hash function has to be changed, and too many keys have to be redistributed.



* Consistent hashing:

Every value to sit on the next server in clockwise / anticlockwise (choose any one) direction



map output value if
hashfun \rightarrow points on ring

Also map the values of server $(0-100)$ to points on ring

~~Problem: removing a server, adds too much load on the immediate next server~~
~~Sol: create replicas of servers and map each of them to form circle~~

Functional & Non-functional requirements

Functional requirements: functionalities provided by system to fulfill user requirements

- ex: Twitter : post tweet, delete tweet, favourite a tweet, follow/unfollow
- Outcome: system / components of system
- distribute the responsibilities of functional requirements to API's, workers, events, messaging (System design diagram, component & architecture diagram)

Non-functional:

- latency of posting tweet, availability of system
- Traffic, Users, Availability, throughput, latency
- choice of resources → No. of resources → Capacity Estimating
- Tech choices , resource utilization, data storages, services or hardware