

**SHARDING****TYPES****Range Based / Dynamic**

- Each shard holds a range of values
- Ex: date ranges, time ranges

→ when ordering + ranges are there in data

**Hash-Based / Algorithmic / Key based**

→ use a hash function on data to get which shard

→ used when there is no natural order (ids)

**Directory Based**

→ you manually maintain a data → shard mapping

**Geo-Based**

→ shard based on geo locations or nearest to locations

**IMPROVEMENTS****SHARD + INDEXING**

Ex: shard on city + index on age

→ here first shard on city + internally each shard is indexed based on age

Query - give me all persons in DELHI with age between 20-40

dividing data into smaller fragments + storing them on different distributed machines

**Partitioning** → same logic but logically within same system

**PROS**

- faster and efficient queries (given you pick good sharding strategy)
- eliminate single pt of failure

**CONS**

→ outage possible as it is distributed

(soln) **master-slave** design

handles writes handles reads

→ fixed number of shards → need rebalancing / rehashing

(soln) **consistent hashing**

→ skewed shards → better sharding for approach

→ joins across shards are costly

**CONSISTENT HASHING**

\* Normal hashing - lets say you have  $N$  servers

$(h(\text{data}) \% N)$  → you can't change this  
hashfn% no. of servers if done needs rebalancing everything

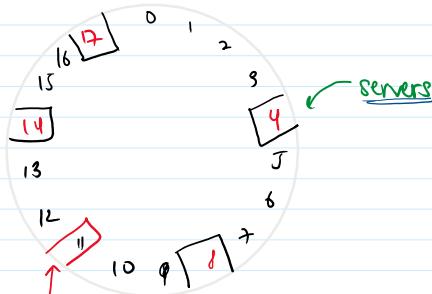
+ consistent hashing

→ consider your slots are in circular form

→  $x \rightarrow \text{hash fn}(x) \rightarrow (0-360)$  example

→ Each result of hash fn will lie on this circle

→ some of these slots will be servers also

**HOW IT WORKS**

→ data  $\rightarrow \text{hash fn}(\text{data}) \rightarrow \text{slot-no} = 3$

→ move clockwise and see next server to 3

→ 3 gets stored in 4 i.e. 4

**New server gets added**

→ 11 is now made a server

→ New indexing doesn't change (9,10 will go to 11 instead of 14)

**Rebalancing is minimal**

→ only 14 which was previously handled by 11 now needs to be rendered to 11

\* Load factor =  $1/N$  (ideal case) equal load to each server

→ But this isn't possible when you have sparse

**(soln) Virtual servers**

→ Add more server slots which brings load factor to  $1/N$

→ Internally a single server handles  $(2^f)$  slots

Ex: 11 might be handled by a single physical server

**ACID**

Maintained in **Relational DBs**

- Good at vert scale (Bad at vertical)
- Fixed schema
- Generally not distributed

↳ Maintained in **Relational DBs**

- Good at vert scale (Bad at vertical)
- Fixed schema
- Generally not distributed

### Atomicity

→ A transaction can either happen completely or not happen at all

Ex - Bank amt transfer

① deduct from ur acnt → ② Add to the des' acnt

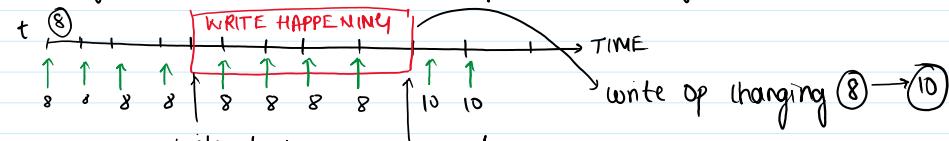
Now this shud happen completely. Ex lets say it got deducted from acnt but then failed

→ In such case revert back everything approach  
everything or nothing approach

**Consistency** → multiple reads at same time should give same results  
why easy? since its not distributed no issue of out-of-sync replicas

### Isolation

→ Every transaction happens independently without caring abt other parallel transactions



→ Here write op was changing val from 8 → 10 but during the process any read doesn't know or care abt this and still reads at 8

→ Only after write operation is completed and state is updated, Read returns 10

**Durability** Every transaction is logged and data is being persisted

### BASE

↳ Shown in **NOSQL Databases**

→ No fixed schema

→ Vertically scalable

→ Distributed (data stored across machines)

### Basically Available

→ Since data is distributed across servers it can handle outages

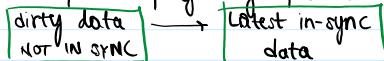
→ If one server crashes another replica/machine can give u the data

### Salt States

→ States/values can auto change without txns or user queries

How? SYNC across replicas

↳ since data is stored across machines/replicas it may happen that data copies keep getting updated



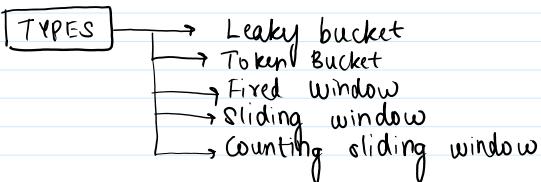
### Eventually Consistent

→ Not immediately consistent → All nodes are immediate consistent

→ 2 simultaneous reads may give different results sometimes

why? SYNC might have been happening b/w 2 reads

∴ TXNS are not always consistent due to distributed nature  
BUT will become consistent after some time (once sync/updates are completed)



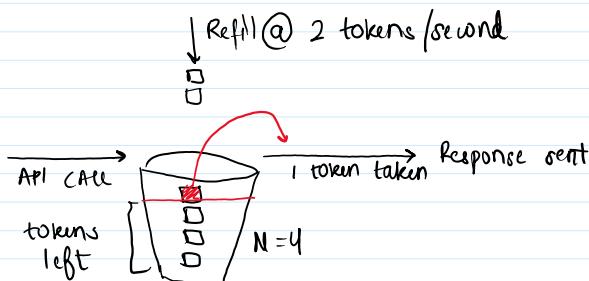
### TOKEN BUCKET

→ You have  $(N)$  tokens in a Bucket (Consider this as a currency)

→ whenever a request comes  
① Check if bucket has token

    → (No) → send 429 response code

    → (Yes) → Take one token from the bucket



★ Refiller : Refill the tokens in the bucket at certain rate

$$\text{refill\_rate} = \text{TPS}/\text{speed of your api}$$

\* Configurations → ①  $N$  = Bucket size  
② Refill rate

### LEAKY BUCKET

→ You have a bucket of size  $(N)$

→ A request comes in, if bucket

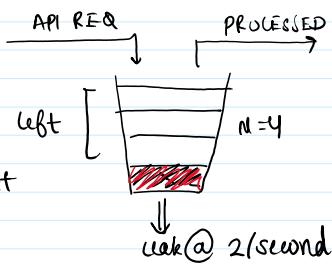
    → (Full) → 429 code

    → (else) → Add 1 token to the bucket

→ Leak rate → You keep on removing tokens at this given rate

Similar to refill rate in token bucket

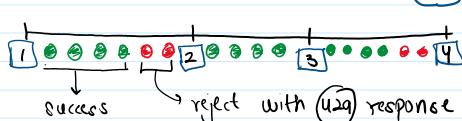
\* IDEAL CASE → Rate of API hits  $\leq$  LEAK RATE



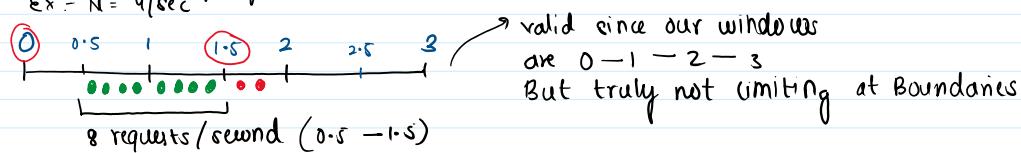
### FIXED WINDOW

→ You maintain certain windows (INTERVALS) and each window can only contain  $(N)$  items

→ Ex - Window of 1 min each



CON → Doesn't always guarantee rate limit  
Ex -  $N = 4/\text{sec}$



### SLIDING LOG WINDOW

→ Instead of fixed window you use timestamps to dynamically calculate it

$$\text{rate} = 4 \text{ req/sec}$$

Ex → Request comes at 1.5, you count no of req b/w 0.5  $\leftrightarrow$  1.5 if num  $< 4$

(Ex)

Request comes at 1.5, you count no of req. b/w 0.5  $\leftrightarrow$  1.5  
if num < 4

Insert this log (timestamp + reqid)

\* Ensures proper dynamic rate limiting

CONS

→ you store each request as one entry

∴ Space needed = N

lets say your Rate is 1000 rps  $\Rightarrow$  you will store 1000 log entries

SOLUTION

### SLIDING COUNTER WINDOW

→ Instead of pushing each request as log item maintain a counter

→ Rate = 1000 rps

ts	counter
0.1	1
0.2	100
0.3	200
:	
0.9	400

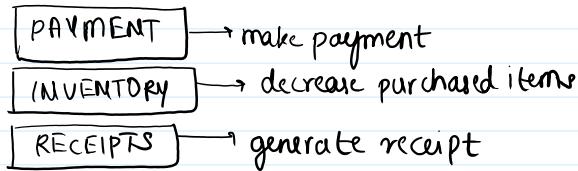
→ Here instead of 1000 logs you  
will only have 10 logs

→ space = 1000  $\rightarrow$  10

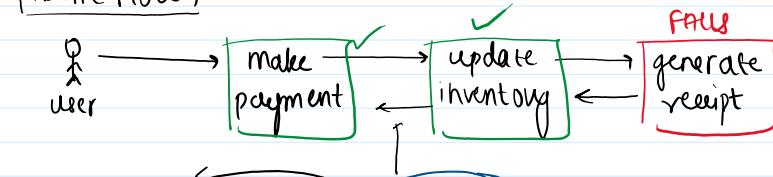
→ No of ts steps

## DISTRIBUTED TRANSACTIONS

**PREMISE** → You have 3 services



### IDEAL FLOW



How will you handle this in microservice

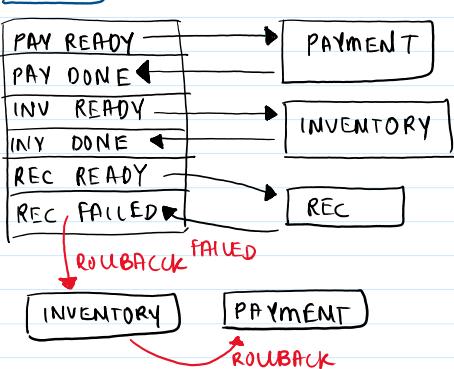
- ① 2PC [2/3 phase commits] BLOCKING
- ② 3PC
- ③ Saga (ASYNC)

### SAGA

\* Asynchronous and non-blocking but complex to implement

→ Queues are used for Inter Process Communication

#### STEPS



#### Implementation

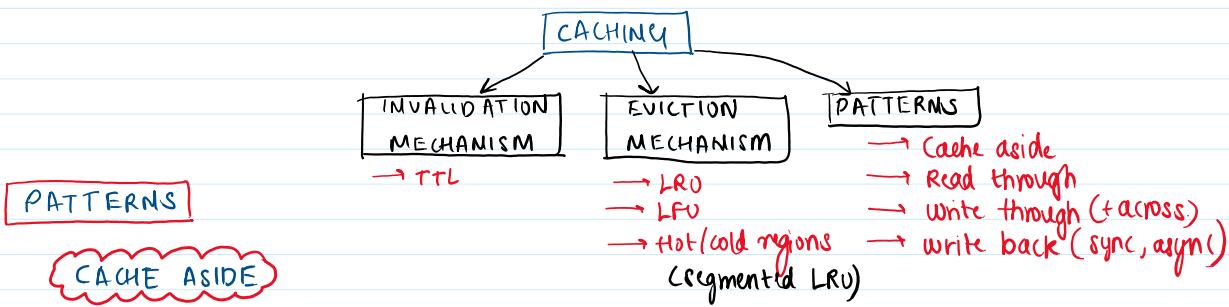
→ 3 microservices  
→ (3 \* 2) Queues      [ 1 success ] one bus each mservice  
                        [ 1 failure ]

→ 2 Queues

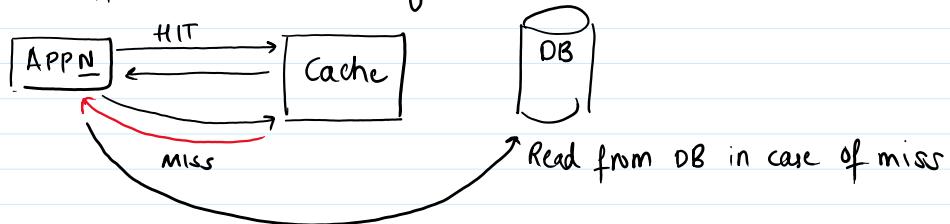
- ① Handles all success (all mservice listen to this)
- ② Handles all failures OR listen to a specified partition

## Caching

07 June 2024 23:04



→ Cache can't communicate directly with DB  
Application does that for it



- Initially data is loaded into the cache by app?  
On every cache miss, the app directly fetches data from DB  
→ On new data, it will always be MISS  
    until you update cache

### PROS

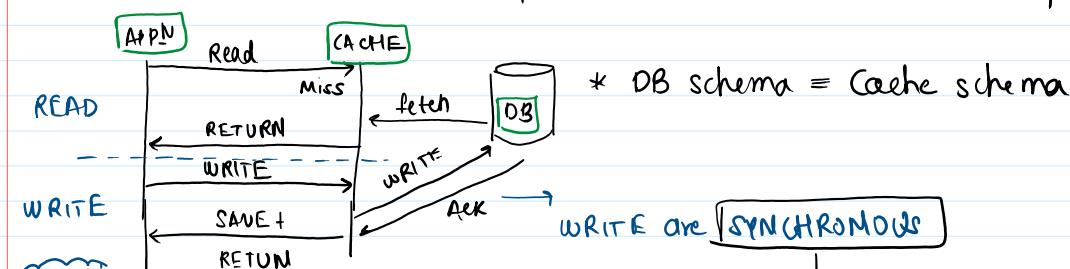
- \* Good for HEAVY-READ apps like dashboards
- \* If cache fails → you can directly hit the down DB each time
- \* Since data stored in cache and DB is different, you can have different schemas

### CONS

- \* New data will be MISS
- \* Inconsistency in data if it gets updated in DB

### READ THROUGH

- Cache can communicate directly with DB App! only accesses Cache  
→ READ MISS: Cache itself updates value from DB  
    → New val, update and then returns
- Writing new data: - you will write to cache + cache will write to DB  
    → if this was update invalidate cache which will be updated on next miss



### PROS

- Great for heavy reads why? (on write cache writes to DB, waits for ACK saves it in cache and then sends it back)
- You don't worry abt miss and fetch from DB

### CONS

- Cache can't fail
- Cache miss for new data (PRE-HEATING can solve)
- Cache schema ≠ DB schema

- Cache miss for new data (PRE-HEATING can solve)
- Cache schema = DB schema
- Data inconsistency

**WRITE THROUGH** → used with READ THROUGH (100% CONSISTENT)

- Here for writes you directly write to DB and invalidate cache (if updated)
- Next hit will be cache miss and cache will fetch from DB  
(saves latency of app → cache → DB for write operations)

**WRITE AROUND** → You ONLY WRITE TO DB and forget

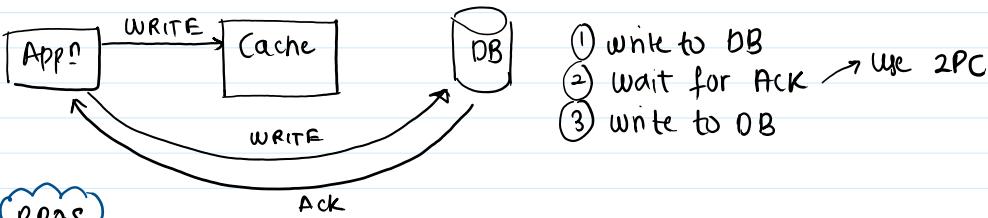
\* when you do a get not time; CACHE miss and cache will fetch from DB

**WRITE BACK** → SYNCHRONOUS — same as write through (BUT cache updates DB)

**SYNCHRONOUS** → ASYNCHRONOUS — NOT 100% consistent

→ Not app itself

- App! itself writes to both DB and cache whenever a write happens



### PROS

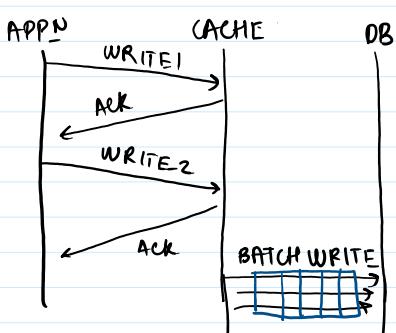
- Very high hits and low misses
- Very rare chance of inconsistency each time. No lag

### CONS

- High latency due to 2 writes and ACK wait
- Need 2PC to make sure DB update is ACID

## ASYNCHRONOUS

- App! only writes to cache and then forgets
- Cache then asynchronously writes data into the database independently
- done via QUEUES
- you can publish 1 event for each write OR Batch (① write and at once write to DB)



### PROS

- WRITE HEAVY applications
- Lesser write latency and lower inconsistency
- Better performance

### CONS

- Sync issues if you immediately read again before DB write is done

## Start with READ THROUGH CACHE

why? → when write happens u just write to cache and forget

→ Now you can directly read that data from cache (even though its still not in DB)

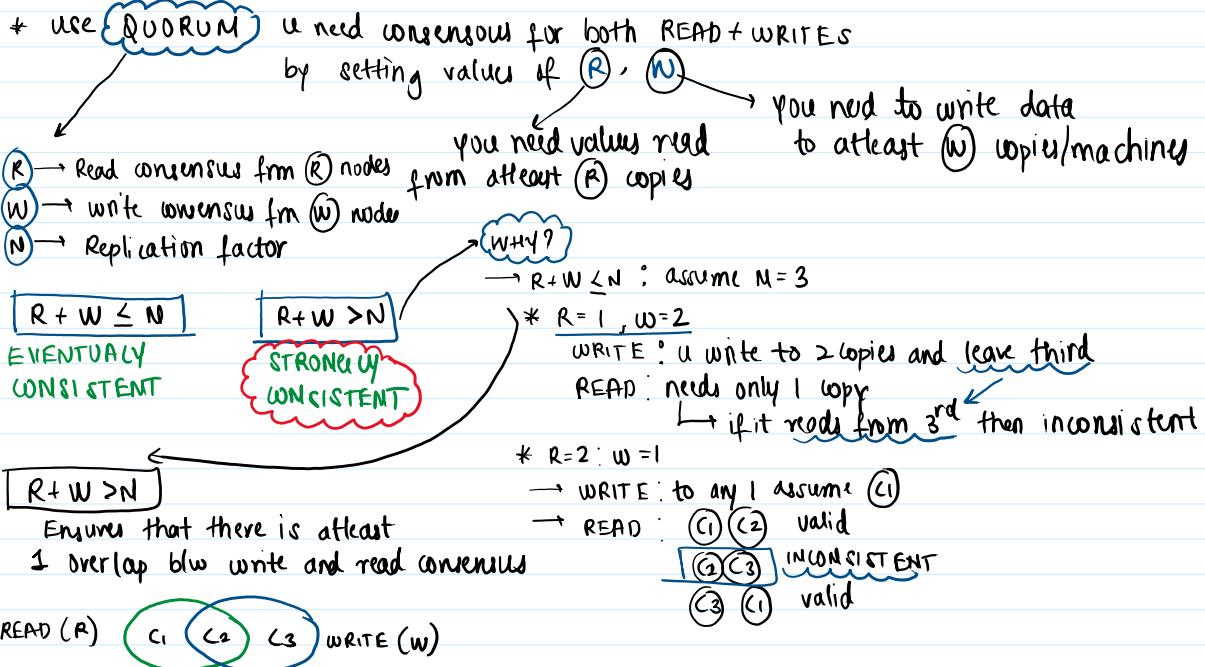
## GLOBALLY DISTRIBUTED CACHE

U: → now you can directly read that data from cache (even though it's still not in DB)

## GLOBALLY DISTRIBUTED CACHE

- ① There can be  $(N)$  different cache servers distributed across clusters
  - ↳ A cache gateway redirects requests to appropriate servers
  - ② Default uses **consistent hashing** for allocation
  - ③ Geographically based - route/store in nearest store
  - ④ Identifier based - Ex: Based on country-id group and store in same server or multiple countries into 1 continent server

## REDUNDANCY



## Isolation Levels

08 June 2024 14:54



### LINEARIZABLE

- single thread sequential execution of Queries
- No chance of any mismatches
- **PERFECT** consistency

### EVENTUAL

- DB might not be always consistent BUT eventually consistent
- multiple queries/read/writes can happen of multiple threads/servers
  - ↳ parallel process so is faster
- \* But sync takes some time after transactions

### CASUAL

- Allows parallel transactions BUT operations in single transaction are linear/serialize
- Consistent within transaction
- Better than causal but not as good as linear

#### CONS

Aggregates/joins are inconsistent

### QUORUM

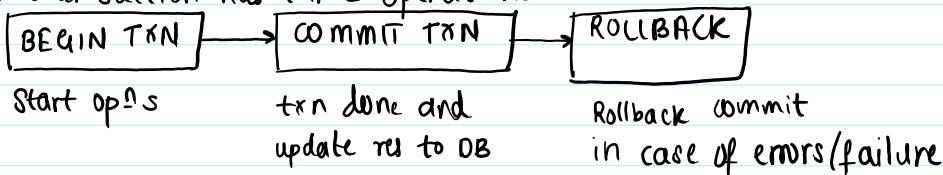
Configurable consistency vs performance

$R + W > N$  Highly consistent

$R + W \leq N$  Eventually consistent

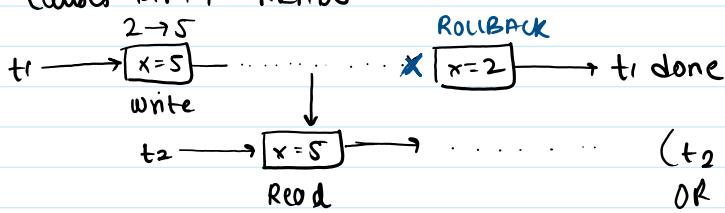
## TRANSACTION ISOLATION LEVELS

A transaction has three operations



### READ UNCOMMITTED

- Each trn directly writes to DB ≡ No concept of commits (No data replication)
- Causes **DIRTY READS**



( $t_2$  never gets rolled back value of  $x$ )  
OR some other op! might also have updated it

### READ COMMITTED

**READ UNCOMMITTED**

- Trn can only read values which are committed by other trns
- slower since you have commit step → it wont roll back later and dont directly update DB

\* In both read committed/uncommitted you cant ensure REPEATABLE READS  
 → if you have 2 reads of same var in the trn  
 you can get 2 different values (someone else must have updated it)

**REPEATABLE READS** → SNAPSHOT ISOLATION

- when any row is read by a transaction its LOCKED until that transaction ends
- other trns which update this locked value can locally commit and keep it as new version and update in DB once lock is released

**SERIALIZABLE**

- All operations in trn are executed sequentially in order
- if 2 trns dont have any common ops then they can run in parallel
- use serialized locks ↓  
 both writers and readers

least isolation  
most efficient

Least efficiency  
most isolation

