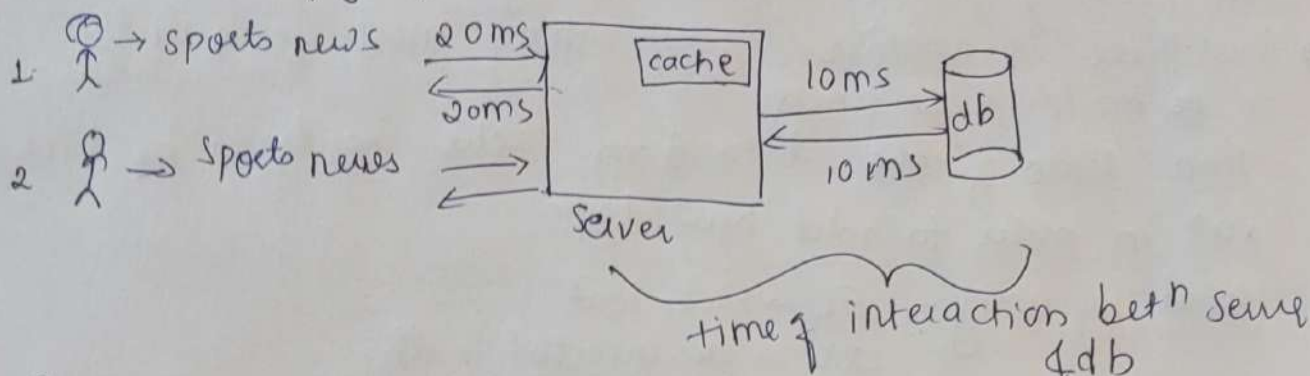


LECTURE 8 - CACHING

process of storing frequently accessed data in temporary storage location so it can be accessed faster.
cache is a copy of database.

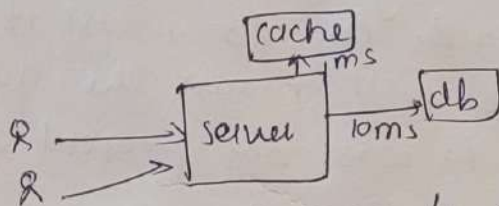


Person 1 requests for sports news $\rightarrow 20\text{ms}$

req goes to db, queries sql & gives results $\rightarrow 20\text{ms}$ (going & coming)

if person 2 also asks for same data, then instead of querying same sql again, we can store the 1st result in our cache memory & when some new user (person 2) asks for same request, we can directly send data that is stored in cache memory.

\rightarrow this will reduce the time of interaction between server and db.



it is much fast and easy to access the data from cache rather than accessing / querying the req & access the data from db.

Ques Then why cannot we put our entire db in cache
Ans ~~It is~~ if db is small then its possible. But if db is large it is not possible cause our cache size is already small and limited.
So, our cache would only store the frequently requested query results.

CACHE POLICY

least recent used \rightarrow LRU
least frequently used \rightarrow LFU

Suppose our cache is already filled with most frequently requested query.

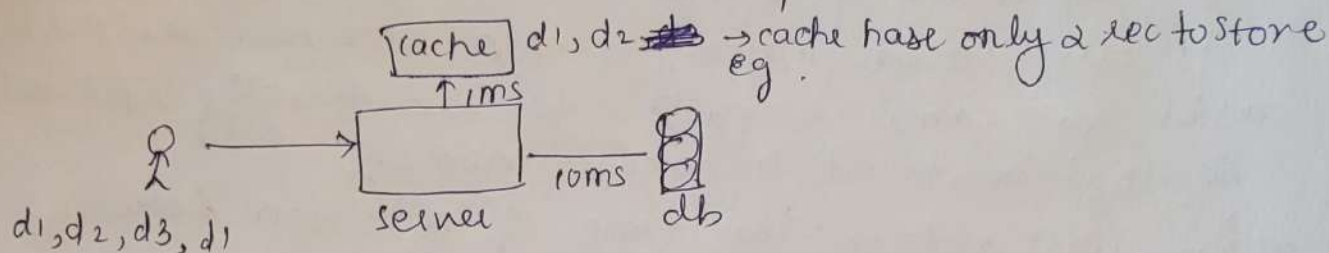
\rightarrow Now there is another request coming most frequently i.e. movie release news

Then some of the data from cache needs to be kicked out in order to add new data.

~~used~~ by using - LFU \rightarrow kicked out
LRU \rightarrow inserted in db.

DRAWBACKS

① Trashing \rightarrow when system spends more time on non productive task rather than on productive task



eg ~~used~~ cache is already filled with d1 & d2 data now user requests for d3, first it will check on cache then it will process the query & give result (as d3 was not present in cache it took extra 1ms) i.e. this was waste of time thrashing

\rightarrow insertion & removing & then inserting the prev data is time consuming

② Eventual consistency

eg cache is updated every one hour, our db & cache data are same at this pt eg (o likes in yt)

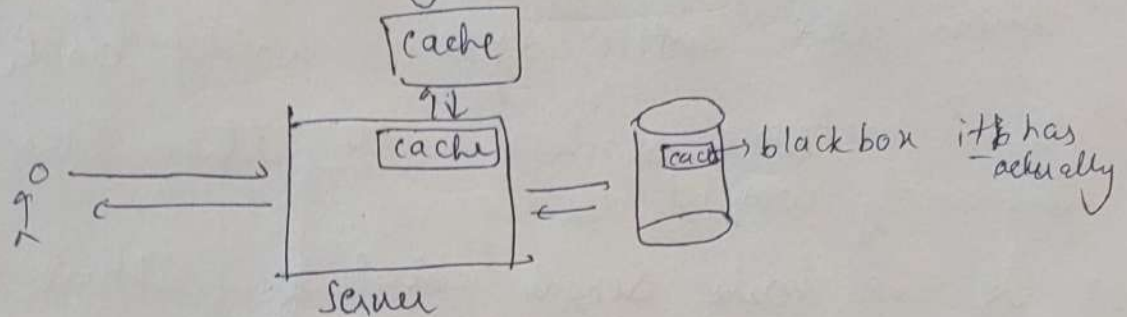
\rightarrow now there is update in data within 5 min eg no of youtube likes got increased in 5 min

\rightarrow now our db is having actual likes & cache is having the old data (i.e. data is not updated)

here we can see false data and true data is only seen after one hour.

in case of transaction user don't want to see false data
→ This cache update is dependent on the policy that we use.

Where is cache memory stored.



ideally ~~the~~ cache ~~is~~ is present Globally, inside db, inside server
independent in memory cache

ideal cache to choose is Global cache

→ as it can scale Independently (kabhi bhi kitna bhi bada kar sakte hai)

→ if we change algo of cache, then there is no need to redeploy the code on servers. i.e. deployment is independent

→ Multiple servers can use the same cache. (using same logic & executing query) that was executed earlier by some other server.

Eg brother asks 15 square one morning → it will take a some time
if he asks same in evening we couldn't be able to response quickly