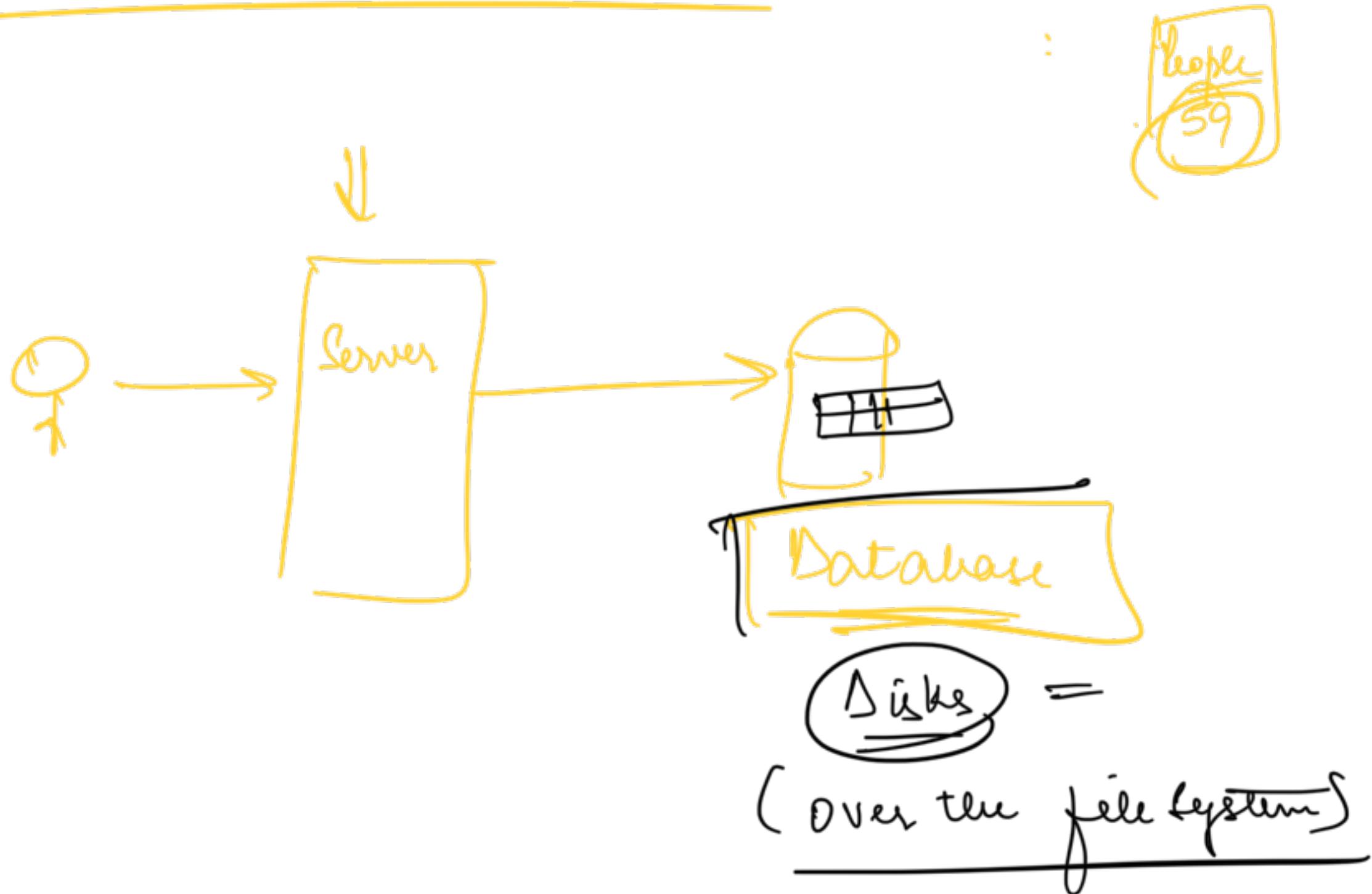


# Design a Distributed Cache

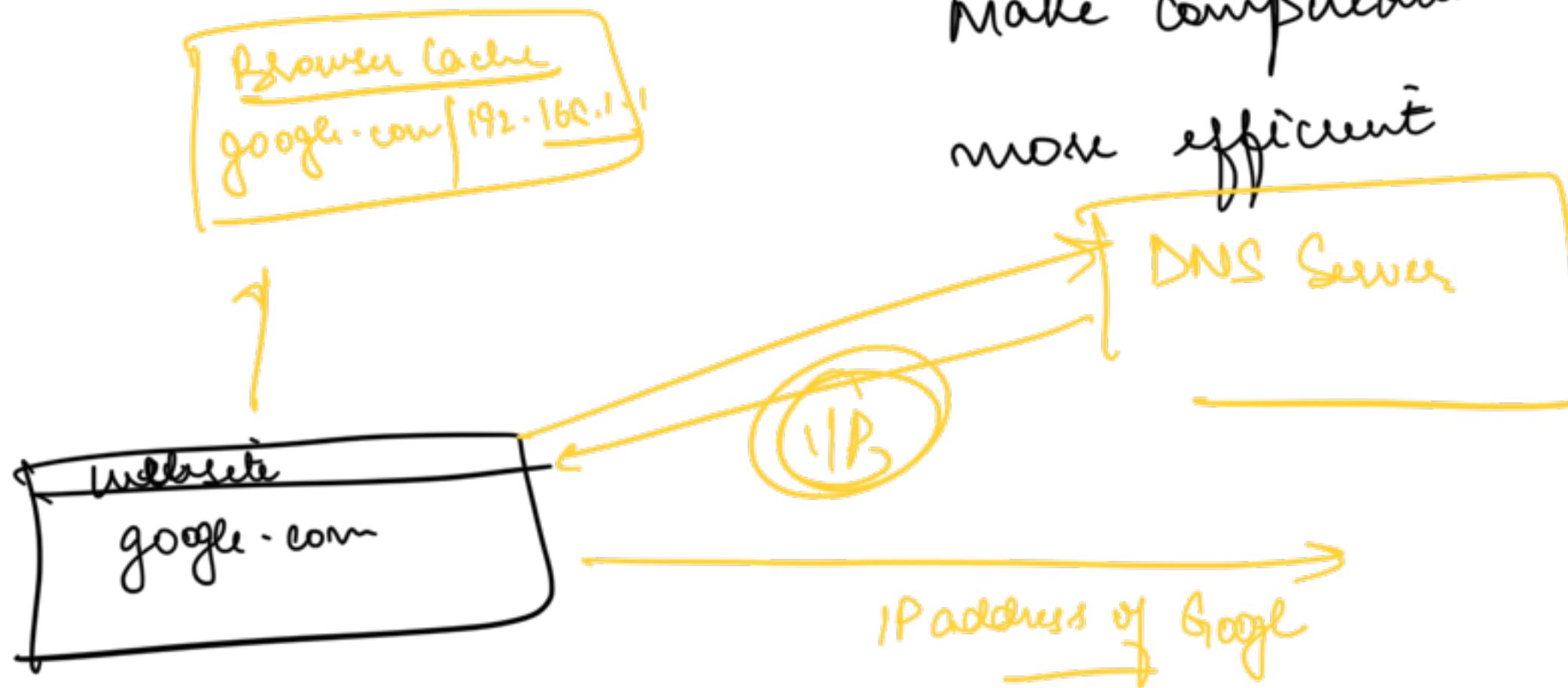


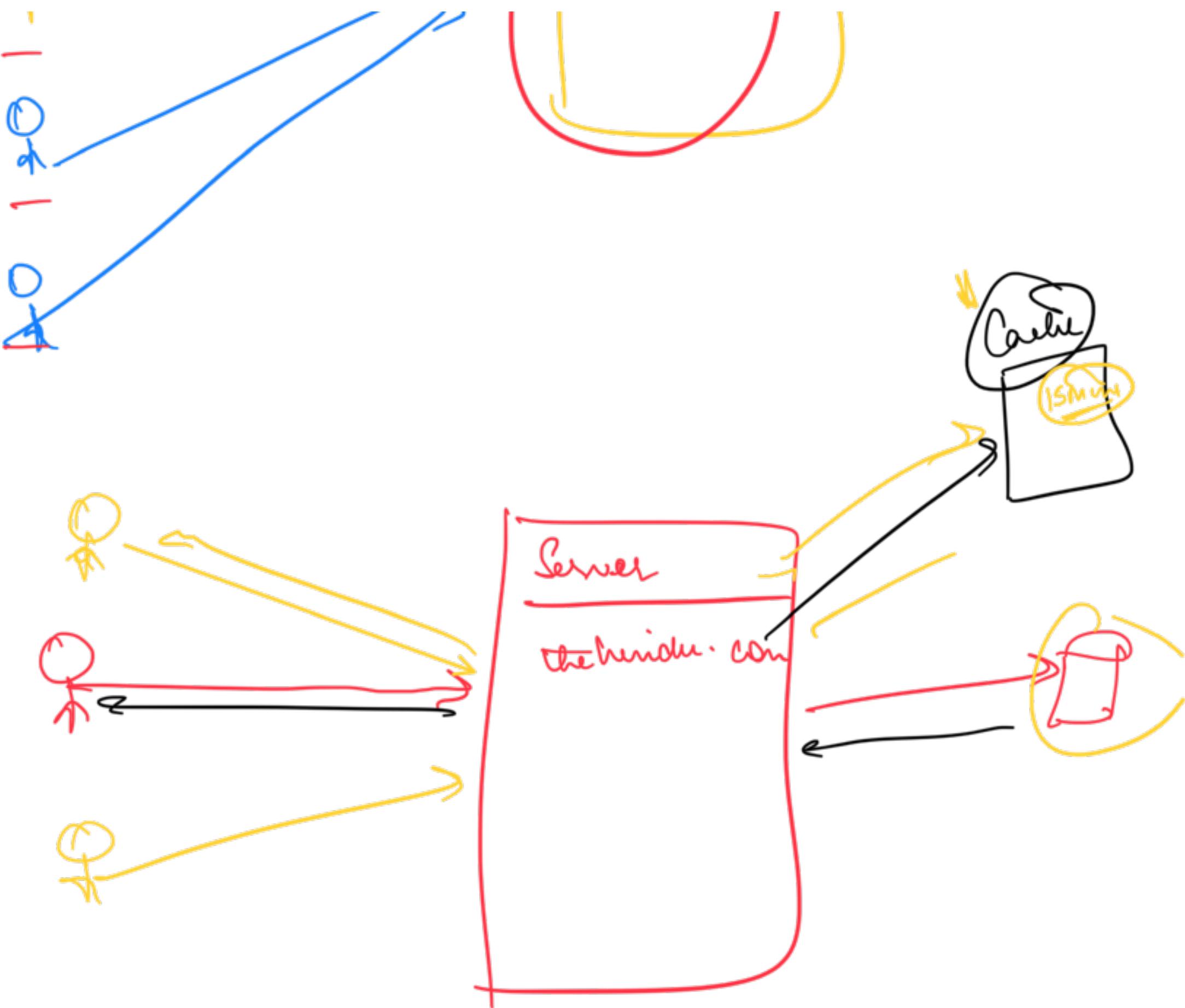
mat

Caches  $\Rightarrow$  temporary storage  $\rightarrow$

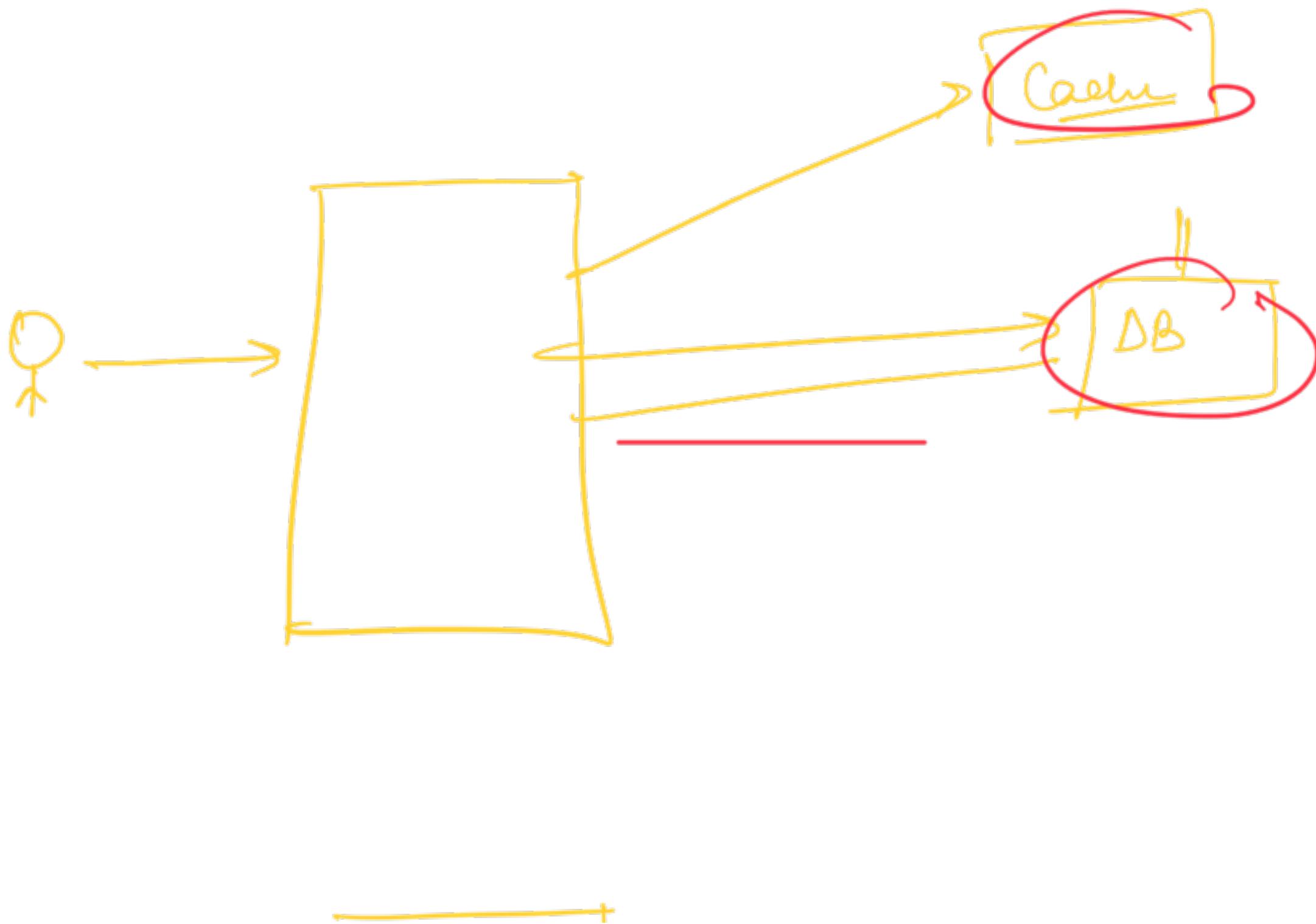
make computations

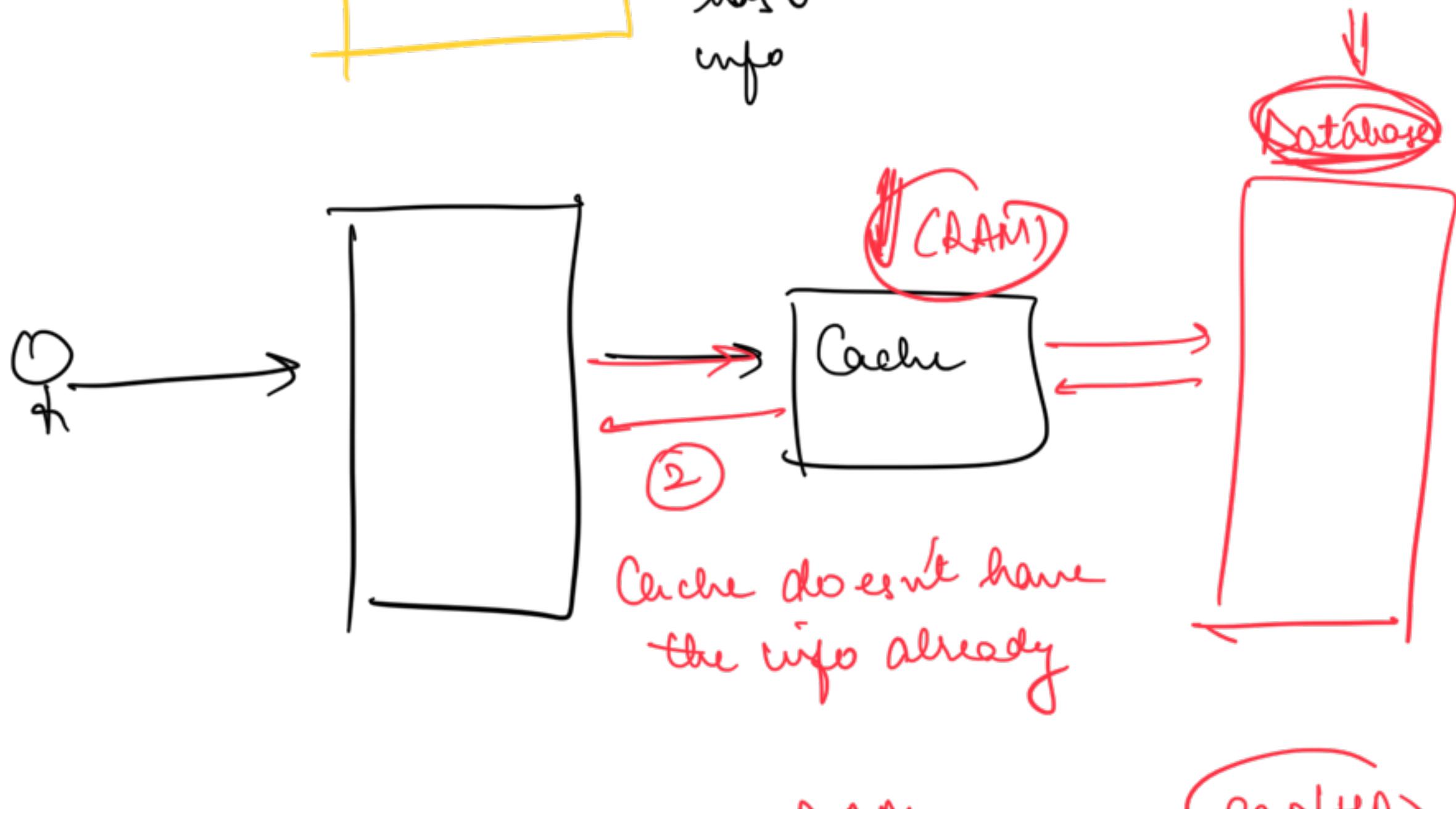
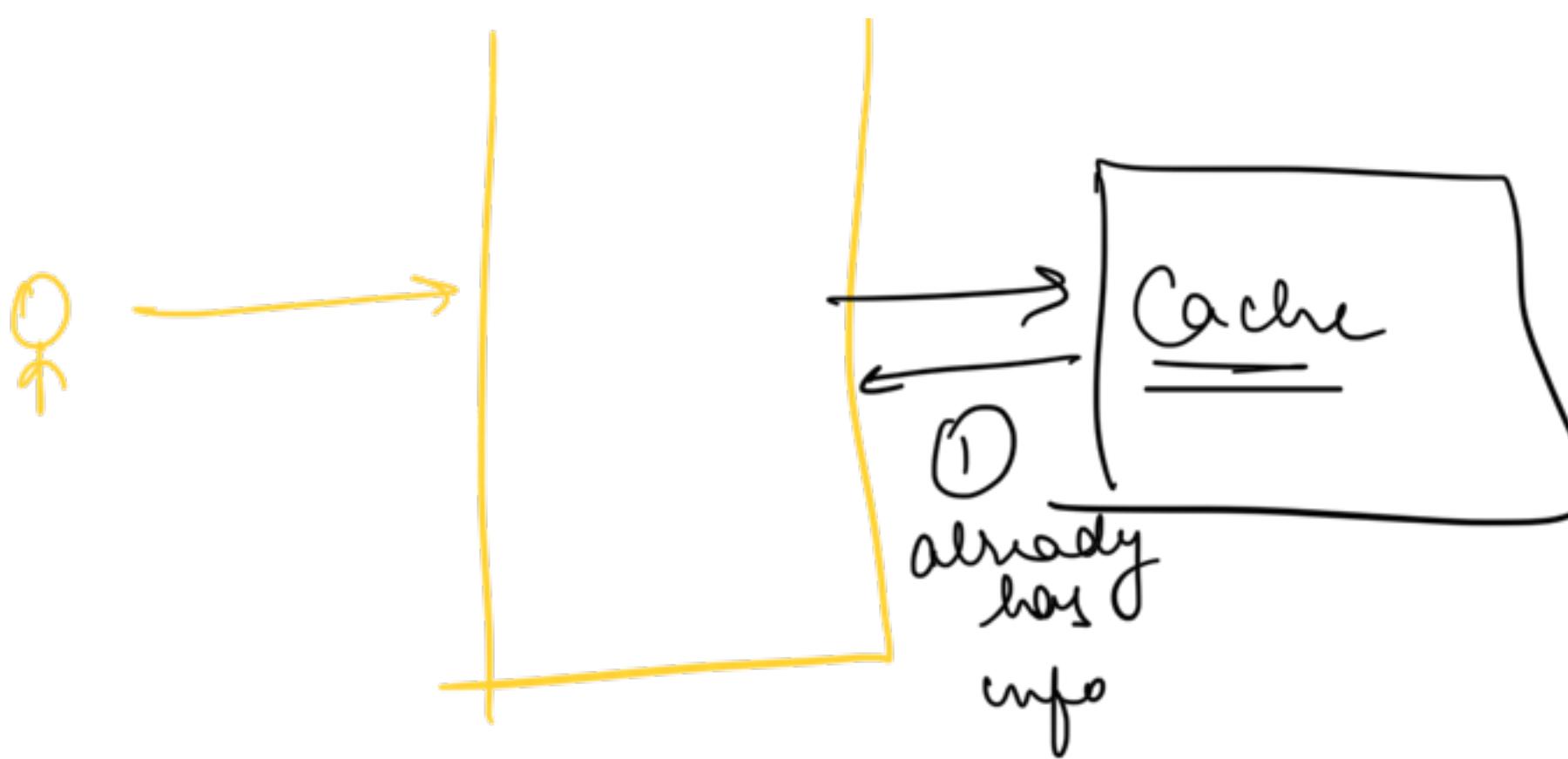
more efficient





## Distributed Cache





RAM

16 GB,

[Expensive)

low storage

(large)  $\Rightarrow$  high mem  
medium cpu

(SDI)

512 GB

[Cheap)

high storage

low cpu

low mem

portals

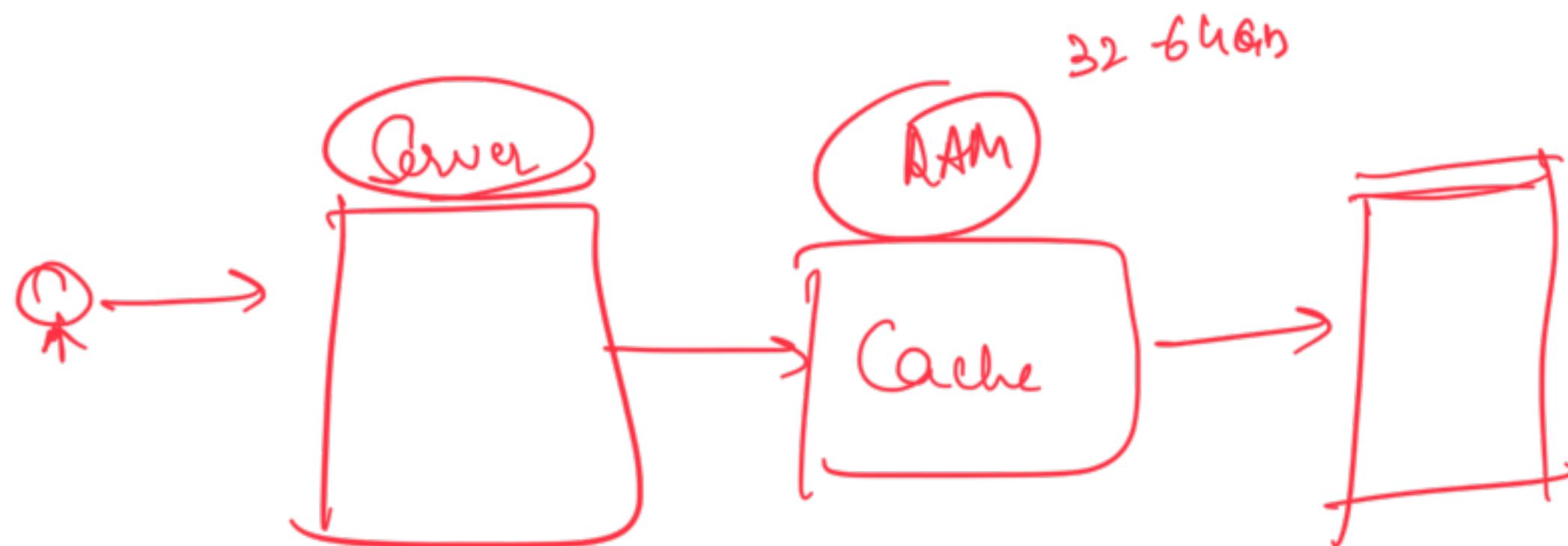
NM

RAM  $\Rightarrow$  512 GB

RAM  $\Rightarrow$  128 GB

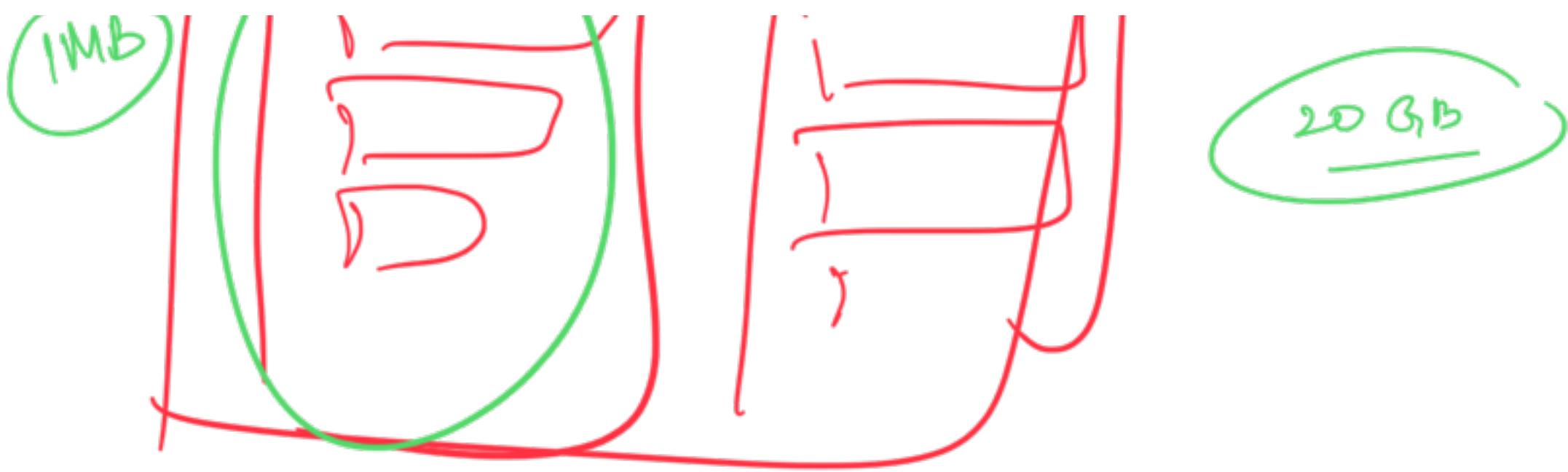
Storage = 2TB

CPU



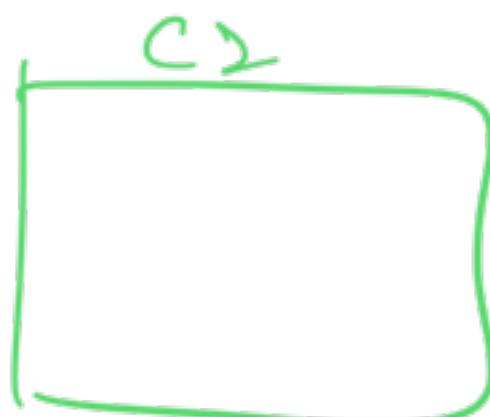
Distributed Cache







(A - E)



(F - J)



(K - O)



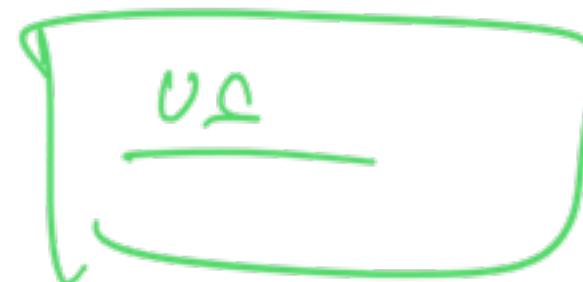
(P - T)



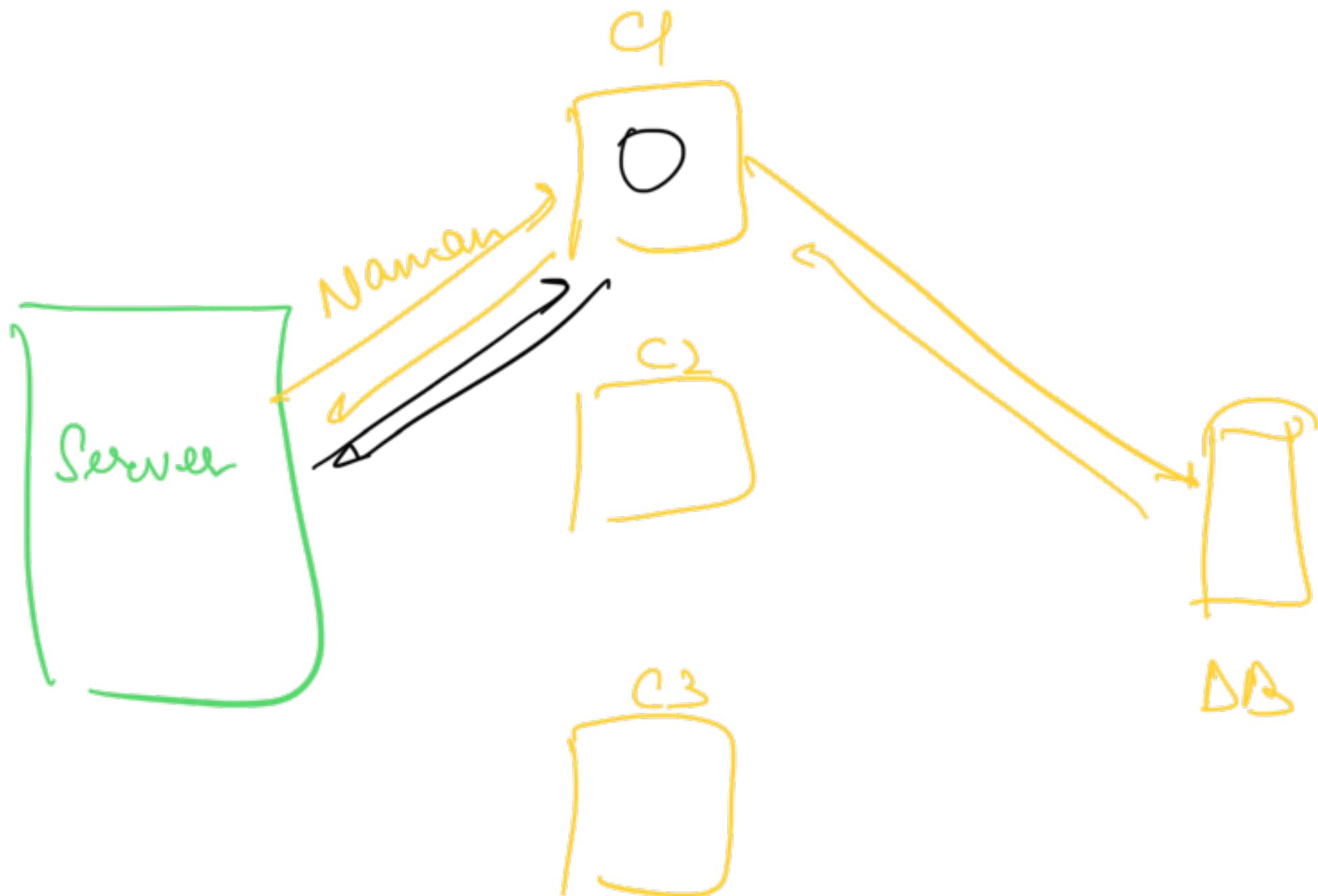
(U - Z)

geo sharding

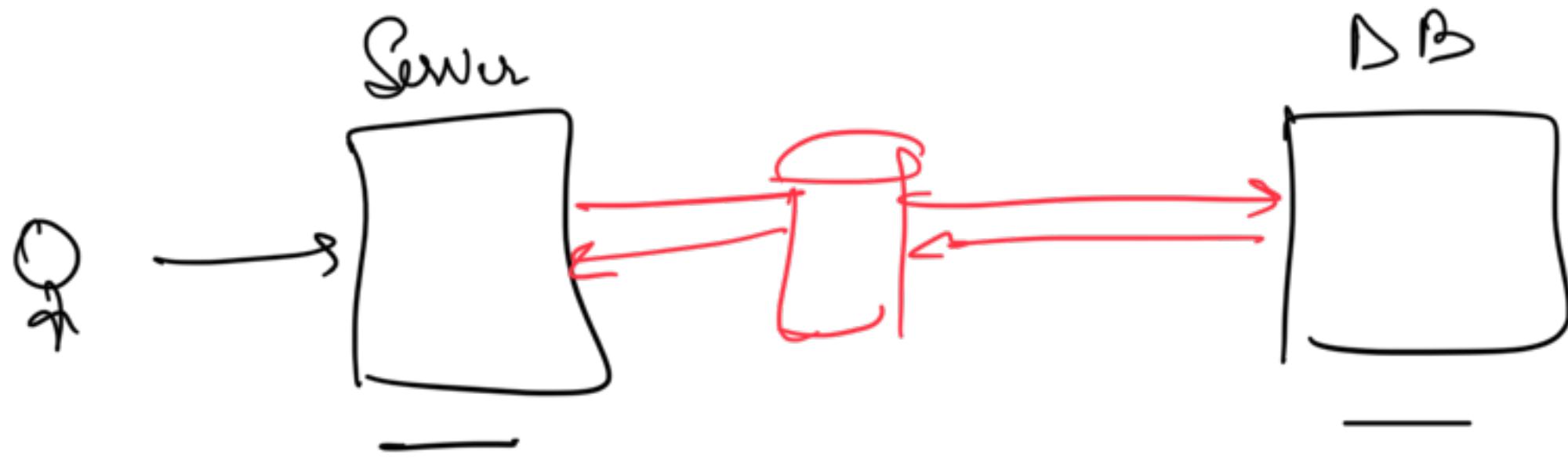
distribute based on location



UK



# Design a Distributed Cache



Requirements

- handle
- ① Should support ~~CRUD~~ in a unit till write policies
- C → Create / Add  
R → Read / fetch  
U → update  
D → Delete

- (2) Shovel support  $\rightarrow$
- (2) Limit on max size of cache & support for diff eviction strategy

(1) Support for TTL

(5) Hot loading (Warming up the cache)

(6) Handle Concurrency & Thread Affinity

(7) Request Collapsing



@Getter

@Setter

@DI

RAM

$\Rightarrow$  POJO

Create

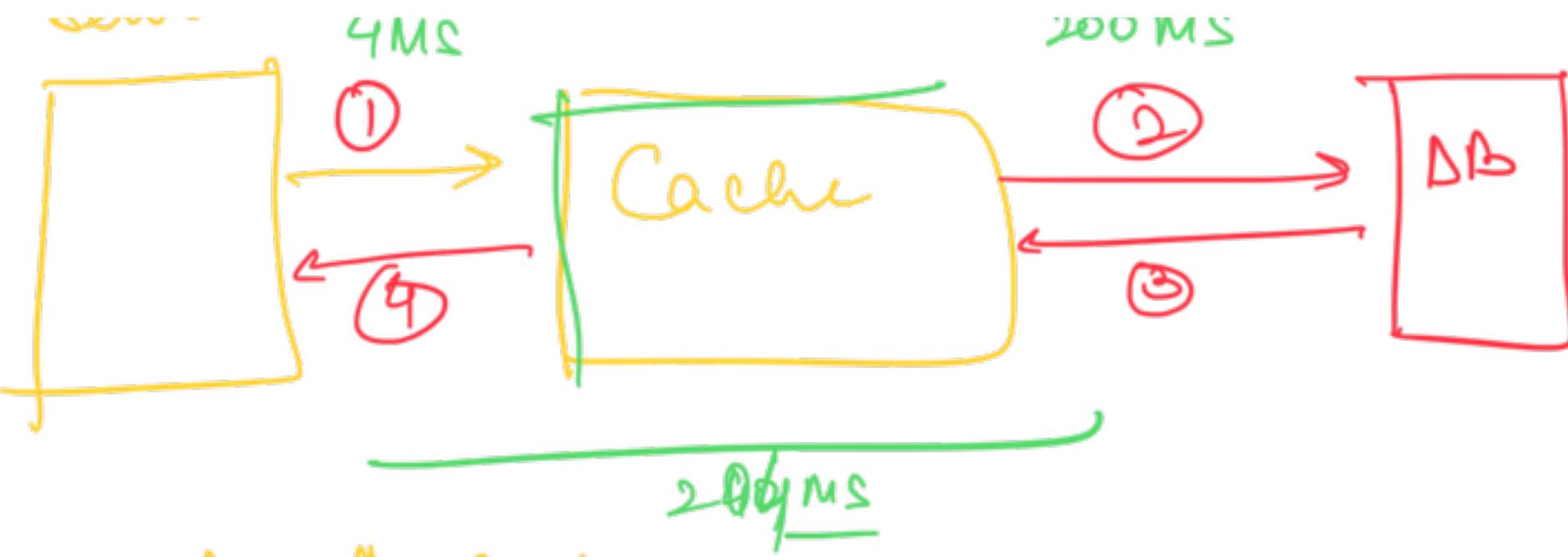
2 Types of Cache Based on How a Write is Stored

①

Write Through Cache

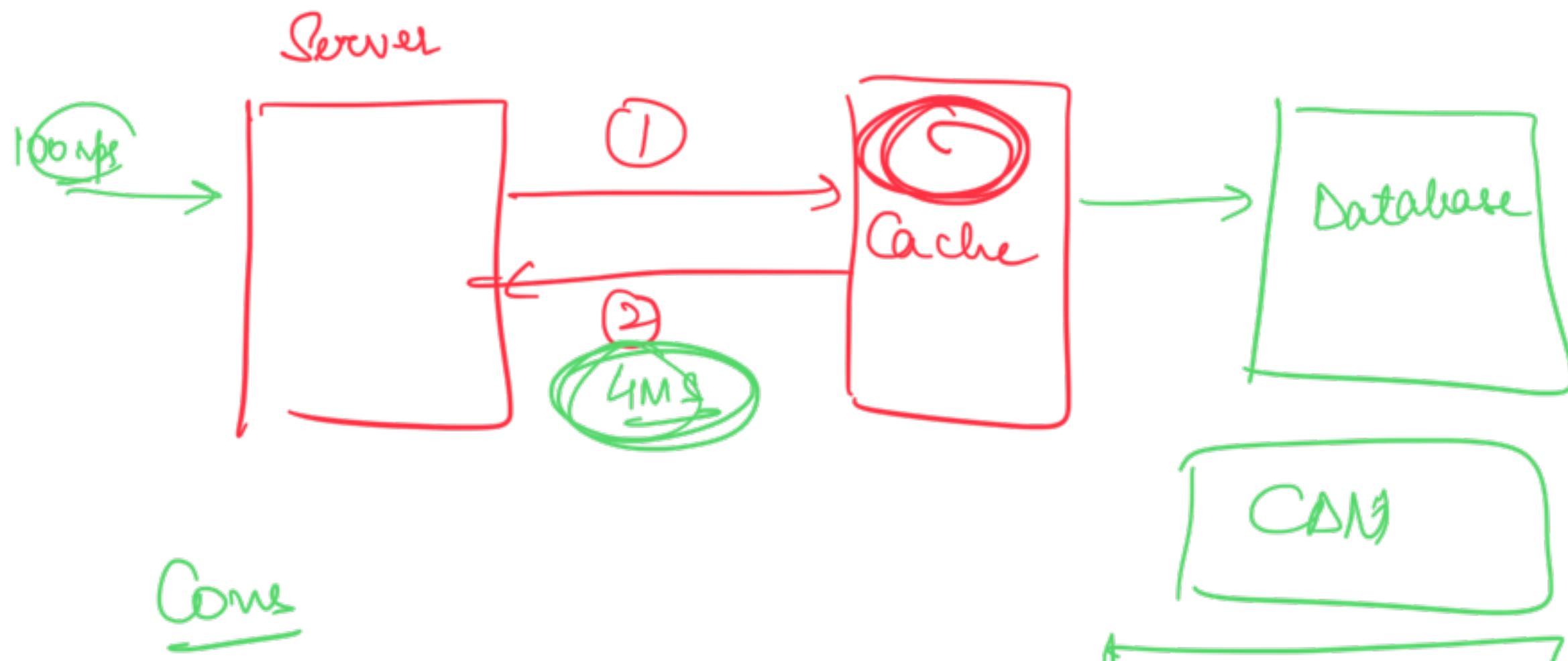
Concurrent

Memory



② Write back Cache

Twitter Newsfeed



Cdns

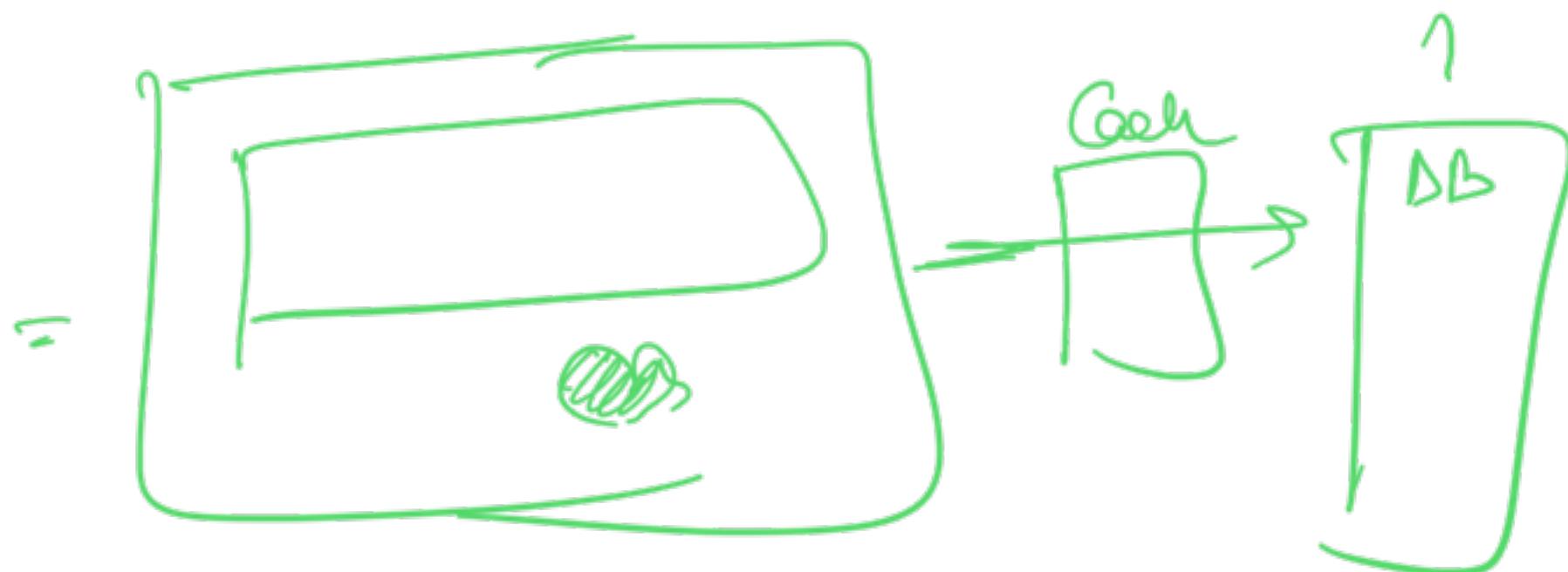
① Data Loss

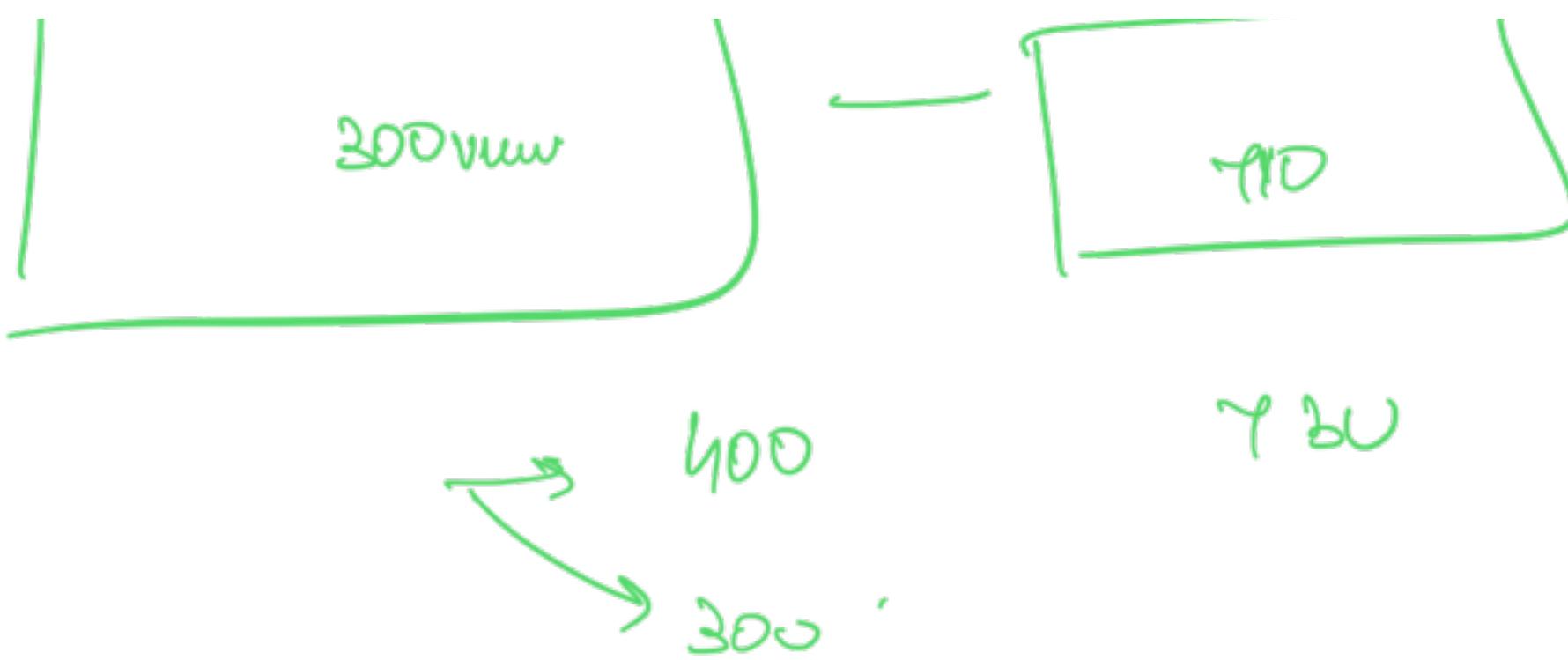
Pros

① Efficiency (Fast)



Write Heavy

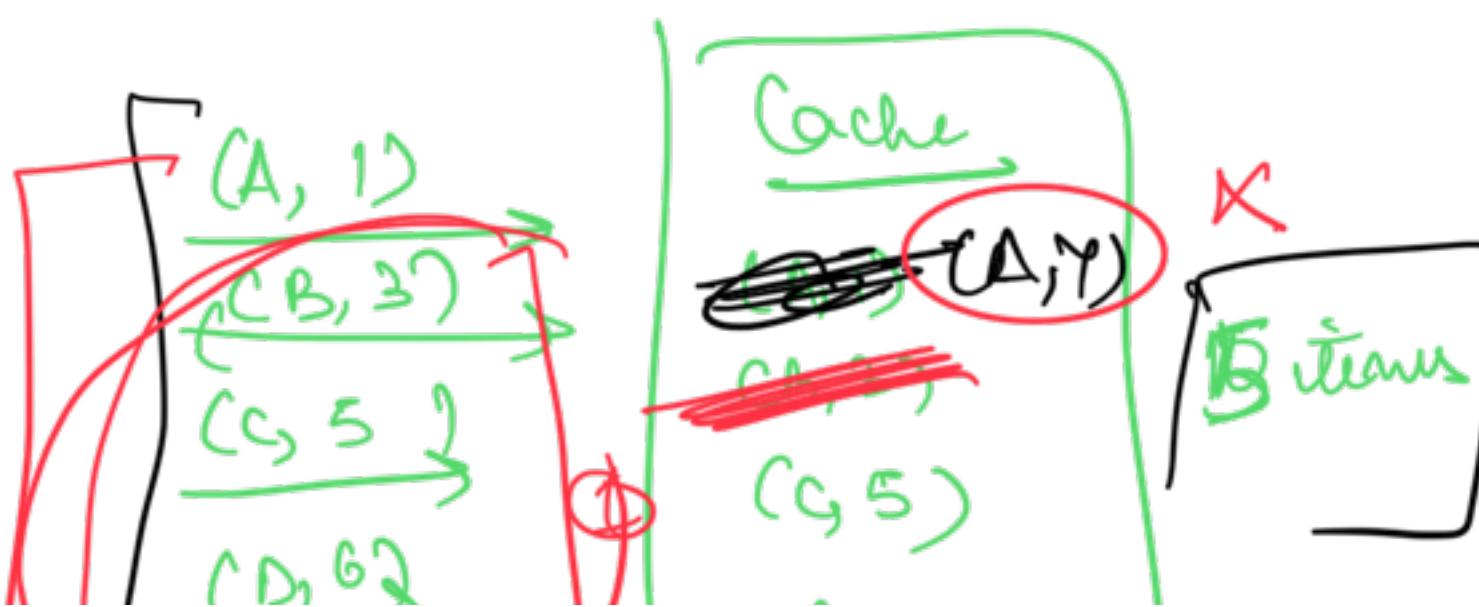
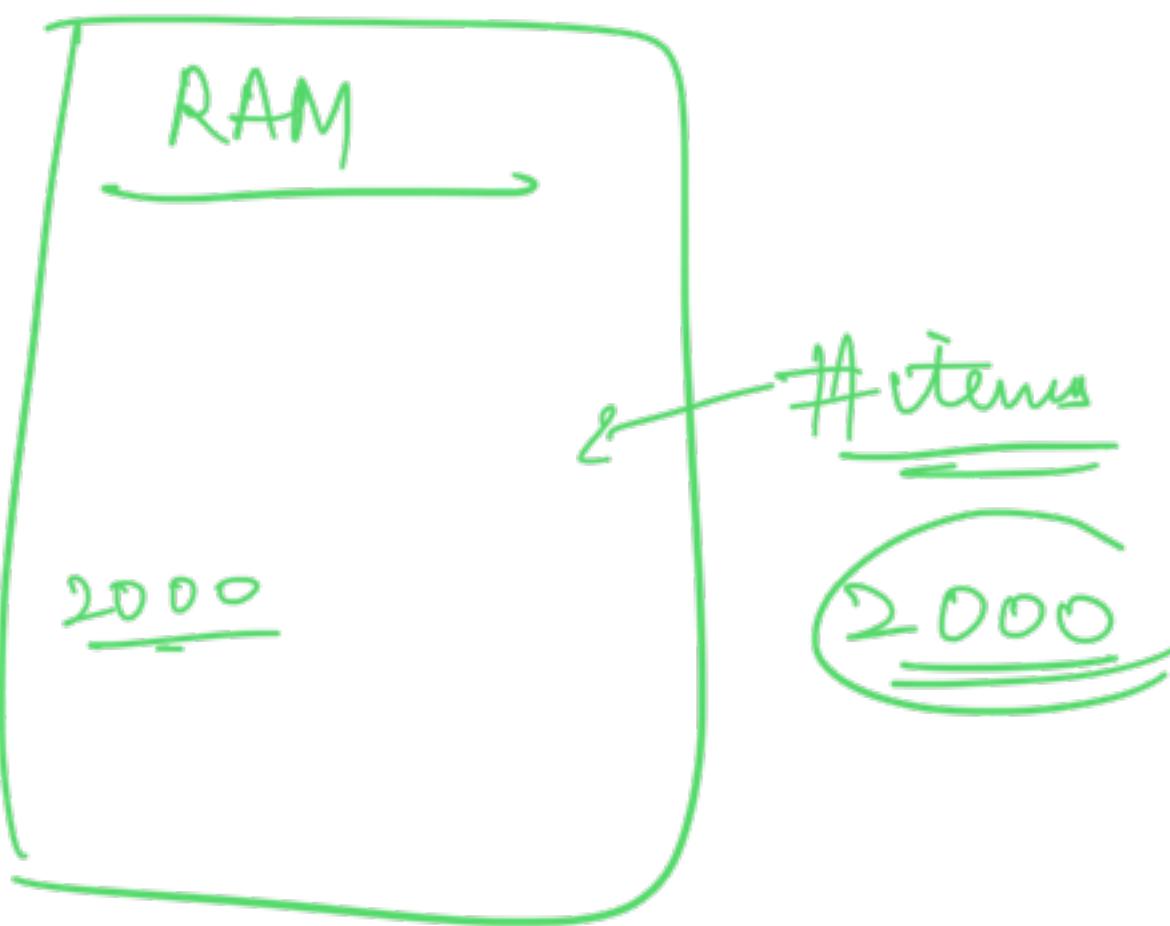




Twitter → Delete Your Account

+1 day Twitter → Reactwall







What to do when you add something after cache  
is full  $\Rightarrow$

- ① Remove one of the existing items
- ② Add the new item

① Random

② LRU (Least Recently Used)



LFU

Remove the item that is used  
the least # of times

Quick Sort

Pivot: left, middle, right

Randomized Quick Sort

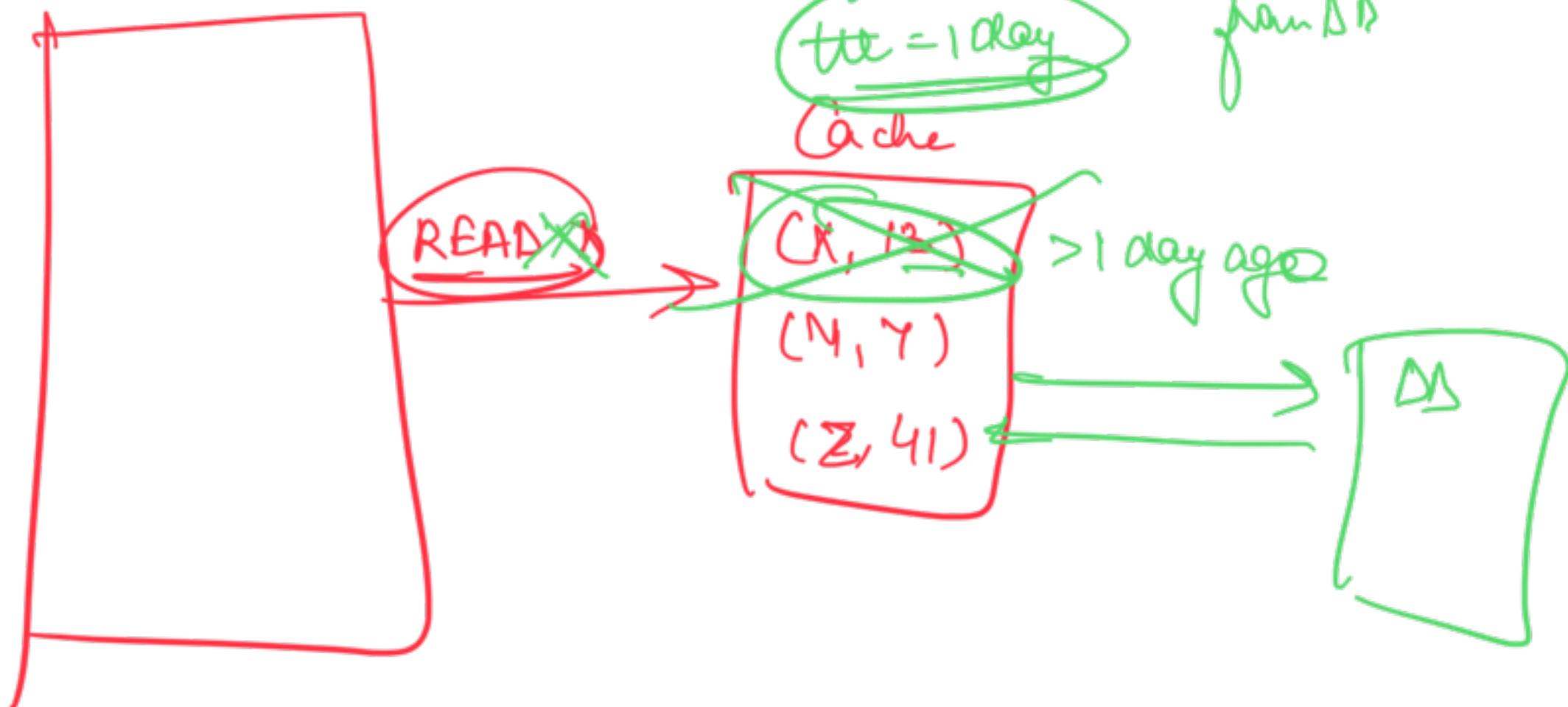
Pivot: random

Read the proof of probabilistic TC of randomized  
 $\Theta(\log N) \rightarrow O(\log N)$

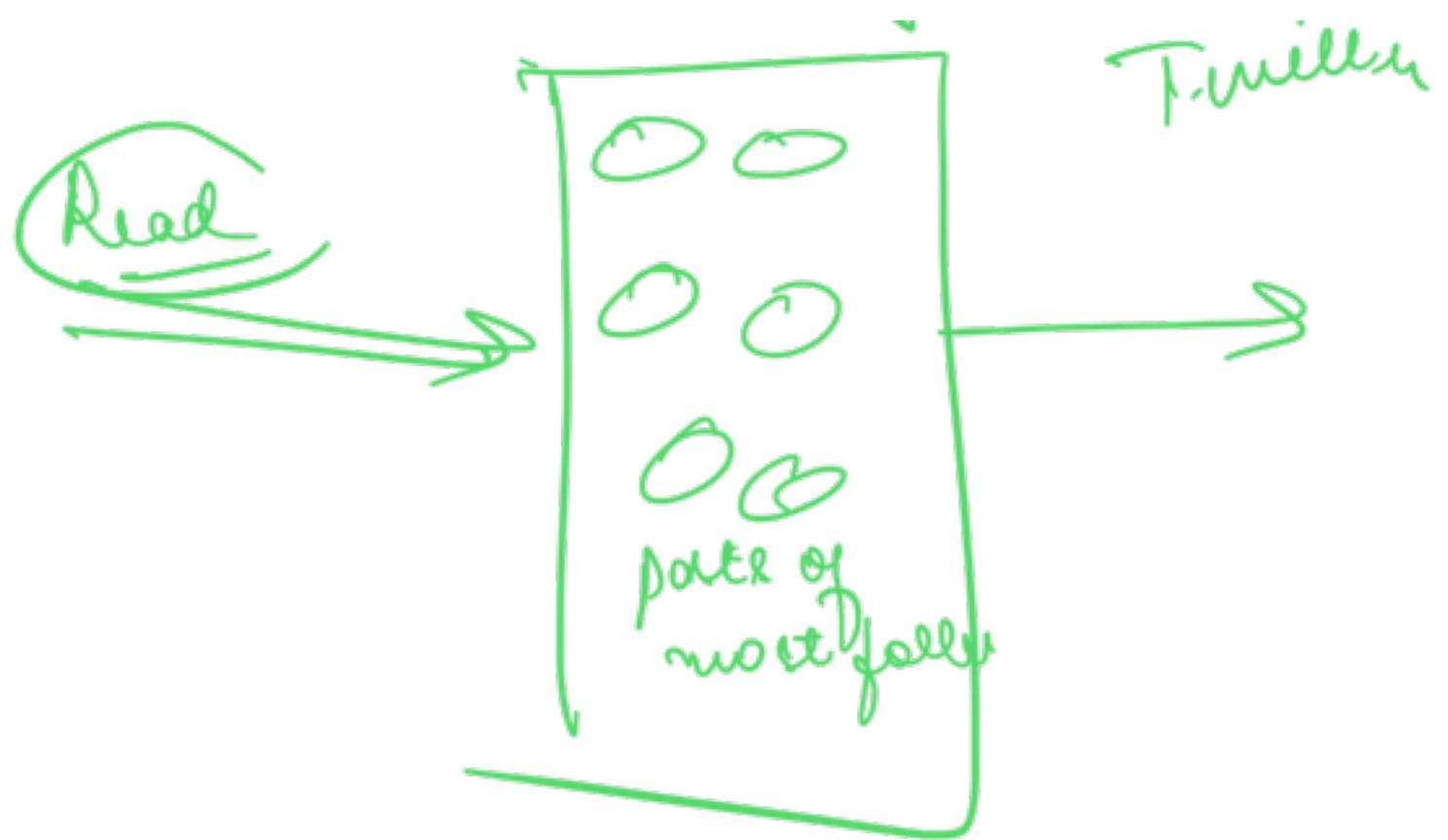
Cache

$O, \text{ttl}$  → time to live

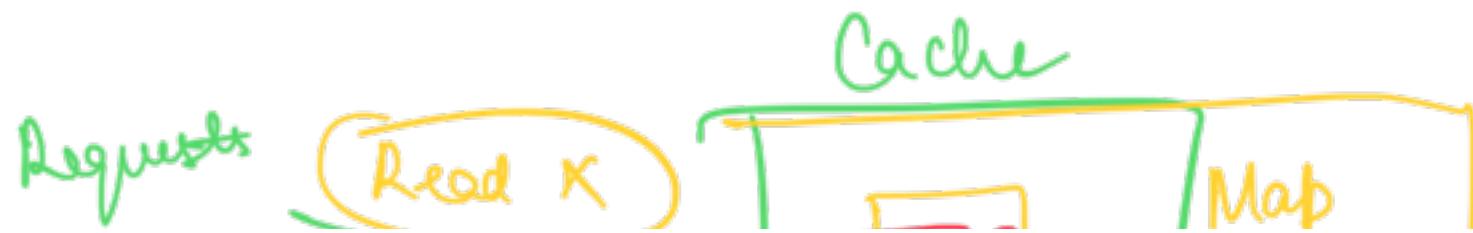
if data in cache is older than ttl, invalidate (remove) from cache and get new data from DB

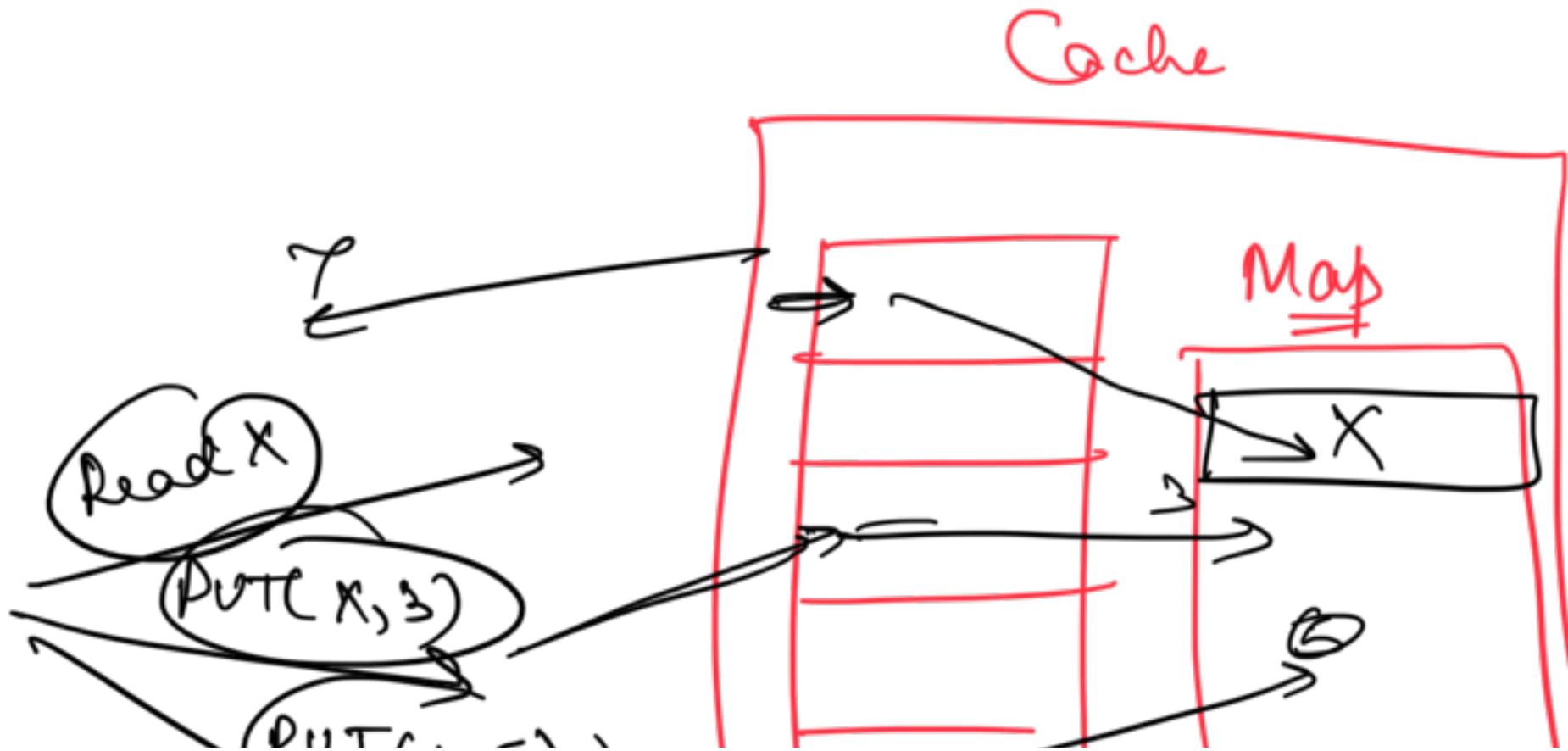
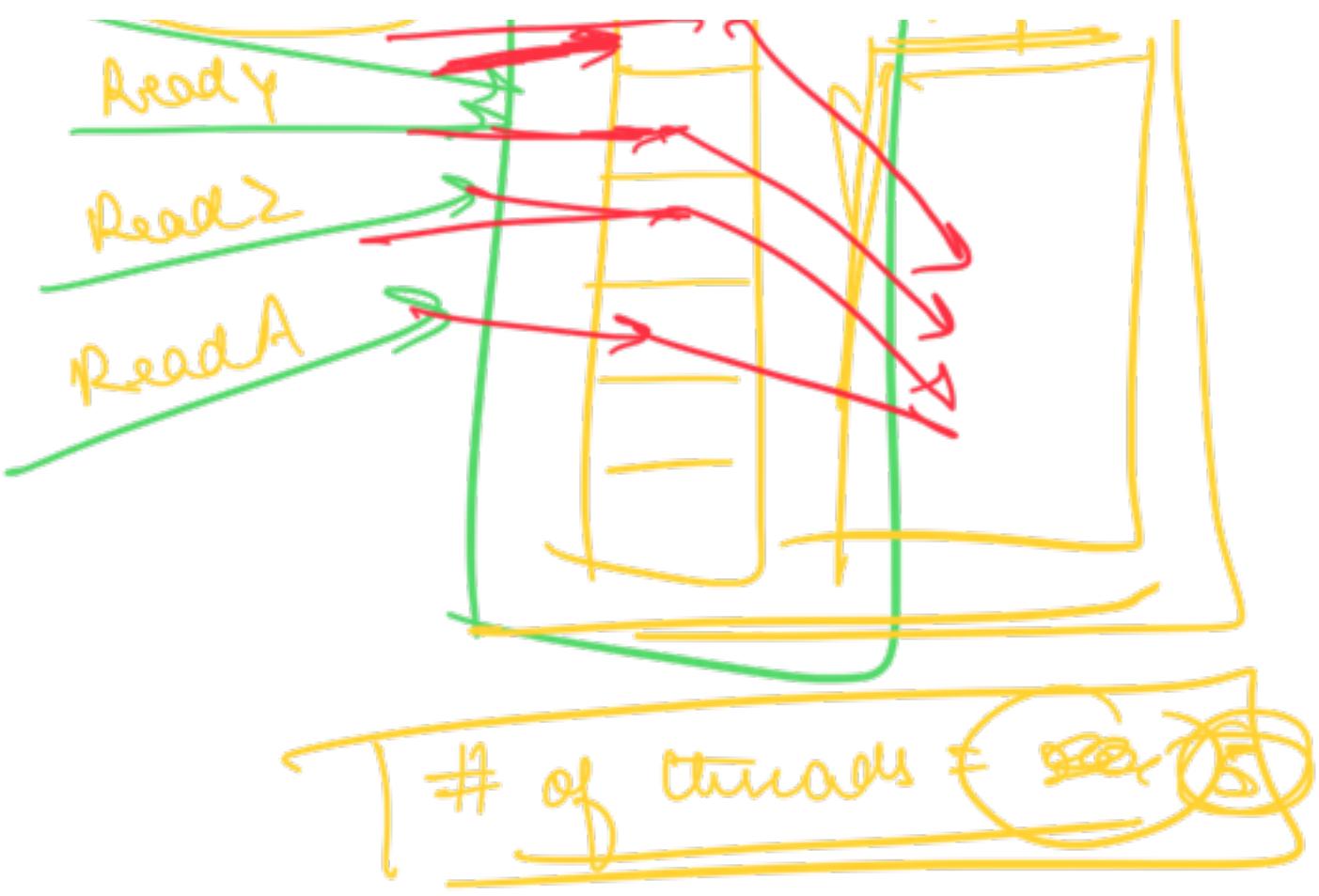


initially cache will be slower



Hot loading  
Cache Warming







sol

①

lock → flow

②

Cache





Concurrency issues because multiple threads same key

What if I say that all the access to a particular key will happen from 1 thread

→ 1 Thread

→ Executor

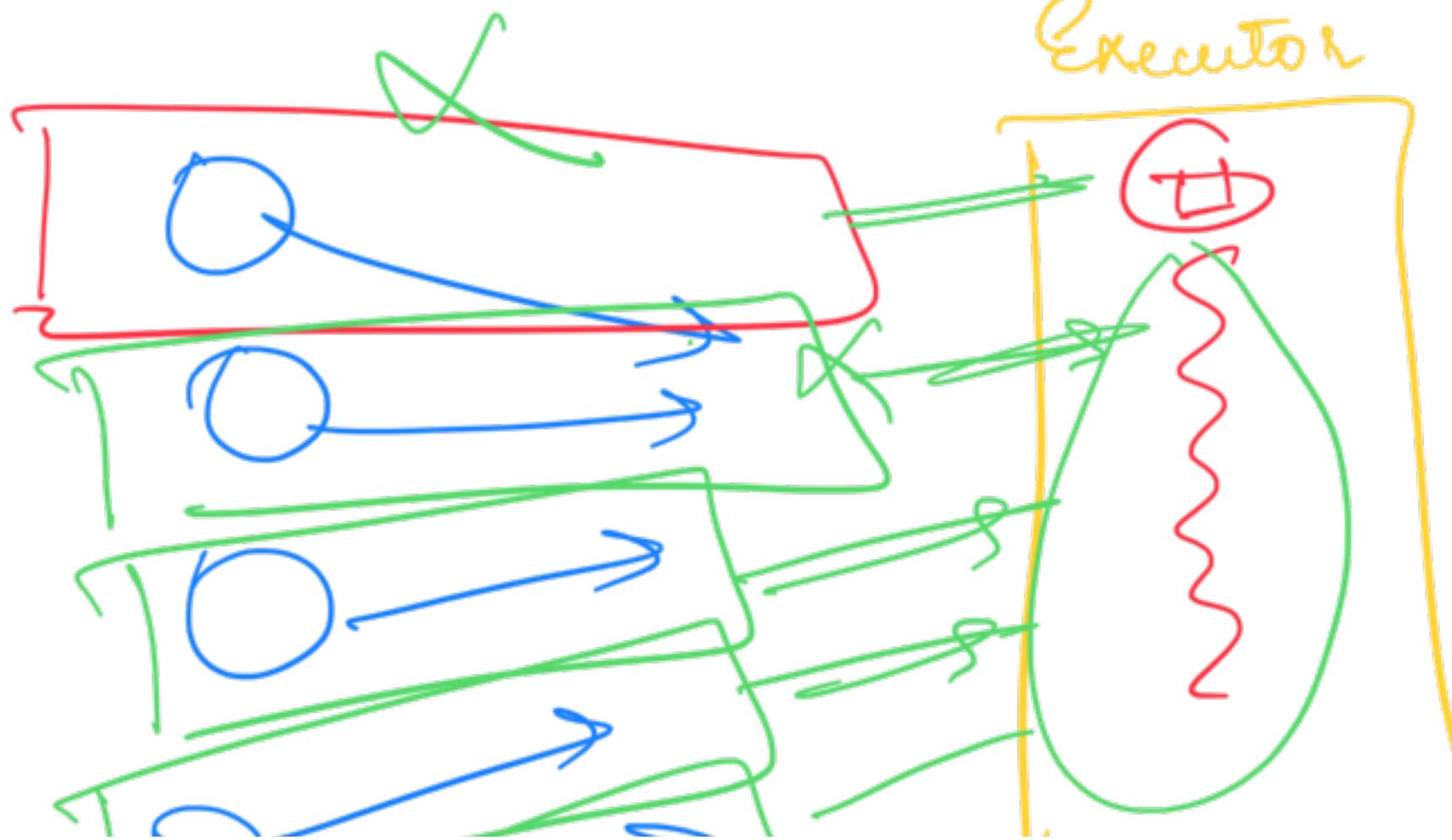


fixed Thread Pool

Cached Thread Pool

Single Thread Executor

Executor





No of Register = 16

Cache {

1Byte       $10^9$       1GB

Executor [15]  $\Rightarrow$  [STE, STE, STE, STE, STE]

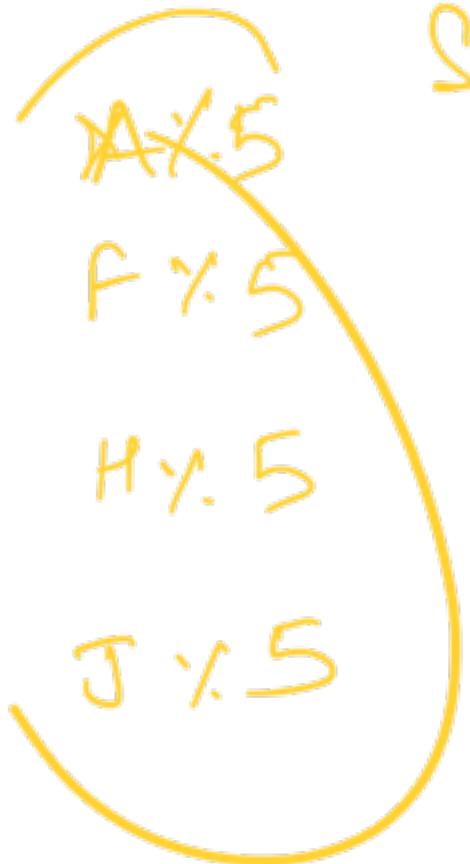
Save (X)  $\rightarrow$  init  
20 h

Executor [ $X \% 5$ ]. execute(Save(i))

3

3

Still faster than lock



Thread Affinity



delete (\*)

get (\*)



Box create  $\langle K, \text{ID} \rangle$

get  $\langle C, K \rangle$

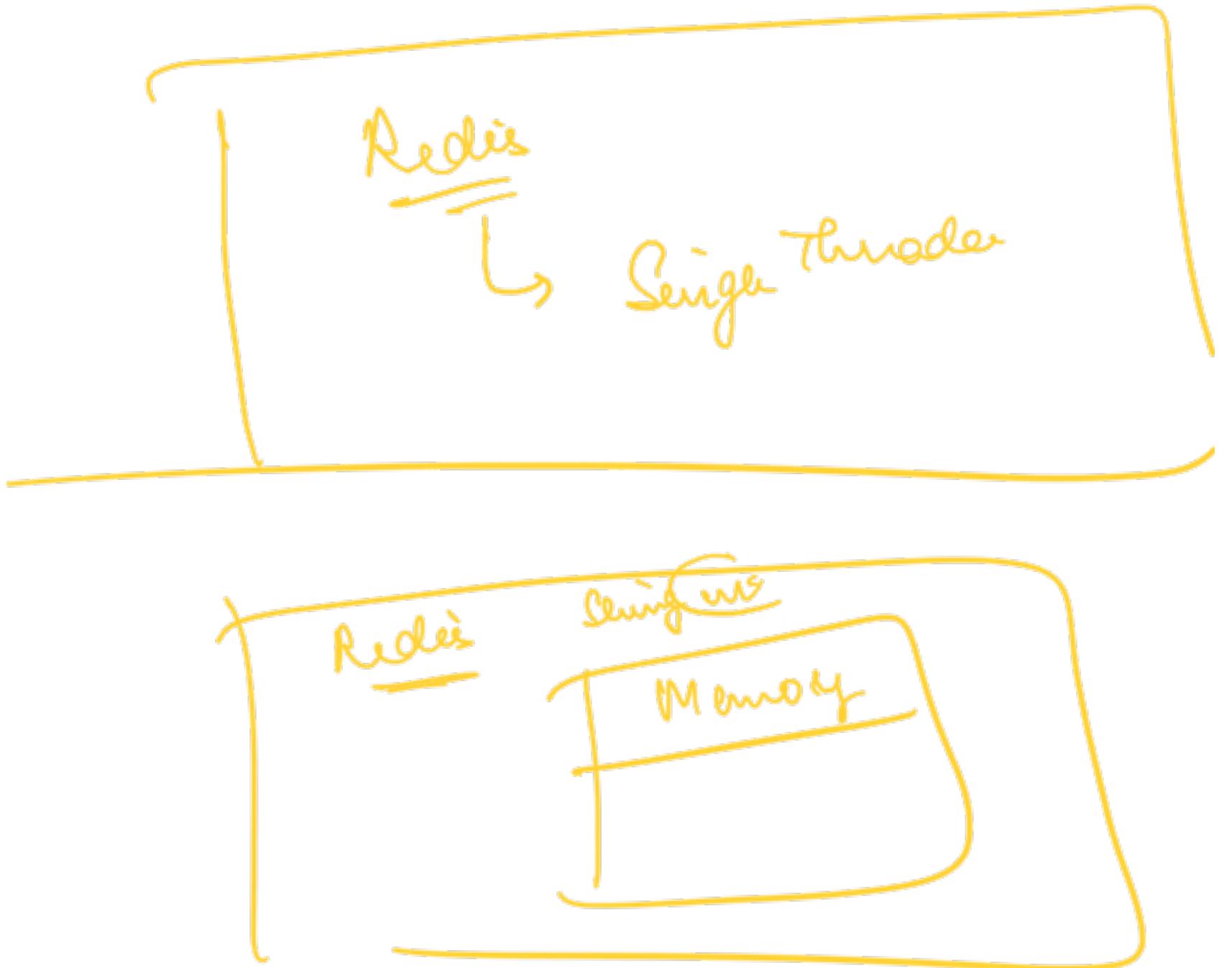
update  $\langle K, V \rangle$

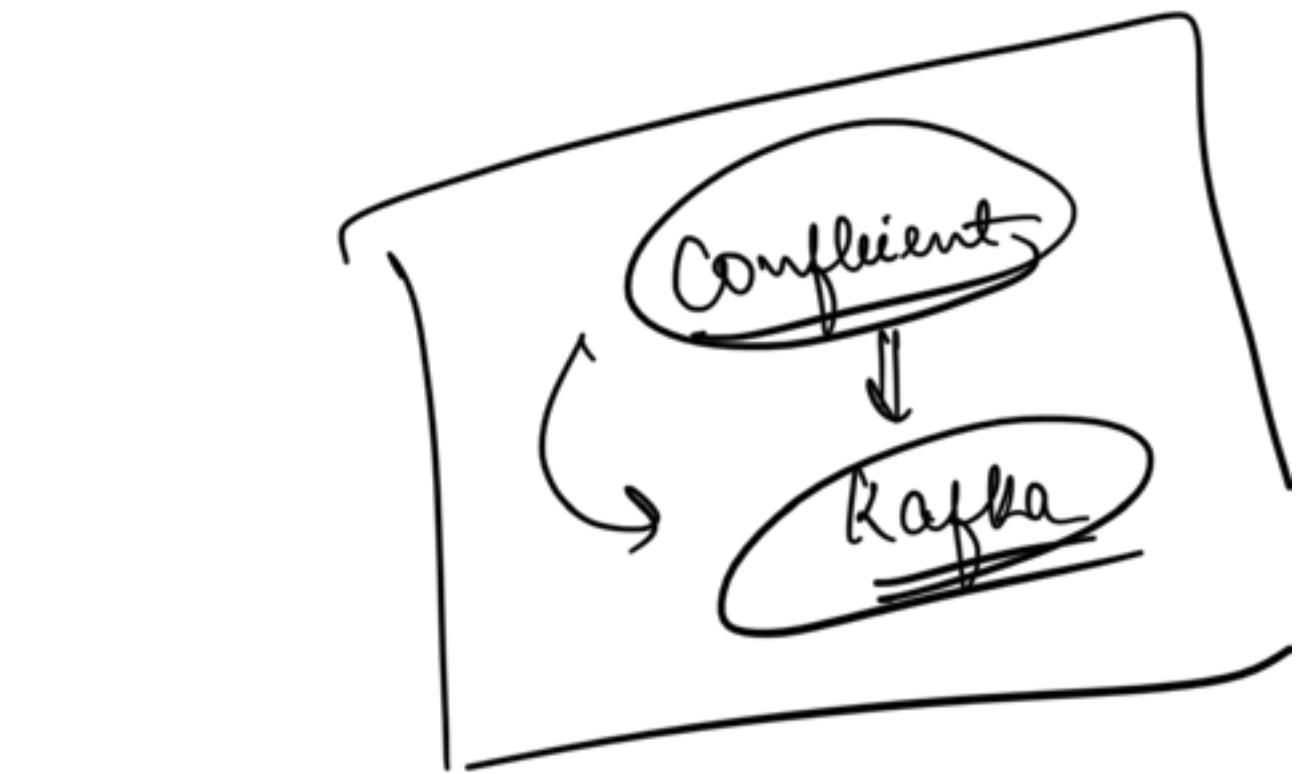
delete  $\langle K \rangle$



If your op<sup>v</sup> are ID intensive  $\Rightarrow$  More

If your op<sup>v</sup> are CPU intensive  $\Rightarrow$  less





Dist Cache

Concurrent

Dist Cache

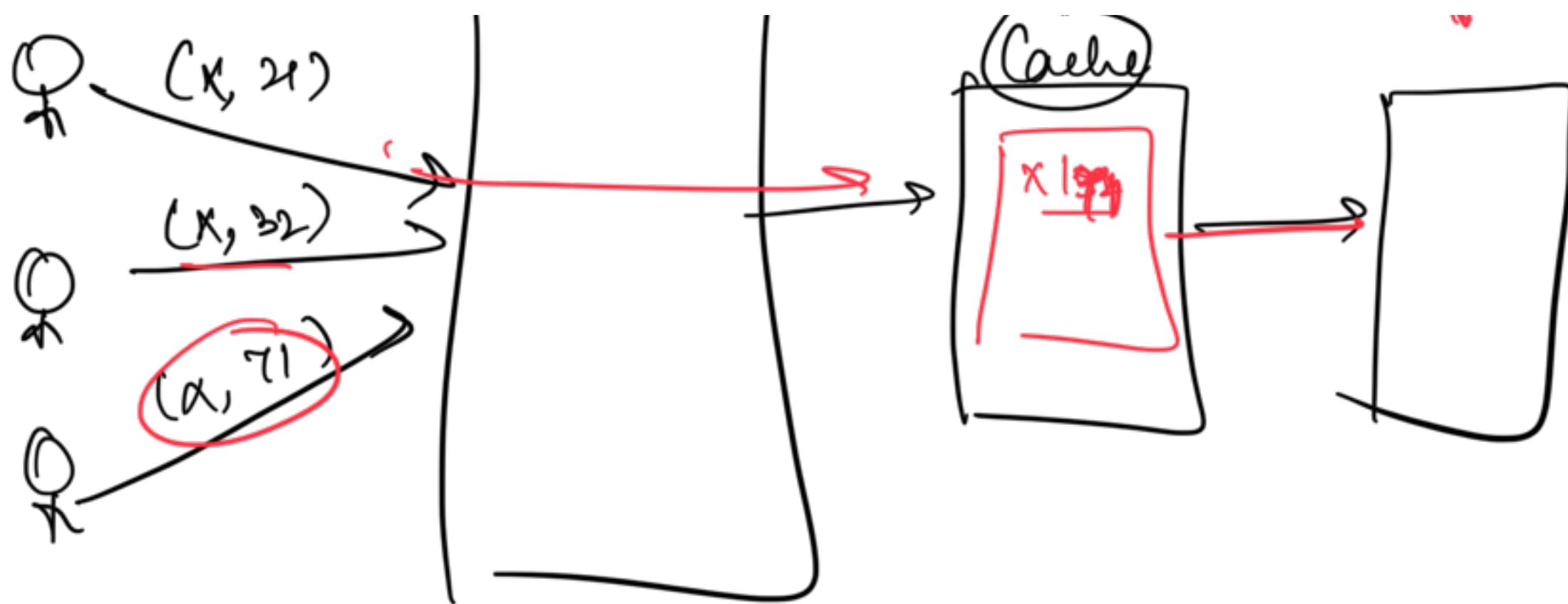
each op<sup>n</sup> → fetching from dist<sup>n</sup>

Request Collapsing



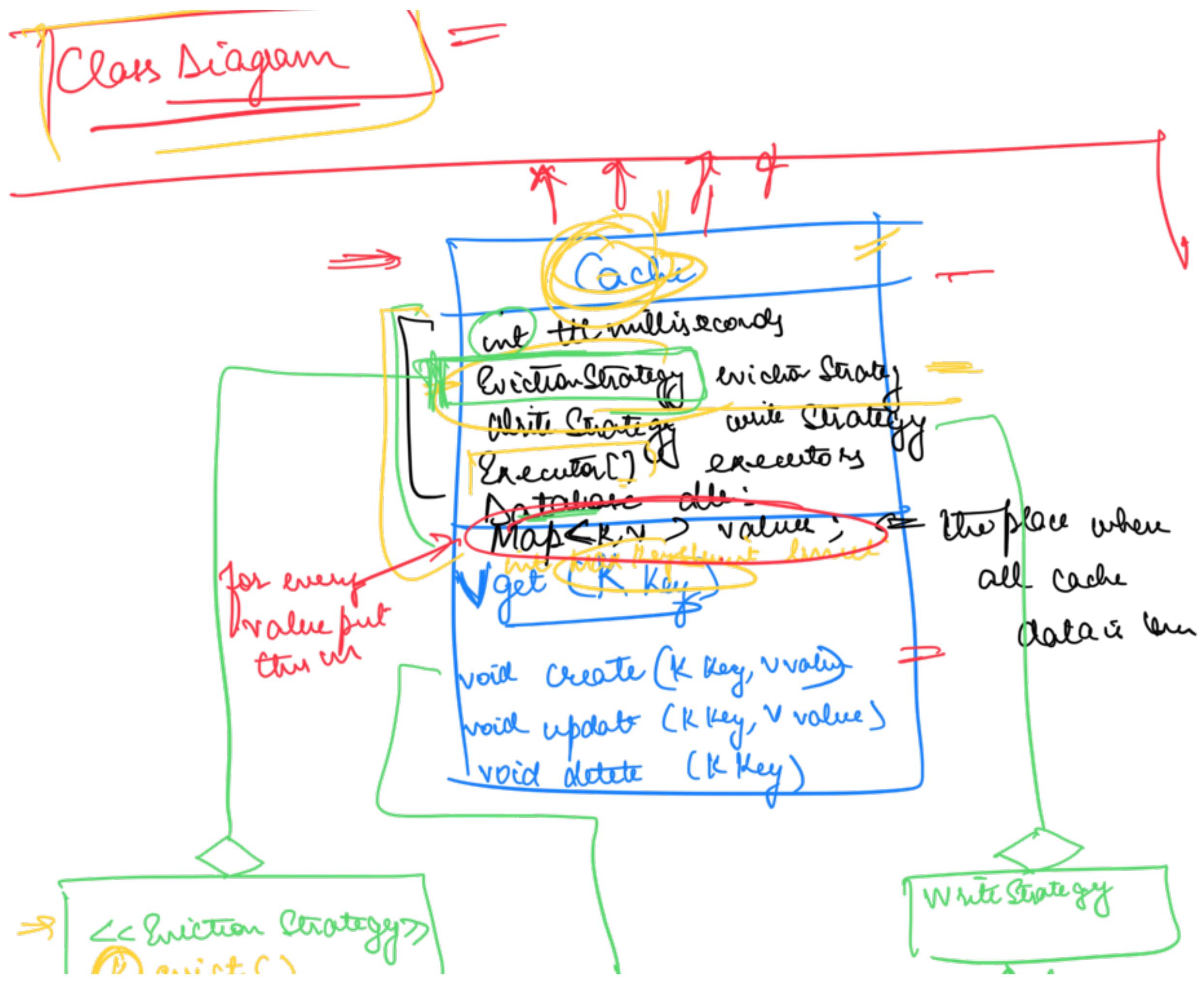
3 write  
↓

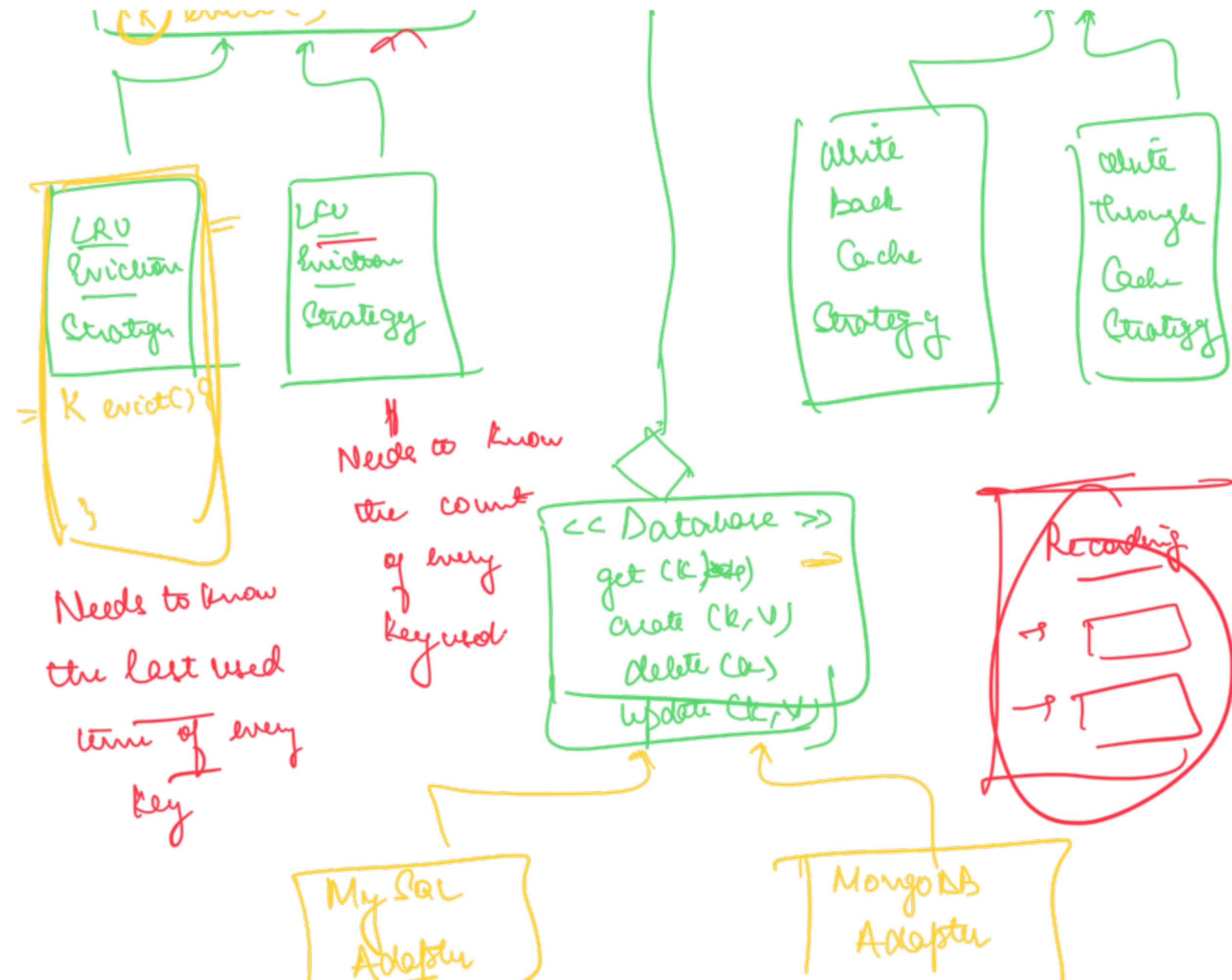
1 write  
↓



### Request Collapsing

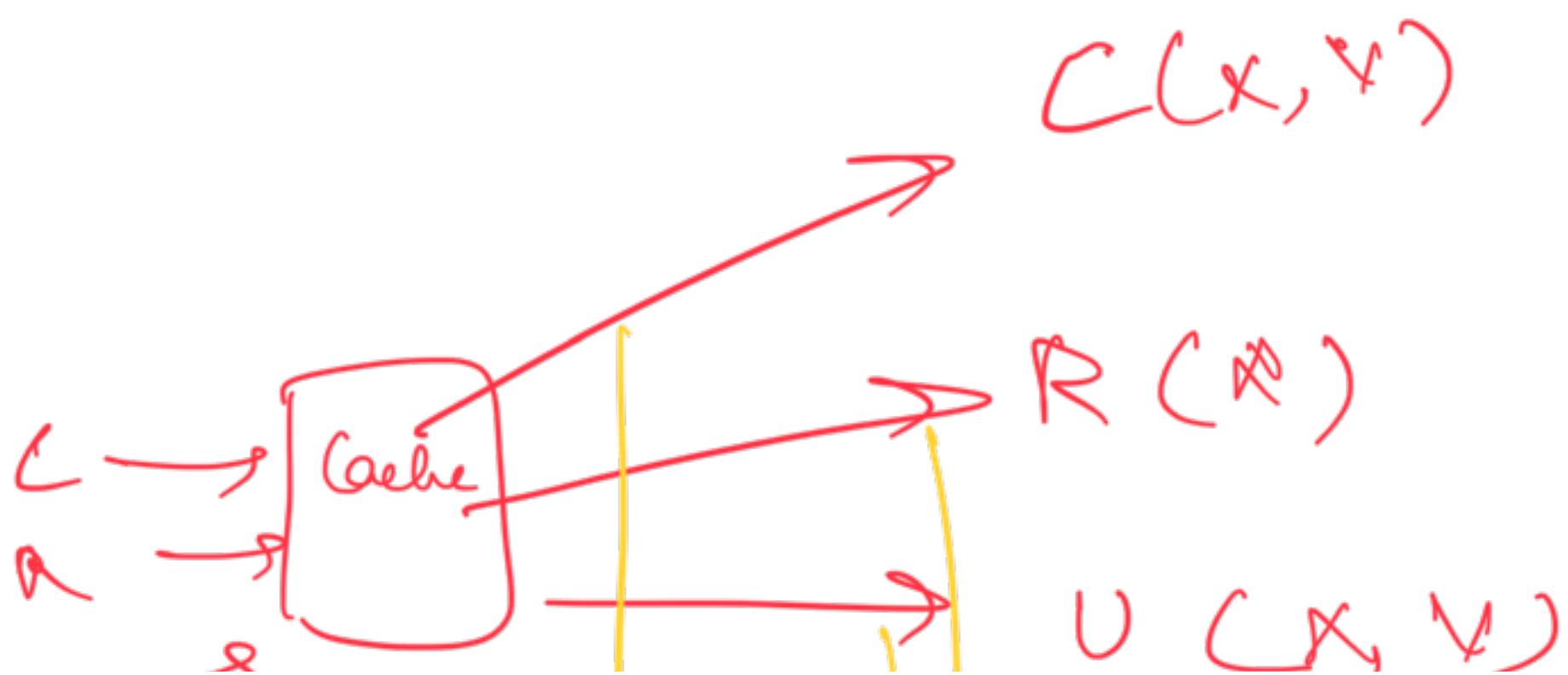
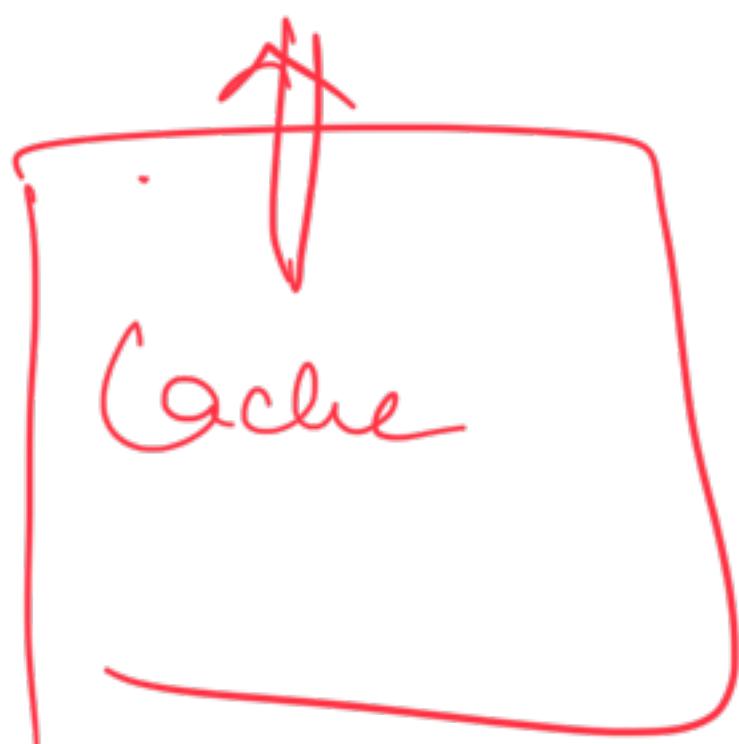
- Write Back Cache
- Reduced # of Write Back



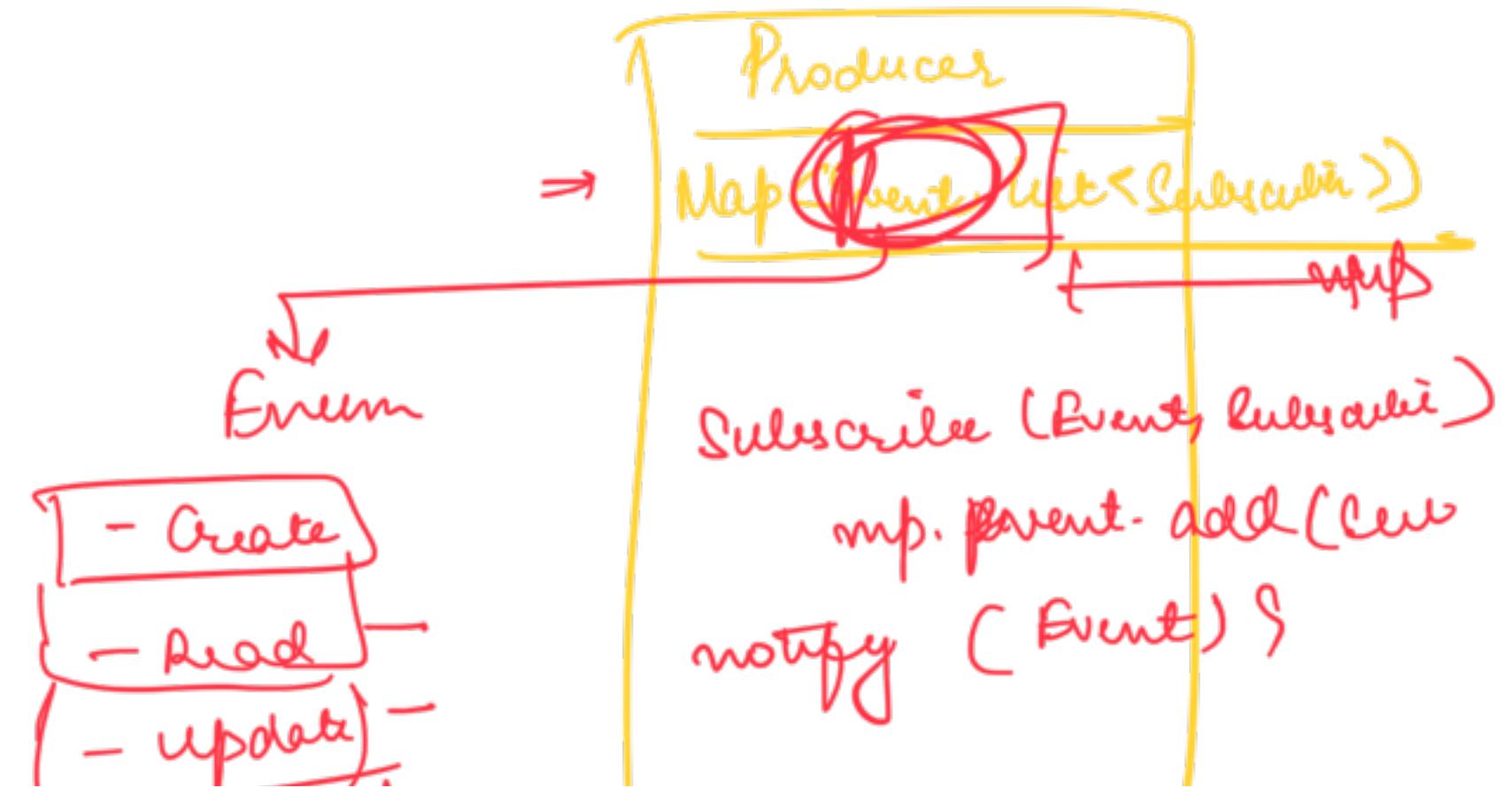




- ① Cache class shouldn't know what  
every eviction strategy needs
- ② Eviction strategy should get the data that  
is there with cache



V →



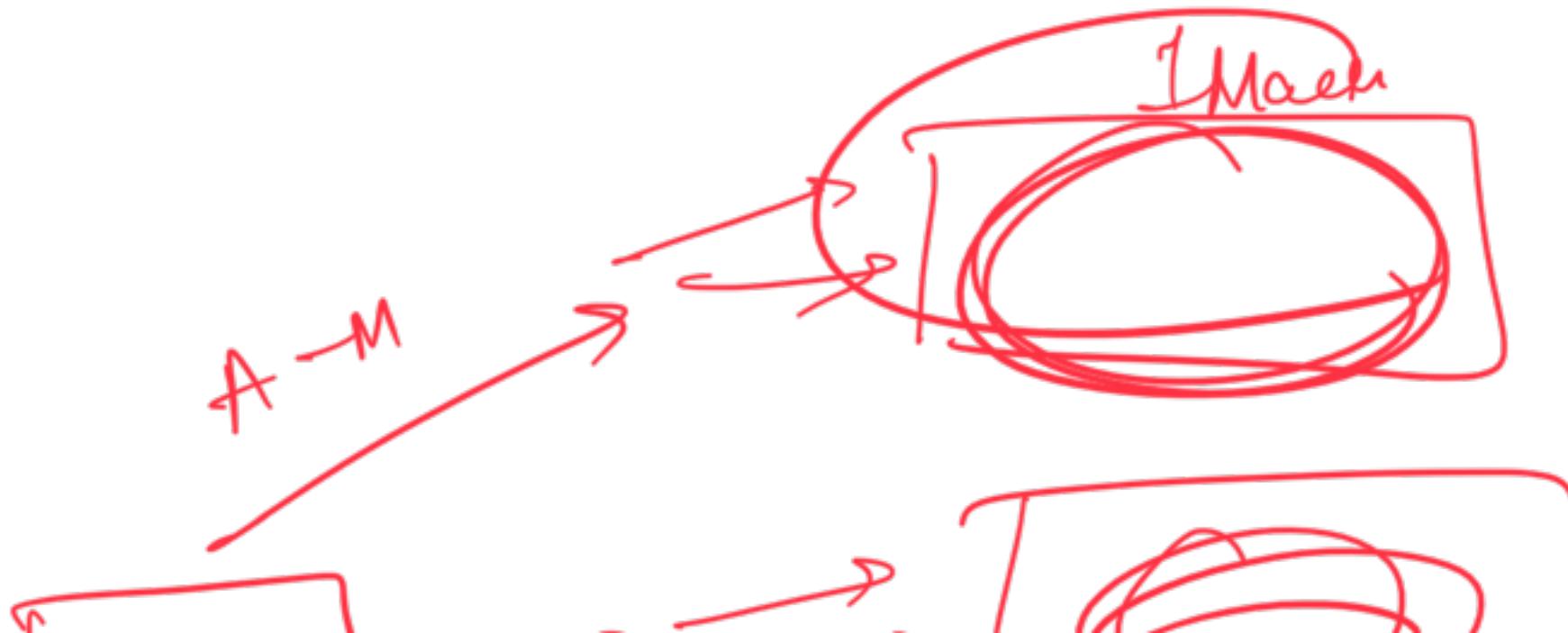
~~- Delete) -~~ ↗

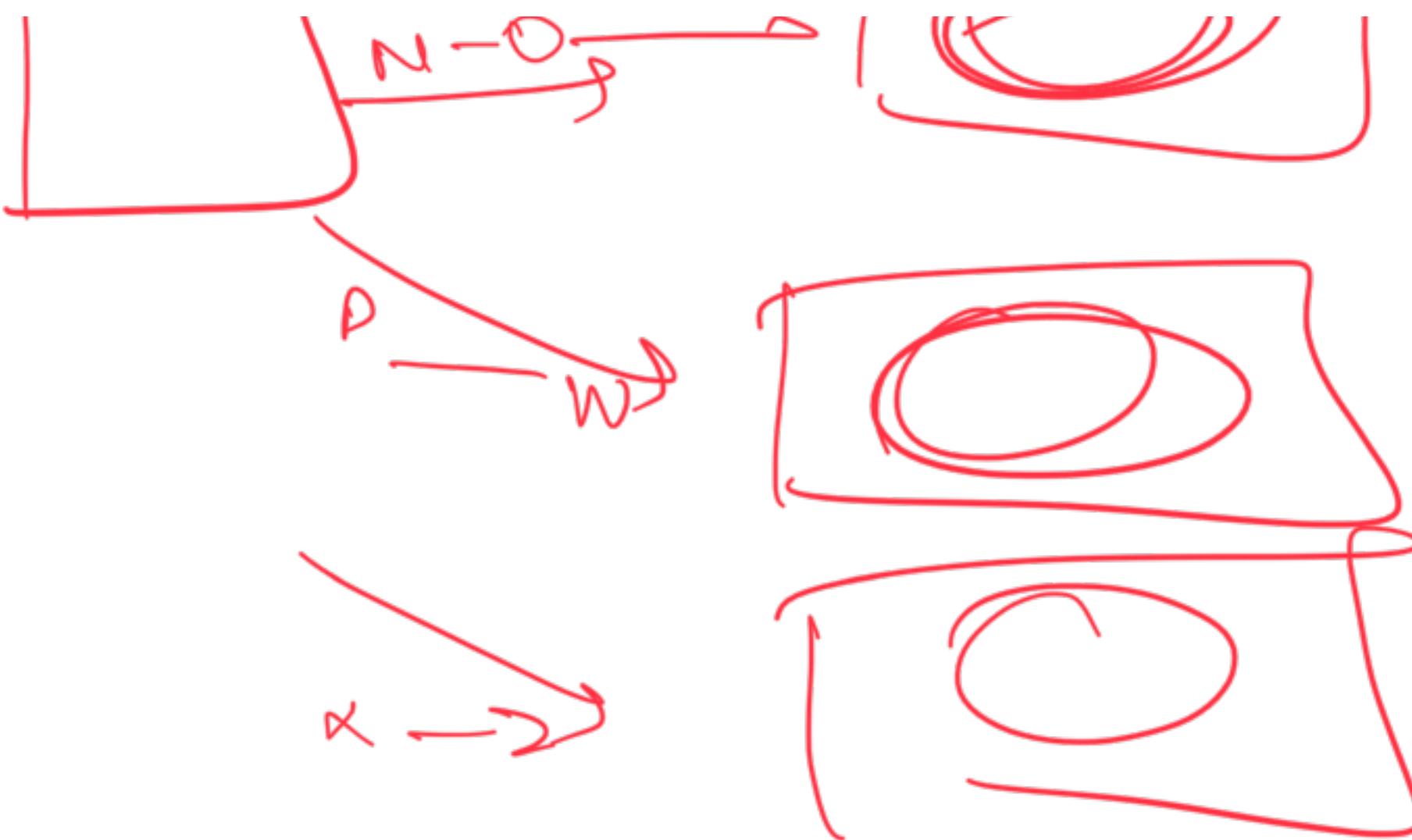
## Code

→ All the knowledge of conc classes

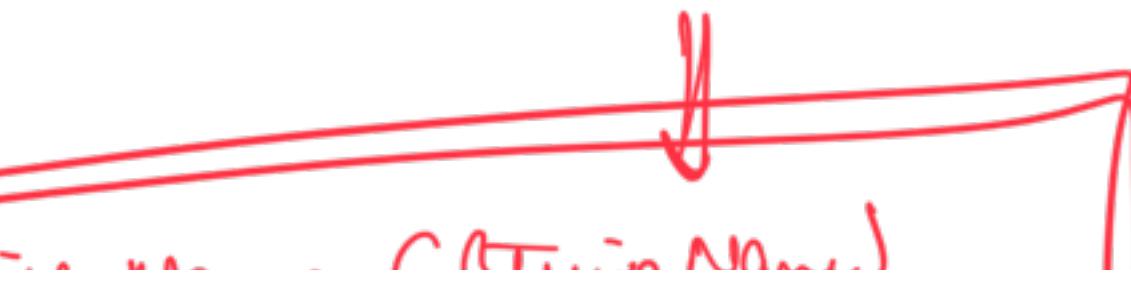
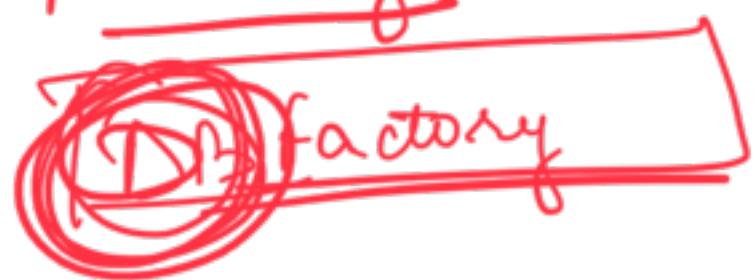
↳ Concurrent DT

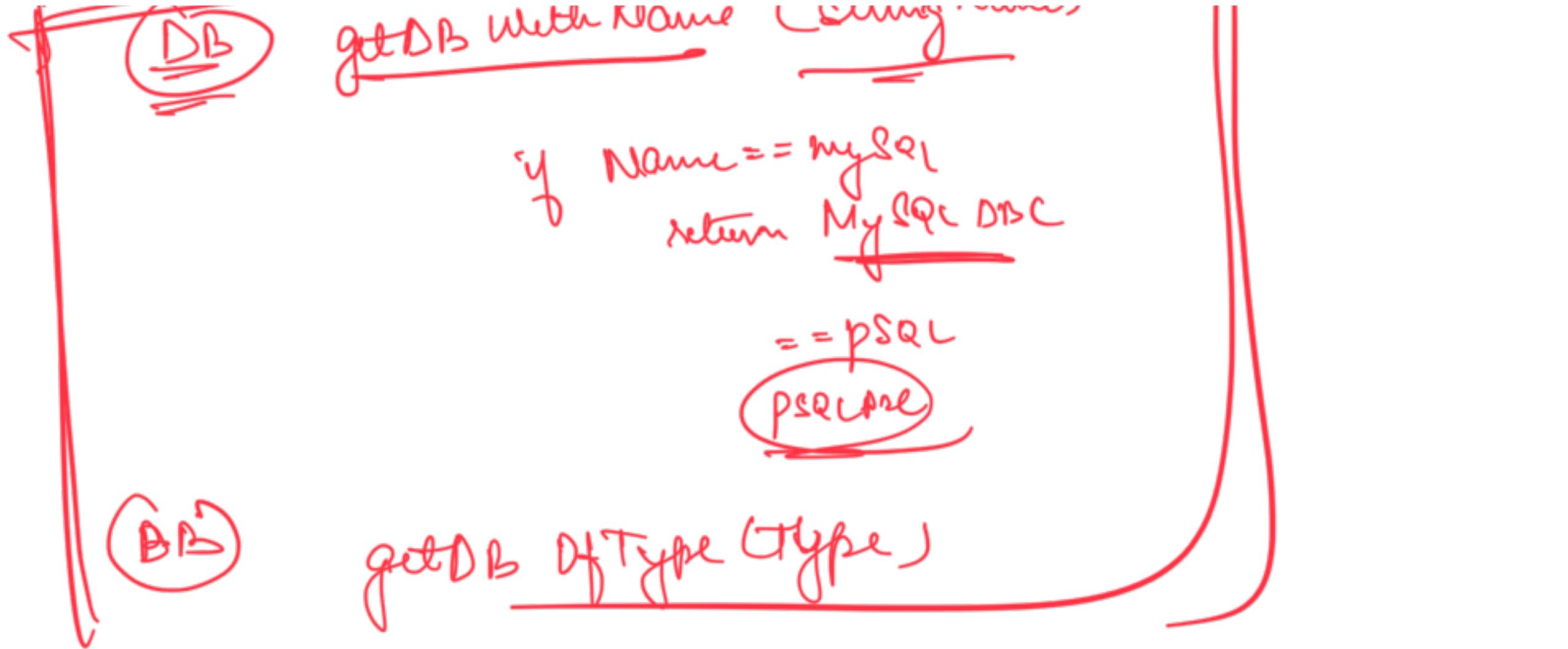
↳ futures





Factory





Abs Factory : related object

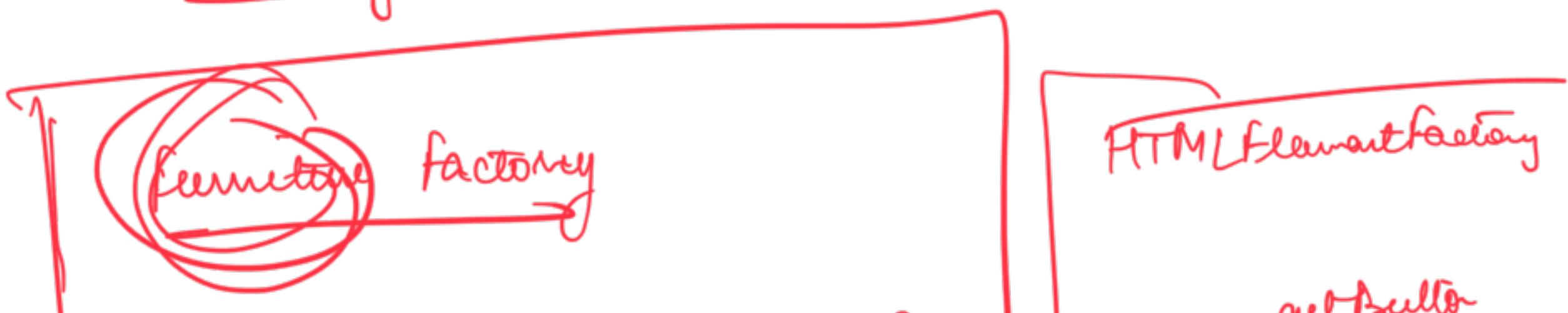


Table getTable()

)

Chain getChain()

)

getMany

getDropdown

Concurrency 3

→ Concurrent  
DS

→ Future

Tom

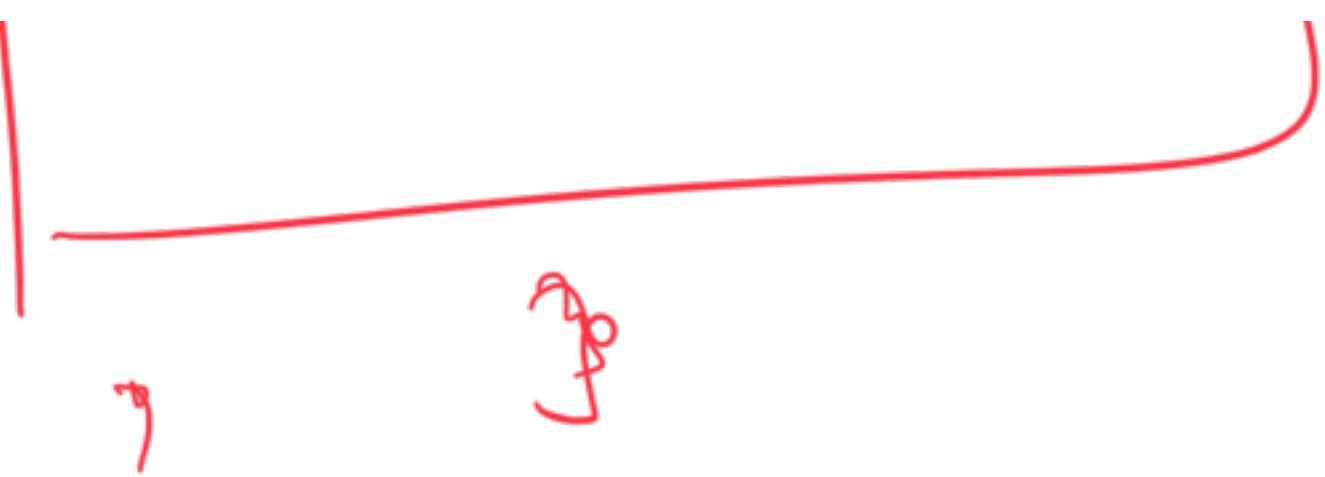
Recordings of 2 SDP

Dec  
Adapter  
Facade  
Flyweight

Code Search Controller {  
  
~~Locate~~ Search() }

+  
Search

Search Controller



Search Service

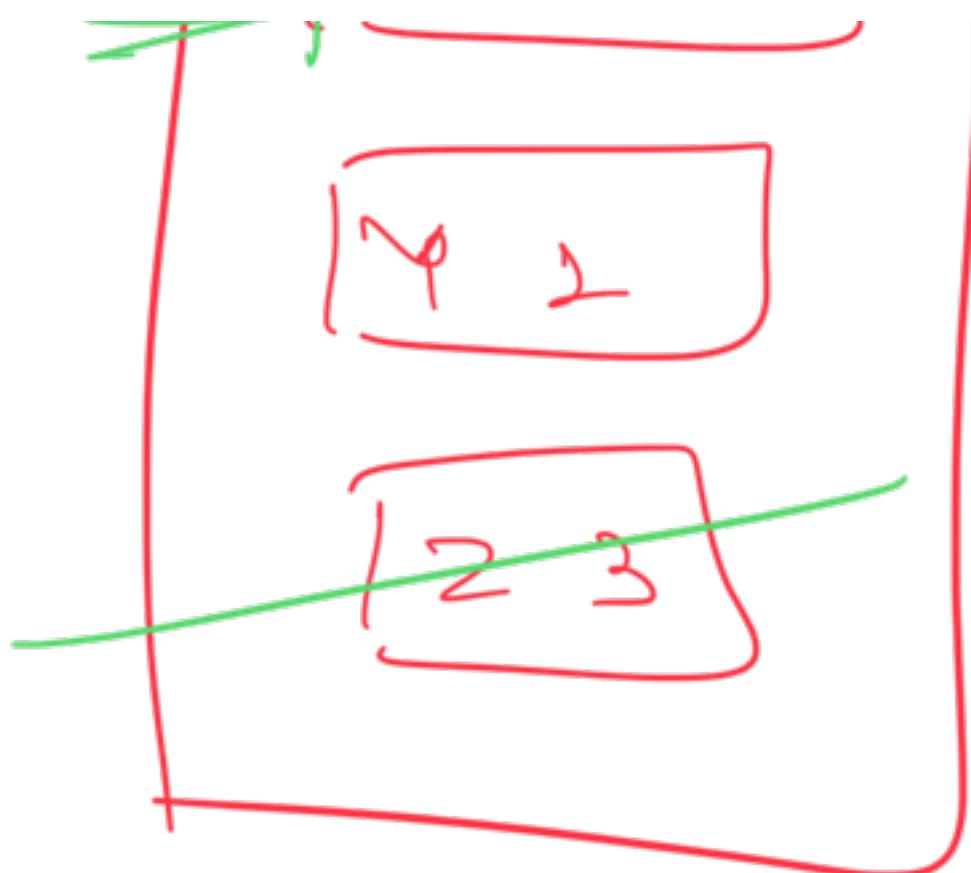
}

Request Pooling



X





① Check at time of access if the key  
is expired

② Have a parallel thread which periodically removes expired keys.