

# CACHING

- 1] Caching in algos → To avoid doing some ops especially computationally complex ops. that take a lotta time multiply times  
∴ caching used to improve time complexity (T.C.)
- 2] Caching in sys. → Used to reduce (improve) the latency of a sys. (speeds up the sys.)
- 3] CACHING → Storing data in a loc<sup>n</sup> that's diff from the original data loc<sup>n</sup> s.t. it's faster to access this data from memory.  
It's a way ~~to~~ to design a sys. s.t. if we were originally using ~~some~~ ops. or data transfers that take a lotta time (e.g.: Network requests) the sys. is designed is s.t. we do diff types of ops. or data transfers that are faster.

## I

### Caching in diff places in a sys.

e.g:-



Regular op. → client makes net. req. to server to get some data from dB and then data transferred from dB → server → client

1) Caching at client level - Client caches a data value so it no longer needs to go to server to retrieve it

2) Caching at server level - You want client to always interact with server but maybe the server doesn't need to go to dB for data retrieval each time

- Go the dB once → get data → cache in server → transfer data to client as per net. req.

3) Cache in b/w components - eg:- cache in b/w server and dB

not in our control when we des. sys. → 4) Caching at hardware level - eg:- CPU cache - make it faster to retrieve data from memory.

5) Caching occurs by default, at many other such kinda b/w. of sys.

## II

### Advantages / Utility of Caching

Caching speeds up sys. → 1) Avoid large magnitudes of net. reqs. → If the client wants to access the same data again and again, no new net req. (client → server → dB → server → client) done each time as we cache reqd. data so client doesn't need to go to dB each time.

Caching done at client or server lvl.

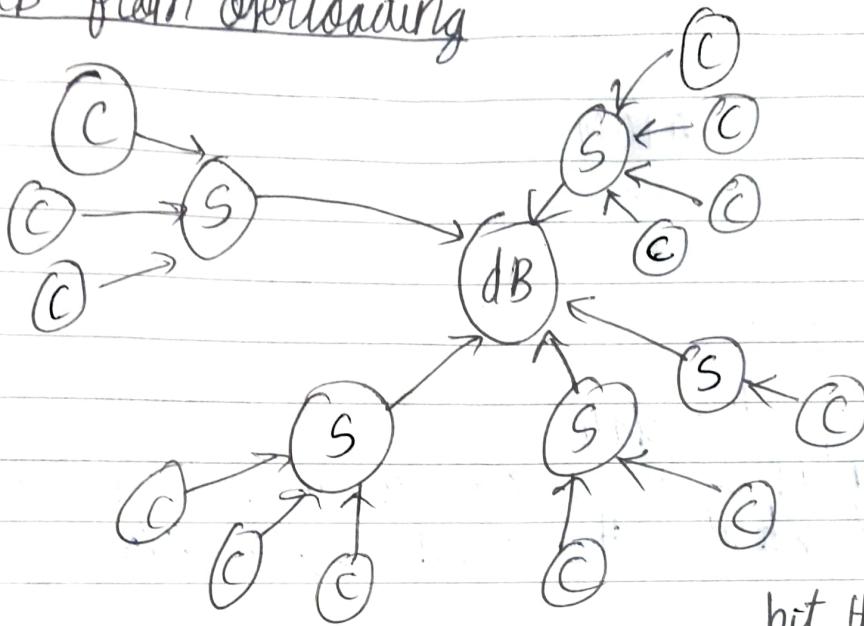
Caching speeds up sys. → 2) Avoid redoing a computationally, long, bad T.C., complex operatn each time  
eg:- say on net. req. from client → server

everytime, this time - consuming algo runs.  
To avoid performin this op. multiple times,  
caching is done.

### 3) Prev. dB from overloading

dB  
optimized

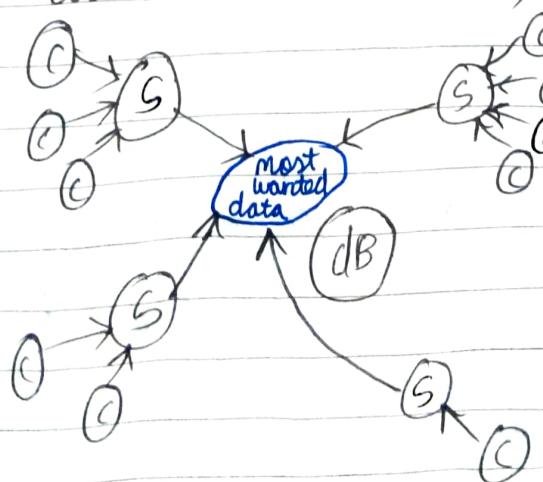
e.g:-



hit the dB

Multip clients across multip server, with the same net. req. e.g:- Each client makes <sup>net.</sup> req. to server to get Instagram profile of a popular celeb.

Indir. req. are fast, so speed is not the concern.  
We ~~use~~ use caching to prev. same data being read from dB 10<sup>6</sup> times (or somethin like that)  
∴ a cache is created, outside dB, for that data



If we don't cache,  
too many req.s  
to dB for this data +  
in gen. other data  
accumulate and  
overload dB

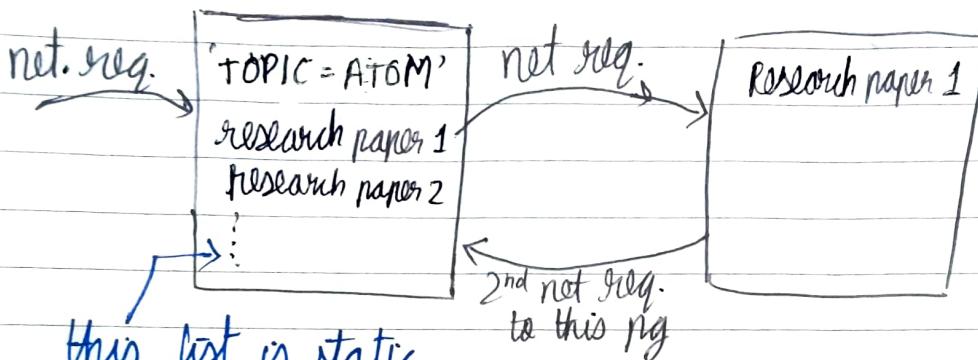
$k-v$  = key value

Date: \_\_\_\_\_  
M T W T F S S

III

## CACHING in ac<sup>n</sup>

- 1) Reasons state components <sup>immutable</sup> <sup>or content</sup> in a website  
 e.g. - ResearchGate : every once a yr papers published on site



this list is static content so if we cache it after 1<sup>st</sup> net req. to this pg of website, the 2<sup>nd</sup> net req. to list pg will be a lot faster

### 2) Runnng code on algoswp. (AE)

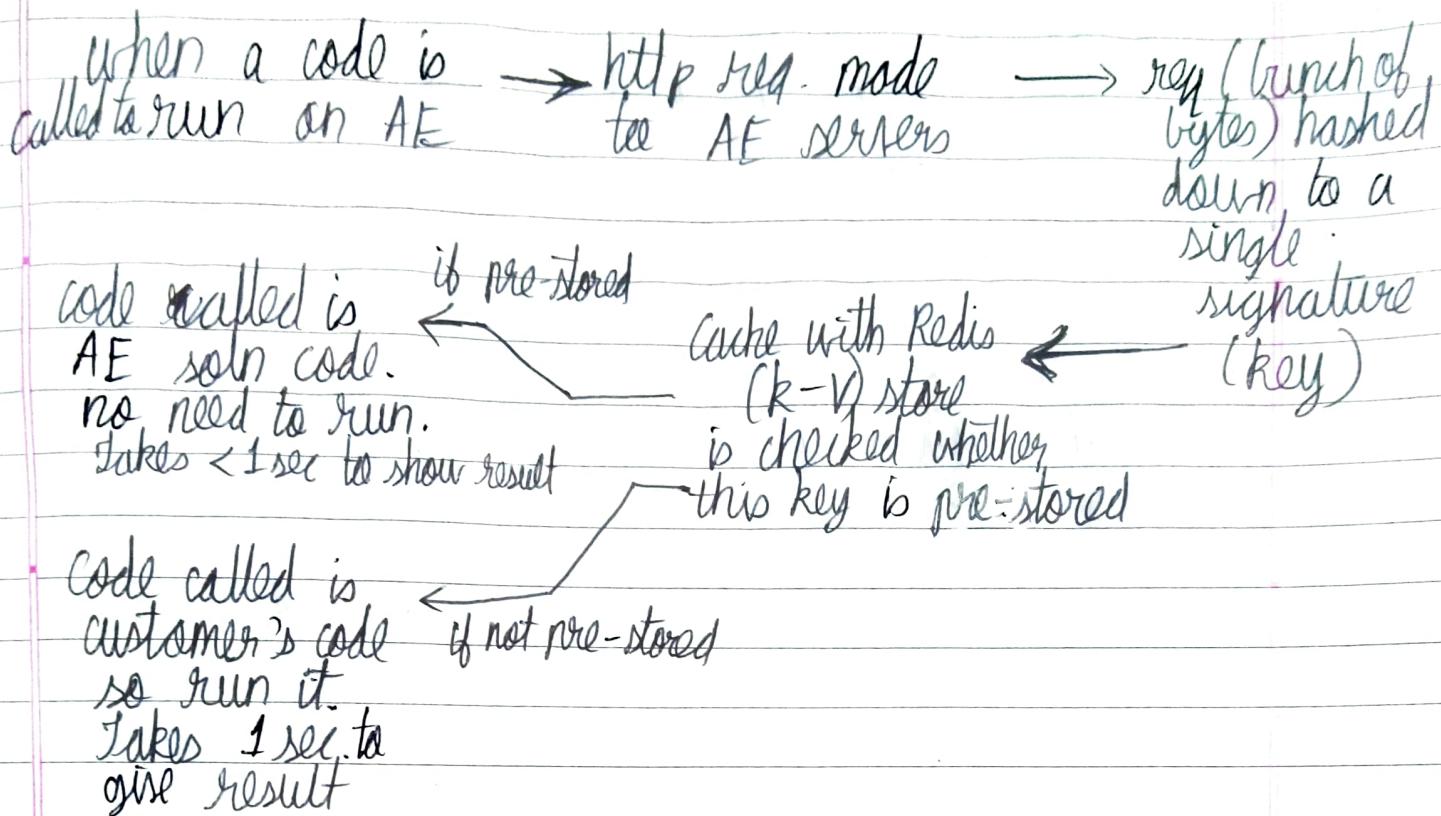
- i) Runnng, testin code takes  $\approx 1$  sec.
- ii) for customer code  $\rightarrow$  runnng, testin code reqd.  
 for AE soln code  $\rightarrow$  w.r.t this code is correct code i.e. it'll pass all test cases so no need to run, test

Cache stored  $\xrightarrow{\text{server level}}$  OR  $\xrightarrow{\text{detached from server is an indep. comp.}}$

Redis (popular in-memory dB)  $\Rightarrow$  a k-v store used for caching.

All AE solns. codes pre-stored with a unique 'key' ~~in~~ in this Redis cache

## Caching in ac<sup>n</sup>



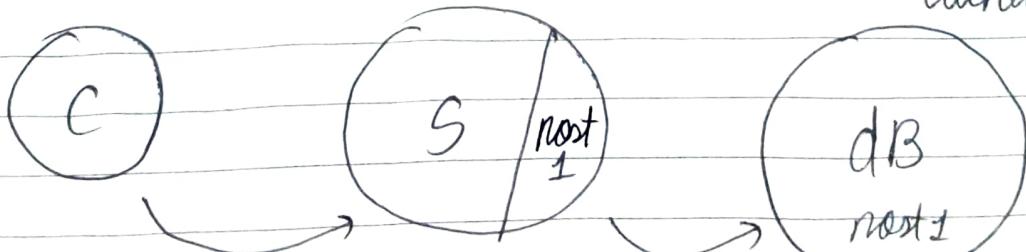
Until now we discussed caching for a sys. that needs to read data.

IV Caching and Sys. where users can read, write, edit data.

e.g.: - Readin, writin, editin posts of FB or linkedIN

1) When you write a post to FB, say that using not reqs. that post is stored in dB

→ cached in server



2 sources of truth for given post → cache in server → dB

2) When you now wanna update your post, depending on type of caching in server post stored in the

```
js server.js > app.get('/nocache/index.html') callback > database.get("index.html")
1  const database = require('./database');
2  const express = require('express');
3
4  const app = express();
5  const cache = {};
6
7  app.get('/nocache/index.html', (req, res) => {
8    database.get('index.html', page => {
9      res.send(page);
10   });
11 });
12
13 app.get('/withcache/index.html', (req, res) => {
14   if ('index.html' in cache) {
15     res.send(cache['index.html']);
16     return;
17   }
18
19   database.get('index.html', page => {
20     cache['index.html'] = page;
21     res.send(page);
22   });
23 });
24
25 app.listen(3001, function() {
26   console.log('Listening on port 3001!');
27 });
```

```
js database.js > ...
1  const database = {
2    ['index.html']: '<html>Hello World!</html>',
3  };
4
5  module.exports.get = (key, callback) => {
6    setTimeout(() => {
7      callback(database[key]);
8    }, 3000);
9};
```

DB → updated instantly  
DB → updated async. ly.

Date: \_\_\_\_\_  
M T W T H S S

### Write-thru cache

- 1) dB and cache both updated instantly when you edit post

### Write-back cache

- 1<sup>st</sup> cache is updated instantly  
Then dB is updated async.

→ updated after a period of time e.g. 5 hrs, 30 days etc.

→ once cache is full and you gotta evict stuff outta cache

- 2) Takes time as you still gotta go to dB and come back for an edit

Takes less time as only cache update at the instant of edit

- 3) Edit stored in 2 sources of truth at the instant of edit

Edit stored only in cache at instant of edit  
∴ if somethin happens (e.g.: delete) to our cache before dB is async. ly updated, we lose <sup>the</sup> edit.

V

## Staleness

eg:- YouTube comments see?

You don't wanna update data asynch as it might happen that:

- someone posts a comment (SERVER 1)
- other person reads it from diff. server (SERVER 2)
- 1st person edits comment (but dB not updated) and his server's cache updated asynch yet
- other person still sees original comment (unedited) and replies to unedited comment (STALE DATA)  
as SERVER 1's cache ↓  
by SERVER 2

Cache in server 2 is STALE : it hasn't been updated properly

Soln: Move cache outside servers into a common comp. for them all outside dB  
∴ Each time cache is updated from server 1, server 2 accesses this latest, updated cached data

NOTE: While sys. design just see if that part of sys. is ok with staleness or not, and accordingly choose how you gonna cache the data

eg:- YouTube → view count → lit bit of staleness on video → is fine  
comment sec → even a lit bit of staleness is NOT fine

VI

Date: \_\_\_\_\_  
M T W T F S S

NOTE:

## WHEN TO USE CACHING?

- ✓ 1) Static (immutable) data → Caching is easy to implement, fast while using  
eg:- Research papers website
- ✗ 2) Dynamic (mutable) data → Caching is tricky to implement as data stored at 2 diff loca's and you gotta make sure data in both loca's is in sync.
- eg:- <sup>(FB)</sup> Facebook wall → not cached
- 1<sup>st</sup> net req. → my FB wall → friend's profile  
net req.  
2<sup>nd</sup> net req.
- On Facebook website, the data you see on your wall after 1<sup>st</sup> and 2<sup>nd</sup> net req. to FB wall page is diff.
- ✓ You're doing only 1 of reading or writing data  
single thing eg: either
- ✗ You're doing multiple things → reading, writing, editing data
- ✓ You don't care abt staleness of data or if you can invalidate (get rid of) stale data in caches in a distributed manner when dealing with a distributed sys.

## VII

# CACHING - EVICTION POLICIES

We might wanna evict (delete) data from cache for several reasons. Eg : data is stale, cache's storage is full.

To evict data from cache, we follow certain rules a.k.a "evic" policies

- 1) LRU policy → Least Recently Used (LRU)  
pieces of data in a cache are removed if ~~you~~ you have some way of tracking what pieces of data are LRU. cuz we ~~care~~ care least abt em.
- 2) LFU policy → Least Freq. ly Used (LFU) pieces of data evicted
- 3) LIFO or FIFO → Data ejected on a Last In First Out or First In First Out basis
- 4) other ways.

Choice of evic" policy depends on kind of product you wanna do sys. des. for