# High Level Design

neighborhood DNS

spread out across the world

Latency ↓↓↓L

CLIENT

Browser
WebApp
Mobile App

Domain Name server

↑ I.P

OS

TLD
COM
.us
.in

$T_1 + T_3$

$T_1$

A

$T_3$

$T_2$

Latency (A)
$= \Sigma T_i$

$T_1 + T_3$

DNS

FB

ICANN

Client

DNS

IP

ISPs
Google
FB
US Army

FB.com
IP'
hdf.com.

newFB.com

Client

DNS

Waterfall

IP₁ IP₂ IP₃ IP₄ IP₅

IP

shuffle

IP₄
IP₂
IP₁
IP₅

System

Proxying
Proxy Servers

IP

GATEWAY ≡ Load Balancer

Horizontal Scaly

APP
Servers

I        II

Vertical

→ much easier
to get some
gainers

Horizontal

→ flexibility
with changing load

→ cost effective

CAP theorem



Client

DNS

Gateway
Proxy

( Load  Balancer ) ( Δ ~

App
Server

Computational
magic

Storage
Layer

cold
storage

internet                 client        open internet

R|509    L      B    ≈   dc.crypt

↓           ↑           i

→ prevent load on down boxes
intranet → rejoin new boxes
→ balance out

— X.509 —

**Layer 4** ≡ Transport layer

addition security

Source IP
Source Port
Dest IP
Dest Port

fast simple

---

**Layer 7 LB** ≡ Application

use if ≡ 707
usid = 107

much better flexibility

Slower

---

**Stateless**          **Stateful**

17+15    |  4+5
         |  a+b

LB 🙂           Calculate

stateless

App Serv

A₁     A₂     A₃

---

↓ Stushom        Chatbot

*** 

LB 🙂

stateful

App Serv

Hi       hello

pr+t

data

$Q_1$

A1

$Q_2$

(A1)

---

## Load Balancing ✓

(stateles)

→ Round Robin

R R



```
LB
```

1 6    2 7    3 8    4 9    5 10    6

---

```
→        LB =
```



---

### WRR
→ (Least Connection first)

→ WLR ALR

## Load Balancy in Statfull

$U = 10^7$    ℓ

2 B    uid bid      LB        Mapping

Sol ① Maintain Map (BAD)

map-size = order of key-size

Load Balancing func will slow down

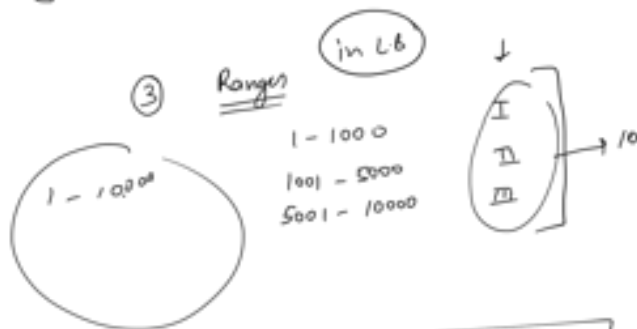Sol ② Hash

BAD!!!

10 boxes ≡ 10 App Serv

$uid \% 10 = \Box$ 0-9

Apush = 102

$uid / 5$     LB → 2

0-9

0    1    2    3    4    5

⊢ S

→ gaurav of stickiness using hashing

→ LB ☺

**Downside**

On server inc/dec
you will have          %
to move state of all boxes
internally

| | | |
|---|---|---|
| 1 %.5 ⟶ 1 | %.6 | 1 |
| | | 2 |
| 2 ⟶ 2 | | |
| | | 3 |
| 3 ⟶ 3 | | |
| 4 ⟶ 4 | | 4 |
| | | 5 |
| 5 ⟶ 5 | | |
| | | 6 |
| 6 ⟶ ℓ 1 | | |
| 7 ⟶ 2 | | 1 |
| 8 ⟶ 3 | | 2 |
| 9 ⟶ 4 | | 3 |

③ Ranges          (in LB)          ↓

1 - 1000

1 - 10000          1001 - 5000          I
                   5001 - 10000         II  ⟶ 10
                                        III

Bottleneck: Add/Removal
of servers is        ☺ ☹
inefficient

━━━━━━━━━━━━━━━━━

④  Consistent Hashing

LB for Stateful Application

4 servers

$(\overline{H_R}(l_1) = v_1 = 10,000$

cont addition/removal of boxes

$[0 - 10^{18}]$

$S_2$

$10^{18}\ 0/2,5$

$R_1$

$10,000$

$S_1$

$H_S$ (server id)

$R_2$

$R_3$

piggyback on people q random

$S_4$

$H_S$ (s-id -1)

$R_4$

$S_3$

$H_S$ (s-id 2)

$0 - 10^{18}$

allmet

database $\equiv$ (Sharding) key

$ry\text{-}id = $ user-id

$ry\text{-}id = $ location id

$\left(H_R\right)(ry\text{-}id)$
$=$

$[0 - 10^{18}]$

S5        S1

S4

$S_2$

$S_3$

$R_3$  $R_2$   $S_1$      $R_1$  → Again

$R_{10}$

$S_4$

$S_2$

sgjar

Add($S_5$)

S5

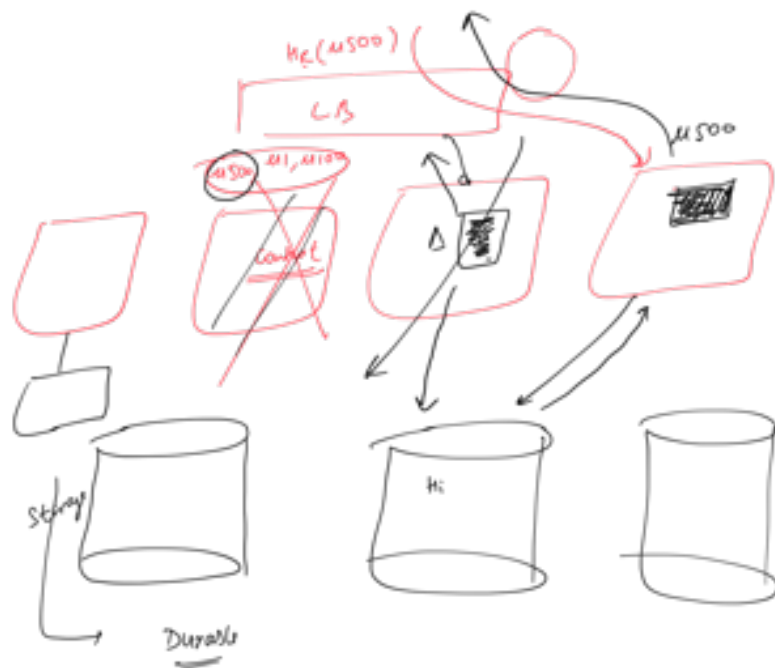H.(s-id 5)

Lokesh

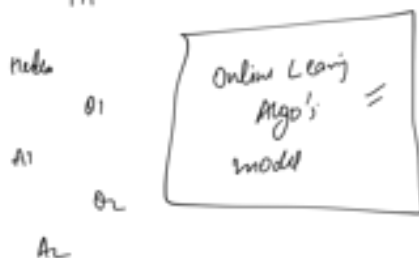✓ (I)  Amount of stale transfer that happens
        has been _reduced_ / a lot

😊 😊

        😊 😊  by a fact of # servers  😊

✓ (II)  downtime will also be low.....  small set of
✓ (III) you don't have downtime for every serv.  users get impacted
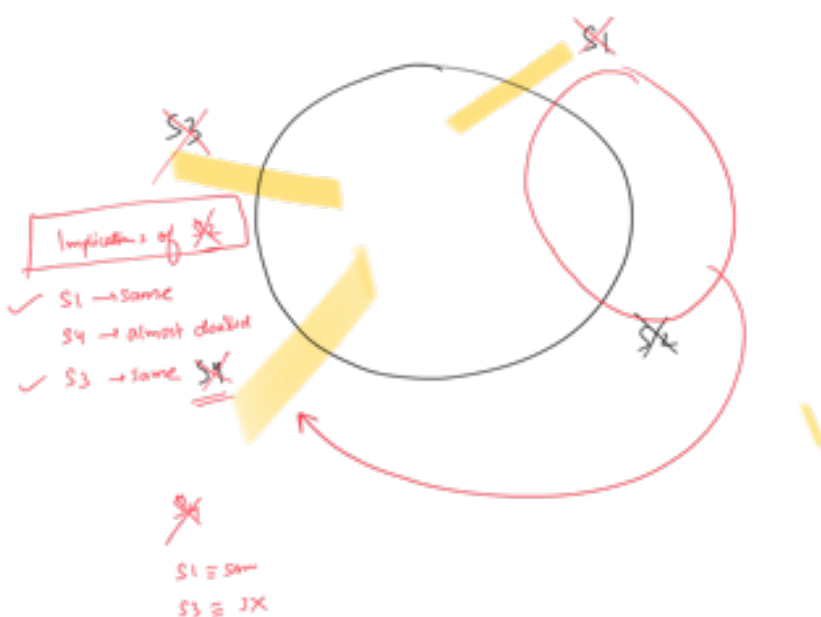
S1

S3    S4

HR(u500)

LB

u500 u1, u100

conflict

Δ

u500

Δ

strong

Hi

Durable

→ we are not losing the ultimate
source of truth

(A) cold start ≡ requests latency inc for some time

Hi

node

θ1

A1

θ2

A2

┌─────────────────┐
│ Online Learning │
│ Algo's     ≡   │
│ model          │
└─────────────────┘

(B) pre-compute copying

(I) ← Storage layer → new App Server

(II) ← M slave

┌──────────────────────────────────────┐
│ Solved problem                        │
│     impact of server addition or removal │
│  'is'  minimized                      │
└──────────────────────────────────────┘

Implication of S1

✓ S1 → same
S4 → almost doubled
✓ S3 → same

S1 ≡ same
S3 ≡ 2X

"  ─────────────→  "

# Cascading Failure



— evryth will go down

$H_R (r - id) \rightarrow [0, 10^{18}]$

$10^8$ 0

$S1$

$H_S^1 (s - id) \rightarrow [0, 10^{18}]$

$H_S^2 (s - id)$

$H_S^3 (s - id)$

$|servers| \times |H_S \; func|$
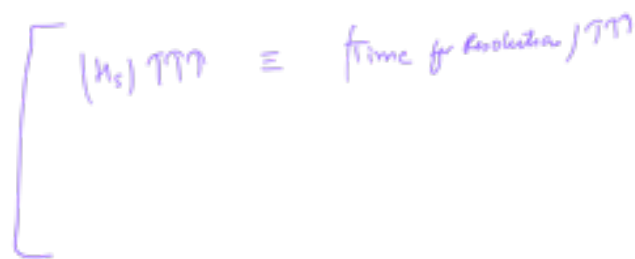
---

$S_1 \uparrow$

$SS \Uparrow$

$S4 \leftrightarrow$



Solved

Cascading Failures

( Now, the load will get divided almost equally.

$H_R$ (107)

S2

S3

S1

R1 ——— 100 720

S3

100, 720

S1

S2

S3

S2  S1

$$\Bigg[ |H_s| \uparrow\uparrow \equiv \text{Randomness} \uparrow\uparrow$$

$$\downarrow$$

Equitable Distribution $\uparrow\uparrow$

$$\Bigg[ |n_s| \uparrow\uparrow\uparrow \equiv (\text{Time for Resolution}) \uparrow\uparrow$$

A.SG



Client

$n_s^A$ $n_s^B$, $n_s^C$

$H_R$ $(r \cdot \text{id})$   Reverse Proxy

Load Balancer $\equiv$ GATEWAY

I          II          III          IV

10 servers

$H_s^A (1\text{-}id) \longrightarrow$

$\longrightarrow$

$[0, 10^{18}] \quad \longrightarrow$

$\longrightarrow$

$r\_id = 107 \quad H_R(107) = \boxed{6222} \quad \longrightarrow \qquad H_s^A$

| (100, 2) | (500, 1) | (5000, 3) | (10,000, 4) |

$\vdash N \longrightarrow$

$*$
$O(\log N)$
Map a $x_j$ to
a server
$*$

/

$H_s^A$

$H_s^B$

$h_s^C$

$\vdash \quad 4 \times 3 \longrightarrow$

| (560, 2) | (2000, 1) | (5700, 4) | (6000, 1) | (200,000, 3) | — —

$H_R(107) = 6220$

$\longrightarrow 3$

## High Level Design

OSI Layers

Internet ≡ network of networks



Ethernet

INTERNET

Intranet

GATEWAY

Outside

X42    II

P3

private IP    private IP    private IP

A    B    C    P R O X Y    IP

V    II    E    F    PROXY    VPN

PL    Reverse Proxy    → Load Balancer

Gateway

Proxy ≡ Gateway to the outside world

Reverse Proxy ≡ Gateway to the inside world

Cisco conn    II

domain

un-id
IP

netflow

netfix.com
Ikmf    UPNs    IP

India

HTTP

GET header = 

≡ Application Layer — origination of request happens ✓

- Encode your request
- Encryption
- Compression

Presentation Layer

7th
6th
https

tagged by a session-id =

Session Layer

5th

s-p    d p-n.o    d

Transport Layer

4th

3rd = | Network Layer |

OSI Layers ≡ Open Systems Interconnection

Intranet

Proxy

10.0.0(2)    MAC A

Laptop A

Router

MAC E

Laptop C

WiFi

M1    MAC M1

static

s-p MAC E

10.0.0(3)

M2

10.0.0(4)

M3

Laptop B    MAC B

hostname service

WiFi

Router

S-ip | ar.p

S_ip | d_ip

S+p | d ip

S | d

S | d

# 2nd    Data Link Layer

Source MAC Add | 0 | 0 | 0 | 1 | 1 - 1 - | dest MAC Address

ARP
Address Res Protocal

1st

## Physical Layer

d | 1 | 0 | 0

elec / light / Radio

M-id

Machine

7 — Application

6 — Presentation

5 — Session

Client

4 — Transport

3 — Network

2 — Data link

1 — Physical Layer

Physical Layer

devices

X    N

Data Link Layer

X N

Dest MAC No
=
Your MAC No

Network Layer

X 1

Transport layer

Session

App

Prnt.

App.





$s\_M = t\text{-}M = A_{MAC}$
$d\text{-}M$

I
Source → A
destination → B

"Same network"
→ ARP (B)
MAC B

II
Source → A
destination → B
"NOT in same"

RoutyTable (A)
→ M
nexhop (B)

ARP → M
MAC - M

s-P | in M   mpp A ad
of next hop

III
.hop — A
dest — B

A,B NOT in same network

s-M | d-M
M'

IV hap → A
ds → B
A, B
ARF
INV
TME
SAVYE
NGW

ARP

---

# (CAP) THEOREM

[ only gaurantee $\frac{2}{3}$ propertes simultaneously ]

- C ≡ Consistency

- A ≡ Availability

- P ≡ Partition Tolerance

**I) Consistency** → Immediate Consistency

Any read that you do, must return the value after the latest write

**II) Availability**

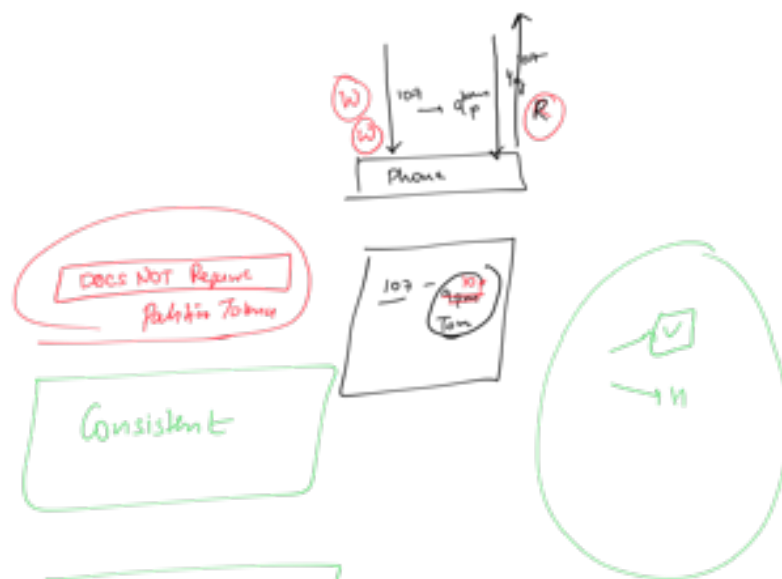Every and any request that goes to a non-dead server, must be responded to

**III) Partition Tolerance**

The system as a whole is allowed to have partitions.

OR    [ A ] / [ B ]—[ C ]

You have accepted the fact
that network partitions can happen at any
time
and the system as a whole will NOT
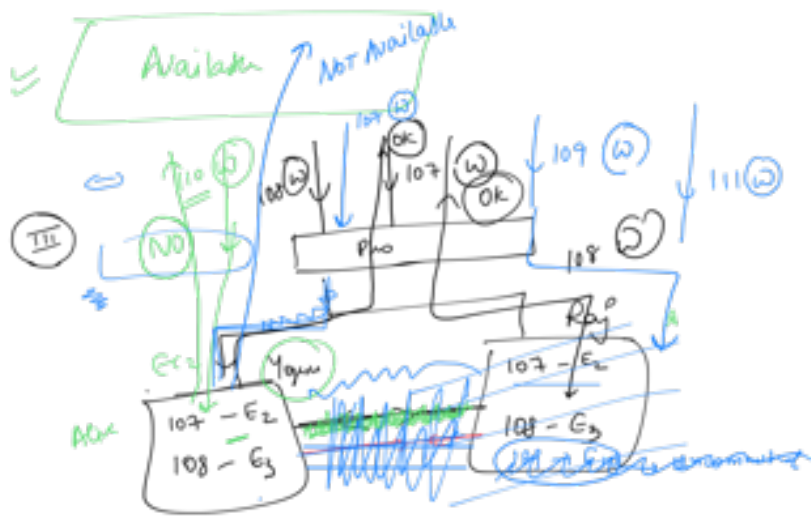stop functioning

---

Event Reminder Service



DOES NOT Require
Partition Tolerance

Consistent

Available

Position Tolerance

$$\frac{2}{?}$$

Available    Not Available

NO
TI

110    D
100 (w)    107 (w)    OK
OK    107
W
OK    109 (w)    111 (v)

Pro    108    D

Rdy    0    107 — E₂

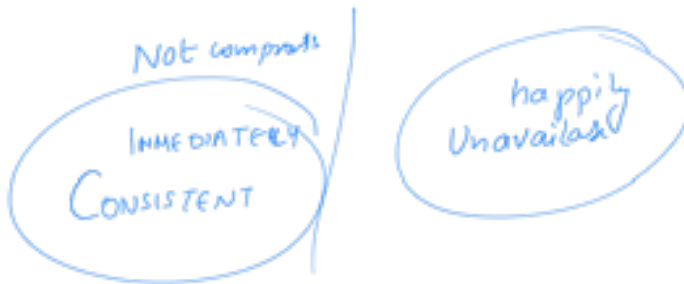E₂    Your    108 — E₃

Alive    107 — E₂
108 — E₃

CONSISTENT
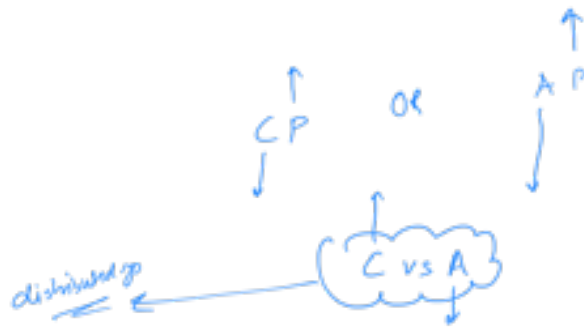
PARTITION TOLERANT

NOT AVAILABLE

CP    OR    AP

— CA — if you are doing
vertical scaling

L1    L2

Done    NOT AV

CP or AP

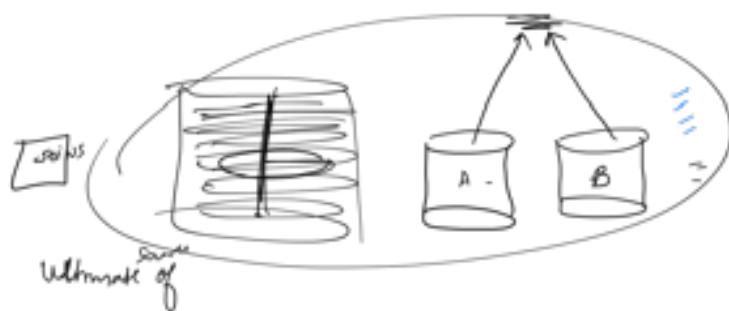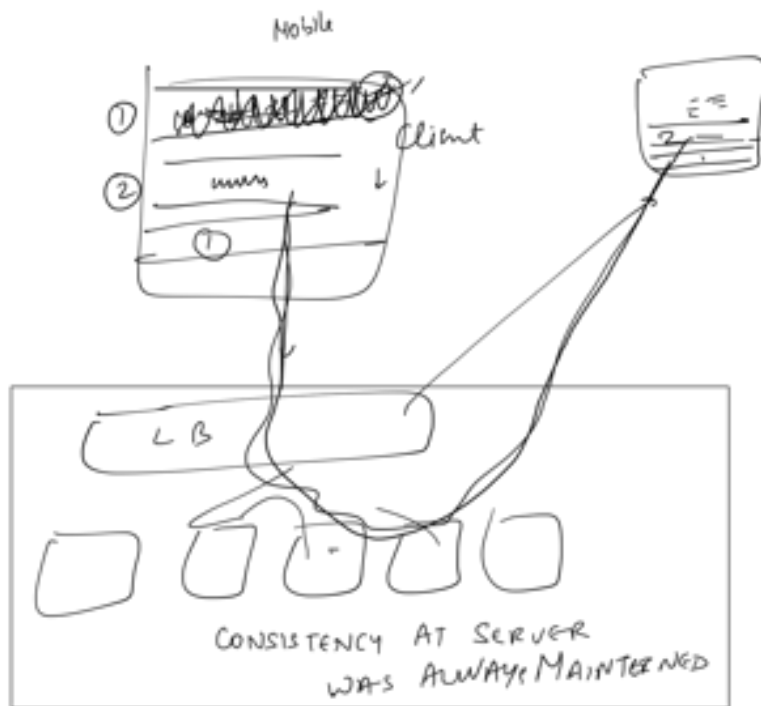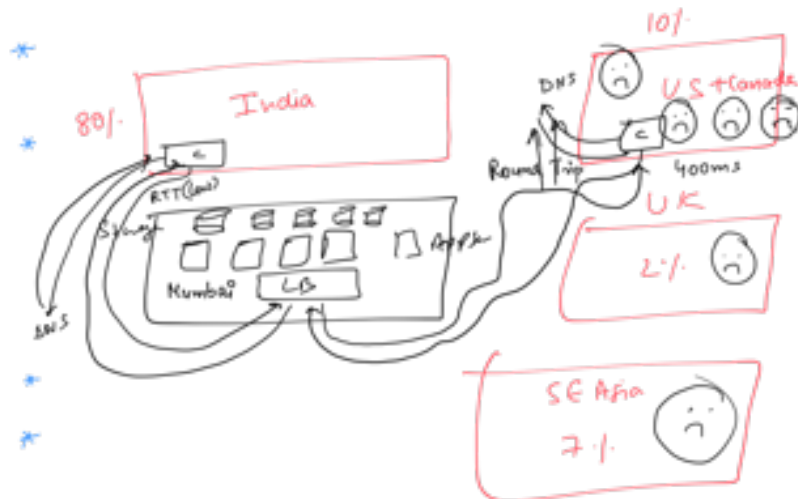C vs A

distributed sp

CA   X not real

Not compromis

IMMEDIATELY CONSISTENT

happily Unavailable

Eventually Consistent

120mm
60mins

A B    A B

Never confuse latency
with consistency

Mobile

① 

② ~~~

Client

E ...

L B

CONSISTENCY AT SERVER
WAS ALWAYS MAINTERNED

CACHING

Understand Caches

faster access of data
(Small)

Cache

Copy (Quick)
(Retrival)
① Memory   ② Ready To Use

Source

Ultimate of Source

① I
② II

Profile Info
PSP
Slow
Mid data

KB
KB
KB

II

Images
Videos
Another pages

MBS
GBS

lighter

heavier

dynamic

non sharable

static

sharable

India

80%

RTT(2ms)

Spray

Kumbai    LB    App

DNS

DNS

Round Trip

10%

US +Canada

400ms

UK

2%

SE Asia

7%

Content Delivery Network

CDN

Akamai

CDNs

Img = cdn-url

RTT

Client    cdn-url    cdn-url

DNS    cdn-url

RTT ↑↑↑

Load Balance

API

Storage

SQL    NoSQL    Obj

File
=
Obj

img1 ⟷ cdn-url

→ network is NOT getting checked

→ latency solver

Deals

x TBs

latency < 10ms
for 99% of the traffic

→ USA
→ Canada
→ Europe

for 90%
< 4ms

Company → CDN

CDN Hub → CDN Subhubs

User → CDN

cdn

{scheme net}

→

DNS

IP

cbcdn.net

DNS

Client

Focus

CDN mas

= $x_2$  IP
$url_1$
$url_2$

LB

$I_1$

$url_1$

CDN Distributed System

---

Cache (Invalidation) *

(A) → Time to Live    (TTL)
   = url should automatically
     become void
     after this bed

(B) → start sending a new URL

(C) Versioning
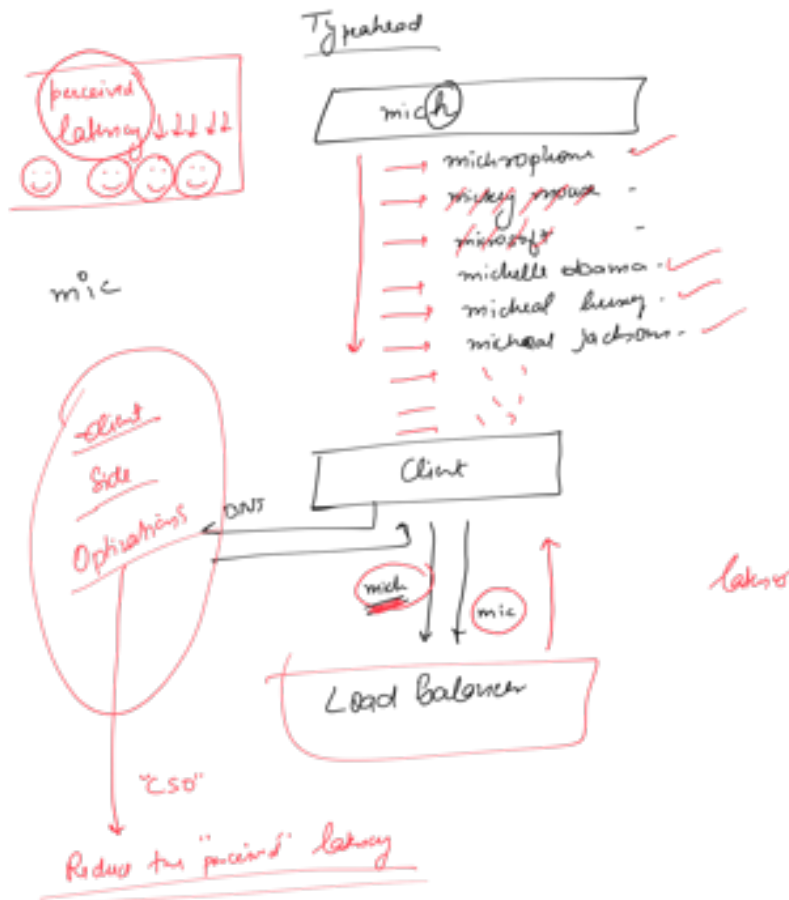
Privacy in CDNs
(1) CDN urls are by design not guessable
(2) encrypted format
(3) Authentication

② Client Side Caching

You can choose to stow some amount of "important" data on the client's side

Browser
Mail App
Desktop App

Typeahead



mic

perceived latency ↓ ↓↓ ↓↓

mic(h)
→ michrophone ✓
→ mickey mouse –
→ microsoft –
→ michelle obama ✓
→ micheal hussey ✓
→ michael jackson ✓

Client Side Optimisations

Client

DNS

mich    mic    lakso

Load balancer
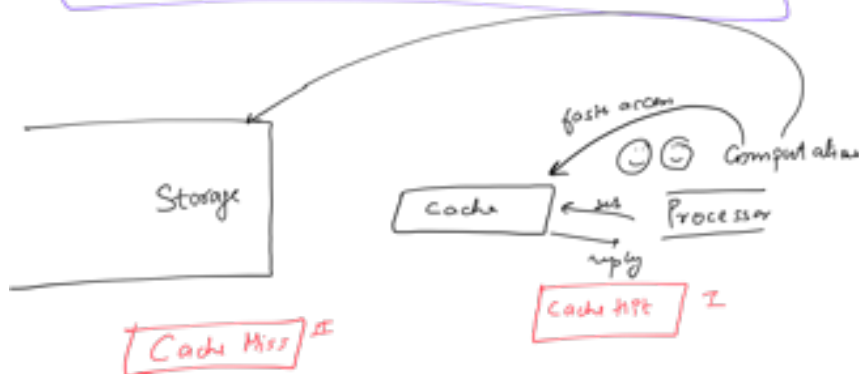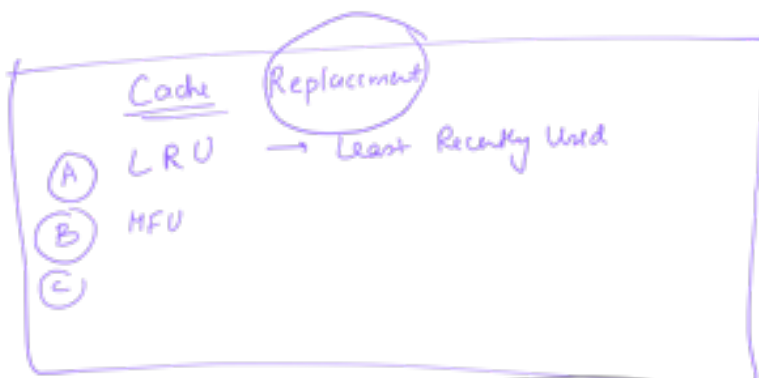
"CSO"

Reduce the "perceived" latency

Browser Optimisation

Automatically identify a few static resources for you and cache them

Hard Reload / Forced Relet

'ctrl + shift + R'

(III) App Server Cache

client Side Cache (II)

Client

DNS ← →    →  CDN (I)

lb?

Stateful          Load Balancers

App Server

107
[A, B, C]
110
2000
2701

Cache (Replacement)

(A)  LRU  → Least Recently Used

(B)  MFU

(C)

Storage

fast access

😊 😊  Computer alias

Cache  ← hit  Processor

reply

Cache Miss  II          Cache Hit  I

Access Time for Cache = X
      "      "   DB = 100X

k calls

Without Cache
          k × Avg DB Access Time

$$\text{Time} = \begin{pmatrix} k \times \text{avg} \\ + \\ k \times \text{Avg Network Access Time} \end{pmatrix}$$

80-20

**With Cache**

$$\begin{bmatrix} 70\% & \underline{C} & \text{Hit Rate} \\ 30\% & \text{Cache Miss Rate} \end{bmatrix}$$

$$\text{Time} = \left[\frac{70}{100} \times k\right] \times \text{Cache Access Time}$$

$$+$$

$$\left[\frac{30}{100} \times k\right] \times \text{Cache Access Time}$$

$$+$$

$$0.3k \times \text{DB Access Time} =$$

$$+$$

$$0.3k \quad \text{Ne A.Time} =$$

(A) How dangit are your incoming requests

(B) Cache Eviction Policy

(C) Amount of Cache Storage

---

④ Global Cache

Client

→ Replication
→ Re-lect

Cache!

LB

A B C

---

(Q1)

Client

pusher

Raj
Jay
Royble

Global
Cache

LB

Stateful

Pusher
C
Raj

B
Raj

Jay

B

---

~~~

LB

App
Server

---

Tinder

Cache Invalidation

Recent Message

Profile Information

Client Side Cache

CONS

Client

Android
iOS

set Profile info Profile
set Profile info

get Profile

Global Mc | Load Balancer | Storage

APP

---

1. Client side Cache = traffic light

Design a Facebook News Feed → Timeline

① Enlist all the features
   → Minimum Viable Product

② Estimate the scale

③ Trade Offs

④ Architecture Design
   → APIs
   → Data Storage
   → Data Flow

Linked In
Instagram
Quora
Twitter

---

Ⓘ  Ability to view posts of interest

Ⓘ  Like/ Comment /share a post

Ⓘ  Multimedia Support → Share Image/Videos

Ⓘ  Bookmark a post

(V) Making a friend / follow ———→ ancillary feature

(VI) Supporting friends

searches
ban ads
→ hashtags

(VII) Repost a post
(VIII) Support hashtags

(IX) Relevance ———→ Time Decay
                 ———→ Connection strength

(X) Ads

(XI) Automatic Refresh
(XII) (Write a post)
(XIII) Tagging users

---

(B) Estimation of scale

① #users = 2 B

② DAU = 400 M ———→ new feeds

③ #people who = 40 M
   produce content
   or
   made a post

④ # new posts /day = 2.5 × 40 M
                    = 100 M

COMPUTE LOAD

(A) writes/day = 100M

(B) Reads/day = [ ] → X mul fact
                                As fetch many
                                feed
             = 400M × 25
             = 10,000 M

10,000 M → 1,10,000        = 10,000 M
                                reads/day

← 30H →

$= \left( \underline{10\,B} \quad \text{xrows} \right)$

Read **OPS**

$$= \frac{10\,\cancel{BS}\,10^9 \quad 10^9 \times 7}{\cancel{7}\,4 \times 60 \times 60}$$

$= 115\,000$ OPS

$= \boxed{115k\ OPS}$

Write **OPS**

$$\frac{100M}{86400} = \boxed{1000\ \text{writes /sec}}$$

---



$2 \times avg$

$avg = 115k$

Read OPS $= 115k \times 2$

$= \boxed{230k}$

Write OPS $= 2000$

How OPS one machine can support!

|box? or /server



8 queries/sec    50 – 1000    1M/s

— hardware
— computation
— data fetch

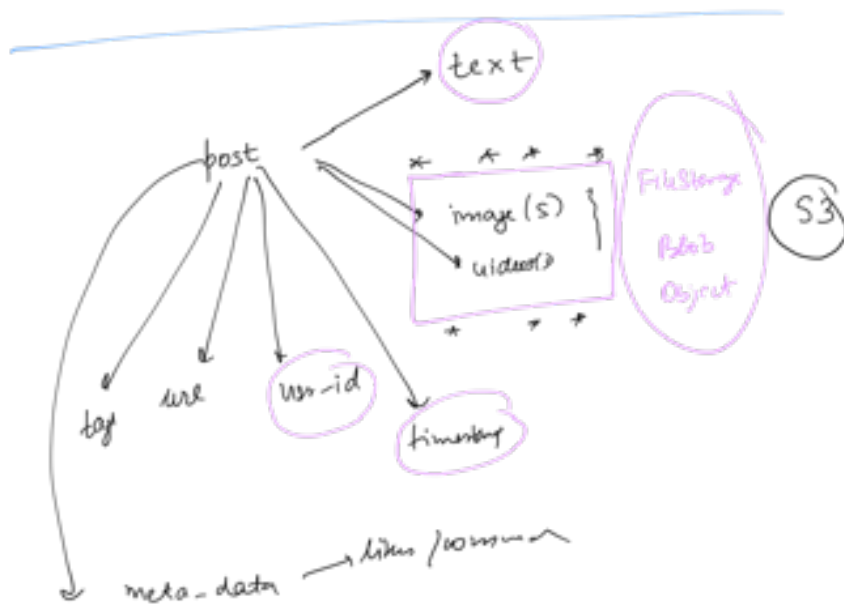$\longrightarrow$ | 100 OPS/box |

1000

$$\frac{2\,30\,k}{100} = \#\ boxes$$

2 300 boxes

Load Factor $\Rightarrow$ 65% - 75%
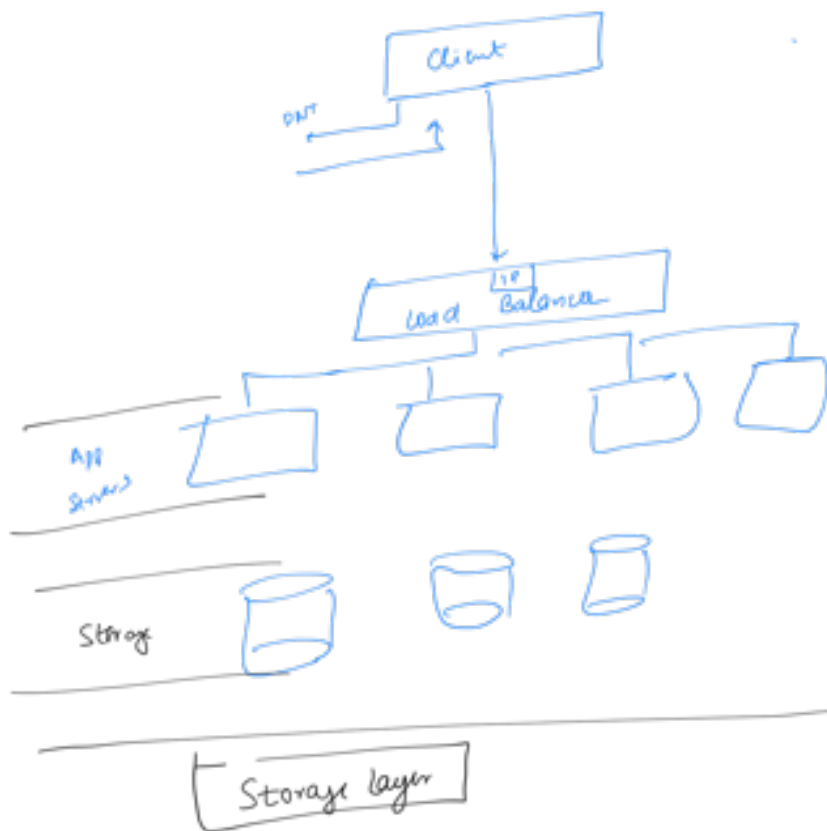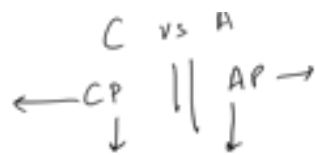
$$\frac{2\,300}{0.7} = \boxed{\sim 4000}$$
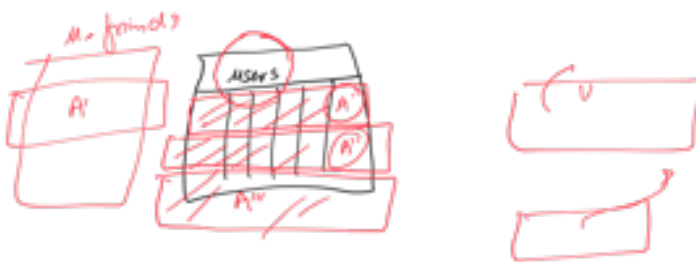
EC2 $\equiv$ App Server

---

text

post

image(s)
videos

FileStorage
Blob
Object

S3

tag    url    user-id

timestamp

meta-data $\longrightarrow$ likes /comment

---

(III) Trade Off.

CAP Theorem

C vs A

← CP || AP →

Available + Partition Tolerant ☆ ☆

Stateful OR Stateless

Client

DNS

Load Balancer [IP]

App Servers

Storage

Storage layer

(2TB) MySQL

Sharding

M. friends

users

A'

A''

A'''

U

users

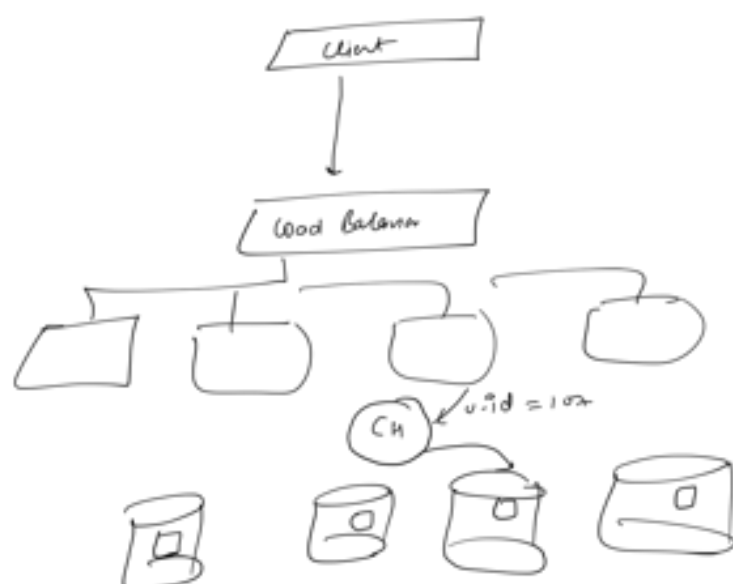| u id | f-nam | l-nam | age | score I | score II | | |
|------|-------|-------|-----|---------|----------|--|--|
|      |       |       |     |         |          |  |  |
|      |       |       |     |         |          |  |  |

normalization

(I) Consistent Hashing

key

CH

(II) Range band shard

$[90,000 - 98,000]$    $[98,000 - 110,000]$

(A - D)    [G - H]

Ⅶ    Hot    Shards    ≡

Client

Load Balancer

u.id = 10a

CH

Ⅷ    Time Based Shards

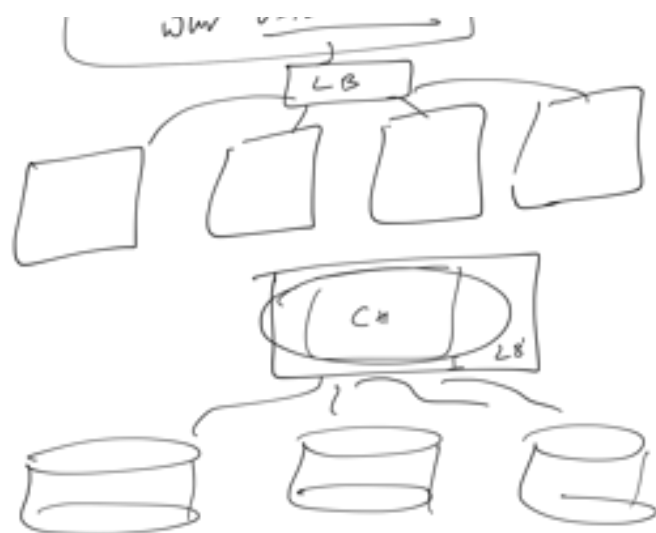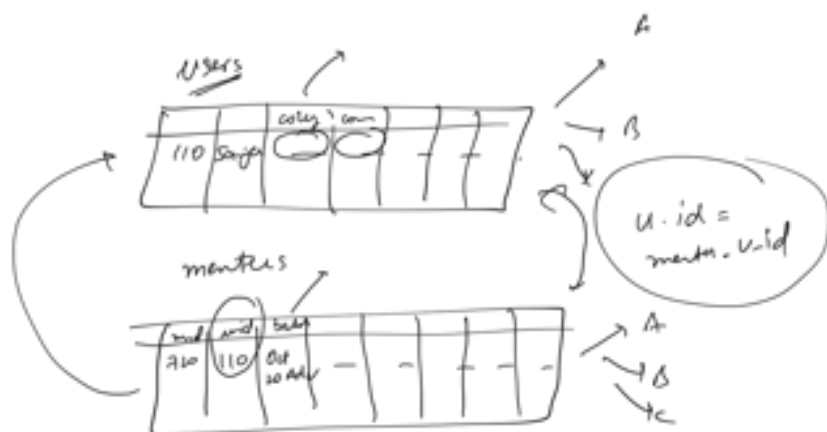30-60    0-2    2day - 30 days

4-10p    12-4    10p - 11:59a

Sat

M-T          W-t





Single box | Few Boxes ≡ Optimal ∦
                         ↓ uur

Evy shard / most of the ≡ Un optimal
               shards        Ouey

v-id = 107

users

via
  I

  II

→ timestamp

users—friend

| u-id | f-id | timestamp |
|------|------|-----------|
| 107  | 20   | (   )     |
| 107  | 22   | (   )     |
| 107  | 110  | (   )     |

select *
from us-friends
    u-id = 107          X N

WWW

LB

C#

LB'

JOINS

Users

| 110 | Sanjay | col<sub>x</sub> | col<sub>4</sub> | | | |
|---|---|---|---|---|---|---|

A
B

$$u.id = mentor.u\_id$$

mentors

| 710 | 110 | Out 10 AM | | | | |
|---|---|---|---|---|---|---|

A
B
C

Intra Shard → JOIN

Inter Shard JOIN

Client

LS

LB

A    B'

A''
B''

B'    B

□ X ▱

Inter Shard JOINs        Avoid

10>

10>%
users → (u_id)

10>

menke → user_id

(I) denormalized way

Ⅱ ( duplicated data ) = 😊 😊 😊

Ranguli ⟶ [ 120 ]        > U-id

USN_submisst
u-id        problem-id
107        11110007        10th Mnue

Recent_Submissontalk

| u-id | p-id | t- |
|------|------|----|
| —    |      |    |
| —    |      |    |
| —    |      |    |
| —    |      |    |

8 hours

Sharding

50:7P5
128gps

2TB



→ scaled my storage
8 TB

→ scaled concurrent load

↓↓↓↓↓↓

1TB    2TB

↑ ▢ ◯ ◯ ◯

500GB    2TB    load

FAIL OVER

Mastr Slove Repilcation

A

Read ↑    M    30 sec

## MS Replication
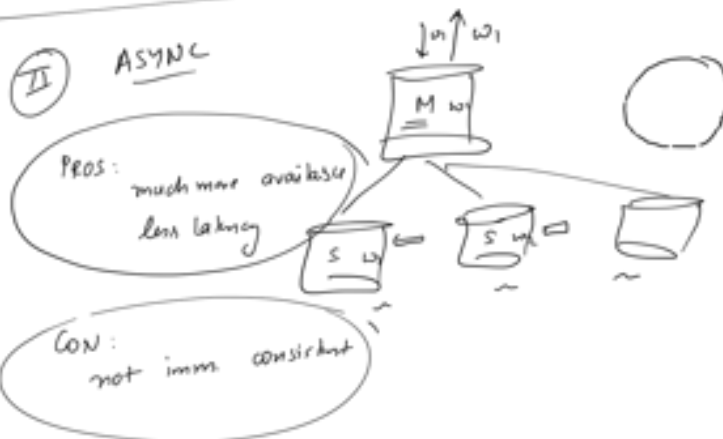
:)

- allows to scale READS
- requirement of backup has reduced
- increased READ Availability

---



① WRITE THROUGH SYNC

$w_1$ ↑ $w_1$ done

M $w_1$ ①

Pro:
- strongly consistent

Con:
- less available
- write latency

$w_1$ ① $S_1$

$w_1$ ① $S_2$

$w_1$ ① $S_3$

---

② ASYNC

$w_1$ ↓↑ $w_1$

M $w_1$

PROS: much more available
less latency

S $w_1$ ⇐ S $w_1$ ⇐

CON: not imm consistent

---

③ QUORUM

↑ $w_2$
↓ $w_1$

M

$\frac{N}{2}+1$  w1 ↗  ↓ w1  → w1



P    S    S

available
not consistent

latency M

---

M — S    Replication

READ

→ scale READ
→ give your
  surface of
  bug

→ Mc. availability
→ less latency latency

M

S1    S2    S3

---



A    B    C    D

A1  A2  A3

---

Master - Mess Replication

P2P

↓W1