## Consistent Hashing

Stateful L B

```
┌─────────────────────────┐
│      Load Balancing      │
└─────────────────────────┘
```

```
┌──────┐   ┌──────┐   ┌──────┐   ┌──────┐
│      │   │      │   │      │   │      │
│      │   │      │   │      │   │      │
└──────┘   └──────┘   └──────┘   └──────┘
```

$10^{18}$ 0 1 2 3 4

S4 @ $10^{16}$

S2 @ 10,200

S3 @ $10^{10}$

R @ 24,000

S1 @ $10^6$

sorted array of server positions

$A =$

| 10,200 | $10^6$ | $10^{10}$ | $10^{16}$ |
|--------|--------|-----------|-----------|
| S2     | S1     | S3        | S4        |

Rid = 107

Load Balancer

Step 1 :    $H_R(107) = 24,000$

Step 2 :    Binary Search (A) to find

ceiling 😃 🙂

$\log(|s|)$    time 🙂

---

PROBLEM:

```
            A  B  C D
        S1  R  R  R R
```

O  R

N  R

M  R

S4

L  R

S2 ↔

S3 ↔

S1  ⇈

K

J  R      I  R      S3

R  E

S2

R  F

R  G

R  H

Now  S4  dies  😐  😏

With this approach of C.N, we have solved

the problem of asking every state to move,
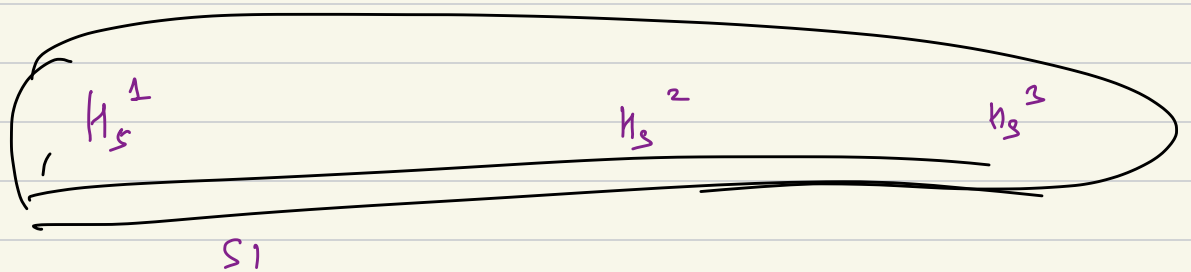
BUTT.TTT.----

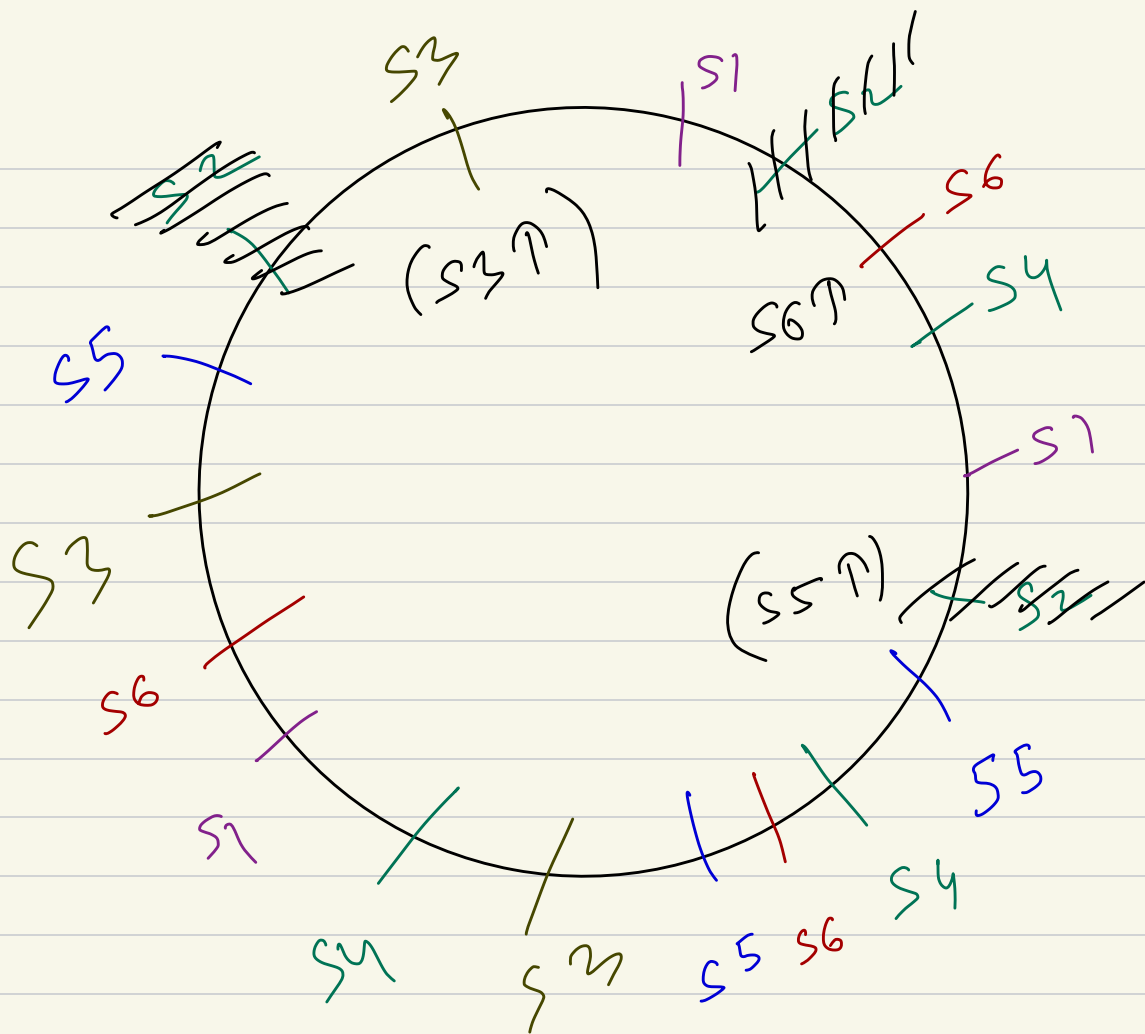We are being unfair with

re-distribution

Now S1 dies 😕😟

S2 ↑↑↑

S3 ⟷

Cascading Failure 😐😐😐

( domino falling )
ripple failure

S3

S2

S1

S6N

S6

S4

S5

(S3N)

S6N

S1

S3

(S5N)

S2

S6

S5

S1

S4

S4

S2

S5  S6

$H_s^1$

$h_s^2$

$h_s^3$

S1

$H_s^1(s_1) = value$

$H_s^2(s_F) = value\ 2$

$H_s^3(s_1) = value\ 3$

| | 10 | 20 | 32 | 48 | | 60 | 72 | 80 |
|---|---|---|---|---|---|---|---|---|

$NGI \ (49)$

$H_R \left( R_{107} \right) = 49$

---

C.4     Redis

Rafka

etc

etc

$H_R \left( R_{107} \right) = \underline{650}$

$H_R \left( R_{111} \right) = 950$

$NGI \Rightarrow S1$

$A =$

| | 100 S1 | | 500 S2 | | 600 S3 | S1 700 | | 800 S5 | 900 S4 | 1000 S1 | | 2000 S7 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|

S1 dies

$H_R \left( R_{107} \right) = 650$     $NGI \Rightarrow S5$

| | 500 S2 | | 600 S3 | | | 800 S5 | 900 S4 | | 2000 S7 |
|---|---|---|---|---|---|---|---|---|---|

Client

Web Browser
Android App
IOS App
Desktop APP

DNS

Gateway + Load Balancer

Client
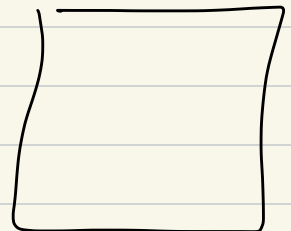
Web Browser
Android App
IOS App
Desktop App

DNS

Gateway + Load Balancer

App servers

App Server

Storage

In the backend, we don't just have one kind of machines doing everything for us

[ Jack of all , master of none !!! 😑 😑 ]

Instead,

we could want to employ different machines with different resources for different jobs.
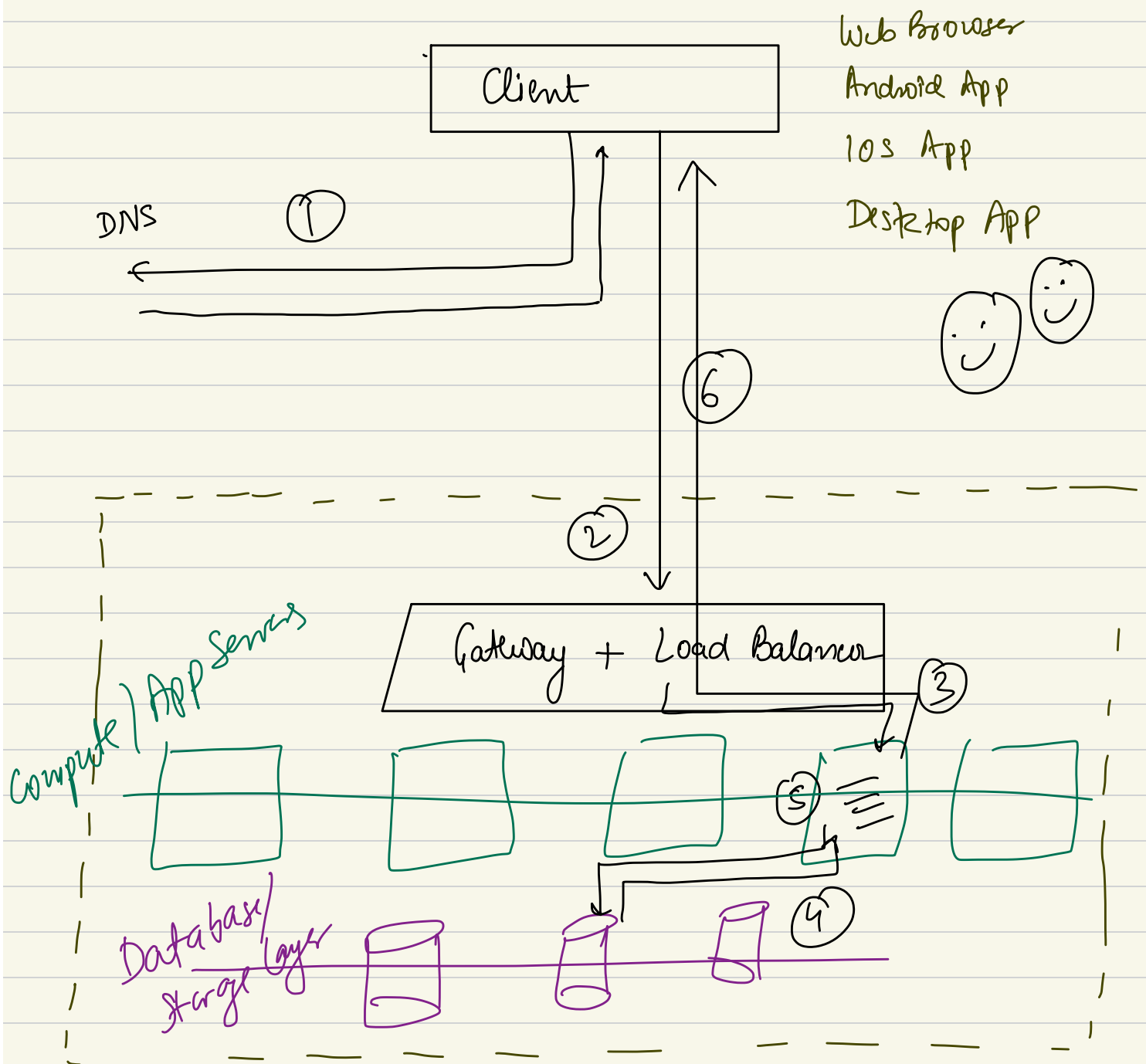
For ex.

① compute layer | App Server layer

App Server / compute machines which are going to have good RAM, good CPU so that code and computations can run better.

② Storage Layer / DB Layer

HDD ↗↗
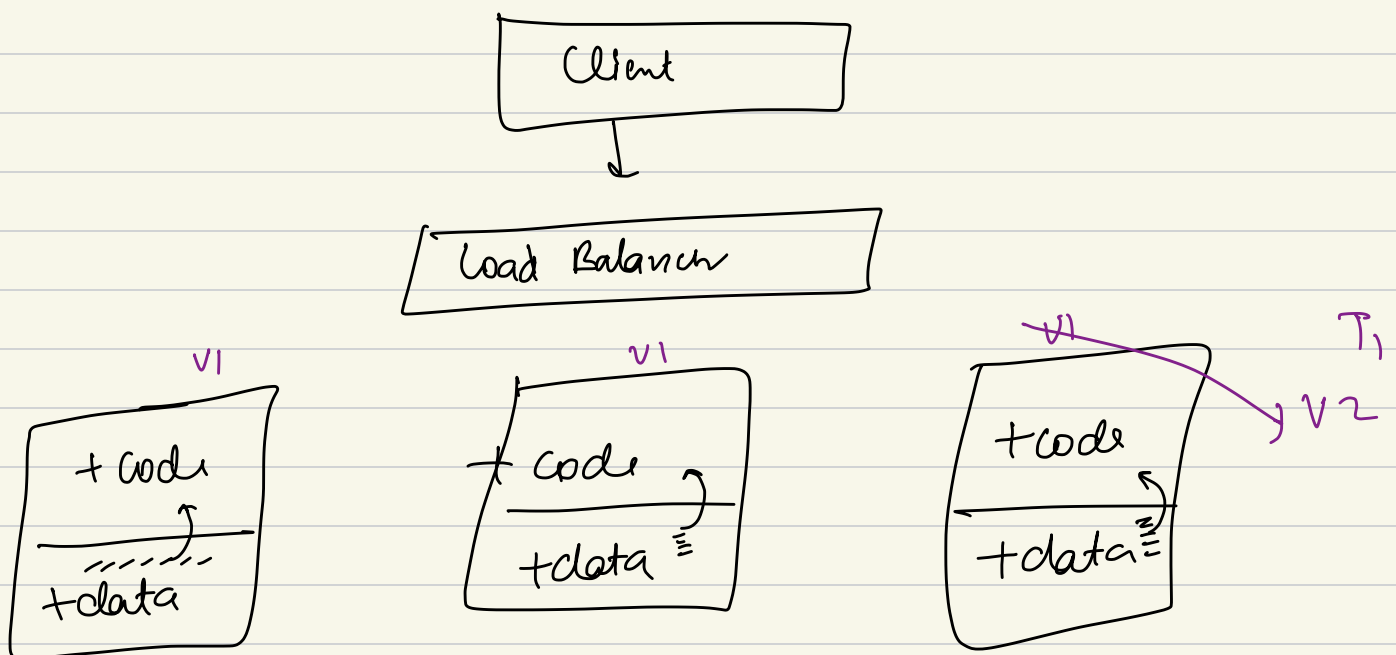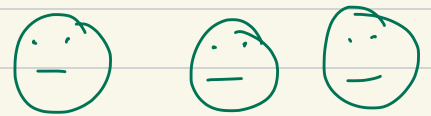
SSD ↗↗

Client

Web Browser
Android App
IOS App
Desktop APP

DNS ①

②

⑥

Gateway + Load Balancer ③

Compute / APP Servers

⑤

④

Database / Storage Layer

## Why decoupling

① Specialized resources.

② Single Responsibility Principle  ☺ ☺

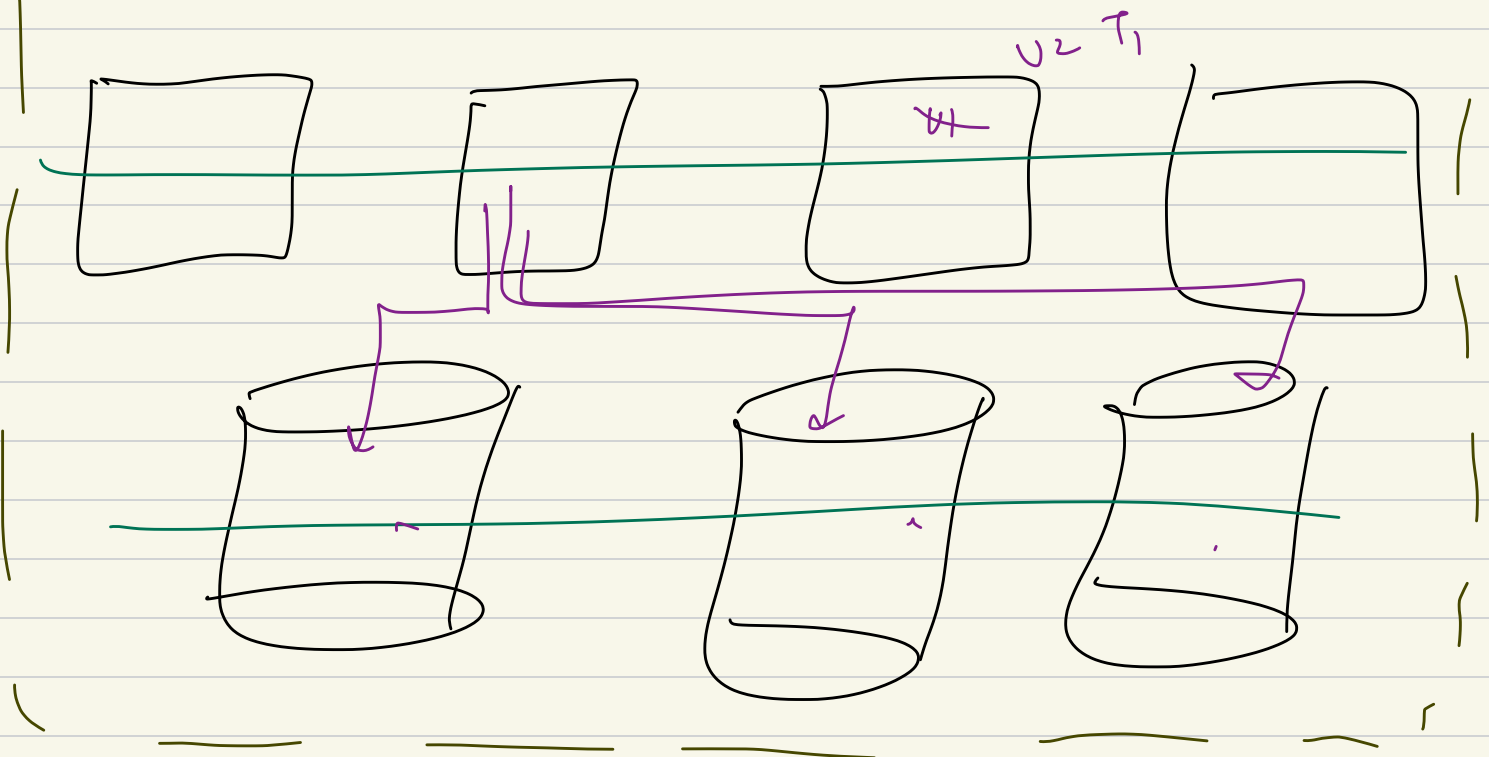③ Scale your systems better.

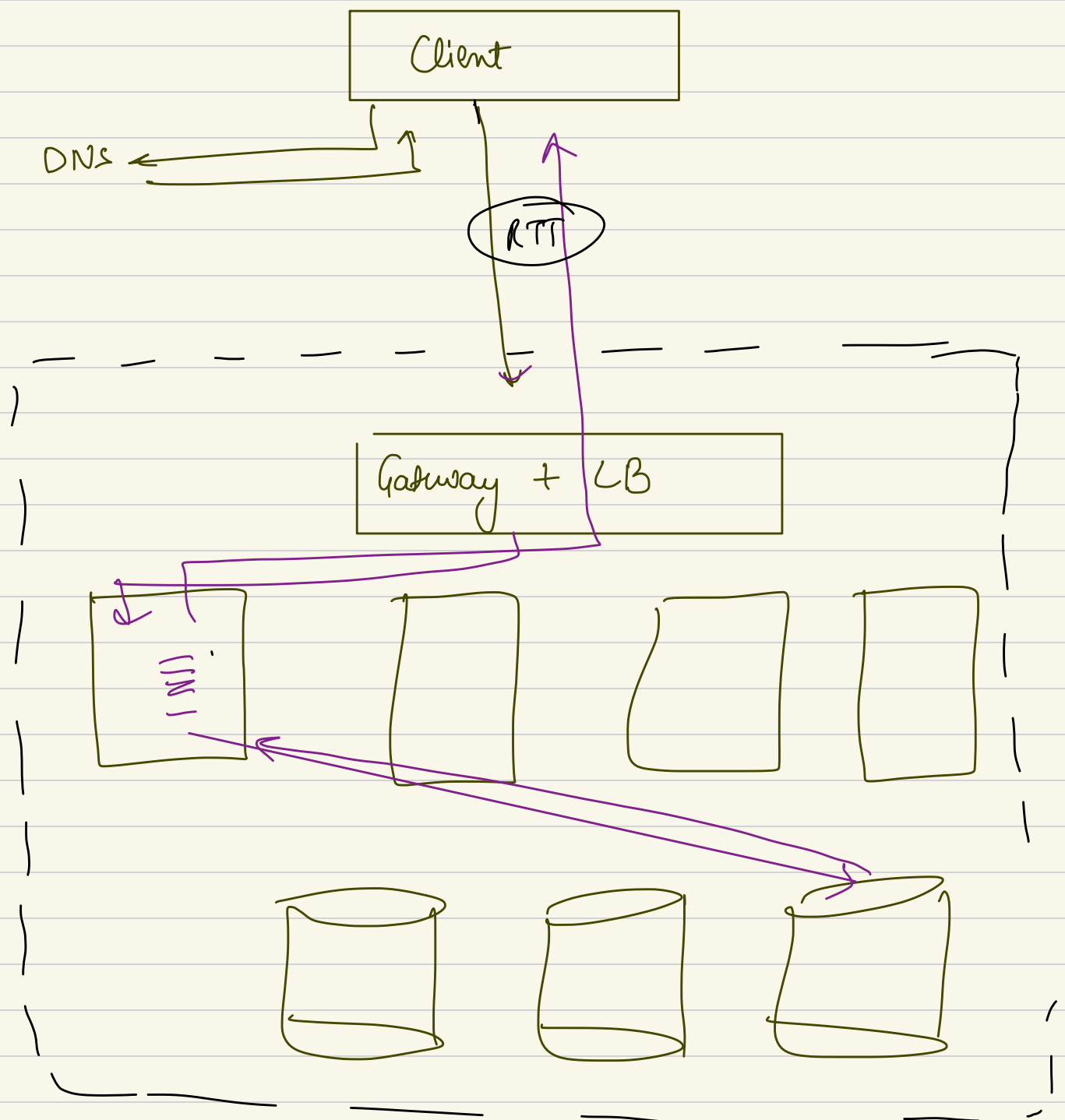④ Code deployments become smoother

---

Case 1 : Coupled System  😐 😐 😐

```
          ┌──────────────┐
          │    Client    │
          └──────────────┘
                 │
                 ▼
        ┌──────────────────┐
        │  Load Balancer   │
        └──────────────────┘
```

V1                        V1                      V1          $T_1$
┌──────────┐          ┌──────────┐          ┌──────────┐      V2
│ + code   │          │ + code   │          │ + code   │
│──────────│          │──────────│          │──────────│
│ + data   │          │ + data   │          │ + data   │
└──────────┘          └──────────┘          └──────────┘

# Case 2: decoupling compute machine and storage machine



Client

Load Balancer

$U_2 T_1$

$\mathcal{H}$

# New Reality

$$\text{Latency} = \left[ \begin{array}{c} \text{Time at which the} \\ \text{request got completed} \\ - \\ \text{Time at which request} \\ \text{started} \end{array} \right]$$

[ Turaround Time ]

User Latency

$$\text{User Latency} = \text{RTT (Round Trip Time)}$$
$$+$$
$$\text{Network Call b/w App server and DB}$$
$$+$$
$$\text{DB Access Time}$$
$$+$$
$$\text{Computation Time}$$

hops

Caching: The idea of creating partial copies of underlying data so that we can access that data quickly. This is called caching :) :)
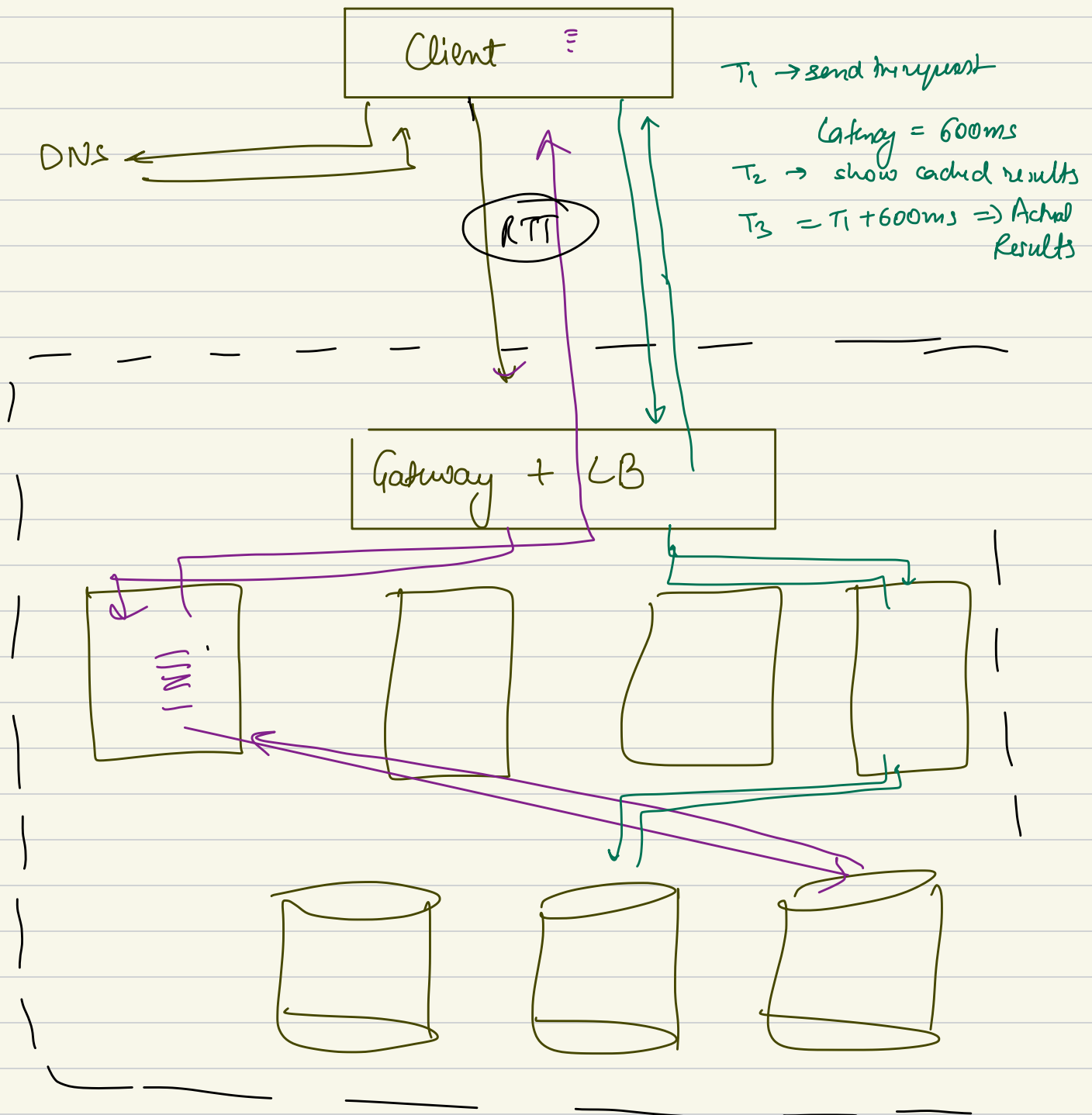
for example in DNS

(1) Browser Cache

(2) OS Cache

(3) Router Cache

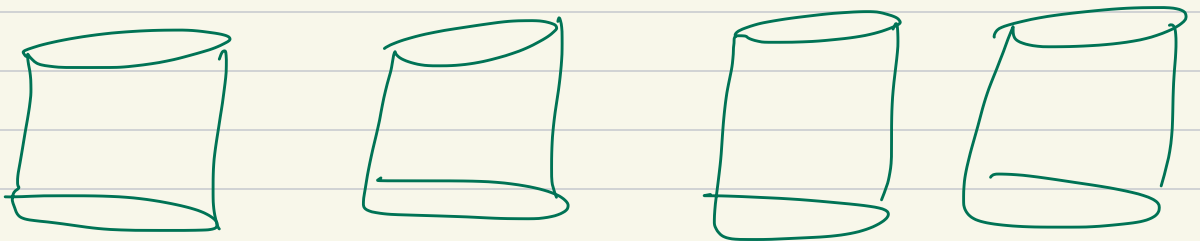(4) Neighbourhood DNS Cache

① Browser / Client Side cache 😊 😊

Client ⋮

$T_1 \rightarrow$ send the request

DNS ←

RTT

$T_2 \rightarrow$ show cached results

$T_3 = T_1 + 600ms \Rightarrow$ Actual Results
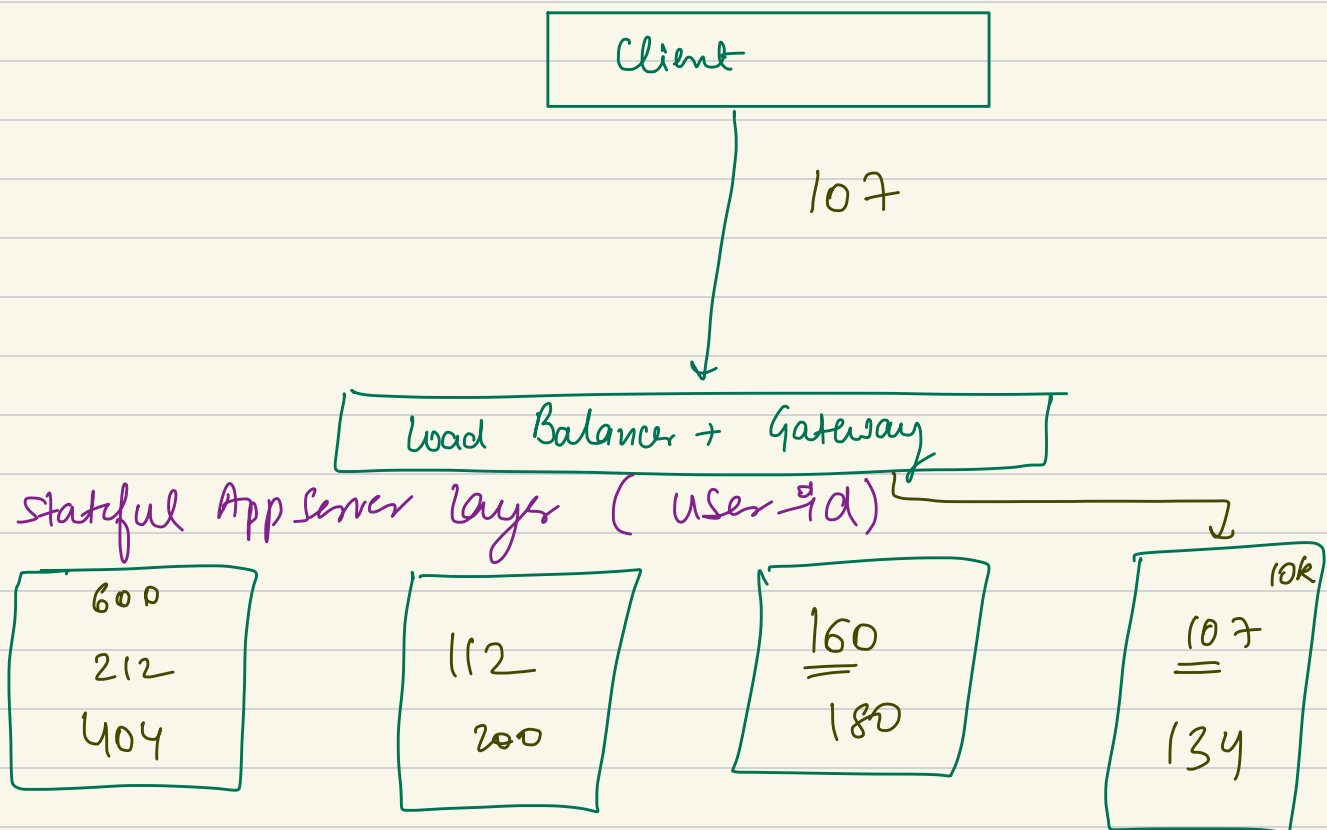
Latency = 600ms

Gateway + LB

Using client side caching ( app data, browser cache)
The client engages the user by showing
cached data, so that the perceived latency
of the user reduces.

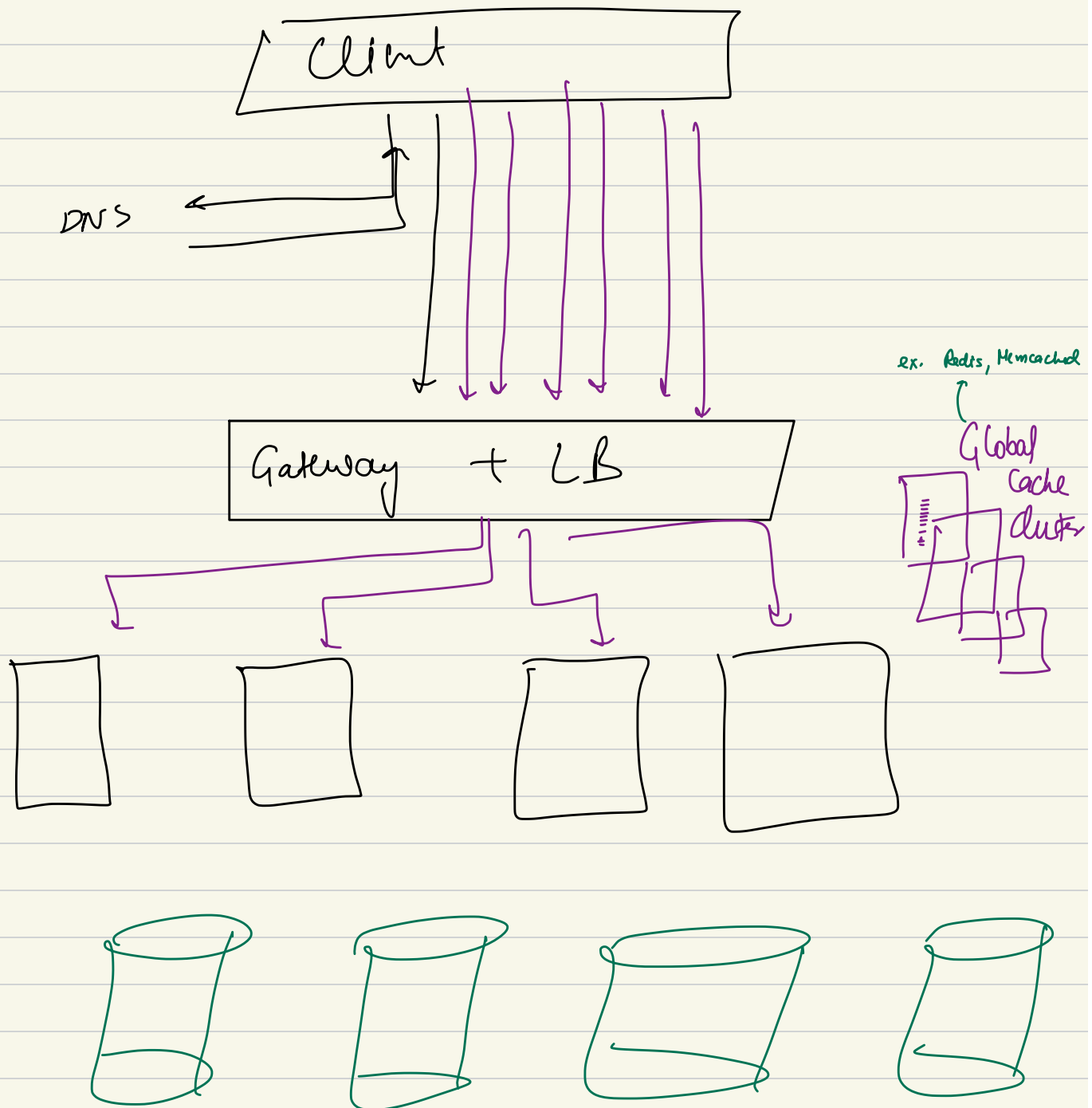"Client Side Optimizations"

## ② App Server Cache

Fb.com

```
Client
```
↓ 107

```
Load Balancer + Gateway
```

Statful App Server Layer ( user-id)

| 600 | 112 | $\underline{160}$ | 10k |
|---|---|---|---|
| 212 | | 180 | $\underline{107}$ |
| 404 | 200 | | 134 |

## ② App Server Caching / Local Caching  ☺

Given we know that in statful load balancing, a particular request is always going to land on the same machine, hence, we can make use of this and cache the relevant data on the relevant machines.
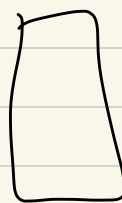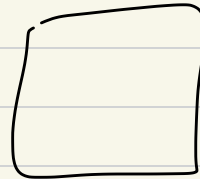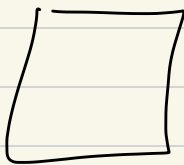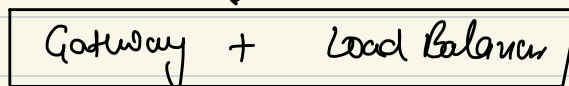
③ Global Caching

Client

DNS

Gateway + LB

ex. Redis, Memcached

Global
Cache
clusters

4　　　　　CDN　　　　　:)

Questions

Session stickyness ≡ Session-id

Load Balancer

Gateway + Load Balancer

Region 1
north central us

proxy

clustP

Region
Paris south

IP1  IP2

IP3  IP4  IP5

The same consistent hashing logic that we discussed for App Servers, the same is applicable for database layer as well; in fact the same is applicable for the caching cluster as well.

App Server → Stateful
App Server → Stateless

storage layer → Stateful