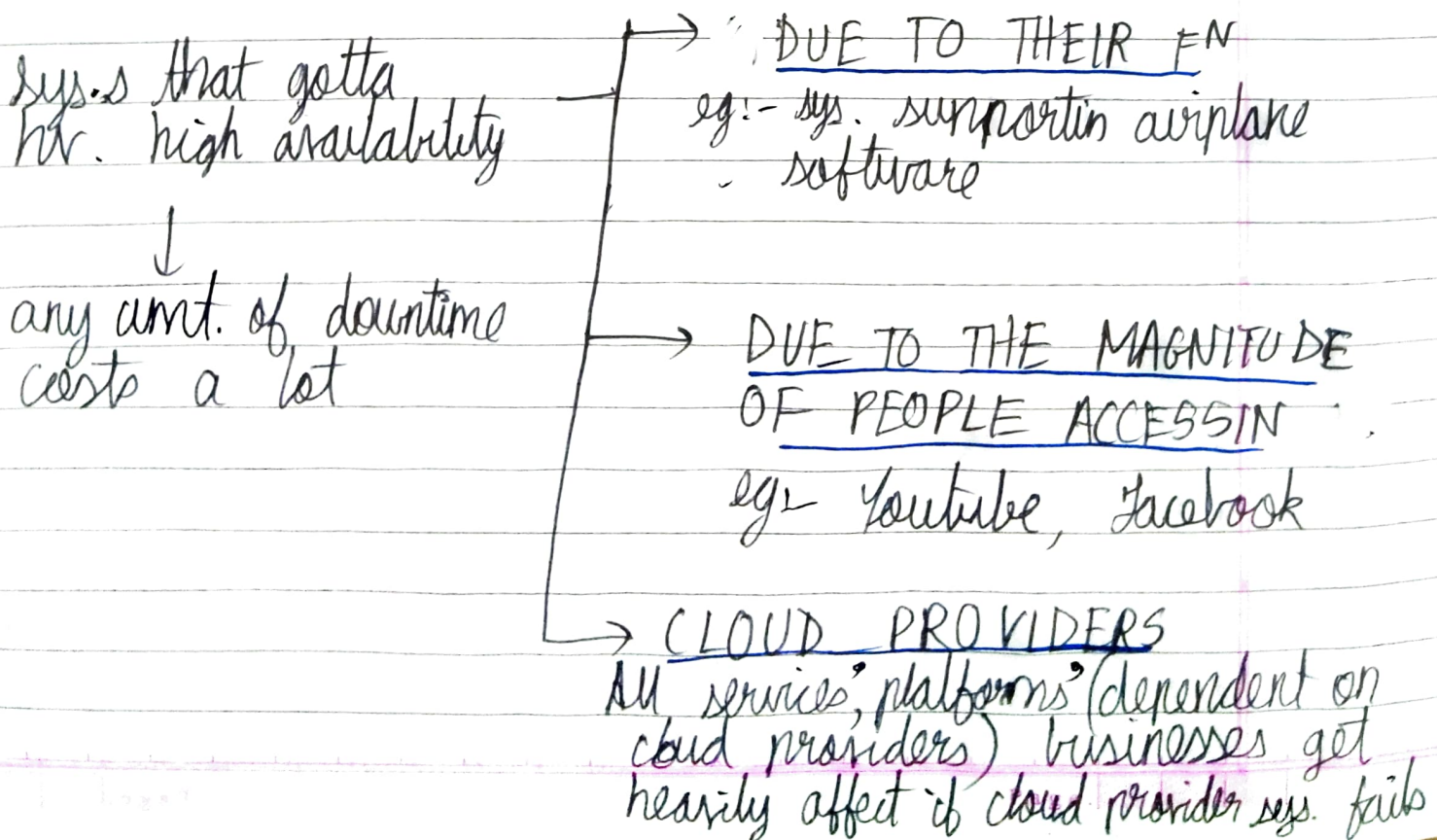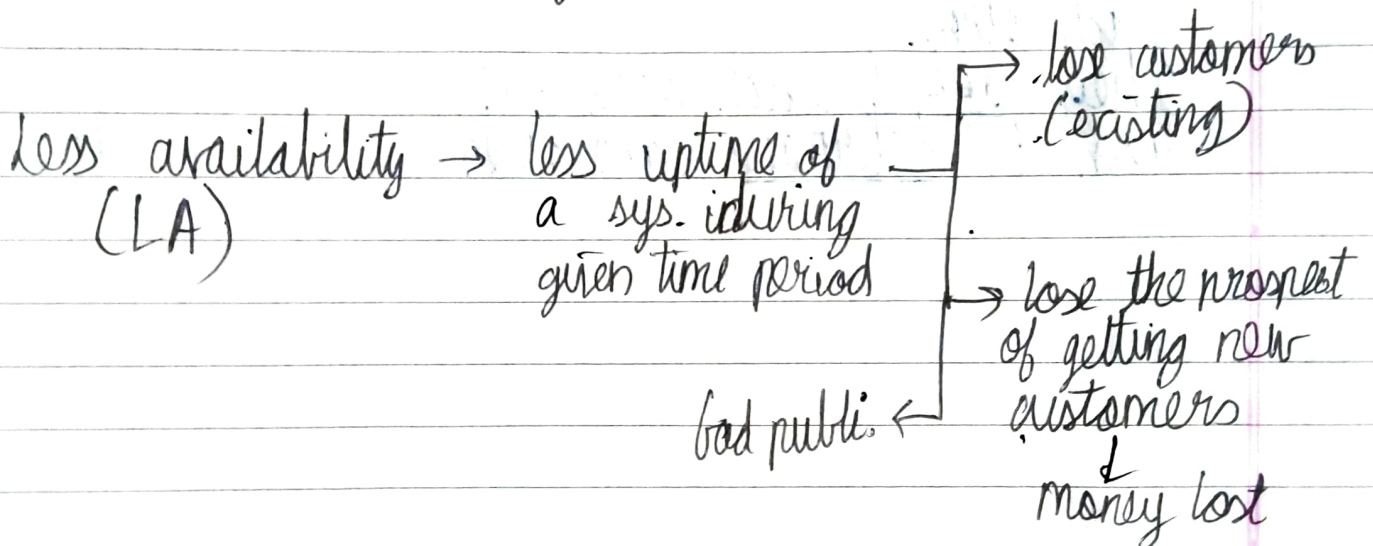# AVAILABILITY

Fault Tolerance :- How ~~tolerant~~ resistant is a sys. to failures. What happens if a server in sys fails? dB fails ? → is sys gonna still be opera"al

Availability → % of a time (eg:- month, year) given period of for which the sys. is atleast opera"al enough to get its primary f"s satisfied

Less availability → less uptime of (LA) a sys. during given time period

→ lose customers (existing)

→ lose the prospect of getting new customers ↓ money lost

bad publi. ←

sys..s that gotta hv. high availability ↓ any amt. of downtime costs a lot

→ DUE TO THEIR F^N eg:- sys. supportin airplane software

→ DUE TO THE MAGNITUDE OF PEOPLE ACCESSIN eg:- Youtube, Facebook

→ CLOUD PROVIDERS All services, platforms (dependent on cloud providers) businesses get heavily affect if cloud provider sys. fails

service ≡ sys.

eg of cloud providers : Azure, AWS, Google cloud, Oracle Platform.

abbr.

## Measuring availability (albt)

Availability is measured in terms of % of sys's uptime in a given year. (90% albt → sys. in up. 90% of the time in a yr)

Most powerful services, sys.s have ~~to be~~ to be really HA. .. i.e. 99.9%, 99.99% etc.

Nines → : Percentages with the no. 9

eg:- 99% albt → 2 nines of albt

| Albt | Downtime /yr |
|---|---|
| 2 nines (99%) | 87.7 hr |
| 3 nines (99.9%) | 8.8 hr |
| 4 nines (99.99%) | 52.6 mins |
| 5 nines (99.999%) | 5.3 mins |

Highly available sys → albt > 5 nines albt.

### SLAs and SLOs (albt. guaranteed explicitly

SLA → Service lvl Agreement → Agreement b/w service b/w service provider and customers /end users of service that guarantee customers amount of albt, error free utility and some other objectives (known as SLOs → service lvl objective), failing which service the service provider pays back the customer some % of their fees

# SLA

eg:-

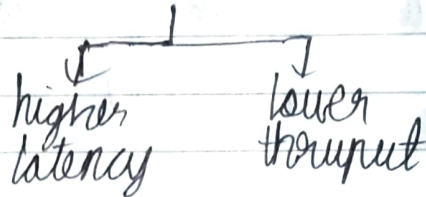AWS (provider)                    Netflix (user)

SLO 1 → 1) Netflix    get 99.9999% albt.
SLO 2 → 2) Netflix    get x errors while usin my service
⇓
If AWS fails to meet SLOs, it returns
10% of the service fees Netflix pays
to it every month

NOTE: Achievin higher albt is really difficult
and may accompany tradeoffs like

higher          lower
latency         thruput

As a sys. des.er, you gotta (HA) decide, which part of
sys you want "high albt." and which part you
don't really need HA.

eg:- Stripe (payment service for business

→ Core services → handlin payments, charging
customers

they gotta have high HA else Stripe, as
well as its client businesses, platforms
lose lotta money

→ Business monitor dashboard → used by client platform
to use business stats.

doesn't really require HA, something,
that if fails, doesn't cause
catastrophic losses

# How a sys. is made "AVAILABLE"

**1]** single pts of failure (parts of the sys. to which, on failing cause entire sys. to fail) should be removed → For this we use **redundancy** while ~~sy~~ designing the sys.

__Passive Redundancy__ → act of duplicatin, or triplicatin or multiplyin even more, certain parts of sys.

C- client
LB- load balancer
S- server
dB- database

**Simple sys** ⇒

(C)
(C)
(C)

(S)    (dB)
↑
single point of failure as, if it fails client can never access database so purpose of sys. has failed

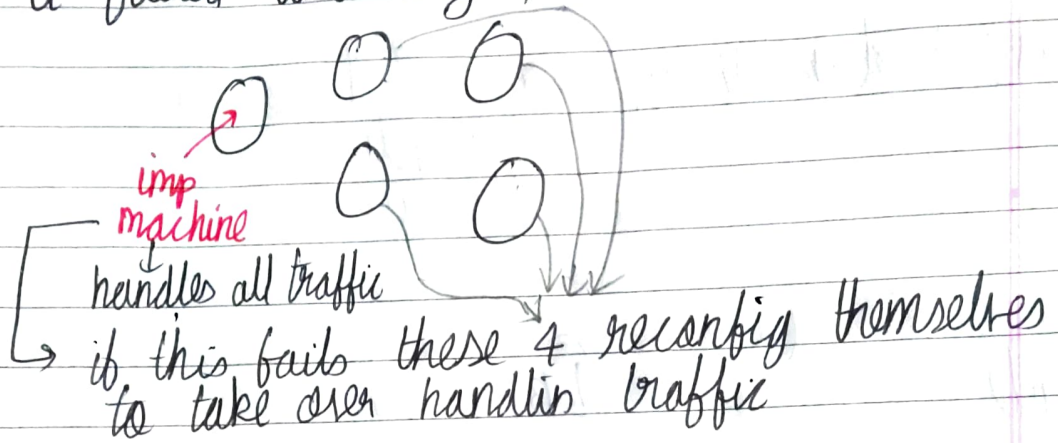**Redundant sys** ⇒

(C)
(C)
(C)

(LB)
(LB)
(LB)

(S)
(S)

(dB)

Load balancer → to ~~reduce~~ equally distri. load across all servers (prvr. + server from becomin too overloaded & actin as single pt. of failure
↓
Multiplied load balancers in this sys to avoid 1 LB from gettin too overloaded

→ duplicated

eg:- twin engine sys. in airplane

Active redundancy – Multip. machines that work together in such a way that only 1 or few of the machines are gonna be typically handlin traffic or doin work (imp. machine). If an 'imp. machine' fails, the others machines know it failed and they'll take over

eg :-



imp machine
- handles all traffic
→ if this fails these 4 reconfig themselves to take over handlin traffic

2] Have rigorous processes in place to handle to handle sys failures cuz its possib. that these failures might req. human interven"

∴ Have these backup op"s in mind while sys. des. to get a crashed sys., up in runnin in proper timeframe.