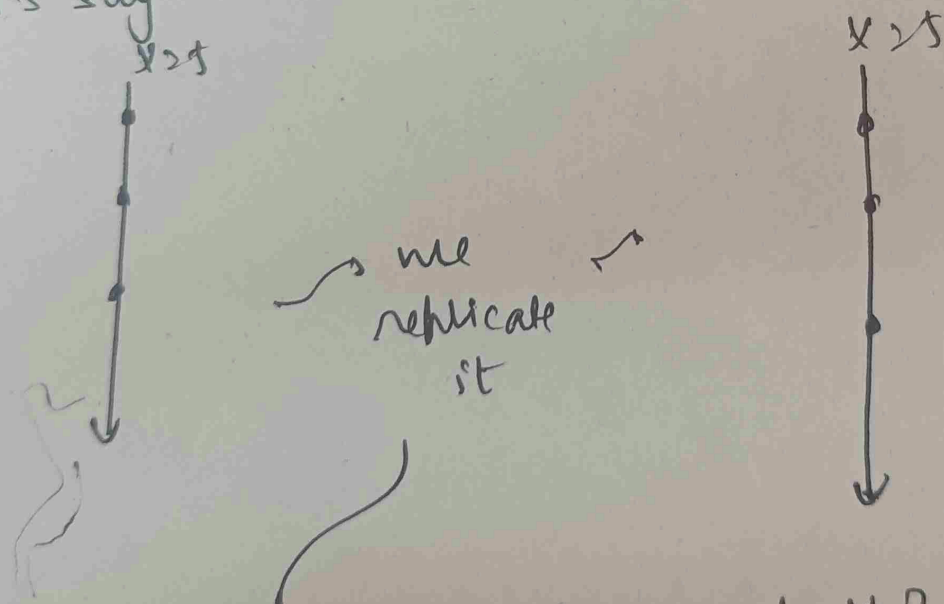# Lec-4 :- Replication, total order v/s determinism

So, last time we began to talk about replication

It is often the case that we mitigate loss of something
by making copies.
→ So we mitigate message loss by making
  copies of messages
→ we mitigate loss of state by making
  copies of state.

Let's say we have a process with some events.
and memo

$x = 5$                                  $x = 5$



→ we
replicate
it
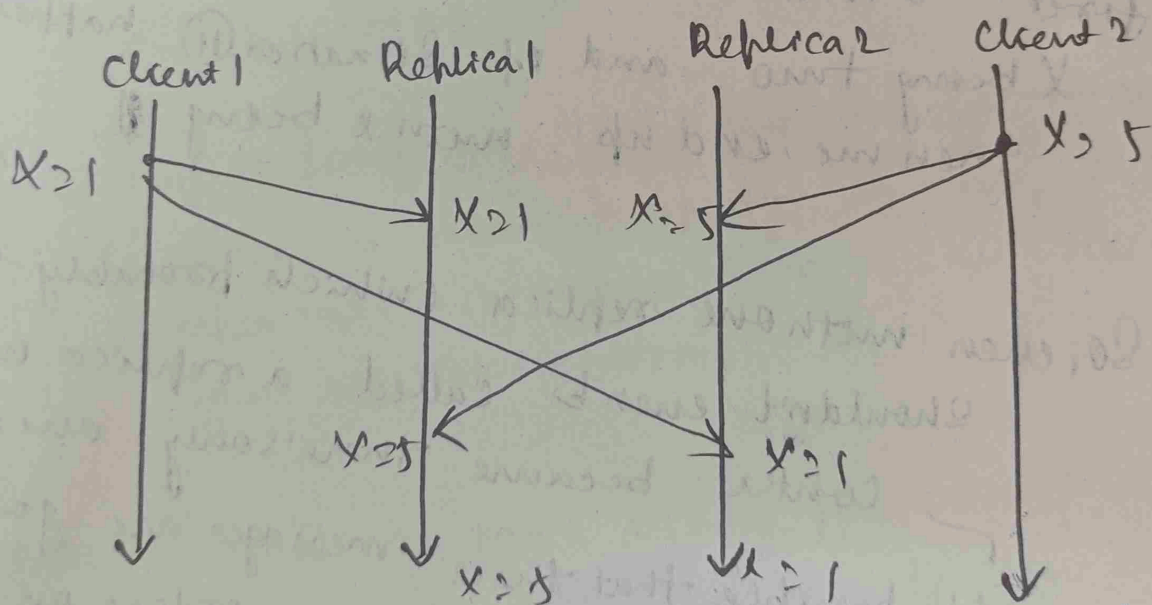
and why would we do it?

⌐→ If one there process crashes, then
   the fact that $x = 5$ isn't lost forever.
⌐→ that's one reason to do replication
   # mitigating data loss

Reasons to do replication:

— mitigating data loss (fault tolerance)

— data locality (you want the copies to be close to clients, for faster response)

— Scalability (serve more requests)
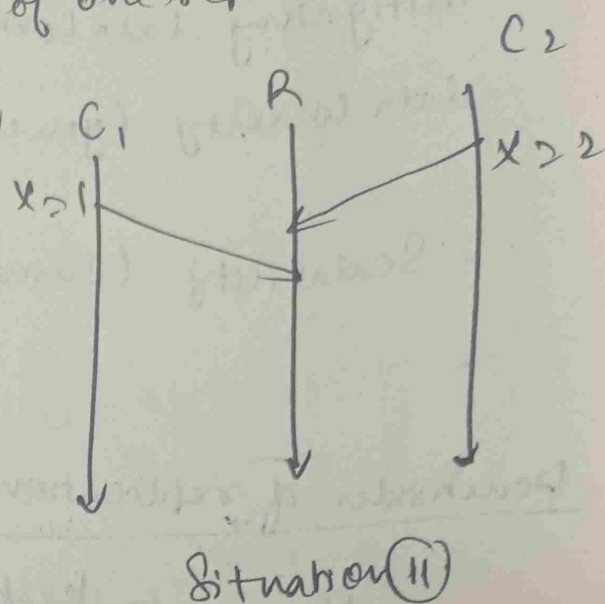
Downsides of replication:

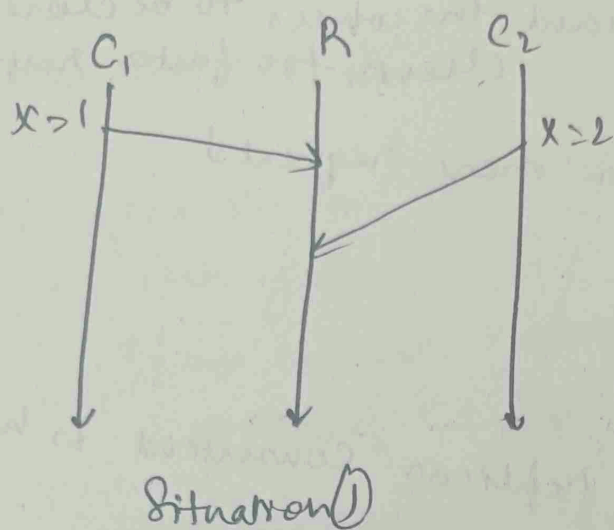→ Having to keep replicas consistent is huge problem

→ Its expensive



This illustrates, one of the challenges with replication. Issue here is that different clients happened in different order on two replicas

∴ Replicas end up being inconsistent

Let's just consider the case of one replica



Situation ①

Situation ⑪

If first situation happens, then we end with $x$ being two, and if Situation ⑪ happens, then we end up with $x$ being 1
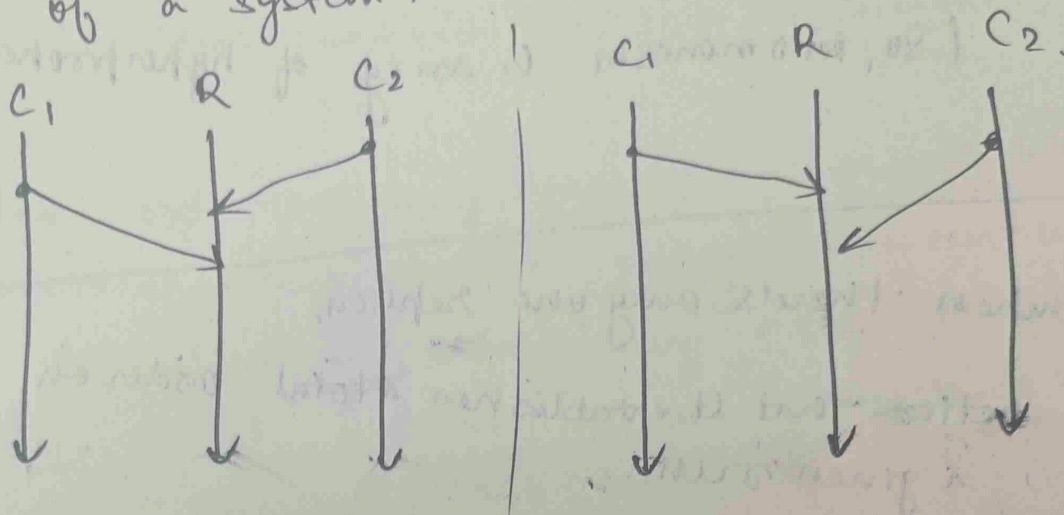
$13!$ N.

So, even with one replica (which probably shouldn't even be called a replica in this context because there is only one of it)

↳ It's possible that there messages are going to arrive in different orders on different runs.

Situation ① ya ⑪ mai kor galat ni keh sakte naa hi dono mai kor Total order delorus ko violate ker raha

# Determinism

- It is a property that relates multiple _runs_ of a system to each other.



It is a violation of determinism,
it is not a violation of Total Order.

It is imp to point out, that even though total order delivery is a nice property to have, it's important to be aware that it still doesn't give you determinism.

**Note:** The totally order delivery property is a property that can only be true or not true of a single run.

But determinism is a property that relates multiple runs.

So, you can't look at a single run and say is this system deterministic. You have to be looking at multiple runs.

A property that relates multiple runs to each other
is called a hyper property.

(So, determinism is an eg of hyperproperty)

---

So, when there's only one replica,

notice that it establishes a total order on
a given run.

what we would like to do is make it so that
there can be multiple replicas but this
nice Total ordered delivery property.

→ In other words, we want clients to think that

they are only dealing with one replica,

we want them to think that they are dealing with
a system that does not have multiple replicas.

## Strong consistency (Informal def):

→ A replicated storage system is strongly consistent if clients can't tell that it is replicated.

It turns out that every strongly consistent replication protocol, that we are going to discuss is going to work by establishing a total order on events, but they all are going to do in different ways.
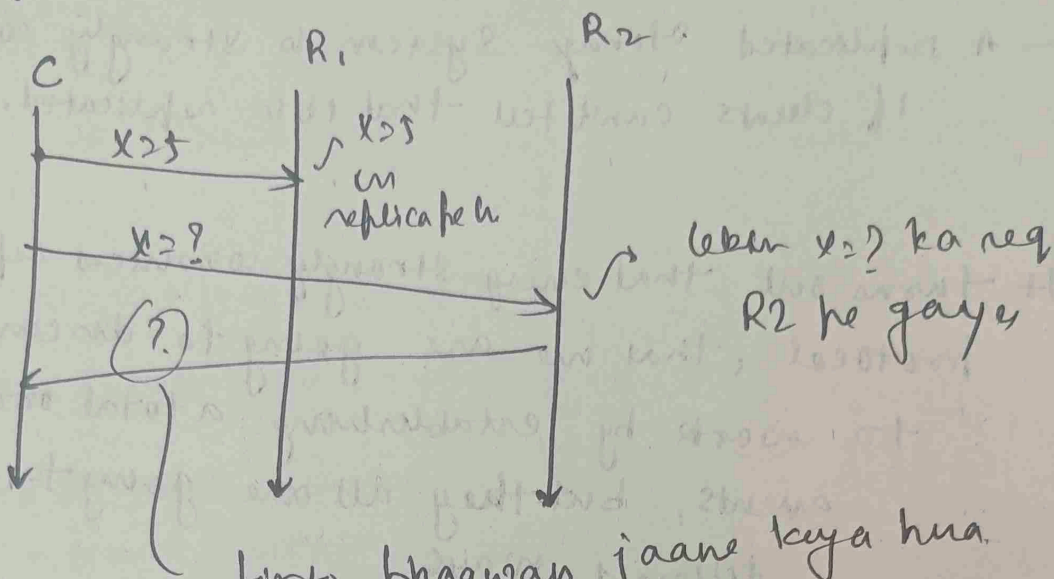
Before that, let's get into particular approaches for establishing this total order and implementing strong consistency.

Let's first talk about some of the ways where a client might be able to tell that data has been replicated.

In other ways, different ways in which replicas might disagree.
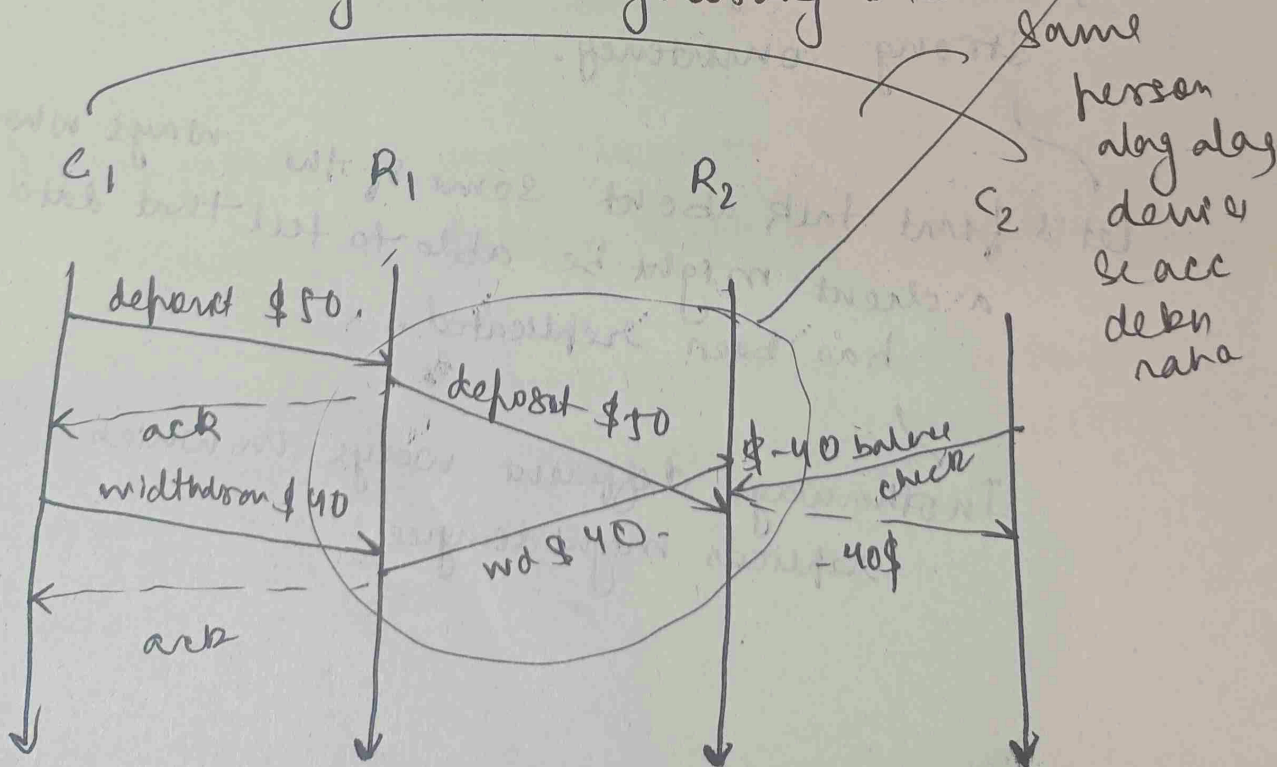
One of the worst cases that can happen k:

C          $R_1$                    $R_2$

$x>5$ → $x>5$ in replicate h

$x=?$ → leken $x=?$ ka req $R_2$ pe gaye

(?)

tirto bhagwan jaane kya hua.

# Read your writes anamoly.

Another thing that can go wrong k:-

same person alag alag device se acc dekh raha

$C_1$        $R_1$                $R_2$                $C_2$

deposit $50.

deposit $50

$-40 balance check

ack

widthdraw $40

wd $40:        -40$

ark

Case ① : mar to $R_1$ bothered hi mi thg
$R_2$ ko update
karne ko.

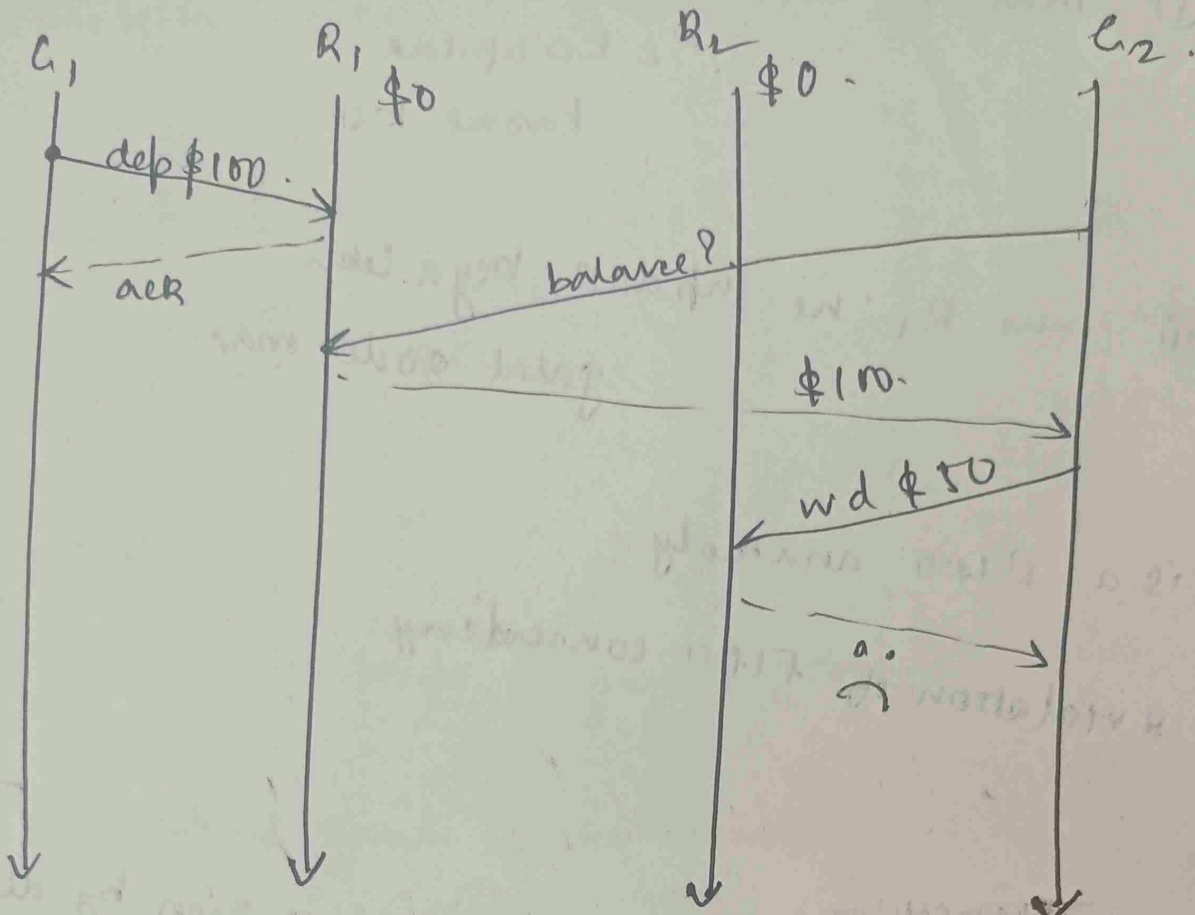Case ② mar $R_1$ ne update kiya leken
galat order mar

→ It is a FIFO anamoly.

# violation of FIFO consistevey.

[
FIFO consistency :—
   writes done by a single procens are seen by all
   procenes in the order, they were ensued.
]

Another bad thing that can happen : PTO

A sequence diagram with vertical lifelines labeled $C_1$, $R_1$ ($\$0$), $R_2$ ($\$0$), $C_2$.

- $C_1 \to R_1$: dep $\$100$.
- $R_1 \to C_1$: ack
- $C_2 \to R_1$: balance?
- $R_1 \to$ : $\$100$.
- $C_2 \to$ : wd $\$50$
- : a.

why did client②  think that they can do withdrawl?

→ withdraw is being done because of the earlier deposit

and Client② find out about this deposit because of "balance?"

so, if I follow chain of events

deposit happened before withdrawl.

So client 2 is not doing anything wrong.

The problem is that Replica ② sees the withdrawl but doesn't see the deposit which was in the causal history of withdraw.

So, because R② sees wd but doesn't see everything that was in the causal history of wd, this is called violation of causal consistency.

Causal consistency :

writes that are potentially related (ie. related by →) must be seen by all processes in the same order.

# Consistency Hierarchy

Read-your-writes consistency
$\downarrow$

PRO consistency
$\downarrow$

Causal "
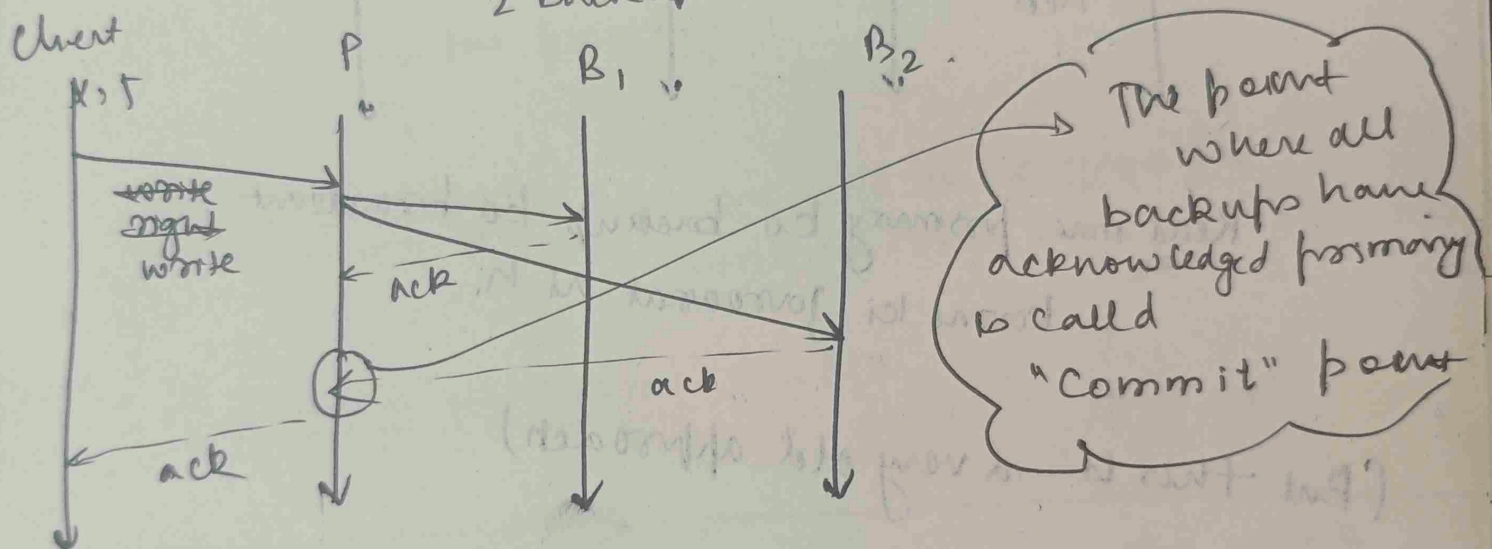$\downarrow$

Strong consistency

Bad

# Primary- backup replication

So, the idea is pretty straightforward:

we pick a particular replica to be what's called the primary, and other replicas are backups.
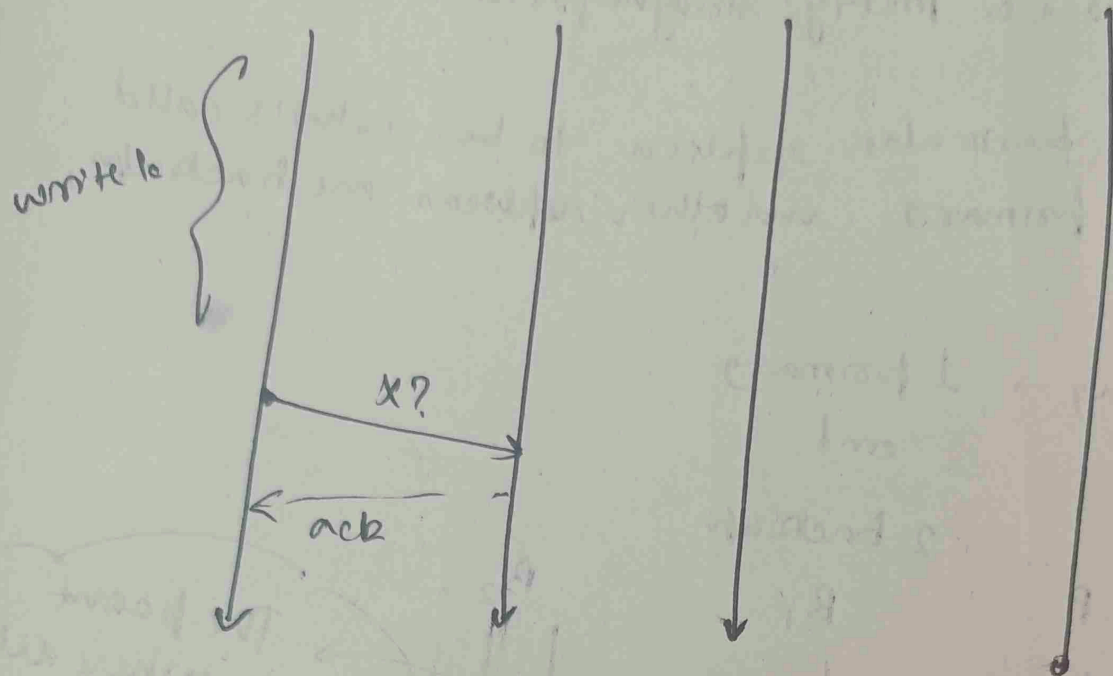
let say → 1 primary
and
2 backup



Client
X, 5

P

B₁

B₂

write
write
write

ack

ack

ack

The point where all backups have acknowledged primary is called "commit" point

clients only interact with the primary.
So, when the primary get a write from a client, (~~right~~)
it then broadcast that right to all the backups, and all backups send an ack
(P)
to the primary,
At the point when all of the backups have acknowledged the right to the primary, then the primary can tell the client that the right succeeded.

what about  read operation?



write to {

X?

ack

Read mai primary ko backup ko broadcant
karone ki Jaroorat ni h.

(But this is a very old approach)

Drawback of this:-

— It is slow

Does Primary backup replication help with:

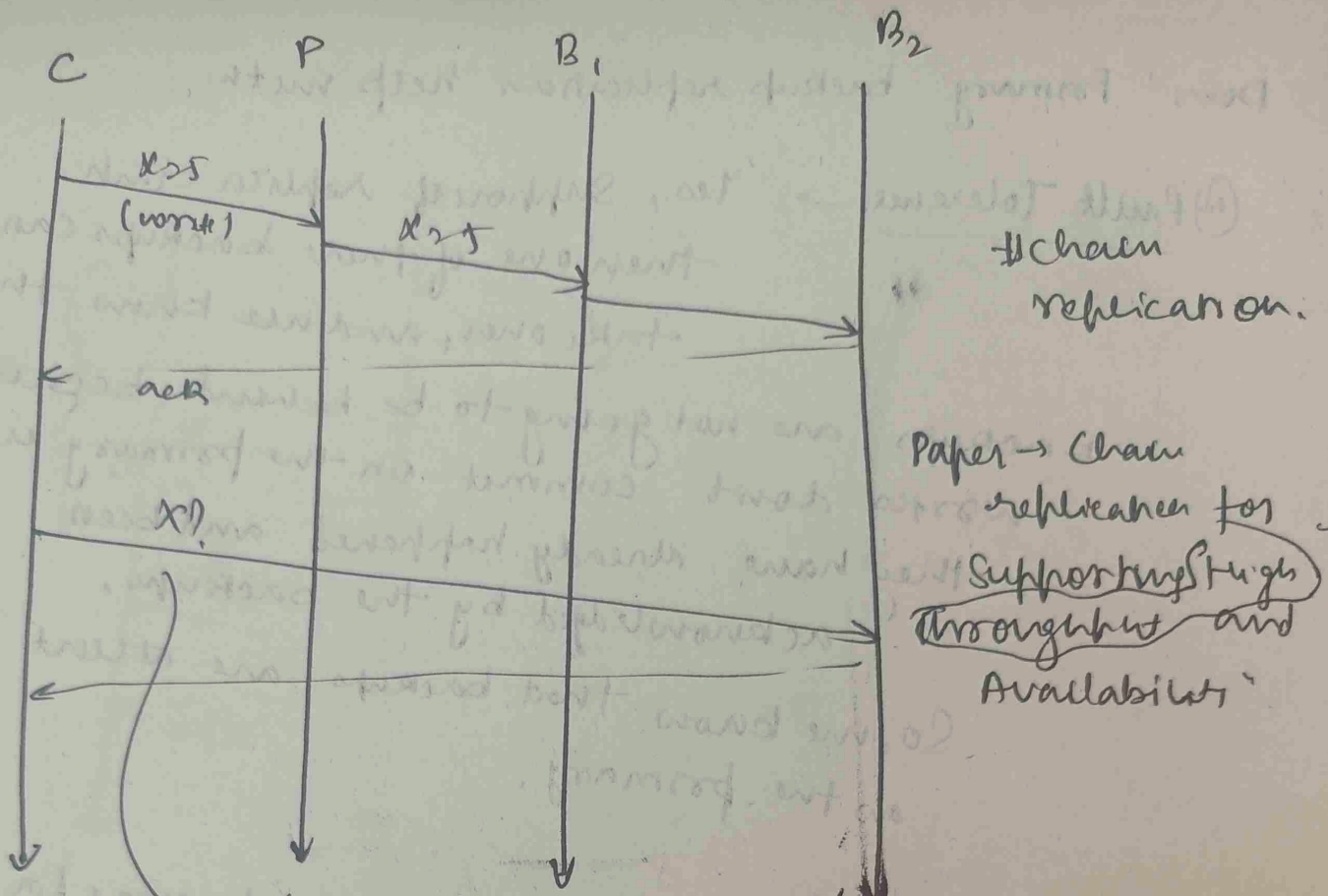(a) Fault Tolerence → Yes, Suppose if replica crash then one of there backups can take over, and we know that backups are not going to be behind, because writes don't commit on the primary until they have already happened and been acknowledged by the backups.

So, we know that backups are atleast uptodate as the primary.

(b) Data Locality :- No, as read to kogage hai write mi korh naa hi read.

(c) Dividing up the work → No

Could we do better?, Could we spread out the work any better?

# chain replication.

Paper → Chain replication for Supporting high throughput and Availability

**Diagram labels (sequence diagram):**

Columns: C, P, $B_1$, $B_2$

- $x=5$ (write) — from C to P
- $x=5$ — from P to $B_1$
- " — from $B_1$ to $B_2$
- ack — back to C
- $x$? — from C across
- (reply back to C)

for read operation we choose one of the backup

Head → that only handles write.
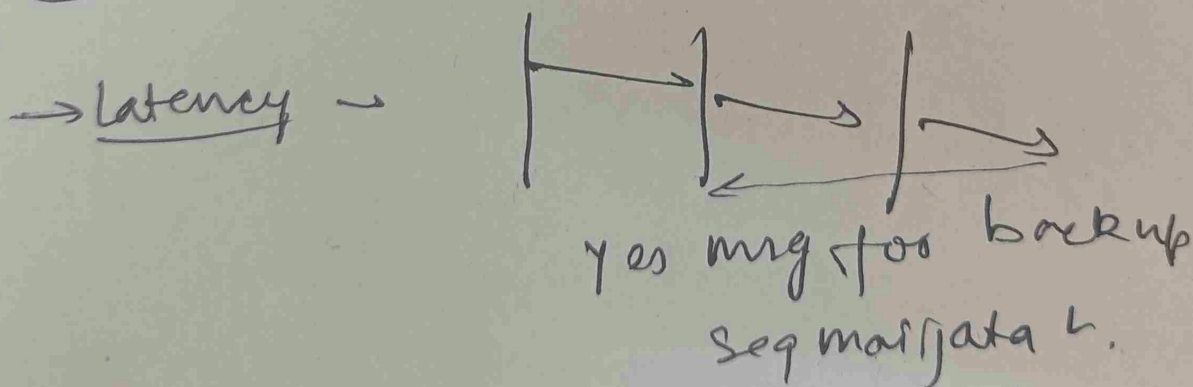
Tail → that only handles read.

→ How are they claiming High Throughput?

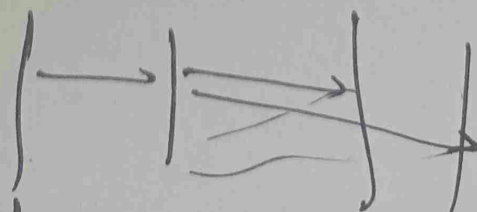[Throughput - no. of actions per unit of time.]

If we have a more or less equal mix of writes and reads, then in theory, chain replication is going to be a better choice than primary backup replication.

As we have one node handling writes and one handling reads then we can process more requests in a given amt of time.

Downside of Chain Replication

→ Latency →

Yes msg for backup
Seq mai jata 4.

bkan primary me

parallel mai jata tha, to faster tha.