

GSCI 1380 : Day 11

Time



Today

Ordering

- ① Total ordering
- ② FIFO ordering
- ③ Causal ordering

Clocks

- ① Logical clocks
- ② Vector clock

(Cockroach DB / YugardB)

(Amazon's Dynamo)

Ordering based on clocks

Global Snapshots

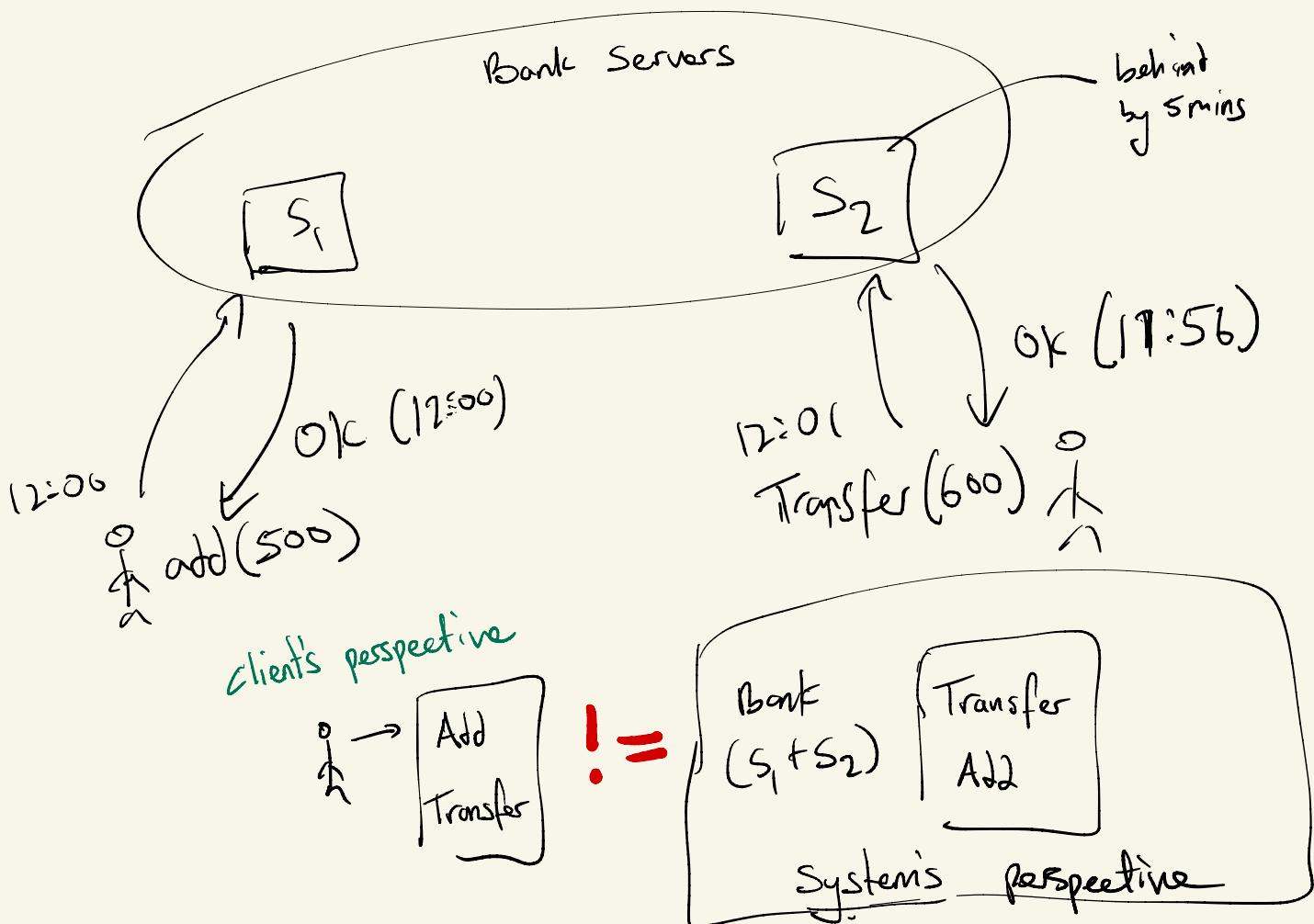
- ① Identifying inconsistent snapshots with
Vector clocks

Real Time is Bad

① hard to synchronize servers

② time is not monotonic

lookup the "leap second" bug

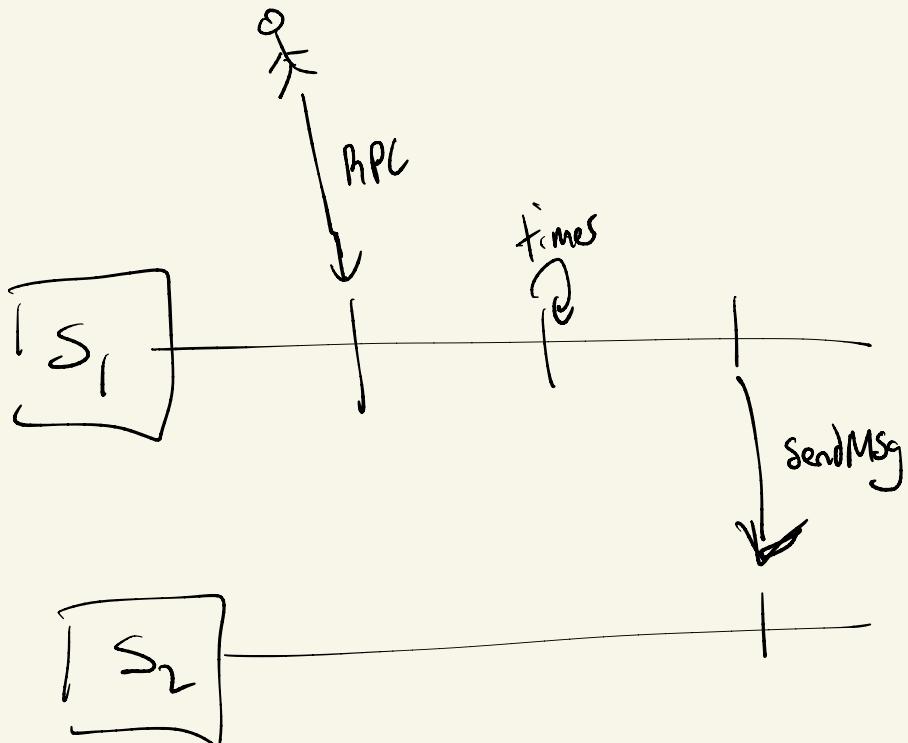


We want a monotonic clock

Event types / Assumptions

event Types

- (1) Receive
- (2) local (timer)
- (3) Send msg



Assumption

- (1) all servers can dump events to external store (S_3)

Type of order (total / FIFO / causal)

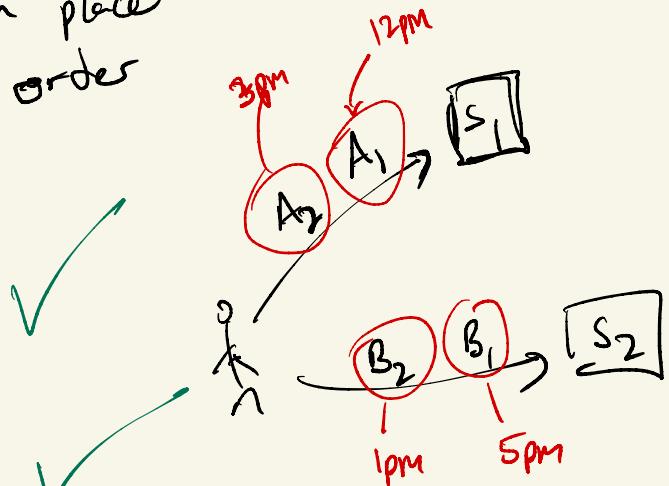
Total = every server can place events in same order

S_1
 A_1, B_1, A_2, B_2

S_2
 A_1, B_1, A_2, B_2

S_1
 B_2, A_1, A_2, B_1

S_2
 B_2, A_1, A_2, B_1



FIFO = every server orders event in order that it was received by the local server: they may not agree on ordering

the servers which got them
 $A_1, A_2 \in S_1$

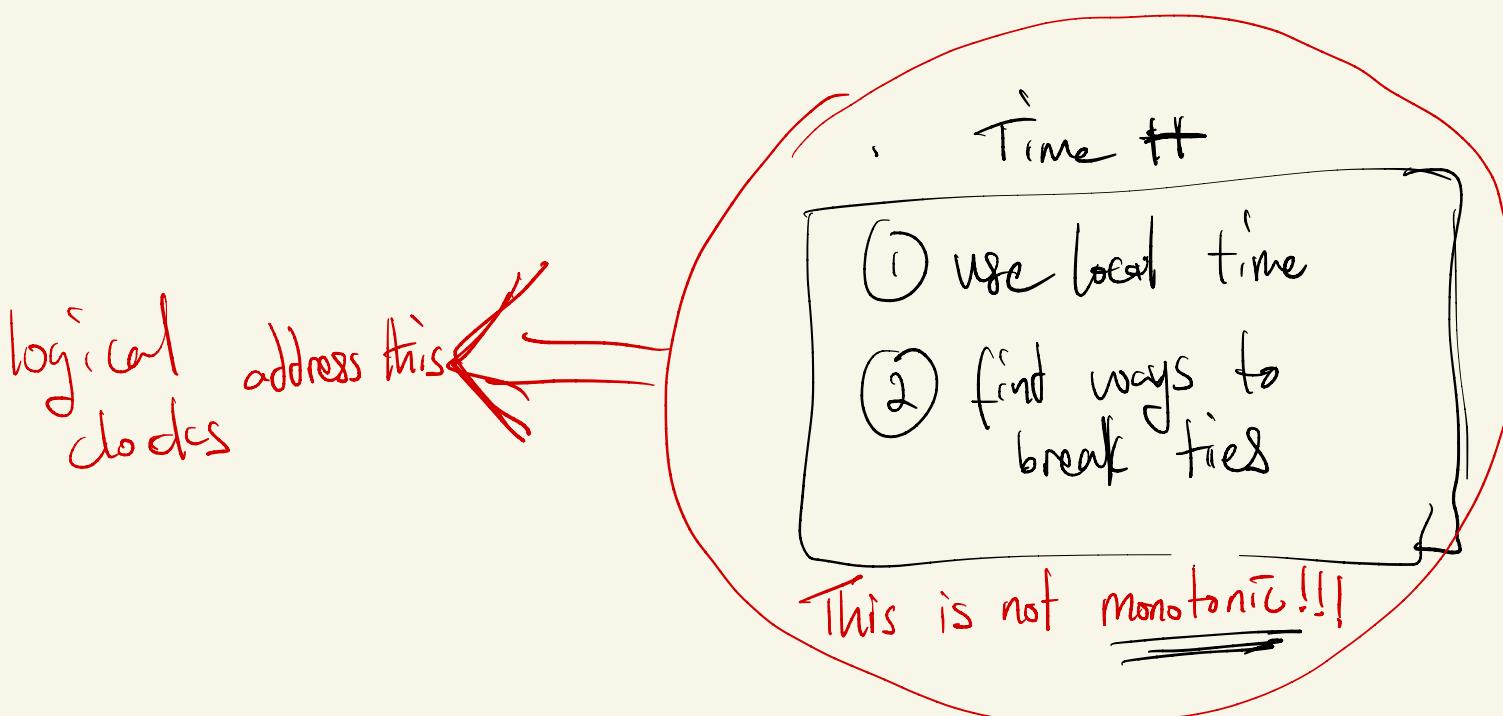
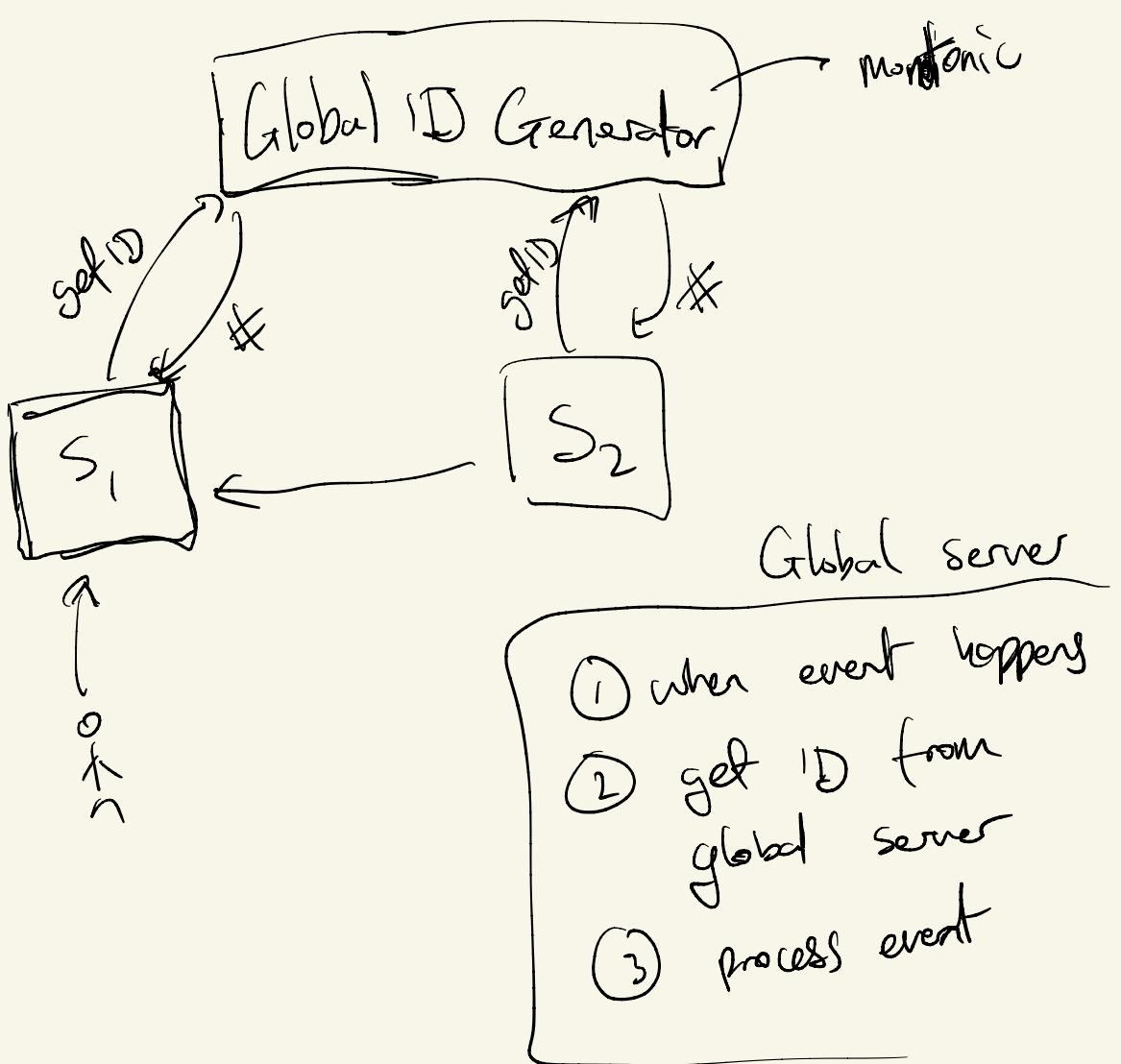
Both are FIFO
because A_2 is always after A_1 ,
& B_2 is always after A_1

S_1
 A_1, A_2, B_1, B_2

S_1
 A_1, B_1, A_2, B_2

S_2
 B_1, B_2, A_1, A_2

S_2
 A_1, B_1, B_2, A_2



rules for logical clocks

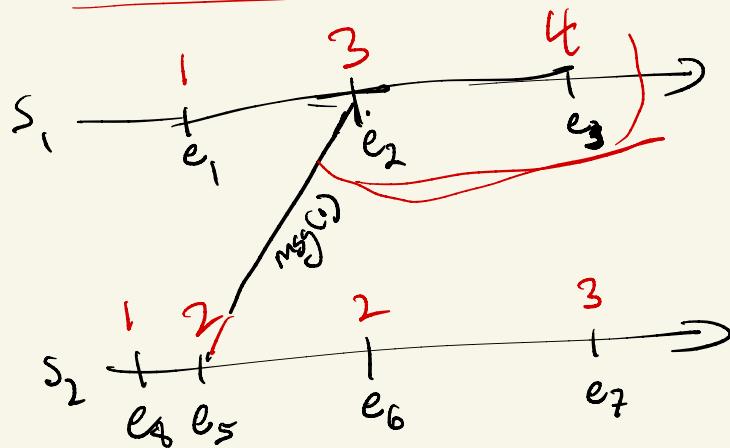
each S maintains a clock ($C=0$)

if event == (local | SendMsg)

$C++$ → include clock in msg

if event == recv

$$C = \max(C, \text{msg}.C) + 1$$



$e_1:1$
$e_2:2$
$e_3:3$
$e_5:1$
$e_6:2$
$e_7:3$

sort based on clock & break ties using server IPs

$e_1, e_2, e_3, e_4, \dots$

Vector Clock

Rules for Vector Clocks

each server maintains a vector of clock ($VC_{[N]}$)

of servers
in system

initialize $VC[] = 0$

if (event == (Send | local))

$VC[i]++$

$SendMsg(VC)$

for Server 5
 $VC[5]++$

else if (event == recv)

for j to N :

$VC[j] = \max(VC[j], msg.VC[j])$

}

$VC[i]++$

Performance problems

① maintain VC

② send VC in each msg

How to determine N (when cluster is dynamic)

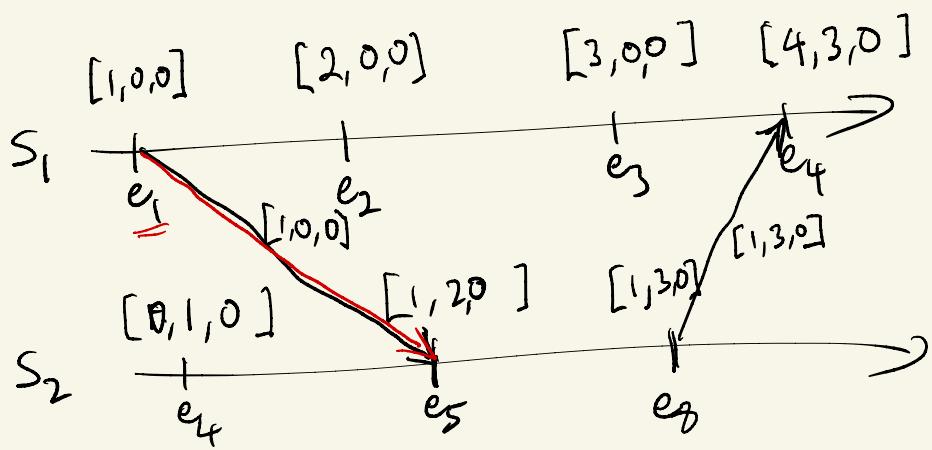
Inherently

① virtual nodes (issues with nodes leaving)

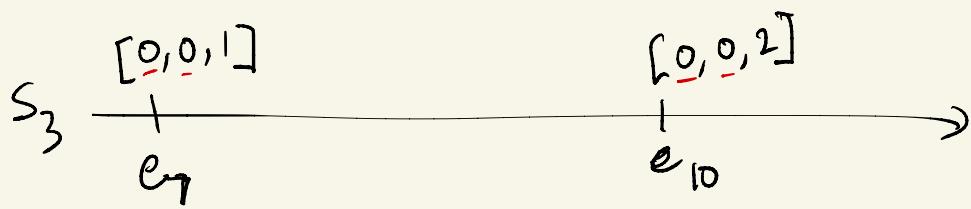
② place a limit on your system (< 1000 nodes)

③ fix the number of nodes (config with all)

Dynamically



$$e_5 = \frac{\max([0, 1, 0], [1, 0, 0])}{\downarrow} \\ \boxed{[1, 1, 0]} \\ \downarrow \\ [1, 2, 0]$$



$$e_4 = \frac{\max(e_3, \text{msg}_{e_8}) = \max([1, 3, 0], [3, 0, 0])}{\downarrow} \\ \boxed{[3, 3, 0]} \\ \downarrow \\ e_4 = [4, 3, 0]$$

Observations

without msgs you can't learn
about another server's clock

Causal ordering (benefits of VC)

"happens before"

A happens before B

] A & B are events

$$VC_A \leq VC_B$$

all clocks in VC_B must be equal / greater than VC_A
& at least one of VC_B elements must
be greater than VC_A

$$\begin{bmatrix} 1, 0, 1 \end{bmatrix} \leq \begin{bmatrix} 1, 0, \underline{2} \end{bmatrix}$$

↑

happens before
the other

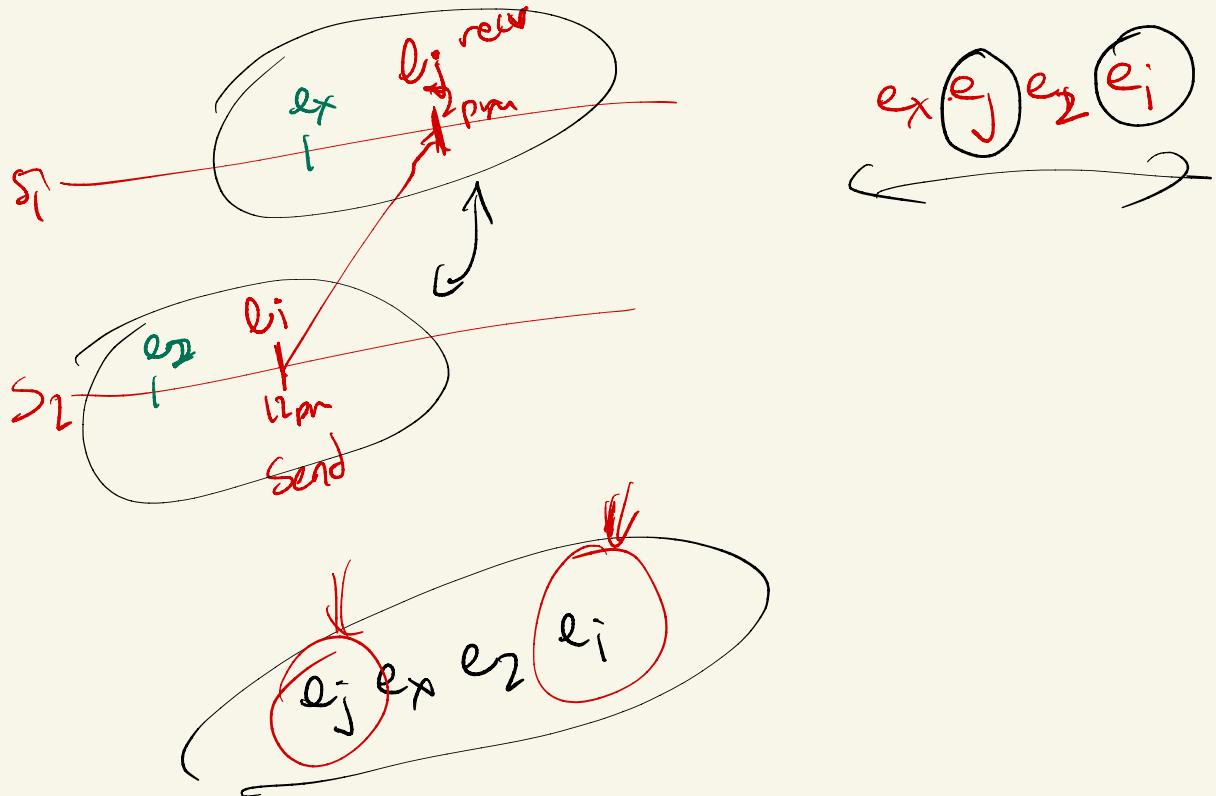
greater

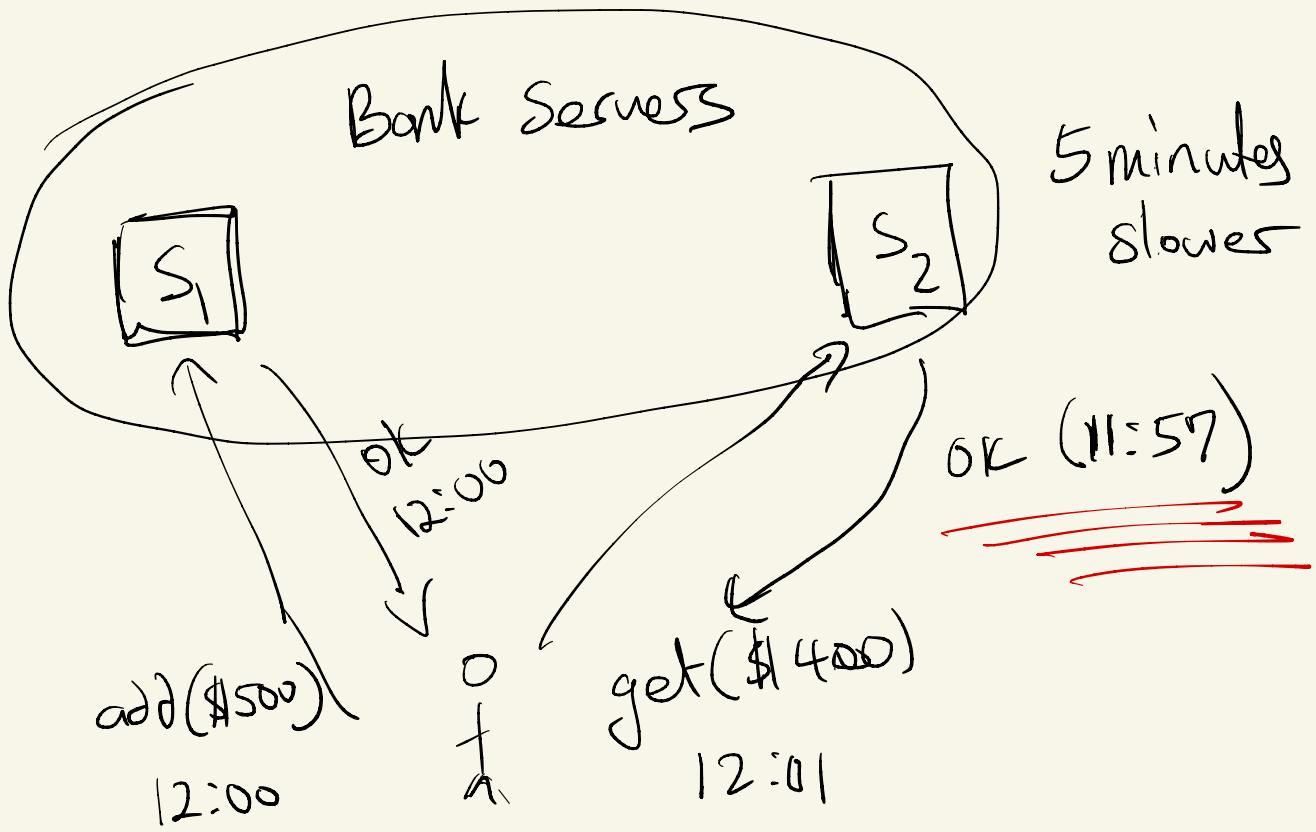
Summary

real time is bad (not monotonic)

FIFO / CAUSAL / TOTAL

Logical / Vector Clocks





People care about "time"

Distributed systems care about ordering :

add → get } what is the
or
get → add } right ordering