

CSCI 1380 : Day 12

Time + Global state



Last Class

Real time is bad (not monotonic & hard to synch)

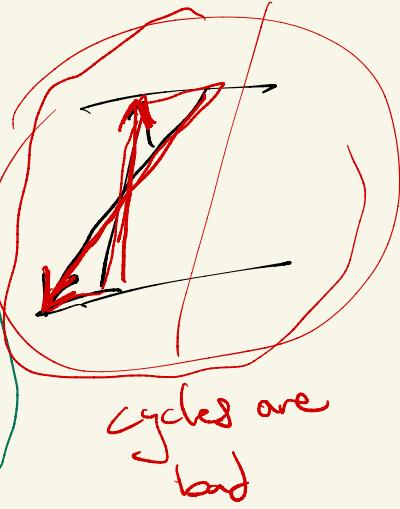
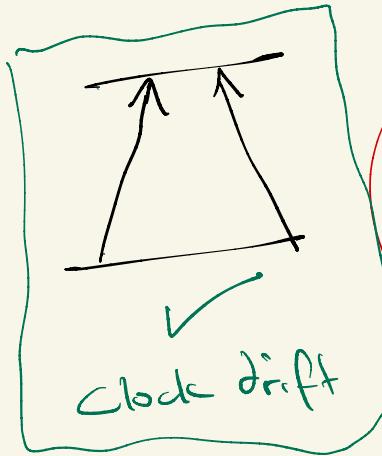
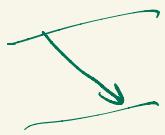
Systems want ordering that is monotonic

Ordering

Total] → logical clock

FIFO

Causal] → VC



Today

Global State

- * consistent snapshots
- * chandy Lamport

Replication (Passive/Active)

- * quorum ∨ consensus
- * Raft = passive

Global state (global snapshot)

Global state (distributed snapshot) = distributed backup of the cluster

Timebased snapshot

① some time T (2pm) all servers make backup (snapshot)

Potential problems

① 2pm all servers stop running

② 2pm is different on every server

manual snapshot

manually log into each server to create a snapshot

Problems

① takes too long

② inconsistency

Continuous snapshot

every N hours each server takes a snapshot
Don't pause during snapshot & don't try to coordinate

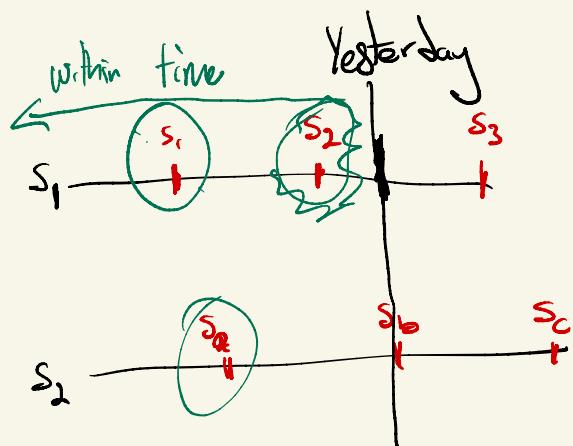
The trick is when you rollback (restore)

You look across servers for a consistent snapshot

① pick a time that you want to roll back to

② figure all snapshots within that time

③ is latest snapshot from s_1 consistent with s_2 if yes \Rightarrow done
else go back a snapshot



Problems

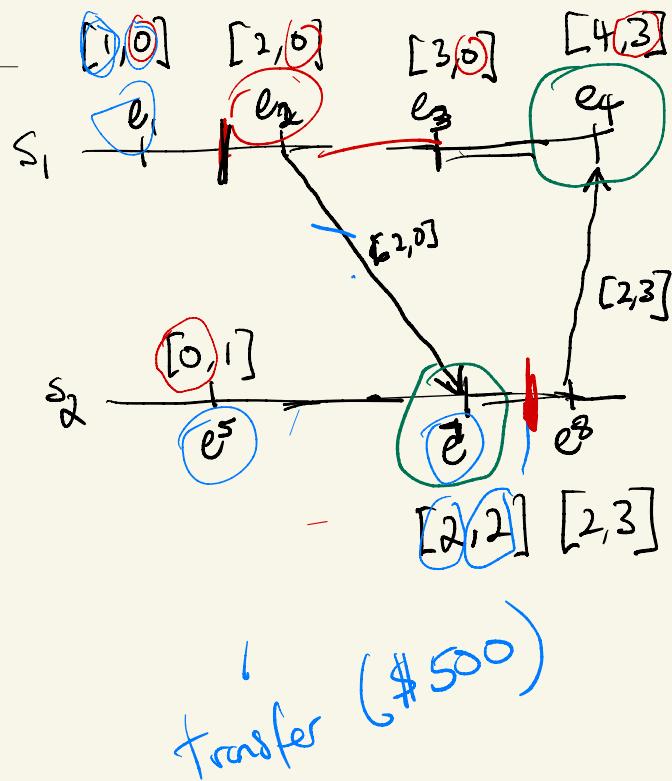
① still have some notion of real time

main problem ② Cascading rollback \rightsquigarrow the only consistent snapshot is ② $T=0$ when system starts
no guarantee to find a consistent snapshot

Consistent Snapshots

only way for consistency or to intertwine state on different servers is via msg

consistent = recv event also have then sent event in the snapshot



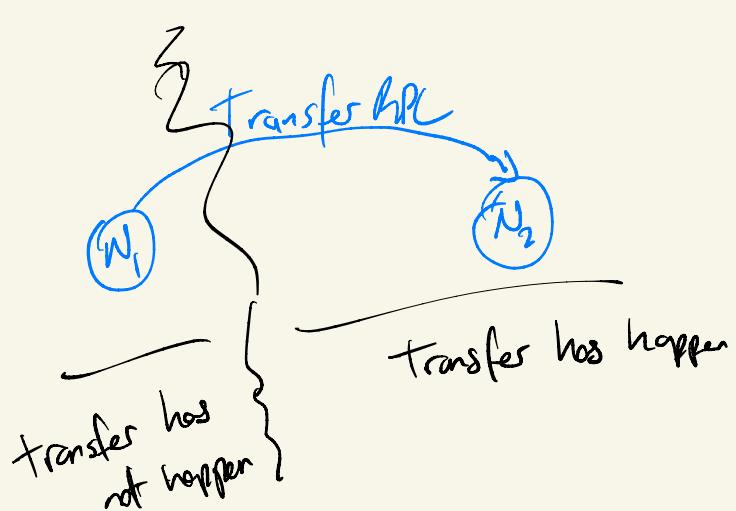
Pick random snapshot (e^5, e^2)



if pick e_i from S_i

$$VC_{e_i}[j] > VC_{S_j}[j]$$

$$j \neq i$$



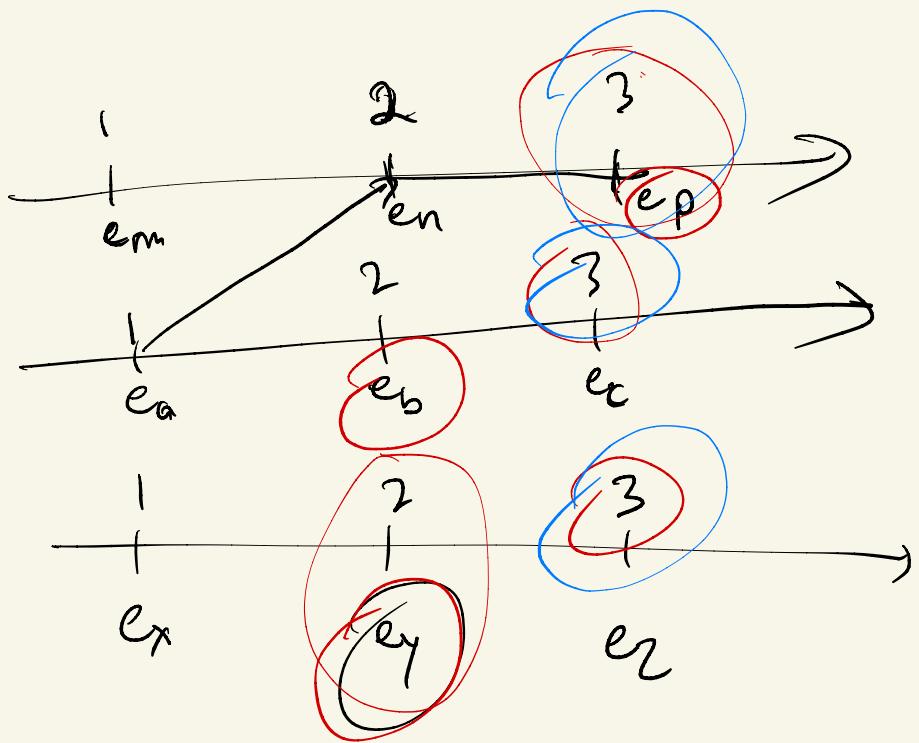
Detecting Inconsistent Snapshot

If an event at S_i has a vector clock VC_i

and an event at S_j ($S_i \neq S_j$) has a vector clock VC_j

then the snapshot is inconsistent if $VC_i[j] > VC_j[j]$

Thus S_i knows about more events at S_j than S_j does hence inconsistency

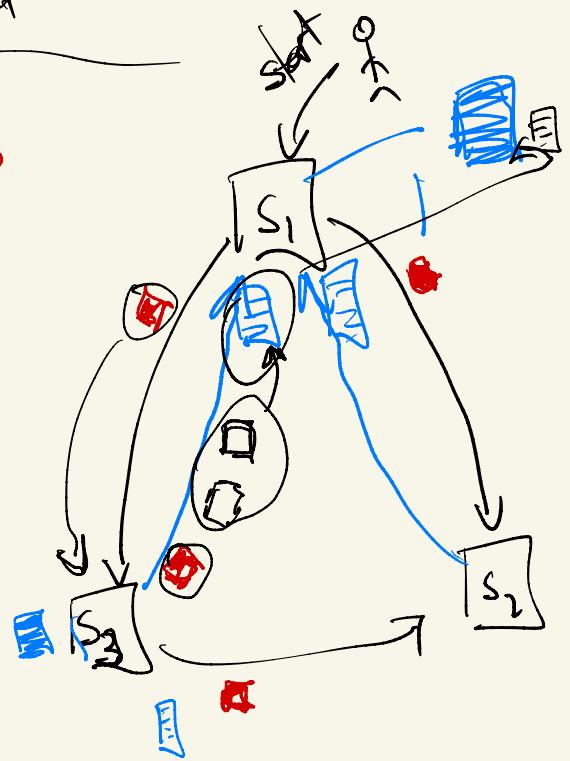


logical clocks don't capture causality
 & can't help infer msg b/w servers
 so they can't help with consistent snapshot

Chandy - Lamport

Rules

■ markers



① any node can initiate a snapshot

- * Sends a marker to all other nodes

- * Start watching for markers for all other nodes

- * Create a snapshot of state
 - \$ maintain queues for each server

② when you get your first marker

- * checkpoint of state

- * Send a new to every one

- * Maintain queues for every one but the person that sent you a marker

Marker \Rightarrow Signal
Msg

Start a
node has
done a
checkpoint

③ each server considers the checkpoint complete after receiving markers for everyone else

Assumptions

① Servers never crash (if they do the system stops)

② external entity starts the snapshot process

③ msg are delivered in order (FIFO order)

④ also msg

don't get lost

↳ if marker is lost
then protocol never ends

out of order
introduces inconsistency

Today

Global state (distribute
checkpoint/
snapshot)

- * approaches for capturing global state
- * determining consistent snapshots
 - using vector clock
- * Chandy - Lamport algorithm for distributed snapshot