

Distributed Transactions with Two-Phase Commit

Alvin Cheung

Fall 2022

Reading: R&G Chapter 22



Distributed vs. Parallel?



- Earlier we discussed Parallel DBMSs
 - Shared-memory
 - Shared-disk
 - Shared-nothing
- **Distributed** is basically shared-nothing parallel
 - With machines geographically distributed over network

What's Special About Distributed Computing?



- Inherited from shared-nothing parallel computation
 - Parallel computation
 - No shared memory/disk
- **Unreliable Networks**
 - Delay, reordering, loss of packets
- **Unsynchronized clocks**
 - Impossible to have perfect synchrony
- Partial failure: can't know what's up, what's down

Distributed Database Systems



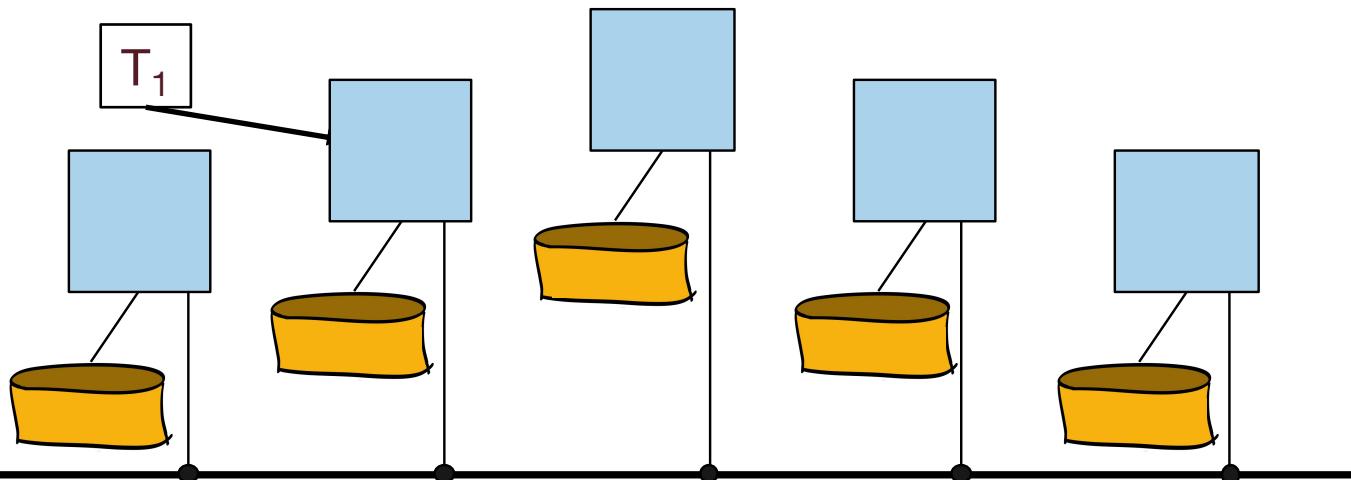
- DBMS is an influential special case of distributed computing
 - The trickiest part of distributed computing is state, i.e., Data
 - Transactions provide an influential model for concurrency/parallelism
 - DBMSs worried about fault handling early on
- But also a special-case as not all distributed programs are written transactionally
 - And if not, database techniques may not apply
- Many of today's most complex distributed systems are databases
 - Cloud SQL databases like Google BigQuery, AWS Aurora, Azure SQL
 - NoSQL databases like DynamoDB, Cassandra, MongoDB, Couchbase...
- We'll focus on transactional concurrency control and recovery
 - You already know many lessons of distributed query processing

Distributed Locking

Distributed Concurrency Control



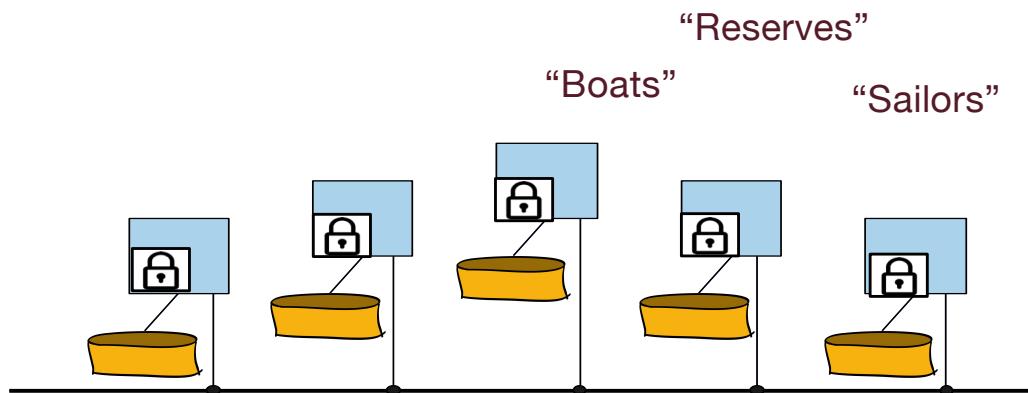
- Consider a **shared-nothing distributed DBMS**
- For today, assume partitioning but no replication of data
- **Each transaction arrives at some node:** *some machine's memory*
 - The “coordinator” for the transaction
 - Can be designated or assigned on the fly



Where is the Lock Table



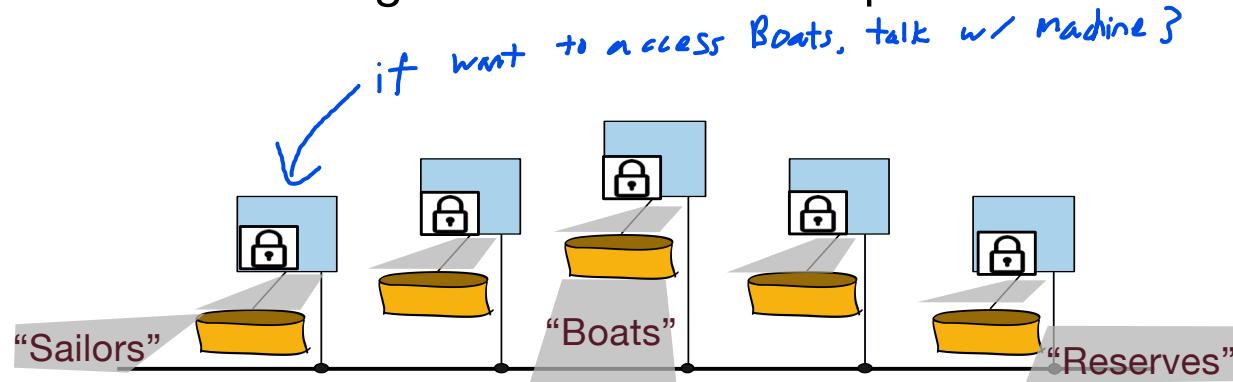
- Typical design: Locks partitioned with the data
 - Independent: each node manages its own lock table
 - Works for objects that fit on one node (pages, tuples)
- For coarser-grained locks, assign a “home” node
 - Object being locked (table, DB) exists across nodes



Where is the Lock Table, Pt 2



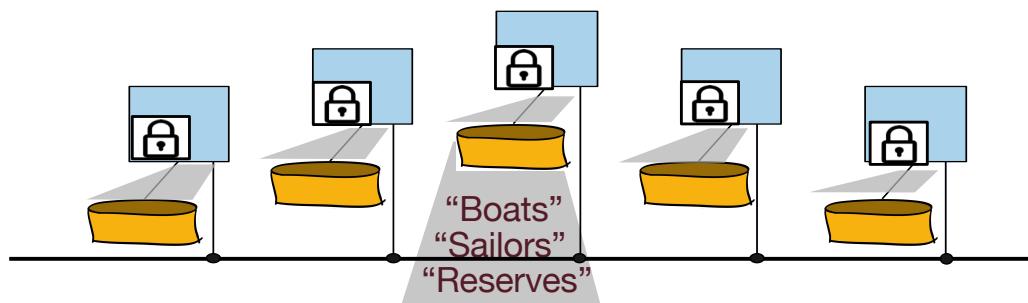
- Typical design: Locks partitioned with the data
 - Independent: each node manages its own lock table
 - Works for objects that fit on one node (pages, tuples)
- For coarser-grained locks, assign a “home” node
 - Object being locked (table, DB) exists across nodes
 - These coarse-grained locks can be partitioned across nodes



Where is the Lock Table, Pt 3



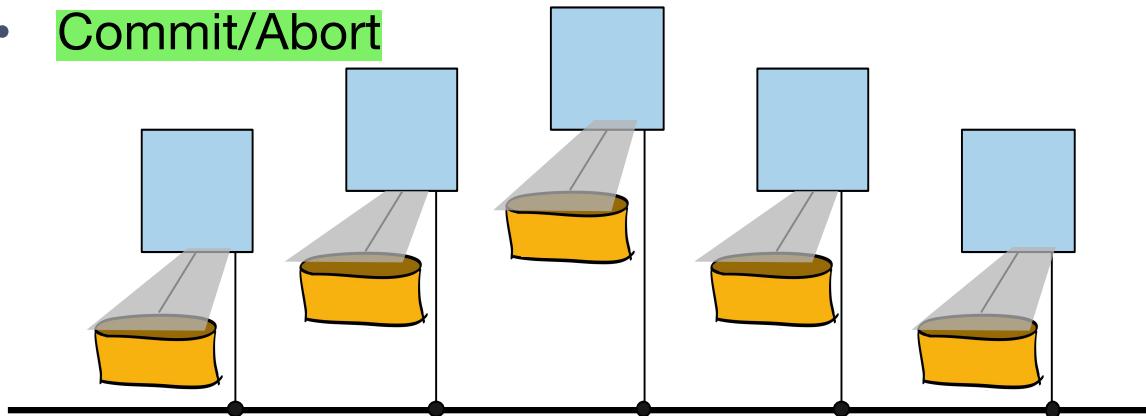
- Typical design: Locks partitioned with the data
 - Independent: each node manages its own lock table
 - Works for objects that fit on one node (pages, tuples)
- For coarser-grained locks, assign a “home” node
 - Object being locked (table, DB) exists across nodes
 - These coarse-grained locks can be partitioned across nodes
 - Or centralized at a master node



Ignore global coarse-grained locks for a moment...



- Every node does its own locking
 - Clean and efficient
 - Nicely generalizes the single-node setting
- “Global” issues remain:
 - Deadlock
 - Commit/Abort

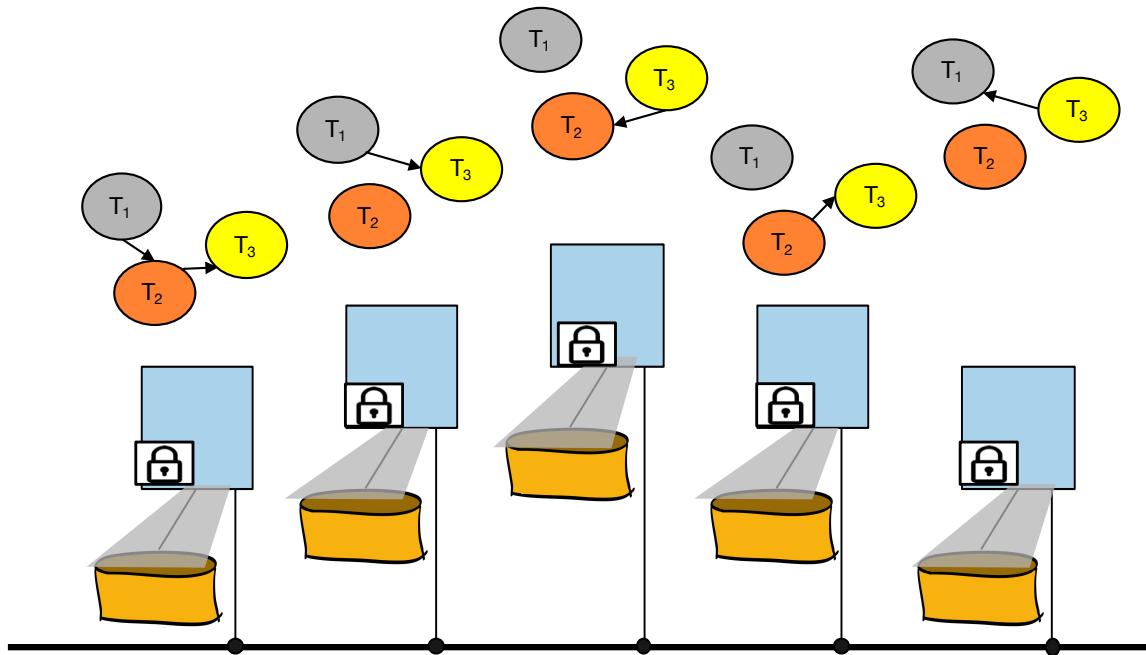


Distributed Deadlock Detection

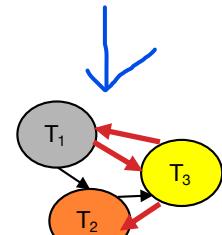
What Could Go Wrong? #1



- Deadlock detection via waits-for graphs



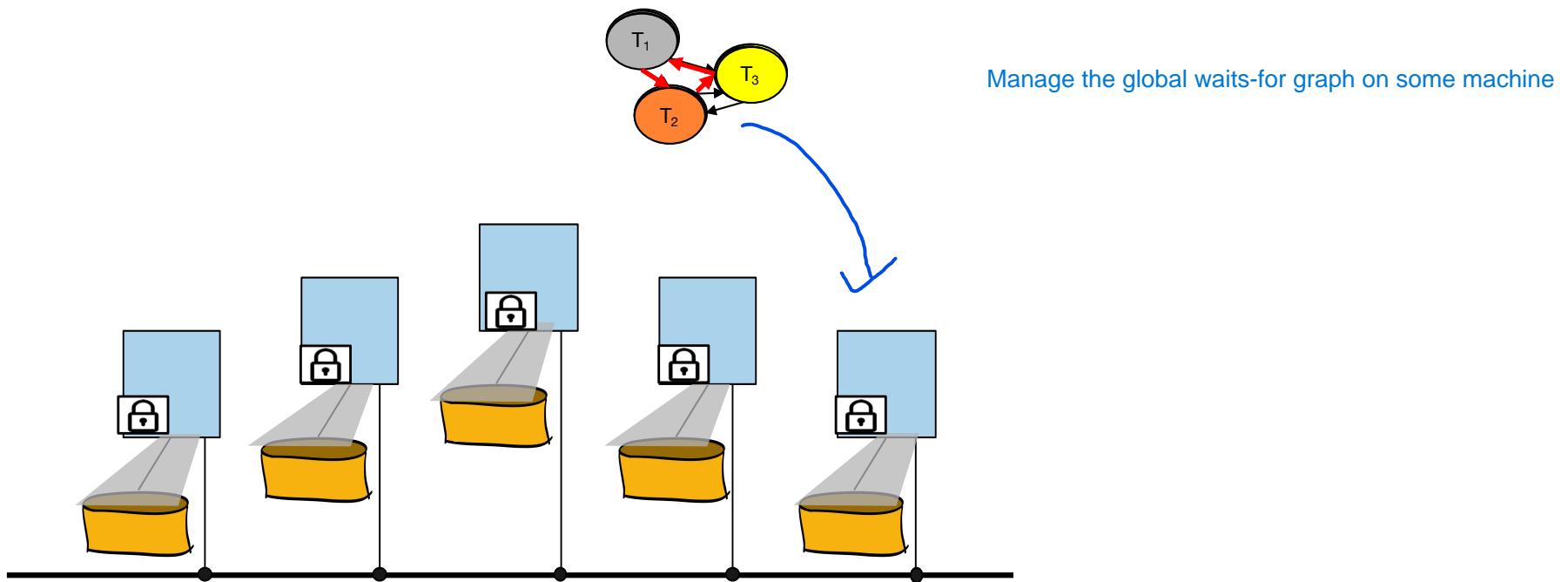
Each machine doesn't have a cycle, but there is a global cycle



What Could Go Wrong? #1 Part 2



- Deadlock detection via waits-for graphs
 - Easy fix: periodically union at designated coordinator



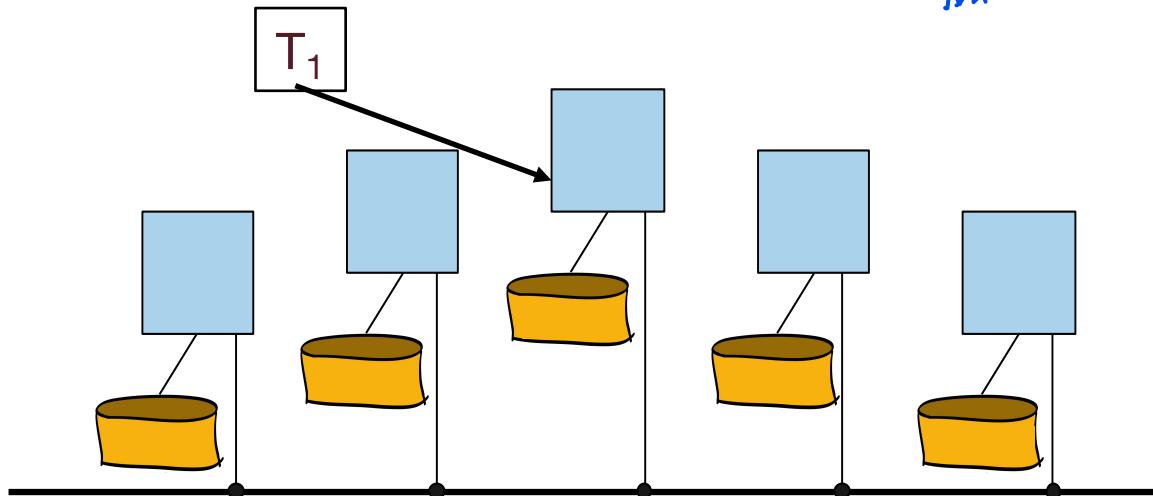
Distributed Commit: 2PC

Strawman: Coordinator makes Decision



- Recall that every txn has a coordinator node
- Coordinator decides if the txn is going to commit or abort.
- Lets all the other nodes know.
- Q: Why is this scheme problematic?

so all machines
will rollback the
txn

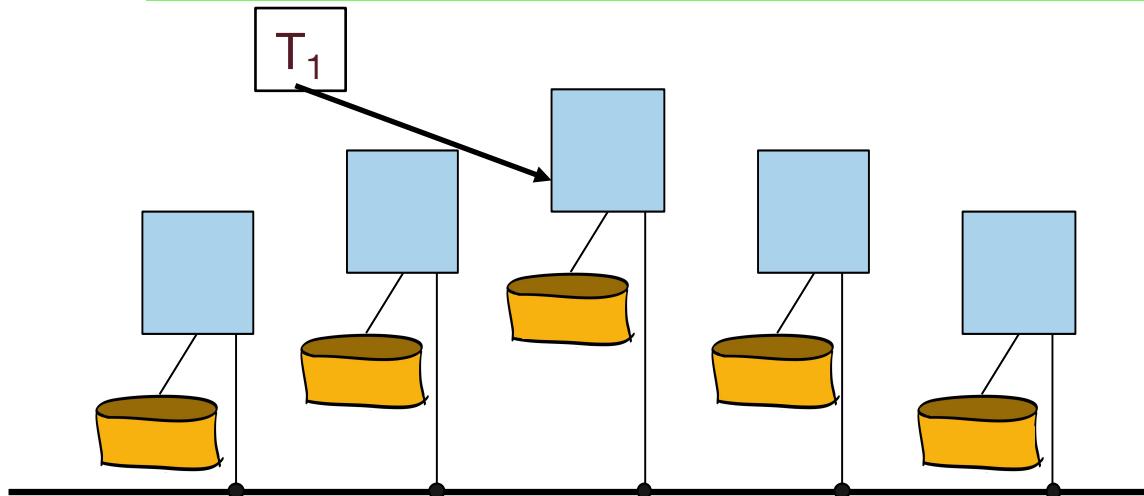
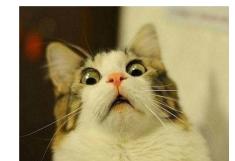


Strawman: Coordinator makes Decision



- Recall that every txn has a coordinator node
- Coordinator decides if the txn is going to commit or abort.
- Lets all the other nodes know.
- Q: Why is this scheme problematic?
 - Among other things, one of the nodes may want to abort, even if the coordinator wants to commit
 - Some nodes may actually be down (so any txn touching their data shouldn't proceed)

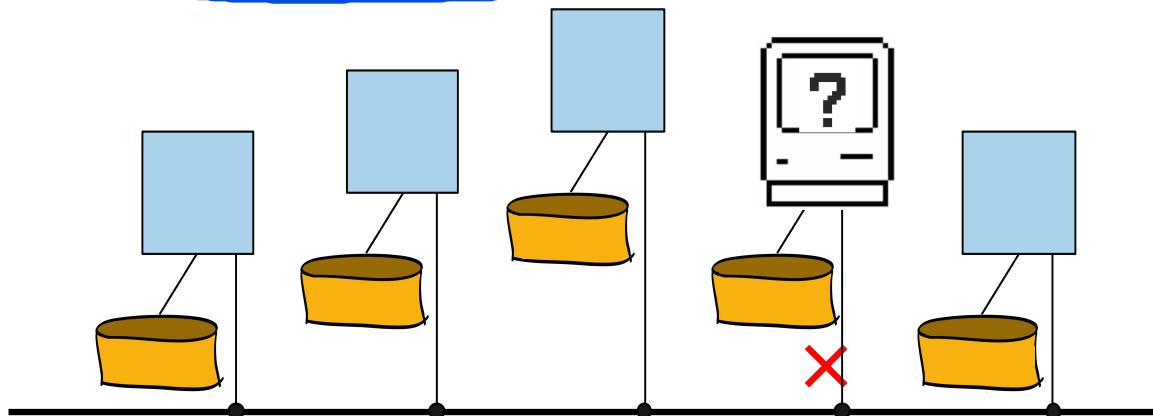
So do we commit or abort?



In General, What Could Go Wrong? #2



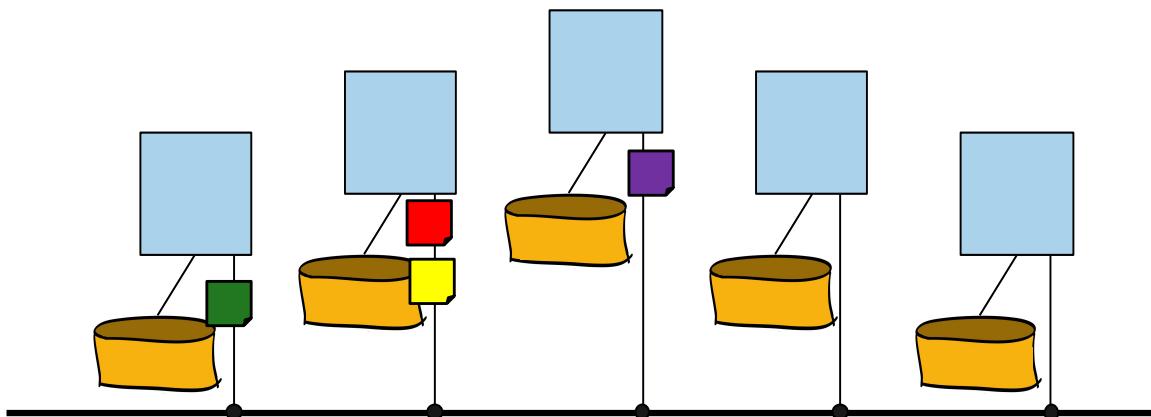
- Failures/Delays: Nodes
 - Commit? Abort?
 - If we haven't heard from a node, we don't know if it is alive or dead.
 - The decision may hinge on this node (imagine a FK violation at that node)
 - When the node comes back, how does it recover in a world that moved forward?



What Could Go Wrong? #2, Part 2



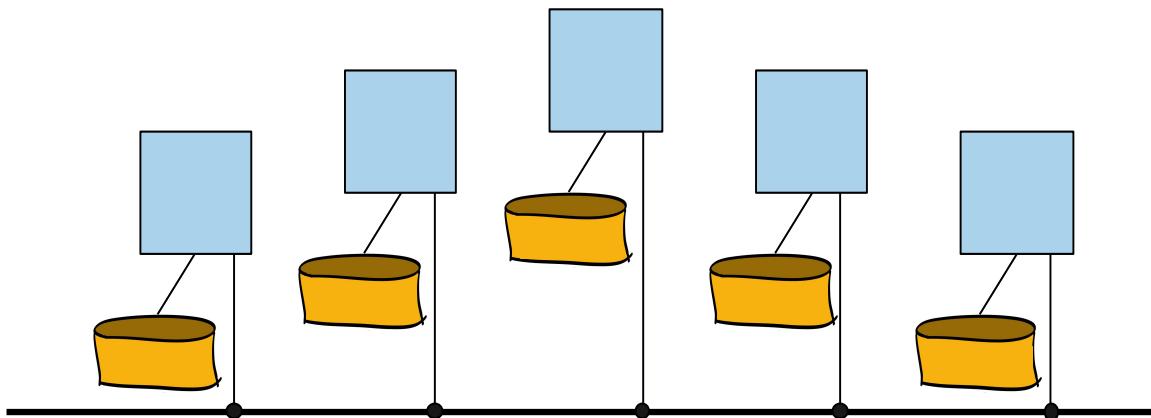
- Failures/Delays: Nodes
- Failures/Delays: Messages
 - Non-deterministic reordering per channel, interleaving across channels
 - “Lost” (very delayed) messages
 - How long should we wait for this?



What Could Go Wrong? #2, Part 3



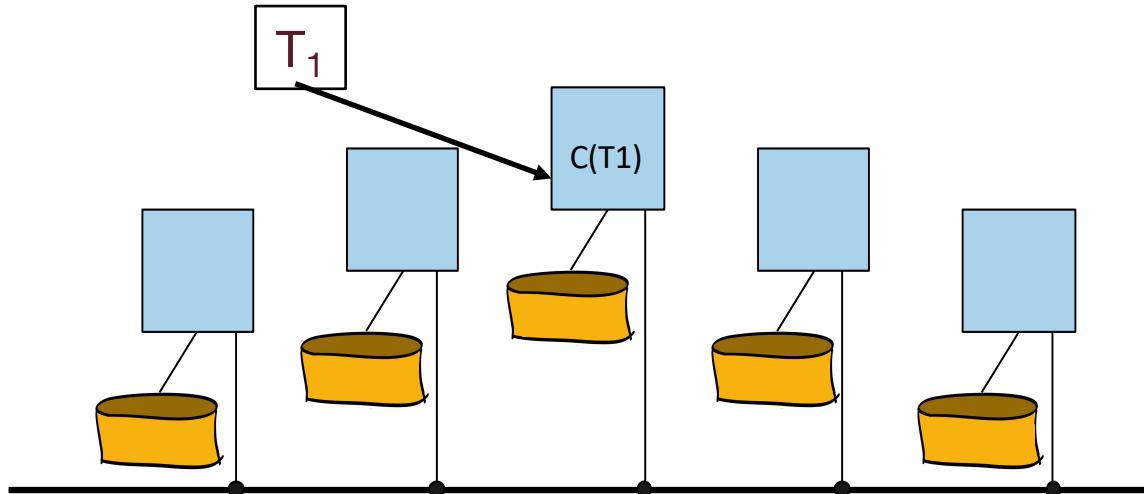
- Failures/Delays: Nodes
- Failures/Delays: Messages
 - Non-deterministic reordering per channel, interleaving across channels
 - “Lost” (very delayed) messages
- Given this, how do all nodes agree on Commit vs. Abort?



Basic Idea: Distributed Voting



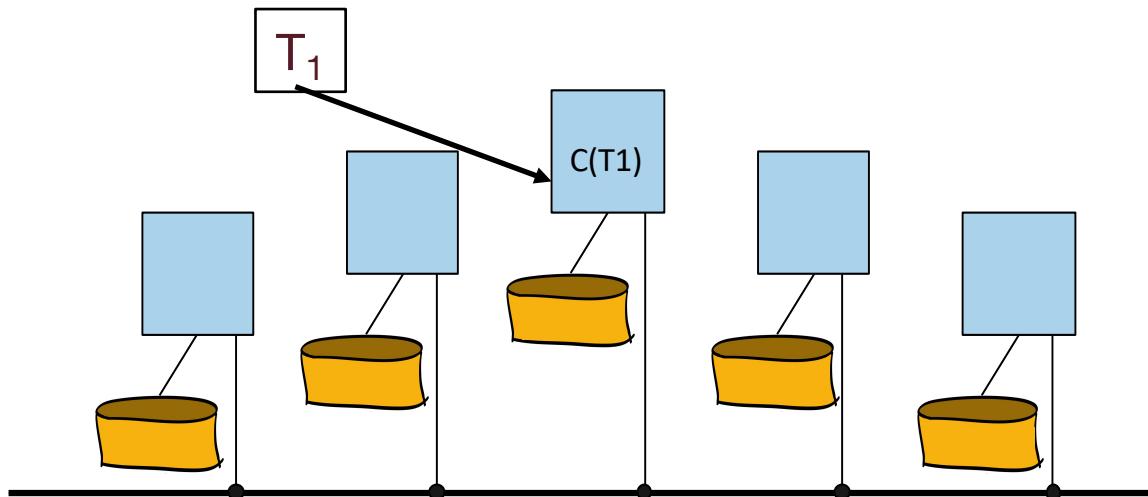
- Vote for Commitment
 - How many votes does a commit need to win?
 - Any single node could observe a problem (e.g., deadlock, constraint violation)
 - Hence must be unanimous. *in order to win commit*



Distributed voting? How?



- How do we implement distributed voting?!
 - In the face of message/node failure/delay?



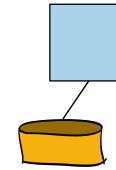
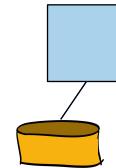
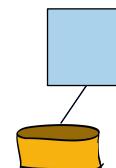
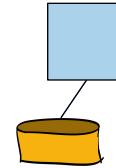
2-Phase Commit



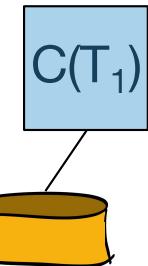
- A.k.a. 2PC. (Not to be confused with 2PL!)
- Like a wedding ceremony!
- **Phase 1:** “do you take this person...”
 - Coordinator tells participants to “prepare”
 - Participants respond with yes/no votes
 - Unanimity required for yes!
- **Phase 2:** “I now pronounce you...”
 - Coordinator disseminates result of the vote
- Need to do some logging for failure handling....

2-Phase Commit, Part 1

- **Phase 1:**
 - Coordinator tells participants to “prepare”
make sure all writes persist aka all operations done
 - Participants respond with yes/no votes
 - Unanimity required for commit!
- Phase 2:
 - Coordinator disseminates result of the vote
 - Participants respond with Ack



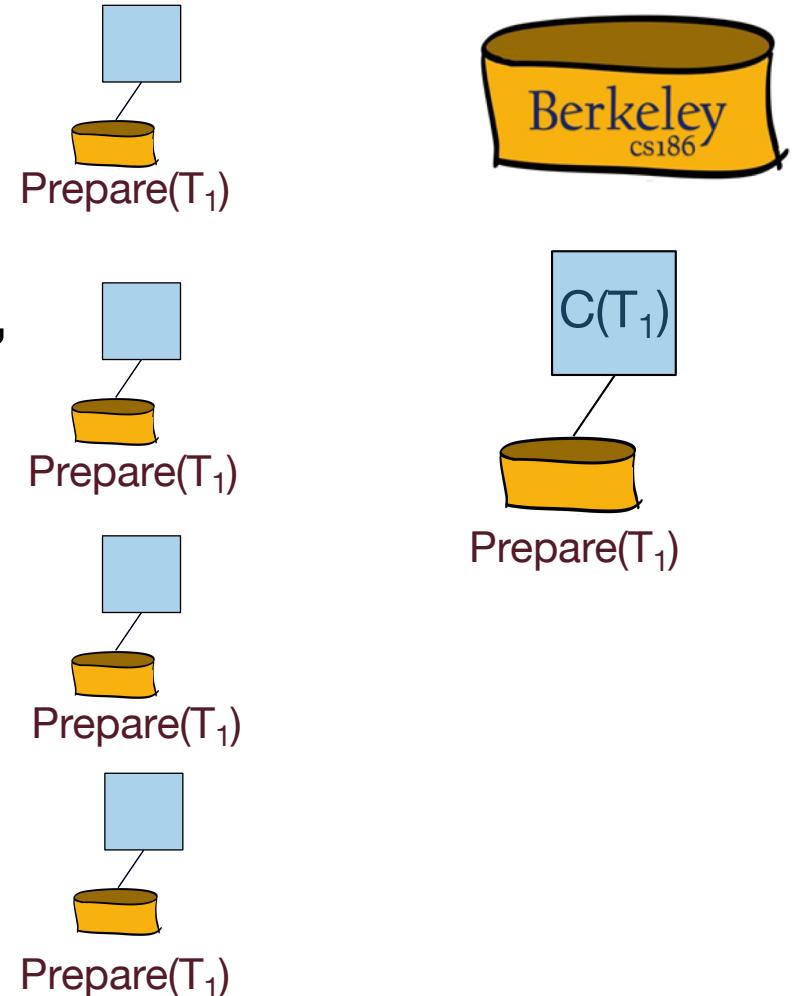
COORD



Prepare(T_1)

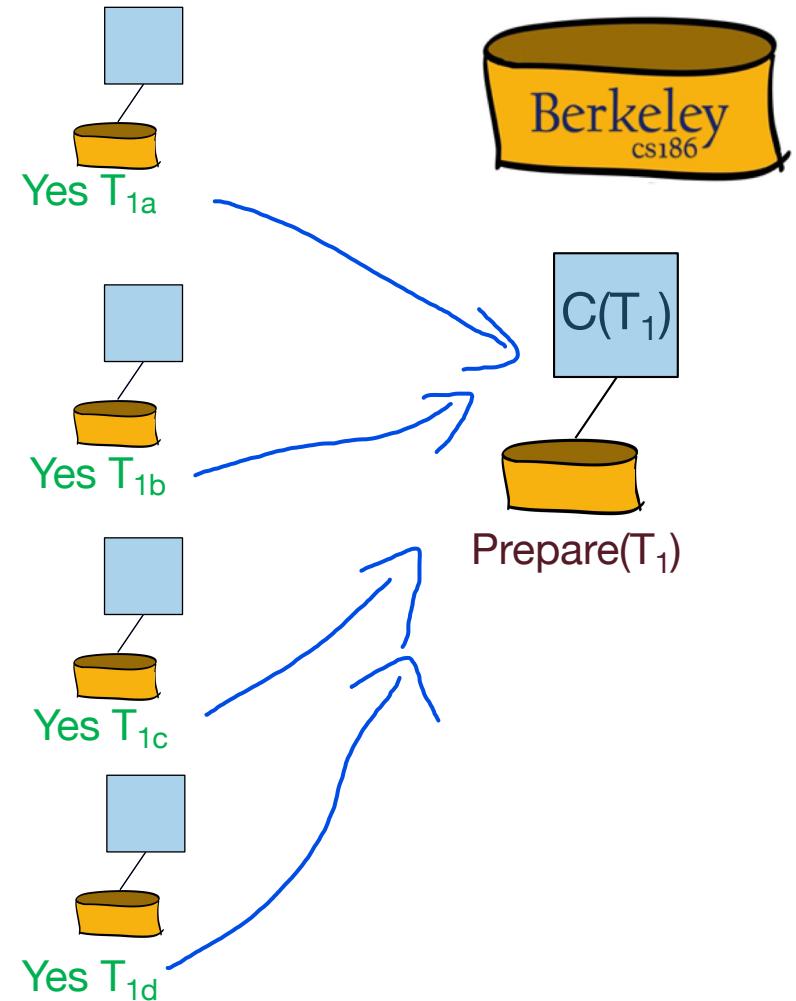
2-Phase Commit, Part 2

- **Phase 1:**
 - Coordinator tells participants to “prepare”
 - Participants respond with yes/no votes
 - Unanimity required for commit!
- Phase 2:
 - Coordinator disseminates result of the vote
 - Participants respond with Ack



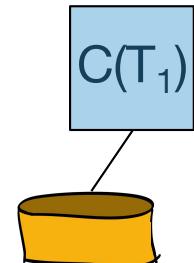
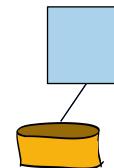
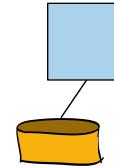
2-Phase Commit, Part 3

- **Phase 1:**
 - Coordinator tells participants to “prepare”
 - **Participants respond with yes/no votes**
 - Unanimity required for commit!
- Phase 2:
 - Coordinator disseminates result of the vote
 - Participants respond with Ack

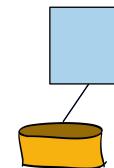
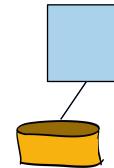


2-Phase Commit, Part 4

- **Phase 1:**
 - Coordinator tells participants to “prepare”
 - **Participants respond with yes/no votes**
 - **Unanimity required for commit!**
- Phase 2:
 - Coordinator disseminates result of the vote
 - Participants respond with Ack

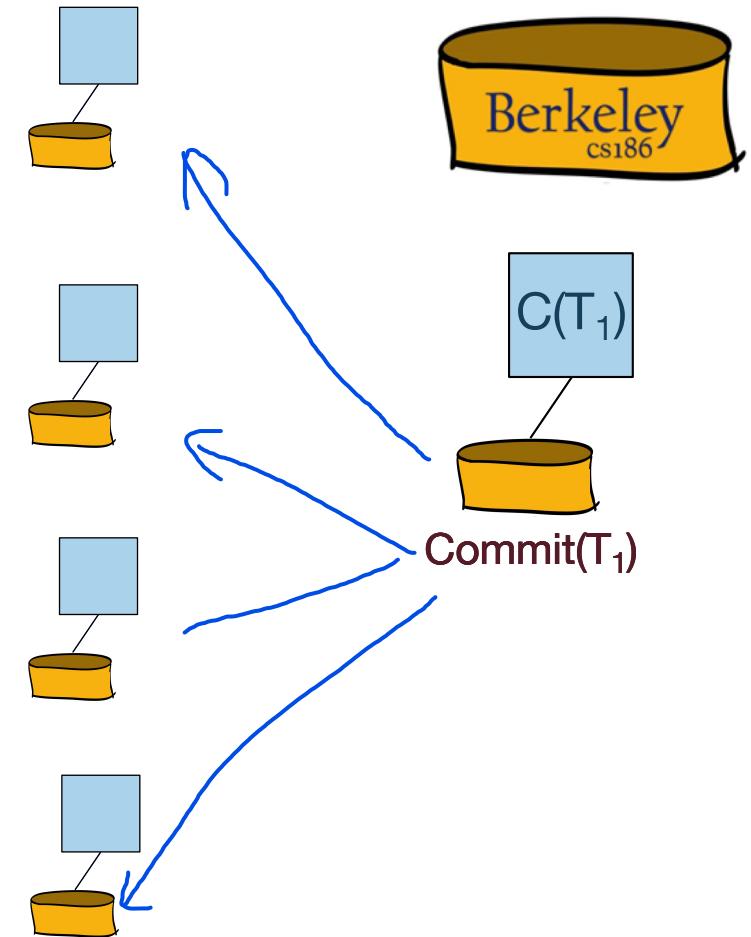


Yes T_{1b}



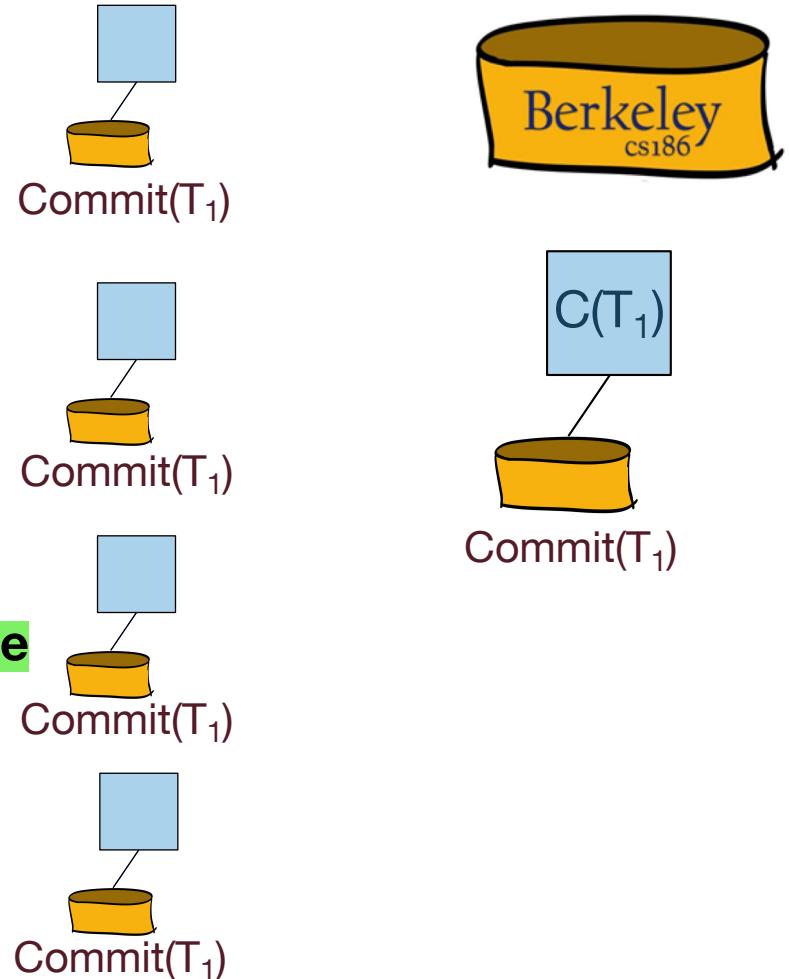
2-Phase Commit, Part 5

- Phase 1:
 - Coordinator tells participants to “prepare”
 - Participants respond with yes/no votes
 - Unanimity required for commit!
- **Phase 2:**
 - **Coordinator disseminates result of the vote**
 - Participants respond with Ack



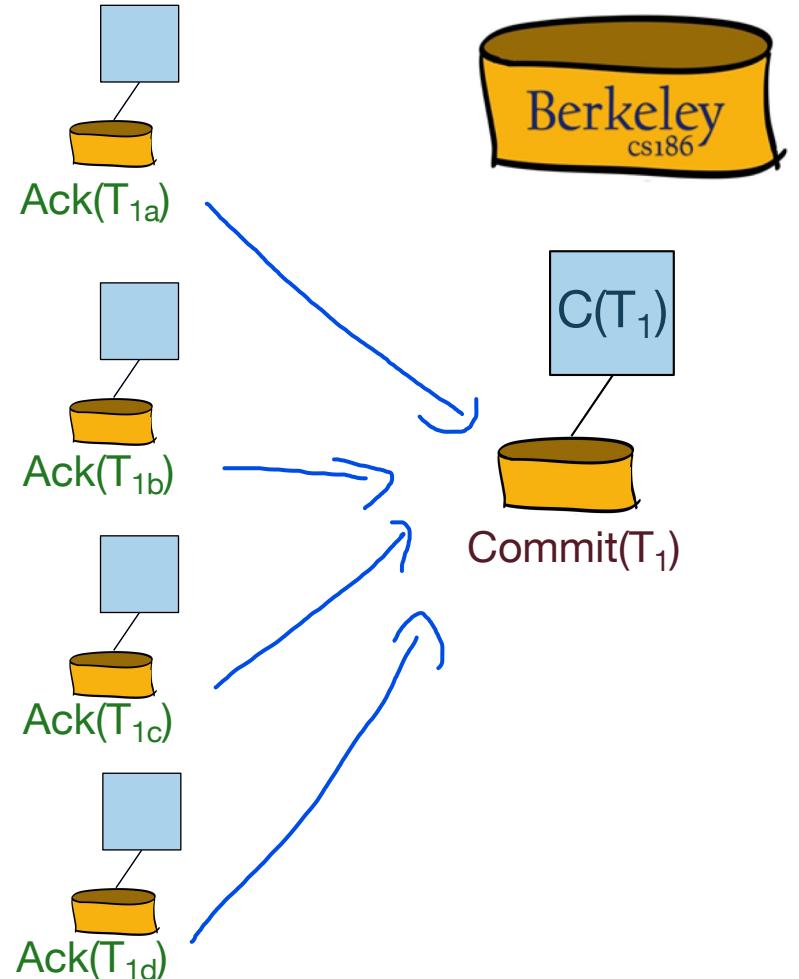
2-Phase Commit, Part 6

- Phase 1:
 - Coordinator tells participants to “prepare”
 - Participants respond with yes/no votes
 - Unanimity required for commit!
- Phase 2:
 - **Coordinator disseminates result of the vote**
 - Participants respond with Ack



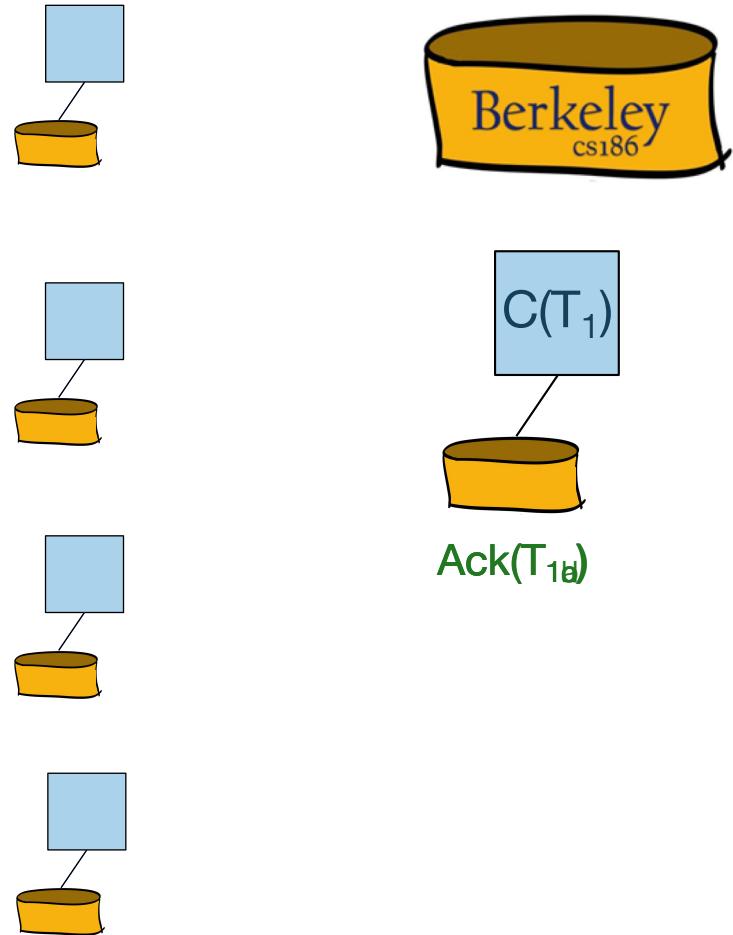
2-Phase Commit, Part 7

- Phase 1:
 - Coordinator tells participants to “prepare”
 - Participants respond with yes/no votes
 - Unanimity required for commit!
- **Phase 2:**
 - Coordinator disseminates result of the vote
 - **Participants respond with Ack**



2-Phase Commit, Part 8

- Phase 1:
 - Coordinator tells participants to “prepare”
 - Participants respond with yes/no votes
 - Unanimity required for commit!
- **Phase 2:**
 - Coordinator disseminates result of the vote
 - **Participants respond with Ack**



When the coordinator receives messages from all participants, txn is complete

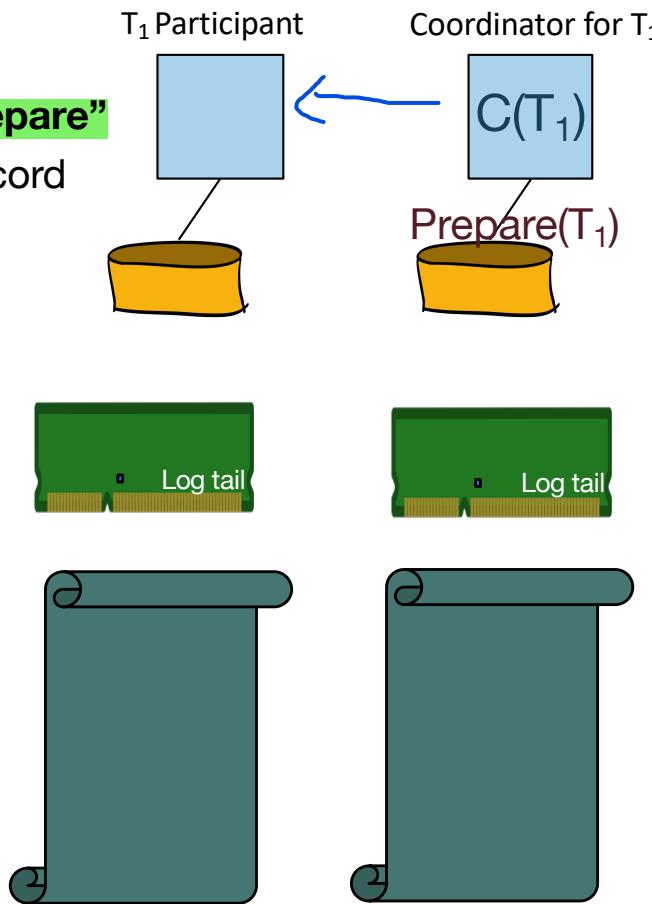
Recovery and Locking

Phase 1

One More Time, With Logging



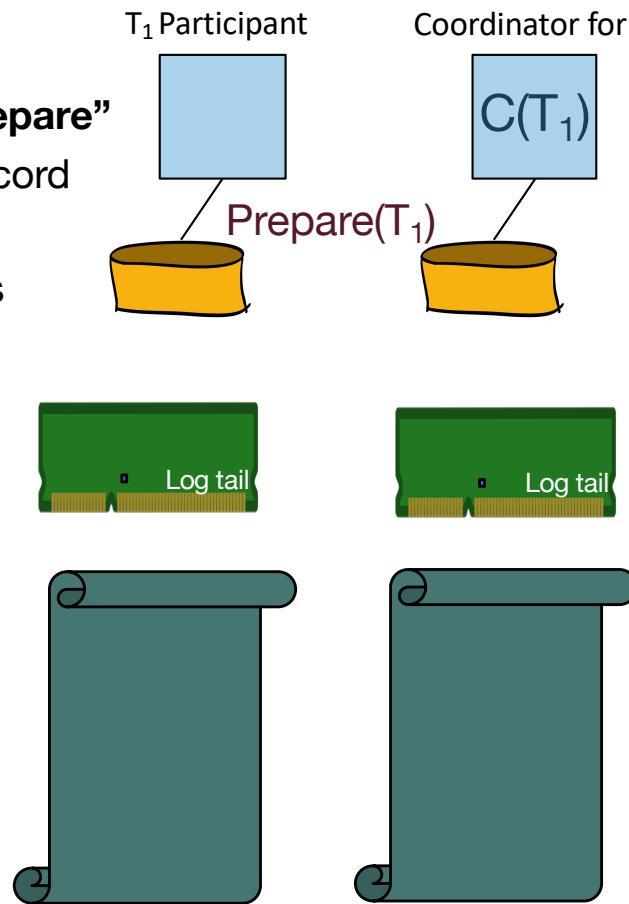
- Phase 1
- Coordinator tells participants to “prepare”
- Participants generate prepare/abort record
- Participants flush prepare/abort record
- Participants respond with yes/no votes
- Coordinator generates commit record
- Coordinator flushes commit record



One More Time, With Logging, Part 2



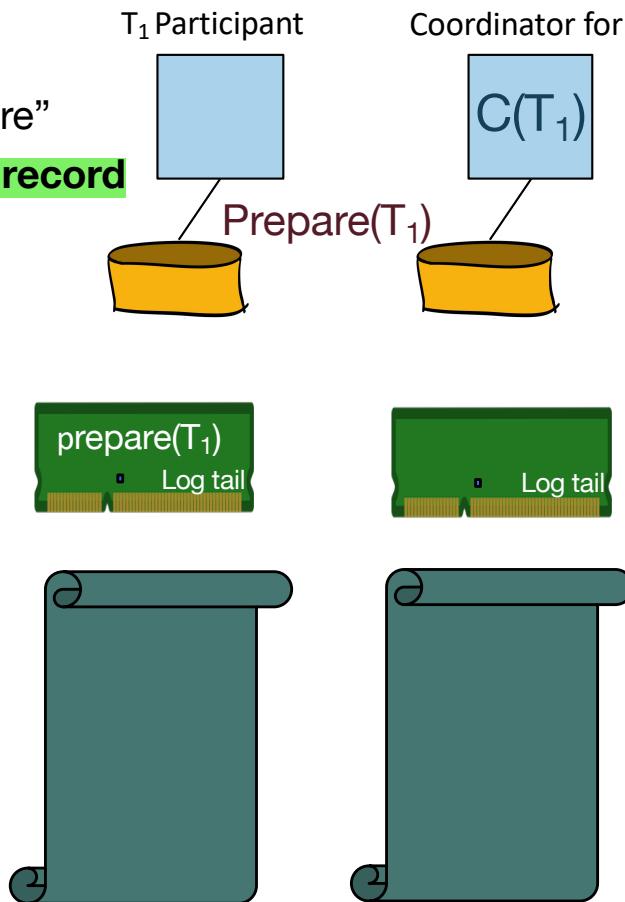
- **Phase 1**
- **Coordinator tells participants to “prepare”**
- Participants generate prepare/abort record
- Participants flush prepare/abort record
- Participants respond with yes/no votes
- Coordinator generates commit record
- Coordinator flushes commit record



One More Time, With Logging, Part 3



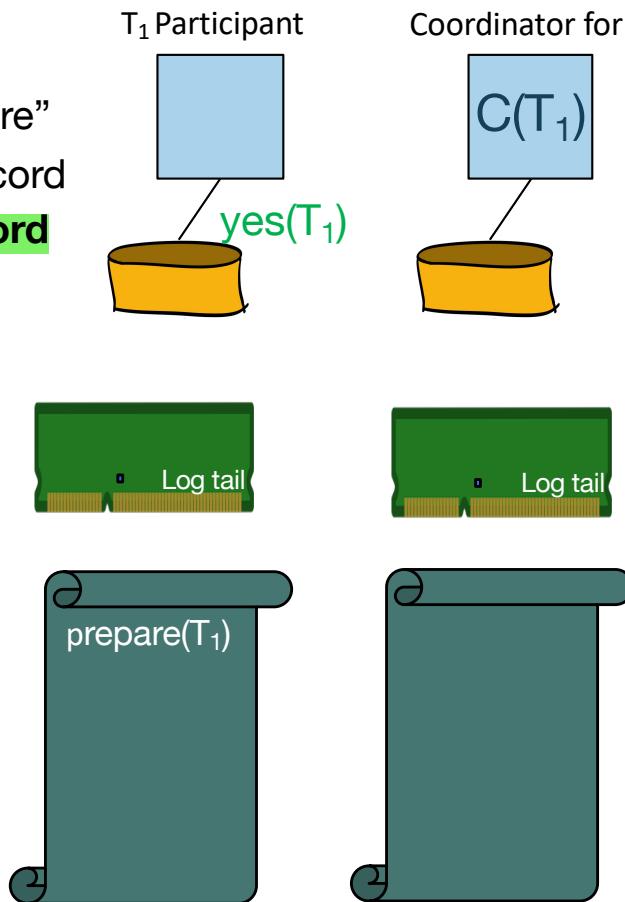
- **Phase 1**
- Coordinator tells participants to “prepare”
- **Participants generate prepare/abort record**
- Participants flush prepare/abort record
- Participants respond with yes/no votes
- Coordinator generates commit record
- Coordinator flushes commit record



One More Time, With Logging, Part 4



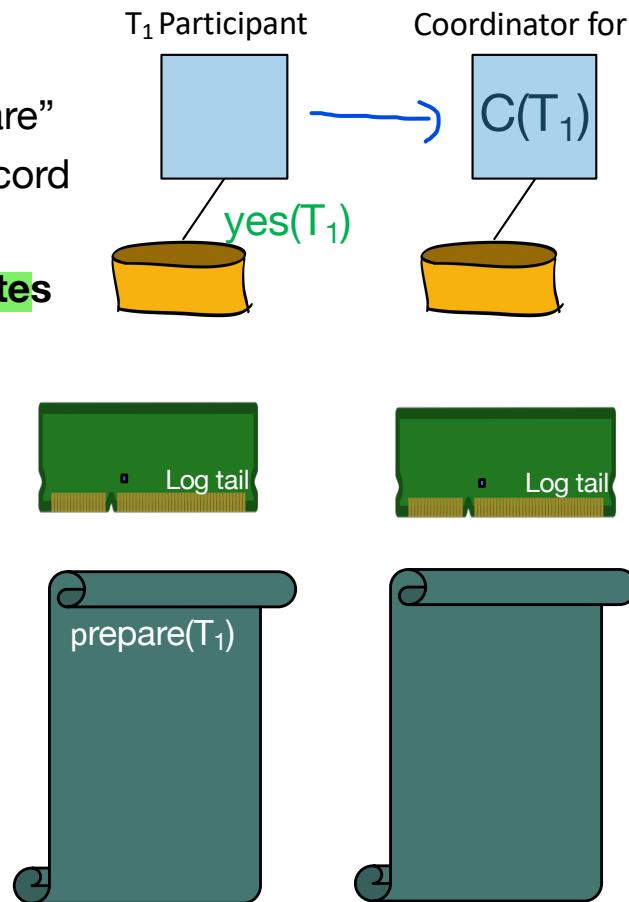
- **Phase 1**
- Coordinator tells participants to “prepare”
- Participants generate prepare/abort record
- **Participants flush prepare/abort record**
- Participants respond with yes/no votes
- Coordinator generates commit record
- Coordinator flushes commit record



One More Time, With Logging, Part 5



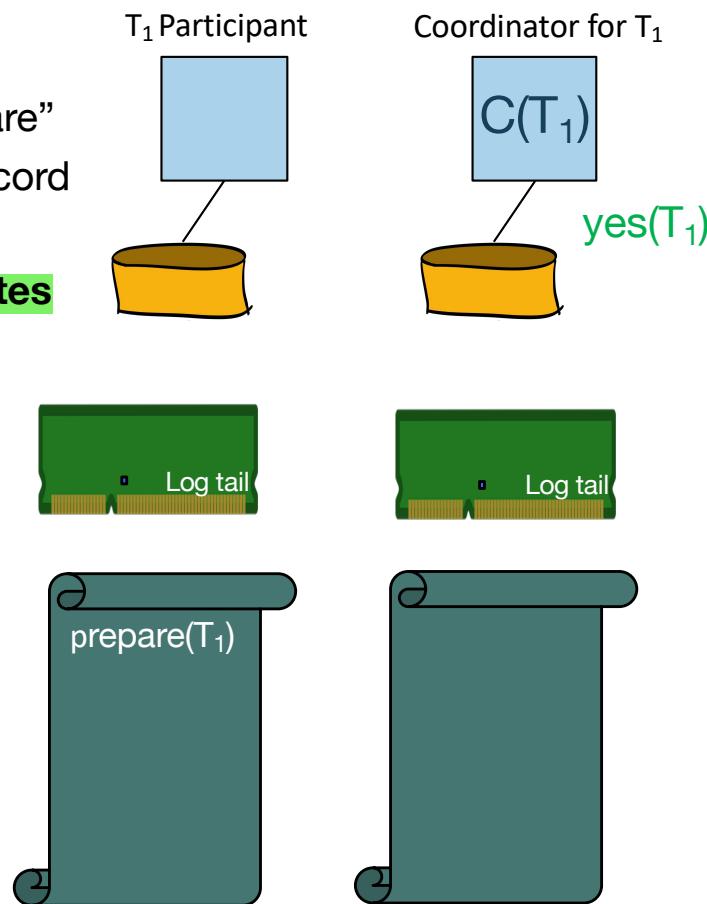
- **Phase 1**
- Coordinator tells participants to “prepare”
- Participants generate prepare/abort record
- Participants flush prepare/abort record
- **Participants respond with yes/no votes**
- Coordinator generates commit record
- Coordinator flushes commit record



One More Time, With Logging, Part 6



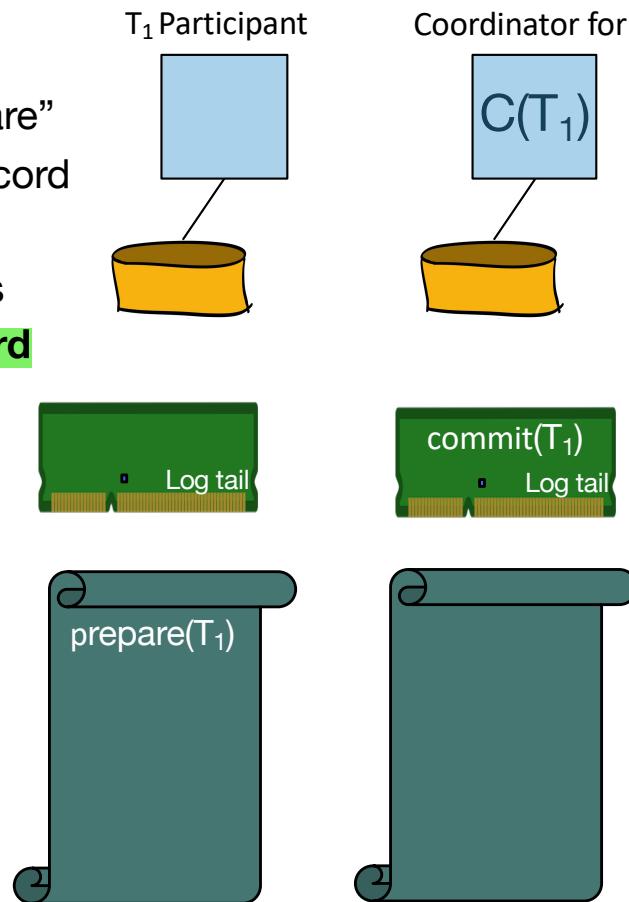
- **Phase 1**
- Coordinator tells participants to “prepare”
- Participants generate prepare/abort record
- Participants flush prepare/abort record
- **Participants respond with yes/no votes**
- Coordinator generates commit record
- Coordinator flushes commit record



One More Time, With Logging, Part 7



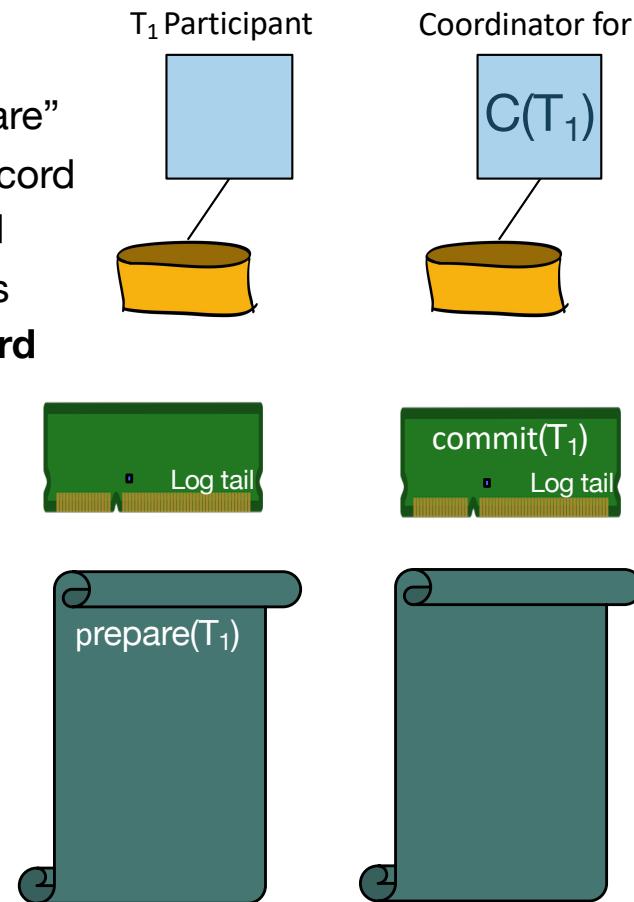
- **Phase 1**
- Coordinator tells participants to “prepare”
- Participants generate prepare/abort record
- Participants flush prepare/abort record
- Participants respond with yes/no votes
- **Coordinator generates commit record**
- Coordinator flushes commit record



One More Time, With Logging, Part 8



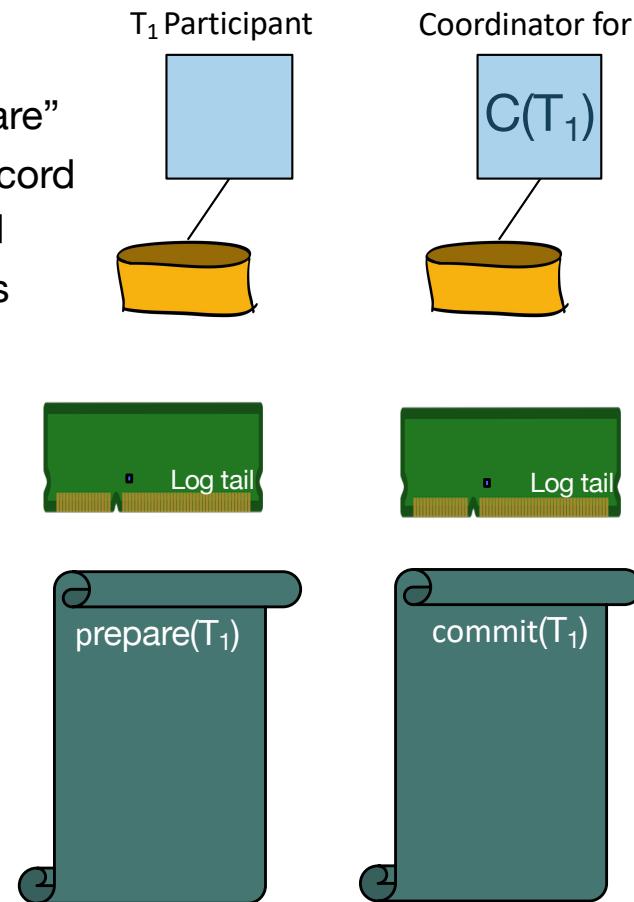
- **Phase 1**
- Coordinator tells participants to “prepare”
- Participants generate prepare/abort record
- Participants flush prepare/abort record
- Participants respond with yes/no votes
- **Coordinator generates commit record**
- Coordinator flushes commit record



One More Time, With Logging, Part 9



- **Phase 1**
- Coordinator tells participants to “prepare”
- Participants generate prepare/abort record
- Participants flush prepare/abort record
- Participants respond with yes/no votes
- Coordinator generates commit record
- **Coordinator flushes commit record**

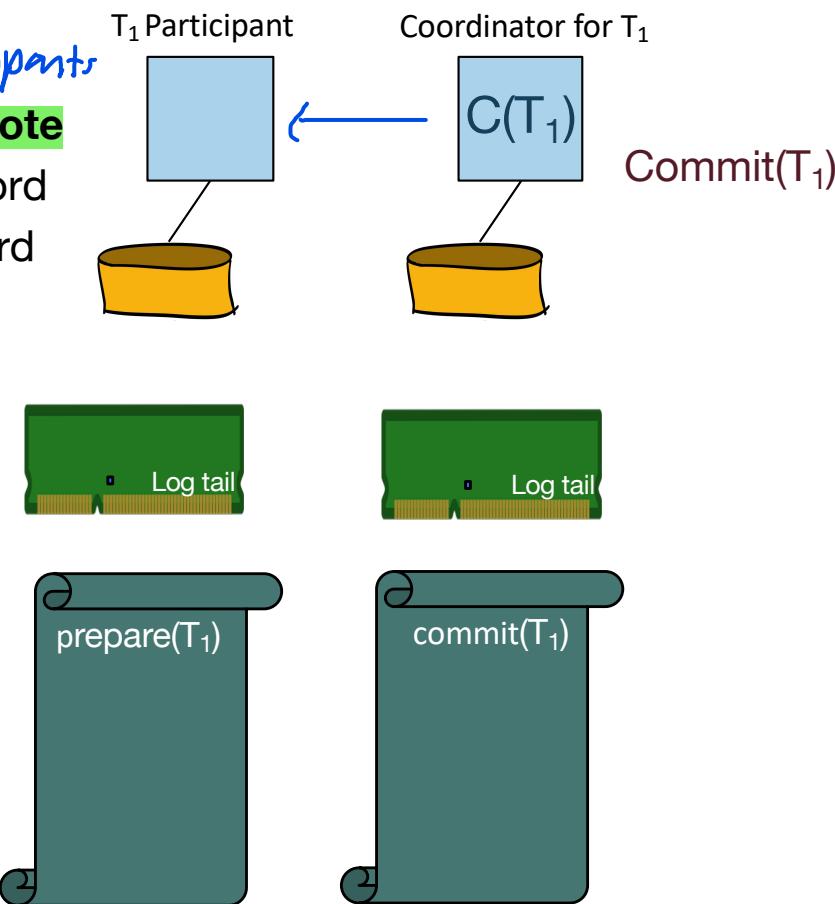


Phase 2

One More Time, With Logging, Part 10



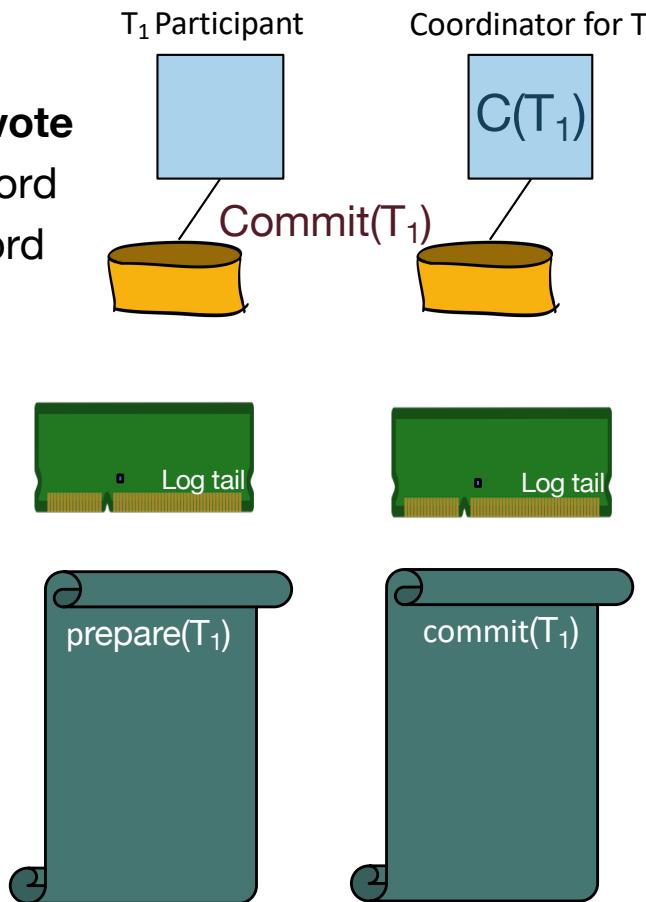
- Phase 2: *to participants*
- Coordinator broadcasts result of vote
- Participants make commit/abort record
- Participants flush commit/abort record
- Participants respond with Ack
- Coordinator generates end record
- Coordinator flushes end record



One More Time, With Logging, Part 11



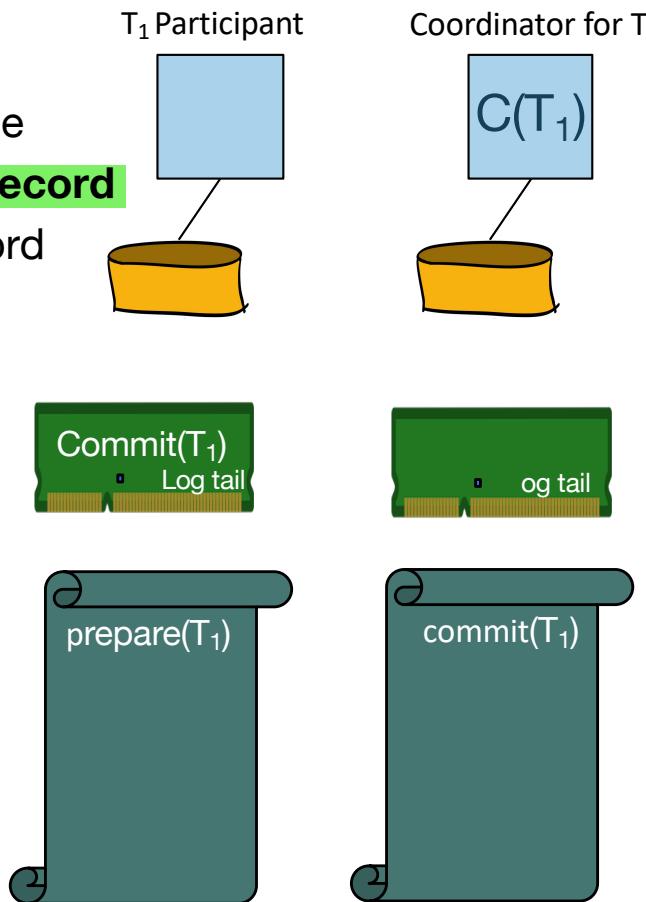
- **Phase 2:**
- **Coordinator broadcasts result of vote**
- Participants make commit/abort record
- Participants flush commit/abort record
- Participants respond with Ack
- Coordinator generates end record
- Coordinator flushes end record



One More Time, With Logging, Part 12



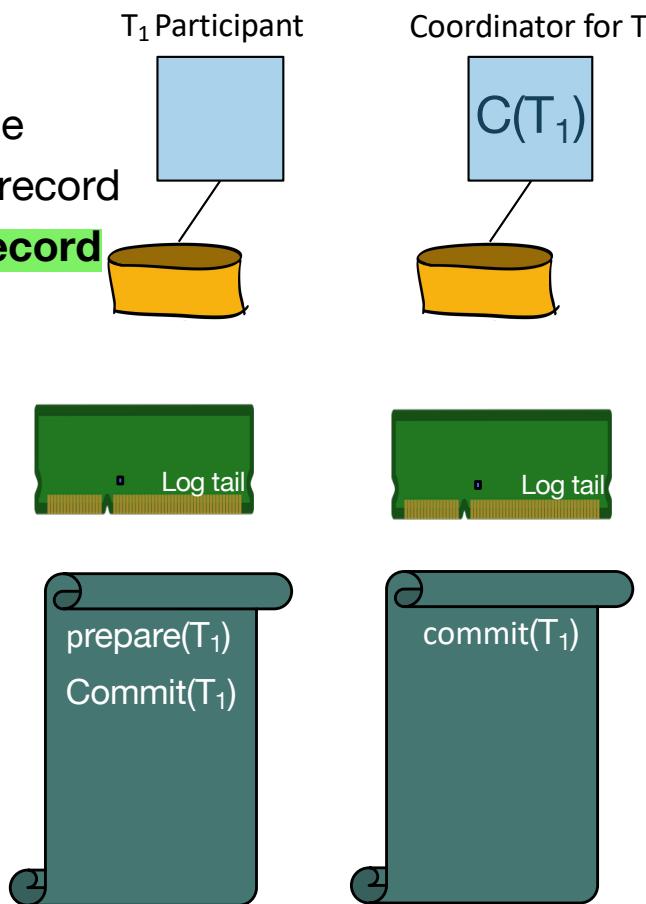
- **Phase 2:**
- Coordinator broadcasts result of vote
- **Participants make commit/abort record**
- Participants flush commit/abort record
- Participants respond with Ack
- Coordinator generates end record
- Coordinator flushes end record



One More Time, With Logging, Part 13



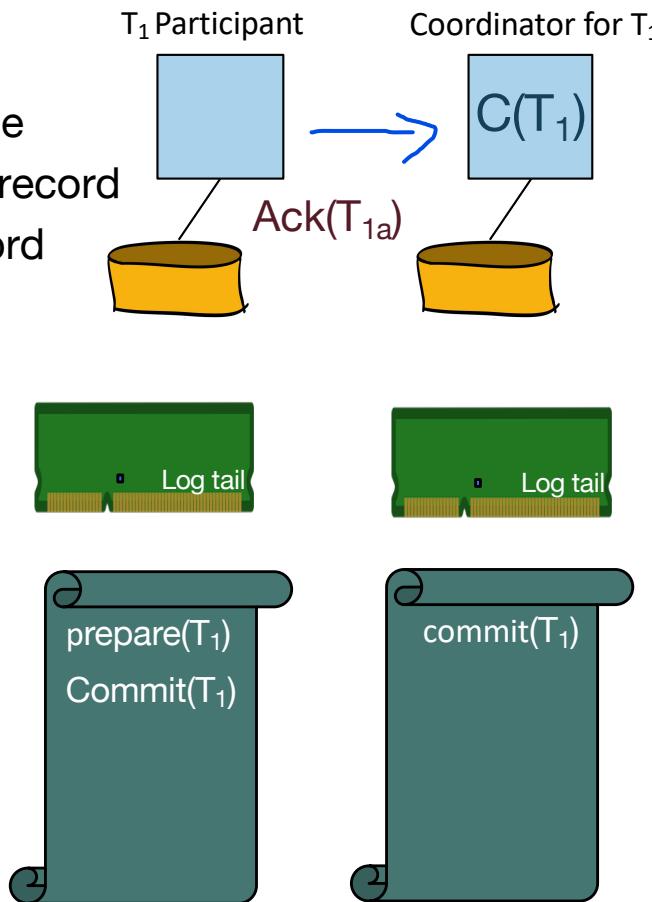
- **Phase 2:**
- Coordinator broadcasts result of vote
- Participants generate commit/abort record
- **Participants flush commit/abort record**
- Participants respond with Ack
- Coordinator generates end record
- Coordinator flushes end record



One More Time, With Logging, Part 14



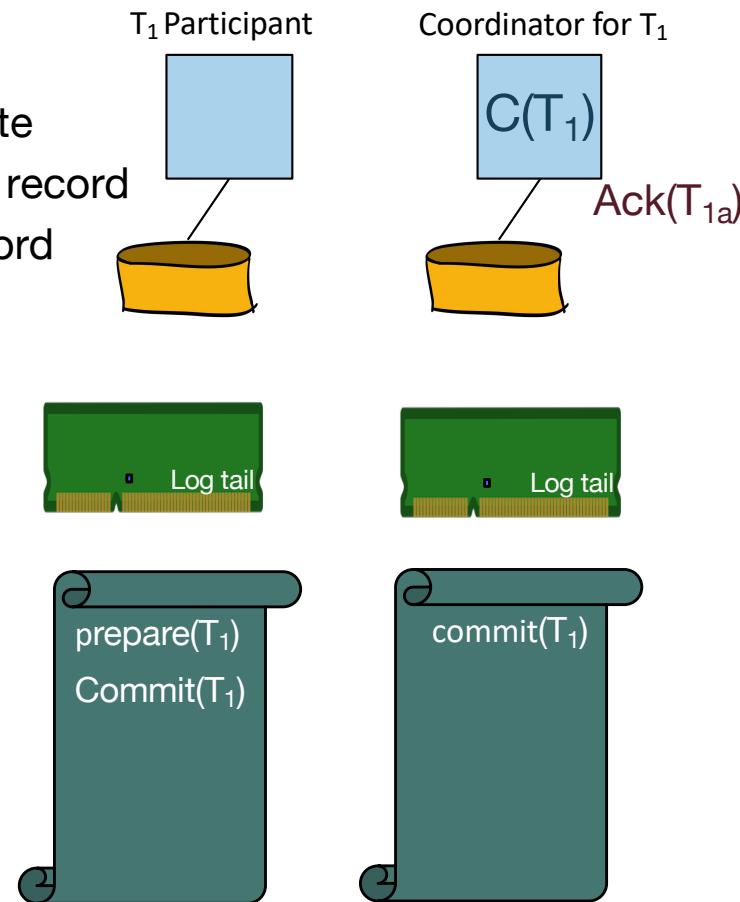
- **Phase 2:**
- Coordinator broadcasts result of vote
- Participants generate commit/abort record
- Participants flush commit/abort record
- **Participants respond with Ack**
- Coordinator generates end record
- Coordinator flushes end record



One More Time, With Logging, Part 15



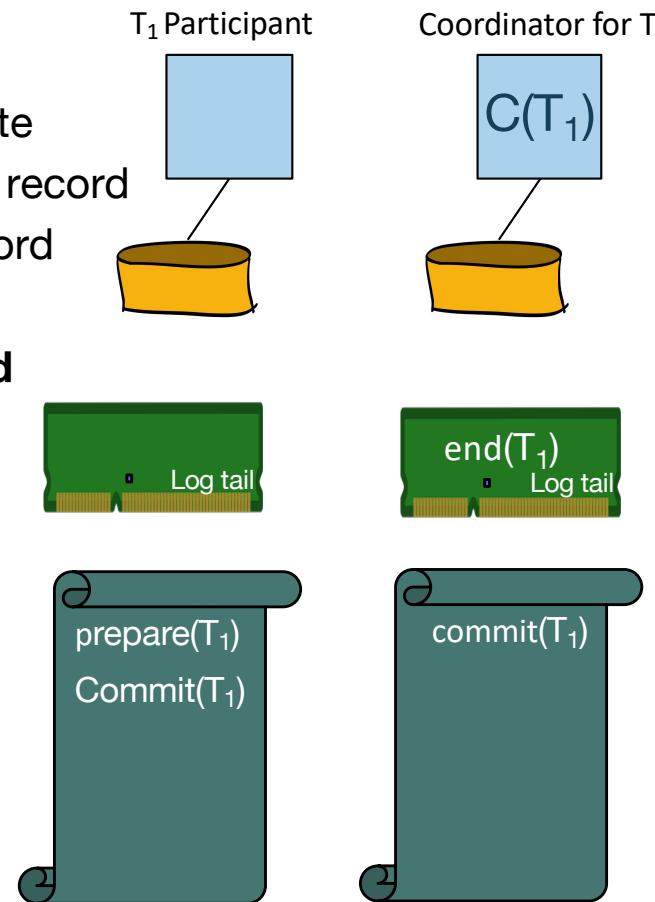
- **Phase 2:**
- Coordinator broadcasts result of vote
- Participants generate commit/abort record
- Participants flush commit/abort record
- **Participants respond with Ack**
- Coordinator generates end record
- Coordinator flushes end record



One More Time, With Logging, Part 16



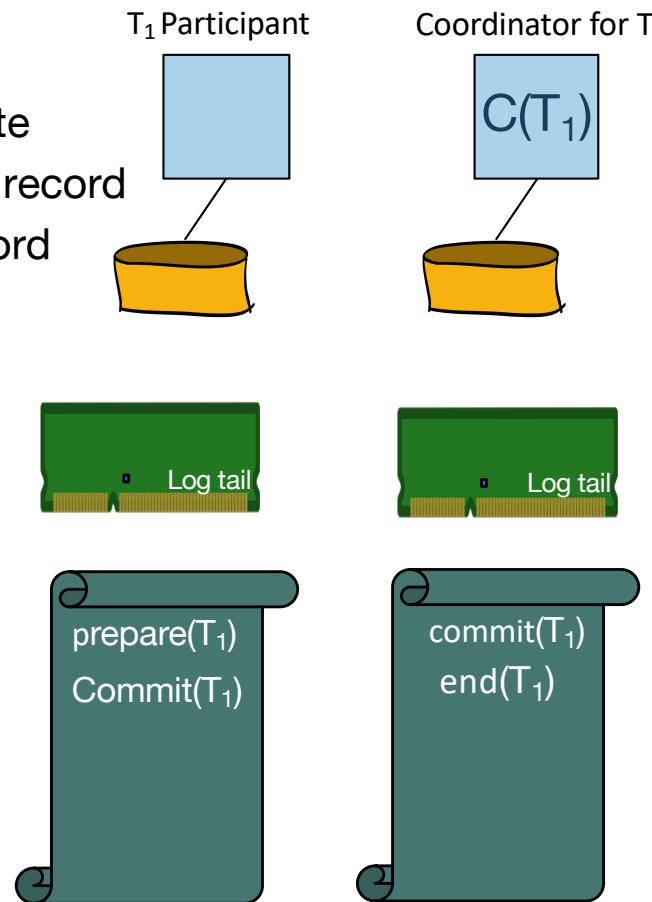
- **Phase 2:**
- Coordinator broadcasts result of vote
- Participants generate commit/abort record
- Participants flush commit/abort record
- Participants respond with Ack
- **Coordinator generates end record**
- Coordinator flushes end record



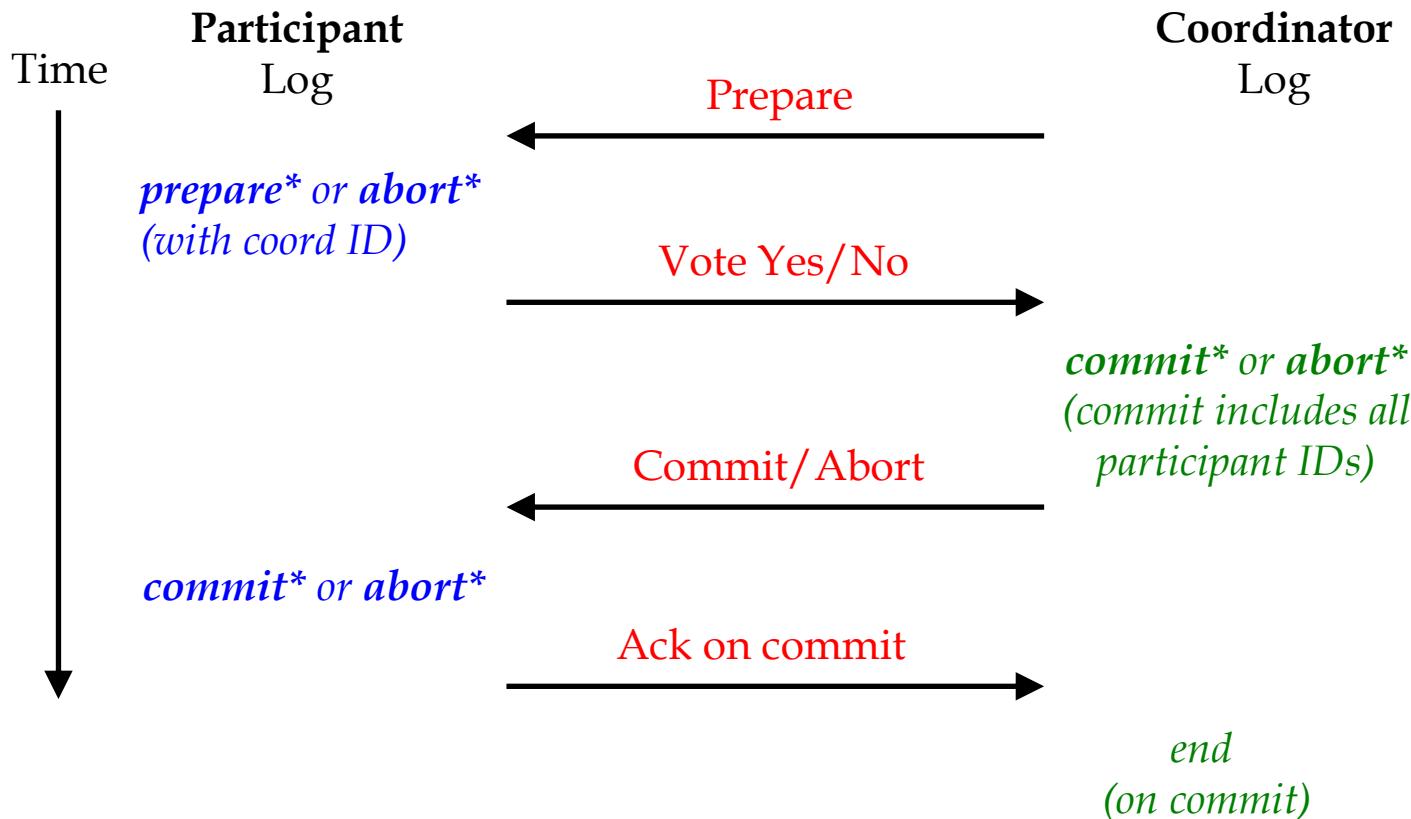
One More Time, With Logging, Part 17



- **Phase 2:**
- Coordinator broadcasts result of vote
- Participants generate commit/abort record
- Participants flush commit/abort record
- Participants respond with Ack
- Coordinator generates end record
- **Coordinator flushes end record**



2PC In a Nutshell



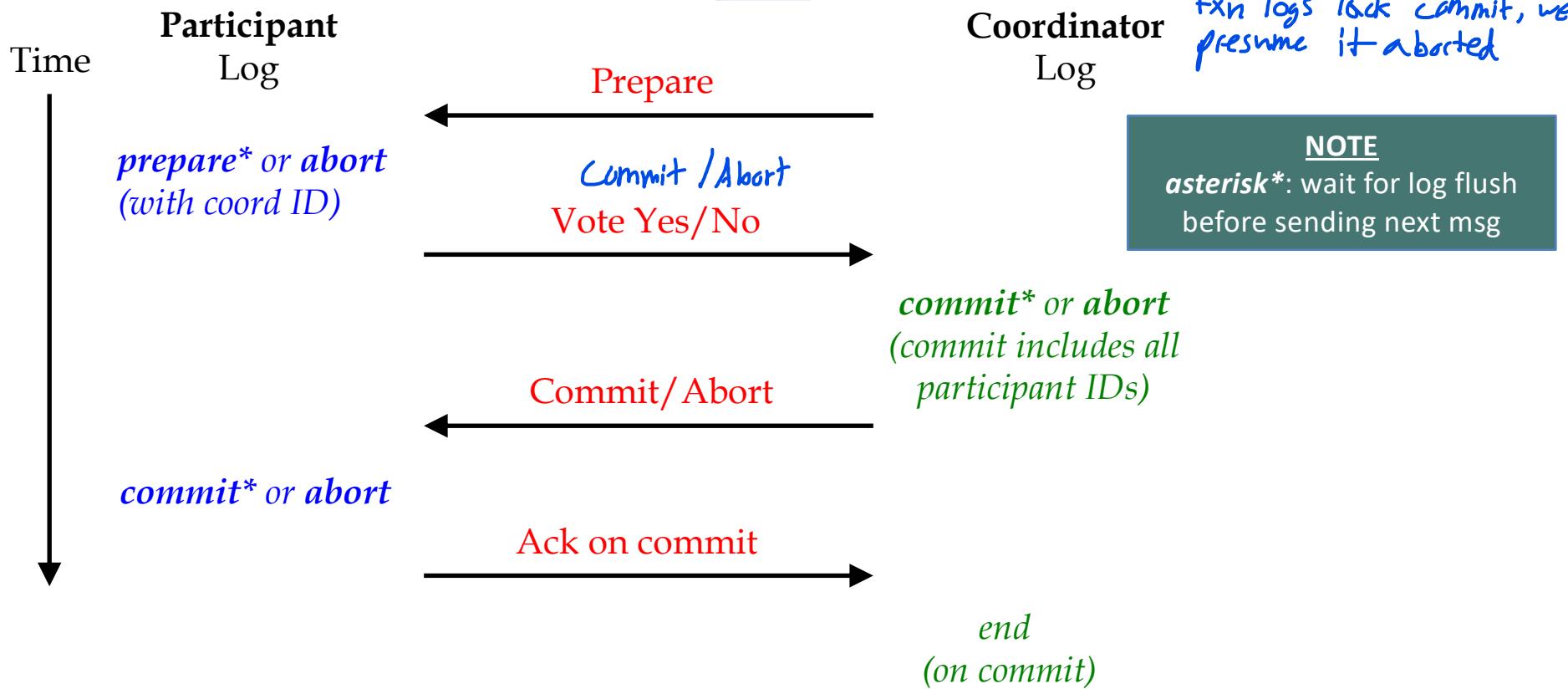
NOTE
*asterisk**: wait for log flush
before sending next msg

Presumed Abort

- Previous scheme needs prepare / abort / commit logged before sending msgs
- Presumed abort: Don't flush abort records to log before sending next msg



Don't need b/c if
txn logs lock commit, we
presume it aborted



Recovery and 2PC

Failure Handling



- Assume everybody recovers eventually
 - Big assumption!
 - Assume short downtimes!
- Coordinator notices a Participant is down?
 - If participant hasn't voted yet, coordinator aborts transaction
 - If waiting for a commit Ack, hand to “recovery process”
- Participant notices Coordinator is down?
 - If it hasn't yet logged prepare, then abort unilaterally
 - If it has logged prepare, hand to “recovery process”
- Note
 - Thinking a node is “down” may be incorrect!

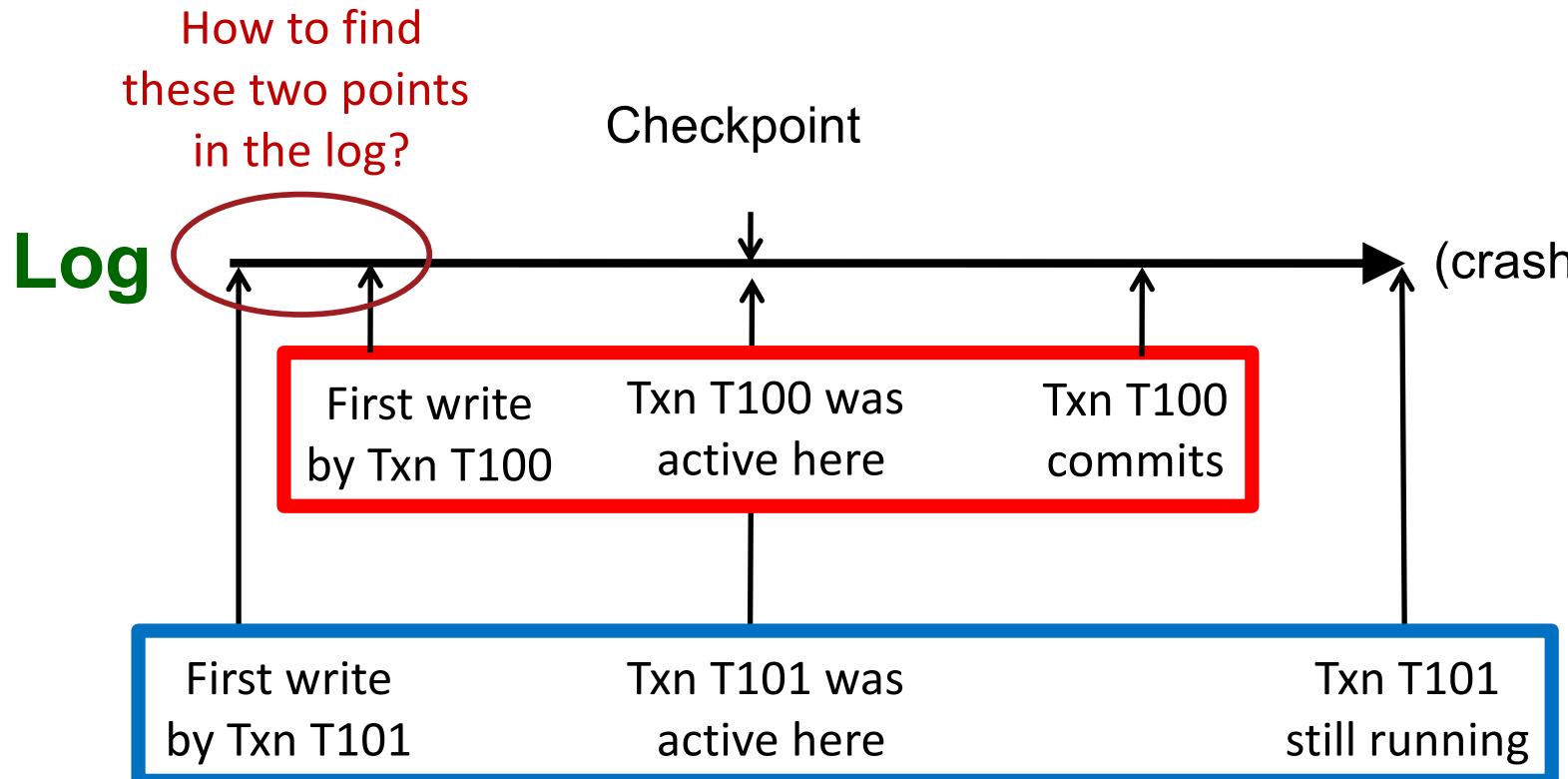
May just be lagging behind

Integration with ARIES Recovery



- On recovery
 - Assume there's a “Recovery Process” at each node
 - It will be given tasks to do by the Analysis phase of ARIES
 - These tasks can run in the background (asynchronously)
- Note: multiple roles on a single node
 - Coordinator for some xacts, Participant for others

Recall: ARIES Recovery



Recall: ARIES Recovery



Recovery from a system crash is done in 3 passes:

1. Analysis pass

- Recreate list of dirty pages and active transactions

2. Redo pass

- Redo all operations, even for those that were incomplete before crash
- Goal is to replay DB to the state at the moment of the crash

3. Undo pass

- Unroll effects of all incomplete transactions at time of crash
- Log changes during undo in case of another crash during undo

Integration with ARIES Recovery



- Recall transaction table states
 - Running, Committing, Aborting
- Same goal as before: recover the state of txns at crash
 - Caveat: 2PC introduces a new phase

prepare phases before commit

Integration with ARIES: Participant Analysis



- On seeing Prepare log record
 - Change state to committing
 - Tell recovery process to ask coordinator recovery process for status
 - When coordinator responds, recovery process handles commit/abort as usual
- On seeing Commit or Abort log record
 - Process the same way as single node ARIES
 - Send Ack back to coordinator when done as usual

Integration with ARIES: Coordinator Analysis



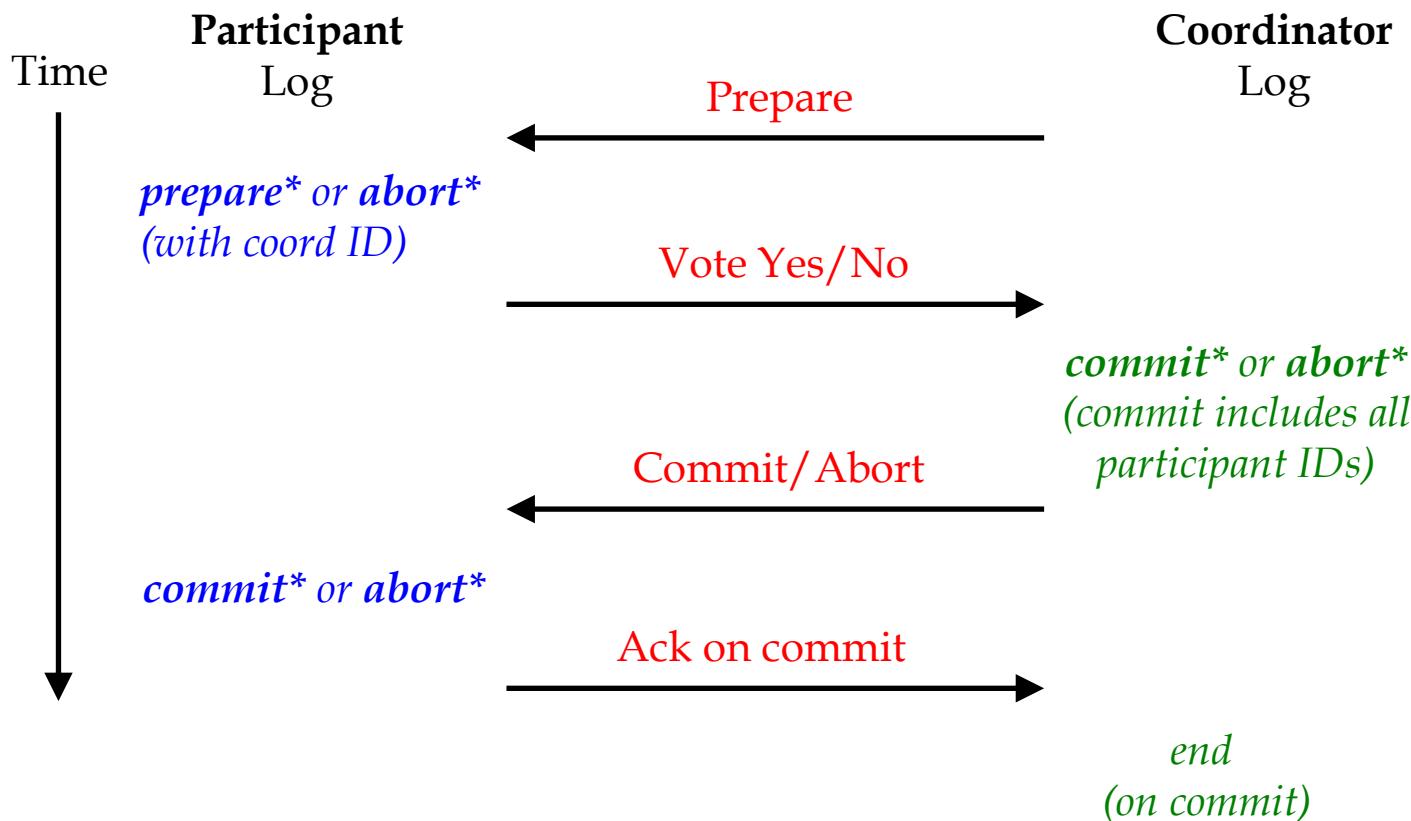
- On seeing Commit/Abort log record
 - Change state to committing/aborting respectively
 - Tell recovery process to send commit/abort msgs to participants
 - Once all participants ack commit, recovery process writes End record
- If at end of analysis there's no 2PC log records for xact X
 - This means we have not enter phase 2 or txn has aborted (if presumed abort)
 - Simply set to Aborting locally, and let ToUndo handle it.
 - Same for participant and coordinator

How Does Recovery Process Work?



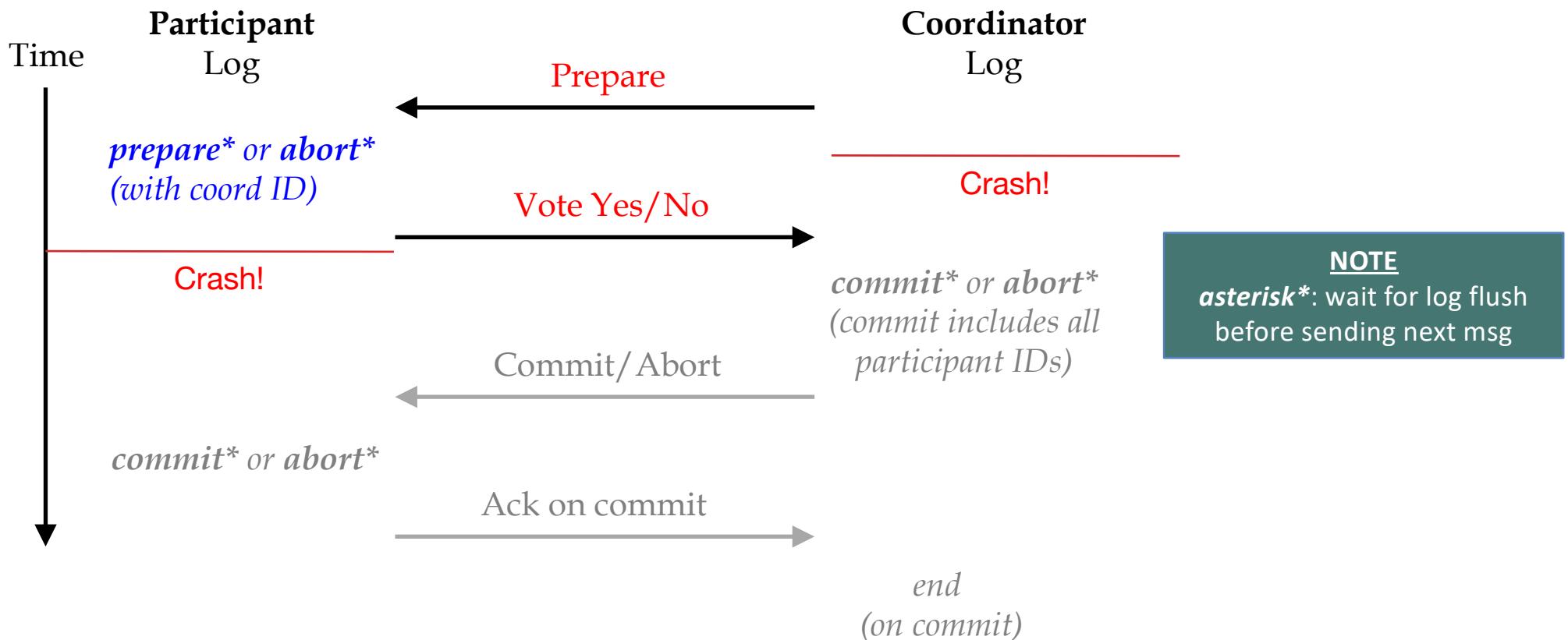
- Coordinator recovery process gets inquiry from a “prepared” participant
 - If transaction table at coordinator says aborting/committing
 - send appropriate response and continue protocol on both sides
 - If transaction table at coordinator says nothing: send ABORT
 - Only happens if coordinator had also crashed before writing commit/abort
 - Inquirer does the abort on its end

2PC In a Nutshell



NOTE
asterisk*: wait for log flush
before sending next msg

2PC In a Nutshell





Recovery: Think it through

- What happens when coordinator recovers?
 - With “commit” and “end”? **Nothing**
 - With just “commit”? **Rerun Phase 2!**
 - With “abort”? **Nothing (Presumed Abort)**
- What happens when participant recovers:
 - With no prepare/commit/abort? **Nothing (Presumed Abort)**
 - With “prepare” & “commit”? **Send Ack to coordinator.**
 - With just “prepare”? **Send inquiry to Coordinator**
 - With “abort”? **Nothing (Presumed Abort)**

2PC + Strict 2PL



- Ensure point-to-point messages are densely ordered
 - 1,2,3,4,5...
 - Dense per (sender/receiver/XID)
 - Receiver can detect anything missing or out-of-order
 - Receiver buffers message $k+1$ until $[1..k]$ received
 - Effect: receiver handles messages in order
- Commit:
 - When a participant processes Commit request, it has all the locks it needs
 - Flush log records and drop locks atomically
- Abort:
 - It's safe to abort autonomously, locally: no cascade.
 - Write abort to log, flush if not presumed abort (otherwise don't flush)
 - Perform local Undo, drop locks atomically

Availability Concerns



- What happens while a node is down?
 - Other nodes may be in limbo, holding locks
 - So certain data is unavailable
 - This may be bad...
- Dead Participants? Respawned by coordinator
 - Recover from log
 - And if the old participant comes back from the dead, just ignore it and tell it to recycle itself
- Dead Coordinator?
 - This is a problem!
 - 3-Phase Commit was an early attempt to solve it
 - Paxos Commit provides a more comprehensive solution
 - Gray+Lamport paper! Out of scope for this class.

Summing Up



- Distributed Databases
 - A central aspect of Distributed Systems
- Partitioning provides Scale-Up
- Can also partition lock tables and logs
- But need to do some global coordination:
 - Deadlock detection: easy
 - Commit: trickier
- Two-phase commit is a classic distributed consensus protocol
 - Logging/recovery aspects unique:
 - many distributed protocols gloss over
 - But 2PC is unavailable on any single failure
 - This is bad news for scale-up,
 - because odds of failure go up with #machines
 - Paxos Commit (Gray+Lamport) addresses that problem