# Distributed Transactions
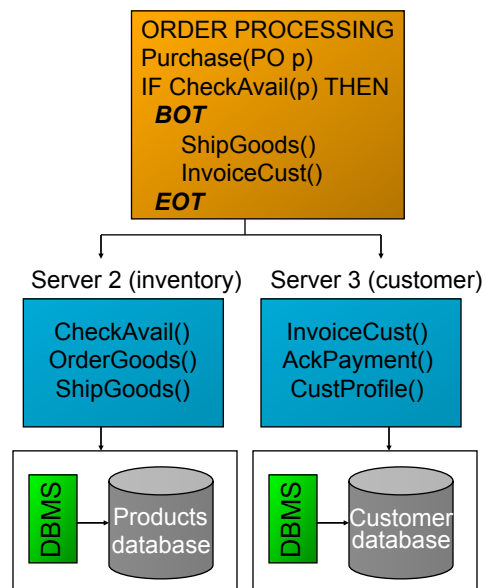
Dominic Duggan

Stevens Institute of Technology

Based in part on materials by K. Birman
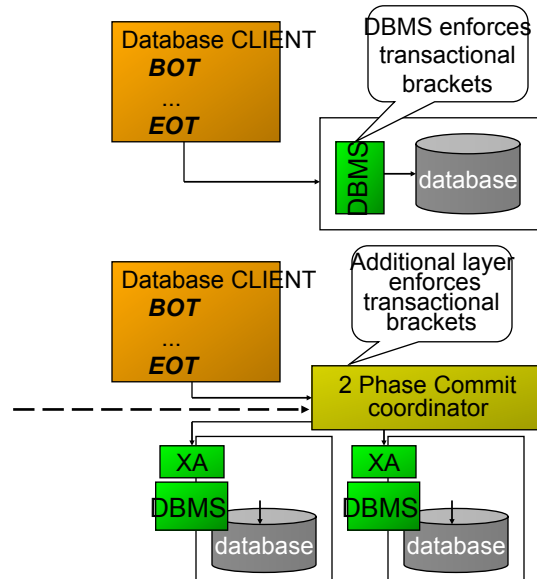
1

# Client, server, and databases
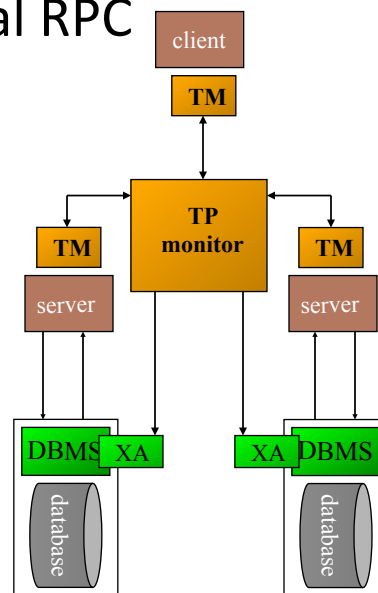
ORDER PROCESSING
Purchase(PO p)
IF CheckAvail(p) THEN
  **BOT**
    ShipGoods()
    InvoiceCust()
  **EOT**

Server 2 (inventory)

CheckAvail()
OrderGoods()
ShipGoods()

Server 3 (customer)

InvoiceCust()
AckPayment()
CustProfile()

DBMS → Products database

DBMS → Customer database

2

# Multiple Databases

Database CLIENT
**BOT**
...
**EOT**

DBMS enforces transactional brackets

DBMS

database

Database CLIENT
**BOT**
...
**EOT**

Additional layer enforces transactional brackets

2 Phase Commit coordinator

XA

DBMS

database

XA

DBMS

database

3

---

# Transactional RPC

- Transaction Managers
- Transaction Processing Monitor (TPM)
- XA API

client

**TM**

**TM**

**TP monitor**

**TM**

server

server

DBMS XA

XA DBMS

database

database

4

## Basic TRPC (making calls)

**Client**
**BOT**
**...**

**1**

**Client stub**

5

## Basic TRPC (making calls)

**Client**
**BOT**
**...**

**1**

**Client stub**
**Get tid**
**from TM**

**2**

**Transaction Manager (TM)**
**Generate tid**
**store context for tid**

6

3

# Basic TRPC (making calls)

**Client**
  BOT
  …

  Service_call
  …

**1**

**Client stub**
  Get tid
  from TM

  Add tid to
  call

**2**

**Transaction Manager (TM)**
  Generate tid
  store context for tid

**3**

**Server stub**
  Get tid

**Server**

7

---

# Basic TRPC (making calls)

**Client**
  BOT
  …

  Service_call
  …

**1**

**Client stub**
  Get tid
  from TM

  Add tid to
  call

**2**

**Transaction Manager (TM)**
  Generate tid
  store context for tid

  Associate server to tid

**4**

**3**

**Server stub**
  Get tid
  register with
    the TM
  Invoke service

**Server**

  Service
  procedure

8

4

# Basic TRPC (making calls)

| Client | Client stub | | Transaction Manager (TM) |
|---|---|---|---|
| BOT | Get tid | 2 | Generate tid |
| ... | from TM | | store context for tid |
| Service_call | Add tid to | | Associate server to tid |
| ... | call | | |

1

4

3

| Server stub | Server |
|---|---|
| Get tid | |
| register with | |
| the TM | |
| Invoke service | Service |
| return | procedure |

5

9

# Basic TRPC (committing calls)

| Client | Client stub | | TP Monitor (TPM)) |
|---|---|---|---|
| ... | | 1 | Look up tid |
| Service_call | | | |
| ... | | | |
| EOT | Send to TM | | |

10

5

# Basic TRPC (committing calls)

**Client**
 ...
 Service_call
 ...
 EOT

**Client stub**

Send to TM

**1**

**TP Monitor (TPM)**
Look up tid

**2** Run 2PC with all servers
 associated with tid

**Server stub** **Server**
Participant
 in 2PC

11

---

# Basic TRPC (committing calls)

**Client**
 ...
 Service_call
 ...
 EOT

**Client stub**

Send to TM
commit(tid)

**1**

**TP Monitor (TPM)**
Look up tid

**2** Run 2PC with all servers
 associated with tid

**3** Confirm commit

**Server stub** **Server**
Participant
 in 2PC

12

# TWO PHASE COMMIT (2PC)

# Atomic Commitment

- Given a set of processes
- **Coordinator** (i.e. TPM) wants to initiate an action (commit)
- **Participants** may vote for or against the action
- Perform the action only if all vote in favor
- Otherwise abort
- Goal is *all-or-nothing* outcome

# Non-triviality

- Avoid solutions that do nothing
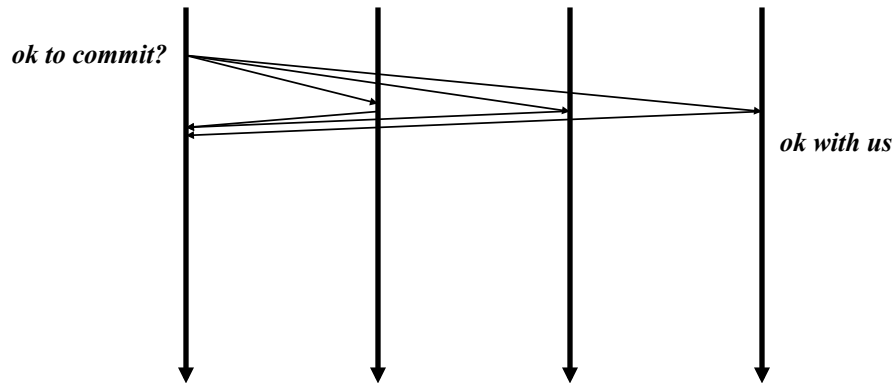- What is a trivial solution?
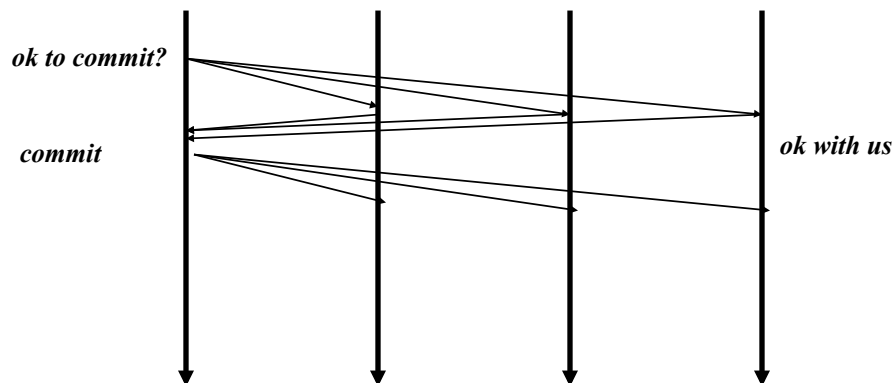- What is a validity condition?

15

# Commit protocol illustrated

*ok to commit?*

16

# Commit protocol illustrated

*ok to commit?*

*ok with us*

17

# Commit protocol illustrated

*ok to commit?*

*commit*

*ok with us*

*Note: garbage collection protocol not shown here*

18

# Two-phase Commit

- Phase 0: Flush caches on Web, app server
- Phase 1: Coordinator asks participants for vote
  - Data managers force updates to the log
  - Then say "ok to commit"
- Phase 2: If all are ok to commit, then coordinator tells participants to commit.  Otherwise, abort.
- Data managers then make updates permanent or rollback to old values, and release locks

19

---

# Commit protocol illustrated



*ok to commit?*

*… times out abort!*

*crashed!*

*ok with us*

*Note: garbage collection protocol not shown here*

20

# Non-triviality

- Avoid solutions that do nothing
- Commit validity: if all vote for commit, protocol must commit
- ...but what if participant vote is lost?
- "Non-triviality" condition hard to capture

# Unilateral abort

- Any data manager can unilaterally abort a transaction until it has said "prepared"
- Implication: even a data manager where only reads were done must participate in 2PC protocol!

**PROBLEMS WITH 2PC**

23

# Non-blocking Commit

- Goal: a protocol that allows all operational processes to terminate the protocol even if some subset crash

24

# Commit with
# unreliable failure detectors

- Assume processes fail by crashing
  - No Byzantine failures
- Coordinator detects failures (unreliably) using timouts
- Challenge: terminate the protocol if the coordinator fails

25

# 2 phase commit

**THREE PHASE COMMIT (3PC)**

# Three-phase commit

- Seeks to increase availability
- Makes an unrealistic assumption that failures are accurately detectable
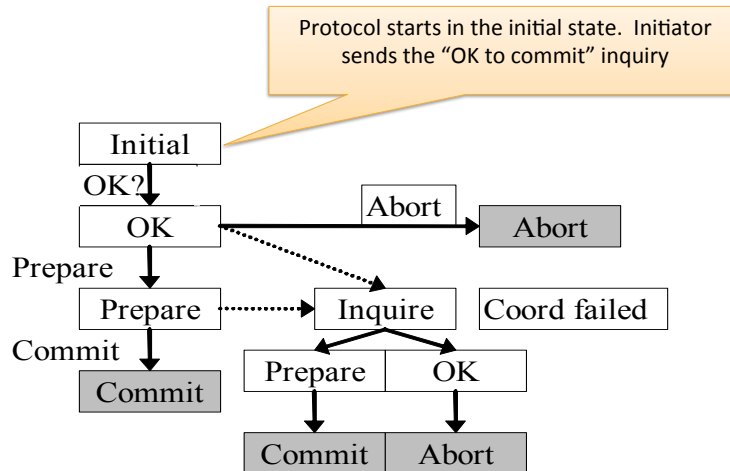- With this, can terminate the protocol even if a failure does occur

# Three-phase commit

- Coordinator starts protocol by sending request
- Participants vote to commit or to abort
- Coordinator collects votes, decides on outcome
- Coordinator can abort immediately
- To commit, coordinator first sends a "prepare to commit" message
- Participants acknowledge, commit occurs during a final round of "commit" messages
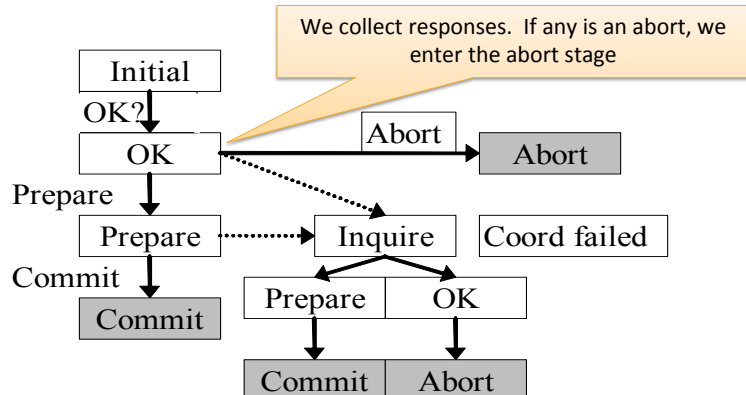
29

# Three-phase commit

| Phase 1 | Phase 2 | Phase 3 |
|---------|---------|---------|
| Vote? | Prepare to commit | Commit! |

*3PC initiator*

p
q
r
s
t

| All vote "commit" | All say "ok" | They commit |

# State diagram for non-faulty member

Protocol starts in the initial state. Initiator sends the "OK to commit" inquiry

Initial → OK? → OK → Abort → Abort

OK → (dotted) → Inquire    Coord failed

Prepare → Prepare → (dotted) → Inquire

Inquire → Prepare / OK

Prepare → Commit    OK → Abort

Commit → Commit

# State diagram for non-faulty member

We collect responses. If any is an abort, we enter the abort stage

Initial → OK? → OK → Abort → Abort

Prepare → Prepare → (dotted) → Inquire    Coord failed

Inquire → Prepare / OK

Commit → Commit

Prepare → Commit    OK → Abort

## State diagram for non-faulty member

Initial

OK?

OK

Prepare

Prepare

Commit

Commit

Inquire

Coord failed

Prepare

OK

Commit

Abort

Abort

Otherwise send prepare-to-commit messages out

33

## State diagram for non-faulty member

Initial

OK?

OK

Prepare

Prepare

Commit

Commit

Inquire

Coord failed

Prepare

OK

Commit

Abort

This state corresponds to the coordinator sending out the commit messages. We enter the state when *all* members receive them

34

## State diagram for non-faulty member

Coordinator failure sends us into an inquiry mode in which someone (anyone) tries to figure out the situation

Initial

OK?

OK

Abort → Abort

Prepare

Prepare → Inquire → Coord failed

Commit

Commit

Prepare | OK

Commit | Abort

35

## State diagram for non-faulty member

Here, we "finish off" the prepare state if a crash interrupted it, by resending the prepare message (needed in case only some processes saw the coordinator's message before it crashed)

Initial

OK?

OK

Prepare

Prepare → Inquire → Coord failed

Commit

Commit

Prepare | OK

Commit | Abort

36

18

State diagram for non-faulty member

Initial
OK?
OK
Prepare
Prepare
Commit
Commit
Inquire
Coord failed
Prepare
OK
Commit
Abort

We get here if all surviving processes are still in the initial "OK to commit?" stage

37



State diagram for non-faulty member

Initial
OK?
OK
Prepare
Prepare
Commit
Commit
Abort
Inquire
Coord failed
Prepare
OK
Commit
Abort

In this case it is safe to abort, and we do so

38

## Observations about 3PC

- Key point: Extra buffer state

- What if none of surviving participants have heard from coordinator?
  - After voting phase
  - 2PC: Some crashed processes may have committed
  - 3PC: No crashed process has committed yet (Why?)

## Observations about 3PC

- If any process is in "prepare to commit" all voted for commit
- Protocol commits only when all surviving processes have acknowledged prepare to commit
- After coordinator fails, it is easy to run the protocol forward to commit state (or back to abort state)

# Problems with 3PC

- Assumes reliable failure detectors
- But even with realistic failure detectors (that can make mistakes), protocol still blocks!
  - "Network partitioning"
- Can prove that this problem is not avoidable

41

# Situation in practical systems?

- Most use protocols based on 2PC
- Need to extend garbage collection
  - protocol state information
- Some systems accept the risk of blocking
- Others reduce the consistency property to make progress

42

## Example: Primary with Backup

Primary          Backup

$T_1$

$T_1$

Logged updates

$T_k$

$T_{k+1}$

$T_k$

Lost updates

$T_1'$

$T_{k+m}$

To be reconciled

$T_n'$

43

# PAXOS AND ATOMIC COMMIT

44

# Database Commit

client      TM      RM

RM

RM

Commit      Ready?

Ready

Commit      Commit

TM: Transaction Manager
RM: Resource Manager

45

---

# Two Phase Commit

- *N* Resource Managers (RMs)
- Want all RMs to commit or all abort.
- Coordinated by Transaction Manager (TM)
  TM sends Prepare, Commit-Abort
- RM responds Prepared, Aborted
- *3N+1* messages
- *N+1* stable writes
- Delay
  - 4 message
  - 2 stable write
- Blocking:
  if TM fails,
  Commit-Abort stalls

RequestCommit

Prepare

Prepared

Commit

**Resource Manager**

working

prepared

committed    aborted

**Transaction Manager**

working

committed    aborted

46

23

# The Problem With 2PC

- Atomicity – all or nothing
- Consistency – does right thing
- Isolation – no concurrency anomalies
- Durability / Reliability – state survives failures
- Availability: always up

**Blocks if TM fails**

# Problem Statement

- ACID Transactions make error handling easy.

- One fault can make 2-Phase Commit block.

- Goal: ACID and Available.
  Non-blocking despite *F* faults.

# Fault-Tolerant Two Phase Commit

client → RequestCommit → TM → Prepare → RM

Prepare → RM

RequestCommit    Prepare

49



# Fault-Tolerant Two Phase Commit

client → RequestCommit → TM    ← Prepared    RM

← Prepared

RM

RequestCommit    Prepare →    ← Prepared

If the 2PC Transaction Manager (TM) Fails, transaction blocks.

50

# Fault-Tolerant Two Phase Commit

client → RequestCommit → TM ⊘

RM

TM

RM

RequestCommit →

Prepare →

← Prepared

If the 2PC Transaction Manager (TM) Fails, transaction blocks.

Solution: Add a "spare" transaction manager
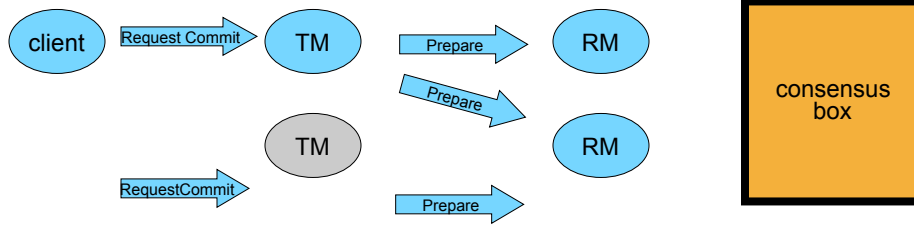*(non blocking commit, 3 phase commit)*

51

# Fault-Tolerant Two Phase Commit

client → RequestCommit → TM ⊘

RM

TM

RM

RequestCommit →

Prepare →

Prepare →

Prepare →

← Prepared

Prepare →

If the 2PC Transaction Manager (TM) Fails, transaction blocks.

Solution: Add a "spare" transaction manager
*(non blocking commit, 3 phase commit)*

52

# Fault-Tolerant Two Phase Commit

client → RequestCommit → TM (failed)

RM

Prepared ←

TM          Prepared ←          RM

RequestCommit →

Prepare →
Prepared ←
Prepare →
Prepared ←

If the 2PC Transaction Manager (TM) Fails,  transaction blocks.

Solution: Add a "spare" transaction manager
*(non blocking commit, 3 phase commit)*

53



# Fault-Tolerant Two Phase Commit

client → RequestCommit → TM (failed)

RM

commit →

TM          commit →          RM

RequestCommit →

Prepare →
Prepared ←
Prepare →
Prepared ←
commit →

If the 2PC Transaction Manager (TM) Fails,  transaction blocks.

Solution: Add a "spare" transaction manager
*(non blocking commit, 3 phase commit)*

54

## Fault-Tolerant Two Phase Commit



If the 2PC Transaction Manager (TM) Fails,  transaction blocks.

Solution: Add a "spare" transaction manager
*(non blocking commit, 3 phase commit)*

55

## Paxos Consensus Box



- Collects proposed values
- Picks one proposed value
- Remembers it forever

56

28

## Consensus for Commit
## The Obvious Approach

client —Request Commit→ TM —Prepare→ RM

TM

RequestCommit→ —Prepare→ RM

—Prepare→

consensus box

- Get consensus on TM's decision.
- TM just learns consensus value.
- TM is "stateless"
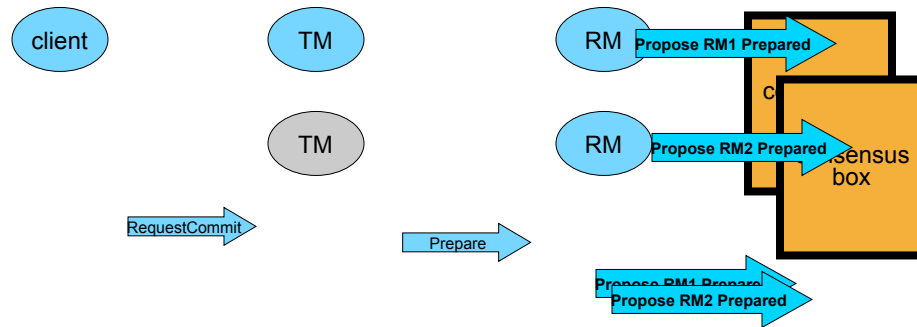
57

## Consensus for Commit
## The Obvious Approach

client    TM ←Prepared— RM

TM ←Prepared— RM

RequestCommit→ —Prepare→

←Prepared—

consensus box

- Get consensus on TM's decision.
- TM just learns consensus value.
- TM is "stateless"

58

Consensus for Commit
The Obvious Approach

- Get consensus on TM's decision.
- TM just learns consensus value.
- TM is "stateless"

59



Consensus for Commit
The Obvious Approach

- Get consensus on TM's decision.
- TM just learns consensus value.
- TM is "stateless"

60

## Consensus for Commit
## The Obvious Approach



- Get consensus on TM's decision.
- TM just learns consensus value.
- TM is "stateless"

61

## Consensus for Commit
## The Paxos Commit Approach



- Get consensus on each RM's choice.
- TM just combines consensus values.
- TM is "stateless"

62

Consensus for Commit
The Paxos Commit Approach

- Get consensus on each RM's choice.
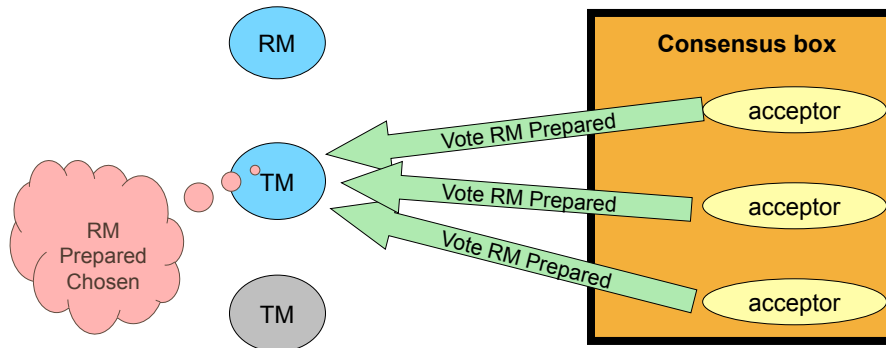- TM just combines consensus values.
- TM is "stateless"

63



Consensus for Commit
The Paxos Commit Approach

- Get consensus on each RM's choice.
- TM just combines consensus values.
- TM is "stateless"

64

# Consensus for Commit
## The Paxos Commit Approach

client · TM · RM

Commit ← (green arrow)

Commit → (green arrow)

TM (grey)

Commit → (green arrow)

RM

consensus box (orange)

RequestCommit → (blue arrow)

Prepare → (blue arrow)

Propose RM1 Prepare
Propose RM2 Prepared (blue arrow)

RM1 Prepared Chosen
RM2 Prepared Chosen (green arrow)

Commit ← (green arrow)

Commit → (green arrow)

- Get consensus on each RM's choice.
- TM just combines consensus values.
- TM is "stateless"

65

---

# Consensus in Action

RM

Propose RM Prepared
Propose RM Prepared
Propose RM Prepared

**Consensus box**

acceptor

acceptor

TM

acceptor

TM (grey)

- The normal (failure-free) case

66

# Consensus in Action



- The normal (failure-free) case
- Two message delays
- Can optimize

67

# Consensus in Action



TM can always learn what was chosen,
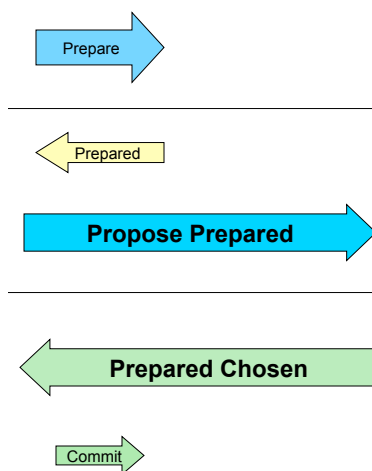or get *Aborted* chosen if nothing chosen yet;

.

68

34

# Consensus in Action

RM

TM

TM

**Consensus box**

acceptor

acceptor

acceptor

TM can always learn what was chosen,
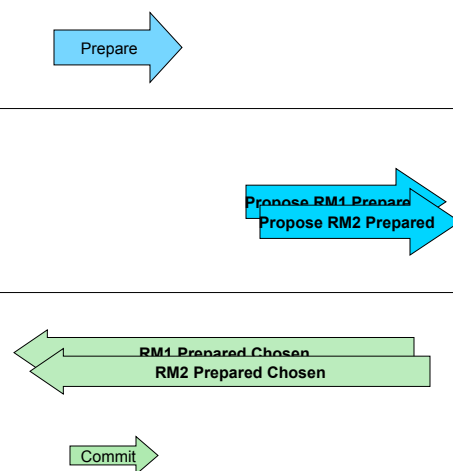or get *Aborted* chosen if nothing chosen yet;
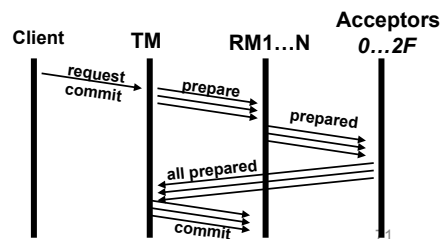if majority of acceptors working .

69

---

# The Obvious Approach

# Paxos Commit

**One fewer message delay**

Prepare

Prepare

Prepared

**Propose Prepared**

Propose RM1 Prepare
Propose RM2 Prepared

**Prepared Chosen**

RM1 Prepared Chosen
RM2 Prepared Chosen

Commit

Commit

70

# Paxos Commit

- *N* RMs

- *2F+1* acceptors (*~2F+1* TMs)

- If *F+1* acceptors see all RMs prepared, then transaction committed.

- *(N+1)(F+3) - 2* messages

  – 1 request to commit

  – N - 1 prepare

  – N (F + 1) prepared

  – F + 1 all prepared

  – N commit

- 5 message delays
  2 stable write delays.

Diagram labels: Client, TM, RM1...N, Acceptors 0...2F, request commit, prepare, prepared, all prepared, commit

---

| **Two-Phase Commit** | **Paxos Commit** |
|---|---|
| | tolerates *F* faults |
| • 3*N*+1 messages | • 3*N*+ 2*F(N+1)* +1 messages |
| • *N*+1 stable writes | • *N*+2*F*+1 stable writes |
| • 4 message delays | • 5 message delays |
| • 2 stable-write delays | • 2 stable-write delays |

## Same algorithm when *F*=0 and TM = Acceptor

72

# EXTENDED TRANSACTION MODELS

# Savepoints

**Transaction Execution**

- Idea: allow the state of a transaction to be rolled back to a certain point in execution
- Not necessary to roll back the entire transaction

Do some processing A

↓

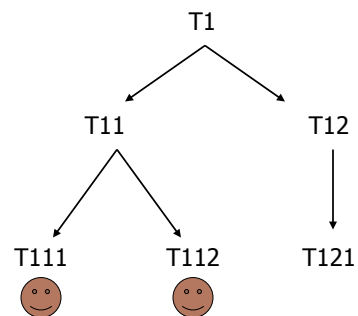Create a savepoint

↓

Do some more processing B

↓

Rollback to last savepoint (undoing any changes done by B)

## Nested Transactions (1)

- Idea: Allow transactions to be nested inside other transactions
- Child transaction's changes not visible to parent, siblings until commit
- Failure of a child transaction does not imply failure of the parent
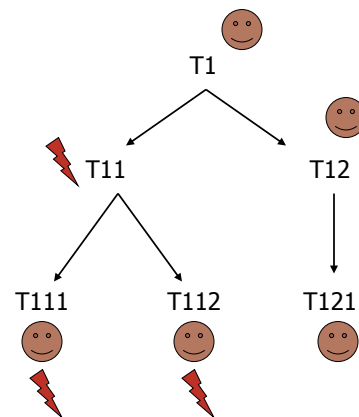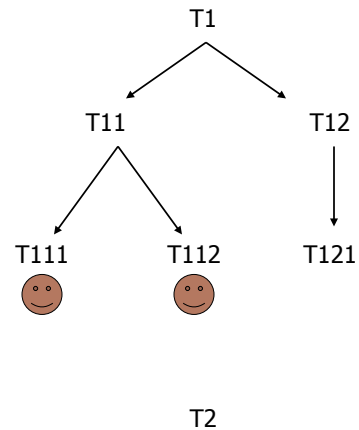
## Nested Transactions (2)

- Abort of a transaction forces all of its descendants to abort
- So commit of a non-root transaction is *tentative*
- Child can obtain locks from parent (*lock inheritance*)
- Child's locks released to parent (*lock anti-inheritance*)

# Nested Transactions (2)

- Abort of a transaction forces all of its descendants to abort
- So commit of a non-root transaction is *tentative*
- Child can obtain locks from parent (*lock inheritance*)
- Child's locks released to parent (*lock anti-inheritance*)



77

# Nested Transactions (2)

- Abort of a transaction forces all of its descendants to abort
- So commit of a non-root transaction is *tentative*
- Child can obtain locks from parent (*lock inheritance*)
- Child's locks released to parent (*lock anti-inheritance*)
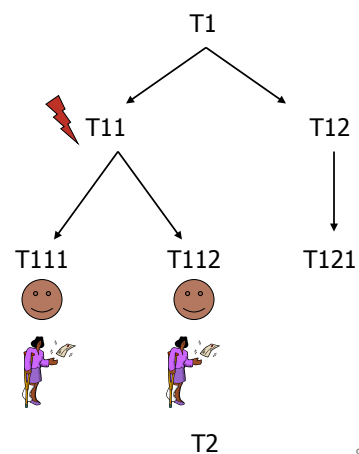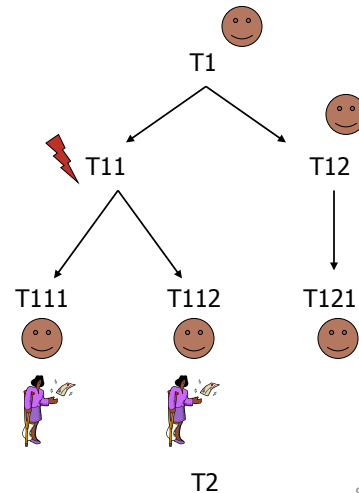


78

# Multi-Level Transactions (1)

- Unlike Nested, Multi-Level allows a transaction's changes to be made visible to *everyone* when it commits
  - Nested: visible to parent
- A compensating transaction is run if the parent transaction subsequently aborts
  - Each transaction has its own specific compensating transaction
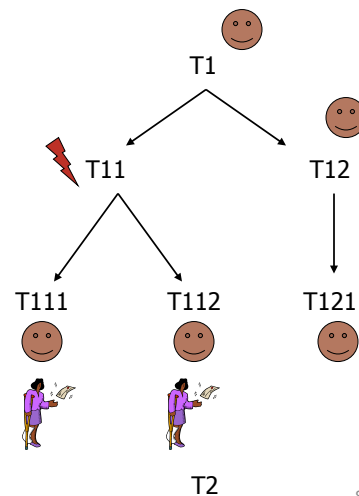


79

---

# Multi-Level Transactions (1)

- Unlike Nested, Multi-Level allows a transaction's changes to be made visible to *everyone* when it commits
  - Nested: visible to parent
- A compensating transaction is run if the parent transaction subsequently aborts
  - Each transaction has its own specific compensating transaction



80

40

# Multi-Level Transactions (1)

- Unlike Nested, Multi-Level allows a transaction's changes to be made visible to *everyone* when it commits
  - Nested: visible to parent
- A compensating transaction is run if the parent transaction subsequently aborts
  - Each transaction has its own specific compensating transaction

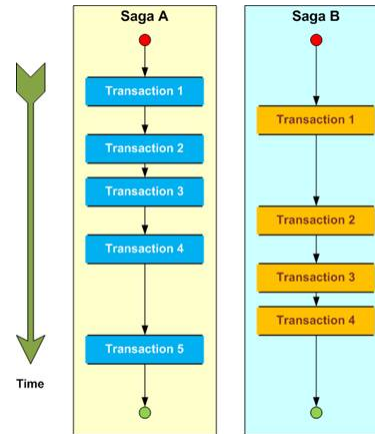

81

---

# Multi-Level Transactions (2)

- Motivation for multi-level: sharing of low-level resources in layered software systems
- Ex: insert records into database:
  - Find page on disk with free space
  - Insert records
  - Update indexes
  - Commit
- Nested would require page to be locked to further allocation until commit
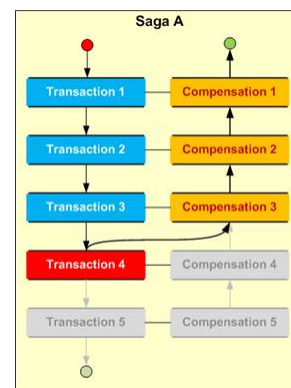


82

41

# Sagas (1)

- Transactions are incompatible with long-lived applications
  - Locks must be held to ensure ACID properties
- Sagas relax ACID properties, allow intermediate states to be visible to outside
- A saga is a sequence of transactions
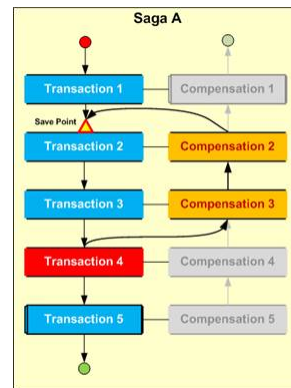  - "long-lived transaction"

# Sagas (2)

- Each transaction $T_i$ has a corresponding compensating transaction $CT_i$
- If $T_k$ aborts in a transaction, then the compensations

$$CT_{k-1},...,CT_2,CT_1$$

  are executed, in that order
  - Explicit rollback

# Sagas (3)

- Sagas also support forward recovery using savepoints
- If state is persisted between transactions, then compensating transactions roll back to that savepoint
- Saga continues from that point forward