# Distributed Transactions

Dominic Duggan
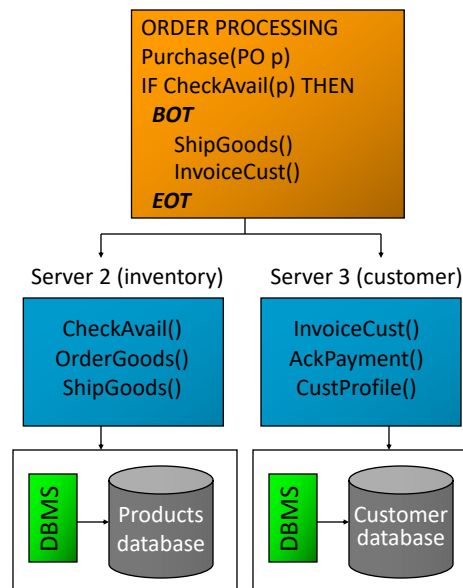
Stevens Institute of Technology

Based in part on materials by K. Birman

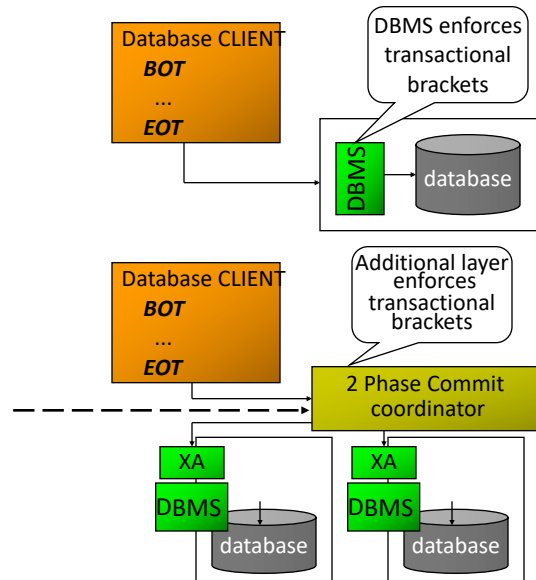1

1

---

# Client, server, and databases

ORDER PROCESSING
Purchase(PO p)
IF CheckAvail(p) THEN
  **BOT**
    ShipGoods()
    InvoiceCust()
  **EOT**

Server 2 (inventory)

CheckAvail()
OrderGoods()
ShipGoods()

DBMS → Products database

Server 3 (customer)

InvoiceCust()
AckPayment()
CustProfile()

DBMS → Customer database

2

2

1

# 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

3

# Transactional RPC

client

TM

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

TP monitor

TM
server

TM
server

DBMS XA
database

XA DBMS
database

4

4

2

# 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
...

**1**

**Client stub**
→ Get tid ←
from TM

**Service_call** → Add tid to
... call

**2**

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

**3**

**Server stub** | **Server**
→ Get tid

# Basic TRPC (making calls)

**Client**
BOT
...

**1**

**Client stub**
→ Get tid ←
from TM

**Service_call** → Add tid to
... call

**2**

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

Associate server to tid

**4**

**3**

**Server stub** | **Server**
→ Get tid
→ register with
the TM
Invoke service → Service
procedure

## Basic TRPC (making calls)

**Client**
BOT
…

Service_call
…

**Client stub**
Get tid
from TM

Add tid to
call

1

2

3

4

5

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

Associate server to tid

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

**Server**

Service
procedure

9

9

## Basic TRPC (committing calls)

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

**Client stub**

Send to TM

1

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

10

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
Participant
in 2PC

Server

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
Participant
in 2PC

Server

12

6

# 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
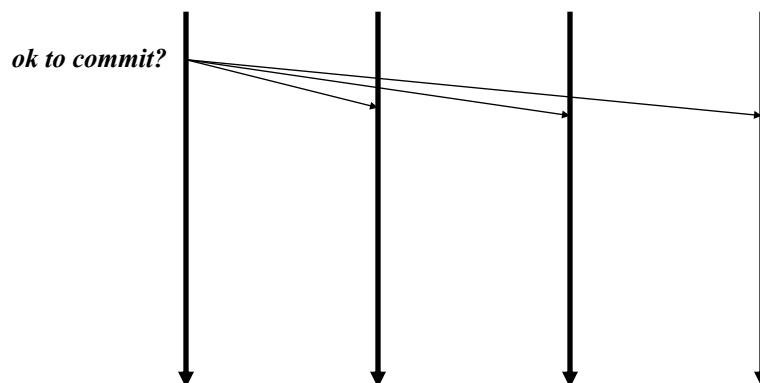- Otherwise abort
- Goal is *all-or-nothing* outcome

# Non-triviality

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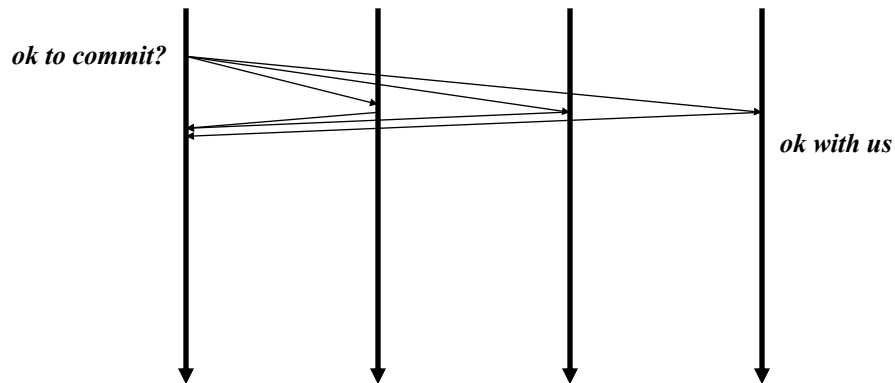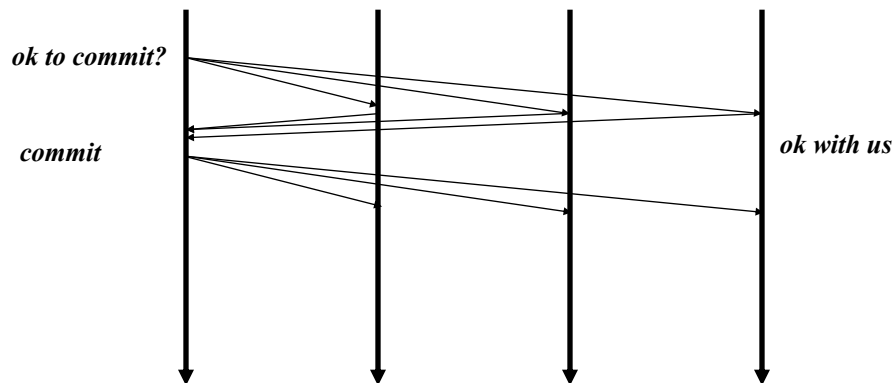- What is a validity condition?

# Commit protocol illustrated

*ok to commit?*

# Commit protocol illustrated

*ok to commit?*

*ok with us*

17

17

# Commit protocol illustrated

*ok to commit?*

*commit*

*ok with us*

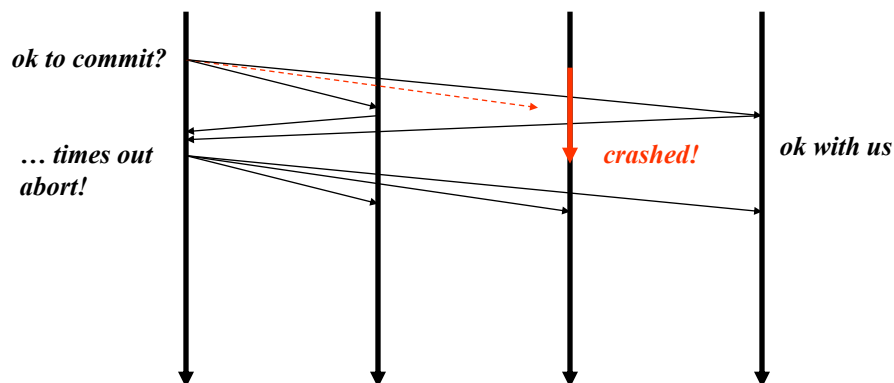*Note: garbage collection protocol not shown here*

18

18

9

# 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

19

# Commit protocol illustrated



*ok to commit?*

*… times out abort!*

*crashed!*

*ok with us*

*Note: garbage collection protocol not shown here*

20

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

21

21

# 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!

22

22

**PROBLEMS WITH 2PC**

# Non-blocking Commit

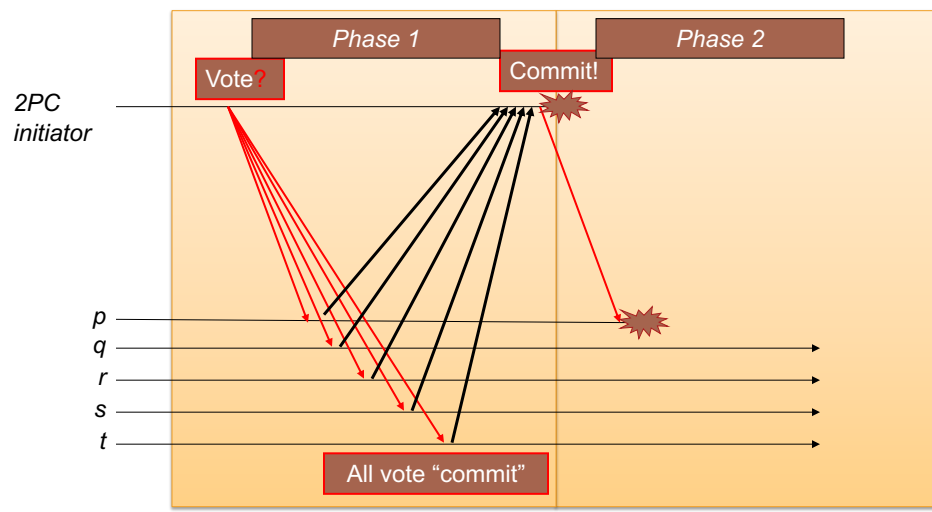- Goal: a protocol that allows all operational processes to terminate the protocol even if some subset crash

# 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
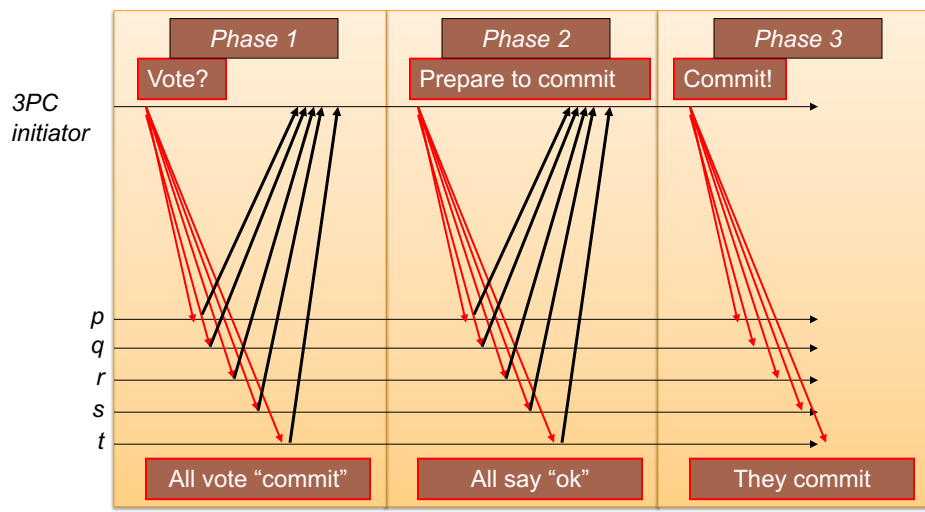
# Three-phase commit

- 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
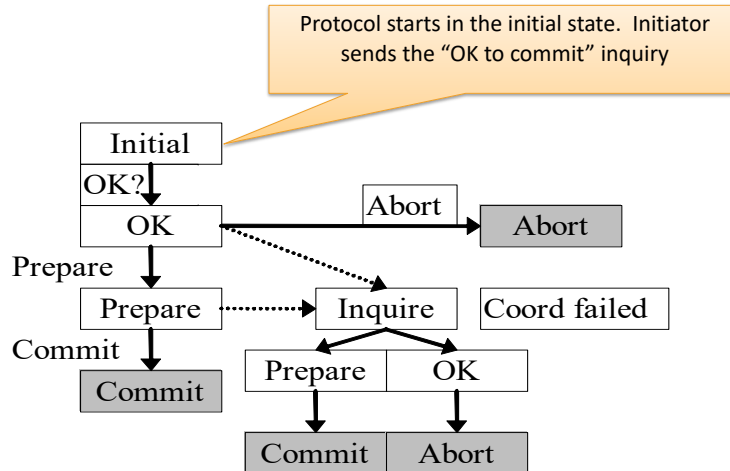- Participants acknowledge, commit occurs during a final round of "commit" messages
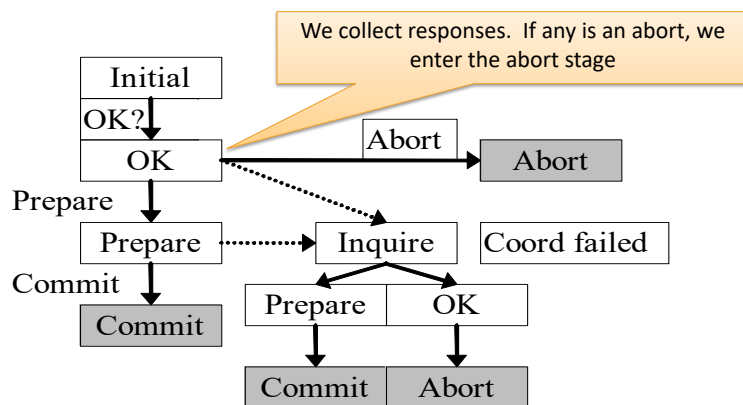
29

29

# Three-phase commit



*3PC initiator*

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

*p*
*q*
*r*
*s*
*t*

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

30

State diagram for non-faulty member

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

31



State diagram for non-faulty member

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

32

## State diagram for non-faulty member

Initial

OK?

OK

Prepare

Prepare

Commit

Commit

Abort

Inquire

Coord failed

Prepare

OK

Commit

Abort

Otherwise send prepare-to-commit messages out

33

33

## State diagram for non-faulty member

Initial

OK?

OK

Prepare

Prepare

Commit

Commit

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

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

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)

36
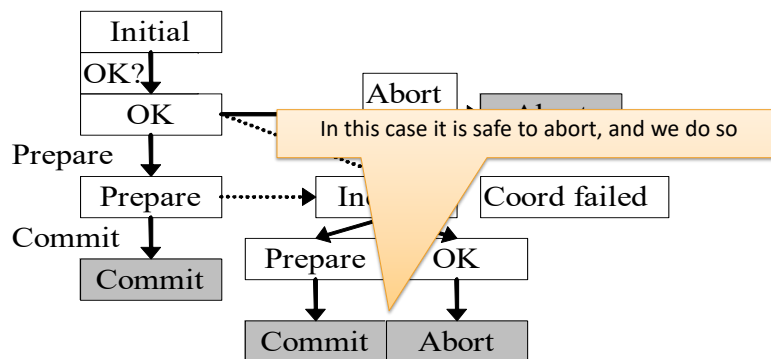
State diagram for non-faulty member



State diagram for non-faulty member

## 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"
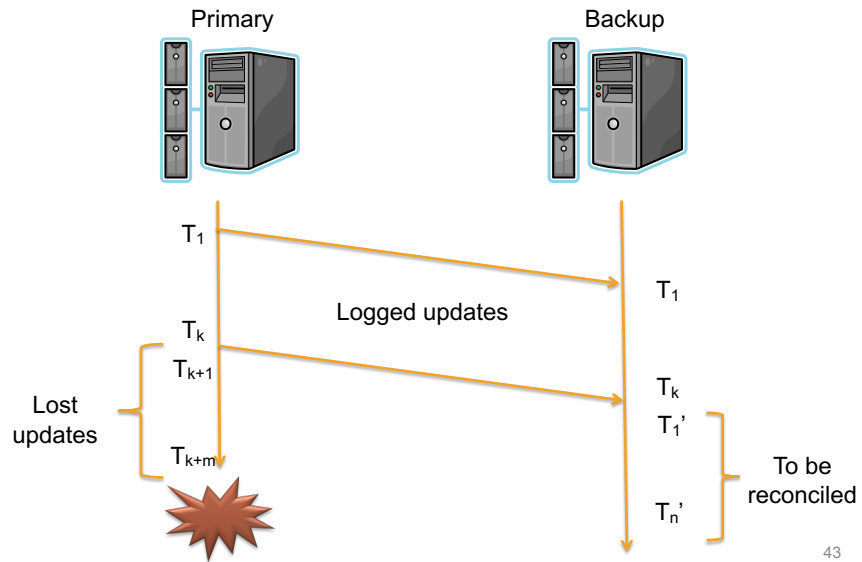- Can prove that this problem is not avoidable

41

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

42

# Example: Primary with Backup

Primary

Backup

$T_1$

$T_1$

Logged updates

$T_k$

$T_{k+1}$

$T_k$

$T_1$'

Lost
updates

To be
reconciled

$T_{k+m}$

$T_n$'

43

43

22