**Chapter 20**

## Distributed DBMSs - Advanced Concepts

## Transparencies

---

## Chapter - Objectives

- u **Distributed transaction management.**
- u **Distributed concurrency control.**
- u **Distributed deadlock detection.**

## Distributed Transaction Management

- u **Distributed transaction accesses data stored at more than one location.**
- u **Divided into a number of *sub-transactions*, one for each site that has to be accessed, represented by an *agent*.**
- u **Indivisibility of distributed transaction is still fundamental to transaction concept.**
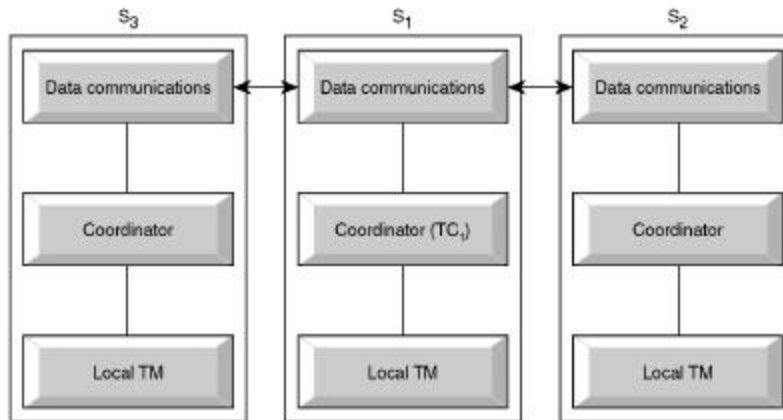- u **DDBMS must also ensure indivisibility of each sub-transaction.**

3

## Distributed Transaction Management

- u **Thus, DDBMS must ensure:**
  - – **synchronization of subtransactions with other local transactions executing concurrently at a site;**
  - – **synchronization of subtransactions with global transactions running simultaneously at same or different sites.**
- u **Global transaction manager (transaction coordinator) at each site, to coordinate global and local transactions initiated at that site.**

4

## Coordination of Distributed Transaction

## Distributed Locking

- u **Look at four schemes:**
  - – **Centralized locking**
  - – **Primary Copy 2PL**
  - – **Distributed 2PL**
  - – **Majority Locking**

## Centralized Locking

- u **Single site that maintains all locking information.**
- u **One lock manager for whole of DDBMS.**
- u **Local transaction managers involved in global transaction request and release locks from lock manager.**
- u **Or transaction coordinator can make all locking requests on behalf of local transaction managers.**
- u **Advantage - easy to implement.**
- u **Disadvantages - bottlenecks and lower reliability.**

7

## Primary Copy 2PL

- u **Lock managers distributed to a number of sites.**
- u **Each lock manager responsible for managing locks for set of data items.**
- u **For replicated data item, one copy is chosen as *primary copy*, others are *slave copies***
- u **Only need to write-lock primary copy of data item that is to be updated.**
- u **Once primary copy has been updated, change can be propagated to slaves.**

8

## Primary Copy 2PL

- u **Disadvantages - deadlock handling is more complex due; still a degree of centralization in system.**

- u **Advantages - lower communication costs and better performance than centralized 2PL.**

9

## Distributed 2PL

- u **Lock managers distributed to every site.**

- u **Each lock manager responsible for locks for data at that site.**

- u **If data not replicated, equivalent to primary copy 2PL.**

- u **Otherwise, implements a Read-One-Write-All (ROWA) replica control protocol.**

10

## Distributed 2PL

- u **Using ROWA protocol:**
  - – **Any copy of replicated item can be used for read.**
  - – **All copies must be write-locked before item can be updated.**
- u **Disadvantages - deadlock handling more complex; communication costs higher than primary copy 2PL.**

11

## Majority Locking

- u **Extension of distributed 2PL.**
- u **To read or write data item replicated at $n$ sites, sends a lock request to more than half the $n$ sites where item is stored.**
- u **Transaction cannot proceed until majority of locks obtained.**
- u **Overly strong in case of read locks.**

12

## Distributed Deadlock

- u **More complicated if lock management is not centralized.**
- u **Local Wait-for-Graph (LWFG) may not show existence of deadlock.**
- u **May need to create GWFG, union of all LWFGs.**
- u **Look at three schemes:**
  - – **Centralized Deadlock Detection**
  - – **Hierarchical Deadlock Detection**
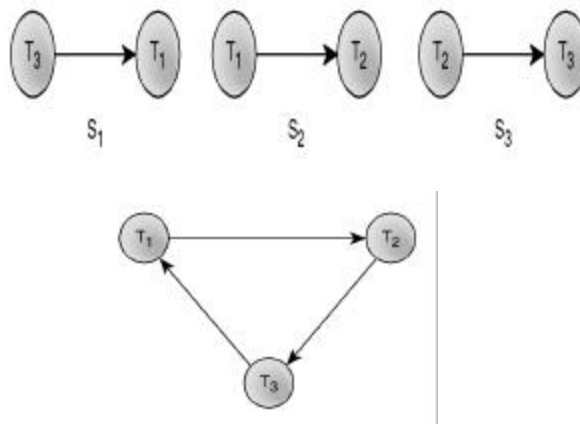  - – **Distributed Deadlock Detection.**

## Example - Distributed Deadlock

- $T_1$ initiated at site $S_1$ and creating agent at $S_2$,
- $T_2$ initiated at site $S_2$ and creating agent at $S_3$,
- $T_3$ initiated at site $S_3$ and creating agent at $S_1$.

| Time | $S_1$ | $S_2$ | $S_3$ |
|------|-------|-------|-------|
| $t_1$ | read_lock($T_1$, $x_1$) | write_lock($T_2$, $y_2$) | read_lock($T_3$, $z_3$) |
| $t_2$ | write_lock($T_1$, $y_1$) | write_lock($T_2$, $z_2$) | |
| $t_3$ | write_lock($T_3$, $x_1$) | write_lock($T_1$, $y_2$) | write_lock($T_2$, $z_3$) |

## Example - Distributed Deadlock

## Centralized Deadlock Detection

- u **Single site appointed deadlock detection coordinator (DDC).**
- u **DDC has responsibility of constructing and maintaining GWFG.**
- u **If one or more cycles exist, DDC must break each cycle by selecting transactions to be rolled back and restarted.**
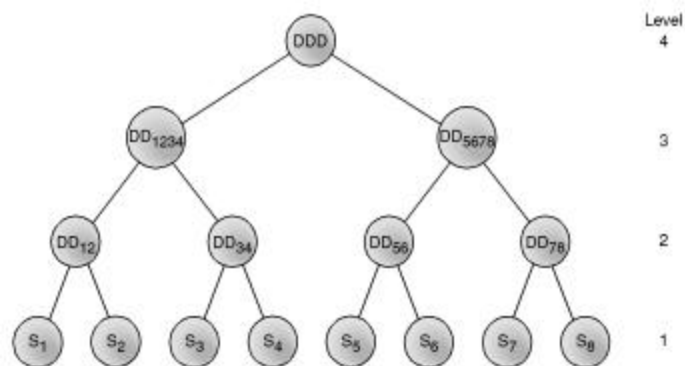
## Hierarchical Deadlock Detection

- u **Sites are organized into a hierarchy.**
- u **Each site sends its LWFG to detection site above it in hierarchy.**
- u **Reduces dependence on centralized detection site.**

## Hierarchical Deadlock Detection

## Two-Phase Commit (2PC)

- u **Two phases: a *voting phase* and a *decision phase*.**
- u **Coordinator asks all participants whether they are prepared to commit transaction.**
  - – **If one participant votes abort, or fails to respond within a timeout period, coordinator instructs all participants to abort transaction.**
  - – **If all vote commit, coordinator instructs all participants to commit.**
- u **All participants must adopt global decision .**

29

## Two-Phase Commit (2PC)

- u **If participant votes abort, free to abort transaction immediately**
- u **If participant votes commit, must wait for coordinator to broadcast global-commit or global-abort message.**
- u **Protocol assumes each site has its own local log and can rollback or commit transaction reliably.**
- u **If participant fails to vote, abort is assumed.**
- u **If participant gets no vote instruction from coordinator, can abort.**
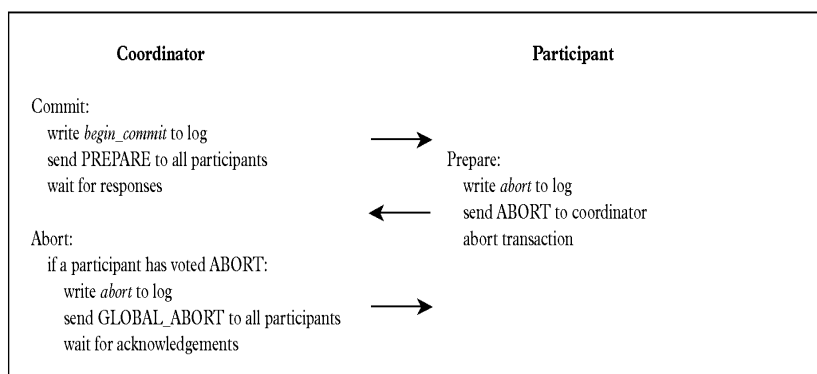
30

# 2PC Protocol for Participant Voting Commit

Coordinator                                        Participant

Commit:
    write *begin_commit* to log
    send PREPARE to all participants      ⟶     Prepare:
    wait for responses                            write *ready_commit* to log
                                    send READY_COMMIT to coordinator
                    ⟵          wait for GLOBAL_COMMIT or GLOBAL_ABORT
Ready_commit:
    if all participants have voted READY:
        write *commit* to log
        send GLOBAL_COMMIT to all participants    ⟶
        wait for acknowledgements            Global_commit:
                                    write *commit* record to log
                                    commit transaction
Ack:
    if all participants have acknowledged:    ⟵     send acknowledgements
        write *end_of_transaction* to log

(a)

31

# 2PC Protocol for Participant Voting Abort

Coordinator                                        Participant

Commit:
    write *begin_commit* to log
    send PREPARE to all participants     ⟶    Prepare:
    wait for responses                          write *abort* to log
                                    send ABORT to coordinator
                    ⟵     abort transaction
Abort:
    if a participant has voted ABORT:
        write *abort* to log
        send GLOBAL_ABORT to all participants    ⟶
        wait for acknowledgements

(b)

32

## Termination Protocols

- u **Invoked whenever a coordinator or participant fails to receive an expected message and times out.**

**Coordinator**

- u **Timeout in WAITING state**
  - – **Globally abort the transaction.**

- u **Timeout in DECIDED state**
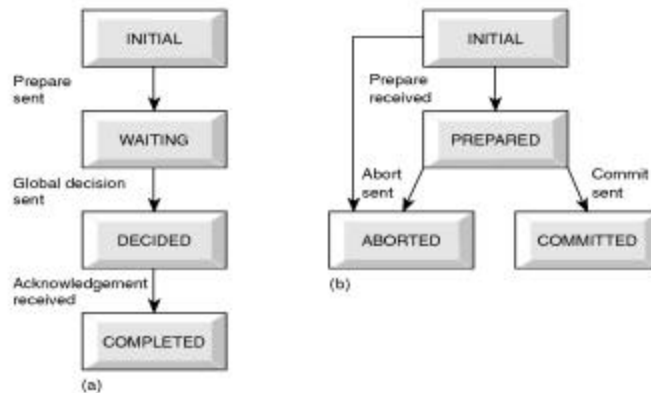  - – **Send global decision again to sites that have not acknowledged.**

33

## Termination Protocols - Participant

- u **Simplest termination protocol is to leave participant blocked until communication with the coordinator is re-established. Alternatively:**

- u **Timeout in INITIAL state**
  - – **Unilaterally abort the transaction.**

- u **Timeout in the PREPARED state**
  - – **Without more information, participant blocked.**
  - – **Could get decision from another participant .**

34

## State Transition Diagram for 2PC



35

---

## Recovery Protocols

u **Action to be taken by operational site in event of failure. Depends on what stage coordinator or participant had reached.**

**Coordinator Failure**

u **Failure in INITIAL state**
  – **Recovery starts the commit procedure.**

u **Failure in WAITING state**
  – **Recovery restarts the commit procedure.**

36

## 2PC - Coordinator Failure

- u **Failure in DECIDED state**
  - – **On restart, if coordinator has received all acknowledgements, it can complete successfully. Otherwise, has to initiate termination protocol discussed above.**

## 2PC - Participant Failure

- u **Objective to ensure that participant on restart performs same action as all other participants and that this restart can be performed independently.**
- u **Failure in INITIAL state**
  - – **Unilaterally abort the transaction.**
- u **Failure in PREPARED state**
  - – **Recovery via termination protocol above.**
- u **Failure in ABORTED/COMMITTED states**
  - – **On restart, no further action is necessary.**

## Three-Phase Commit (3PC)

- u **2PC *is not* a non-blocking protocol.**

- u **For example, a process that times out after voting commit, but before receiving global instruction, is blocked if it can communicate only with sites that do not know global decision.**

- u **Probability of blocking occurring in practice is sufficiently rare that most existing systems use 2PC.**

40

## Three-Phase Commit (3PC)

- u **Alternative non-blocking protocol, called *three-phase commit (3PC)* protocol.**

- u **Non-blocking for site failures, except in event of failure of all sites.**

- u **Communication failures can result in different sites reaching different decisions, thereby violating atomicity of global transactions.**

- u **3PC removes uncertainty period for participants who have voted commit and await global decision.**
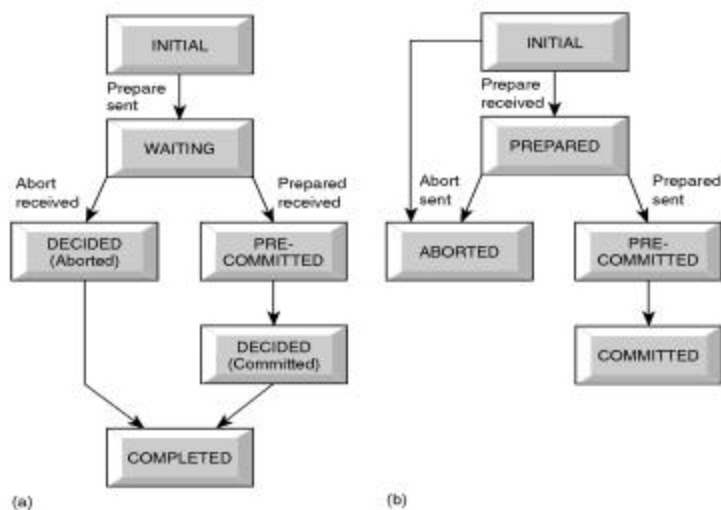
41

## Three-Phase Commit (3PC)

- u **Introduces third phase, called *pre-commit*, between voting and global decision.**
- u **On receiving all votes from participants, coordinator sends global pre-commit message.**
- u **Participant who receives global pre-commit, knows all other participants have voted commit and that, in time, participant itself will definitely commit.**

42

## State Transition Diagram for 3PC



(a)    (b)

43