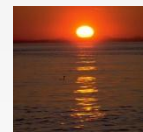




# Chapter 22: Distributed Databases

**Database System Concepts, 5th Ed.**

©Silberschatz, Korth and Sudarshan  
See [www.db-book.com](http://www.db-book.com) for conditions on re-use





# Chapter 22: Distributed Databases

- Heterogeneous and Homogeneous Databases
- Distributed Data Storage
- Distributed Transactions
- Commit Protocols
- Concurrency Control in Distributed Databases
- Availability
- Distributed Query Processing
- Heterogeneous Distributed Databases
- Directory Systems





# Distributed Database System

- A distributed database system consists of loosely coupled sites that share no physical component
- Database systems that run on each site are independent of each other
- Transactions may access data at one or more sites





# Homogeneous Distributed Databases

- In a homogeneous distributed database
  - All sites have identical software
  - Are aware of each other and agree to cooperate in processing user requests.
  - Each site surrenders part of its autonomy in terms of right to change schemas or software
  - Appears to user as a single system
- In a heterogeneous distributed database
  - Different sites may use different schemas and software
    - 4 Difference in schema is a major problem for query processing
    - 4 Difference in software is a major problem for transaction processing
  - Sites may not be aware of each other and may provide only limited facilities for cooperation in transaction processing





# Data Replication

- A relation or fragment of a relation is **replicated** if it is stored redundantly in two or more sites.
- **Full replication** of a relation is the case where the relation is stored at all sites.
- Fully redundant databases are those in which every site contains a copy of the entire database.





# Data Replication (Cont.)

- Advantages of Replication
  - **Availability:** failure of site containing relation  $r$  does not result in unavailability of  $r$  if replicas exist.
  - **Parallelism:** queries on  $r$  may be processed by several nodes in parallel.
  - **Reduced data transfer:** relation  $r$  is available locally at each site containing a replica of  $r$ .
- Disadvantages of Replication
  - Increased cost of updates: each replica of relation  $r$  must be updated.
  - Increased complexity of concurrency control: concurrent updates to distinct replicas may lead to inconsistent data unless special concurrency control mechanisms are implemented.
    - 4 One solution: choose one copy as **primary copy** and apply concurrency control operations on primary copy





# Data Fragmentation

- Division of relation  $r$  into fragments  $r_1, r_2, \dots, r_n$  which contain sufficient information to reconstruct relation  $r$ .
- **Horizontal fragmentation**: each tuple of  $r$  is assigned to one or more fragments
- **Vertical fragmentation**: the schema for relation  $r$  is split into several smaller schemas
  - All schemas must contain a common candidate key (or superkey) to ensure lossless join property.
  - A special attribute, the tuple-id attribute may be added to each schema to serve as a candidate key.
- Example : relation account with following schema
- $Account = (account\_number, branch\_name, balance)$





# Horizontal Fragmentation of *account* Relation

<i>account_numbe</i> <i>r</i>	<i>branch_name</i>	<i>balance</i>
A-305	Hillside	500
A-226	Hillside	336
A-155	Hillside	62

$$account_1 = \sigma_{branch\_name="Hillside"}(account)$$

<i>account_numbe</i> <i>r</i>	<i>branch_name</i>	<i>balance</i>
A-177	Valleyview	205
A-402	Valleyview	10000
A-408	Valleyview	1123
A-639	Valleyview	750

$$account_2 = \sigma_{branch\_name="Valleyview"}(account)$$







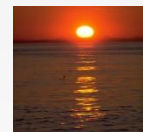
# Vertical Fragmentation of *employee\_info* Relation

<i>branch_name</i>	<i>customer_name</i>	<i>tuple_id</i>
Hillside	Lowman	1
Hillside	Camp	2
Valleyview	Camp	3
Valleyview	Kahn	4
Hillside	Kahn	5
Valleyview	Kahn	6
Valleyview	Green	7

$$deposit_1 = \Pi_{branch\_name, customer\_name, tuple\_id}(employee\_info)$$

<i>account_numbe</i>	<i>balance</i>	<i>tuple_id</i>
A-305	500	1
A-226	336	2
A-177	205	3
A-402	10000	4
A-155	62	5
A-408	1123	6
A-639	750	7

$$deposit_2 = \Pi_{account\_number, balance, tuple\_id}(employee\_info)$$





# Advantages of Fragmentation

- Horizontal:
  - allows parallel processing on fragments of a relation
  - allows a relation to be split so that tuples are located where they are most frequently accessed
- Vertical:
  - allows tuples to be split so that each part of the tuple is stored where it is most frequently accessed
  - tuple-id attribute allows efficient joining of vertical fragments
- Vertical and horizontal fragmentation can be mixed.
  - Fragments may be successively fragmented to an arbitrary depth.
- Replication and fragmentation can be combined
  - Relation is partitioned into several fragments: system maintains several identical replicas of each such fragment.





# Data Transparency

- **Data transparency:** Degree to which system user may remain unaware of the details of how and where the data items are stored in a distributed system
- Consider transparency issues in relation to:
  - Fragmentation transparency
  - Replication transparency
  - Location transparency
- Naming of data items: criteria
  1. Every data item must have a system-wide unique name.
  2. It should be possible to find the location of data items efficiently.
  3. It should be possible to change the location of data items transparently.
  4. Each site should be able to create new data items autonomously.





# Centralized Scheme - Name Server

- Structure:
  - name server assigns all names
  - each site maintains a record of local data items
  - sites ask name server to locate non-local data items
- Advantages:
  - satisfies naming criteria 1-3
- Disadvantages:
  - does not satisfy naming criterion 4
  - name server is a potential performance bottleneck
  - name server is a single point of failure





# Use of Aliases

- Alternative to centralized scheme: each site prefixes its own site identifier to any name that it generates i.e., *site 17.account*.
  - Fulfills having a unique identifier, and avoids problems associated with central control.
  - However, fails to achieve network transparency.
- Solution: Create a set of **aliases** for data items; Store the mapping of aliases to the real names at each site.
- The user can be unaware of the physical location of a data item, and is unaffected if the data item is moved from one site to another.





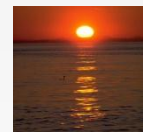
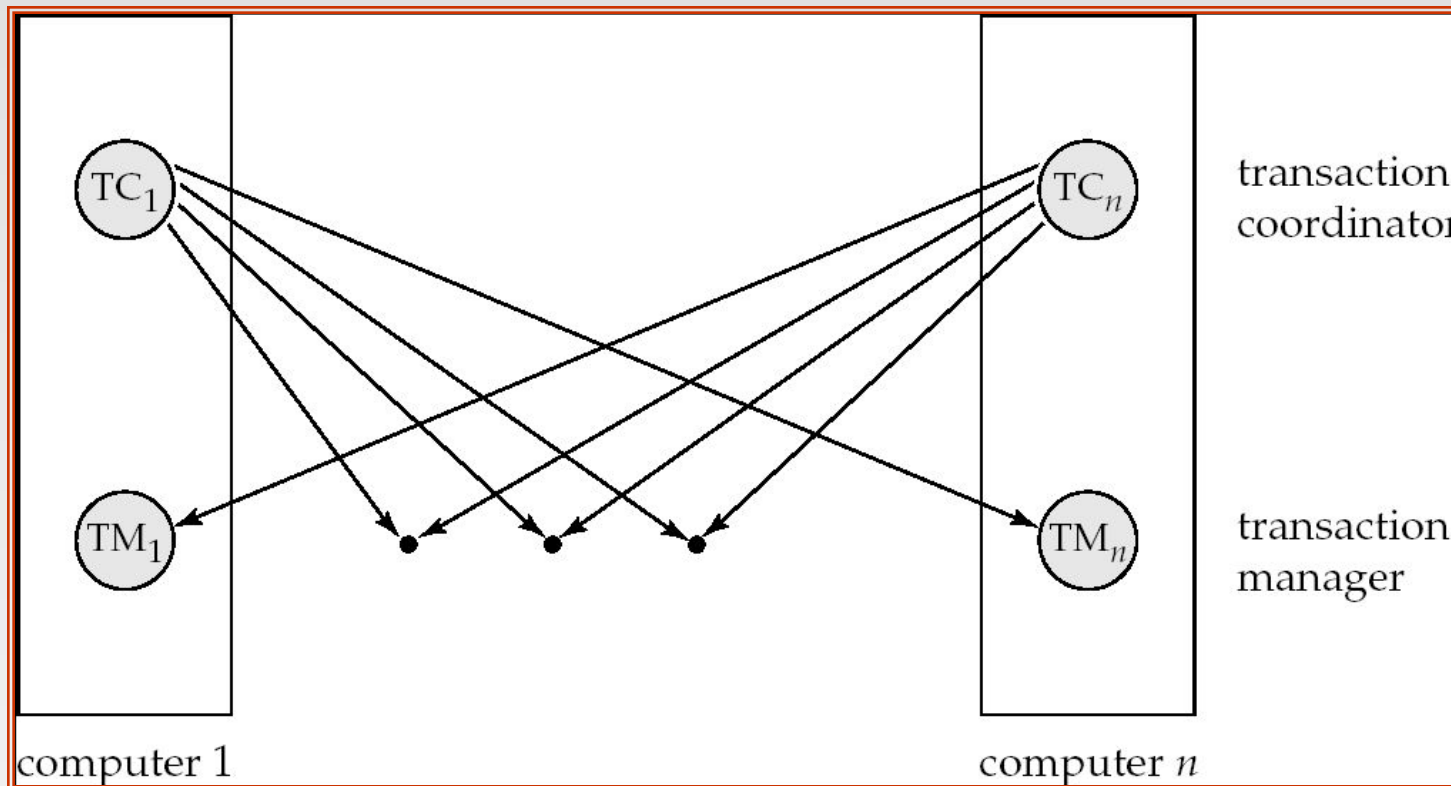
# Distributed Transactions

- Transaction may access data at several sites.
- Each site has a local **transaction manager** responsible for:
  - Maintaining a log for recovery purposes
  - Participating in coordinating the concurrent execution of the transactions executing at that site.
- Each site has a **transaction coordinator**, which is responsible for:
  - Starting the execution of transactions that originate at the site.
  - Distributing subtransactions at appropriate sites for execution.
  - Coordinating the termination of each transaction that originates at the site, which may result in the transaction being committed at all sites or aborted at all sites.





# Transaction System Architecture





# System Failure Modes

- Failures unique to distributed systems:
  - Failure of a site.
  - Loss of messages
    - 4 Handled by network transmission control protocols such as TCP-IP
  - Failure of a communication link
    - 4 Handled by network protocols, by routing messages via alternative links
  - **Network partition**
    - 4 A network is said to be **partitioned** when it has been split into two or more subsystems that lack any connection between them
      - Note: a subsystem may consist of a single node
- Network partitioning and site failures are generally indistinguishable.







# Commit Protocols

- Commit protocols are used to ensure atomicity across sites
  - a transaction which executes at multiple sites must either be committed at all the sites, or aborted at all the sites.
  - not acceptable to have a transaction committed at one site and aborted at another
- The *two-phase commit* (2PC) protocol is widely used
- The *three-phase commit* (3PC) protocol is more complicated and more expensive, but avoids some drawbacks of two-phase commit protocol. This protocol is not used in practice.





# Two Phase Commit Protocol (2PC)

- Assumes **fail-stop** model – failed sites simply stop working, and do not cause any other harm, such as sending incorrect messages to other sites.
- Execution of the protocol is initiated by the coordinator after the last step of the transaction has been reached.
- The protocol involves all the local sites at which the transaction executed
- Let  $T$  be a transaction initiated at site  $S_p$  and let the transaction coordinator at  $S_i$  be  $C_i$





# Phase 1: Obtaining a Decision

- Coordinator asks all participants to *prepare* to commit transaction  $T_i$ .
  - $C_i$  adds the records  $\langle \mathbf{prepare} T \rangle$  to the log and forces log to stable storage
  - sends  $\mathbf{prepare} T$  messages to all sites at which  $T$  executed
- Upon receiving message, transaction manager at site determines if it can commit the transaction
  - if not, add a record  $\langle \mathbf{no} T \rangle$  to the log and send  $\mathbf{abort} T$  message to  $C_i$
  - if the transaction can be committed, then:
    - add the record  $\langle \mathbf{ready} T \rangle$  to the log
    - force *all records* for  $T$  to stable storage
    - send  $\mathbf{ready} T$  message to  $C_i$





## Phase 2: Recording the Decision

- $T$  can be committed if  $C_i$  received a **ready**  $T$  message from all the participating sites: otherwise  $T$  must be aborted.
- Coordinator adds a decision record, **<commit  $T$ >** or **<abort  $T$ >**, to the log and forces record onto stable storage. Once the record stable storage it is irrevocable (even if failures occur)
- Coordinator sends a message to each participant informing it of the decision (commit or abort)
- Participants take appropriate action locally.





# Handling of Failures - Site Failure

When site  $S_i$  recovers, it examines its log to determine the fate of transactions active at the time of the failure.

- Log contain **<commit  $T$ >** record: site executes **redo** ( $T$ )
- Log contains **<abort  $T$ >** record: site executes **undo** ( $T$ )
- Log contains **<ready  $T$ >** record: site must consult  $C_i$  to determine the fate of  $T$ .
  - If  $T$  committed, **redo** ( $T$ )
  - If  $T$  aborted, **undo** ( $T$ )
- The log contains no control records concerning  $T$ 
  - implies that  $S_k$  failed before responding to the **prepare**  $T$  message from  $C_i$
  - $S_k$  must execute **undo** ( $T$ )





# Handling of Failures- Coordinator Failure

- If coordinator fails while the commit protocol for  $T$  is executing then participating sites must decide on  $T$ 's fate:
  1. If an active site contains a **<commit  $T$ >** record in its log, then  $T$  must be committed.
  2. If an active site contains an **<abort  $T$ >** record in its log, then  $T$  must be aborted.
  3. If some active participating site does not contain a **<ready  $T$ >** record in its log, then the failed coordinator  $C_i$  cannot have decided to commit  $T$ .
    1. Can therefore abort  $T$ .
  4. If none of the above cases holds, then all active sites must have a **<ready  $T$ >** record in their logs, but no additional control records (such as **<abort  $T$ >** or **<commit  $T$ >**).
    - In this case active sites must wait for  $C_i$  to recover, to find decision.
- **Blocking problem:** active sites may have to wait for failed coordinator to recover.





# Handling of Failures - Network Partition

- If the coordinator and all its participants remain in one partition, the failure has no effect on the commit protocol.
- If the coordinator and its participants belong to several partitions:
  - Sites that are not in the partition containing the coordinator think the coordinator has failed, and execute the protocol to deal with failure of the coordinator.
    - 4 No harm results, but sites may still have to wait for decision from coordinator.
- The coordinator and the sites are in the same partition as the coordinator think that the sites in the other partition have failed, and follow the usual commit protocol.
  - 4 Again, no harm results





# Coordinator Selection

- Backup coordinators
  - site which maintains enough information locally to assume the role of coordinator if the actual coordinator fails
  - executes the same algorithms and maintains the same internal state information as the actual coordinator fails executes state information as the actual coordinator
  - allows fast recovery from coordinator failure but involves overhead during normal processing.
- Election algorithms
  - used to elect a new coordinator in case of failures
  - Example: Bully Algorithm - applicable to systems where every site can send a message to every other site.







# Bully Algorithm

- If site  $S_i$  sends a request that is not answered by the coordinator within a time interval  $T$ , assume that the coordinator has failed  $S_i$  tries to elect itself as the new coordinator.
- $S_i$  sends an election message to every site with a higher identification number,  $S_i$  then waits for any of these processes to answer within  $T$ .
- If no response within  $T$ , assume that all sites with number greater than  $i$  have failed,  $S_i$  elects itself the new coordinator.
- If answer is received  $S_i$  begins time interval  $T'$ , waiting to receive a message that a site with a higher identification number has been elected.





# Bully Algorithm (Cont.)

- If no message is sent within  $T$ , assume the site with a higher number has failed;  $S_i$  restarts the algorithm.
- After a failed site recovers, it immediately begins execution of the same algorithm.
- If there are no active sites with higher numbers, the recovered site forces all processes with lower numbers to let it become the coordinator site, even if there is a currently active coordinator with a lower number.





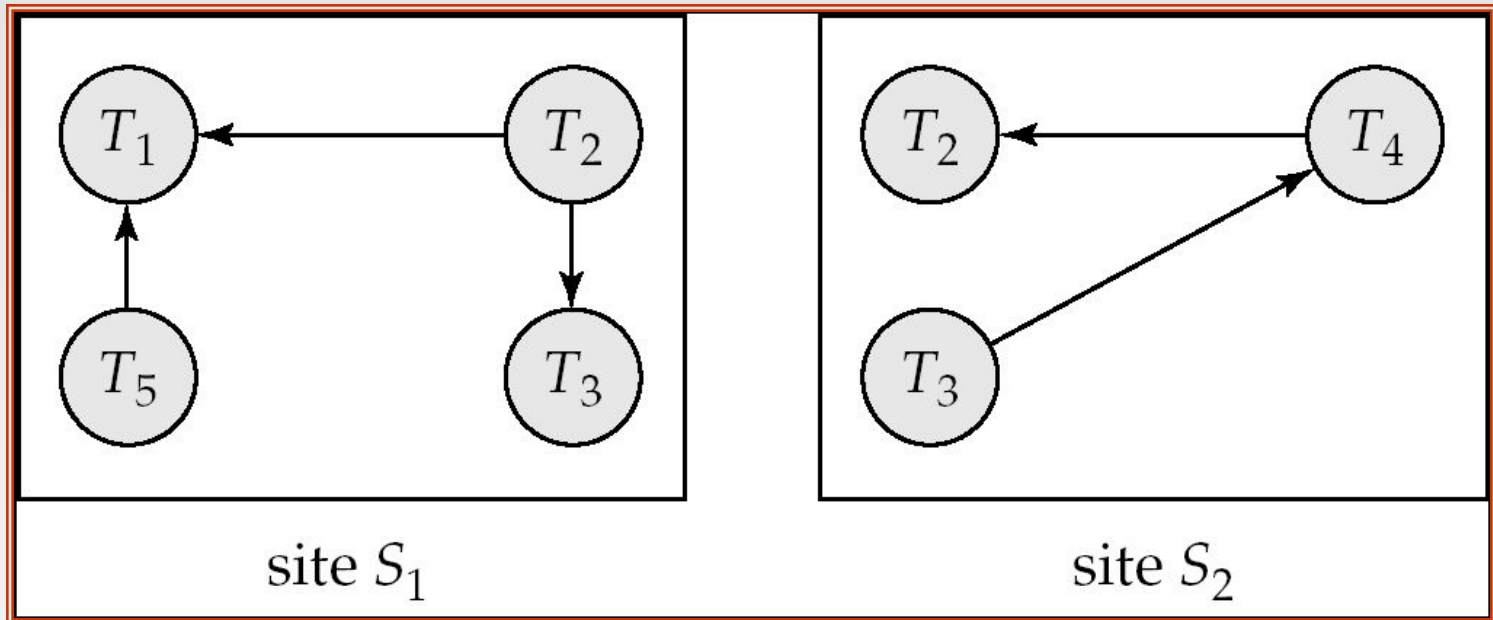
# Semijoin Strategy

- Let  $r_1$  be a relation with schema  $R_1$  stores at site  $S_1$   
Let  $r_2$  be a relation with schema  $R_2$  stores at site  $S_2$
- Evaluate the expression  $r_1 \bowtie r_2$  and obtain the result at  $S_1$ .
- 1. Compute  $temp_1 \leftarrow \Pi_{R_1 \cap R_2}(r_1)$  at  $S_1$ .
- 2. Ship  $temp_1$  from  $S_1$  to  $S_2$ .
- 3. Compute  $temp_2 \leftarrow r_2 \bowtie temp_1$  at  $S_2$
- 4. Ship  $temp_2$  from  $S_2$  to  $S_1$ .
- 5. Compute  $r_1 \bowtie temp_2$  at  $S_1$ . This is the same as  $r_1 \bowtie r_2$ .



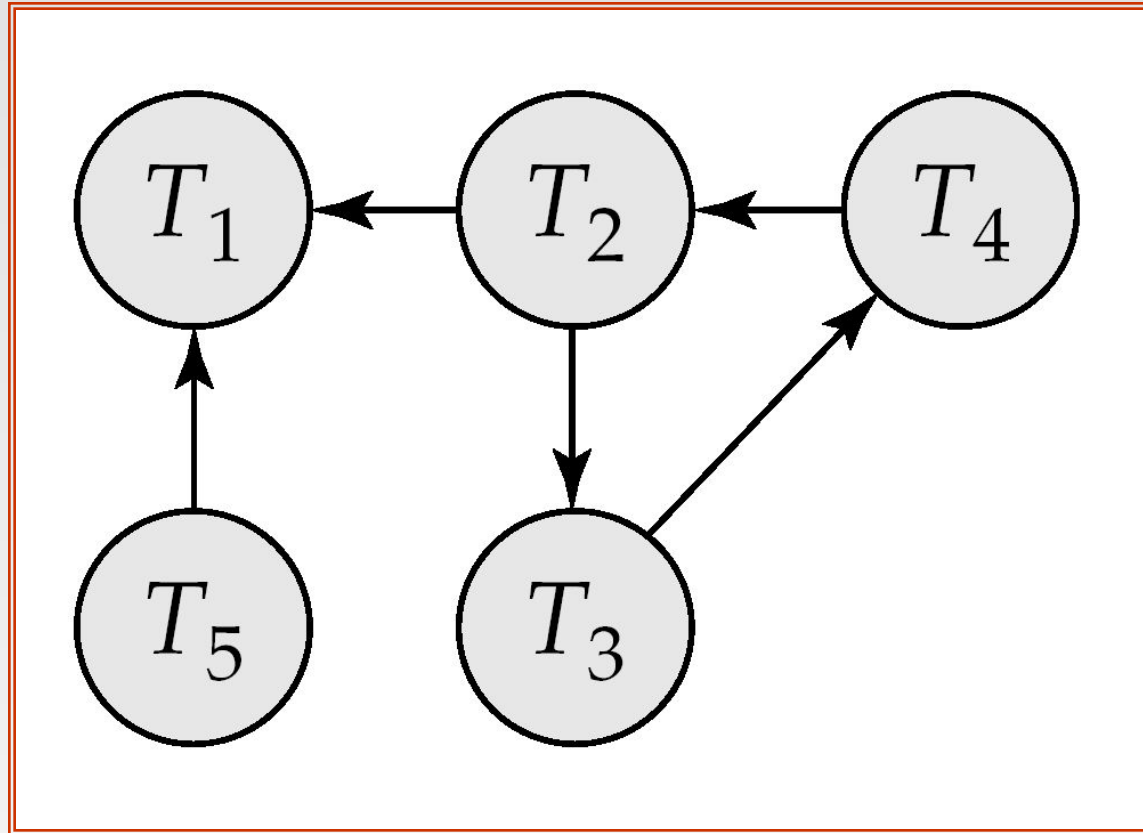


# Figure 22.3



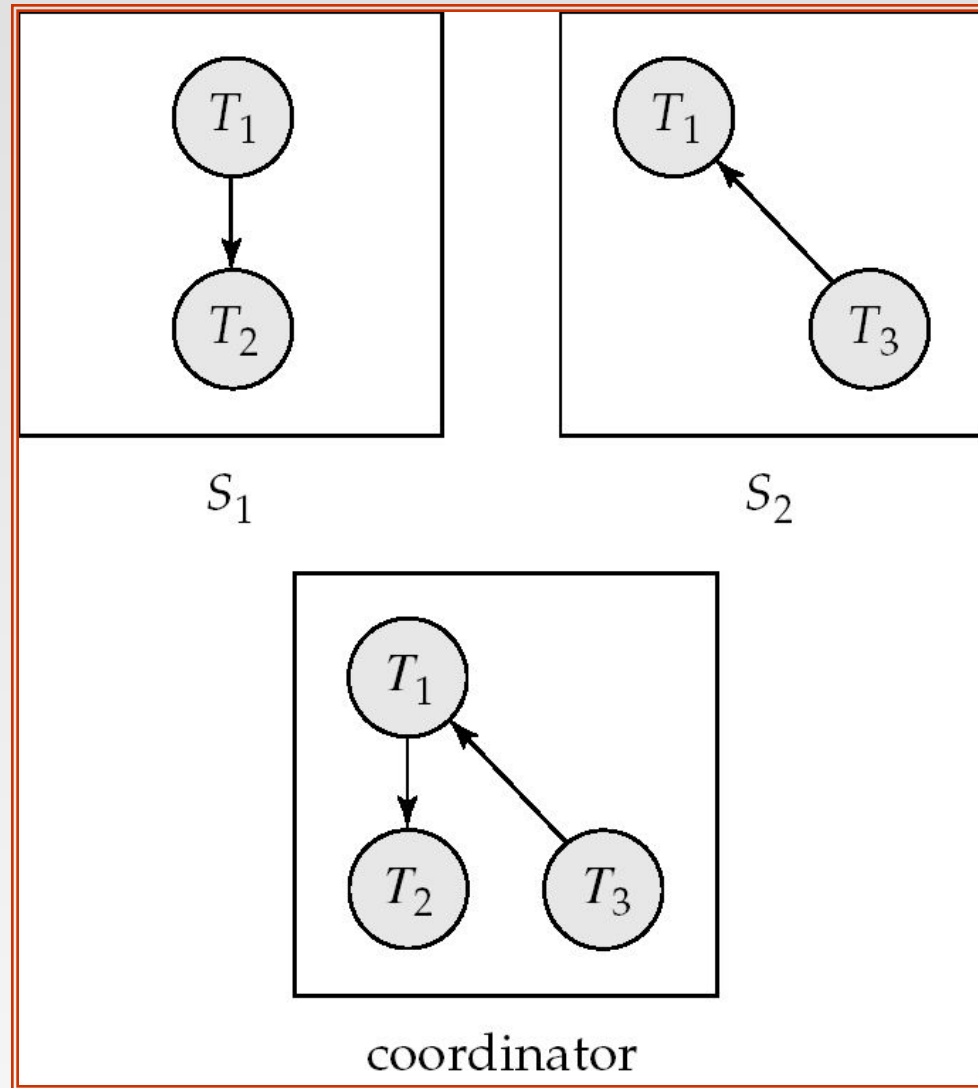


# Figure 22.4





# Figure 22.5





# Figure 22.7

<i>A</i>	<i>B</i>	<i>C</i>		<i>C</i>	<i>D</i>	<i>E</i>
1	2	3		3	4	5
4	5	6		3	6	8
1	2	4		2	3	2
5	3	2		1	4	1
8	9	7		1	2	3
<i>r</i>				<i>s</i>		

