

## **Chapter 19: Distributed Databases**

**Database System Concepts, 6th Ed.** 

©Silberschatz, Korth and Sudarshan See <a href="https://www.db-book.com">www.db-book.com</a> for conditions on re-use



### **Chapter 19: 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
    - Difference in schema is a major problem for query processing
    - 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



## **Distributed Data Storage**

- Assume relational data model
- Replication
  - System maintains multiple copies of data, stored in different sites, for faster retrieval and fault tolerance.
- Fragmentation
  - Relation is partitioned into several fragments stored in distinct sites
- Replication and fragmentation can be combined
  - Relation is partitioned into several fragments: system maintains several identical replicas of each such fragment.



## **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 is 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.
    - One solution: choose one copy as primary copy and apply concurrency control operations on primary copy



### **Data Fragmentation**

- Division of relation r into fragments  $r_1, r_2, ..., r_n$  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.



## Horizontal Fragmentation of account Relation

branch_name	account_number	balance
Hillside	A-305	500
Hillside	A-226	336
Hillside	A-155	62

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

branch_name	account_number	balance
Valleyview	A-177	205
Valleyview	A-402	10000
Valleyview	A-408	1123
Valleyview	A-639	750

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



#### Vertical Fragmentation of employee\_info Relation

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

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

account_number	balance	tuple_id
A-305 A-226 A-177 A-402 A-155 A-408	500 336 205 10000 62 1123	1 2 3 4 5 6
A-639	750	

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



### **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
- allows parallel processing on a relation
- Vertical and horizontal fragmentation can be mixed.
  - Fragments may be successively fragmented to an arbitrary depth.



## **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 and 2 Phase Commit

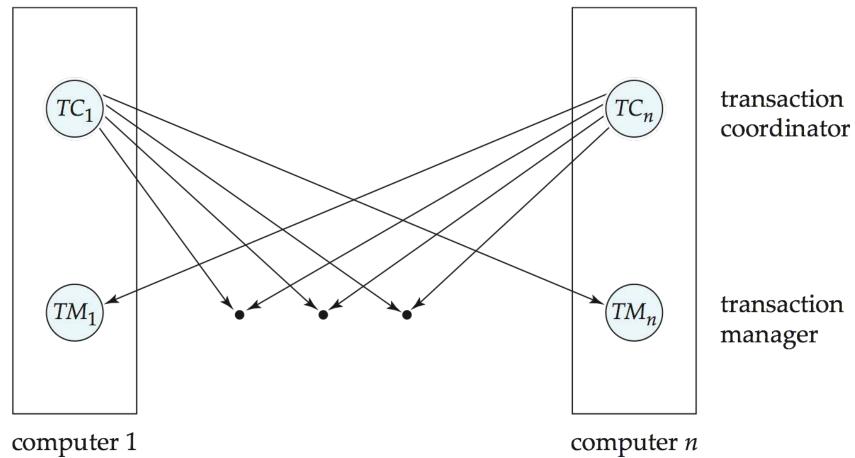


#### **Distributed Transactions**

- 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 massages
    - Handled by network transmission control protocols such as TCP-IP
  - Failure of a communication link
    - Handled by network protocols, by routing messages via alternative links
  - Network partition
    - 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  $S_i$ , 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<sub>i</sub>.
  - C<sub>i</sub> adds the records prepare T> to the log and forces log to stable storage
  - sends 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 <no T> to the log and send abort T message to C<sub>i</sub>
  - if the transaction can be committed, then:
  - add the record <ready T> to the log
  - force all records for T to stable storage
  - send ready T message to C<sub>i</sub>



## **Phase 2: Recording the Decision**

- T can be committed of  $C_i$  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: txn had completed, nothing to be done
- Log contains <abort T> record: txn had completed, nothing to be done
- Log contains <ready T> record: site must consult C<sub>i</sub> to determine the fate of T.
  - If T committed, redo (T); write < commit T> record
  - If T aborted, undo (T)
- The log contains no log records concerning T:
  - Implies that S<sub>k</sub> failed before responding to the prepare T message from C<sub>i</sub>
  - since the failure of S<sub>k</sub> precludes the sending of such a response, coordinator C₁ must abort T
  - $S_k$  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  $\langle ready \ T \rangle$  record in its log, then the failed coordinator  $C_i$  cannot have decided to commit T.
    - Can therefore abort T; however, such a site must reject any subsequent prepare T> message from C;
  - 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 <abord T> of <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.
    - 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.
  - Again, no harm results



## **Recovery and Concurrency Control**

- In-doubt transactions have a <ready T>, but neither a <commit T>, nor an <abort T> log record.
- The recovering site must determine the commit-abort status of such transactions by contacting other sites; this can slow and potentially block recovery.
- Recovery algorithms can note lock information in the log.
  - Instead of <ready T>, write out <ready T, L> L = list of locks held by T when the log is written (read locks can be omitted).
  - For every in-doubt transaction T, all the locks noted in the
     <ready T, L> log record are reacquired.
- After lock reacquisition, transaction processing can resume; the commit or rollback of in-doubt transactions is performed concurrently with the execution of new transactions.



## **Three Phase Commit (3PC)**

- Assumptions:
  - No network partitioning
  - At any point, at least one site must be up.
  - At most K sites (participants as well as coordinator) can fail
- Phase 1: Obtaining Preliminary Decision: Identical to 2PC Phase 1.
  - Every site is ready to commit if instructed to do so
- Phase 2 of 2PC is split into 2 phases, Phase 2 and Phase 3 of 3PC
  - In phase 2 coordinator makes a decision as in 2PC (called the pre-commit decision) and records it in multiple (at least K) sites
  - In phase 3, coordinator sends commit/abort message to all participating sites,
- Under 3PC, knowledge of pre-commit decision can be used to commit despite coordinator failure
  - Avoids blocking problem as long as < K sites fail</li>
- Drawbacks:
  - higher overheads
  - assumptions may not be satisfied in practice



# Alternative Models of Transaction Processing

- Notion of a single transaction spanning multiple sites is inappropriate for many applications
  - E.g. transaction crossing an organizational boundary
  - No organization would like to permit an externally initiated transaction to block local transactions for an indeterminate period
- Alternative models carry out transactions by sending messages
  - Code to handle messages must be carefully designed to ensure atomicity and durability properties for updates
    - Isolation cannot be guaranteed, in that intermediate stages are visible, but code must ensure no inconsistent states result due to concurrency
  - Persistent messaging systems are systems that provide transactional properties to messages
    - Messages are guaranteed to be delivered exactly once
    - Will discuss implementation techniques later



## **Alternative Models (Cont.)**

- Motivating example: funds transfer between two banks
  - Two phase commit would have the potential to block updates on the accounts involved in funds transfer
  - Alternative solution:
    - Debit money from source account and send a message to other site
    - Site receives message and credits destination account
  - Messaging has long been used for distributed transactions (even before computers were invented!)
- Atomicity issue
  - once transaction sending a message is committed, message must guaranteed to be delivered
    - Guarantee as long as destination site is up and reachable, code to handle undeliverable messages must also be available
      - e.g. credit money back to source account.
  - If sending transaction aborts, message must not be sent



# Error Conditions with Persistent Messaging

- Code to handle messages has to take care of variety of failure situations (even assuming guaranteed message delivery)
  - E.g. if destination account does not exist, failure message must be sent back to source site
  - When failure message is received from destination site, or destination site itself does not exist, money must be deposited back in source account
    - Problem if source account has been closed
      - get humans to take care of problem
- User code executing transaction processing using 2PC does not have to deal with such failures
- There are many situations where extra effort of error handling is worth the benefit of absence of blocking
  - E.g. pretty much all transactions across organizations



## **Persistent Messaging and Workflows**

- Workflows provide a general model of transactional processing involving multiple sites and possibly human processing of certain steps
  - E.g. when a bank receives a loan application, it may need to
    - Contact external credit-checking agencies
    - Get approvals of one or more managers
       and then respond to the loan application
  - We study workflows in Chapter 25
  - Persistent messaging forms the underlying infrastructure for workflows in a distributed environment



## Implementation of Persistent Messaging

#### Sending site protocol.

- When a transaction wishes to send a persistent message, it writes a record containing the message in a special relation messages\_to\_send; the message is given a unique message identifier.
- A message delivery process monitors the relation, and when a new message is found, it sends the message to its destination.
- The message delivery process deletes a message from the relation only after it receives an acknowledgment from the destination site.
  - If it receives no acknowledgement from the destination site, after some time it sends the message again. It repeats this until an acknowledgment is received.
  - If after some period of time, that the message is undeliverable, exception handling code provided by the application is invoked to deal with the failure.
- Writing the message to a relation and processing it only after the transaction commits ensures that the message will be delivered if and only if the transaction commits.



# Implementation of Persistent Messaging (Cont.)

#### Receiving site protocol.

- When a site receives a persistent message, it runs a transaction that adds the message to a received\_messages relation
  - provided message identifier is not already present in the relation
- After the transaction commits, or if the message was already present in the relation, the receiving site sends an acknowledgment back to the sending site.
  - Note that sending the acknowledgment before the transaction commits is not safe, since a system failure may then result in loss of the message.
- In many messaging systems, it is possible for messages to get delayed arbitrarily, although such delays are very unlikely.
  - Each message is given a timestamp, and if the timestamp of a received message is older than some cutoff, the message is discarded.
  - All messages recorded in the received messages relation that are older than the cutoff can be deleted.



# **Concurrency Control**



## **Concurrency Control**

- Modify concurrency control schemes for use in distributed environment.
- We assume that each site participates in the execution of a commit protocol to ensure global transaction automicity.
- We assume all replicas of any item are updated
  - Will see how to relax this in case of site failures later



## Single-Lock-Manager Approach

- System maintains a single lock manager that resides in a single chosen site, say S<sub>i</sub>
- When a transaction needs to lock a data item, it sends a lock request to S<sub>i</sub> and lock manager determines whether the lock can be granted immediately
  - If yes, lock manager sends a message to the site which initiated the request
  - If no, request is delayed until it can be granted, at which time a message is sent to the initiating site



# Single-Lock-Manager Approach (Cont.)

- The transaction can read the data item from any one of the sites at which a replica of the data item resides.
- Writes must be performed on all replicas of a data item
- Advantages of scheme:
  - Simple implementation
  - Simple deadlock handling
- Disadvantages of scheme are:
  - Bottleneck: lock manager site becomes a bottleneck
  - Vulnerability: system is vulnerable to lock manager site failure.



#### **Distributed Lock Manager**

- In this approach, functionality of locking is implemented by lock managers at each site
  - Lock managers control access to local data items
    - But special protocols may be used for replicas
- Advantage: work is distributed and can be made robust to failures
- Disadvantage: deadlock detection is more complicated
  - Lock managers cooperate for deadlock detection
    - More on this later
- Several variants of this approach
  - Primary copy
  - Majority protocol
  - Biased protocol
  - Quorum consensus



## **Primary Copy**

- Choose one replica of data item to be the primary copy.
  - Site containing the replica is called the primary site for that data item
  - Different data items can have different primary sites
- When a transaction needs to lock a data item Q, it requests a lock at the primary site of Q.
  - Implicitly gets lock on all replicas of the data item
- Benefit
  - Concurrency control for replicated data handled similarly to unreplicated data - simple implementation.
- Drawback
  - If the primary site of Q fails, Q is inaccessible even though other sites containing a replica may be accessible.



#### **Majority Protocol**

- Local lock manager at each site administers lock and unlock requests for data items stored at that site.
- When a transaction wishes to lock an unreplicated data item Q residing at site S<sub>i</sub>, a message is sent to S<sub>i</sub> 's lock manager.
  - If Q is locked in an incompatible mode, then the request is delayed until it can be granted.
  - When the lock request can be granted, the lock manager sends a message back to the initiator indicating that the lock request has been granted.



## **Majority Protocol (Cont.)**

- In case of replicated data
  - If Q is replicated at n sites, then a lock request message must be sent to more than half of the n sites in which Q is stored.
  - The transaction does not operate on Q until it has obtained a lock on a majority of the replicas of Q.
  - When writing the data item, transaction performs writes on all replicas.
- Benefit
  - Can be used even when some sites are unavailable
    - details on how handle writes in the presence of site failure later
- Drawback
  - Requires 2(n/2 + 1) messages for handling lock requests, and (n/2 + 1) messages for handling unlock requests.
  - Potential for deadlock even with single item e.g., each of 3 transactions may have locks on 1/3rd of the replicas of a data.



#### **Biased Protocol**

- Local lock manager at each site as in majority protocol, however, requests for shared locks are handled differently than requests for exclusive locks.
- **Shared locks**. When a transaction needs to lock data item *Q*, it simply requests a lock on *Q* from the lock manager at one site containing a replica of *Q*.
- Exclusive locks. When transaction needs to lock data item Q, it requests a lock on Q from the lock manager at all sites containing a replica of Q.
- Advantage imposes less overhead on read operations.
- Disadvantage additional overhead on writes



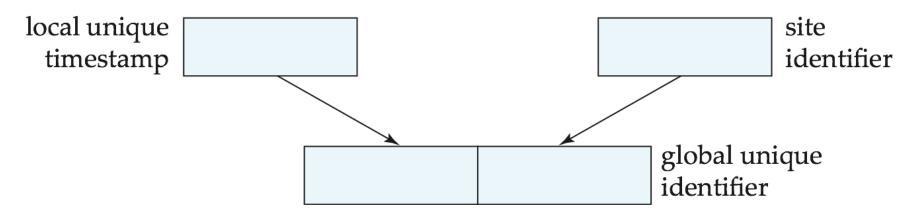
#### **Quorum Consensus Protocol**

- A generalization of both majority and biased protocols
- Each site is assigned a weight.
  - Let S be the total of all site weights
- Choose two values read quorum Q<sub>r</sub> and write quorum Q<sub>w</sub>
  - Such that  $Q_r + Q_w > S$  and  $2 * Q_w > S$
  - Quorums can be chosen (and S computed) separately for each item
- Each read must lock enough replicas that the sum of the site weights is >= Q<sub>r</sub>
- Each write must lock enough replicas that the sum of the site weights is >= Qw
- For now we assume all replicas are written
  - Extensions to allow some sites to be unavailable described later



#### **Timestamping**

- Timestamp based concurrency-control protocols can be used in distributed systems
- Each transaction must be given a unique timestamp
- Main problem: how to generate a timestamp in a distributed fashion
  - Each site generates a unique local timestamp using either a logical counter or the local clock.
  - Global unique timestamp is obtained by concatenating the unique local timestamp with the unique identifier.





# **Timestamping (Cont.)**

- A site with a slow clock will assign smaller timestamps
  - Still logically correct: serializability not affected
  - But: "disadvantages" transactions
- To fix this problem
  - Define within each site S<sub>i</sub> a logical clock (LC<sub>i</sub>), which generates the unique local timestamp
  - Require that S<sub>i</sub> advance its logical clock whenever a request is received from a transaction Ti with timestamp < x,y> and x is greater that the current value of LC<sub>i</sub>.
  - In this case, site S<sub>i</sub> advances its logical clock to the value x + 1.



#### Replication with Weak Consistency

- Many commercial databases support replication of data with weak degrees of consistency (I.e., without a guarantee of serializabiliy)
- E.g.: master-slave replication: updates are performed at a single "master" site, and propagated to "slave" sites.
  - Propagation is not part of the update transaction: its is decoupled
    - May be immediately after transaction commits
    - May be periodic
  - Data may only be read at slave sites, not updated
    - No need to obtain locks at any remote site
  - Particularly useful for distributing information
    - ▶ E.g. from central office to branch-office
  - Also useful for running read-only queries offline from the main database



#### Replication with Weak Consistency (Cont.)

- Replicas should see a transaction-consistent snapshot of the database
  - That is, a state of the database reflecting all effects of all transactions up to some point in the serialization order, and no effects of any later transactions.
- E.g. Oracle provides a create snapshot statement to create a snapshot of a relation or a set of relations at a remote site
  - snapshot refresh either by recomputation or by incremental update
  - Automatic refresh (continuous or periodic) or manual refresh



#### **Multimaster and Lazy Replication**

- With multimaster replication (also called update-anywhere replication) updates are permitted at any replica, and are automatically propagated to all replicas
  - Basic model in distributed databases, where transactions are unaware of the details of replication, and database system propagates updates as part of the same transaction
    - Coupled with 2 phase commit
- Many systems support lazy propagation where updates are transmitted after transaction commits
  - Allows updates to occur even if some sites are disconnected from the network, but at the cost of consistency



## **Deadlock Handling**

Consider the following two transactions and history, with item X and transaction  $T_1$  at site 1, and item Y and transaction  $T_2$  at site 2:

$$T_1$$
: write (X) write (Y)

 $T_2$ : write (Y) write (X)

X-lock on X write (X)	X-lock on Y write (Y) wait for X-lock on X	
Wait for X-lock on Y		

Result: deadlock which cannot be detected locally at either site

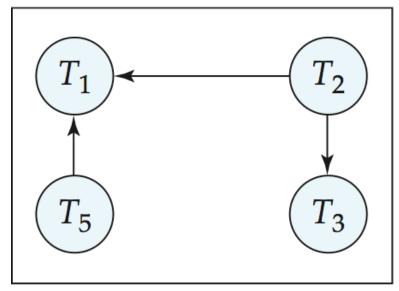


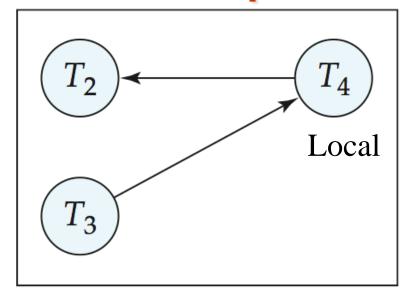
## **Centralized Approach**

- A global wait-for graph is constructed and maintained in a single site; the deadlock-detection coordinator
  - Real graph: Real, but unknown, state of the system.
  - Constructed graph: Approximation generated by the controller during the execution of its algorithm.
- the global wait-for graph can be constructed when:
  - a new edge is inserted in or removed from one of the local waitfor graphs.
  - a number of changes have occurred in a local wait-for graph.
  - the coordinator needs to invoke cycle-detection.
- If the coordinator finds a cycle, it selects a victim and notifies all sites.
  The sites roll back the victim transaction.

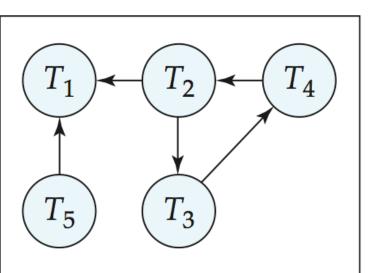


#### **Local and Global Wait-For Graphs**





site  $S_1$ 



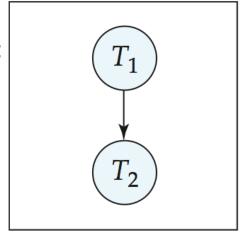
site  $S_2$ 

Global



#### **Example Wait-For Graph for False Cycles**

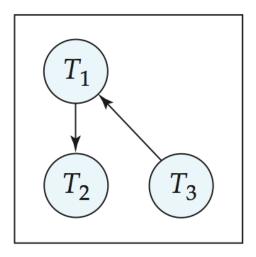
Initial state:



 $T_1$   $T_3$ 

 $S_1$ 

 $S_2$ 



coordinator



# False Cycles (Cont.)

- Suppose that starting from the state shown in figure,
  - 1.  $T_2$  releases resources at  $S_1$ 
    - resulting in a message remove  $T_1 \rightarrow T_2$  message from the Transaction Manager at site  $S_1$  to the coordinator)
  - 2. And then  $T_2$  requests a resource held by  $T_3$  at site  $S_2$ 
    - resulting in a message insert  $T_2 \rightarrow T_3$  from  $S_2$  to the coordinator
- Suppose further that the insert message reaches before the delete message
  - this can happen due to network delays
- The coordinator would then find a false cycle

$$T_1 \rightarrow T_2 \rightarrow T_3 \rightarrow T_1$$

- The false cycle above never existed in reality.
- False cycles cannot occur if two-phase locking is used.



#### **Unnecessary Rollbacks**

- Unnecessary rollbacks may result when deadlock has indeed occurred and a victim has been picked, and meanwhile one of the transactions was aborted for reasons unrelated to the deadlock.
- Unnecessary rollbacks can result from false cycles in the global waitfor graph; however, likelihood of false cycles is low.