

These slides are a modified version of the slides provided with the book Özsu and Valduriez, *Principles of Distributed Database Systems* (3rd Ed.), 2011 The original version of the slides is available at: extras.springer.com

Outline (distributed DB)

Introduction (Ch. 1) *

- Distributed Database Design (Ch. 3) *
- Distributed Query Processing (Ch. 6-8) *
- Distributed Transaction Management (Ch. 10-12) *
 Introduction to transaction management (Ch. 10) *
 - → Distributed Concurrency Control (Ch. 11) *
 - → Distributed DBMS Reliability (Ch. 12) *

* Özsu and Valduriez, Principles of Distributed Database Systems (3rd Ed.), 2011

Outline (today)

Distributed Concurrency Control (Ch. 11) *

- ➡ Serializability Theory
 - Formalization/Abstraction of Transactions
 - + Formalization/Abstraction of Concurrent Transactions (Histories)
 - Serial Histories
- ➡ Locking-based
- (strict) 2-phase Locking (2PL)
- Deadlock management

* Özsu and Valduriez, Principles of Distributed Database Systems (3rd Ed.), 2011

Concurrency Control

- The problem of synchronizing concurrent transactions such that the consistency of the database
 is maintained while, at the same time, maximum degree of concurrency is achieved
- This has to do with C(onsistency) and I(solation) from the ACID properties
- Consistency: assuming that each transaction is internally consistent (no integrity constraint violations) it is obtained by guaranteeing the right level of isolation (serializability)
 Isolation is obtained on the DR Mana
- Isolation: isolating transactions from one another in terms of their effects on the DB. More precisely, in terms of the effect on the DB of intermediate operations (before commit)
 Tradeoff between isolation and parallel execution (concurrency)
- Assumptions
 - System is fully reliable (no failures) we deal with reliability in Ch. 12*
 No data replication discussion on data replication is in Ch. 13* (we do not cover this chapter)
- Possible anomalies
- Lost updates
- The effects of some transactions are not reflected on the database
 Inconsistent retrievals
- + A transaction, if it reads the same data item more than once, should always read the same value

* Özsu and Valduriez, Principles of Distributed Database Systems (3rd Ed.), 2011

POSET's to Model Transactions

• We treat a transaction as a POSET (aka partially ordered set, partial order)

- POSET's are pairs $< \Sigma$, < > where
 - $\Rightarrow \Sigma$ is a set (domain)
 - $\Rightarrow \prec$ is a binary relation over $\sum (\prec \subseteq \sum x \sum)$ that is
 - irreflexive (not a < a, for all a)
 - asymmetric (a ≺ b implies not b ≺ a, for all a,b)
 - transitive (a ≤ b and b ≤ c implies a ≤ c, for all a,b,c)
- Operations are
 - DB operations (read or write): R(x), W(x) (where x is a data entity, e.g., a tuple) or
 - ➡ termination conditions (abort or commit): A, C

 A transaction is modeled as a partially ordered set of operations containing exactly one termination condition

Formalization/Abstraction of Transactions

- A transition is a POSET $T = \langle \Sigma, \langle \rangle$ where
 - Σ is finite: it is the set of operations of T
 O is the set of DB operations in Σ (operation that are not termination conditions)
 i.e., elements of O are of the kind R(x), W(x) where x is a data entity
 - Thus, $\Sigma = O \cup \{N\}$, where N ∈ $\{A, C\}$
 - 2 DB operations (elements of O) conflict iff they act on the same data entity x and one of them is a write W operation

(exactly 1 termination condition)

- W(x), R(x) are in conflict, W(x), W(x) are in conflict
 W(x), R(y) are NOT in conflict, W(x), W(y) are NOT in conflict, R(x), R(x) are NOT in conflict
- NOTICE: it is possible to have 2 distinct W(x) operations
- We assume implicit indices to make every operation uniqu
 Operations are atomic (indivisible units)
- Operations
- → < is s.t.</p>
 - order of conflicting operation is specified
 ✓ for all o₁, o₂ ∈ O: if o₁ and o₂ conflict, then either o₁ < o₂ or o₂ < o₁
 - all DB operations precede the unique termination condition
 - for all $o \in O$: $o \prec N$

Formalization/Abstraction of Transactions - cont'd • POSET's are DAG (directed acyclic graphs)

- We represent a transaction either way (as a POSET or as a DAG)
- The order of 2 conflicting operations is important and MUST be specified it specify the execution order between the 2 operations
- Operations that are not related can be executed in parallel
- A transaction might force other precedence order relations besides the ones between conflicting operations

These depend on application semantics













Conflict equivalence

- Definition. Two histories over the same set of transactions are conflict equivalent (or, simply equivalent) iff they agree on the execution order of the conflicting operations
 H₁ = < ∑₁, <₁ > and H₂ = < ∑₂, <₂ > with O and O' conflicting operations
 then O <₁O' if and only if O <₂O'
- (we are ignoring abort transaction to keep definition simpler)
- $H' = \{ W_2(x) \prec_{H'} W_2(y) \prec_{H'} R_2(z) \prec_{H'} R_1(x) \prec_{H'} W_1(x) \prec_{H'} R_3(x) \prec_{H'} R_3(y) \prec_{H'} R_3(z) \}$ = is NOT equivalent to $H_1 = \{ W_2(x) \prec_{H_1} R_1(x) \prec_{H_2} R_3(x) \prec_{H_1} W_1(x) \prec_{H_2} W_2(y) \prec_{H_1} R_3(y) \prec_{H_1} R_3(z) \prec_{H_1} R_3(z) \}$ + because $W_1(x) \prec_{H'} R_3(x)$ in H' but $R_3(x) \prec_{H_1} W_1(x)$ in H_1
- is equivalent to H₂ = { W₂(x) < H₂ R₁(x) < H₂ W₁(x) < H₂ R₃(x) < H₂ W₂(y) < H₂ R₃(y) < H₂ R₂(y) < H₂ R₂(z) < H₂ R₃(z) } (actually H', H₂, H₂ are all concurrent transaction executions)
 Definition. A history is serializable iff it is equivalent to a serial execution
- Therefore, H₂ is *serializable* (because H₂ is equivalent to a serial execution
 Therefore, H₂ is *serializable* (because H₂ is equivalent to H', which is a serial history)
 Primary function of a concurrency controller is to produce a serializable history over the set of pending transactions

Serializability in Distributed DBMS

- Somewhat more involved. Two histories have to be considered:
- local histories: histories over sets of transactions at the same site
 global history: union of local histories
- For global transactions (i.e., global history) to be serializable, two conditions are necessary:
- ⇒ Each local history should be serializable
- Identical local serialization order

* Özsu and Valduriez, Principles of Distributed Database Systems (3rd Ed.), 2011

























Deadlock Detection Transactions are allowed to wait freely Wait-for graphs and cycles Topologies for deadlock detection algorithms Centralized Distributed Hierarchical

Centralized Deadlock Detection

- One site is designated as the deadlock detector for the system. Each scheduler periodically sends its local WFG to the central site which merges them to a global WFG to determine cycles.
- How often to transmit?
 - Too often ⇒ higher communication cost but lower delays due to undetected deadlocks
- Too late ⇒ higher delays due to deadlocks, but lower communication cost
 Would be a reasonable choice if the concurrency control algorithm is also centralized (centralized 2PL)



Distributed Deadlock Detection

- Each site has a DD that maintain an LWFG
- Sites cooperate in detection of global deadlocks
- One example:
 - The local WFGs are formed at each site and passed on to other sites. Each local WFG is modified as follows:

 - Ocal WFG is modified as follows:
 Since each site receives the potential deadlock cycles from other sites, these edges are added to the local WFGs
 The edges in the local WFG which show that local transactions are waiting for transactions at other sites are joined with edges in the local WFGs which show that remote transactions are waiting for local ones
 - ➡ Each local deadlock detector:
 - ach local deadlock defector:
 looks for a cycle that does not involve the external edge. If it exists, there is a local deadlock which can be handled locally
 looks for a cycle involving the external edge. If it exists, it indicates a potential global deadlock. Pass on the information to the next site