

Relational ALGEBRA

↳ basic set of operations for relational model

UNARY OPERATIONS

↳ Select $\sigma_{\text{selection condition}}$

LIST \leftarrow ↳ Project $\pi_{\text{projection list}}$

↳ Rename $\rho_{\text{new name}}$

SET THEORY OPERATIONS

↳ Union \cup

↳ Intersection \cap

↳ Difference $-$

↳ Cartesian Product \times

need to be
TYPE COMPATIBLE

BINARY OPERATIONS

↳ JOIN \bowtie

↳ Division \div $\rightarrow \text{all, every, each}$

ADDITIONAL OPERATIONS

↳ OUTER JOINS

↳ OUTER UNION

↳ Aggregate function



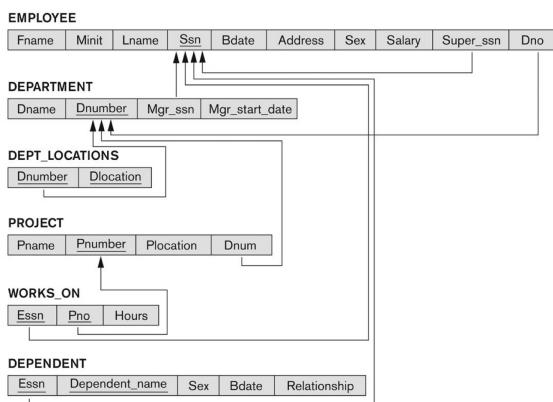
TYPE COMPATIBLE

↳ same no of attributes

↳ same domain for attributes

$$\text{dom}(A_i) = \text{dom}(B_j) \text{ for } i=1,2,\dots,n$$

The set of operations including SELECT σ , PROJECT π , UNION \cup , DIFFERENCE $-$, RENAME ρ , and CARTESIAN PRODUCT \times is called a *complete set* because any other



UNARY OPERATIONS

1. Select (δ) $\delta <\text{condition}_1>, <\text{condition}_2>, \dots, <\text{condition}_n>$ (R)

↳ filter out tuples

- Select the EMPLOYEE tuples whose department number is 4:

$\delta_{DNO=4}$ (EMPLOYEE)

- Select the employee tuples whose salary is greater than \$30,000:

$\delta_{SALARY > 30,000}$ (EMPLOYEE)

2. Project (π) $\pi <\text{attribute list}>$ (R)

↳ keeps certain columns

↳ removes duplicate tuples

P

π

δ

δ) Select fname, lname, salary from Employee E where Dno = 5 AND salary > 300k

1. Single relational Algebra expression

$\pi_{fname, lname, salary}(\delta_{DNO=5 \text{ AND } SALARY > 300k}$ (EMPLOYEE))

2. Immediate result relations expression

DEP5_Emp $\leftarrow \delta_{DNO=5 \text{ AND } SALARY > 300k}$ (EMPLOYEE)

final result $\leftarrow \pi_{fname}$ (DEP5_Emp)

3. Rename (P) $P_s(B_1, B_2, \dots, B_n)$ (R)

↳ changes relation name to S
↳ changes column names to B_1, B_2, \dots, B_n

$P(A_1, A_2)$ (EMPLOYEE (A₁, A₂)) \rightarrow A renamed to A₁,
B renamed to B₁

- If we write:
 - RESULT $\leftarrow \pi_{FNAME, LNAME, SALARY}$ (DEP5_EMPS)
 - RESULT will have the same attribute names as DEP5_EMPS (same attributes as EMPLOYEE)
- If we write:
 - RESULT (F, M, L, S, B, A, SX, SAL, SU, DNO) $\leftarrow P_{RESULT(F, M, L, S, B, A, SX, SAL, SU, DNO)}$ (DEP5_EMPS)
 - The 10 attributes of DEP5_EMPS are renamed to F, M, L, S, B, A, SX, SAL, SU, DNO, respectively

Note: the \leftarrow symbol is an assignment operator

A ₁	B ₁	C	D
A ₁	B ₁		

Set Theory Operations

1. Union (\cup) $R \cup S$

↳ combine R and S

↳ no duplicates

↳ must be type compatible

- Example:
 - To retrieve the social security numbers of all employees who either work in department 5 (RESULT1 below) or directly supervise an employee who works in department 5 (RESULT2 below)
 - We can use the UNION operation as follows:
$$\begin{aligned} DEP5_EMPS &\leftarrow \sigma_{DNO=5}(\text{EMPLOYEE}) \\ \text{RESULT1} &\leftarrow \pi_{SSN}(DEP5_EMPS) \\ \text{RESULT2}(SSN) &\leftarrow \pi_{\text{SUPERSSN}}^{\text{EMPLOYEE}}(DEP5_EMPS) \\ \text{RESULT} &\leftarrow \text{RESULT1} \cup \text{RESULT2} \end{aligned}$$
 - The union operation produces the tuples that are in either RESULT1 or RESULT2 or both

Some properties of UNION, INTERSECT, and DIFFERENCE

- Notice that both union and intersection are commutative operations; that is
 - $R \cup S = S \cup R$, and $R \cap S = S \cap R$
- Both union and intersection can be treated as n-ary operations applicable to any number of relations as both are associative operations; that is
 - $R \cup (S \cup T) = (R \cup S) \cup T$
 - $(R \cap S) \cap T = R \cap (S \cap T)$
- The minus operation is not commutative; that is, in general
 - $R - S \neq S - R$

2. Intersection (\cap) $R \cap S$

↳ common tuples of R and S

↳ must be type compatible

3. Difference - $R - S$

↳ all tuples that are in R but not S

↳ must be type compatible

4. Cartesian Product \times $R(A_1, A_2, \dots, A_n) \times S(B_1, B_2, \dots, B_n)$

↳ combine tuples from 2 diff relations

↳ dont need to be type compatible

- Generally, CROSS PRODUCT is not a meaningful operation
 - Can become meaningful when followed by other operations
- Example (not meaningful):
 - $FEMALE_EMPS \leftarrow \sigma_{SEX='F'}(\text{EMPLOYEE})$
 - $EMPNAMESS \leftarrow \pi_{\text{FNAME}, \text{LNAME}, \text{SSN}}(FEMALE_EMPS)$
 - $EMP_DEPENDENTS \leftarrow EMPNAMESS \times \text{DEPENDENT}$
 - $EMP_DEPENDENTS$ will contain every combination of EMPNAMESS and DEPENDENT
 - whether or not they are actually related
- To keep only combinations where the DEPENDENT is related to the EMPLOYEE, we add a SELECT operation as follows
- Example (meaningful):
 - $FEMALE_EMPS \leftarrow \sigma_{SEX='F'}(\text{EMPLOYEE})$
 - $EMPNAMESS \leftarrow \pi_{\text{FNAME}, \text{LNAME}, \text{SSN}}(FEMALE_EMPS)$
 - $EMP_DEPENDENTS \leftarrow EMPNAMESS \times \text{DEPENDENT}$
 - $ACTUAL_DEPS \leftarrow \sigma_{\text{SSN}=\text{ESN}}(\text{EMP_DEPENDENTS})$
 - $RESULT \leftarrow \pi_{\text{FNAME}, \text{LNAME}, \text{DEPENDENT_NAME}}(ACTUAL_DEPS)$
 - RESULT will now contain the name of female employees and their dependents

Binary Operations

1. Division $\div R(A) \div S(B)$

L applied on 2 relations

DIVISION Operation

- The division operation is applied to two relations
- $R(Z) \div S(X)$, where X subset Z . Let $Y = Z - X$ (and hence $Z = X \cup Y$); that is, let Y be the set of attributes of R that are not attributes of S .
- The result of DIVISION is a relation $T(Y)$ that includes a tuple t if tuples t_R appear in R with $t_R[Y] = t$, and with $t_R[X] = t_s$ for every tuple t_s in S .
- For a tuple t to appear in the result T of the DIVISION, the values in t must appear in R in combination with every tuple in S .

$$R \div S \\ R \leftarrow \pi_{SSN} (\text{EMPLOYEE})$$

- Q1: Retrieve the name and address of all employees who work for the 'Research' department.

RESEARCH_DEPT $\leftarrow \sigma_{DNAME='Research'} (\text{DEPARTMENT})$
 RESEARCH_EMPS $\leftarrow (\text{RESEARCH_DEPT} \bowtie_{\text{DNUMBER}=\text{DNUMBER}} \text{EMPLOYEE})$
 RESULT $\leftarrow \pi_{FNAME, LNAME, ADDRESS} (\text{RESEARCH_EMPS})$

- Q6: Retrieve the names of employees who have no dependents.

ALL_EMPS $\leftarrow \pi_{SSN} (\text{EMPLOYEE})$
 EMPS_WITH_DEPS(SSN) $\leftarrow \pi_{ESSN} (\text{DEPENDENT})$
 EMPS_WITHOUT_DEPS $\leftarrow (\text{ALL_EMPS} - \text{EMPS_WITH_DEPS})$
 RESULT $\leftarrow \pi_{LNAME, FNAME} (\text{EMPS_WITHOUT_DEPS} * \text{EMPLOYEE})$

As a single expression, these queries become:

- Q1: Retrieve the name and address of all employees who work for the 'Research' department.

$\pi_{Fname, Lname, Address} (\sigma_{Dname='Research'} (\text{DEPARTMENT} \bowtie_{\text{Dnumber}=\text{Dno}} \text{EMPLOYEE}))$

- Q6: Retrieve the names of employees who have no dependents.

$\pi_{Lname, Fname} ((\pi_{SSN} (\text{EMPLOYEE}) - \rho_{SSN} (\pi_{ESSN} (\text{DEPENDENT}))) * \text{EMPLOYEE})$

2. JOIN $\bowtie R \bowtie_{\text{join condition}} S$

Binary Relational Operations (contd.)

- Example: Suppose that we want to retrieve the name of the manager of each department.

To get the manager's name, we need to combine each DEPARTMENT tuple with the EMPLOYEE tuple whose SSN value matches the MGRSSN value in the department tuple.

- We do this by using the join \bowtie operation.

$\text{DEPT_MGR} \leftarrow \text{DEPARTMENT} \bowtie_{\text{MGRSSN}=SSN} \text{EMPLOYEE}$

MGRSSN=SSN is the join condition

- Combines each department record with the employee who manages the department

The join condition can also be specified as
 $\text{DEPARTMENT.MGRSSN} = \text{EMPLOYEE.SSN}$

Equi Join $R \bowtie_{\text{join condition}} S$

Natural Join $R \bowtie_{\text{conditions}} S$

- Example: To apply a natural join on the DNUMBER attributes of DEPARTMENT and DEPT_LOCATIONS, it is sufficient to write:

$\text{DEPT_LOC} \leftarrow \text{DEPARTMENT} * \text{DEPT_LOCATIONS}$

- Only attribute with the same name is DNUMBER

An implicit join condition is created based on this attribute:
 $\text{DEPARTMENT.DNUMBER} = \text{DEPT_LOCATIONS.DNUMBER}$

- Another example: $Q \leftarrow R(A,B,C,D) * S(C,D,E)$

- The implicit join condition includes each pair of attributes with the same name, AND'd together:

$R.C=S.C \text{ AND } R.D=S.D$

- Result keeps only one attribute of each such pair:

$Q(A,B,C,D,E)$

Use of the Aggregate Functional operation \mathcal{F}

- $\mathcal{F}_{MAX \text{ Salary}} (\text{EMPLOYEE})$ retrieves the maximum salary value from the EMPLOYEE relation

- $\mathcal{F}_{MIN \text{ Salary}} (\text{EMPLOYEE})$ retrieves the minimum salary value from the EMPLOYEE relation

- $\mathcal{F}_{SUM \text{ Salary}} (\text{EMPLOYEE})$ retrieves the sum of the salary from the EMPLOYEE relation

- $\mathcal{F}_{COUNT \text{ SSN}, AVERAGE \text{ Salary}} (\text{EMPLOYEE})$ computes the count (number) of employees and their average salary

- Note: count just counts the number of rows, without removing duplicates

- Grouping attribute placed to left of symbol
- Aggregate functions to right of symbol
- $\text{DNO } \mathcal{F}_{COUNT \text{ SSN}, AVERAGE \text{ Salary}} (\text{EMPLOYEE})$

Single User DBMS

- ↳ one user at a time
- ↳ e.g. home computer

Multi User DBMS

- ↳ many users can access system at a time
- ↳ e.g. airline reservation system, banks, supermarket

Multiprogramming: execute multiple processes concurrently

→ execute some commands of P_1 , then suspend
execute some commands of P_2 , then suspend
resume P_1

Interleaved Processing

- ↳ concurrent execution of processes
- ↳ keep CPU busy, when I/O operation required by switching to another process instead of remaining idle
- ↳ prevents long process from delaying other processes

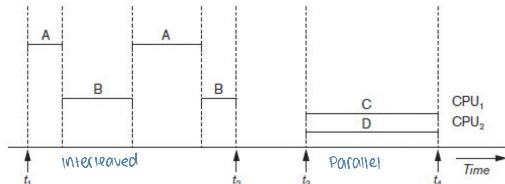


Figure 20.1 Interleaved processing versus parallel processing of concurrent transactions

Transaction

↳ atomic operation
either be completed or not done at all

- ↳ an executing program that forms a logical unit of DB processing
- ↳ has many operations
 1. Insertion
 2. Deletion
 3. Modification
 4. Retrieval

Specifying transaction boundaries

1. begin transaction statement
2. end transaction statement

Types of Transaction

1. Read Only Transaction → only retrieve no update
2. Read Write Transaction → retrieve + update

Parallel Processing

- ↳ uses multiple CPUs
- ↳ parallel processing of multiple processes

Database

- ↳ a collection of named data items
 - ↳ its size is called Granularity
- ↳ data items
 1. DB Record
 2. Disk block
 3. Attribute value of a DB record

Transaction Processing Systems

- ↳ systems with large databases
- ↳ has many concurrent users
- ↳ requires high availability, fast response time
- ↳ these are independent of item granularity

Read-item(x)

- Reads a database item named X into a program variable named X

Steps:

- Find the address of the disk block that contains item X .
- Copy that disk block into a buffer in main memory (if that disk block is not already in some main memory buffer). The size of the buffer is the same as the disk block size.
- Copy item X from the buffer to the program variable named X .

Write-item(x)

- Writes the value of program variable X into the database item named X
- Process includes finding the address of the disk block, copying to and from a memory buffer, and storing the updated disk block back to disk

- Find the address of the disk block that contains item X .
- Copy that disk block into a buffer in main memory (if that disk block is not already in some main memory buffer).
- Copy item X from the program variable named X into its correct location in the buffer.
- Store the updated disk block from the buffer back to disk (either immediately or at some later point in time).

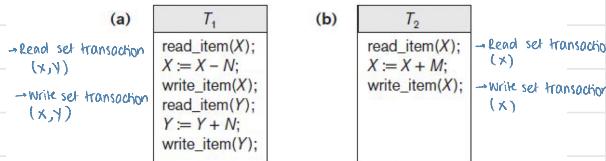


Figure 20.2 Two sample transactions (a) Transaction T_1 (b) Transaction T_2

Concurrency Control

- transactions by multiple users may execute concurrently leading to access and update of same DB items causing an inconsistent DB

Problems

1. LOST UPDATE Problem
2. TEMPORARY UPDATE Problem
3. INCORRECT SUMMARY Problem
4. UNREPEATABLE READ Problem

DBMS Buffers

- DBMS will maintain several main memory data buffers in the database cache
- When buffers are occupied, a buffer replacement policy is used to choose which buffer will be replaced
 - Example policy: least recently used

1. LOST UPDATE PROBLEM → nullifying update of first transaction

- ↳ T_1, T_2 perform interleaved operation on same data item

- ↳ resulting in incorrect value of DB item

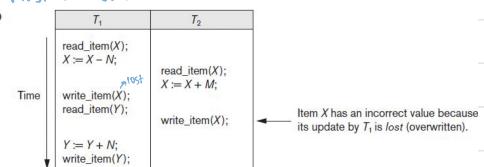


Figure 20.3 Some problems that occur when concurrent execution is uncontrolled (a) The lost update problem

2. TEMPORARY UPDATE / DIRTY READ PROBLEM → reading uncommitted data

- ↳ T_1 updates a X and the transaction fails *for some reason*
 - ↳ Meanwhile updated transaction is read by T_2 before its changed back to its old value
- roll back

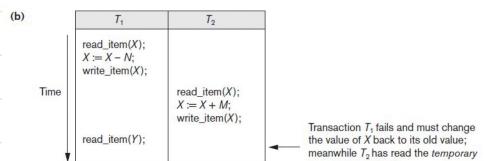


Figure 20.3 (cont'd.) Some problems that occur when concurrent execution is uncontrolled (b) The temporary update problem

3. INCORRECT SUMMARY PROBLEM → invalid result on aggregating data

- ↳ T_1 calculating an aggregate summary function on multiple DB items
- ↳ while other transactions are updating some of these items
- ↳ the aggregate function may calculate some values before and some values after they are updated

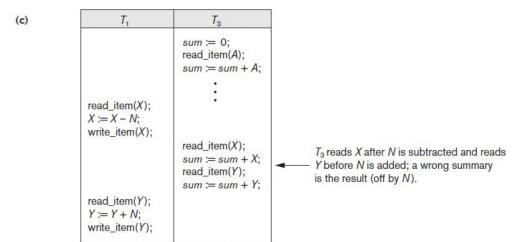


Figure 20.3 (cont'd.) Some problems that occur when concurrent execution is uncontrolled (c) The incorrect summary problem

4. UNREPEATABLE READ PROBLEMS → diff value for same data set variable

- ↳ T_1 reads same DB item twice
- ↳ while T_2 changed its value before T_1 read it the second time
- ↳ so T_1 receive 2 diff values of same item

WHY IS RECOVERY needed

↳ To make sure

1. Committed Transactions: effects are permanently recorded in DB
2. Aborted Transactions: don't affect the DB

Types of Transaction Failures

1. Computer failure, system crash → during transaction execution
hardware/software/network error
2. Transaction/system error → integer overflow, bad casts, logical programming error
3. Local errors → delete item not found
4. Concurrency control enforcement → insufficient account balance in banking
may cause withdraw fund be canceled
aborts transactions
serializing violations in denotates
→ restarted automatically later
5. Disk failure → block disks lost their data
due to read/write head crash
6. Physical Problems, catastrophes → power conditioning failure
fire, fire, theft

* Failure 1,4 occurs
System must keep sufficient info to recover quickly

common

long recovery time

Transaction States

↳ BEGIN_TRANSACTION

↳ READ/WRITE

↳ END_TRANSACTION

↳ COMMIT_TRANSACTION

↳ ROLLBACK/ABORT

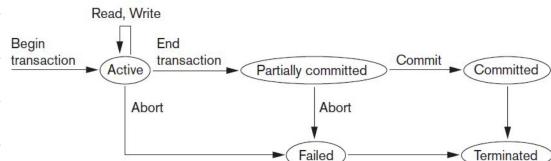


Figure 20.4 State transition diagram illustrating the states for transaction execution

SYSTEM log

↳ to be able to recover from system failures

- ↳ keeps track of transaction operations
- ↳ its a sequential, append only file
- ↳ not affected by failure → except disk/catastrophic
- ↳ Periodically blocked up
- ↳ allows undo/redo operations based on log

log buffer

- ↳ main memory buffers
- ↳ when full, appends to the end of system log file

Commit Point of Transaction

- ↳ when all operations completed successfully
- ↳ and effect of operation stored in log
- Transaction writes a commit record into the log
 - If system failure occurs, can search for transactions with recorded start_transaction but no commit record
- Force-writing the log buffer to disk
 - Writing log buffer to disk before transaction reaches commit point

DBMS Specific Buffer Replacement Policy

- * DBMS Cache: hold disk pages that are currently being processed in main memory buffers
- * Domains: DBMS cache divided into domains

- Page replacement policy
 - Selects particular buffers to be replaced when all are full
- Domain separation (DS) method → LRU page replacement
 - DBMS CACHE DIVIDED INTO DOMAINS
 - Each domain handles one type of disk pages
 - Index pages
 - Data file pages
 - Log file pages
 - Number of available buffers for each domain is predetermined
- Hot set method → PAGE REPLACEMENT ALSO
 - Useful in queries that scan a set of pages repeatedly
 - JOIN OPERATION IN NESTED LOOPS
 - DISK PAGES THAT WILL BE REFERENCES REPEATEDLY
 - Does not replace the set in the buffers until processing is completed
- The DBMIN method → PAGE REPLACEMENT POLICY
 - Predetermines the pattern of page references for each algorithm for a particular type of database operation
 - Calculates locality set using query locality set model (QLSM)

ACID PROPERTIES: Transactions Several Properties

→ enforced by concurrency control, recovery methods

- Atomicity
 - Transaction performed in its entirety or not at all
 - Consistency preservation
 - ↑ A transaction takes database from one consistent state to another, if there's no interference from any other transaction
 - Isolation
 - ↑ A transaction must not interfere with/by other transactions
 - Durability or permanency
 - Changes must persist in the database of commit transaction.
These changes must not be lost
- ...
...
- Levels of isolation
 - Level 0 isolation does not overwrite the dirty reads of higher-level transactions
 - Level 1 isolation has no lost updates
 - Level 2 isolation has no lost updates and no dirty reads
 - Level 3 (true) isolation has repeatable reads
 - In addition to level 2 properties
 - Snapshot isolation → another type of isolation

Schedule / History

- Schedule or history of n transactions
 - Order of operations of transactions
 - Operations from different transactions can be interleaved in the schedule
- Total ordering of operations in a schedule
 - For any two operations in the schedule, one must occur before the other

Conflicting operations in a schedule

- Two conflicting operations in a schedule are in conflict, if they satisfy these conditions
 - Operations belong to different transactions
 - Operations access the same item X
 - At least one of the operations is a write item(X)
- Two operations conflict if changing their order results in a different outcome
- Read-write conflict
- Write-write conflict

CHARACTERIZING SCHEDULES BASED ON RECOVERABILITY

Recoverable schedules → no abort

↳ a committed T should never be rolled back

↳ recoverable

Non Recoverable schedules → abort used

↳ a committed T is rolled back

↳ should not be permitted by DBMS

T ₁	T ₂
r(x)	r(x)
w(x)	w(x)
r(y)	commit
w(y)	
commit	

T ₁	T ₂
r(x)	r(x)
w(x)	r(x)
r(y)	w(x)
abort	commit
	so now value of x is no longer valid

T ₁	T ₂
r(x)	r(x)
w(x)	r(y)
r(y)	w(x)
abort	commit

CASCADING ROLLBACK

↳ uncommitted T, may need to be rolled back

↳ This is time consuming →

SOLUTION

CASCADELESS SCHEDULES

↳ every T only reads items written by committed T

↳ read needs to be committed

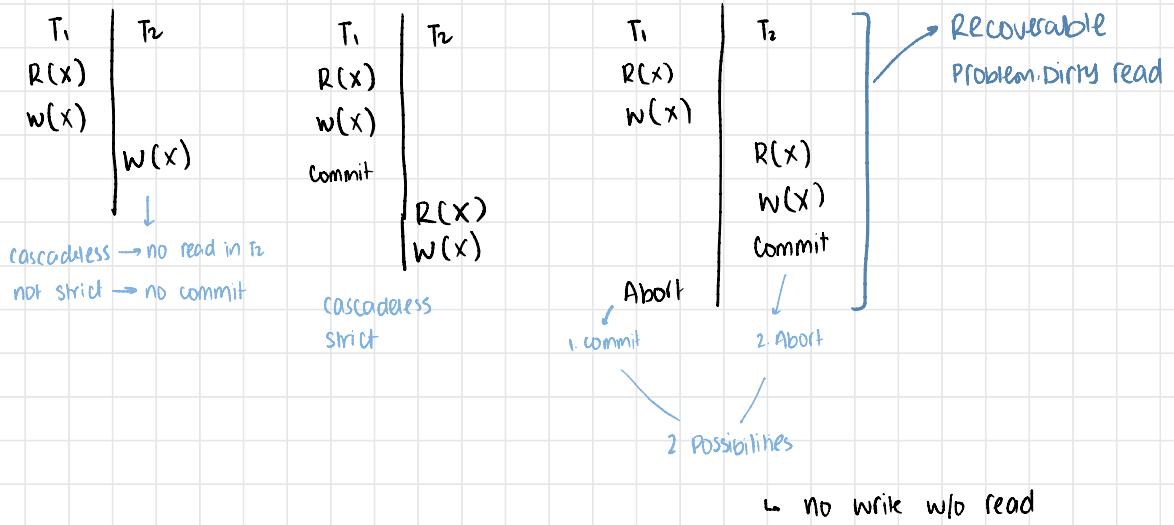
STRICK SCHEDULES

↳ every T only reads, writes items written by committed T

↳ read, write committed

↳ simpler recovery process → restore the before image

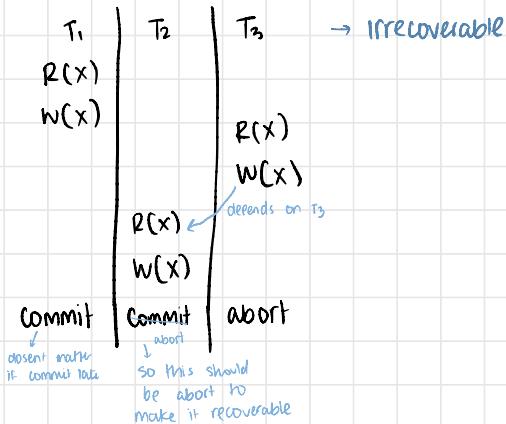
recoverable
cascadeless by default



cascading abort

↳ if T_1 aborts so do the rest of the transaction

↳ concurrency failure transaction



T ₁	T ₂	T ₃
R(x)		
R(z)	R(z)	
		R(x)
		R(y)
w(x) Commit		
		w(y) Commit
		R(y) w(z) w(y) Commit

cascadeless
strict
recoverable

T ₁	T ₂	T ₃
R(x)		
R(z)	R(z)	
		R(x)
		R(y)
	w(x)	
		R(y)
		w(z)
		w(y)
		commit
		commit
		commit

last transaction
and committed in no
end so its committed

is committed read
is committed write

cascadeless
strict
recoverable

T ₁	T ₂	T ₃
R(x)		
R(z)	R(z)	R(x)
		R(y)
R(z)	R(y)	
		R(y)
w(x) Commit		
		w(y)
		Commit
		w(z)
		w(y)
		Commit

cascadeless
strict
recoverable

T ₁	T ₂	T ₃
R(x)		
W(x)		
		R(x)
		W(x)
		abort
		R(x)
		W(x)
		Commit
		abort
		Commit

now its recoverable

irrecoverable

CHARACTERIZING SCHEDULES BASED ON SERIALIZABILITY

Serializable schedules

↳ always correct when concurrent T are executing

↳ places simultaneous T in series $\rightarrow T_1 T_1 \text{ before } T_2 / T_2 \text{ before } T_1$

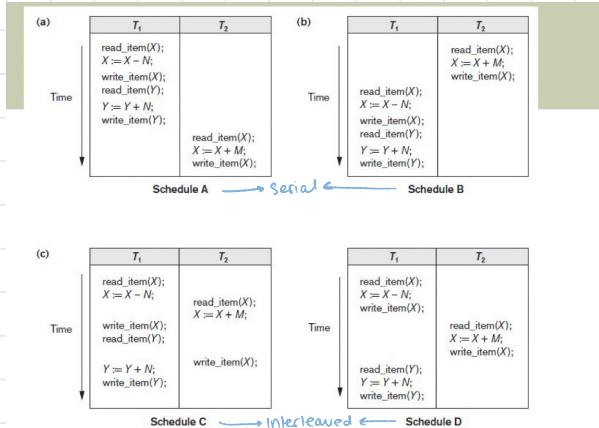


Figure 20.5 Examples of serial and nonserial schedules involving transactions T_1 and T_2 (a) Serial schedule A: T_1 followed by T_2 (b) Serial schedule B: T_2 followed by T_1 (c) Two nonserial schedules C and D with interleaving of operations

Copyright © 2016 Ramez Elmasri and Shamkant B. Navathe

Slide 20-3:

Serial Schedule

↳ T_1 complete, then T_2 start

↳ no concurrency

Interleaved Schedule

↳ T_1, T_2 concurrency

Problem with serial schedules

- Limit concurrency by prohibiting interleaving of operations
- Unacceptable in practice
- Solution: determine which schedules are equivalent to a serial schedule and allow those to occur

Serializability

↳ interleaved make same transaction changes as serial

↳ interleaved schedule is equivalent to serial schedule

↳ interleaved, serial both produce same final state

Conflicting Operations

- ↳ W(x), R(x) no R(x), R(x)
- ↳ W(x), W(x)
- ↳ R(x), W(x)

Conflict equivalence

- ↳ relative order of conflicting operations
is same in both schedules

CHECK Serializability

If conflict equivalence

↳ make nodes of all transactions

↳ check for conflicting

↳ $W_1(x), R_2(x)$ $T_1 \rightarrow T_2$

↳ $R_1(x), W_2(x)$ $T_1 \rightarrow T_2$

↳ $W_1(x), W_2(x)$ $T_1 \rightarrow T_2$

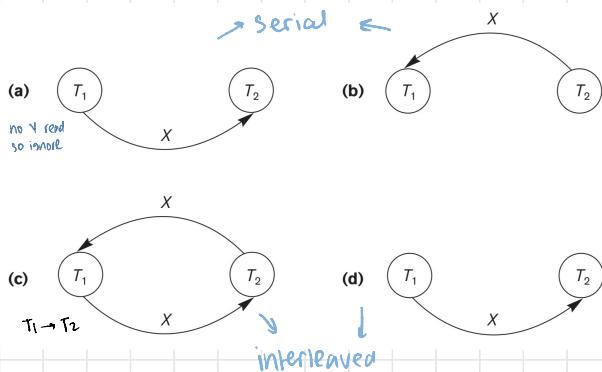
* min incoming edges runs first

cycle → not equivalent to serial

no cycle → equivalent to serial

(a)	T_1	T_2	(b)	T_1	T_2
Time	<pre>read_item(X); X := X - N; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>		Time	<pre>read_item(X); X := X + M; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>	
Schedule A	→ serial ←		Schedule B	→ serial ←	

(c)	T_1	T_2	(d)	T_1	T_2
Time	<pre>read_item(X); X := X - N; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>	<pre>read_item(X); X := X + M; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>	Time	<pre>read_item(X); X := X + M; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>	<pre>read_item(X); X := X + M; write_item(X); read_item(Y); Y := Y + N; write_item(Y);</pre>
Schedule C	→ interleaved ←		Schedule D	→ interleaved ←	



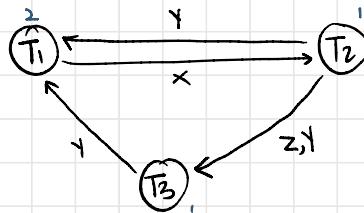
(e)	Transaction T_1	Transaction T_2	Transaction T_3
Time	<pre>read_item(X); write_item(X);</pre>		<pre>read_item(Y); read_item(Z);</pre>
	<pre>read_item(Y); write_item(Y);</pre>	<pre>read_item(Z);</pre>	<pre>write_item(Y); write_item(Z);</pre>

Schedule F

(b)

Transaction T_1	Transaction T_2	Transaction T_3
$\text{read_item}(X); \text{write_item}(X);$ $\text{read_item}(Y); \text{write_item}(Y);$	$\text{read_item}(Z); \text{read_item}(Y); \text{write_item}(Z);$ $\text{read_item}(X);$	$\text{read_item}(Y); \text{read_item}(Z);$ $\text{write_item}(Y); \text{write_item}(Z);$

Schedule E



$T_3 \rightarrow T_2 \rightarrow T_1$

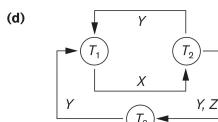
OR

$T_2 \rightarrow T_3 \rightarrow T_1$

(b)

Transaction T_1	Transaction T_2	Transaction T_3
$\text{read_item}(X); \text{write_item}(X);$ $\text{read_item}(Y); \text{write_item}(Y);$	$\text{read_item}(Z); \text{read_item}(Y); \text{write_item}(Y);$ $\text{read_item}(X);$	$\text{read_item}(Y); \text{read_item}(Z);$ $\text{write_item}(Y); \text{write_item}(Z);$

Schedule E



Equivalent serial schedules

None

Reason

Cycle $X(T_1 \rightarrow T_2), Y(T_2 \rightarrow T_1)$

Cycle $X(T_1 \rightarrow T_2), YZ(T_2 \rightarrow T_3), Y(T_3 \rightarrow T_1)$

HOW serialibility is used for concurrency control

Serializable schedule

↳ diff from being serial as

↳ it gives benefit of concurrent execution

without giving up any correctness

- Difficult to test for serializability in practice

- Factors such as system load, time of transaction submission, and process priority affect ordering of operations

- DBMS enforces protocols

- Set of rules to ensure serializability

View equivalence of 2 schedules

- Same set of T and operations of those T
- read operation value reads result of same write operation in both schedule
write operation produces same results in both schedule
- Last operation same in both schedules

View Resizability schedule

↳ view equivalent to a serial schedule

Conflict Resizability

↳ blind writes allowed

↳ by default view resizable

↳ write without a preceding read operation

Constrained write assumption

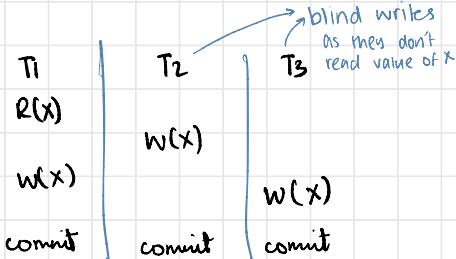
↳ read operation, then write operation only

↳ no blind writes

Unconstrained write assumption

↳ value written by an operation

can be independent of its old value



View Resizable

not conflict Resizable → as no conflict equivalence



OTHER TYPES OF EQUIVALENCE SCHEDULES

Debit card transactions

can use non serializable schedules *as order doesn't matter*

TRANSACTION SUPPORT IN SQL

- No explicit Begin_Transaction statement
- Every transaction must have an explicit end statement
 - COMMIT
 - ROLLBACK
- Access mode is READ ONLY or READ WRITE
- Diagnostic area size option
 - Integer value indicating number of conditions held simultaneously in the diagnostic area
- Isolation level option
 - Dirty read
 - Nonrepeatable read
 - Phantoms

		Type of Violation		
Isolation Level	Dirty Read	Nonrepeatable Read	Phantom	
READ UNCOMMITTED	Yes	Yes	Yes	
READ COMMITTED	No	Yes	Yes	
REPEATABLE READ	No	No	Yes	
SERIALIZABLE	No	No	No	

Table 20.1 Possible violations based on isolation levels as defined in SQL

- Snapshot isolation
 - Used in some commercial DBMSs
 - Transaction sees data items that it reads based on the committed values of the items in the database snapshot when transaction starts
 - Ensures phantom record problem will not occur

again, T_1 will see a different value.

3. **Phantoms.** A transaction T_1 may read a set of rows from a table, perhaps based on some condition specified in the SQL WHERE-clause. Now suppose that a transaction T_2 inserts a new row r that also satisfies the WHERE-clause condition used in T_1 , into the table used by T_1 . The record r is called a **phantom record** because it was not there when T_1 starts but is there when T_1 ends. T_1 may or may not see the phantom, a row that previously did not exist. If the equivalent serial order is T_1 followed by T_2 , then the record r should not be seen; but if it is T_2 followed by T_1 , then the phantom record should be in the result given to T_1 . If the system cannot ensure the correct behavior, then it does not deal with the phantom record problem.

CONCURRENCY CONTROL PROTOCOLS

- ↳ set of rules to guarantee serializability

Timestamp

uses

- ↳ unique identifier for each T

TWO PHASE LOCKING PROTOCOL

- ↳ locks data items to prevent multiple T from accessing items concurrently
- ↳ has high overhead

MULTIVERSION CONCURRENCY PROTOCOLS

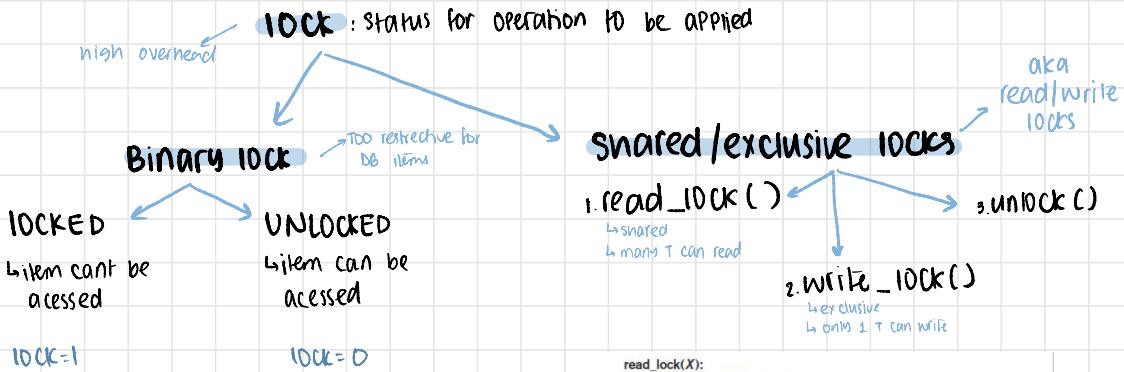
- ↳ uses multiple versions of a data item

LOWER OVERHEAD

OPTIMISTIC PROTOCOLS

- ↳ Validation of a T
- ↳ Certification of a T
- ↳ Then also assume multiple versions of a data item

TWO PHASED LOCKING TECHNIQUES FOR CONCURRENCY CONTROL



- Transaction requests access by issuing a `lock_item(X)` operation

```
lock_item(X):
B: if LOCK(X) = 0           (*item is unlocked*)
    then LOCK(X) ← 1        ("lock the item")
else
begin
    wait (until LOCK(X) = 0
        and the lock manager wakes up the transaction);
    go to B
end;
unlock_item(X):
    LOCK(X) ← 0;           (* unlock the item *)
    if any transactions are waiting
        then wakeup one of the waiting transactions;
```

Figure 21.1 Lock and unlock operations for binary locks

```
read_lock(X):
B: if LOCK(X) = "unlocked"
    then begin LOCK(X) ← "read-locked";
        no_of_reads(X) ← 1
    end
else if LOCK(X) = "read-locked"
    then no_of_reads(X) ← no_of_reads(X) + 1
else begin
    wait (until LOCK(X) = "unlocked"
        and the lock manager wakes up the transaction);
    go to B
end;

write_lock(X):
B: if LOCK(X) = "unlocked"
    then LOCK(X) ← "write-locked"
else begin
    wait (until LOCK(X) = "unlocked"
        and the lock manager wakes up the transaction);
    go to B
end;

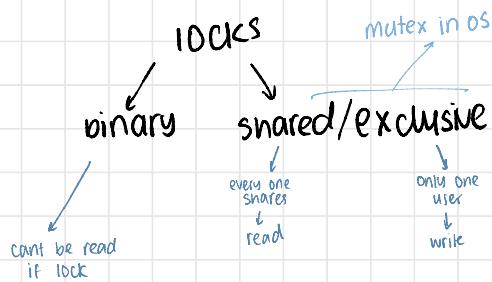
unlock(X):
if LOCK(X) = "write-locked"
    then begin LOCK(X) ← "unlocked";
        wakeup one of the waiting transactions, if any
    end
else if LOCK(X) = "read-locked"
    then begin
        no_of_reads(X) ← no_of_reads(X) - 1;
        if no_of_reads(X) = 0
            then begin LOCK(X) = "unlocked";
                wakeup one of the waiting transactions, if any
            end
    end;
```

LOCK TABLE

- Specifies items that have locks

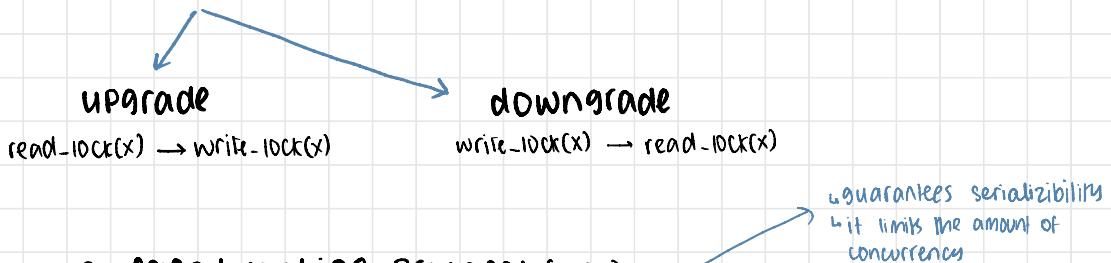
LOCK manager

- tracks and controls access to locks
- rules enforced by lock manager



lock conversion

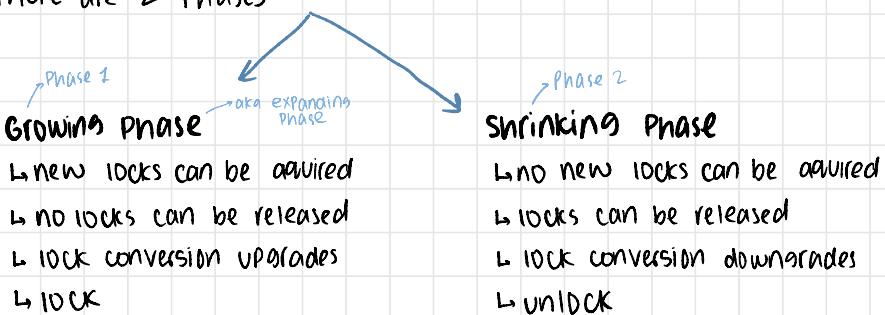
- ↳ a T already holding a lock can convert lock to another state



TWO FACED LOCKING PROTOCOL (TPL)

- ↳ once T releases its first lock, it can no longer acquire new locks

- ↳ There are 2 Phases



* Once gone shrink, can't go back to growing Phase

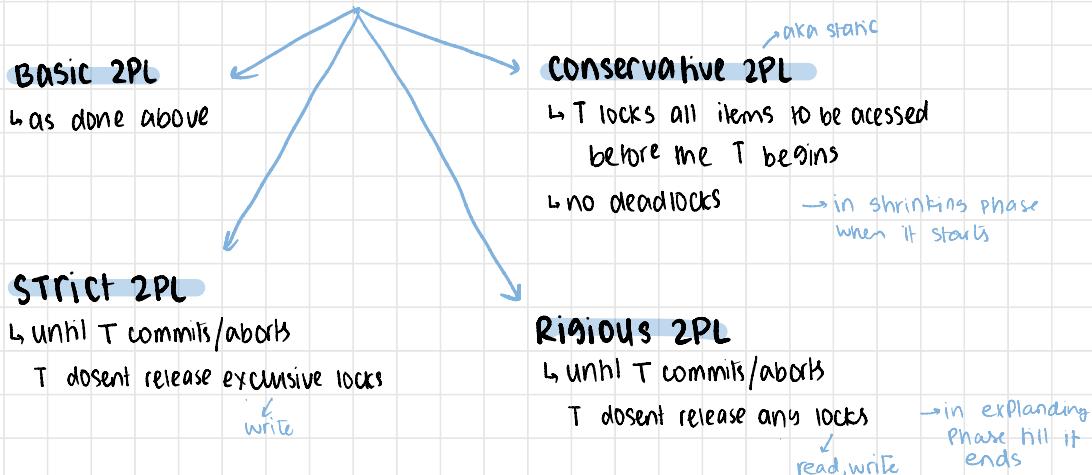
Figure 21.3 Transactions that do not obey two-phase locking (a) Two transactions T1 and T2 (b) Results of possible serial schedules of T1 and T2 (c) A nonserializable schedule S that uses locks

(a)	<table border="1"> <tr> <th>T_1</th> <th>T_2</th> </tr> <tr> <td> <code>read_lock(X); read_item(Y); unlock(X); write_lock(Y); read_item(Z); X := X + Y; write_item(X); unlock(X);</code> </td> <td> <code>read_lock(); read_item(); unlock(); write_lock(); write_lock(); read_item(); Y := X + Y; write_item(); unlock();</code> </td> </tr> </table>	T_1	T_2	<code>read_lock(X); read_item(Y); unlock(X); write_lock(Y); read_item(Z); X := X + Y; write_item(X); unlock(X);</code>	<code>read_lock(); read_item(); unlock(); write_lock(); write_lock(); read_item(); Y := X + Y; write_item(); unlock();</code>	(b)	<p>Initial values: $X=20, Y=30$</p> <p>Result serial schedule T_1, followed by T_2: $X=50, Y=80$</p> <p>Result of serial schedule T_2, followed by T_1: $X=70, Y=60$</p>				
T_1	T_2										
<code>read_lock(X); read_item(Y); unlock(X); write_lock(Y); read_item(Z); X := X + Y; write_item(X); unlock(X);</code>	<code>read_lock(); read_item(); unlock(); write_lock(); write_lock(); read_item(); Y := X + Y; write_item(); unlock();</code>										
(c)	<table border="1"> <tr> <th>T_1</th> <th>T_2</th> </tr> <tr> <td> <code>read_lock(Y); read_item(); unlock();</code> </td> <td> <code>read_lock(); read_item(); unlock(); write_lock(); write_lock(); Y := X + Y; write_item(); unlock();</code> </td> </tr> </table>	T_1	T_2	<code>read_lock(Y); read_item(); unlock();</code>	<code>read_lock(); read_item(); unlock(); write_lock(); write_lock(); Y := X + Y; write_item(); unlock();</code>	Time	<table border="1"> <tr> <th>T_1'</th> <th>T_2'</th> </tr> <tr> <td> <code>read_lock(Y); read_item(); write_lock(); unlock(); read_item(); X := X + Y; write_item(); unlock();</code> </td> <td> <code>read_lock(); read_item(); write_lock(); unlock(); read_item(); Y := X + Y; write_item(); unlock();</code> </td> </tr> </table>	T_1'	T_2'	<code>read_lock(Y); read_item(); write_lock(); unlock(); read_item(); X := X + Y; write_item(); unlock();</code>	<code>read_lock(); read_item(); write_lock(); unlock(); read_item(); Y := X + Y; write_item(); unlock();</code>
T_1	T_2										
<code>read_lock(Y); read_item(); unlock();</code>	<code>read_lock(); read_item(); unlock(); write_lock(); write_lock(); Y := X + Y; write_item(); unlock();</code>										
T_1'	T_2'										
<code>read_lock(Y); read_item(); write_lock(); unlock(); read_item(); X := X + Y; write_item(); unlock();</code>	<code>read_lock(); read_item(); write_lock(); unlock(); read_item(); Y := X + Y; write_item(); unlock();</code>										

Figure 21.4

Transactions T_1' and T_2' , which are the same as T_1 and T_2 in Figure 21.3 but follow the two-phase locking protocol. Note that they can produce a deadlock.

VARIATIONS OF TWO PHASED LOCKING



hence no other T can read/write an item that is written by T unless T is committed

CONCURRENCY CONTROL SYSTEM

- ↳ generates read-lock, write-lock requests

Deadlock

- ↳ set of locked T, each holding a data item and waiting to acquire a data item held by another T

- ↳ process stuck in circular waiting for the resources

* deadlock implies starvation ↳ starvation does not imply deadlock

Starvation

- ↳ when a T waits because it requires a data item to run, that is never allocated to this T

↳

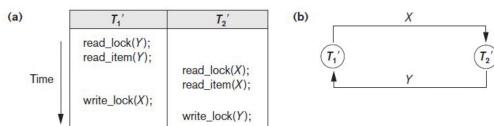


Figure 21.5 Illustrating the deadlock problem (a) A partial schedule of T_1' and T_2' that is in a state of deadlock (b) A wait-for graph for the partial schedule in (a)

Deadlock Prevention Protocols

1

- ↳ conservative 2PL
- ↳ lock all data items it needs in advance

2

- ↳ Order all items in DB
- ↳ T that needs several items will lock them in that order

both limit concurrency

Wait die

- ↳ If T_y wants to access item locked by T_o
 T_y is aborted and restarted with same priority
- ↳ If T_o wants to access item locked by T_y
 T_o will have to wait

Wound wait

- ↳ If T_y wants to access item locked by T_o
 T_y is has to wait
- ↳ If T_o wants to access item locked by T_y
 T_o will have access

avoids starvation

No Waiting algo (NW)

- ↳ If T can't obtain a lock
it is aborted and restarted later
without checking if deadlock
occurs or not

needlessly aborts, restarts T

Cautious waiting algo (CW)

- ↳ If T_1 is not blocked
then T_2 blocked and allowed to wait
- ↳ If T_1 is blocked
abort T_1

solution
waiting for a dead item

Deadlock Detection

- ↳ System checks if deadlock state exists
by making a wait-for graph

Victim selection: which T to abort due to deadlock

Timeouts: if a T waits longer than a predefined time, abort the T

Starvation: when T can't proceed for an indefinite time, while other continue normally

↳ solution: first come first serve queue

Concurrency Control Based on Timestamp Ordering (cont'd.)

- **Timestamp**
 - Unique identifier assigned by the DBMS to identify a transaction
 - Assigned in the order submitted
 - Transaction start time
- Concurrency control techniques based on timestamps do not use locks
 - Deadlocks cannot occur

- **Generating timestamps**
 - Counter incremented each time its value is assigned to a transaction
 - Current date/time value of the system clock
 - Ensure no two timestamps are generated during the same tick of the clock
- **General approach**
 - Enforce equivalent serial order on the transactions based on their timestamps

- **Timestamp ordering (TO)**
 - Allows interleaving of transaction operations
 - Must ensure timestamp order is followed for each pair of conflicting operations
- **Each database item assigned two timestamp values**
 - `read_TS(X)`
 - `write_TS(X)`

- **Basic TO algorithm**
 - If conflicting operations detected, later operation rejected by aborting transaction that issued it
 - Schedules produced guaranteed to be conflict serializable
 - Starvation may occur
- **Strict TO algorithm**
 - Ensures schedules are both strict and conflict serializable

- **Thomas's write rule**
 - Modification of basic TO algorithm
 - Does not enforce conflict serializability
 - Rejects fewer write operations by modifying checks for `write_item(X)` operation

- Recovery process restores database to most recent consistent state before time of failure
- Information kept in system log

TYPICAL RECOVERY STRATEGIES

- Restore backed up copy of DB → when extensive damage
- Identify any changes that may cause inconsistency → when noncatastrophic failure
 - Some operations may require redo

3. DEFERRED UPDATE TECHNIQUE

- When T commits
- then physically update DB
- No undo needed
- redo might be needed from log
- No undo/redo also
no steal

- Can only be used for short transactions and transactions that change few items
 - Buffer space an issue with longer transactions
- Transaction does not reach its commit point until all its REDO-type log entries are recorded in log and log buffer is force-written to disk

4. IMMEDIATE UPDATE TECHNIQUE

- DB might be updated by some operations before T commits
- Operations recorded in logs
- undo/redo also
- steal/no force strategy

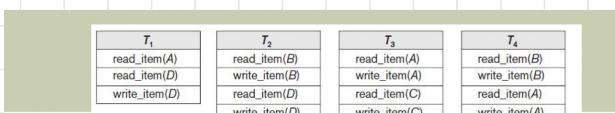


Figure 22.3 An example of recovery using deferred update with concurrent transactions (a) The READ and WRITE operations of four transactions (b) System log at the point of crash

[start_transaction, T_1]
[write_item, T_1 , D, 20]
[commit, T_1]
[checkpoint]
[start_transaction, T_2]
[write_item, T_2 , B, 15]
[write_item, T_4 , A, 20]
[commit, T_4]
[start_transaction, T_3]
[write_item, T_3 , B, 12]
[start_transaction, T_3]
[write_item, T_3 , A, 30]
[write_item, T_3 , D, 25]

T_2 and T_3 are ignored because they did not reach their commit points.
 T_4 is redone because its commit point is after the last system checkpoint.

UNDO AND READ OPERATIONS

↳ idempotent → executing operation = executing just once
multiple times

↳ entire recovery process should be idempotent

CACHING OF DISK BLOCKS

DBMS Cache: a collection of in-memory buffers

cache directory: tracks which DB items are in buffers

flush: cache buffers flushed to make space for new items

if dirty bit = 1 → modified buffer

↳ contents written back to disk before flush

If pin-unpin bit = 1 → page can't be written back to disk yet

Main Strategies when pushing a modified buffer back to disk

dirty bit = 1

Inplace updating

- ↳ writes in same db disk location
- ↳ overwrites the changed values

log for recovery
inorder to
undo/redo

- Before-image: old value of data item → BFIM
- After-image: new value of data item → AFIM

Write ahead logging

- ↳ ensure before image is recorded in log and
- ↳ log entry is flushed to disk before AFIM overwrites
- UNDO type log entries → old value → BFIM
- redo type log entries → new value → AFIM

Shadowing

- ↳ writes in diff disc location
- to maintain multiple versions of data items

Steal/no steal, Force/no force

- ↳ rules for when a page from DB cache can be written to disk

NO Steal approach

- ↳ T commits, only then
- ↳ write updated buffer to disk

Steal approach

- ↳ write updated buffer to disk
- ↳ men T commits

NO Force approach

- ↳ T commits, only then
- ↳ write updated buffer to disk

Force approach

- ↳ write all updated buffer to disk immediately
- ↳ men T commits
- ↳ eliminates need for redo

Typical database systems employ a steal/no-force strategy

- Avoids need for very large buffer space → steal
- Reduces disk I/O operations for heavily updated pages → no force

Write-ahead logging protocol for recovery algorithm requiring both UNDO and REDO

- BFIM of an item cannot be overwritten by its after image until all UNDO-type log entries have been force-written to disk
- Commit operation of a transaction cannot be completed until all REDO-type and UNDO-type log records for that transaction have been force-written to disk

check Point

↳ another type of entry in log

- Taking a checkpoint
 - Suspend execution of all transactions temporarily
 - Force-write all main memory buffers that have been modified to disk
 - Write a checkpoint record to the log, and force-write the log to the disk
 - Resume executing transactions
- DBMS recovery manager decides on checkpoint interval

time needed to
force write may
delay T

solution

Fuzzy checkpointing

- System can resume transaction processing after a begin_checkpoint record is written to the log
- Previous checkpoint record maintained until end_checkpoint record is written

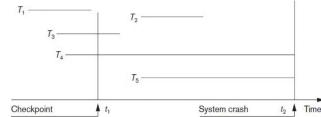


Figure 22.2 An example of a recovery timeline to illustrate the effect of checkpointing

Presented by Ramez Elmasri and Shamkant B. Navati

Slide 22-0

Transaction Rollback

- Transaction failure after update but before commit
 - Necessary to roll back the transaction
 - Old data values restored using undo-type log entries
- Cascading rollback
 - If transaction T is rolled back, any transaction S that has read value of item written by T must also be rolled back
 - Almost all recovery mechanisms designed to avoid this

Figure 22.1 illustrating cascading rollback (a process that never occurs in strict or cascadeless schedules) (a) The read and write operations of three transactions (b) System log at point of crash (c) Operations before the crash

(a)

T_1	T_2	T_3
read_item(A) read_item(C) write_item(C)	read_item(B) write_item(B)	read_item(B) read_item(A) write_item(A)

(b)

	A	B	C	D
(start transaction, T_1)	90	15	40	20
(read_item, T_1 , A)	12			
(write_item, T_1 , C, 12, 12)		12		
(read_item, T_2 , B)		18		
(write_item, T_2 , B, 12, 18)		18		
(read_item, T_3 , A)			25	
(read_item, T_3 , C)			25	
(write_item, T_3 , D, 25, 25)			25	
(read_item, T_3 , C)			25	
(write_item, T_3 , D, 25, 24)			25	
(write_item, T_3 , C, 25, 24)			25	

(c)

Copyright © 2016 Ramez Elmasri and Shamkant B. Navati

Slide 22-1

Transactions that Do Not Affect the Database

- Example actions: generating and printing messages and reports
- If transaction fails before completion, may not want user to get these reports
 - Reports should be generated only after transaction reaches commit point
- Commands that generate reports issued as batch jobs executed only after transaction reaches commit point
 - Batch jobs canceled if transaction fails

ishma hafeez
notes
refresh tree