

Transaction Models and Concurrency Control

5DV120 — Database System Principles

Umeå University

Department of Computing Science

Stephen J. Hegner

`hegner@cs.umu.se`

`http://www.cs.umu.se/~hegner`

The Issue of Concurrency in the DBMS Context

- It is often the case that a database system will be accessed by many users simultaneously.
- If this access is read-only, then there are no serious integrity problems; only ones of performance.
- If the access includes writing the database, then serious problems will arise if the interaction is not regulated.
- It is therefore necessary to characterize correct behavior in the context of concurrent transactions.

The ACID Characterization

- The properties which a set of concurrent transactions should exhibit is often expressed via the acronym *ACID*:

Atomicity: For each transaction, either the complete result of its execution is recorded in the database, or else nothing about its results is recorded.

Consistency: The execution of any transaction in isolation preserves the integrity of the database.

Isolation: The execution of one running transaction must not affect the execution of another concurrently running transaction.

Durability: The results of the transactions are permanent in the database.

- These slides will focus primarily upon *isolation*.
- A subsequent set of slides will focus upon *atomicity* and *durability*.
- *Consistency* is a property of a single transaction and will not be the focus here.

Example Transactions

Example (simplified bank transactions) : Two transactions T_1 and T_2 .

- R_i and W_i are local variables for transaction i with $i \in \{1, 2\}$.
- There are the following operations:
 - $R_Bal_i\langle a \rangle$ means that transaction T_i reads the balance of account a into a local variable R_i : $R_i \leftarrow Bal\langle a \rangle$.
 - $W_Bal_i\langle a \rangle$ means that transaction T_i writes the balance of account a from variable W_i to the database: $Bal\langle a \rangle \leftarrow W_i$.
 - $Cpd_Bal_i\langle X \rangle$ is a local operation that adds $X\%$ interest to R_i and places the result in W_i : $W_i \leftarrow R_i \times (1 + X/100)$.
 - $Wthd_i\langle X \rangle$ means that X Euros are subtracted from the local value R_i and placed in W_i : $W_i \leftarrow R_i - X$.

T_1 Compound 10% on account 15:

$R_Bal_1\langle 15 \rangle; Cpd_Bal_1\langle 10 \rangle; W_Bal_1\langle 15 \rangle$.

T_2 Withdraw 2000 from account 15:

$R_Bal_2\langle 15 \rangle; Wthd_2\langle 2000 \rangle; W_Bal_2\langle 15 \rangle$.

Order of Execution

- Shown below are two possibilities for schedules for these transactions.

T_1	T_2	Bal $\langle 15 \rangle$	T_1	T_2	Bal $\langle 15 \rangle$
R_Bal ₁ $\langle 15 \rangle$		10000		R_Bal ₂ $\langle 15 \rangle$	10000
Cpd_Bal ₁ $\langle 10 \rangle$		10000		Wthd ₂ $\langle 2000 \rangle$	10000
W_Bal ₁ $\langle 15 \rangle$		11000		W_Bal ₂ $\langle 15 \rangle$	8000
	R_Bal ₂ $\langle 15 \rangle$	11000	R_Bal ₁ $\langle 15 \rangle$		8000
	Wthd ₂ $\langle 2000 \rangle$	11000	Cpd_Bal ₁ $\langle 10 \rangle$		8000
	W_Bal ₂ $\langle 15 \rangle$	9000	W_Bal ₁ $\langle 15 \rangle$		8800

- Both schedules are *serial* and both are correct ...
 - ... even though the results differ.
- The order of *serial* execution does not affect correctness.
- The system cannot and should not decide which order is better.

Lost Updates

- If the steps of the transactions are interleaved in certain ways, updates may be lost. Shown below are two possibilities for schedules for these transactions.

T_1	T_2	Bal $\langle 15 \rangle$	T_1	T_2	Bal $\langle 15 \rangle$
R_Bal ₁ $\langle 15 \rangle$		10000		R_Bal ₂ $\langle 15 \rangle$	10000
Cpd_Bal ₁ $\langle 10 \rangle$		10000		Wthd ₂ $\langle 2000 \rangle$	10000
	R_Bal ₂ $\langle 15 \rangle$	10000	R_Bal ₁ $\langle 15 \rangle$		10000
	Wthd ₂ $\langle 2000 \rangle$	10000	Cpd_Bal ₁ $\langle 10 \rangle$		10000
	W_Bal ₂ $\langle 15 \rangle$	8000	W_Bal ₁ $\langle 15 \rangle$		11000
W_Bal ₁ $\langle 15 \rangle$		11000		W_Bal ₂ $\langle 15 \rangle$	8000

- In the schedule on the left, the result of T_2 is lost.
- In the schedule on the right, the result of T_1 is lost.

Basic Steps and Transactions

- To study the issues surrounding concurrency systematically, some formal notions are necessary.

Basic steps: A *basic step* for a transaction T is either a read $r\langle x \rangle$ or a write $w\langle x \rangle$ of a data object x .

- The actual values of x which are read and written are not important to the model.
- The internal steps (e.g., $R_Bal_i\langle x \rangle$, $R_Bal_i\langle x \rangle$, $Wthd_i\langle n \rangle$, $Cpd_Bal_i\langle n \rangle$) are not represented.
- Only the fact that T read or wrote that object is important.
- For T_i , these are usually written $r_i\langle x \rangle$ and $w_i\langle x \rangle$, respectively.

Transaction: A *transaction* $T = \langle t_1, t_2, \dots, t_n \rangle$ is modelled by a finite sequence of steps, with each t_i a basic step for T .

Example: $T_1 = r_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle w_1\langle z \rangle$ is a transaction.

- $Steps\langle T \rangle$ denotes the set of basic steps of T .

Example: $Steps\langle T_1 \rangle = \{r_1\langle x \rangle, r_1\langle y \rangle, w_1\langle y \rangle, w_1\langle z \rangle\}$.

Simplifying Assumptions for Transactions

- There are some simplifying assumptions which are made for the basic steps of a transaction.

Single read and write: For a given data object x , a transaction T reads x at most once and writes x at most once.

Read before write: If a transaction T both reads and writes a data object x , then it reads x before it writes x .

- $r\langle x \rangle \dots w\langle x \rangle$, not $w\langle x \rangle \dots r\langle x \rangle$.
- From the point of view of transaction semantics, these are not significant limitations, and they simplify the modelling of concurrency greatly.

Reason: In most cases, other transactions will not be able to see what T does until it finishes (commits), so the internal order of operations within T is not of significance to other transactions.

Schedules

- A *schedule* for a set of transactions is a specification of the order in which the basic steps will be executed.
- Formally, let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions, with

$$T_i = \langle t_{i1}, t_{i2}, \dots, t_{in_i} \rangle$$

for $1 \leq i \leq m$.

The steps of a schedule: Define $\text{Steps}\langle\mathbf{T}\rangle = \bigcup_{i=1}^m \text{Steps}\langle T_i \rangle$.

Schedule: A *schedule* S for \mathbf{T} is any total ordering \leq_S of the set $\text{Steps}\langle\mathbf{T}\rangle$ with the property that $t_{ij} \leq_S t_{ik}$ whenever $j \leq_S k$.

- In other words, the order of elements within each T_i is preserved.

Serial Schedules

Serial schedules: A schedule S for the set $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ of transactions is *serial* if there is a total ordering \leq of \mathbf{T} with the property that if $T_i < T_j$, then all elements of T_i occur before any element of T_j in the ordering \leq_s .

Examples: Let

$$T_1 = r_1\langle x \rangle r_1\langle y \rangle w_1\langle x \rangle w_1\langle y \rangle$$

$$T_2 = r_2\langle z \rangle w_2\langle z \rangle w_2\langle y \rangle$$

$$T_3 = r_3\langle z \rangle w_3\langle z \rangle r_3\langle x \rangle w_3\langle x \rangle$$

- Then

$$r_2\langle z \rangle w_2\langle z \rangle w_2\langle y \rangle \quad r_1\langle x \rangle r_1\langle y \rangle w_1\langle x \rangle w_1\langle y \rangle \quad r_3\langle z \rangle w_3\langle z \rangle r_3\langle x \rangle w_3\langle x \rangle$$

is the schedule corresponding to $T_2 < T_1 < T_3$, while

$$r_1\langle x \rangle r_1\langle y \rangle \quad r_3\langle z \rangle w_3\langle z \rangle \quad r_2\langle z \rangle \quad w_1\langle x \rangle w_1\langle y \rangle \quad w_2\langle z \rangle w_2\langle y \rangle \quad r_3\langle x \rangle w_3\langle x \rangle$$

is not a serial schedule.

Serializability

- A serial schedule exhibits a correct semantics of concurrency, as there is no undesirable intertwining of actions of different transactions.
- Allowing only serial schedules is too restrictive.
 - It prohibits any form of concurrency whatever.
 - Performance would be compromised greatly in many situations.
- The solution is to allow *serializable* schedules – ones which are equivalent to serial schedules.
 - Parallelism is allowed.
 - The correctness of transactions is not compromised.

Question: How is *serializability* defined?

- It turns out that there are (at least) three reasonable definitions.

Three Notions of Serializability

View serializability: In view serializability, it is ensured that the reads and subsequent writes of each data object occur in the same order as in some serial schedule.

- This is the most important theoretical notion of serializability.
- It is the “correct” theoretical notion of serializability.
- Testing a schedule for view serializability is NP-complete.

Final-state serializability: In final-state serializability, it is ensured that the final result (*i.e.*, the final values of the data objects) is the same as in some serial schedule.

- This form of serializability is strictly weaker than view serializability and not widely used.
- It will not be considered further in this course.

Conflict serializability: In conflict serializability, specific forms of conflict are ruled out.

- Conflict serializability is strictly stronger than view serializability.
- It is of interest because there exist efficient algorithms for testing conflict serializability.

The Three Conditions Surrounding View Equivalence

- Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions, and let S be a schedule for \mathbf{T} .
- Let $r_i\langle x \rangle \in \text{Steps}\langle T_i \rangle$ and $w_j\langle x \rangle \in \text{Steps}\langle T_j \rangle$.

Read from: $r_i\langle x \rangle$ *reads from* $w_j\langle x \rangle$ in S if $w_j\langle x \rangle \leq_s r_i\langle x \rangle$ and there is no $k \neq j$ for which $w_j\langle x \rangle \leq_s w_k\langle x \rangle \leq_s r_i\langle x \rangle$.

Initial read: $r_i\langle x \rangle$ is an *initial read* in S if there is no k for which $w_k\langle x \rangle \leq_s r_i\langle x \rangle$.

Final write: $w_j\langle x \rangle$ is a *final write* in S if there is no $k \neq j$ for which $w_j\langle x \rangle \leq_s w_k\langle x \rangle$.

Example: In

$r_1\langle x \rangle r_1\langle y \rangle$	$r_3\langle z \rangle w_3\langle z \rangle$	$r_2\langle z \rangle$	$w_1\langle x \rangle w_1\langle y \rangle$	$w_2\langle z \rangle w_2\langle y \rangle$	$r_3\langle x \rangle w_3\langle x \rangle$
---	---	------------------------	---	---	---

- $r_2\langle z \rangle$ reads from $w_3\langle z \rangle$.
- $r_3\langle x \rangle$ reads from $w_1\langle x \rangle$.
- $r_1\langle x \rangle$, $r_1\langle y \rangle$, and $r_3\langle z \rangle$ are initial reads.
- $w_2\langle z \rangle$, $w_2\langle y \rangle$, and $w_3\langle x \rangle$ are final writes.

View Equivalence and View Serializability

- Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions, and let S and S' be schedules for \mathbf{T} .

View equivalence: S and S' are *view equivalent*, written $S \approx_V S'$, if:

- (ve-i) Every read action of the form $r_i\langle x \rangle$ is, in each schedule, either an initial read or else reads from the same write action $w_j\langle x \rangle$.
- (ve-ii) The two schedules have the same final-write steps.

View serializability: S is said to be *view serializable* if there is a serial schedule S' such that $S \approx_V S'$.

Examples: The following schedules are view equivalent and view serializable.

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle w_3\langle z \rangle$ $r_2\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $r_3\langle x \rangle w_3\langle x \rangle$

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_3\langle z \rangle$ $r_2\langle z \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $r_3\langle x \rangle w_3\langle x \rangle$

$r_1\langle x \rangle r_1\langle y \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $r_3\langle z \rangle$ $w_3\langle z \rangle r_3\langle x \rangle w_3\langle x \rangle$ $r_2\langle z \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$

- The following schedule is not view serializable.

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle w_3\langle z \rangle r_3\langle x \rangle$ $r_2\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $w_3\langle x \rangle$

Motivation for Conditions of View Equivalence

- The motivation for condition (ve-i) is clear.

Question: Is (ve-ii) (equivalence of final writes) really necessary?

Example to motivate (ve-ii): Let

$$T_1 = w_1\langle x \rangle w_1\langle y \rangle \quad T_2 = w_2\langle x \rangle w_2\langle y \rangle$$

- Note that:
 - In any serial schedule, the first transaction has no effect.
 - Since there are no reads, no schedule can violate (ve-i).

Example: The following schedule is not equivalent to a serial schedule, yet satisfies (ve-i):

$$w_1\langle x \rangle w_2\langle x \rangle w_2\langle y \rangle w_1\langle y \rangle$$

- In general, both (ve-i) and (ve-ii) are necessary to obtain a satisfactory notion of serializability.

Blind Writes

Examples Each of the two schedules contains write operations which do not first read the associated data object.

$$T_1 = w_1\langle x \rangle w_1\langle y \rangle \quad T_2 = w_2\langle x \rangle w_2\langle y \rangle$$

Blind write: Let $T = \langle t_1, t_2, \dots, t_n \rangle$ be a transaction. The operation $t_j = w\langle x \rangle$ is called a *blind write* (of x) if for no $i < j$ is it the case that $t_i = r\langle x \rangle$.

- Without blind writes, condition (ve-ii) would not be necessary.

Theorem: In general, the problem of deciding whether a given schedule is view equivalent to some serial schedule is NP-complete. \square

But...: There is a polynomial-time algorithm to decide view serializability for the special case that none of the transactions involves blind writes.

Direct Serialization Conditions

- Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions, let S be a schedule for \mathbf{T} , and let T_i and T_j be distinct transactions in \mathbf{T} .

rw-conflict: There is a *direct read-write conflict* (or *rw-conflict*) on x from T_i to T_j , written $T_i \xrightarrow{rw\langle x \rangle} T_j$, if $i \neq j$, $r_i\langle x \rangle \leq_s w_j\langle x \rangle$ and $r_i\langle x \rangle \leq_s w_k\langle x \rangle \leq_s w_j\langle x \rangle$ implies $k = j$.

ww-conflict: There is a *direct write-write conflict* (or *ww-conflict*) on x from T_i to T_j , written $T_i \xrightarrow{ww\langle x \rangle} T_j$, if $i \neq j$, $w_i\langle x \rangle \leq_s w_j\langle x \rangle$ and $w_i\langle x \rangle \leq_s w_k\langle x \rangle \leq_s w_j\langle x \rangle$ implies $k = i$ or $k = j$.

wr-conflict: There is a *direct write-read conflict* (or *wr-conflict*) on x from T_i to T_j , written $T_i \xrightarrow{wr\langle x \rangle} T_j$, if $i \neq j$, $w_i\langle x \rangle \leq_s r_j\langle x \rangle$ and $w_i\langle x \rangle \leq_s w_k\langle x \rangle \leq_s r_j\langle x \rangle$ implies $k = i$.

Direct Serialization Graph: The *direct-serialization graph*, or *DSG*, for S is the directed graph whose vertices are the elements of \mathbf{T} and whose edges are exactly those defined by and labelled with the three forms of conflict above.

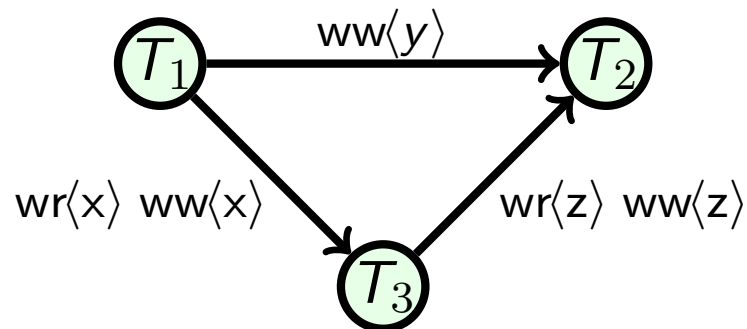
Examples of the DSG

Examples: The DSG for all of these schedules

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle w_3\langle z \rangle$ $r_2\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $r_3\langle x \rangle w_3\langle x \rangle$

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_3\langle z \rangle$ $r_2\langle z \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $r_3\langle x \rangle w_3\langle x \rangle$

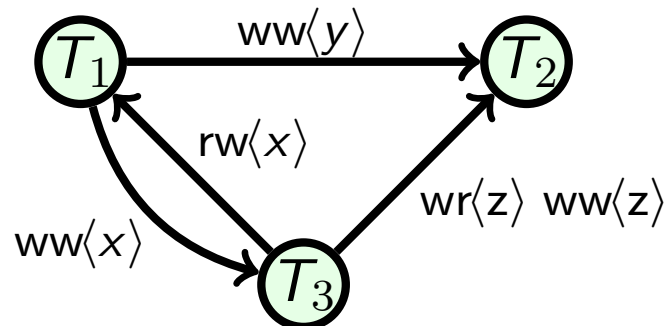
$r_1\langle x \rangle r_1\langle y \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $r_3\langle z \rangle$ $w_3\langle z \rangle r_3\langle x \rangle w_3\langle x \rangle$ $r_2\langle z \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$



The edges are labelled with the data names which induce them, as well as the *conflict types* (rw, wr, and ww).

Example: The DSG for

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle w_3\langle z \rangle r_3\langle x \rangle$ $r_2\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $w_3\langle x \rangle$



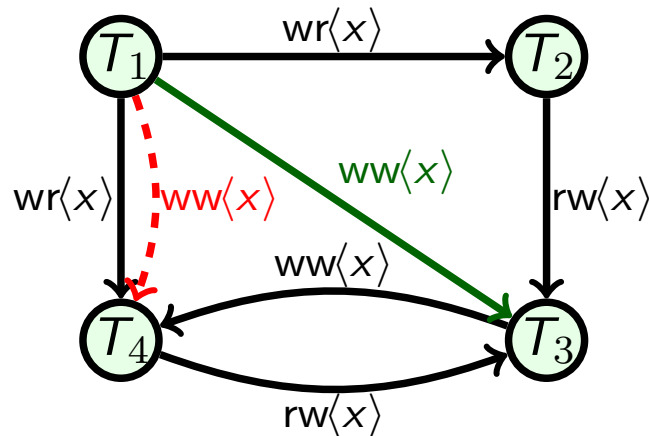
Next and Previous Operations and the DSG

- In the DSG, for any edge of the form $T_i \xrightarrow{\alpha w\langle x \rangle} T_j$, for $\alpha \in \{r, w\}$, T_j must perform the *next* write operation on x after T_i performs the α -operation.
- Similarly, for any $T_i \xrightarrow{wr\langle x \rangle} T_j$, T_i must perform the *last* write operation on x before T_j performs the read operation.

Example:



DSG:



- $T_1 \xrightarrow{ww\langle x \rangle} T_4$ is not a true ww-edge, since T_4 is not the *next* writer of x , after T_1 .
- However, $T_1 \xrightarrow{ww\langle x \rangle} T_3$ is a true ww-edge, since T_3 is the next writer of x , even though there are two reads sandwiched between these writes.

Conflict Serializability

- Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions.

Conflict equivalence: Two schedules S and S' for \mathbf{T} are *conflict equivalent*, denoted $S \approx_C S'$, if they have the same DSG.

Conflict serializability: The schedule S for \mathbf{T} is *conflict serializable* if there is a serial schedule S' for \mathbf{T} with $S \approx_C S'$.

Theorem: A schedule S for \mathbf{T} is conflict serializable iff the associated DSG is acyclic. \square

Theorem: Every conflict-serializable schedule is also view serializable. \square

- The advantage of conflict serializability (over view serializability) is that that there are simple and highly tractable algorithms for construction of the DSG, and testing it for cycles,

The Relationship Between View and Conflict Equivalence

Recall: Every conflict-serializable schedule is also view-serializable. \square

Restricted model: In the *restricted model* of transactions, there are no blind writes.

Theorem: In the restricted model, a schedule is view serializable iff it is conflict serializable. \square

Corollary: It is blind writes which force the decision problem for view serializability to be NP-complete. \square

Example: For $i \in \{1, 2, 3\}$, let $T_i = w_i\langle x \rangle w_i\langle y \rangle$.
Then

$w_1\langle x \rangle \quad w_2\langle x \rangle w_2\langle y \rangle \quad w_1\langle y \rangle \quad w_3\langle x \rangle w_3\langle y \rangle$

is view serializable and with $T_1 < T_2 < T_3$, but it is not conflict serializable since $T_1 \xrightarrow{ww\langle x \rangle} T_2 \xrightarrow{ww\langle y \rangle} T_1$ occurs in the DSG.

Realizing Serializable Schedules

- It is not reasonable to generate candidate schedules and then test for serializability.
- Rather, what is needed is a systematic way of guaranteeing that constructed schedules are serializable.
- There are several approaches in practice:
 - Locking:** Locks are used to prevent more than one transaction from writing the same data object concurrently, and also to prevent reads of objects which are being written.
 - Pure optimism:** Nothing is locked; conflicts are detected when transactions commit, and conflicts are resolved by aborting one or more transactions.
 - Multiversioning:** Each write operation generates a new version of the data object which is written. The versions are consolidated when the transactions finish.
- Many “real” approaches combine aspects of all three.

Locks

- In a lock-based approach, for a transaction to access a data object, it must request and be granted a *lock* on that object.
- There are two basic forms of lock:

Write lock: A write lock permits a transaction both to read and to write a data object.

- Only one transaction may hold a write lock on a data object at any given point in time.
- Also called an *exclusive lock* or *X-lock*.

Read lock: A read lock permits a transaction to read a data object, but not to write it.

- Several transactions may hold read locks on a data object concurrently.
- Also called a *shared lock* or *S-lock*.

Lock Requests and Releases:

- The following three basic lock operations are defined for a data object x by transaction T_i .

$rlk_i\langle x \rangle$: Request a read lock on x . This request may be granted provided there are no current write locks on x .

$wlk_i\langle x \rangle$: Request a write lock on x . This request may be granted provided there are no locks on x .

$unlk_i\langle x \rangle$: Dissolve the lock on x held by T_i .

- There are also two operations which upgrade and downgrade locks:

$upgr_i\langle x \rangle$: Convert a read lock by T_i on x to a write lock. This request may only be granted in the case that no other transaction holds a read lock on x .

$dngr_i\langle x \rangle$: Convert a write lock by T_i on x to a read lock.

- For various reasons, upgrades and downgrades are sometimes excluded from a modelling situation.

Transactions with Locks

- Informally, a *transaction with locks* is a transaction with lock commands interspersed.

Transaction with locks A *transaction with locks* is a sequence T_i of elements of the form $r_i\langle x \rangle$, $w_i\langle x \rangle$, $rlk_i\langle x \rangle$, $wlk_i\langle x \rangle$, $unlk_i\langle x \rangle$, $upgr_i\langle x \rangle$, and $dngr_i\langle x \rangle$, where x may be any data object and need not be the same for each element in the sequence, such that:

- If the operations of the form $rlk_i\langle x \rangle$, $wlk_i\langle x \rangle$, $unlk_i\langle x \rangle$, $upgr_i\langle x \rangle$, and $dngr_i\langle x \rangle$ are removed, the result is an ordinary transaction.
- The sequence must obey the *locking protocol* given on the next slide.

Locking Requirements

- A transaction with locks T_i must obey the following locking rules:
 - Before a data object x is read by T_i , a lock (read or write) must be requested and granted.
 - Before a data object x is written by T_i , a write lock must be requested and granted.
 - All reads on x must be performed before the corresponding lock on x is released.
 - All writes on x must be performed before the corresponding lock on x is released or downgraded.
 - All locks must be released (via unlock) before the transaction finishes (commits).
 - It is usually (but not always) assumed that transactions do not request redundant locks.
 - This makes analyses simpler.

Locking protocol: A scheduler operates according to a *locking protocol* just in case these conventions are followed.

Examples of Transactions with Locks

- Consider the transaction $T_1 = r_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle w_1\langle z \rangle$.
- The following are schedules with locks for T_1 .

$wlk_1\langle x \rangle wlk_1\langle y \rangle wlk_1\langle z \rangle r_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle w_1\langle z \rangle unl_1\langle x \rangle unl_1\langle y \rangle unl_1\langle z \rangle$

$rlk_1\langle x \rangle wlk_1\langle y \rangle wlk_1\langle z \rangle r_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle w_1\langle z \rangle unl_1\langle x \rangle unl_1\langle y \rangle unl_1\langle z \rangle$

$rlk_1\langle x \rangle r_1\langle x \rangle wlk_1\langle y \rangle r_1\langle y \rangle w_1\langle y \rangle wlk_1\langle z \rangle w_1\langle z \rangle unl_1\langle x \rangle unl_1\langle y \rangle unl_1\langle z \rangle$

$rlk_1\langle x \rangle r_1\langle x \rangle unl_1\langle x \rangle wlk_1\langle y \rangle r_1\langle y \rangle w_1\langle y \rangle unl_1\langle y \rangle wlk_1\langle z \rangle w_1\langle z \rangle unl_1\langle z \rangle$

$rlk_1\langle x \rangle r_1\langle x \rangle unl_1\langle x \rangle rlk_1\langle y \rangle r_1\langle y \rangle upgr_1\langle y \rangle w_1\langle y \rangle unl_1\langle y \rangle wlk_1\langle z \rangle w_1\langle z \rangle unl_1\langle z \rangle$

Schedules with Locks

- Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions, and let S be a schedule for \mathbf{T} .
- A *schedule with locks* S' is a schedule S which has been augmented with lock operations.
- More precisely, it is a sequence of operations of the form $r_i\langle x \rangle$, $w_i\langle x \rangle$, $rlk_i\langle x \rangle$, $wlk_i\langle x \rangle$, $unlk_i\langle x \rangle$, $upgr_i\langle x \rangle$, and $dngr_i\langle x \rangle$ which satisfies:
 - If the lock, unlock, upgrade, and downgrade operations are removed, the result is a schedule.
 - The rules given on the previous slide which define when these locking operations may be applied are followed.
 - The locking protocol is followed.
 - Informally, this means that objects must be locked appropriately before they are accessed.
 - This idea is expanded on the next slide.

Locking schedule: In this case, S' is said to be a locking schedule for S .

Example of a Schedule with Locks

- Here is a nonserializable schedule considered earlier.

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle w_3\langle z \rangle r_3\langle x \rangle$ $r_2\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $w_3\langle x \rangle$

- Here is one valid schedule of locks for it:

$rlk_1\langle x \rangle wlk_1\langle y \rangle$ $r_1\langle x \rangle r_1\langle y \rangle$ $rlk_3\langle z \rangle$ $r_3\langle z \rangle$ $upgr_3\langle z \rangle rlk_3\langle x \rangle$
 $w_3\langle z \rangle r_3\langle x \rangle$ $unlk_3\langle z \rangle unl k_3\langle x \rangle wlk_2\langle z \rangle$ $r_2\langle z \rangle$ $upgr_1\langle x \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$
 $unlk_1\langle y \rangle wlk_2\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $unlk_1\langle x \rangle wlk_3\langle x \rangle$ $w_3\langle x \rangle$
 $unlk_3\langle x \rangle unl k_2\langle x \rangle unl k_2\langle y \rangle$

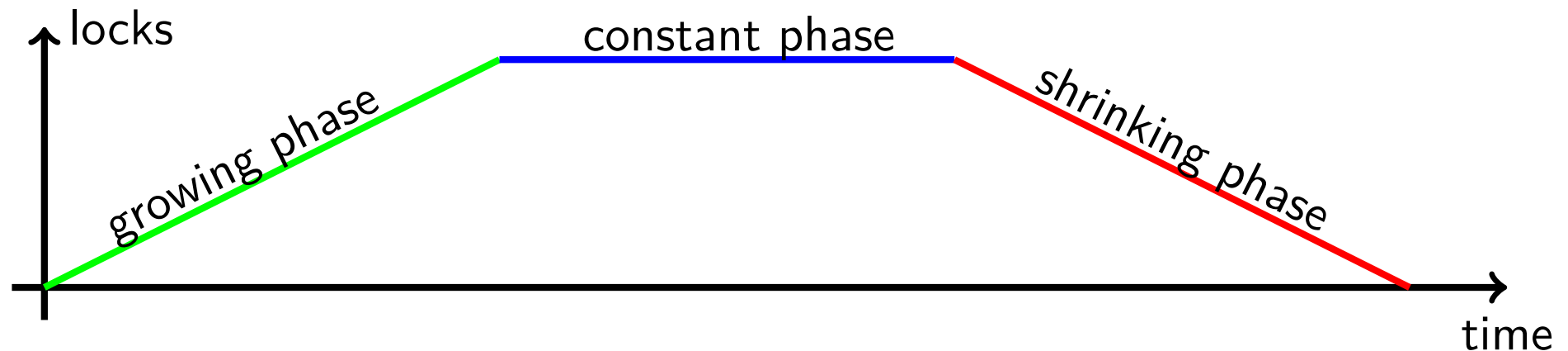
- Note that it is necessary for T_3 to lock x , release it, and then lock it again.
- This is for illustration only; it is not a reasonable schedule.

The Two-Phase Locking Protocol

- The *two-phase locking protocol* is defined for each transaction T_i individually.
- Let T_i be a transaction with locks.

Condition for two-phase locking (2PL): T_i satisfies the *two-phase locking protocol (2PL)* if no lock or upgrade operation comes after an unlock or downgrade operation in the ordering.

- All lock and upgrade operations precede all unlock and downgrade operations.



Definition: A schedule with locks is defined to be 2PL if each of its transactions with locks has that property.

Examples of 2PL

- Here is a nonserializable schedule considered earlier.

$r_1\langle x \rangle r_1\langle y \rangle$ $r_3\langle z \rangle w_3\langle z \rangle r_3\langle x \rangle$ $r_2\langle z \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $w_3\langle x \rangle$

- Here is one valid schedule of locks for it:

$rlk_1\langle x \rangle wlk_1\langle y \rangle$ $r_1\langle x \rangle r_1\langle y \rangle$ $rlk_3\langle z \rangle$ $r_3\langle z \rangle$ $upgr_3\langle z \rangle rlk_3\langle x \rangle$
 $w_3\langle z \rangle r_3\langle x \rangle$ $unlk_3\langle z \rangle unl k_3\langle x \rangle wlk_2\langle z \rangle$ $r_2\langle z \rangle$ $upgr_1\langle x \rangle$ $w_1\langle x \rangle w_1\langle y \rangle$
 $unlk_1\langle y \rangle wlk_2\langle y \rangle$ $w_2\langle z \rangle w_2\langle y \rangle$ $unlk_1\langle x \rangle wlk_3\langle x \rangle$ $w_3\langle x \rangle$
 $unlk_3\langle x \rangle unl k_2\langle x \rangle unl k_2\langle y \rangle$

- In this schedule with locks, T_1 and T_2 are 2PL, but T_3 is not.
- Hence, the schedule is not 2PL.

2PL Schedules with Locks

- Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions, let S be a schedule for \mathbf{T} , and let S' be a locking schedule for S .

Theorem: If S' is 2PL, then S is conflict serializable. \square

- Call a schedule with locks S' *view serializable* (resp. *conflict serializable*) iff the underlying schedule without locks has that property.

Corollary: If S' is 2PL, then it is view serializable. \square

Remark: There exist schedules with locks which are conflict serializable but not 2PL.

Example: Let

$$T_1 = w_1\langle x \rangle w_1\langle y \rangle$$

$$T_2 = r_2\langle x \rangle r_2\langle z \rangle$$

$$T_3 = r_3\langle y \rangle$$

and let $S = w_1\langle x \rangle r_2\langle x \rangle r_3\langle y \rangle r_2\langle z \rangle w_1\langle y \rangle$.

- S is conflict serializable with $T_3 < T_1 < T_2$.
- There is no 2PL locking schedule for S .

Assessment of 2PL

Question: To what extent is 2PL useful in real systems?

- The answer is not a simple one.
- There are at least three issues which must be considered.

Recoverability: If a transaction does not finish normally, that is, if it *aborts*, it must be handled in such a way that preserves the integrity of the remaining transactions.

Management of deadlock: Transactions can *deadlock* in their requests for resources. If they occur, these deadlocks must be resolved.

Implications of locking: Locking entails significant costs, and can reduce parallelism immensely.

- Each of these issues will be considered in turn.

Termination of Transactions

- The *atomicity* requirement of ACID demands that a transaction either run to completion or else have no effect on the database.

Commit: When a transaction *commits*, its results are irrevocably entered into the database, and the transaction ceases to exist.

Abort: When a transaction *aborts*, it is terminated without entering any updates into the database.

- The model of transactions which has been considered so far does not take the possibility of abort into account.

Problem: What if a transaction aborts after executing at least one write operation?

- A second transaction may have read from that write.
- The effects of that second transaction must be reversed.

Question: What if that second transaction has already committed?

- This process can lead to *cascading aborts* of many transactions.
- A more detailed analysis of this phenomenon is required.

The Commit Operation

- When a transaction has completed its operations successfully, it *commits*.
 - The results of its operations are made a permanent part of the database.
 - The transaction ceases to exist and so cannot be aborted any more.
- The commit operation is, by definition, the last thing that a successful transaction does.
- It is useful to express the commit operation explicitly.
- Write *cmt_i* to indicate that transaction T_i commits.

Example: $T_1 = r_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle w_1\langle z \rangle$ with explicit commit is written
 $T_1 = r_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle w_1\langle z \rangle \text{cmt}_1$.

- Call such a representation a *transaction with explicit commit*.

Schedules with Explicit Commits

- It is often useful to write *schedules with explicit commits* for its transactions.

Examples: In this example, the respective commit operations occur immediately after the end of each transaction.

$r_1\langle x \rangle r_1\langle y \rangle \quad r_3\langle z \rangle w_3\langle z \rangle \text{ cmt}_3 \quad r_2\langle z \rangle \quad w_1\langle x \rangle w_1\langle y \rangle \text{ cmt}_1 \quad w_2\langle z \rangle w_2\langle y \rangle \text{ cmt}_2$

- However, this is not required.
- Each of the following is also admissible.

$r_1\langle x \rangle r_1\langle y \rangle \quad r_3\langle z \rangle w_3\langle z \rangle \quad r_2\langle z \rangle \quad w_1\langle x \rangle w_1\langle y \rangle \quad w_2\langle z \rangle w_2\langle y \rangle \text{ cmt}_1 \text{ cmt}_2 \text{ cmt}_3$
 $r_1\langle x \rangle r_1\langle y \rangle \quad r_3\langle z \rangle w_3\langle z \rangle \quad r_2\langle z \rangle \quad w_1\langle x \rangle w_1\langle y \rangle \text{ cmt}_3 \quad w_2\langle z \rangle \text{ cmt}_1 \quad w_2\langle y \rangle \text{ cmt}_2$

Nonrecoverable Schedules

Example: Consider the following two simple transactions:

$$T_1 = r_1\langle x \rangle w_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle$$

$$T_2 = r_2\langle x \rangle w_2\langle x \rangle$$

- and the following schedule:



in which T_1 aborts before completion.

- Note that T_2 read the value of x from T_1 .
- Since T_1 aborts, this value is invalid.
- Thus, T_2 must be aborted as well.
- But it has committed.
- This is an example of a *nonrecoverable schedule*.

Cascading Nonrecoverability

Example: Consider the following three simple transactions:

$$T_1 = r_1\langle x \rangle w_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle$$

$$T_2 = r_2\langle x \rangle w_2\langle x \rangle r_2\langle z \rangle w_2\langle z \rangle$$

$$T_3 = r_3\langle z \rangle w_3\langle z \rangle$$

- and the following schedule:



- T_2 reads x from T_1 and then commits.
- T_3 reads z from T_2 and then commits.
- Both T_2 and T_3 must be aborted even though they have committed.
- This illustrates *cascading nonrecoverability*.
- It may clearly be extended to any finite number of transactions.

Recoverable Schedules

- A schedule is *recoverable* if T_j reads from T_i implies that T_i commits before T_j .

Example: Consider again the following two simple transactions:

$$T_1 = r_1\langle x \rangle w_1\langle x \rangle r_1\langle y \rangle w_1\langle y \rangle$$

$$T_2 = r_2\langle x \rangle w_2\langle x \rangle$$

- The first schedule below is recoverable, while the other two are not.

$r_1\langle x \rangle w_1\langle x \rangle$	$r_2\langle x \rangle w_2\langle x \rangle$	$r_1\langle y \rangle w_1\langle y \rangle$	$\text{cmt}_1 \text{ cmt}_2$
$r_1\langle x \rangle w_1\langle x \rangle$	$r_2\langle x \rangle w_2\langle x \rangle$	$r_1\langle y \rangle w_1\langle y \rangle$	$\text{cmt}_2 \text{ cmt}_1$
$r_1\langle x \rangle w_1\langle x \rangle$	$r_2\langle x \rangle w_2\langle x \rangle$	cmt_2	$r_1\langle y \rangle w_1\langle y \rangle \text{ cmt}_1$

2PL and Recoverability

- It is easy to see that 2PL does not guarantee recoverability.

Example: The following schedule is 2PL but not recoverable.

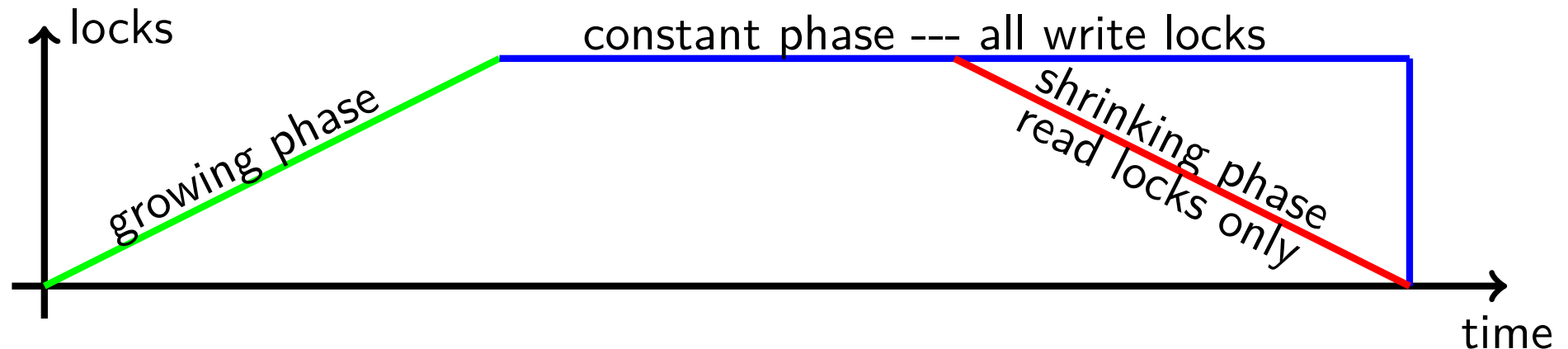
$wl k_1 \langle x \rangle wl k_1 \langle y \rangle$ $r_1 \langle x \rangle w_1 \langle x \rangle$ $unlk_1 \langle x \rangle wl k_2 \langle x \rangle$
 $r_2 \langle x \rangle w_2 \langle x \rangle$ $unlk_2 \langle x \rangle$ cmt_2 $r_1 \langle y \rangle w_1 \langle y \rangle$ $unlk_1 \langle y \rangle$ cmt_1

- To guarantee recoverability, transactions must not release locks too early.

Strict 2PL

- One way to ensure recoverability is to require that each transaction retain all of its write locks until it commits.
- Let T_i be a transaction with locks and explicit commit.

Condition for strict two-phase locking (S2PL): T_i satisfies the *strict two-phase locking protocol (S2PL)* if it satisfies 2PL and all write locks are held until the transaction commits.



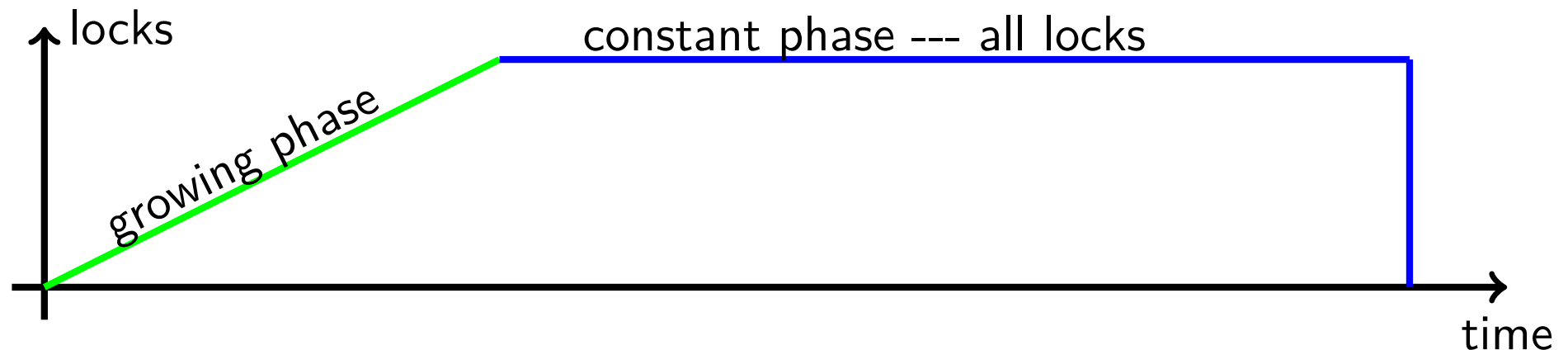
Definition: A schedule with locks and explicit commits is defined to be *strict 2PL (S2PL)* if each of its transactions has that property.

Theorem: Every S2PL schedule is recoverable. \square

Strong Strict 2PL

- Let T_i be a transaction with locks and explicit commit.

Condition for strong strict two-phase locking (SS2PL): T_i satisfies the *strong strict two-phase locking protocol (SS2PL)* if it satisfies 2PL and all locks (read and write) are held until the transaction commits.



Definition: A schedule with locks and explicit commits is defined to be *strong strict 2PL (SS2PL)* if each of its transactions has that property.

Theorem: Every SS2PL schedule is recoverable. \square

- SS2PL is also called *rigorous 2PL*.

2PL in Real Systems

- Textbooks on database systems often state that SS2PL is widely used in practice.
- The degree to which this is true will be discussed later in these lectures.
- What can be stated is the following:
 - To the extent that 2PL is used in real systems, it is of the form SS2PL.
- Nonrecoverable schedules are almost never acceptable in real systems.

Question: Why SS2PL and not S2PL?

- I do not have a good answer to that question.
- Possibly complexity of implementation is an issue.

Remark: In early literature, SS2PL was sometimes called S2PL.

- This terminology is no longer used.

The Problem of Deadlock

Motivating example: Consider the following two transactions:

$$T_1 = r_1\langle x \rangle r_1\langle y \rangle w_1\langle x \rangle$$

$$T_2 = r_2\langle y \rangle r_2\langle x \rangle w_2\langle y \rangle$$

- Suppose that scheduling of execution begins as follows:

$$wlk_1\langle x \rangle \quad r_1\langle x \rangle \quad wlk_2\langle y \rangle \quad r_2\langle y \rangle$$

- To continue, either T_1 must acquire at least a read lock on y , or else T_2 must acquire at least a read lock on x .
- Neither is possible without forcing the other transaction to release a lock, which it still needs.
- A *deadlock* has occurred.
- This can happen even if T_1 and T_2 begin with read locks.

Detection of Deadlock

- Let $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be a set of transactions.
- A *lock set* for \mathbf{T} is any subset of
$$\{\text{wlk}_i\langle x \rangle \mid 1 \leq i \leq m \text{ and } x \text{ is a data object}\}.$$
 - For simplicity, only write locks are considered.
- A *lock situation* for \mathbf{T} is a pair (L, R) in which L and R are lock sets.
 - L is the set of locks which are currently held.
 - R is the set of locks which must be obtained in order to continue.

Wait-for graph: The (directed) *wait-for graph* for (L, R) has:

Vertices: \mathbf{T} .

Edges: $T_i \xrightarrow{x} T_j$ iff $\text{wlk}_i\langle x \rangle \in L$ and $\text{wlk}_j\langle x \rangle \in R$.

Theorem: (L, R) represents a deadlock situation iff the wait-for graph has a (directed) cycle. \square

Example of the Wait-For Graph

- Return to the motivating example:

$$T_1 = r_1\langle x \rangle r_1\langle y \rangle w_1\langle x \rangle$$

$$T_2 = r_2\langle y \rangle r_2\langle x \rangle w_2\langle y \rangle$$

- The scheduling of execution begins as follows:

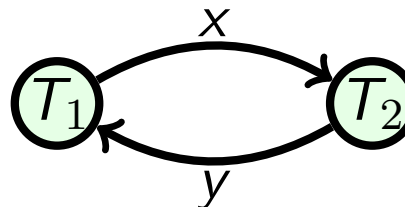
$$wlk_1\langle x \rangle \text{ } r_1\langle x \rangle \text{ } wlk_2\langle y \rangle \text{ } r_2\langle y \rangle$$

- The lock sets are:

$$L = \{wlk_1\langle x \rangle, wlk_2\langle y \rangle\}$$

$$R = \{wlk_1\langle y \rangle, wlk_2\langle x \rangle\}$$

- The wait-for graph is:



Resolution of Deadlock

- To resolve deadlock, an edge (or edges) must be removed from the wait-for graph to render it acyclic.
- There are two main approaches to managing deadlock:

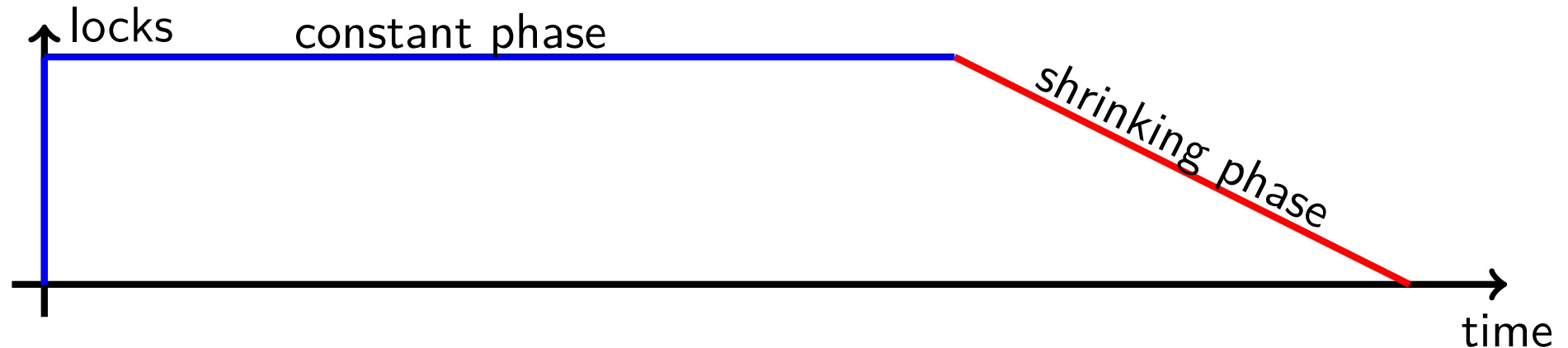
Pessimistic resolution: Do not allow a transaction to begin until it is guaranteed that it can acquire all of the locks that it needs.

Optimistic resolution: Allow transactions to proceed unhindered.

- When a deadlock is detected, abort one or more transactions in order to render the wait-for graph acyclic.

Pessimistic Resolution of Deadlock

- Pessimistic resolution may be guaranteed via *conservative 2PL*, in which all locks are acquired before the transaction is allowed to proceed.



- Pessimistic resolution is seldom employed in the DBMS context.
- Typically, conflicts due to lock contention far outnumber conflicts due to deadlock.
- Also, when a transaction begins, it is not always known which resources it will need.
 - ⇒ Conservative 2PL may result in many unnecessary locks.

Bottom line: The performance penalty imposed by conservative 2PL outweighs the advantages gained.

Optimistic Resolution of Deadlock

- Optimistic resolution of deadlock proceeds by choosing a *victim* transaction to abort when a deadlock is detected.

Livelock: Livelock (also called *starvation*) occurs when a given transaction is chosen to be the victim over and over, and so never is able to complete.

- Livelock may be avoided by timestamping each transaction with the time of its initial begin.
- When a transaction is restarted after an abort, it is restarted with its timestamp.
- When a victim is chosen due to deadlock, it is the youngest transaction in the cycle.
- In this way, transactions which have been aborted repeatedly receive increasing priority and will eventually complete.

Caution: *Optimistic resolution of deadlock* and *optimistic concurrency control* are two entirely different things which address two completely different issues.

Granularity of Locks

Question: What size of objects should be locked? (*lock granularity*)

- At first thought, it might seem best to lock the smallest possible objects.
- Smaller lock objects (finer granularity) have the advantage of allowing increased parallelism due to lesser contention for data objects.
- However, finer granularity of locks implies greater overhead from lock management.

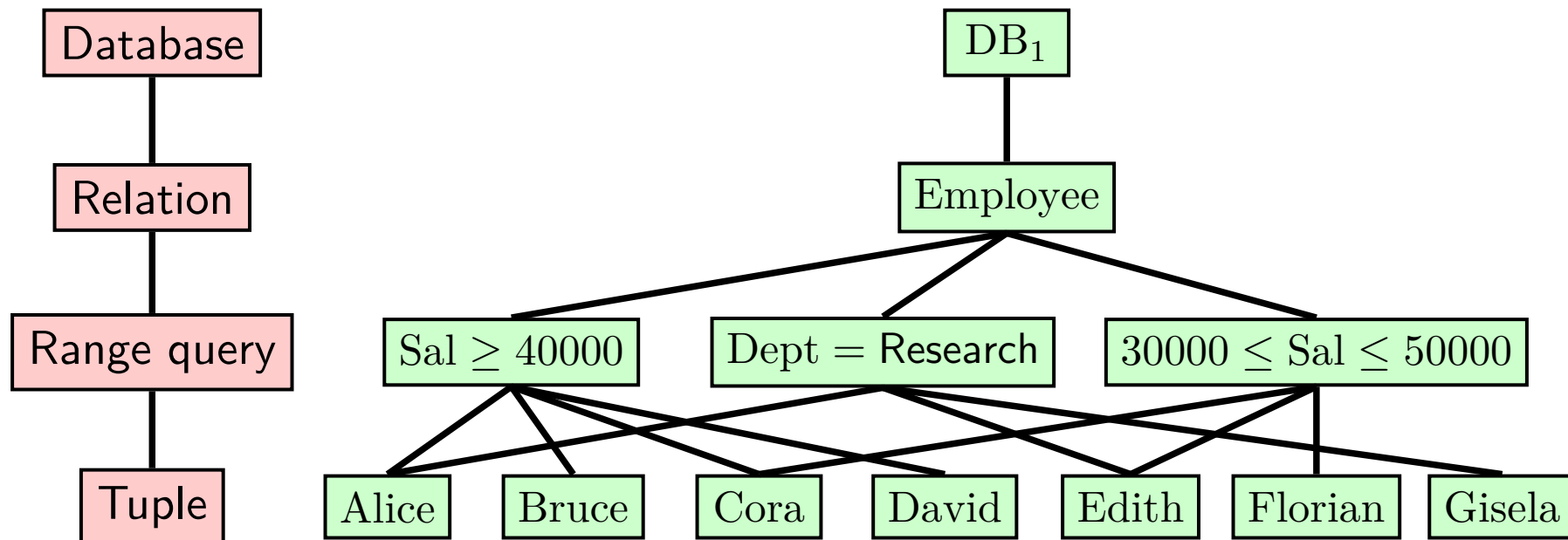
Observation: Different transactions may require different lock granularities.

- Transaction A processes a whole relation or a large part of a relation and so works best with coarse-grained locks.
- Transaction B processes only a few tuples at a time and so will interfere less with other transactions if its locks are fine grained.
- Transaction C needs to read lock an entire relation but then updates only a small part of it.

⇒ It is thus advantageous to allow read and write locks of differing granularity.

Multigranular Locking

- The classes of objects to be locked are arranged in a hierarchy.
- A simple example is shown in pink to the left below.
- An example of object instances is shown to the right in green.
- This hierarchy need not be a tree, but its (directed from parent to child) graph must be acyclic.
- When an object is (read/write) locked, all objects below it are also (read/write) locked.
- Note also there is a hierarchy within the range queries (not shown).



Intention Locks

- Locking all objects below a given object provides correct semantics of multigranular locking.
- Additional efficiency may be obtained by propagating a certain type of lock to those objects above the object to be locked.
- When an object is locked, all objects above it are assigned an *intention lock* of the same type.

Example: Consider the following two transactions:

- T_1 updates information on employee Alice.
- T_2 updates information on employees in the Research department.
- Before T_1 is allowed to write lock the Alice tuple, it must obtain an intention write lock on all employees in the Research department, as well as the Employees relation and the whole database.
- This intention lock ensures that no other transaction will be able to lock those objects which subsume the Alice tuple.

Intention Locks — Formalization

Intention-to-read locks: Before a transaction T_i may obtain a read lock $rlk_i\langle x \rangle$ on object x , it must obtain an intention-to-read lock $irlk_i\langle x' \rangle$ on every data object x' above and including x in the hierarchy.

- Also called a *intention-shared lock* or *IS-lock*.

Intention-to-write locks: Before a transaction T_i may obtain a write lock $wlk_i\langle x \rangle$ on object x , it must obtain an intention-to-write lock $iwlk_i\langle x' \rangle$ on every data object x' above and including x in the hierarchy.

- Also called an *intention-exclusive lock* or *IX-lock*.

Compatibility matrix: Shows which types of locks are compatible.

Type of Lock Held		Type of Lock Requested			
		$rlk_i\langle x \rangle$	$wlk_i\langle x \rangle$	$irlk_i\langle x \rangle$	$iwlk_i\langle x \rangle$
	$rlk_j\langle x \rangle$	yes	no	yes	no
	$wlk_j\langle x \rangle$	no	no	no	no
	$irlk_j\langle x \rangle$	yes	no	yes	yes
	$iwlk_j\langle x \rangle$	no	no	yes	yes

Read with Intention to Write

- It is often the case that a transaction will require a write lock on some object x as well as a read lock on some object x' which is above x in the hierarchy.

Example: Execute a large range query (read only), and then update (write) just a few tuples.

RIW-locks: A *read-with-intention-to-write lock* $\text{riw}l_k_i\langle x \rangle$ is equivalent to $\text{r}l_k_i\langle x \rangle$ and $\text{i}w_l_k_i\langle x \rangle$ together.

- Also called a *shared and intention-exclusive lock* or *SIX-lock*.

Compatibility matrix: Shows which types of locks are compatible.

Type of Lock Held		Type of Lock Requested				
		$\text{r}l_k_i\langle x \rangle$	$\text{w}l_k_i\langle x \rangle$	$\text{i}r_l_k_i\langle x \rangle$	$\text{i}w_l_k_i\langle x \rangle$	$\text{riw}l_k_i\langle x \rangle$
	$\text{r}l_k_j\langle x \rangle$	yes	no	yes	no	no
	$\text{w}l_k_j\langle x \rangle$	no	no	no	no	no
	$\text{i}r_l_k_j\langle x \rangle$	yes	no	yes	yes	yes
	$\text{i}w_l_k_j\langle x \rangle$	no	no	yes	yes	no
	$\text{riw}l_k_j\langle x \rangle$	no	no	yes	no	no

Concurrency Control in Real Systems

- SS2PL provides the degree of transaction correctness and recoverability needed for true isolation.
- However, it comes with a price.

Example: Consider an update which first requires a range query on an attribute which is not indexed.

- Give all employees with $29000 \leq \text{Salary} \leq 30000$ a 10% raise.
- To execute that query, the entire Employee relation must be read locked in order to identify those employee which meet the range condition.
- Those parts which are to be updated must then be write locked, and this can cause a huge delay if the read lock is shared.
- Additionally, according to SS2PL, this read lock cannot be released until the entire transaction has completed.
- Locking the entire relation, even for reading, can have a serious impact on performance.
- Other methods of implementing isolation are therefore often used in practice.

Standard Degrees of Isolation

- The SQL standard specifies four degrees of isolation.
- Each is described in terms of certain anomalies.
- They are summarized in the table below.

Degree of Isolation	Level of isolation	Anomaly allowed		
		Dirty read	Nonrepeatable read	Phantom
1	Read uncommitted (RU)	Yes	Yes	Yes
2	Read committed (RC)	No	Yes	Yes
3	Repeatable read (RR)	No	No	Yes
4	Serializable (SER)	No	No	No

- Each of these modes requires write locks, but not necessarily following 2PL.
- The differences lie in the degrees of read locks required.
- Because of their importance in practice, each will be discussed briefly.

Degree 1 Isolation — Read Uncommitted

- Under the *read-uncommitted (RU)* isolation level, data which are not committed may be read by a transaction.
- In the context of locking, no locks are required for reading.
 - Thus, the usual locking protocol need not be followed.
- This mode allows *dirty reads*.

- An example in which T_2 reads dirty data is shown below.

$w_{l_1}\langle x \rangle$ $r_1\langle x \rangle w_1\langle x \rangle$ $w_{l_2}\langle y \rangle$ $r_2\langle x \rangle r_2\langle y \rangle w_2\langle y \rangle$ $r_1\langle z \rangle$ abort_1

- Since T_1 aborts, the value which it wrote for x is invalid.
- However, T_2 uses it anyway.
- Degree 1 isolation is useful in computing summary results, where small errors are not an issue.

Degree 2 Isolation — Read Committed

- Under the *read-committed (RC)* isolation level, only committed data may be read.
 - There are no other guarantees, however.
- In the setting of a locking protocol, this isolation level is often implemented by requiring that a transaction acquire a read lock before reading a given data item.
 - But the lock may be released as soon as the item has been read.
 - Another transaction may alter that data before the original reader commits.
- This mode allows *nonrepeatable reads*.
- An example is shown on the next slide.

Example of Nonrepeatable Read

- The classical form of a nonrepeatable read occurs when a transaction reads the same data object x twice, with different results (because another transaction wrote x between the two reads).
- A more general form of nonrepeatable read occurs when a transaction T_1 reads two data objects x and y , T_2 writes both data objects between the reads of T_1 , as illustrated below.

$\text{rlk}_1\langle x \rangle$ $\text{r}_1\langle x \rangle$ $\text{unlk}_1\langle x \rangle$ $\text{wlk}_2\langle x \rangle$ $\text{wlk}_2\langle y \rangle$ $\text{r}_2\langle x \rangle$ $\text{r}_2\langle y \rangle$ $\text{w}_2\langle x \rangle$ $\text{w}_2\langle y \rangle$
 $\text{unlk}_2\langle x \rangle$ $\text{unlk}_2\langle y \rangle$ cmt_2 $\text{rlk}_1\langle y \rangle$ $\text{wlk}_1\langle z \rangle$ $\text{r}_1\langle y \rangle$ $\text{w}_1\langle z \rangle$ $\text{unlk}_1\langle y \rangle$ $\text{unlk}_1\langle z \rangle$ cmt_1

Concrete interpretation: x and y hold account balances.

- T_1 computes $z \leftarrow x + y$.
- T_2 transfers 100€ from x to y .
- Note that this is not possible with 2PL.
- Note that the read uncommitted error illustrated on Slide 57 is not possible under read committed.

Degree 3 Isolation — Repeatable Read

- Under the *repeatable read (RR)* isolation level, the value read from a single data item must be the same over multiple reads by the transaction.
- However, the set of tuples in a range query may change.
- In the context of locks, this issue is whether read locks are allowed on nonexistent tuples.
 - With RR, read locks are not necessary on nonexistent tuples.
 - With serializable isolation (SER), read locks are necessary for all *possible* tuples in the range of the query.
 - Such nonexistent tuples are called *phantoms*.
- An example is presented on the next slide.

Example of Phantom with Repeatable Read

- An example in which T_1 performs a repeatable read may be expressed informally as follows.
 - T_1 computes the sum of the salaries of all employees.
 - T_2 inserts a new employee with a positive salary.
 - T_1 computes again the sum of the salaries of all employees.
- The second read by T_1 may return a different value than the first.
- The inserted tuple is called a *phantom*.
- This is not possible with 2PL, since in 2PL all tuples in the range of the query must be locked.
- Note, however, that the example of nonrepeatable read given previously is not possible with repeatable read isolation.
 - Reads of existing data items are by definition repeatable under RR isolation.

Full Isolation — Serializability?

Obvious fact: The serializable isolation SER mode of the SQL standard is view serializability. Right?

- Wrong! Nothing is ever easy with the SQL standard.
- The SQL standard defines serializability as the absence of dirty reads, nonrepeatable reads, and phantoms...
 - ... and there are anomalies which pass those three tests yet violate the SER isolation level.

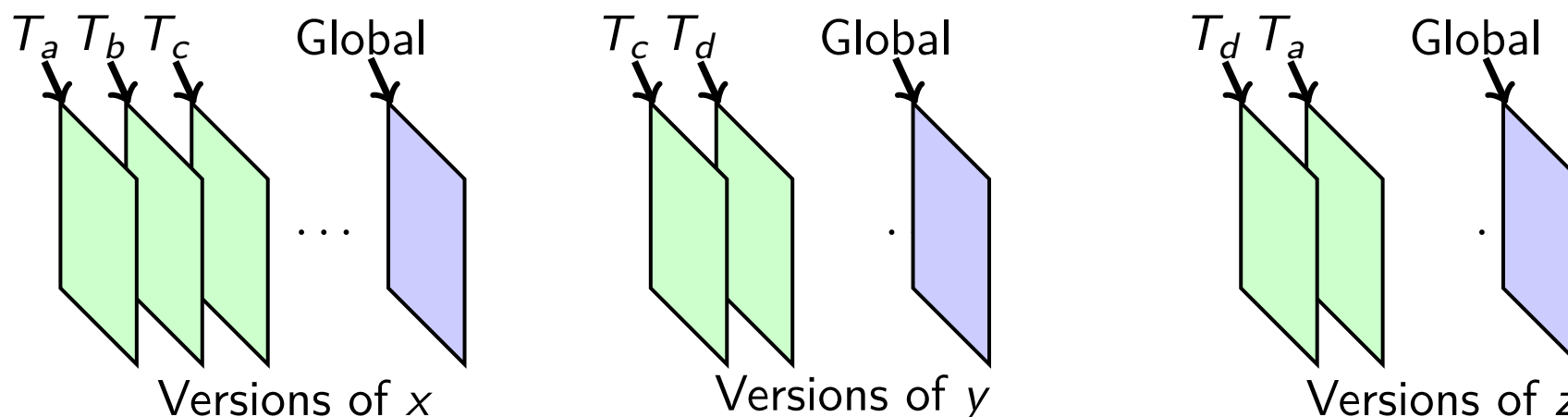
Question: But surely the major DBMS vendors implement the standard SQL serializable isolation level as SS2PL?

Answer: Until very recently, of the five major DBMSs Oracle, IBM DB2, Microsoft SQL Server, PostgreSQL, and MySQL/InnoDB...

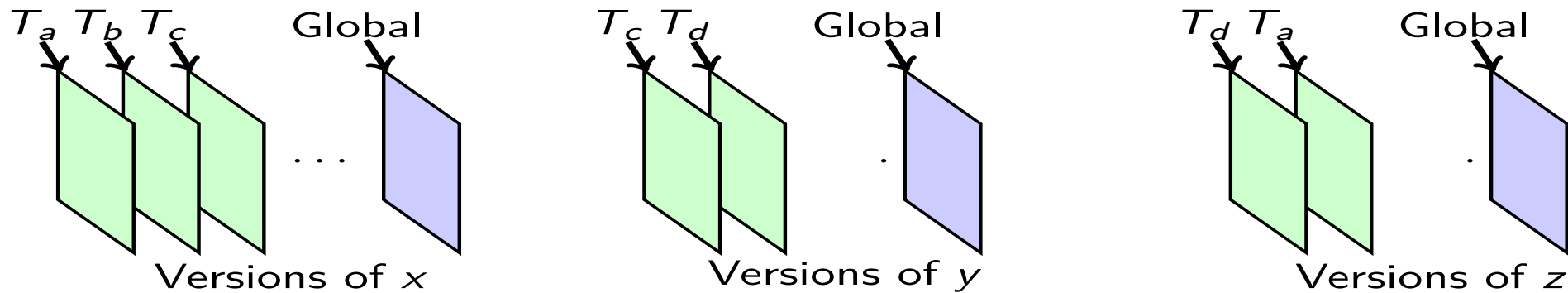
- ...only SQL Server and IBM DB2 even provided true serializable mode.
- The others provided something called *snapshot isolation* as the strongest level of isolation.
- These properties still hold for system versions in common use.

Multiversion Concurrency Control

- Historically, systems which work with just a single version of the database have have been widely used in DBMSs.
 - This model is called *single-version concurrency control (SVCC)*.
- Nowadays, however, a much more common approach is *multiversion concurrency control (MVCC)*.
- In MVCC, there may be several versions of each primitive data object.
- Exactly one version for each data object is committed, and belongs to the *global DB* – the “true” database which other transactions can see.
- The others are local copies for transactions.
- All updates by transactions are to local copies.

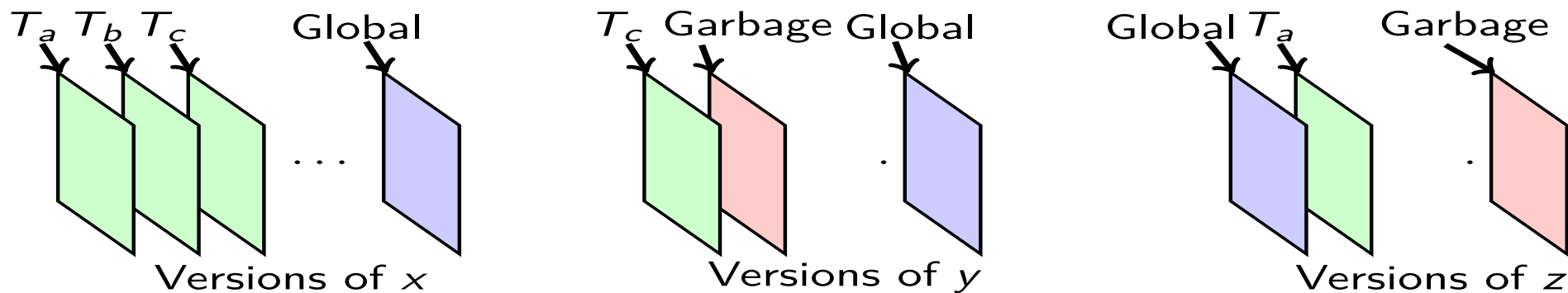


Multiversion Concurrency Control — Modify and Commit



Example: T_d reads y and writes z .

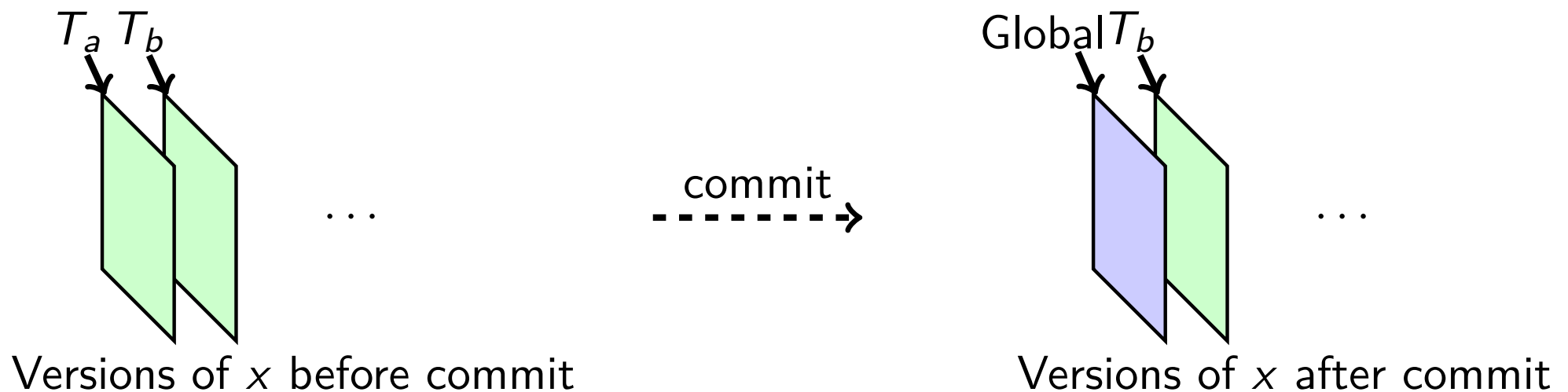
- When a transaction commits:
 - Its local writes become the new committed versions.
 - The old committed versions become garbage.
 - Its local read-only copies become garbage.



- The garbage versions are recycled via the *vacuuming process*.

Multiversion Concurrency Control — Insertion

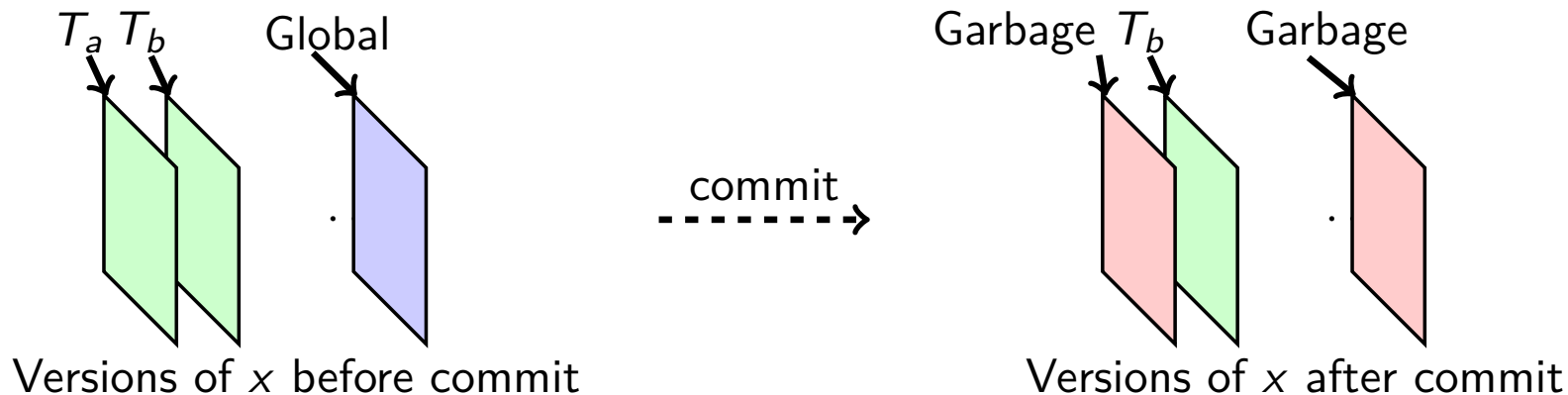
- For an insertion operation, a current value of the data object does not exist in the global database.
- Otherwise, the situation is very similar to that of a modification update.
- In the example illustrated below, transaction T_a inserts a record for data object x .



- T_b must also intend to write x , since that data object did not exist when the transaction started.
- In most reasonable isolation levels, T_b cannot commit its write.

Multiversion Concurrency Control — Deletion

- For a deletion operation, a record corresponding to the data object will not exist after the transaction commits.
- Otherwise, the situation is very similar to that of a modification update.
- In the example illustrated below, transaction T_a deletes a record for data object x .



- T_b may be just a reader, or else it may intend either to delete x also, or else to modify it.
- If T_b intends to modify x , that write cannot be committed.
- If T_b also intends to delete x , in principle, there is no conflict.
 - In practice, the redundant deletion would likely be blocked.

Motivation for Studying Lock-Based Concurrency Control

Question: Given that MVCC has become the dominant form of concurrency control in DBMSs, why study SVCC lock-based approaches?

- SVCC provides a firm theoretical foundation for understanding what a concurrency-control mechanism should do.
- It thus forms something of a reference model.
- All MVCC implementations use some locking and thus share features with SVCC lock-based approaches.
- This will become clearer when recovery techniques are studied.
- All that having been said, given its importance, most DBMS textbooks unfortunately do not provide anything close to adequate coverage of MVCC.

The Models of MVCC Used in this Presentation

- It is relatively straightforward to implement classical lock-based concurrency control, and in particular 2PL, S2PL, and SS2PL, within MVCC.
- For reasons of limited time, the details will not be given in these lectures.
- Rather, a higher-level *version-based* model will be used, which builds upon the idea that each transaction sees a distinct version of the database.
- The details how this version-based model is mapped to the lower-level, data-item-based model will not be given in these lectures.
 - However, the general principles follow the MVCC model presented in the previous slides.
- The focus here is upon concepts, not implementation details.

The Version-Based Model of MVCC

- The following summarizes the version model used for the data objects which are used by transactions.
- The basic idea behind the version-based model of MVCC is that there are many *versions* of the database.
- One of these versions is, of course, the global database, reflecting just the committed updates.
- (Obviously), the whole database is not replicated in each copy.
- Non-committed versions of a data object are private to a single transaction.
- Uncommitted updates by a transaction T_i are usually (always?) reflected only in a private version of the database which is associated with T_i .
- Reads by T_i are made from a version which is determined by the isolation level.

Two Fundamental Modes for MVCC

- There are two fundamental modes for the version-based model of MVCC, which correspond to distinct levels of isolation.

Snapshot mode: In *snapshot mode*, a “snapshot” of the database is taken for transaction T_i when it begins execution.

- Throughout the life of the transaction, it reads from and writes to this snapshot.
- This mode is used to define a new and very important isolation mode called *snapshot isolation*.

Read-Committed dynamic mode (RC dynamic mode): In this mode, reads of data objects which the transaction has not yet written are always made from the latest committed version of that object.

- However, once a transaction has written data item x , that value becomes part of its private version and it no longer see values of x which are committed afterwards by other transactions.
- Often used to implement the classical RC mode within MVCC.

Read-Committed Isolation in MVCC

- Suppose that transaction T_i is operating with the RC isolation level.
- It will then use RC dynamic mode of MVCC.
- All reads are made from the (unique) current committed version of the database.
 - Note that this version can (and usually will) change during the lifetime of T_i .
- Writes are always made to a private version of the database which is associated with T_i .
- If the transaction is to be aborted, this private version is simply discarded.
 - Since it is invisible to the other transactions, it is not necessary to “undo” anything which T_i has written to this private copy.
- Before T_i can commit, its private copy must be integrated into the database.

Managing Update Conflicts in MVCC

- Before moving on to snapshot isolation, it is important to sketch how update conflicts and transaction commits are managed in MVCC.

Update conflict: Say that T_i and T_j are in *update conflict* if they run concurrently (overlap in time) and their private versions contain at least one update on a common data object.

- There are many ways to resolve such conflicts.
- One natural one is...

First Committer Wins (FCW): Let T_i be a transaction.

- No locks are required.
- No action is taken until a T_i is ready to commit.
- When it is, the global version of the database is checked to see whether any committed updates have been made to data items which T_i has updated in its private version.
- If there is a conflict, T_i is not allowed to commit.
- Otherwise, its updates are committed to the global database.
- FCW may lead to livelock-like problems.

Managing Update Conflicts in MVCC via Locks

- Update conflicts in MVCC may also be managed using write locks.
- A transaction must write lock all data objects which it intends to write, but there are no read locks.
- If T_i and T_j each hold a write lock on the same data object x , they may proceed to update their local copies.
- However, only one will be allowed to commit.
 - The choice of which is to commit may be made in many ways.
- A particular case, which is widely used in practice, is the following.

First updater wins (FUW): Let T_i and T_j be transactions.

- In this protocol, a transactions still write locks the data objects which it intends to update.
- If T_j already holds a lock on some data object x which T_i also wishes to write, then T_i must wait until T_j commits or aborts.
- If T_j commits, then T_i must abort.
- If T_j aborts, then T_i may obtain a write lock on x and continue.
- Livelock-like problems may be resolved via timestamping together with the modification *oldest updater wins (OUW)*.

Managing Update Conflicts under RC Isolation

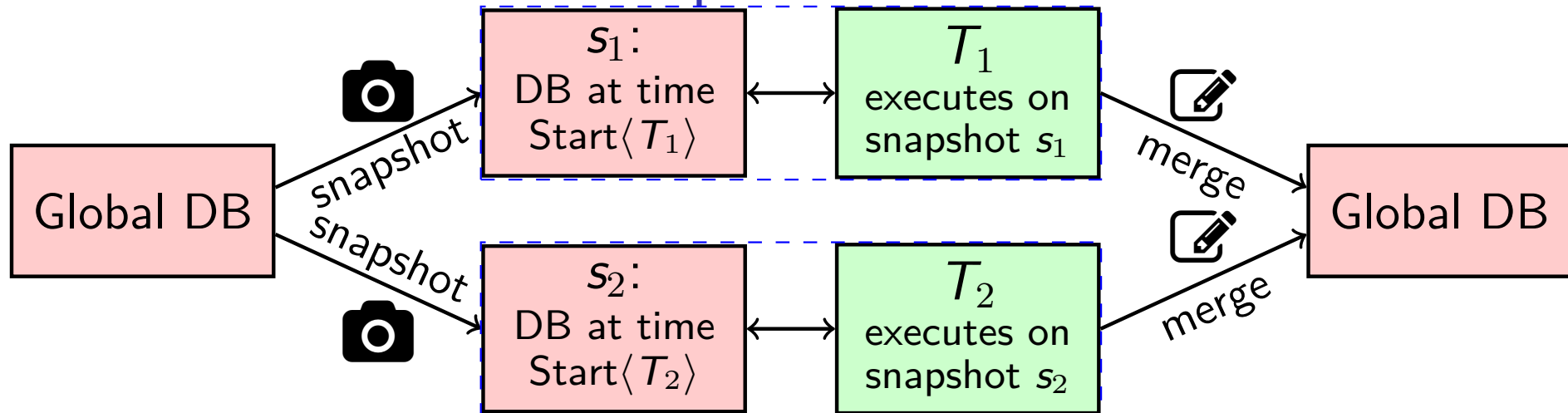
- FCW and FUW are natural choices for MVCC when operating in snapshot mode, since all external reads occur at the beginning of the transaction.
- For RC dynamic mode, things are much more complex, and the result can depend upon the winner of a race condition.
- In contrast to SI, there is no prohibition of concurrent writes to the same data object under RC.
- However, the second writer must generally wait for the first to commit or roll back, as in FUW.
- The second writer will overwrite the value which the first transaction wrote.
- However, if the two writers update different fields of the same data object, both updates will be visible after both commit.

Managing Update Conflicts under RC Isolation — 2

- PostgreSQL has two special clauses for the SELECT statement:
 - `SELECT ... FOR SHARE` of <table1>, <table2>, ...
 - `SELECT ... FOR UPDATE` of <table1>, <table2>, ...
- These two statements issue read locks (resp. write locks) on the retrieved rows of the designated tables for the lifetime of the SELECT.
- If these statements are not used under RC isolation, there is a risk that the SELECT statements will not see values from a single, consistent database during its execution.
- These are not standard SQL, although other systems may provide similar ways to accomplish the same effect.
- If this sort of statement is not used, unexpected behavior may result, which may even vary from system to system.

Guideline: Use RC isolation committed with care, because unexpected results may otherwise arise.

Snapshot Isolation



Snapshot Isolation (SI): Let T_i be a transaction.

- The private version of the database for T_i is a “snapshot” of the database at the time at which T_i begins.
- This version cannot (normally) be accessed by other transactions.
- Thus, it appears to T_i that it is executing without any concurrent operations from other transactions.
- When a transaction is ready to commit, the updates to its local version must be integrated into the main database.
- This is typically realized via the *first-updater-wins* protocol.

Anomalies of Snapshot Isolation

- With SI, dirty reads and nonrepeatable reads cannot occur.
- Thus, the isolation level is at least as great as those provided by RU (read uncommitted) and RC (read committed).
- Whether SI is as strong as repeatable read depends upon technical details of the model.
 - See [Berenson et al. 1995] for a detailed discussion.
- However, with the anomaly model presented here, it is strictly stronger than RR (repeatable read).
- Therefore it satisfies the conditions for the Serializable mode of the SQL standard.
- There are nevertheless other anomalies which cannot occur in true serializable mode.
- The two principal ones are known as *write skew* and *SI read-only anomaly*.

Write Skew

Example: Let x and y represent the balances of two distinct accounts.

Integrity constraint: $x + y \geq 500\text{€}$.

Initial state: $x = 300\text{€}$, $y = 300\text{€}$.

T_1 : Withdraw 100€ from x .

T_2 : Withdraw 100€ from y .

- To verify that its update will not violate the integrity constraint, each transaction reads both x and y .
- Suppose that the two transactions run concurrently, so that they see the same initial state “snapshot”.
- Each runs without knowledge of what the other does.
- The final state is $(x, y) = (200\text{€}, 200\text{€})$, which violates the constraint.
- This schedule clearly does not involve dirty reads, nonrepeatable reads, or phantoms, so it passes the test for the isolation levels RU, RC, and RR.
- Write skew cannot occur with true serializability.
- Thus, SI is strictly weaker than SER isolation.

Schedules for SI

- Transaction schedules for SI have a special form, due to the snapshot-at-start, write-at-commit nature of SI.
- Specifically, in a schedule for SI:

Read for the snapshot: All reads for a transaction occur at its beginning.

Write at commit: All writes for a transaction occur at its end.

Example:

bgn₁ r₁⟨x⟩w₁⟨x⟩ cmt₁ bgn₂ r₂⟨x⟩ bgn₃ r₃⟨z⟩ bgn₄ r₄⟨x⟩w₄⟨x⟩ cmt₄ w₂⟨y⟩ cmt₂ w₃⟨x⟩ cmt₃

SI Read-Only Anomaly

- This anomaly is interesting in that two transactions produce a final result which is consistent with a serializable schedule.
 - However, a read-only transaction sees a state which is not possible in any serial schedule.

- Let

$$T_1 = r_1\langle x \rangle w_1\langle x \rangle \quad T_2 = r_2\langle x \rangle r_2\langle y \rangle w_2\langle y \rangle \quad T_3 = r_3\langle x \rangle r_3\langle y \rangle$$

- If the SI schedule is

bgn₂ r₂⟨x⟩r₂⟨y⟩ bgn₁ r₁⟨x⟩w₁⟨x⟩ cmt₁ bgn₃ r₃⟨x⟩r₃⟨y⟩ cmt₃ w₂⟨y⟩ cmt₂

then the result which T_3 sees need not be the result of a serializable schedule.

- Note that this schedule becomes serializable if T_3 is removed.
- These properties are most easily seen via a concrete interpretation, given on the next slide.

SI Read-Only Anomaly 2

- Let x and y represent the balances of bank accounts.
- If a transaction forces $x + y < 0$, then a 10% interest charge is imposed.
- Let the initial balances be $(x, y) = (0\text{€}, 0\text{€})$.
- Let T_1 read and then add 20€ to the balance of x .
- Let T_2 read both balances and then deduct 10€ from the balance of y .
 - Note that if the snapshot for T_2 is taken before T_1 commits, then 1€ in interest is also deducted.
- The final result of

$r_2\langle x \rangle r_2\langle y \rangle$ $r_1\langle x \rangle w_1\langle x \rangle$ cmt₁ $r_3\langle x \rangle r_3\langle y \rangle$ cmt₃ $w_2\langle y \rangle$ cmt₂

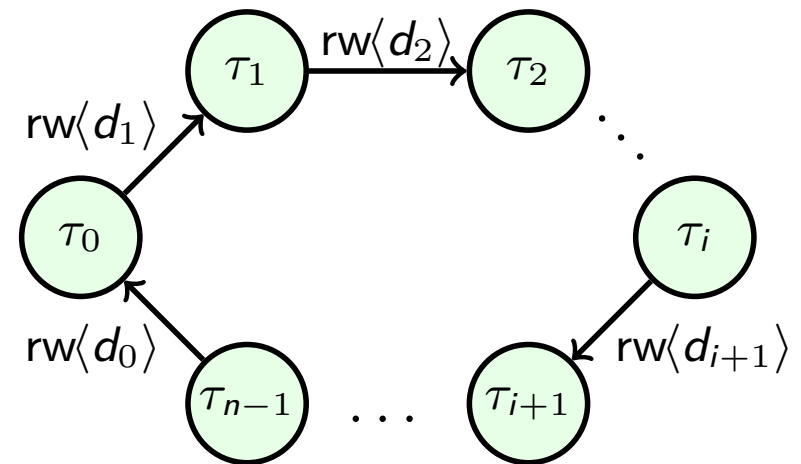
is $(x, y) = (20\text{€}, -11\text{€})$, while T_3 sees $(x, y) = (20\text{€}, 0\text{€})$.

- The values seen by T_3 cannot be the result of reading during any serializable execution of $\{T_1, T_2\}$.
 - T_3 sees the result of running T_1 on the initial data, so the serial execution would have to be $\langle T_1, T_2 \rangle$.
 - Without T_3 , the schedule is view equivalent to $\langle T_2, T_1 \rangle$.

Large Sets of Non-Serializable Transactions under SI

Example (SI permutation): $n \in \mathbb{N}$;

- d_0, d_1, \dots, d_{n-1} data objects.
- $\tau_0, \tau_1, \dots, \tau_{n-1}$ transactions with $\tau_i: d_i \leftarrow d_{(i+1) \bmod n}$.
- The n transactions, run concurrently under SI, effect a permutation of the values of the d_i 's (shift to the left).



- $\tau_i \xrightarrow{rw\langle d_i \rangle} \tau_{(i+1) \bmod n}$ denotes that τ_i reads d_i and $\tau_{(i+1) \bmod n}$ writes it.
- This behavior cannot be view serializable since if τ_i is run first, the old value of d_i is lost.
- However, if any transaction (say τ_i) is removed, the result of running all transactions concurrently under SI is serializable.
 - Run them in this order: $\tau_{i+1} \dots \tau_{n-1} \tau_0 \dots \tau_{i-1}$.

Observation: For any $n \in \mathbb{N}$, there is a set of n transaction which, when run concurrently under SI, results in nonserializable behavior, yet any proper subset produces serializable behavior under SI. No integrity constraints are involved. \square

Serializable Snapshot Isolation

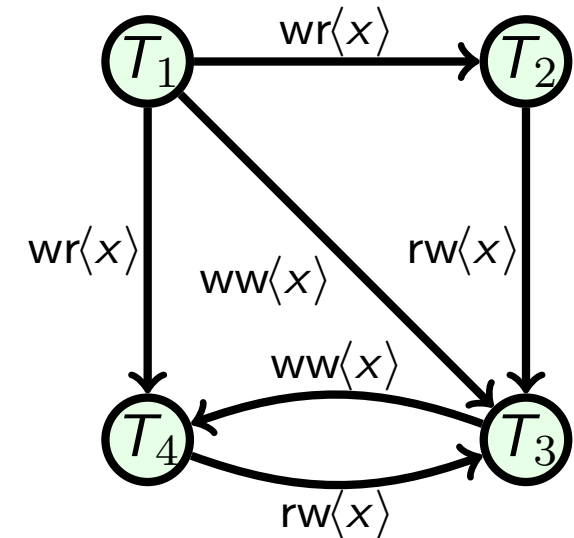
- The weakness of SI, relative to true serializable isolation, is that SI only addresses write conflicts, while conflicts which arise when transactions only reads, but do not write, a data object x lie outside of conflict detection.

Serializable SI: Very recently, a method for augmenting SI so that it always produces serializable isolation has been developed.

- This method is called *serializable SI* or *SerSI* or *SSI*.
- It is an *optimistic approach* in that it looks for certain triples, called *dangerous structures* in the DSG.
- Such dangerous structures are a necessary (but not sufficient) condition for SI to be non-serializable.
- When such a triple is found, a victim transaction is aborted to ensure serializability.
- It has recently been introduced into real systems, including PostgreSQL, and so is worth a closer look.

The Direct Serialization Graph Revisited

Recall the DSG: Example:



- Here it is very important to distinguish the three *conflict types*, or *dependency types*, characterizing an edge $T_i \longrightarrow T_j$.
 - wr-dependency:** $(T_i \xrightarrow{wr\langle x \rangle} T_j)$ T_i writes a data item x , and then T_j reads a that version.
 - ww-dependency:** $(T_i \xrightarrow{ww\langle x \rangle} T_j)$ T_i writes a data item x , and then T_j writes the next version.
 - rw-dependency:** (or *antidependency*) $(T_i \xrightarrow{rw\langle x \rangle} T_j)$ T_i reads a data item x , and then T_j writes the next version.

Concurrency and Dependency Types under SI

Context: $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ be set of transactions, S a schedule on \mathbf{T} , $T_i, T_j \in \mathbf{T}$, each running under SI.


Fact: Of the three types of dependency, rw, ww, and wr, the only type which may occur between concurrent transactions is an rw-dependency.


- Put another way, if x is a data object, and at least one of $T_i \xrightarrow{wr\langle x \rangle} T_j$ and $T_i \xrightarrow{ww\langle x \rangle} T_j$ is an edge in the DSG, then T_i and T_j cannot be concurrent.

Proof sketch: Consider the two cases:

ww: Immediate, since concurrent transaction running under SI may not write the same data object.

wr: Since the read of T_j occurs when it starts, and the write of T_j occurs when it commits, no overlap is possible. \square

 For T_j running RC dynamic mode and T_i SI, $T_i \xrightarrow{ww\langle x \rangle} T_j$ is not possible for concurrent transactions; however, $T_i \xrightarrow{wr\langle x \rangle} T_j$ is quite possible.

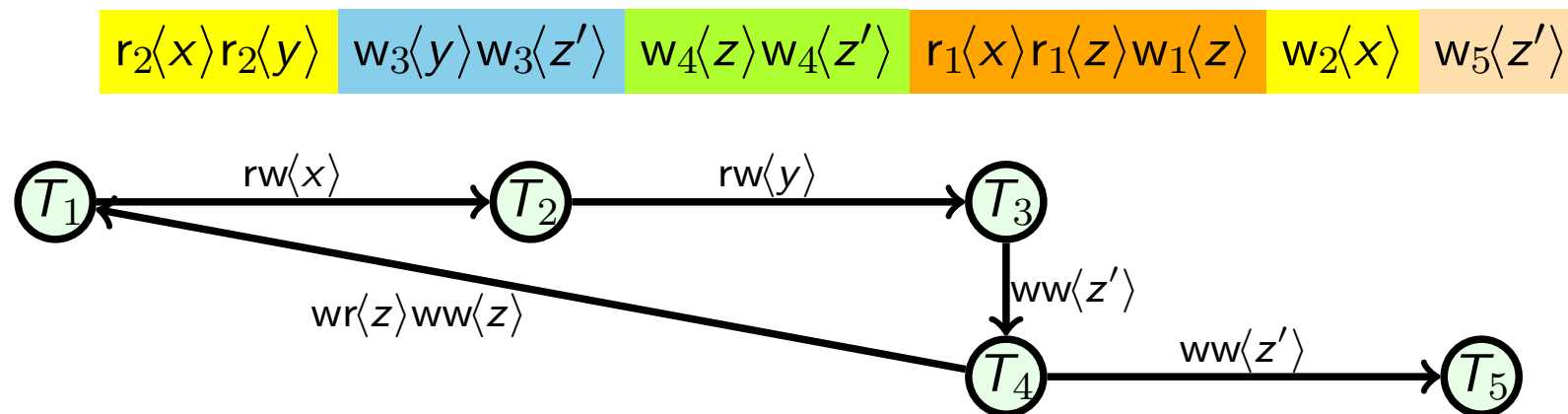
 If both T_i and T_j run RC dynamic mode, all three types of edges are possible between concurrent transactions.

View Serializability and Conflict Serializability under SI

⚠ View serializability and conflict serializability may differ even when all transactions run under SI, if there are blind writes.

Example: $T_1 = r_1\langle x \rangle r_1\langle z \rangle w_1\langle z \rangle$ $T_3 = w_3\langle y \rangle w_3\langle z' \rangle$ $T_5 = w_5\langle z' \rangle$
 $T_2 = r_2\langle x \rangle r_2\langle y \rangle w_2\langle x \rangle$ $T_4 = w_4\langle z \rangle w_4\langle z' \rangle$

- The schedule is:



- It is view equivalent to run T_4 first.

$w_4\langle z \rangle w_4\langle z' \rangle$ $r_2\langle x \rangle r_2\langle y \rangle$ $w_3\langle y \rangle w_3\langle z' \rangle$ $r_1\langle x \rangle r_1\langle z \rangle w_1\langle z \rangle$ $w_2\langle x \rangle$ $w_5\langle z' \rangle$

- This breaks the cycle in the DSG.
- The view-equivalent serialization is $\langle T_4, T_1, T_2, T_3, T_5 \rangle$.

Convention for Non-Serializability under SI

Convention: In the characterization of non-serializability, the notion of serializability which will be used will always be *conflict* serializability.

- As illustrated by the example on the previous slide, this may result in some false positives.
- However, examples which are view serializable, but not conflict serializable, are almost always pathological, so this convention makes very little difference in practice.

Characterization of Non-Serializability under SI

Concurrency: The two transactions T_i and T_j are *concurrent* if they overlap in time: either T_i begins before T_j ends or else T_j begins before T_i ends.

Observation: If there is a ww- or wr-edge from T_i to T_j , then these transactions cannot be concurrent under SI.

- T_i must commit before T_j can see its write.
- Thus, rw-edges are of primary interest in SI.

Vulnerable edge: An rw-dependency between concurrent transactions

$$T_i \xrightarrow{rw\langle x \rangle} T_j.$$

Dangerous structure: In the DSG, a *dangerous structure* is a sequence

$T_i \xrightarrow{rw\langle - \rangle} T_j \xrightarrow{rw\langle - \rangle} T_k$ of two consecutive vulnerable edges which occur in a cycle.

- T_i and T_k may be the same transaction, but need not be.

Characterization Non-Serializability Under SI — 2

Theorem: (A sufficient condition for serializability)

Given: A set S for a set $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ of transactions.

Assumed: All transactions run under isolation level SI.

Conclusion: If every cycle of the DSG for S is free of dangerous structures, then S is view serializable. \square

Important: The theorem gives a sufficient condition, but not a necessary one, for serializability.

- The condition can be sharpened a bit to give a necessary condition under a certain interpretation.
- This is described on Slide 93.

Detection of Non-Serializable SI – Example 1

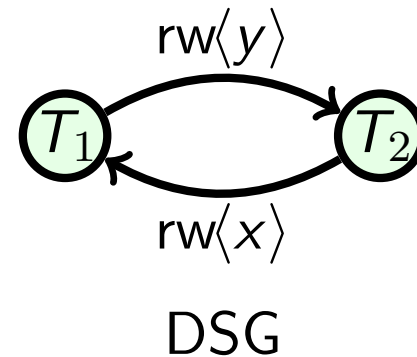
- Consider again the example of write skew: x and y represent the balances of two distinct accounts.

Integrity constraint: $x + y \geq 500\text{€}$.

Initial state: $x = 300\text{€}$, $y = 300\text{€}$.

T_1 : Withdraw 100€ from x .

T_2 : Withdraw 100€ from y .



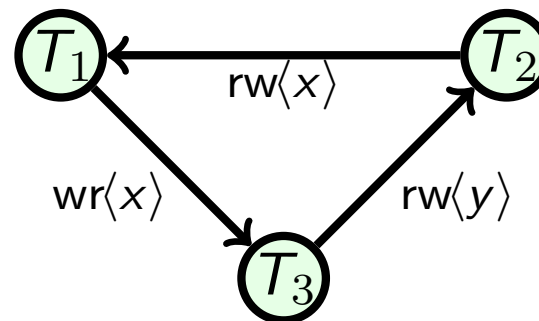
- Both transactions must read both x and y in order to check the integrity constraint.
- $T_1 \xrightarrow{rw\langle y \rangle} T_2 \xrightarrow{rw\langle x \rangle} T_1$ forms a dangerous structure, so the possibility of non-serializability exists.

Detection of Non-Serializable SI – Example 2

- Consider again the example of SI read-only anomaly with begin points marked explicitly.

$\text{bgn}_2 \text{ } r_2\langle x \rangle r_2\langle y \rangle \text{ } \text{bgn}_1 \text{ } r_1\langle x \rangle w_1\langle x \rangle \text{ } \text{cmt}_1 \text{ } \text{bgn}_3 \text{ } r_3\langle x \rangle r_3\langle y \rangle \text{ } \text{cmt}_3 \text{ } w_2\langle y \rangle \text{ } \text{cmt}_2$

- DSG:



- $T_3 \xrightarrow{rw\langle y \rangle} T_2 \xrightarrow{rw\langle x \rangle} T_1$ forms a dangerous structure in the cycle
 $T_1 \xrightarrow{wr\langle x \rangle} T_3 \xrightarrow{rw\langle y \rangle} T_2 \xrightarrow{rw\langle x \rangle} T_1$.
- Thus, this schedule is potentially not serializable.
- Note that for $T_1 \xrightarrow{wr\langle x \rangle} T_3$ to hold, T_1 must commit before T_3 begins.

Implementation of Serializable SI

- In principle, the implementation strategy is very simple.

Principle: Look for dangerous structures in the DSG.

- When such a structure is found, choose one of its member transactions as a victim, abort it, and then re-run it.

Basic Solution: A basic solution is to look for potentially dangerous structures, and abort one of the component transactions whenever such a pair occurs.

Potentially Dangerous Structure: Two consecutive vulnerable edges (need not occur in a cycle).

- This basic solution does not require building the entire DSG, nor does it require knowing the commit order.

Conservative: Due to its simple approach, the number of false positives (transactions which are unnecessarily aborted) is relatively high.

- For this reason, improvements on this basic approach have been the topic of recent investigations.

Essential Dangerous Structures

Essential dangerous structures: In the DSG, a dangerous structure

$T_i \xrightarrow{rw\langle - \rangle} T_j \xrightarrow{rw\langle - \rangle} T_k$ is *essential* if T_k is the first of the three transactions to commit.

Theorem: (Improved sufficient condition for serializability)

Given: A schedule S for a set $\mathbf{T} = \{T_1, T_2, \dots, T_m\}$ of transactions.

Assumed: All transactions run under isolation level SI.

Conclusion: If every cycle of the DSG for S is free of essential dangerous structures, then S is view serializable. \square

Necessity: This condition is also *necessary* for serializability if only the presence of a cycle, and not the identity of the data objects which are arguments of the $rw\langle - \rangle$, $ww\langle - \rangle$, and $wr\langle - \rangle$ labels, are considered.

An Improved Implementation of Serializable SI — ESSI

ESSI: Fine-tune the implementation of SSI by regarding only essential dangerous structures, and not all dangerous structures, as that which must be avoided.

- Properties of ESSI:
 - + Result in fewer false positives.
 - Requires, for each dangerous structure, waiting until one of its component transactions is ready to commit before a decision to abort is made.

Implementation of Serializable SI — PSSI

Precisely Serializable SI (PSSI): This approach augments ESSI by checking for cycles in the DSG which contain essential dangerous structures.

- It aborts a transaction only if:
 - (a) It is part of an essential dangerous structure; and
 - (b) That essential dangerous structure is part of a cycle in the DSG.
- Properties of PSSI:
 - + *Subject to the necessity condition on Slide 93*, it does not result in false positives of conflict non-serializability; it only causes an abort if non-serializability would otherwise result.
 - In addition to the added costs of ESSI, it requires maintenance of the DSG in order to check for cycles.
- It has been implemented as a modification of MySQL/InnoDB, and performance studies have been positive [Revilak et al, ICDE 2011].
 - But the transaction mix used in the study must always be taken into account.

SSI in “Real” Systems

- Until very recently, of the five major DBMSs Oracle, IBM DB2, Microsoft SQL Server, PostgreSQL and MySQL/InnoDB, only SQL Server and IBM DB2 even provided true serializable mode.
 - And these are direct, lock-based SS2PL implementations, with the associated compromise on concurrency.

PostgreSQL 9.1: As of version 9.1, the serializable isolation level of PostgreSQL is implemented as SSI.

- Thus PostgreSQL now provides true serializable isolation.
- The repeatable read isolation level is (ordinary) SI.
- Prior to version 9.1, serializable and repeatable read were both implemented as (ordinary) SI.

Performance: There have been several studies of the performance of SSI and its relatives (ESSI, PSSI), and the results are generally positive.

- It provides serializable isolation without the limitations and overhead of huge locks.
- At least in the transaction mixes which were studied, the number of aborted transactions was not limiting factor.

SI and Integrity Constraints

Selling point of SI: Readers are never blocked by writers.

- In the snapshot model, all data to a transaction are private.
- Thus, there are never any delays until commit time, when constraints must be checked.

A significant exception: Internal integrity constraints (primary keys, foreign-keys, check clauses, for example) are generally checked immediately against the committed database, not the snapshot.

- If this were not done, all integrity constraints would need to be checked at every commit of a transaction.
- This applies only to constraints specified in the DDL.
- It does not apply to constraints which are defined via triggers, for example.

Locks: This immediate checking is generally implemented via locks, even for MVCC in which locks are not otherwise used.

- This locking means that readers can be forced to wait for writers.

The Complexity of SSI – A Tradeoff

Observation: Although performance studies on “typical” transaction mixes have been positive, SSI and its variants are *optimistic* strategies, and depend upon terminating a transaction when a (potentially) nonserializable situation is found.

Tradeoff: This involves a tradeoff between higher computational complexity and a higher number of false positives:

SSI and ESSI: Test at most three transactions at a time, so the complexity is lower, but may result in false positives.

PSSI: No false positives, but must test a larger number of transactions.

Question: How large can the test set become with PSSI?

Answer: It can be as large as the largest number of (concurrent) transactions, as illustrated by the example on Slide 82.

Experimenting with Concurrency in PostgreSQL

- The general structure of a transaction:

```
BEGIN TRANSACTION;  
SET TRANSACTION ISOLATION LEVEL SERIALIZABLE;  
-- Other options (instead of SERIALIZABLE) are:  
-- REPEATABLE READ, READ COMMITTED, READ UNCOMMITTED.  
-- READ UNCOMMITTED is the same as READ COMMITTED.  
-- For versions < 9.1, SERIALIZABLE is the same as REPEATABLE READ.  
  
<SQL CODE>  
  
-- Delays can be realized by sleeping; the argument is in seconds.  
-- Delays are useful for studying the interaction of transactions.  
SELECT pg_sleep(10);  
  
-- The COMMIT directive ends the transaction.  
COMMIT TRANSACTION;
```

- To study the interaction of several transactions, run each one in a separate psql client window, on the same database.
- Use delays to enable human interaction and nontrivial concurrency.
- Some simple examples may be found on the course Web site.

SVCC vs. MVCC in Current Systems

- While SVCC used to be the norm, virtually all current-generation DBMSs use MVCC.

Why?

- MVCC offers superior support for concurrency control.
 - But it used to be too expensive to implement effectively.
 - Memory (both primary and secondary) has become much less expensive and available in much larger sizes.
 - MVCC requires lots of memory to store the versions.
- The sole holdout seems to be IBM DB2, which is still primarily based upon SVCC.
 - However, even that system now offers a mode which behaves in a way similar in many ways to MVCC.

Isolation Levels in Current Systems

- Most systems with MVCC offer RC and SI as options.
 - This is understandable since these two have natural implementations with MVCC.
- Most systems have RC as the default isolation level.
 - This is despite the fact that the SQL standard specifies SER as the default isolation level.
- SQL Server and DB2 offer SER isolation level via SS2PL with locks.
- In other systems, with the notable exception of PostgreSQL ≥ 9.1 , the isolation level which is identified as SER is generally something else (SI or something close)!

Note: In DB2, the SER isolation level is called *Repeatable Read*.

- The isolation level which is called RR in these slides is called *Read Stability* in DB2.
- The locking granularity for DB2 SER isolation is the table.
- In other words, if any part of a relation is involved in a transaction, the entire relation is locked, regardless of available indices.

Isolation Levels in Current Systems 2

- In PostgreSQL, RU is the same as RC, and RR is the same as SI.
 - This is consistent with the standard, since RU and RR isolation levels forbid certain anomalies but do not require that others be possible.
- As will be discussed shortly, RU and RR are not natural modes in MVCC.
- This PostgreSQL convention is likely used in many other systems as well.
- Many of the systems offer other non-standard modes as well.
 - In particular, locking of physical entities such as pages and files is sometimes supported.
- It is very easy to develop applications which are not portable because they use choices of isolation level which are not used by other systems.
- The best choices for portability are RC and SI.

Read-Uncommitted Isolation in MVCC

Interesting question: Does RU isolation make sense in MVCC?

- It could be implemented in a manner similar to that used for RC, except that for a data item x which T_i reads but has not yet written, the version of the database for T_i would contain the latest available version of x , regardless of whether or not it has been committed.

Question: Why would this be easier to implement than RC?

Answer: In the general MVCC context, it would not be.

- Since RC provides “superior” data to RU, there is no apparent advantage to RU over RC in MVCC.

Consequence: At least in some real DBMSs (e.g., PostgreSQL), RC and RU isolation levels are identical and behave as RC.

However: It might be possible to implement RU to advantage in the context of the DB cache.

- This possibility will be discussed in the context of recovery.

Repeatable-Read Isolation in MVCC

- Snapshot isolation already delivers a strong version of repeatable read.
- Since all reads are from the same snapshot, those reads are always repeatable.
- It is not clear how SI could be weakened within MVCC while maintaining RR.
- Thus, repeatable read as a fundamental mode of isolation, differing from SI, also seems a bit problematic in MVCC.

Optimistic and Pessimistic Concurrency Control

- In the discussion of the resolution of deadlock for lock-based SVCC, notions of *optimistic* and *pessimistic* methods for the resolution of deadlocks were presented.
- These concepts make sense in a more general context, including but not limited to MVCC, possibly without locks and without deadlocks.

Optimistic concurrency control: refers to an approach in which transactions are allowed to proceed, with conflicts resolved as late as possible, often at commit time.

Pessimistic concurrency control: refers to an approach in which potential conflicts are detected and resolved early on.

Example: Within MVCC, First Committer Wins is an example of optimistic concurrency control.

Example: First Updater Wins is an example which has both optimistic and pessimistic aspects.

Issues with Concurrent Isolation Modes

- In a system with multiple concurrency models, each transaction is allowed to choose its own concurrency model.
- The issue of isolation level is one of interacting transactions, and not just a single transaction.
- Thus, even if transaction T_1 chooses true serializable isolation, if transaction T_2 chooses read uncommitted, it can compromise the results of T_1 .
- This underlines the necessity of having a policy for transactions which support the overall goals of the enterprise.

Transactions in Current Systems

- Major DBMSs do not in general follow SQL standard specifications in regards to directives surrounding transactions.
 - Each follows its own conventions.
 - Thus, the SQL standard will not be discussed here.
- The general convention is that transaction initiation is *implicit*.
 - It is not necessary (and in some cases not possible) to give an explicit `Begin Transaction` statement or the like.
 - On the other hand, it is generally possible to give a `Commit` or `Rollback` statement.
- It is also possible to give directives to set the isolation level.

Transaction Initiation and Commit in Current Systems

- There are two general models of transaction initiation:

Session based: SQL statements are executed one after the other, but a commit occurs only at the end of the session or when an explicit Commit directive is issued.

- Default for Oracle.

Statement based (autocommit): A Commit occurs immediately after each SQL statement.

- Default for the other four systems if block markers are not used (Begin Transaction, Commit, Rollback).
- In all cases, there is a directive to choose which of these models applies to a given session.

Long-Running Transactions

- As the name suggests, *long-running transactions* are those which take a “long” time to complete.
 - They often access many different data objects, although they may need some such objects for only short a short interval.
- Long-running transactions pose a particularly difficult problem for concurrency control.
 - If an optimistic strategy is employed, then the risk is that transactions which have been running for a long time and are near completion must be aborted.
 - If a pessimistic strategy is employed, then the risk is that execution will be nearly serial and so there will be unacceptably long waits before a transaction is allowed to run.
- Solutions for dealing with long-running transactions must often be customized for the given application area.

Interactive Transactions and Conflict Resolution via Negotiation

- A special case of long-running transactions are *interactive transactions*, in which interaction with a human is integral to execution.
- In this situation, a strategy of abort and rerun to resolve conflicts is almost never acceptable.

Negotiation: An alternative is to manage conflicts via *negotiation*.

- When a conflict occurs, the conflicting transactions attempt to negotiate changes in their behavior in order to avoid the conflict.
- This is particularly attractive in interactive situations, since the humans involved can make decisions about which compromises to make.
- However, it requires a way to identify conflicts (preferably) between only two transactions, or at least between very few transactions.

Constraint Preservation and SI

- Complete serializability is often not required.
- However, preservation of integrity constraints is always necessary.
 - As illustrated by the write skew example, among others, SI is not always constraint preserving.

CPSI: A compromise solution is *constraint-preserving snapshot isolation*.

- Full serializability is not guaranteed, but it is ensured that constraints are always satisfied.
- This includes constraints implemented using triggers, which are not maintained automatically by the system under SI.
- Interactive business process often involve such constraints.
- Testing for compliance is pairwise.
- False positives are possible in certain situations (just as they are for SSI), but the method nevertheless appears to be quite effective.
- Further details will not be presented here.

Index Maintenance under Transaction Concurrency — Crabbing

- Because an index must be modified as indexed entries are inserted, deleted, and updated, it is necessary to have a concurrency policy.

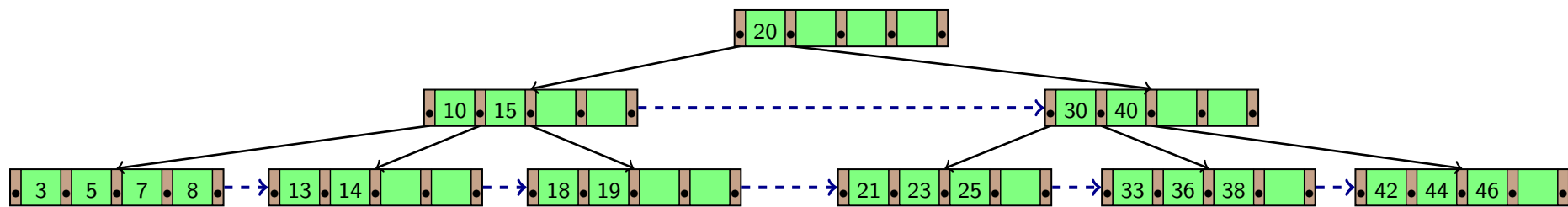
Crabbing: One simple method for B^+ -trees is to traverse the index using a technique called *crabbing*.

- To find a record, instead of read locking the entire index, proceed as follows:
 - Read lock the root v_0 .
 - Identify the appropriate child v_1 and read lock it as well.
 - Release the lock on v_0 .
 - Identify the appropriate child v_2 of v_1 and read lock it as well.
 - Release the lock on v_1 .
 - Repeat this procedure until finding the appropriate leaf vertex.
- This finds the appropriate leaf vertex without forcing a lock on the entire B^+ -tree.
- If the operation is just a read, keep the leaf vertex locked until the operation is complete; then, release the lock.
- If the operation involves a write, proceed as described on the next slide.

Write Operations under Index Traversal via Crabbing

- This process assumes that the appropriate leaf vertex has been identified using the crabbing technique of the previous slide.
- If the leaf vertex to be written has not already been write locked, then upgrade the lock to write.
- If the operation does not require any changes to the index, or any movement of records to a sibling, then complete the operation and release the lock.
- If an operation requires modification of the parent vertex (in the index) and/or a sibling, then these must be write locked as well.
- This process may propagate up the tree if further modifications to the index are necessary.
 - It is therefore important to maintain a list of the vertices visited on the way down the tree.
- All locks are held until the entire operation is complete.
- The general process may involve deadlocks and consequent livelocks, which may be resolved via timestamping and giving preference to older operations.

Index Maintenance under Transaction Concurrency — Linking

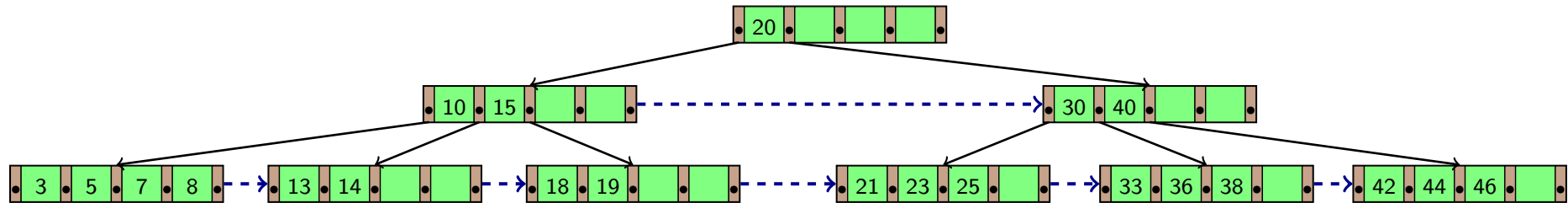


- In a *B-linked* B⁺-tree, all of the index vertices have a special link to the next vertex to the right, on the same level.
 - For reasons of space, only the index vertices are shown above.
 - The leaf (data) vertices always contain two-way links in B⁺-tree.

Insertions (both key and data): Implemented so that when a vertex becomes full, new items are always place in the (possibly new) vertex to the right.

- The link(s) may be followed to find a sought value, without the need to back up to the parent.
- An item which is not found where it “should” be will always be to the right.
- A further access to the parent (which may be locked by another process) is not necessary.

Index Maintenance with Linking — 2



Deletions (both key and data): Unfortunately, if a coalesce operation is involved, a path may become invalid and the query may need to be re-run.

- However, insertions are typically much more common than deletions, so this tends not to be a major problem.
- There are many technical details regarding how concurrency is managed, which may be found in the paper noted at the end of these slides.

Key-Value Locking in B⁺-Trees

- The most straightforward approach to locking data items has the granularity of a page (leaf vertex).

Key-value locking: To achieve greater concurrency, it is possible to lock one record at a time.

- This allows concurrent access to the same leaf vertex, on different records.
- When the lock granularity is the page, phantoms are excluded because all the keys are locked which lie in the range of possible key values for that page, as identified by the index.
- However, when key-value locking is used, phantom tuples may create problems.

Next-key locking: To address this issue, when a key is locked in key-value locking, all (possibly phantom) tuples in the range from the current key up to and including the next physical key are also locked.

- Any concurrent operation on (phantom) tuples in this range is thus blocked.

Classical Reference Books on Concurrency Control

- This classical reference is available online at <http://research.microsoft.com/en-us/people/philbe/ccontrol.aspx>. It contains a detailed presentation of classical MVCC.

Bernstein, P. A., V. Hadzilacos, and N. Goodman, *Concurrency Control and Recovery in Database Systems*, Addison-Wesley, 1987.

- The following is still one of the best references on the basic theory of concurrency control. It is concise and well written.

Papadimitriou, C., *The Theory of Database Concurrency Control*, Computer Science Press, 1986.

Recent Reference Books on Concurrency Control

- This book is very current and presents an application-oriented perspective without going into detailed theory. It is a great book for obtaining the overall picture of transaction processing in the real world.

Philip Bernstein and Eric Newcomer.

Principles of Transaction Processing.

Morgan Kaufmann, second edition, 2009.

- This book is a comprehensive reference on the theory of concurrency control.

Gerhard Weikum and Gottfried Vossen.

Transactional Information Systems.

Morgan Kaufmann, 2002.

Paper on the DSG

- Material on the DSG is based upon ideas and notation from the following paper:

Atul Adya, Barbara Liskov, and Patrick E. O'Neil,
[Generalized isolation level definitions](#),

In David B. Lomet and Gerhard Weikum, editors, *Proceedings of the 16th International Conference on Data Engineering, San Diego, California, USA, February 28 - March 3, 2000*, pp. 67-78, 2000.

Papers on Snapshot Isolation

- The following now-classical paper presents a simple yet formal model of modelling of transaction anomalies. It is the first paper to discuss snapshot isolation from a formal perspective and illustrate write skew. It is available for free download at <http://arxiv.org/abs/cs/0701157>.

Hal Berenson, Philip A. Bernstein, Jim Gray, Jim Melton,
Elizabeth J. O'Neil, and Patrick E. O'Neil.

A critique of ANSI SQL isolation levels.

In Proceedings of the 1995 ACM SIGMOD International Conference on Management of Data, San Jose, California, May 22-25, 1995, pages 1–10, 1995.

Papers on Snapshot Isolation 2

- Augmenting SI to provide serializable isolation.
Alan Fekete, Dimitrios Liarokapis, Elizabeth J. O’Neil, Patrick E. O’Neil, and Dennis Shasha. [Making snapshot isolation serializable](#). *ACM Trans. Database Syst.*, 30(2):492–528, 2005.
- ESSI and more recent developments:
Michael J. Cahill, Uwe Röhm, and Alan David Fekete. [Serializable isolation for snapshot databases](#). *ACM Trans. Database Syst.*, 34(4), 2009.
- PSSI:
Stephen Revilak, Patrick E. O’Neil, and Elizabeth J. O’Neil. [Precisely serializable snapshot isolation \(PSSI\)](#). In *Proceedings of the 27th International Conference on Data Engineering, ICDE 2011, April 11-16, 2011, Hannover, Germany*, pages 482–493, 2011.
- SSI in PostgreSQL:
Dan R. K. Ports and Kevin Grittner. [Serializable snapshot isolation in PostgreSQL](#). *Proc. VLDB Endowment*, 5(12):1850–1861, 2012.

Papers on Snapshot Isolation 3

- The following paper examines SI in the context of classical schedules.

Ragnar Normann and Lene T. Østby.

[A theoretical study of 'snapshot isolation'.](#)

In Luc Segoufin, editor, *Database Theory - ICDT 2010, 13th International Conference, Lausanne, Switzerland, March 23-25, 2010, Proceedings*, pages 44–49, ACM, 2010.

<http://www.edbt.org/Proceedings/2010-Lausanne/icdt/papers/p0044-Normann.pdf>

- The above paper is based upon the following MSc thesis at the University of Oslo.

Lene T. Østby.

[En teoretisk studie av "snapshot isolation".](#)

Masteroppgave, Institutt for informatikk, Universitetet i Oslo, 2008.

<http://www.duo.uio.no/sok/work.html?WORKID=74076>

Paper on Linking in B-trees

- The following paper is the original, main presentation of B-linking.

Philip L. Lehman and S. Bing Yao,

[Efficient Locking for Concurrent Operations on B-trees.](#)

ACM Trans. Database Syst. 6(4):650-670, 1981.

DBTech Resources on Transactions

- DBTech EXT is a consortium, funded by the EU, which develops educational materials in the DBMS area, with a focus upon hands-on use of real systems.
- Their main portal is here: <http://dbtech.uom.gr/>
- Of particular interest is the materials which they have developed for concurrency control and recovery, which may be found by clicking on the appropriate link of the above site.
- These material include not only papers, but also exercises and even a downloadable VBox image which contains Linux with the free versions of both Oracle and DB2 installed.
- Two of the participants, Martti Laiho and Fritz Laux, have written a paper entitled “On SQL Concurrency Technologies for Application Developers”, which covers in detail how real systems handle concurrency control.
- It is available for free downloaded at:
http://www.dbtechnet.org/papers/SQL_ConcurrencyTechnologies.pdf.