Transactions

<□ > < @ > < E > < E > E のQ @

Transactions

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Transaction

Execution of a user program in a DBMS.

Transactions

Transaction

Execution of a user program in a DBMS.

Transaction properties

- Atomicity: all-or-nothing execution
- Consistency: database consistency is preserved
- Isolation: concurrently executing transactions have no effect on one another
- Durability: results survive failures.

Schedules

Schedules

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Transaction (DBMS view)

- list of actions (read or write)
- terminated by commit or abort

Schedules

Transaction (DBMS view)

- list of actions (read or write)
- terminated by commit or abort

Schedule

- interleaving of multiple transactions
- action order within transaction preserved
- complete: commit/abort for every transaction
- serial: no interleaving of actions from different transactions
- serializable: equivalent to a serial schedule (assuming all transactions commit).

Conflicts

(日) (個) (目) (日) (日) (の)

Conflicts

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 善臣 - のへで

Conflict

- a pair of actions of different transactions on the same object
- one action is a write
- a conflict orders the transactions

Conflicts

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Conflict

- a pair of actions of different transactions on the same object
- one action is a write
- a conflict orders the transactions

Conflicts influence serializability

- WR: reading uncommitted data
- RW: unrepeatable reads
- WW: overwriting uncommitted data.

Reading uncommitted data

Reading uncommitted data

<□ > < @ > < E > < E > E のQ @

| T_1 | debit(A,1000), credit(B,1000) |
|-----------------------|---|
| <i>T</i> ₂ | increase A by 10%, increase B by 10% |

Reading uncommitted data

| T_1 | debit(A,1000), credit(B,1000) |
|-------|---|
| T_2 | increase A by 10%, increase B by 10% |

| T_1 | T ₂ |
|--------|----------------|
| R(A) | |
| W(A) | |
| | R(A) |
| | W(A) |
| | R(B) |
| | W(B) |
| | Commit |
| R(B) | |
| W(B) | |
| Commit | |

Unrepeatable read

Unrepeatable read

| <i>T</i> ₃ | credit(A,1000) |
|-----------------------|----------------|
| <i>T</i> ₄ | credit(A,2000) |

Unrepeatable read

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

| <i>T</i> ₃ | credit(A,1000) |
|-----------------------|----------------|
| <i>T</i> ₄ | credit(A,2000) |

| <i>T</i> ₃ | T ₄ |
|-----------------------|----------------|
| R(A) | |
| | R(A) |
| | W(A) |
| | Commit |
| W(A) | |
| Commit | |

Overwriting uncommitted data

<□ > < @ > < E > < E > E のQ @

Overwriting uncommitted data

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

 T_5 book(F1,AA), book(F2,AA)

 T_6 book(F1,Delta), book(F2,Delta)

Overwriting uncommitted data

| <i>T</i> ₅ | book(F1,AA), book(F2,AA) |
|-----------------------|--------------------------------|
| T_6 | book(F1,Delta), book(F2,Delta) |

| T_5 | T_6 |
|--------|--------|
| W(F1) | |
| | W(F1) |
| | W(F2) |
| | Commit |
| W(F2) | |
| Commit | |

(4日) (個) (目) (目) (目) (の)

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

The effect of aborted transactions has to be completely undone.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

The effect of aborted transactions has to be completely undone.

Problems

The effect of aborted transactions has to be completely undone.

Problems

 a transaction depending on an aborted transaction may have already committed (unrecoverable schedule)

The effect of aborted transactions has to be completely undone.

Problems

- a transaction depending on an aborted transaction may have already committed (unrecoverable schedule)
- aborting a transaction requires aborting other transactions (cascading aborts).

Unrecoverable schedule

▲□▶ ▲圖▶ ▲≧▶ ▲≣▶ = 目 - のへで

Unrecoverable schedule

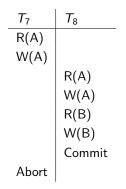
▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ = ● ● ●

| <i>T</i> ₇ | debit(A,100) |
|-----------------------|---|
| <i>T</i> ₈ | increase A by 10%, increase B by 10% |

Unrecoverable schedule

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ ―臣 … のへで

| <i>T</i> ₇ | debit(A,100) |
|-----------------------|---|
| <i>T</i> ₈ | increase A by 10%, increase B by 10% |



Strict two-phase locking

▲□▶ <圖▶ < ≧▶ < ≧▶ = のQ@</p>

Strict two-phase locking

Rules

- before an object is accessed, an appropriate lock on the object(read: shared mode, write: exclusive mode) needs to be obtained
- 2 lock in exclusive mode: no other transaction can lock the object in any mode
- Olock in shared mode: other transactions can lock the object in shared mode
- 4 a transaction cannot lock an object more than once
- 5 all the locks are held until the end of transaction.

Strict two-phase locking

▲口 → ▲圖 → ▲ 臣 → ▲ 臣 → □ 臣 □

Rules

- before an object is accessed, an appropriate lock on the object(read: shared mode, write: exclusive mode) needs to be obtained
- 2 lock in exclusive mode: no other transaction can lock the object in any mode
- Olds in shared mode: other transactions can lock the object in shared mode
- 4 a transaction cannot lock an object more than once
- 5 all the locks are held until the end of transaction.

Guarantees

- schedule serializability
- schedule recoverability
- no cascading aborts

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Locks are stored in a lock table (managed by DBMS), lock requests are queued.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Locks are stored in a lock table (managed by DBMS), lock requests are queued.

Lock/unlock: atomic operations.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

Locks are stored in a lock table (managed by DBMS), lock requests are queued.

Lock/unlock: atomic operations.

Problems

- deadlocks
- starvation.

Deadlocks

Deadlocks

Deadlock

A set of transactions such that each waits for a lock held by another one.

Deadlocks

▲ロト ▲帰ト ▲ヨト ▲ヨト - ヨ - の々ぐ

Deadlock

A set of transactions such that each waits for a lock held by another one.

Handling deadlocks

- prevention:
 - object ordering
 - transaction priorities
 - obtaining all the locks at the beginning
- detection:
 - · identifying cycles in the waits-for graph or timeout, and
 - abort transaction.

Deadlocks

Deadlock

A set of transactions such that each waits for a lock held by another one.

Handling deadlocks

- prevention:
 - object ordering
 - transaction priorities
 - obtaining all the locks at the beginning
- detection:
 - · identifying cycles in the waits-for graph or timeout, and
 - abort transaction.

Handling starvation

FIFO lock queues.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Types of failures

- transaction abort
- system crash
- media failure

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Types of failures

- transaction abort
- system crash
- media failure

Memory levels

- disk blocks
- main memory buffers
- local variables
- the same object may have a copy at each level

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Types of failures

- transaction abort
- system crash
- media failure

Memory levels

- disk blocks
- main memory buffers
- local variables
- the same object may have a copy at each level

▲□▶ ▲圖▶ ★園▶ ★園▶ - 園 - のへで

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Operations

• INPUT(X): Copy the disk block containing the object X to a buffer.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Operations

- INPUT(X): Copy the disk block containing the object X to a buffer.
- READ(X,A): Copy the object X to a local variable A (preceded by INPUT(X) if necessary).

Operations

- INPUT(X): Copy the disk block containing the object X to a buffer.
- READ(X,A): Copy the object X to a local variable A (preceded by INPUT(X) if necessary).
- WRITE(X,A): Copy the value of the local variable A to the object X (preceded by INPUT(X) if necessary).

Operations

- INPUT(X): Copy the disk block containing the object X to a buffer.
- READ(X,A): Copy the object X to a local variable A (preceded by INPUT(X) if necessary).
- WRITE(X,A): Copy the value of the local variable A to the object X (preceded by INPUT(X) if necessary).
- OUTPUT(X): Copy the block containing X from buffer to disk.

Operations

- INPUT(X): Copy the disk block containing the object X to a buffer.
- READ(X,A): Copy the object X to a local variable A (preceded by INPUT(X) if necessary).
- WRITE(X,A): Copy the value of the local variable A to the object X (preceded by INPUT(X) if necessary).
- OUTPUT(X): Copy the block containing X from buffer to disk.

Assumption: each object fits into one block.

Logging

Logging

Recording all the operations in an append-only log (also stored on disk).

Logging

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 三臣 - のへで

Recording all the operations in an append-only log (also stored on disk).

Log records

- <START T>
- <COMMIT T>
- <ABORT T>
- <T,X,old,new>

UNDO/REDO logging

(4日) (個) (目) (目) (目) (の)()

UNDO/REDO logging

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

UNDO/REDO rule

Before modifying an object X on disk on behalf of a transaction T, a log update record <T,X,old,new> needs to be written to disk.

$\mathsf{UNDO}/\mathsf{REDO}\ \mathsf{logging}$

UNDO/REDO rule

Before modifying an object X on disk on behalf of a transaction T, a log update record <T,X,old,new> needs to be written to disk.

Recovery

Redo all the committed transactions in the order earliest-first.
Undo all the incomplete transactions in the order latest-first.

$\mathsf{UNDO}/\mathsf{REDO}\ \mathsf{logging}$

UNDO/REDO rule

Before modifying an object X on disk on behalf of a transaction T, a log update record <T,X,old,new> needs to be written to disk.

Recovery

Redo all the committed transactions in the order earliest-first.
Undo all the incomplete transactions in the order latest-first.

Checkpointing

- Write <START CKPT (T1,...Tk)> log record, where T1,...Tk are all the active transactions, and flush the log.
- 2 Flush all dirty buffers.
- 3 Write <END CKPT> log record, and flush the log.

Distributed transactions

▲□▶ ▲圖▶ ▲国▶ ▲国▶ - 国 - のへで

Distributed transactions

▲□▶ ▲圖▶ ★ 国▶ ★ 国▶ - 国 - のへで

Transactions

- subtransactions executing at different sites
- all subtransactions commit or none does (commit protocol)
- site and link failures.

Distributed transactions

Transactions

- subtransactions executing at different sites
- all subtransactions commit or none does (commit protocol)
- site and link failures.

Two-phase commit

A site is designated as a coordinator, other participating sites are subordinates.

▲□▶ ▲圖▶ ▲≣▶ ▲≣▶ = = の�?

1 Coordinator: send a PREPARE message to each subordinate

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

1 *Coordinator:* send a PREPARE message to each subordinate

2 Subordinate: receive PREPARE and decide to commit or abort:

- **1** *Coordinator:* send a PREPARE message to each subordinate
- 2 Subordinate: receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

- **1** *Coordinator:* send a PREPARE message to each subordinate
- *Subordinate:* receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- **1** *Coordinator:* send a PREPARE message to each subordinate
- *Subordinate:* receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.
- 3 Coordinator:

- **1** *Coordinator:* send a PREPARE message to each subordinate
- 2 Subordinate: receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.
- 3 Coordinator:
 - all subordinates reply YES: write a commit log record, flush log, send a COMMIT message to each subordinate;

- **1** *Coordinator:* send a PREPARE message to each subordinate
- 2 Subordinate: receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.
- 3 Coordinator:
 - all subordinates reply YES: write a commit log record, flush log, send a COMMIT message to each subordinate;
 - one replies NO or times out: write an abort log record, flush log, send an ABORT message to each subordinate.

- **1** *Coordinator:* send a PREPARE message to each subordinate
- 2 Subordinate: receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.
- 8 Coordinator:
 - all subordinates reply YES: write a commit log record, flush log, send a COMMIT message to each subordinate;
 - one replies NO or times out: write an abort log record, flush log, send an ABORT message to each subordinate.

4 Subordinate:

- **1** *Coordinator:* send a PREPARE message to each subordinate
- 2 Subordinate: receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.
- 3 Coordinator:
 - all subordinates reply YES: write a commit log record, flush log, send a COMMIT message to each subordinate;
 - one replies NO or times out: write an abort log record, flush log, send an ABORT message to each subordinate.

4 Subordinate:

 receive COMMIT: write a commit log record, flush log, send ACK to coordinator, commit;

- **1** *Coordinator:* send a PREPARE message to each subordinate
- 2 Subordinate: receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.
- 3 Coordinator:
 - all subordinates reply YES: write a commit log record, flush log, send a COMMIT message to each subordinate;
 - one replies NO or times out: write an abort log record, flush log, send an ABORT message to each subordinate.

4 Subordinate:

- receive COMMIT: write a commit log record, flush log, send ACK to coordinator, commit;
- receive ABORT: write an abort log record, flush log, send ACK, abort.

- **1** *Coordinator:* send a PREPARE message to each subordinate
- 2 Subordinate: receive PREPARE and decide to commit or abort:
 - commit: write a prepare log record, flush log, reply YES;
 - abort: write an abort log record, flush log, reply NO.
- 3 Coordinator:
 - all subordinates reply YES: write a commit log record, flush log, send a COMMIT message to each subordinate;
 - one replies NO or times out: write an abort log record, flush log, send an ABORT message to each subordinate.

4 Subordinate:

- receive COMMIT: write a commit log record, flush log, send ACK to coordinator, commit;
- receive ABORT: write an abort log record, flush log, send ACK, abort.
- 6 Coordinator: receive ACK from all subordinates: write end log record.