# DBMS IMPLEMENTATION CONCURRENCY CONTROL

### Transaction

- transaction is the basic unit of work in a DBMS
- 4 basic or acid properties of a transaction:

atomicity

consistency

independency

durability

Concurrency Ctrl -Section 6

© Vincent P. Wade

# **DB Processing**

• Execution of application processing consists of a series of atomic transactions with non-DB processing taking place in between

Figure 6.1 - Program executing several transactions

Concurrency Ctrl -Section 6

© Vincent P. Wade

## **Concurrency Control Properties**

atomicity

the all or nothing property; a transactions is an indivisible unit of work

consistency

transactions transform the DB from one consistent state to another consistent state

independence

transactions execute independently of one another i.e. partial effects of one transaction are not visible to other transactions

durability ( also called persistence)

the effects of a successfully completed (committed) transaction are permanently recorded in the DB and cannot be undone

Concurrency Ctrl -Section 6

© Vincent P. Wade

2

# DB Processing (contd.)

- transaction manager oversees execution of a transactions and co-ordinates DB requests on behalf of the transaction
- scheduler implements a particular strategy for transaction execution
- objective of the scheduler is to maximise concurrency without allowing transactions to interfere with one another and compromise the integrity of the DB

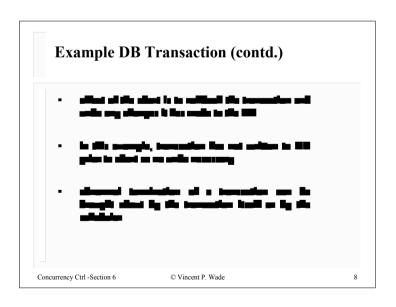
Concurrency Ctrl -Section 6

© Vincent P. Wade

# Transactions can be interleaved in 2 ways: • end-to-end txn execution Figure 6.2 • concurrent txn execution Figure 6.3 Concurrency Ctrl - Section 6 © Vincent P. Wade 5

# Example DB Transaction Figure 6.4 Funds transfer transaction Concurrency Ctrl - Section 6 © Vincent P. Wade 7

# DB Processing (contd.) • start of a transactions signalled by begin transaction • end of a transaction signalled by either commit (= successful termination) or abort (= unsuccessful termination) (rollback) Concurrency Ctrl - Section 6 ♥ Vincent P. Wade 6



## **Transaction Interference**

There are three different ways in which concurrently executing transactions can interfere with one another:

- lost update problem
- violation of integrity constraints
- · inconsistent retrieval

Concurrency Ctrl -Section 6

Concurrency Ctrl -Section 6

C Vincent P. Wade

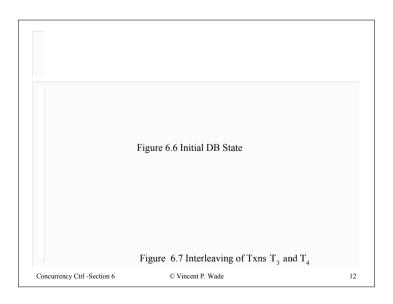
9

11

# • Apparently successful updates can be overwritten by other transactions Figure 6.5 Lost Update Problem Concurrency Ctrl-Section 6 © Vincent P. Wade 10

# Violation of Integrity Constraints Example DB: SCHEDULE (surgeon\_name, operation, date) SURGEON (surgeon\_name, operation) SCHEDULE specifies which surgeon is to perform a particular operation on a particular day. SURGEON records the qualifications by operation for eachoperation integrity constraint: surgeons must be qualified to perform the operations for which they are scheduled Figure 6.6 shows the initial state of the DB, Figure 6.7 shows the interleaving of 2 transactions, T<sub>3</sub> and T<sub>4</sub>

© Vincent P. Wade



### **Violation of Integrity Constraints (Contd.)**

• Figure 6.8 shows the invalid state of the DB after execution of  $T_3$  and  $T_4$ .

### Figure 6.8

Note: Neither transaction is aware of the action of the other transaction as they are both updating different data

Concurrency Ctrl -Section 6

© Vincent P. Wade

13

# Example: Inconsistent Retrieval Problem Figure 6.9 Concurrency Ctrl - Section 6 © Vincent P. Wade 15

### **Inconsistent retrieval (dirty read)**

- Most Concurrency Control work concentrates on transactions which update DB since only they can corrupt the DB
- If transactions are allowed to read the partial results of incomplete transactions, they can obtain an inconsistent view of the DB (dirty or unrepeatable read)
- Example given in Figure 6.9

Concurrency Ctrl -Section 6

© Vincent P. Wade

14

## **Schedules and serialisation**

- A transaction consists of a sequence of reads and writes to the database
- The entire sequence of reads and writes by all concurrent transactions in a database taken together is known as a *schedule*.

Concurrency Ctrl -Section 6

© Vincent P. Wade

### Schedules and serialisation contd.

• a schedule *S* is generally written:

$$S = [O^1, O^2, O^3, \dots, O^m]$$

where O<sup>i</sup> indicates either a read (R) or write (W) operation executed by a transaction on a data item

 O<sup>1</sup> precedes O<sup>2</sup>, which in turn precedes O<sup>3</sup>, and so on. This is generally denoted

$$O^1 < O^2 < O^3 < \dots < O^m$$

Concurrency Ctrl -Section 6

© Vincent P. Wade

17

19

### Schedules and serialisation contd

Most concurrency control algorithms assume that transactions read a data item *before* they update it (*constrained write rule*) i.e.

$$R_{i}(x_{j}) < W_{i}(x_{j}).$$

 The order of interleaving of operations from different transactions is crucial to maintaining the consistency of the database

Concurrency Ctrl -Section 6

© Vincent P. Wade

### Schedules and serialisation contd

• The schedule S for transactions T<sub>1</sub> and T<sub>2</sub> in Figure 6.5 would be

$$S = [R_2(balance_x), R_1(balance_x), W_2(balance_x), W_1(balance_x), R_1(balance_y), W_1(balance_y)]$$

where  $R_i$  and  $W_i$  denote read and write operations, respectively, by transaction  $T_i$ .

Concurrency Ctrl -Section 6

© Vincent P. Wade

18

## **Serial schedule**

A *serial* schedule is one in which all the reads and writes of each transaction are grouped together so that the transactions are run sequentially one after the other, as in Figure 6.2(a)

Schedule, S, is said to be serialisable if all the reads and writes
of each transaction can be reordered in such a way that when
they are grouped together as in a serial schedule, the net effect
of executing this serial schedule is the same as that of the
original schedule S

Concurrency Ctrl -Section 6

© Vincent P. Wade

### Serial schedule contd.

- This reorganised schedule is called the *equivalent serial* schedule of the original serialisable schedule
- A serialisable schedule will therefore be *equivalent* to and have the same effect on the DB as *some* serial schedule

Example

$$S_1 = [R_7(x), R_8(x), W_8(x), R_6(y), W_6(y), R_7(y), W_7(y)]$$

Concurrency Ctrl -Section 6

© Vincent P. Wade

21

23

### Serial schedule contd.

• schedule S<sub>1</sub> for Figure 6.10 is

$$S_1 = [R_7(x), R_8(x), W_8(x), R_6(y), W_6(y), R_7(y), W_7(y)]$$

and the equivalent serial schedule,  $SR_1$  is

$$SR_1 = [R_6(y), W_6(y), R_7(x), R_7(y), W_7(y), R_8(x), W_8(x)]$$

Concurrency Ctrl -Section 6

© Vincent P. Wade

### Serial schedule contd.

- chronologically, by start-time, the order of execution of transactions is T<sub>7</sub>, then T<sub>8</sub> and followed by T<sub>6</sub>
- logically however T<sub>6</sub> precedes T<sub>7</sub> which in turn precedes T<sub>8</sub>, since T<sub>6</sub> reads the pre-T<sub>7</sub> value of y, while T<sub>7</sub> sees the pre-T<sub>8</sub> value of x,

i.e. 
$$T_6 < T_7 < T_{8}$$

in spite of the fact that chronologically  $\mathrm{T}_8$  finishes before  $\mathrm{T}_6$  begins!

Concurrency Ctrl -Section 6

C Vincent P. Wade

22

### Serial schedule contd.

- note that a serialisable schedule is not the same as a serial schedule
- Serialisability is taken as proof of correctness since a serial schedule cannot contain transactions which interfere with each other and a serialisable schedule is a schedule which has the same net effect as executing a serial schedule
  - => one possible approach is to examine the schedule produced and see if it is serialisable

Concurrency Ctrl -Section 6

© Vincent P. Wade

### Serial schedule contd.

- However, deciding whether a schedule is equivalent to some serial schedule is a very difficult computational problem (with constrained write rule: polynomial complexity; without: NPcomplete problem)
- better to design schedulers in such a way that they are guaranteed to generate only serialisable and hence correct schedules

Concurrency Ctrl -Section 6

© Vincent P. Wade

25

27

### **Conflicting operations**

 Read operations cannot conflict with one another and the order of execution of R<sub>1</sub>(x) and R<sub>2</sub>(x) does not matter, i.e.

$$[R_1(x), R_2(x)]$$
 o  $[R_2(x), R_1(x)]$ 

but

$$[R_1(x), W_1(x), R_2(x)]^{-1} [R_1(x), R_2(x), W_1(x)]$$

Concurrency Ctrl -Section 6

© Vincent P. Wade

### Serial schedule contd.

Formally, we can state the rules for equivalence of schedules as:

Rule 1: Each read operation reads the same values in both schedules; this effectively means that those values must have been produced by the same write operations in both schedules

Rule 2: The final database state is the same for both schedules; thus the final write operation on each data item is the same in both schedules.

Concurrency Ctrl -Section 6

© Vincent P. Wade

# **Conflicting operations contd.**

- In terms then of schedule equivalence, it is the ordering of conflicting operations which must be the same in both schedules
- The conflict between a read and a write operation is called a read-write conflict, and a conflict between two write operations a write-write conflict.

Concurrency Ctrl -Section 6

© Vincent P. Wade

28

### **CONCURRENCY CONTROL TECHNIQUES**

Three basic concurrency control techniques:

- locking methods (conservative)
- timestamp methods (conservative)
- optimistic methods

Concurrency Ctrl -Section 6

© Vincent P. Wade

29

31

# Locking methods contd.

- Only one transaction at a time can hold a write lock on a data item
- Transaction holds a lock until it explicitly releases it
- Effects of the write operation are not visible to other transactions until the write lock has been released

Concurrency Ctrl -Section 6

© Vincent P. Wade

## **Locking methods**

- Most widely used approach
- A transaction must claim a *read (shared)* or *write (exclusive)* lock on a data item prior to the execution of the corresponding read or write operation on that data item
- More than one transaction can hold read locks simultaneously on the same data item

Concurrency Ctrl -Section 6

© Vincent P. Wade

30

# Locking methods contd.

- in Figure 6.10, the request by transaction T1 for a write lock on balance<sub>x</sub> would be denied since  $T_5$  already holds a read lock on balance<sub>x</sub>.
- some systems allow upgrading and downgrading of locks

i.e. a read lock on a data item can be *upgraded* to a write lock, if it is the *only* transaction holding a read lock on that data item; similarly a write-lock can be *downgraded* to a read lock (potentially allows greater concurrency)

Concurrency Ctrl -Section 6

© Vincent P. Wade

### Granularity of locks

- promis sin in halling now may fine a single topic, in the sales had no
- ألمان الدولية المنافع المنافعة الم

Concurrency Ctrl -Section 6

© Vincent P. Wade

33

# 2 Phase Locking

- n most common locking protocol is known as *two-phase locking* (2PL)
- n transactions operate in 2 distinct phases:
  - a growing phase during which the transaction acquires locks and a
  - shrinking phase during which it releases those locks

Concurrency Ctrl -Section 6

© Vincent P. Wade

2.4

# Rules for transactions obeying 2PL

- transactions are well-formed, thus a transaction must acquire a lock on a data object before operating on it and all locks held by a transaction must be released when the transaction is finished
- compatibility rules for locking are observed, thus no conflicting locks are held (write-write and read-write conflicts are forbidden)
- 3. once the transaction has released a lock, no new locks are acquired
- 4. all write locks are released together when the transaction commits (to ensure atomicity).

Concurrency Ctrl -Section 6

© Vincent P. Wade

35

# Rules for transactions obeying 2PL (cont)

### Note:

- n upgrading and downgrading of locks are possible under 2PL, with the restriction that downgrading is only permitted during the shrinking phase.
- n it can be proved that the schedules produced by transactions which obey 2PL are guaranteed to be serialisable.

Concurrency Ctrl -Section 6

© Vincent P. Wade

# Deadlock

n Under 2PL, if a transaction T<sub>1</sub> requests a write-lock on data item x<sub>i</sub>, which is currently locked by another transaction, T<sub>2</sub>, there are two possibilities:

Approach 1:

Place  $T_1$  on a queue for  $x_i$  awaiting release of lock by  $T_2$  Approach 2:

Abort and rollback T<sub>1</sub>

Concurrency Ctrl -Section 6

© Vincent P. Wade

37

39

# Deadlock (cont 3)

- n With Approach 2, T<sub>1</sub> must release all its locks and restart (this could involve high overhead)
- n Approach 2 is referred to as a **deadlock prevention protocol** because transaction releases all its locks on becoming blocked, so deadlock cannot occur

Concurrency Ctrl -Section 6

© Vincent P. Wade

Deadlock (cont)

- n with Approach 1, T<sub>1</sub> retains all the locks it currently holds and just enters a wait state
- n Approach 1 can lead to deadlock

deadlock occurs when one transaction is waiting for a lock to be released by a second transaction, which is in turn waiting for a lock currently held by the first transaction

n Therefore, in approach 1, a **deadlock detection protocol** is needed

Concurrency Ctrl -Section 6

© Vincent P. Wade

38

# LiveLock

- n transactions can be repeatedly rolled back or left in a wait state indefinitely, unable to acquire their locks, even though the system is not deadlocked => livelock, e.g. summary transactions
- n priority system is needed to avoid livelock, whereby the longer a transaction has to wait, the higher its priority.

Concurrency Ctrl -Section 6

© Vincent P. Wade

# **Deadlock Detection**

- n deadlock is generally detected by means of wait-for graphs.
- n transactions are represented by nodes and blocked requests for locks are represented by labelled, directed edges (see Figure 6.11)
- n deadlock is represented in the graph by the existence cycle between  $T_1$  and  $T_2$
- n informally, a cycle in a graph forms a closed loop by which it is possible to start at one node, traverse the edges of the graph according to the directions of the arrows and get back to where you started

Concurrency Ctrl -Section 6

© Vincent P. Wade

41

# Figure 6.11 Concurrency Ctrl -Section 6 © Vincent P. Wade 42

# **DeadLock Detection (cont)**

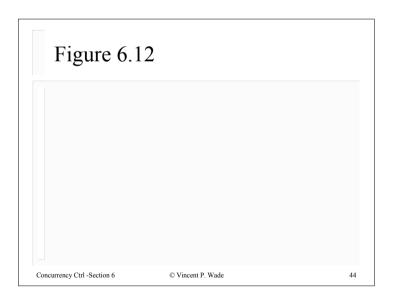
n wait-for graph, G, can also be represented symbolically as:

$$G = T_1 \otimes T_2 \otimes T_1$$

- n where  $T_1 \otimes T_2$  indicates that transaction  $T_1$  is waiting for transaction  $T_2$ ; the  $\otimes$  represents the wait-for relationship)
- n deadlock can occur between two transactions indirectly via a chain of intermediate transactions as shown in Figure 6.12

Concurrency Ctrl -Section 6

© Vincent P. Wade



# Deadlock Detection Graph

$$G' = T_1 \otimes T_2 \otimes T_3 \otimes T_4 \otimes T_5 \otimes T_1$$

n once detected, deadlock must be resolved by preemptying, i.e. aborting and rolling back, one of the transactions (any one will do since this will always break the cycle)

Concurrency Ctrl -Section 6

© Vincent P. Wade

45

Wait Die - Fig 6.13

### **BEGIN**

T1 requests lock on data item currently held by T2 if T1 is older than T2 (I.e. ts(T1)<ts(T2)) then T1 waits for T2 to commit or rollback else T1 is rolled back endif

**END** 

Concurrency Ctrl -Section 6

© Vincent P. Wade

# **Deadlock Prevention**

- n how does the transaction manager decide whether or not to allow a transaction  $T_1$ , which has requested a lock on data item  $x_i$  currently held by transaction  $T_2$ , to wait and to *guarantee* that this waiting *cannot* give rise to deadlock?
- n can force locks to be acquired in a certain datadependent order (difficult to implement)
- n alternatively, transactions can be ordered and ensure that all conflicting operations are executed in sequence according to this order;

Concurrency Ctrl -Section 6

© Vincent P. Wade

46

# **Deadlock Prevention (Cont)**

- n deadlock is prevented by only allowing blocked transactions to wait under certain circumstances which will maintain this ordering
- n ordering mechanism is generally based on **timestamps** (unique identifier assigned to transaction when it is launched)
- n can ensure that *either* older transactions wait for younger ones (Wait-die) *or* vice versa (Woundwait)
- n figure 6.13 gives the algorithm for wait-die, while figure 6.14 gives the algorithm for wound-wait.

Concurrency Ctrl -Section 6

© Vincent P. Wade

47

# Wound Wait Fig 6.14

### **BEGIN**

T1 requests lock on data item currently held by T2 if T1 is older than T2 (I.e. ts(T1)<ts(T2)) then T2 is rolled back else T1 waits for T2 to commit or rollback endif END

Concurrency Ctrl -Section 6

© Vincent P. Wade

49

### 50

# The deadlock spectrum

- n at one end of this spectrum, techniques are deadlock-free, i.e. deadlocks can never occur (require no run-time support)
- n at the other end, techniques detect and recover from deadlock
- n deadlock prevention methods lie in the middle,

Concurrency Ctrl -Section 6

© Vincent P Wade

51

## **Deadlock Prevention (Cont)**

- n if a transaction is rolled back, it retains its original timestamp (otherwise it could be repeatedly rolled back)
- n wait-die and wound-wait use locks as the primary concurrency control mechanism and are therefore classified as *lock-based* rather than timestamp

Concurrency Ctrl -Section 6

© Vincent P. Wade

## TIMESTAMP METHODS

- n no locks are involved and there can therefore be no deadlock.
- n no waiting; transactions involved in conflict are simply rolled back and restarted.
- n transactions are ordered globally in such a way that older transactions, transactions with *smaller* timestamps, get priority in the event of conflict
- n read or write only be allowed if the last update on that data item was carried out by an older transaction; otherwise the requesting transaction is restarted *and given a new timestamp*.

Concurrency Ctrl -Section 6

© Vincent P. Wade

# **Timestamps (cont)**

### n Thus:

Timestamp methods produce serialisable schedules, which are equivalent to the serial schedule defined by the timestamps of successfully committed transactions.

- n each transaction is assigned a unique timestamp, when it is launched
- n timestamps can use the system clock, a global counter or sequence number generator

Concurrency Ctrl -Section 6

© Vincent P. Wade

53

55

# **Basic timestamping**

n To implement basic timestamping, the following variables are required:

for each data item x

 $ts(read\ x)$  = the timestamp of the transaction which last read data item x and

ts(write x) = the timestamp of the transaction which last updated data item x

and for each transaction  $T_i$ 

 $ts(T_i)$  = the timestamp assigned to transaction  $T_i$  when it is launched

- transactions actually issue pre-writes to buffers rather than writes to the DB
- physical writes to the DB are only performed at commit

Concurrency Ctrl -Section 6 © Vincent P. Wade

## **Atomicity of transactions with timestamps**

- n with timestamp protocols we must prevent other transactions from seeing partial updates since there are no locks
- n done by using pre-writes (deferred update)
- n updates of uncommitted transactions are not written out to the database, but instead are written to a set of buffers, which are only flushed out to the DB when the transaction commits.
- n has the advantage that when a transaction is aborted and restarted, no physical changes need to be made to the DB

Concurrency Ctrl -Section 6

© Vincent P. Wade

5.4

# Pre-Write Algorithm Fig 6.18

### **BEGIN**

Ti attempts to pre-write data item x

if x has been read or written to by a younger transaction (I.e. ts(Ti)<ts(read x) or ts(Ti)<ts(write x)

then reject Ti and restart as Ti

else accept pre-write: buffer (pre) write together with ts(Ti)

end if

**END** 

Concurrency Ctrl -Section 6

© Vincent P. Wade

# Write Algorithm Fig 6.19

### **BEGIN**

Ti attempts to update (I.e. write) data item x
if there is an update pending on x by an older
transaction Tj (I.e. ts(Tj) < ts(Ti)
then Ti waits until Tj is committed or restarted
else Ti committs update and sets ts(write x) = ts(Ti)
end if
END

Concurrency Ctrl -Section 6

© Vincent P. Wade

57

# **Timestamping (cont)**

- n Figures 6.18, 6.19 and 6.20 illustrate the algorithms for the pre-write, write and read operations, respectively, under basic timestamping
- n Timestamping method is equivalent to applying exclusive locks on the data items between the prewrite and write operations
- n Since all waits consist of younger transactions waiting for older transactions "deadlock" is not possible.

Concurrency Ctrl -Section 6

© Vincent P. Wade

59

# Read Operation using Basic Timestamping Fig 6.20

```
BEGIN
```

```
Ti attempts a read operation on data item x

if x has been updated by a younger transaction (I.e.

ts(Ti)<ts(write x)

then reject read operation and restart Ti

else if there is an update pending on x by an older
transaction Tj (I.e. ts(Tj)<ts(Ti))

then Ti waits for Tj to commit or restart
else accept read operation and set
ts(read x) = max(ts(read x), ts(Ti))

endif
enf if

cend
Concurrency Ctrl-Section 6

© Vincent P. Wade

58
```