

Transactions

Distributed Software Systems

Transactions

✍ Motivation

- ✍ Provide atomic operations at servers that maintain shared data for clients
- ✍ Provide recoverability from server crashes

✍ Properties

- ✍ Atomicity, Consistency, Isolation, Durability (ACID)

✍ Concepts: commit, abort

Operations of the *Account* interface

deposit(amount)
deposit amount in the account
withdraw(amount)
withdraw amount from the account
getBalance() -> *amount*
return the balance of the account
setBalance(amount)
set the balance of the account to amount

Operations of the Branch interface

create(name) -> *account*
create a new account with a given name
lookUp(name) -> *account*
return a reference to the account with the given
name
branchTotal() -> *amount*
return the total of all the balances at the branch

Transactions 3

A client's banking transaction

Transaction T:
a.withdraw(100);
b.deposit(100);
c.withdraw(200);
b.deposit(200);

Transactions 4

Operations in Coordinator interface

openTransaction() -> *trans*;

starts a new transaction and delivers a unique TID *trans*. This identifier will be used in the other operations in the transaction.

closeTransaction(trans) -> (*commit*, *abort*);

ends a transaction: a *commit* return value indicates that the transaction has committed; an *abort* return value indicates that it has aborted.

abortTransaction(trans);

aborts the transaction.

Transaction life histories

<i>Successful</i>	<i>Aborted by client</i>		<i>Aborted by server</i>
<i>openTransaction</i>	<i>openTransaction</i>		<i>openTransaction</i>
<i>operation</i>	<i>operation</i>		<i>operation</i>
<i>operation</i>	<i>operation</i>		<i>operation</i>
⋮	⋮	server aborts transaction	⋮
<i>operation</i>	<i>operation</i>	→	<i>operation ERROR reported to client</i>
<i>closeTransaction</i>	<i>abortTransaction</i>		

Concurrency control

- ✦ Motivation: without concurrency control, we have lost updates, inconsistent retrievals, dirty reads, etc. (see following slides)
- ✦ Concurrency control schemes are designed to allow two or more transactions to be executed correctly while maintaining serial equivalence
 - ✦ Serial Equivalence is correctness criterion
 - ✦ Schedule produced by concurrency control scheme should be equivalent to a serial schedule in which transactions are executed one after the other
- ✦ Schemes: locking, optimistic concurrency control, time-stamp based concurrency control

Transactions 7

The lost update problem

Transaction T:	Transaction U:
<code>balance = b.getBalance();</code> <code>b.setBalance(balance*1.1);</code> <code>a.withdraw(balance/10)</code>	<code>balance = b.getBalance();</code> <code>b.setBalance(balance*1.1);</code> <code>c.withdraw(balance/10)</code>
<code>balance = b.getBalance();</code> \$200 <code>b.setBalance(balance*1.1);</code> \$220 <code>a.withdraw(balance/10)</code> \$80	<code>balance = b.getBalance();</code> \$200 <code>b.setBalance(balance*1.1);</code> \$220 <code>c.withdraw(balance/10)</code> \$280

Transactions 8

The inconsistent retrievals problem

Transaction V		Transaction W:	
<i>a.withdraw(100)</i>		<i>aBranch.branchTotal()</i>	
<i>b.deposit(100)</i>			
<i>a.withdraw(100);</i>	\$100	<i>total = a.getBalance()</i>	\$100
		<i>total = total+b.getBalance()</i>	\$300
		<i>total = total+c.getBalance()</i>	
<i>b.deposit(100)</i>	\$300	⋮	

Transactions 9

A serially equivalent interleaving of T and U

Transaction T		Transaction U	
<i>balance = b.getBalance()</i>		<i>balance = b.getBalance()</i>	
<i>b.setBalance(balance*1.1)</i>		<i>b.setBalance(balance*1.1)</i>	
<i>a.withdraw(balance/10)</i>		<i>c.withdraw(balance/10)</i>	
<i>balance = b.getBalance()</i>	\$200	<i>balance = b.getBalance()</i>	\$220
<i>b.setBalance(balance*1.1)</i>	\$220	<i>b.setBalance(balance*1.1)</i>	\$242
<i>a.withdraw(balance/10)</i>	\$80	<i>c.withdraw(balance/10)</i>	\$278

Transactions 10

A serially equivalent interleaving of V and W

TransactionV:		TransactionW:
<i>a.withdraw(100);</i>		<i>aBranch.branchTotal()</i>
<i>b.deposit(100)</i>		
<i>a.withdraw(100);</i>	\$100	
<i>b.deposit(100)</i>	\$300	
		<i>total = a.getBalance()</i> \$100
		<i>total = total+b.getBalance()</i> \$400
		<i>total = total+c.getBalance()</i>
		...

Transactions 11

A dirty read when transaction T aborts

TransactionT:	TransactionU:
<i>a.getBalance()</i>	<i>a.getBalance()</i>
<i>a.setBalance(balance + 10)</i>	<i>a.setBalance(balance + 20)</i>
<i>balance = a.getBalance()</i> \$100	
<i>a.setBalance(balance + 10)</i> \$110	
	<i>balance = a.getBalance()</i> \$110
	<i>a.setBalance(balance + 20)</i> \$130
	<i>commit transaction</i>
<i>abort transaction</i>	

Transactions 12

Serializability

BEGIN_TRANSACTION
 x = 0;
 x = x + 1;
 END_TRANSACTION

(a)

BEGIN_TRANSACTION
 x = 0;
 x = x + 2;
 END_TRANSACTION

(b)

BEGIN_TRANSACTION
 x = 0;
 x = x + 3;
 END_TRANSACTION

(c)

Schedule 1	x = 0; x = x + 1; x = 0; x = x + 2; x = 0; x = x + 3	Legal
Schedule 2	x = 0; x = 0; x = x + 1; x = x + 2; x = 0; x = x + 3;	Legal
Schedule 3	x = 0; x = 0; x = x + 1; x = 0; x = x + 2; x = x + 3;	Illegal

(d)

- a) – c) Three transactions T_1 , T_2 , and T_3
 d) Possible schedules

Transactions 13

Read and write operation conflict rules

<i>Operations of different transactions</i>	<i>Conflict</i>	<i>Reason</i>
<i>read read</i>	No	Because the effect of a pair of <i>read</i> operations does not depend on the order in which they are executed
<i>read write</i>	Yes	Because the effect of a <i>read</i> and a <i>write</i> operation depends on the order of their execution
<i>write write</i>	Yes	Because the effect of a pair of <i>write</i> operations depends on the order of their execution

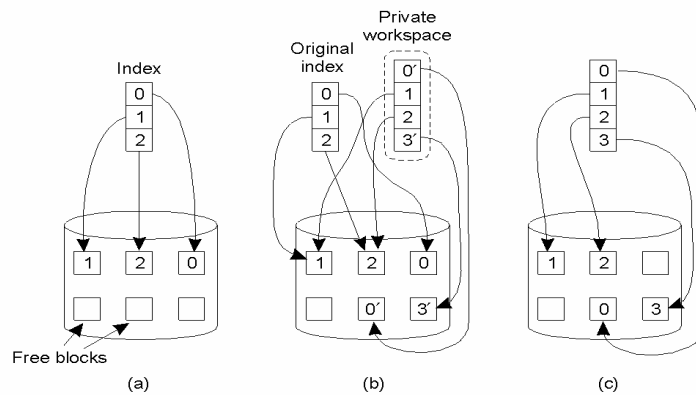
Transactions 14

A non-serially equivalent interleaving of operations of transactions *T* and *U*

Transaction <i>T</i> :	Transaction <i>U</i> :
<i>x</i> = read(<i>i</i>)	
write(<i>i</i> , 10)	<i>y</i> = read(<i>j</i>)
	write(<i>j</i> , 30)
write(<i>j</i> , 20)	
	<i>z</i> = read(<i>i</i>)

Transactions 15

Implementing Transactions: Private Workspace



- a) The file index and disk blocks for a three-block file
- b) The situation after a transaction has modified block 0 and appended block 3
- c) After committing

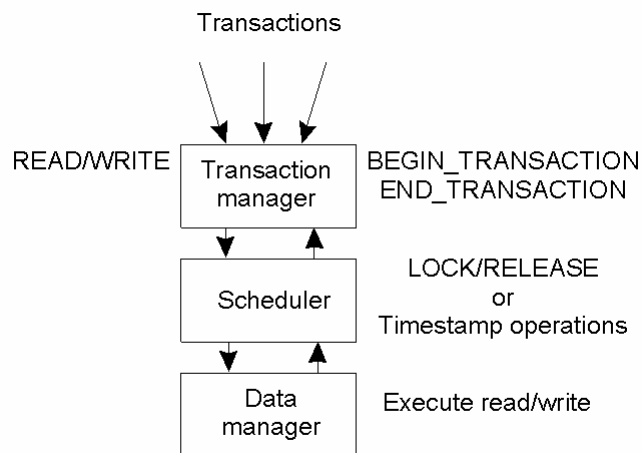
Transactions 16

Implementing Transactions: Writeahead Log

<pre>x = 0; y = 0; BEGIN_TRANSACTION; x = x + 1; y = y + 2; x = y * y; END_TRANSACTION;</pre>	<p>Log</p> <p>[x = 0 / 1]</p> <p>(b)</p>	<p>Log</p> <p>[x = 0 / 1] [y = 0/2]</p> <p>(c)</p>	<p>Log</p> <p>[x = 0 / 1] [y = 0/2] [x = 1/4]</p> <p>(d)</p>
---	--	--	--

- a) A transaction
- b) – d) The log before each statement is executed

Concurrency Control



General organization of managers for handling transactions.

Transactions T and U with exclusive locks

Transaction T		Transaction U	
<i>balance = b.getBalance()</i>		<i>balance = b.getBalance()</i>	
<i>b.setBalance(bal*1.1)</i>		<i>b.setBalance(bal*1.1)</i>	
<i>a.withdraw(bal/10)</i>		<i>c.withdraw(bal/10)</i>	
Operations	Locks	Operations	Locks
<i>openTransaction</i>		<i>openTransaction</i>	
<i>bal = b.getBalance()</i>	lock B	<i>bal = b.getBalance()</i>	waits for T 's lock on B
<i>b.setBalance(bal*1.1)</i>		•••	
<i>a.withdraw(bal/10)</i>	lock A		lock B
<i>closeTransaction</i>	unlock A, B	<i>b.setBalance(bal*1.1)</i>	
		<i>c.withdraw(bal/10)</i>	lock C
		<i>closeTransaction</i>	unlock B, C

Transactions 19

Lock compatibility

<i>For one object</i>		<i>Lock requested</i>	
		<i>read</i>	<i>write</i>
<i>Lock already set</i>	<i>none</i>	OK	OK
	<i>read</i>	OK	wait
	<i>write</i>	wait	wait

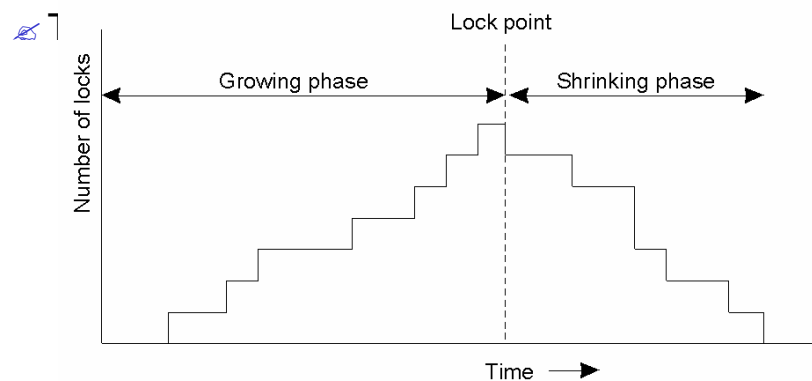
Transactions 20

Use of locks in strict two-phase locking

1. When an operation accesses an object within a transaction:
 - (a) If the object is not already locked, it is locked and the operation proceeds.
 - (b) If the object has a conflicting lock set by another transaction, the transaction must wait until it is unlocked.
 - (c) If the object has a non-conflicting lock set by another transaction, the lock is shared and the operation proceeds.
 - (d) If the object has already been locked in the same transaction, the lock will be promoted if necessary and the operation proceeds. (Where promotion is prevented by a conflicting lock, rule (b) is used.)
2. When a transaction is committed or aborted, the server unlocks all objects it locked for the transaction.

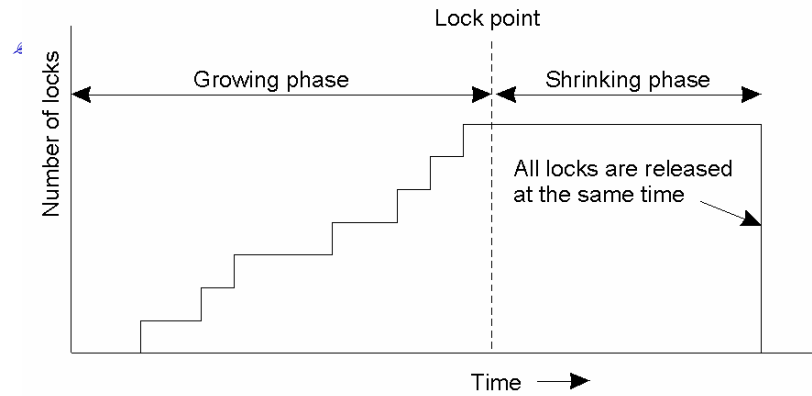
Transactions 21

Two-Phase Locking (1)



Transactions 22

Strict Two-Phase Locking (2)



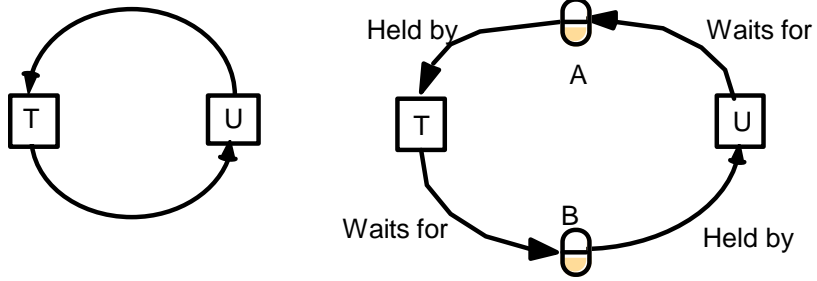
Transactions 23

Deadlock with write locks

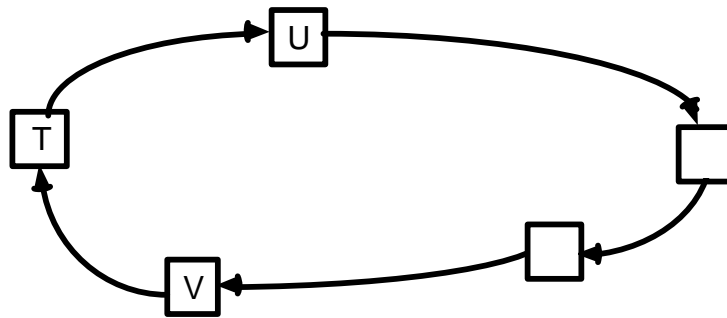
Transaction <i>T</i>		Transaction <i>U</i>	
Operations	Locks	Operations	Locks
<i>a.deposit(100);</i>	write lock <i>A</i>	<i>b.deposit(200)</i>	write lock <i>B</i>
<i>b.withdraw(100)</i>	waits for <i>U</i> 's	<i>a.withdraw(200);</i>	waits for <i>T</i> 's
•••	lock on <i>B</i>	•••	lock on <i>A</i>
•••		•••	
•••		•••	

Transactions 24

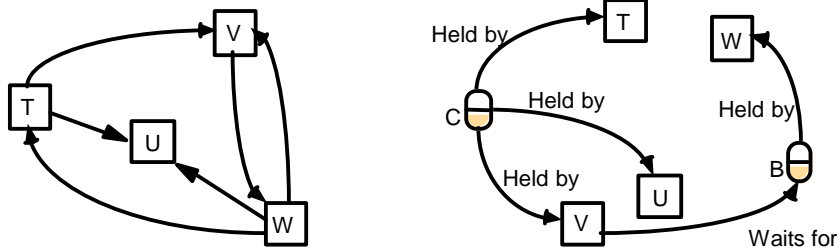
The wait-for graph



A cycle in a wait-for graph



Another wait-for graph



Transactions 27

Resolution of deadlock

Transaction T		Transaction U	
Operations	Locks	Operations	Locks
<i>a.deposit(100);</i>	write lock <i>A</i>	<i>b.deposit(200)</i>	write lock <i>B</i>
<i>b.withdraw(100)</i>		<i>a.withdraw(200);</i>	waits for T's
...	waits for U's	...	lock on <i>A</i>
	lock on <i>B</i>	...	
	(timeout elapses)	...	
<i>T's lock on A becomes vulnerable,</i>		<i>a.withdraw(200);</i>	write locks <i>A</i>
unlock <i>A</i> , abort T			unlock <i>A, B</i>

Transactions 28

Optimistic Concurrency Control

- ✍ Drawbacks of locking
 - ✍ Overhead of lock maintenance
 - ✍ Deadlocks
 - ✍ Reduced concurrency
- ✍ Optimistic Concurrency Control
 - ✍ In most applications, likelihood of conflicting accesses by concurrent transactions is low
 - ✍ Transactions proceed as though there are no conflicts
 - ✍ Three phases
 - ✍ Working Phase – transactions read and write private copies of objects
 - ✍ Validation Phase – each transaction is assigned a transaction number when it enters this phase
 - ✍ Update Phase

Transactions 29

Optimistic Concurrency Control: Serializability of transaction T_v with respect to transaction T_i

T_v and T_i are overlapping transactions

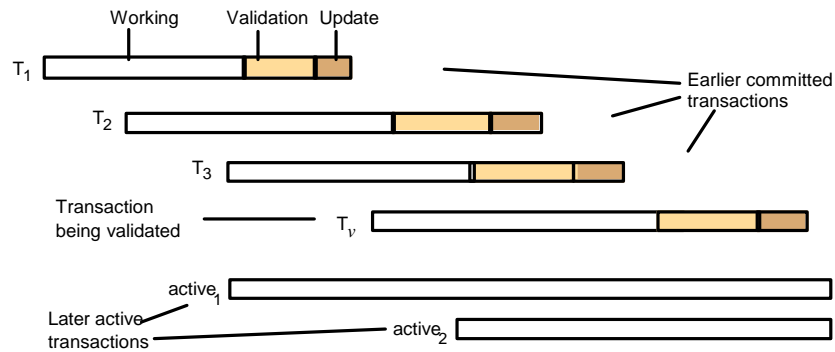
For T_v to be serializable wrt T_i the following rules must hold

T_v	T_i	Rule
<i>write</i>	<i>read</i>	1. T_i must not read objects written by T_v
<i>read</i>	<i>write</i>	2. T_v must not read objects written by T_i
<i>write</i>	<i>write</i>	3. T_i must not write objects written by T_v and T_v must not write objects written by T_i

If simplification is made that only one transaction may be in its validation or write phases at one time, then third rule is always satisfied

Transactions 30

Validation of transactions



Transactions 31

Validation of Transactions

Backward validation of transaction T_v

```

boolean valid = true;
for (int  $T_i = startTn+1$ ;  $T_i \leq finishTn$ ;  $T_i++$ ){
    if (read set of  $T_v$  intersects write set of  $T_i$ ) valid = false;
}

```

Forward validation of transaction T_v

```

boolean valid = true;
for (int  $T_{id} = activeI$ ;  $T_{id} \leq activeN$ ;  $T_{id}++$ ){
    if (write set of  $T_v$  intersects read set of  $T_{id}$ ) valid = false;
}

```

Transactions 32

Timestamp based concurrency control

- ✍ Each timestamp is assigned a unique timestamp at the moment it starts
 - ✍ In distributed transactions, Lamport's timestamps can be used
- ✍ Every data item has a timestamp
 - ✍ Read timestamp = timestamp of transaction that last read the item
 - ✍ Write timestamp = timestamp of transaction that most recently changed an item

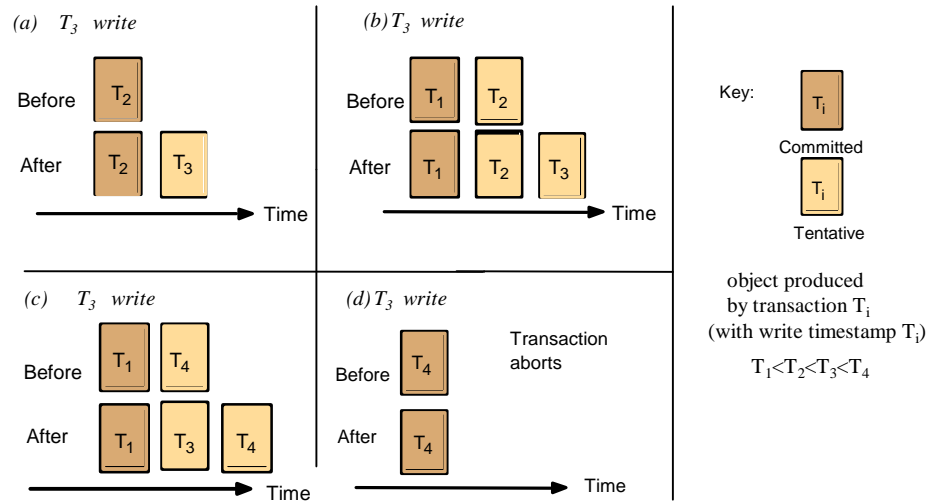
Operation conflicts for timestamp ordering

Rule	T_c	T_i	
1.	<i>write</i>	<i>read</i>	T_c must not <i>write</i> an object that has been <i>read</i> by any T_i where $T_i > T_c$. this requires that $T_c =$ the maximum read timestamp of the object.
2.	<i>write</i>	<i>write</i>	T_c must not <i>write</i> an object that has been <i>written</i> by any T_i where $T_i > T_c$. this requires that $T_c >$ write timestamp of the committed object.
3.	<i>read</i>	<i>write</i>	T_c must not <i>read</i> an object that has been <i>written</i> by any T_i where $T_i > T_c$. this requires that $T_c >$ write timestamp of the committed object.

Timestamp ordering write rule

if (T_c = maximum read timestamp on D &&
 $T_c >$ write timestamp on committed version of D)
 perform write operation on tentative version of D with write timestamp T_c
 else /* write is too late */
 Abort transaction T_c

Write operations and timestamps



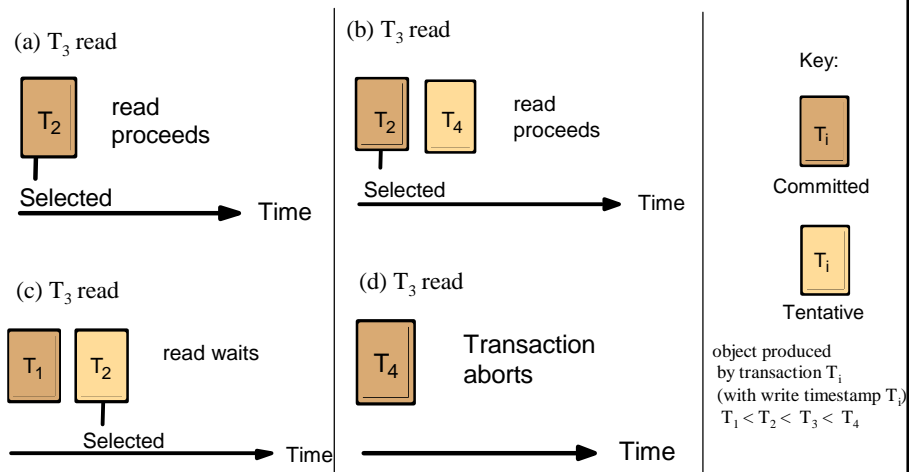
Timestamp ordering read rule

```

if ( $T_c >$  write timestamp on committed version of  $D$ ) {
  let  $D_{\text{selected}}$  be the version of  $D$  with the maximum write timestamp =  $T_c$ 
  if ( $D_{\text{selected}}$  is committed)
    perform read operation on the version  $D_{\text{selected}}$ 
  else
    Wait until the transaction that made version  $D_{\text{selected}}$  commits or aborts
    then reapply the read rule
} else
  Abort transaction  $T_c$ 
  
```

Transactions 37

Read operations and timestamps



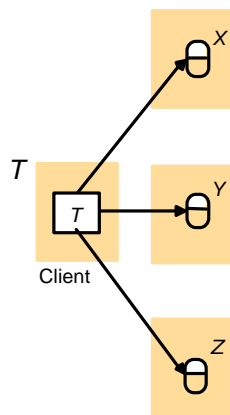
Transactions 38

Timestamps in transactions T and U

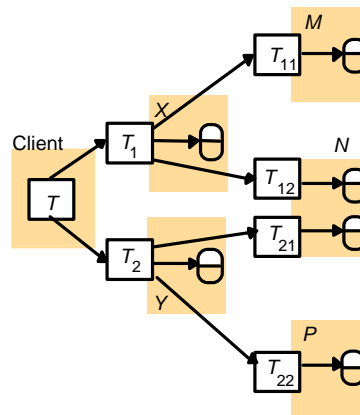
		<i>Timestamps and versions of objects</i>					
T	U	A		B		C	
		<i>RTS</i>	<i>WTS</i>	<i>RTS</i>	<i>WTS</i>	<i>RTS</i>	<i>WTS</i>
<i>openTransaction</i>		{ }	S	{ }	S	{ }	S
$bal = b.getBalance()$		{ T }					
$b.setBalance(bal*1.1)$	<i>openTransaction</i>	S, T					
$a.withdraw(bal/10)$	$bal = b.getBalance()$	S, T					
<i>commit</i>	<i>wait for T</i>	T					
	•••	T					
	•••	T					
	$bal = b.getBalance()$	{ U }					
	$b.setBalance(bal*1.1)$	T, U					
	$c.withdraw(bal/10)$	S, U					

Distributed transactions

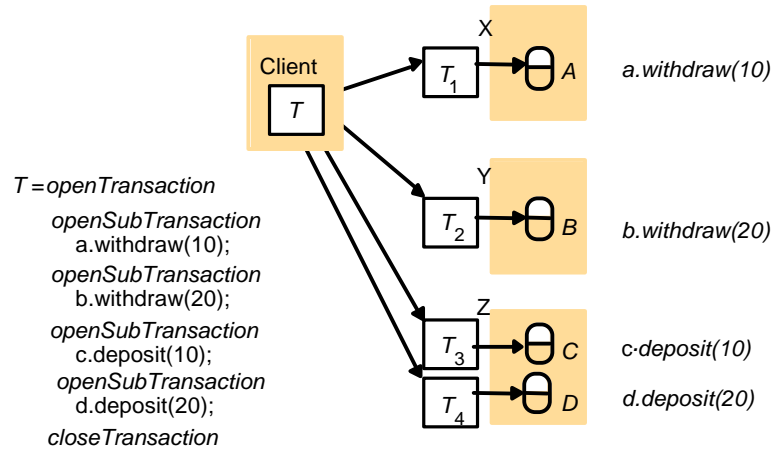
(a) Flat transaction



(b) Nested transactions

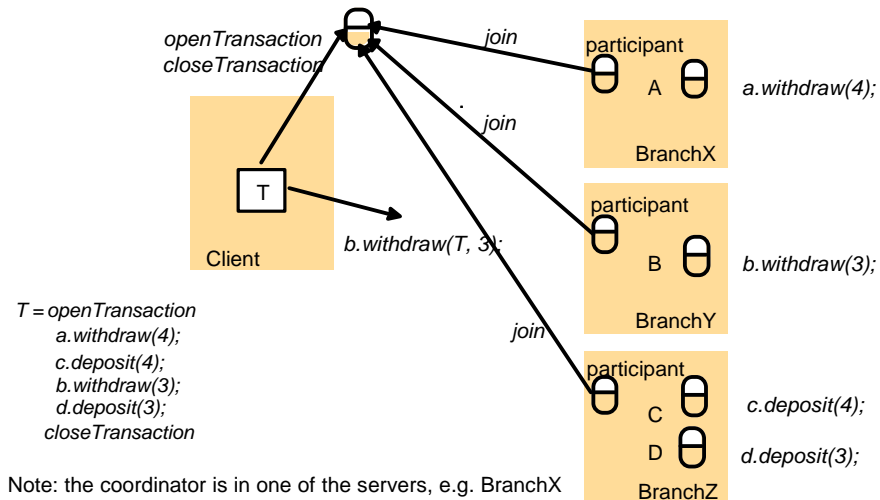


Nested banking transaction



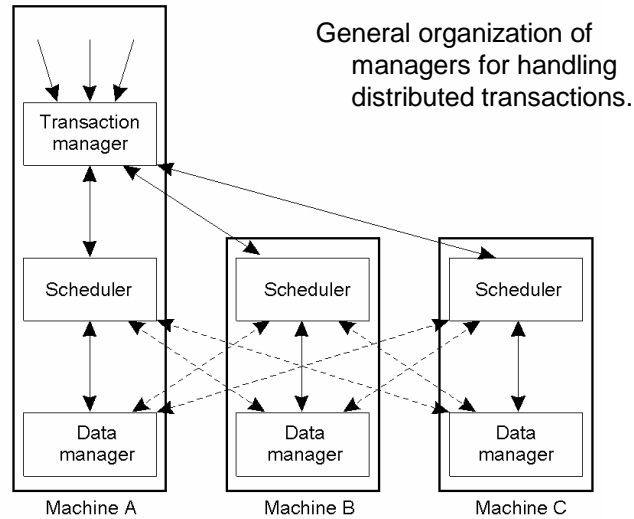
Transactions 41

A distributed banking transaction



Transactions 42

Concurrency Control for Distributed Transactions



Transactions 43

Concurrency Control for Distributed Transactions

✍ Locking

- ✍ Distributed deadlocks possible

✍ Timestamp ordering

- ✍ Lamport time stamps

✍ for efficiency it is required that timestamps issued by coordinators be roughly synchronized

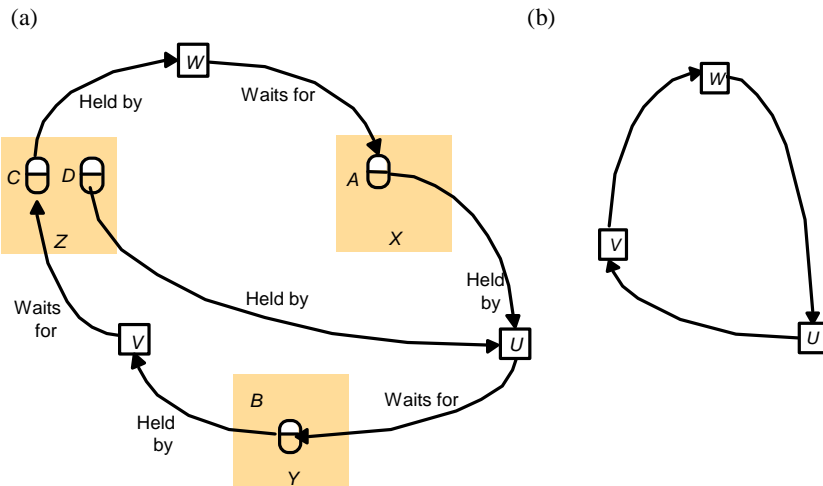
Transactions 44

Interleavings of transactions *U*, *V* and *W*

<i>U</i>	<i>V</i>	<i>W</i>
<i>d.deposit(10)</i> lock <i>D</i>	<i>b.deposit(10)</i> lock <i>B</i> at <i>Y</i>	
<i>a.deposit(20)</i> lock <i>A</i> at <i>X</i>		<i>c.deposit(30)</i> lock <i>C</i> at <i>Z</i>
<i>b.withdraw(30)</i> wait at <i>Y</i>	<i>c.withdraw(20)</i> wait at <i>Z</i>	<i>a.withdraw(20)</i> wait at <i>X</i>

Transactions 45

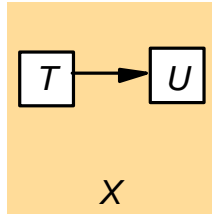
Distributed deadlock



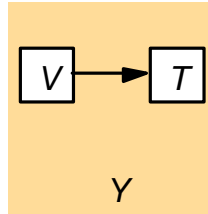
Transactions 46

Local and global wait-for graphs

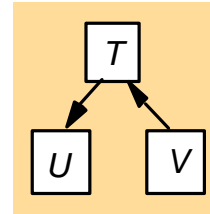
local wait-for graph



local wait-for graph



global deadlock detector



Transactions 47

Atomic Commit Protocols

- ✍ The atomicity of a transaction requires that when a distributed transaction comes to an end, either all of its operations are carried out or none of them
- ✍ One phase commit
 - ✍ Coordinator tells all participants to commit
 - ✍ If a participant cannot commit (say because of concurrency control), no way to inform coordinator
- ✍ Two phase commit (2PC)

Transactions 48

The two-phase commit protocol

Phase 1 (voting phase):

1. The coordinator sends a *canCommit?* request to each of the participants in the transaction.
2. When a participant receives a *canCommit?* request it replies with its vote (*Yes* or *No*) to the coordinator. Before voting *Yes*, it prepares to commit by saving objects in permanent storage. If the vote is *No* the participant aborts immediately.

Phase 2 (completion according to outcome of vote):

3. The coordinator collects the votes (including its own).
 - (a) If there are no failures and all the votes are *Yes* the coordinator decides to commit the transaction and sends a *doCommit* request to each of the participants.
 - (b) Otherwise the coordinator decides to abort the transaction and sends *doAbort* requests to all participants that voted *Yes*.
4. Participants that voted *Yes* are waiting for a *doCommit* or *doAbort* request from the coordinator. When a participant receives one of these messages it acts accordingly and in the case of commit, makes a *haveCommitted* call as confirmation to the coordinator.

Transactions 49

Operations for two-phase commit protocol

canCommit?(trans) -> Yes / No

Call from coordinator to participant to ask whether it can commit a transaction. Participant replies with its vote.

doCommit(trans)

Call from coordinator to participant to tell participant to commit its part of a transaction.

doAbort(trans)

Call from coordinator to participant to tell participant to abort its part of a transaction.

haveCommitted(trans, participant)

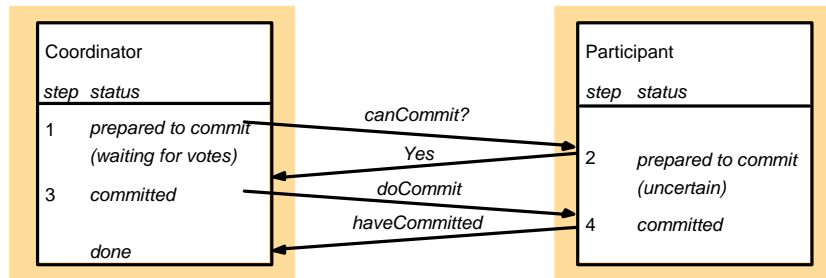
Call from participant to coordinator to confirm that it has committed the transaction.

getDecision(trans) -> Yes / No

Call from participant to coordinator to ask for the decision on a transaction after it has voted *Yes* but has still had no reply after some delay. Used to recover from server crash or delayed messages.

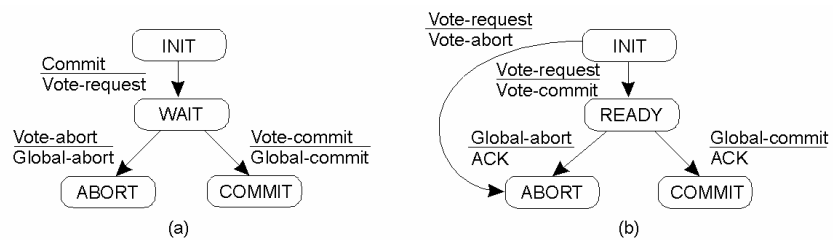
Transactions 50

Communication in two-phase commit protocol



Transactions 51

Two-Phase Commit (1)



- a) The finite state machine for the coordinator in 2PC.
- b) The finite state machine for a participant.

Transactions 52

Two-Phase Commit (2)

State of Q	Action by P
COMMIT	Make transition to COMMIT
ABORT	Make transition to ABORT
INIT	Make transition to ABORT
READY	Contact another participant

Actions taken by a participant *P* when residing in state *READY* and having contacted another participant *Q*.

Transactions 53

Two-Phase Commit (3)

actions by coordinator:

```
while START_2PC to local log;
multicast VOTE_REQUEST to all participants;
while not all votes have been collected {
  wait for any incoming vote;
  if timeout {
    write GLOBAL_ABORT to local log;
    multicast GLOBAL_ABORT to all participants;
    exit;
  }
  record vote;
}
if all participants sent VOTE_COMMIT and coordinator votes COMMIT{
  write GLOBAL_COMMIT to local log;
  multicast GLOBAL_COMMIT to all participants;
} else {
  write GLOBAL_ABORT to local log;
  multicast GLOBAL_ABORT to all participants;
}
```

Outline of the steps taken by the coordinator in a two phase commit protocol

Transactions 54

Two-Phase Commit (4)

Steps taken by
participant
process in
2PC.

actions by participant:

```
write INIT to local log;
wait for VOTE_REQUEST from coordinator;
if timeout {
  write VOTE_ABORT to local log;
  exit;
}
if participant votes COMMIT {
  write VOTE_COMMIT to local log;
  send VOTE_COMMIT to coordinator;
  wait for DECISION from coordinator;
  if timeout {
    multicast DECISION_REQUEST to other participants;
    wait until DECISION is received; /* remain blocked */
    write DECISION to local log;
  }
  if DECISION == GLOBAL_COMMIT
    write GLOBAL_COMMIT to local log;
  else if DECISION == GLOBAL_ABORT
    write GLOBAL_ABORT to local log;
} else {
  write VOTE_ABORT to local log;
  send VOTE_ABORT to coordinator;
}
```

Transactions 55

Two-Phase Commit (5)

actions for handling decision requests: /* executed by separate thread */

```
while true {
  wait until any incoming DECISION_REQUEST is received; /* remain blocked */
  read most recently recorded STATE from the local log;
  if STATE == GLOBAL_COMMIT
    send GLOBAL_COMMIT to requesting participant;
  else if STATE == INIT or STATE == GLOBAL_ABORT
    send GLOBAL_ABORT to requesting participant;
  else
    skip; /* participant remains blocked */
}
```

Steps taken for handling incoming decision requests.

Transactions 56

Three Phase Commit

✍ Problem with 2PC

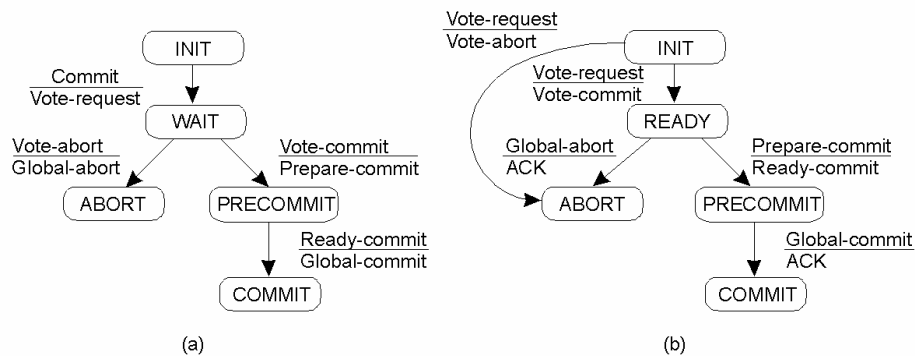
- ✍ If coordinator crashes, participants cannot reach a decision, stay blocked until coordinator recovers

✍ 3PC

- ✍ There is no single state from which it is possible to make a transition directly to either COMMIT or ABORT states
- ✍ There is no state in which it is not possible to make a final decision, and from which a transition to COMMIT can be made

Transactions 57

Three-Phase Commit



- a) Finite state machine for the coordinator in 3PC
 b) Finite state machine for a participant

Transactions 58