

Universität Stuttgart
Master of Science

Distributed Transactions

Marco Aiello

Transactions

- A *transaction* is a group of operations on an object to be executed coherently according to some policy, typically, atomicity (i.e., all or nothing)

Single server, single client

- If one unique “client” of an “object” at one server then locking (e.g. synchronized in java) is enough

ACIDity

A tomicity
C onsistency
I solation
D urability

For atomicity and durability, objects must be *recoverable*

Transactions facts

- The goal of any server supporting transactions is to maximize concurrency
- *Serializability* is often the requirement
- Transaction support can be part of the middleware (e.g, CORBA's Object Transaction Service, the TID are implicit)
- Transactions are created and managed by a *coordinator*

Figure 16.3

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.

Figure 16.4
Transaction life histories

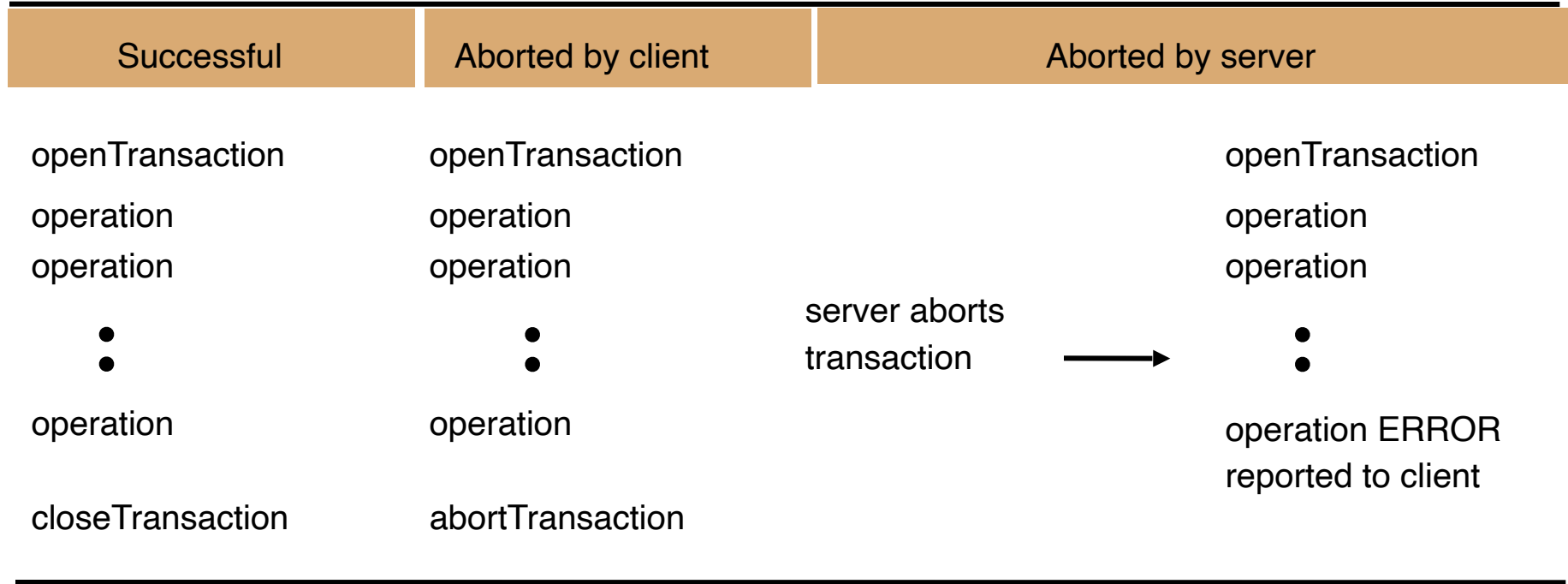


Figure 16.1

Banking example: 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

Figure 16.2

A client's banking transaction

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

Figure 16.5
The lost update problem

a,b,c are \$100, \$200, \$300, resp.

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

Figure 16.7
 A serially equivalent interleaving of *T* and *U*

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

Figure 16.6
The inconsistent retrievals problem

a,b are \$200

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

Figure 16.8
 A serially equivalent interleaving of *V* and *W*

Transaction V:		Transaction W:	
<i>a.withdraw(100);</i> <i>b.deposit(100)</i>		<i>aBranch.branchTotal()</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>	
		...	

Conflicting operations

- A pair of operations is *conflicting* when their combined effect depends on the order of execution

Figure 16.9
Read and write operation conflict rules

<i>Operations of different transactions</i>		<i>Conflict</i>	<i>Reason</i>
<i>read</i>	<i>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</i>	<i>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</i>	<i>write</i>	Yes	Because the effect of a pair of <i>write</i> operations depends on the order of their execution

Figure 16.10

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

Transaction T :	Transaction U :
$x = \text{read}(i)$	
$\text{write}(i, 10)$	$y = \text{read}(j)$
	$\text{write}(j, 30)$
$\text{write}(j, 20)$	
	$z = \text{read}(i)$

Serially equivalence

Serially equivalence requires either one of:

1. T accesses i before U and T access j before U
2. U accesses i before T and U access j before T

Can be obtained with:

1. locks
2. optimistic concurrency
3. timestamping

Recoverability from aborts

- How to take into account that transactions may abort while executing?
 - *Recoverability*: any possible dirty read forces delayed commit
 - *Avoid cascading aborts*: transactions are allowed only to read values which are committed

Figure 16.11

A dirty read when transaction T aborts

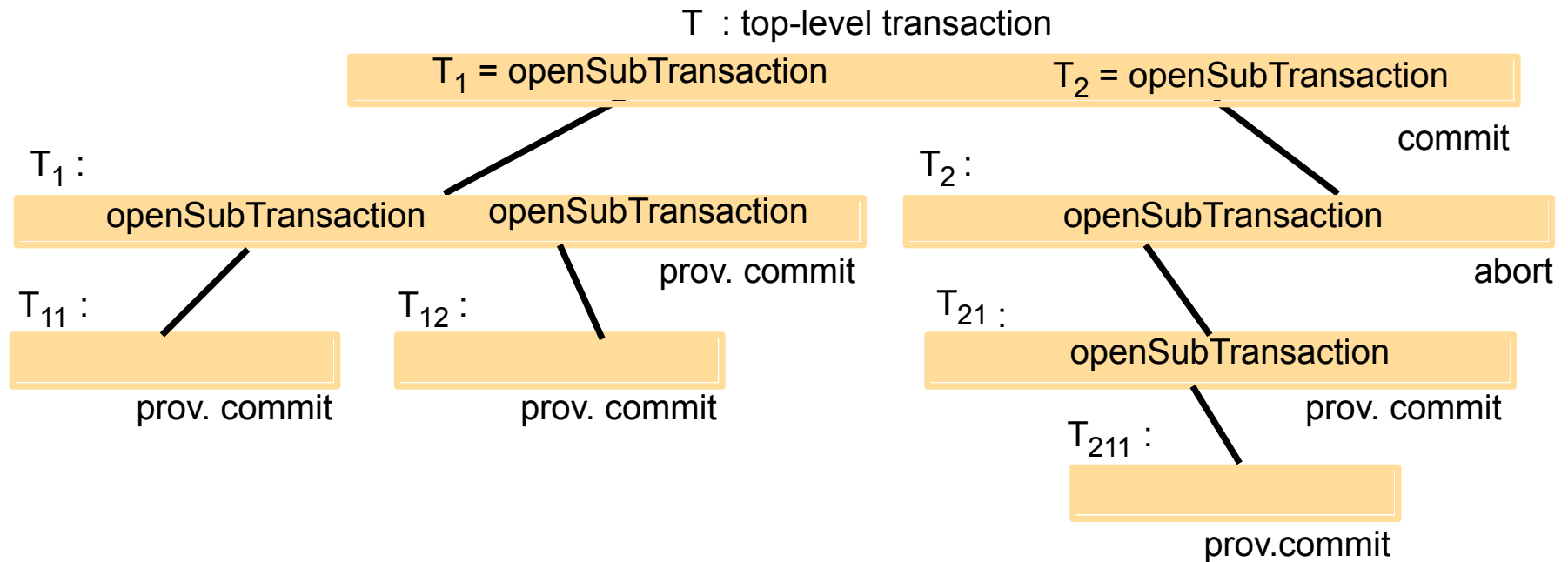
Transaction T :	Transaction U :
$a.getBalance()$ $a.setBalance(balance + 10)$	$a.getBalance()$ $a.setBalance(balance + 20)$
$balance = a.getBalance()$ \$100 $a.setBalance(balance + 10)$ \$110	$balance = a.getBalance()$ \$110 $a.setBalance(balance + 20)$ \$130 <i>commit transaction</i>
<i>abort transaction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

Figure 16.12
Overwriting uncommitted values

Transaction <i>T</i> :		Transaction <i>U</i> :	
<i>a.setBalance(105)</i>		<i>a.setBalance(110)</i>	
	\$100		
<i>a.setBalance(105)</i>	\$105		
		<i>a.setBalance(110)</i>	\$110

Figure 16.13

Nested transactions



Nested Transactions

- Subtransactions at one level may run concurrently with other subtransactions at the same level
- Subtransactions can commit or abort independently

Nested transaction commit rules

1. A transaction may commit or abort only after its children transactions have completed
2. When a subtransaction completes, it decides independently to provisionally commit or abort
3. When a parent aborts, all its subtransactions abort
4. When a subtransaction aborts, the parent decides whether to abort or not
5. When the top-level transaction commits, all of the subtransaction that have provisionally committed can commit too

How to guarantee serial equivalence?

1. Locks (seen in Operating Systems)
2. Optimistic concurrency control
3. Timestamping

Figure 16.14
Transactions *T* and *U* with exclusive locks

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

Figure 16.15
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

Figure 16.16

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.
-

Figure 16.19
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 lock on <i>B</i>	<i>a.withdraw(200);</i>	waits for <i>T</i> 's lock on <i>A</i>
• • •		• • •	
• • •		• • •	
		• • •	

Figure 16.20
The wait-for graph for Figure 16.19

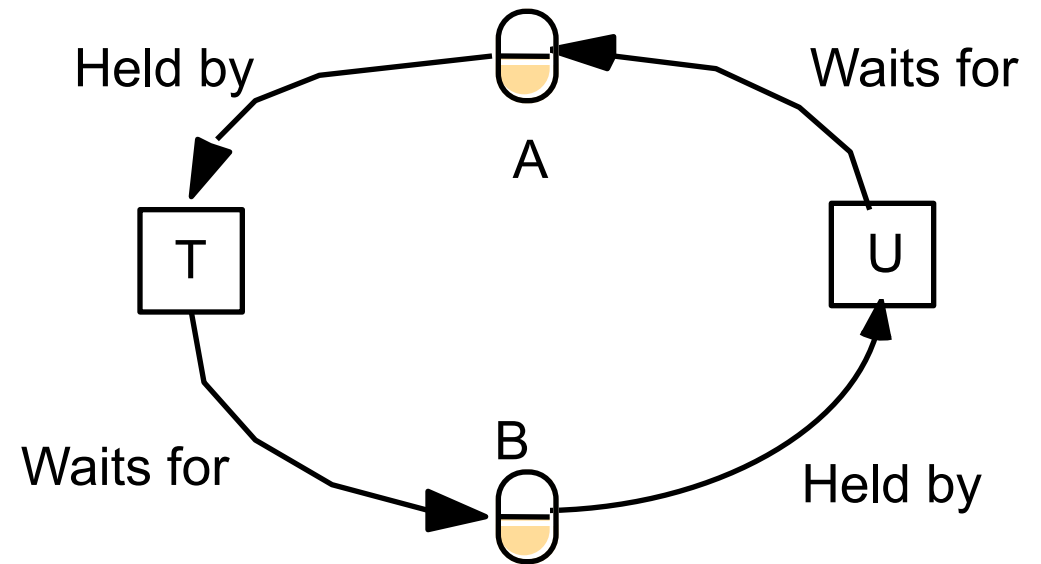
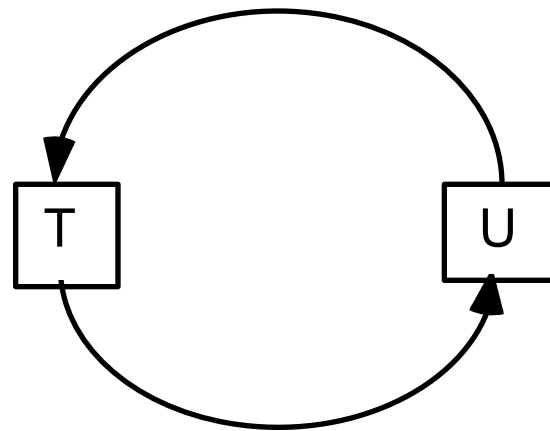


Figure 16.21
A cycle in a wait-for graph

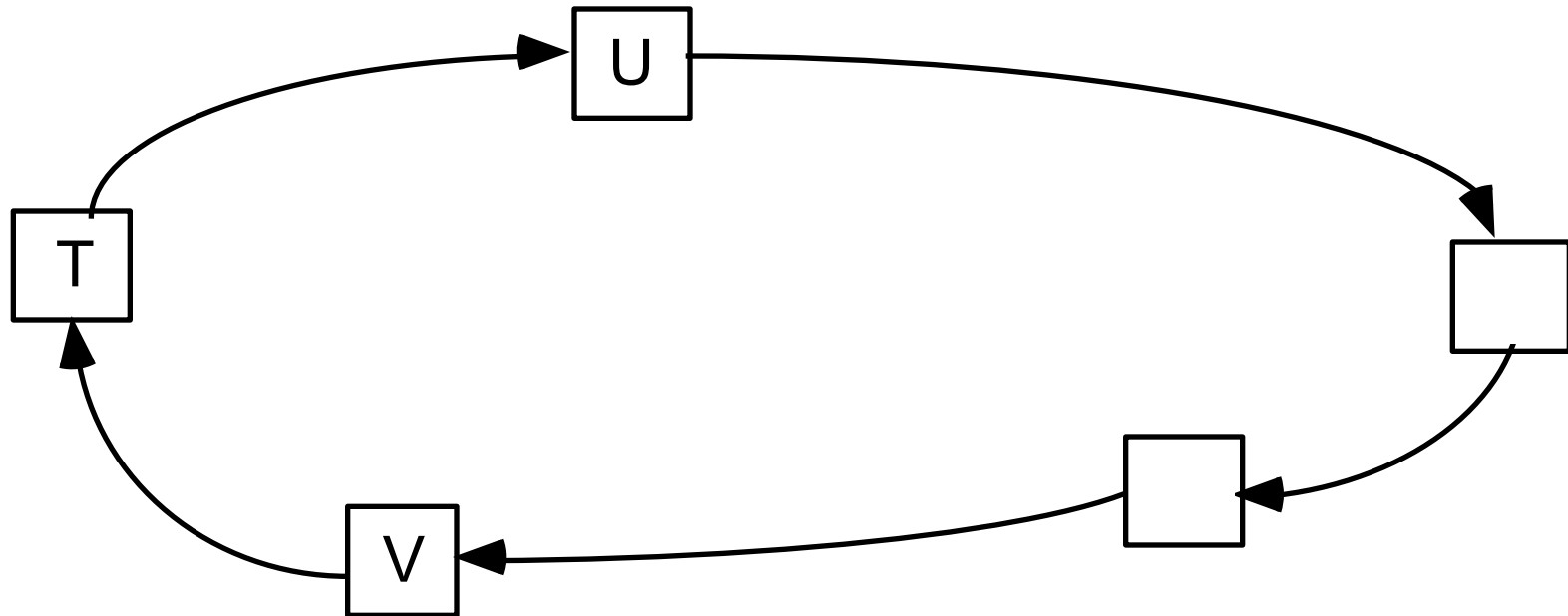


Figure 16.22
Another wait-for graph

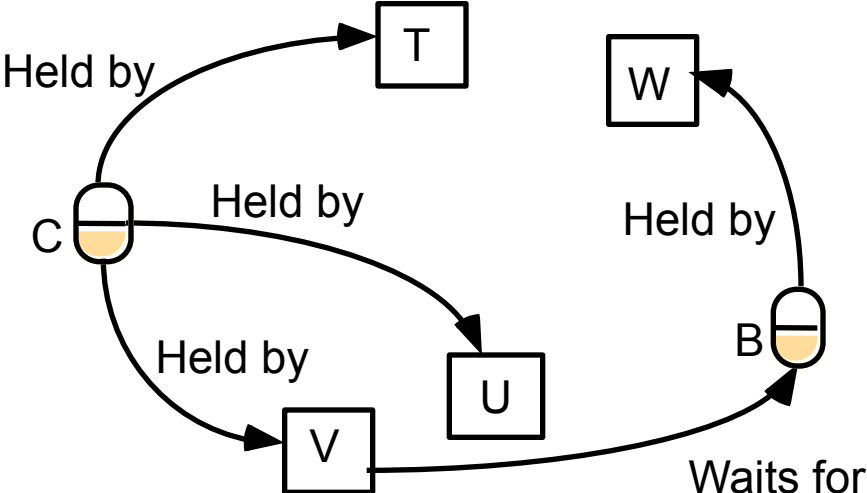
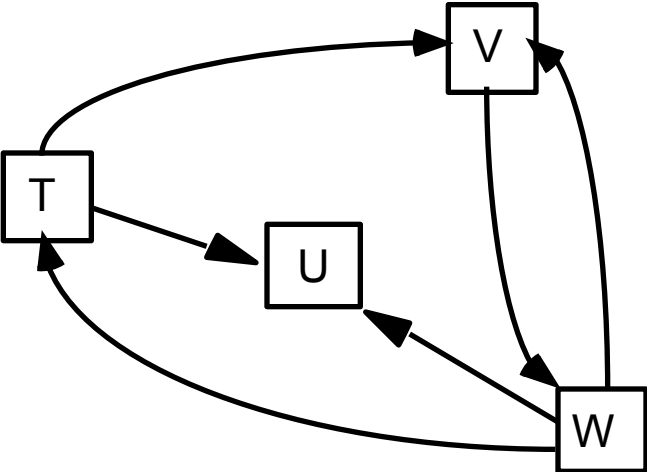


Figure 16.23
Resolution of the deadlock in Figure 15.19

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

Optimistic concurrency control

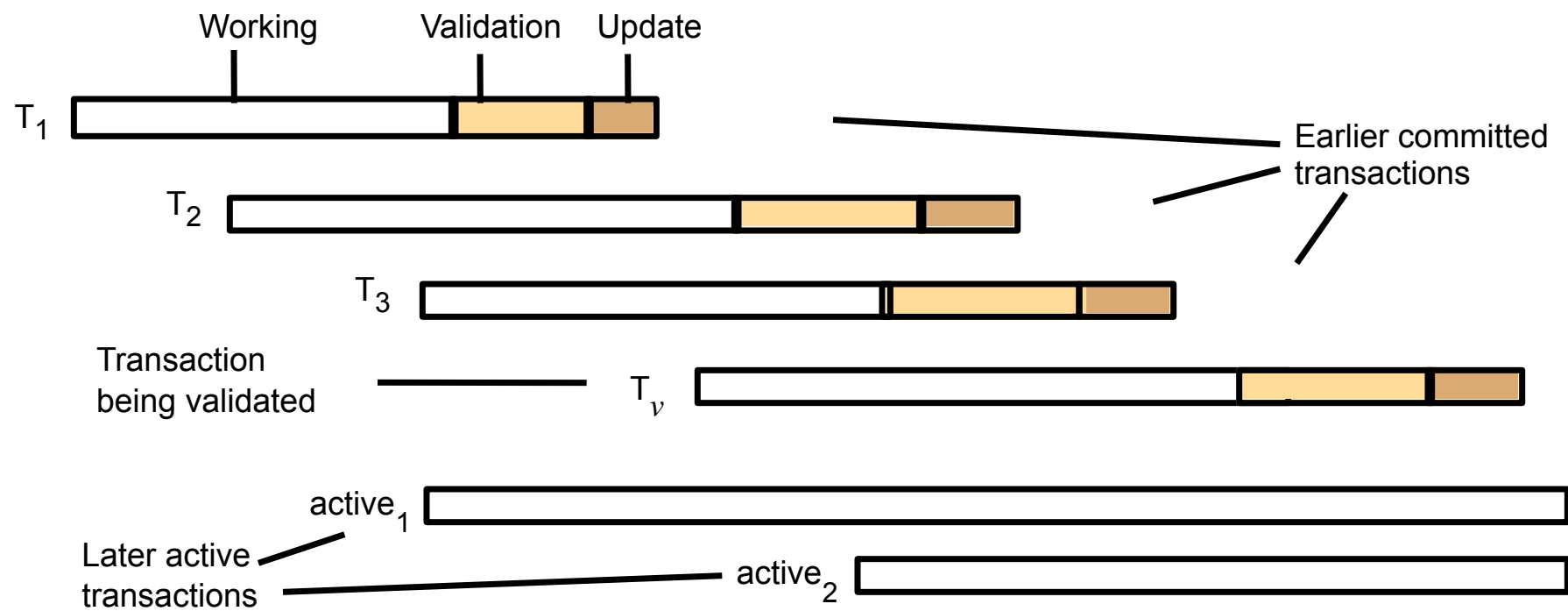
- ◆ Locking has drawbacks: overheads, deadlocks, limited concurrency to avoid cascading aborts
- ◆ Optimistic: transactions are allowed to proceed until ready to commit. If a conflict occurs, then abort.
- ◆ How? In phases:
 - ◆ Working phase: execute operations. If concurrent transactions, many values of an object may coexist. Keep *write set* and *read set* for each transaction.
 - ◆ Validation phase: either no conflict or call conflict resolution policy.
 - ◆ Update phase: if validated, make changes permanent

Table on page 708

Serializability of transaction T with respect to transaction T_i

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

Figure 16.28
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;
}
```

◆ the read set of T_v must be compared with the write sets of $T_2 T_3$

Forward validation of transaction T_v

```
boolean valid = true;
for (int  $T_{id} = active1$ ;  $T_{id} \leq activeN$ ;  $T_{id}++$ ) {
    if (write set of  $T_v$  intersects read set of  $T_{id}$ ) valid = false;
}
```

◆ the write set of T_v must be compared with the read sets of $active_{1,2}$

starvation?

Optimistic concurrency control

- ◆ Alternatives to conflict resolution:
 - ◆ defer validation until all conflicting transactions have finished
 - ◆ abort all the conflicting transactions except the one being validated
 - ◆ abort the transaction being validated

Timestamp ordering

- ♦ A transaction's request to write an object is valid only if that object was last read and written by earlier transactions.
- ♦ A transaction's request to read an object is valid only if that object was last written by an earlier transaction.
- ♦ I.e., tentative version of each object are committed in the order determined by the timestamps of the transactions.

Figure 16.29

Operation conflicts for timestamp ordering

<i>Rule</i>	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 $T_c \geq$ 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.

Examples



- **Dropbox:** uses an optimistic form of concurrency control, keeps track of conflicts, granularity of files



- **Google apps:** finer granularity, user resolve conflict manually



- **Wikipedia:** editing concurrency is optimistic, manual conflict resolution



- **Amazon Dynamo:** key value storage. No isolation guarantee, optimistic concurrency, conflict resolution via application logic or timestamping (last write wins)

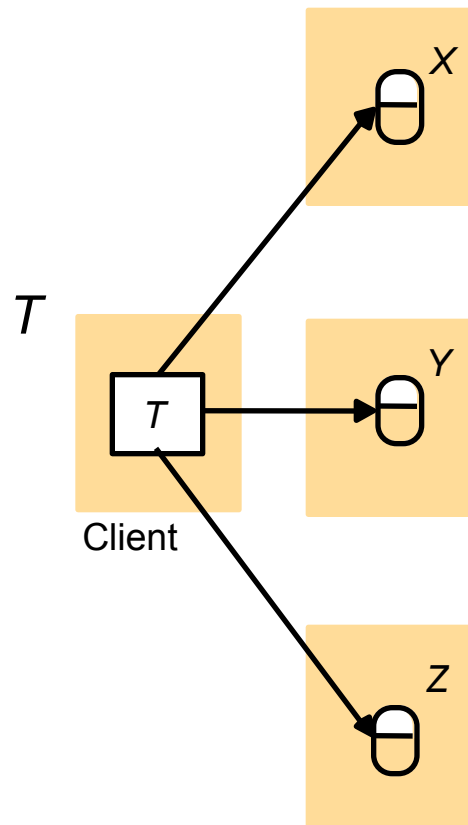
Distributed transactions

- ♦ What if more servers hold the objects accessed by the operations of a transaction? *Distributed transactions*.
- ♦ *Coordinator*: ensuring the same outcome at all of the servers

Figure 17.1

Distributed transactions

(a) Flat transaction



(b) Nested transactions

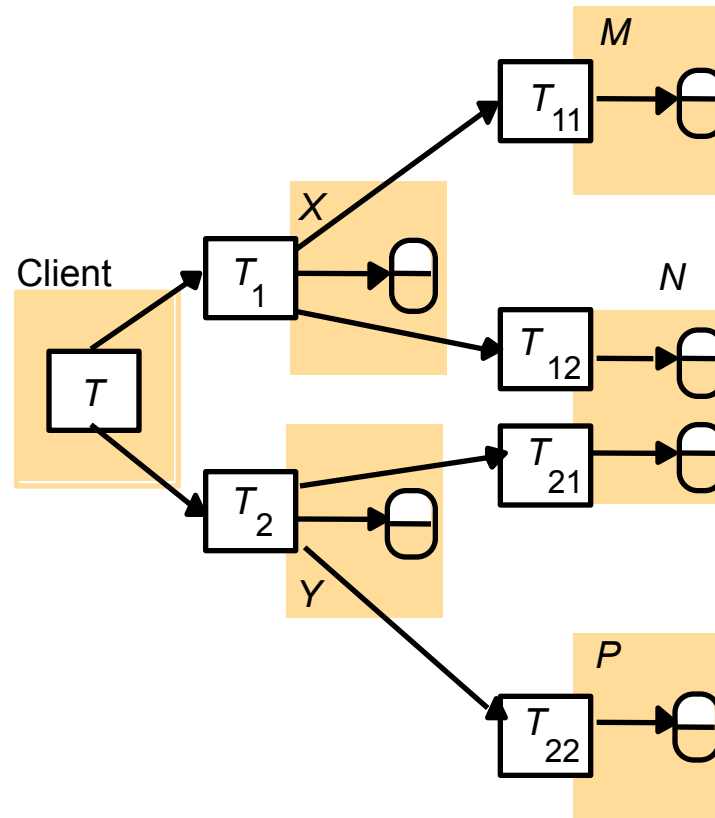


Figure 17.2
Nested banking transaction

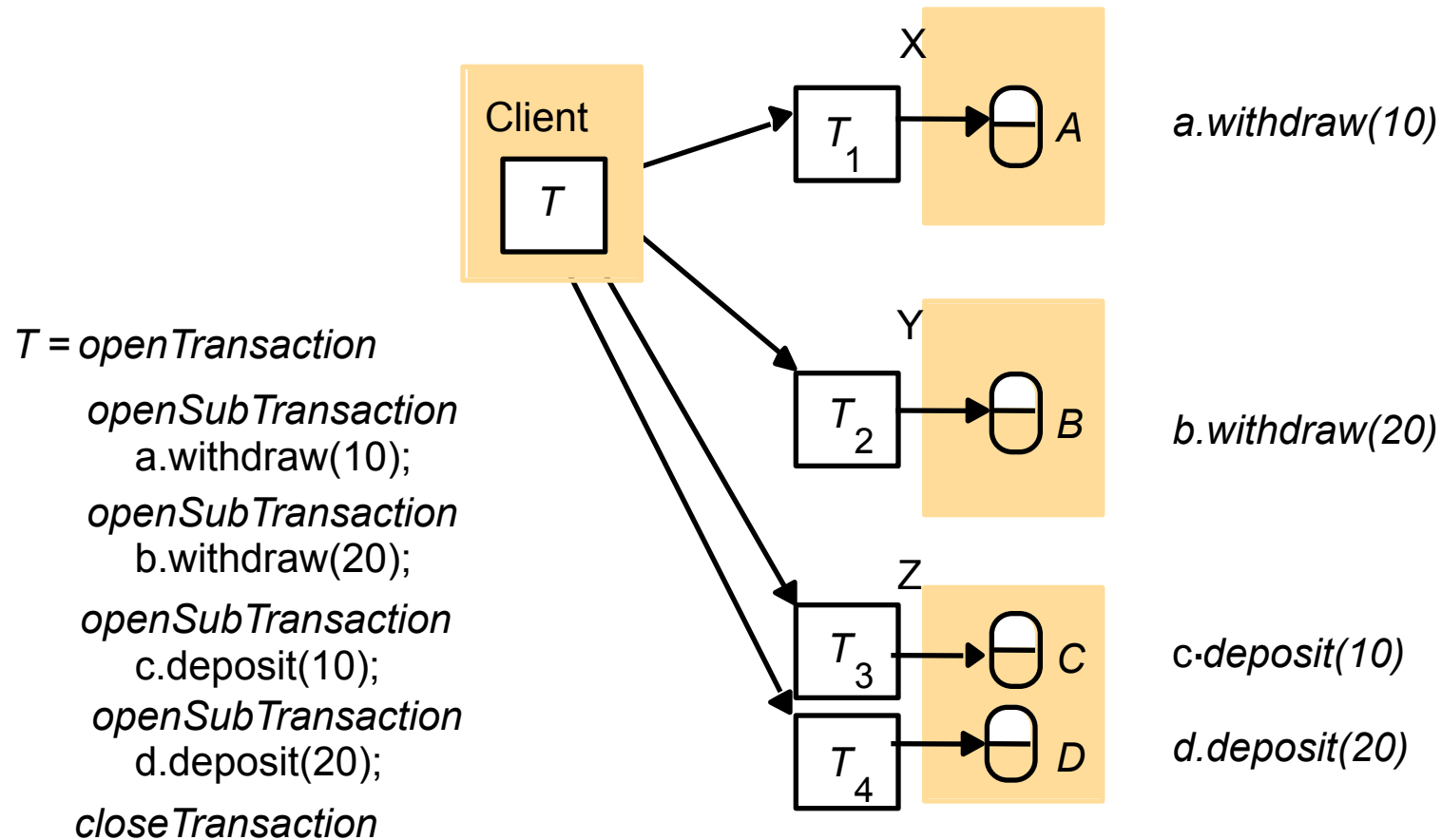


Figure 17.3
A distributed banking transaction

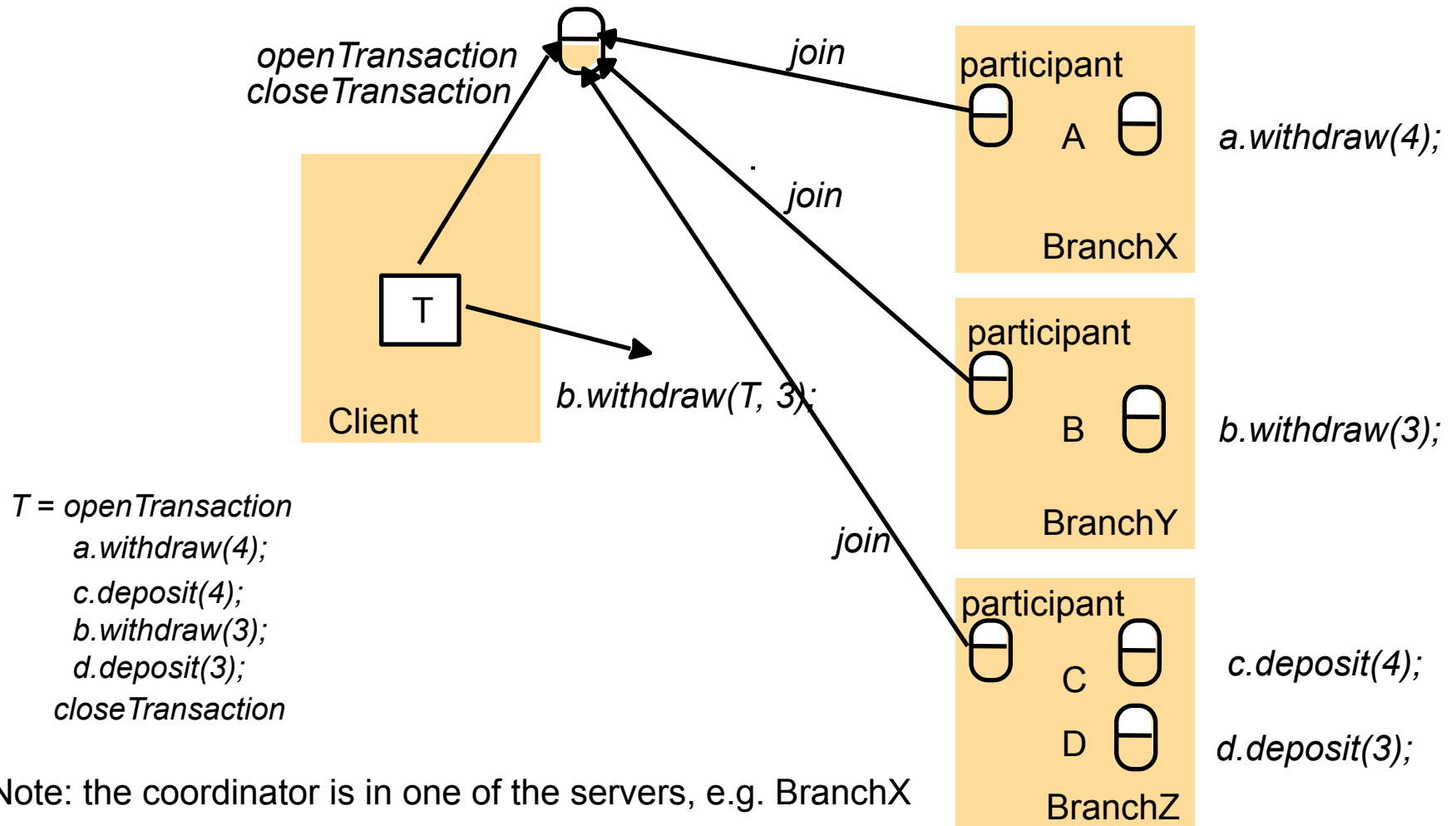


Figure 17.4

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.

Figure 17.5

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.

Figure 17.6
Communication in two-phase commit protocol

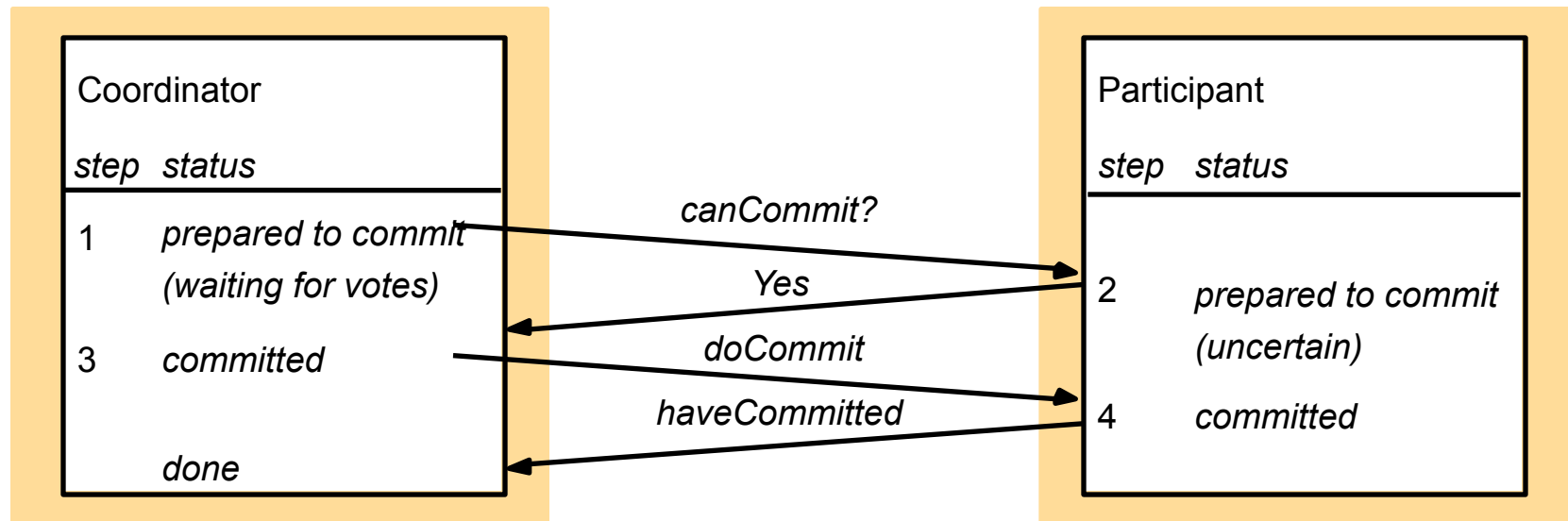


Figure 17.7

Operations in coordinator for nested transactions

openSubTransaction(trans) -> subTrans

Opens a new subtransaction whose parent is *trans* and returns a unique subtransaction identifier.

getStatus(trans)-> committed, aborted, provisional

Asks the coordinator to report on the status of the transaction *trans*. Returns values representing one of the following:
committed, aborted, provisional.

Complexity

- ◆ In case of no failures
 - ◆ Cost in messages is proportional to $3N$ (with N servers)
 - ◆ Cost in time: 3 rounds of messages
- ◆ In case of failures (server replacement)
 - ◆ guaranteed to complete eventually, but with no bound on the complexity
- ◆ Three phase commit protocol exist to avoid waiting for uncertain outcome

Figure 17.8
Transaction T decides whether to commit

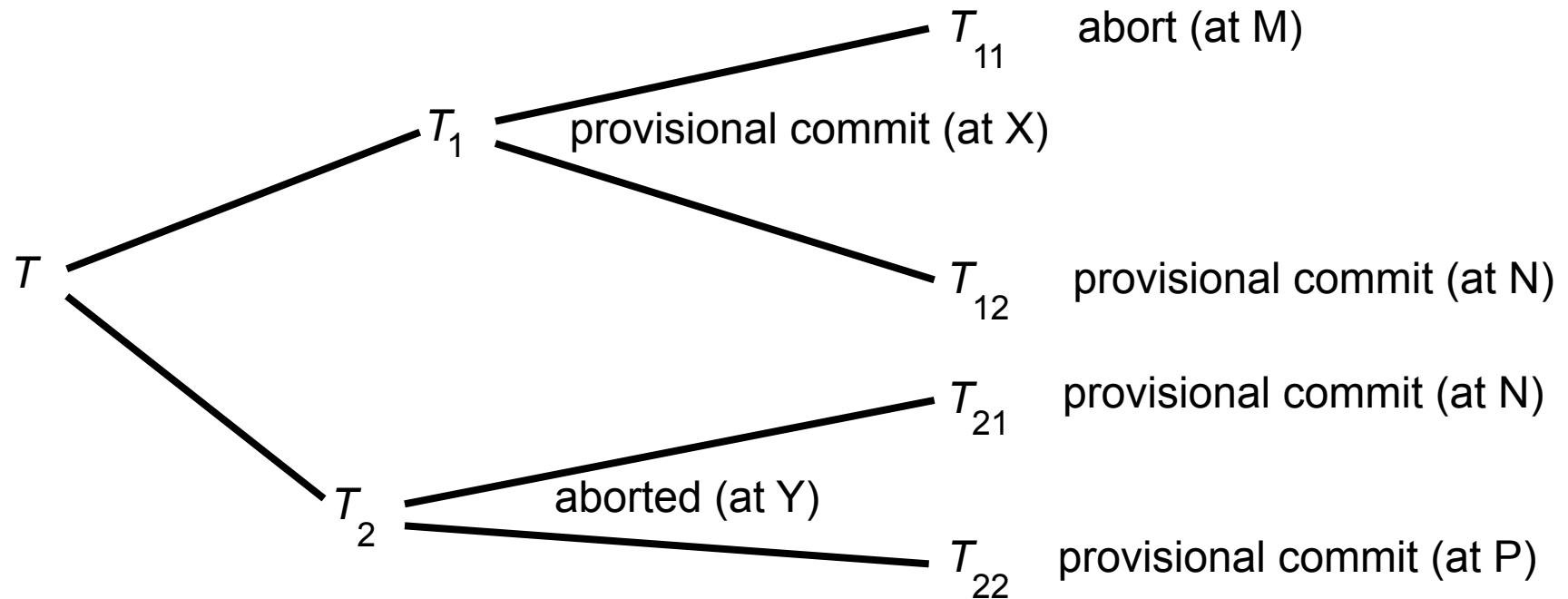


Figure 17.9
Information held by coordinators of nested transactions

<i>Coordinator of transaction</i>	<i>Child transactions</i>	<i>Participant</i>	<i>Provisional commit list</i>	<i>Abort list</i>
<i>T</i>	T ₁ , T ₂	yes	T ₁ , T ₁₂	T ₁₁ , T ₂
T ₁	T ₁₁ , T ₁₂	yes	T ₁ , T ₁₂	T ₁₁
T ₂	T ₂₁ , T ₂₂	no (aborted)		T ₂
T ₁₁		no (aborted)		T ₁₁
T ₁₂ , T ₂₁		T ₁₂ but not T ₂₁ *	T ₂₁ , T ₁₂	
T ₂₂		no (parent aborted)	T ₂₂	

*T₂₁'s parent has aborted

Figure 17.10

canCommit? for hierarchic two-phase commit protocol

canCommit?(trans, subTrans) -> Yes / No

Call a coordinator to ask coordinator of child subtransaction whether it can commit a subtransaction *subTrans*. The first argument *trans* is the transaction identifier of top-level transaction. Participant replies with its vote *Yes / No*.

Figure 17.11 *canCommit?* for flat two-phase commit protocol

canCommit?(trans, abortList) -> Yes / No

Call from coordinator to participant to ask whether it can commit a transaction. Participant replies with its vote *Yes / No*.

Other transaction models

- **ACID** is the traditional model (consistency)
- **BASE** often with NoSQL databases (availability, scalability)

CAP Theorem

- **B**asically **A**vailable: via replication
- **S**oft state: with possibly inconsistent values
- **E**ventually consistent: with reads possible before consistency

NoSQL (not only SQL) databases:

- ♦ key value, wide column, graph
- ♦ horizontal scaling
- ♦ partial replication