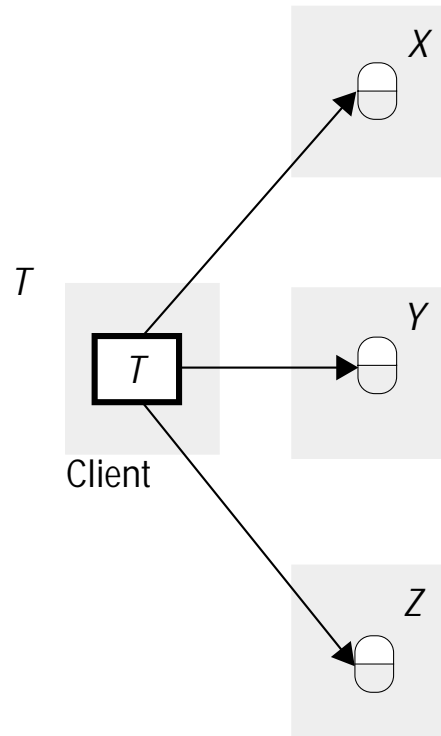


Figure 13.1 Distributed transactions

(a) Flat transaction



(b) Nested transactions

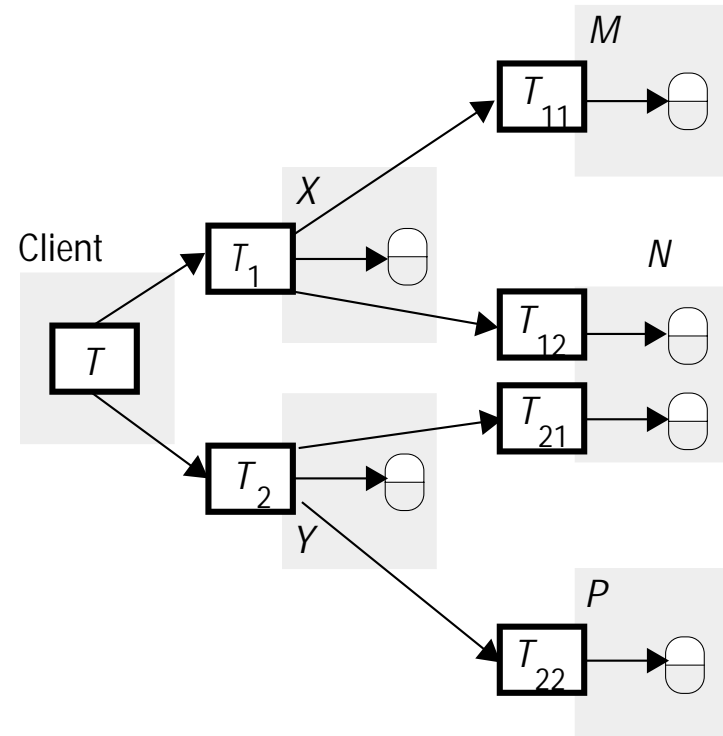


Figure 13.2 Nested banking transaction

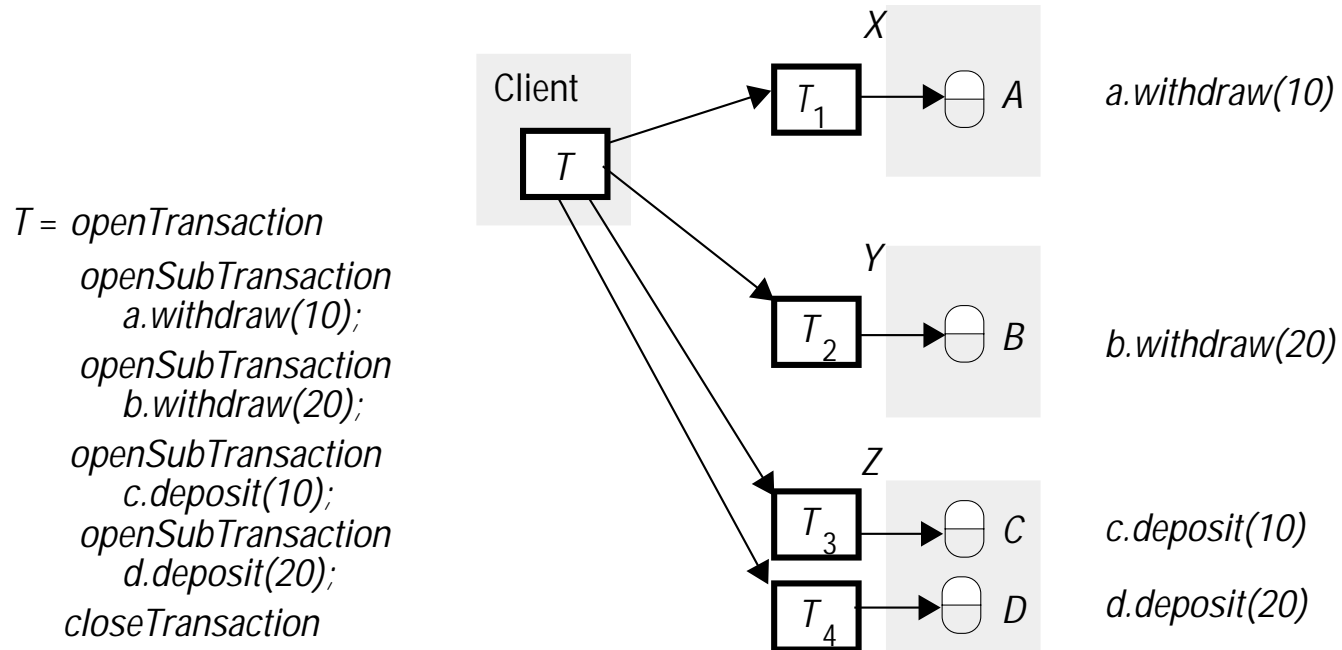


Figure 13.3 A distributed banking transaction

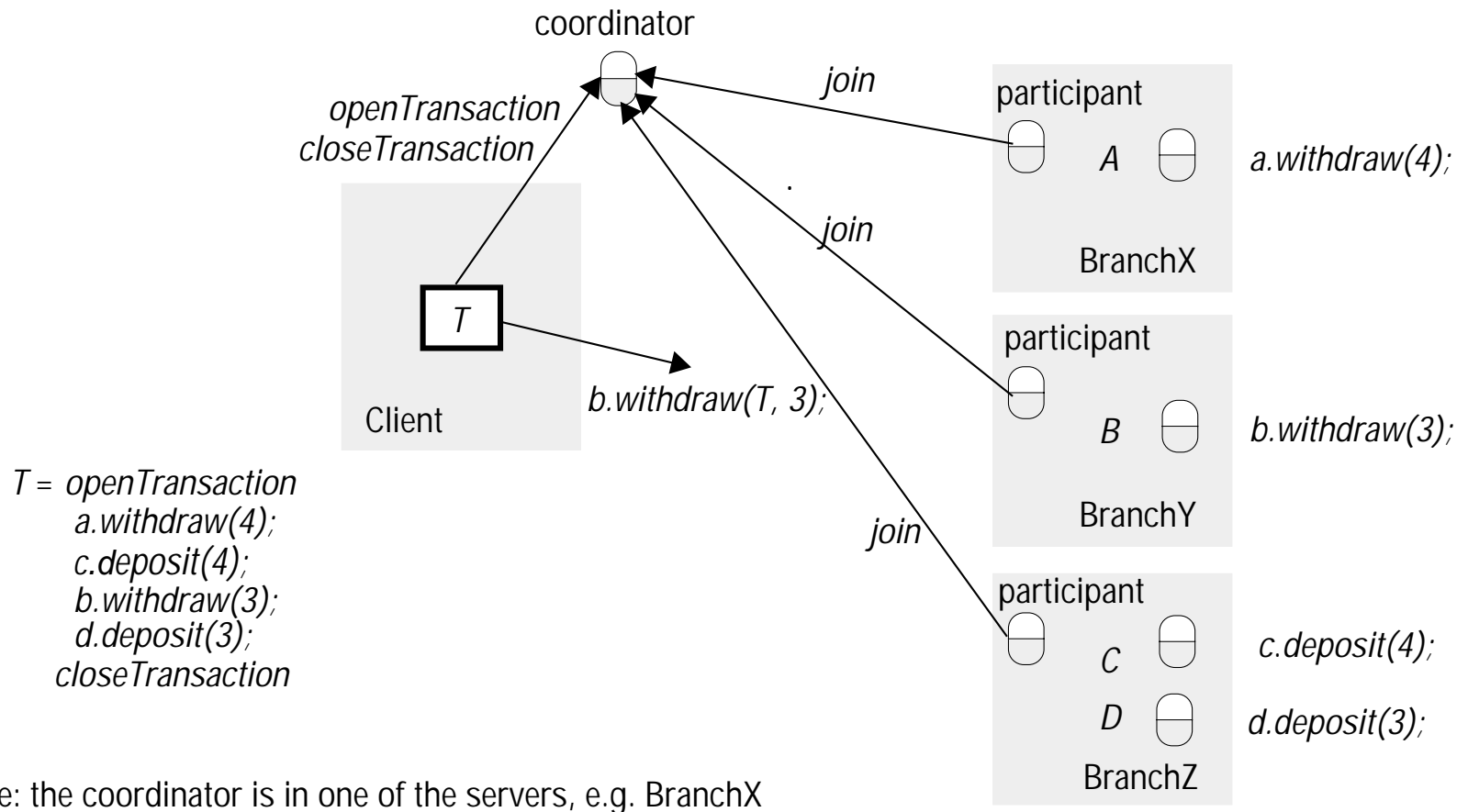


Figure 13.4 Operations for two-phase commit protocol

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 13.5 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 13.6 Communication in two-phase commit protocol

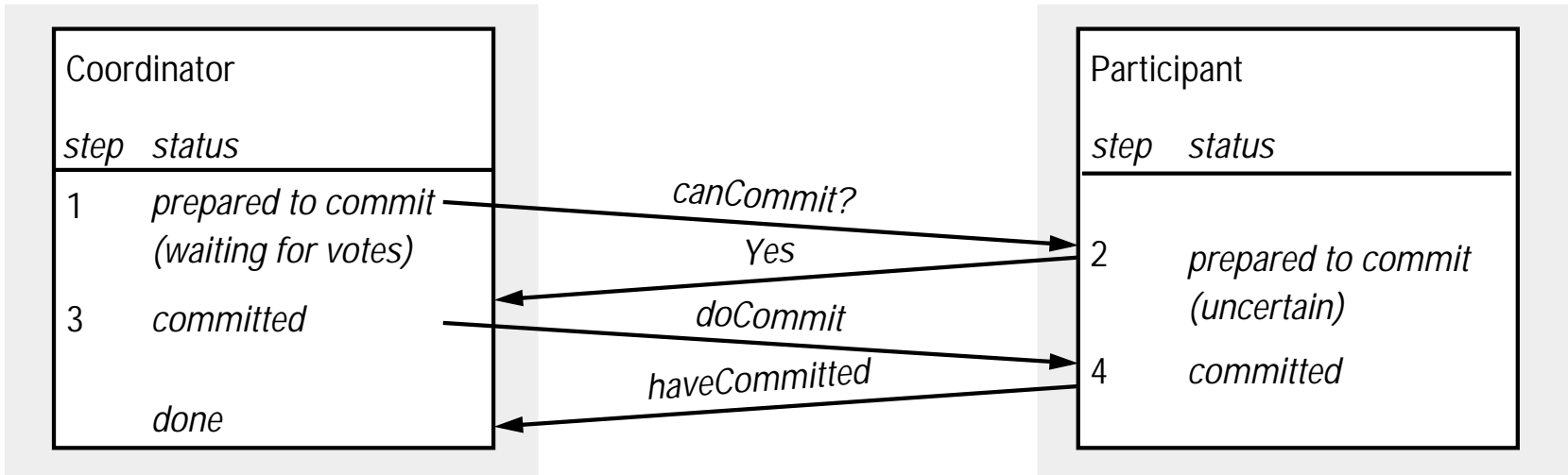


Figure 13.7 Operations in coordinator for nested transactions

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*.

Figure 13.8 Transaction T decides whether to commit

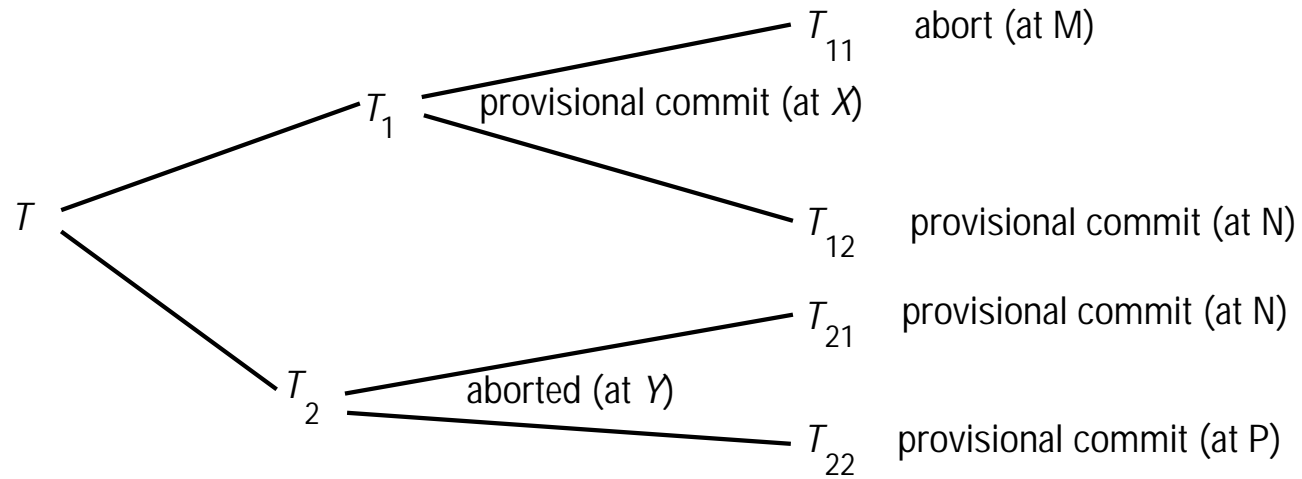


Figure 13.9 Information held by coordinators of nested transactions:

<i>Coordinator of transaction</i>	<i>Child transactions</i>	<i>Participant</i>	<i>Provisional commit list</i>	<i>Abort list</i>
T	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}	

Figure 13.10 *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 13.11 *canCommit?* for flat two-phase commit protocol

canCommit?(trans, abortList) → Yes / No

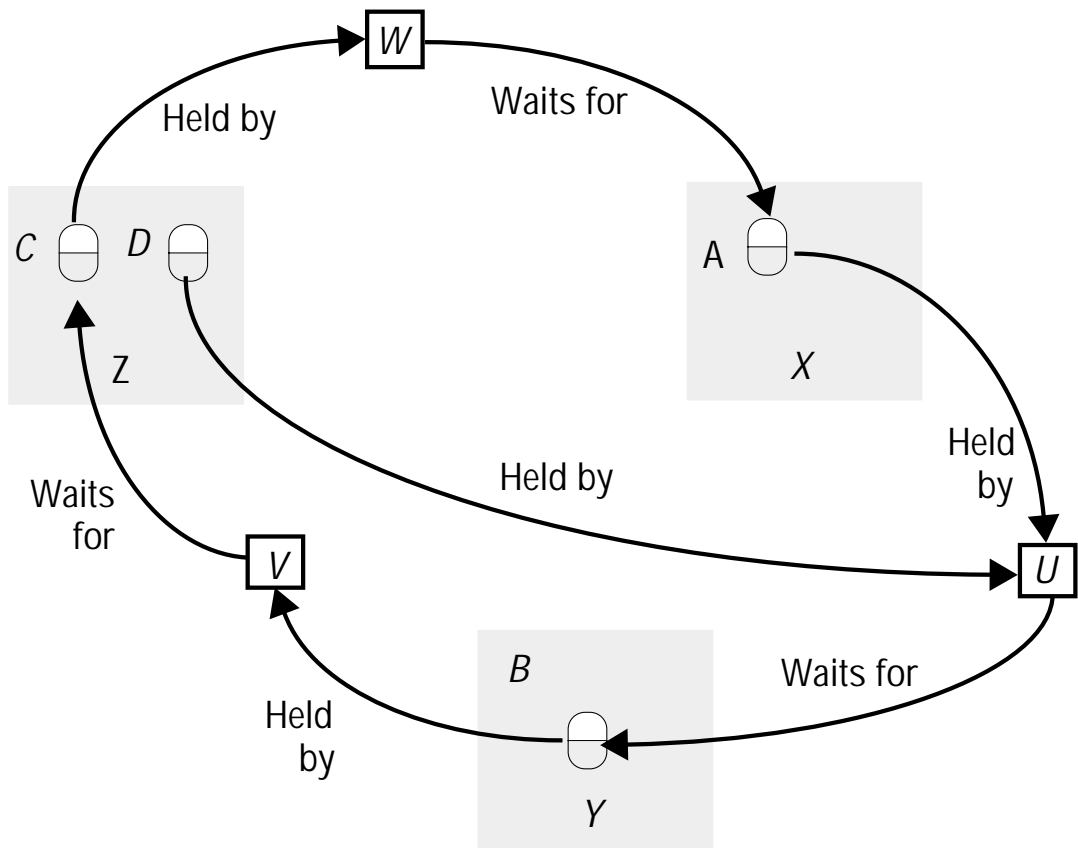
Call from coordinator to participant to ask whether it can commit a transaction. Participant replies with its vote *Yes / No*.

Figure 13.12 Interleavings of transactions U , V and W

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

Figure 13.13 Distributed deadlock

(a)



(b)

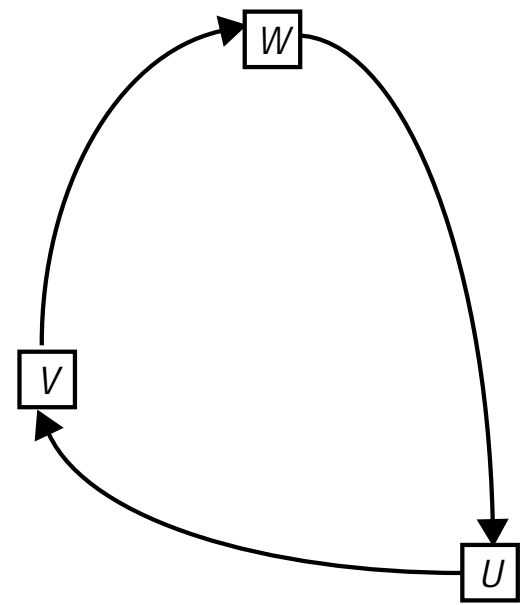
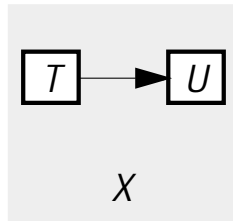
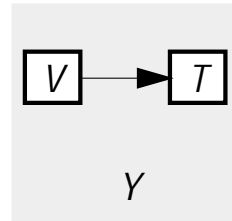


Figure 13.14 Local and global wait-for graphs

local wait-for graph



local wait-for graph



global deadlock detector

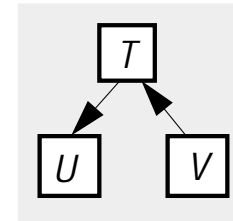


Figure 13.15 Probes transmitted to detect deadlock

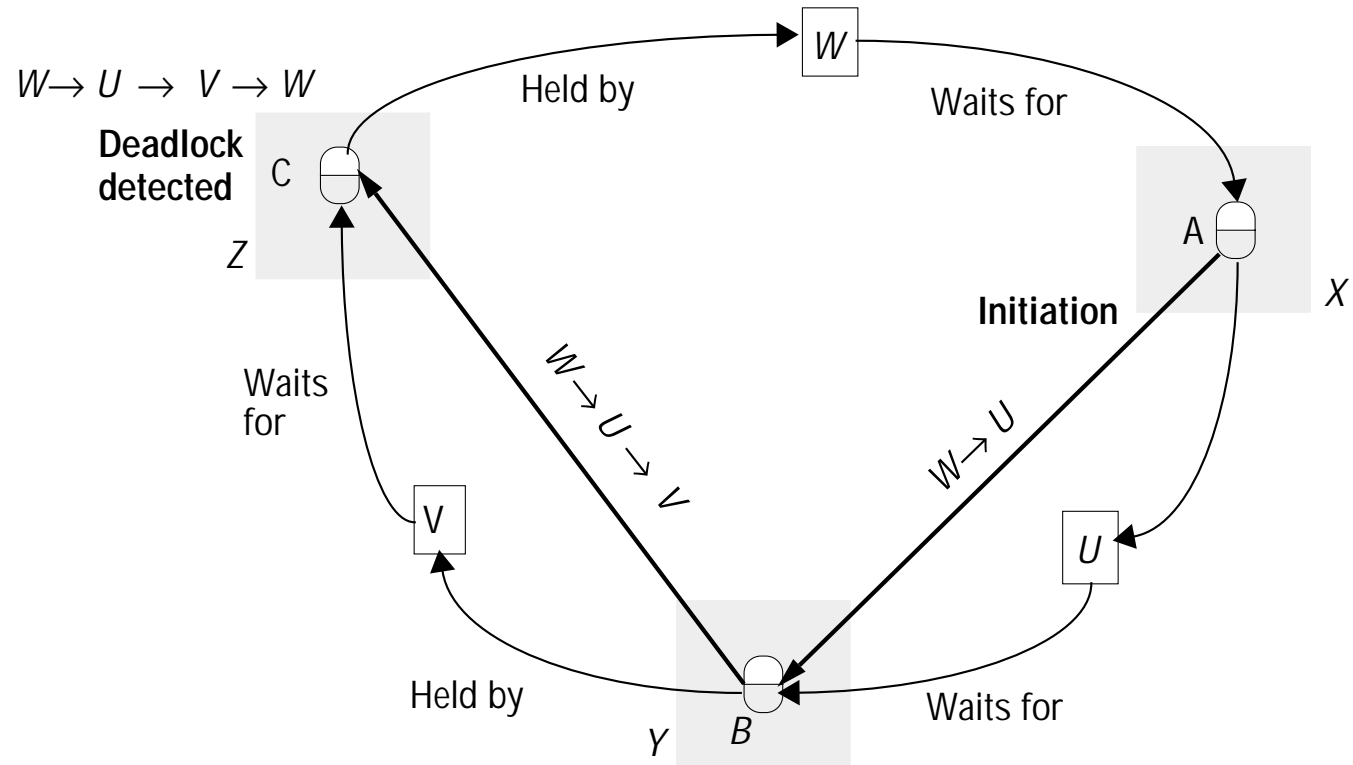
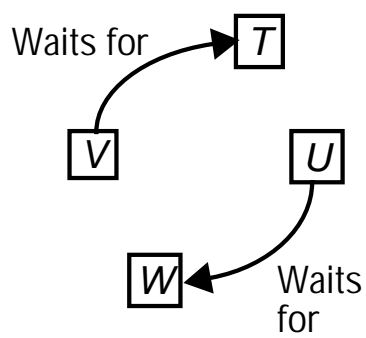
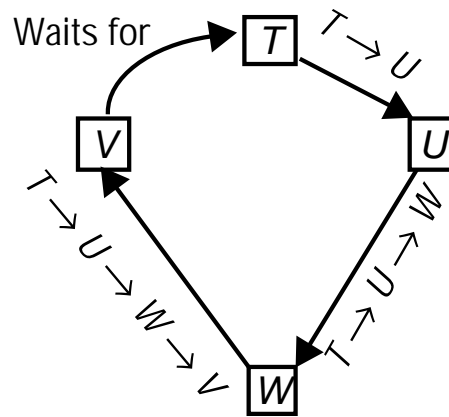


Figure 13.16 Two probes initiated

(a) initial situation



(b) detection initiated at object requested by T



(c) detection initiated at object requested by W

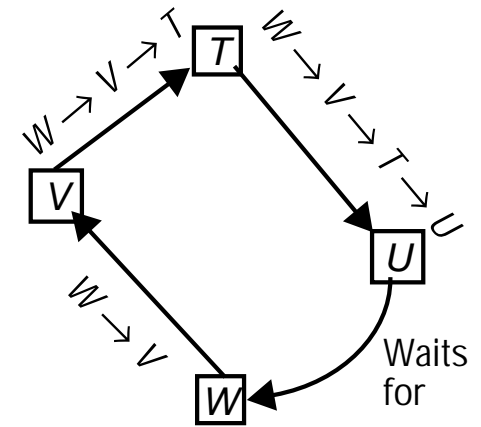


Figure 13.17 Probes travel downhill

(a) V stores probe when U starts waiting (b) Probe is forwarded when V starts waiting

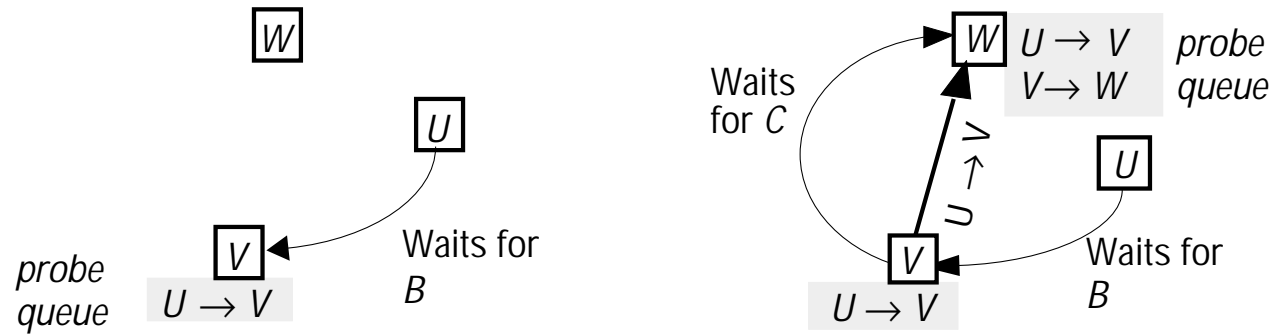


Figure 13.18 Types of entry in a recovery file

<i>Type of entry</i>	<i>Description of contents of entry</i>
Object	A value of an object.
Transaction status	Transaction identifier, transaction status (<i>prepared, committed, aborted</i>) – and other status values used for the two-phase commit protocol.
Intentions list	Transaction identifier and a sequence of intentions, each of which consists of <identifier of object>, <position in recovery file of value of object>.

Figure 13.19 Log for banking service

P_0			P_1		P_2	P_3	P_4	P_5	P_6	P_7
Object:A	Object:B	Object:C	Object:A	Object:B		Trans:T	Trans:T	Object:C	Object:B	Trans:U
100	200	300	80	220		prepared	committed	278	242	prepared
						$\langle A, P_1 \rangle$				$\langle C, P_5 \rangle$
						$\langle B, P_2 \rangle$				$\langle B, P_6 \rangle$
						P_0	P_3			P_4

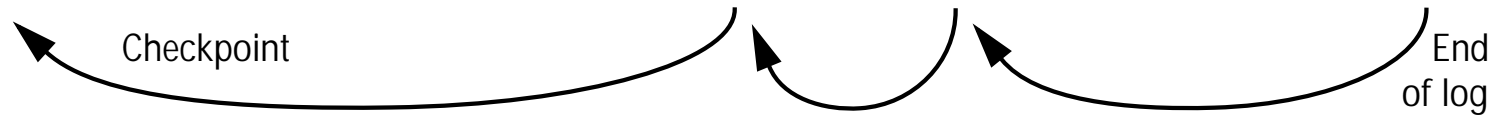


Figure 13.20 Shadow versions

		<i>Map at start</i>			<i>Map when T commits</i>			
		$A \rightarrow P_0$			$A \rightarrow P_1$			
		$B \rightarrow P_0'$			$B \rightarrow P_2$			
		$C \rightarrow P_0''$			$C \rightarrow P_0''$			
		P_0	P_0'	P_0''	P_1	P_2	P_3	P_4
<i>Version store</i>		100	200	300	80	220	278	242
		<i>Checkpoint</i>						

Figure 13.21 Log with entries relating to two-phase commit protocol

Trans: <i>T</i> prepared intentions list	Coord'r: <i>T</i> part'pant list: . . .	•	•	Trans: <i>T</i> committed	Trans: <i>U</i> prepared intentions list	•	•	Part'pant: <i>U</i> Coord'r: . . .	Trans: <i>U</i> uncertain	Trans: <i>U</i> committed
---	---	---	---	------------------------------	---	---	---	---------------------------------------	------------------------------	------------------------------

Figure 13.22 Recovery of the two-phase commit protocol

<i>Role</i>	<i>Status</i>	<i>Action of recovery manager</i>
Coordinator	<i>prepared</i>	No decision had been reached before the server failed. It sends <i>abortTransaction</i> to all the servers in the participant list and adds the transaction status <i>aborted</i> in its recovery file. Same action for state <i>aborted</i> . If there is no participant list, the participants will eventually timeout and abort the transaction.
Coordinator	<i>committed</i>	A decision to commit had been reached before the server failed. It sends a <i>doCommit</i> to all the participants in its participant list (in case it had not done so before) and resumes the two-phase protocol at step 4 (Fig 13.5).
Participant	<i>committed</i>	The participant sends a <i>haveCommitted</i> message to the coordinator (in case this was not done before it failed). This will allow the coordinator to discard information about this transaction at the next checkpoint.
Participant	<i>uncertain</i>	The participant failed before it knew the outcome of the transaction. It cannot determine the status of the transaction until the coordinator informs it of the decision. It will send a <i>getDecision</i> to the coordinator to determine the status of the transaction. When it receives the reply it will commit or abort accordingly.
Participant	<i>prepared</i>	The participant has not yet voted and can abort the transaction.
Coordinator	<i>done</i>	No action is required.

Figure 13.23 Nested transactions

