

Transactions

Distributed Software Systems

Transactions 1

Transactions

- ❑ A *transaction* is a sequence of server operations that is guaranteed by the server to be atomic in the presence of multiple clients and server crashes.
 - Free from interference by operations being performed on behalf of other concurrent clients
 - Either all of the operations must be completed successfully or they must have no effect at all in the presence of server crashes
- ❑ Properties(ACID)
 - Atomicity – all or nothing
 - Consistency – a transaction must move the system from one consistent state to another
 - Isolation – intermediate effects must be invisible to other clients
 - Durability
- ❑ Concepts: commit, abort

Transactions 2

Bank Operations

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

```

A client's banking transaction

Transaction T:

```

a.withdraw(100);
b.deposit(100);
c.withdraw(200);
b.deposit(200);

```

Transactions 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.

```

Transactions 4

Transaction life histories

Successful	Aborted by client	Aborted by server
openTransaction	openTransaction	openTransaction
operation	operation	operation
operation	operation	operation
⋮	⋮	⋮
operation	operation	operation ERROR reported to client
closeTransaction	abortTransaction	

server aborts transaction →

Transactions 5

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 6

Lost Update Problem - 1

Consider the following transactions T and U with the initial account balances for a, b, and c as 100, 200, and 300 respectively.

T: *balance = b.getbalance()*
*b.setbalance(balance*1.1)*
a.withdrawn(balance/10)

U: *balance = b.getbalance()*
*b.setbalance(balance*1.1)*
c.withdrawn(balance/10)

If we execute T followed by U, we get balances of 80, 242 and 278

If we execute U followed by T, we get balances of 78, 242 and 280

Transactions 7

Lost Update Problem - 2

TransactionT:	TransactionU:
<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

In the above, we get balances of 80, 220 and 280 – U's update of b was lost

Transactions 8

Lost Update Problem - 3

T: <i>balance = b.getbalance()</i> <i>b.setbalance(balance*1.1)</i> <i>a.withdrawn(balance/10)</i>	U: <i>balance = b.getbalance()</i> <i>b.setbalance(balance*1.1)</i> <i>c.withdrawn(balance/10)</i>
---	---

balance = b.getbalance()
*b.setbalance(balance*1.1)*

a.withdrawn(balance/10)

balance = b.getbalance()
*b.setbalance(balance*1.1)*

c.withdrawn(balance/10)

We get balances of 80, 242 and 278, same as T followed by U

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		

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>total = a.getBalance()</i>	\$100
<i>b.deposit(100)</i>	\$300	<i>total = total+b.getBalance()</i>	\$400
		<i>total = total+c.getBalance()</i>	
		...	

Serializability

BEGIN_TRANSACTION x = 0; x = x + 1; END_TRANSACTION	BEGIN_TRANSACTION x = 0; x = x + 2; END_TRANSACTION	BEGIN_TRANSACTION x = 0; x = x + 3; END_TRANSACTION
(a)	(b)	(c)

Schedule	Order of Operations	Legal
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 T₁, T₂, and T₃

d) Possible schedules

Read and write operation conflict rules

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

Transactions 13

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

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

Transactions 14

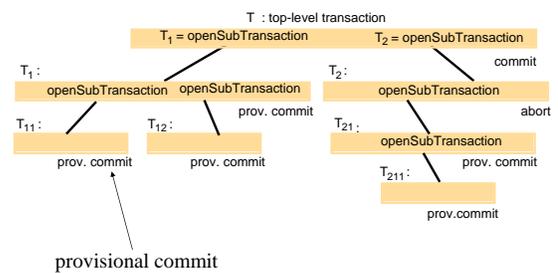
A dirty read when transaction T aborts

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

uses result of uncommitted transaction!

Transactions 15

Nested transactions

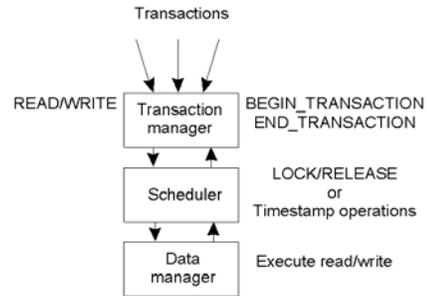


Transactions 16

Committing Nested Transactions

- ❑ A transaction may commit or abort only after its child transactions have completed
- ❑ When a sub-transaction completes, it makes an independent decision either to commit provisionally or to abort. Its decision to abort is final.
- ❑ When a parent aborts, all of its sub-transactions are aborted.
- ❑ When a sub-transaction aborts, the parent can decide whether to abort or not.
- ❑ If a top-level transaction commits, then all of the sub-transactions that have provisionally committed can commit too, provided that non of their ancestors has aborted.

Concurrency Control



General organization of managers for handling transactions.

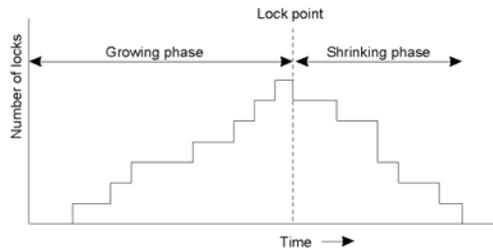
Schemes for Concurrency control

- ❑ Locking
 - Server attempts to gain an exclusive 'lock' that is about to be used by one of its operations in a transaction.
 - Can use different lock types (read/write for example)
 - Two-phase locking
- ❑ Optimistic concurrency control
- ❑ Time-stamp based concurrency control

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 lock on B
<i>b.setBalance(bal*1.1)</i>		...	lock B
<i>a.withdraw(bal/10)</i>	lock A	<i>c.withdraw(bal/10)</i>	lock C
<i>closeTransaction</i>	unlock A, B	<i>closeTransaction</i>	unlock B, C

Two-Phase Locking (1)

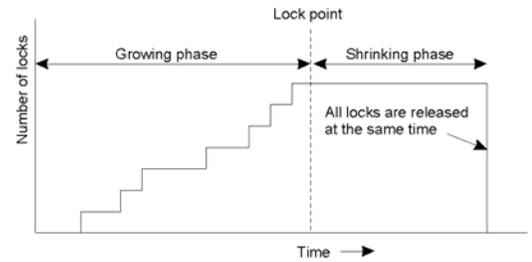


In two-phase locking, a transaction is not allowed to acquire any new locks after it has released a lock

Transactions 21

Strict Two-Phase Locking (2)

□ Strict two-phase locking.



Transactions 22

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 23

Lock compatibility

For one object		Lock requested	
		read	write
Lock already set	none	OK	OK
	read	OK	wait
	write	wait	wait

Operation Conflict rules:

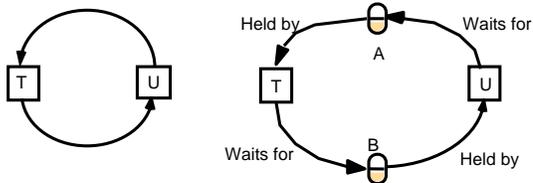
1. If a transaction T has already performed a read operation on a particular object, then a concurrent transaction U must not write that object until T commits or aborts
2. If a transaction T has already performed a read operation on a particular object, then a concurrent transaction U must not read or write that object until T commits or aborts

Transactions 24

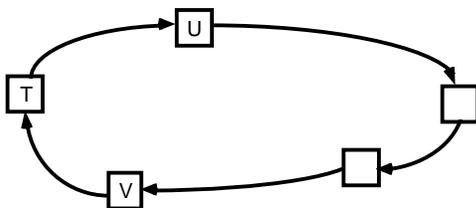
Deadlock with write locks

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

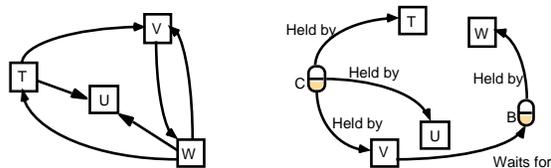
The wait-for graph



A cycle in a wait-for graph



Another wait-for graph



Dealing with Deadlock in two-phase locking

- Deadlock prevention
 - Acquire all needed locks in a single atomic operation
 - Acquire locks in a particular order
- Deadlock detection
 - Keep graph of locks held. Check for cycles periodically or each time an edge is added
 - Cycles can be eliminated by aborting transactions
- Timeouts
 - Aborting transactions when time expires

Transactions 29

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

Transactions 30

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

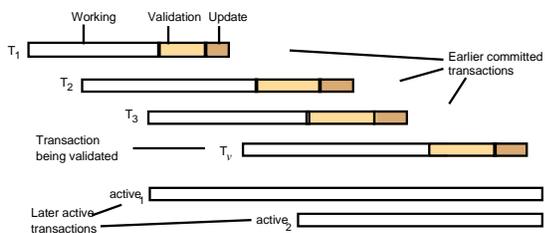
Transactions 31

Optimistic Concurrency Control

- Three phases:
 - **Working Phase** – transactions read and write private copies of objects (most recently committed)
 - **Validation Phase** – Once transaction is done, the transaction is validated to establish whether or not its operations on objects conflict with operations of other transactions on the same object. If not conflict, can commit; else some form of conflict resolution is needed and the transaction may abort.
 - **Update Phase** – if commit, private copies are used to make permanent change.

Transactions 32

Validation of transactions



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
write	read	1. T_i must not read objects written by T_v
read	write	2. T_v must not read objects written by T_i
write	write	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

Validation of Transactions

Backward validation of transaction T_v

```

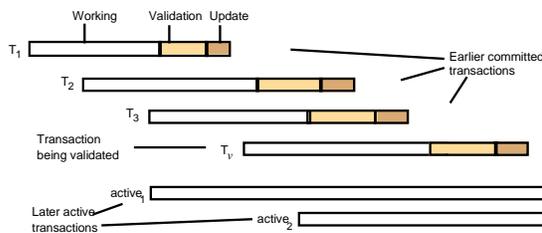
boolean valid = true;
for (int  $T_i = startTn+1; T_i \leq finishTn; T_i++$ ) {
    if (read set of  $T_v$  intersects write set of  $T_i$ ) valid = false;
}
    
```

Forward validation of transaction T_v

```

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

Validation of transactions



Schemes for Concurrency control

- ❑ Locking
- ❑ Optimistic concurrency control
- ❑ Time-stamp 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
 - A request to write to an object is valid only if that object was last read and written by earlier (wrt timestamps) transactions.
 - A request to read an object is valid only if that object was last written by an earlier transaction.

Transactions 37

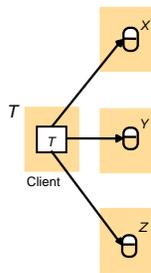
Distributed Transactions

- ❑ Motivation
 - Provide distributed atomic operations at multiple servers that maintain shared data for clients
 - Provide recoverability from server crashes
- ❑ Properties
 - Atomicity, Consistency, Isolation, Durability (ACID)
- ❑ Concepts: commit, abort, distributed commit

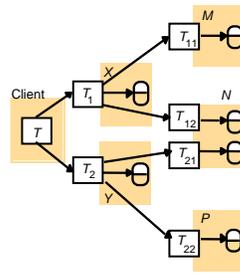
Transactions 38

Distributed Transactions

(a) Flat transaction

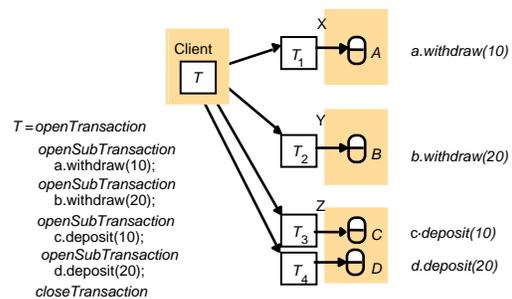


(b) Nested transactions



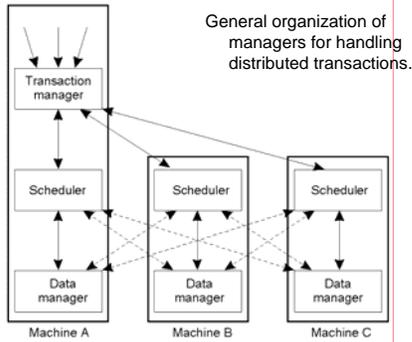
Transactions 39

Nested banking transaction



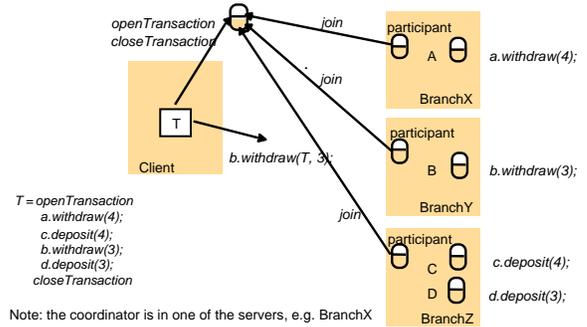
Transactions 40

Concurrency Control for Distributed Transactions



Transactions 41

A distributed banking transaction



Transactions 42

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

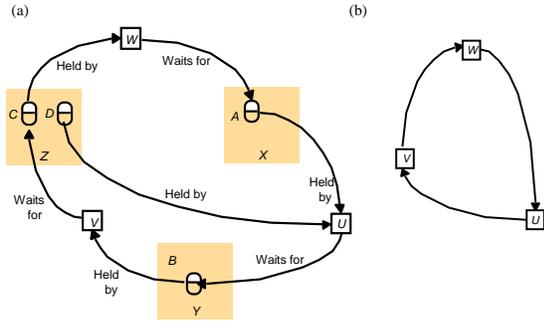
Transactions 43

Interleavings of transactions U, V and W

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

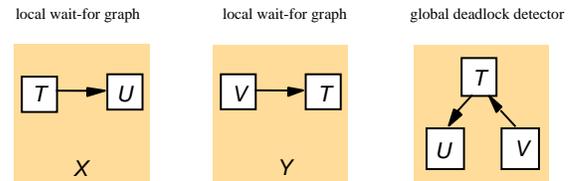
Transactions 44

Distributed deadlock



Transactions 45

Local and global wait-for graphs



Transactions 46

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 47

The two-phase commit protocol - 1

Phase 1 (voting phase):

1. The coordinator sends a *canCommit?* (*VOTE_REQUEST*) request to each of the participants in the transaction.
2. When a participant receives a *canCommit?* request it replies with its vote *Yes* (*VOTE_COMMIT*) or *No* (*VOTE_ABORT*) 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.

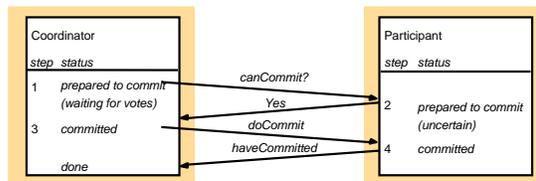
Transactions 48

The two-phase commit protocol - 2

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* (*GLOBAL_COMMIT*) request to each of the participants.
 - (b) Otherwise the coordinator decides to abort the transaction and sends *doAbort* (*GLOBAL_ABORT*) 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.

Communication in two-phase commit protocol



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.

Two-Phase Commit protocol - 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

Two-Phase Commit protocol - 4

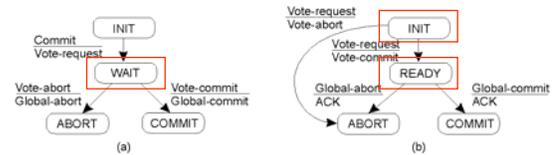
Steps taken by participant process in 2PC.

```

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 53

Two-Phase Commit protocol - 5



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

If a failure occurs during a 'blocking' state (red boxes), there needs to be a recovery mechanism.

Transactions 54

Two Phase Commit Protocol - 6

Recovery

- 'Wait' in Coordinator – use a time-out mechanism to detect participant crashes. Send GLOBAL_ABORT
- 'Init' in Participant – Can also use a time-out and send VOTE_ABORT
- 'Ready' in Participant P – abort is not an option (since already voted to COMMIT and so coordinator might eventually send GLOBAL_COMMIT). Can contact another participant Q and choose an action based on its state.

State of Q	Action by P
COMMIT	Transition to COMMIT
ABORT	Transition to ABORT
INIT	Both P and Q transition to ABORT (Q sends VOTE_ABORT)
READY	Contact more participants. If all participants are 'READY', must wait for coordinator to recover

Transactions 55

Two-Phase Commit protocol - 7

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 protocol - 1

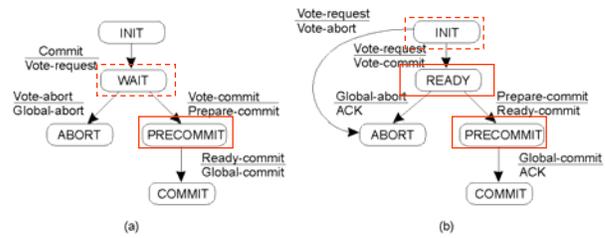
❑ Problem with 2PC

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

❑ Three Phase Commit3PC

- 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

Three-Phase Commit protocol - 2



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

Three Phase Commit Protocol - 3

Recovery

- ❑ 'Wait' in Coordinator – same
- ❑ 'Init' in Participant – same
- ❑ 'PreCommit' in Coordinator – Some participant has crashed but we know it wanted to commit. GLOBAL_COMMIT the application knowing that once the participant recovers, it will commit.
- ❑ 'Ready' or 'PreCommit' in Participant P – (i.e. P has voted to COMMIT)

State of Q	Action by P
PRECOMMIT	Transition to PRECOMMIT. If all participants in PRECOMMIT, can COMMIT the transaction
ABORT	Transition to ABORT
INIT	Both P (in READY) and Q transition to ABORT (Q sends VOTE_ABORT)
READY	Contact more participants. If can contact a majority and they are in 'Ready', then ABORT the transaction. If the participants contacted in 'PreCommit' it is safe to COMMIT the transaction

Note: if any participant is in state PRECOMMIT, it is impossible for any other participant to be in any state other than READY or PRECOMMIT.