

Last Class

- Distributed Snapshots
 - Termination detection
- Election algorithms
 - Bully
 - Ring



Today: Still More Canonical Problems

- Distributed synchronization and mutual exclusion
- Distributed transactions



Distributed Synchronization

- Distributed system with multiple processes may need to share data or access shared data structures
 - Use critical sections with mutual exclusion
- Single process with multiple threads
 - Semaphores, locks, monitors
- How do you do this for multiple processes in a distributed system?
 - Processes may be running on different machines
- Solution: lock mechanism for a distributed environment
 - Can be centralized or distributed

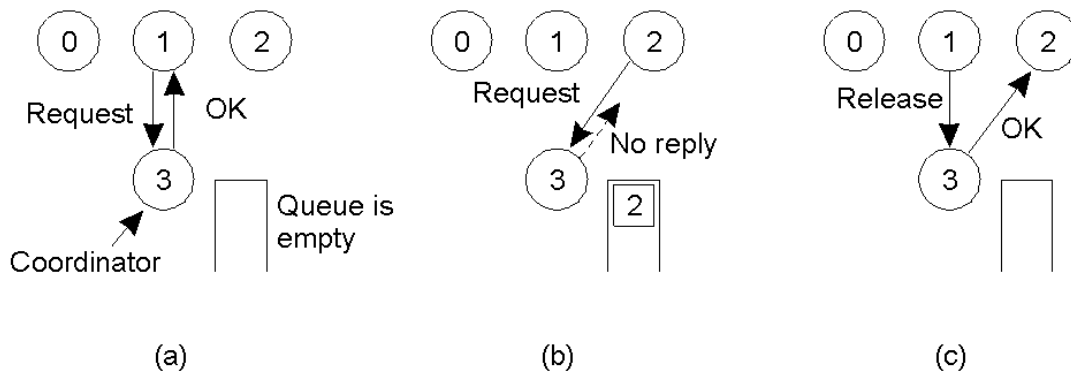


Centralized Mutual Exclusion

- Assume processes are numbered
- One process is elected coordinator (highest ID process)
- Every process needs to check with coordinator before entering the critical section
- To obtain exclusive access: send request, await reply
- To release: send release message
- Coordinator:
 - Receive *request*: if available and queue empty, send grant; if not, queue request
 - Receive *release*: remove next request from queue and send grant



Mutual Exclusion: A Centralized Algorithm



- Process 1 asks the coordinator for permission to enter a critical region. Permission is granted
- Process 2 then asks permission to enter the same critical region. The coordinator does not reply.
- When process 1 exits the critical region, it tells the coordinator, when then replies to 2



Properties

- Simulates centralized lock using blocking calls
- Fair: requests are granted the lock in the order they were received
- Simple: three messages per use of a critical section (request, grant, release)
- Shortcomings:
 - Single point of failure
 - How do you detect a dead coordinator?
 - A process can not distinguish between “lock in use” from a dead coordinator
 - No response from coordinator in either case
 - Performance bottleneck in large distributed systems

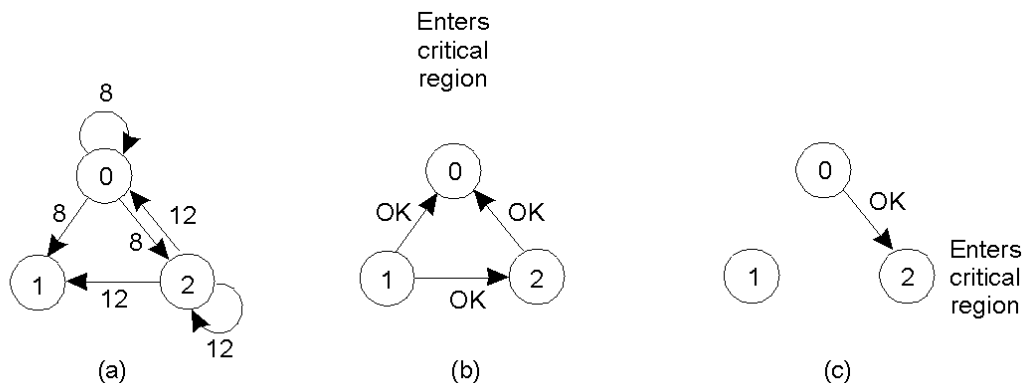


Distributed Algorithm

- [Ricart and Agrawala]: needs $2(n-1)$ messages
- Based on event ordering and time stamps
 - Assumes total ordering of events in the system (Lamport's clock)
- Process k enters critical section as follows
 - Generate new time stamp $TS_k = TS_k + 1$
 - Send $request(k, TS_k)$ all other $n-1$ processes
 - Wait until $reply(j)$ received from all other processes
 - Enter critical section
- Upon receiving a $request$ message, process j
 - Sends $reply$ if no contention
 - If already in critical section, does not reply, queue request
 - If wants to enter, compare TS_j with TS_k and send reply if $TS_k < TS_j$, else queue



A Distributed Algorithm



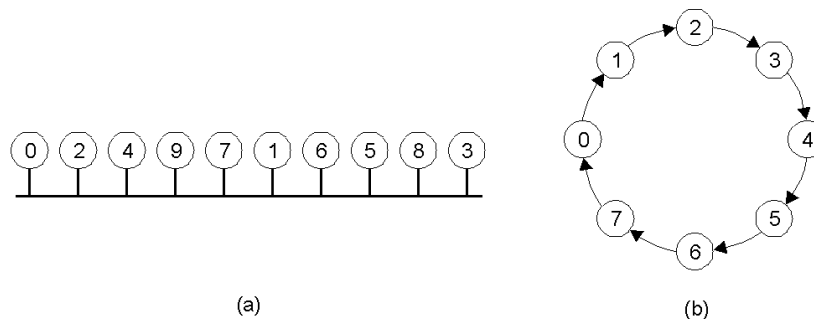
- Two processes want to enter the same critical region at the same moment.
- Process 0 has the lowest timestamp, so it wins.
- When process 0 is done, it sends an OK also, so 2 can now enter the critical region.



Properties

- Fully decentralized
- N points of failure!
- All processes are involved in all decisions
 - Any overloaded process can become a bottleneck

A Token Ring Algorithm



- An unordered group of processes on a network.
 - A logical ring constructed in software.
- Use a token to arbitrate access to critical section
 - Must wait for token before entering CS
 - Pass the token to neighbor once done or if not interested
 - Detecting token loss in non-trivial

Comparison

Algorithm	Messages per entry/exit	Delay before entry (in message times)	Problems
Centralized	3	2	Coordinator crash
Distributed	$2(n - 1)$	$2(n - 1)$	Crash of any process
Token ring	1 to ∞	0 to $n - 1$	Lost token, process crash

- A comparison of three mutual exclusion algorithms.

Transactions

- Transactions provide higher level mechanism for *atomicity* of processing in distributed systems
 - Have their origins in databases
- Banking example: Three accounts A:\$100, B:\$200, C:\$300
 - Client 1: transfer \$4 from A to B
 - Client 2: transfer \$3 from C to B
- Result can be inconsistent unless certain properties are imposed on the accesses

Client 1	Client 2
Read A: \$100	
Write A: \$96	
	Read C: \$300
	Write C:\$297
Read B: \$200	
	Read B: \$200
	Write B:\$203
Write B:\$204	

ACID Properties

- *Atomic*: all or nothing
- *Consistent*: transaction takes system from one consistent state to another
- *Isolated*: Immediate effects are not visible to other (serializable)
- *Durable*: Changes are permanent once transaction completes (commits)

Client 1	Client 2
Read A: \$100	
Write A: \$96	
Read B: \$200	
Write B:\$204	
	Read C: \$300
	Write C:\$297
	Read B: \$204
	Write B:\$207

Transaction Primitives

Primitive	Description
BEGIN_TRANSACTION	Make the start of a transaction
END_TRANSACTION	Terminate the transaction and try to commit
ABORT_TRANSACTION	Kill the transaction and restore the old values
READ	Read data from a file, a table, or otherwise
WRITE	Write data to a file, a table, or otherwise

Example: airline reservation

Begin_transaction

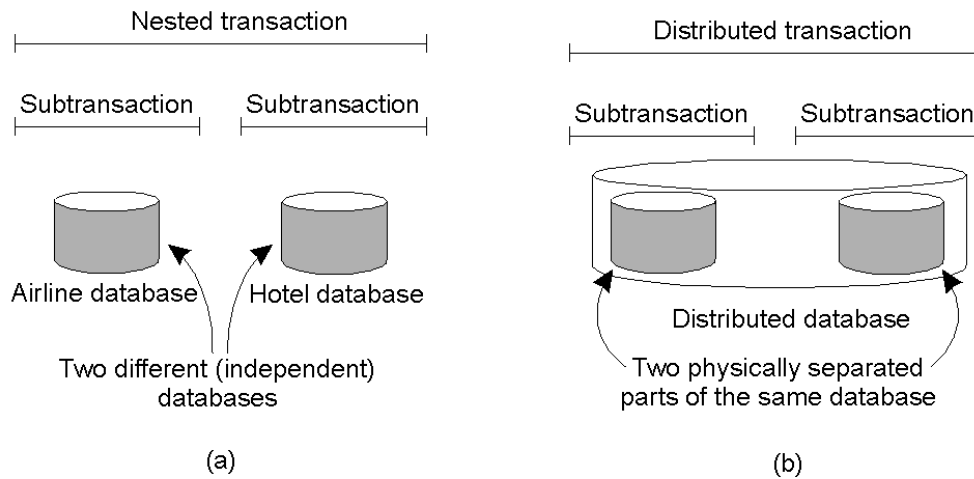
if(reserve(NY,Paris)==full) Abort_transaction

if(reserve(Paris,Athens)==full)Abort_transaction

if(reserve(Athens,Delhi)==full) Abort_transaction

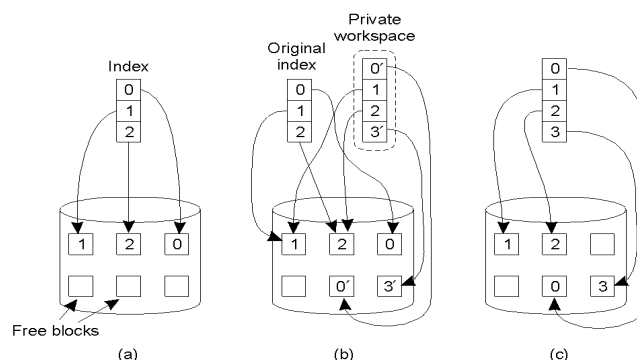
End_transaction

Distributed Transactions



Implementation: Private Workspace

- Each transaction get copies of all files, objects
- Can optimize for reads by not making copies
- Can optimize for writes by copying only what is required
- Commit requires making local workspace global



Option 2: Write-ahead Logs

- *In-place updates*: transaction makes changes *directly* to all files/objects
- *Write-ahead log*: prior to making change, transaction writes to log on *stable storage*
 - Transaction ID, block number, original value, new value
- Force logs on commit
- If abort, read log records and undo changes [*rollback*]
- Log can be used to rerun transaction after failure

- Both workspaces and logs work for distributed transactions
- Commit needs to be *atomic* [will return to this issue in Ch. 7]



Writeahead Log Example

x = 0;	Log	Log	Log
y = 0;			
BEGIN_TRANSACTION;			
x = x + 1;	[x = 0 / 1]	[x = 0 / 1]	[x = 0 / 1]
y = y + 2		[y = 0/2]	[y = 0/2]
x = y * y;			[x = 1/4]
END_TRANSACTION;			
(a)	(b)	(c)	(d)

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

