

Chapter 7 Consistency And Replication

- Data are replicated to increase the reliability of a system.
- Replication for performance
- Scaling in numbers
- Scaling in geographical area
- Caveat
 - Gain in performance
 - Cost of increased bandwidth for maintaining replication

Reasons for Replication

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

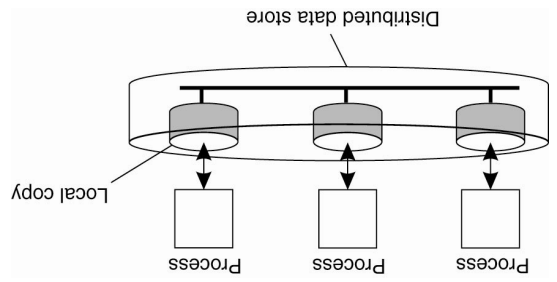


Figure 7-1. The general organization of a logical replicated across multiple processes.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Continuous Consistency Models

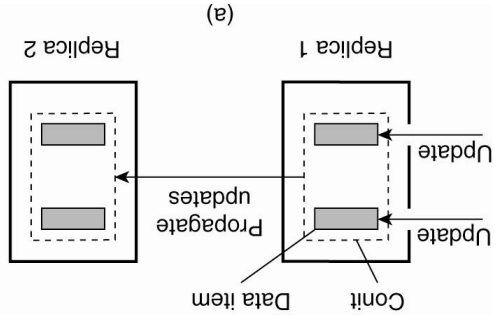


Figure 7-3. Choosing the appropriate granularity for a cont. (a) Two updates lead to update propagation.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Continuous Consistency (1)

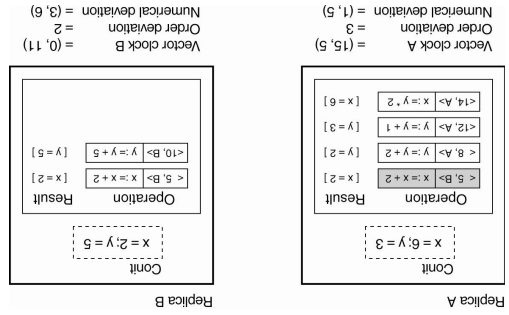


Figure 7-2. An example of keeping track of consistency deviations [adapted from (Yu and Vahdat, 2002)].

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Continuous Consistency (3)

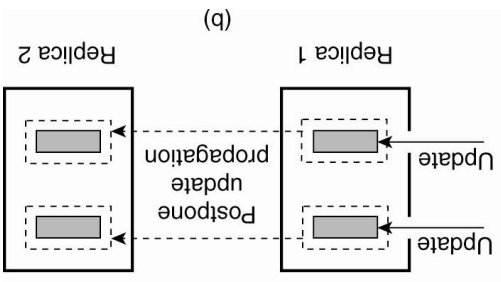


Figure 7-3. Choosing the appropriate granularity for a cont. (b) No update propagation is needed (yet).

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Sequential Consistency (1)

P1: W(x)a

P2: R(x)NIL R(x)a

Figure 7-4. Behavior of two processes operating on the same data item. The horizontal axis is time.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Sequential Consistency (3)

P1: W(x)a

P2: W(x)b

P3: R(x)b R(x)a

P4: R(x)b R(x)a

P1: W(x)a

P2: W(x)b

P3: R(x)b R(x)a

P4: R(x)a R(x)b

Figure 7-5. (a) A sequentially consistent data store. (b) A data store that is not sequentially consistent.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Sequential Consistency (5)

(a) P1: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P2: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P3: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P4: print(x, y);

Signature: 001011
Prints: 101011
Signature: 101011

(b) P1: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P2: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P3: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P4: print(x, y);

Signature: 010111
Prints: 010111
Signature: 110101

(c) P1: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P2: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P3: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P4: print(x, y);

Signature: 111111
Prints: 111111
Signature: 111111

(d) P1: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P2: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P3: print(x, z); y ← 1; y ← 1; x ← 1; print(x, y);

P4: print(x, y);

Signature: 111111
Prints: 111111
Signature: 111111

Figure 7-7. Four valid execution sequences for the processes of Fig. 7-6. The vertical axis is time.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Sequential Consistency (2)

A data store is sequentially consistent when: The result of any execution is the same as if the (read and write) operations by all processes on the data store ... were executed in some sequential order and ... the operations of each individual process appear ... in this sequence

- in the order specified by its program.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Sequential Consistency (4)

Process P1	x ← 1; y ← 1; print(x, z);	Process P2	y ← 1; print(x, z);	Process P3	z ← 1; print(x, y);
-------------------	----------------------------------	-------------------	------------------------	-------------------	------------------------

Figure 7-6. Three concurrently-executing processes.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Causal Consistency (1)

For a data store to be considered causally consistent, it is necessary that the store obeys the following condition:
Writes that are potentially causally related ... must be seen by all processes in the same order.
Concurrent writes ... may be seen in a different order on different machines.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Figure 7-10. A valid event sequence for entry consistency.

P1:	Acq(Lx) W(x)a Acq(Ly) W(y)b Rel(Lx) Rel(Ly)
P2:	Acq(Lx) R(x)a R(y) NIL
P3:	Acq(Ly) R(y)b

Grouping Operations (2)

Figure 7-9. (b) A correct sequence of events in a causally-consistent store.

P1:	W(x)a
P2:	W(x)b
P3:	R(x)b R(x)a
P4:	R(x)a R(x)b

(b)

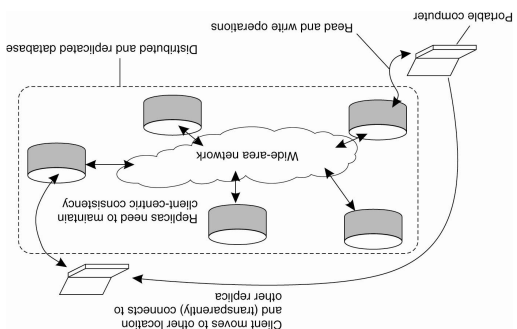
Causal Consistency (4)

Figure 7-8. This sequence is allowed with a causally-consistent store, but not with a sequentially consistent store.

P1:	W(x)a
P2:	R(x)a W(x)b
P3:	R(x)a R(x)c
P4:	R(x)a R(x)b R(x)c

Causal Consistency (2)

Figure 7-11. The principle of a mobile user accessing different replicas of a distributed database.



Eventual Consistency

- An acquire access of a synchronization variable, not allowed to perform until all updates to guarded shared data have been performed with respect to that process.
 - Before exclusive mode access to synchronization variable by other process may hold synchronization variable, not even in nonexclusive mode.
 - After exclusive mode access to synchronization variable has been performed, any other process' next nonexclusive mode access to that synchronization variable may not be performed until it has performed with respect to that variable's owner.
- Necessary criteria for correct synchronization:

Grouping Operations (1)

Figure 7-9. (a) A violation of a causally-consistent store.

P1:	W(x)a
P2:	R(x)a W(x)b
P3:	R(x)b R(x)a
P4:	R(x)a R(x)b

(a)

Causal Consistency (3)

Monotonic Reads (1)

- A data store is said to provide monotonic-read consistency if the following condition holds:
 - If a process reads the value of a data item x ...
 - any successive read operation on x by that process
 - will always return that same value
 - or a more recent value.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Monotonic Reads (2)

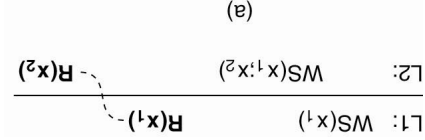


Figure 7-12. The read operations performed by a single process P at two different local copies of the same data store. (a) A monotonic-read consistent data store.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Monotonic Reads (3)

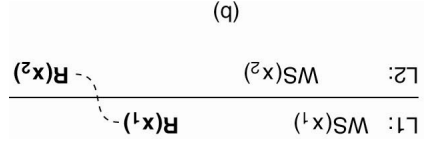


Figure 7-12. The read operations performed by a single process P at two different local copies of the same data store. (b) A data store that does not provide monotonic reads.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Monotonic Writes (2)

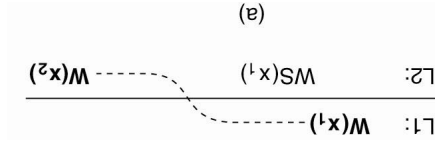


Figure 7-13. The write operations performed by a single process P at two different local copies of the same data store. (a) A monotonic-write consistent data store.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Monotonic Writes (3)

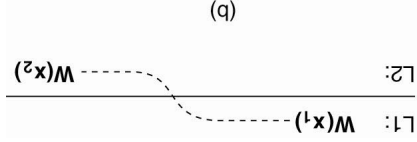


Figure 7-13. The write operations performed by a single process P at two different local copies of the same data store. (b) A data store that does not provide monotonic-write consistency.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

In a monotonic-write consistent store, the following

condition holds:

- A write operation by a process on a data item x ...
- is completed before any successive write operation on x
- by the same process.

Monotonic Writes (1)

Read Your Writes (1)

- A data store is said to provide read-your-writes consistency, if the following condition holds: The effect of a write operation by a process on data item x ...
- will always be seen by a successive read operation on x
- by the same process.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Read Your Writes (2)

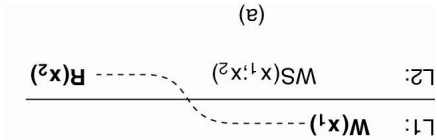


Figure 7-14. (a) A data store that provides read-your-writes consistency.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Read Your Writes (3)

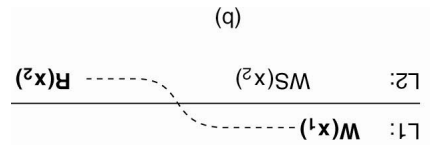


Figure 7-14. (b) A data store that does not.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Writes Follow Reads (1)

- A data store is said to provide writes-follow-reads consistency, if the following holds: A write operation by a process ...
- on a data item x following a previous read operation on x by the same process ...
- is guaranteed to take place on the same or a more recent value of x that was read.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Writes Follow Reads (2)

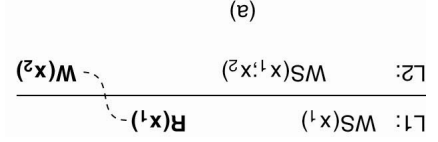


Figure 7-15. (a) A writes-follow-reads consistent data store.

Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Writes Follow Reads (3)

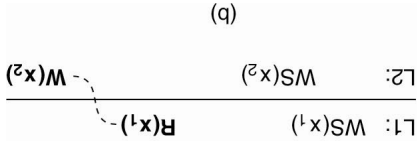


Figure 7-15. (b) A data store that does not provide writes-follow-reads consistency.

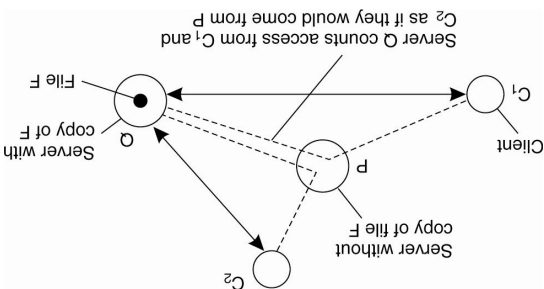
Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Figure 7-19. A comparison between push-based and pull-based protocols in the case of multiple-client, single-server systems.

Issue	Push-based	Pull-based
State at server	List of client replicas and caches	None
Messages sent	Update (and possibly fetch update later)	Poll and update
Response time at client	Immediate (or fetch-update time)	Fetch-update time

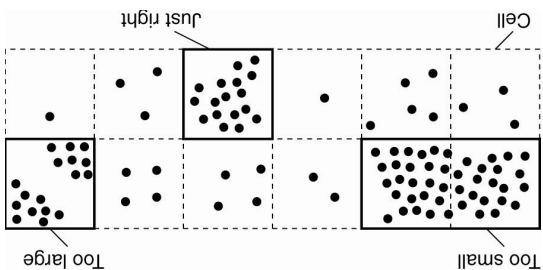
Pull versus Push Protocols

Figure 7-18. Counting access requests from different clients.



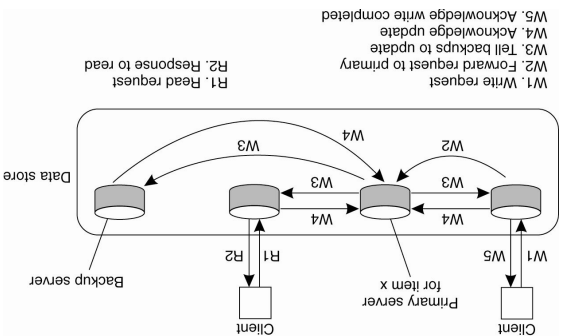
Server-Initiated Replicas

Figure 7-16. Choosing a proper cell size for server placement.



Replica-Server Placement

Figure 7-20. The principle of a primary-backup protocol.



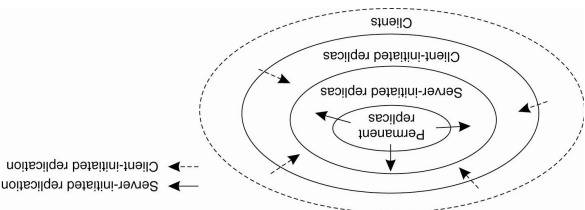
Remote-Write Protocols

1. Propagate only a notification of an update.
2. Transfer data from one copy to another.
3. Propagate the update operation to other copies.

Possibilities for what is to be propagated:

State versus Operations

Figure 7-17. The logical organization of different kinds of copies of a data store into three concentric rings.



Content Replication and Placement

Local-Write Protocols

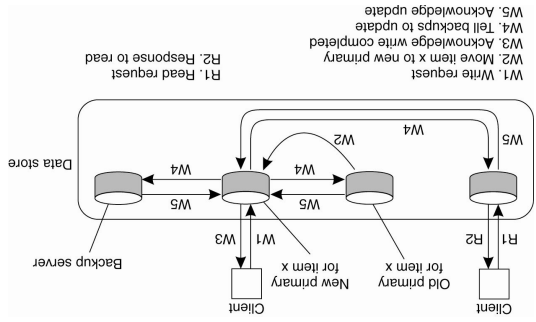


Figure 7-21. Primary-backup protocol in which the primary migrates to the process wanting to perform an update. Tanenbaum & Van Steen. Distributed Systems: Principles and Paradigms, 2e. (c) 2007

Quorum-Based Protocols

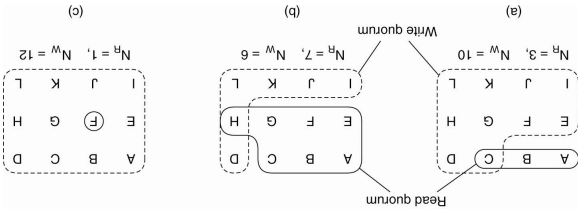


Figure 7-22. Three examples of the voting algorithm. (a) A correct choice of read and write set. (b) A choice that may lead to write-write conflicts. (c) A correct choice, known as ROWA (read one, write all).