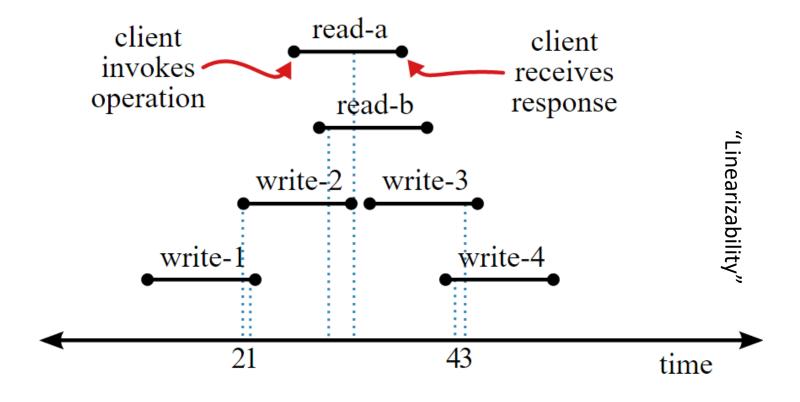
Consistency



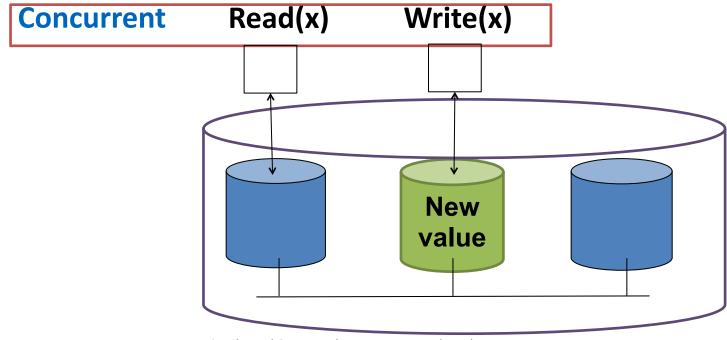
Consistency – according to dictionary.com

 steadfast adherence to the same principles, course, form, etc.

CONSISTENCY MODELS

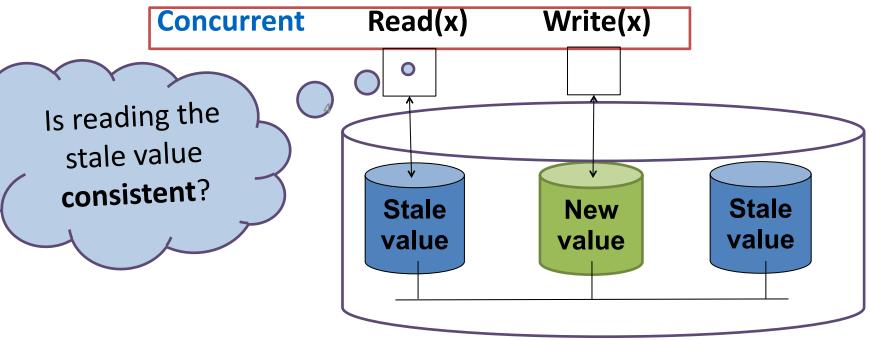
Concurrent Operations

- All operations must be applied in a specific order to all replicas
- Global ordering is too costly and not scalable (e.g., using a consensus algorithm)



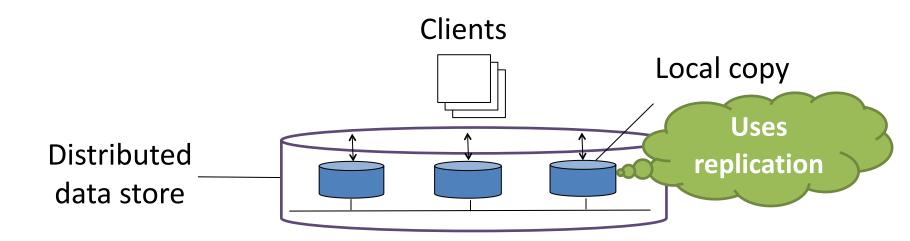
Concurrent Operations

- All operations must be applied in a specific order to all replicas
- Global ordering is too costly and not scalable (e.g., using a consensus algorithm)
- **Solution**: Avoid global ordering using **weaker** consistency requirements **suitable** for application



Consistency Models

- Definition (data-centric consistency model)
 - A contract between a distributed data store and a set of clients which specifies what results of concurrent read/ write operations can be



Distributed data store as synonym for replicas, distributed database, shared memory, shared files, etc.

Data-Centric Consistency Models

- **Data-centric** consistency models dictate the outcome of concurrent reads and writes:
 - Strict consistency
 - Sequential consistency
 - Linearizable consistency
 - Causal consistency
 - FIFO consistency
 - Weak consistency

— ...

Strict Consistency

 Definition: Any read on a data item x returns a value corresponding to the result of the most recent write on x

 Uni-processor systems have traditionally observed strict consistency, ...

$$=$$
 a = 1; a = 2; print(a);

Output?

• Definition assumes existence of **absolute global time** for unambiguous determination of "most recent".

Strict Consistency

 Definition: Any read on a data item x returns a value corresponding to the result of the most recent write on x

Uni-processor systems have traditionally observed strict consistency, ...
 but what about replicated systems?

• Definition assumes existence of **absolute global time** for unambiguous determination of "most recent".



```
x = 0;
y = 0;
x = 1;
r = y;
```

```
Result: x = 1 and r = 0
```

H.-A. Jacobsen 8



```
x = 0;

y = 0;
```

$$y = 1;$$

 $s = x;$

Result:
$$y = 1$$
 and $s = 0$



```
Shared variables
```

$$x = 0;$$

$$y = 0;$$

$$x = 1;$$

$$r = y;$$

$$y = 1;$$

$$s = x;$$

Result:

$$r = 0$$
 and $s = 1$

$$r = 1$$
 and $s = 0$

$$r = 1$$
 and $s = 1$



```
Shared variables x = 0; y = 0;
```

$$x = 1;$$
 $y = 1;$ $s = x;$



```
Shared variables
```

$$x = 0;$$

$$y = 0;$$

$$x = 1;$$

$$r = y;$$

$$y = 1;$$

$$s = x;$$

Result:

Write buffer x86

$$r = 0$$
 and $s = 1$

$$r = 1$$
 and $s = 0$

$$r = 1$$
 and $s = 1$

$$r = 0$$
 and $s = 0$

H.-A. Jacobsen



```
Shared variables
```

$$x = 0;$$

$$y = 0;$$

$$x = 1;$$

$$r = y;$$

$$y = 1;$$

$$s = x;$$

Result:

Write buffer

x86

r = 0 and s = 1

r = 1 and s = 0

r = 1 and s = 1

$$r = 0$$
 and $s = 0$

Compiler optimization / reordering

Interpretation of Strict Consistency

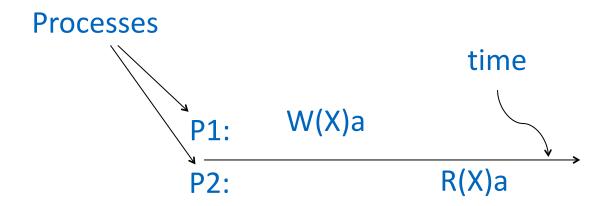
- Under strict consistent, all writes are instantaneously visible to all nodes and absolute global time order is maintained
- If a **replica is updated**, ...
 - all subsequent reads, see the new value, no matter how soon after the update the reads are done
 - and no matter which node is doing the reading and where it is located
- Similarly, if a read is done, then it gets the most recent value,
 no matter how quickly the next write is done

Notation

W(X)a: Represents **writing** value **a** to data **X** (memory)

R(X)a: Represents reading data X, which returns value a

Initial value of X is NIL



Strict Consistency Example



P1: W(X)a

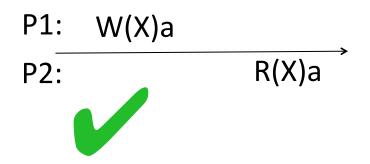
P2: R(X)a

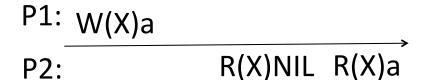
P1: W(X)a

P2: R(X)NIL R(X)a

Strict Consistency Example

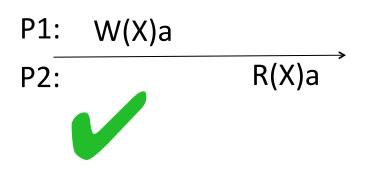


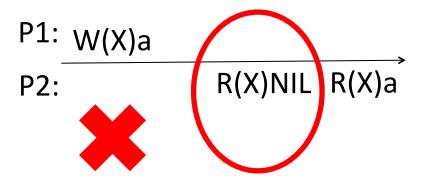




Strict Consistency Example

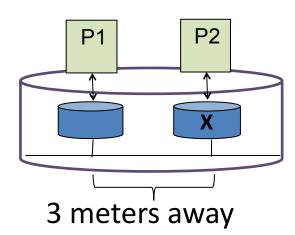






Strict Consistency: Thought Experiment

@P1: W(X)a at T_1 **@P2:** R(X)a at $T_1 + 1$ ns



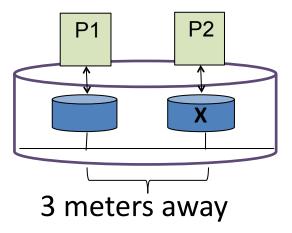
Strict Consistency: Thought Experiment

- To satisfy strict consistency, laws of physics may have to be violated! Obviously, not an option!!!
- Example:

———————→

@P1: W(X)a at T₁

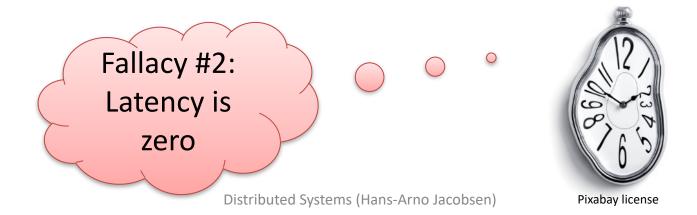
@P2: R(X)a at $T_1 + 1ns$



To realize strict consistency in this case, W(X)a would have to travel at 10 times the speed of light!

Strict Consistency: The Bad News

- It is impossible to instantaneously replicate operations
- It is **impossible to perfectly** synchronize clocks
 - How to accurately determine the time of each operation?



So, what are we to do?

 Relax consistency and introduce models with weaker consistency requirements

Data-Centric Consistency Models

- Data-centric consistency models dictate the outcome of concurrent reads and writes:
 - Strict consistency
 - Sequential consistency
 - Linearizable consistency
 - Causal consistency
 - FIFO consistency
 - Weak consistency

— ...

Pixabay license

Distributed Systems (Hans-Arno Jacobsen)

Self-study questions

- Verify the timing assumptions in the thought experiment on strict consistency?
- What aspect of a consistency model could be relaxed and why?
- Create some concurrent read and write sequences that conform to strict consistency.
- Create some concurrent operation sequences that violate strict consistency.
- See if the projected outcome of r=0, s=0 can occur on an x86 or due to compiler optimizations.
- Would a quantum tunnel help us achieve stric consistency in a replicated system?

of each individual processor does not guarantee that the multissor computer is sequentially consistent. In this brief note we describe a method of interconnecting sequential processors with memory modules that insures the sequential consistency of the resulting multiprocessor.

We assume that the computer consists of a collection of proces sors and memory modules, and that the processors communicate with one another only through the memory modules. (Any special communication registers may be regarded as separate memory modules.) The only processor operations that concern us are the operations of sending fetch and store requests to memory modules. We assume that each processor issues a sequence of such requests. (It must sometimes wait for requests to be executed, but

that does not concern us.)

We illustrate the problem by considering a simple two-process mutual exclusion protocol. Each process contains a critical section, and the purpose of the protocol is to insure that only one process may be executing its critical section at any time. The protocol is as follows.

process 1 a := 1: if b = 0 then critical section;process 2 b := 1: if a = 0 then critical section; else ··· fi

The else clauses contain some mechanism for guaranteeing even tual access to the critical section, but that is irrelevant to the discussion. It is easy to prove that this protocol guarantees mu-tually exclusive access to the critical sections. (Devising a proof provides a nice exercise in using the assertional techniques of [2] and [3], and is left to the reader.) Hence, when this two-process program is executed by a sequentially consistent multiprocessor computer, the two processors cannot both be executing their criti-

We first observe that a sequential processor could execute the

condition will be called sequentially consistent. The sequentially consistent with the called sequentially consistent condition will be called sequentially consistent. The sequentially consistent is consistent to misming waiting, the processor can issue the store request to the memory module without specifying the value to be

0018-9340/79/0900-0690\$00.75 © 1979 IEEE

How to Make a Multiprocessor Computer That Correctly Executes Multiprocess Programs

LESLIE LAMPORT

Abstract—Many large sequential computers execute operations in a faitent order than is specified by the program. A correct execution is a faitened if the results produced and the same as would be produced by executing the program steps in order. For a multiprocessor computer, such a correct execution by each processor does not guarantee the correct execution of the entire program. Additional conditions are given which be guarantee that a computer correctly conditions are given which be guarantee that a computer correctly

Index Terms—Computer design, concurrent computing, hardware correctness, multiprocessing, parallel processing.

A high-speed processor may execute operations in a different cal sections at the same time. order than is specified by the program. The correctness of the "b = 1" and "fetch b" operations of process 1 in either order. order than is specimed by the processor satisfies the following (When process 1's program is considered by itself, it does not condition: the result of an execution is the same as if the opera-tions had been executed in the order specified by the program. A processor satisfying this condition will be called sequential. Con-sider a computer composed of several such processors accessing a common memory. The customary approach to designing and proving the correctness of multiprocess algorithms [11-13] for such a computer assumes that the following condition is satisfied: the result of any execution is the same as if the operations of all the result of any execution is the same as if the operations of all the processors were executed in some sequential order, and the ing a value is possible only after the value has been computed. A the processor were second in control of the processor appear in this sequence in the order specified by its program. A multiprocessor satisfying this processor will often be ready to issue a memory fetch request before it knows the value to be stored by a preceding store

Manuscript received September 28, 1977; revised May 8, 1979.
The author is with the Computer Science Laboratory, SRI International, Meris
Act, CA 94025.

SEQUENTIAL CONSISTENCY

FORMALIZED IN 1979

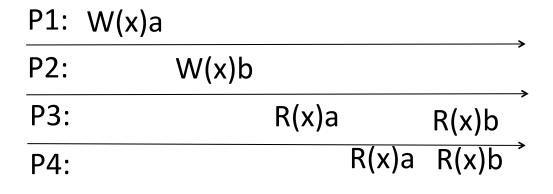
Definition of Sequential Consistency

 The result of any execution is the same as if the operations by all processes on the data store were executed in some sequential order and ...

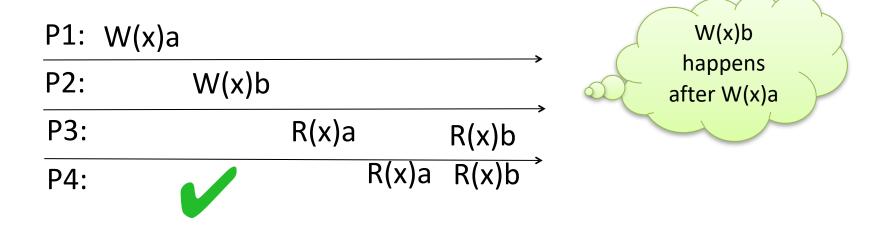
 ... the operations of each individual process appear in this sequence in the order specified by its program

 Weaker than strict consistency: logical time instead of physical time

Sequential Consistency Example I



Sequential Consistency Example I



Sequential Consistency Example II

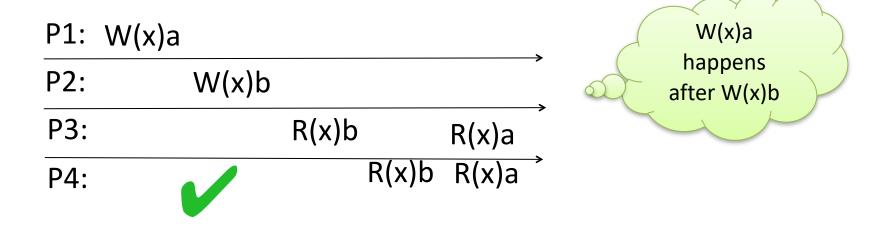
P1: W(x)a P2: W(x)b P3: R(x)b R(x)a P4: R(x)b R(x)a

P1: W(x)a
P2: W(x)b

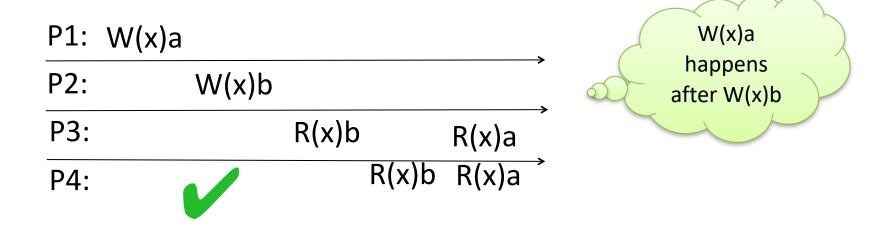
P3: R(x)b R(x)a

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

Sequential Consistency Example II



Sequential Consistency Example II



Cyclic dependency between the writes!

P1: W(x)a
P2: W(x)b

R(x)b
R(x)a
R(x)a
R(x)b

Self-study questions

- Create some concurrent read and write sequences that conform to sequential consistency.
- Create some concurrent operation sequences that violate sequential consistency.
- Is a strictly consistent system, sequentially consistent and vice versa, discuss?
- How could sequential consistency be implemented?

LINEARIZABILITY

FORMALIZED IN 1990

Linearizability: A Correctness Condition for Concurrent Objects

MAURICE P. HERLIHY and JEANNETTE M. WING Carnegie Mellon University

A concurrent object is a data object chared by concurrent processes. Linearisability is a correctness condition for concurrent objects that exploits the semantics of abstract data types. It permits a high degree of concurrency, yet it permits programmers to specify and reason about concurrent objects using known techniques from the sequential domain. Linearizability provides the illusion that each operation applied by concurrent processes takes effect instantaneously at some point between its invocation and its response, implying that the meaning of a concurrent object's operations can be given by pre- and post-conditions. This paper defines linearizability, compares it to other correctness conditions, presents and demonstrates a method for proving the correctness of implementations, and shows how to reason about concurrent objects, given they are linearizable.

Catagories and Subject Descriptors: D.1.3 [Programming Techniques]: Concurrent Programming; D.2.1 [Software Engineering]: Requirements/Specifications; D.3.3 [Programming Languages]: Language Construct—abstract data types concurrent programming structures, data types and structures; F.1.2 [Computation by Abstract Devices]: Modes of Computation—parallelism; F.3.1 [Logics and Meanings of Programs]: Specifying and Verifying and Reasoning about Programs—pre- and post-conditions, specification techniques

General Terms: Theory, Verification

Additional Key Words and Phrases: Concurrency, correctness, Larch, linearizability, multiprocessing, serializability, shared memory, specification

1. INTRODUCTION

1.1 Overview

Informally, a concurrent system consists of a collection of sequential processes that communicate through shared typed objects. This model encompasses both message-passing architectures in which the shared objects are message queues,

A preliminary version of this paper appeared in the Proceedings of the 14th ACM Symposium on Principles of Programming Languages, January 1987 [21].

This research was sponsored by IBM and the Defense Advanced Research Projects Agents (DOD), ARI'A order 4976 (Amendment 20), under contract F33015-07-C-1499, monitored by the Avionics Laboratory, Air Force Wright Aeronautical Laboratories, Wright-Patterson AFB. Additional spport for J. M. Wing was provided in part by the National Science Foundation under grant CCR-8620027. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of the Defense Advanced Research Projects Agency or the US Government.

Authors' address: Department of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213-3890.

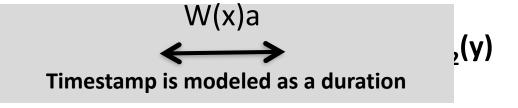
Permission to copy without fee all or part of this material is granted provided that the copies are not made or distributed for direct commercial advantage, the ACM copyright notice and the title of the publication and its date appear, and notice is given that copyring is by permission of the Association for Computing Machinery. To copy otherwise, or to republish, requires a fee and/or specific permission.

© 1990 ACM 0164-0925/90/0700-0463 \$01.50

ACM Transactions on Programming Languages and Systems, Vol. 12, No. 3, July 1990, Pages 463-492.

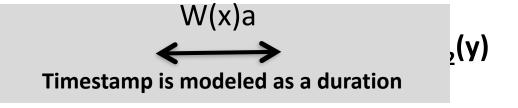
Definition of Linearizability

- 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.
- Here, $ts_{OP}(x)$ denotes **timestamp** assigned to operation OP that is performed on data item x, and OP is either read (R) or write (W)
- If ts_{OP1}(x) < ts_{OP2}(y), then c in this sequence



Definition of Linearizability

- 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 sequen process appear process appear program.
- Here, $ts_{OP}(x)$ denotes **timestamp** assigned to operation OP that is performed on data item x, and OP is either read (R) or write (W)
- If ts_{OP1}(x) < ts_{OP2}(y), then c in this sequence

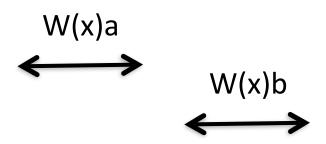


 Linearizability is weaker than strict consistency, but stronger than sequential consistency.

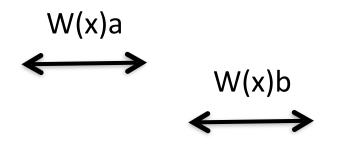
 Like strict consistency, assumes global time, but not absolute global time

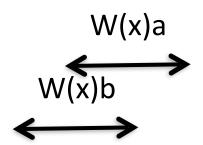
 Assumes processes in the system have physical clocks synchronized to within an bounded error (captured by timestamp's duration)

• If *W(x)b* was the **most**recent write operation
and there is no other write
operation overlapping
with *W(x)b*, then any later
read should return *b*

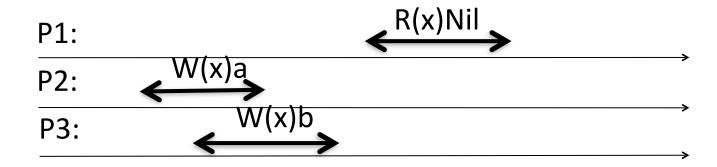


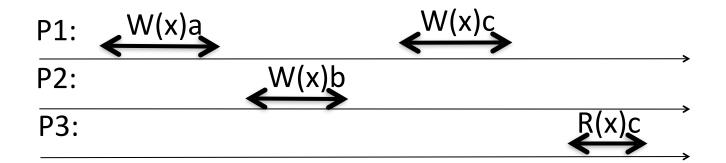
- If W(x)b was the most recent write operation and there is no other write operation overlapping with W(x)b, then any later read should return b
- If W(x)a and W(x)b were two most-recent overlapping write operation, then any later read should return either a or b, not something else

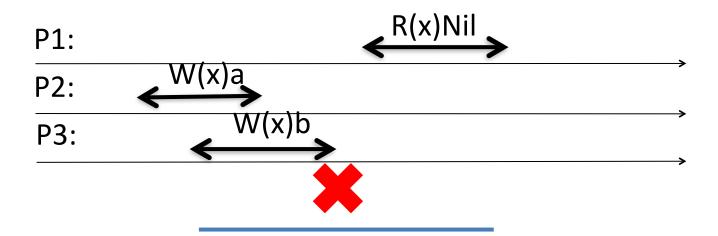


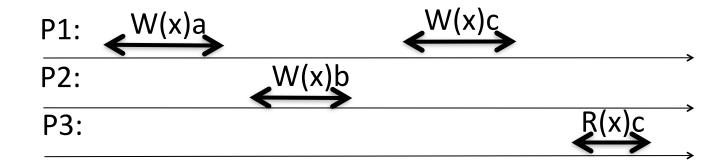


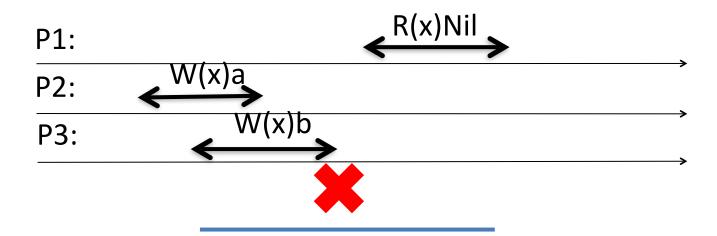
- A linearizable data store is also sequentially consistent
- I.e., linearizability is more restrictive (stronger)
- Difference is ordering according to a set of synchronized clocks
- Linearizability prevents stale reads, since it guarantees correct reads after a write is completed

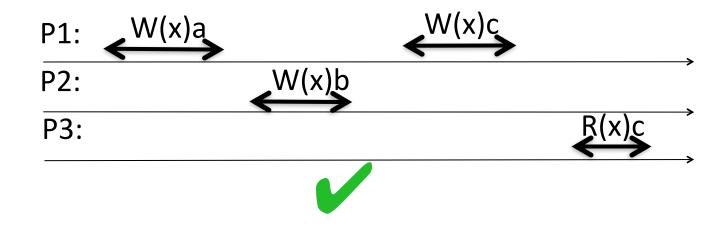


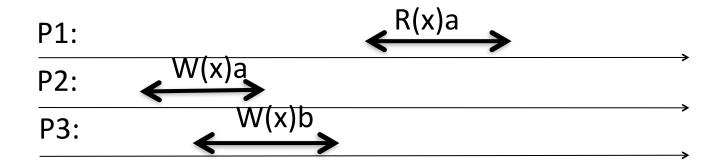


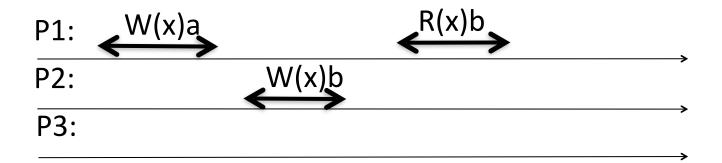


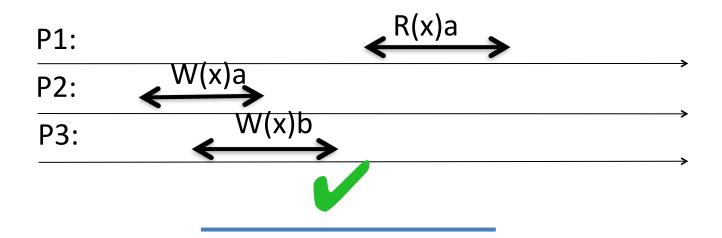


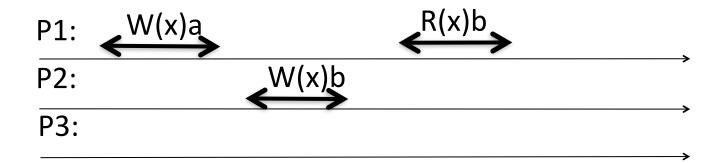


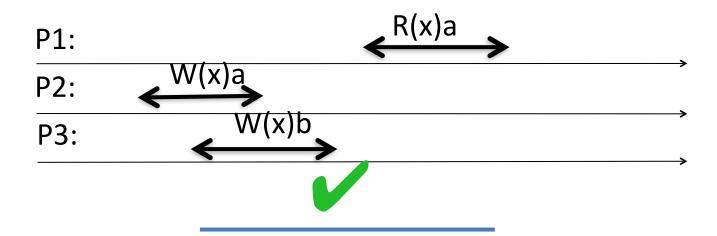


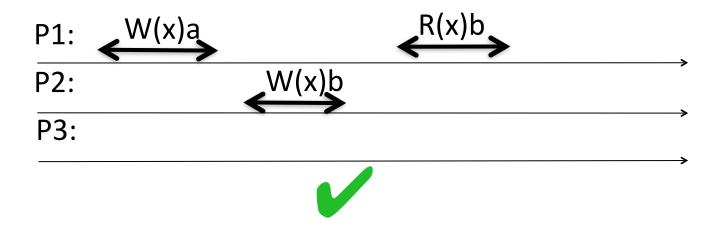


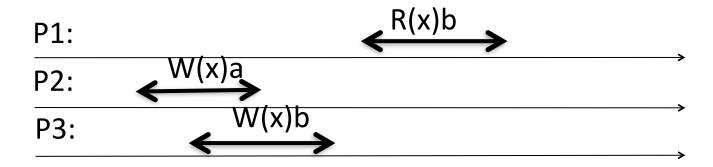


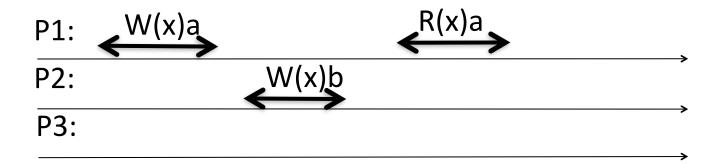


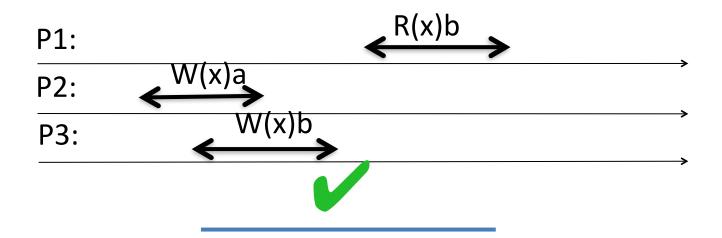


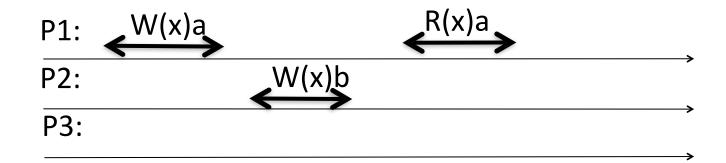


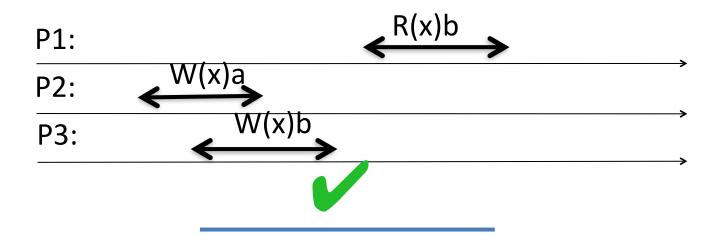


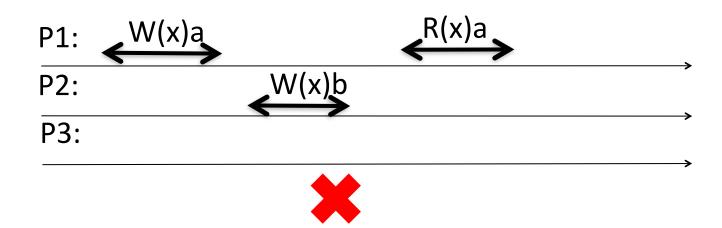


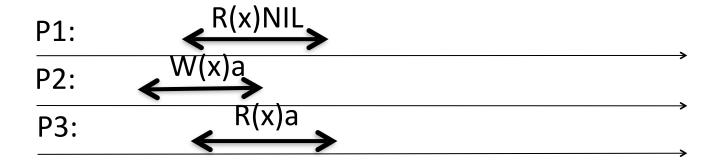


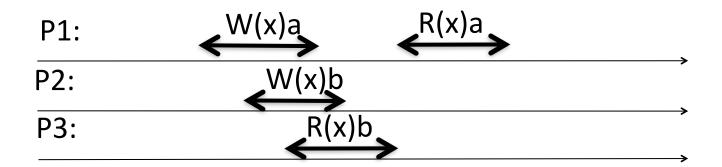


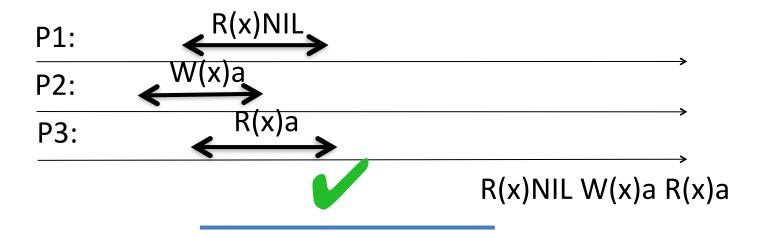


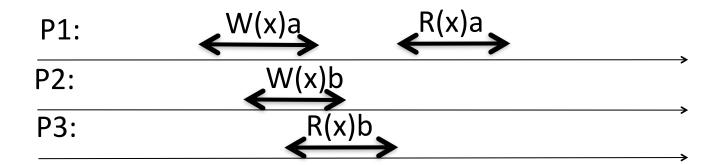


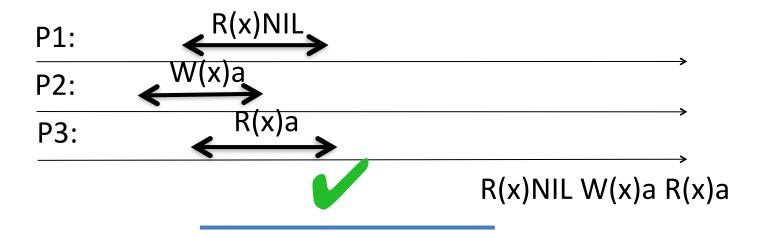


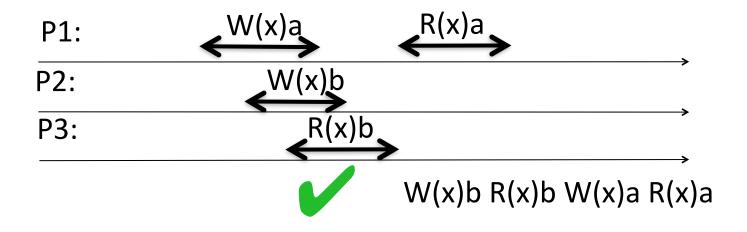


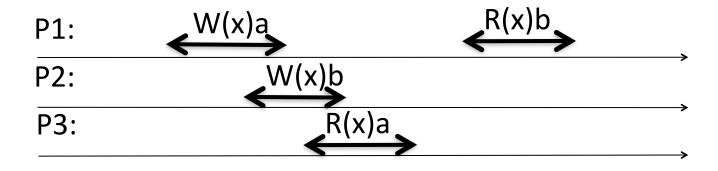


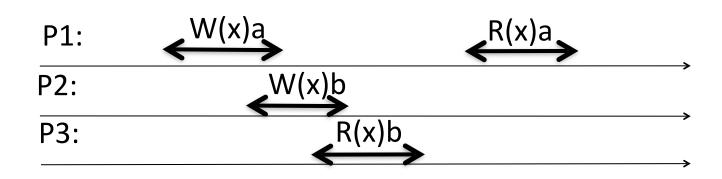


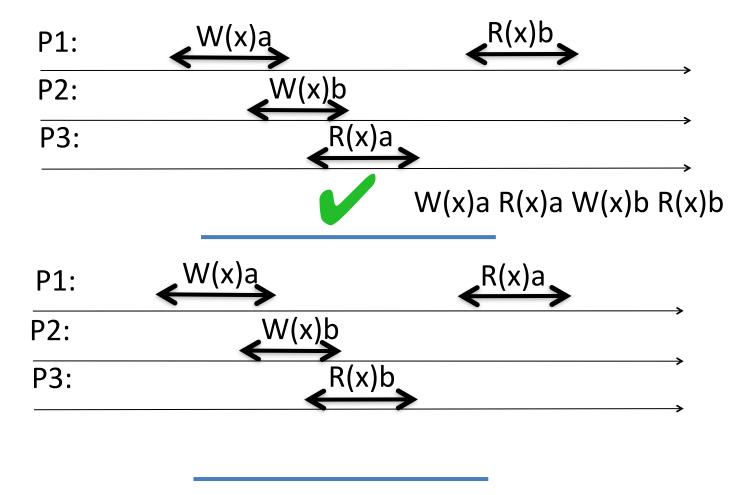


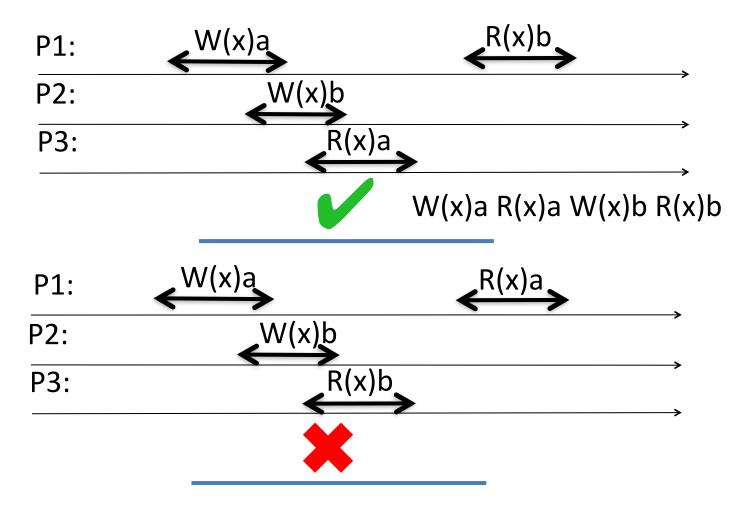


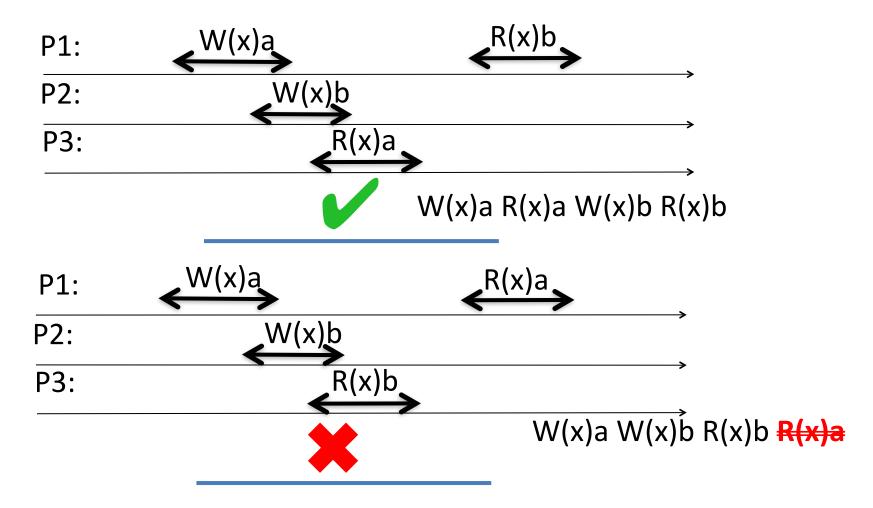


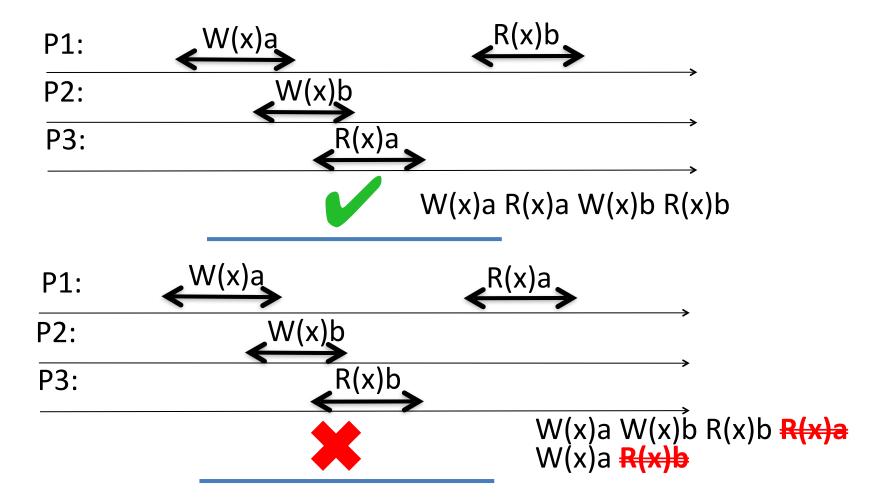


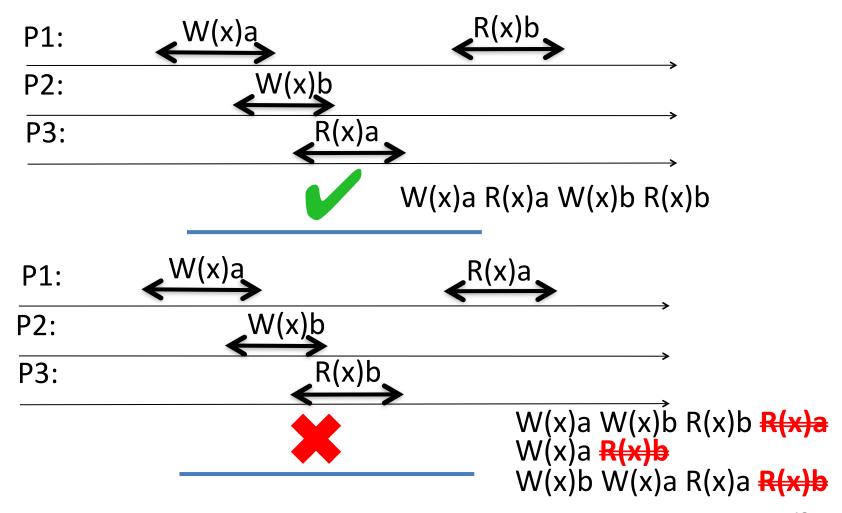


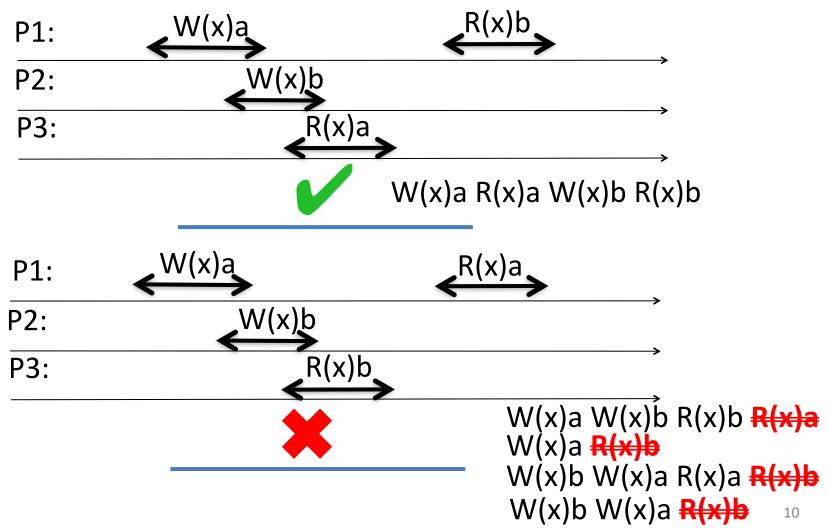












Self-study questions

- Create some concurrent read and write sequences that conform to linearizability.
- Create some concurrent operation sequences that violate linearizability.
- Is a sequentially consistent data store linearizable, argue for or against?
- Is a strictly consistent system, sequentially consistent and vice versa, discuss?
- How could linearizability be implemented?



Pixabay.com

CAUSAL CONSISTENCY

Intuition for Causal Consistency

- Weaker than sequential consistency
- Distinguish events (writes) that are potentially causally related and those that are not (concurrent writes)
- Similar to happened-before relation (cf. Lamport clock), but with read and write operations instead of events
- If Write B is influenced (caused) by an earlier Write A (A→B),
 causal consistency requires that every process first sees effect
 of A then B
- Concurrent writes may be seen in a different order by different processes

Causal Relationship

- Read followed by write in same process, both are causally related
 - **Example**: $R(x)a \rightarrow W(y)b$ (e.g., it may be that y=f(x))
- Read is causally related to write that provided value read got (different processes)
 - **Example**: W(x)a \rightarrow R(x)a
- Transitivity: if $Op_1 \rightarrow Op_2$, $Op_2 \rightarrow Op_3$, then $Op_1 \rightarrow Op_3$
- Independent writes by different processes are not causally related - they are concurrent
 - **Example:** W(x)a || W(x)b

"Potentially" Causally Related

```
P1: W(x)a
R(x)a \qquad W(y)b
```

- Writing of x in P1 and reading of x in P2 are causally related
- Reading of x in P2 and writing of y in P2 are potentially causally related
- Computation of y may have depended on value of x read by P2 (written by P1); e.g., y = f(x)
- On the other hand, y may not have depended on x, yet potential causality still holds in our formalization!
- There is a causal chain between writing of x in P1 and writing of y in P2 (via reading of x in P2)

"Potentially" Causally Related

```
P1: W(x)a
R(x)a \rightarrow W(y)b
```

- Writing of x in P1 and reading of x in P2 are causally related
- Reading of x in P2 and writing of y in P2 are potentially causally related
- Computation of y may have depended on value of x read by P2 (written by P1); e.g., y = f(x)
- On the other hand, y may not have depended on x, yet potential causality still holds in our formalization!
- There is a causal chain between writing of x in P1 and writing of y in P2 (via reading of x in P2)

Definition of Causal Consistency

- "Potentially" causally related writes must be seen by all processes in the same order
- Concurrent writes may be seen in a different order by different processes

```
\begin{array}{ccc}
P1: & W(x)a \\
\hline
P2: & W(x)b
\end{array}
```

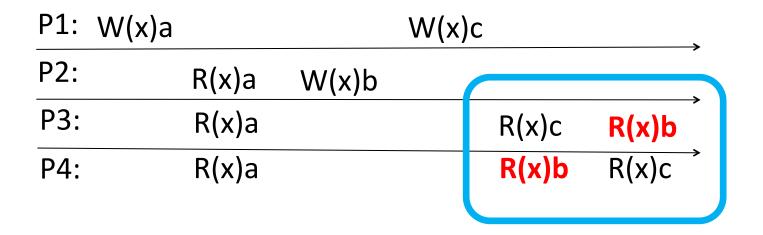
Program order must be met

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

Strictly, sequentially or causally consistent?

P1:	W(x)a			W(x)c			~
P2:		R(x)a	W(x)b				
P3:		R(x)a			R(x)c	R(x)b	7
P4:		R(x)a			R(x)b	R(x)c	>

Neither strictly, nor sequentially consistent, ...



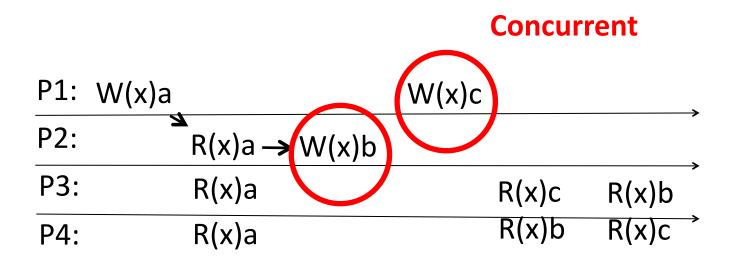
Neither strictly, nor sequentially consistent, ...

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

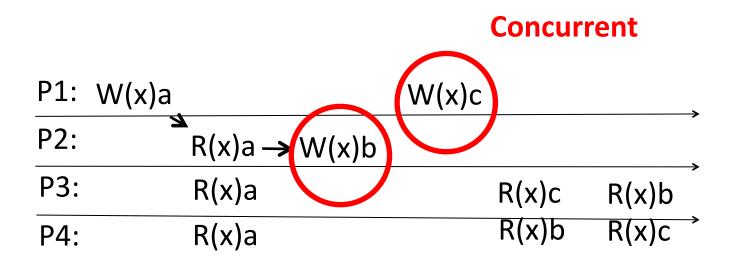
Causal Consistency Example I

P1:	W(x)a		W(x)c		•	>
P2:	2	$R(x)a \rightarrow W(x)b$				
P3:		R(x)a		R(x)c	R(x)b	>
P4:		R(x)a		R(x)b	R(x)c	>

Causal Consistency Example I



Causal Consistency Example I



Neither strictly, nor sequentially consistent, but causally consistent!

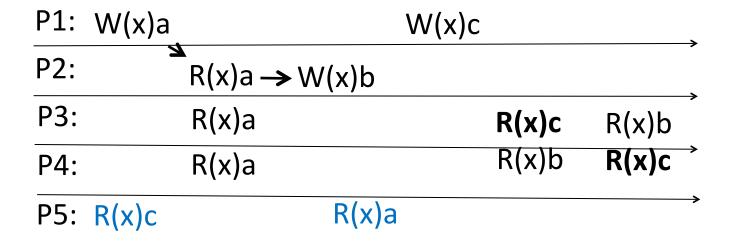
Causal Consistency Example I'

P1: W(x)a			W(x)c		
P2:	R(x)a	W(x)b			
P3:	R(x)a			R(x)c	R(x)b
P4:	R(x)a			R(x)b	R(x)c
P5: R(x)c		R(x)a			

Causal Consistency Example I'

P1:	W(x)a		W(x)c		_
P2:	2	$R(x)a \rightarrow W(x)b$			
P3:		R(x)a		R(x)c	R(x)b
P4:		R(x)a		R(x)b	R(x)c →
P5:	R(x)c	R(x)a	9		

Causal Consistency Example I'



Not causally consistent, W(x)a happened-before W(x)c in P1

Self-study questions

- Create some concurrent read and write sequences that conform to causal consistency.
- Create some concurrent operation sequences that violate causal consistency.
- Is a sequentially consistent data store causally consistent, argue for or against?
- Is a causally consistent data store, sequentially consistent, argue for or against?
- How could causal consistency be implemented?

FIFO CONSISTENCY

FIRST IN – FIRST OUT

Definition of FIFO Consistency

 Writes by a single process are seen by all other processes in the order in which they were issued

 Writes from different processes may be seen in a different order by different processes

 Easy to maintain: Simply send writes in FIFO order from each process to each replica (e.g., using TCP)

FIFO Consistency Example I



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

P1: W(x)aP2: R(x)a W(x)b W(x)cP3: R(x)c R(x)a R(x)b

FIFO Consistency Example I



```
P1: W(x)a

P2: R(x)a W(x)b W(x)c

P3: R(x)bR(x)a R(x)c

P1: W(x)a

P2: R(x)a W(x)b W(x)c

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

FIFO Consistency Example I



```
P1:
      W(x)a
             R(x)a
                    W(x)b W(x)c
P2:
                                   R(x)bR(x)aR(x)c
P3:
P1:
     W(x)a
                    W(x)b W(x)c
P2:
             R(x)a
                                   R(x)c R(x)a R(x)b
P3:
```

FIFO Consistency Example II



```
P1: W(x)a
```

P2: R(x)a W(x)b W(x)c

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

P1: W(x)a

P2: R(x)a W(x)b W(x)c

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

FIFO Consistency Example II



P1: W(x)a

P2: R(x)a W(x)b W(x)c

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

P1: W(x)a

P2: R(x)a W(x)b W(x)c

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

FIFO Consistency Example II



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

Self-study Questions

- Create some concurrent read and write sequences
 Pixabay.com
 that conform to FIFO consistency.
- Create some concurrent operation sequences that violate FIFO consistency.
- Is a sequentially consistent system, FIFO consistent, argue for or against?
- Is a causally consistent system, FIFO consistent, argue for or against?
- Find sample applications of FIFO consistency.

WEAK CONSISTENCY

Weak Consistency

- Not all processes need to see all writes, let alone in same order
- Based on synchronization variable S with single operation
 Synchronize(S)
- Any process can perform read/write operations and synchronize
- Order of operations before synchronization is not consistent
- After synchronization, all processes see same outcome of operations preceding synchronization point

Weak Consistency Interpretation

- Process forces the just written value out to all synchronizing replicas
- Process can be sure to get the most recent value written before it reads

- Model enforces consistency on a group of operations as opposed to individual reads and writes
- Care about the effect of a group of reads and writes



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

P1: W(x)a W(x)b S
P2: S R(x)a



P2:

Weak Consistency Ex

 $\frac{1}{\sqrt{x}}$

P2 & P3 have yet to synchronize, no guarantees about values read

R(x)a

S

P1:	vv(x)a	vv(x)b	3			
P2:				R(x)a	R(x)b	S (
P3:				R(x)b	R(x)a	S
		V				
P1:	W(x)a	W(x)b	S			

C



Weak Consistency Ex

P2 & P3 have yet to synchronize, no guarantees about values read

P1:	W(x)a	W(x)b	S				
P2:				R(x)a	R(x)b	S	
P3:				R(x)b	R(x)a	S	
		V		F	P1 propa	agate	es
P1:	W(x)a	W(x)b	S	,00	value	e b	
P2:				S	R(x)	a	
		×					•



Agreement on value is implementation dependent.

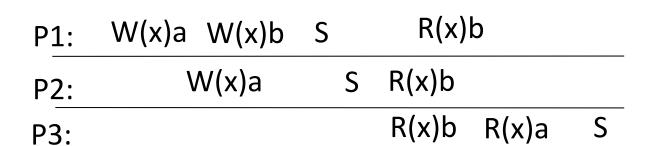
P1:	W(x)a W(x)b	S	R(x)a	
P2:	W(x)a	S	R(x)a	
P3:			R(x)b R(x)a	S

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



Agreement on value is implementation dependent.

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

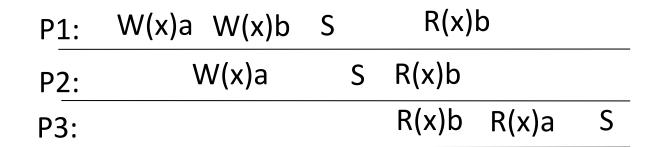




Agreement on value is implementation dependent.

P1 & P2 agree on value a

P1:	W(x)a W(x)b	S	R(x)a	00	
P2:	W(x)a	S	R(x)a		
P3:			R(x)b	R(x)a	S





Agreement on value is implementation dependent.

P1 & P2 agree on value a

P1:	W(x)a W(x)b	S	R(x)a	00	
P2:	W(x)a	S	R(x)a		
P3:			R(x)b R	k(x)a	S



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





Agreement on value is implementation dependent.

P1 & P2 agree on value a

P1:	W(x)a W(x)b	S	R(x)a	000	
P2:	W(x)a	S	R(x)a		
P3:			R(x)b	R(x)a	S



 $\frac{1}{2}$

P1 & P2 agree on value b

LT:	vv(x)a vv(x)D	3	I (X)D	0 0	
P2:	W(x)a	S	R(x)b		
P3:			R(x)b	R(x)a	S

R(x)h





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

```
P1: W(x)a W(x)b S R(x)b

P2: W(x)a S R(x)a

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



P1: W(x)a W(x)b S R(x)a
P2: W(x)a S R(x)b
R(x)b R(x)a S



P1: W(x)a W(x)b S R(x)b

P2: W(x)a S R(x)a

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



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



P1: W(x)a W(x)b S R(x)b

P2: W(x)a S R(x)a

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



Self-study Questions

- Create some concurrent read and write sequences that conform to weak consistency.
- Can weak consistency emulate any of the other consistency models we study?
- How can a weakly consistent data store be implemented?
- Find sample applications of weak consistency.

DATA-CENTRIC CONSISTENCY MODELS SUMMARY

Strongest

Summary of Consistency Models

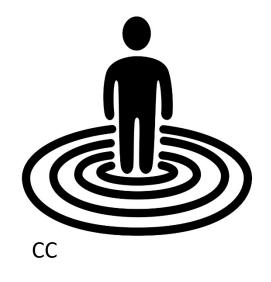
Consistency	Description
Strict	Absolute time ordering of all writes
Linear-	All processes see all writes in same order.
izability	Accesses ordered according to a (non-unique)
	global timestamp
Sequential	All processes see all writes in same order
	Accesses are not ordered in time
Causal	All processes see causally-related writes in same
	order; concurrent writes must not be ordered
FIFO	All processes see writes from each other in the
	order they were issued; writes from different
	processes must not be seen in that order
Weak	Shared data can be counted on to be consistent
	only after synchronizing

Where does ... come into the picture?

- Linearizability: Strongest distributed solution, possible with eager replication (synchronous), chain replication, supported in Hbase, Bigtable
- **Sequential**: System-wide consistent reads, e.g., everyone sees replies to a post in same order
- Causal: Everyone sees posts before replies
- FIFO: Reading all messages from each friend in order but not across friends
- Weak: Responsibility left to developer who must explicitly enforce synchronization

Self-study Questions

- Create some concurrent read and write sequences that conform to different consistency models.
- How do consistency models relate to one another, if at all?
- Find sample applications for each consistency model.



CLIENT-CENTRIC CONSISTENCY

Client-centric Consistency

- So far, goal was to maintain consistency in presence of concurrent read and write operations
- There are use cases with no (few) concurrent writes to the same key, or consistency of write operations are secondary
 - DNS: No write-write conflicts since there is a single authority updating each domain (disjoint partitions)
 - **Key-value stores**: **Usually no write-write conflicts** since updates partitioned by keys, e.g., Dynamo, Cassandra (also, optimistic concurrency control)
 - WWW: Heavy use of client-side caching, reading stale pages is acceptable in many cases

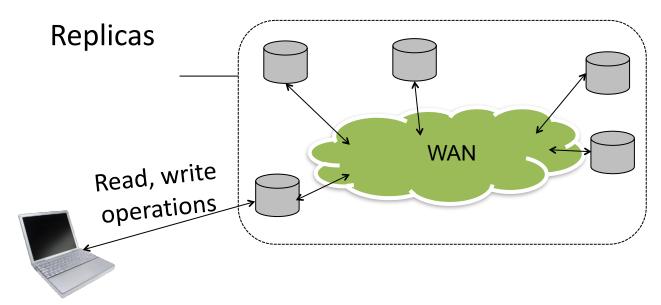
Client-centric Consistency

- Client-centric consistency puts emphasis on maintaining a consistent view for a single process, instead of on data stored in system
- Emphasis on read-write conflicts, and we assume write-write conflicts do not exist (i.e., no (few) concurrent writes)
- Client-centric consistency models describe what happens when a single client writes/reads from multiple replicas

Eventual Consistency

- Eventual consistency states that all replicas **eventually** converge when write operations stop
 - E.g., lazy replication using gossiping (cf. replication)
- Very weak form of consistency with no time bound, but highly available (i.e., always return a value, but could be stale)
- Works fine if a client always reads from same replica...
- ...but gives "weird" results if client reads from multiple replicas, due to:
 - Client mobility
 - Replica failure

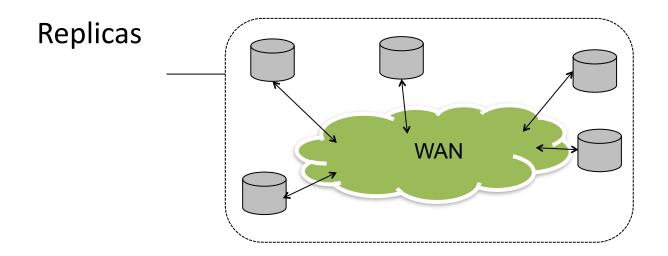
Mobile Users Use Case (Bayou System)



Replicas maintain client-centric consistency

- System developed by Xerox PARC in mid 90's
- Designed to work under network partitions
- Replicas may take some time to converge

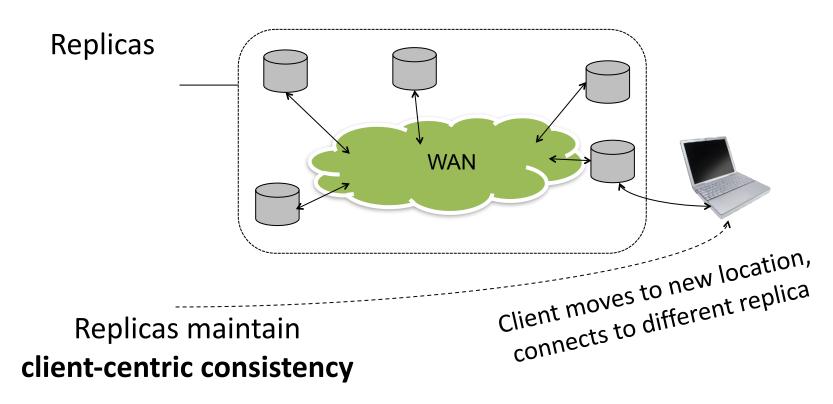
Mobile Users Use Case (Bayou System)



Replicas maintain client-centric consistency

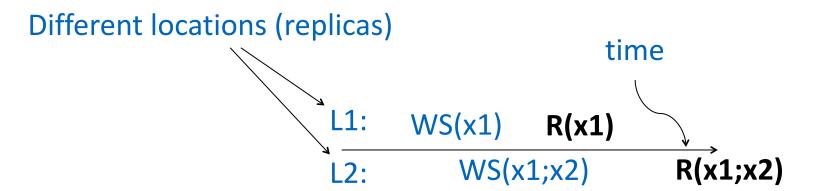
- System developed by Xerox PARC in mid 90's
- Designed to work under network partitions
- Replicas may take some time to converge

Mobile Users Use Case (Bayou System)



- System developed by Xerox PARC in mid 90's
- Designed to work under network partitions
- Replicas may take some time to converge

Notation



- Reads and writes by one client at different locations
- WS represents a series of write operations at L_i
- E.g., WS(x1) denotes that only x1 has been applied;
 WS(x1;x2) denotes that x1 and then later x2 have been applied
- R(x1;x2) means read returns a value formed by x1 followed by x2
- W(x1), R(x1) are write and read operations

Client-centric Consistency Models

- Read your writes consistency
- Monotonic reads consistency
- Writes follow reads consistency
- Monotonic writes consistency

Read-Your-Writes Consistency

The following consistency models always state:

A data store provides "XYZ" consistency if

• • •

in the same process potentially issueing reads or writes against different replicas.

Read Your Writes Consistency

Informally, ... W ... R ... (by same process!)

 If a read follows a write then the effect of the write is included in the set of writes read by R

 A write operation is always completed before a successive read operation by the same process, no matter where the read operation takes place

Read Your Writes



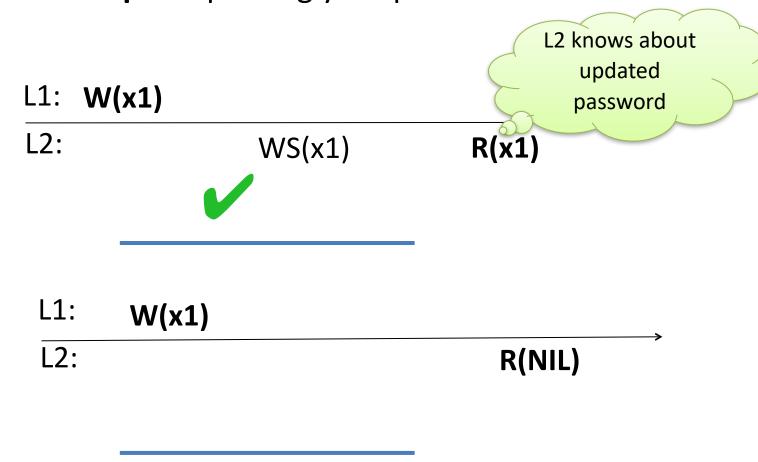
Example: Updating your password on a cluster

L1: W	V(x1)			
L2:		WS(x1)	R(x1)	7
L1:	W(x1)			
L2:			R(NIL)	~
			-	

Read Your Writes



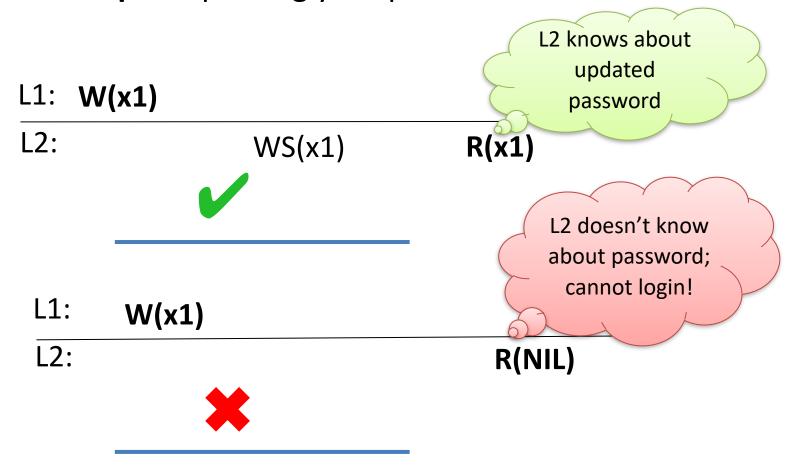
Example: Updating your password on a cluster



Read Your Writes



Example: Updating your password on a cluster



Monotonic-Reads Consistency

- ... $WS_1(..) R ... WS_2(..) R ... (s.t. <math>WS_1(..) \subseteq WS_2(..)$)
- If a process reads a value formed by a set of operations, WS₁, any successive read operation of that process always returns a value formed from a superset of WS₂
- A process always sees more recent data (but not necessarily fresh!)
- If a process reads again from another replica that replica must have already received the relevant older operations

Example: Monotonic-Reads



Example: E-mail client

Each read returns list of emails

Each write adds one e-mail

L1: WS(x1) R(x1)L2: WS(x1;x2) R(x1;x2)

L1: WS(x1) R(x1)
L2: WS(x2) R(x2)

Example: Monotonic-Reads



Example: E-mail client

Each read returns list of emails

Each write adds one e-mail

L1: WS(x1) R(x1)

L2: WS(x1;x2)

 $S(x1;x2) \qquad \qquad R(x1;x2)$

L2 knows about all emails

L1: WS(x1) R(x1)

L2: WS(x2) R(x2)

Example: Monotonic-Reads



Example: E-mail client

Each read returns list of emails

Each write adds one e-mail

L1: WS(x1)

R(x1)

L2:

WS(x1;x2)

R(x1;x2)

L2 knows about all e-mails



L1: WS(x1)

R(x1)

L2:

WS(x2)

R(x2)

L2 doesn't know about e-mail x1!



Writes Follow Reads Consistency

• ... WS(X1, ..., Xn) **R1** ... WS(X1, ..., Xn) **W2** ...

If a read R1 precedes a write W2 then at any replica, if W2 is written, it is preceded by the write set read by R1

 Any successive write operation by a process will be performed on a state that is up to date with the value most recently read by that process

Writes Follow Reads



Example: Facebook Wall

Reading a post

Replying to a post

L1: WS(x1) R(x1)
L2: WS(x1) W(x2)

 $\begin{array}{ccc}
L1: & WS(x1) & R(x1) \\
\hline
L2: & WS(NIL) & W(x2)
\end{array}$

Writes Follow Reads



Example: Facebook Wall

Reading a post

Replying to a post

L2 knows post (x1), writes reply (x2) to known post

L1: WS(x1) R(x1)

L2: WS(x1) W(x2)



L1: WS(x1) R(x1)
L2: WS(NIL) W(x2)

Writes Follow Reads



Example: Facebook Wall

Reading a post

Replying to a post

L2 knows post (x1), writes reply (x2) to known post

L1: WS(x1) R(x1)

L2: WS(x1) W(x2)



L2 doesn't know about post to reply to!

L1: WS(x1) R(x1)

L2: WS(NIL) W(x2)



Monotonic-Writes Consistency

Informally, ... W1 ... W2 ...

• If a write W1 precedes a write W2 then on any replica, W1 must precede W2

• In other words, writes by the same process are performed in the same order at every replica

Resembles FIFO consistency model

Monotonic Writes



Example: Replicating software code W(x1) adds a new method to a program W(x2) calls new method

L1: W(x1)

L2: WS(x1) W(x2)

L1: W(x1)

L2: W(x2)

Monotonic Writes

Example: Replicating software code

W(x1) adds a new method to a program

W(x2) calls new method

L2 knows about new method

L1: W(x1)

L2:

WS(x1)



W(x2)

L1: **W(x1)**

L2:

W(x2)

Monotonic Writes

Example: Replicating software code

W(x1) adds a new method to a program

W(x2) calls new method

L2 knows about new method

L1: W(x1)

L2:

WS(x1)

W(x2)

L1: **W(x1)**

L2:

W(x2) (

L2 doesn't know new method! Code does not compile

Monotonic writes

- Informally, writes are applied in the order in which they arrive
- For example,
 - x = 100
 - W1 is x = x + x * 10 %
 - W2 is x = x + 100
- W1, W2
 - 210
- W2, W1
 - 220

Summary: Client-centric consistency

- Eventual consistency can be used when:
 - No or few concurrent writes on the same key occur
 - Used in Dynamo, Cassandra, Riak, et al.
- Eventual consistency provides high availability and scalability
- Client-centric consistency dictates how reads and writes should look from the client's perspective
 - Read your writes
 - Monotonic-reads
 - Monotonic-writes
 - Writes follow reads

Self-study Questions

- Write some concrete operation sequences that illustrate write-write and read-write conflicts in a data store.
- Create some concurrent read and write sequences that conform to different client-centric consistency models.
- Is there a correspondence among data-centric and clientcentric consistency models? Could one be used to emulate the other, argue for or against.
- How do client-centric consistency models relate to one another, if at all?
- Find sample applications for each consistency model.
- How can client-centric consistency models be implemented?