

Stronger properties for low-latency geo-replicated storage.

BY WYATT LLOYD, MICHAEL J. FREEDMAN,
MICHAEL KAMINSKY, AND DAVID G. ANDERSEN

Don't Settle for Eventual Consistency

GEO-REPLICATED STORAGE PROVIDES copies of the same data at multiple, geographically distinct locations. Facebook, for example, geo-replicates its data (profiles, friends lists, likes, and so on) to datacenters on the East and West coasts of the U.S. and in Europe. In each datacenter, a tier of separate Web servers accepts

browser requests and then handles those requests by reading and writing data from the storage system, as shown in Figure 1.

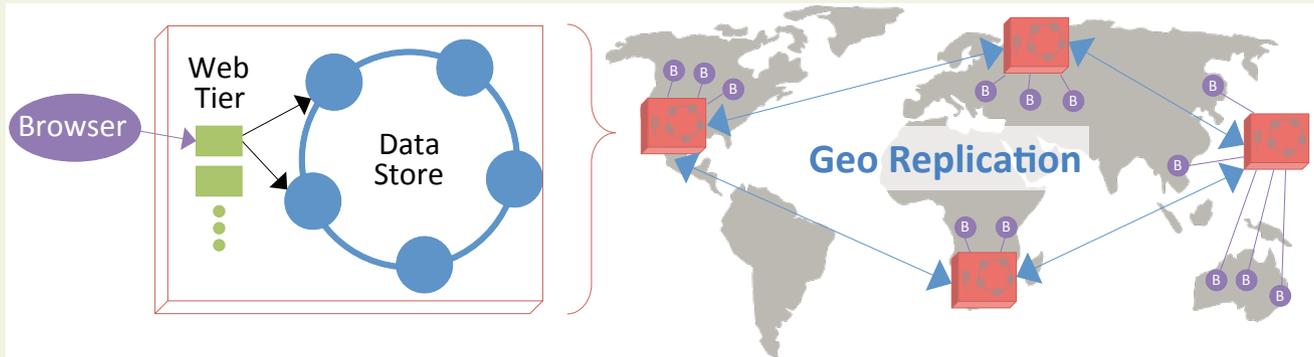
Geo-replication brings two key benefits to Web services: fault tolerance and low latency. It provides fault tolerance through redundancy: if one datacenter fails, others can continue to provide the service. It provides low latency through proximity: clients can be directed to and served by a nearby datacenter to avoid speed-of-light delays associated with cross-country or round-the-globe communication.

Geo-replication brings its challenges, however. The famous CAP theorem, conjectured by Brewer¹ and proved by Gilbert and Lynch,⁷ shows it is impossible to create a system that has strong *consistency*, is always *available* for reads and writes, and is able

to continue operating during network *partitions*. Each of these properties is highly desirable. Strong consistency—more formally known as *linearizability*—makes programming easier. Availability ensures front-end Web servers can always respond to client requests. Partition tolerance ensures the system can continue operating even when datacenters cannot communicate with one another. Faced with the choice of at most two of these properties, many systems^{5,8,16} have chosen to sacrifice strong consistency to ensure availability and partition tolerance. Other systems—for example, those that deal with money—sacrifice availability and/or partition tolerance to achieve the strong consistency necessary for the applications built on top of them.^{4,15}

The former choice of availability and partition tolerance is not surpris-

Figure 1. Geo-replicated storage.



ing, however, given it also enables the storage system to provide low latency—defined as latency for reads and writes that is less than half the speed-of-light delay between the two most distant datacenters. A proof that predates the

CAP theorem by 14 years¹⁰ shows it is impossible to guarantee low latency and provide strong consistency at the same time. Front-end Web servers read or write data from the storage system potentially many times to answer a

single request; therefore, low latency in the storage system is critical for enabling fast page-load times, which are linked to user engagement with a service—and, thus, revenue. An always-available and partition-tolerant system can provide low latency on the order of milliseconds by serving all operations in the local datacenter. A strongly consistent system must contact remote datacenters for reads and/or writes, which takes hundreds of milliseconds.

Thus, systems that sacrifice strong consistency gain much in return. They can be always available, guarantee responses with low latency, and provide partition tolerance. In Clusters of Order-preserving Servers (COPS),¹¹ developed for our original work on this subject, we coined the term ALPS for systems that provide these three properties—*always available*, *low latency*, and *partition tolerance*—and one more: *scalability*. Scalability implies that adding storage servers to each datacenter produces a proportional increase in storage capacity and throughput. Scalability is critical for modern systems because data has grown far too large to be stored or served by a single machine.

The question remains as to what consistency properties ALPS systems can provide. Before answering this, let's consider the consistency offered by existing ALPS systems. For systems such as Amazon's Dynamo, LinkedIn's Project Voldemort, and Facebook/Apache's Cassandra, the answer is *eventual consistency*.

Eventual Consistency

Eventual consistency is a widely used

Figure 2. Anomaly 1: Comment reordering.

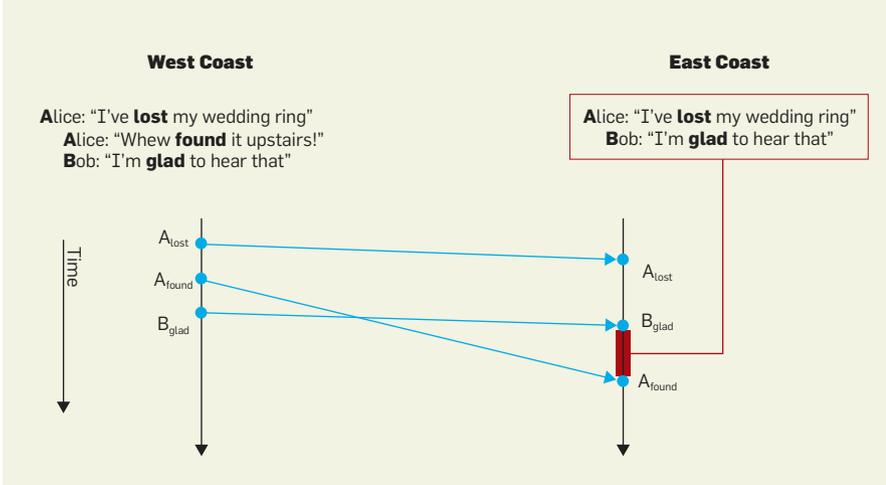
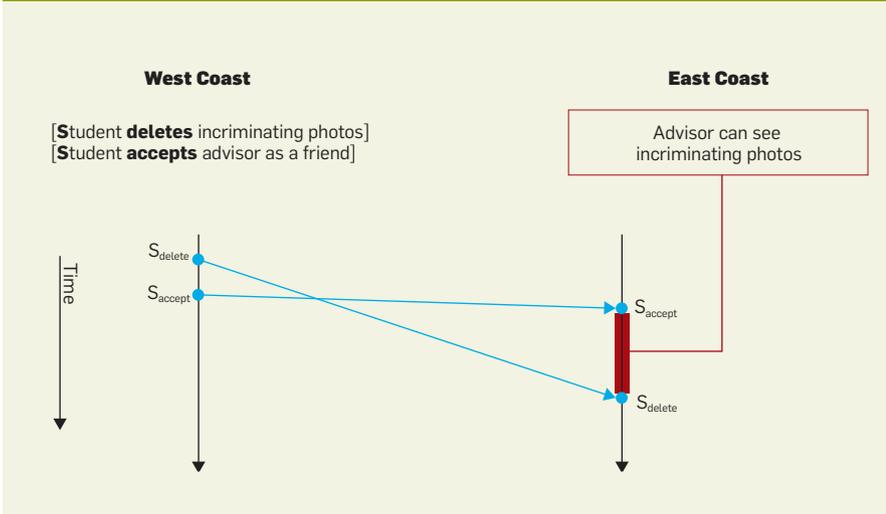


Figure 3. Anomaly 2: Photo privacy violation.



term that can have many meanings. Here, it is defined as the strongest property provided by all systems that claim to provide it: namely, writes to one datacenter will eventually appear at other datacenters, and if all datacenters have received the same set of writes, they will have the same values for all data.

Contrast this with the following part of the definition of strong consistency (linearizability): a total order exists over all operations in the system. This makes programming a strongly consistent storage system simple, or at least simpler: it behaves as a single entity. Eventual consistency does not say anything about the ordering of operations. This means that different datacenters can reflect arbitrarily different sets of operations. For example, if someone connected to the West Coast datacenter sets $A=1$, $B=2$, and $C=3$, then someone else connected to the East Coast datacenter may see only $B=2$ (not $A=1$ or $C=3$), and someone else connected to the European datacenter may see only $C=3$ (not $A=1$ or $B=2$). This makes programming eventually consistent storage complicated: operations can appear out of order.

The out-of-order arrival leads to many potential anomalies in eventually consistent systems. Here are a few examples for a social network:

Figure 2 shows that in the West Coast datacenter, Alice posts, she comments, and then Bob comments. In the East Coast datacenter, however, Alice’s comment has not appeared, making Bob look less than kind. Figure 3 shows that in the West Coast datacenter, a grad student carefully deletes incriminating photos before accepting an advisor as a friend. Unfortunately, in the East Coast datacenter, the friend-acceptance appears before the photo deletions, allowing the advisor to see the photos.³

Figure 4 shows that in the West Coast datacenter, Alice uploads photos, creates an album, and then adds the photos to the album, but in the East Coast datacenter, the operations appear out of order and her photos do not end up in the album. Finally, in Figure 5, Cindy and Dave have \$1,000 in their joint bank account. Concurrently, Dave withdraws \$1,000 from the East Coast datacenter and Cindy withdraws

\$1,000 from the West Coast datacenter. Once both withdrawals propagate to each datacenter, their account is in a consistent state ($-\$1,000$), but it is too late to prevent the mischievous couple from making off with their ill-gotten gains.

Is eventual consistency the only option? Given that theoretical results show the ALPS properties are incompatible with strong consistency, do we have to settle for eventual consistency? Are we stuck with all the anomalies that come with eventual consistency? No!

Our research systems, COPS¹¹ and Eiger,¹² have pushed on the properties that ALPS systems can provide. In particular, they provide *causal* consistency instead of eventual, which prevents the first three anomalies. (The fourth anomaly in Figure 5 is unfortunately unavoidable in a system that accepts writes in every location and guarantees

low latency.) In addition, they provide limited forms of read-only and write-only transactions that allow programmers to consistently read or write data spread across many different machines in a datacenter.

Causal Consistency

Causal consistency ensures that operations appear in the order the user intuitively expects. More precisely, it enforces a partial order over operations that agree with the notion of potential causality. If operation A happens before operation B, then any datacenter that sees operation B must see operation A first.

Three rules define potential causality:⁹

1. **Thread of Execution.** If a and b are two operations in a single thread of execution, then $a \rightarrow b$ if operation a happens before operation b .
2. **Reads-From.** If a is a write operation

Figure 4. Anomaly 3: Broken photo album.

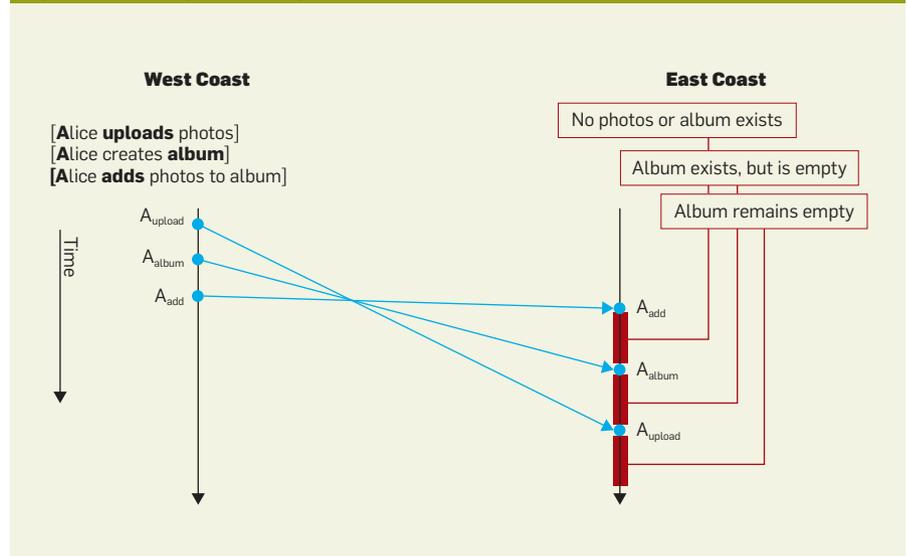


Figure 5. Anomaly 4: Double money withdrawal.

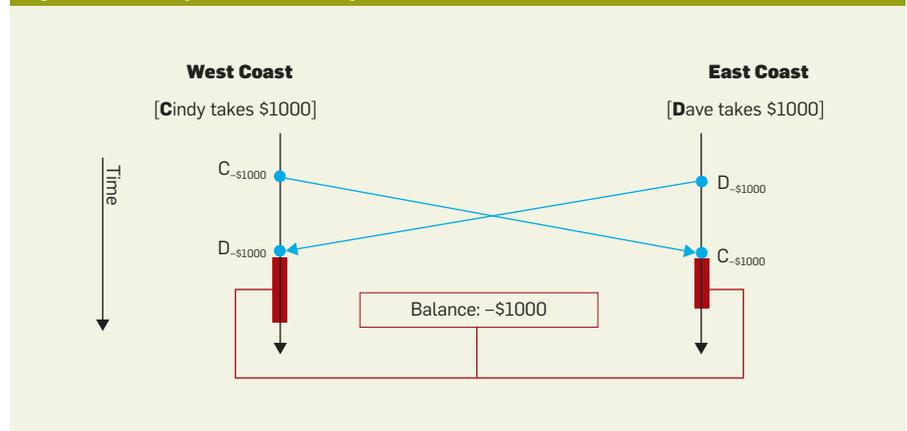
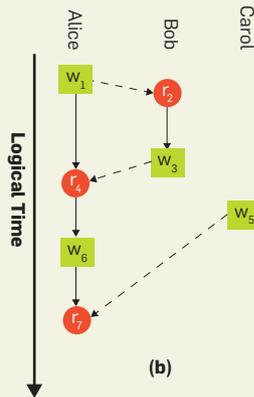


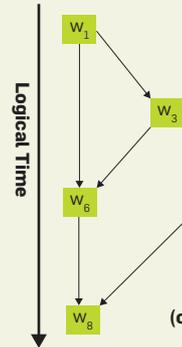
Figure 6. Causality and dependency.

User	Op ID	Operation
Alice	w_1	write(Alice:Town, NYC)
Bob	r_2	read(Alice:Town)
Bob	w_3	write(Bob:Town, LA)
Alice	r_4	read(Bob:Town))
Carol	w_5	write(Carol:Likes, ACM, 8/31/12)
Alice	w_6	write(Alice:Likes, ACM, 9/1/12)
Alice	r_7	read(Carol:Likes, ACM)
Alice	w_8	write(Alice:Friends, Carol, 9/2/12)

(a)



(b)



(c)

Op ID	Dependencies
w_1	—
w_3	w_1
w_5	—
w_6	w_5, w_1
w_8	w_6, w_5, w_3, w_1

(d)

tion and b is a read operation that returns the value written by a , then $a \rightarrow b$.

3. **Transitivity.** For operations a , b , and c , if $a \rightarrow b$ and $b \rightarrow c$, then $a \rightarrow c$. Thus, the causal relationship between operations is the transitive closure of the first two rules.

Causal consistency ensures operations appear in an order that agrees with these rules. This makes users happy because their operations are applied everywhere in the order they intended. It makes programmers happy because they no longer have to reason about out-of-order operations.

Causal consistency prevents each of our first three anomalies, turning them into *regularities*.

Regularity 1: No Missing Comments

In the West Coast datacenter, Alice posts, and then she and Bob comment:

- Op_1 Alice: I've lost my wedding ring.
- Op_2 Alice: Whew, found it upstairs.
- Op_3 [Bob reads Alice's post and comment.]
- Op_4 Bob: I'm glad to hear that.

$Op_1 \rightarrow Op_2$ by the thread-of-execution rule; $Op_2 \rightarrow Op_3$ by the reads-from rule; $Op_3 \rightarrow Op_4$ by the thread-of-execution rule.

Write operations are only propagated and applied to other datacenters, so the full causal ordering that is enforced is $Op_1 \rightarrow Op_2 \rightarrow Op_4$.

Now, in the East Coast datacenter, operations can appear only in an order that agrees with causality. Thus:

- Op_1 Alice: I've lost my wedding ring.
- Then

- Op_1 Alice: I've lost my wedding ring.
- Op_2 Alice: Whew, found it upstairs.
- Then
- Op_1 Alice: I've lost my wedding ring.
- Op_2 Alice: Whew, found it upstairs.
- Op_4 Bob: I'm glad to hear that.

but never the anomaly that makes Bob look unkind.

Regularity 2: No Leaked Photos³

In the West Coast datacenter, a grad student carefully deletes incriminating photos before accepting an advisor as a friend:

- Op_1 [Student deletes incriminating photos.]
- Op_2 [Student accepts advisor as a friend.]
- $Op_1 \rightarrow Op_2$ by the thread-of-execution rule.

Now, in the East Coast datacenter, operations can appear only in an order that agrees with causality, which is:

- [Student deletes incriminating photos.]
- Then
- [Student deletes incriminating photos.]
- [Student accepts advisor as a friend.]

but never the anomaly that leaks photos to the student's advisor.

Regularity 3: Normal Photo Album

In the West Coast datacenter, Alice uploads photos and then adds them to her Summer 2013 album:

- Op_1 [Alice uploads photos.]
- Op_2 [Alice creates an album.]
- Op_3 [Alice adds photos to the album.]
- $Op_1 \rightarrow Op_2 \rightarrow Op_3$ by the thread-of-

execution rule.

Now, in the East Coast datacenter, the operations can appear only in an order that agrees with causality:

- Op_1 [Alice uploads photos.]
- Then
- Op_1 [Alice uploads photos.]
- Op_2 [Alice creates an album.]
- Then
- Op_1 [Alice uploads photos.]
- Op_2 [Alice creates an album.]
- Op_3 [Alice adds photos to the album.]

but never in a different order that results in an empty album or complicates what a programmer must think about.

What causal consistency cannot do.

Anomaly 4 represents the primary limitation of causality consistency: it cannot enforce global invariants. Anomaly 4 has an implicit global invariant—that bank accounts cannot go below \$0—that is violated. This invariant cannot be enforced in an ALPS system. Availability dictates that operations must complete, and low latency ensures they are faster than the time it takes to communicate between datacenters. Thus, the operations must return before the datacenters can communicate and discover the concurrent withdrawals.

True global invariants are quite rare, however. E-commerce sites, where it seems inventory cannot go below 0, have back-order systems in place to deal with exactly that scenario. Even some banks do not enforce global \$0 invariants, as shown by a recent concurrent withdrawal attack on ATMs that extracted \$40 million from only 12 account numbers.¹⁴

Another limitation of causal consistency also stems from the possibility of concurrent operations. Programmers must decide how to deal with concurrent write operations to the same data at different datacenters. A common strategy is the last-writer-wins rule in which one concurrent update overwrites the other. For example, a social-network user can have only one birthday. Some situations, however, require a more careful approach. Consider a scenario where Alice has two pending friend requests being accepted concurrently at different datacenters. Each accepted friend request should increase Alice's friend count by one. With the last-writer-wins rule, however, one of the increments will overwrite the other. Instead, the two increments must be merged to increase Alice's total friend count by two. With causally consistent storage (as with eventually consistent storage), programmers must determine if the last-writer-wins rule is sufficient, or if they have to write a special function for merging concurrent updates.

The final limitation of causal consistency is it cannot see or enforce causality outside of the system. The classic example is a cross-country phone call. If Alice on the West Coast updates her profile, calls Bob on the East Coast, and then Bob updates his profile, the system will not see the causal relationship between the two updates and will not enforce any ordering between them.

Providing causal consistency. At a high level, our systems, COPS and Eiger capture causality through a client library and then enforce the observed ordering when replicating writes to other datacenters. The ordering is enforced by delaying the application of a write until all causally previous operations have been applied. This delay is necessary only in remote datacenters; all causally previous operations have already been applied at the datacenter that accepts the write. The client library that tracks causality sits between the Web servers and the storage tiers in each datacenter. (In current implementations it is on the Web servers.) Individual clients are identified through a special `actor_id` field in the API to the client library that allows the operations of different users on the same Web server to be

disentangled. For example, in a social network the unique user ID could be used as the `actor_id`.

Let's first describe an inefficient system that provides causality and then explain how to refine it to make it efficient.

Our systems operate by tracking and enforcing the ordering only between write operations. Read operations establish causal links between write operations by different clients, but they are not replicated to other datacenters and thus do not need to have an ordering enforced on them. For example, in anomaly/regularity 1, Bob's read (Op_3) of Alice's post (Op_1) and comment (Op_2) creates the causal link that orders Bob's later comment (Op_4) after Alice's post and comment. A causal link between two write operations is called a *dependency*—the later operation depends on the earlier operation.

Figure 6 shows the relationship between the graph of causality and the graph of dependencies. A dependency is a small piece of metadata that uniquely identifies a write operation. It has two fields: a key, which is the data location that is updated by the write; and a timestamp, which is a globally unique logical timestamp assigned by the logical clock of the server in the datacenter where it was originally written. Figure 6 illustrates (a) a set of example operations; (b) the graph of causality between them; (c) the corresponding dependency graph; and (d) a table listing dependencies with one-hop dependencies shown in bold.

In the initial design the client library tracks the full set of dependencies for

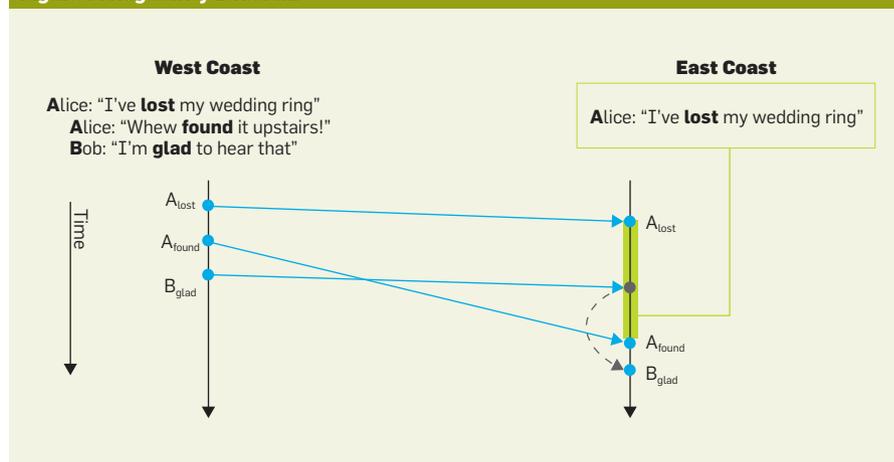
each client. Tracking all dependencies for a client requires tracking three sets of write operations:

1. All of the client's previous write operations, because of the thread-of-execution rule.
2. All of the operations that wrote values it read, because of the reads-from rule.
3. All of the operations that the operations in 1 and 2 depend on, because of the transitivity rule.

Tracking the first set is straightforward: servers return the unique timestamp assigned to each write to the client library, which then adds a dependency on that write. Tracking the second set is also straightforward: servers return the timestamp of the write that wrote the value when they respond to reads, and then the client library adds a dependency on that write. The third set of operations is a bit trickier: it requires that every write carry with it all of its dependencies, that these dependencies are stored with the value, returned with reads of that value, and then added to the reader's set of dependencies by the client library.

With the full set of dependencies for each client stored in its client library, all of these dependencies can be attached to each write operation the client issues. Now when a server in a remote datacenter receives a write with its full set of dependencies, it blocks the write and verifies each dependency is satisfied. Blocking these replicated write operations is acceptable because they are not client-facing and do not block reads to whatever data they update. Here, we have explicitly chosen to delay these write

Figure 7. Regularity 1-Redux.



operations until they can appear in the correct order, as shown in Figure 7. The dependency check for B_{grad} does not return until after A_{found} is applied on the East Coast, which ensures Bob is never misunderstood.

The system described thus far provides causal consistency and all of the ALPS properties. Causal consistency is provided by tracking causality with a client library and enforcing the causal order with dependency checks on replicated writes. Availability and partition tolerance are ensured by keeping all operations inside the local datacenter. Low latency is guaranteed by keeping all operations local, nonblocking, and lock free. Finally, a fully decentralized design ensures the system has scalability.

The current system, however, is inefficient. It has a huge amount of dependency metadata that travels around with write operations and a huge number of dependency checks

to execute before applying them. Both of these factors steal throughput from user-facing operations and reduce the utility of the system. Luckily, our systems can exploit the transitivity inherent in the graph of causality to drastically reduce the dependencies that must be tracked and enforced. The subset of dependencies being tracked is the *one-hop* dependencies, which have an arc to the current operation in the graph of causality. (Note that in graph theoretic terms, the one-hop dependencies subset is the direct predecessor set of an operation.) In Figure 6, the one-hop dependencies are shown in bold. They transitively capture all of the ordering constraints on an operation. In particular, because all other dependencies are depended upon by at least one of the one-hop dependencies by definition, if this current operation occurs after the one-hop dependencies, then by transitivity it will occur after all others as well.

Limited Transactions

In addition to causal consistency, our systems provide limited forms of transactions. These include read-only transactions, which transactionally read data spread across many servers in a datacenter, and write-only transactions, which transactionally update data spread across many servers in a datacenter.

These limited transactions are necessitated—and complicated—by the current scale of data. Data for many services is now far too large to fit on a single machine and instead must be spread across many machines. With data resident on many machines, extracting a consistent view of that data becomes tricky. Even though a data store itself may always be consistent, a client can extract an inconsistent view because the client’s reads will be served at different times by different servers. This, unfortunately, can reintroduce many of the anomalies inherent in eventual consistency. In Figure 8, for example, in the West Coast datacenter, a grad student removes photo-viewing permissions from an advisor and uploads incriminating photos. The advisor concurrently tries to view the student’s photos and, incorrectly, is shown the incriminating photos. To avoid these anomalies, causal consistency must be extended from the storage system to the Web servers and then on to users of the service. This can be done using read-only transactions.

Read-only transactions allow programmers to transactionally read data spread across many servers, yielding a consistent view of the data. The interface for a read-only transaction is simple: a list of data locations. Instead of issuing many individual reads for different data locations, a programmer issues a single read for all those locations. This is similar to batching these operations, which is often done to make dispatching reads more efficient—except it also ensures the results are isolated.

With read-only transactions, anomaly 5 can now be converted into a regularity as well. Figure 9 shows that with read-only transactions, the permissions and photos are read together transactionally, yielding any of the three valid states shown, but never the anomaly that leaks the incriminating photos to the student’s advisor.

Figure 8. Anomaly 5: Leaked photos.

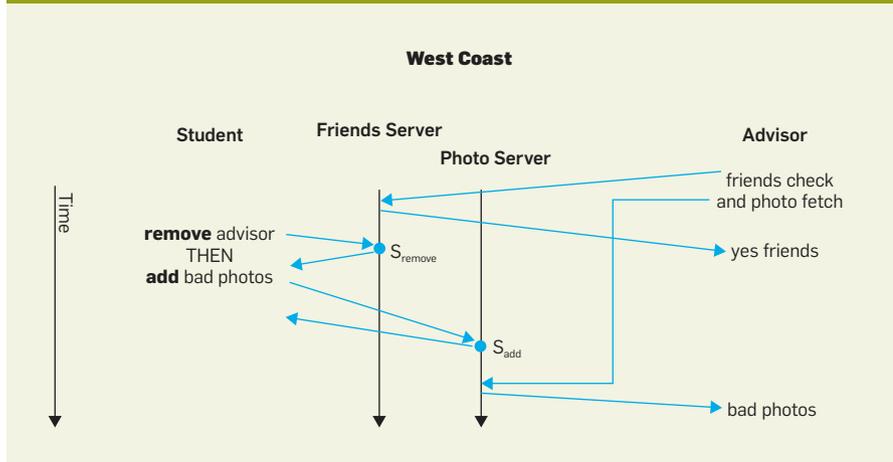
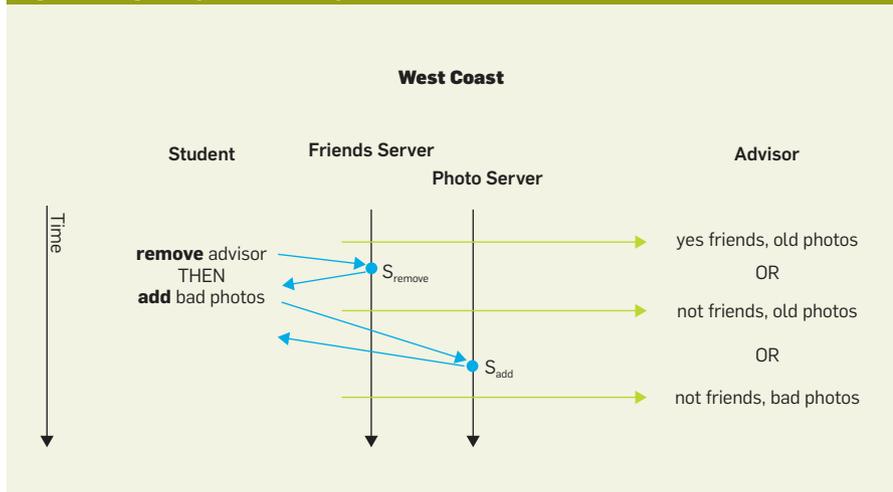


Figure 9. Regularity 5: No leaked photos.



Providing read-only transactions.

There are many techniques for ensuring isolated access to data locations spread across many machines. The most popular of these include two-phase locking (2PL), using a transaction manager (TM) to schedule when reads are applied, and maintaining multiple versions of data with multiversion concurrency control (MVCC). The first two approaches are at odds with the ALPS properties. All forms of locking, and 2PL in particular, can encounter locks that are already acquired and then either fail the operation, which gives up on availability, or block the operation until it can acquire the lock, which gives up on low latency. A TM is a centralized entity, and directing all reads through it is a bottleneck that inhibits scalability. This leaves MVCC, and our approach may be viewed as a particularly aggressive variant of it that is possible because our transactions are limited.

The basic idea behind our read-only transaction algorithm is that we want to read the entire distributed data store at a single logical time. (For logical time, each node in the system keeps a logical clock that is updated every time an event occurs (for example, writing a value or receiving a message). When sending a message, a node includes a timestamp t set to its logical clock c ; when receiving a message, a node sets $c \leftarrow \max(c, t + 1)$.⁹ Logical time LT provides a progressing logical view of the system even though it is distributed and there is no centralized coordination of it. If event a happens before event b , then $LT(a) < LT(b)$. Thus, if distributed data is read at a single logical time LT for all events seen at time t , we know all events that happen before them have lower logical times and thus are reflected in the results. Figure 10 shows an example of this graphically, with validity periods for values, represented by letters, written to different locations.

You can determine if values within a set are consistent with one another by annotating them with the logical time they became visible and then were overwritten. For example, in Figure 10 consistent sets include $\{A, J, X\}$, $\{B, K, X\}$, $\{B, K, Y\}$, $\{B, L, Y\}$ and inconsistent sets include $\{A, K, X\}$, $\{B, J, Y\}$, and $\{C, L, X\}$, among others. Our serv-

Figure 10. Logical time.

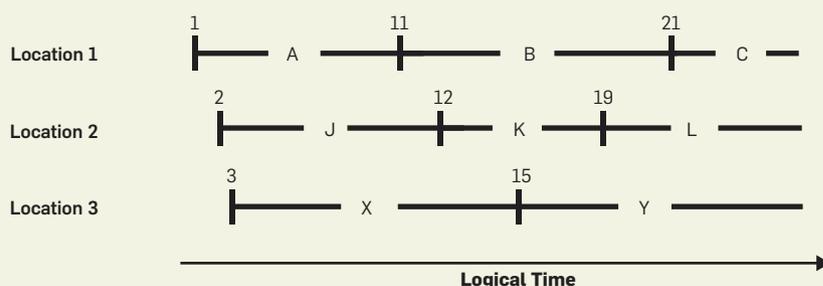
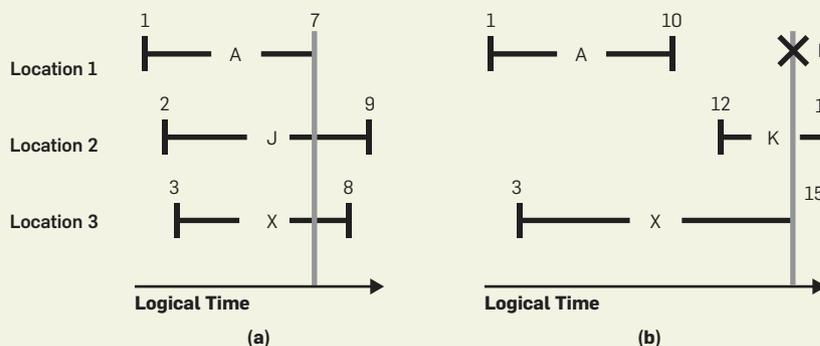


Figure 11. Read-only transactions.



ers annotate values in this way and include them when returning results to the client library so it can determine if values are mutually consistent.

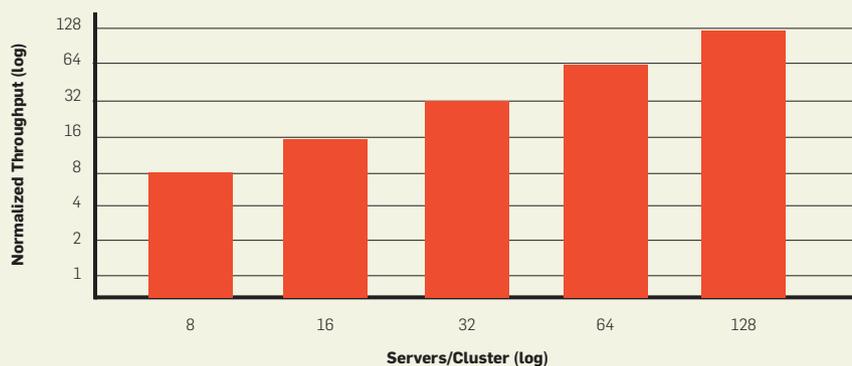
Our read-only transaction algorithm is run by the client library and takes at most two rounds of parallel reads. In the first round, the client library sends out parallel requests for all requested data. Servers respond with their current visible values and validity intervals, which is the logical time the value was written and the current logical time at the server. The value may be valid at future logical times as well, but conservatively we know it is valid for at least this interval. Once the client receives all responses, it determines if all the values are mutually consistent by checking for intersection in their validity intervals. If there is intersection—which is almost always the case unless some of the values are overwritten concurrently with the read-only transaction—then the client library knows the values are consistent and returns them to the client.

If the validity intervals do not all intersect, then the process moves to

the second round of the algorithm. The second round begins by calculating a single logical time at which to read values, called the *effective time*. It is calculated by choosing a time that ensures an up-to-date view of the data instead of being stuck on an old consistent cut of it, and it allows the use of many of the values retrieved in the first round. The client library then issues a second round of parallel reads for all data for which it does not have a valid value at the effective time. These second-round reads ask for the value of the data at the effective time, and servers answer these reads by traversing the older version of a value until they find the one that was valid at the effective time. Figure 11 shows the second round in action. Figure 11a is a read-only transaction that completes in a single round, while Figure 11b is a read-only transaction that requires a second round and requests data location 1 at the effective time 15 and receives value B in response.

This read-only transaction algorithm is specifically designed to maintain all the ALPS properties and provide high performance. It is available

Figure 12. Eiger scales linearly.



because all reads ask for a current value or an old value. It is low latency because it requires at most two non-blocking rounds of parallel reads. It is partition tolerant because all reads are in the local datacenter (partitions are assumed to occur only in the wide area, not in the local datacenter). It is scalable because it is fully decentralized. Finally, it is performant because it normally takes only a single round of parallel reads and only two rounds of reads in the worst case.

Our previous work on Eiger¹² has more details on how to choose the effective time, how to limit server-side storage of old versions, and an algorithm for write-only transactions that also maintains all the ALPS properties.

The cost of causal consistency and limited transactions. For one realistic view of Eiger's overhead, we parameterized a synthetic workload based upon Facebook's production The Associations and Objects (TAO) system.² We then compared Eiger's throughput with that of eventually consistency Cassandra from which it was forked in an experiment with clusters of eight servers each in Washington and California. The Cassandra setup achieved 23,657 operations per second that touched 498,239 data locations per second on average. The Eiger setup, with causal consistency and all inconsistent batch operations converted to read or write transactions, achieved 22,891 operations per second that touched 480,904 data locations per second on average. This experiment shows that for this real-world workload Eiger's causal consistency

and stronger semantics do not impose significant overhead.

To demonstrate the scalability of Eiger, we ran the Facebook TAO workload on N client machines that fully loaded an N -server cluster that is replicating writes to another N -server cluster (that is, the $N=128$ experiment involves 384 machines). This experiment was run on PROBE's Kodiak testbed,⁶ which provides an Emulab with exclusive access to hundreds of machines. Figure 12 shows the throughput for Eiger as N scales from eight to 128 servers/cluster. The bars show throughput normalized against the throughput of the eight-server cluster. Eiger scales out as the number of servers increases; each doubling of the number of servers increases cluster throughput by 96% on average.

More information is available in our papers on COPS¹¹ and Eiger,¹² and Wyatt Lloyd's dissertation.¹³ The code for Eiger is available from <https://github.com/wlloyd/eiger>.

Acknowledgements

This work was supported by funding from National Science Foundation Awards CSR-0953197 (CAREER), CCF-0964474, CNS-1042537 (PROBE), and CNS-1042543 (PROBE); and by Intel via the Intel Science and Technology Center for Cloud Computing (ISTC-CC). □



Related articles
on queue.acm.org

Proving the Correctness of Nonblocking Data Structures

Mathieu Desnoyers

<http://queue.acm.org/detail.cfm?id=2490873>

Eventual Consistency Today: Limitations, Extensions, and Beyond

Peter Bailis, Ali Ghodsi

<http://queue.acm.org/detail.cfm?id=2462076>

Structured Deferral: Synchronization via Procrastination

Paul E. McKenney

<http://queue.acm.org/detail.cfm?id=2488549>

References

- Brewer, E. Towards robust distributed systems. In *Proceedings of the 19th Annual ACM Symposium on Principles of Distributed Computing*, (2000)
- Bronson, N., et al. TAO: Facebook's distributed data store for the social graph. In *Proceedings of the Usenix Annual Technical Conference*, (2013).
- Cham, J. PhD Comics (June 2013); <http://www.phdcomics.com/comics.php?f=1592>.
- Corbett, J.C. et al. Spanner: Google's globally distributed database. *ACM Transactions on Computer Systems* 31, 3 (2013).
- DeCandia, G., et al. Dynamo: Amazon's highly available key-value store. In *Proceedings of the 21st ACM Symposium on OS Principles*, (2007), 205–220.
- Gibson, G., Grider, G., Jacobson, A. and Lloyd, W. PROBE: A thousand-node experimental cluster for computer systems research. *Usenix :login:* 38, 3 (2013).
- Gilbert, S. and Lynch, N. Brewer's conjecture and the feasibility of consistent, available, partition-tolerant Web services. *ACM SIGACT News* 33, 2 (2002), 51–59.
- Lakshman, A. and Malik, P. Cassandra—a decentralized structured storage system. In *The 3rd ACM SIGOPS International Workshop on Large-scale Distributed Systems and Middleware*, (2009).
- Lamport, L. Time, clocks, and the ordering of events in a distributed system. *Commun. ACM* 21, 7 (July 1978), 558–565.
- Lipton, R.J. and Sandberg, J.S. PRAM: A scalable shared memory. Technical Report TR-180-88. Princeton University, Department of Computer Science, 1988.
- Lloyd, W., Freedman, M. J. and Kaminsky, M. and Andersen, D.G. Don't settle for eventual: Scalable causal consistency for wide-area storage with COPS. In *Proceedings of the 23rd Symposium on Operating Systems Principles*: (2011), 401–416.
- Lloyd, W., Freedman, M.J., Kaminsky, M. and Andersen, D.G. Stronger semantics for low-latency geo-replicated storage. In *Proceedings of the 10th Usenix Conference on Networked Systems Design and Implementation* (2013), 313–328.
- Lloyd, W. Stronger consistency and semantics for low-latency geo-replicated storage. Ph.D. Dissertation, 2013, Princeton University.
- Santora, M. In hours, thieves took \$45 million in ATM scheme. *New York Times* (May 9, 2013).
- Sovran, Y., Power, R., Aguilera, M.K. and Li, J. Transactional storage for geo-replicated systems. In *Proceedings of the 23rd Symposium on Operating Systems Principles*, (2011) 385–400.
- Voldemort. 2013; <http://project-voldemort.com>.

Wyatt Lloyd is a postdoctoral researcher at Facebook and is to begin a position as an assistant professor at the University of Southern California. His research interests include the distributed systems and networking problems that underlie the architecture of large-scale websites, cloud computing, and big data.

Michael J. Freedman is an associate professor of computer science at Princeton University, with a research focus on distributed systems, networking, and security.

Michael Kaminsky is a senior research scientist at Intel Labs and is an adjunct faculty member of the computer science department at Carnegie Mellon University. He is part of the Intel Science and Technology Center for Cloud Computing based in Pittsburgh, PA.

David G. Andersen is an associate professor of computer science at Carnegie Mellon University. In 1995, he cofounded an Internet service provider in Salt Lake City, Utah.

Copyright held by Owner(s)/Author(s). Publication rights licensed to ACM. \$15.00.