

Disciplined Inconsistency

Double-blind submission

Abstract

Distributed systems are often wrong. Type systems can protect programmers from making mistakes, force them to handle error cases, and expose additional information, such as the quality or correctness of values.

1. Introduction

- Applications have performance requirements
 - Sometimes explicit in the form of SLAs, promising a certain latency or availability
 - Sometimes more implicit (i.e. every additional ms of latency reduces revenue)
- Constantly balancing performance vs correctness / programmability
 - If it isn't scaling well, or latencies are too high, then relax consistency in some places and hope...
- This is error prone: every time you change consistency, there are new reorderings and conditions to consider
 - new edge cases to handle, *implicit* in the consistency model
 - accidentally leak into places that weren't intended to be weakened
- Worse: conditions can change at any moment; node goes down, network unreliable, traffic surges
 - In test environment, inconsistency is typically unlikely
 - Adverse conditions in production can cause errors that never appeared in testing, or are very difficult to test for
 - No way to know if you've caught them all
- Furthermore, when conditions are good, there's no need to resort to weak consistency
- It would be great if we had a way to:
 - Express performance bounds
 - Have the system help achieve them
 - Make inconsistency explicit and restricted
 - handle different cases in a disciplined way
 - restrict possible values, and where they can be used
- So the question is: *where to introduce this abstraction?*
 - As part of the data type!
 - Couples the effects of mutating operations with reads
 - Concise and modular: re-use data types, no annotations on individual operations
 - Safe: inconsistency expressed as return types

2. Type System

- IPA type lattice
 - Inconsistent (\perp)
 - Rushed | Interval | Leased
 - Consistent (\top)
- Performance bounds
 - `Set[ID] with Consistency(Strong)`
 - `Set[ID] with LatencyBound(50 ms) -> contains(ID): Rushed[Boolean]`
 - `Counter with ErrorTolerance(5%) -> read(): Interval[Long]`
- Semantics of mixed consistency levels?
- Futures
 - (talk about how everything is implemented with futures, or just elide that?)

3. Implementation

We demonstrate one possible instantiation of the Disciplined Inconsistency model with an implementation of a Scala client, using Cassandra as the backing store. Most of the functionality required to implement the model is relatively datastore-agnostic; most Dynamo-style datastores support some form of tunable consistency, so porting our implementation to another backing datastore such as Riak should be possible.

[explain the basics of how Cassandra's consistency levels work] see: [Cassandra Consistency](#)

3.1. Latency bounds

As discussed earlier, a common desire is to be able to guarantee a certain response time, for example in order to meet an SLA. However, within that window of time, we would like to provide the strongest guarantees possible, so that users typically observe consistent, up-to-date data.

Conceptually, any Dynamo-style datastore implements configurable consistency levels by adjusting the number of replicas that a client request waits for a response from. [explain first how it would work conceptually by sending read requests to a quorum of replicas, and then proceeding with whatever we have when time is up; then explain how in Cassandra we have to cheat by issuing reads in parallel]

3.2. Reservations

In order to implement the `Interval` bounds, we build on the concept of *escrow* and *reservations* [10, 17–19].

We implement reservations as a middleware layer: a reservation server runs alongside each Cassandra server. Any operations with error tolerance bounds are routed to a reservation server, using the Cassandra client's knowledge of which replicas are up.

3.3. Leases

[???

4. Related Work

4.1. Consistency Models

A vast number of consistency models have been proposed over the years. From Lamport's *sequential consistency* [13] and Herlihy's *linearizability* [12] on the strong side, to *eventual consistency* [25] at the other extreme. A variety of intermediate models fit elsewhere in the spectrum, each making different trade-offs balancing high performance and availability against ease of programming. For example, a family of models including *read-your-writes* and *monotonic reads* use *sticky sessions* [23], which reduces availability in a small way, but provides users with a bit more certainty about what values they will observe.

A single global consistency model for an entire database or application is restrictive; some datastores support configuring consistency at a finer granularity: Cassandra [3] per operation, Riak [5] on an object or namespace granularity, as well as others [14, 22].

4.2. Explicit performance bounds

It is difficult for programmers to determine the correct consistency level for each operation. Ideally, everything would be as consistent as possible, but in some situations, performance needs (such as availability) force inconsistency.

[will probably have to introduce this earlier when explaining Rushed, but putting the text here for now] With *consistency-based SLAs* in Pileus [24], programmers can explicitly trade off consistency for latency. A consistency SLA specifies a target latency and a consistency level (e.g. 100 ms with read-my-writes). In this programming model, operations specify a set of desired SLAs, each associated with a *utility*. Using a prediction mechanism similar to PBS, Pileus attempts to determine which SLA to target to maximize utility, typically to achieve the best consistency possible within a certain latency.

In Pileus, SLAs are specified on each *read* operation, which returns both the value it got and the achieved consistency level. This allows programs to behave different depending on changing conditions. Our Rushed IPA types, which were inspired by Pileus, provide a more disciplined way to let programmers express how behavior should depend on consistency, protecting them from inadvertently misusing the returned value. In addition, Pileus's SLAs are assigned only to individual reads; writes are all assumed to be the same, and data type is not considered. Working with latency

bounds at the ADT level allows reads and writes to be coupled, enabling more potential optimizations.

[are there other systems with explicit performance bounds enforced by the system?]

4.3. Controlling staleness

Most eventually consistent models provides no guarantees about how long it will take for updates to propagate. However, there are several techniques to help bound the staleness of reads.

Leases are an old technique that essentially gives reads an *expiration date*: the datastore promises not to modify the value that was just read until the lease term is over. First proposed to avoid explicit invalidations in distributed file system caches [11], leases have since been used in a multitude of ways: in Facebook's Memcache system [16] for invalidations, Google's Chubby [7] and Spanner [9] to adjust the frequency of heartbeat messages, and on mobile clients with *exo-leases* [21]. Warranties [15] are a generalization of leases, allowing arbitrary assertions over state or behavior.

[explain how our leases relate (if they get implemented)]

[Probabilistically bounded staleness [4]]

4.4. Types for distributed systems

Convergent (or *conflict-free*) *replicated data types* (CRDTs) [20] are data types designed for eventual consistency. Similar to how IPA types express weakened semantics which allow for implementation on weak consistency, CRDTs guarantee that they will converge on eventual consistency by forcing all update operations to commute. For example, Set add and remove typically do not commute, but a CRDT called an OR-Set re-defines them so that add wins over remove, making them commute again. CRDTs can be enormously useful because they allow concurrent updates with sane semantics, but they are still only eventually (or causally) consistent, so users must still deal with temporary divergence and out-of-date reads, and they do not incorporate performance bounds or variable accuracy.

Bloom [1, 2, 8] is a language and runtime system for defining whole applications that are guaranteed to converge. Based around a conceptual monotonically growing set of facts, the language encourages coordination-free computation, but automatically creates synchronization points where necessary.

[Session types?]

4.5. Approximate types / Trading off correctness

[cite some approximate computing papers?]

[Something something Uncertain<T> 6]

[Conit-based Continuous Consistency Model 26]

References

- [1] Peter Alvaro, Neil Conway, Joe Hellerstein, and William R Marczak. Consistency analysis in bloom: a calm and col-

- lected approach. In *Conference on Innovative Data Systems Research (CIDR)*, CIDR, pages 249–260. Citeseer, 2011.
- [2] Peter Alvaro, Neil Conway, Joseph M. Hellerstein, and David Maier. Blazes: Coordination analysis for distributed programs. In *IEEE International Conference on Data Engineering*. Institute of Electrical & Electronics Engineers (IEEE), March 2014. doi:[10.1109/icde.2014.6816639](https://doi.org/10.1109/icde.2014.6816639).
 - [3] Apache Software Foundation. Cassandra. <http://cassandra.apache.org/>, 2015.
 - [4] Peter Bailis, Shivaram Venkataraman, Michael J. Franklin, Joseph M. Hellerstein, and Ion Stoica. Probabilistically bounded staleness for practical partial quorums. *Proceedings of the VLDB Endowment*, 5 (8): 776–787, April 2012. doi:[10.14778/2212351.2212359](https://doi.org/10.14778/2212351.2212359).
 - [5] Basho Technologies, Inc. Riak. <http://docs.basho.com/riak/latest/>, 2015.
 - [6] James Bornholt, Todd Mytkowicz, and Kathryn S. McKinley. Uncertain<T>: A First-Order Type for Uncertain Data. In *Proceedings of the 19th International Conference on Architectural Support for Programming Languages and Operating Systems - ASPLOS 14*, ASPLOS. Association for Computing Machinery (ACM), 2014. doi:[10.1145/2541940.2541958](https://doi.org/10.1145/2541940.2541958).
 - [7] Mike Burrows. The chubby lock service for loosely-coupled distributed systems. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation*, pages 335–350. USENIX Association, 2006.
 - [8] Neil Conway, William R. Marczak, Peter Alvaro, Joseph M. Hellerstein, and David Maier. Logic and lattices for distributed programming. In *Proceedings of the Third ACM Symposium on Cloud Computing - SoCC 12*, SoCC. ACM Press, 2012. doi:[10.1145/2391229.2391230](https://doi.org/10.1145/2391229.2391230).
 - [9] James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, J. J. Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson Hsieh, Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David Mwaura, David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Szymaniak, Christopher Taylor, Ruth Wang, and Dale Woodford. Spanner: Google’s globally-distributed database. In *USENIX Conference on Operating Systems Design and Implementation*, OSDI, pages 251–264, 2012. ISBN 978-1-931971-96-6. URL <http://dl.acm.org/citation.cfm?id=2387880.2387905>.
 - [10] Dieter Gawlick and David Kinkade. Varieties of Concurrency Control in IMS/VS Fast Path. *IEEE Database Engineering Bulletin*, 8 (2): 3–10, 1985.
 - [11] C. Gray and D. Cheriton. Leases: an efficient fault-tolerant mechanism for distributed file cache consistency. In *ACM Symposium on Operating Systems Principles (SOSP)*, SOSP. Association for Computing Machinery (ACM), 1989. doi:[10.1145/74850.74870](https://doi.org/10.1145/74850.74870).
 - [12] Maurice P. Herlihy and Jeannette M. Wing. Linearizability: a correctness condition for concurrent objects. *ACM Transactions on Programming Languages and Systems*, 12 (3): 463–492, July 1990. doi:[10.1145/78969.78972](https://doi.org/10.1145/78969.78972).
 - [13] Lamport. How to make a multiprocessor computer that correctly executes multiprocess programs. *IEEE Transactions on Computers*, C-28 (9): 690–691, September 1979. doi:[10.1109/tc.1979.1675439](https://doi.org/10.1109/tc.1979.1675439).
 - [14] Cheng Li, Daniel Porto, Allen Clement, Johannes Gehrke, Nuno Preguiça, and Rodrigo Rodrigues. Making geo-replicated systems fast as possible, consistent when necessary. In *Presented as part of the 10th USENIX Symposium on Operating Systems Design and Implementation (OSDI 12)*, pages 265–278, Hollywood, CA, 2012. USENIX. ISBN 978-1-931971-96-6. URL <https://www.usenix.org/conference/osdi12/technical-sessions/presentation/li>.
 - [15] Jed Liu, Tom Magrino, Owen Arden, Michael D. George, and Andrew C. Myers. Warranties for faster strong consistency. In *USENIX Symposium on Networked Systems Design and Implementation (NSDI 14)*, pages 503–517, Seattle, WA, April 2014. USENIX Association. ISBN 978-1-931971-09-6. URL https://www.usenix.org/conference/nsdi14/technical-sessions/presentation/liu_jed.
 - [16] Rajesh Nishtala, Hans Fugal, Steven Grimm, Marc Kwiakowski, Herman Lee, Harry C. Li, Ryan McElroy, Mike Paleczny, Daniel Peek, Paul Saab, David Stafford, Tony Tung, and Venkateshwaran Venkataramani. Scaling memcache at facebook. In *Presented as part of the 10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*, pages 385–398, Lombard, IL, 2013. USENIX. ISBN 978-1-931971-00-3. URL <https://www.usenix.org/conference/nsdi13/technical-sessions/presentation/nishtala>.
 - [17] Patrick E. O’Neil. The escrow transactional method. *ACM Transactions on Database Systems*, 11 (4): 405–430, December 1986. doi:[10.1145/7239.7265](https://doi.org/10.1145/7239.7265).
 - [18] Nuno Preguiça, J. Legatheaux Martins, Miguel Cunha, and Henrique Domingos. Reservations for conflict avoidance in a mobile database system. In *Proceedings of the 1st international conference on Mobile systems, applications and services - MobiSys 03*, MobiSys. Association for Computing Machinery (ACM), 2003. doi:[10.1145/1066116.1189038](https://doi.org/10.1145/1066116.1189038).
 - [19] Andreas Reuter. *Concurrency on high-traffic data elements*. ACM, New York, New York, USA, March 1982.
 - [20] Marc Shapiro, Nuno Preguiça, Carlos Baquero, and Marek Zawirski. Conflict-free Replicated Data Types. In *Proceedings of the 13th International Conference on Stabilization, Safety, and Security of Distributed Systems*, SSS, pages 386–400, 2011. ISBN 978-3-642-24549-7.
 - [21] Liuba Shrira, Hong Tian, and Doug Terry. Exo-leasing: Escrow synchronization for mobile clients of commodity storage servers. In *Middleware 2008*, Middleware, pages 42–61. Springer Science & Business Media, 2008. doi:[10.1007/978-3-540-89856-6_3](https://doi.org/10.1007/978-3-540-89856-6_3).
 - [22] Yair Sovran, Russell Power, Marcos K. Aguilera, and Jinyang Li. Transactional storage for geo-replicated systems. In *ACM Symposium on Operating Systems Principles - SOSP’11*, SOSP. Association for Computing Machinery (ACM), 2011. doi:[10.1145/2043556.2043592](https://doi.org/10.1145/2043556.2043592).
 - [23] D.B. Terry, A.J. Demers, K. Petersen, M.J. Spreitzer, M.M. Theimer, and B.B. Welch. Session guarantees for weakly

- consistent replicated data. In *Proceedings of 3rd International Conference on Parallel and Distributed Information Systems*, PDIS. Institute of Electrical & Electronics Engineers (IEEE), 1994. doi:[10.1109/pdis.1994.331722](https://doi.org/10.1109/pdis.1994.331722).
- [24] Douglas B. Terry, Vijayan Prabhakaran, Ramakrishna Kotla, Mahesh Balakrishnan, Marcos K. Aguilera, and Hussam Abu-Libdeh. Consistency-based service level agreements for cloud storage. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles - SOSP 13*. ACM Press, 2013. doi:[10.1145/2517349.2522731](https://doi.org/10.1145/2517349.2522731).
- [25] Werner Vogels. Eventually consistent. *Communications of the ACM*, 52 (1): 40, January 2009. doi:[10.1145/1435417.1435432](https://doi.org/10.1145/1435417.1435432).
- [26] Haifeng Yu and Amin Vahdat. Design and evaluation of a conit-based continuous consistency model for replicated services. *ACM Transactions on Computer Systems (TOCS)*, 20 (3): 239–282, 2002.