Understanding Eventual Consistency

Hongseok Yang University of Oxford

Joint work with Sebastian Burckhardt (MSR), Alexey Gotsman (IMDEA)

Monday, 11 February 13

C Reader 0 4 1 + a www.amazon.co.uk — Amazon.co.uk: Low Prices in Electronics, Books, Sports Equipment & more ** **January Deals** amazon > Shop now HONGSEOK's Amazon Today's Deals Gift Cards Sell Help Shop by Department Wish Hello, HONGSEOK Go Search AI -Basket -Your Account -List -MP3s & Cloud Player New Year. Amazon MP3 Cloud Player Kindle LOVEFILM Appstore for Android 20 million songs, play anywhere Audible **New Resolutions** Amazon Cloud Drive 5 GB of free storage > See more Appstore for Android Meet the Get a paid app for free every day Kindle Family Join amazonfamily Books Music, Games, Film & TV Protect and personalise Kindle Kindle Paperwhite Kindle Fire HD Kindle your Kindle >£69 > from £109 > from £159 Electronics > Kindle Accessories Save up to £50 every month Computers & Office and get three months' FREE Home, Garden & Pets One-Day Delivery Toys, Children & Baby Amazon Family Trade-In Clothing Store Amazon Prime Subscribe & Save Clothes, Shoes & Watches FASHION > Clothing Sports & Outdoors Just some of the benefits of Amazon Family up to 60% Off Become a member Grocery, Health & Beauty Shoes & Bags DIY, Tools & Car up to 60% Off Advertisement Jewellery > Full Shop Directory up to 70% Off Customised Watches Photo Calendar up to 75% Off

More Items to Consider



You viewed

The Signal and the Noise: The Art and ... Nate Silver Hardcover ***** (9) 00.40





Customers who viewed this also viewed

Antifragile: How to Live in a

World.. Nassim Taleb Hardcover ******* (26)

AL



companies... Ben Goldacre Paperback ****** (97)

LOOK INSIDE



The Theory That Would Not Die: How.... Sharon Bertsch McGrayne Paperback ******** (9)



LOOK INSIDE

SIGNAL

NOISE

***** ----





Geo-replicated systems

Every data centre stores a complete replica of data
Purpose: fault tolerance, minimising latency

Geo-replicated systems

Every data centre stores a complete replica of data
Purpose: fault tolerance, minimising latency

Geo-replicated systems

Every data centre stores a complete replica of data
Purpose: fault tolerance, minimising latency



- Any user can access any replica
- Updating only one replica leads to weak consistency - allows strange behaviours



- Any user can access any replica
- Updating only one replica leads to weak consistency - allows strange behaviours



- Any user can access any replica
- Updating only one replica leads to weak consistency - allows strange behaviours



- Any user can access any replica
- Updating only one replica leads to weak consistency - allows strange behaviours



- Strong consistency requires updating multiple replicas before completing a user request
- Problematic in geo-replicated systems: links may go down \Rightarrow the network partitions



- Strong consistency requires updating multiple replicas before completing a user request
- Problematic in geo-replicated systems: links may go down \Rightarrow the network partitions



- Strong consistency requires updating multiple replicas before completing a user request
- Problematic in geo-replicated systems: links may go down \Rightarrow the network partitions



Also an issue with mobile devices









- Update other replicas later, when the connectivity is restored: give up strong consistency
- Wait until the connectivity is restored: give up availability
- Assume network partitions don't happen: give up partition tolerance



- Update other replicas later, when the connectivity is restored: give up strong Consistency
- Wait until the connectivity is restored: give up Availability
- Assume network partitions don't happen: give up Partition tolerance

CAP Theorem [Lynch&Gilbert]: Achieving all three at once is impossible



- Update other replicas later, when the connectivity is restored: give up strong Consistency eventual consistency
- Wait until the connectivity is restored: give up Availability
- Assume network partitions don't happen: give up Partition tolerance

CAP Theorem [Lynch&Gilbert]: Achieving all three at once is impossible

"If no new updates are made to the object, eventually all accesses will return the last updated value"

practice

DOI:10.1145/1435417.1435432

Building reliable distributed systems at a worldwide scale demands trade-offs between consistency and availability.

BY WERNER VOGELS

Eventually Consistent

AT THE FOUNDATION of Amazon's cloud computing are infrastructure services such as Amazon's S3 (Simple Storage Service), SimpleDB, and EC2 (Elastic Compute Cloud) that provide the resources for constructing Internet-scale computing platforms and a great variety of applications. The requirements placed on these infrastructure services are very strict; they need to score high marks in the areas of security, scalability, availability, performance, and cost-effectiveness, and they need to meet these requirements while serving millions of customers around the globe, continuously.

Under the covers these services are massive distributed systems that operate on a worldwide scale. This scale creates additional challenges, because when a system processes trillions and trillions of requests, events that normally have a low probability of occurrence are now guaranteed to happen and must be accounted for upfront in the design and architecture of the system. Given the worldwide scope of these systems, we use replication techniques ubiquitously to guarantee consistent performance and high availability. Although replication brings us closer to our goals, it cannot achieve them in a perfectly

transparent manner; under a number of conditions the customers of these services will be confronted with the consequences of using replication techniques inside the services.

One of the ways in which this manifests itself is in the type of data consistency that is provided, particularly when many widespread distributed systems provide an eventual consistency model in the context of data replication. When designing these largescale systems at Amazon, we use a set of guiding principles and abstractions related to large-scale data replication and focus on the trade-offs between high availability and data consistency. Here, I present some of the relevant background that has informed our approach to delivering reliable distributed systems that must operate on a global scale. (An earlier version of this article appeared as a posting on the "All Things Distributed" Weblog and was greatly improved with the help of its readers.)

Historical Perspective

In an ideal world there would be only one consistency model: when an update is made all observers would see that update. The first time this surfaced as difficult to achieve was in the database systems of the late 1970s. The best "period piece" on this topic is "Notes on Distributed Databases" by Bruce Lindsay et al.5 It lays out the fundamental principles for database replication and discusses a number of techniques that deal with achieving consistency. Many of these techniques try to achieve distribution transparency-that is, to the user of the system it appears as if there is only one system instead of a number of collaborating systems. Many systems during this time took the approach that it was better to fail the complete system than to break this transparency.2

In the mid-1990s, with the rise of larger Internet systems, these practices were revisited. At that time people began to consider the idea that availability was perhaps the most impor"If no new updates are made to the object, eventually all accesses will return the last updated value"

But updates never stop!

So what does this tell to me as a client?

practice

DOI:10.1145/1435417.1435432

Building reliable distributed systems at a worldwide scale demands trade-offs between consistency and availability.

BY WERNER VOGELS

Eventually Consistent

AT THE FOUNDATION of Amazon's cloud computing are infrastructure services such as Amazon's S3 (Simple Storage Service), SimpleDB, and EC2 (Elastic Compute Cloud) that provide the resources for constructing Internet-scale computing platforms and a great variety of applications. The requirements placed on these infrastructure services are very strict; they need to score high marks in the areas of security, scalability, availability, performance, and cost-effectiveness, and they need to meet these requirements while serving millions of customers around the globe, continuously.

Under the covers these services are massive distributed systems that operate on a worldwide scale. This scale creates additional challenges, because when a system processes trillions and trillions of requests, events that normally have a low probability of occurrence are now guaranteed to happen and must be accounted for upfront in the design and architecture of the system. Given the worldwide scope of these systems, we use replication techniques ubiquitously to guarantee consistent performance and high availability. Although replication brings us closer to our goals, it cannot achieve them in a perfectly

transparent manner; under a number of conditions the customers of these services will be confronted with the consequences of using replication techniques inside the services.

One of the ways in which this manifests itself is in the type of data consistency that is provided, particularly when many widespread distributed systems provide an eventual consistency model in the context of data replication. When designing these largescale systems at Amazon, we use a set of guiding principles and abstractions related to large-scale data replication and focus on the trade-offs between high availability and data consistency. Here, I present some of the relevant background that has informed our approach to delivering reliable distributed systems that must operate on a global scale. (An earlier version of this article appeared as a posting on the "All Things Distributed" Weblog and was greatly improved with the help of its readers.)

Historical Perspective

In an ideal world there would be only one consistency model: when an update is made all observers would see that update. The first time this surfaced as difficult to achieve was in the database systems of the late 1970s. The best "period piece" on this topic is "Notes on Distributed Databases" by Bruce Lindsay et al.5 It lays out the fundamental principles for database replication and discusses a number of techniques that deal with achieving consistency. Many of these techniques try to achieve distribution transparency-that is, to the user of the system it appears as if there is only one system instead of a number of collaborating systems. Many systems during this time took the approach that it was better to fail the complete system than to break this transparency.2

In the mid-1990s, with the rise of larger Internet systems, these practices were revisited. At that time people began to consider the idea that availability was perhaps the most impor-

50 shades of eventual consistency

Session Guarantees for Weakly Consistent Replicated Data

Douglas B. Terry, Alan J. Demers, Karin Petersen, Mike J. Spreitzer, Marvin M. Theimer, and Brent B. Welch

Don't Settle for Eventual: Scalable Causal Consistency for Wide-Area Storage with COPS

Wyatt Lloyd*, Michael J. Freedman*, Michael Kaminsky[†], and David G. Andersen[‡] *Princeton University, [†]Intel Labs, [‡]Carnegie Mellon University

Conflict-free Replicated Data Types *

Marc Shapiro^{1,5}, Nuno Preguiça^{2,1}, Carlos Baquero³, and Marek Zawirski^{1,4}

¹ INRIA, Paris, France
 ² CITI, Universidade Nova de Lisboa, Portugal
 ³ Universidade do Minho, Portugal
 ⁴ UPMC, Paris, France
 ⁵ LIP6, Paris, France

Transactional storage for geo-replicated systems

Yair Sovran* Russell Power* Marcos K. Aguilera† Jinyang Li* *New York University †Microsoft Research Silicon Valley Strengthen consistency (somewhat)

Add features that make coping with weak consistency easier

Problem

- Different formalisms/levels of abstraction: how do I compare systems?
- Tied to implementation: what do I tell the programmer/verification person?



 How do I combine different features/ explore the design space?

Need a formal and declarative definition for the semantics of eventually consistent systems



Update the local replica now, propagate to others later





Replicas diverge \Rightarrow eventual consistency violated



Updates have to commute \Rightarrow convergence guaranteed



Timestamp updates if they don't commute Not only physical time [Lamport 78]

Replicated data types

- Lots more interesting ways to make operations commutative and specify a conflict resolution policy (an example shortly)
- More complicated data types than counters or read/write registers: sets, graphs, sequences, file systems [Shapiro+ 2011]
- Every object in the database can be assigned a replicated data type → conflict resolution policy







access.write(all)



access.write(noboss)

post.write(photo)

Monday, 11 February 13





Formalising eventual consistency

- I. How do we resolve conflicts resulting from concurrent updates at different replicas?
- 2. Does a user's knowledge change in line with the system evolution?

- First, define a notion of an execution
- Weak consistency

 straightforward
 interleaving semantics doesn't work
- Similar to axiomatic memory models





What happens on the interface client/database





The order of submission to the database by a session



Execution: (A, so, vis, ar)





Affects the results of operations





The order of timestamps, used for conflict resolution



System specification = set of valid executions

- I. How do we resolve conflicts resulting from concurrent updates at different replicas?
 Data type specification
- 2. Does a user's knowledge change in line with the system evolution?

Consistency axioms



System specification = set

- I. How do we resolve conflicts resul concurrent updates at different re
 Data type specification
- Does a user's knowledge change in the system evolution?
 Consistency axioms



Figure 1. Axioms of eventual consistency WELL-FORMEDNESS AXIOMS SOWF: so is the union of transitive, irreflexive and total orders on actions by each session VISWF: $\forall a, b. a \xrightarrow{\text{vis}} b \implies \text{obj}(a) = \text{obj}(b)$ ARWF: $\forall a, b, a \xrightarrow{ar} b \implies \mathsf{obj}(a) = \mathsf{obj}(b),$ ar is transitive and irreflexive, and $\operatorname{ar}_{\operatorname{vis}^{-1}(a)}$ is a total order for all $a \in A$ AUXILIARY RELATIONS Per-object session order: $soo = (so \cap sameobj)$ Per-object causality order: $hbo = (soo \cup vis)^+$ Causality order: $hb = (so \cup vis)^+$ BASIC EVENTUAL CONSISTENCY AXIOMS RVAL: $\forall a \in A$. $\operatorname{rval}(a) = F_{\operatorname{type}(a)}(\operatorname{cone}(a))$ EVENTUAL: $\forall a \in A. \neg (\exists infinitely many \ b \in A. sameobj(a, b) \land \neg (a \xrightarrow{\mathsf{vis}} b))$ THINAIR: so ∪ vis is acyclic SESSION GUARANTEES RYW (Read Your Writes): soo \subseteq vis MR (Monotonic Reads): (vis; soo) \subseteq vis WFRV (Writes Follow Reads in Visibility): (vis; soo^{*}; vis) \subseteq vis WFRA (Writes Follow Reads in Arbitration): (vis; soo^{*}) \subseteq ar MWV (Monotonic Writes in Visibility): $(soo; vis) \subseteq vis$ MWA (Monotonic Writes in Arbitration): $soo \subseteq ar$ CAUSALITY AXIOMS POCV (Per-Object Causal Visibility): hbo \subseteq vis POCA (Per-Object Causal Arbitration): $hbo \subseteq ar$ COCV (Cross-Object Causal Visibility): $(hb \cap sameobj) \subseteq vis$ COCA (Cross-Object Causal Arbitration): hb ∪ ar is acyclic

	Figure 1. Axioms of eventual consistency
	WELL-FORMEDNESS AXIOMS
System specification $=$ set	SOWE: so is the union of transitive, irreflexive and total orders
Jocenn specification set	on actions by each session
	$\mathbf{VISWF}: \forall a, b, a \xrightarrow{vis} b \implies obj(a) = obj(b)$
I How do we resolve conflicts resul	$ARWF: \forall a, b, a \xrightarrow{ar} b \implies obj(a) = obj(b),$
	ar is transitive and irreflexive, and
concurrent updates at different re	$\operatorname{ar} _{\operatorname{vis}^{-1}(a)}$ is a total order for all $a \in A$
	AUXILIARY RELATIONS
Data type specification	Per-object session order: $soo = (so \cap sameobj)$
	Per-object causality order: $hbo = (soo \cup vis)^+$
	Causality order: $hb = (so \cup vis)^+$
2 Doos a usor's knowledge change in	PASIC EVENTUAL CONSISTENCY AVIONS
Z. DOES à user s'Knowledge change n	$\mathbf{P}_{\mathbf{A}\mathbf{A}}: \forall \mathbf{a} \in \mathcal{A} \mathbf{p}_{\mathbf{C}\mathbf{A}}(\mathbf{a}) = F_{\mathbf{C}\mathbf{A}} \cdot (\mathbf{c} \mathbf{c} \mathbf{p}_{\mathbf{C}}(\mathbf{a}))$
the system evolution?	EVENTUAL: $\forall a \in A$. $\forall \forall a(a) = r_{type(a)}(cone(a))$
	$\forall a \in A, \neg(\exists \text{ infinitely many } b \in A, \text{ sameobi}(a, b) \land \neg(a \xrightarrow{\text{vis}} b))$
Consistency axioms	THINAIR: so \cup vis is acyclic
	SESSION GUARANTEES
	RYW (Read Your Writes): soo \subseteq vis
	MR (Monotonic Reads): (vis; soo) \subseteq vis
access.write(all)	WFRV (Writes Follow Reads in Visibility): (vis; soo [*] ; vis) \subseteq vis
	WFRA (Writes Follow Reads in Arbitration): $(vis; soo^*) \subseteq ar$
ar: so	MWV (Monotonic Writes in Visibility): $(soo; vis) \subseteq vis$
access write (nobess) post read() : photo	MWA (Monotonic Writes in Arbitration): soo \subseteq ar
access.write(noboss) post.read() : photo	
vis	CAUSALITY AXIOMS
vis vis	POCV (Per-Object Causal Visibility): hbo \subseteq vis
post write(photo)	POCA (Per-Object Causal Arbitration): hbo \subseteq ar
	COCV (Cross-Object Causal Visibility): (hb \cap sameobj) \subseteq vis
	COCA (Cross-Object Causal Arbitration): hb ∪ ar is acyclic

- How do I compute the return value of an action a?
- Only actions visible to a are important: have been delivered to the replica performing a



- How do I compute the return value of an action a?
- Only actions visible to a are important: have been delivered to the replica performing a



F: context of a \rightarrow return value of a

QUERY. $\forall a \in A$. $\mathsf{retval}(a) = F_{\mathsf{type}(a)}(\mathsf{ctxt}(a))$



F: context of a \rightarrow return value of a

QUERY. $\forall a \in A$. $\mathsf{retval}(a) = F_{\mathsf{type}(a)}(\mathsf{ctxt}(a))$



Counter data type

F: context of a \rightarrow return value of a



F: apply standard counter operations in any order Commutative -> ar is not used





Remove wins:s.add(x); s.remove(x) \rightarrow s = {}Add wins:s.remove(x); s.add(x) \rightarrow s = {x}



Remove wins: s.add(x); s.remove(x) \rightarrow s = {}

Add wins: $s.remove(x); s.add(x) \rightarrow s = \{x\}$

- Time-stamping is not always the best solution
- Application semantics might prefer a given outcome: add wins for a shopping cart

Add-wins set

F: context of a \rightarrow return value of a



F: remove cancels out only vis-preceding adds If you saw it, it's not a conflict

Add-wins set

F: context of a \rightarrow return value of a



F: remove cancels out only vis-preceding adds If you saw it, it's not a conflict

System specification = set

- I. How do we resolve conflicts resul concurrent updates at different re
 Data type specification
- 2. Does a user's knowledge change ir the system evolution?

Consistency axioms



Figure 1. Axioms of eventual consistency WELL-FORMEDNESS AXIOMS SOWF: so is the union of transitive, irreflexive and total orders on actions by each session VISWF: $\forall a, b. a \xrightarrow{\text{vis}} b \implies \text{obj}(a) = \text{obj}(b)$ ARWF: $\forall a, b, a \xrightarrow{ar} b \implies \mathsf{obj}(a) = \mathsf{obj}(b),$ ar is transitive and irreflexive, and $\operatorname{ar}_{\operatorname{vis}^{-1}(a)}$ is a total order for all $a \in A$ AUXILIARY RELATIONS Per-object session order: $soo = (so \cap sameobj)$ Per-object causality order: $hbo = (soo \cup vis)^+$ Causality order: $hb = (so \cup vis)^+$ BASIC EVENTUAL CONSISTENCY AXIOMS RVAL: $\forall a \in A$. $\operatorname{rval}(a) = F_{\operatorname{type}(a)}(\operatorname{cone}(a))$ EVENTUAL: $\forall a \in A. \neg (\exists infinitely many \ b \in A. sameobj(a, b) \land \neg (a \xrightarrow{\mathsf{vis}} b))$ THINAIR: so ∪ vis is acyclic SESSION GUARANTEES RYW (Read Your Writes): soo \subseteq vis MR (Monotonic Reads): (vis; soo) \subseteq vis WFRV (Writes Follow Reads in Visibility): (vis; soo^{*}; vis) \subseteq vis WFRA (Writes Follow Reads in Arbitration): (vis; soo^{*}) \subseteq ar MWV (Monotonic Writes in Visibility): $(soo; vis) \subseteq vis$ MWA (Monotonic Writes in Arbitration): $soo \subseteq ar$ CAUSALITY AXIOMS POCV (Per-Object Causal Visibility): hbo \subseteq vis POCA (Per-Object Causal Arbitration): $hbo \subseteq ar$ COCV (Cross-Object Causal Visibility): $(hb \cap sameobj) \subseteq vis$ COCA (Cross-Object Causal Arbitration): hb ∪ ar is acyclic

Axioms

Basic eventual consistency

Session guarantees

Per-object causal consistency

Causal consistency

Strong consistency

Figure 1. Axioms of eventual consistency
WELL-FORMEDNESS AXIOMSWELL-FORMEDNESS AXIOMSSOWF: so is the union of transitive, irreflexive and total orders
on actions by each sessionVISWF: $\forall a, b. a \xrightarrow{vis} b \implies obj(a) = obj(b)$ ARWF: $\forall a, b. a \xrightarrow{ar} b \implies obj(a) = obj(b)$,
ar is transitive and irreflexive, and
 $ar|_{vis^{-1}(a)}$ is a total order for all $a \in A$ AUXILIARY RELATIONSPer-object session order: soo = (so \cap sameobj)Per-object causality order: hbo = (soo \cup vis)⁺Causality order: hbo = (soo \cup vis)⁺

BASIC EVENTUAL CONSISTENCY AXIOMS RVAL: $\forall a \in A$. $\operatorname{rval}(a) = F_{\operatorname{type}(a)}(\operatorname{cone}(a))$ EVENTUAL: $\forall a \in A$. $\neg(\exists \text{ infinitely many } b \in A. \operatorname{sameobj}(a, b) \land \neg(a \xrightarrow{\operatorname{vis}} b))$ THINAIR: so \cup vis is acyclic

SESSION GUARANTEES RYW (Read Your Writes): $soo \subseteq vis$ MR (Monotonic Reads): (vis; $soo) \subseteq vis$ WFRV (Writes Follow Reads in Visibility): (vis; soo^* ; $vis) \subseteq vis$ WFRA (Writes Follow Reads in Arbitration): (vis; soo^*) \subseteq ar MWV (Monotonic Writes in Visibility): (soo; vis) $\subseteq vis$ MWA (Monotonic Writes in Arbitration): $soo \subseteq ar$

CAUSALITY AXIOMS POCV (Per-Object Causal Visibility): hbo \subseteq vis POCA (Per-Object Causal Arbitration): hbo \subseteq ar COCV (Cross-Object Causal Visibility): (hb \cap sameobj) \subseteq vis COCA (Cross-Object Causal Arbitration): hb \cup ar is acyclic

Formalising "eventual"

EVENTUAL. $\forall a \in A. \neg (\exists \text{ infinitely many } b \in A. \text{ sameobj}(a, b) \land \neg (a \xrightarrow{\mathsf{vis}} b))$

An action cannot be invisible to infinitely many actions on the same object

Formalising "eventual"

EVENTUAL. $\forall a \in A. \neg (\exists \text{ infinitely many } b \in A. \text{ sameobj}(a, b) \land \neg (a \xrightarrow{\mathsf{vis}} b))$

An action cannot be invisible to infinitely many actions on the same object



Formalising "eventual"

EVENTUAL. $\forall a \in A. \neg (\exists \text{ infinitely many } b \in A. \text{ sameobj}(a, b) \land \neg (a \xrightarrow{\mathsf{vis}} b))$

An action cannot be invisible to infinitely many actions on the same object



Strengthening consistency

General principle: mandate that certain actions be visible by including additional edges into vis

READ YOUR WRITES: so \cap sameobj \subseteq vis



How do we disallow this?



How do we disallow this?



write(noboss) happened-before the read, but the read didn't see it

How do we disallow this?



Mandate that all actions that happened before have to be visible

Causal consistency



Happens-before relation:

$$\mathsf{hb} = (\mathsf{so} \cup \mathsf{vis})^+$$

CAUSAL CONSISTENCY:

 $(\mathsf{hb} \cap \mathsf{sameobj}) \subseteq \mathsf{vis}$

 $(\mathsf{hb} \cap \mathsf{sameobj}) \subseteq \mathsf{ar}$



Monday, 11 February 13

Causal consistency



Happens-before relation:

$$\mathsf{hb} = (\mathsf{so} \cup \mathsf{vis})^+$$

CAUSAL CONSISTENCY:

 $(\mathsf{hb} \cap \mathsf{sameobj}) \subseteq \mathsf{vis}$

 $(\mathsf{hb} \cap \mathsf{sameobj}) \subseteq \mathsf{ar}$

Technical take-aways

- Ways in which the database processes requests specified declaratively by vis and ar
- Conflict resolution policies specified by functions of vis and ar



 Consistency strengthened by including additional edges into vis or ar

READ YOUR WRITES:

 $\mathsf{so} \cap \mathsf{sameobj} \subseteq \mathsf{vis}$

CAUSAL CONSISTENCY: (hb \cap sameobj) \subseteq vis (hb \cap sameobj) \subseteq ar

What does this buy us?

 Exploiting testing and verification technology developed for weak memory models

Throw the axiomatic model at a SAT solver to do bounded model checking

- Letting the programmer switch between different types of eventual consistency within the same system implementation
 - E.g., using fences, like in weak memory
- Compositional reasoning about eventually consistent systems

The End

Draft paper available at

http://www.cs.ox.ac.uk/people/ hongseok.yang/Public/Publications.html