Design and Validation of Cloud Storage Systems using Maude

Peter Csaba Ölveczky

University of Oslo University of Illinois at Urbana-Champaign

Based on joint work with Jon Grov and members of UIUC's Center for Assured Cloud Computing

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 1 / 58



Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 2 / 58



Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

◆□▶ ◆□▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ◆ □ ▶ ○ へ ○ UCM, February 20, 2017 3 / 58

Global 500 2017

The most valuable brands of 2017



Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 4 / 58

Cloud Computing Data Stores

Cloud computing systems store/retrieve large amounts of data



Are Planning To

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

5 / 58

(同) (三) (三)

Availability



- Data should always be available
 - network/site failures, network congestion, scheduled upgrades
 - ightarrow data must be replicated

A B > A B >

Availability



- Data should always be available
 - network/site failures, network congestion, scheduled upgrades
 - → data must be replicated
- Large and growing data
 - Facebook (2014): 300 petabytes data; 350M photos uploaded every day
 - \longrightarrow data must be partitioned

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 6 / 58

Consistency in Replicated Systems



Figure by Jiaqing Du

Consistency: All replicas of a data item should have same value

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

"CAP Theorem"

Data consistency + partition tolerance + availability impossible



(Figure from http://flux7.com/blogs/nosql/cap-theorem-why-does-it-matter/)

8 / 58

Peter Csaba Ölveczky (U. Oslo/UIUC) Cloud Storage Systems in Maude UCM, February 20, 2017

Slightly Different View

Trade-off

$\mathsf{consistency} \; \mathsf{level} \;\; \longleftrightarrow \;\; \mathsf{latency}$

Peter Csaba Ölveczky (U. Oslo/UIUC) Cloud Storage Systems in Maude UCM, February 20, 2017 9 / 58

(日) (同) (三) (三)

- B

Eventual Consistency

Weak consistency OK for some applications



Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

10 / 58

Eventual Consistency

- Weak consistency OK for some applications
- ... but not others:







Electronic medical records

Electronic patient registry Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

10 / 58

Designing Data Stores

- Complex systems
 - size
 - replication
 - concurrence
 - fault tolerance

- B

(日) (同) (三) (三)

Designing Data Stores

- Complex systems
 - size
 - replication
 - concurrence
 - fault tolerance

• Many hours of "whiteboard analysis"



Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 11 / 58

Validating Data Store Designs

• Correctness: "hand proofs"

- error prone
- informal
- key assumptions implicit
- does not scale to nontrivial systems

∃ → < ∃ →</p>

Validating Data Store Designs

- Correctness: "hand proofs"
 - error prone
 - informal
 - key assumptions implicit
 - does not scale to nontrivial systems
- Performance: simulation tools, real implementations
 - additional artifact
 - cannot be used to reason about correctness

Use formal methods to develop and validate designs

- define mathematical model of system
- use mathematical rules to analyze system

A B + A B +

Use formal methods to develop and validate designs

- define mathematical model of system
- use mathematical rules to analyze system

Find errors early!

Using Formal Methods (I): Validation Perspective

• Formal system model S

- precise mathematical model
- makes assumptions precise and explicit
- amenable to mathematical analysis

Using Formal Methods (I): Validation Perspective

• Formal system model S

- precise mathematical model
- makes assumptions precise and explicit
- amenable to mathematical analysis
- Formal property specification P
 - precise description of consistency model
 - can check whether $S \models P$

化原因 化原因

Using Formal Methods (I): Validation Perspective

• Formal system model S

- precise mathematical model
- makes assumptions precise and explicit
- amenable to mathematical analysis
- Formal property specification P
 - precise description of consistency model
 - can check whether $S \models P$
- What about performance analysis?

Need:

• expressive and intuitive modeling language

Peter Csaba Ölveczky (U. Oslo/UIUC) Cloud Storage S

Cloud Storage Systems in Maude

► < ☐ ► < 2 ► < 2 ► 2</p>
UCM, February 20, 2017

15 / 58

Need:

- expressive and intuitive modeling language
- expressive and intuitive property specification language

・ 同 ト ・ ヨ ト ・ ヨ ト

Need:

- expressive and intuitive modeling language
- expressive and intuitive property specification language
- automatically check whether design satisfies property
 - quick and extensive feedback
 - saves days of whiteboard analysis
 - "extensive and automatic test suite"

- 4 B b 4 B b

Need:

- expressive and intuitive modeling language
- expressive and intuitive property specification language
- automatically check whether design satisfies property
 - quick and extensive feedback
 - saves days of whiteboard analysis
 - "extensive and automatic test suite"
- design model also for performance analysis!
 - no new artifact for performance analysis

Difficult challenges:

- intuitive
- expressive
- useful automatic analyses
- both correctness and performance analysis
- complex properties to check
- mature tool support
- real-time and probabilistic features

3 N

Our Framework: Rewriting Logic

• Rewriting logic: equations and rewrite rules

- expressive
- simple/intuitive
- object-oriented

Our Framework: Rewriting Logic

• Rewriting logic: equations and rewrite rules

- expressive
- simple/intuitive
- object-oriented
- Maude tool:
 - simulation
 - temporal logic model checking
 - ★ expressive property specification language

Our Framework: Rewriting Logic

• Rewriting logic: equations and rewrite rules

- expressive
- simple/intuitive
- object-oriented
- Maude tool:
 - simulation
 - temporal logic model checking
 - ★ expressive property specification language
- Extensions:
 - real-time systems
 - probabilistic systems

Maude: Software Engineering Perspective I

• Models can be developed quickly

Maude: Software Engineering Perspective I

- Models can be developed quickly
- Simulation gives quick feedback (rapid prototyping)

Maude: Software Engineering Perspective I

- Models can be developed quickly
- Simulation gives quick feedback (rapid prototyping)
- Model checking: analyze all behaviors from one initial state



http://embsys.technikum-wien.at/projects/decs/verification/formalmethods.php

formal test-driven development: "test-driven development approach where many complex scenarios can be quickly tested by model checking"

Peter Csaba Ölveczky (U. Oslo/UIUC) Cle

Cloud Storage Systems in Maude

UCM, February 20, 2017 18 / 58

Maude: Software Engineering Perspective (cont.)

What about performance analysis?

Peter Csaba Ölveczky (U. Oslo/UIUC)

Maude: Software Engineering Perspective (cont.)

What about performance analysis?

(Randomized) simulations

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

19 / 58

What about performance analysis?

- (Randomized) simulations
- Probabilistic analysis (using PVeStA)
 - statistical model checking

Same artifact for:

- precise system description
- rapid prototyping
- extensive testing
- correctness analysis
- performance estimation
Case Study I

Modeling, Analyzing, and Extending Megastore



Joint work with Jon Grov (U. Oslo)

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 21 / 58

Megastore:

- Google's wide-area replicated data store
- 3 billion write and 20 billion read transactions daily (2011)





Android Market

Megastore: Key Ideas (I)



(Figure from http://cse708.blogspot.jp/2011/03/megastore-providing-scalable-highly.html)

Data divided into entity groups

- Peter's email
- Books on rewriting logic
- Narciso's documents

・ 同 ト ・ ヨ ト ・ ヨ ト

- Consistency for transactions accessing a single entity group
 - no guarantee if transaction reads multiple entity groups

• [Developed and] formalized [our version of the] Megastore [approach] in Maude

3

(a)

- [Developed and] formalized [our version of the] Megastore [approach] in Maude
 - first (public) formalization/detailed description of Megastore

(日) (同) (三) (三)

- [Developed and] formalized [our version of the] Megastore [approach] in Maude
 - first (public) formalization/detailed description of Megastore
- 56 rewrite rules (37 for fault tolerance features)

・ 同 ト ・ ヨ ト ・ ヨ ト …

- Key performance measures:
 - average transaction latency
 - number of committed/aborted transactions

< ≣ > <

- Key performance measures:
 - average transaction latency
 - number of committed/aborted transactions
- Randomly generated transactions (rate 2.5 TPS)
- Network delays:

	30%	30%	30%	10%
$Madrid\leftrightarrowParis$	10	15	20	50
$Madrid \leftrightarrow New \; York$	30	35	40	100
$Paris \leftrightarrow New \; York$	30	35	40	100

26 / 58

- Key performance measures:
 - average transaction latency
 - number of committed/aborted transactions
- Randomly generated transactions (rate 2.5 TPS)
- Network delays:

	30%	30%	30%	10%
$Madrid \leftrightarrow Paris$	10	15	20	50
$Madrid \leftrightarrow New \; York$	30	35	40	100
$Paris \leftrightarrow New \; York$	30	35	40	100

• Simulating for 200 seconds:

	Avg. latency (ms)	Commits	Aborts
Madrid	218	109	38
New York	336	129	16
Paris	331	116	21

Megastore-CGC: extending Megastore

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

글 🖌 🔺 글 🕨

27 / 58

• Some transactions must access multiple entity groups

3

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

- Some transactions must access multiple entity groups
- Our work: extend Megastore with consistency for transactions accessing multiple entity groups

・ 同 ト ・ ヨ ト ・ ヨ ト

- Some transactions must access multiple entity groups
- Our work: extend Megastore with consistency for transactions accessing multiple entity groups
- Megastore-CGC piggybacks ordering and validation onto Megastore's coordination protocol
 - no additional messages for validation/commit!
 - maintains Megastore's performance and fault tolerance

Performance Comparison using Real-Time Maude

- Simulating for 1000 seconds (no failures)
- Megastore:

	Commits	Aborts	Avg. latency (ms)
Madrid	652	152	126
Paris	704	100	118
New York	640	172	151

Megastore-CGC:

	Commits	Aborts	Val. aborts	Avg.latency (ms)
Madrid	660	144	0	123
Paris	674	115	15	118
New York	631	171	10	150

Model checking scenarios

- 5 transactions , no failures, message delay 30 ms or 80 ms $\rightarrow 108,279$ reachable states, 124 seconds
- 3 transactions, one site failure and fixed message delay
 → 1,874,946 reachable states, 6,311 seconds
- 3 transactions, fixed message delay and one message failure → 265,410 reachable states, 858 seconds

Case Study II



Work by Si Liu, Muntasir Raihan Rahman, Stephen Skeirik, Indranil Gupta, José Meseguer, Son Nguyen, Jatin Ganhotra (ICFEM'14, QEST'15)

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

- Key-value data store originally developed at Facebook
- Used by Amadeus, Apple, CERN, IBM, Netflix, Facebook/Instagram, Twitter, . . .
- Open source

Cassandra Overview

Read consistency either one, quorum, or all



3

< 6 >

Cassandra Overview

Read consistency either one, quorum, or all



Write consistency either zero, one, quorum, or all



[Figures from http://www.slideshare.net/nuboat/cassandra-distributed-data-store]

Formal model from 345K LOC

- allows experimenting with different optimizations/variations
- Analyze basic property: eventual consistency
- When/how often does Cassandra give stronger guarantees?
 - strong consistency
 - read-your-writes
- Performance evaluation:
 - compare PVeStA analyses with real implementations

• • = • • = •

34 / 58

Formal Analysis with Multiple Clients

	Consistency Lv. Latency	ONE	QUORUM	ALL
Strong	L1 (L1 <d1)< td=""><td>×</td><td>×</td><td>×</td></d1)<>	×	×	×
	L2 (D1 <l2<d2)< td=""><td>×</td><td>×</td><td>×</td></l2<d2)<>	×	×	×
Relation between	L3 (D2 <l3)< th=""><th>\checkmark</th><th>\checkmark</th><th>\checkmark</th></l3)<>	\checkmark	\checkmark	\checkmark
issuing latency and				
message delays	Consistency Lv. Latency	ONE	QUORUM	ALL
Eventual	L1 (L1 <d1)< td=""><td>\checkmark</td><td>\checkmark</td><td>\checkmark</td></d1)<>	\checkmark	\checkmark	\checkmark
	L2 (D1 <l2<d2)< td=""><td>\checkmark</td><td>\checkmark</td><td>\checkmark</td></l2<d2)<>	\checkmark	\checkmark	\checkmark
	L3 (D2 <l3)< td=""><td>\checkmark</td><td>\checkmark</td><td>\checkmark</td></l3)<>	\checkmark	\checkmark	\checkmark

Conclusion

- strong consistency depends on the latency between requests
- eventual consistency is guaranteed

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 35 / 58



- (X axis =) Issuing Latency = time difference between the given read request and the latest write request
- (Y axis =) Probability of a request satisfying that model

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

P-Store [N. Schiper, P. Sutra, and F. Pedone; IEEE SRDS'10]

- Replicated and partitioned data store
- Serializability

э

▲ □ ▶ ▲ □ ▶ ▲ □ ▶

P-Store [N. Schiper, P. Sutra, and F. Pedone; IEEE SRDS'10]

- Replicated and partitioned data store
- Serializability

Atomic multicast orders concurrent transactions

∃ >

P-Store [N. Schiper, P. Sutra, and F. Pedone; IEEE SRDS'10]

- Replicated and partitioned data store
- Serializability
- Atomic multicast orders concurrent transactions
- Group commitment for atomic commit

Atomic Multicast

Definition

Atomic Multicast: Consistent reception order of messages

- (a): any pair of nodes receive the same atomic-multicast messages in the same order
- (b): induced "global read order" must be acyclic

Atomic Multicast

Definition

Atomic Multicast: Consistent reception order of messages

- (a): any pair of nodes receive the same atomic-multicast messages in the same order
- (b): induced "global read order" must be acyclic

Example

- A reads *m*₁ < *m*₂
- *B* reads *m*₂ < *m*₃
- *C* reads *m*₃ < *m*₁

satisfies (a) but not (b)

< ロ > < 同 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < 回 > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ >

- Fundamental problem in distributed systems
- Impose order on conflicting concurrent transactions

Atomic Multicast in Maude (I)

- Fundamental problem in distributed systems
- Impose order on conflicting concurrent transactions
- Many algorithms for atomic multicast

Atomic Multicast in Maude (I)

- Fundamental problem in distributed systems
- Impose order on conflicting concurrent transactions
- Many algorithms for atomic multicast
- Define generic atomic multicast primitive in Maude
 - abstract
 - covers all possible receiving orders

化原因 化原因

Atomic Multicast in Maude (I)

- Fundamental problem in distributed systems
- Impose order on conflicting concurrent transactions
- Many algorithms for atomic multicast
- Define generic atomic multicast primitive in Maude
 - abstract
 - covers all possible receiving orders
- Infrastructure stores (un)read AM messages

My Work: Atomic Multicast in Maude (II)

• Atomic-multicast message M:

```
rl [atomic-multicast] :
    < 0 : Node | msgToSend : M, receivers : OS >
    =>
        < 0 : Node | ... >
        (atomic-multicast M from 0 to OS) .
```

My Work: Atomic Multicast in Maude (II)

• Atomic-multicast message M:

```
rl [atomic-multicast] :
    < 0 : Node | msgToSend : M, receivers : OS >
=>
    < 0 : Node | ... >
    (atomic-multicast M from 0 to OS) .
```

Read:

```
crl [receiveAtomicMulticast] :
    (msg M from 02 to 0)
    < 0 : Node | ... >
    AM-TABLE
=>
    < 0 : Node | ... >
    updateAM(MC, 0, AM-TABLE)
if okToRead(MC, 0, AM-TABLE) .
```

Peter Csaba Ölveczky (U. Oslo/UIUC)

Analyzing P-Store

Find all reachable final states from init3:

```
Maude> (search init3 =>! C:Configuration .)
```

```
Solution 1
C:Configuration --> ...
< c1 : Client | pendingTrans : t1, txns : emptyTransList >
< c2 : Client | pendingTrans : t2, txns : emptyTransList >
< r1 : PStoreReplica | aborted : none,
                       committed : < t1 : Transaction | ... >
< r2 : PStoreReplica | aborted : none,
                       committed : < t2 : Transaction | ... >
```

. . .

41 / 58

Analyzing P-Store

Find all reachable final states from init3:

```
Maude> (search init3 =>! C:Configuration .)
```

sites validate transactions

. . .

but client never gets result

◆□ ▶ ◆□ ▶ ◆ □ ▶ ◆ □ ● ● ● ● ● ●
Solution 5

Solution 5



Host does not validate t1 even when needed info known

- Found the source of the errors
 - all replicas must be involved in voting and notification
 - ★ not just write replicas
- Modeled and analyzed proposed corrected version

A B > A B >

P-Store Summary

Algorithm A_{ae} A Genuine Certification Protocol - Code of site s 1. Initialization 2: $Votes \leftarrow \emptyset$ 3: function ApplyUpdates(T) foreach $\forall (k, v) \in T.up : k \in Items(s)$ do 4: 5: let ts be Version(k, s) $w_T[k, v, ts + 1]$ {write to the database} 6. 7: function Certify(T)**return** $\forall (k, ts) \in T.rs \text{ s.t. } k \in Items(s) : ts = Version(k, s)$ 9: To submit transaction T {Task 1} 10: A-MCast(T) to Replicas(T) ${Executing \rightarrow Submitted}$ {Task 2} 11: When receive(VOTE, T.id, vote) from s' 12: $Votes \leftarrow Votes \cup (T.id, s', vote)$ 13: When A-Deliver(T) {Task 3} 14: if T is local then 15: if Certify(T) then 16: ApplyUpdates(T)17: commit T{Submitted \rightarrow Committed} 18: else abort T {Submitted \rightarrow Aborted} 19: else 20: if $\exists (k, -) \in T.rs : k \in Items(s)$ then 21: $Votes \leftarrow Votes \cup (T.id, s, Certify(T))$ 22: send(VOTE, T.id, Certify(T)) to all s' in WReplicas(T) s.t. $s' \notin group(s)$ 23: if $s \in WReplicas(T)$ then wait until $\exists VQ \in VQ(T)$: 24: $\forall s' \in VQ : (T.id, s', -) \in Votes$ 25: if $\forall s' \in VQ : (T.id, s', yes) \in Votes$ then ApplyUpdates(T)26: 27: {Submitted \rightarrow Committed} commit T28: {Submitted \rightarrow Aborted} else abort T29: if $s \in WReplicas(T)$ then send T's outcome to Proxy(T)

- "P-Store verified"
- 3 significant errors found
- one confusing definition
- key assumption missing

(日) (同) (三) (三)

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017 44 / 58

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

UCM, February 20, 2017

< ロ > < 同 > < 回 > < 回 >

45 / 58

э

Developed formal models of large industrial data stores

- Google's Megastore (from brief description)
- Apache Cassandra (from 345K LOC and description)
- P-Store (academic)

Developed formal models of large industrial data stores

- Google's Megastore (from brief description)
- Apache Cassandra (from 345K LOC and description)
- P-Store (academic)
- Automatic model checking analysis of consistency properties

A B + A B +

Developed formal models of large industrial data stores

- Google's Megastore (from brief description)
- Apache Cassandra (from 345K LOC and description)
- P-Store (academic)
- Automatic model checking analysis of consistency properties
- Designed own transactional data stores
 - Megastore-CGC
 - variation of Cassandra

A = A = A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A
A

Developed formal models of large industrial data stores

- Google's Megastore (from brief description)
- Apache Cassandra (from 345K LOC and description)
- P-Store (academic)
- Automatic model checking analysis of consistency properties
- Designed own transactional data stores
 - Megastore-CGC
 - variation of Cassandra
- Errors, ambiguities, missing assumptions found in "verified" P-Store

Developed formal models of large industrial data stores

- Google's Megastore (from brief description)
- Apache Cassandra (from 345K LOC and description)
- P-Store (academic)
- Automatic model checking analysis of consistency properties
- Designed own transactional data stores
 - Megastore-CGC
 - variation of Cassandra
- Errors, ambiguities, missing assumptions found in "verified" P-Store
- Maude/PVeStA performance estimation close to real implementations

- Quickly develop formal models/prototypes of complex systems
 - experiment with different design choices

→ □ → → □ →

47 / 58

- Quickly develop formal models/prototypes of complex systems
 - experiment with different design choices
- Simulation and model checking throughout design phase
 - model-checking-based-testing for subtle "corner cases"
 - replaces days of whiteboard analysis
 - too many scenarios for standard test-based development
 - catch bugs early!

- Quickly develop formal models/prototypes of complex systems
 - experiment with different design choices
- Simulation and model checking throughout design phase
 - model-checking-based-testing for subtle "corner cases"
 - replaces days of whiteboard analysis
 - too many scenarios for standard test-based development
 - catch bugs early!

Single artifact for

- system description
- rapid prototyping
- model checking
- performance estimation

- Quickly develop formal models/prototypes of complex systems
 - experiment with different design choices
- Simulation and model checking throughout design phase
 - model-checking-based-testing for subtle "corner cases"
 - replaces days of whiteboard analysis
 - too many scenarios for standard test-based development
 - catch bugs early!
- Single artifact for
 - system description
 - rapid prototyping
 - model checking
 - performance estimation
- Megastore and Megastore-CGC modeler had no formal methods experience

47 / 58

DOI:10.1145/2699417

Engineers use TLA+ to prevent serious but subtle bugs from reaching production.

BY CHRIS NEWCOMBE, TIM RATH, FAN ZHANG, BOGDAN MUNTEANU, MARC BROOKER, AND MICHAEL DEARDEUFF

How Amazon Web Services Uses Formal Methods

Peter Csaba Ölveczky (U. Oslo/UIUC)

Cloud Storage Systems in Maude

• Amazon Web Services (AWS):

- world's largest cloud computing service provider
- more profitable than Amazon's retail business

- Amazon Web Services (AWS):
 - world's largest cloud computing service provider
 - more profitable than Amazon's retail business
- Amazon Simple Storage Service (S3)
 - stores > 3 trillion objects
 - 99.99% availability of objects
 - > 1 million requests per second
- DynamoDB data store

Amazon Web Services and Formal Methods

- Formal methods used extensively at AWS during design of S3, DynamoDB, ...
- Used Lamports TLA+
 - model checking

Model checking finds "corner case" bugs that would be hard to find with standard industrial methods:

イロト イポト イヨト イヨト 三日

Model checking finds "corner case" bugs that would be hard to find with standard industrial methods:

• "We have found that standard verification techniques in industry are necessary but not sufficient. We routinely use deep design reviews, static code analysis, stress testing, and fault-injection testing but still find that subtle bugs can hide in complex fault-tolerant systems."

・ 同 ト ・ ヨ ト ・ ヨ ト …

Model checking finds "corner case" bugs that would be hard to find with standard industrial methods:

 "the model checker found a bug that could lead to losing data [...]. This was a very subtle bug; the shortest error trace exhibiting the bug included 35 high-level steps. [...] The bug had passed unnoticed through extensive design reviews, code reviews, and testing."

・ロト ・ 一日 ・ ・ 日 ・ ・ 日 ・ ・ 日

Experiences at Amazon WS II

A formal specification is a valuable precise description of an algorithm:

Experiences at Amazon WS II

A formal specification is a valuable precise description of an algorithm:

 "the author is forced to think more clearly, helping eliminating "hand waving," and tools can be applied to check for errors in the design, even while it is being written. In contrast, conventional design documents consist of prose, static diagrams, and perhaps psuedo-code in an ad hoc untestable language."

Experiences at Amazon WS II

A formal specification is a valuable precise description of an algorithm:

 "Talk and design documents can be ambiguous or incomplete, and the executable code is much too large to absorb quickly and might not precisely reflect the intended design. In contrast, a formal specification is precise, short, and can be explored and experimented on with tools."

Experiences at Amazon WS III

Formal methods are surprisingly feasible for mainstream software development and give good return on investment:

Formal methods are surprisingly feasible for mainstream software development and give good return on investment:

"In industry, formal methods have a reputation for requiring a huge amount of training and effort to verify a tiny piece of relatively straightforward code. Our experience with TLA+ shows this perception to be wrong. [...] Amazon engineers have used TLA+ on 10 large complex real-world systems. In each, TLA+ has added significant value. [...] Engineers have been able to learn TLA+ from scratch and get useful results in two to three weeks."

Experiences at Amazon WS III

Formal methods are surprisingly feasible for mainstream software development and give good return on investment:

 "Using TLA+ in place of traditional proof writing would thus likely have improved time to market, in addition to achieving greater confidence in the system's correctness."

く 戸 と く ヨ と く ヨ と …

Quick and easy to experiment with different design choices:

Peter Csaba Ölveczky (U. Oslo/UIUC) **Cloud Storage Systems in Maude**

Quick and easy to experiment with different design choices:

• "We have been able to make innovative performance optimizations [...] we would not have dared to do without having model-checked those changes. A precise, testable description of a system becomes a what-if tool for designs."

• • • • • • • • •

Experiences at Amazon WS: Limitations

TLA+ did/could not analyze performance degradation

3

() < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < () < ()

Maude should be better suited!

- more intuitive and expressive specification language
 - ► 00
 - hierarchical states
 - dynamic object/message creation/deletion
 - ▶ ...
- Support for real-time and probabilistic systems
- Also for performance estimation!

→ Ξ → < Ξ →</p>

56 / 58

Conclusions at Amazon

» key insights

- Formal methods find bugs in system designs that cannot be found through any other technique we know of.
- Formal methods are surprisingly feasible for mainstream software development and give good return on investment.
- At Amazon, formal methods are routinely applied to the design of complex real-world software, including public cloud services.

57 / 58

Peter Csaba Ölveczky (U. Oslo/UIUC) Cloud Storage Systems in Maude UCM, February 20, 2017

Take Away from Talk

• Formal methods can be an efficient way to

- design
- test
- describe
- validate correctness and performance
- experiment with different design choices

industrial state-of-the-art fault-tolerant distributed systems also for non-experts

・ 同 ト ・ ヨ ト ・ ヨ ト

Take Away from Talk

• Formal methods can be an efficient way to

- design
- test
- describe
- validate correctness and performance
- experiment with different design choices

industrial state-of-the-art fault-tolerant distributed systems also for non-experts

• Maude suitable modeling language and analysis toolset

・ 同 ト ・ ヨ ト ・ ヨ ト