

MODELING, ANALYZING, AND EXTENDING MEGASTORE USING REAL-TIME MAUDE

Jon Grov¹ and Peter Ölveczky^{1,2}

¹University of Oslo

²University of Illinois at Urbana-Champaign

Thanks to Indranil Gupta (UIUC) and UIUC's Center for Assured
Cloud Computing



GOAL

Cloud computing also when *data consistency* critical

- *e-commerce, banking, medical systems*

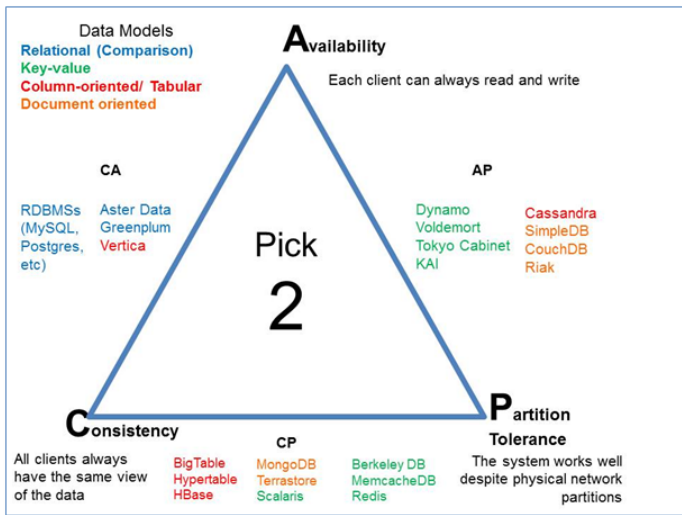
Part I: Formalizing Megastore

DATA IN THE CLOUD

- **Availability** and **scalability**: data must be **replicated**

DATA IN THE CLOUD

- **Availability** and **scalability**: data must be **replicated**



EVENTUAL CONSISTENCY

“Eventual consistency” OK for some applications

EVENTUAL CONSISTENCY

“Eventual consistency” OK for some applications

The image shows a screenshot of a Facebook profile page for Csaba Ölveczky. The page layout includes a top navigation bar with a search bar, the user's name 'Csaba', and links for 'Home' and 'Find Friends'. On the left side, there is a navigation menu with sections for 'News Feed', 'Messages', 'Events', 'Find Friends', 'GROUPS', and 'APPS'. The main content area shows a post from Denise West to Csaba Ölveczky, dated '1 hr · 1 like'. The post text reads: 'Happy Birthday!!! Hope your day is wonderful, and your personal new year fantastic!! BIG HUG!'. Below the text is a large image of lit birthday candles. On the right side, there are sponsored advertisements for 'SF Film Norge' and 'G-Sport', each with a 'Like Page' button. At the bottom of the page, there are navigation icons for back, forward, and search.

EVENTUAL CONSISTENCY

“Eventual consistency” OK for some applications

U.S. INTERNATIONAL 中文網

Introducing **nytnow**
From The New York Times

The New York Times

Thursday, April 10, 2014 | Today's Paper | Personalize Your Weather | f | t

STORIES AT THE SPEED OF LIFE.

WORLD U.S. NEW YORK BUSINESS OPINION SPORTS SCIENCE ARTS FASHION & STYLE VIDEO All Sections

International New York Times

GLOBAL COVERAGE THAT GOES EVERYWHERE YOU DO. 99¢ FOR 12 WEEKS

GET IT NOW >

Secretary Resigns After HealthCare.gov Missteps

By MICHAEL D. SHEAR 6:31 PM ET

Ending a stormy five-year tenure marred by the disastrous rollout of President Obama's signature legislative achievement, Kathleen Sebelius is resigning as secretary of health and human services.

90 Comments



Todd Heislert/The New York Times

Near the Anzalduas Dam, two guides led a raft full of migrants from Honduras and Guatemala across the Rio Grande.

The Opinion Pages

TAKING NOTE

Will 'House of Cards' Deal Elsewhere?

By FRANCIS X. CLINES

The show's producers want Maryland to increase filming subsidies.



- Editorial: **The Truth About the Pay Gap**
- Private Lives: **Parenting the Non-Girlie Girl**

Abusing Both Medicare and Politics

By THE EDITORIAL BOARD

Does anyone doubt that there is a quid pro quo when doctors earning millions of Medicare dollars write big checks to senators and political parties?

OP-ED COLUMNISTS

- **Blow: We Should Be in a Rage**
- **Kristof: Where the G.O.P. Gets It Right**

Today's Times Insider

Michael Sokolove spent several long days with Oscar Pistorius in 2012, including time at a firing range. for a



EVENTUAL CONSISTENCY

“Eventual consistency” OK for some applications

... but not others:

- banking
- stock exchange
- electronic commerce
 - online auctions (eBay)
 - plane tickets ...
- medical (information) systems

TRANSACTIONS

Databases traditionally provide **ACID transactions**

- atomicity
- consistency
- isolation (“serializability”)
- durability

TRANSACTIONS

Databases traditionally provide **ACID transactions**

- atomicity
- consistency
- isolation (“serializability”)
- durability

Little support for ACID transactions in cloud computing

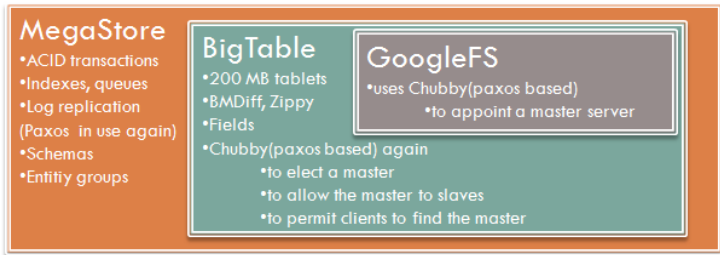
MEGASTORE

Megastore:

- Google's wide-area **replicated data store**
- Key part of Google's cloud infrastructure
- **3 billion write** and **20 billion read** transactions daily (2011)
- Adds (limited) **transactions** to wide-area replicated data stores



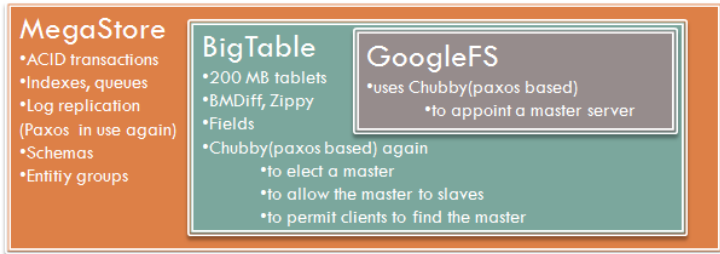
MEGASTORE: KEY IDEAS (I)



(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
 - Jon's documents

MEGASTORE: KEY IDEAS (I)



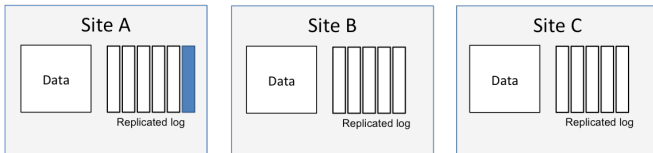
(Figure from <http://cse708.blogspot.jp/2011/03/megastore-providing-scalable-highly.html>)

- Data divided into **entity groups**
 - Peter's email
 - Books on rewriting logic
 - Jon's documents
- **Replicated transaction log** for each entity group

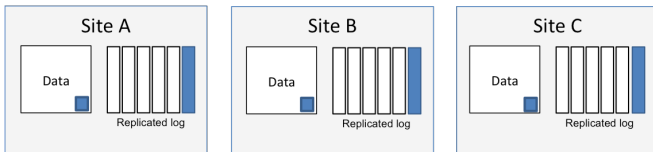
MEGASTORE: KEY IDEAS (II)

- Node suggests next log entry for an entity group
- **Paxos**: agree on next log entry if concurrent updates

Before coordination



After (successful) coordination



MEGASTORE: KEY IDEAS (III)

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

MEGASTORE: KEY IDEAS (III)

- **Consistency** for transactions accessing **a single** entity group
 - no guarantee if transaction reads **multiple** entity groups
- **ElasTraS**, **Spinnaker**, **Calvin**, and Microsoft's **Azure**: consistency **within** each data partition

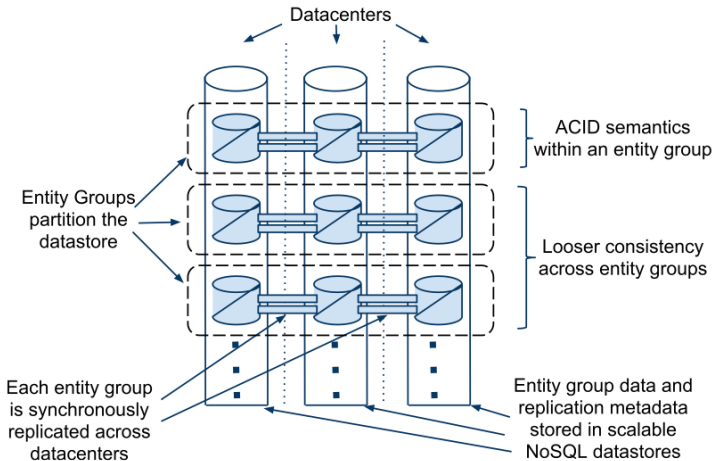


Figure 1: Scalable Replication

(figure from J. Baker et al., "Megastore: Providing Scalable, Highly Available Storage for Interactive Services")

OUR WORK: MOTIVATION

- 12 page informal overview paper of **really complex system**
- Understand system and guarantees
- Basis for further research and extensions

OUR WORK

- [Developed and] formalized [our version of the] Megastore [approach] in Real-Time Maude

OUR WORK

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

OUR WORK

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

OUR WORK

- [Developed and] formalized [our version of the] Megastore [approach] in Real-Time Maude
 - first (public) formalization/detailed description of Megastore
- 56 rewrite rules (37 for fault tolerance features)
- Real-Time Maude simulation and model checking used extensively throughout development of (our) Megastore

BACKGROUND: REAL-TIME MAUDE

Real-Time Maude: rewriting-logic-based language and analysis tool

- **Expressive** formalism
 - algebraic specification specify data types
 - rewrite rules specify local transitions
- **Object-oriented** specification of distributed real-time systems

BACKGROUND: REAL-TIME MAUDE

Real-Time Maude: rewriting-logic-based language and analysis tool

- **Expressive** formalism
 - algebraic specification specify data types
 - rewrite rules specify local transitions
- **Object-oriented** specification of distributed real-time systems
- Formal analysis:
 - simulation
 - reachability analysis
 - **LTL** model checking
 - **Timed CTL** model checking

A REWRITE RULE (IN FAILURE-FREE SETTING)

When all operations in the operations list are completed (reads) or buffered (writes), the transaction is ready to commit. All buffered updates are merged into a candidate log entry. If the transaction updates entities from several entity groups, one log entry is created for each group. For each such entity group, the first step is to send the candidate log entry to the leader for the *next* log position, which was selected during the previous coordination round. The rule for initiating Paxos is modeled as follows:

```
cr1 [initiateCommit] :
  < SID' : Site |
    entityGroups EGROUPS,
    localTransactions : LOCALTRANS
      < TID : Transaction | operations : emptyOpList,
        writes : WRITEOPS, status : idle
        readState : RSTATE, paxosState : PSTATE > >
=>
  < SID : Site |
    localTransactions : LOCALTRANS
      < TID : Transaction | paxosState : NEW-PAXOS-STATE,
        status : in-paxos > >
  ACC-LEADER-REQ-MSG
  if EIDSET := getEntityGroupIds(WRITEOPS) /\
    NEW-PAXOS-STATE := initiatePaxosState(EIDSET, TID, WRITEOPS,
      SID, RSTATE, EGROUPS)
  /\ (createAcceptLeaderMessages(SID, NEW-PAXOS-STATE)) => ACC-LEADER-REQ-MSG
```

A REWRITE RULE (CONT.)

`getEntityGroupIds(WRITEOPS)` contains entity groups accessed by operations in `WRITEOPS`, and `NEW-PAXOS-STATE` contains one record for each entity group. These records contain the log position that `TID` requests to update and the candidate log entry `le`. The operator `createAcceptLeaderMessages` generates an `acceptLeaderReq` message to the leader of each entity group containing the transaction id `TID` and candidate log entry.

ANOTHER REWRITE RULE (NO FAILURES)

The following rule [...] where a replicating site receives an `acceptAllReq` message. The site verifies that it has not already granted an accept for this log position (since messages could be delayed for a long time, it checks both the transaction log and received proposals). If there are no such conflicts, the site responds with an accept message, and stores its accept in proposals for this entity group. The record (TID' LP SID OL) represents the candidate log entry, containing the transaction identifier TID', the log position LP, the proposed leader site SID, and the list of update operations OL.

```
cr1 [rcvAcceptAllReq] :
  (msg acceptAllReq(TID, EG, (TID' LP SID OL), PROPNUM) from SENDER to THIS)
  < THIS : Site |
    entityGroups : EGROUPS
      < EG : EntityGroup | proposals : PROPSET, transactionLog : LEL > >
=>
  < THIS : Site |
    entityGroups : EGROUPS
      < EG : EntityGroup |
        proposals : accepted(SENDER, (TID' LP SID OL), PROPNUM) ;
          removeProposal(LP, PROPSET) > >
  dly(acceptAllResp(TID, EG, LP, PROPNUM) from THIS to SENDER), T)
if not (containsLPos(LP, LEL) or hasAcceptedForPosition(LP, PROPSET))
  /\ T ; TS := possibleMessageDelay(THIS, SENDER) .
```

FORMAL ANALYSIS

- “Monte-Carlo” simulations for performance estimation
- LTL model checking
 - highly nondeterministic setting (**many conflicts**)
 - limited scenarios
 - super useful to discover many bugs not found during simulation

PERFORMANCE ESTIMATION

- Key performance measures:
 - average **transaction latency**
 - number of **committed/aborted transactions**
- 2 entity groups
- Randomly generated transactions (rate 2.5 TPS)

- Network delays:

	30%	30%	30%	10%
London ↔ Paris	10	15	20	50
London ↔ New York	30	35	40	100
Paris ↔ New York	30	35	40	100

PERFORMANCE ESTIMATION (CONT.)

- Simulating for 200 seconds (no failures):

	Avg. latency (ms)	Commits	Aborts
London	122	149	15
New York	155	132	33
Paris	119	148	18

PERFORMANCE ESTIMATION (CONT.)

- Simulating for 200 seconds (no failures):

	Avg. latency (ms)	Commits	Aborts
London	122	149	15
New York	155	132	33
Paris	119	148	18

- Site failures:
 - mean-time-to-failure 10 seconds per site
 - mean-time-to repair 2 seconds

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

MODEL CHECKING (I)

- Desired properties (for **finite** number of transactions):
 - all transactions **finish** execution
 - **all replicas** of an entity must **eventually** have **same value**
 - **all logs** for entity group must **eventually** have **same entries**
 - execution was **serializable**

MODEL CHECKING (I)

- Desired properties (for **finite** number of transactions):
 - all transactions **finish** execution
 - **all replicas** of an entity must **eventually** have **same value**
 - **all logs** for entity group must **eventually** have **same entries**
 - execution was **serializable**
- **Serializability** tricky property
 - construct serialization graph **during** execution
 - check that graph has no cycles

MODEL CHECKING (II)

- All replicas same unless coordinator invalidated:

op entityGroupsEqualOrInvalid : -> Prop [ctor] .

```
ceq {< S1 : Site | coordinator : eglp(EG1, LP) ; EGLP,
      entityGroups :
        < EG1 : EntityGroup | entitiesState : ES1 > EGS1 >
    < S2 : Site | coordinator : eglp(EG1, LP) ; EGLP,
      entityGroups :
        < EG1 : EntityGroup | entitiesState : ES2 > EGS2 >
    REST} |= entityGroupsEqual = false if ES1 /= ES2 .
```

eq {SYSTEM} |= entityGroupsEqualOrInvalid = true [owise] .

MODEL CHECKING (II)

- All replicas same unless coordinator invalidated:

`op entityGroupsEqualOrInvalid` : \rightarrow Prop [ctor] .

`ceq` {< S1 : Site | coordinator : eglp(EG1, LP) ; EGLP,
entityGroups :
 < EG1 : EntityGroup | entitiesState : ES1 > EGS1 >
 < S2 : Site | coordinator : eglp(EG1, LP) ; EGLP,
 entityGroups :
 < EG1 : EntityGroup | entitiesState : ES2 > EGS2 >
 REST} |= `entityGroupsEqual` = false if ES1 \neq ES2 .

`eq` {SYSTEM} |= `entityGroupsEqualOrInvalid` = true [owise] .

- Execution serializable:

`op isSerializable` : \rightarrow Prop [ctor] .

`eq` {< th : TransactionHistory | graph : GRAPH > REST}
 |= `isSerializable` = not hasCycle(GRAPH) .

MODEL CHECKING (II)

- All replicas same unless coordinator invalidated:

```
op entityGroupsEqualOrInvalid : -> Prop [ctor] .
```

```
ceq {< S1 : Site | coordinator : eglp(EG1, LP) ; EGLP,  
      entityGroups :  
        < EG1 : EntityGroup | entitiesState : ES1 > EGS1 >  
    < S2 : Site | coordinator : eglp(EG1, LP) ; EGLP,  
      entityGroups :  
        < EG1 : EntityGroup | entitiesState : ES2 > EGS2 >  
    REST} |= entityGroupsEqual = false if ES1 /= ES2 .
```

```
eq {SYSTEM} |= entityGroupsEqualOrInvalid = true [owise] .
```

- Execution serializable:

```
op isSerializable : -> Prop [ctor] .
```

```
eq {< th : TransactionHistory | graph : GRAPH > REST}  
    |= isSerializable = not hasCycle(GRAPH) .
```

- Desired property:

```
<> [] (allTransFinished /\ entityGroupsEqualOrInvalid  
      /\ transLogsEqualOrInvalid /\ isSerializable)
```

MODEL CHECKING (III)

- Without fault tolerance:
 - model checked **untimed** Maude model
 - covers **all possible** message delays, transaction start times, etc.
 - 3 sites and 3 transactions

MODEL CHECKING (III)

- Without fault tolerance:
 - model checked **untimed** Maude model
 - covers **all possible** message delays, transaction start times, etc.
 - 3 sites and 3 transactions
- With fault tolerance: model checked **timed** model

Msg. delay	#Trans	Trans. start time	#Fail.	Fail. time	Run (s)
{20, 100}	4	{19, 80} and {50, 200}	0	-	1367
{20, 100}	3	{10, 50, 200}	1	60	1164
{20, 40}	3	20, 30, and {10, 50}	2	{40, 80}	872
{20, 40}	4	20, 20, 60, and 110	2	70 and {10, 130}	241
{20, 40}	4	20, 20, 60, and 110	2	{30, 80}	DNF
{10, 30, 80},and {30, 60, 120}	3	20, 30, 40	1	{30, 80}	DNF
{10, 30, 80},and {30, 60, 120}	3	20, 30, 40	1	60	DNF

Part II: Megastore-CGC: extending Megastore

OUR WORK: MOTIVATION

GOAL

*Cloud computing also when **data consistency** critical*

- Sometimes transactions **must** access **multiple** entity groups

OUR WORK: MOTIVATION

GOAL

*Cloud computing also when **data consistency** critical*

- Sometimes transactions **must** access **multiple** entity groups
- **Megastore**: consistency only if transactions access **single** entity group

OUR WORK: MOTIVATION

GOAL

*Cloud computing also when **data consistency** critical*

- Sometimes transactions **must** access **multiple** entity groups
- **Megastore**: consistency only if transactions access **single** entity group
- Our approach: extend **Megastore** with consistency for transactions accessing **multiple** entity groups

OUR WORK: MOTIVATION

GOAL

Cloud computing also when data consistency critical

- Sometimes transactions **must** access **multiple** entity groups
- **Megastore**: consistency only if transactions access **single** entity group
- Our approach: extend **Megastore** with consistency for transactions accessing **multiple** entity groups
 - must maintain **Megastore's**
 - performance
 - strong **fault tolerance**

OUR IDEA

- Key observations:
 1. A Megastore site replicating different entity groups participates in all updates on those entity groups
 2. → site has implicit local ordering on those updates

OUR IDEA

- **Key observations:**
 1. A **Megastore** site replicating different entity groups participates in all updates on those entity groups
 2. → site has **implicit local** ordering on those updates
- **Our idea:**
 1. Select **one** such site and make its order **explicit**
 2. This site validates transactions on multiple entity groups before commit
 - **ordering class** = set of entity groups
 - **ordering site** for each ordering class

MEGASTORE-CGC

Megastore-CGC piggybacks ordering and validation onto Megastore's coordination protocol

- no additional messages for validation/commit!
- maintains Megastore's performance and fault tolerance

MEGASTORE-CGC

Megastore-CGC piggybacks ordering and validation onto Megastore's coordination protocol

- no additional messages for validation/commit!
- maintains Megastore's performance and fault tolerance
- failover protocol when ordering site fails

DEVELOPING MEGASTORE-CGC

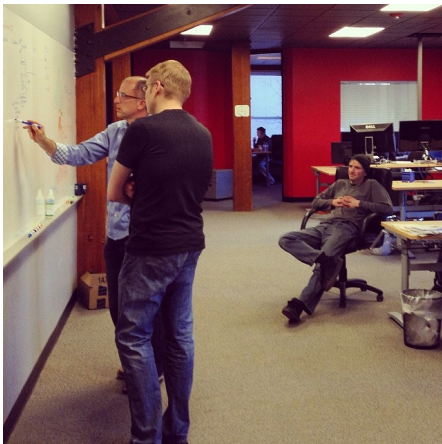
- Developed **Megastore-CGC** using **Real-Time Maude**
- **72** (fairly complex) rewrite rules
- **Real-Time Maude simulation** and **model checking** used extensively **throughout** development of **Megastore-CGC**

SOFTWARE ENGINEERING PERSPECTIVE (I)

- Developing a **fault-tolerant** protocol is **difficult**
- Typically: pseudo-code description → Java prototype

SOFTWARE ENGINEERING PERSPECTIVE (I)

- Developing a **fault-tolerant** protocol is **difficult**
- Typically: pseudo-code description → Java prototype
 - **many** hours of “**whiteboard analysis**”
 - still incorrect



SOFTWARE ENGINEERING PERSPECTIVE (II)

- Our approach: **Formal test-driven development:**
 1. Express requirements as LTL formulas
 2. Develop Real-Time Maude model
 3. Test model using Real-Time Maude simulation and model checking
 4. For failing tests, analyse the problem and go to step (2)

SOFTWARE ENGINEERING PERSPECTIVE (II)

- Our approach: **Formal test-driven development**:
 1. Express requirements as LTL formulas
 2. Develop Real-Time Maude model
 3. Test model using Real-Time Maude simulation and model checking
 4. For failing tests, analyse the problem and go to step (2)
- Allowed **one person** to develop and validate complex protocol in short time

SOFTWARE ENGINEERING PERSPECTIVE (III)

- With **much less effort**, the fault-tolerance features of Megastore-CGC are **significantly more mature** than similar protocols the first author has been working on
- **Model checking** replaces most of the **whiteboard analysis**

PERFORMANCE COMPARISON USING REAL-TIME MAUDE

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

	Commits	Aborts	Avg. latency (ms)
Site 1	652	152	126
Site 2	704	100	118
RSite	640	172	151

- Megastore-CGC:

	Commits	Aborts	Val. aborts	Avg.latency (ms)
Site 1	660	144	0	123
Site 2	674	115	15	118
RSite	631	171	10	150

MODEL CHECKING (III)

Model checking scenarios

- 5 transactions (3 fixed, 2 with 2 start times), **no failures**, message delay 30 ms or 80 ms
→ **108,279** reachable states, 124 seconds
- 3 transactions (all with 2 start times), **one site failure** and fixed message delay
→ **1,874,946** reachable states, 6,311 seconds
- 3 transactions (all with 2 start times), fixed message delay and **one message failure**
→ **265,410** reachable states, 858 seconds

CONCLUSION

- Developed model of Megastore from limited information
 - **first** formal analysis of transactional data stores

CONCLUSION

- Developed model of Megastore from limited information
 - **first** formal analysis of transactional data stores
- Extended the Megastore approach to **Megastore-CGC**
 - consistency for transactions accessing multiple entity groups
 - **performance** and **fault tolerance** on par with Megastore

CONCLUSION

- Developed model of Megastore from limited information
 - **first** formal analysis of transactional data stores
- Extended the Megastore approach to **Megastore-CGC**
 - consistency for transactions accessing multiple entity groups
 - **performance** and **fault tolerance** on par with Megastore
- Detailed formal specification in **Real-Time Maude**
 - shorter development time
 - increased quality of specification
 - analyze **performance** and **correctness**

CONCLUSION

- Developed model of Megastore from limited information
 - **first** formal analysis of transactional data stores
- Extended the Megastore approach to **Megastore-CGC**
 - consistency for transactions accessing multiple entity groups
 - **performance** and **fault tolerance** on par with Megastore
- Detailed formal specification in **Real-Time Maude**
 - shorter development time
 - increased quality of specification
 - analyze **performance** and **correctness**
- First generic (?) cloud data management protocol with consistency for transactions across partitions
 - Google's **Spanner**: similar features with advanced infrastructure

CONCLUSION

- Developed model of Megastore from limited information
 - **first** formal analysis of transactional data stores
- Extended the Megastore approach to **Megastore-CGC**
 - consistency for transactions accessing multiple entity groups
 - **performance** and **fault tolerance** on par with Megastore
- Detailed formal specification in **Real-Time Maude**
 - shorter development time
 - increased quality of specification
 - analyze **performance** and **correctness**
- First generic (?) cloud data management protocol with consistency for transactions across partitions
 - Google's **Spanner**: similar features with advanced infrastructure
- Real-Time Maude **expressiveness** key
 - model complex system
 - express complex requirements