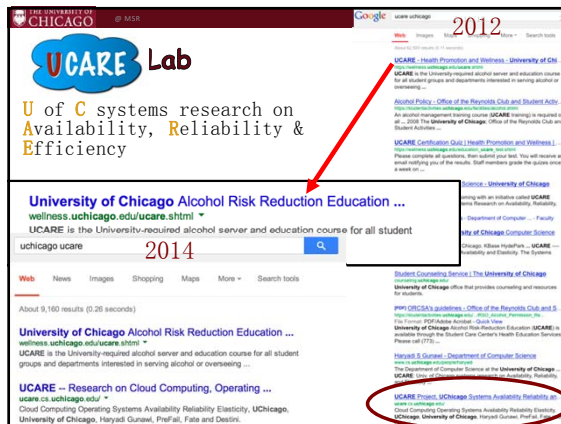


What New Bugs Live in the Cloud? (and how to exterminate them)

Haryadi Gunawi

UCARE Lab
U of C systems research on Availability, Reliability & Efficiency

University of Chicago Alcohol Risk Reduction Education ...
wellness.uchicago.edu/ucare.shtml *
UCARE is the University-required alcohol server and education course for all student groups and departments interested in serving alcohol or overseeing ...

UCARE - Research on Cloud Computing, Operating ...
ucare.cs.uchicago.edu/*
Cloud Computing Operating Systems Availability Reliability Elasticity, UChicago, University of Chicago, Harvard Gunawi, Prafail, Fate and Destin.

What new bugs live in the cloud?

Datacenter distributed systems

# of Bug Reports	Jan 2014	Jan 2016
Hadoop+MR+Yarn	17454	23811
HDFS	5710	9605
HBase	10263	15062
Cassandra	6535	10960
ZooKeeper	1854	2350

We studied 3000+ issues

"New" classes of bugs

- Distributed concurrency bugs
 - + Timings of multiple failures
- Non-deterministic performance bugs
- Scalability bugs
- Other outage-causing bugs:
 - SPOF/cascading bugs
 - Cross-layer upgrade bugs

TaxDC [ASPLOS '16]
SACM [OSDI '14]
FATE & DESTINI [NSDI '11]
The Tail at Store [FAST '16]
SPV [HotCloud '15]
Limpware [SoCC '13]
Tiny Tail [In Subm.]
Path-Based Spec. Exec. [In Subm.]
SCK [In Subm.]
Cloud Bug Study [SoCC '14]
Cloud Outage Study [In Subm.]

"New" classes of bugs

- Distributed concurrency (DC) bugs
 - TaxDC [ASPLOS '16]
 - SACM [OSDI '14]
- Non-deterministic performance bugs
- Scalability bugs

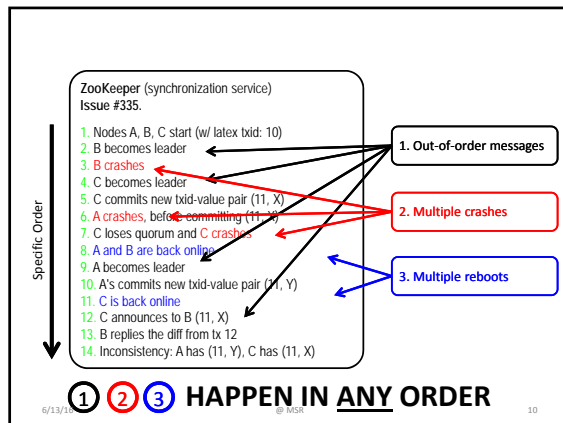
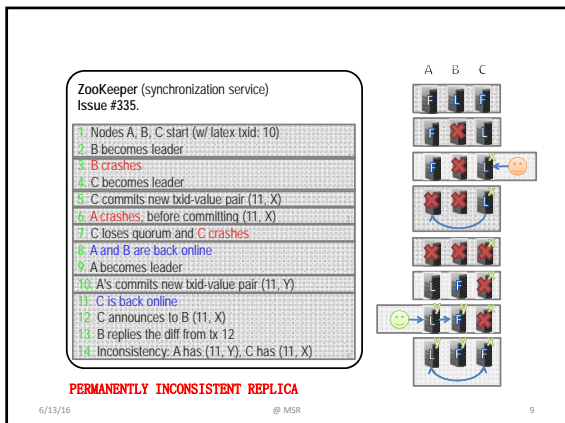
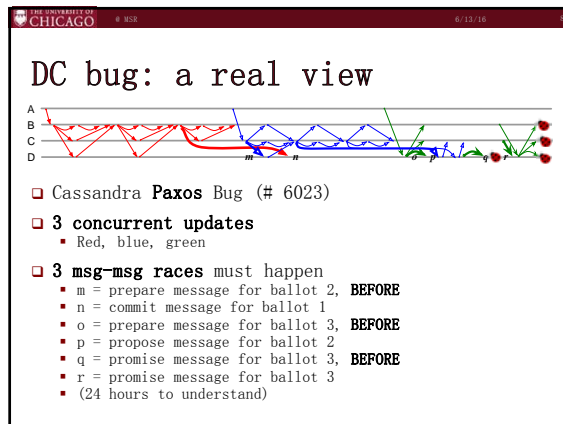
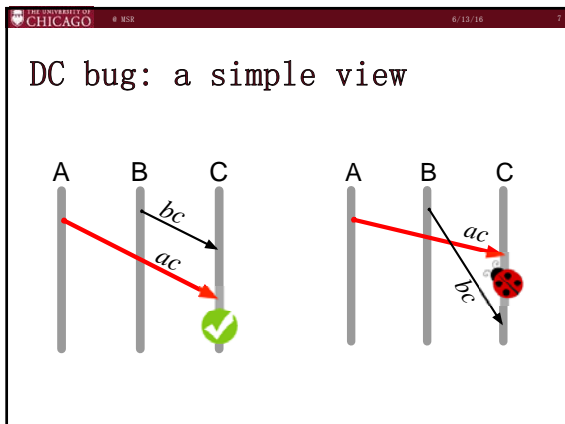
Distributed concurrency (DC) bug

- Caused by non-deterministic timing of concurrent events involving more than one node
- Events: Messages, crashes, reboots, timeouts, local computations

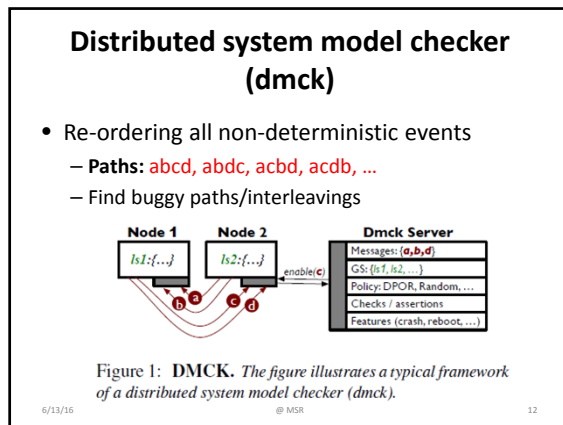
6% of the bugs in our study

CAUTION

Data loss, downtimes, inconsistent replicas, hanging jobs, etc.



How can we catch deep concurrency bugs in distributed systems?



Event re-orderings by dmck

ZooKeeper (synchronization service)
Issue #335.
Permanent inconsistent data
 1. Nodes A, B, C start (w/ latest bid: 10)
 2. B becomes leader
 3. B crashes
 4. C becomes leader
 5. C commits new bid-value pair ((1, X))
 6. A crashes, before committing the new bid !!
 7. C loses quorum and C crashes
 8. A and B are back online after C crashes
 9. A becomes leader
 10. A's commits new bid-value pair ((1, Y))
 11. C is back online after A's new tx commit
 12. C announce to B ((1, X))
 13. B replies diff starting with tx 12
 14. Inconsistency: A has ((1, Y)), C has ((1, X))

3
4
5
6
7
8
9
10
11
12
13

✓

2
7
1
4
5
6
7
8
11
12
10
9
10
11
12
14
13

✓

6
9
3
4
5
1
7
8
2
10
10
11
13
12
14

✓

1
2
3
4
5
6
7
8
9
10
11
12
13
14

✗

2
1
3
4
5
6
7
8
9
10
11
12
13

✓

6/13/16 @ MSR 13

SAMC: Semantic-Aware Model Checking

for Fast Discovery of Deep DC Bugs

with Tanakorn Leesatapornwongsa,
Mingzhe Hao, Pallavi Joshi, and Jeffrey F. Lukman
[OSDI '14]

What's Wrong with Existing Model Checkers?

- Last 7 years
 - MaceMC [NSDI '07], Modist [NSDI '09], dBug [SSV '10], Demeter [SOSP '13], etc.
- BUT**
 - Too many events to permute
 - Must add **crashes** and **reboots**
 - State-space explosion!
 - (skipped in existing checkers)
 - Cannot find deep bugs!

ZooKeeper (synchronization service)
Issue #335.
Permanent inconsistent data
 1. Nodes A, B, C start (w/ latest bid: 10)
 2. B becomes leader
 3. B crashes
 4. C becomes leader
 5. C commits new bid-value pair ((1, X))
 6. A crashes, before committing the new bid !!
 7. C loses quorum and C crashes
 8. A and B are back online after C crashes
 9. A becomes leader
 10. A's commits new bid-value pair ((1, Y))
 11. C is back online after A's new tx commit
 12. C announce to B ((1, X))
 13. B replies diff starting with tx 12
 14. Inconsistency: A has ((1, Y)), C has ((1, X))

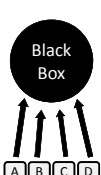
100 events

6/13/16 @ MSR 15

How can we catch deep bugs **REALLY FAST?**

6/13/16 @ MSR 16

- Why are existing checkers slow?
- They treat target system as a **black box**
 - Must re-order everything



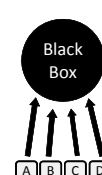
Black Box

Black Box Model Checker

ABCD
ABDC
ACBD
ACDB
ADBC
...
(24 total)


6/13/16 @ MSR 17

- How can we make model checkers fast?
 - Exploit **semantic knowledge**
 - E.g. knowledge of how messages are processed
 - Reduce unnecessary re-orderings



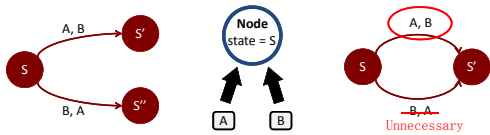
Black Box

Semantic Awareness



6/13/16 @ MSR 18

Dependency vs. Independence



A, B = Dependent

A, B = Independent

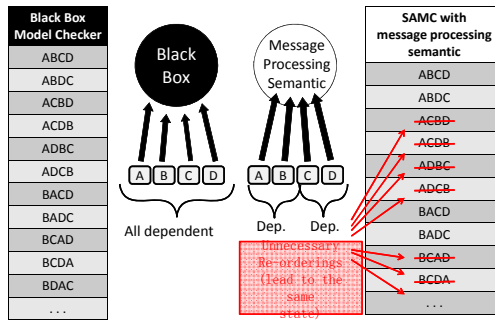
Independent = No need to reorder

6/13/16

@ MSR

19

Black Box vs. SAMC

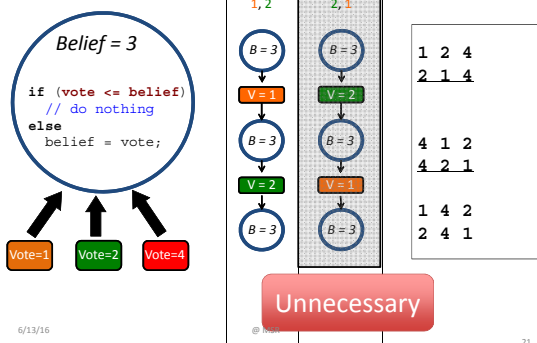


6/13/16

@ MSR

20

Message Processing Semantic in a Leader Election



6/13/16

@ MSR

21

Discard pattern

```

    MESSAGE PROCESSING SEMANTIC
    if (msg.vote <= state.belief)
      // do nothing
    else
      belief = vote;

    DISCARD PATTERN
    if (isDiscard(msg, state)) {
      // do nothing;
    }

    DISCARD PREDICATE
    boolean isDiscard(msg, state) {
      if (msg.vote <= state.belief)
        return true;
      else
        return false;
    }
  
```

6/13/16

@ MSR

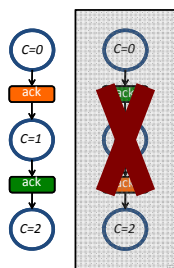
22

- Discard pattern
- Increment pattern
- Constant pattern

```

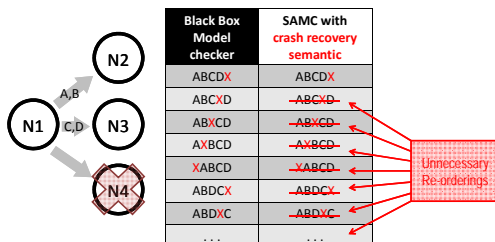
    if (msg.type == ack) {
      node.ackCount++;
    }

    boolean isIncrement(msg, ls) {
      if (msg.type == ack)
        return true;
      else
        return false;
    }
  
```



Local-Message Independence (LMI)

SAMC with Crashes



6/13/16

@ MSR

24

Crash-Msg Independence

```

void handleCrash() {
  if (X == follower &&
      isQuorum())
    followerCount--;
  // No new messages!!
}
    
```

Black Box

- ABCDX
- ABXC D
- ABXCD
- AXB CD
- XABCD
- ABDCX
- ...

Crash a **follower**
 → **Local Impact**
 (no new messages & only state changes in leader L)

6/13/16 @ MSR 25

Crash-Msg Independence

```

void handleCrash() {
  if (X == leader || !isQuorum())
    electLeader();
  // New messages created
}
    
```

Black Box

- ABCDX
- ABXC D
- ABXCD
- AXB CD
- XABCD
- ABDCX
- ...

Crash the **leader**
 → **Global Impact**
 (cannot prune re-orderings)

6/13/16 @ MSR 26

SAMC Architecture

SAMC

Protocol Specific Rules: Leader Election, Atomic Broadcast, ...

Generic Reduction Policies:

- Local-Message Indep. (LMI)
- Crash-Message Indep. (CMI)
- Crash Recovery Symmetry (CRS)
- Reboot Sync. Symmetry (RSS)

6/13/16 27

Protocol-specific predicates (extra)

(e.g. ZooKeeper Leader Election)

Local-Message Independence (LMI)	Crash-Message Independence (CMI)	Crash Recovery Symmetry (CRS)
<pre> bool pl : !newVote(a, s); bool pm : newVote(a, s); bool newVote(a, s) : ret i; if (a.sp > s.sp) ret i; else if (a.sp == s.sp) if (a.tx > s.tx) ret i; else if (a.tx == s.tx && m.lid > s.lid) ret i; ret 0; </pre>	<pre> bool pg (s, X) : if (s.r1 == F && X.r1 == L) ret i; if (s.r1 == L && X.r1 == F) ret i; if (s.r1 == S && X.r1 == S) ret i; bool pl (s, X) : if (s.r1 == L && X.r1 == F) ret i; bool quorumAfter(s) : ret ((s.fol-1) >= s.allF2); </pre>	<pre> bool pr1(s,C) : if (s.r1 == L && C.r1 == F) ret i; bool pr2(s,C) : if (s.r1 == L && C.r1 == F) ret i; bool pr3(s,C) : if (s.r1 == F && C.r1 == L) ret i; bool pr4 : if (s.r1 == S) ret i; </pre>

- 35 LOC on average per protocol

6/13/16 @ MSR 28

Speed in Reaching Old Bugs

#executions/paths to reach the bugs (e.g., 2 paths = abcd, abdc)

Bug#	SAMC	Black-Box DPOR	Random	Random DPOR
ZooKeeper-335				
ZooKeeper-790				
ZooKeeper-975				
ZooKeeper-1075				
ZooKeeper-1419				
ZooKeeper-1492				
ZooKeeper-1653				
MapReduce-4748				
MapReduce-5489				
MapReduce-5505				
Cassandra-3395				
Cassandra-3626				

6/13/16 @ MSR 29

Speed in Reaching Old Bugs

#executions/paths to reach the bugs (e.g., 2 paths = abcd, abdc)

Bug#	SAMC	Black-Box DPOR	Random	Random DPOR
	#exe	#exe speedup	#exe speedup	#exe speedup
ZooKeeper-335	117	5000+	43+	1057
ZooKeeper-790	7	14	2	225
ZooKeeper-975	53	967	18	71
ZooKeeper-1075	16	1081	68	86
ZooKeeper-1419	100	924	9	2514
ZooKeeper-1492	576	5000+	9+	5000+
ZooKeeper-1653	11	945	86	3756
MapReduce-4748	4	22	6	6
MapReduce-5489	53	5000+	94+	5000+
MapReduce-5505	40	1212	30	5000+
Cassandra-3395	104	2552	25	191
Cassandra-3626	96	5000+	52+	5000+

6/13/16 @ MSR 30

Summary

- Distributed concurrency bugs → hard to catch
- **Semantic-awareness** for model checking is powerful
 - Find bugs **2 - 340x** faster, **49x** on average

6/13/16

@ MSR

31

TaxDC:

Taxonomy of Non-Deterministic Concurrency Bugs in Datacenter Distributed Systems

with Tanakorn Leesatapornwongsa,
Jeffrey F. Lukman and Shan Lu
[ASPLOS '16]

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 33

Google local concurrency bug

(LC bug: multi-threaded single machine software)

Learning from mistakes: a comprehensive study on real world concurrency bug characteristics
S. Lu, S. Park, E. Seo, Y. Zhou - ACM Sigplan Notices, 2008 - dl.acm.org
Cited by 558 **Top 10 most cited ASPLOS paper**

≠

Google distributed concurrency bug


Learning from mistakes: a comprehensive study on real world concurrency bug characteristics
S. Lu, S. Park, E. Seo, Y. Zhou - ACM Sigplan Notices, 2008 - dl.acm.org

PDF TaxDC: A Taxonomy of Non-Deterministic Concurrency Bugs in Datacenter Distributed Systems
T. Leesatapornwongsa, J.F. Lukman, S. Lu, H.S. Gunawi - ucare.cs.uchicago.edu
Cited by 1

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 34

TaxDC

- Taxonomy of distributed concurrency bugs
- **104** bugs
- **4** varied distributed systems



- Bugs in **2011-2014**
- Study description, source code, patches

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 35

Detailed Characteristics

Input:
4 Protocol initiations

ZooKeeper-1264

1. Follower F crashes, reboots, and joins cluster
2. Leader L sync snapshot with F
3. Client requests new update, F applies this only in memory
4. Sync finishes
5. Client requests other update, F writes this to disk correctly
6. F crashes, reboots, and joins cluster again
7. This time L sends only diff after update in step 5.
8. F loses update in step 3.

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 36

Detailed Characteristics

Input:
4 Protocols
- 2 faults
- 2 reboots

ZooKeeper-1264

1. Follower F **crashes, reboots,** and joins cluster
2. Leader L sync snapshot with F
3. Client requests new update, F applies this only in memory
4. Sync finishes
5. Client requests other update, F writes this to disk correctly
6. F **crashes, reboots,** and joins cluster again
7. This time L sends only diff after update in step 5.
8. F loses update in step 3.

Detailed Characteristics

ZooKeeper-1264

1. Follower F crashes, reboots, and joins cluster
2. Leader L sync snapshot with F
3. Client requests new update, F applies this only in memory
4. Sync finishes
5. Client requests other update, F writes this to disk correctly
6. F crashes, reboots, and joins cluster again
7. This time L sends only diff after update in step 5.
8. F loses update in step 3.

Input:
- 4 Protocols
- 2 faults
- 2 reboots

Timing:
- Atomicity violation
- Fault Timing

Detailed Characteristics

ZooKeeper-1264

1. Follower F crashes, reboots, and joins cluster
2. Leader L sync snapshot with F
3. Client requests new update, F applies this only in memory
4. Sync finishes
5. Client requests other update, F writes this to disk correctly
6. F crashes, reboots, and joins cluster again
7. This time L sends only diff after update in step 5.
8. F loses update in step 3.

Input:
- 4 Protocols
- 2 reboots

Fix:
Delay msg.

Timing:
- Atomicity violation
- Fault Timing

Error:
- Global

Failure:
Data inconsistency

Detailed Characteristics

ZooKeeper-1264

1. Follower F crashes, reboots, and joins cluster
2. Leader L sync snapshot with F
3. Client requests new update, F applies this only in memory
4. Sync finishes
5. Client requests other update, F writes this to disk correctly
6. F crashes, reboots, and joins cluster again
7. This time L sends only diff after update in step 5.
8. F loses update in step 3.

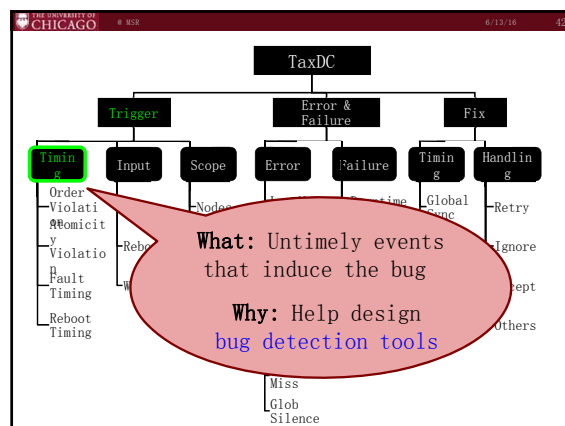
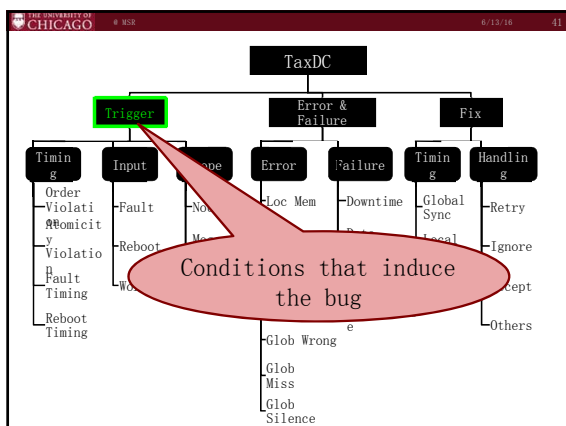
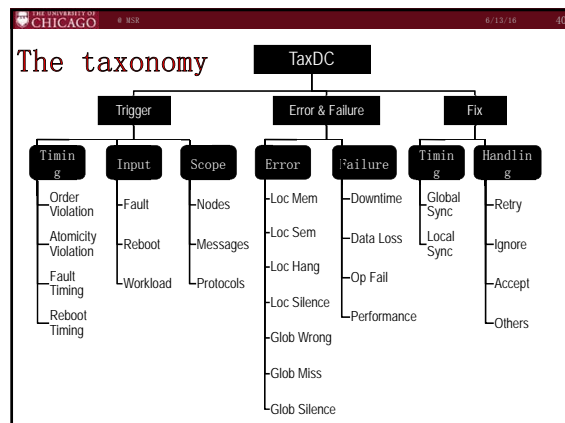
Input:
- 4 Protocols
- 2 faults
- 2 reboots

Timing:
- Atomicity violation
- Fault Timing

Error:
- Global

Failure:
Data inconsistency

Fix:
Delay msg.



Trigger

- Timing
- Message

Messages arrive in untimely order

Trigger

- Timing
- Message
- Order violation (44%)

Y must happen **after** X
But Y happens **before** X

Trigger

- Timing
- Message
- Order violation (44%)
- Msg-msg race

Y must happen **after** X
But Y happens **before** X

Ex: MapReduce-3274

Trigger

- Timing
- Message
- Order violation (44%)
- Msg-msg race

MapReduce-3274 HBase-5780 MapReduce-5358

Trigger

- Timing
- Message
- Order violation (44%)
- Msg-msg race
- Msg-compute race

Ex: MapReduce-4157

Trigger

- Timing
- Message
- Order violation (44%)
- Atomicity violation (20%)

A message comes in the **middle** of "atomic" operation

Ex: Cassandra-1011, HBase-4729, MapReduce-5009, Zookeeper-1496

Trigger

- Timing
- Message
- Fault (21%)

Fault at specific timing

No fault timing in LC bugs
Only in DC bugs

Ex: Cassandra-6415, Hbase-5806, MapReduce-3858, Zookeeper-1653

Trigger

- Timing
- Message
- Fault
- Reboot (11%)

Reboot at specific timing

Ex: Cassandra-2083, Hadoop-3186, MapReduce-5489, Zookeeper-975

Trigger

- Timing

Implication: simple patterns can inform pattern-based bug detection tools, etc.

Message timing Fault timing Reboot

TaxDC

- Trigger
 - Timing
 - Input
 - Scope
- Error & Failure
 - Error
 - Failure
- Fix
 - Timing
 - Handling

What: Input to exercise buggy code

Why: Improve testing coverage

Trigger

- Timing
- Input
- Fault

"How many bugs require fault injection?"

37% = No fault 63% = Yes

"What kinds of fault? & How many times?"

88% = No timeout 12%

53% = No crash 35% = 1 crash 12%

Real-world DC bugs are NOT just about message re-ordering, but faults as well

Trigger

- Timing
- Input
- Fault
- Reboot

"How many reboots?"

73% = No reboot 20% = 1 7%

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 55

Cassandra Paxos bug

(Cassandra-6023)

3 concurrent user requests!

"How many protocol initiations to run as input?"

20% = 1	29% = 2	24% = 3	27% = 4+
---------	---------	---------	----------

Implication: **multiple protocols** for DC testing

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 56

TaxDC

- Trigger
 - Timing
 - Order Violation
 - Atomicity Violation
 - Fault
 - Reboot
 - Timing
- Error & Failure
 - Scope
 - Error
 - Failure
- Fix
 - Timing
 - Handling

What: How developers fix bugs

Why: Help design runtime prevention and automatic patch generation

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 57

- Trigger
- Error
- Fix
- Complex

Add new states & transitions

Similar to fixing LC bugs: add synchronization e.g. lock()

Add Global Synchronization

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 58

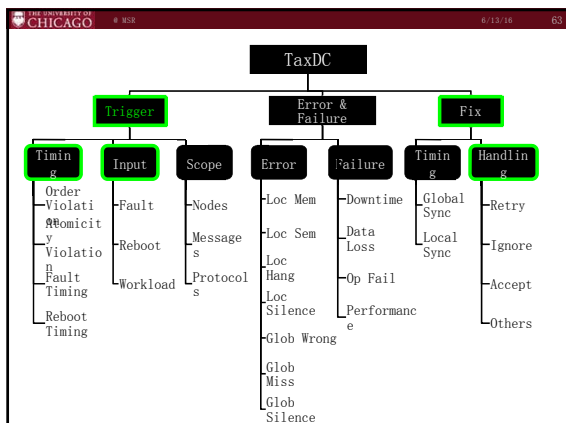
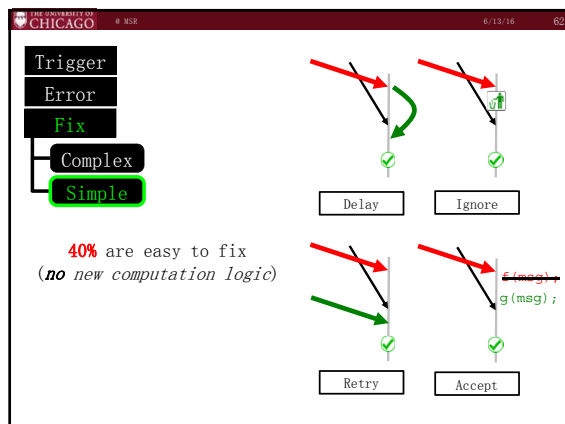
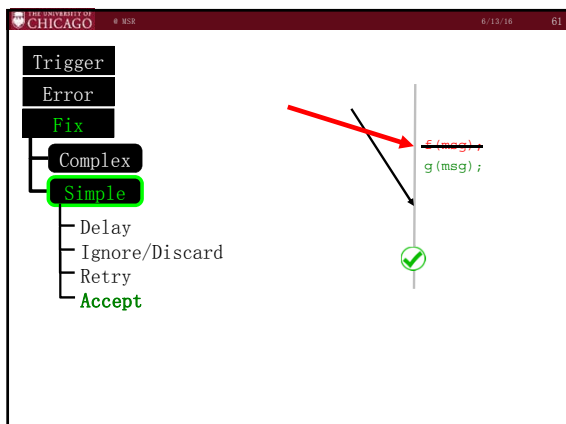
- Trigger
- Error
- Fix
- Complex
- Simple
- Delay

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 59

- Trigger
- Error
- Fix
- Complex
- Simple
- Delay
- Ignore/discard

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 60

- Trigger
- Error
- Fix
- Complex
- Simple
- Delay
- Ignore/Discard
- Retry



- ### Challenges & Opportunities in ...
- ❑ Distributed system model checker
 - ❑ Formal verification
 - ❑ DC bug detection
 - ❑ Runtime prevention

Distributed System Model Checkers

Event	Modist <i>NSDI'11</i>	Demeter <i>SOSP'11</i>	MaceMC <i>NSDI'07</i>	SAMC <i>OSDI'14</i>	Reality
Message	✓	✓	✓	✓	✓
Crash	✓	✓	✓	✓	✓
Multiple crashes	✗	✗	✗	✓	✓
Reboot	✗	✗	✗	✓	✓
Multiple reboots	✓	✓	✓	✗	✓
Timeout	✓	✓	✗	✗	✓
Computation	✗	✗	✗	✗	✓
Disk fault					

Formal Verification

- ❑ State-of-the-art
 - Verdi [PLDI'15]
 - Raft update protocol
 - IronFleet [SOSP'15]
 - Paxos update protocol
 - Lease-based read/write
- ❑ Challenges
 - Foreground & Background
 - 19% = FG 52% = BG 29% = Mix
 - #Protocol interactions
 - 20% = 1 80% = 2+ Protocols

Only verify foreground protocols

Foreground & background

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 67

DC Bug Detection

- State-of-the-art:
 - LC bug detection
 - Pattern-based detection
 - Error-based detection
 - Statistical bug detection
- Opportunities:
 - DC bug detection?
 - Pattern-based detection

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 68

Runtime Failure Prevention

- State-of-the-art:
 - LC bug prevention
 - Deadlock Immunity [OSDI '08]
 - Aviso [ASPLOS '13]
 - ConAir [ASPLOS '13]
 - (many more)
- Opportunities:
 - DC bug prevention

Fixes

60% = Complex 40% = Simple

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 69

Dev' s comments on DC bugs

- "Do we have to rethink this entire [HBase] root and meta ' huh hah' ? There isn' t a week going by without some new bugs about races between splitting and assignment [distributed protocols]. " — hbase4397
- "That is one monster of a race!" — mr3274
- "This has become quite messy, we didn' t foresee some of this during design, sigh. " — mr4819
- "Great catch, Sid! Apologies for missing the race condition " — mr4099
- "We have already found and fix many cases ... however it seems exist many other cases. " — hb6147

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 70

"New" classes of bugs ...

- Distributed concurrency bugs
- Non-deterministic performance bugs
 - Limpware [SoCC '13]
 - Detect performance bugs [HotCloud '15]
 - Path-Based Spec. Exec. [In Subm.]
- Scalability bugs

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 71

A "limpware" anecdote

(limping hardware)

- "... 1Gb NIC card on a machine that suddenly only transmits at 1 kbps,
- this slow machine caused a chain reaction upstream
- in such a way that the 100 node cluster began to crawl at a snail's pace,
- making the system non-availible for all practical purposes. " — Ben Facebook

Limping NIC!

Cascading impact!

THE UNIVERSITY OF CHICAGO # MSR 6/13/16 72

Limpware, really?

- "In 2011, one of the DDN 9900 units had 4 servers having high wait times on I/O for a certain set of disk LUNs. The maximum wait time was 103 seconds. This was left uncorrected for 50 days. " - Kasick of CMU, Harms of Argonne
- "The disk attempts to re-read each block multiple times before responding. " - Baptist of Cleversafe
- "On Intrepid, we had a bad batch of optical transceivers with an extremely high error rate. That results in an effective throughput of 1-2 Kbps. " - Harms of Argonne
- Many others: "Yes, we've seen that in production"

Limpware impacts?

- Modern distributed systems are ...
 - ... fault tolerant
 - ... limpware tolerant?
- Limpware-injection experiments
 - Run HDFS, Hadoop, ZooKeeper, Cassandra, Hbase
 - Run load-intensive workload + inject limpware
 - E.g. slow a NIC to 1 Mbps, 0.1 Mbps, etc.

An example

- Run a distributed protocol
 - E.g., write pipeline in HDFS
- Measure slowdowns under:
 - No failure, crash, a limping NIC

The diagram shows a write pipeline with three nodes. The second node is highlighted in yellow, and the third node is highlighted in red. To the right, a bar chart shows execution slowdowns on a logarithmic scale (1x, 10x, 100x, 1000x slower) for different NIC speeds: 10 Mbps NIC (green bar, ~10x slower), 1 Mbps NIC (red bar, ~100x slower), and 0.1 Mbps NIC (black bar, ~1000x slower).

Benchmarks

ID	Protocol	Limp-ware	Injected Node	Workload	Base Latency
F1	Logging	Disk	Master	Create 8000 empty files	12
F2	Write	Disk	Data	Create 30 64-MB files	182
F3	Read	Disk	Data	Read 30 64-MB files	120
F4	Metadata Read/Logging	Disk	Master	Stats 1000 files + heavy updates	9
F5	Checkpoint	Disk	Secondary	Checkpoint 60K transactions	39
F6	Write	Net	Data	Create 30 64-MB files	208
F7	Read	Net	Data	Read 30 64-MB files	104
F8	Regeneration	Net	Data	Regenerate 90 blocks	432
F9	Regeneration	Net	Data-S/Data-D	Scale replication factor from 2 to 4	11
F10	Balancing	Net	Data-O/Data-U	Move 3.47 GB of data	4105
F11	Decommission	Net	Data-L/Data-R	Decommission a node having 90 blocks	174
H1	Speculative execution	Net	Mapper	WordCount: 512 MB dataset	80
H2	Speculative execution	Net	Reducer	WordCount: 512 MB dataset	80
H3	Speculative execution	Net	Job Tracker	WordCount: 512 MB dataset	80
H4	Speculative execution	Net	Task Node	1000-task Facebook workload	4320
Z1	Get	Net	Leader	Get 7000 1-KB znodes	4
Z2	Get	Net	Follower	Get 7000 1-KB znodes	5
Z3	Set	Net	Leader	Set 7000 1-KB znodes	23
Z4	Set	Net	Follower	Set 7000 1-KB znodes	26
Z5	Set	Net	Follower	Set 20KB data 6000 times to 100 znodes	64
C1	Put (quorum)	Net	Data	Put 240K KeyValues	66
C2	Get (quorum)	Net	Data	Get 45K KeyValues	73
C3	Get (one) + Put (all)	Net	Data	Get 45K KeyValues + heavy puts	36
B1	Put	Net	Region Server	Put 300K KeyValues	61
B2	Get	Net	Region Server	Get 300K KeyValues	151
B3	Scan	Net	Region Server	Scan 300K KeyValues	17
B4	Cache Get/Put	Net	Data-H	Get 100 KeyValues + heavy puts	4
B5	Compaction	Net	Region Server	Compact + 100-MB sstables	122

Fail-stop tolerant, but not limpware tolerant

(no failover)

The grid shows 20 small bar charts, each representing a different task from the benchmarks table. Each chart compares execution slowdowns under four conditions: 'No failure', 'Node crash', '10 Mbps NIC', and '0.1 Mbps NIC'. Red arrows indicate significant slowdowns when the NIC is slowed down, demonstrating that the system is not limpware tolerant.

(The root causes are in Limpware paper [SOCC '13]; this talk focuses on Hadoop MapReduce)

Hadoop MapReduce

- Supposedly tail tolerant
- Why not limpware tolerant?
- Why Speculative Execution fails:

The four bar charts show that speculative execution (Spec. Exec.) for all components (Mapper, Reducer, Job Tracker, Task Node) experiences significant slowdowns when the NIC is slowed down, indicating that MapReduce is not limpware tolerant.

Loophole #1

- Backup task reads from the same slow datanode
 - Hadoop and HDFS don't cooperate
 - No history of bad "paths"

The diagram shows two Datanodes, A and B, and two Mappers, M1 and M2. M1 is connected to A, and M2 is connected to B. A red arrow points from M2 to B, indicating that M2 is reading from the slow datanode B, even though A is faster. This is the 'loophole' where Hadoop and HDFS do not cooperate to avoid slow paths.

Loophole #2

- All reducers fetch from a mapper with a slow NIC
 - All reducers slow → no straggler
 - M2 reads data locally (not slow)
- (many other loopholes in the paper)

Cascading failures

- A limping NIC → limping tasks
 - (Limping tasks are slower by orders of magnitude)
- Limping tasks use up slots → limping node
 - If all slots are used → node is “unavailable”
- All nodes in limp mode → limping cluster

Cluster collapse

- Macrobenchmark: Facebook Hadoop workload
 - 30-node cluster
 - One node w/ limping NIC (0.1 Mbps)

Formalizing the problem

- A job = various deployment scenarios
- Untriggered speculative execution
 - (DSR_i & FTY_i & FPL_i & DLC_i) or (JCH_i & TPL_i & FTY_i & FPL_i) or ...
 - Unanticipated scenario

Scenario Type	Possible Conditions
DLC: Data Locality	(1) Read from remote disk, (2) read from local disk, ...
DSR: Data Source	(1) Some tasks read from same datanode, (2) all tasks read from different datanodes, ...
JCH: Job Characteristic	Map-reduce is all-to-all, (2) all-to-many, (3) large fan-in, (4) large fan-out, ...
JSZ: Job Size	(1) 1 GB jar file, (2) 1 MB jar file, ...
LSZ: Load Size	(1) Thousands of tasks, (2) small number of tasks, ...
FTY: Fault Type	(1) Slowdown fault injection at the mapper, (2) Node disconnect/racklet drop, (3) Disk error/out of space, (4) Rack switch, ...
FPL: Fault Placement	Slowdown fault injection at the mapper, (2) mapper, (3) reducer, ...
FGF: Fault Granularity	(1) Single disk/NIC, (2) single node (deadnode), (3) entire rack (network switch), ...
FTM: Fault Timing	(1) During shuffling, (2) during 95% of task completion, ...
TOP: Topology Scenario	(1) 30 nodes per rack, (2) 3 nodes per rack, ...
TPL: Task Placement	(1) Mappers and reducers are in different nodes, (2) AM and reducers in different nodes, (3) Mappers are in the same node, (4) Most of reducers placed in the same rack, ...

Non-deterministic performance bugs

- Untriggered Speculative Execution
 - MR-70001 = JCH_i & TPL_i & FPL_i & FTY_i
 - MR-70002 = DSR_i & DLC_i & FPL_i & FTY_i
 - MR-5533 = FTY_i & FPL_i & TPL_i
 - ...
- 0(n) Recovery
 - MR-5251 = FTY_i & FPL_i & FTM_i
 - MR-5080 = TPL_i & TPL_i & FTY_i & FPL_i
 - MR-1800 = TPL_i & TPL_i & FTY_i & TOP_i
 - ...
- Long lock contention
 - MR-9191 = FTY_i & FPL_i & FTM_i
 - MR-9292 = TPL_i & TPL_i & FTY_i & FPL_i
 - MR-9393 = TPL_i & TPL_i & FTY_i & TOP_i
 - ...

Scenario Type	Possible Conditions
DLC: Data Locality	(1) Read from remote disk, (2) read from local disk, ...
DSR: Data Source	(1) Some tasks read from same datanode, (2) all tasks read from different datanodes, ...
JCH: Job Characteristic	Map-reduce is all-to-all, (2) all-to-many, (3) large fan-in, (4) large fan-out, ...
JSZ: Job Size	(1) 1 GB jar file, (2) 1 MB jar file, ...
LSZ: Load Size	(1) Thousands of tasks, (2) small number of tasks, ...
FTY: Fault Type	(1) Slowdown fault injection at the mapper, (2) Node disconnect/racklet drop, (3) Disk error/out of space, (4) Rack switch, ...
FPL: Fault Placement	Slowdown fault injection at the mapper, (2) mapper, (3) reducer, ...
FGF: Fault Granularity	(1) Single disk/NIC, (2) single node (deadnode), (3) entire rack (network switch), ...
FTM: Fault Timing	(1) During shuffling, (2) during 95% of task completion, ...
TOP: Topology	(1) 30 nodes per rack, (2) 3 nodes per rack, ...
TPL: Task Placement	(1) Mappers and reducers are in different nodes, (2) AM and reducers in different nodes, (3) Mappers are in the same node, (4) Most of reducers placed in the same rack, ...

Perf. Model Checking [HotCloud '15]

- Goal: Permute many topological/failure/placement scenarios
- Real Java code → Colored Petri Nets (CPN) model
 - Automated conversion (“compiler”)
 - Abstract system-level constructs
 - E.g., queues, tasks, resources, locks
- Permute the scenarios in CPN
- Abstract performance faults
 - Boolean result: limping or not
 - No need for precise latency/bandwidth predictions
- Test the buggy scenarios in real runs

THE UNIVERSITY OF CHICAGO MSR 6/13/16 85

Path Based Spec. Exec. [In Subm.]

- Hadoop SE:
 - **Straggler**: if task T's progress is **slower** than the rest
 - Task T is just a **progress score** → fundamental flaw
- Our observation:
 - Task T is a **path**
 - **Map path**: source datanode → map node
 - **Shuffle path**: map node → reduce node
 - **Output path**: reduce node → pipeline of datanodes
- **PBSE**: Path-based speculative execution
 - It's about the progress of individual "paths"
 - SE algorithm is based on path progress
 - Diverse paths: no single point of path failure

THE UNIVERSITY OF CHICAGO MSR 6/13/16 86

Conclusion

- Distributed concurrency bugs
- Non-deterministic performance bugs
- Scalability bugs
- Other outage-causing bugs:
 - SPOF/cascading bugs
 - Cross-layer upgrade bugs

The complexity of cloud-scale hardware and software ecosystem has outpaced existing testing, debugging, and verification tools.

Many new classes of bugs to hunt!

THE UNIVERSITY OF CHICAGO MSR 6/13/16 87

Thank you! Questions?



ucare.cs.uchicago.edu



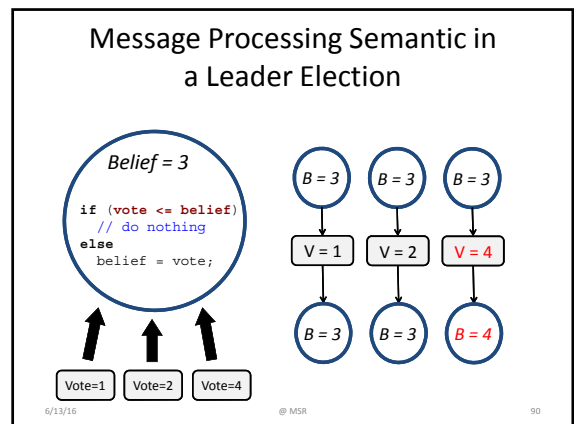
ceres.cs.uchicago.edu

THE UNIVERSITY OF CHICAGO MSR 6/13/16 88

EXTRA

THE UNIVERSITY OF CHICAGO MSR 6/13/16 89

Extra -- SAMC

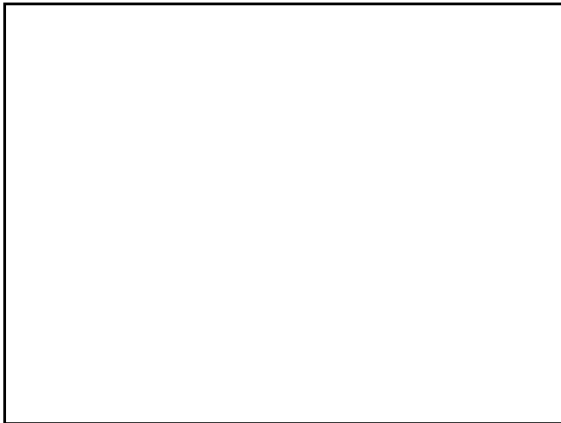
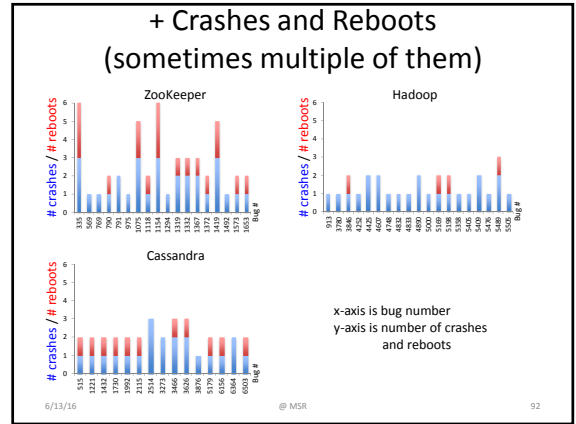


THE UNIVERSITY OF CHICAGO @ MSR 6/13/16 91

SAMC server logic (extra)

vote	belief	isDiscard
1	3	true
2	3	true
4	3	false

m_x	m_y	discard(m_x)	discard(m_y)	Independent
1	2	true	true	✓
1	4	true	false	x
2	4	true	false	x



THE UNIVERSITY OF CHICAGO @ MSR Re-orderings via Message Processing Semantic 6/13/16

Message Processing Semantic

```

if (vote <= belief)
  // do nothing
else
  belief = vote;
    
```

Errors, Faults, Failure

- To quote the [Software Engineering Body of Knowledge](#)
- Different cultures and standards may use somewhat different meanings for these terms, which have led to attempts to define them.
- Partial definitions taken from standard (IEEE610.12-90) are:
- Error: "A difference...between a computed result and the correct result"
- Fault: "An incorrect step, process, or data definition in a computer program"
- Failure: "The [incorrect] result of a fault"
- Mistake: "A human action that produces an incorrect result"

6/13/16 @ MSR 95