

THE UNIVERSITY OF CHICAGO

UNEARTHING CONCURRENCY AND SCALABILITY BUGS IN CLOUD-SCALE
DISTRIBUTED SYSTEMS

A DISSERTATION SUBMITTED TO
THE FACULTY OF THE DIVISION OF THE PHYSICAL SCIENCES
IN CANDIDACY FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY

DEPARTMENT OF COMPUTER SCIENCE

BY
TANAKORN LEESATAPORNWONGSA

CHICAGO, ILLINOIS

JUNE 2017

Copyright © 2017 by Tanakorn Leesatapornwongsa
All Rights Reserved

To my family: father, mother, Louise, Fon, Nuch, and Nid

“Debugging is twice as hard as writing the code in the first place. Therefore, if you write the code as cleverly as possible, you are, by definition, not smart enough to debug it.” — Brian Kernigham

TABLE OF CONTENTS

LIST OF FIGURES	viii
LIST OF TABLES	ix
ACKNOWLEDGMENTS	x
ABSTRACT	xi
1 INTRODUCTION	1
1.1 Distributed Concurrency Bugs	2
1.1.1 Distributed Concurrency Bug Study & Taxonomy	3
1.1.2 Semantic-Aware Model Checking	4
1.2 Scalability Bugs	5
1.3 Summary of Contributions and Outline	8
2 BACKGROUND AND RELATED WORK	9
2.1 Cloud Computing and Cloud-Scale Distributed Systems	9
2.1.1 Cloud Computing	9
2.1.2 Cloud-Scale Distributed Systems	10
2.2 Distributed Concurrency	11
2.2.1 Local Concurrency and Distributed Concurrency	11
2.2.2 Distributed Systems Model Checker (DMCK)	13
2.2.3 Symbolic Execution	14
2.2.4 Fault Injector	15
2.3 Scalability	15
2.3.1 Vertical Scaling vs Horizontal Scaling	15
2.3.2 Scalability Testing	16
2.3.3 Scalability Benchmarking	17
2.4 Conclusion	17
3 TAXDC: A TAXONOMY OF NON-DETERMINISTIC CONCURRENCY BUGS IN CLOUD DISTRIBUTED SYSTEMS	18
3.1 TaxDC	19
3.2 Methodology	20
3.2.1 Basic Definitions	20
3.2.2 Target Systems and Dataset	20
3.2.3 Taxonomy	21
3.2.4 Threats to Validity	21
3.2.5 TaxDC Database	23
3.2.6 Detailed Terminologies	23
3.3 Trigger	24
3.3.1 Timing Conditions (TC)	24

3.3.2	Input Preconditions (IP)	28
3.3.3	Triggering Scope (TS)	30
3.4	Errors and Failures	33
3.4.1	Error Symptoms	33
3.4.2	Failure Symptoms	34
3.5	Fixes	35
3.5.1	Message Timing Bug Fixes	35
3.5.2	Fault/Reboot Timing Bug Fixes	36
3.6	Root Causes	38
3.7	Other Statistics	39
3.8	LC bugs vs. DC bugs	39
3.9	Conclusion	40
4	SAMC: SEMANTIC-AWARE MODEL CHECKING FOR FAST DISCOVERY OF DC BUGS IN CLOUD DISTRIBUTED SYSTEMS	42
4.1	Background	42
4.1.1	DMCK Framework and Terms	42
4.1.2	State-of-the-Art DMCKs	44
4.1.3	Does State of the-Art Help?	46
4.2	SAMC	47
4.2.1	Semantic-Awareness Example	47
4.2.2	Architecture	48
4.2.3	Semantic-Aware Reduction Policies	49
4.2.4	Pattern Extraction	56
4.3	Implementation and Integration	57
4.3.1	SAMPRO	57
4.3.2	Integration to Target Systems	58
4.4	Evaluation	58
4.4.1	Speed in Finding Old Bugs	60
4.4.2	Ability of Finding New Bugs	63
4.4.3	Reduction Ratio	63
4.5	Discussion	65
4.5.1	Simplicity	65
4.5.2	Generality	66
4.5.3	Soundness	66
4.6	Conclusion	67
5	SCK: A SINGLE-MACHINE APPROACH FOR FINDING SCALABILITY BUGS IN CLOUD-SCALE DISTRIBUTED SYSTEMS	68
5.1	Scalability Bug	69
5.1.1	A Sample Cassandra Bug	70
5.1.2	Observations	72
5.2	State of the Art of Large-Scale Emulation	74
5.3	SCK	75

5.3.1	Processing Illusion (PIL)	75
5.3.2	Single-Process Cluster (SPC)	80
5.3.3	Global Event Driven Arch. (GEDA)	80
5.3.4	Memory Footprint Reduction (MFR)	81
5.3.5	Putting It All Together	82
5.4	Implementation and Integration	83
5.5	Evaluation	85
5.5.1	Colocation Factor	85
5.5.2	Accuracy	87
5.5.3	Bugs Reproduced	90
5.5.4	New Bugs	91
5.6	Conclusion	93
6	CONCLUSIONS AND FUTURE WORK	94
6.1	Conclusion	94
6.1.1	Distributed Concurrency Bugs	94
6.1.2	Scalability Bugs	95
6.2	Future Work	96
6.2.1	Distributed Concurrency Bugs	96
6.2.2	Scalability Bugs	103
6.3	Closing Words	104
	REFERENCES	106

LIST OF FIGURES

2.1	ZooKeeper-1264 bug	12
2.2	DMCK	14
3.1	Triggering patterns	27
3.2	A DC bug in ZooKeeper	29
3.3	Statistical overview of TaxDC	31
3.4	A Cassandra's Paxos bug	32
4.1	How dmck works	44
4.2	SAMC Architecture	48
4.3	LMI and CMI	51
4.4	Crash Recovery in Leader Election	54
4.5	Complexity of Deep Bugs	61
5.1	The problem of gossip-based failure detection in Cassandra	69
5.2	Maximum colocation factor	85
5.3	Cassandra internal metrics	88
5.4	Accuracy in reproducing CA-6127	89
5.5	Accuracy in reproducing other bugs	92
6.1	Race of HBase's messages to ZooKeeper	102

LIST OF TABLES

3.1	Taxonomy of DC Bugs	22
3.2	#DC bugs triggered by timing conditions	25
3.3	First error symptoms of DC bugs	34
3.4	Fix strategies for message timing bugs	36
3.5	Fix strategies for fault/reboot timing bugs	37
4.1	Protocol-Specific Reduction Rules for ZLE	59
4.2	SAMC Speed in Finding Old Bugs	60
4.3	SAMC Reduction Ratio	65
5.1	Integration (LOC)	83
5.2	Reproduced bugs	90

ACKNOWLEDGMENTS

This Ph.D. could not be accomplished, if I did not get supports from faculty, family, and friends, which I would like to thank these individuals here.

The first person I need to thank is my adviser (and of course, also one of the dissertation committee), Prof. Haryadi Gunawi. He guided me since the beginning to the end. I have learned a lot from him from “*how to survive Ph.D.?*” to “*how to find a job?*”. It is my great pleasure to have him as my adviser (and it is also a great experience to be his *first* student).

Next, I need to thank the other two committee members, Professor Shan Lu and Professor Ravi Chugh that kindly help to be my committee. I appreciate their time and their suggestions on my presentation (which is my job talk). I also had a chance to work with Prof. Lu in one project which is one part in this dissertation. Working with Prof. Lu taught me many great lessons.

I also need to thank my colleagues (aka co-authors), Tiratat Patana-anake, Mingzhe Hao, Pallavi Joshi, Riza Suminto, Thanh Do, Jeffrey Lukman, Huan Ke, Cesar Stuardo, Daniar Kurniawan, and Bo Fu for their hard working; thank to other UCARE students Shiqin Yan, Michael Tong, and Huaicheng Li to make UCARE group lively; and thank to all my friends, department faculty and staff that helped me many things when I was working on the dissertation.

And also I want to give big thanks to my family. The first most important one is my mother, the woman who always supports me throughout my life, without her, there would not be this Tanakorn. The second one is Louise, my wife; she always helps and supports me during my hard time. The other three are my lovely sisters, Fon, Nuch, and Nid; they always make me feel good every time I talk with them. Lastly, I want to thank my father, a man who is my motivations for many things including this Ph.D.

ABSTRACT

In the era of cloud computing, users move their data and computation from local machines to cloud, thus the services are expected to be 24/7 dependable. Cloud services must be accessible anytime and anywhere, not lose or corrupt users data, and scale as user base continues to grow. Unfortunately, guaranteeing cloud services' dependability is challenging because these cloud services are backed by large sophisticated distributed systems such as scalable data stores, data-parallel frameworks, and cluster management systems. Such cloud-scale distributed systems remain difficult to get right because they need to address data races among nodes, complex failures in commodity hardware, tremendous user requests, and much more. Addressing these cloud-specific challenges makes the systems more complex and new intricate bugs continue to create dependability problems.

This dissertation tries to answer a vital question of cloud dependability: *“how can we make cloud-scale distributed systems more dependable?”* We try to answer this question by focusing on the problems of distributed concurrency bugs and scalability bugs. We focus on these two problems because they are novel issues that occur in cloud-scale environment only and not many works addressing them.

Distributed concurrency bug (DC bug) is one unsolved reliability problem in cloud systems. DC bugs are caused by non-deterministic order of distributed events such as message arrivals, machine crashes, and reboots. Cloud systems execute multiple complicated distributed protocols concurrently. The possible interleavings of the distributed events are beyond developer's anticipations and some interleavings might not be handled properly that can lead to catastrophic failures. To combat DC bugs, we make two contributions. First, we conduct a formal study on DC bugs to gain foundation knowledge for DC-bug combating research. We study 104 DC bugs from various widely-deployed cloud-scale distributed systems in many characteristics along several axes of analysis such as the triggering timing condition, input preconditions, error and failure symptoms, and fix strategies. We present the first complete taxonomy of DC bugs, TaxDC, along with many

findings on DC bugs that can guide future research.

Second, we advance state of the art of distributed system model checking by introducing “*semantic-aware model checking*” (SAMC). Distributed system model checkers (dmck) are used to test system reliability of real systems. Existing dmcks however rarely exercise multiple faults due to the state-space explosion problem, and thus do not address present reliability challenges of cloud systems in dealing with complex faults. SAMC pushes the boundary of dmcks by introducing a white-box principle that takes simple semantic information of the target system and incorporates that knowledge into state-space reduction policies. We show that SAMC can find deep bugs one to two orders of magnitude faster compared to state-of-the-art techniques.

And for the second aspect of system dependability, we focus on scalability bugs. Scale surpasses the limit of a single machine in meeting users’ increasing demands for computing and storage. On the negative side, scale creates new development and deployment issues. Developers must ensure that their algorithms and protocol designs to be scalable. However, until real deployment takes place, unexpected bugs in the actual implementations are unforeseen. This new era of cloud-scale distributed systems has given birth to “scalability bugs”, latent bugs that are scale-dependent, and only surface in large scale.

To address scalability bugs, we conduct a study on scalability bugs to understand how they manifest and what their root causes are, and introduce SCK, a methodology that enables developers to *scale-check* distributed systems and find scalability bugs economically on one machine. SCK helps developers identify potential buggy code and allows developers to colocate a large number of nodes to test the potential buggy code without sacrificing accuracy. We remove a problem of hardware contentions (*i.e.*, CPU, memory, and thread) with four novel strategies, and we successfully integrate SCK to Cassandra, Riak, and Voldemort. With SCK, we achieve a high colocation factor (500 nodes), and can reproduce six scalability bugs and identify two new hidden bugs.

CHAPTER 1

INTRODUCTION

“*Cloud computing*” has been given many definitions from many companies and experts [32, 33, 34, 70]. These definitions are different in details, but they have some common characteristics; they are on-demand internet-based services that can scale to fit increasing users, and users pay only for their use. Cloud computing help users (from end users to organizational users) reduce the capital investment in hardware that is mostly underutilized [89] and help business moves faster [124]. We see a trend that users are moving their data and computation from local machines and in-house datacenters to the cloud [2, 8, 19, 21, 26, 30, 31, 121].

This trend makes client-side software get thinner and more heavily rely on the cloud services, thus the services are expected to be 24/7 dependable. Cloud services must be accessible any-time and anywhere, not lose or corrupt users data, and scale as user base continues to grow [45]. Unfulfilled dependability is costly. Some researchers estimate that 568 hours of downtime at 13 well-known cloud services since 2007 to 2012 had an economic impact of more than \$70 million [66]. Others predict worse: for every hour it is not up and running, a cloud service can take a hit between \$1 to 5 million [112]. A study shows that service disruptions are widespread and still continue to occur [79].

Unfortunately, proving cloud services’ dependability is challenging. Behind cloud computing is backed by large sophisticated distributed software stack [44, 48, 49, 56, 58, 59, 73, 94, 105, 125, 162] that is running on top of large-scale cluster [3, 35, 36]. Such cloud distributed systems remain difficult to get right because they need to address data races among machines, complex failures that randomly happen, tremendous user requests, and much more issues that caused from cloud computing infrastructure.

Data races are known to be a core problem in any concurrent software systems. Unlike non-distributed software, cloud distributed systems are subject to not only local concurrency bugs, which basically come from thread interleaving, but also distributed concurrency bugs, which come

from inter-node message interleaving. Moreover, cloud hardware is built from commodity hardware such that failures can happen at anytime and can be very complex. The timing of these hardware failures plus message interleaving makes it hard to handle the concurrency correctly.

Moreover, the size of cloud users is tremendous and cloud service providers need to guarantee service quality (*i.e.*, availability and performance) to their users. The providers need to ensure that their capabilities can satisfy the current users and also make sure there is no glitch when users are growing. Cloud providers normally employ large-scale systems to achieve high aggregate capabilities, but large-scale systems are challenging to build and costly to test their correctness.

Addressing these challenges makes the systems getting more complex. New intricate bugs continue to happen and create dependability problems. Guaranteeing dependability has proven to be challenging in these systems [77, 81, 113, 151, 157]. This raises a vital question: “*how can we make cloud-scale distributed systems more dependable?*” We try to answer this question by focusing on the problems of distributed concurrency bugs and scalability bugs. These two are critical problems because they are novel issues that occur in cloud environment only and not many works addressing them. The following sections discuss our contributions to address the challenges.

1.1 Distributed Concurrency Bugs

Distributed concurrency bugs (DC bugs) are bugs that caused by nondeterministic orders of distributed events. Distributed events could be message arrivals, hardware crashes/reboots, network timeout, etc. Cloud systems execute multiple complicated distributed protocols concurrently (*e.g.*, serving users’ requests, operating background tasks, and combined with untimely hardware failures), and possible interleavings of the distributed events are beyond developers’ anticipation, which some interleavings might not be handled properly, and can cause catastrophic failures such as data loss/inconsistency and downtimes. Compared to the “countless” of efforts in combating “local” concurrency bugs in multi-threaded software, DC bugs have not received the same amount of attention within the research community.

Here are our contributions in combating DC bugs in systematic and comprehensive manners:

1. Bug study and taxonomy for DC bugs (TaxDC): We perform an in-depth study of more than 100 real-world DC bugs and build a first complete taxonomy of DC bugs. This study can give insight to guide many future research work on DC bugs.
2. Semantic-aware model checking (SAMC): We advance the state of the art of model checking for distributed systems by adopting white-box approach to tackle state-space explosion which is the current limitation of model checking.

The brief detail of these two works are discussed below.

1.1.1 Distributed Concurrency Bug Study & Taxonomy

Bug and failure studies can significantly guide many aspects of dependability research. Many researchers have recently employed formal studies on bugs and failures [96, 108, 109, 118, 137, 146, 155, 158]. However, we are not aware of any public large-scale DC-bug study, a recent study from Microsoft analyzed the effect of distributed concurrency of workload and only studied five DC bugs in MapReduce [155], and researchers from NEC Labs dissected only network-failure-related DC bugs to study and did not publicly release it [99].

In this dissertation, we fill the void by performing large-scale DC-bug study. We study 104 real-world DC bugs from four various popular cloud-scale distributed systems: Cassandra, HBase, Hadoop MapReduce/Yarn, and ZooKeeper. We study DC bugs in all aspects including trigger, errors and failures, and fixes.

For triggering conditions, we study DC bugs from two perspectives:

1. Timing conditions: For every DC bug, we identify the smallest set of concurrent events E , so that a specific ordering of E can guarantee the bug manifestation. This is similar to the interleaving condition for local concurrency bugs.

2. Input preconditions: In order for those events in E to happen, regardless of the ordering, certain inputs or fault conditions (*e.g.*, node crashes) must occur. This is similar to the input condition for local concurrency bugs.

Understanding the triggering can help the design of testing tools that can proactively trigger DC bugs, bug detection tools that can predict which bugs can be triggered through program analysis, and failure prevention tools that can sabotage the triggering conditions at run time.

Other than the trigger, we also look into errors and failures. From the triggering conditions, we then scrutinize the first error that happens immediately after. First errors are the pivotal point that bridges the triggering and error-propagation process. And after the first errors, we track down to system failures that are noticeable to users such as downtimes, lost/corrupted/inconsistent data, failed operations, and degraded performance. Identifying errors and failures help failure diagnosis get closer to disclosing bug triggering and root causes and help bug detection get closer to accurately predict failures.

Lastly, we study how developers fix DC bugs to understand their fix strategies. We want to see how different DC bug fixes compared to local concurrency bugs. Understanding the fix strategies will help research on runtime failure prevention and automatic bug fixing.

Our contribution from the study is the first complete taxonomy of DC bugs which named TaxDC. TaxDC contains in-depth characteristics of DC bugs, stored in the form of 2,083 classification labels and 4,528 lines of re-enumerated steps to the bugs that we manually added. And as mentioned above, TaxDC can guide various future research on combating DC bugs such as model checking, bug detections, failure diagnosis, and failure prevention and fixing.

1.1.2 *Semantic-Aware Model Checking*

One powerful method for discovering hidden DC bugs is the use of an *implementation-level (software) distributed system model checker* (dmck). A dmck can discover buggy interleavings that lead to DC bugs by reordering every possibility of nondeterministic distributed events. The last ten

years have seen a rise of dmcks such as MaceMC, MODIST, or Demeter. One big challenge faced by a dmck is the state-space explosion problem (*i.e.*, there are too many distributed events to re-order). To address this, existing dmcks adopt a random walk or basic reduction techniques such as dynamic partial order reduction (DPOR). Despite these early successes, existing approaches cannot unearth many real-world DC bugs, so we advance state of the art of dmck in this dissertation.

We start by addressing two limitations of existing dmcks. First, existing dmcks treat every target system as a complete *black box*, and perform unnecessary reorderings of distributed events that would lead to the same states (*i.e.*, redundant executions). Second, they do not incorporate complex multiple fault events (*e.g.*, crashes, reboots, etc.) into their exploration strategies, as such inclusion would exacerbate the state-space-explosion problem.

To address these limitations, we introduce *Semantic-Aware Model Checking* (SAMC), a novel white-box model checking approach that takes *semantic knowledge* of how distributed events (specifically, messages, crashes, and reboots) are processed by the target system and incorporates that to create reduction policies. The policies are based on sound reduction techniques such as DPOR and symmetry. The policies tell SAMC not to re-order some pairs of events such as message-message pairs, and message-crash pairs, yet preserves soundness, because those cut out re-orderings are redundant, and unnecessary to check.

SAMC can reproduce twelve old bugs in three cloud distributed systems (Cassandra, Hadoop MapReduce, and ZooKeeper) involving 30-120 distributed events and multiple crashes and reboots. Some of these bugs cannot be unearthed by non-SAMC approaches, even after two days. SAMC can find the bugs up to 340 (49x on average) faster compared to state-of-the-art techniques, it found two new bugs in Hadoop MapReduce and ZooKeeper.

1.2 Scalability Bugs

Scalability bug is a type of bug that newly born in the era of cloud computing. These bugs are latent such that they do not surface in small/medium-scale deployments, but only surface in large

scale. They threaten systems reliability and availability at scale. As we discussed above, cloud backend needs to be scalable; algorithms and protocols in cloud distributed systems are designed to be scalable. However, until real deployment takes place, if developers do not have a large cluster to test their actual implementations, unexpected bugs are unforeseen.

To unearth latent scalability bugs, we need an effective and economic approach to test the systems prior to deployments, but in order to do that, we need to understand the nature of scalability bugs first. Unfortunately, we are not aware of any study on scalability bugs at all, so in this dissertation, we perform a study of scalability bugs to gain some foundational knowledge about them. We study 41 bugs in seven systems including Cassandra, Couchbase, Hadoop MapReduce, HBase, HDFS, Riak, and Voldemort. And here is our brief observations from the study:

- Scalability bugs only appear at extreme scale (*e.g.*, hundreds node).
- Systems can be scalable in design, but not in practice.
- Scalability bugs could be implementation specific and hard to predict.
- Scalability bugs are caused from cascading impacts of “not independent” nodes.
- It is long and difficult to debug large-scale.
- Not all developers have large cluster to test the systems, especially in open-source project.

These observations accentuate the need for scale-checking distributed system *implementations* at *real scale*, not via simulation nor extrapolation. The challenge of large-scale emulation is resource contention problem that is nodes compete to consume resources (*e.g.*, CPU, memory, and threads) and make test outcome inaccurate. In this context, we start a pilot work, SCK, a large-scale emulation that allows developers to colocate hundreds nodes in one machines to test system scalability, yet still get accurate testing results. SCK contains four techniques to mitigate resource contention which we briefly describe below.

First, we introduce *processing illusion* (PIL), which replaces scale-dependent CPU-intensive computations with `sleep()` without changing the cluster behavior. The insight behind PIL is that the key to computation is not the intermediate results, but rather the execution time and eventual output. To make PIL feasible, we analyze the characteristics of functions that can take PIL. We employ pre-memoization and order determinism to record the output data and execution time of PIL-replaceable functions.

In addition to PIL, we introduce other colocation strategies that reduce unnecessary CPU and memory contentions, strategies such as *single process cluster* (SPC), which runs the whole cluster in a single process, *global even driven architecture* (GEDA), which replaces hundreds of threads in SPC with a few event-handler threads shared by all nodes, and *memory footprint reduction* (MFR), which removes high system-specific memory footprints in our target systems.

We created SCK tools for Cassandra [105], Riak [28], and Voldemort [27]. We scale-checked a total of 5 protocols; 3 Cassandra (bootstrap, scale-out, decommission), 1 Riak (bootstrap+rebalance), and 1 Voldemort (rebalance) protocols. To show the simplicity of developing SCK, we have migrated SCK to a total of 9 old and new releases (5 Cassandra, 2 Riak, and 2 Voldemort versions). Across these versions, we have colocated 500 nodes and reproduced 7 (old and new) scalability bugs (5 Cassandra, 1 Riak, and 1 Voldemort bugs).

In summary, our contributions are:

1. We present a method for scale-checking distributed systems and reproducing the scalability bugs within.
2. We uncover the reasons why existing distributed systems are not easily scale-checkable (*i.e.*, the colocation bottlenecks).
3. We show the generality of SCK by applying the concept to three real-world cloud-scale distributed systems.

Overall, we believe that scalability bugs are new-generation bugs to combat in modern cloud-

scale distributed systems and SCK is one of the pilot solutions in this new area of research.

1.3 Summary of Contributions and Outline

We summarize our contributions and present the outline for the rest of dissertation below.

- **Background and related work:** Chapter 2 discusses a background of cloud computing and cloud-scale distributed systems, and how concurrency and scalability affect system development. We also discuss related work in this chapter.
- **Distributed concurrency bug study and taxonomy:** Chapter 3 introduces TaxDC, a DC-bug study and taxonomy that provides insights how DC bugs are triggered, how they affect systems, and how they are fixed which can help future research on DC bugs including model checking, bug detection, formal verification, failure diagnosis, and failure prevention.
- **Semantic-aware model checking:** Chapter 4 presents SAMC, a white-box approach to model check distributed systems. The approach advances the state-of-the-art of leveraging semantic knowledge to prune out redundant executions to tackle state space explosion. SAMC introduces some principles of semantic awareness that can mitigate the problem and this chapter also shows SAMPRO, a prototype of how SAMC approach can be adopted to Cassandra, Hadoop, and ZooKeeper.
- **Scalability checking methodology:** Chapter 5 shows our observations on scalability bugs that accentuate the need for scale checking, and presents SCK, a methodology to check scalability of systems in an economical manners. SCK introduces an approach for developers to colocate multiple CPU-intensive nodes on one machine to test and debug scalability bugs. With SCK, developers can mitigate a resource contention problem and yield the similar test and debug outcome as deploying systems on real large scale.

- **Conclusion and future work:** Chapter 6 concludes this dissertation, and discuss research avenues to make research work in this dissertation more complete.

CHAPTER 2

BACKGROUND AND RELATED WORK

This chapter provides the background of this dissertation. We briefly discuss about cloud computing and cloud-scale distributed systems in Section 2.1. Section 2.2 and 2.3 discuss about distributed concurrency and scalability respectively. These are two aspects of distributed systems that could threaten system dependability. And in those two sections (Section 2.2 and Section 2.3), we also briefly discuss how system community address issues from distributed concurrency and scalability.

2.1 Cloud Computing and Cloud-Scale Distributed Systems

2.1.1 Cloud Computing

In the past decade, *cloud computing* has become widespread buzzword that IT people often talk about. However, there is few agreement on what it really means; many companies and experts have given many definitions [32, 33, 34, 70] that are different in details, but they share some common characteristics: cloud computing is on-demand internet-based services that can scale to serve growing users' requests, and users *pay as they go*.

Cloud computing can attract a number of users to move their data and computation from local machines and private datacenters to the cloud [2, 8, 19, 21, 26, 30, 31, 121]. It provides many benefits as shown below:

- Users can access their computing resources and data any time and any where. For example, Google Doc allows users to create/modify their documents on one machine and access them later on mobile devices. This advantage also enables new classes of applications, especially, mobile applications [17, 18, 20, 24, 29]
- Users can cut the cost of hardware investment, but can get immediate access to computing resources. Moreover, this will help improve hardware utilization because users do not need

to pay while they are not using the resources [124].

- Most importantly, cloud computing help users to scale their services in a convenient manner. When their workload is growing up, they can just purchase more computing power from their cloud service providers.

2.1.2 *Cloud-Scale Distributed Systems*

Behind the cloud computing, it is sophisticated distributed systems running on large-scale clusters. Unlike traditional computing, when the number of users or workload increase, we do not upgrade hardware specification, but add more machines to the systems. Thus, cloud-scale distributed systems are distributed systems that are designed to be scalable to handle tremendous users' requests. Moreover, the cluster is built from commodity machines which hardware failures are not optional and can be very complex; machines in the cluster can stop working at any time; switches and routers can fail, and make network partitions; disks can be corrupted but still serve I/O requests; hardware can be “*limping*” which functions in degraded mode [62, 147]. Cloud systems must be able to handle these complex failures [37, 76, 85]. We will discuss some types of cloud-scale distributed systems below as these systems are systems that we work on in this dissertation.

- *Distributed file systems*: these are file systems that stores files across machines in clusters (however, some systems cannot be mounted as other traditional file systems, and some users consider them as data store [9]). Distributed file systems do replication or erasure coding across machines in order to be fault tolerant, and to increase performance via data aggregation. Examples of these systems are GFS from Google [73], and HDFS from Apache [141] which is an open-source version of GFS.
- *Data-parallel framework*: these frameworks process *big data* by leveraging parallelism. It eases parallel computing by enabling users to increase computing power by just adding more machines without changing their programs. One of well-known data-parallel frameworks is

MapReduce from Google [58]. As its name states, MapReduce is a programming model that consist of *map* and *reduce* functions. Map functions process key-value pairs of data and generate intermediate key-value pairs, which reduce functions will process values of data with the same keys to generate final results. Apache also has open-source framework, which is similar to Google MapReduce, called Hadoop MapReduce [5].

- *NoSQL data stores*: These are data stores that are not relational database. Storing and accessing data do not have strict tabular relations like in relational database, and are not done by SQL query (some systems support SQL query but not fully), such as key-value stores, document stores, and object stores. Most NoSQL stores adopt the concept of “eventual consistency” to improve availability during network partition (favoring “C” in CAP theorem). An example of NoSQL is Dynamo from Amazon [59]. Dynamo is decentralized distributed key-value store. Its open-source counterpart from Apache is Cassandra [105].
- *Synchronization services*: These are utility services that support other large distributed systems. They help nodes in other systems synchronize some metadata such as global locking, configuration maintaining, and naming. Examples of the synchronization services are Chubby from Google [44] and ZooKeeper from Apache [94].

2.2 Distributed Concurrency

2.2.1 Local Concurrency and Distributed Concurrency

A well-known concurrency that developers most familiar with is from thread execution interleaving in multi-threaded software or what we call in this dissertation “*local concurrency*”. Multi-threaded software has become common in the age of multi-core processor, however, building multi-threaded software is hard. Developers need to handle all possible interleaving of multiple threads that are accessing to same data properly, otherwise concurrency bugs will happen. These bugs are timing-related and non-deterministic, and they are extremely difficult to test and debug.

ZooKeeper Bug #1264:

1. *Follower F crashed* in the past,
2. *F reboots* and joins the cluster; then *F synchronizes data* with Leader L
3. F sends FOLLOWERINFO message to L [synchronization message]
4. L sends LEADERINFO message to F [synchronization message]
5. F sends ACKEPOCH message to L [synchronization message]
6. L sends SNAP message to F [synchronization message]
7. L sends data tree snapshot to F [synchronization message]
8. L sends NEWLEADER message to F [synchronization message]
9. *Client C sends a request* to update data with Tx-#15 to L; L does atomic broadcast to update all followers
10. L sends update proposal message for Tx-#15 to F [broadcast message]
11. F sends update ack message for Tx-#15 to L [broadcast message]
12. *L sends update commit message* for Tx-#15 to F [broadcast message]
13. *F applies the update* for Tx-#15 to in-memory data tree, but not to on-disk log (because F has not received UPTODATE message)
14. *L sends UPTODATE message* to F [synchronization message]
15. C sends a request to update data with Tx-#16 to L
16. L sends update proposal for Tx-#16 to F
17. F sends update ack for Tx-#16 to L
18. L sends update commit for Tx-#16 to F
19. F applies the update for Tx-#16 to in-memory data tree and on-disk log
20. *F crashes* (before *F does snapshot*)
21. F reboots and joins the cluster again
22. L synchronized data with F by sending update starting from Tx-#17
23. F loses the update for Tx-15 C did in step 9

Figure 2.1: **ZooKeeper-1264 bug.** A concurrency bug in ZooKeeper that is caused from a mix of untimely message arrivals and crash timing. This bug surfaces when a follower receives update commit message (step 12) in the middle of an atomic operation (step 3-14) and the follower crashes before it does snapshot (step 20)

For distributed systems, other than local concurrency, the systems are also subject to “*distributed concurrency*” that is caused from interleaving of computations in multiple nodes. Nodes in distributed systems do not have shared memory and they access data in other nodes via network communication, so concurrency in distributed systems can also come from concurrent message arrivals and internal computations in running nodes.

Other than timing of message arrivals and local computations, concurrent bugs in distributed systems are caused from timing of failures as well. Cloud systems are built on top of commodity hardware for horizontal scaling purpose (Section 2.3.1). This commodity hardware is unreliable, and hardware failures are not an option [37, 76, 85]. Cloud distributed systems need to response to these failures, they need to detect and recover from the failures and makes sure that users’ data will not be lost or corrupted. Guaranteeing this correctness is proven to be hard [61, 77], cloud-scale distributed systems need to handle failures that can happen at any time and at any state of the systems. Some ordering of message arrivals could make systems into state that developers never anticipate and is prone to error when failure happens.

Considering hardware failures, concurrency bugs in cloud distributed systems is not only about interleaving of message arrivals and local computations, but also timing of hardware failures as well. Figure 2.1 shows an example of a concurrency bug that happens because of untimely ordering of message arrivals and node crashes. This bugs surfaces only if a follower receives an UPTODATE message (step 14) after a commit message (step 12), and the follower crashes before it does snapshot (step 20); only untimely message arrivals or the timing of follower crash is not enough for bug to surface.

2.2.2 Distributed Systems Model Checker (DMCK)

In order to unearth DC bugs the question we have to answer is: “*can we exercise necessary conditions (i.e. workloads and faults) and test different event re-ordering to hit the bugs?*”. This is the job of distributed system model checkers (dmck), which are gaining popularity recently

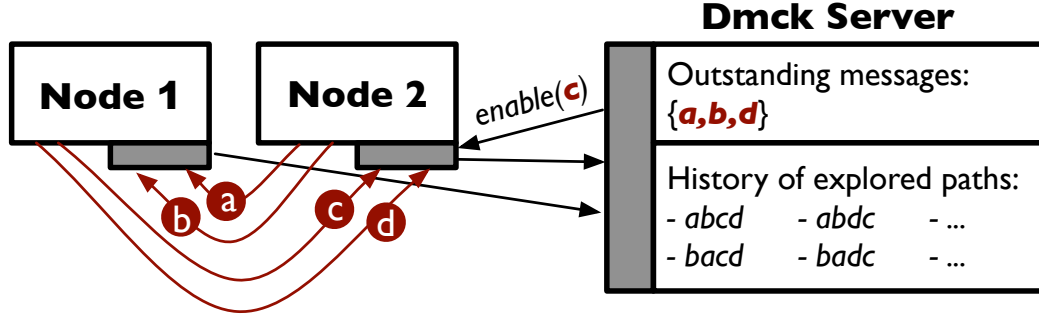


Figure 2.2: **DMCK**. A typical framework of a distributed system model checker (*dmck*).

[81, 103, 143, 157]. Dmck works by intercepting distributed events and permuting their ordering, and hereby pushing the target system into corner-case situations and unearthing hard-to-find bugs. However, the more events included, the more scalability issues will arise due to state-space explosion.

The last ten years have seen a rise of software model checker that checks distributed systems directly at the implementation level. Figure 2.2 illustrates a dmck integration to a target distributed system, a simple representation of existing dmck frameworks [81, 103, 143, 157]. The dmck inserts an interposition layer in each node of the target system with the purpose of controlling all important events (*e.g.*, network messages, timeouts) and preventing the target system to process the events until the dmck enables them. A main dmck mechanism is the permutation of events; the goal is to push the target system into all possible ordering scenarios. For example, the dmck can enforce *abcd* ordering in one execution, *bcad* in another, and so on.

2.2.3 Symbolic Execution

Symbolic execution is another powerful formal method to verify systems correctness. Symbolic execution also faces an explosion problem, specifically the path explosion problem. A huge body of work has successfully addressed the problem and made symbolic execution scale to large (non-distributed) software systems [41, 46, 50, 57, 163]. Symbolic execution and model checking can formally be combined into a more powerful method [42], however this concept has not permeated

the world of distributed systems; it is challenging to track symbolic values across distributed nodes.

2.2.4 Fault Injector

Reliability bugs are often caused by incorrect handling of failures [77, 80]. Fault-injection testing however is challenging due to the large number of possible failures to inject. This challenge led to the development of efficient fault-injection testing frameworks. For example, AFEX [38] and LFI [123] automatically prioritize “high-impact targets” (*e.g.*, unchecked system calls). These novel frameworks target non-distributed systems and thus the techniques are different than ours.

Similarly, recent work highlights the importance of testing faults in cloud systems (*e.g.*, FATE [77], SETSUDO [99], PREFAIL [100], and OpenStack fault-injector [101]). However, these frameworks are not a dmck; they cannot re-order concurrent messages and failures and therefore cannot catch distributed concurrency bugs systematically.

2.3 Scalability

2.3.1 Vertical Scaling vs Horizontal Scaling

When systems’ workload grows (the number of users raises or individual users’ requests increase), developers need to scale the systems to add more capability and keep users satisfied. Two traditional approaches to scale system are used as we show below [126]:

- **Vertical scaling or scale-up:** this approach expands system capabilities by adding more resources (*e.g.*, CPU, memory, and storage) to a single node to boost its performance, and make software to leverage additional resources. For example, run more processes of applications in the node.
- **Horizontal scaling or scale-out:** this approach enhance the capability by adding more nodes to current distributed systems to yield higher aggregate capability; mostly, the nodes that we are adding are low-cost machines.

In the past, vertical scaling was widely favored by many companies. Multiprocessor with higher clock rate can satisfy computing power need of largest companies [126]. Vertical scaling requires less human effort than horizontal scaling; it does not need more administrative effort because the number of machines and systems administrators need to handle is still the same. The disadvantages of vertical scaling are the upgradability is limited by existing hardware manufacturing, and the upgrade cost is expensive.

Because of the upgradability and price issues, nowadays, the trend goes to horizontal scaling. Many cloud service companies adopt this approach (*e.g.*, Google, Facebook, Amazon, etc.). The cost of horizontal scaling is much more cheaper than the vertical scaling and there is not limitation for hardware to scale out infinitely (the limitations are posed by software stack) [16]. In addition, hardware manufacturers try to facilitate scale-out approach [126].

2.3.2 Scalability Testing

As we discussed in 2.3.1, horizontal scaling or scale-out is a trend now because it does not expose hardware limitation. The limitation is in software stack so developers need to invent scalable algorithms and protocols, however, before real deployment, if they do not have a large cluster to test their implementations, there could be “*scalability bugs*” hide there.

We now discuss popular approaches (simulation, extrapolation, and emulation) for unearthing scalability bugs that avoid acquiring a number of machines, because testing on such deployments is costly. First, simulation approaches test system/application models in different scales [47, 104]. However, a model can look scalable but the actual implementation can contain unforeseen bugs. Later in Section 5.1.1, we will show our observations from *real-world* scalability bugs that simulation cannot detect them.

Second, extrapolation monitors system behaviors in “mini clusters” and extrapolates them to larger scales (Section 2.1 in [151]). However, mini clusters tend to be order(s) of magnitude smaller than real deployments. Most importantly, system behaviors do not always extrapolate

linearly [151].

Finally, real-scale emulation checks real implementations in an emulated environment [83, 151]. This approach emulates real large-scale system in a single machine. For example, a naive way to achieve this is just colocating multiple processes/VMs on one machine. The limitation here is emulation consumes real resources (CPU, memory, and storage), so with limited resources, we can not emulate *really large* deployment (*e.g.*, can test up to 50-node deployment). Some works try to resolve this resource contention issue [83, 151]. We will discuss them in Section 5.2

2.3.3 Scalability Benchmarking

Scalability benchmarking (*e.g.*, YCSB [55]) is a standard method to check throughput/latency scalability. The results are useful for advertising system capabilities, and thus acquiring a large number of machines is justifiable. But this is different from testing and debugging code to check system scalability.

2.4 Conclusion

In this chapter, we discuss about cloud-scale distributed systems, software backend for the cloud computing. We show what the current trend of the systems is and how they are design. We also discuss about distributed concurrency and scalability, the two important aspects of cloud-scale distributed systems that could threaten dependability of the systems. We briefly discuss how system community address issues from the concurrency and scalability. Unfortunately, distributed concurrency bugs are still an unsolved problem, and scalability bugs are novel and not many works address about them.

CHAPTER 3

TAXDC: A TAXONOMY OF NON-DETERMINISTIC CONCURRENCY BUGS IN CLOUD DISTRIBUTED SYSTEMS

Concurrency bugs are one notorious type of software bugs that happen in concurrency systems. These timing-related bugs manifest non-deterministically, and hence are extremely difficult to detect, diagnose, and fix. A huge body of work exists in this space that focuses on “local” concurrency bugs (LC bugs) in single-machine multi-threaded software, caused by incorrect interleaving of memory accesses. And for cloud-scale distributed systems, the reliability is also severely threatened by non-deterministic concurrency bugs as well, which we refer as *distributed concurrency bugs* (DC bugs). Distributed systems execute many complicated distributed protocols on hundreds/thousands of machines with no common clocks, and must face a variety of random hardware failures [62, 78]. This combination makes distributed systems prone to DC bugs caused by non-deterministic timing of distributed events such as message arrivals, node crashes, reboots, and timeouts. These DC bugs cannot be directly tackled by LC bug techniques, and they cause fatal implications such as operation failures, downtimes, data loss and inconsistencies.

Fighting DC bugs is challenging, particularly given the preliminary understanding of real-world DC bugs. To make progress, a comprehensive bug study is needed. Past studies have closely examined bugs in various software systems [51, 117, 131], which have motivated and guided many aspects of reliability research. There are few bug studies on cloud-scale distributed systems [78, 108], but they did not specifically dissect DC bugs. There was an internal bug study dissecting network-failure-related DC bugs to be a foundation to combat those bugs, but it was not published [99], and one recent work analyzed non-determinism in MapReduce programs but only discussed five bugs [155]. Thorough studies have also been conducted for LC bugs [69, 118] with many follow-up work to date, yet *there is no comprehensive study on real-world distributed concurrency bugs*.

In this chapter, we fill the void by presenting our in-depth analysis of real-world DC bugs in

well-known cloud distributed systems, and introducing TaxDC, the largest comprehensive taxonomy of DC bugs that covers several axes. We briefly give an overview of TaxDC in Section 3.1, and present our analysis in Section 3.3-3.7.

3.1 TaxDC

In this formal study on DC bugs, we do in-depth analysis of 104 DC bugs. The bugs came from four popular cloud distributed systems: Cassandra [4], HBase [6], Hadoop MapReduce [5], and ZooKeeper [7]. We introduce TaxDC, a comprehensive taxonomy of real-world DC bugs across several axes of analysis such as the triggering timing condition and input preconditions, error and failure symptoms, and fix strategies, as shown in detail in Table 3.1.

As the main contribution, TaxDC will be the first large-scale DC-bug benchmark. In the last six years, bug benchmarks for LC bugs have been released [95, 160], but no large-scale benchmarks exist for DC bugs. Researchers who want to evaluate the effectiveness of existing or new tools in combating DC bugs do not have a benchmark reference. TaxDC provides researchers with more than 100 thoroughly taxonomized DC bugs to choose from. Practitioners can also use TaxDC to check whether their systems have similar bugs. The DC bugs we studied are considerably general, representing bugs in popular types of distributed systems.

As a side contribution, TaxDC can help open up new research directions. In the past, the lack of understanding of real-world DC bugs has hindered researchers to innovate new ways to combat DC bugs. The state of the art focuses on three lines of research: monitoring and postmortem debugging [71, 115, 116, 136], testing and model checking [81, 103, 143, 157], and verifiable language frameworks [60, 154]. We hope our study will not only improve these lines of research, but also inspire new research in bug detection tool design, runtime prevention, and bug fixing, as elaborated more in Chapter 6.

3.2 Methodology

3.2.1 Basic Definitions

A *distributed concurrency bug* (DC bug) is a concurrency bug in distributed systems caused by distributed events that can occur in non-deterministic order. An *event* can be a message arrival/sending, local computation, fault, and reboot. A *local concurrency bug* (LC bug) is a concurrency bug that happens locally within a node due to thread interleaving. In our model, a *distributed system* is a collection of shared-nothing nodes. Each node can run multiple protocols in multiple threads.

3.2.2 Target Systems and Dataset

Our study examined bugs from four widely-deployed open-source datacenter distributed systems that represent a diverse set of system architectures: Hadoop MapReduce (including Yarn) [5] representing distributed computing frameworks, HBase [6] and Cassandra [4] representing distributed key-value stores (also known as NoSQL systems), and ZooKeeper [7] representing synchronization services. They are all fully complete systems containing many complex concurrent protocols. Throughout the chapter, we will present short examples of DC bugs in these systems. Some detailed examples are illustrated in Figure 3.2, 3.4 and 6.1.

The development projects of our target systems are all hosted under Apache Software Foundation wherein organized issue repositories (named “JIRA”) are maintained. To date, across the four systems, there are over 30,000 issues submitted. One major challenge is that issues pertaining to DC bugs do not always contain plain terms such as “concurrency”, “race”, “atomicity”, etc. Scanning all the issues is a daunting task. Thus, we started our study from an open source cloud bug study (CBS) database [1], which already labels issues related to concurrency bugs. However, beyond simple labeling, the CBS work did not differentiate DC from LC bugs and did not dissect DC bugs further.

From CBS, we first filtered out LC bugs, then exclude DC bugs that do not contain clear description, and finally randomly picked 104 samples from the remaining detailed DC bugs, specifically 19 Cassandra, 30 HBase, 36 Hadoop MapReduce, and 19 ZooKeeper DC bugs, reported in January 2011-2014 (the time range of CBS work). We have seen much fewer clearly explained DC bugs in CBS from Cassandra and ZooKeeper than those from HBase and Hadoop MapReduce, which may be related to the fact that they are different types of distributed systems. For example, ZooKeeper, as a synchronization service, is quite robust as it is built on the assumption of event asynchrony since day one. Cassandra was built on eventual consistency, and thus did not have many complex transactions, until recently when Cassandra adopts Paxos. We still see new DC bugs throughout 2014-2015 (some pointed to us by the developers); they can be included into TaxDC in the future.

3.2.3 Taxonomy

We study the characteristics of DC bugs along three key stages: triggering, errors & failures, and fixing (Table 3.1). *Triggering* is the process where software execution states deviate from correct to incorrect under specific conditions. At the end of this process, the manifestation of DC bugs changes from non-deterministic to deterministic. *Errors and failures* are internal and external software misbehaviors. *Fixing* shows how developers correct the bug. We will discuss in detail these categories in their respective sections.

3.2.4 Threats to Validity

For every bug, we first ensure that the developers marked it as a real bug (not a false positive). We also check that the bug description is clear. Finally, We then *re-enumerate* the full sequence of operations (the “*steps*”) to a clearer and more concise description such as the ones in Figure 3.2. Our study cannot and does not cover DC bugs not fixed by the developers. Even for fixed bugs, we do not cover those that are not described clearly in the bug repositories, a sacrifice we had to make

Triggering
<i>What is the triggering timing condition?</i>
Message arrives unexpectedly late/early
Message arrives unexpectedly in the middle
Fault (component failures) at an unexpected state
Reboot at an unexpected state
<i>What are the triggering inputs preconditions?</i>
Fault, reboot, timeout, background protocols, and others
<i>What is the triggering scope?</i>
<i>How many nodes/messages/protocols are involved?</i>
Errors & Failures
<i>What is the error symptom?</i>
Local memory exceptions
Local semantic error messages & exceptions
Local hang
Local silent errors (inconsistent local states)
Global missing messages
Global unexpected messages
Global silent errors (inconsistent global states)
<i>What is the failure symptom?</i>
Node downtimes, data loss/corruption, operation failures, slowdowns
Fixing
<i>What is the fix strategy?</i>
Fix Timing: add global synchronization
Fix Timing: add local synchronization
Fix Handling: retry message handling at a later time
Fix Handling: ignore a message
Fix Handling: accepting a message without new computation logics
Fix Handling: others

Table 3.1: **Taxonomy of DC Bugs.**

to maintain the accuracy of our results.

Readers should be cautioned not to generalize the statistics we report as each distributed system has unique purpose, design and implementation. For example, we observe 2:1 overall ratio between order and atomicity violations (Section 3.3.1), however the individual ratios are different across the four systems (*e.g.* 1:2 in ZooKeeper and 6:1 in MapReduce). Like all empirical studies, our findings have to be interpreted with our methodology in mind.

3.2.5 TaxDC Database

We name the product of our study TaxDC database. TaxDC contains in total 2,083 classification labels and 4,528 lines of clear and concise re-description of the bugs (our version, that we manually wrote) including the re-enumeration of the steps, triggering conditions, errors and fixes. We release TaxDC to the public¹. We believe TaxDC will be a rich “bug benchmark” for researchers who want to tackle distributed concurrency problems. They will have sample bugs to begin with, advance their work, and do not have to repeat our multi-people-year effort.

3.2.6 Detailed Terminologies

Below are the detailed terminologies we use in this chapter. We use the term “state” to interchangeably imply *local state* (both in-memory and on-disk per-node state) or *global state* (a collection of local states and outstanding messages). A *protocol* (e.g., read, write, load balancing) creates a chain of events that modify system state. User-facing protocols are referred as *foreground* protocols while those generated by daemons or operators are referred as *background* protocols.

We consider four types of *events*: message, local computation, fault and reboot. The term *fault* represents component failures such as crashes, timeouts, and disk errors. A *timeout* (system-specific) implies a network disconnection or busy peer node. A *crash* usually implies the node experiences a power failure. A *reboot* means the node comes back up.

Throughout the chapter, we present bug examples by abstracting system-specific names. As shown in Figure 3.1, we use capital letters for nodes (e.g., A, B), two small letters for a message between two nodes (*ab* is from A to B). Occasionally, we attach system-specific information in the subscript (e.g., A_{AppMaster} sends *ab_{taskKill}* message to B_{NodeManager}). We use “/” to imply concurrency (*ac/bc* implies the two messages can arrive at C in different orders, *ac* or *bc* first). A dash, “–”, means causal relation of two events (*ab-bc* means *ab* causally precedes *bc*). Finally, we use “N*” to represent crash, “N!” reboot, and “N+” local computation at N.

1. <http://ucare.cs.uchicago.edu/project/taxDC>

We cite bug examples with clickable hyperlinks (*e.g.*, [MR-3274](#)). To keep most examples uniform, we use MapReduce examples whenever possible. We use the following abbreviations for system names: “c/CA” for Cassandra, “h/HB” for HBase, “m/MR” for Hadoop MapReduce, and “z/ZK” for ZooKeeper; and for system-specific components: “AM” for application master, “RM” for resource manager, “NM” for node manager, “RS” for region server, and “ZAB” for ZooKeeper atomic broadcast.

3.3 Trigger

DC bugs often have a long triggering process, with many local and global events involved. To better reason about this complicated process, we study them from two perspectives:

1. *Timing conditions* (Section [3.3.1](#)): For every DC bug, we identify the smallest set of concurrent events E , so that a specific ordering of E can guarantee the bug manifestation. This is similar to the interleaving condition for LC bugs.
2. *Input preconditions* (Section [3.3.2](#)): In order for those events in E to happen, regardless of the ordering, certain inputs or fault conditions (*e.g.*, node crashes) must occur. This is similar to the input condition for LC bugs.

Understanding the triggering can help the design of testing tools that can proactively trigger DC bugs, bug detection tools that can predict which bugs can be triggered through program analysis, and failure prevention tools that can sabotage the triggering conditions at run time.

3.3.1 Timing Conditions (TC)

Most DC bugs are triggered either by untimely delivery of messages, referred to as *message timing bugs*, or by untimely faults or reboots, referred to as *fault timing bugs*. Rarely DC bugs are triggered by both untimely messages and untimely faults, referred to as *message-fault bugs*. Table [3.2](#) shows

	Ordering	Atomicity	Fault	Reboot
CA	4	4	6	5
HB	13	9	8	1
MR	25	4	5	3
ZK	4	8	7	5
All	46	25	26	14

Table 3.2: **#DC bugs triggered by timing conditions.** *The total is more than 104 because some bugs require more than one triggering condition. More specifically, 46 bugs (44%) are caused only by ordering violations, 21 bugs (20%) only by atomicity violations, and 4 bugs (4%) by multiple timing conditions (as also shown in Figure 3.3a).*

the per-system breakdown and Figure 3.3a (TC) the overall breakdown. Since a few bugs are triggered by more than one type of timing conditions (Section 3.3.3), the sum of numbers in Table 3.2 is slightly larger than the total number of DC bugs.

Message Timing Bugs. The timing conditions can be abstracted to two categories:

- a. *Order violation* (44% in Table 3.2) means a DC bug manifests whenever a message comes earlier (later) than another event, which is another message or a local computation, but not when the message comes later (earlier).
- b. *Atomicity violation* (20% in Table 3.2) means a DC bug manifests whenever a message comes in the middle of a set of events, which is a local computation or global communication, but not when the message comes either before or after the events.

LC and DC bugs are similar in that their timing conditions can both be abstracted into the above two types. However, the subjects in these conditions are different: shared memory accesses in LC and message deliveries in DC. The ratio between order violation and atomicity violation bugs are also different: previous study of LC bugs showed that atomicity violations are much more common than order violations in practice [118]; our study of DC bugs shows that this relationship does not apply or even gets reversed in several representative distributed systems.

An order violation can originate from a race between two messages (*message-message race*) at one node. The race can happen between two message arrivals. For example, Figure 3.1a illustrates *ac/bc* race at node C in MR-3274. Specifically, B_{RM} sends to C_{NM} a task-init message (*bc_{init}*), and soon afterwards, A_{AM} sends to C_{NM} a task-kill preemption message (*ac_{kill}*), however *ac_{kill}* arrives *before* *bc_{init}* and thus is incorrectly ignored by C. The bug would not manifest if *ac_{kill}* arrives *after* *bc_{init}* (Figure 3.1b). Message-message race can also happen between a message arrival and a message sending. For example, the *ab/bc* race in Figure 3.1c depicts HB-5780. In this bug, B_{RS} sends to C_{Master} a cluster-join request (*bc_{join}*) unexpectedly *before* a security-key message (*ab_{key}*) from A_{ZK} arrives at B, causing the initialization to abort.

Interestingly, message-message race can also occur concurrently across two nodes. For example, Figure 3.1d illustrates *ab/ba* race crisscrossing two nodes A and B in MR-5358. Specifically, A_{AM} sends *ab_{kill}* to a backup speculative task at B_{NM} because the job has completed, but concurrently the backup task at B sends *ba_{complete}* to A, creating a double-complete exception at A. If *ab_{kill}* arrives early at B, *ba* will not exist and the bug will not manifest (Figure 3.1e).

An order violation can also originate from a race between a message and a local computation (*message-compute race*). For example, Figure 3.1f illustrates *ab/b+* race in MR-4157. First, B_{AM} was informed that a task has finished and B plans to close the job and remove its local temporary files (*b+*). However, just *before* *b+*, A_{RM} sends to B a kill message (*ab*) and hence the files are never removed, eventually creating space issues. To prevent the failure, the kill message has to arrive after the local cleanup (Figure 3.1g).

An atomicity violation, as defined above, originates when a message arrives in the middle of a supposedly-atomic local computation or global communication. For example, Figure 3.1h illustrates MR-5009. When B_{NM} is in the middle of a commit transaction, transferring task output data (*bc*) to C_{HDFS}, A_{RM} sends a kill preemption message (*ab*) to B, preempting the task without resetting commit states on C. The system is never able to finish the commit — when B later reruns the task and tries to commit to C (*bc'*), C throws a double-commit exception. This failure would

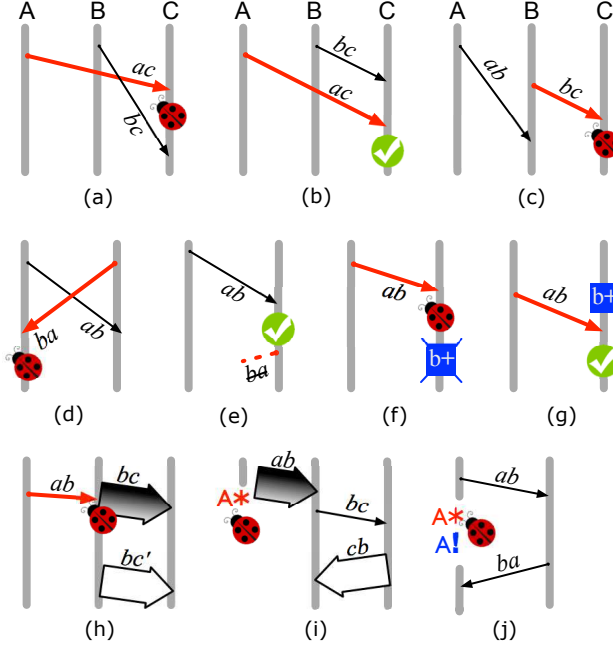


Figure 3.1: **Triggering patterns (Section 3.3.1).** The three vertical lines represent the timeline of nodes A, B and C. An arrow with xy label implies a message from X to Y. A square box with label $x+$ implies a local state-modifying computation at node X. A thick arrow implies a set of messages performing an atomic operation. X^* and $X!$ implies a crash and reboot at node X respectively (Section 3.2.6). All figures are discussed in Section 3.3.1

not happen if the kill message (ab) comes before or after the commit transaction (bc).

Fault and Reboot Timing Bugs. Fault and reboot timing bugs (32% in Table 3.2) manifest when faults and/or reboots occur at specific global states S_i ; the bugs do not manifest if the faults and reboots happen at different global states S_j .

Figure 3.1i illustrates a fault-timing bug in MR-3858. Here, A_{NM1} is sending a task's output to B_{AM} (ab) but A crashes in the middle (A^*) leaving the output half-sent. The system is then unable to recover from this untimely crash — B detects the fault and reruns the task at C_{NM2} (via bc) and later when C re-sends the output (cb), B throws an exception. This bug would not manifest, if the crash (A^*) happens before/after the output transfer (ab).

Figure 3.1j depicts a reboot-timing bug in MR-3186. Here, A_{RM} sends a job (ab) to B_{AM} and while B is executing the job, A crashes and reboots (A^* , $A!$) losing all its in-memory job

description. Later, B sends a job-commit message (*ba*) but A throws an exception because A does not have the job information. The bug would not manifest if A reboots later: if A is still down when B sends *ba_{commit}* message, B will realize the crash and cancel the job before A reboots and A will repeat the entire job assignment correctly.

Message-Fault Bugs. Four DC bugs are caused by a combination of messages and faults. For example, in Figure 3.2, a message (step 12) arrives in the middle of some atomic operation (step 3-14). This message atomicity violation leads to an error that further requires a fault timing (step 20) to become an externally visible failure.

Finding #1: DC bugs are triggered mostly by *untimely messages* (64% in Table 3.2) and sometimes by *untimely faults/reboots* (32%), and occasionally by a *combination* of both (4%). Among untimely messages, two thirds commit order violations due to message-message or message-computation race on the node they arrive; the others commit atomicity violations.

3.3.2 Input Preconditions (IP)

The previous section presents simple timing conditions that can be understood in few simple steps. In practice, many of the conditions happen “deep” in system execution. In other words, the triggering path is caused by complex input preconditions (IP) such as faults, reboots, and multiple protocols. Let’s use the same example in Figure 3.2. First, a fault and a reboot (step 1-2) and a client request (step 9) must happen to create a path to the message atomicity violation (step 9 interfering with step 3-14). Second, conflicting messages from two different protocols (ZAB and NodeJoin initiated in step 2 and 9) have to follow specific bug-triggering timing conditions. Even after the atomicity violation (after step 14), the bug is not guaranteed to lead to any error yet (*i.e.*, a benign race). Finally, the follower experiences an untimely fault (step 20), such that after it re-

ZK-1264:

1. *Follower F crashed* in the past,
2. *F reboots* and joins the cluster; then *F synchronizes data* with Leader L
3. F sends FOLLOWERINFO message to L [synchronization message]
4. L sends LEADERINFO message to F [synchronization message]
5. F sends ACKEPOCH message to L [synchronization message]
6. L sends SNAP message to F [synchronization message]
7. L sends data tree snapshot to F [synchronization message]
8. L sends NEWLEADER message to F [synchronization message]
9. *Client C sends a request* to update data with Tx-#15 to L; L does atomic broadcast to update all followers
10. L sends update proposal message for Tx-#15 to F [broadcast message]
11. F sends update ack message for Tx-#15 to L [broadcast message]
12. *L sends update commit message* for Tx-#15 to F [broadcast message]
13. *F applies the update* for Tx-#15 to in-memory data tree, but not to on-disk log (because F has not received UPTODATE message)
14. *L sends UPTODATE message* to F [synchronization message]
15. C sends a request to update data with Tx-#16 to L
16. L sends update proposal for Tx-#16 to F
17. F sends update ack for Tx-#16 to L
18. L sends update commit for Tx-#16 to F
19. F applies the update for Tx-#16 to in-memory data tree and on-disk log
20. *F crashes* (before *F does snapshot*)
21. F reboots and joins the cluster again
22. L synchronized data with F by sending update starting from Tx-#17
23. F loses the update for Tx-15 C did in step 9

Figure 3.2: **A DC bug in ZooKeeper.** This figure shows Figure 2.1 again. It shows a DC bug in ZooKeeper that is caused from a mix of untimely message arrivals and crash timing. This bug surfaces when a follower receives update commit message (step 12) in the middle of an atomic operation (step 3-14) and the follower crashes before it does snapshot (step 20)

boots (step 21), a global replica-inconsistency error will happen (step 23). Put it in a reverse way, before step 20, the global state is S_i and $S_i + \text{crash} \rightarrow \text{error}$, and the only way for the system to reach S_i is from complex preconditions such as a fault, a reboot, and some foreground and background protocols.

Statistically, Figure 3.3b (FLT) shows that 63% of DC bugs must have at least one fault. In more detail, Figure 3.3c-e (TO, CR, RB) shows the percentage of issues that require timeouts, crashes and reboots respectively, including how many instances of such faults must be there; the rest is other faults such as disk errors (not shown).

Figure 3.3f (PR) shows how many “protocol initiations” mentioned in the bug description. For example, if the system needs to perform one execution of background protocol and also three concurrent calls to the `write` protocol, then we label it with four protocol initiations. Up to 3 protocol initiations covers three quarters of DC bugs. When we count the number of *unique* protocols involved in all the bugs we study, we record 10 Cassandra, 13 HBase, 10 MapReduce, 6 ZooKeeper unique protocols, or 39 protocols in total. This again highlights the complexity of fully complete systems. Figure 3.3g (B/F) shows our categorization of protocols that are concurrently running into foreground only, background only, and foreground-background (mix) categories. More than three quarters of the bugs involve some background protocols and about a quarter involves a mix of foreground and background protocols.

Finding #2: Many DC bugs need *complex input preconditions*, such as faults (63% in Figure 3.3b), multiple protocols (80% in Figure 3.3f), and background protocols (81% in Figure 3.3g) .

3.3.3 Triggering Scope (TS)

We now analyze the triggering scope (TS), which is a complexity measure of DC-bug timing conditions. We use four metrics to measure the scope: message count (TS-MSG), node (TS-ND),

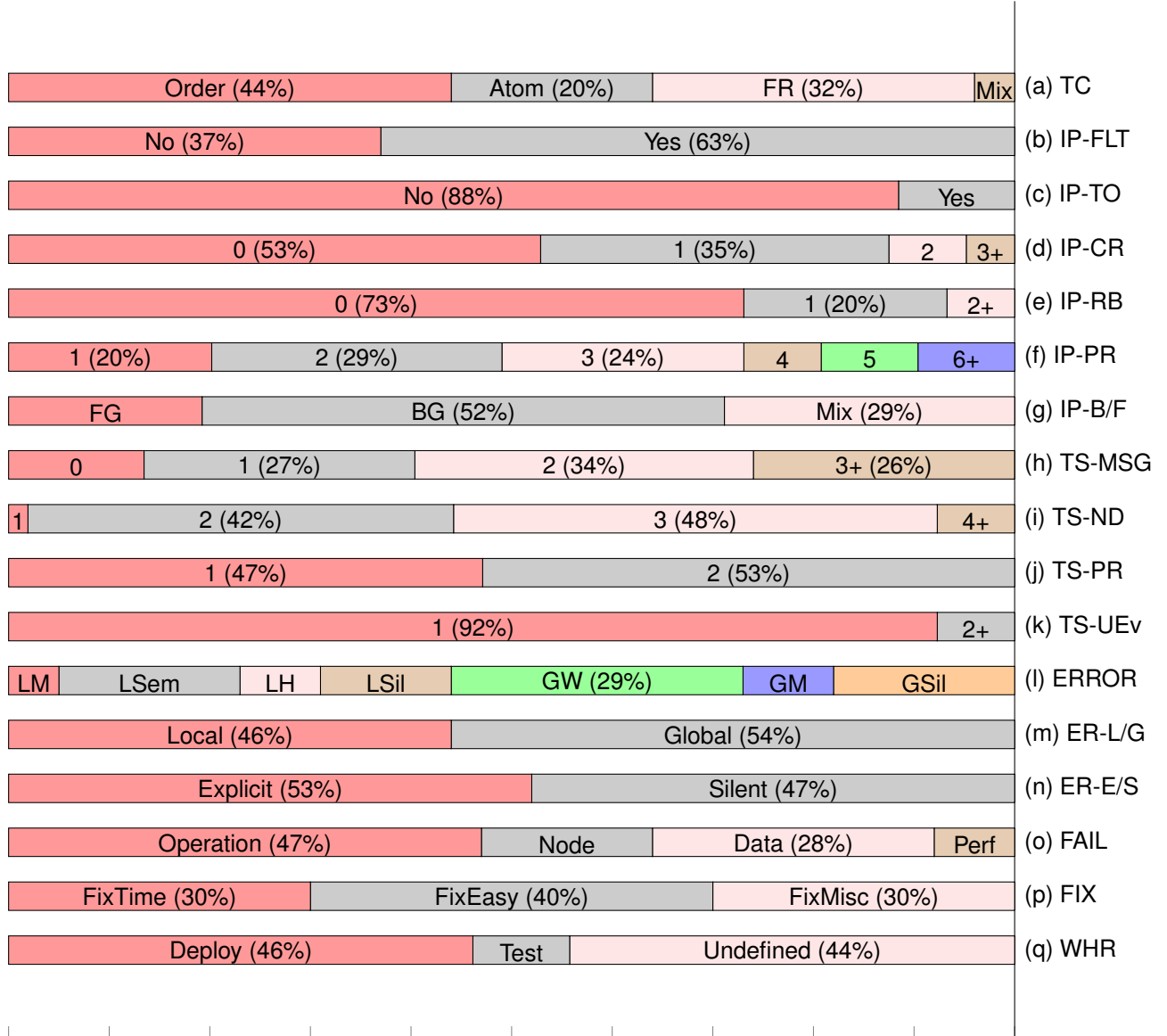


Figure 3.3: **Statistical overview of TaxDC.** *Timing Conditions (TC)* is discussed in Section 3.3.1, *Input Preconditions (IP)* in Section 3.3.2, *Triggering Scope (TS)* in Section 3.3.3, *Errors (ER)* in Section 3.4.1, *Failures (FAIL)* in Section 3.4.2, *Fixes (FIX)* in Section 3.5, and *Where Found (WHR)* in Section 3.7.

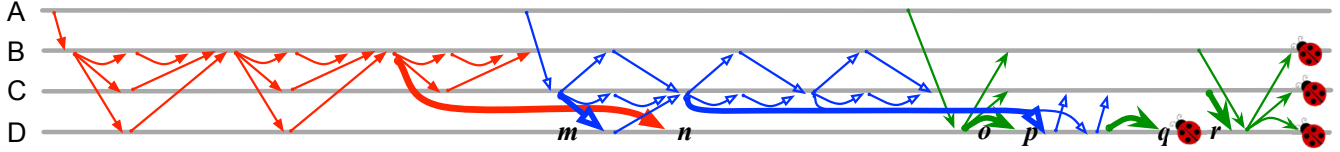


Figure 3.4: **A Cassandra's Paxos bug.** In *CA-6023*, three key-value updates (different arrow types) concurrently execute the Paxos protocol on four nodes (we simplify from the actual six nodes). The bug requires three message-message race conditions: (1) *m* arrives before *n*, (2) *o* before *p*, and (3) *q* before *r*, which collectively makes *D* corrupt the data and propagate the corruption to all replicas after the last broadcast. Note that the bug would not surface if any of the conditions did not happen. It took us one full day to study this bug.

protocol (TS-PR), and untimely event (TS-UEv) counts as shown in Figure 3.3h-k. This statistic is important with respect to the scalability of model checking, bug detection and failure diagnostic tools (Section 6.2.1-6.2.1).

Message count implies the minimum number of messages involved in *E* as defined in the beginning of section 3.3. Figure 3.3h (TS-MSG) shows that one or two triggering messages are the most common, with 7 messages as the maximum. Informally, zero implies fault timing bugs without any message-related races, one implies message-compute race, two implies message-message as in Figure 3.1a, and three implies a scenario such as *ac/(ab-bc)* race where *ab* and *ac* are concurrent or non-blocking message sending operations.

The node and protocol scopes present how many nodes and protocols are involved within the message scope. Figure 3.3i-j (TS-ND and TS-PR) shows that the scale of node and protocol triggering scope is also small, mostly two or three nodes and one or two protocols.

The untimely events count implies the total number of order violations, atomicity violations, untimely faults and reboots in the triggering timing condition of a bug. Figure 3.3k (TS-UEv) shows that only eight bugs require more than one untimely events. Four of them are message-fault bugs, each requiring one untimely message and one untimely fault to trigger (e.g., step 9 and 20 in Figure 3.2). Three are fault-reboot timing bugs, each requiring one untimely fault and one untimely reboot. The last one is *CA-6023*, shown in Figure 3.4, requiring three message-message order violations to happen.

Finding #3: The *timing conditions* of most DC bugs only involve *one to three* messages, nodes, and protocols (>90% in Figure 3.3h-j). Most DC bugs are mostly triggered by only *one* untimely event (92% in Figure 3.3k).

3.4 Errors and Failures

3.4.1 Error Symptoms

From the triggering conditions, we then scrutinize the *first* error that happens immediately after. First errors are the pivotal point that bridges the triggering and error-propagation process. Identifying first errors help failure diagnosis get closer to disclosing bug triggering and root causes and help bug detection get closer to accurately predict failures.

We categorize first errors into *local* errors and *global* errors, based on whether they can be observed from the triggering node N_T alone. Here, N_T is the node where triggering ends. It is the receiver node of untimely messages (e.g., node C in Figure 3.1a) or the node with untimely fault (e.g., node A in Figure 3.1i). For each error, we also check whether it is an *explicit* or *silent* error. Table 3.3 and Figure 3.3l (ERR) show the per-system and overall breakdowns respectively. Some MapReduce bugs caused multiple concurrent first errors of different types.

First, DC bugs can manifest into both local explicit errors and local silent errors. The former includes *memory exceptions* such as null-pointer exceptions (5% in Table 3.3) and *semantic errors* such as wrong state-machine transition exceptions thrown by the software (19%). Local silent errors include *hangs*, such as forever waiting for certain states to change or certain messages to arrive which are typically observed implicitly by users (9%), and *local silent* state corruption, such as half-cleaned temporary files (13%).

When local error is non-obvious in N_T , we analyze if the error is observable in other nodes communicating with N_T . Many DC bugs manifest into explicit global errors through *wrong mes-*

	Local Errors				Global Errors		
	Mem	Sem	Hang	Sil	Wrong	Miss	Sil
CA	2	0	0	4	3	3	7
HB	1	2	1	2	15	3	6
MR	2	13	7	4	14	4	0
ZK	0	6	2	5	1	0	5
All	5	21	10	15	33	10	18

Table 3.3: **First error symptoms of DC bugs.** *Some bugs cause multiple concurrent first errors.*

sages (29% in Table 3.3). Specifically, the communicating node receives an incorrect message from N_T , and throws an exception during the message handling. However, a few DC bugs still lead to silent global errors. These include *missing messages*, where N_T never sends a reply that the communicating node is waiting for in the absence of timeout (9%), and *global silent* state corruption such as replica inconsistencies between N_T and the other nodes (16%).

Finding #4: *Local and global first errors are about equally common; 46% vs. 54% in Figure 3.3m (ER-L/G). About half of the DC bugs generate explicit first errors (53%), including local exceptions and global wrong messages, and the remaining DC bugs lead to silent errors (47%), as shown in Figure 3.3n (ER-E/S). Some of them immediately lead to hangs in the triggering node N_T (9%) or a node communicating with N_T (9%).*

3.4.2 Failure Symptoms

Figure 3.3o (FAIL) shows that errors from DC bugs will eventually lead to a wide range of fatal failures including node downtimes (17%), data loss/corruption/inconsistencies (28%), operation failures (47%), and performance degradation (8%). A node downtime happens when the node either crashes or hangs (*i.e.*, it may still be heartbeating). It happens to both master/leader nodes and worker/follower nodes in our study. Data-related failures and performance problems are an

artifact of incorrect state logic induced from DC bugs. For example, in HBase, concurrent region updates and log splittings can cause data loss. In Cassandra, some dead nodes are incorrectly listed as alive causing unnecessary data movement that degrades performance. Node downtimes and data-related failures could also cause some operations to fail. To avoid double counting, we consider a bug as causing operation failures only when it does not cause node downtimes or data-related failures.

3.5 Fixes

We next analyze bug patches to understand developers' fix strategies. In general, we find that DC bugs can be fixed by either disabling the triggering timing or changing the system's handling to that timing (*fix timing* vs. *fix handling*). The former prevents concurrency with extra synchronization and the latter allows concurrency by handling untimely events properly. Since message timing bugs are fixed quite differently from fault timing bugs, we separate them below.

3.5.1 Message Timing Bug Fixes

The left half of Table 3.4 shows that only one fifth of message timing bugs are fixed by disabling the triggering timing, through either global or local synchronization. Only a couple of bugs are fixed through extra *global synchronization*, mainly due to its complexity and communication cost. For example, to prevent a triggering pattern $b+/ab$ in Figure 3.1f, MR-5465's fix *adds* a monitor on A_{RM} to wait for ba_{done} message from B_{AM} after B finishes with its local computation ($b+$); the result is $b+-ba-ab$ global serialization. More often, the buggy timing is disabled through *local synchronization*, such as re-ordering message sending operations within a single node. For example, HB-5780's fix for ab/bc race in Figure 3.1c forces the sending of bc request at B to wait for the receipt of ab ; the result is $ab-bc$ local serialization at B.

The right half of Table 3.4 shows that fix handling is more popular. Fix handling fortunately can be simple; many fixes do *not* introduce brand-new computation logic into the system, which

	Fix Timing		Fix Handling			
	Glob	Loc	Ret	Ign	Acc	Misc
CA	0	0	0	1	3	4
HB	2	7	2	1	7	3
MR	2	8	2	7	8	3
ZK	0	4	0	3	0	1
All	4	19	4	12	18	11

Table 3.4: **Fix strategies for message timing bugs.** *Some bugs require more than one fix strategy.*

can be done in three ways. First, the fix can handle the untimely message by simply *retrying* it at a later time (as opposed to ignoring or accepting it incorrectly). For example, to handle *bc/ac* race in Figure 3.1a, [MR-3274](#) retries the unexpectedly-early *ac_{kill}* message at a later time, right after the to-be-killed task starts. Second, the fix can simply *ignore* the message (as opposed to accepting it incorrectly). For example, to handle *ab/ba* race in Figure 3.1d, [MR-5358](#) simply ignores the unexpectedly-late *ba_{finish}* message that arrives after *A_{AM}* sends an *ab_{kill}* message. Finally, the patch can simply *accept* the untimely message by *re-using* existing handlers (as opposed to ignoring it or throwing an error). For example, [MR-2995](#)’s fix changes the node *AM* to accept an unexpectedly-early expiration message using an existing handler that was originally designed to accept the same message at a later state of *AM*. [MR-5198](#)’s fix handles the atomicity violation by using an existing handler and simply cancels the atomicity violating local operation. The rest of the fix-handling cases require new computation logic to fix bugs.

3.5.2 Fault/Reboot Timing Bug Fixes

Table 3.5 summarizes fix strategies for fault/reboot timing bugs. Unlike message timing, only rare bugs can be fixed by controlling the triggering timing either globally or locally (*e.g.*, by controlling the timing of the fault recovery actions). A prime example is an HBase cluster-wide restart scenario ([HB-3596](#)). Here, as *A* shuts down earlier, *B* assumes responsibility of *A*’s regions (via a region-takeover recovery protocol), but soon *B* shuts down as well with the regions still locked in

	Fix Timing		Fix Handling			
	G	L	Detect		Recover	
			TO	Msg	Canc	Misc
CA	1	0	3	2	4	6
HB	0	1	3	1	6	1
MR	2	1	1	1	2	1
ZK	0	3	0	1	1	7
All	3	5	7	5	13	15

Table 3.5: **Fix strategies for fault/reboot timing bugs.** *Some bugs require more than one fix strategy.*

ZooKeeper and the takeover cannot be resumed after restart. The patch simply adds a delay before a node starts region takeover so that it will likely get forced down before the takeover starts.

For the majority of fault timing bugs, their patches conduct two tasks: (1) detect the local/global state inconsistency caused by the fault and (2) repair/recover the inconsistency. The former is accomplished through timeouts, additional message exchanges, or others (omitted from Table 3.5).

The latter can be achieved by simply canceling operations or adding new computation logic.

Finding #5: *A small number of fix strategies have fixed most DC bugs. A few DC bugs are fixed by *disabling* the triggering timing (30% in Figure 3.3p), occasionally through extra messages and mostly through local operation re-orderings. Most DC bugs are fixed by better handling the triggering timing, most of which do not introduce new computation logic — they *ignore* or *delay* messages, *re-use* existing handlers, and *cancel* computation (40%).*

3.6 Root Causes

It is difficult to know for sure why many DC-bug triggering conditions were not anticipated by the developers (*i.e.*, the root causes). In this section, we postulate some possible and common misbeliefs behind DC bugs.

- “*One hop is faster than two hop.*” Some bugs manifest under scenario $bc/(ba-ac)$, similar to Figure 3.1a. Developers may assume that bc (one hop) should arrive earlier than $ba-ac$ (two hops), but ac can arrive earlier and hit a DC bug.
- “*No hop is faster than one hop.*” Some bugs manifest under scenario $ba-(b+/ab)$, similar to Figure 3.1f. Developers may incorrectly expect $b+$ (local computation with no hop) to always finish before ab arrives (one hop).
- “*Atomic blocks cannot be broken.*” Developers might believe that “atomic” blocks (local or global transactions) can only be broken unintentionally by some faults such as crashes. However, we see a few cases where atomic blocks are broken inadvertently by the system itself, specifically via untimely arrival of kill/preemption messages in the middle of an atomic block. More often, the system does not record this interruption and thus unconsciously leaves state changes half way. Contrary, in fault-induced interruption, some fault recovery protocol typically will handle it.
- “*Interactions between multiple protocols seem to be safe.*” In common cases, multiple protocols rarely interact, and even when they do, non-deterministic DC bugs might not surface. This can be unwittingly treated as normally safe, but does not mean completely safe.
- “*Enough states are maintained.*” Untimely events can unexpectedly corrupt system states and when this happens the system does not have enough information to recollect what had happened in the past, as not all event history is logged. We observe that some fixes add new in-memory/on-disk state variables to handle untimely message and fault timings.

Finding #6: Many DC bugs are related with a few common misconceptions that are unique to distributed systems.

3.7 Other Statistics

We now present other quantitative findings not included in previous discussions. We attempted to measure the complexity of DC bugs using four metrics: (a) the number of “re-enumerated steps” as informally defined in Section 3.2.4, (b) the patch LOC including new test suites for the corresponding bug, (c) the time to resolve (TTR), and (d) the number of discussion comments between the bug submitter and developers. The 25th percentile, median, and 75th percentile for the four metrics are (a) 7, 9, and 11 steps, (b) 44, 172, and 776 LOC, (c) 4, 14, and 48 days to resolve, (d) 12, 18, and 33 comments.

In terms of where the bugs were found, Figure 3.3r (WHR) highlights that 46% were found in deployment and 10% from failed unit tests. The rest, 44%, are not defined (could be manually found or from deployment). Some DC bugs were reported from large-scale deployments such as executions of thousands of tasks on hundreds of machines.

3.8 LC bugs vs. DC bugs

We now compare LC bugs and DC bugs to identify their similarities and differences. There are clearly similarities between LC bugs and DC bugs, as, by definition, they are both timing-related non-deterministic bugs. Many DC bugs contain LC components: untimely messages may lead to unsynchronized accesses from multiple threads or multiple event-handlers [135, 92] in a single machine. It is probably not a surprise that atomicity violations and order violations are two dominant triggering timing conditions for both LC and DC bugs (Finding #1). Our observation of the small triggering scope of most DC bugs (Finding #3) is similar with that for LC bugs, which may

be related to the nature of the bug sets — more complicated bugs may be more difficult to fix, and hence less likely to be included in empirical studies.

There are also many differences between LC bugs and DC bugs, as they originate from different programming paradigms and execution environments. For example, order violations are much more common in DC bugs than those in LC bugs (Finding #1); faults and reboots are much more common in DC bugs than those in LC bugs (Finding #2); the diagnosis of many DC bugs will have to reason beyond one node, clearly different from that of LC bugs (Finding #4); the fix strategies for DC bugs are very different from those of LC bugs, because enforcing global synchronization is difficult (Finding #5).

3.9 Conclusion

In this chapter, we have shown our in-depth study result on DC bugs and introduce TaxDC, the largest and most comprehensive taxonomy of DC bugs. TaxDC categorizes DC bugs in three major aspects including triggering, error and failure, and fix. This chapter show some complexity of DC bugs such as:

- We see a lack of effective testing, verification, and analysis tools to detect DC bugs during development process.
- Many DC bugs hide in complex concurrent executions of *multiple* protocols and not only user-facing foreground protocols, but also background and operational.
- Majority of DC bugs surface in the presence of hardware faults such as machine crashes (and reboots), network delay and partition (timeouts), and disk errors.
- Half of DC bugs lead to silent failures and hence are hard to debug in production and reproduce offline.

However, our results also bring fresh and positive insights:

- From the triggering patterns, we see opportunities to build DC bug detection to focus on timing-specification inference and violation detection.
- Half of DC bugs lead to *explicit* local or global errors which allow inferring timing specifications based on local correctness specifications, in the form of error checking already provided by developers.
- Most DC bugs are fixed through a small set of strategies and some are simple which implies research opportunities for automated in-production fixing for DC bugs.

We believe that our observations in this chapter will be a foundation to help us advance the state of the art to combat DC bugs.

CHAPTER 4

SAMC: SEMANTIC-AWARE MODEL CHECKING FOR FAST DISCOVERY OF DC BUGS IN CLOUD DISTRIBUTED SYSTEMS

As we discuss in Section 2.2.2, a recent popular approach to unearth DC bugs is adopting distributed system model checker or *dmck*. However, due to the complexity of DC bugs that we discuss in Chapter 3 (*e.g.*, interactions of multiple protocols, and various and multiple faults), the state of the art of *dmcks* still cannot work effectively [81, 103, 143, 157]. We observe that the existing systematic reduction policies cannot find bugs quickly, and cannot scale with the inclusion of fault events.

In this chapter, we discuss how to advance the state of the art by leveraging semantic awareness to assist model checking, and introduce “*Semantic-Aware Model Checking*” (SAMC), a white-box model checking approach that takes semantic knowledge of how events (*e.g.*, messages, crashes, and reboots) are processed by the target system and incorporates that information in reduction policies. We first discuss background of *dmck* and the state of the art in Section 4.1, then we discuss about the concept of SAMC and its implementation in Section 4.2 and 4.3 respectively. We evaluate SAMC by comparing with the state of the art in Section 4.4, and lastly, we discuss other issues in adopting SAMC to check distributed systems in Section 4.5.

4.1 Background

This section gives more background on *dmck* and related terms, followed with a detailed overview of the state of the art.

4.1.1 DMCK Framework and Terms

As mentioned before, we define *dmck* as software model checker that checks distributed systems directly at the implementation level. Figure 4.1 illustrates a *dmck* integration to a target distributed

system, a simple representation of existing dmck frameworks [81, 103, 143, 157]. The dmck inserts an interposition layer in each node of the target system with the purpose of controlling all important events (*e.g.*, network messages, timeouts) and preventing the target system to process the events until the dmck enables them. A main dmck mechanism is the permutation of events; the goal is to push the target system into all possible ordering scenarios. For example, the dmck can enforce *abcd* ordering in one execution, *bcad* in another, and so on.

We now provide an overview of basic dmck terms we use in this dissertation and Figure 4.1. Each node of the target system has a *local state* (*ls*), containing many variables. An *abstract local state* (*als*) is a subset of the local state; dmck decides which *als* is important to check. The collection of all (and abstract) local states is the *global state* (*gs*) and the *abstract global state* (*ags*) respectively. The *network state* describes all the *outstanding messages* currently intercepted by dmck (*e.g.*, *abd*). To model check a specific protocol, dmck starts a *workload driver* (which restarts the whole system, runs specific workloads, etc.). Then, dmck generates many (typically hundreds/thousands) executions; an *execution* (or a *path*) is a specific ordering of events that dmck enables (*e.g.*, *abcd*, *dbca*) from an initial state to a termination point. A *sub-path* is a subset of a path/execution. An *event* is an action by the target system that is intercepted by dmck (*e.g.*, a network message) or an action that dmck can inject (*e.g.*, a crash/reboot). Dmck enables one event at a time (*e.g.*, *enable(c)*). To permute events, dmck runs *exploration methods* such as brute-force (*e.g.*, depth first search) or random. As events are permuted, the target system enters hard-to-reach states. Dmck continuously runs state *checks* (*e.g.*, safety checks) to verify the system's correctness. To reduce the state-space explosion problem, dmck can employ *reduction policies* (*e.g.*, DPOR or symmetry). A policy is *systematic* if it does not use randomness or bug-specific knowledge. In this chapter, we focus on advancing systematic reduction policies.

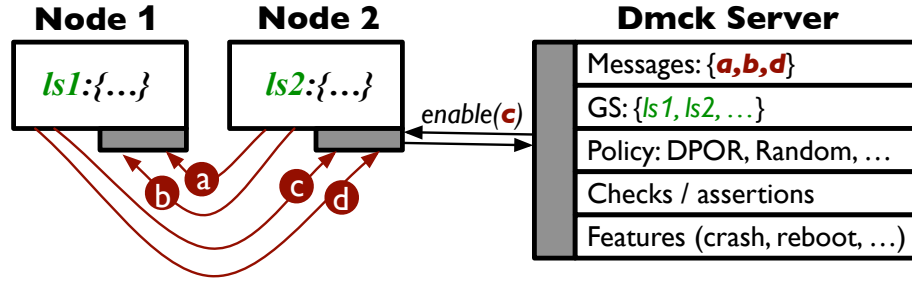


Figure 4.1: **How dmck works.** The figure illustrates a typical framework of a distributed system model checker (dmck).

4.1.2 State-of-the-Art DMCKs

MODIST [157] is arguably one of the most powerful dmcks that comes with systematic reduction policies. MODIST has been integrated to real systems due to its exploration scalability. At the heart of MODIST is *dynamic partial order reduction (DPOR)* [68] which exploits the *independence* of events to reduce the state explosion. Independent events mean that it does not matter in what order the system execute the events, as their different orderings are considered equivalent.

To illustrate how MODIST adopts DPOR, let's use the example in Figure 4.1, which shows four concurrent outstanding messages abcd (a and b for N1, c and d for N2). A brute-force approach will try all possible combinations (abcd, abdc, acbd, acdb, cabd, and so on), for a total of $4!$ executions. Fortunately, the notion of event independence can be mapped to distributed system properties. For example, MODIST specifies this reduction policy: a message to be processed by a given node is independent of other concurrent messages destined to other nodes (based on vector clocks). Applying this policy to the example in Figure 4.1 implies that a and b are dependent¹ but they are independent of c and d (and vice versa). Since only dependent events need to be reordered, this reduction policy leads to only 4 executions (ab-cd, ab-dc, ba-cd, ba-dc), giving a 6x speed-up ($4!/4$).

Although MODIST's speed-up is significant, we find that one scalability limitation of its DPOR

1. In model checking, “dependent” events mean that they must be re-ordered. “Dependent” does not mean “causally dependent”.

application is within its *black-box* approach; it only exploits general properties of distributed systems to define message independence. It does not exploit any semantic information from the target system to define more independent events. We will discuss this issue later (Section 4.2.1).

Dynamic interface reduction (DIR) [81] is the next advancement to MODIST. This work suggests that a complete dmck must re-order not only messages (global events) but also thread interleavings (local events). The reduction intuition behind DIR is that different thread interleavings often lead to the same global events (*e.g.*, a node sends the same messages regardless of how threads are interleaved in that node). DIR records local exploration and replays future incoming messages without the need for global exploration. In our work, SAMC focuses only on global exploration (message and failure re-orderings). We believe DIR is orthogonal to SAMC, similar to the way DIR is orthogonal to MODIST.

MODIST and DIR are examples of dmcks that employ advanced systematic reduction policies. LMC [75] is similar to DIR; it also decouples local and global exploration. dBug [143] applies DPOR similarly to MODIST. There are other dmcks such as MACEMC [103] and CrystalBall [156] that use basic exploration methods such as depth first (DFS), weight-based, and random searches.

Other than the aforementioned methods, *symmetry* is another foundational reduction policy [63, 145]. Symmetry-based methods exploit the architectural symmetry present in the target system. For example, in a ring of nodes, one can rotate the ring without affecting the behavior of the system. Symmetry is powerful, but we find no existing dmcks that adopt symmetry.

Besides dmcks, there exists sophisticated testing frameworks for distributed systems (*e.g.*, FATE [77], PREFAIL [100], SETSUDO [99], OpenStack fault-injector [101]). This set of work emphasizes the importance of multiple failures, but their major limitation is that they are not a dmck. That is, they cannot systematically control and permute non-deterministic choices such as message and failure reorderings.

4.1.3 Does State of the-Art Help?

We now combine our observations in the previous section and our insight from TaxDC (Chapter 3), and describe why state-of-the-art dmcks do not address present reliability challenges of cloud systems.

First, *existing systematic reduction policies often cannot find bugs quickly*. Experiences from previous dmck developments suggest that significant savings from sound reduction policies do not always imply high bug-finding effectiveness [81, 157]. To cover deep states and find bugs, many dmcks revert to non-systematic methods such as randomness or manual checkpoints. For example, MODIST combines DPOR with random walk to “jump” faster to a different area of the state space (Section 4.5 of [157]). DIR developers find new bugs by manually setting “interesting” checkpoints so that future state explorations happen from the checkpoints (Section 5.3 of [81]). In our work, although we use different target systems, we are able to reproduce the same experiences above (Section 4.4.1).

Second, *existing dmcks do not scale with the inclusion of failure events*. Given the first problem above, exercising multiple failures will just exacerbate the state-space explosion problem. Some frameworks that can explore multiple failures such as MACEMC [103] only do so in a random way; however, in our experience (Section 4.4.1), randomness many times cannot find deep bugs quickly. MODIST also enabled only one failure. In reality, multiple failures is a big reliability threat, and thus must be exercised.

We conclude that finding systematic (no random/checkpoint) policies that can find deep bugs is still an open dmck research problem. We believe without semantic knowledge of the target system, dmck hits a scalability wall (as also hinted by DIR authors; Section 8 of [81]). In addition, as crashes and reboots need to be exercised, we believe recovery semantics must be incorporated into reduction policies. All of these observations led us to SAMC, which we describe next.

4.2 SAMC

4.2.1 Semantic-Awareness Example

In a simple leader election protocol, every node broadcasts its vote to reach a quorum and elect a leader. Each node begins by voting for itself (*e.g.*, N2 broadcasts `vote=2`). Each node receives vote broadcasts from other peers and processes every vote with this simplified code segment below. As depicted in the code segment below, if an incoming vote is less than the node's current vote, it is simply discarded. If it is larger, the node changes its vote and broadcasts the new vote.

```
if (msg.vote < myVote) {  
    discard;  
} else {  
    myVote = msg.vote; broadcast(myVote);  
}
```

Let's assume N4 with `vote=4` is receiving three concurrent messages with votes 1, 2, and 3 from its peers. Here, a dmck with a black-box DPOR approach must perform 6 (3!) orderings (123, 132, and so on). This is because a black-box DPOR does *not* know the *message processing semantic* (*i.e.*, how messages will be processed by the receiving node). Thus, a black-box DPOR must treat all of them as dependent (Section 4.1.2); they must be re-ordered for soundness. However, by knowing the processing logic above, a dmck can soundly conclude that all orderings will lead to the same state; all messages will be discarded by N4 and its local state will not change. Thus, a semantic-aware dmck can reduce the 6 redundant executions to just 1 execution.

The example above shows a scalability limitation of a black-box dmck. Fortunately, simple semantic knowledge has a great potential in removing redundant executions. Furthermore, semantic knowledge can be incorporated on top of sound model checking foundations such as DPOR and symmetry, as we describe next.

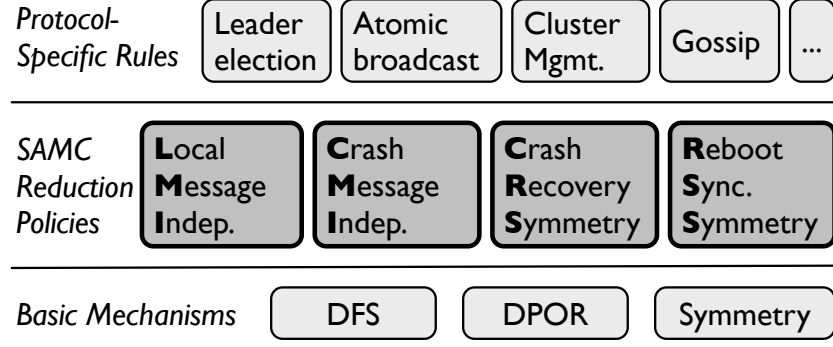


Figure 4.2: **SAMC Architecture.**

4.2.2 Architecture

Figure 4.2 depicts the three levels of SAMC: sound exploration mechanisms, reduction policies, and protocol-specific rules. SAMC is built on top of sound model checking foundations such as DPOR [68, 74] and symmetry [53, 145]. We name these foundations as mechanisms because a *dmck* must specify accordingly what events are dependent/independent and symmetrical, which in SAMC will be done by the reduction policies and protocol-specific rules.

Our main contribution lies within our four novel *semantic-aware reduction policies*: local-message independence (LMI), crash-message independence (CMI), crash recovery symmetry (CRS), and reboot synchronization symmetry (RSS). To the best of our knowledge, none of these approaches have been introduced in the literature. At the heart of these policies are *generic event processing patterns* (*i.e.*, patterns of how messages, crashes, and reboots are processed by distributed systems). Our policies and patterns are simple and powerful; they can be applied to many different distributed systems. Testers can extract the patterns from their target protocols (*e.g.*, leader election, atomic broadcast) and write protocol-specific rules in few lines of code.

In the next section, we first present our four reduction policies along with the processing patterns. Later, we will discuss ways to extract the patterns from target systems (Section 4.2.4) and then show the protocol-specific rules for our target systems (Section 4.3.2).

4.2.3 Semantic-Aware Reduction Policies

We now present four semantic-aware reduction policies that enable us to define fine-grained event dependency/independency and symmetry beyond what black-box approaches can do.

Local-Message Independence (LMI)

We define *local messages* as messages that are concurrently in flight to a given node (*i.e.*, intra-node messages). As shown in Figure 4.3a, a black-box DPOR treats the message processing semantic inside the node as a black box, and thus must declare the incoming messages as dependent, leading to $4!$ permutation of $abcd$. On the other hand, with white-box knowledge, local-message independence (LMI) can define *independency relationship among local messages*. For example, illustratively in Figure 4.3b, given the node's local state ($1s$) and the processing semantic (embedded in the `if` statement), LMI is able to define that a and b are dependent, c and d are dependent, but the two groups are independent, which then leads to only 4 re-orderings. This reduction illustration is similar to the one in Section 4.1.2, but this time LMI enables DPOR application on local messages.

LMI can be easily added to a `dmck`. A `dmck` server typically has a complete view of the local states (Section 4.1.1). What is needed is the *message processing semantic*: how will a node (N) process an incoming message (m) given the node's current local state ($1s$)? The answer lies in these four simple *message processing patterns* (discard, increment, constant, and modify):

<u>Discard:</u>	<u>Increment:</u>
if (pd(m,ls))	if (pi(m,ls))
(noop);	ls++;
<u>Constant:</u>	<u>Modify:</u>
if (pc(m,ls))	if (pm(m,ls))
ls = Const;	ls = modify(m,ls);

In practice, ls and m contain many fields. For simplicity, we treat them as integers. The functions with prefix p are boolean read-only functions (predicates) that compare an incoming message (m) with respect to the local state (ls); for example, pd can return true if $m < s$. The first pattern is a *discard* pattern where the message is simply discarded if pd is true. This pattern is prevalent in distributed systems with votes/versions; old votes/versions tend to be discarded (*e.g.*, our example in Section 4.2.1). The *increment* pattern performs an increment-by-one update if pi is true, which is also quite common in many protocols (*e.g.*, counting commit acknowledgments). The *constant* pattern changes the local state to a constant whenever pc is true. Finally, the *modify* pattern changes the local state whenever pm is true.

Based on these patterns, we can apply LMI in the following ways.

1. x is independent of y if pd is true on any of x and y . That is, if x (or y) will be discarded, then it does not need to be re-ordered with other messages.
2. x is independent of y if pi (or pc) is true on both x and y . That is, the re-orderings do not matter because the local state is monotonically increasing by one (or changed to the same constant).
3. x and y are dependent if pm is true on x and pd is not true on y . That is, since both messages modify the local state in unique ways, then the re-orderings can be “interesting” and hence should be exercised.

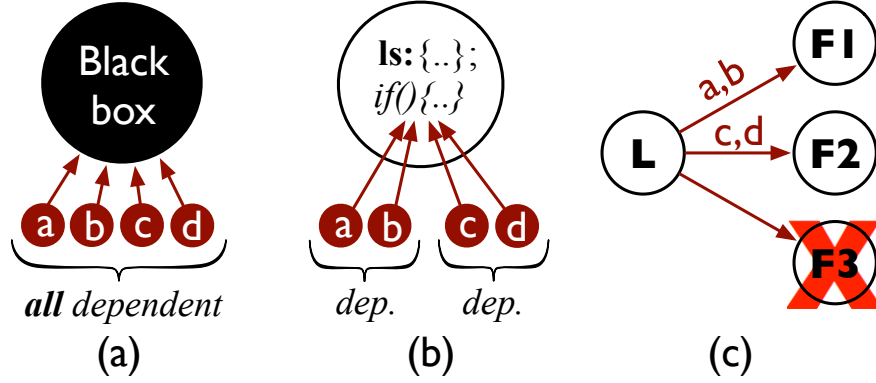


Figure 4.3: **LMI and CMI.** The figures illustrate (a) a black-box approach, (b) local-message independence with white-box knowledge, and (c) crash-message independence.

All these rules are continuously evaluated before every event is enabled. If multiple cases are true, dependency has higher precedence than independency.

Overall, LMI allows dmck to smartly skip redundant re-orderings by leveraging simple patterns. The job of the tester is to find the message processing patterns from a target protocol and write *protocol-specific rules* (i.e., filling in the content of the four LMI predicate functions (pd, pi, pc, and pm) specific to the target protocol). As an example, for our simple leader election protocol (Section 4.2.1), pd can be as simple as: `return m.vote < ls.myVote`.

Crash-Message Independence (CMI)

Figure 4.3c illustrates the motivation behind our next policy. The figure resembles an atomic broadcast protocol where a leader (L) sends commit messages to the followers (Fs). Let's assume commit messages ab to F1 and cd to F2 are still in flight (i.e., currently outstanding in the dmck; not shown). In addition, the dmck would like to crash F3, which we label as a crash event X. The question we raise is: how should X be re-ordered with respect to other outstanding messages (a, b, c, and d)?

As we mentioned before, we find *no* single dmck that incorporates crash semantics into reduction policies. As an implication, in our example, the dmck must reorder X with respect to other outstanding messages, generating executions Xabcd, aXbcd, abXcd, and so on. Worse, when abcd

are reordered, X will be reordered again. We find this as one major fundamental problem why existing dmcks do not scale with the inclusion of failures.

To solve this, we introduce crash-message independence (CMI) which defines *independency relationship between a to-be-injected crash and outstanding messages*. The key lies in these two crash reaction patterns (global vs. local impact) running on the surviving nodes (*e.g.*, the leader node in Figure 4.3c).

<u>Global impact:</u>	<u>Local impact:</u>
if (pg(X, ls))	if (pl(X, ls))
modify(ls);	modify(ls);
sendMsg();	

The functions with prefix p are predicate functions that compare the crash event X with respect to the surviving node's local state (*e.g.*, the leader's local state). The pg predicate in the *global-impact* pattern defines that the crash X during the local state ls will lead to a local state change *and* new outgoing messages (*e.g.*, to other surviving nodes). Here, no reduction can be done because the new crash-induced outgoing messages must be re-ordered with the current outstanding messages. On the other hand, reduction opportunities exist within the *local-impact* pattern, wherein the pl predicate specifies that the crash will just lead to a local state change but not new messages, which implies that the crash does not need to be re-ordered.

Based on the two crash impact patterns, we apply CMI in the following ways. Given a local state ls at node N , a peer failure X , and outstanding messages ($m_1 \dots m_n$) from N to other surviving peers, CMI performs:

1. If pl is true, then X and $m_1 \dots m_n$ are independent.
2. If pg is true, then X and $m_1 \dots m_n$ are dependent.

In Figure 4.3c for example, if pl is true in node L , then X does not impact outstanding messages to $F1$ and $F2$, and thus X is independent to $abcd$; exercising $Xabcd$ is sufficient.

An example of CMI application is a quorum-based write protocol. If a follower crash occurs and quorum is still established, the leader will just decrease the number of followers (local state change only). Here, for the protocol-specific rules, the tester can specify `p1` with `#follower >= majority` and `pg` with the reverse. Overall, CMI helps `dmck` scale with the inclusion of failures, specifically by skipping redundant re-orderings of crashes with respect to outstanding messages.

Crash Recovery Symmetry (CRS)

Before we discuss our next reduction policy, we emphasize again the difference between message event and crash/reboot event. Message events are generated by the target system, and thus `dmck` can only reduce the number of re-orderings (but it cannot reduce the events). Contrary, crash events are generated by `dmck`, and thus there exists opportunities to reduce the number of injected crashes. For example, in Figure 4.3c, in addition to crashing `F3`, the `dmck` can also crash `F1` and `F2` in different executions, but that might not be necessary.

To omit redundant crashes, we develop crash recovery symmetry (CRS). The intuition is that some crashes often lead to symmetrical recovery behaviors. For example, let's assume a 4-node system with node roles `FFFL`. At this state, crashing the first or second or third node perhaps lead to the same recovery since all of them are followers, and thereby injecting one follower crash could be enough. Further on, if the system enters a slightly different state, `FFLF`, crashing any of the followers might give the same result as above. However, crashing the leader in the two cases (`N4` in the first case and `N3` in the second) should perhaps be treated differently because the recovery might involve the dead leader ID. The goal of CRS is to help `dmck` with crash decision.

The main question in implementing CRS is: how to incorporate crash recovery semantics into `dmck`? Our solution is to compute *recovery abstract global state* (`rags`), a simple and concise representation of crash recovery. CRS builds `rags` with the following steps:

First, we define that two recovery actions are symmetrical if they produce the same messages and change the same local states in the same way.

```

broadcast()  sendMsgToAll(role, leaderId);
quorumOkay()  return (follower > nodes / 2);

// pr1
if (role == L && C.role == F && quorumOkay())
    follower--;

// pr2
if (role == L && C.role == F && !quorumOkay())
    follower = 0;
    role = S;
    broadcast();

// pr3
if (role == F && C.role == L)
    leaderId = myId;
    broadcast();

```

Figure 4.4: Crash Recovery in Leader Election. *The figure shows a simplified example of crash recovery in a leader election protocol. The code runs in every node. C implies the crashing node; each node typically has a view of the states of its peers. Three predicate-recovery pairs are shown (pr₁, pr₂, and pr₃). In the first, if quorum still exists, the leader simply decrements the follower count. In the second, if quorum breaks, the leader falls back to searching mode (S). In the third, if the leader crashes, the node (as a follower) votes for itself and broadcasts the vote to elect a new leader.*

Second, we extract recovery logic from the code by flattening the predicate-recovery pairs (*i.e.*, recovery-related `if` blocks). Figure 4.4 shows a simple example. Different recovery actions will be triggered based on which recovery predicate (pr₁, pr₂, or pr₃) is true. Each predicate depends on the local state and the information about the crashing node. Our key here is to map each predicate-recovery pair to this formal pattern:

```

if (pri(ls, C.ls))
    modify(ralsi);
    (and/or)
    sendMsg(ralsi);

```

Here, pr_i is the recovery predicate for the *i*-th recovery action, and rals_i is the recovery abstract

local state (*i.e.*, a subset of all fields of the local state involved in recovery). That is, each recovery predicate defines what recovery abstract local state that matters (*i.e.*, $pr_i \rightarrow \{rals_i\}$). For example, in Figure 4.4, if pr_1 is true, then $rals_1$ only contains the `follower` variable; if pr_3 is true, $rals_3$ contains `role` and `leaderId` variables.

Third, before we crash a node, we check which pr_i will be true on each surviving node and then obtain the $rals_i$. Next, we combine $rals_i$ of all surviving nodes and *sort* them into a recovery abstract global state ($rags$); sorting $rags$ helps us exploit topological symmetry (*e.g.*, individual node IDs often do not matter).

Fourth, given a plan to crash a node, the algorithm above gives us the $rags$ that represents the corresponding recovery action. We also maintain a history of $rags$ of previous injected crashes. If the $rags$ already exists in the history, then the crash is skipped because it will lead to a symmetrical recovery of the past.

To recap with a concrete example, let's go back to the case of FFFL where we plan to enable `crash(N1)`. Based on the code in Figure 4.4, the $rags$ is $\{*, \emptyset, \emptyset, \#follower=3\}$; $*$ implies the crashing node, \emptyset means there is no true predicate at the other two follower nodes, and $\#follower=3$ comes from $rals_1$ of pr_1 of $N4$ (the leader). CRS will sort this and check the history, and assuming no hit, then `crash(N1)` will be enabled. In another execution, SAMC finds that `crash(N2)` at FFFL will lead to $rags:\{\emptyset, *, \emptyset, \#follower=3\}$, which after sorting will hit the history, and hence `crash(N2)` is skipped. If the system enters a different state FFLE, no follower crash will be injected, because the $rags$ will be the same as above. In terms of leader crash, crashing the leader in the two cases will be treated differently because in a leader crash, pr_3 is true on followers and pr_3 involves `leaderId` which is different in the two cases.

In summary, the foundation of CRS is the computation of recovery abstract global state ($rags$) from the crash recovery logic extracted from the target system via the $pr_i \rightarrow \{rals_i\}$ pattern. We believe this extraction method is simple because CRS does *not* need to know the specifics of crash recovery; CRS just needs to know what variables are involved in recovery (*i.e.*, the $rals$).

Reboot Synchronization Symmetry (RSS)

Reboots are also essential to exercise (Section ??), but if not done carefully, will introduce more scalability problems. Reboot reduction policy is needed to help dmck inject reboots “smartly”. The intuition behind reboot synchronization symmetry (RSS) is similar to that of CRS. When a node reboots, it typically *synchronizes* itself with the peers. However, a reboot will not lead to a new scenario if the current state of the system is similar to the state when the node crashed. To implement RSS, we extract reboot-synchronization predicates and the corresponding actions. Since the overall approach is similar to CRS, we omit further details.

In our experience RSS is extremely powerful. For example, it allows us to find deep bugs involving multiple reboots in the ZooKeeper atomic broadcast (ZAB) protocol. RSS works efficiently here because reboots in ZAB are only interesting if the live nodes have seen new commits (*i.e.*, the dead node falls behind). In contrast, a black-box dmck without RSS initiates reboots even when the live nodes are in similar states as in before the crash, prolonging the discovery of deep bugs.

4.2.4 Pattern Extraction

We have presented four general, simple, and powerful semantic-aware reduction policies along with the generic event processing patterns. With this, testers can write protocol-specific rules by extracting the patterns from their target systems. Given the patterns described in previous sections, a tester must perform what we call as “extraction” phase. Here, the tester must extract the patterns from the target system and write protocol-specific rules specifically by filling in the predicates and abstractions as defined in previous sections; in Section 4.3.2, we will show a real extraction result (*i.e.*, real rules). Currently, the extraction phase is manual; we leave automated approaches as a future work (Section 4.5). Nevertheless, we believe manual extraction is bearable for several reasons. First, today is the era of DevOps [111] where developers are testers and vice versa; testers know the internals of their target systems. This is also largely true in cloud system development.

Second, the processing patterns only cover high-level semantics; testers just fill in the predicates and abstractions but no more details. In fact, simple semantics are enough to significantly help dmck go faster to deeper states.

4.3 Implementation and Integration

In this section, we first describe our SAMC prototype, SAMPRO, which we built from scratch because existing dmcks are either proprietary [157] or only work on restricted high-level languages (*e.g.*, Mace [103]). We will then describe SAMPRO integration to three widely popular cloud systems, ZooKeeper [94], Hadoop/Yarn [148], and Cassandra [105]. Prior to SAMPRO, there was no available dmck for these systems; they are still tested via unit tests, and the test code size is bigger than the main code, but the tests are far from reaching deep bugs.

4.3.1 SAMPRO

SAMPRO is written in 10,886 lines of code in Java, which includes all the components mentioned in Section 4.1.1 and Figure 2.2. The detailed anatomy of dmck has been thoroughly explained in literature [75, 81, 103, 143, 157], and therefore for brevity, we will not discuss many engineering details. We will focus on SAMC-related parts.

We design SAMPRO to be highly portable; we do not modify the target code base significantly as we leverage a mature interposition technology, AspectJ, for interposing network messages and timeouts. Our interposition layer also sends local state information to the SAMPRO server. SAMPRO is also equipped with crash and reboot scripts specific to the target systems. The tester can specify a budget of the maximum number of crashes and reboots to inject per execution. SAMPRO employs basic reduction mechanisms and advanced reduction policies as described before. We deploy safety checks at the server (*e.g.*, no two leaders). If a check is violated, the trace that led to the bug is reported and can be deterministically replayed in SAMPRO. Overall, we have built all the necessary features to show the case of SAMC. Other features such as intra-node thread

interleavings [81], scale-out parallelism [144], and virtual clock for network delay [157] can be integrated to SAMPRO as well.

4.3.2 *Integration to Target Systems*

In our work, the target systems are ZooKeeper, Hadoop 2.0/Yarn, and Cassandra. ZooKeeper [94] is a distributed synchronization service acting as a backbone of many distributed systems such as HBase and High-Availability HDFS. Hadoop 2.0/Yarn [148] is the current generation of Hadoop that separates cluster management and processing components. Cassandra [105] is a distributed key-value store derived from Amazon Dynamo [59].

In total, we have model checked 7 protocols: ZooKeeper leader election (ZLE) and atomic broadcast (ZAB), Hadoop cluster management (CM) and speculative execution (SE), and Cassandra read/write (RW), hinted handoff (HH) and gossip (GS). These protocols are highly asynchronous and thus susceptible to message re-orderings and failures.

Table 4.1 shows a real sample of protocol-specific rules that we wrote. Rules are in general very short; we only wrote 35 lines/protocol on average. This shows the simplicity of SAMC's integration to a wide variety of distributed system protocols.

4.4 **Evaluation**

We now evaluate SAMC by presenting experimental results that answer the following questions:

1. How fast is SAMC in finding deep bugs compared to other state-of-the-art techniques?
2. Can SAMC find new deep bugs?
3. How much reduction ratio does SAMC provide?

To answer the first question, we show SAMC's effectiveness in finding old bugs. For this, we have integrated SAMPRO to old versions of our target systems that carry deep bugs: ZooKeeper

Local-Message Independence (LMI)	Crash-Message Independence (CMI)	Crash Recovery Symmetry (CRS)	Reboot Synchronization Symmetry (RSS)
<pre> bool pd : !newVote(m, s); bool pm : newVote(m, s); bool newVote(m, s) : if (m.ep > s.ep) return true; else if (m.ep == s.ep) if (m.tx > s.tx) return true; else if (m.tx == s.tx && m.lid > s.lid) return true; </pre>	<pre> bool pg (s, X) : if (s.rl == F && X.rl == L) return true; if (s.rl == L && X.rl == F && !quorumAfterX(s)) return true; if (s.rl == S && X.rl == S) return true; bool pl (s, X) : if (s.rl == L && X.rl == F && quorumAfterX(s)) return true; bool quorumAfterX(s) : ret ((s.fol-1) >= s.all/2); </pre>	<pre> bool pr1(s,C): if (s.rl == L && C.rl == F && quorumAfterX(s)) return true; rals1:{rl,fol,all}; bool pr2(s,C): if (s.rl == L && C.rl == F && !quorumAfterX(s)) return true; rals2: {rl,fol,lid,ep,tx,clk} bool pr3(s,C): if (s.rl == F && c.rl == L) return true; rals3: {rl,fol,lid,ep,tx,clk} bool pr4: if (s.rl == S) return true; rals4: {rl,lid,ep,tx,clk} </pre>	<pre> bool ps1(s,R): if (s.rl == L) return true; sals1: {rl,lid,ep,tx,clk} bool ps2(s,R): if (s.rl == F) return true; sals2: {rl,lid,ep,tx,clk} bool ps3(s,R): if (s.rl == S && s.clk > R.clk) return true; sals3: {rl,lid,ep,tx,clk} bool ps4(s,R): if (s.rl == S && moreUpdated(s, R)) return true; sals4: {rl,lid,ep,tx,clk} bool moreUpdated(s, R): if (R.ep > s.ep) return true; else if (R.ep == s.ep) if (R.tx > s.tx) return true; else if (R.tx == s.tx) if (R.lid > s.lid) return true; </pre>

Table 4.1: Protocol-Specific Reduction Rules for ZLE. *The code above shows the actual protocol-specific rules for ZLE protocol. These rules are the inputs to the four reduction policies. Many variables are abbreviated (ep: epoch, tx: latest transaction ID, lid: leader ID, rl: role, fol: follower count, all: total node count, clk: logical clock, L: leading, F: following, S: searching, X/C: crashing node, R: rebooting node). LMI pc and pi predicates are not used for ZLE, but used for other protocols.*

v3.1.0, v3.3.2, v3.4.3, and v3.4.5, Hadoop v2.0.3 and v2.2.0, and Cassandra v1.0.1 and v1.0.6. To answer the second question, we have integrated SAMPRO to two recent stable versions: ZooKeeper v3.4.6 (released March 2014) and Hadoop v2.4.0 (released April 2014). In total, we have integrated SAMPRO to 10 versions, showing the high portability of our prototype. Overall, our extensive evaluation exercised more than 100,000 executions and used approximately 48 full machine days.

4.4.1 Speed in Finding Old Bugs

This section evaluates the speed of SAMC vs. state-of-the-art techniques in finding old deep bugs. In total, we have reproduced 12 old deep bugs (7 in ZooKeeper, 3 in Hadoop, and 2 in Cassandra). Figure 4.5 illustrates the complexity of the deep bugs that we reproduced.

Old Issue#	Protocol				#Executions				Speed-up of SAMC vs.		
		E	C	R	bDP	RND	rDP	SAMC	bDP	RND	rDP
ZK-335	ZAB	120	3	3	↑5000	1057	↑5000	117	↑43	9	↑43
ZK-790	ZLE	21	1	1	14	225	82	7	2	32	12
ZK-975	ZLE	21	1	1	967	71	163	53	18	1	3
ZK-1075	ZLE	25	3	2	1081	86	250	16	68	5	16
ZK-1419	ZLE	25	3	2	924	2514	987	100	9	25	10
ZK-1492	ZLE	31	1	0	↑5000	↑5000	↑5000	576	↑9	↑9	↑9
ZK-1653	ZAB	60	1	1	945	3756	3462	11	86	341	315
MR-4748	SE	25	1	0	22	6	6	4	6	2	2
MR-5489	CM	20	2	1	↑5000	↑5000	↑5000	53	↑94	↑94	↑94
MR-5505	CM	40	1	1	1212	↑5000	1210	40	30	↑125	30
CA-3395	RW+HH	25	1	1	2552	191	550	104	25	2	5
CA-3626	GS	15	2	1	↑5000	↑5000	↑5000	96	↑52	↑52	↑52

Table 4.2: **SAMC Speed in Finding Old Bugs.** “E”, “C” and “R” represent the number of events, crashes, and reboots necessary to hit the bug. The numbers in the middle four columns represent the number of executions to hit the bug across different policies. “bDP”, “RND”, and “rDP” stand for black-box DPOR (in MODIST), random, and random + black-box DPOR respectively. We stop at 5000 executions (around 2 days) if the bug cannot be found (labeled with “↑”). Thus, speed-up numbers marked with “↑” are potentially much higher.

Table 4.2 shows the result of our comparison. We compare SAMC with basic techniques (DFS and Random) and advanced state-of-the-art techniques such as black-box DPOR (“bDP”) and Random+bDP (“rDP”). Black-box DPOR is the MODIST-style of DPOR (Section 4.1.2). We

MR-5505:

1. A job finishes,
2. Application manager (AM) sends a “remove-app” message to Resource Manager (RM),
3. RM receives the message,
4. AM is unregistering,
5. *RM crashes* before completely processes the message,
6. AM finishes unregistering,
7. *RM reboots* and reads the old state file,
8. RM thinks the job has never started and runs the job again.

CA-3395

1. Three nodes N1-3 started and formed a ring,
2. Client writes data,
3. *N3 crashes*,
4. Client updates the data via N1; N3 misses the update,
5. *N3 reboots*,
6. N1 begins the hinted handoff process, (7) Another client reads the data with strong consistency via N1 as a coordinator,
7. N1 and N2 provide the updated value, but N3 still provides the stale value,
8. The coordinator gets “confused” and returns the stale value to the client!

Figure 4.5: **Complexity of Deep Bugs.** Above are two sample deep bugs in Hadoop and Cassandra. A sample for ZooKeeper was shown in the introduction (Section ??). Deep bugs are complex to reproduce; crash and reboot events must happen in a specific order within a long sequence of events (there are more events behind the events we show in the bug descriptions above). To see the high degree of complexity of other old bugs that we reproduced, interested readers can click the issue numbers (hyperlinks) in Table 4.2.

include Random+DPOR to mimic the way MODIST authors found bugs faster (Section 4.1.3). The table shows the number of executions to hit the bug. As a note, software model checking with the inclusion of failures takes time (back-and-forth communications between the target system and the dmck server, killing and restarting system processes multiple times, restarting the whole system from a clean state, etc.). On average, each execution runs for 40 seconds and involves a long sequence of 20-120 events including the necessary crashes and reboots to hit the bug. We do not show the result of running DFS because it never hits most of the bugs.

Based on the result in Table 4.2, we make several conclusions. First, with SAMC, we prove that smart systematic approaches can reach to deep bugs quickly. We do not need to revert to randomness or incorporate checkpoints. As a note, we are able to reproduce every deep bug that we picked; we did not skip any of them. (Hunting more deep bugs is possible, if needed).

Second, SAMC is one to two orders of magnitude faster compared to state-of-the-art techniques. Our speed-up is up to 271x (33x on average). But most importantly, there are bugs that other techniques cannot find even after 5000 executions (around 2 days). Here, SAMC’s speed-up factor is potentially much higher (labeled with “↑”). Again, in the context of dmck (a process of hours/days), large speed-ups matter. In many cases, state-of-the-art policies such as bDP and rDP cannot reach the bugs even after very long executions. The reasons are the two problems we mentioned earlier (Section 4.1.3). Our micro-analysis (not shown) confirmed our hypothesis that non-SAMC policies frequently make redundant crash/reboot injections and event re-orderings that anyway lead to insignificant state changes.

Third, Random is truly “random”. Although many previous dmcks embrace randomness in finding bugs [103, 157], when it comes to failure-induced bugs, we have a different experience. Sometimes Random is as competitive as SAMC (e.g., ZK-975), but sometimes Random is much slower (e.g., ZK-1419), or worse Random sometimes did not hit the bug (e.g., ZK-1492, MR-5505). We find that some bugs require crashes and/or reboots to happen at very specific points, which is probabilistically hard to reach with randomness. With SAMC, we show that being systematic and

semantic aware is consistently effective.

4.4.2 Ability of Finding New Bugs

The previous section was our main focus of evaluation. In addition to this, we have integrated SAMPRO to recent stable versions of ZooKeeper (v3.4.6, released March 2014) and Hadoop (v2.4.0, released April 2014). In just hours of deployment, we found 1 new ZLE bug involving 2 crashes, 2 reboots, and 52 events, and 1 new Hadoop speculative execution bug involving 2 failures and 32 events. These two new bugs are distributed data race bugs. The ZLE bug causes the ZooKeeper cluster to create two leaders at the same time. The Hadoop bug causes a speculative attempt on a job that is wrongly moved to a scheduled state, which then leads to an exception and a failed job. We can deterministically reproduce the bugs multiple times and we have reported the bugs to the developers. Currently, the bugs are still marked as major and critical, the status is still open, and the resolution is still unresolved.

We also note that in order to unearth more bugs, a dmck must have several complete features: workload generators that cover many protocols, sophisticated perturbations (*e.g.*, message re-ordering, fault injections) and detailed checks of specification violations. Further discussions can be found in our previous work [77]. Currently, SAMPRO focuses on speeding up the perturbation part. By deploying more workload generators and specification checks in SAMPRO, more deep bugs are likely to be found. As an illustration, the 94 deep bugs we mentioned in Section ?? originated from various protocols and violated a wide range of specifications.

4.4.3 Reduction Ratio

Table 4.3 compares the reduction ratio of SAMC over black-box DPOR (bDP) with different budgets (#crashes and #reboots). This evaluation is slightly different than the bug-finding speed evaluation in Section 4.4.1. Here, we measure how many executions in bDP are considered redundant based on our reduction policies and protocol-specific rules. Specifically, we run bDP for 3000

executions and run SAMC policies on the side to mark the redundant executions. The reduction ratio is then 3000 divided by the number of non-redundant executions. Table 4.3 shows that SAMC provides between 37x-166x execution reduction ratio in model checking ZLE and ZAB protocols across different crash/reboot budgets.

Table 4.3b shows that with each policy the execution reduction ratio increases when the number of crashes and reboots increases. With more crashes and reboots, the ZLE protocol generates more messages and most of them are independent, and thus the LMI policy has more opportunities to remove redundant message re-orderings. Similarly, the crash and reboot symmetry policies give better benefits with more crashes and reboots. The table also shows that LMI provides the most reduction. This is because the number of message events is higher than crash and reboot events (as also depicted in Table 4.2).

We now discuss our reduction ratio with that of DIR [81]. As discussed earlier (Section 4.1.2), DIR records local exploration (thread interleavings) and replays future incoming messages whenever possible, reducing the work of global exploration. If the target system does not have lots of thread interleavings, DIR’s reduction ratio is estimated to be between 10^1 to 10^3 (Section 5 of [81]). As we described earlier (Section 4.1.2), DIR is orthogonal to SAMC. Thus, the reduction ratios of SAMC and DIR are complementary; when both methods are combined, there is a potential for a higher reduction ratio. The DIR authors also hinted that domain knowledge can guide dmcks (and also help their work) to both scale and hit deep bugs (Section 8 of [81]). SAMC has successfully addressed such need.

Finally, we note that in evaluating SAMC, we use execution reduction ratio as a primary metric. Another classical metric to evaluate a model checker is state coverage (*e.g.*, a dmck that covers more states can be considered a more powerful dmck). However, in our observation state coverage is not a proper metric for evaluating optimization heuristics such as SAMC policies. For example, if there are three nodes ABC that have the same role (*e.g.*, follower), a naive black-box dmck will crash each node and covers three distinct states: $*BC$, $A*C$ and $AB*$. However, with a semantic-

C	R	Execution Reduction Ratio in		C	R	Execution Reduction Ratio in ZLE with				
		ZLE	ZAB			All	LMI	CMI	CRS	RSS
1	1	37	93	1	1	37	18	5	4	3
2	2	63	107	2	2	63	35	6	5	5
3	3	103	166	3	3	103	37	9	9	14

(a)
(b)

Table 4.3: **SAMC Reduction Ratio.** The table (a) shows the execution reduction ratio of SAMC over black-box DPOR (bDP) in checking ZLE and ZAB under different crash/reboot budgets. “C” and “R” are the number of crashes and reboots. The table (b) shows the execution reduction ratio in ZLE with individual policies over black-box DPOR (bDP).

aware approach (e.g., symmetry), we know that covering one of the states is sufficient. Thus, less state coverage does not necessarily imply a less powerful dmck.

4.5 Discussion

In this section, we discuss SAMC’s simplicity, generality and soundness. We would like to emphasize that the main goal of this thesis is to show the power of SAMC in finding deep bugs both quickly and systematically, and thus we intentionally leave some subtasks (e.g., automated extraction, soundness proofs) for future work.

4.5.1 Simplicity

In previous sections, we mentioned that policies can be written in few lines of code. Besides LOC, simplicity can be measured by how much time is required to understand a protocol implementation, extract the patterns and write the policies. This time metric is unfortunately hard to quantify. In our experience, the bulk of our time was spent in developing SAMPRO from scratch and integrating policies to dmck mechanisms (Section 4.1.1). However, the process of understanding protocols and crafting policies requires a small effort (e.g., few days per protocol to the point where we feel the policies are robust). We believe that the actual developers will be able to perform this process much faster than we did as they already have deeper understandings of their code.

4.5.2 Generality

Our policies contain patterns that are common in distributed systems. One natural question to ask is: how much semantic knowledge should we expose to dmck? The ideal case is to expose as much information as possible as long as it is sound. Since proving soundness and extracting patterns automatically are our future work, in this thesis we only propose exposing high-level processing semantics. With advanced program analysis tools that can analyze deep program logic, we believe more semantic knowledge can be exposed to dmck in a sound manner. For example, LMI can be extended to include commutative modifications. This is possible if the program analysis can verify that the individual modification does not lead to other state changes. This will perhaps be the point where symbolic execution and dmck blend in the future (Section ??).

Nevertheless, we find that high-level semantics are powerful enough. Beyond the three cloud systems we target in this thesis, we have been integrating SAMC to MaceMC [103]. MACEMC already employs random exploration policies to model check Mace-based distributed systems such as Mace-based Chord and Pastry. To integrate SAMC, we first must re-implement DPOR in MaceMC (existing DPOR implementation in MaceMC is proprietary [81]). Then, we have written 18 lines of LMI protocol-specific rules for Chord and attain two orders of magnitude of reduction in execution. This shows the generality of SAMC to many distributed systems.

4.5.3 Soundness

SAMC policies only skip re-orderings and crash/reboot events that are redundant by definition, however currently our version of SAMC is not sound; the unsound phase is the manual extraction process. For example, if the tester writes a wrong predicate definition (*e.g.*, `pd`) that is inconsistent with what the target system defines, then soundness (and correctness) is broken. Advanced program analysis tools can be developed to automate and verify this extraction process and make SAMC sound. Currently, the fact that protocol-specific rules tend to be short might also help in reducing human errors. Our prototype, SAMPRO, is no different than other testing/verification

tools; full correctness requires that such tools to be free of bugs and complete in checking all specifications, which can be hard to achieve. Nevertheless, we want to bring up again the discussion in Section 4.1.3 that dmck’s scalability and ability to find deep bugs in complex distributed systems are sometimes more important than soundness. We leave soundness proofs for future work, but we view this as a small limitation, mainly because we have successfully shown the power of SAMC.

4.6 Conclusion

In this chapter, we present semantic-aware model checking, a white-box principle that takes simple semantic information of the target system and incorporates that knowledge in state-space reduction policies by skipping redundant executions. We find that simple semantic knowledge can assist dmck greatly; the semantic are from *event processing semantic* (i.e., how messages, crashes, and reboots are processed by the target system). To help testers extract such information from the target system, we provide *generic event processing patterns*, patterns of how messages, crashes, and reboots are processed by distributed systems in general.

We introduce four novel semantic-aware reduction policies:

1. **Local-Message Independence** (LMI) reduces re-orderings of intra-node messages.
2. **Crash-Message Independence** (CMI) reduces re-orderings of crashes among messages.
3. **Crash Recovery Symmetry** (CRS) skips crashes that lead to same recovery behaviors.
4. **Reboot Synchronization Symmetry** (RSS) skips reboots that lead to same synchronization actions.

Our reduction policies are *generic*; they are applicable to many distributed systems. SAMC users (i.e., testers) only need to feed the policies with short *protocol-specific rules* that describe event independence and symmetry specific to their target systems. And SAMC is purely systematic; it does not incorporate randomness or bug-specific knowledge.

CHAPTER 5

SCK: A SINGLE-MACHINE APPROACH FOR FINDING SCALABILITY BUGS IN CLOUD-SCALE DISTRIBUTED SYSTEMS

In the previous two chapters, we discuss about distributed concurrency bugs that how we make sense of them and how we advance the state of the art of model checking to unearth them faster. In this chapter, we will focus on scalability aspect of cloud backend. As we mention in Section 2.3, the trend of cloud distributed systems goes to horizontal scaling or scale-out systems. On the positive side, scale surpasses the limit of a single machine in meeting users’ increasing demands of computing and storage, which led to many inventions of “cloud-scale” distributed systems [48, 58, 59, 73, 91, 150]. The field has witnessed a phenomenal deployment scale of such systems; Netflix runs tens of 500-node Cassandra clusters [35], Apple deploys a total of 100,000 Cassandra nodes [3], and Yahoo! recently revealed the use of 40,000 Hadoop servers, with a 4500-node cluster as the largest one [36].

On the negative side, scale creates new development and deployment issues. Developers must ensure that their algorithms and protocol designs to be scalable. However, until real deployment takes place, scalability bugs in the actual implementations are unforeseen. These scalability bugs are latent bugs that are scale-dependent; they only surface in large-scale deployments, but not in small/medium-scale ones. Their presence jeopardizes systems reliability and availability at scale.

In this chapter, we show our observations from our initial study on scalability bugs that highlights an urgency in tackling scalability bugs and our pilot work to introduce a single-machine testing methodology to check scalability of the systems. Section 5.1 discusses about our initial study and shows our observations toward scalability bugs; Section 5.2 discusses about the state of the art; Section 5.3-5.5 explains our pilot work on scalability bug checking named SCK¹ and some evaluations.

1. SCK is a joint work with Cesar A. Stuardo.

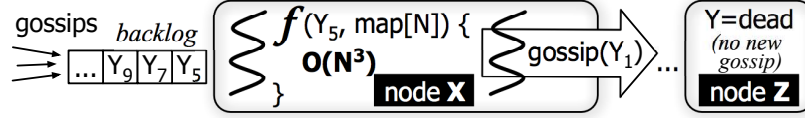


Figure 5.1: **The problem of gossip-based failure detection in Cassandra.**

5.1 Scalability Bug

To begin our work on combating scalability bugs, we start by conducting a study to see their nature and how critical they are. We perform a study of 41 scalability bugs reported from the deployments of popular large-scale systems such as Hadoop, HBase, HDFS, Cassandra, Couchbase, Riak, and Voldemort. From the study, we observed many challenges in finding, reproducing, and debugging scalability bugs. For example, let us consider a bug in Cassandra, a highly-scalable peer-to-peer key-value store. If a user decommissions a single node from a cluster of 50 nodes, the operation can be done smoothly. However, if the user decommissions a node from a cluster of 200 nodes, the protocol that re-calculate the key-range partitions (which nodes should own which key ranges) becomes CPU intensive as the calculation has an $O(N^3)$ complexity where N is the number of nodes. This combined with the gossiping and failure detection logic leads to a scalability bug that makes the cluster unstable (many live nodes are declared as dead, making some data not reachable by the users). We give a full detail of the Cassandra bug in Section 5.1.1.

As in the example above, bug symptoms sometimes surface only in large deployment scales (e.g., $N > 100$ nodes), hence small/medium-scale testing is not enough. Yet, not all developers have large test budgets, and even when they do, debugging on hundreds of nodes is time consuming and difficult. Furthermore, protocol algorithms can be scalable in the design sketches, but not necessarily in the real deployments; there are specific implementation details whose implications at scale are hard to predict. We discuss more about our observations on scalability bugs in Section 5.1.2.

5.1.1 A Sample Cassandra Bug

We now describe in detail a scalability bug in Cassandra, which we use as a sample bug. Our journey in understanding scalability bugs began when we observed repeated “flapping” problems in large-scale Cassandra deployments (*i.e.*, hundreds of nodes). Flapping is a cluster instability problem where node’s up/down status continuously flaps. A “flap” is when a node X marks a peer node Y as down (and soon marks Y as alive again). We rigorously study a series of Cassandra bugs below that surfaced as the code evolved.

To understand this bug, we need to understand the following protocols.

1. **Bootstrap:** Each node first creates partition keys (*e.g.*, 32 random numbers) and gossips this information to peer nodes.
2. **Gossip broadcast:** *Every second*, each node gossips to one random node about a list of nodes and partitions it knows (including itself) and their *version* numbers. Each node also increments its version number (“I’m still alive”) before gossiping.
3. **Gossip processing:** The receiving node then finds any state (metadata) differences between the two nodes to synchronize their views of the ring. Eventually, all nodes know about each other.
4. **Failure detection:** *Every second*, a failure detection daemon runs [105]. Put simply, if a node X has not received a new gossip about Y *from anyone* (Y’s version has not changed after some period of time), X will declare Y dead (a flap). When X receives a new gossip about Y, it marks Y alive.

There are two factors that induce the bug. The first is the *long latency of scale-dependent state-update gossip processing during bootstrapping* (“f” in Figure 5.1). While gossip processing is usually fast in a stable cluster, it is expensive during bootstrapping as the gossips carry many new state changes about the ring; the state-update processing time is scale-dependent (*i.e.*, greater

than $O(N^3)$); the larger the cluster (N), the larger the ring map, the longer the processing time is. This long latency is caused by (1) state-update checkpoint to on-disk database and (2) multi-map cloning and updates. The first one is needed for fast fault tolerance; after a node crashes, it can reboot fast as it knows the latest view of the ring. The second one is preferred for simplicity; Cassandra clones its `MultiMap` ring table and applies changes one by one to alleviate long write locks.

The second factor is the *single threaded* implementation of gossip processing. As shown in Figure 5.1, this inability to process multiple gossips/state updates concurrently (for the sake of preventing concurrency bugs) creates a *backlog* of new gossips. For example, in *every second*, Y tells someone it's alive with increasing version number (e.g., Y_7), but the receiving nodes are still busy processing state changes and only forward Y's old version number (e.g., Y_1). As Y's new gossip is not propagated on time, other nodes (e.g., Z) will mark Y as dead. This happens to all nodes, not just Y.

The journey starts with Bug #CA-3831 [10], when a node D is decommissioned from a cluster ring, D initiates a gossip telling that all other nodes must rebalance the ring's key-ranges. This scale-dependent "pending key-range calculation" is CPU intensive with $O(MN^3 \log^3(N))$ complexity; M is the list of key-range changes in the gossip message. This in turn leaves many gossips not propagated on time, creating flapping symptoms that only appear at scale (at 200+ nodes). The developers then optimized the code to $O(MN^2 \log^2(N))$ complexity.

Soon afterwards (Bug #CA-3881 [11]), Cassandra added the concept of virtual partitions/nodes (e.g., $P=256$ per physical node). As an implication, the fix above did not scale as " N " becomes $N \times P$. The bug was fixed with a complete redesign of the pending key-range calculation, making it $O(MNP \log^2(NP))$.

About a year later (CA-5456 [12]), Cassandra code employs multi-threading between the pending key-range calculation and the gossip processing with a coarse-grained lock to protect sharing of the ring table. Unbeknownst to the developers, at scale, the key-range calculation can acquire

the lock for a long time, causing flapping to reappear again. The fix clones the ring table for the key-range calculation, to release the lock early.

Later on (CA-6127 [13]), a similar bug reappeared. In the above cases, the problems appeared when the cluster grows/shrinks gradually. However, if user bootstrap a large cluster (*e.g.*, 500+ nodes) from scratch (*i.e.*, all nodes do not know each other, with no established key ranges), the execution traverses a different code path that performs a fresh ring-table/key-range construction with $O(MN^2)$ complexity.

The story continues on (CA-6345, CA-6409, etc.). Fast forward today, Cassandra developers recently started a new umbrella ticket for discussing “Gossip 2.0,” supposedly scalable to 1000+ nodes [22, 23]. Similar to Cassandra, other large-scale systems are prone to the same problem. So far, we have collected and analyzed 9 Cassandra, 5 Couchbase, 2 Hadoop, 10 HBase, 13 HDFS, 1 Riak, and 1 Voldemort scalability bugs, all caused user-visible impacts. This manual mining was arduous because there is no searchable jargon for “scalability bugs”; we might have missed other bugs.

5.1.2 Observations

From all the bugs we studied, we make several important observations.

- *Only appear at extreme scale:* This Cassandra bug does not surface in 30-node deployment. In 128-node cluster, the symptom appears mildly (tens of flaps). From 200-500 nodes, flapping skyrockets from hundreds to thousands of flaps. Testing in small/medium scales is not sufficient, which is also true for other bugs we studied.
- *Scalable in design, but not in practice.* Related to the Cassandra bug, the accrual failure detector/gossiper [88] was interestingly adopted by Cassandra as it is scalable in design [105]. However, the design proof does not account gossip processing time during bootstrap, which can be long. To understand the bug, the developers tried to “do the [simple] math”

[13] but failed. In practice, the assumption that new gossips are propagated every second is not met (due to the backlog). The actual implementations overload gossips with many other purposes (*e.g.*, announcing boot/rebalance changes) beyond their original design sketch.

- *Implementation specific and hard to predict.* The backlog-induced flapping in Cassandra was caused specifically by Cassandra’s implementation choice: metadata checkpoint, multi-map cloning, and its single-threaded implementation. State-update processing time is hard to predict (ranges from 0.001 to 4 seconds) as it depends on a 2-dimensional input: the receiving node’s ring table size and the number of new state changes.
- *Cascading impacts of “not-so-independent” nodes.* In cluster-wide control protocols, distributed nodes are not necessarily independent; nodes must communicate with each other to synchronize their views of cluster metadata. As the cluster grows, the cluster metadata size increases. Thus, unpredictable processing time in individual nodes can create cascading impacts to the whole cluster.
- *Long and difficult large-scale debugging:* Some bug reports like CA-6127 generated over 40 back-and-forth discussion comments and took 2 months to fix. It is apparent [13] that there were many hurdles of deploying and debugging the buggy protocol at real scale. Important to note is that debugging is *not* a single iteration; developers must *repeatedly* instrument the system (add more logs) and re-run the system at scale to find and fix the bug, which is not trivial. The scalability bugs we studied took 6 to 157 days to fix (27 on average).
- *Not all developers have large test budgets:* Another factor of delayed fixes is the lack of budget for large test clusters. Such luxury tends to be accessible to developers in large companies, but not to open-source developers. When CA-6127 was submitted by a customer who had hundreds of nodes, the Cassandra developers did not have an instant access to a test cluster of the same scale.

- *Quick fixes and repeated bugs:* Bugs are often fixed with quick patches (development pressures), but the new fix might not eradicate the problem completely [159]. For example, for CA-6127, the patch simply disables failure detection during bootstrap. As the protocol was not redesigned, the bug still appeared in another workload (*e.g.*, scaling out from 128 to 256 nodes). In the latest Cassandra, the simple fix has been removed and the gossip protocol has been redesigned. We also found that old fixes can become obsolete in protocol re-designs, which then can give birth to new scalability bugs. For example, the fix for CA-3831 became obsolete as “vnodes” was introduced, which then gave rise to a new vnode-related scalability bug (CA-3881). A scale-check could have ensured that new fixes remove old scalability bugs entirely and similar bugs do not re-surface in new designs.

5.2 State of the Art of Large-Scale Emulation

Our observations above accentuate the need for scale-checking distributed system *implementations* at *real scale*, not via simulation nor extrapolation. In this context, we now discuss the state of the art of large-scale emulation..

DieCast [83], invented for network emulation, can colocate many processes/VMs on a single machine as if they run individually without contention. The trick is adding “time dilation factor” (TDF) support [84] into the VMM (*e.g.*, Xen). For example, TDF=5 implies that for every second of wall-clock time, each emulated VM on the VMM believes that time has advanced by only 200 ms. The most significant drawback of DieCast is that high colocation factor (*e.g.*, TDF=100) is likely not desirable, for two reasons: prolonged testing time (TDF=100 implies 100x longer run) and memory overcapacity. DieCast was only evaluated with TDF=10.

Exalt [151] targets I/O-intensive scalability bugs. With a custom data compression, users’ data is compressed to zero byte on disk (but the size is recorded) while metadata is not compressed. With this, Exalt can colocate 100 emulated HDFS datanodes on one machine. In its evaluation, most of the bugs reproduced are in the HDFS namenode which runs alone on one machine. As the

authors stated, their approach “may not discover scalability problems that arise at the nodes that are being emulated [the datanodes]” (Section 4.1 in [151]). Thus, Exalt is not suitable for finding scalability bugs in CPU-intensive distributed systems.

In summary, we did not find a fast and accurate single-machine approach that can scale-check CPU-intensive protocols in cloud systems. The scalability bugs could be caused by the scale-dependent processing time, not network or I/O bottlenecks. As DieCast targets *network* emulation via time dilation and Exalt targets *storage* space emulation via compression.

5.3 SCK

We (along with my colleague, Cesar A. Stuardo) now present the four SCK techniques to achieve high colocation factor in one machine (Section 5.3.1-5.3.4) and summarize how to use these techniques to scale check the systems (Section 5.3.5).

The enabler of our methods is the fact that SCK is an offline check, which provides us the liberty to re-architect distributed systems to be scale-checkable offline, but without introducing performance overhead online. When N is large, every small optimization in SCK will lead to a multitude of benefits. When the four techniques are combined, SCK can achieve a colocation factor of around 500 nodes with high result fidelity.

5.3.1 Processing Illusion (PIL)

While I/O and memory bottlenecks are often blamed for many scalability problems [130, 142, 151], scalability bugs also could be caused by cascading impacts of CPU-intensive processing as the example of Cassandra bugs. To emulate CPU-intensive processing, we introduce *processing illusion* (PIL), an approach that *replaces an actual processing with* `sleep()`. For example, in the sample Cassandra bug, we can replace the expensive ring-table update with `sleep(τ)` where τ is an accurate timing of how long the update takes.

The intuition behind PIL is similar to the intuition behind other emulation techniques. For

example, Exalt provides an illusion of storage space; their insight was “how data is processed is not affected by the content of the data being written, but only by its size” [151]. PIL provides an illusion of compute processing; our insight is that “*the key to computation is not the intermediate results, but rather the execution time and eventual output*”.

PIL might sound outrageous in the first place, but it is feasible only if the following concerns are addressed: How can a function be safely replaced with `sleep()` *without* changing the whole processing semantic? And how to find specific functions that should be replaced with `sleep()`? How can we produce the output if the actual compute is skipped? How can we predict the actual compute time (τ) accurately?

PIL-Safe Functions

Our first challenge is to ensure that functions (or code blocks) can be safely replaced with `sleep()`, but still retain the cluster-wide behavior and unearth the bug symptoms. We name such functions as *PIL-safe functions*. We identify two main characteristics of sleep-safe functions:

1. Memoizable output: a PIL-safe function must have a memoizable (deterministic) output based on the input of the function.
2. Non-pertinent I/Os: if a function performs disk/network I/Os that are not pertinent to the correctness of the corresponding protocol, the function is PIL-safe. For example, in the sample Cassandra bugs, there is a ring-table checkpoint (not shown) needed for fault tolerance but is irrelevant (never read) during bootstrapping.

Depending on the modularity of the target system, manually finding such target functions can be challenging, primarily because scale-dependent nested loops can span across multiple functions. Right now, we rely on developers to identify such that functions, and we leave automatic PIL-taking function extracting for future work (Chapter 6).

Pre-Memoization with Order Determinism

As sleep-safe functions no longer perform the actual computation, the next question to address is: how do we manufacture the output? We find there are sleep-safe functions with non-pertinent outputs (Section 5.3.1). For these functions, only time profiling is needed (Section 5.3.1) but not output recording. However, there are also sleep-safe functions with non-pertinent intermediate data but with outputs that are needed. For this latter case, we need to manufacture the outputs such that the global behavior is not altered (*e.g.*, cluster bootstrapping or rebalancing should terminate successfully). Our solution is *pre-memoization*: given a sleep-safe function, we identify all the possible inputs, and for every input, run the function and pre-memoize the output. When SCK runs, it will use the pre-memoized outputs.

Unfortunately, pre-memoization in the context of large-scale, decentralized, non-deterministic distributed systems requires an “*infinite*” time and storage space. The issue is that the state of each node (the input) depends on the *order* in which messages arrive (which can be random). As an example, let’s consider Riak’s bootstrap+rebalance protocol where eventually all nodes own a similar number of partitions. A node initially has an unbalanced partition table, receives another partition table from a peer node, then inputs it to a rebalance function, and finally sends the output to a *random* node via gossiping. *Every* node repeats the same process until the cluster is balanced. In a Riak cluster with $N=256$ and $P^2=64$, there are in total 2489 rebalance iterations with a set of specific inputs in *one* run. *Another* run of the protocol will result in a *different* set of inputs due to gossip randomness. Our calculation shows that there are $(N^{NP})^2$ possible inputs.

To address this problem, we pre-memoize with *order determinism*. Thus, repeated runs of the same workload in SCK mode will use the *same* global message ordering (akin to deterministic record and replay [71]). For example, across different runs, a Riak node now receives gossips from the same sequence of nodes. With order determinism, pre-memoization and SCK work as follow:

1. We run all nodes on one machine without PIL (more details in Section 5.3.5) and interpose

2. P : the number of key-partitions per node; A key-partition is typically a random integer within 2^{64} keyrange.

sleep-safe functions.

2. When sleep-safe functions are executed, we record the inputs and corresponding outputs to a *memoization database* (SSD-backed files).
3. During this pre-memoization phase, we *record message non-determinism* (e.g., gossip send-receive pairs and their orderings).
4. After pre-memoization completes, we can repeatedly run SCK wherein order determinism is enforced (e.g., no randomness), sleep-safe functions replaced with PIL, and their outputs retrieved from the memoization database.

We omit some details above but will summarize the steps again later along with other features (Section 5.3.5).

With order determinism, the memoization database is kept small as we only record the possible inputs within *one* deterministic order. In the 256-node Riak’s case above, the database only needs to store around 2500 input-output pairs (the number of rebalance iterations observed) in 1.3 GB of memoized data (and 5.3 GB for the 512-node setup). We also note that while the idea of deterministic systems has been made popular recently, the concept of deterministic distributed systems is still not practical due to the excessive runtime overhead (e.g., 10x slower [93]). However, in the context of offline methodology such as SCK, order determinism can be exploited in a fitting manner.

Time Profiling

As sleep-safe functions are replaced with `sleep(τ)`, we need to accurately predict the actual compute time (τ). There are two different approaches we take, depending on the target protocols.

1. The first approach is to profile compute time in situ with the pre-memoization phase (*in-situ time profiling*). That is, for each input observed during pre-memoization, we also record how long the processing takes.

2. Another approach is to profile compute time with an *offline time profiling*, which is feasible for functions with non-pertinent outputs (Section 5.3.1), which do not need pre-memoized outputs.

With offline time profiling, we simply profile the expensive function exclusively by itself with the possible input space. If faster profiling is needed, we can sample the input space. For example, in the sample Cassandra bug, the expensive function depends on a 2-dimensional input (#commit states and current ring table size), each ranges from 1 to N . With N^2 profiles, the profiling time can take more than one day without sampling when N is large (e.g., 512 nodes). When SCK runs, given an input, we normalize t based on the sampled profile.

We now address two further questions. First, is time profiling necessary? In other words, is static prediction sufficient (e.g., extrapolation based on a `for-loop` timing)? In our case, static prediction is hard to achieve; nested loops can span across multiple functions with many `if-else` conditions. For example, state-update processing in the sample Cassandra bug can range from 0.001 to 4 seconds depending on the multi-dimensional input (Section 5.1.2).

Second, is it obvious from the profiled time that a scalability bug will appear, hence obviating the need for SCK? As suggested earlier, every implementation is unique (Section 5.1.2). In the sample Cassandra bug for example, *if* Cassandra processes gossips in a multi-threaded manner, long processing time would not lead to the scalability bug. In fact, patches for scalability bugs do not always remove the expensive computation. Scalability bugs are not merely about the expensive functions, but rather their potential cascading impacts, hence it is essential to run SCK in addition to time profiling. We want to emphasize that our profiling approach is not the same as extrapolation, which tends to stop profiling at a certain scale (e.g., 100 nodes) and extrapolates the behavior for larger scales. Our profiling and SCK phases run at real scales.

Overall, PIL significantly removes processing contention and reduces CPU utilization. Interestingly, as we colocate more nodes, before we hit 100% CPU utilization, we hit other major colocation bottlenecks such as memory exhaustion and process/thread context-switching delays.

For this reason, we re-architect our target systems to make them scale-checkable with the next three optimization techniques (Section 5.3.2-Section 5.3.4).

5.3.2 *Single-Process Cluster (SPC)*

Many distributed systems today are implemented in managed languages (*e.g.*, Java, Erlang) whose runtimes consume non-negligible memory overhead. Java and Erlang VMs for example use around 70 and 64 MB of memory per process respectively. As we target 3-digit colocation factor, this memory overhead becomes an unnecessary limitation. Furthermore, a managed-language VM can contain advanced services. For example, Erlang VM contains a DNS service which sends heartbeat messages to other connected VMs. As hundreds of Erlang VMs (one for each Riak node) run on one machine, the heartbeat messages cause a “network” overflow that disconnects Erlang VMs.

To address this, we create Single-Process Cluster (SPC) support wherein the whole cluster runs as threads in a single process. Surprisingly, our target systems do not have this simple support; there is no scalability check in the unit tests (mainly feature correctness tests). Thus, we need to slightly re-design the code to support SPC (*e.g.*, adding arrays of per-node global data structures, removing static-synchronized functions that lock the whole cluster). As all nodes run in one process, user-kernel switching to send messages becomes unnecessary. Thus, we create a shim layer in our target systems to bypass OS network calls when they run in SCK mode.

5.3.3 *Global Event Driven Arch. (GEDA)*

With SPC, a node still runs multiple daemon threads (gossiper, failure detector, etc.). With high colocation factor, there are more than one thousand threads that cause severe context switching and long queuing delays. Because of this overhead, we noticed that events become late (*e.g.*, gossips are not sent every second) even though CPU utilization has not reached 100%.

To address this, we leverage the staged event-driven architecture (SEDA) [153] common in distributed system implementations. With SEDA, each service/stage in each node exclusively has

an event queue and a handler thread. In SCK mode, we convert SEDA to *global-event driven architecture* (GEDA). That is, for every stage, there is only *one* queue and one handler for the *whole* cluster.

As an example, let's consider a periodic gossip service. With 500-node colocation, there are 500 gossip threads in SPC, each sending a gossip every second. With GEDA, we only have a few threads (matched with the number of available cores). These global handler threads are shared among all the nodes for sending gossips. Here, unless a high CPU utilization is reached (*e.g.*, 90%), GEDA guarantees no late event. As another example, for gossip processing, there is only one global gossip-receiving queue shared among all nodes. Overall, GEDA removes thread context switching and queuing delays that should have never existed in the first place and does so *without* changing the processing logic, as if the nodes run exclusively on independent machines.

5.3.4 Memory Footprint Reduction (MFR)

Finally, to achieve a high colocation factor, we must perform memory footprint reduction (MFR) to prevent out-of-memory exceptions that originate from system-specific root causes.

First, relevant services in the target protocol can “over-allocate” memory. For example, in Riak's bootstrap+rebalance protocol, each node creates $N \times P$ partition services although at the end only retain P partitions and never use (remove) the other $(N-1) \times P$ partitions (as reclaimed by other nodes). Worse, each partition service is an Erlang process (1.3 MB of memory overhead); collocating 30 nodes ($N=30$ with default $P=64$) will directly consume 75 GB of memory ($30 \times 30 \times 64 \times 1.3$ MB) from the start. In SCK, we must modify Riak to remove this unoptimized memory usage.

Second, some libraries can cause high memory footprints. For example, Voldemort nodes use Java NIO [?] which is fast but contains buffers and connection metadata that take up memory space. In SCK, we address this with network bypass from SPC (Section 5.3.2).

The lesson learned here is that modern distributed systems are implemented without scale-

checkability in mind. In our target systems, we must address the specific memory issues above; other systems can potentially face other root causes. This system-specific memory optimization is crucial in SCK; as we colocate hundreds of nodes, a small memory footprint reduction per node will bring orders of magnitude reduction globally.

5.3.5 *Putting It All Together*

Integration steps: We now summarize how our four SCK techniques can be integrated to a target system.

(1) We first reduce memory footprints with SPC (Section 5.3.2) and MFR (Section 5.3.4) to avoid out-of-memory exceptions. We then modify the target system with GEDA (Section 5.3.3) to remove excessive thread context switching.

(2) Next, before running SCK with PIL, we must find all PIL-safe functions (Section 5.3.1) and interpose them to record the inputs, outputs, processing time, and non-deterministic events (Section 5.3.1, Section 5.3.1).

(3a) After the preparation, if the target protocol exhibits non-pertinent outputs, we can simply perform offline time profiling without pre-memoization (Section 5.3.1).

(3b) Otherwise, we execute the pre-memoization step (Section 5.3.1). For example, if we want to scale-check an N -node deployment, we record pre-memoized data and in-situ profiled time with colocation factor of N .

(4) Finally, we begin scale-checking the target protocol with all the features enabled, including PIL. This scale-check phase will use the recorded output and time profiles and run in deterministic order.

Note that all these features are only enabled in SCK mode. In online deployment, the system runs normally as if without any changes and modification overhead.

Debugging efficiency: We now emphasize how SCK eases scale-checking and large-scale debugging efforts.

	Cassandra	Riak	Voldemort
System	95757	18250	63485
Unit Tests	18891	1341	15820
(a) + SCK-able	2874	217	800
(b) + sck (tool)	2308	429	613
Total changes (%)	5%	4%	2%

Table 5.1: **Integration (LOC).** *The table quantifies our efforts in integrating SCK. Riak’s integration is smaller due to the conciseness of Erlang. Voldemort’s integration does not include PIL for the old bug we reproduced*

First, the only step that consumes time is the time profiling and pre-memoization phase (step 3).³ This is because nodes compete for CPU resources (PIL is still disabled). However, this is only a *one-time* overhead.

Second and most importantly, developers can *repeatedly re-run* the scale-check phase (step 4) as many times as needed (tens/hundreds of iterations) until the bug’s root cause is found. In this phase, the target protocol runs in a *similar* duration and behavior as if all the nodes run on independent machines.

Finally, developers can quickly apply and test new fixes. Some fixes can be tested by only re-running the last phase; for example, fixes such as changing the failure detector Φ threshold (for CA-6127), caching slow methods (CA-3831), changing lock management (CA-5456), and enabling parallel processing (VO-1212). However, if the fixes involve a complete redesign (*e.g.*, optimized gossip processing in CA-3881, decentralized to centralized rebalancing in RI-3926), the integration and profiling/pre-memoization steps (2 and 3) must be repeated.

5.4 Implementation and Integration

We integrate our SCK methodology to three popular distributed key-value stores: Cassandra [105], Riak [28], and Voldemort [27]. The implementation of SCK involves two parts: **(a)** changes to make the target system “SCK-able” on one machine and **(b)** the sck tool itself (specifically:

3. Ranges 1-34 hours for 100-500 nodes.

`sck.cass`, `sck.vold`, and `sck.riak`).⁴ Part (a) involves the integration of SPC, MFR, and GEDA, as well as PIL interpositioning points. Part (b) is the code for pre-memoization, time profiling, and other `sck` setups. Table 5.1 quantifies the integration efforts.

Generality: We show the generality of SCK with two major efforts. First, we picked three systems and scale-checked various control-path protocols within them, for a total of 5 protocols: 3 Cassandra (bootstrap, scale-out, decommission), 1 Riak (bootstrap+rebalance), and 1 Voldemort (rebalancing) protocols. A protocol can be built on top of other protocols (*e.g.*, bootstrap on gossip and failure detection protocols).

Second, we migrated SCK to a total of 9 old and new releases: 5 Cassandra (v0.8.9, v1.1.10, v1.2.0, v1.2.9, v2.2.5), 2 Riak (v0.14.2, v2.1.3), and 2 Voldemort (v0.90.1, v1.10.21). This effort is also important to show how SCK can find old and new bugs.

Simplicity: Table 5.1 shows SCK requires thousands of LOC, which we believe is a justified cost for supporting scale-checkability. We want to emphasize that this is a *one-time* cost; subsequent migrations are fast. Our first complete integration to Cassandra (v1.2.9) took almost a year; we needed to learn from scratch about Cassandra and its scalability bugs and also design SCK. However, after SCK techniques were solid, migration to other versions (Cassandra-v0.8.9, v1.1.10, v2.2.5) only took *1 week* each. Next, our first integration to Riak (v0.14.2) only took *4 weeks* (although Riak was completely new to us). A subsequent migration (Riak-v2.1.3) only took *4 days*. The efforts for Voldemort is also similar. Overall, we expect SCK integration can be done more seamlessly with today’s DevOps practice [111], where developers are testers and testers are developers.

4. As an analogy, part (b) is similar to specific `fscks` (*e.g.*, `fsck.ext3`, `fsck.xfs`) while part (a) is similar to how file system code is modified to make `fsck` fast [90, 122].

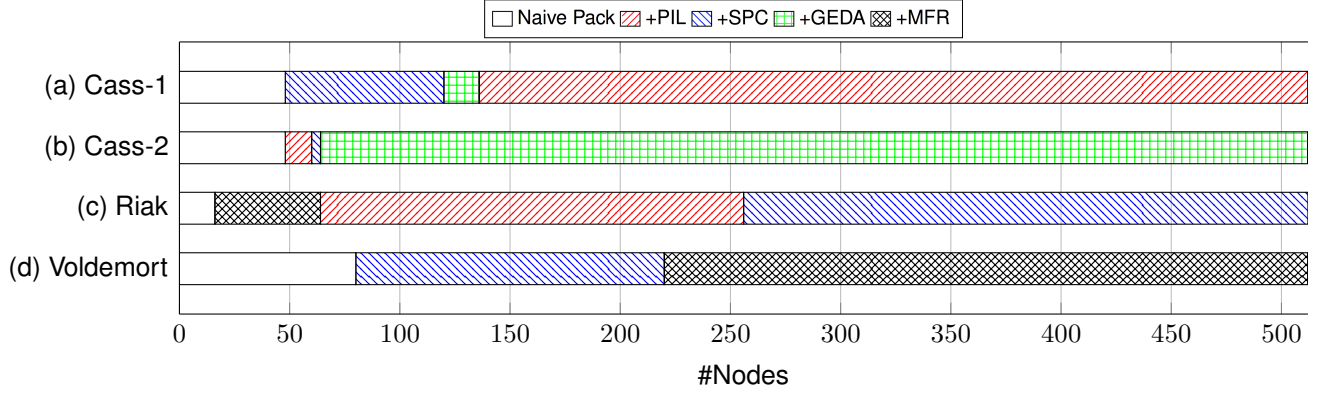


Figure 5.2: **Maximum colocation factor.**

5.5 Evaluation

We now evaluate our SCK integration to Cassandra, Riak, and Voldemort. Our evaluation answers the following questions. Section 5.5.1: What is the maximum colocation factor achieved? Section 5.5.2: How accurate is SCK compared to real deployments? Section 5.5.3: Can SCK find old scalability bugs? Section 5.5.4: Can SCK reveal new bugs?

We use the Nome cluster; each machine has 16-core AMD Opteron(tm) 8454 processors with 32-GB DRAM [25]. To measure SCK accuracy (Section 5.5.2), we compare it with real deployments of 32, 64, 128, 256, and 512 nodes, deployed on at most 128 Nome machines.⁵

5.5.1 Colocation Factor

We first show the maximum colocation factor SCK can achieve as each feature is added *one at a time* on top of the other. To recap, the features are: processing illusion (PIL; Section 5.3.1), single-process cluster (SPC; Section 5.3.2), global event driven architecture (GEDA; Section 5.3.3), and memory footprint reduction (MFR; Section 5.3.4). The results are based on a 16-core Nome machine [25].

Maximum colocation factor (MaxCF): A maximum colocation factor is reached when the target

⁵. Our target protocols only make at most 2 busy cores per node, which justifies why we run 8 nodes per one 16-core machine for the real deployment.

system’s behavior in SCK mode starts to deviate from the real deployment behavior. Deviation happens when one or more of the following bottlenecks are reached: (1) high average CPU utilization ($>90\%$), (2) memory exhaustion (nodes receive out-of-memory exceptions and crash), and (3) high event “lateness”; queuing delays from thread context switching can make events late to be processed, although CPU utilization is not high. We instrument our target systems to measure *event lateness* of relevant events (*e.g.*, gossip sending, gossip processing, and failure detection events). For example, if gossips should be sent every 1 second, but they are sent every 1.5 second on average, then the lateness is 50%. We use 10% as the maximum acceptable event lateness. Note that the residual limiting bottlenecks above come from the main logic of the target protocols, which cannot be removed with general methods.

Results and observations: Figure 5.2 shows *different sequences of integration* to our three target systems and the resulting maximum colocation factors. We discuss four important findings from this figure.

First, when multiple techniques are combined, they collectively achieve a high colocation factor (up to 512 nodes for the three systems respectively). For example, in Figure 5.2a, with just adding SPC+GEDA to Cassandra, MaxCF only reaches 136. But with SPC+GEDA+PIL, MaxCF significantly jumps to 512. When we increase the colocation factor (+100 nodes) beyond the maximum, we hit the residual bottlenecks mentioned before; at this point, we do not test MaxCF with small increments (*e.g.*, +1 node) as pre-memoization and profiling (step 3 in Section 5.3.5) takes time. The bug in Voldemort’s rebalancing protocol (VO-1212) involve sequential operations (no parallel CPU-intensive computations), hence GEDA and PIL are not necessary.

Second, distributed systems of the same type (*e.g.*, P2P key-value stores) are implemented in uniquely different ways. Thus, integrations to different systems face different sequences of bottleneck. To show this, we tried different *integration sequences*. For example, for reproducing CA-6127 in Cassandra (Figure 5.2a), our integration sequence is: +SPC, +GEDA, and +PIL (as we continuously hit CPU contention). For RI-3926 (Figure 5.2c), we began with MFR as we hit

a memory bottleneck first in Riak (the excessive Erlang processes; Section 5.3.4). For VO-1212 (Figure 5.2d), we began with SPC to reduce Java VM memory overhead in Voldemort.

Third, not all features get the chance to show their benefits as the fundamental bottlenecks are reached. For example, for Cassandra CA-6127 (Figures 5.2a-b), MFR is unnecessary as we will hit CPU contention in >512 nodes. For Riak RI-3926 (Figure 5.2c), GEDA is not needed as we will hit a memory bottleneck in >512 nodes, and similarly for Voldemort VO-1212 (Figure 5.2d).

Fourth, an SCK technique can hit a different bottleneck before showing its full potential. For example, for Cassandra, we tried two different integration sequences (Figure 5.2a-b). With Naive Packing⁶, we initially hit a MaxCF of 48 nodes due to CPU contention. At this point, there are two choices: add SPC+GEDA (to reduce process/thread context switching) or PIL (to reduce expensive processing). In Figure 5.2b, we tried +PIL first and we found that it does not help much as process/thread queuing delays are still the bottleneck. Similarly, in Figure 5.2a, SPC+GEDA also can only reach a certain maximum. This again highlights that it is the *combination* of the techniques that make SCK powerful.

So far, SCK is limited by the single machine's resources. To increase colocation factor, a higher-end machine can be used. SCK can also be extended to run on multiple machines (a future work).

5.5.2 Accuracy

Next, we provide a detailed accuracy evaluation of SCK. Due to space constraints, this section only focuses on one bug (CA-6127 [13]) while the next section briefly discuss other bugs we reproduced.

Figure 5.3a-d presents the internal metrics within Cassandra failure detection protocol that we measured for *every pair* of nodes. That is, the algorithm runs on every node A for every peer B. Figure 5.4a-d compare in detail the accuracy of SCK compared to real deployments. For example,

6. All nodes run as processes on one machine without modification.

a) $\#flaps = f(\Phi > 8)$
b) $\Phi = f(T_{avgGossip}, T_{lastGossip})$ $T_{avgGossip} = \text{avg. of last 1000 } T_{lastGossip}$
c) $T_{lastGossip} = f(\#hops, T_{gossipExec})$ $\#hops = \log(N)$ on average $T_{gossipExec} = T_{stateUpdate}$ (if new state changes)
d) $T_{stateUpdate} = f(Size_{ringTable}, Size_{newStates})$ $Size_{ringTable} \leq N \times P$ and $Size_{newStates} \leq N$

Figure 5.3: **Cassandra internal metrics.** Above are the metrics we measured within the Cassandra bootstrap protocol for measuring SCK accuracy (Figure 5.4). “ f ” represents “a function of” (i.e., an arbitrary function).

$x=512$ implies the comparison of 512-node colocation in SCK versus a real deployment of 512 nodes. Note that for CA-6127, we only need time profiling with offline sampling (Section 5.3.1) and no pre-memoized data (Section 5.3.1). We use pre-memoization to reproduce Riak RI-3926 (next section).

Figure 5.4a shows the total number of flaps (alive-to-dead transitions) observed in the whole cluster during bootstrapping. As shown, SCK closely mimics real deployment scenarios. Most importantly here, a significant $\#flaps$ does not appear until 256-node deployment, hence mini-cluster extrapolation techniques will not work (Section 5.2). Figure 5.3a defines that $\#flaps$ depends on Φ [88]. Every node A maintains a Φ value for a peer node B (a total of $N \times (N-1)$ variables to monitor). If $\Phi > 8$ for B, A will declare B dead (a flap).

Figure 5.4b shows the maximum Φ values observed for every peer node. For example, for the 512-node setup, the whisker plots show the distribution of the maximum Φ values observed for each of the 512 nodes. As shown, the larger the cluster, more Φ values exceeds the threshold value of 8, hence the flapping. Figure 5.3b points that Φ depends on the average inter-arrival time of when new gossips about B arrives at A ($T_{avgGossip}$) and the time since A heard the last gossip about B ($T_{lastGossip}$); the “last gossip” is the last version number received (Section 5.1.1). The

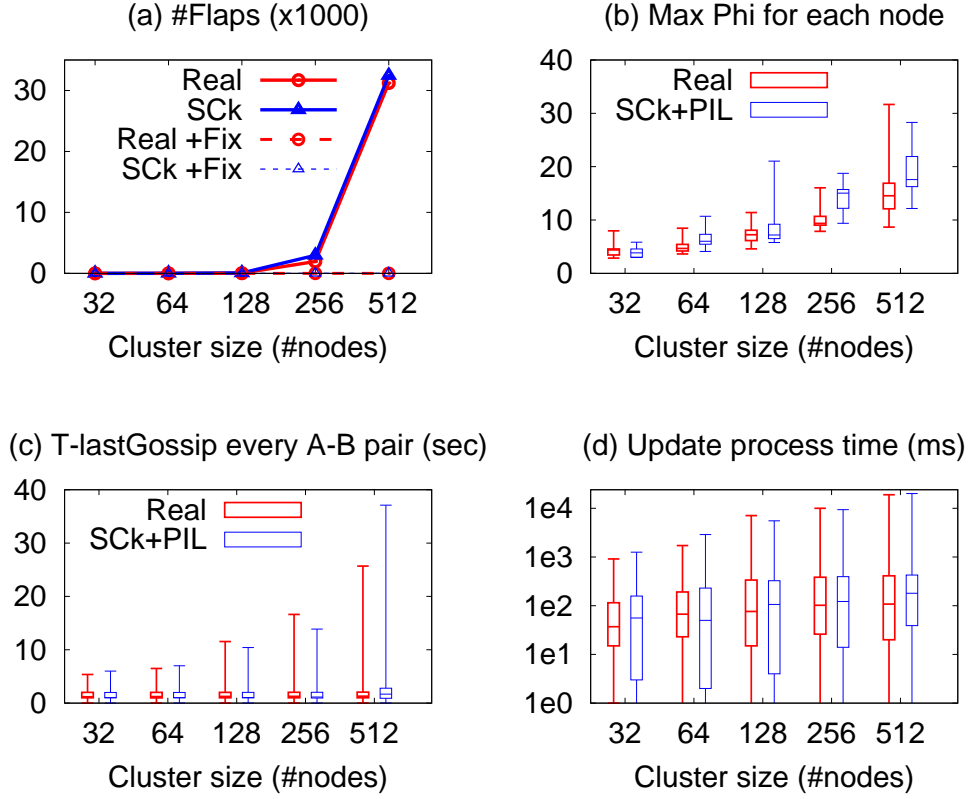


Figure 5.4: **Accuracy in reproducing CA-6127.** The figures represent the metrics presented in Figure 5.3, measured in real deployment (“Real”) and SCk with different cluster sizes (32, 64, 128, 256, and 512). Figure title represents the y-axis.

point is that $T_{lastGossip}$ should not be much higher than $T_{avgGossip}$.

Figure 5.4c shows the whisker plots of gossip inter-arrival times ($T_{lastGossip}$) that we collected for every A-B pair. For example, for the 512-node setup, the whisker plots represent the distribution of around 41 million gossip inter-arrival times; this large number is because a message contains gossips of many peer nodes. The figure shows that in larger clusters, new gossips do not arrive as fast as in smaller clusters, especially at high percentiles. Figure 5.3c shows that $T_{lastGossip}$ depends on how far B’s new gossips propagate through other nodes to A ($\#hops$) and the gossip processing time in each hop ($T_{gossipExec}$). The $\#hops$ is stable at $\log(N)$ on average in SCk and real deployment (not shown). The latter ($T_{gossipExec}$) is essentially state-update processing time ($T_{stateUpdate}$) whenever there are state changes, which is the culprit.

Bug#	Surface	Protocol	Metric
CA-6127 [13]	$N \geq 256$	Bootstrap	#flaps
CA-3831 [10]	≥ 256	Decommission	#flaps
CA-3881 [11]	≥ 64	Add new nodes	#flaps
CA-5456 [12]	≥ 256	Add new nodes	#flaps
RI-3926 [14]	≥ 128	Boot+rebalance	$T_{Complete}$
VO-1212 [15]	≥ 128	Rebalance	$T_{Complete}$

Table 5.2: **Reproduced bugs.** “Surface” implies the number of nodes needed for the bug symptom to surface. “c” stands for Cassandra, “r” for Riak, and “v” for Voldemort.

Figure 5.4d (in log scale) shows the whisker plots of state-update processing time ($T_{stateUpdate}$); in the 512-node setup, we measured around 25,000 state-update invocations. The figure shows that at high percentiles, $T_{stateUpdate}$ is scale-dependent. As explained in Figure 5.3d (and Section 5.1.1), $T_{stateUpdate}$ complicatedly depends on a scale-dependent 2-dimensional input ($Size_{ringTable}$ and $Size_{newStates}$); a node’s $Size_{ringTable}$ depends on how many nodes it knows, including the partition arrangement ($\leq N \times P$) and $Size_{newStates}$ ($\leq N$) increases as cluster size increases. Note that the $T_{stateUpdate}$ in SCK comes from the sampling-based time profiling (Section 5.3.1), which is relatively accurate as the figure shows.

We conclude that SCK mimics similar behaviors as in real deployments and is accurate for reproducing scalability bugs. As an additional note, we have applied the bug patch in both SCK and real deployment modes; Figure 5.4a shows #flaps is always zero in both modes.

5.5.3 Bugs Reproduced

Table 5.2 lists all the 7 bugs we have reproduced (4 Cassandra, 1 Riak, and 1 Voldemort bugs). We chose these 7 bugs (among the 41 bugs we studied) because the reports contain more detailed descriptions about the bugs, the affected protocols, the affected code version numbers, configuration setups, and the patches. Table 5.2 also shows the number of nodes needed for the bug symptoms to surface and the quantifiable metrics of the symptoms. Our first target system was Cassandra, hence the more bugs reproduced compared to Riak and Voldemort; the latter two were added for stronger proof of concept. Figure 5.5 shows the accuracy of SCK in reproducing the 6 bugs using

the metrics shown in Table 5.2. All Cassandra bugs have been described in Section 5.1.1 so we now briefly discuss the other two bugs in Voldemort and Riak.

Figure 5.5d: In VO-1212 [15], Voldemort’s rebalancing was not optimized for large clusters; it led to more stealer-donor partition transitions as the cluster size grows (128+ nodes). To fix this, the developers completely changed the stealer-donor partition transition algorithm.

Figure 5.5e: In RI-3926 [14], Riak bootstrapping employed a complex 3-stage rebalancing algorithm (claim-target, claim-hole, full-rebalance) that each node runs to eventually converge and achieve a perfect balance of the ring. Each node runs this CPU-intensive algorithm on *every* bootstrap gossip received. The larger the cluster, the longer time perfect balance is achieved (observed in 128+ nodes). For Riak, we profile the rebalancing time along with pre-memoization (with order determinism; Section 5.3.1-5.3.1). Figure 5.5f (similar to Figure 5.4d) compares the execution time of the rebalance function invocations in SCK and real deployments. The figure shows that SCK’s PIL exhibits a high accuracy.

We make several remarks from this experience. First, if SCK had existed in the first place, it might have *prevented* the Cassandra bugs; they all involve the same protocols (gossip, rebalance, and failure detector) and create the same symptom (high *#flaps*). These bugs highlight that code evolution can introduce new bugs in the same protocols. In this context, SCK is highly useful. Second, reproducing scalability bugs is relatively *easy* as we achieve a high colocation factor. Unlike non-deterministic bugs which require complex timing reordering to reproduce [81, 107], symptoms of scalability bugs are “deterministically scale-dependent.” Third, different systems of the same type (*e.g.*, key-value store) implement similar protocols. The *generality* of SCK methods in scale-checking the protocols above can be useful to many other distributed key-value stores.

5.5.4 New Bugs

We also scale-checked the latest stable versions of Cassandra (v2.2.5), Riak (v2.1.3), and Voldemort (v1.10.21). In Cassandra, SCK shows that cluster-wide flapping resurfaces again but only

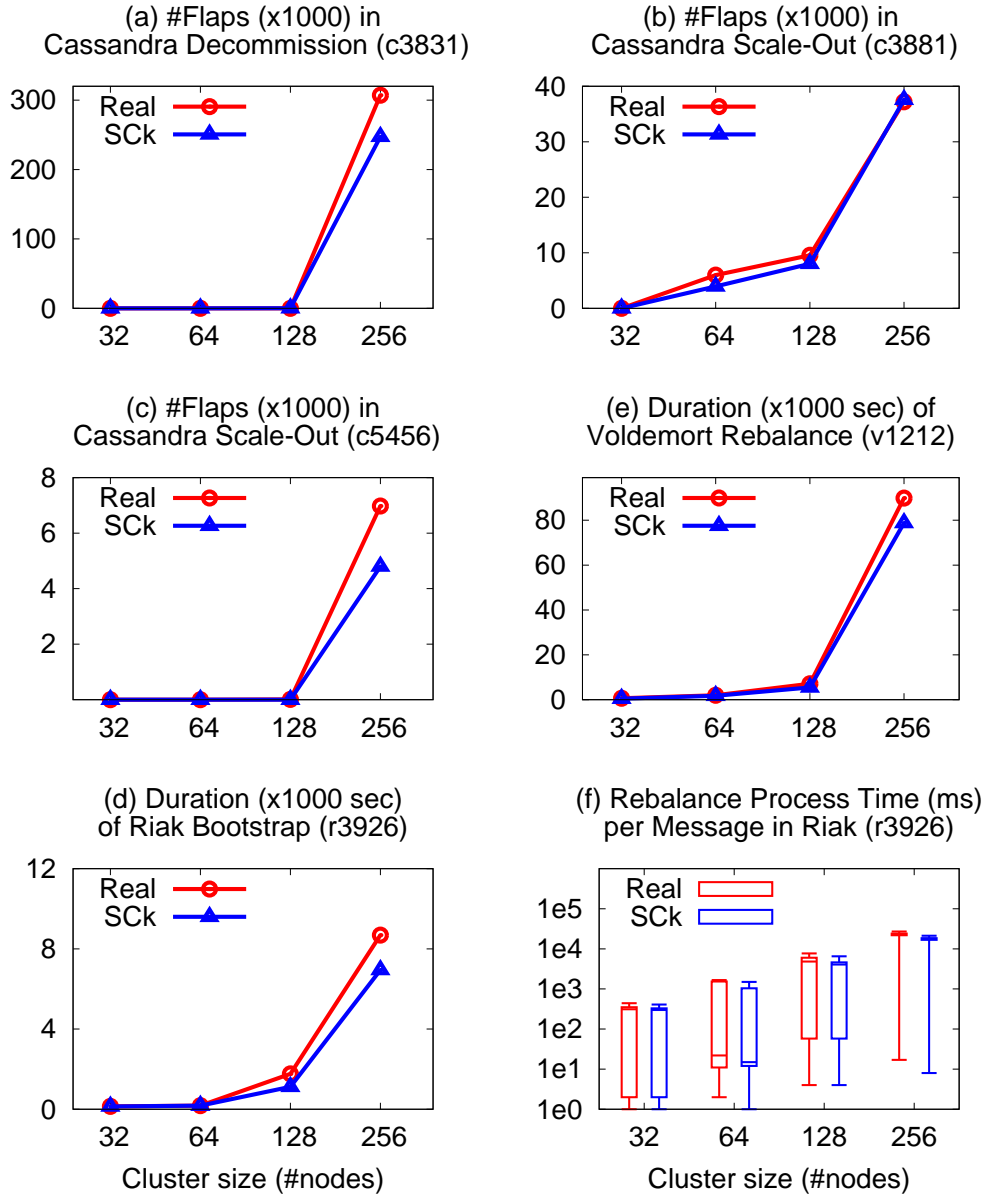


Figure 5.5: **Accuracy in reproducing other bugs.** The figures represent the bugs described in Table 5.2. The title represents the y-axis. We cap the y-axis to show the scale at which the bug symptoms start to appear.

observable in 512-node deployment (*e.g.*, decommissioning only one node caused almost 100,000 flaps). We submitted the bug few months back and it is still unresolved (the fix might require new design). Meanwhile, the developers suggested us to add/remove node one at a time with 2-minute separation, which means scaling-out/down 100 nodes will take over 3 hours; instant elasticity is not achievable. For Riak and Voldemort, we found that their latest-stable bootstrap/rebalance protocols do not exhibit any scalability bug, up to 512 nodes.

5.6 Conclusion

In this chapter, we have presented our observations on scalability bugs that highlight a need of an attention to combat them. Scalability bugs are latent bugs that are scale-dependent that only manifest in large scale. We have also presented our pilot work SCK, a methodology to enable developers to colocate hundreds of nodes on one machine to emulate large-scale deployments; this helps developers save cost of testing and speed up testing process. We have introduced four techniques in this chapter:

1. **Processing Illusion** (PIL) helps reduce CPU contention so CPU-intensive nodes can be more colocated on one machine.
2. **Single Process Cluster** (SPC) helps reduce memory consumption and context switching.
3. **Global Event Driven Architecture** (GEDA) helps reduce the number of threads we need to run the systems.
4. **Memory Footprint Reduction** (MFR) helps reduce memory consumption further.

SCK is a pilot work, it still needs a lot of manual efforts from developers to do scale check, but we hope that this work can raise awareness from system community to pay attention to this new class of bugs in cloud-scale distributed systems.

CHAPTER 6

CONCLUSIONS AND FUTURE WORK

In this dissertation, we aim to strengthen dependability of cloud-scale distributed systems by addressing two new types of bugs in cloud systems, distributed concurrency bugs (DC bugs) and scalability bugs. We have performed bug studies to gain insights about the nature of these bugs, and we have also advanced state of the art of system testing. This chapter concludes this dissertation work and discusses future work in combating DC bugs and scalability bugs.

6.1 Conclusion

6.1.1 Distributed Concurrency Bugs

The first problem we focus in this dissertation is DC bugs. We have conducted in-depth study and created the largest and most comprehensive of DC bugs named TaxDC. We categorize DC bugs in three dimensions. The first dimension is triggering which is conditions that makes bugs happen. We studied timing conditions and found four main timing patterns: order violation, atomicity violation, fault timing, and reboot timing. We also studied input conditions that are ingredients for bugs to surface. We found that most DC bugs will surface only systems execute multiple protocols, and more than 50% of bugs surface in recovery protocols (*i.e.*, the bugs surface only when there are hardware failures).

The second dimension that we studied is errors and failures. We studied the first errors that happen immediately after the bugs are triggered. We see half of the bugs have local errors that we can see the errors by observing only triggering node, but half of them have global errors that require us to observe the whole systems to notice the errors. Moreover, we looked into failure symptom induced by DC bugs and found that the bugs can lead to severe failures like system downtime, operation failures, data loss/corruption/inconsistencies, and performance degradation.

Lastly, the third dimension that we studied is fixes that are how developers fix the DC bugs. We

saw two main strategies to fix the bugs that are fixing the timing and fixing the handling. For timing fixes, developers can do it globally or locally (global synchronization or local synchronization). For handling fixes, developers change the logic of message handling or fault handling such that the systems still behave correctly.

Other than the bug study, we have introduced semantic-aware model checking (SAMC) that is a white-box approach to model check the systems. SAMC prunes out some executions because it knows that those executions are redundant with previous executions it already tested by using semantic knowledge of target systems. We show a strong case that SAMC can elegantly address state-space-explosion problem. We have introduced four novel semantic-aware reduction policies, and built SAMPRO and integrated it to three systems including Hadoop MapReduce, Cassandra and ZooKeeper. On average, SAMC can find bugs 49x faster than other states of the art.

6.1.2 Scalability Bugs

The second problem we focus in this dissertation is scalability bugs. Scalability bugs are bugs that specific to cloud systems and we found there is not much attention paid on them. We observed that scalability bugs are scale dependent and only surface at extreme scale (*e.g.*, hundreds of nodes). We noticed that although the systems are designed to be scalable, but actual implementations can introduces the bugs. We also saw that not all developers can afford large clusters to check scalability of their code, and make bugs linger until the systems are deployed on large scale.

Our observations highlight the need for scale checking that check the implementation of the systems. Hence, we have introduced SCK, a scale-checking methodology to allow developers colocate hundreds nodes on a single machine to check their systems. We introduced four techniques to mitigate resource contentions issue (*i.e.*, CPU, memory, and threads) regarding to colocating several nodes in one machine. We adopted SCK to three systems including Cassandra, Riak, and Voldemort, and were able to reproduce six old scalability bugs with high accuracy.

6.2 Future Work

We now discuss our research impact and future research revenues in combating DC bugs and scalability bugs. Regarding DC-bug combating, although many high-level directions can be adopted from work on LC bugs, there are many interesting challenges and opportunities unique to DC bugs. For scalability bugs, our pilot work just targets a subset of scalability bugs, there are many classes of scalability bugs that we have not touched yet.

6.2.1 *Distributed Concurrency Bugs*

Fault Paths and Multi-Protocol Interactions

Individual protocols tend to be robust in general. Only 18 DC bugs occur in *individual* protocols *without* any input fault condition; only 8 of them are in foreground protocols. On the other hand, a large majority of DC bugs happen due to concurrent executions of multiple protocols and/or different fault timings (Finding #2). This has a tremendous implication to input testing: *all types of verification, testing, and analysis approaches must consider fault injections and multiple protocols as input conditions*. Although recent work has paid attention to this [77, 100, 161], we emphasize that all forms of faults (Section 3.2.6) must be exercised.

Distributed Systems Model Checkers

Assuming the necessary input conditions are exercised, the next question is: can we test different event re-orderings to hit the triggering timing (Section 3.3.1)? This is the job of distributed system model checkers (dmck), which are gaining popularity recently [81, 103, 143, 157]. Dmck works by intercepting distributed events and permuting their ordering. The more events included, the more scalability issues will arise due to state-space explosion. To date, *no dmck completely controls the timings of all necessary events* that might contribute to the triggering timing (Finding #1). MaceMC [103] only reorders messages and network disconnections. MoDist [81] exercises

timeouts and Demeter [81] intercepts messages and local computation but they do not explore different timing of multiple crashes and reboots. Also, none of the above include storage faults or timing issues [86]. Therefore, continued research on scalable exploration algorithms is needed, specifically when *all* the necessary events need to be controlled. This could be helped by DC bugs’ triggering scope characteristics (Finding #3), just like that in LC model checkers [127].

Domain-Specific Specifications

Now, assuming the necessary events are controlled, the next question is: do we have the specification to judge the manifestation of a bug? This is a plague for many tools. For example, Demeter does not find new bugs [81]. Conversations with the authors suggest that their target systems do not deploy detailed specifications, and thus some bugs are left uncaught. Deploying generic “textbook” specifications (*e.g.*, “only one leader exists”) does not help as they could lead to false positives (*e.g.*, ZooKeeper allows two leaders at a single point in time). Many research papers on specifications only deploy few of them [77, 115, 136]. Developers also bemoan the hard-to-debug fail-silent problems MR-3634 and prefer to see easier-to-debug fail-stop bugs.

On the positive side, 53% of DC bugs lead to explicit first errors (Finding #4), implying that sanity checks already in software can be harnessed as specifications (more in Section 6.2.1). On the other side, compared to single-machine systems, distributed systems are much more capable of masking errors. Therefore, these error specifications have to be used with caution to avoid false positives. Furthermore, 47% of DC bugs lead to silent first errors (Finding #4). Many of them proceed to “silent failures”, such as data loss, node hangs, etc. Even if they become explicit errors later, these explicit errors could be far away from the initial triggering conditions (*e.g.*, Figure 3.4). In short, *no matter how sophisticated the tools are, they are ineffective without accurate specifications*. This motivates the creation or inference of local specifications that can show early errors or symptoms of DC bugs.

Bug Detection Tools

We now discuss bug detection tools, which are unfortunately rare for DC bugs, although very popular for LC bugs [40, 67, 92, 119, 127, 138]. Bug detection tools look for bugs that match specific patterns. They cannot provide bug-free proof, but can be efficient in discovering bugs when guided by the right patterns. Our study provides guidance and patterns that can be exploited by future DC bug detection.

Generic detection framework. Finding #1 implies that detecting DC bugs, particularly message-timing DC bugs, should focus on two key tasks: (1) obtaining timing specifications, including order and atomicity specifications among messages and computation; and (2) detecting violations to these specifications through dynamic or static analysis.

Invariant-guided detection. Likely program invariants can be learned from program behaviors, and used as specifications in bug detection [64, 65, 119]. The key challenge is to design simple and suitable invariant templates. For example, “function F_1 should always follow F_2 ” is a useful template for API-related semantic bugs [64]; “the atomicity of accesses a_1 and a_2 should never be violated” is effective for LC bugs [119]. Finding #1 about triggering timing and Finding #4 about error patterns provide empirical evidence that these templates can be effective for DC bugs: “message bc should arrive at C before message ac (ca) arrives (leaves)”; “message ab should never arrive in the middle of event e on node B ”; and “message ab should always be replied”.

Misconception-guided bug detection. Knowing programmers’ misconceptions can help bug detectors focus on specifications likely to be violated. LC bug researchers have leveraged misconceptions such as “two near-by reads of the same variable should return the same value” [119] and “a condition checked to be true should remain true when used” [134]. Finding #6 reveals that common misconceptions, such as “a single hop is faster than double hops”, “local computation is faster than remote computation”, “atomic blocks cannot be broken” can help DC bug detection.

Error-guided bug detection. Finding #4 shows that many DC bugs lead to explicit local/global errors, which implies that timing specifications for many DC bugs can be inferred backward based on explicit errors. For example, program analysis may reveal that a state-machine exception e will arise whenever C receives message ac before bc , which provides a timing specification (ac arrives before bc) whose violation leads to a *local error*; or, the analysis may reveal that exception e arises whenever node B receives a message cb from node C and C only sends cb when ac arrives at C before bc , which provides a timing specification whose violation leads to a *wrong-message global error*; and so on.

Software testing. Testing takes a quarter of all software development resources, and is crucial in exposing bugs before code release. Although many testing techniques have been proposed for LC bugs [43, 132, 139], there have been few for DC bugs [140]. Finding #2 implies that test input design has to consider faults, concurrent protocols, and background protocols. Finding #3 implies that *pairwise testing*, which targets every pair of message ordering, every pair of protocol interaction, and so on, will work much more effectively than *all combination testing*, which exercises all possible total orders and interactions of all messages and all protocols. For example, a large number of DC bugs (Figure 3.3d-f) can be found with inputs of at most two protocols, crashes and reboots.

Failure Diagnosis

Given failure symptoms, distributed systems developers have to reason about many nodes to figure out the triggering and root cause of a failure. Our study provides guidance to this challenging process of failure diagnosis.

Slicing/dependence analysis. Identifying which instructions can affect the outcome of an instruction i is a widely used debugging technique for deterministic sequential bugs. However, it cannot scale to the whole distributed systems, and hence is rarely used. Finding #3 indicates that

most DC bugs have deterministic error propagation; Finding #4 shows that many DC bugs have their errors propagate through missing or wrong messages. Therefore, per-node dependence analysis that can quickly identify whether the generation of a local error depends on any incoming messages would help DC bug failure diagnosis to get closer and closer to where the triggering events happen.

Error logging. Error logging is crucial in failure diagnosis. If the first error of a DC bug is an explicit local error, the error log can help developers quickly identify the triggering node and focus their diagnosis on one node. Finding #4 unfortunately shows that only 23% of DC bugs lead to explicit local errors. This finding motivates future tool to help make more DC bugs lead to explicit local errors.

Statistical debugging. Comparing success-run traces with failure-run traces can help identify failure predictors for semantic bugs [110] and concurrency bugs [97] in single-machine software. The key design question is what type of program properties should be compared between failure and success runs. For example, branch outcomes are compared for diagnosing semantic bugs but not for LC bugs. Finding #1 and #3 about triggering timing conditions provide guidance for applying this approach for DC bugs. We can collect all message sending/arrival time at runtime, and then find rare event orderings that lead to failures by contrasting them with common “healthy” orderings (*e.g.*, Figure 3.1b happens 99.99% of the time while Figure 3.1a happens 0.01% of the time). Of course, there are challenges. Finding #2 and #3 show that many DC bugs come from the interactions of many protocols. Thus, it is not sufficient to only log a chain of messages originated from the same request, a common practice in request logging [52]. Furthermore, some DC bugs are triggered by message-computation ordering. Therefore, logging messages alone is not sufficient.

Record and Replay. Debugging LC concurrency bugs with record and deterministic replay is a popular approach [133, 149]. However, such an approach has not permeated practices in dis-

tributed systems debugging. A ZooKeeper developer pointed us to a fresh DC bug that causes a whole-cluster outage but has not been fixed for months because the deployment logs do not record enough information to replay the bug (ZK-2172). There has been 9 back-and-forth log changes and attachments with 72 discussion comments between the bug submitter and the developers. More studies are needed to understand the gap between record-replay challenges in practice and the current state of the art [72, 116].

Failure Prevention and Fixing

Runtime Prevention. The manifestation of concurrency bugs can sometimes be prevented by injecting delays at runtime. This technique has been successfully deployed to prevent LC bugs based on their timing conditions [102, 120, 164]. Finding #1 shows that many DC bugs are triggered by untimely messages and hence can potentially be prevented this way. For example, none of the bugs shown in Figure 3.1a–h would happen if we delay a message arrival/sending or local computation. Of course, different from LC bugs, some of these delays have to rely on a network interposition layer; similar with LC bugs, some delays may lead to hangs, and hence cannot be adopted.

Bug Fixing. Recent work automatically fixes LC bugs by inserting lock/unlock or signal/wait to prohibit buggy timing [98, 114, 152]. Finding #5 shows that the same philosophy is promising for 30% of studied DC bugs. Our study shows that this approach has to be tweaked to focus on using global messages (*e.g.*, ACKs) or local operation re-ordering, instead of lock or signal, to fix DC bugs. Finding #5 indicates that 40% of those DC bugs are fixed by shifting message handlers, ignoring messages, and canceling computation, without adding new computation logic. This presents a unique opportunity for developing new and more aggressive fixing techniques.

HB-9095:

1. RS successfully OPENED region R,
2. RS notifies ZK that region R is OPENED,
3. ZK continues region R state msg to Master,
4. Master starts processing OPENED msg,
5. Meanwhile RS CLOSED region R (asked by Client),
6. RS notifies ZK that region R is CLOSED,
7. Master asks ZK to delete znode for region R, *concurrently racing with step 6!*,
8. ZK deletes region R's znode,
9. Master never assigns region R to any RS. *R becomes an orphan!*

Figure 6.1: **Race of HBase's messages to ZooKeeper.**

Distributed Transactions

In the middle of our study, we ask ourselves: if DC bugs can be theoretically solved by distributed transactions, why doesn't such technique eliminate DC bugs in practice? Our answers are: first, the *actual implementations of theoretically-proven distributed transactions are not always correct* (as also alluded in other work [44, 129]). For example, new DC bugs continue to surface in complex distributed transactions such as ZooKeeper's ZAB and Cassandra's Paxos as they are continuously modified. Second, *distributed transactions are only a subset of a full complete system*. A prime example is the use of ZooKeeper in HBase for coordinating and sharing states between HBase masters and region servers. Although ZooKeeper provides linearization of updates, HBase must handle its concurrent operations to ZooKeeper, for example, step 6 and 7 in Figure 6.1; there are many other similar examples. Put simply, there are many protocols that do not use distributed transactions, instead they use domain-specific finite state machines, which should be tested more heavily.

Another approach to eliminate non-deterministic bugs in distributed protocols is by building deterministic distributed systems. However, the technique is still in its infancy, at least in terms of

the impact to performance (*e.g.*, an order of magnitude of overhead [93]).

Verifiable Frameworks

Recently there is a growing work on new programming language frameworks for building verifiable distributed systems [60, 87, 154], but they typically focus on the main protocols and not the full system including the background protocols. One major challenge is that just for the basic read and write protocols, the length of the proofs can reach thousands of lines of code, potentially larger than the protocol implementation. Unfortunately, our study shows that the complex interaction between foreground and background protocols can lead to DC bugs. Therefore, for complete real-world systems, verification of the entire set of the protocols is needed.

6.2.2 Scalability Bugs

SCK is an initial effort to combat scalability bugs. It focuses on scale-dependent CPU/processing time. However, there are other scaling problems that lead to I/O and memory contentions [78, 130, 142], usually caused by the scale of load [39, 82] or data size [128]. SCK cannot reproduce such issues on one machine as we do not address memory and IO emulation. Moreover, we find that some bugs are caused by scale of failures [78] (*e.g.*, a great number of machines fails at the same time) and these bugs are hard to catch during testing process because failures in cloud-scale distributed systems can be very complex. We believe there are many open problems to solve in this new research area. We will discuss some possible research directions here.

Program Analysis

Another approach to scale check systems on just one machine is program analysis. There is a previous work adopting static analysis to check if software is scalable on multicore processors [54]. However, cloud-scale distributed systems are different from multicore software, and the analysis

for multicore software is not applicable for cloud systems. In this dissertation, we show that scalability bugs in cloud distributed systems are caused from scale-dependent CPU/processing time. Building a program analysis that covers all paths and understands the cascading impacts without false positives is challenging. Not all scale-dependent loops imply buggy code. For example, in Cassandra gossip protocol, if Cassandra processes gossips in a multi-threaded manner, the long processing time might not cascade to failures.

However, building program analysis to point out potential buggy functions is still useful. It can reduce manual efforts in order to identify PIL-safe and offending functions which makes SCk more automatic.

In-Production Checking

If scale checking on one machine is hard and checking on real scale is expensive, *can we piggyback scale checking on production cluster?* The idea of this in-production checking is that we already have ready-to-check environment (*i.e.*, large scale setup and big data stored in clusters), we should be able to test critical scenarios such as machine decommissions, a surge of requests, and datacenter failures. Scalability bugs in these critical scenarios are unlikely to be covered in offline testing, but are possible to be detected in production systems, however, in-production checking remains a “controversial” idea, mainly because of the soundness of the checking process. No service provider would like to report to their clients “an in-production checking that we scheduled has caused an outage/data loss/performance disruption.” The risk is too high. An initial idea of in-production checking is proposed [106] and discussed, but there is still no formal research on this topic.

6.3 Closing Words

Cloud-scale distributed systems are the leading actors in the age of cloud computing. The dependability of cloud services mainly rely on these distributed systems. This dissertation addresses an essential question of cloud dependability: “*how can we make cloud-scale distributed systems more*

dependable?” Our work focuses on two unsolved problems, distributed concurrency bugs and scalability bugs. These two types of bugs are novel and unique to cloud-scale environment, and we found not many works addressing them. We address these two types of bugs by conducting formal studies to understand the nature of the bugs and advancing state of the art of system checking. We hope this dissertation brings new insights on combating distributed concurrency bugs and draws more attentions for scalability bugs. However, there are still many new classes of cloud-unique bugs waiting for system community to address in the future.

REFERENCES

- [1] <http://ucare.cs.uchicago.edu/projects/cbs/>.
- [2] Adobe Creative Cloud Adoption Grows to 9 Million Paid Members. <http://prodesigntools.com/creative-cloud-one-million-paid-members.html>.
- [3] Apache Cassandra. https://en.wikipedia.org/wiki/Apache_Cassandra.
- [4] Apache Cassandra. <http://cassandra.apache.org>.
- [5] Apache Hadoop. <http://hadoop.apache.org>.
- [6] Apache HBase. <http://hbase.apache.org>.
- [7] Apache ZooKeeper. <http://zookeeper.apache.org>.
- [8] AWS Customer Success. <https://aws.amazon.com/solutions/case-studies/>.
- [9] Big data storage: Hadoop storage basics. <http://www.computerweekly.com/feature/Big-data-storage-Hadoop-storage-basics>.
- [10] BUG: CASSANDRA-3831: scaling to large clusters in GossipStage impossible due to calculatePendingRanges. <https://issues.apache.org/jira/browse/CASSANDRA-3831>.
- [11] BUG: CASSANDRA-3881: reduce computational complexity of processing topology changes. <https://issues.apache.org/jira/browse/CASSANDRA-3881>.
- [12] BUG: CASSANDRA-5456: Large number of bootstrapping nodes cause gossip to stop working. <https://issues.apache.org/jira/browse/CASSANDRA-5456>.
- [13] BUG: CASSANDRA-6127: vnodes don't scale to hundreds of nodes. <https://issues.apache.org/jira/browse/CASSANDRA-6127>.
- [14] BUG: RIAK: Large ring_creation_size. http://lists.basho.com/pipermail/riak-users_lists.basho.com/2011-April/003895.html.
- [15] BUG: VOLDEMORT: Number of partitions. <https://groups.google.com/forum/#!msg/project-voldemort/3vrZfZgQp2Y/Uqt8NgJHg4AJ>.
- [16] Database Scalability: Vertical Scaling vs Horizontal Scaling. <http://www.vcloudnews.com/database-scalability-vertical-scaling-vs-horizontal-scaling/>.
- [17] Dropbox. <https://www.dropbox.com/>.
- [18] Gmail. <https://www.google.com/gmail/about>.
- [19] Gmail Now Has More Than 1B Monthly Active Users. <https://techcrunch.com/2016/02/01/gmail-now-has-more-than-1b-monthly-active-users/>.

- [20] Google Drive. <https://www.google.com/intl/en/drive/>.
- [21] Google Drive claims one million paying customers, er, organizations. <http://fortune.com/2015/09/21/google-drive-1m-paid-users/>.
- [22] Gossip 2.0. <https://issues.apache.org/jira/browse/CASSANDRA-12345>.
- [23] Gossip 2.0. http://mail-archives.apache.org/mod_mbox/cassandra-dev/201609.mbox/%3CCAHjqPuJMkfZwp9DDX45PNBNhkoGXsPW4TFT6Zxv%2BTTz_Pg3Y%2Bg%40mail.gmail.com%3E.
- [24] iCloud. <https://www.icloud.com/>.
- [25] NMC PRObE Nome Nodes. <https://www.nmc-probe.org/wiki/Nome:Nodes>.
- [26] Number of registered Dropbox users from April 2011 to March 2016 (in millions). <https://www.statista.com/statistics/261820/number-of-registered-dropbox-users/>.
- [27] Project Voldemort. <http://www.project-voldemort.com/voldemort/>.
- [28] Riak. <http://basho.com/products/riak-kv>.
- [29] Siri. <https://www.apple.com/ios/siri/>.
- [30] Taking astronomy to the cloud. <https://www.inthefieldstories.net/taking-astronomy-to-the-cloud/>.
- [31] The Top 20 Valuable Facebook Statistics - Updated May 2017. <https://zephoria.com/top-15-valuable-facebook-statistics/>.
- [32] Twenty-One Experts Define Cloud Computing. <http://cloudcomputing.sys-con.com/node/612375/print>.
- [33] What is cloud computing? <https://www.ibm.com/cloud-computing/learn-more/what-is-cloud-computing/>.
- [34] What Is Cloud Computing? <http://www.pcmag.com/article2/0,2817,2372163,00.asp>.
- [35] Running Netflix on Cassandra in the Cloud. <https://www.youtube.com/watch?v=97VBdgIgcCU>, 2013.
- [36] Why the world's largest Hadoop installation may soon become the norm. <http://www.techrepublic.com/article/why-the-worlds-largest-hadoop-installation-may-soon-become-the-norm/>, 2014.
- [37] Daniel Abadi. Data Management in the Cloud: Limitations and Opportunities. *IEEE Data Engineering Bulletin*, 32(1):3–12, March 2009.

- [38] Radu Banabic and George Candea. Fast black-box testing of system recovery code. In *Proceedings of the 2012 EuroSys Conference (EuroSys)*, 2012.
- [39] Peter Bodik, Armando Fox, Michael Franklin, Michael Jordan, and David Patterson. Characterizing, Modeling, and Generating Workload Spikes for Stateful Services. In *Proceedings of the 1st ACM Symposium on Cloud Computing (SoCC)*, 2010.
- [40] Michael D. Bond, Katherine E. Coons, and Kathryn S. McKinley. PACER: Proportional Detection of Data Races. In *PLDI*, 2010.
- [41] Stefan Bucur, Vlad Ureche, Cristian Zamfir, and George Candea. Parallel Symbolic Execution for Automated Real-World Software Testing. In *Proceedings of the 2011 EuroSys Conference (EuroSys)*, 2011.
- [42] J. R. Burch, E. M. Clarke, K L. McMillan, D L. Dill, and L J. Hwang. Symbolic model checking: 10^{20} states and beyond. *Information and Computation*, 98(2):142–170, June 1992.
- [43] Sebastian Burckhardt, Pravesh Kothari, Madanlal Musuvathi, and Santosh Nagarakatte. A Randomized Scheduler with Probabilistic Guarantees of Finding Bugs. In *ASPLOS*, 2010.
- [44] Mike Burrows. The Chubby lock service for loosely-coupled distributed systems. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI)*, 2006.
- [45] Rajkumar Buyya, Srikumar Venugopal Chee Shin Yeo and, James Broberg, and Ivona Brandic. Cloud Computing and Emerging IT Platforms: Vision, Hype, and Reality for Delivering Computing As the 5th Utility. In *Future Generation Computer Systems*, 2009.
- [46] Cristian Cadar, Daniel Dunbar, and Dawson R. Engler. KLEE: Unassisted and Automatic Generation of High-Coverage Tests for Complex Systems Programs. In *Proceedings of the 8th Symposium on Operating Systems Design and Implementation (OSDI)*, 2008.
- [47] Alexandru Calotoiu, Torsten Hoefler, Marius Poke, and Felix Wolf. Using Automated Performance Modeling to Find Scalability Bugs in Complex Codes. In *Proceedings of International Conference on High Performance Computing, Networking, Storage and Analysis (SC)*, 2013.
- [48] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Michael Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. Bigtable: A Distributed Storage System for Structured Data. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation (OSDI)*, 2006.
- [49] John Chapin, Mendel Rosenblum, Scott Devine, Tirthankar Lahiri, Dan Teodosiu, and Anoop Gupta. Hive: Fault Containment for Shared-Memory Multiprocessors. In *Proceedings of the 15th ACM Symposium on Operating Systems Principles (SOSP)*, 1995.

- [50] Vitaly Chipounov, Volodymyr Kuznetsov, and George Candea. S2E: a platform for in-vivo multi-path analysis of software systems. In *Proceedings of the 16th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2011.
- [51] Andy Chou, Junfeng Yang, Benjamin Chelf, Seth Hallem, and Dawson Engler. An Empirical Study of Operating System Errors. In *Proceedings of the 18th ACM Symposium on Operating Systems Principles (SOSP)*, 2001.
- [52] Michael Chow, David Meisner, Jason Flinn, Daniel Peek, and Thomas F. Wenisch. The Mystery Machine: End-to-end Performance Analysis of Large-scale Internet Services. In *Proceedings of the 11th Symposium on Operating Systems Design and Implementation (OSDI)*, 2014.
- [53] Edmund M. Clarke, E. Allen Emerson, Somesh Jha, and A. Prasad Sistla. Symmetry reductions in model checking. In *10th International Conference on Computer Aided Verification (CAV)*, 1998.
- [54] Austin T. Clements, M. Frans Kaashoek, Nickolai Zeldovich, Robert T. Morris, , and Eddie Kohler. The Scalable Commutativity Rule: Designing Scalable Software for Multicore Processors. In *Proceedings of the 24th ACM Symposium on Operating Systems Principles (SOSP)*, 2013.
- [55] Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears. Benchmarking Cloud Serving Systems with YCSB. In *Proceedings of the 1st ACM Symposium on Cloud Computing (SoCC)*, 2010.
- [56] James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, JJ Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, Wilson Hsieh, Sebastian Kanthak, Eugene Kogan, Hongyi Li, Alexander Lloyd, Sergey Melnik, David Mwaura, David Nagle, Sean Quinlan, Rajesh Rao, Lindsay Rolig, Yasushi Saito, Michal Szymaniak, Christopher Taylor, Ruth Wang, and Dale Woodford. Spanner: Google’s Globally-Distributed Database. In *Proceedings of the 10th Symposium on Operating Systems Design and Implementation (OSDI)*, 2012.
- [57] Heming Cui, Gang Hu Columbia, Jingyue Wu, and Junfeng Yang. Verifying Systems Rules Using Rule-Directed Symbolic Execution. In *Proceedings of the 18th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2013.
- [58] Jeffrey Dean and Sanjay Ghemawat. MapReduce: Simplified Data Processing on Large Clusters. In *Proceedings of the 6th Symposium on Operating Systems Design and Implementation (OSDI)*, 2004.
- [59] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Voshall, and Werner Vo-

- gels. Dynamo: Amazon’s Highly Available Key-value Store. In *Proceedings of 21st ACM SIGOPS symposium on Operating systems principles (SOSP)*, 2007.
- [60] Ankush Desai, Vivek Gupta, Ethan Jackson, Shaz Qadeer, Sriram Rajamani, and Damien Zufferey. P: Safe Asynchronous Event-Driven Programming. In *Proceedings of the ACM SIGPLAN 2013 Conference on Programming Language Design and Implementation (PLDI)*, 2013.
 - [61] Thanh Do. Towards Reliable Cloud Systems. Doctoral Dissertation, University of Wisconsin-Madison (UMI 3632947), 2014.
 - [62] Thanh Do, Mingzhe Hao, Tanakorn Leesatapornwongsa, Tiratat Patana-anake, and Haryadi S. Gunawi. Limplock: Understanding the Impact of Limpware on Scale-Out Cloud Systems. In *Proceedings of the 4th ACM Symposium on Cloud Computing (SoCC)*, 2013.
 - [63] E. Allen Emerson, Somesh Jha, and Doron Peled. Combining Partial Order and Symmetry Reductions. In *The 3rd International Workshop on Tools and Algorithms for Construction and Analysis of Systems (TACAS)*, 1997.
 - [64] Dawson Engler, David Yu Chen, Seth Hallem, Andy Chou, and Benjamin Chelf. Bugs as Deviant Behavior: A General Approach to Inferring Errors in Systems Code. In *SOSP*, 2001.
 - [65] Michael D. Ernst, Adam Czeisler, William G. Griswold, and David Notkin. Quickly Detecting Relevant Program Invariants. In *ICSE*, 2000.
 - [66] Loek Essers. Cloud Failures Cost More Than \$70 Million Since 2007, Researchers Estimate. <http://www.pcworld.com>, 2012.
 - [67] Cormac Flanagan and Stephen N. Freund. FastTrack: Efficient and Precise Dynamic Race Detection. In *PLDI*, 2009.
 - [68] Cormac Flanagan and Patrice Godefroid. Dynamic Partial-Order Reduction for Model Checking Software. In *Proceedings of the 33rd ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL)*, 2005.
 - [69] Pedro Fonseca, Cheng Li, Vishal Singhal, and Rodrigo Rodrigues. A Study of the Internal and External Effects of Concurrency Bugs. In *DSN*, 2010.
 - [70] Ian Foster, Young Zhao, Ioan Raicu, and Shiyong Lu. Cloud Computing and Grid Computing 360-Degree Compared. In *2008 Grid Computing Environments Workshop*, 2008.
 - [71] Dennis Geels, Gautam Altekar, Petros Maniatis, Timothy Roscoe, and Ion Stoica. Friday: Global Comprehension for Distributed Replay. In *Proceedings of the 4th Symposium on Networked Systems Design and Implementation (NSDI)*, 2007.

- [72] Dennis Geels, Gautam Altekar, Scott Shenker, and Ion Stoica. Replay Debugging for Distributed Applications. In *Proceedings of the USENIX Annual Technical Conference (ATC)*, 2006.
- [73] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The Google File System. In *Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP)*, 2003.
- [74] Patrice Godefroid. Partial-Order Methods for the Verification of Concurrent Systems: An Approach to the State-Explosion Problem. volume 1032, 1996.
- [75] Rachid Guerraoui and Maysam Yabandeh. Model Checking a Networked System Without the Network. In *Proceedings of the 8th Symposium on Networked Systems Design and Implementation (NSDI)*, 2011.
- [76] Haryadi S. Gunawi, Thanh Do, Joseph M. Hellerstein, Ion Stoica, Dhruba Borthakur, and Jesse Robbins. Failure as a Service (FaaS): A Cloud Service for Large-Scale, Online Failure Drills. UC Berkeley Technical Report UCB/EECS-2011-87.
- [77] Haryadi S. Gunawi, Thanh Do, Pallavi Joshi, Peter Alvaro, Joseph M. Hellerstein, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, Koushik Sen, and Dhruba Borthakur. FATE and DESTINI: A Framework for Cloud Recovery Testing. In *Proceedings of the 8th Symposium on Networked Systems Design and Implementation (NSDI)*, 2011.
- [78] Haryadi S. Gunawi, Mingzhe Hao, Tanakorn Leesatapornwongsa, Tiratat Patana-anake, Thanh Do, Jeffry Adityatama, Kurnia J. Eliazar, Agung Laksono, Jeffrey F. Lukman, Vincentius Martin, and Anang D. Satria. What Bugs Live in the Cloud? A Study of 3000+ Issues in Cloud Systems. In *Proceedings of the 5th ACM Symposium on Cloud Computing (SoCC)*, 2014.
- [79] Haryadi S. Gunawi, Mingzhe Hao, Riza O. Suminto, Agung Laksono, Anang D. Satria, Jeffry Adityatama, and Kurnia J. Eliazar. Why Does the Cloud Stop Computing? Lessons from Hundreds of Service Outages. In *Proceedings of the 7th ACM Symposium on Cloud Computing (SoCC)*, 2016.
- [80] Haryadi S. Gunawi, Cindy Rubio-Gonzalez, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, and Ben Liblit. EIO: Error Handling is Occasionally Correct. In *Proceedings of the 6th USENIX Symposium on File and Storage Technologies (FAST)*, 2008.
- [81] Huayang Guo, Ming Wu, Lidong Zhou, Gang Hu, Junfeng Yang, and Lintao Zhang. Practical Software Model Checking via Dynamic Interface Reduction. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP)*, 2011.
- [82] Zhenyu Guo, Sean McDirmid, Mao Yang, Li Zhuang, Pu Zhang, Yingwei Luo, Tom Bergan, Madan Musuvathi, Zheng Zhang, and Lidong Zhou. Failure Recovery: When the Cure Is Worse Than the Disease. In *The 14th Workshop on Hot Topics in Operating Systems (HotOS XIV)*, 2013.

- [83] Diwaker Gupta, Kashi Venkatesh Vishwanath, and Amin Vahdat. DieCast: Testing Distributed Systems with an Accurate Scale Model. In *Proceedings of the 5th Symposium on Networked Systems Design and Implementation (NSDI)*, 2008.
- [84] Diwaker Gupta, Kenmeth Yocum, Marvin McNett, Alex C. Snoeren, Amin Vahdat, and Geoffrey M. Voelker. To Infinity and Beyond: Time-Warped Network Emulation. In *Proceedings of the 3rd Symposium on Networked Systems Design and Implementation (NSDI)*, 2006.
- [85] James Hamilton. On Designing and Deploying Internet-Scale Services. In *Proceedings of the 21st Large Installation System Administration Conference (LISA)*, 2007.
- [86] Mingzhe Hao, Gokul Soundararajan, Deepak Kenchammana-Hosekote, Andrew A. Chien, and Haryadi S. Gunawi. The Tail at Store: A Revelation from Millions of Hours of Disk and SSD Deployments. In *Proceedings of the 14th USENIX Symposium on File and Storage Technologies (FAST)*, 2016.
- [87] Chris Hawblitzel, Jon Howell, Manos Kapritsos, Jacob R. Lorch, Bryan Parno, Michael L. Roberts, Srinath Setty, and Brian Zill. IronFleet: Proving Practical Distributed Systems Correct. In *Proceedings of the 25th ACM Symposium on Operating Systems Principles (SOSP)*, 2015.
- [88] Naohiro Hayashibara, Xavier Defago, Rami Yared, and Takuya Katayama. The Phi Accrual Failure Detector. In *The 23rd Symposium on Reliable Distributed Systems (SRDS)*, 2004.
- [89] Brian Hayes. Cloud Computing. In *Communications of the ACM*, 2008.
- [90] Val Henson, Arjan van de Ven, Amit Gud, and Zach Brown. Chunkfs: Using divide-and-conquer to improve file system reliability and repair. In *IEEE 2nd Workshop on Hot Topics in System Dependability (HotDep)*, 2006.
- [91] Benjamin Hindman, Andy Konwinski, Matei Zaharia, Ali Ghodsi, Anthony D. Joseph, Randy Katz, Scott Shenker, and Ion Stoica. Mesos: A Platform for Fine-Grained Resource Sharing in the Data Center. In *Proceedings of the 8th Symposium on Networked Systems Design and Implementation (NSDI)*, 2011.
- [92] Chun-Hung Hsiao, Cristiano L. Pereira, Jie Yu, Gilles A. Pokam, Satish Narayanasamy, Peter M. Chen, Ziyun Kong, and Jason Flinn. Race Detection for Event-Driven Mobile Applications. In *PLDI*, 2014.
- [93] Nicholas Hunt, Tom Bergan, Luis Ceze, and Steven D. Gribble. DDOS: Taming Non-determinism in Distributed Systems. In *Proceedings of the 18th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2013.
- [94] Patrick Hunt, Mahadev Konar, Flavio P. Junqueira, and Benjamin Reed. ZooKeeper: Wait-free coordination for Internet-scale systems. In *Proceedings of the 2010 USENIX Annual Technical Conference (ATC)*, 2010.

- [95] Nicholas Jalbert, Cristiano Pereira, Gilles Pokam, and Koushik Sen. Radbench: A concurrency bug benchmark suite. In *The 3rd USENIX Workshop on Hot Topics in Parallelism (HOTPAR)*, 2011.
- [96] Guoliang Jin, Linhai Song, Xiaoming Shi, Joel Scherpelz, and Shan Lu. Understanding and Detecting Real-World Performance Bugs. In *Proceedings of the ACM SIGPLAN 2012 Conference on Programming Language Design and Implementation (PLDI)*, 2012.
- [97] Guoliang Jin, Aditya Thakur, Ben Liblit, and Shan Lu. Instrumentation and Sampling Strategies for Cooperative Concurrency Bug Isolation. In *OOPSLA*, 2010.
- [98] Guoliang Jin, Wei Zhang, Dongdong Deng, Ben Liblit, and Shan Lu. Automated Concurrency-Bug Fixing. In *OSDI*, 2012.
- [99] Pallavi Joshi, Malay Ganai, Gogul Balakrishnan, Aarti Gupta, and Nadia Papakonstantinou. SETSUDO : Perturbation-based Testing Framework for Scalable Distributed Systems. In *Conference on Timely Results in Operating Systems (TRIOS)*, 2013.
- [100] Pallavi Joshi, Haryadi S. Gunawi, and Koushik Sen. PREFAIL: A Programmable Tool for Multiple-Failure Injection. In *Proceedings of the 26th Annual ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*, 2011.
- [101] Xiaoen Ju, Livio Soares, Kang G. Shin, Kyung Dong Ryu, and Dilma Da Silva. On Fault Resilience of OpenStack. In *Proceedings of the 4th ACM Symposium on Cloud Computing (SoCC)*, 2013.
- [102] Horatiu Julia, Daniel Tralamazza, Cristian Zamfir, and George Candea. Deadlock Immunity: Enabling Systems To Defend Against Deadlocks. In *OSDI*, 2008.
- [103] Charles Killian, James W. Anderson, Ranjit Jhala, and Amin Vahdat. Life, Death, and the Critical Transition: Finding Liveness Bugs in Systems Code. In *Proceedings of the 4th Symposium on Networked Systems Design and Implementation (NSDI)*, 2007.
- [104] Ignacio Laguna, Dong H. Ahn, Bronis R. de Supinski, Todd Gamblin, Gregory L. Lee, Martin Schulz, Saurabh Bagchi, Milind Kulkarni, Bowen Zhou, Zhezhe Chen, and Feng Qin. Debugging High-Performance Computing Applications at Massive Scales. *Communications of the ACM (CACM)*, 58(9), September 2015.
- [105] Avinash Lakshman and Prashant Malik. Cassandra - A Decentralized Structured Storage System. In *The 3rd ACM SIGOPS International Workshop on Large Scale Distributed Systems and Middleware (LADIS)*, 2009.
- [106] Tanakorn Leesatapornwongsa and Haryadi S. Gunawi. The Case for Drill-Ready Cloud Computing. In *Proceedings of the 5th ACM Symposium on Cloud Computing (SoCC)*, 2014.
- [107] Tanakorn Leesatapornwongsa, Mingzhe Hao, Pallavi Joshi, Jeffrey F. Lukman, and Haryadi S. Gunawi. SAMC: Semantic-Aware Model Checking for Fast Discovery of Deep

- Bugs in Cloud Systems. In *Proceedings of the 11th Symposium on Operating Systems Design and Implementation (OSDI)*, 2014.
- [108] Sihan Li, Hucheng Zhou, Haoxiang Lin, Tian Xiao, Haibo Lin, Wei Lin, and Tao Xie. A Characteristic Study on Failures of Production Distributed Data-Parallel Programs. In *Proceedings of the 35th International Conference on Software Engineering (ICSE)*, 2013.
 - [109] Xin Li, Michael C. Huang, and Kai Shen. An Empirical Study of Memory Hardware Errors in A Server Farm. In *The 3rd Workshop on Hot Topics in System Dependability (HotDep)*, 2007.
 - [110] Ben Liblit, Alex Aiken, Alice X. Zheng, and Michael I. Jordan. Bug Isolation via Remote Program Sampling. In *PLDI*, 2003.
 - [111] Thomas A. Limoncelli and Doug Hughe. LISA '11 Theme – DevOps: New Challenges, Proven Values. *USENIX ;login: Magazine*, 36(4), August 2011.
 - [112] David Linthicum. Calculating the true cost of cloud outages. <http://www.infoworld.com>, 2013.
 - [113] Haopeng Liu, Guangpu Li, Jeffrey F. Lukman, Jiaxin Li, Shan Lu, Haryadi S. Gunawi, and Chen Tian. DCatch: Automatically Detecting Distributed Concurrency Bugs in Cloud Systems. In *Proceedings of the 22nd International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2017.
 - [114] Peng Liu, Omer Tripp, and Charles Zhang. Grail: Context-Aware Fixing of Concurrency Bugs. In *FSE*, 2014.
 - [115] Xuezheng Liu, Zhenyu Guo, Xi Wang, Feibo Chen, Xiaochen Lian, Jian Tang, Ming Wu, M. Frans Kaashoek, and Zheng Zhang. D3S: Debugging Deployed Distributed Systems. In *Proceedings of the 5th Symposium on Networked Systems Design and Implementation (NSDI)*, 2008.
 - [116] Xuezheng Liu, Wei Lin, Aimin Pan, and Zheng Zhang. WiDS Checker: Combating Bugs in Distributed Systems. In *Proceedings of the 4th Symposium on Networked Systems Design and Implementation (NSDI)*, 2007.
 - [117] Lanyue Lu, Andrea C. Arpaci-Dusseau, Remzi H. Arpaci-Dusseau, and Shan Lu. A Study of Linux File System Evolution. In *Proceedings of the 11th USENIX Symposium on File and Storage Technologies (FAST)*, 2013.
 - [118] Shan Lu, Soyeon Park, Eunsoo Seo, and Yuanyuan Zhou. Learning from Mistakes — A Comprehensive Study on Real World Concurrency Bug Characteristics. In *Proceedings of the 13th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2008.
 - [119] Shan Lu, Joseph Tucek, Feng Qin, and Yuanyuan Zhou. AVIO: Detecting Atomicity Violations via Access Interleaving Invariants. In *ASPLOS*, 2006.

- [120] Brandon Lucia and Luis Ceze. Cooperative Empirical Failure Avoidance for Multithreaded Programs. In *ASPLOS*, 2013.
- [121] Jake Luo, Min Wu, Deepika Gopukumar, , and Yiqing Zhao. Big Data Application in Biomedical Research and Health Care: A Literature Review. In *Biomedical Informatics Insights*, 2016.
- [122] Ao Ma, Chris Dragga, Andrea C. Arpaci-Dusseau, and Remzi H. Arpaci-Dusseau. ffsck: The Fast File System Checker. In *Proceedings of the 11th USENIX Symposium on File and Storage Technologies (FAST)*, 2013.
- [123] Paul D. Marinescu, Radu Banabic, and George Candea. An Extensible Technique for High-Precision Testing of Recovery Code. In *Proceedings of the 2010 USENIX Annual Technical Conference (ATC)*, 2010.
- [124] Sean Marston, Zhi Li, Subhajyoti Bandyopadhyay, Juheng Zhang, and Anand Ghalsasi. Cloud computing - The business perspective. In *Decision Support Systems*, 2011.
- [125] Sergey Melnik, Andrey Gubarev, Jing Jing Long, Geoffrey Romer, Shiva Shivakumar, Matt Tolton, and Theo Vassilakis. Dremel: Interactive Analysis of Web-Scale Datasets. In *Proceedings of the 36th International Conference on Very Large Data Bases (VLDB)*, 2010.
- [126] Maged Michael, Jose E. Moreira, Doron Shiloach, and Robert W. Wisniewski. Scale-up x Scale-out: A Case Study using Nutch/Lucene. In *2007 IEEE International Parallel and Distributed Processing Symposium*, 2007.
- [127] Madanlal Musuvathi and Shaz Qadeer. Iterative Context Bounding for Systematic Testing of Multithreaded Programs. In *PLDI*, 2007.
- [128] Khanh Nguyen, Lu Fang, Guoqing Xu, Brian Demsky, Shan Lu, Sanazsadat Alamian, and Onur Mutlu. Yak: A High-Performance Big-Data-Friendly Garbage Collector. In *Proceedings of the 12th Symposium on Operating Systems Design and Implementation (OSDI)*, 2016.
- [129] Diego Ongaro and John Ousterhout. In Search of an Understandable Consensus Algorithm. In *Proceedings of the 2014 USENIX Annual Technical Conference (ATC)*, 2014.
- [130] Kay Ousterhout, Ryan Rasti, Sylvia Ratnasamy, Scott Shenker, and Byung-Gon Chun. Making Sense of Performance in Data Analytics Frameworks. In *Proceedings of the 12th Symposium on Networked Systems Design and Implementation (NSDI)*, 2015.
- [131] Nicolas Palix, Gaël Thomas, Suman Saha, Christophe Calvès, Julia Lawall, and Gilles Muller. Faults in Linux: Ten Years Later. In *ASPLOS*, 2011.
- [132] Soyeon Park, Shan Lu, and Yuanyuan Zhou. CTrigger: Exposing Atomicity Violation Bugs from Their Finding Places. In *ASPLOS*, 2009.

- [133] Gilles Pokam, Klaus Danne, Cristiano Pereira, Rolf Kassa, Tim Kranich, Shiliang Hu, Justin Gottschlich, Nima Honarmand, Nathan Dautenhahn, Samuel T. King, and Josep Torrellas. QuickRec: Prototyping an Intel Architecture Extension for Record and Replay of Multi-threaded Programs. In *ISCA*, 2013.
- [134] Shanxiang Qi, Abdullah A. Muzahid, Wonsun Ahn, and Josep Torrellas. Dynamically Detecting and Tolerating IF-Condition Data Races. In *HPCA*, 2014.
- [135] Veselin Raychev, Martin T. Vechev, and Manu Sridharan. Effective race detection for event-driven programs. In *OOPSLA*, 2013.
- [136] Patrick Reynolds, Charles Killian, Janet L. Wiener, Jeffrey C. Mogul, Mehul A. Shah, and Amin Vahdat. Pip: Detecting the Unexpected in Distributed Systems. In *Proceedings of the 3rd Symposium on Networked Systems Design and Implementation (NSDI)*, 2006.
- [137] Swarup Kumar Sahoo, John Criswell, and Vikram Adve. An Empirical Study of Reported Bugs in Server Software with Implications for Automated Bug Diagnosis. In *Proceedings of the 32nd International Conference on Software Engineering (ICSE)*, 2010.
- [138] Stefan Savage, Michael Burrows, Greg Nelson, Patrick Sobalvarro, and Thomas Anderson. Eraser: A Dynamic Data Race Detector for Multithreaded Programs. *ACM TOCS*, 1997.
- [139] Koushik Sen. Race Directed Random Testing of Concurrent Programs. In *PLDI*, 2008.
- [140] Koushik Sen and Gul Agha. Automated Systematic Testing of Open Distributed Programs. In *FSE*, 2006.
- [141] Konstantin Shvachko, Hairong Kuang, Sanjay Radia, and Robert Chansler. The Hadoop Distributed File System. In *Proceedings of the 26th IEEE Symposium on Massive Storage Systems and Technologies (MSST)*, 2010.
- [142] Konstantin V. Shvachko. HDFS Scalability: The Limits to Growth. *USENIX ;login.*, 35(2), April 2010.
- [143] Jiri Simsa, Randy Bryant, and Garth Gibson. dBug: Systematic Evaluation of Distributed Systems. In *5th International Workshop on Systems Software Verification (SSV)*, 2010.
- [144] Jiri Simsa, Randy Bryant, Garth A. Gibson, and Jason Hickey. Scalable Dynamic Partial Order Reduction. In *The 3rd International Conference on Runtime Verification (RV)*, 2012.
- [145] A. Prasad Sistla, Viktor Gyuris, and E. Allen Emerson. SMC: a symmetry-based model checker for verification of safety and liveness properties. *ACM Transactions on Software Engineering and Methodology*, 2010.
- [146] Vilas Sridharan and Dean Liberty. A Study of DRAM Failures in the Field. In *Proceedings of International Conference on High Performance Computing, Networking, Storage and Analysis (SC)*, 2012.

- [147] Riza O. Suminto, Agung Laksono, Anang D. Satria, Thanh Do, and Haryadi S. Gunawi. Towards Pre-Deployment Detection of Performance Failures in Cloud Distributed Systems. In *The 7th USENIX Workshop on Hot Topics in Cloud Computing (HotCloud)*, 2015.
- [148] Vinod Kumar Vavilapalli, Arun C Murthy, Chris Douglas, Sharad Agarwal, Mahadev Konar, Robert Evans, Thomas Graves, Jason Lowe, Hitesh Shah, Siddharth Seth, Bikas Saha, Carlo Curino, Owen O'Malley, Sanjay Radia, Benjamin Reed, and Eric Baldeschwieler. Apache Hadoop YARN: Yet Another Resource Negotiator. In *Proceedings of the 4th ACM Symposium on Cloud Computing (SoCC)*, 2013.
- [149] Kaushik Veeraraghavan, Dongyoon Lee, Benjamin Wester, Jessica Ouyang, Peter M. Chen, Jason Flinn, and Satish Narayanasamy. DoublePlay: Parallelizing Sequential Logging and Replay. In *ASPLOS*, 2011.
- [150] Abhishek Verma, Luis Pedrosa, Madhukar Korupolu, David Oppenheimer, Eric Tune, and John Wilkes. Large-scale cluster management at Google with Borg. In *Proceedings of the 2015 EuroSys Conference (EuroSys)*, 2015.
- [151] Yang Wang, Manos Kapritsos, Lara Schmidt, Lorenzo Alvisi, and Mike Dahlin. Exalt: Empowering Researchers to Evaluate Large-Scale Storage Systems. In *Proceedings of the 11th Symposium on Networked Systems Design and Implementation (NSDI)*, 2014.
- [152] Yin Wang, Terence Kelly, Manjunath Kudlur, Stephane Lafortune, and Scott Mahlke. Gadara: Dynamic Deadlock Avoidance for Multithreaded Programs. In *OSDI*, 2008.
- [153] Matt Welsh, David Culler, and Eric Brewer. SEDA: An Architecture for Well-Conditioned, Scalable Internet Services. In *Proceedings of the 18th ACM Symposium on Operating Systems Principles (SOSP)*, 2001.
- [154] James R. Wilcox, Doug Woos, Pavel Panchekha, Zachary Tatlock, Xi Wang, Michael D. Ernst, and Tom Anderson. Verdi: A framework for formally verifying distributed system implementations. In *Proceedings of the 36th ACM SIGPLAN Conference on Programming Language Design and Implementation (PLDI)*, 2015.
- [155] Tian Xiao, Jiaying Zhang, Hucheng Zhou, Zhenyu Guo, Sean McDirmid, Wei Lin, Wenguang Chen, and Lidong Zhou. Nondeterminism in MapReduce Considered Harmful? An Empirical Study on Non-commutative Aggregators in MapReduce Programs. In *Proceedings of the 36th International Conference on Software Engineering (ICSE)*, 2014.
- [156] Maysam Yabandeh, Nikola Knezevic, Dejan Kostic, and Viktor Kuncak. CrystalBall: Predicting and Preventing Inconsistencies in Deployed Distributed Systems. In *Proceedings of the 6th Symposium on Networked Systems Design and Implementation (NSDI)*, 2009.
- [157] Junfeng Yang, Tisheng Chen, Ming Wu, Zhilei Xu, Xuezheng Liu, Haoxiang Lin, Mao Yang, Fan Long, Lintao Zhang, and Lidong Zhou. MODIST: Transparent Model Checking of Unmodified Distributed Systems. In *Proceedings of the 6th Symposium on Networked Systems Design and Implementation (NSDI)*, 2009.

- [158] Zuoning Yin, Xiao Ma, Jing Zheng, Yuanyuan Zhou, Lakshmi N. Bairavasundaram, and Shankar Pasupathy. An Empirical Study on Configuration Errors in Commercial and Open Source Systems. In *Proceedings of the 23rd ACM Symposium on Operating Systems Principles (SOSP)*, 2011.
- [159] Zuoning Yin, Ding Yuan, Yuanyuan Zhou, Shankar Pasupathy, and Lakshmi Bairavasundaram. How do fixes become bugs? In *Proceedings of the 19th ACM SIGSOFT international symposium on Foundations of software engineering (FSE)*, 2011.
- [160] Jie Yu. A collection of concurrency bugs. <https://github.com/jieyu/concurrency-bugs>.
- [161] Ding Yuan, Yu Luo, Xin Zhuang, Guilherme Renna Rodrigues, Xu Zhao, Yongle Zhang, Pranay U. Jain, and Michael Stumm. Simple Testing Can Prevent Most Critical Failures: An Analysis of Production Failures in Distributed Data-Intensive Systems. In *Proceedings of the 11th Symposium on Operating Systems Design and Implementation (OSDI)*, 2014.
- [162] Matei Zaharia, Mosharaf Chowdhury, Tathagata Das, Ankur Dave, Justin Ma, Murphy McCauley, Michael J. Franklin, Scott Shenker, and Ion Stoica. Resilient Distributed Datasets: A Fault-Tolerant Abstraction for In-Memory Cluster Computing. In *Proceedings of the 9th Symposium on Networked Systems Design and Implementation (NSDI)*, 2012.
- [163] Cristian Zamfir and George Candea. Execution Synthesis: A Technique for Automated Software Debugging. In *Proceedings of the 2010 EuroSys Conference (EuroSys)*, 2010.
- [164] Wei Zhang, Marc de Kruijf, Ang Li, Shan Lu, and Karthikeyan Sankaralingam. ConAir: Featherweight concurrency bug recovery via single-threaded idempotent execution. In *ASPLOS*, 2013.