

Leveraging Software-Defined Networking to Mask Partial Network Partitions

by

Basil Alkhatib

A thesis
presented to the University of Waterloo
in fulfillment of the
thesis requirement for the degree of
Master of Mathematics
in
Computer Science

Waterloo, Ontario, Canada, 2021

© Basil Alkhatib 2021

Author's Declaration

This thesis consists of material all of which I authored or co-authored: see Statement of Contributions included in the thesis. This is a true copy of the thesis, including any required final revisions, as accepted by my examiners

I understand that my thesis may be made electronically available to the public.

Statement of Contributions

The results of this work are included in two papers. The first is a paper at the USENIX Symposium on Operating Systems Design and Implementation (OSDI) 2020, and the second is a main paper under review at ACM Transactions on Computer Systems (TOCS). I am the first author of the main paper and it has been done in collaboration with Ahmad Alquraan, Mohammed Alfatafta, Wael Al-Manasrah, Sreeharsha Udayashankar, and Sara Qunaibi. Ahmad and Mohammed worked with me on the tickets and the systems in Chapter 3. Mohammad also helped with the study of systems design in Chapter 4, Wael has carried out the Redis PubSub experiment in Section 6.1. Sreeharsha and Sara provided technical support for the HDFS experiments in Section 7.1.

Abstract

We present an extensive study focused on partial network partitioning. Partial network partitions disrupt the communication between some but not all nodes in a cluster. First, we conduct a comprehensive study of system failures caused by this fault in 13 popular systems. Our study reveals that the studied failures are catastrophic (e.g., lead to data loss), easily manifest, and are mainly due to design flaws. Our analysis identifies vulnerabilities in core systems mechanisms including scheduling, membership management, and ZooKeeper-based configuration management.

Second, we dissect the design of nine popular systems and identify four principled approaches for tolerating partial partitions. Unfortunately, our analysis shows that implemented fault tolerance techniques are inadequate for modern systems; they either patch a particular mechanism or lead to a complete cluster shutdown, even when alternative network paths exist.

Finally, our findings motivate us to build Nifty, a transparent communication layer that masks partial network partitions. Nifty builds an overlay between nodes to detour packets around partial partitions. Nifty provides an approach for applications to optimize their operation during a partial partition. We demonstrate the benefit of this approach through integrating Nifty with VoltDB and HDFS.

Acknowledgements

I would like to thank my supervisor, Samer Al-Kiswany for all the valuable advice and guidance he provided throughout my Master's journey. To my friends and colleagues at WASL lab who I have had a lot of good times with even though we could not work together in the university for most of my study. I'm thankful for having a supportive family who has always been there whenever I needed them. I would also like to thank my thesis readers, Professor Bernard Wong and Professor Trevor Brown for their valuable feedback and comments.

Table of Contents

List of Figures	viii
List of Tables	x
1 Introduction	1
2 Causes of Partial Network Partitions	5
2.1 Definitions	5
2.2 Causes of Partial Network Partitioning Fault	5
3 Analysis of Partial Network-Partitioning Failures	7
3.1 Methodology	7
3.2 Limitations	8
3.3 Findings	9
3.4 Design Pitfalls	11
3.5 Insights	14
4 Dissecting Modern Fault Tolerance Techniques	16
4.1 Identifying the Surviving Clique	17
4.2 Checking Neighbors' Views	19
4.2.1 RabbitMQ	19
4.2.2 Elasticsearch	21

4.3	Failure Verification	25
4.4	Neutralizing Partitioned Nodes	25
4.5	Summary	26
5	Nifty Design	28
5.1	Implementation	30
6	Evaluation	32
6.1	Overhead Evaluation	32
6.2	Handling Partial Partitions	34
7	Classification API Utility	39
7.1	HDFS	39
7.2	VoltDB	42
8	Related Work	44
9	Concluding Remarks	47
	References	48
	APPENDICES	62

List of Figures

1.1	Partial partition. Groups 1 and 2 are disconnected, while Group 3 can reach both sides of the partition.	2
4.1	VoltDB’s surviving clique. Gray nodes shut down as they are not in the clique.	18
4.2	The probability of a VoltDB cluster shutdown. Different lines represent different cluster sizes. The x-axis shows the number of nodes that are not in the clique.	19
4.3	A scenario for RabbitMQ’s <i>pause</i> policy. Every non-bridge node pauses (gray nodes) as it detects that it cannot reach one node on the other side.	20
4.4	The median number of paused nodes in a cluster of 15 nodes. In all runs, one node is unaffected by the partition. The notation (i, j) shows the number of nodes on each side of the partition.	21
4.5	Elasticsearch unavailability scenario. The master pauses because it cannot reach majority of nodes, and all nodes pause because they cannot reach the master.	22
4.6	A Mesos cluster becomes unavailable when a partial partition isolates the master node and its backups.	23
5.1	A Nifty routing example. A partial network partition isolates node 1 from nodes 3 and 4, and another partial partition isolates node 4 from nodes 1 and 2. Communication between 1 and 4 is routed through nodes 2 and 3.	29
6.1	Nifty’s overhead. The average throughput for HDFS (a) and the average throughput vs. average latency Kafka, ActiveMQ, and MongoDB. (-P) denotes the results with a partial partition.	36

6.2	Nifty's overhead. the average throughput vs. average latency for VoltDB, RabbitMQ, and the average throughput for Redis PubSub. (-P) denotes the results with a partial partition.	37
6.3	Scalability evaluation. Average throughput while increasing the number of nodes.	38
6.4	Tail latency evaluation. Average throughput vs. 99th percentile of latency.	38
7.1	HDFS worst case rerouting. NameNode choosing the replicas to be on 1,2, and 3 will case the data to move across the partition twice.	40
7.2	HDFS write throughput with different optimizations	42
7.3	The impact of MPI placement on VoltDB's performance. Figure shows the average latency (a) and average throughput (b). Standard deviation was less than 2%.	43
1	Probability of a VoltDB system shut down.	64
2	The median number of paused nodes in a cluster of 15 nodes. In all runs, 3 node are unaffected by the partition. The notation (i, j) shows the number of nodes on each side of the partition.	65
3	The median number of paused nodes in a cluster of 15 nodes. In all runs, 5 node are unaffected by the partition. The notation (i, j) shows the number of nodes on each side of the partition.	65

List of Tables

3.1	List of studied systems and the number of studied failures. The shaded rows are systems that implemented a fault tolerance technique specifically for partial network partitioning.	8
4.1	Summary of shortcomings. (D) indicates that the nodes shut down. (P) indicates that the nodes pause until the partition heals. In the worst case, RabbitMQ pauses all nodes except one. We consider this a complete cluster loss (1). Under different RabbitMQ policies, (2) and (3) can occur. (S) indicates a system-wide technique, whereas (M) is a mechanism-specific technique.	24

Chapter 1

Introduction

Modern networks are complex. They use heterogeneous hardware and software [143], deploy diverse middleboxes (e.g., NAT, load balancers, and firewalls) [18, 102, 109], and span multiple data centers [18, 109]. Despite the high redundancy built into modern networks, catastrophic failures are common [143, 102, 101, 144]. Nevertheless, modern cloud systems are expected to be highly available [81, 132] and to preserve stored data despite failures of nodes, networks, or even entire data centers [83, 149, 89].

We focus our investigation on a peculiar type of network fault: *partial network partitions*, which disrupts the communication between some, but not all, nodes in a cluster. Figure 1.1 illustrates how a partial network partition divides a cluster into three groups of nodes, such that two groups (Group 1 and Group 2) are disconnected, but Group 3 can communicate with Groups 1 and 2.

In our previous work [74] we identified this fault and presented examples of how it leads to system failures. Other than our previous preliminary effort, we did not find any in-depth analysis of partial network partition failures and of their fault tolerance techniques. Nevertheless, we found 54 reports of failures caused by partial network partitioning faults¹ in the publicly accessible issue tracking systems of 13 production-quality systems (Chapter 3), numerous blog posts and discussions of this fault on developers' forums (Chapter 2.2), and eight popular systems with fault tolerance techniques specifically designed to tolerate this type of fault (Chapter 4).

Our goal in this work is threefold. First, we aim to study failures caused by partial network partitioning to understand their impact and failure characteristics and, foremost, to identify opportunities to improve systems' resiliency to this type of fault. Second, we

¹A *fault* is the initial root cause. If not properly handled, it may lead to a user-visible system *failure*.

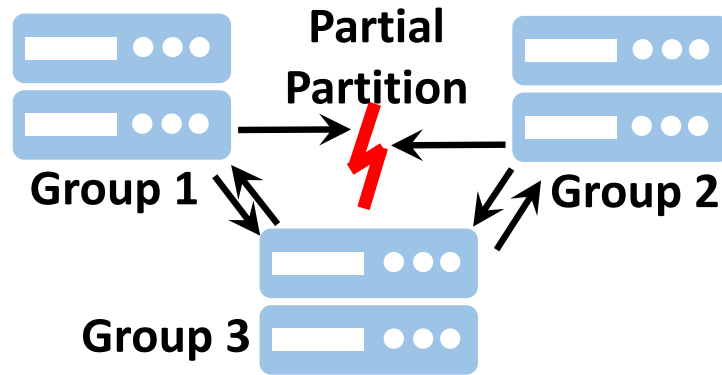


Figure 1.1: Partial partition. Groups 1 and 2 are disconnected, while Group 3 can reach both sides of the partition.

aim to dissect the fault tolerance techniques implemented in popular production systems and identify their shortcomings. Third, we aim to design a generic fault tolerance technique for partial network partitioning. This is the first work to characterize these failures and explore fault tolerance techniques for partial partitioning faults.

It is important to understand that *partial* partitions are fundamentally different from *complete* partitions [74]. Complete partitions split a cluster into two completely disconnected sides and are well studied with known theoretical bounds (CAP theorem [100]) and numerous practical solutions [131, 114, 141, 121]. On the contrary, a cluster experiencing a partial partition is still connected but not all-to-all connected. Consequently, the theoretical bounds of complete partitions do not apply to partial partitions, and fault tolerance techniques for complete partitions are not effective in handling partial partitions (Chapter 8).

An analysis of partial network partitioning failures. We conduct an in-depth study of 54 partial network partitioning failures from 13 cloud systems (Chapter 3). We select a diverse set of systems, including database systems (MongoDB and HBase), file systems (HDFS and MooseFS), an object storage system (Ceph), messaging systems (RabbitMQ, Kafka, and ActiveMQ), a data-processing system (MapReduce), a search engine (Elasticsearch), an in-memory data grid (Hazelcast), and resource managers (Mesos and DKron). For each considered failure, we carefully study the failure report, logs, discussions between users and developers, source code, and code patches.

Failure Impact. Overall, we find that partial network partitioning faults cause silent failures with catastrophic effects (e.g., data loss and corruption) that affect core system mechanisms (e.g., leader election and replication).

Ease of manifestation. Unfortunately, these failures can easily occur. The majority of the failures are deterministic and require less than four events (e.g., read or write request) for the failure to occur. Even worse, all the studied failures can be triggered by partially partitioning a single node. The majority of failures do not require client access or can be triggered by clients only accessing one side of the partition.

Insights. We identify three approaches to improve system resilience: better testing, focused design reviews, and building a generic fault tolerance communication layer. Our analysis of each failure’s manifestation sequence, access patterns, and timing constraints shows that almost all the failures could have been revealed through simple tests and by only using five nodes. Second, the majority of failures are due to design flaws. We posit that design reviews focused on network partitioning could identify these vulnerabilities. Third, building a generic communication layer to mask partial partitions is feasible, simplifies system design, and improves system resiliency.

Finally, we identify that a common deployment approach of Zookeeper introduces a failure vulnerability (Chapter 4). Our analysis shows that system designers need to design additional mechanisms to handle partial partitions when using Zookeeper or other external coordination services.

Dissecting modern fault tolerance techniques. We dissect the implementation of nine popular systems (VoltDB, MapReduce, HBase, MongoDB, Elasticsearch, Mesos, LogCabin, RabbitMQ, and HazelCast) and study the fault tolerance techniques they employ specifically to tolerate partial partitions (Chapter 4). For each system, we study the source code, analyze the fault tolerance technique’s design, extract the design principles, and identify the technique’s shortcomings. We identify four principled approaches for tolerating partial partitions: identifying the surviving clique, checking neighbors’ views, verifying failures announced by other nodes, and neutralizing partially partitioned nodes.

Our analysis reveals that the studied fault tolerance techniques are inadequate. They either patch a specific system mechanism, which leaves the rest of the system vulnerable to failures, or unnecessarily shut down the entire cluster or pause up to half of the cluster nodes (Chapter 4).

Designing a generic fault tolerance technique. Our findings motivate us to build network partitioning fault-tolerance layer (Nifty), a simple, generic, and transparent communication layer that can mask partial network partitions (Chapter 5). Nifty’s approach is simple; it monitors the connectivity in a cluster through all-to-all heart beating, and when it detects a partial partition, it detours the traffic around the partition through intermediate nodes. Nifty overcomes all the shortcomings present in the studied fault tolerance techniques.

The main insight of Nifty is that tolerating partial partitioning does not require elaborate techniques such as the ones adopted by current systems (Chapter 4). Many modern systems already incorporate membership and connectivity monitoring mechanisms based on all-to-all heart beating [58, 66, 13]. Nifty shows that extending these mechanisms with a simple rerouting capability can effectively mask partial partitions.

Nifty reroutes packet between end hosts to mask partial partitions. This approach increases the load on the intermediate nodes and can create a performance bottleneck. To reduce the load on intermediate nodes, system designers may optimize the data or process placement or employ a flow-control mechanism. Nifty provides an API that exposes the network state to the system running atop of it and facilitates building system-specific optimizations.

To demonstrate Nifty’s effectiveness, we deploy it with seven systems: HDFS, Kafka, RabbitMQ, ActiveMQ, MongoDB, VoltDB, and Redis PubSub. We choose these systems because they are data intensive and popular systems. Furthermore, RabbitMQ and VoltDB implement generic techniques to tolerate partial partitions. Our prototype evaluation with synthetic and real-world benchmarks shows that Nifty effectively masks partial partitions while adding negligible overhead.

To demonstrate the utility of the Nifty API, we integrate Nifty with VoltDB and HDFS and explore a number of optimizations. Our evaluation shows that system-specific optimizations can significantly reduce the traffic rerouting overhead during partial partitions.

Chapter 2

Causes of Partial Network Partitions

2.1 Definitions

A *partial network partition* is a network fault that prevents at least one node (e.g., a node in Group 1 in Figure 1.1) from communicating with another node (Group 2) in the system, while a third node (Group 3) can communicate with both affected nodes. Nodes in a partially partitioned cluster are still connected but are not all-to-all connected (i.e., they do not form a complete graph [148]). A partial partition divides a cluster into three groups: two sides and one bridge group. We identify a node as a *bridge* node if it can reach at least one node on each side of a partition. A partial partition has two *sides*, all the nodes on one side of the partition cannot reach all the nodes on the other side of the partition. We note that a cluster may suffer from multiple concurrent partial partitions.

We define a *single-node partial partition* as a partial partition that has a single node on one side of the partition, while the rest of the cluster nodes are bridge nodes or are on the other side of the partition. For instance, a single-node partial partition can be caused by a firewall misconfiguration that prevents a node from communicating with some other nodes.

2.2 Causes of Partial Network Partitioning Fault

Recent reports indicate that network partitioning faults are common and happen at various scales. Connectivity loss between data centers [143] leads to network partitions in geo-replicated systems. Wide area network partitions happen as frequently as once every four

days [144]. Switch failures can cause a network partition in a data center [101]. Switch failures caused 40 network partitions in two years at Google [102] and 70% of the downtime at Microsoft [101]. On a single node, NIC [12] or software failures can partition a node that may host multiple VMs. Finally, network partitions caused by correlated failures are common [109, 101, 144] and are often caused by system-wide maintenance tasks [102, 101].

While we did not find failure reports that detail partial partitioning faults, we found numerous discussions of their impact on production systems. Partial partitions were the cause of service outages at Cloudflare [15], Google [26], Lyft [40], and Amazon AWS [19]. A misbehaving switch caused the failure at Cloudflare. The switch data plane did not process all packets, while the control plane protocols remained operational. This disrupted the communication between some nodes in the cluster and eventually caused a 6-hour outage of Cloudflare. AWS [19] also blame a misbehaving switch for a partial partitioning failure that affected applications that span multiple availability zones. A partial partition also affected Google Compute Engine (GCE) services. When a new VM is added, GCE uses two mechanisms to update the other VMs: one to update VMs in the same zone as the new VM, and another to update VMs in other zones. When the processes responsible for updating the VMs in other zones failed [26], the newly added VM became unreachable from VMs from outside its zone. This created a partial partition since old VMs in the same zone as the new VM could reach all VMs since they are updated through a separate mechanism. Lyft reported cases of partial network partitions while running Kafka at scale at AWS [19]. Finally, an early version of Google’s B4 control plane use a primary master with a standby backup. A partial partition disconnected the primary master from the standby master while both can reach the switches and the system gateway. This led to having two active masters in the infrastructure [109].

Furthermore, we found 54 failure reports detailing system failures due to partial network partitions, and numerous articles and online discussions discussing the fault [124, 122, 127, 49]. Some of these reports and discussions mention the root cause of the partial partition. Partial partitions are caused by a connectivity loss between two data centers [143] while both are reachable by a third center, the failure of additional links between racks [56, 52], network misconfiguration [31], firewall misconfiguration [31], network upgrades [16], and flaky links between switches [65].

Chapter 3

Analysis of Partial Network-Partitioning Failures

We conduct an in-depth study of partial network partitioning failures reported in 13 popular systems (Table 3.1). We aim to understand the impact and characteristics of these failures and to identify opportunities for improving system resilience.

3.1 Methodology

We choose 13 diverse and widely used systems (Table 3.1), including two databases, a data analysis framework, two file systems, three messaging systems, a storage system, a search engine, an in-memory data grid, and two resource managers.

We selected the 54 failures in our study from publicly accessible issue-tracking systems. First, we used the search tools in the issue-tracking systems to find tickets related to partial network partitioning. Users did not classify network partitioning failures based on the partition type, so we had to search for all network partitioning failures and manually identified partial partitioning failures. We used the following keywords: “network partition,” “partial network partition,” “partial partition,” “network failure,” “switch failure,” “isolation,” “split-brain,” and “asymmetric partition.” Second, we considered tickets that were dated 2011 or later. Third, we excluded tickets marked as “Minor.” For each ticket, we studied the failure description, system logs, developers’ and users’ comments, and code patches. For tickets that lacked enough details (e.g., missing output logs or did not have details about the affected mechanism), we manually reproduced them using NEAT [74].

Table 3.1: List of studied systems and the number of studied failures. The shaded rows are systems that implemented a fault tolerance technique specifically for partial network partitioning.

System	Category	Failures	
		Total	Catastrophic
Elasticsearch [23]	Search engine	17	17
MongoDB [45]	Database	9	5
RabbitMQ [58]	Messaging	5	3
MapReduce [4]	Data processing	4	2
HBase [5]	Database	3	2
Mesos [6]	Resource manager	2	1
Hazelcast [27]	In-memory data structures	2	2
Kafka [37]	Messaging	3	3
HDFS [4]	File system	3	1
Ceph [13]	Storage system	2	2
MooseFS [46]	File system	2	2
ActiveMQ [2]	Messaging	1	1
DKron [22]	Resource manager	1	1
Total	-	54	42

Finally, during our evaluation, we found and reported bugs in Kafka and Elasticsearch. We included these failures in our study.

We differentiate failures by their manifestation sequences. In a few cases, the same faulty mechanism leads to two different failure paths. We count these as separate failures, even if they are reported in a single ticket. Similarly, although the exact failure is sometimes reported in multiple tickets, we count it once in our study.

3.2 Limitations

As with any characterization study, our findings may not be generalizable. Here, we list four potential sources of bias and describe our best efforts to address them.

1. *Representativeness of the studied systems.* Although we study 13 diverse systems (Table 3.1), our results may not be generalizable to systems we did not study. The selected systems follow diverse designs from strongly consistent (MongoDB, HBase, and Ceph) to eventually consistent (Elasticsearch) designs and from systems persisting data on disks and replicating data in-memory across nodes to caching systems. They follow a primary-backup or peer-to-peer architecture and use synchronous or

asynchronous replication. The selected systems are widely used: Kafka, ActiveMQ, and RabbitMQ are the most popular open-source messaging systems; MapReduce, HDFS, and HBase are the core of the Hadoop platform; Elasticsearch is a popular search system; and MongoDB is a popular database.

2. *Limited number of tickets.* We study all 54 tickets that we found following our methodology. Statistical inference indicates that 30 samples can sufficiently represent the entire population [107]. More rigorously, if we assume the tickets we found represent a random sample of partial network partition failures in the wild, the central limit theorem predicts that our analysis of 54 tickets has a 13% margin of error at a 95% confidence level. To increase confidence in our findings, we only report findings that apply to at least two-thirds of the studied failures. A third of our findings apply to all failures.
3. *Priority bias.* We include only high-impact tickets and avoid tickets marked by the developers as low-priority. This sampling methodology may bias the results.
4. *Observer error.* To reduce the chance for observer errors, two team members study every failure report using the same classification methodology. Then, we discuss the failure in a group meeting before reaching a verdict.

3.3 Findings

This section presents a summary of our findings. Our study indicates that partial network partitioning leads to catastrophic failures that are easy to manifest. Luckily, our study identified that code reviews and targeted testing can improve systems fault tolerance. We refer the reader to our previous paper for a detailed discussion of our findings [71].

Failure Impact. Overall, *we find that 76.4% of the studied failures lead to catastrophic effects.* A failure is said to be catastrophic if it leads to a system crash or violates the system’s guarantees such as data loss or corruption, system unavailability, and stale or dirty reads. The majority of non-catastrophic failures lead to reducing a system availability such as intermittent disruption of system operation [8].

Data loss and system unavailability are the two most common effects of partial partitions and are the result of 42.5% of failures. For instance, in HBase, region servers store their logs on HDFS. When a log reaches a certain size, the region server creates a new log and informs the master of the new log location. If a partial partition isolates a region server from the master while both can reach HDFS, the master assumes that the region

server has failed and assigns its logs to a new region server. If at this time the old region server creates a new log, the master will not know about it, and the entries in the new log will be lost [57].

The majority of failures (81.5%) are silent, meaning the user is not informed about their occurrence. Some systems return a warning to the user when an operation fails due to partial network partitioning, but these warnings are ambiguous with no clear mechanisms for resolution. For example, in Elasticsearch, if a client sends a request to a replica that is partially isolated from the other replicas, the replica will return “a rejected execution” exception [53]. This confusing warning does not inform the user of the problem’s actual cause nor the steps needed to resolve it. This is unsettling because a lack of error or warning notification delays failure detection.

Ease of manifestation. Unfortunately, the studied failures can easily occur:

- *All the studied failures, except one, are deterministic or have known time constraints*, such as the period before considering a node to have failed.
- *The majority of failures (66.6%) require three or fewer events (other than the partial partition) to manifest.* An event is a user request, a hardware or software fault, or a start of a background operation (e.g., leader election and data rebalancing). This is alarming because in real deployments, many users interact with the system, increasing the probability of failure.
- *Most failures (59.3%) do not require client access or require only that clients access one side of the partition.* To reduce the network partition’s impact, some systems limit client access to one side of the partition [28, 60, 106]. This finding shows that this fault tolerance technique is not sufficient.
- *All the studied failures can be triggered by a single-node partial partition.* Arguably, single-node partial partitions (Section 2.1) are more likely than partitioning more than one node. These partitions could happen due to a single ToR switch malfunction or by misconfiguring a single firewall.

We further study which nodes need to be isolated for a failure to manifest. Of the failures, 33.3% manifest by partitioning any node in the system—regardless of its role. Among the failures that require partitioning a specific node, partitioning the leader replica is most common (44.4%). In real deployments, partitioning a leader is likely because almost every node in the cluster is a leader for some shard.

Failure Characteristics Our study revealed two surprising characteristics of these failures. First, *the majority of the fixed bugs (59.3%) are due design flaws.* We consider a

code patch to be fixing a design flaw if it significantly changes the implemented protocol or logic, such as changing the mechanism to select a master in Elasticsearch. Second, *partial partition fault affects a wide range of system mechanisms* including leader election, configuration change, replication protocol, request routing, scheduling, and data migration. Leader election, configuration change, and replication protocols are the most affected mechanism (affected by 72.6% of failures).

Finally, These failures can be easily reproduced with small clusters of five or fewer nodes, and 75.9% require only three nodes. Furthermore, all the failures except one can be reproduced using a fault-injection framework that can inject partial partitioning faults such as NEAT [74].

3.4 Design Pitfalls

Our study revealed that the majority of failures are due to design flaws. For each design failure we study the code patches and the system design to understand the design flaw. We identified flaws in the following five common designs of core distributed system techniques. Revisiting the design of these techniques to tolerate partial network partitioning is a high impact research frontier that requires further investigation.

Leader election is the most vulnerable mechanism to partial partitions. The following are the most frequent flaws we found in the studied tickets.

- *Two leaders.* Partial network partitioning fault leads to having two active leaders in MongoDB [64] and RabbitMQ [63]. Having more than one leader results in data loss, dirty read, and stale read. This failure typically manifests when two nodes on different sides of the partition start the leader election process. If the bridge node votes for the two candidates, each candidate will get enough votes to become a leader. A common solution to avoid this double voting problem is to divide time into terms or epochs and each node has a single vote in a term [131, 130].
- *No leader in the system.* Some leader election policies may leave a cluster without a leader under partial network partitions. For instance, in an earlier version of Elasticsearch, a live node with the smallest id is the cluster leader. If a node can not reach the leader, it will ask the node with the second smallest id to become a leader. The node with the second smallest id will refuse to become a leader if it can reach the current leader. If a partial partition puts the current leader on one side of the partition and the node with the second smallest id is a bridge node, no node will be elected as a leader and the cluster pauses until the partition heals [55].

- *Leader election thrashing.* Partial partition faults may lead to continuous leader election thrashing if the two sides of the partition keep launching the leader election process. For instance, leader election in MongoDB is based on a majority vote, with an arbiter node included to break ties. Consider a shard that has two replicas (A and B), with A being the leader. If a partial partition disrupts the communication between A and B while both can reach an arbiter, B will detect that A is unreachable and calls for a leader election. Because there is only one candidate in the system, the arbiter votes for it, and B becomes the leader. The arbiter will inform A of the new leader, and A steps down. A will detect that the leader (B) is unreachable, call for a leader election, become a leader, and then B steps down. This leader-election thrashing continues until the network partition heals [8]. The system is unavailable during leader election, so this failure significantly reduces system availability. CloudFlare reported a service outage due to a similar flaw in the leader election mechanism in etcd [15].

Leader election using a coordination service. A common approach for electing a leader in modern systems (e.g., Mesos, Kafka, ActiveMQ, HBase, and Neo4j) is to use a coordination service such ZooKeeper [7] to monitor the nodes and choose a new leader when a leader fails. As this is a common usages of ZooKeeper, the ZooKeeper user guide has a "recipe" [68] for how to use ZooKeeper for leader election that is broadly followed. Unfortunately, this recipe is vulnerable to partial network partitions.

To elect a leader using ZooKeeper, each node creates a "sequence ephemeral" file in a specific shared directory at ZooKeeper. The file has a unique sequence number that is generated by ZooKeeper. The node with the smallest sequence number is the leader. If ZooKeeper misses heartbeats from a node, it deletes all the ephemeral files that are created by the unreachable node, and notifies the other nodes in the clusters. Consequently, if a leader fails, ZooKeeper deletes its sequence ephemeral file and informs the other nodes of this change. Then each node in the cluster will check the shared directory to see if its file has the smallest sequence number. The node with the smallest sequence number becomes the new leader.

If a partial partition isolates the leader from the rest of the cluster while all nodes are reachable from ZooKeeper, the nodes will typically pause their operations because they cannot reach the leader. Because ZooKeeper can reach the current leader, it will not delete its ephemeral file and no new leader will be elected. The cluster remains unavailable until the partial partition heals. This failure manifested in ActiveMQ [1] and Kafka [38].

Scheduling. Resource management and scheduling systems use heartbeating to monitor a cluster's health. If a scheduler missed heartbeats from a worker node, it will suspect that

the node has failed and will reschedule all the tasks that were running on the failed node on other nodes in the cluster.

This fault tolerance technique is vulnerable to partial partitions. If a partial partition isolates the scheduler from one of the nodes, while the affected node can reach the rest of the cluster, the scheduler will reschedule the tasks running on the affected node on other cluster nodes. This leads to double, potentially concurrent, execution. Double execution can corrupt shared state (e.g., data on HDFS) or confuse clients. For instance, in MapReduce, a partial partition leads to a double execution and data corruption of shared data [41]. Mesos [42], and Elasticsearch [21] suffered from a similar failure.

Membership management. Modern systems use membership lists to keep track of nodes in the cluster. Other systems allow for a block/allow list to report slow or unresponsive nodes to avoid performance straggling. If a node detects that another node has failed or might be slow, it will notify the metadata service. The metadata service will update the membership list or the blocklist to avoid using that node in future operations. Our study shows that under partial partitions, these techniques could lead to big availability and performance issues.

MapReduce uses blocklisting to identify slow or unresponsive nodes. If a reducer cannot reach a mapper node, it will report it to the master node. The master will not assign new tasks to the node running that mapper. If a partial partition isolates a reducer from many mappers while all nodes are still reachable by the master, the affected reducer will report and unnecessarily block list many nodes, which leads to a significant drop in cluster performance [65].

RabbitMQ supports message replication for higher availability. RabbitMQ maintains a membership log that lists the current nodes in the cluster. If nodes have conflicting views on which nodes are part of the cluster, the RabbitMQ cluster crashes. For instance, in a cluster with three nodes (A, B, and C), when a partial partition disconnects B and C, B assumes that C crashed and removes it from the membership log, and C assumes that B crashed and removes it from the membership log. This inconsistency in the cluster membership leads to a complete cluster crash [44].

Discovery service. Modern systems often use a metadata or discovery service to direct clients to a node hosting a queue in a messaging system, or to a leader replica in a storage system. If a partial partition isolates a client from some nodes in the cluster while the discovery service can reach all nodes and clients, a discovery service may point the client to a node that the client can not reach due to a partial partition. This problem often leads to system unavailability to some clients. For instance, in Kafka, a client asks the bootstrap service for a list of cluster nodes. If the client cannot reach a topic leader, while

the bootstrap service indicates that the leader is alive, messages to that leader will be lost [36]. Elasticsearch had a similar failure [54].

In HDFS, consider a case when a partial network partition separates a client from, say, rack 0, while the NameNode can reach that rack. If the NameNode allocates replicas for a new data chunk on rack 0, then a client write operation will fail, and the client will ask for a different DataNode to place its replica. The NameNode, following its rack-aware data placement, will likely suggest another node from the same rack. The process repeats five times before the client gives up [30].

3.5 Insights

Surprisingly, partial network partitioning faults trigger silent failures that have catastrophic effects in production-quality systems. It is unsettling to realise how easy it is for these failures to manifest once a partial partitioning fault happens. Isolating a single node, with three or less events, with client access to one side of the partition, deterministically causes over two thirds of the failures.

Fortunately, we identify three approaches for improving system resilience to partial partitions. First, because these faults are deterministic and can be reproduced on a five-node cluster, improved testing can reveal the majority of the studied failures. Our analysis finds timing, client access, and partition characteristics that significantly reduce the number of sufficient test cases. Second, the fact that the majority of failures are due to design flaws, indicates that system designers overlook partial network partitioning failures in the design phase. We posit that design reviews focused on network partitioning could identify these vulnerabilities. Since a large number of failures are triggered without client access, our analysis highlights that system designers should consider the impacts of partial partitioning faults on all operations, including background operations.

Third, partial network partitions have two characteristics that imply that a generic fault tolerance technique is possible. These faults can be detected by exchanging information between the nodes, and by definition, there are alternative paths in the network to reconnect the system. We leverage these two characteristics in building Nifty (Chapter 5).

Finally, we point out design flaws in core system mechanisms including leader election, scheduling, discovery service, and membership management (Section 3.4). Most of the studied failures are caused by the underlying assumption that, if a node can reach a service, all nodes can reach that service, and if a node cannot reach a service then the service is down. Our analysis shows the danger of such assumptions; this leads to a confusing state,

wherein some of the system's parts start executing a fault tolerance mechanism, while others presume the whole system is healthy and carry on normal operations. The mix of these two operation modes is poorly understood and tested.

Chapter 4

Dissecting Modern Fault Tolerance Techniques

We studied the code patches related to the tickets included in our study. Seven of the systems in Table 3.1 (MongoDB, Elasticsearch, RabbitMQ, HBase, MapReduce, Hazelcast, and Mesos) changed the system design to incorporate a fault tolerance technique specific to partial network partitioning faults. The rest of the systems either patched the code with an implementation-specific workaround or did not fix the reported bugs yet.

Furthermore, we found that two additional systems, VoltDB [66, 139] and LogCabin [39] (the original implementation of the Raft [131] consensus protocol), implement fault tolerance techniques for partial partitions. For these systems, we did not find failure reports related to partial partitioning faults in their issue tracking systems, but VoltDB announced that their recent versions tolerates partial partitions [32]. We experimented with LogCabin to understand the impact partial partitions have on strongly consistent systems and found that LogCabin incorporates a technique to tolerate partial partitions. We included VoltDB and LogCabin in our study.

For each of the nine systems, we study the source code, and extract and analyze the design principles of their fault tolerance technique. We identify four approaches for tolerating partial partitions: detecting a surviving clique of nodes, checking neighbors' views, verifying failure reports received from other nodes, and neutralizing one side of the partial partition. Unfortunately, these techniques have severe shortcomings that may lead to a complete system shutdown or to the unavailability of a major part of the system. In this chapter, we detail these techniques and discuss their shortcomings.

4.1 Identifying the Surviving Clique

Main idea. Upon a partial network partition, the system identifies the maximum clique of nodes [90], which is the largest subset of nodes that are completely connected. All nodes that are not part of the maximum clique are shut down. VoltDB and Hazelcast follow this approach.

VoltDB Implementation. VoltDB [66, 139] is a popular ACID, sharded, and replicated relational database. VoltDB follows a peer-to-peer approach to implement this technique. Every node in the system periodically sends a heartbeat to all nodes. If a node loses connectivity to any node, it suspects that a partial network partition occurred and starts the recovery procedure. The recovery procedure has two phases. In the first phase, the node that detects the failure broadcasts a list of nodes it can reach. When a node in the cluster receives this message, it broadcasts its list of reachable nodes to all nodes in the cluster. In phase two, every node independently combines the information from the other nodes into a graph representing the cluster connectivity. Each node analyzes this graph to detect the maximum completely connected clique of nodes. Every node that finds that it is not part of this “surviving” clique shuts itself down. Figure 4.1 shows an example in which a partial partition disrupts the communication between nodes 2, 3, and 4 on one side and nodes 5 and 6 on another. Nodes 5 and 6 are not part of the clique and will shut down.

After identifying the surviving clique, the system verifies that it did not lose any data by verifying that the surviving clique has at least one replica of every data shard. If the clique is missing one shard, such as when all the replicas of a shard are shut down, the entire system shuts down.

Shortcomings. This fault tolerance approach has two severe shortcomings. First, it unnecessarily shuts down up to half of the cluster nodes, reducing the system’s performance and fault tolerance. Second, this approach causes a complete cluster shutdown if the surviving clique is missing a single data shard. To understand how likely a cluster is to shut down, we conduct a probabilistic analysis (Appendix A). Figure 4.2 shows the probability of a complete cluster shutdown while varying the cluster size and the number of nodes that shut down (i.e., nodes that are not part of the surviving clique – the x-axis in Figure 4.2). Each shard has three replicas. Our analysis shows that isolating only 10% of the nodes leads to more than a 50% probability of shutting down the entire cluster, and isolating only 20% of the nodes leads to a staggering 90% chance of a complete cluster shutdown.

Hazelcast Implementation. Hazelcast [27] offers in-memory sharded and replicated data structures. Every node in the system periodically sends a heartbeat to all nodes.

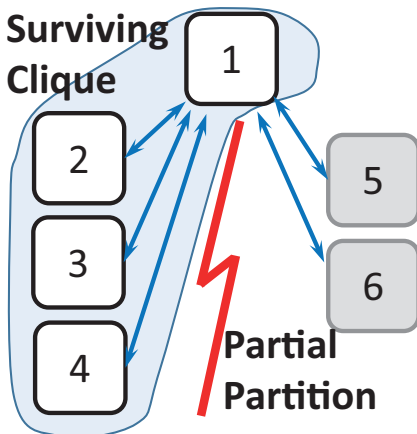


Figure 4.1: VoltDB’s surviving clique. Gray nodes shut down as they are not in the clique.

Hazelcast uses a master node to track the cluster membership, i.e., which nodes are part of the cluster. The master periodically sends a membership list to all nodes. A node will ignore membership updates coming from nodes that are not in the membership list.

Hazelcast escalates partial partitions to complete cluster partitions, such that the cluster is split into completely disconnected sub-clusters. When a partial partition occurs the master node collects connectivity information from all nodes. Nodes that are not reachable by the master are removed from the cluster membership list. The master then constructs a graph representing the cluster connectivity, runs the Bron–Kerbosch algorithm [82] to identify the largest fully connected sub-graph that includes the master node, removes all nodes that are not part of this sub-graph from the membership list, and broadcasts an updated membership list. For nodes that are removed from the membership list, Hazelcast supports two policies: pause or form a new cluster. Having two clusters serving the same application can lead to data inconsistency. When a partial partition heals, Hazelcast merges conflicting versions of the data using automated data consolidation policies (e.g., version with latest access time wins and discarding entries from the smaller cluster). Unfortunately, these policies can lose data or keep an inconsistent version of the data [74].

Shortcomings. This fault tolerance approach offers two undesirable alternatives. The cluster may unnecessarily pause a large number of nodes reducing the system’s performance and fault tolerance. Note that Hazelcast selects the largest subgraph that includes the master which may not include the majority of nodes. Alternatively, Hazelcast may form multiple clusters leading to data loss or inconsistency.

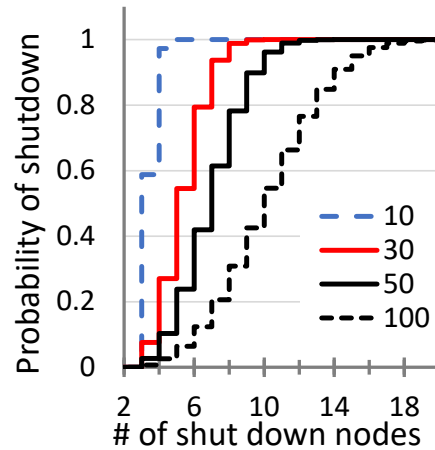


Figure 4.2: The probability of a VoltDB cluster shutdown. Different lines represent different cluster sizes. The x-axis shows the number of nodes that are not in the clique.

4.2 Checking Neighbors' Views

Main idea. When one node (e.g., node S) loses its connection to another node D, it verifies whether the connection is lost due to a partial partition. To this end, S asks all nodes in the cluster whether they can reach D. If a node reports that it can reach D, this indicates that the cluster is suffering a partial network partition.

If S detects a partial network partition, S either disconnects from all nodes that can reach D, which effectively makes the partition a complete partition, or pauses its operation. RabbitMQ and Elasticsearch follow this approach.

4.2.1 RabbitMQ

RabbitMQ [58] is a popular messaging system that replicates message queues for reliability. In RabbitMQ, if a node detects that its communication with another node (e.g., node D) is affected by a partial partition, it applies one of the following policies depending on its configuration.

1. *Escalate to a complete partition.* The node will drop its connection with any node that can reach node D. The goal of this policy is to create a complete partition in which both sides work independently. This configuration leads to data inconsistency and requires running a data consolidation mechanism after the partition heals.

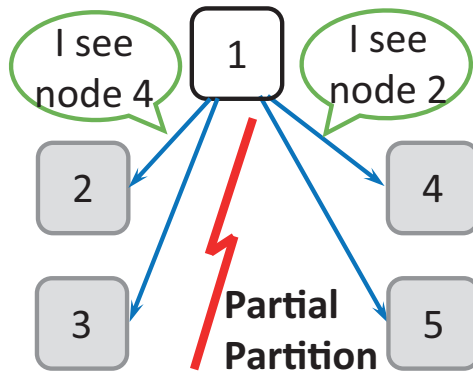


Figure 4.3: A scenario for RabbitMQ’s *pause* policy. Every non-bridge node pauses (gray nodes) as it detects that it cannot reach one node on the other side.

2. *Pause*: To avoid data inconsistency, once a node discovers the partial partition, it pauses its activities. It resumes its activities only when the partition heals. The result of this policy is that a subset of nodes will continue to operate. This subset will be completely connected and will run without sacrificing data consistency.
3. *Pause if anchor nodes are unreachable*: RabbitMQ’s configuration can specify a subset of nodes to act as anchor nodes. If a node cannot reach any of the anchor nodes, it pauses. This may lead to creating multiple complete partitions if the anchor nodes become partially partitioned. This may lead to pausing all nodes if all the anchor nodes are isolated.

After a partition heals, RabbitMQ employs two data consolidation techniques: administrator intervention, in which the administrator decides which side of the partition should become the authoritative version of the data, and auto-heal, in which the system makes this determination based on the number of clients connected to each side. Both techniques may lead to data loss or inconsistency [74].

Shortcomings. RabbitMQ’s policies have serious shortcomings. Changing a partial partition to a complete partition (policies 1 and 3) may lead to multiple inconsistent copies of the data, whereas the *pause* policy (policy 2) may pause the entire system or the majority of the nodes. For instance, in Figure 4.3, if every node except node 1 detects that it cannot reach a node on the other side of the partition, it pauses, leading to a complete cluster pause.

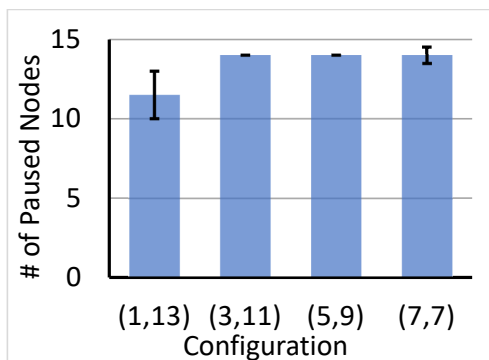


Figure 4.4: The median number of paused nodes in a cluster of 15 nodes. In all runs, one node is unaffected by the partition. The notation (i, j) shows the number of nodes on each side of the partition.

In the case of the *pause* policy (policy 2), to determine how many nodes pause under different partial partition scenarios, we conduct an experiment in which we deploy a 15-node RabbitMQ cluster, introduce a partial partition, and observe how many nodes pause. In all experiments, we inject a partition such that one node remains unaffected and able to reach all nodes. Figure 4.4 shows the median number of paused nodes under various partition configurations. We run each configuration 30 times. Surprisingly, in all configurations almost all the cluster nodes pause because each node detects that it cannot reach at least one node on the other side of the partition. Even isolating a single node (configuration (1,13) in Figure 4.4) leads to pausing 12 nodes. We experimented with additional configurations with a larger number of bridge nodes and noticed a similar behaviour (Appendix B). Our investigation reveals that nodes declare another node unreachable after missing its heartbeats for a timeout period. In RabbitMQ, the default timeout period is 1 minute, which gives enough time for many nodes to detect the partition and pause. Using a shorter timeout periods causes some nodes to declare prematurely that other nodes have failed, even without a partial partition.

4.2.2 Elasticsearch

Elasticsearch [23] is a popular search engine. Its master election protocol uses a fault tolerance technique based on checking neighbors' views. In Elasticsearch, the node with the lowest ID is the master. If a node (e.g., S) cannot reach the master, it contacts all nodes to check whether they can reach the master. If any node reports that it can reach the master, S pauses its operations. If none of the nodes can reach the master, the node with the lowest ID becomes the new master.

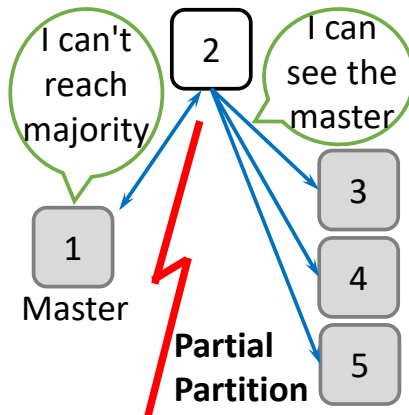


Figure 4.5: Elasticsearch unavailability scenario. The master pauses because it cannot reach majority of nodes, and all nodes pause because they cannot reach the master.

Shortcomings. First, this approach can affect cluster availability quite severely, as all nodes that cannot reach the master pause. In the worst case, it can cause a complete cluster unavailability. For instance, in Figure 4.5, none of the nodes can reach the master except node 2, which refuses to become the new master because it can reach a node with a lower ID (node 1). Consequently, all the nodes in the cluster pause. Furthermore, because the master cannot reach a majority of nodes, it also pauses, which leads to system unavailability [55]. Second, Elasticsearch uses this approach only to fortify the master election protocol, which leaves the rest of the system vulnerable to partial partitions.

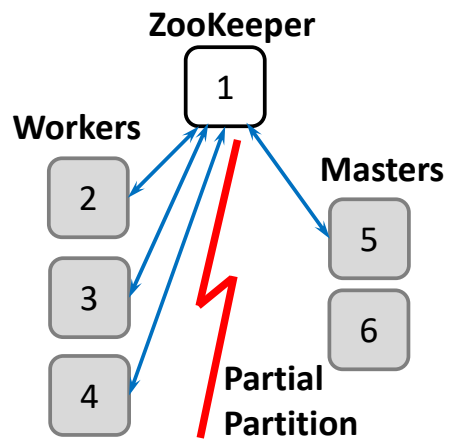


Figure 4.6: A Mesos cluster becomes unavailable when a partial partition isolates the master node and its backups.

	Surviving Clique		Checking w/ Neighbors		Failure Verification	Neutralizing Nodes		Nifty
	VoltDB	Hazelcast	Elasticsearch	RabbitMQ	MongoDB/LogCabin	MapReduce/HBase	Mesos	
Reduced Availability	\times^D	\times^P	\times^P	\times^P	\times^P	\times^D	\times^P	
Complete Unavailability	\times		\times	\times^1				
Complete Partition		\times		\times^2				
Double Execution							\times	
Data Unavailability		\times		\times^3				
Scope (System/Mechanism)	S	S	M	S	M	M	M	S

Table 4.1: Summary of shortcomings. (D) indicates that the nodes shut down. (P) indicates that the nodes pause until the partition heals. In the worst case, RabbitMQ pauses all nodes except one. We consider this a complete cluster loss (1). Under different RabbitMQ policies, (2) and (3) can occur. (S) indicates a system-wide technique, whereas (M) is a mechanism-specific technique.

4.3 Failure Verification

Main idea. If a node (e.g., S) receives a notification from another node that a third node (D) has failed, node S first verifies that it cannot reach D before taking any fault tolerance steps. This approach is used in the leader election protocols of MongoDB [45], and LogCabin [39]. It was also used in an earlier version of Elasticsearch.

In MongoDB and LogCabin, if a leader is on one side of a partial partition but can still reach the majority of nodes, the nodes on the other side of the partition unnecessarily call for leader election. Section 3.4 discusses a scenario in which a partial partition leads to continuous leader election thrashing and to system unavailability [8]. To avoid unnecessary elections, when a node receives a call for election, it first verifies that the current leader is unreachable. A node participates in an election only if it cannot reach the current leader, else it will ignore the failure report.

Shortcomings. This approach has two major shortcomings. First, it leads to the unavailability of a large number of nodes. Second, it is mechanism specific. Designing a system-wide fault tolerance mechanism using this approach is tricky because one cannot ignore every failure notification. For instance, using this approach in an earlier version of Elasticsearch backfired [25]. During data migration from a primary replica of a shard to a secondary replica, if a partial partition isolates the primary replica from the secondary replica while both are reachable from the master node, the primary requests a new secondary replica. Because the master can reach the secondary replica, it ignores the failure report. This leads to the unavailability of the affected shard [25]. Broadly applying this fault tolerance technique is not feasible because designers have to revisit the design of every system mechanism, consider the consequences of ignoring failure reports, and examine the interaction of various mechanisms under partial partitions.

4.4 Neutralizing Partitioned Nodes

Main idea. One challenge related to handling partial network partitions is that nodes may update a shared state that is reachable from both sides of the partition, leading to data loss and inconsistency. To avoid this problem, this approach attempts to neutralize one side of the partition. However, the neutralization method is implementation-specific. HBase, MapReduce, and Mesos use this approach.

HBase Implementation. In HBase, data shards are managed by an HBase node but are stored on HDFS. If the HBase leader cannot reach one of the HBase nodes, it neutralizes

that node by renaming the shard’s directory in HDFS. Renaming a shard’s directory effectively prohibits the old HBase node from making further changes to the shard [57]. The leader then assigns the shards of that node to a new HBase node.

MapReduce Implementation. In MapReduce, a manager node assigns tasks to AppMaster nodes. If the manager cannot reach an AppMaster, it reschedules the tasks assigned to that AppMaster to a new AppMaster. With partial network partitions, this approach may result in two AppMasters working on the same task, which leads to data corruption [41]. To fix this problem, when an AppMaster completes a task, it writes a completion record in a shared log on HDFS. Before an AppMaster executes a new task, it checks the shared log for a completion record. If it finds one, it does not re-execute the task.

Mesos Implementation. In Mesos, a master node assigns tasks to worker nodes. A Zookeeper instance selects the master node. The master sends periodic heartbeats to workers. If a partial partition isolates a worker node from the master, it pauses its operations. Figure 4.6 shows a worst-case scenario in which the partial partition isolates the master and its backup from all workers, which leads to a complete cluster unavailability. Finally, if a master detects that one of the workers is unavailable, it marks the tasks that were running on the unreachable worker as lost and reschedules them on new workers. This may lead to the double execution of a task [20].

Shortcomings. First, it is not practical to use this approach for system-wide fault tolerance, as this approach is specific to a certain protocol and implementation. The presented three systems use this approach for different mechanisms. To use this approach broadly, designers must go through the daunting task of independently designing a fault tolerance technique for every mechanism in the system and understanding the interaction between these mechanisms. Second, this approach leaves the nodes on one side of the partition idle, which reduces system performance and availability.

4.5 Summary

Table 4.1 summarizes the shortcomings of the current fault tolerance techniques, none of which are adequate for modern cloud systems. All current techniques severely affect system availability, as they unnecessarily lose a significant number of nodes. Failure verification and neutralizing partitioned nodes are used to fortify *specific* mechanisms, rather than providing *system-wide* fault tolerance. Using mechanism-specific fault tolerance techniques requires the independent fortification of all system mechanisms and the analysis of the interactions between various mechanisms. This approach complicates system design, fault

analysis, and debugging. An example of a system that uses multiple mechanism-specific techniques to tolerate partial partitions is Elasticsearch, which uses checking neighbors' view, failure verification [25], and neutralizing partitioned nodes [67] in different mechanisms. However, Elasticsearch has the highest number of reported failures due to partial partitions (Table 3.1).

Detecting the surviving clique and checking neighbors' views can be used to build a *system-wide* fault tolerance technique. However, as Table 4.1 shows, these techniques lead to a complete system shutdown or significant loss of system capacity. This realization motivated us to build Nifty (Chapter 5), a system-wide fault tolerance technique that overcomes the aforementioned shortcomings.

Chapter 5

Nifty Design

To overcome the limitations of current fault tolerance techniques, we design a simple, transparent network-partitioning fault-tolerant communication layer (Nifty).

Nifty follows a peer-to-peer design in which every node in the cluster runs a Nifty process. These processes collaborate in monitoring cluster connectivity. When Nifty detects a partial partition, it reroutes the traffic around the partition through intermediate nodes (i.e., bridge nodes). For instance, in Figure 5.1, if two partial partitions isolate node 1 from node 4, Nifty reroutes packets exchanged between nodes 1 and 4 through nodes 2 and 3.

Although Nifty keeps the cluster connected, it may increase the load on the bridge nodes, leading to a lower system performance. System designers who use Nifty may optimize the data or process placement or employ a flow-control mechanism to reduce the load on bridge nodes. To facilitate system-specific optimization, Nifty provides an API to identify bridge nodes and nodes on different sides of the partition, and to help take action when a partial partition occurs or heals.

Connectivity monitoring. Each Nifty process uses heart beating to monitor its connectivity with all other Nifty processes. Each Nifty process maintains a distance vector that includes the distance, in number of hops, to every node in the cluster. If a Nifty process misses five heartbeats from another Nifty process, it assumes that the communication with that process is broken and updates its distance vector. To detect when the communication between nodes recovers, Nifty processes continue to send heartbeats to disconnected nodes.

Recovery. Each Nifty process sends its distance vector (piggybacked on heartbeat messages) to all other nodes. Every Nifty process then uses these vectors to build and maintain a routing table.

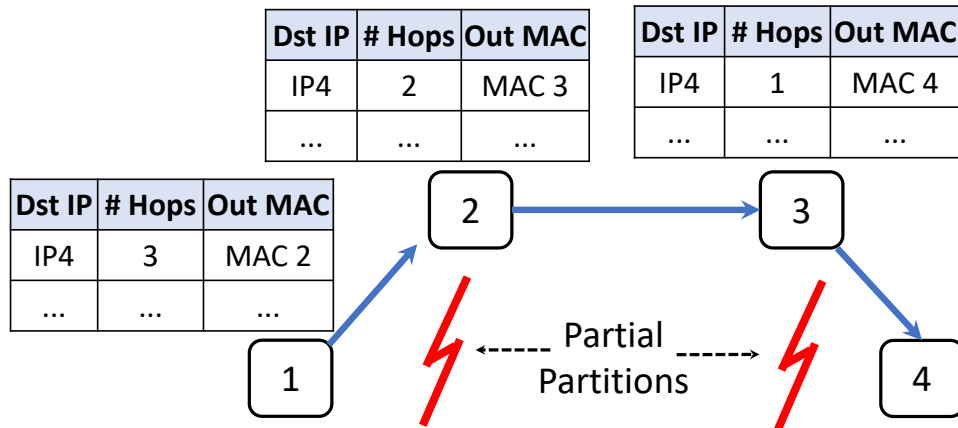


Figure 5.1: A Nifty routing example. A partial network partition isolates node 1 from nodes 3 and 4, and another partial partition isolates node 4 from nodes 1 and 2. Communication between 1 and 4 is routed through nodes 2 and 3.

When a Nifty process detects a change in the cluster (e.g., a node becomes unreachable or reachable), it initiates the route discovery procedure to find new routes. In our prototype, we use the classical Bellman–Ford distance-vector protocol [126, 77]. We use hop count as the link weight. By hop, we mean a hop between end nodes. Using hop count naturally favors direct connections, when they exist, over rerouting through intermediate nodes. Another possible alternative was to use a link state routing. We considered this options and preferred to go with distance vector routing since it is simpler to implement without the need for link state flooding, and is enough to tolerate all the network fault discussed in the tickets.

An entry in the routing table has a destination IP address, hop count, and output MAC address. If a packet is received with a destination IP address that matches an entry in the routing table, Nifty will change the destination MAC address of the packet to equal the output MAC address found in the routing table, then send the packet out.

Route deployment. Nifty uses OpenFlow [51] and Open vSwitch [133] to deploy the new routes. For instance, to reroute packets sent from node 1 to node 4 through nodes 2 and 3 in Figure 5.1, the Nifty process on node 1 installs rules on its local Open vSwitch to change the destination MAC address of any packet destined to node 4 to the MAC address of node 2. Whenever node 2 receives a packet with node 4 IP address as its destination, it changes the destination MAC address to node 3 MAC address and sends the packet out. Finally, when node 3 receives a packet with node 4 IP address, it changes the MAC address to node 4 MAC and sends the packet out.

Node classification. A system using Nifty can be optimized to reduce the amount of data forwarded through bridge nodes. The approach to do so is system-specific and may entail relocating processes in a cluster, dropping client requests, or reducing query result quality [81].

To facilitate the implementation of these mechanisms, Nifty offers an API that informs a system running atop Nifty when a partial partition happens and identifies which nodes are on the same side of the network partition and which nodes serve as bridge nodes. Section 7 demonstrates how this information can facilitate optimizing process placement in a VoltDB cluster.

5.1 Implementation

We implemented Nifty in 575 lines of C++ code. A Nifty process runs as a background process on all cluster nodes. A configuration file lists the IP addresses of all cluster nodes. Each Nifty node heartbeats all the nodes listed in the configuration file. The heartbeat message is sent over UDP packets. The default heartbeat period is 200 ms. A node assumes it cannot reach another node if it misses three heartbeats from that node. We note that this is relatively aggressive heartbeating. The goal is for Nifty to discover the partial partition and create alternative routes before the system atop detects the partition with its own heartbeating mechanism. We found that Raft has the shortest heartbeating periods of 250 ms, hence we choose to heartbeat every 200 ms. Nevertheless, the heartbeat period is configurable. Nifty uses the Bellman-Ford routing algorithm to find routes between end nodes. We used the `ovs-ofctl` program to manipulate Open vSwitch rules.

Nifty API. To facilitate building system specific optimizations, Nifty provides an API. The API mainly notifies the system when a partial partition happens and exposes the cluster connectivity graph to the system running atop Nifty. The current prototype offers a Java wrapper to simplify integrating Nifty with the systems we used in our evaluation.

Listing 5.1 shows the main functions in the NIFTYAPI. The API has two groups of functions. The call back functions are triggered when the network state changes. Systems use the query APIs to find the network topology.

The user must override the abstract methods `atPartialPartition()`, `atHealthyNetwork()` and `atCompletePartition()`. Nifty calls these functions when the network state changes. To identify if a partition is complete or partial, Nifty uses depth-first search to traverse the connectivity graph. Depth first search fails to reach all the nodes in a complete partition.

Listing 5.1: The NiftyAPI

```
// Call back functions
abstract void atPartialPartition ();
abstract void atHealthyNetwork ();
abstract void atCompletePartition ();

// Query APIs
Graph getNetworkGraph ();
boolean isBridgeNode ();
NetState getNetworkState ();
```

Chapter 6

Evaluation

Our evaluation answers three questions. How much overhead does Nifty impose when there are no network partitions? What is a system’s performance with Nifty under a network partition? What is the utility of Nifty’s classification API?

Testbed. We conduct our experiments using 40 x1170 nodes at the Cloudfab Utah cluster. Each node has an Intel Xeon E5 10-core CPU, 64 GB of RAM, and a Mellanox ConnectX-4 25 Gbps NIC. To inject a network partition fault, we modify the Open vSwitch rules on the nodes to drop packets between the affected nodes. In all our experiments, we report the average for 30 runs. We note that the standard deviation in all our experiments is lower than 5%.

6.1 Overhead Evaluation

To evaluate Nifty’s overhead, we measure its impact on the performance of a synthetic benchmark using iperf [34] and seven data-centric systems (i.e., storage, database, and messaging systems). The iperf experiment uses a 100-node cluster to measure Nifty’s impact on larger clusters. The systems we selected are:

- HDFS: We deploy HDFS (v3.3.0) on six nodes (one name node and five data nodes) and with a replication level of three. To avoid disk access, we configure data nodes to use tmpfs. We use the HDFS standard benchmark (TestDFSIO). The benchmark reads and writes 1 GB files.
- Kafka: We deploy Kafka (v2.6.0) on five nodes. We distribute the queues (aka, topics) among nodes to balance the load. Each message is replicated on three nodes.

We use Kafka’s benchmarking tool to generate load on the system. The experiments use a set of producers and consumers. Each producer sends messages to a dedicated queue and each queue has one consumer.

- **ActiveMQ:** We deploy ActiveMQ Artemis (v2.15.0) on five nodes with each queue being replicated on two nodes. The experiments use a set of producers and consumers. Each producer sends messages to a dedicated queue and each queue has one consumer.
- **MongoDB:** We deploy MongoDB (v4.4.1) on six nodes (one config server and five mongod nodes) and with a replication level of three. We discuss our results with the Yahoo benchmark workload B (95% reads and 5% writes) with a uniform distribution [87]. We use 10 million records. The rest of the Yahoo benchmark workloads show similar results.
- **VoltDB:** We deploy VoltDB (v9.0) on nine nodes, with data sharding enabled and a replication level of three. We use the Yahoo benchmark and the TPC-C benchmark. Figure 6.2.a shows the throughput-latency curve under Yahoo benchmark workload B (95% reads and 5% writes) with a uniform distribution. The results using the TPC-C benchmark and the Yahoo benchmark workloads A and C with uniform and skewed loads show similar low overhead.
- **RabbitMQ:** We deploy RabbitMQ (v3.8.2) on three nodes. We use the mirrored mode in which each queue has a leader replica and two backup replicas. We distribute the queue masters among brokers to distribute the load. The experiments use a set of producers and consumers. Each producer sends messages to a dedicated queue and each consumer reads messages from a dedicated queue.
- **Redis PubSub:** We deploy Redis (v6.2) on three nodes. One publisher connects to one node (i.e., root node) and continuously publishes 1 KB messages to one topic. With Redis PubSub, the root Redis node forwards the published messages to the other Redis nodes. The subscribers connect to the other two Redis nodes.

Results. We compare the throughput and average latency of each system with and without Nifty when there is no partial network partition. We evaluate Nifty with a partial partition in Section 6.2.

Figure 6.1 shows the write throughput of HDFS (Figure 6.1.a) and the throughput-latency curve for Kafka (Figure 6.1.b), ActiveMQ (Figure 6.1.c), MongoDB (Figure 6.1.d). Figure 6.2 shows the throughput-latency curve for VoltDB (Figure 6.2.a) and RabbitMQ (Figure 6.2), and the throughput figure for Redis PubSub (Figure 6.2.c). The results show

that Nifty does not add noticeable overhead; for all systems, the curves almost completely overlap. This is because Nifty processes exchange a negligible number of packets. Each Nifty process sends a single UDP heartbeat packet every 200 ms to other nodes in the system. Consequently, in the largest deployment of nine nodes, each node sends only 40 packets every second.

Scalability evaluation. Nifty uses all-to-all heart beating to monitor a cluster’s connectivity. Consequently, Nifty’s overhead increases with the cluster size. To measure Nifty’s scalability, we evaluate its overhead on a 100 m510 nodes at the CloudLab Utah cluster. Each node has an ARMv8 (Atlas/A57) 8-core CPU, 64 GB of RAM, and a Mellanox ConnectX-3 10 Gbps NIC. For this experiment, we limit the throughput of each node to 1 Gbps, as CloudLab can not support a full 10-Gbps connectivity between the 100 nodes we managed to book. To generate network intensive load, we use iperf [34]. Half of the nodes run an iperf server, and the other half run an iperf client. Each client communicates with a single server. Figure 6.3 shows the aggregate throughput of the iperf servers when deployed with and without Nifty. The figure shows that Nifty’s overhead is negligible. When using 100 nodes, Nifty degrades the aggregate throughput by only 3.5%. Nevertheless, this monitoring approach will not scale to clusters with thousands of nodes. We are currently exploring the design of a fault tolerance technique that can scale to larger clusters.

6.2 Handling Partial Partitions

To demonstrate the effectiveness of the proposed approach, we evaluate Nifty’s performance with the seven aforementioned systems under a partial partition fault. We note that RabbitMQ and VoltDB implemented two different techniques for tolerating partial partitions (Chapter 4).

Partial partition setup. We use the same deployment of the seven aforementioned systems. Each system is deployed on an odd number of replicas. We introduce a partial partition that leaves one node as a bridge node and puts an equal number of nodes on each side of the partition. Client nodes are not affected by the partition. We partition the cluster this way to create maximum pressure on the bridge node.

Figure 6.1 and Figure 6.2 show the system performance when the cluster suffers from the partial partition. We notice that all the seven systems are severely effected by the partial partition. ActiveMQ, MongoDB, and VoltDB suffer a complete cluster pause or shutdown when deployed without Nifty. HDFS fails almost all write operations. The VoltDB cluster shuts down because, after detecting the surviving clique, the system misses at least one shard. This confirms our analysis in Section 4.1.

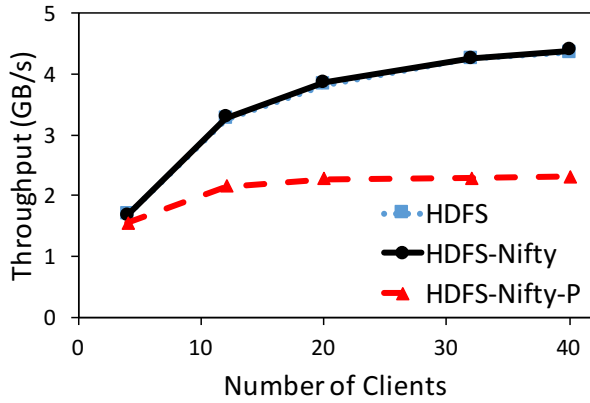
RabbitMQ uses the checking neighbor’s views fault tolerance approach. In our deployment, each queue is mirrored on a backup replica. Due to the strong consistency requirement, we configure RabbitMQ to pause in case of partial partition. We deploy RabbitMQ on three nodes. Unfortunately, we could not use a larger RabbitMQ cluster because partial partitions often lead to the pause of the entire RabbitMQ cluster when Nifty is not used (Figure 4.3). Even with three nodes, partial partitions sometimes lead to pausing two out of three nodes. We discard those results and only include results in which one node pauses. Consequently, our results show the best possible performance of RabbitMQ under partial partitions. Pausing a broker in RabbitMQ leads to more than 50% reduction in throughput (RabbitMQ-P in (Figure 6.2.b)).

In Redis PubSub, if a partial partition isolates a node that receives new messages from another Redis node, Redis will fail to deliver the message to subscribers connected to the isolated node, leading to 50% reduction in throughput (Figure 6.2.c).

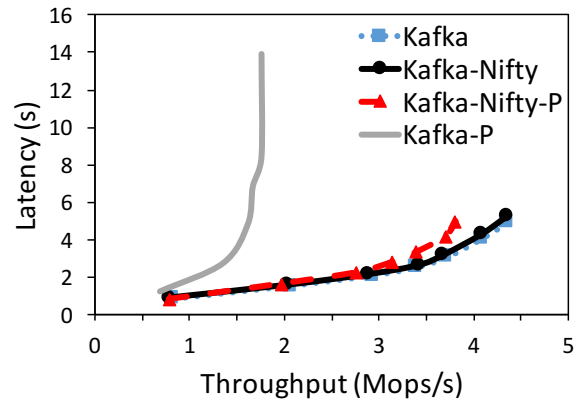
Kafka uses Zookeeper to monitor cluster nodes. If a partial partition isolates a queue leader from the majority of replicas while Zookeeper runs on a bridge node, Zookeeper will not select a new leader and the entire cluster pauses (Finding 1 in Chapter 3). To mitigate this, we made sure that Zookeeper falls on one side of the partition. In this case, all the nodes on the other side of the partition that cannot reach Zookeeper are removed from the cluster. In our experiment, the partial partition causes two nodes to pause, which leads to almost a 50% reduction in system throughput (Figure 6.1.b).

Figure 6.1 and Figure 6.2 show that Nifty effectively masks the partial partition, so none of the nodes shut down or pause. Figure 6.1.a shows the write operation throughput for HDFS. With a replication level of three, each file has replicas on both sides of a partial partition. Consequently, for every 1 GB of data written, up to 2 GB of data are rerouted through the bridge node. This reduces the system throughput by up to 45%. We note that having a partial partition result in a performance degradation is better than a complete system unavailability when HDFS is deployed without Nifty. We present an optimization for HDFS that alleviates this problem in Section 7. For the rest of the systems, during the partial partition, almost 50% of client requests and responses are rerouted through the bridge node. Even so, the system throughput only decreases by 2-6.7% and latency only increases by 3-7.8%. This shows that Nifty can effectively mask partial partitions and is able to utilize remaining connections to reduce the performance impact.

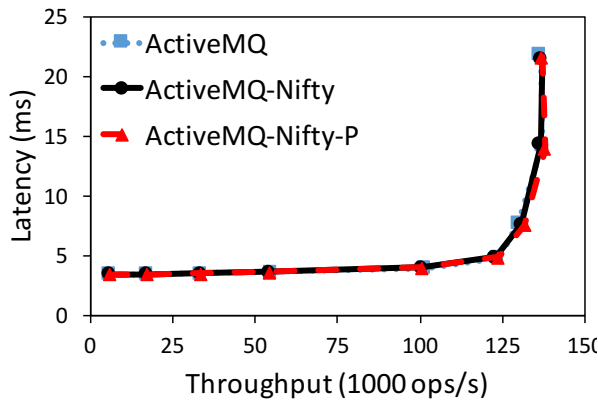
Figure 6.4 shows the tail latency for VoltDB and RabbitMQ for the same experiments presented in Figure 6.2. The figure shows the average throughput and the 99th percentile of latency while increasing the load on the system. The figure shows that Nifty increases the 99th percentile latency by up to 6.8% without a partial partition and by 15% under a partial partition failure.



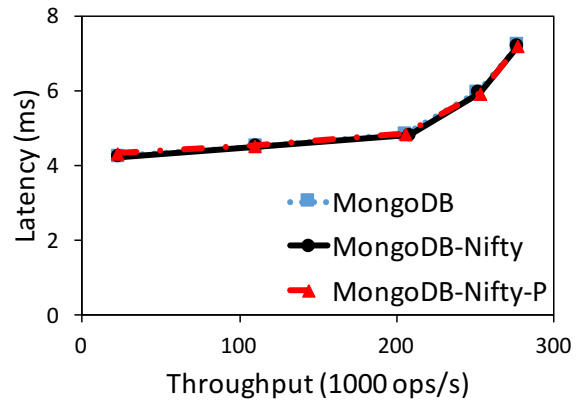
(a) HDFS write throughput



(b) Kafka

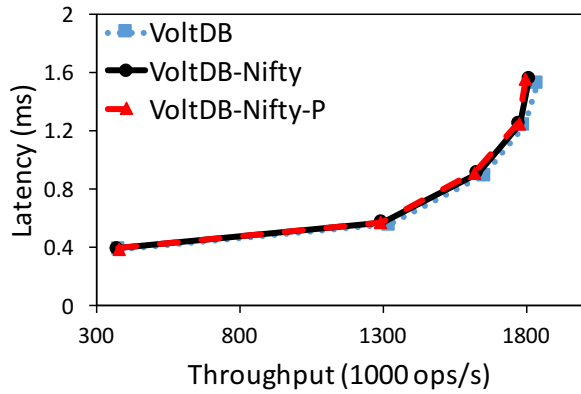


(c) ActiveMQ

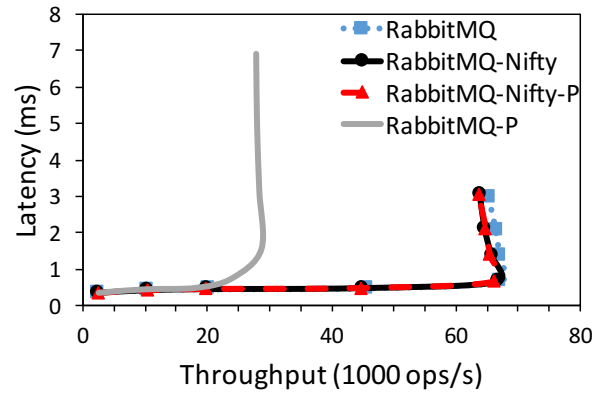


(d) MongoDB

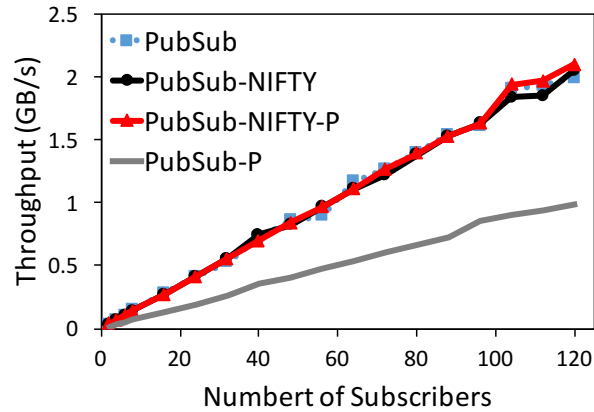
Figure 6.1: Nifty’s overhead. The average throughput for HDFS (a) and the average throughput vs. average latency Kafka, ActiveMQ, and MongoDB. (-P) denotes the results with a partial partition.



(a) VoltDB



(b) RabbitMQ



(c) Redis PubSub

Figure 6.2: Nifty’s overhead. the average throughput vs. average latency for VoltDB, RabbitMQ, and the average throughput for Redis PubSub. (-P) denotes the results with a partial partition.

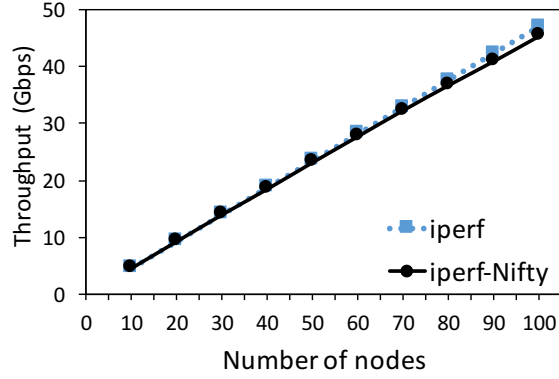
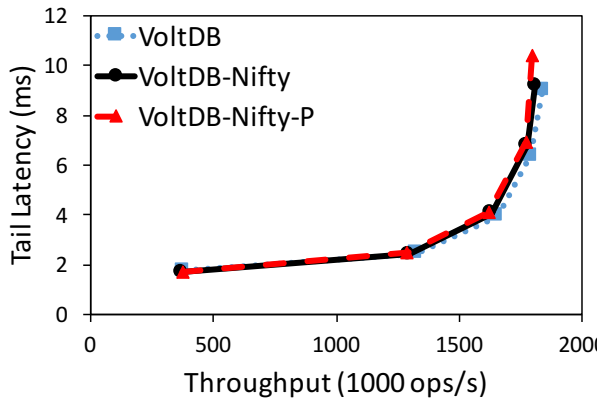
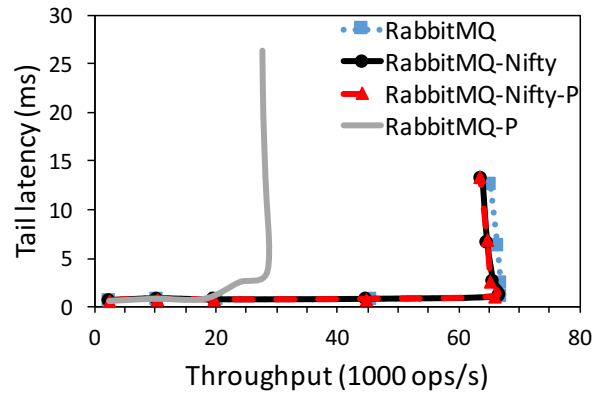


Figure 6.3: Scalability evaluation. Average throughput while increasing the number of nodes.



(a) VoltDB tail latency.



(b) RabbitMQ tail latency.

Figure 6.4: Tail latency evaluation. Average throughput vs. 99th percentile of latency.

Chapter 7

Classification API Utility

While Nifty reroutes packets through bridge nodes to restore a cluster connectivity during a partial partition, it has the potential to significantly increase the load on bridge nodes and significantly reduce performance. Our evaluation with HDFS (Figure 6.1a) shows one scenario in which during a partial partition a system using Nifty experiences close to 50% drop in throughput. Given that network partitions may last hours [143, 144, 101, 102], this performance degradation is highly undesirable.

A system using Nifty can be optimized to reduce the amount of data forwarded through bridge nodes. The approach to do so is system-specific. Nifty offers an API to facilitate the implementation of these mechanisms 5.1.

To demonstrate the benefit of using Nifty’s API we modified the implementation of two mechanisms to reduce the amount of data rerouted through the bridge nodes. We modified the data placement protocol in HDFS, and the processing of multi-shard operations in VoltDB.

7.1 HDFS

HDFS uses chain replication to replicate write operations. Chain replication arranges replicas in a chain. Each node passes the write operations to its successor. When selecting three data nodes for a new data chunk, the name node tries to balance the number of blocks across nodes and select nodes located on more than one rack. Under partial partition large volumes of data can be rerouted through the bridge node. In the worst case, with three-way replication, the same data may traverse bridge nodes twice. Figure 6.1a shows that system

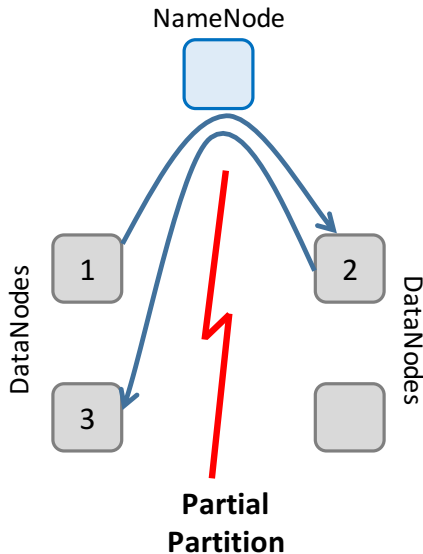


Figure 7.1: HDFS worst case rerouting. NameNode choosing the replicas to be on 1,2, and 3 will cause the data to move across the partition twice.

throughput degrades by up to 45% under partial partitions. During this experiment clients write 48 GB of data, and the bridge node rerouted 39.2 GB of data. Our probabilistic analysis shows that with replication level of 3 the total volume of rerouted traffic through the bridge nodes will be equivalent to 85% of the data written by the clients.

To improve the system performance under partial partitions we used the Nifty’s API to implement three optimizations:

- Optimized chain ordering (Opt.-Chain). In the worst case, when the system is configured with a replication level of three, a newly written data blocks may be rerouted twice through the bridge nodes (Figure 7.1). Our probabilistic analysis shows that in our experimental setup of seven data nodes with one node being a bridge, there is a 17% chance to reroute a block twice through the bridge node.

We modified the HDFS data placement algorithm to avoid the situation in which data is routed twice through bridge nodes (Figure 7.1). After the data placement mechanism picks three replicas, we modified the code to query Nifty to get the network topology. We use the network topology to order the replicas in a chain in such a way to forward data through bridge nodes at most once.

- Two replicas (2-Replicas). Another approach to avoid rerouting the same data twice through bridge nodes is to temporally reduce the replication level to 2 during network partitions. Once the partial partition heals the NameNode can create additional replicas of the affected data chunks.
- Optimized data placement (One-side). In this alternative we modified the HDFS data placement algorithm to query Nifty to identify the cluster topology under partial partitions then for any new data chunk allocate three data nodes on the same side of the partition or bridge nodes. This effectively eliminates any data rerouting through the bridge nodes.

We note that these policies affect data written during a partial partition. When the partition heals, Nifty returns to the original replication factor or placement policy. For the 2-Replicas policy, the system will create an additional replicas after the partition heals.

Results. We deploy HDFS on eight nodes: one name node, and seven data nodes on the same cluster detailed in Chapter 6. The partial partition is injected to put three data nodes on each side of the partition and keep one bridge node. Clients run on dedicated machines and use the TestDFSIO benchmark to write to and read 1 GB files. We use the default replication factor of three. We show the average of 30 runs. The maximum standard deviation of these experiments is 4.7%.

Figure 7.2 shows the system throughput while varying the number of clients. We use the performance of HDFS without a partial partition as a baseline (Baseline in Figure 7.2). The figure shows that when the partial partition is injected, Nifty, without using any optimization, achieves up to 41% of the Baseline throughput. This is mainly because 85% of the client data is rerouted through the bridge node which creates a system bottleneck. In the worst case, a client data will be forwarded twice through the bridge node during the replication step. This scenario accounted for 34% of the forwarded data. The Opt.-Chain optimization guarantees that each write operation is rerouted at most once through the bridge node. This optimization reduces the amount of data rerouted through the bridge node to 68% of the client data and achieves up to 59% of the Baseline throughput. The 2-replicas optimization further reduces the replication overhead, but also reduces the durability guarantees for data written during a partial partition. This optimization achieves 81% of the Baseline throughput. Figure 7.2 shows that the one-side optimization achieves a throughput comparable to the Baseline under the partial partition. This is because this optimization avoids rerouting any client data through the bridge nodes.

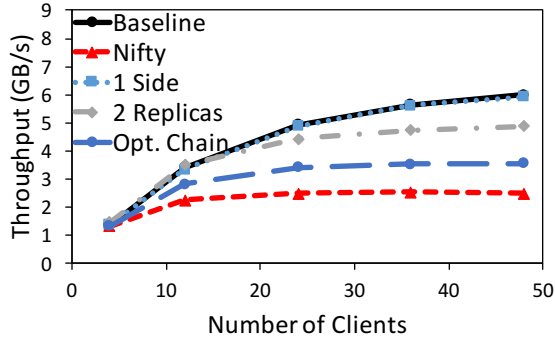


Figure 7.2: HDFS write throughput with different optimizations

7.2 VoltDB

In VoltDB, a single server (aka, multi-data-partition initiator or MPI) processes all multi-shard operations. The MPI divides a multi-shard query (e.g., a join) to sub-queries, such that each sub-query targets a single shard. The MPI forwards each sub-query to its shard leader, gathers the intermediate results, performs final query processing, and sends the result to the client.

When deploying VoltDB atop Nifty, if the MPI node is on one side of the partition, a potentially significant volume of intermediate data passes through the bridge node. In our setup, when the MPI is on one side of the partition, 50% of the intermediate results are rerouted through the bridge node. This increases operation latency and the load on bridge nodes.

To improve the performance of multi-shard operations, the MPI process can be migrated to a bridge node. This effectively eliminates the need to reroute any traffic for multi-shard queries. We modify VoltDB to use Nifty’s API to identify bridge nodes and migrate the MPI to a bridge node.

To evaluate this optimization’s effectiveness, we evaluate the effect of the MPI’s location on system performance. We restrict clients to contacting VoltDB nodes on one side of the partition and compare the system performance of three MPI placements: on clients side of the partition (client side in Figure 7.3), on the bridge node (bridge), and on the side opposite to the clients (opposite side). Bridge placement represents our optimization.

Setup and Workload. We use the same VoltDB configuration and partial partition setup detailed in the previous sections. Unfortunately, VoltDB has limited support for join queries, so it cannot run standard benchmarks such as TPC-H [70]. In our experiments, we

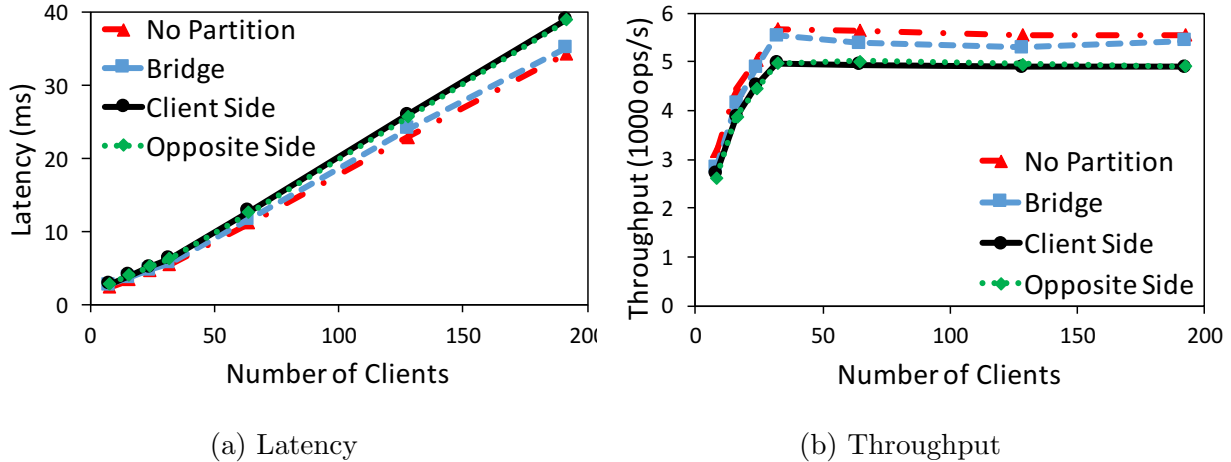


Figure 7.3: The impact of MPI placement on VoltDB’s performance. Figure shows the average latency (a) and average throughput (b). Standard deviation was less than 2%.

use a simple synthetic benchmark that joins two tables. The benchmark has two sharded tables of 20 fields each. Each field is 50 bytes, leading to approximately 1 KB rows. To use multiple shards, clients issue a range query that joins the two tables on the primary key. The client issues a query with a range that includes four primary keys. Consequently, the query result size is limited to four rows, with a total size of almost 8 KB. We populate the database with 20 GB of data before running the experiments. We report the average and standard deviation for 30 runs.

Results. Figure 7.3 shows the system throughput (a) and the average latency (b) for the three possible MPI placements. During a partial partition fault, placing the MPI on a bridge node decreases the latency by up to 11% and improves throughput by 11% compared to client and opposite side placements. Placing the MPI on a bridge node reduces the number of hops the join query must make before the MPI accumulates all the results and sends the query reply. Furthermore, bridge placement achieves throughput and latency within 4% of VoltDB’s performance when there is no partition (“no partition” in Figure 7.3).

We measure the amount of data forwarded through the bridge nodes for each one of those configurations; placing the MPI on the bridge node imposes the least overhead. When using 128 clients, 72 MB, 5 GB, and 6.5 GB of data are forwarded through the bridge node when the MPI is placed on the bridge, client side, and opposite side, respectively. The opposite side reroutes more data than the client side placement, as the client request and the result are also rerouted through the bridge node.

Chapter 8

Related Work

To the best of our knowledge, this is the first study to focus on partial network partitioning, characterize its failures, identify design pitfalls in common distributed systems techniques, dissect modern fault tolerance techniques, and explore the design of a generic fault tolerance technique for this type of fault.

Failure studies. A number of previous efforts analyzed failures in distributed systems, including characterizing specific component failures [101, 144, 145, 78, 97, 110] and characterizing failures in a specific domain such as HPC [94, 120, 136], IaaS clouds [76], data-mining services [154], hosting services [132, 104], data-intensive systems [135, 103, 118], and cloud systems [151]. Our work complements these efforts by focusing on failures triggered by partial network partitions.

Yuan et al. [151] conducted a study on 198 general user-reported failures from six distributed systems. Yuan et al. [151] reports 24% of their failures to be catastrophic failures, while our work shows a much higher percentage in partial network partition failures (76%). Our work also shows that less than 2% of the failures we found are nondeterministic, compared to 26% of general failures they found. This clearly shows that partial network partition failures have more severe effects on the systems than general failures. This also shows that testers should be able to catch partial network partition failures with well written tests since almost all of them are deterministic.

In our previous work [74], we studied 136 network partitioning failures focusing on complete partitions. This previous work identified partial partitions, presented examples of how they can lead to system failures, and presented NEAT, a testing tool that can inject complete and partial network partitioning faults. We use NEAT to reproduce some of the reported failures. This paper presents an in-depth analysis of partial partition failures and fault tolerance techniques and proposes a novel fault-tolerant communication layer.

Complete Network Partitions Comparing the characteristics of partial and complete partitions [74] shows that they have similar catastrophic impact and manifestation and reproducibility characteristics. Partial partitions seem easier to manifest. While all partial partition failures are triggered by a single-node partial partition and almost all of the failures are deterministic, 88% of the complete partitions manifest by isolating a single node and 80% of them are deterministic. Furthermore, we found twice as many failure reports reporting complete partitions than partial partitions.

Despite their similarity in causing catastrophic failures and being easy-to-manifest, partial and complete partitions are fundamentally different faults. Unlike complete partitions, a cluster suffering a partial partition is still connected but not all-to-all connected. Consequently, the CAP theorem bounds [100] do not apply to partial partitions. Furthermore, fault tolerance techniques for complete partitions cannot handle partial partitions or lead to pausing up to half of the cluster nodes. For instance, using majority vote to elect a leader is an effective mechanism to tolerate complete partitions. This approach alone is not effective in handling partial partitions, as there could be multiple completely connected subgroups with each connecting a majority of nodes. Chapter 4 shows how using only majority voting can lead to leader election thrashing and system unavailability.

Overlay Networks. Using the hosts in a cluster of node to create an overlay network like Nifty does to mask partial partitions is a common idea. RON [75] creates an overlay network to recover from path outages in the internet. Unlike RON, Nifty focuses on partial partitions in data center networks, works in the MAC layer, and allows for much faster convergence times. VTrace [95] helps to diagnose persistent packet loss in cloud overlay networks, without taking any action to help with the packets loss like Nifty. DHT Systems like Chord[138] and Symphony[123] use overlay networks to route requests in peer-to-peer distributed storage systems. BDS [153] uses overlay network routing for optimizing inter-datacenter data replication. [84] investigates using overlay networks for better performance in content delivery. Nifty differs from these systems in its purpose, which is to keep a cluster fully connected in the case of partial partitions.

Leveraging SDN for overlay networks. Software-defined networking capabilities have been used to engineer traffic and optimize system operations. Google Andromeda [91] uses SDN for Google Cloud Platform’s network virtualization stack. In a similar fashion [137] shows how overlay networks can be created in OpenStack using SDN.

Other works have focused on using SDN to create or optimize certain systems. [146, 152] use SDN for performing different network measurements tasks like QoS measurements or anomaly detection. [113] has used SDN to implementing an in-network stateful firewalls. [80, 142] show the use of SDN to create load balancers to optimize key-value stores and distributed Memcached deployments. [112, 119] use SDN for key-value-based routing.

Nifty is similar in spirit to these systems, as we use Open vSwitch capabilities to implement an overlay. Our goal however is different: to improve systems fault-tolerance by masking partial network partitions.

Chapter 9

Concluding Remarks

Our work sheds light on a peculiar type of infrastructure fault and highlights the need for further research to understand such faults and explore techniques to improve systems' resiliency.

This is the first work to focus on partial network partitioning fault and present an in-depth analysis of system failures triggered by this fault. We identify characteristics that can facilitate better test design. Our findings highlight that focused design reviews can identify vulnerabilities early in the design process. We dissect the implementation of eight popular systems and study their fault tolerance techniques. In doing so, we identify four main approaches for tolerating partial partitions. Unfortunately, all implemented fault tolerance techniques have severe shortcomings.

We, therefore, build Nifty to overcome the limitations of modern fault tolerance techniques. Nifty is a simple, transparent communication layer that reroutes packets around partial partitions. We note that modern systems already incorporate a membership and connectivity monitoring. We show that extending the current implementations with a detour mechanism is an effective and low overhead fault tolerance technique to partial partitions. The source code for Nifty is available at <https://github.com/UWASL/NIFTY>

References

- [1] Activemq cluster blocks indefinitely in the presence of partial network partition. <https://issues.apache.org/jira/browse/AMQ-7064>. Accessed: June 2021.
- [2] Activemq: Flexible & powerful open source multi-protocol messaging. <http://activemq.apache.org/>. Accessed: June 2021.
- [3] Another journey of chaos engineering (stitch fix). <https://www.datadoghq.com/videos/another-journey-of-chaos-engineering/>. Accessed: June 2021.
- [4] The apache hadoop project. <http://hadoop.apache.org/>. Accessed: June 2021.
- [5] Apache hbase. <https://hbase.apache.org/>. Accessed: June 2021.
- [6] Apache mesos. <http://mesos.apache.org/>. Accessed: June 2021.
- [7] Apache zookeeper. <https://zookeeper.apache.org/>. Accessed: June 2021.
- [8] Arbiters in pv1 should vote no in elections if they can see a healthy primary of equal or greater priority to the candidate. <https://jira.mongodb.org/browse/SERVER-27125>. Accessed: June 2021.
- [9] Asymmetrical network partition can cause the election of two primary nodes. <https://jira.mongodb.org/browse/SERVER-9730>. Accessed: June 2021.
- [10] Autoheal crash during partial partitions and blocks nodes. <https://github.com/rabbitmq/rabbitmq-server/issues/928>. Accessed: June 2021.
- [11] Balancer can cause cascading mongod failures during network partitions. <https://jira.mongodb.org/browse/SERVER-19550>. Accessed: June 2021.
- [12] bnx2 cards intermittantly going offline. <https://www.spinics.net/lists/netdev/msg152880.html>. Accessed: June 2021.

- [13] The ceph object store. <https://ceph.io/>. Accessed: June 2021.
- [14] Cisco data center infrastructure 2.5 design guide. Cisco Systems, Inc., 2011.
- [15] Cloudflare blog: A byzantine failure in the real world. <https://blog.cloudflare.com/a-byzantine-failure-in-the-real-world/>. Accessed: June 2021.
- [16] Cluster broken after switches upgrade. <https://github.com/elastic/elasticsearch/issues/9495>. Accessed: June 2021.
- [17] cluster with pause_minority survive node less than 1/2 after network disconnect. <https://github.com/rabbitmq/rabbitmq-server/issues/1667>. Accessed: June 2021.
- [18] Data center: Load balancing data center, solutions reference network design. Technical report, Cisco Systems, Inc., 2004.
- [19] Datadog: Learning from aws failure. <https://www.datadoghq.com/blog/gray-aws-failures/>. Accessed: June 2021.
- [20] Designing highly available mesos frameworks. <http://mesos.apache.org/documentation/latest/high-availability-framework-guide/>. Accessed: June 2021.
- [21] Disconnect between coordinating node and shards can cause duplicate updates or wrong status code #9967. <https://github.com/elastic/elasticsearch/issues/9967>. Accessed: June 2021.
- [22] Dkron: A distributed cron service. <https://dkron.io/>. Accessed: June 2021.
- [23] Elasticsearch: Distributed search & analytics. <https://www.elastic.co/products/elasticsearch>. Accessed: June 2021.
- [24] Elasticsearch resiliency status. <https://www.elastic.co/guide/en/elasticsearch/resiliency/current/index.html>. Accessed: June 2021.
- [25] Faulty recovery caused by partial network partitions. <https://github.com/elastic/elasticsearch/pull/8720>. Accessed: June 2021.
- [26] Google cloud networking incident #18003. <https://status.cloud.google.com/incident/cloud-networking/18003>. Accessed: June 2021.

- [27] Hazelcast — the leading in-memory computing platform. <https://hazelcast.com/>. Accessed: June 2021.
- [28] Hazelcast: the leading in-memory data grid. <https://hazelcast.com/>. Accessed: June 2021.
- [29] Hbase-3446: Processservershutdown fails if meta moves, orphaning lots of regions. <https://issues.apache.org/jira/browse/HBASE-3446>. Accessed: June 2021.
- [30] Hdfs-1384: Namenode should give client the first node in the pipeline from different rack other than that of excludednodes list in the same rack. <https://issues.apache.org/jira/browse/HDFS-1384>. Accessed: June 2021.
- [31] Healthchecking is not transitive. <https://www.robustperception.io/healthchecking-is-not-transitive>. Accessed: June 2021.
- [32] How does voltdb handle partial network partitions? <https://www.voltdb.com/resources/transaction-consistency-faq#net>. Accessed: June 2021.
- [33] If block report races with closing of file, replica is incorrectly marked corrupt. <https://issues.apache.org/jira/browse/HDFS-2791>. Accessed: June 2021.
- [34] iperf: The ultimate speed test tool for tcp, udp and sctp. <https://iperf.fr/>. Accessed: June 2021.
- [35] Jepsen: A framework for distributed systems verification, with fault injection. <https://github.com/jepsen-io/jepsen>. Accessed: June 2021.
- [36] Kafka-3686: Kafka producer is not fault tolerant. <https://issues.apache.org/jira/browse/KAFKA-3686>. Accessed: June 2021.
- [37] Kafka: A distributed streaming platform. <https://kafka.apache.org/>. Accessed: June 2021.
- [38] Kafka leader election doesn't happen when leader broker port is partitioned off the network. <https://issues.apache.org/jira/browse/KAFKA-8702>. Accessed: June 2021.
- [39] Logcabin. <https://github.com/logcabin/logcabin>. Accessed: June 2021.
- [40] Lyft engineering: Operating apache kafka clusters 24/7 without a global ops team. <https://eng.lyft.com/operating-apache-kafka-clusters-24-7-without-a-global-ops-team-417813a5ce70>. Accessed: June 2021.

- [41] Mapreduce ticket 4832. <https://issues.apache.org/jira/browse/MAPREDUCE-4832>. Accessed: June 2021.
- [42] Mesos-1529: Handle a network partition between master and slave. <https://issues.apache.org/jira/browse/MESOS-1529>. Accessed: June 2021.
- [43] minimum_master_nodes does not prevent split-brain if splits are intersecting. <https://github.com/elastic/elasticsearch/issues/2488>. Accessed: June 2021.
- [44] Mirrored queue crash with out of sync acks. <https://github.com/rabbitmq/rabbitmq-server/issues/749>. Accessed: June 2021.
- [45] MongoDB: The database for modern applications. <https://www.mongodb.com/>. Accessed: June 2021.
- [46] Moosefs: Distributed file system. <https://moosefs.com/>. Accessed: June 2021.
- [47] A network partition can cause in flight documents to be lost. <https://github.com/elastic/elasticsearch/issues/7572>. Accessed: June 2021.
- [48] Nodemangers die on startup if they can't connect to the rm. <https://issues.apache.org/jira/browse/MAPREDUCE-3963>. Accessed: June 2021.
- [49] Observability in paxos clusters. <https://davecturner.github.io/2017/08/18/observability-in-paxos.html>. Accessed: June 2021.
- [50] Onos 1.4 test plan - ha. <https://wiki.onosproject.org/pages/viewpage.action?pageId=7439437>. Accessed: June 2021.
- [51] Openflow switch specification, version 1.5.1 (onf ts-025). Open Networking Foundation, 2015.
- [52] Partial network partition and retries. <https://github.com/elastic/elasticsearch/issues/6105>. Accessed: June 2021.
- [53] Partial network partition and retries. <https://github.com/elastic/elasticsearch/issues/6105>. Accessed: June 2021.
- [54] Partial network partition and retries #6105. <https://github.com/elastic/elasticsearch/issues/6105>. Accessed: June 2021.
- [55] Partial network partitioning leads to cluster unavailability. <https://github.com/elastic/elasticsearch/issues/43183>. Accessed: June 2021.

- [56] Partial network partitions and obstacles to innovation. <https://rachelbythebay.com/w/2012/02/16/partition/>. Accessed: June 2021.
- [57] Possible data loss when rs goes into gc pause while rolling hlog. <https://issues.apache.org/jira/browse/HBASE-2312>. Accessed: June 2021.
- [58] Rabbitmq message broker. <https://www.rabbitmq.com>. Accessed: June 2021.
- [59] Rabbitmq performance testing tool. <https://github.com/rabbitmq/rabbitmq-perf-test>. Accessed: June 2021.
- [60] Redis: in-memory data structure store. <https://redis.io/>. Accessed: June 2021.
- [61] Replicaset servers able to make a putsch due to bad network timeout parameters. <https://jira.mongodb.org/browse/SERVER-26216>. Accessed: June 2021.
- [62] Splitlogmanger async delete node hangs log splitting when zk connection is lost. <https://issues.apache.org/jira/browse/HBASE-5606>. Accessed: June 2021.
- [63] Synchronisation causes crash in duplicated master #1006. <https://github.com/rabbitmq/rabbitmq-server/issues/1006>. Accessed: June 2021.
- [64] Two primaries with network partitioned replica set (non-transient). <https://jira.mongodb.org/browse/SERVER-2544>. Accessed: June 2021.
- [65] Using map output fetch failures to blacklist nodes is problematic. <https://issues.apache.org/jira/browse/MAPREDUCE-1800>. Accessed: June 2021.
- [66] Voltldb in-memory database platform. <https://www.voltldb.com/>. Accessed: June 2021.
- [67] Wait on shard failures. <https://github.com/elastic/elasticsearch/issues/14252>. Accessed: June 2021.
- [68] Zookeeper recipes and solutions. <https://zookeeper.apache.org/doc/current/recipes.html>. Accessed: June 2021.
- [69] TPC-C benchmark standard specification. Transaction Processing Performance Council, February 2010. Revision 5.11.
- [70] TPC-H benchmark (decision support) standard specification. Transaction Processing Performance Council, December 2018. Revision 2.18.0.

- [71] Mohammed Alfatafta, Basil Alkhatib, Ahmed Alquraan, and Samer Al-Kiswany. Toward a generic fault tolerance technique for partial network partitioning. In *14th USENIX Symposium on Operating Systems Design and Implementation (OSDI 20)*, pages 351–368. USENIX Association, November 2020.
- [72] Mohammed Alfatafta, Basil Alkhatib, Ahmed Alquraan, and Samer Al-Kiswany. Understanding partial network partitioning. Technical Report WASL-TR-2020-02, Waterloo Advanced Systems Lab, University of Waterloo, October 2020.
- [73] Mohammed Alfatafta, Zuhair AlSader, and Samer Al-Kiswany. Cool: A cloud-optimized structure for mpi collective operations. In *2018 IEEE 11th International Conference on Cloud Computing (CLOUD)*, pages 746–753. IEEE, 2018.
- [74] Ahmed Alquraan, Hatem Takruri, Mohammed Alfatafta, and Samer Al-Kiswany. An analysis of network-partitioning failures in cloud systems. In *13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18)*, pages 51–68, 2018.
- [75] David Andersen, Hari Balakrishnan, Frans Kaashoek, and Robert Morris. Resilient overlay networks. 35(5):131–145, October 2001.
- [76] Theophilus Benson, Sambit Sahu, Aditya Akella, and Anees Shaikh. A first look at problems in the cloud. *HotCloud*, 10:15, 2010.
- [77] Dimitri P Bertsekas, Robert G Gallager, and Pierre Humblet. *Data networks*, volume 2. Prentice-Hall International New Jersey, 1992.
- [78] Robert Birke, Ioana Giurgiu, Lydia Y Chen, Dorothea Wiesmann, and Ton Engbersen. Failure analysis of virtual and physical machines: patterns, causes and characteristics. In *2014 44th Annual IEEE/IFIP International Conference on Dependable Systems and Networks*, pages 1–12. IEEE, 2014.
- [79] Pat Bosshart, Dan Daly, Glen Gibb, Martin Izzard, Nick McKeown, Jennifer Rexford, Cole Schlesinger, Dan Talayco, Amin Vahdat, George Varghese, et al. P4: Programming protocol-independent packet processors. *ACM SIGCOMM Computer Communication Review*, 44(3):87–95, 2014.
- [80] Anat Bremler-Barr, David Hay, Idan Moyal, and Liron Schiff. Load balancing memcached traffic using software defined networking. In *2017 IFIP Networking Conference (IFIP Networking) and Workshops*, pages 1–9. IEEE, 2017.

- [81] Eric A Brewer. Lessons from giant-scale services. *IEEE Internet computing*, 5(4):46–55, 2001.
- [82] Coen Bron and Joep Kerbosch. Algorithm 457: Finding all cliques of an undirected graph. *Commun. ACM*, 16(9):575–577, September 1973.
- [83] Nathan Bronson, Zach Amsden, George Cabrera, Prasad Chakka, Peter Dimov, Hui Ding, Jack Ferris, Anthony Giardullo, Sachin Kulkarni, Harry Li, et al. TAO: Facebook’s distributed data store for the social graph. In *Presented as part of the 2013 USENIX Annual Technical Conference (USENIX ATC 13)*, pages 49–60, 2013.
- [84] John Byers, Jeffrey Considine, Michael Mitzenmacher, and Stanislav Rost. Informed content delivery across adaptive overlay networks. SIGCOMM ’02, page 47–60, New York, NY, USA, 2002. Association for Computing Machinery.
- [85] Tushar D Chandra, Robert Griesemer, and Joshua Redstone. Paxos made live: an engineering perspective. In *Proceedings of the twenty-sixth annual ACM symposium on Principles of distributed computing*, pages 398–407. ACM, 2007.
- [86] Xin Chen, Charng-Da Lu, and Karthik Pattabiraman. Failure analysis of jobs in compute clouds: A google cluster case study. In *2014 IEEE 25th International Symposium on Software Reliability Engineering*, pages 167–177. IEEE, 2014.
- [87] 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 ’10, page 143–154, New York, NY, USA, 2010. Association for Computing Machinery.
- [88] 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*, pages 143–154. ACM, 2010.
- [89] James C Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, Jeffrey John Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, et al. Spanner: Google’s globally distributed database. *ACM Transactions on Computer Systems (TOCS)*, 31(3):8, 2013.
- [90] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2009.

- [91] Michael Dalton, David Schultz, Jacob Adriaens, Ahsan Arefin, Anshuman Gupta, Brian Fahs, Dima Rubinstein, Enrique Cauich Zermeno, Erik Rubow, James Alexander Docauer, et al. Andromeda: performance, isolation, and velocity at scale in cloud network virtualization. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 373–387, 2018.
- [92] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Vosshall, and Werner Vogels. Dynamo: amazon’s highly available key-value store. In *ACM SIGOPS operating systems review*, volume 41, pages 205–220. ACM, 2007.
- [93] Dmitry Duplyakin, Robert Ricci, Aleksander Maricq, Gary Wong, Jonathon Duerig, Eric Eide, Leigh Stoller, Mike Hibler, David Johnson, Kirk Webb, Aditya Akella, Kuangching Wang, Glenn Ricart, Larry Landweber, Chip Elliott, Michael Zink, Emmanuel Cecchet, Snigdhaswin Kar, and Prabodh Mishra. The design and operation of CloudLab. In *Proceedings of the USENIX Annual Technical Conference (ATC)*, pages 1–14, July 2019.
- [94] Nosayba El-Sayed and Bianca Schroeder. Reading between the lines of failure logs: Understanding how hpc systems fail. In *2013 43rd annual IEEE/IFIP international conference on dependable systems and networks (DSN)*, pages 1–12. IEEE, 2013.
- [95] Chongrong Fang, Haoyu Liu, Mao Miao, Jie Ye, Lei Wang, Wansheng Zhang, Daxiang Kang, Biao Lyv, Peng Cheng, and Jiming Chen. Vtrace: Automatic diagnostic system for persistent packet loss in cloud-scale overlay network. In *Proceedings of the Annual Conference of the ACM Special Interest Group on Data Communication on the Applications, Technologies, Architectures, and Protocols for Computer Communication, SIGCOMM ’20*, page 31–43, New York, NY, USA, 2020. Association for Computing Machinery.
- [96] Daniel Firestone, Andrew Putnam, Sambhrama Mundkur, Derek Chiou, Alireza Dabagh, Mike Andrewartha, Hari Angepat, Vivek Bhanu, Adrian Caulfield, Eric Chung, Harish Kumar Chandrappa, Somesh Chaturmohta, Matt Humphrey, Jack Lavier, Norman Lam, Fengfen Liu, Kalin Ovtcharov, Jitu Padhye, Gautham Popuri, Shachar Raindel, Tejas Sapre, Mark Shaw, Gabriel Silva, Madhan Sivakumar, Nisheeth Srivastava, Anshuman Verma, Qasim Zuhair, Deepak Bansal, Doug Burger, Kushagra Vaid, David A. Maltz, and Albert Greenberg. Azure accelerated networking: Smartnics in the public cloud. In *15th USENIX Symposium on Networked Systems Design and Implementation (NSDI 18)*, pages 51–66, Renton, WA, April 2018. USENIX Association.

- [97] Daniel Ford, François Labelle, Florentina Popovici, Murray Stokely, Van-Anh Truong, Luiz Barroso, Carrie Grimes, and Sean Quinlan. Availability in globally distributed storage systems. 2010.
- [98] Peter Garraghan, Paul Townend, and Jie Xu. An empirical failure-analysis of a large-scale cloud computing environment. In *2014 IEEE 15th International Symposium on High-Assurance Systems Engineering*, pages 113–120. IEEE, 2014.
- [99] Seth Gilbert and Nancy Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *Acm Sigact News*, 33(2):51–59, 2002.
- [100] Seth Gilbert and Nancy Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. *SIGACT News*, 33(2):51–59, June 2002.
- [101] Phillipa Gill, Navendu Jain, and Nachiappan Nagappan. Understanding network failures in data centers: measurement, analysis, and implications. *ACM SIGCOMM Computer Communication Review*, 41(4):350–361, 2011.
- [102] Ramesh Govindan, Ina Minei, Mahesh Kallahalla, Bikash Koley, and Amin Vahdat. Evolve or die: High-availability design principles drawn from googles network infrastructure. In *Proceedings of the 2016 ACM SIGCOMM Conference*, pages 58–72. ACM, 2016.
- [103] Haryadi S Gunawi, Mingzhe Hao, Tanakorn Leesatapornwongsa, Tiratat Patananake, Thanh Do, Jeffry Adityatama, Kurnia J Eliazar, Agung Laksono, Jeffrey F Lukman, Vincentius Martin, et al. What bugs live in the cloud? a study of 3000+ issues in cloud systems. In *Proceedings of the ACM Symposium on Cloud Computing*, pages 1–14. ACM, 2014.
- [104] 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 Seventh ACM Symposium on Cloud Computing*, pages 1–16. ACM, 2016.
- [105] Nikhil Handigol, Mario Flajslik, Srini Seetharaman, Nick McKeown, and Ramesh Johari. Aster* x: Load-balancing as a network primitive. In *9th GENI Engineering Conference (Plenary)*, pages 1–2, 2010.
- [106] A. Herr. Veritas cluster server 6.2 I/O fencing deployment considerations. Technical report, Veritas Technologies, 2016.

- [107] Robert V. Hogg, Elliot Tanis, and Dale Zimmerman. *Probability and Statistical Inference*. Pearson, 9 edition, 2013.
- [108] Patrick Hunt, Mahadev Konar, Flavio Paiva Junqueira, and Benjamin Reed. Zookeeper: Wait-free coordination for internet-scale systems. In *USENIX annual technical conference*, volume 8, 2010.
- [109] Sushant Jain, Alok Kumar, Subhasree Mandal, Joon Ong, Leon Poutievski, Arjun Singh, Subbaiah Venkata, Jim Wanderer, Junlan Zhou, Min Zhu, et al. B4: Experience with a globally-deployed software defined wan. In *ACM SIGCOMM Computer Communication Review*, volume 43, pages 3–14. ACM, 2013.
- [110] Weihang Jiang, Chongfeng Hu, Yuanyuan Zhou, and Arkady Kanevsky. Are disks the dominant contributor for storage failures?: A comprehensive study of storage subsystem failure characteristics. *ACM Transactions on Storage (TOS)*, 4(3):7, 2008.
- [111] Flavio P Junqueira, Benjamin C Reed, and Marco Serafini. Zab: High-performance broadcast for primary-backup systems. In *2011 IEEE/IFIP 41st International Conference on Dependable Systems & Networks (DSN)*, pages 245–256. IEEE, 2011.
- [112] I. Kettaneh, A. Alquraan, H. Takruri, S. Yang, A. S. Dusseau, R. Arpaci-Dusseau, and S. Al-Kiswany. The network-integrated storage system. *IEEE Transactions on Parallel and Distributed Systems*, 2019.
- [113] Pakapol Krongbaramee and Yuthapong Somchit. Implementation of sdn stateful firewall on data plane using open vswitch. In *2018 15th International Joint Conference on Computer Science and Software Engineering (JCSSE)*, pages 1–5. IEEE, 2018.
- [114] Leslie Lamport et al. Paxos made simple. *ACM Sigact News*, 32(4):18–25, 2001.
- [115] Adrian Lara, Anisha Kolasani, and Byrav Ramamurthy. Network innovation using openflow: A survey. *IEEE communications surveys & tutorials*, 16(1):493–512, 2013.
- [116] Jialin Li, Ellis Michael, and Dan RK Ports. Eris: Coordination-free consistent transactions using in-network concurrency control. In *Proceedings of the 26th Symposium on Operating Systems Principles*, pages 104–120. ACM, 2017.
- [117] Jialin Li, Ellis Michael, Naveen Kr Sharma, Adriana Szekeres, and Dan RK Ports. Just say NO to paxos overhead: Replacing consensus with network ordering. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*, pages 467–483, 2016.

- [118] 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 2013 International Conference on Software Engineering*, pages 963–972. IEEE Press, 2013.
- [119] Xiaozhou Li, Raghav Sethi, Michael Kaminsky, David G Andersen, and Michael J Freedman. Be fast, cheap and in control with switchkv. In *13th USENIX Symposium on Networked Systems Design and Implementation (NSDI 16)*, pages 31–44, 2016.
- [120] Yinglung Liang, Yanyong Zhang, Anand Sivasubramaniam, Morris Jette, and Ramendra Sahoo. Bluegene/l failure analysis and prediction models. In *International Conference on Dependable Systems and Networks (DSN'06)*, pages 425–434. IEEE, 2006.
- [121] Barbara Liskov and James Cowling. Viewstamped replication revisited. Technical Report MIT-CSAIL-TR-2012-021, MIT, July 2012.
- [122] Christian Maihofer. A survey of geocast routing protocols. *IEEE Communications Surveys & Tutorials*, 6(2):32–42, 2004.
- [123] G. S. Manku, M. Bawa, and P. Raghavan. Symphony: Distributed hashing in a small world. In *USENIX Symposium on Internet Technologies and Systems*, 2003.
- [124] Simon J Maple and Ian Robinson. Transaction recovery in a transaction processing computer system employing multiple transaction managers, October 20 2015. US Patent 9,165,025.
- [125] Ali José Mashtizadeh, Min Cai, Gabriel Tarasuk-Levin, Ricardo Koller, Tal Garfinkel, and Sreekanth Setty. Xvmotion: Unified virtual machine migration over long distance. In *2014 USENIX Annual Technical Conference (USENIXATC 14)*, pages 97–108, 2014.
- [126] Deep Medhi and Karthik Ramasamy. *Network routing: algorithms, protocols, and architectures*. Morgan Kaufmann, 2017.
- [127] Matthew Milano and Andrew C Myers. Mixt: a language for mixing consistency in geodistributed transactions. In *Proceedings of the 39th ACM SIGPLAN Conference on Programming Language Design and Implementation*, pages 226–241. ACM, 2018.
- [128] Ankur Kumar Nayak, Alex Reimers, Nick Feamster, and Russ Clark. Resonance: dynamic access control for enterprise networks. In *Proceedings of the 1st ACM workshop on Research on enterprise networking*, pages 11–18. ACM, 2009.

- [129] Keita Nomura, Yoshiaki Taniguchi, Nobukazu Iguchi, and Kenzi Watanabe. A system for supporting migration to overlay openflow network using openstack. In *2016 10th International Conference on Complex, Intelligent, and Software Intensive Systems (CISIS)*, pages 595–598. IEEE, 2016.
- [130] Brian M. Oki and Barbara H. Liskov. Viewstamped replication: A new primary copy method to support highly-available distributed systems. In *Proceedings of the Seventh Annual ACM Symposium on Principles of Distributed Computing, PODC '88*, page 8–17, New York, NY, USA, 1988. Association for Computing Machinery.
- [131] Diego Ongaro and John Ousterhout. In search of an understandable consensus algorithm. In *2014 USENIX Annual Technical Conference (USENIX ATC 14)*, pages 305–319, 2014.
- [132] David Oppenheimer, Archana Ganapathi, and David A Patterson. Why do internet services fail, and what can be done about it? In *USENIX symposium on internet technologies and systems*, volume 67. Seattle, WA, 2003.
- [133] Ben Pfaff, Justin Pettit, Teemu Koponen, Ethan Jackson, Andy Zhou, Jarno Rajahalme, Jesse Gross, Alex Wang, Joe Stringer, Pravin Shelar, et al. The design and implementation of open vswitch. In *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI 15)*, pages 117–130, 2015.
- [134] Dan RK Ports, Jialin Li, Vincent Liu, Naveen Kr Sharma, and Arvind Krishnamurthy. Designing distributed systems using approximate synchrony in data center networks. In *12th USENIX Symposium on Networked Systems Design and Implementation (NSDI 15)*, pages 43–57, 2015.
- [135] Ariel Rabkin and Randy Howard Katz. How hadoop clusters break. *IEEE software*, 30(4):88–94, 2012.
- [136] Bianca Schroeder and Garth Gibson. A large-scale study of failures in high-performance computing systems. *IEEE transactions on Dependable and Secure Computing*, 7(4):337–350, 2009.
- [137] Piyush Raman Srivastava and Saket Saurav. Networking agent for overlay l2 routing and overlay to underlay external networks l3 routing using openflow and open vswitch. In *2015 17th Asia-Pacific Network Operations and Management Symposium (APNOMS)*, pages 291–296. IEEE, 2015.

- [138] Ion Stoica, Robert Morris, David Karger, M. Frans Kaashoek, and Hari Balakrishnan. Chord: A scalable peer-to-peer lookup service for internet applications. *SIGCOMM Comput. Commun. Rev.*, 31(4):149–160, August 2001.
- [139] Michael Stonebraker and Ariel Weisberg. The voltdb main memory dbms. *IEEE Data Eng. Bull.*, 36(2):21–27, 2013.
- [140] Douglas B Terry, Vijayan Prabhakaran, Ramakrishna Kotla, Mahesh Balakrishnan, Marcos K Aguilera, and Hussam Abu-Libdeh. Consistency-based service level agreements for cloud storage. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 309–324. ACM, 2013.
- [141] Douglas B Terry, Marvin M Theimer, Karin Petersen, Alan J Demers, Mike J Spreitzer, and Carl H Hauser. Managing update conflicts in bayou, a weakly connected replicated storage system. In *SOSP*, volume 95, pages 172–182, 1995.
- [142] Alex FR Trajano and Marcial P Fernandez. Two-phase load balancing of in-memory key-value storages through nfv and sdn. In *2015 IEEE Symposium on Computers and Communication (ISCC)*, pages 409–414. IEEE, 2015.
- [143] Daniel Turner, Kirill Levchenko, Jeffrey C Mogul, Stefan Savage, Alex C Snoeren, Daniel Turner, Kirill Levchenko, Jeffrey C Mogul, Stefan Savage, and Alex C Snoeren. On failure in managed enterprise networks. *HP Labs HPL-2012-101*, 2012.
- [144] Daniel Turner, Kirill Levchenko, Alex C Snoeren, and Stefan Savage. California fault lines: understanding the causes and impact of network failures. *ACM SIGCOMM Computer Communication Review*, 41(4):315–326, 2011.
- [145] Kashi Venkatesh Vishwanath and Nachiappan Nagappan. Characterizing cloud computing hardware reliability. In *Proceedings of the 1st ACM symposium on Cloud computing*, pages 193–204. ACM, 2010.
- [146] An Wang, Yang Guo, Songqing Chen, Fang Hao, TV Lakshman, Doug Montgomery, and Kotikalapudi Sriram. vprom: Vswitch enhanced programmable measurement in sdn. In *2017 IEEE 25th International Conference on Network Protocols (ICNP)*, pages 1–10. IEEE, 2017.
- [147] Richard Wang, Dana Butnariu, Jennifer Rexford, et al. Openflow-based server load balancing gone wild. *Hot-ICE*, 11:12–12, 2011.
- [148] Robin J. Wilson. *Introduction to Graph Theory*. Prentice Hall/Pearson, New York, 2010.

- [149] Zhe Wu, Michael Butkiewicz, Dorian Perkins, Ethan Katz-Bassett, and Harsha V Madhyastha. Spanstore: Cost-effective geo-replicated storage spanning multiple cloud services. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 292–308. ACM, 2013.
- [150] Nezih Yigitbasi, Matthieu Gallet, Derrick Kondo, Alexandru Iosup, and Dick Epema. Analysis and modeling of time-correlated failures in large-scale distributed systems. In *2010 11th IEEE/ACM International Conference on Grid Computing*, pages 65–72. IEEE, 2010.
- [151] 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 *11th USENIX Symposium on Operating Systems Design and Implementation (OSDI 14)*, pages 249–265, 2014.
- [152] Zili Zha, An Wang, Yang Guo, Doug Montgomery, and Songqing Chen. Instrumenting open vswitch with monitoring capabilities: designs and challenges. In *Proceedings of the Symposium on SDN Research*, page 16. ACM, 2018.
- [153] Yuchao Zhang, Junchen Jiang, Ke Xu, Xiaohui Nie, Martin J. Reed, Haiyang Wang, Guang Yao, Miao Zhang, and Kai Chen. Bds: A centralized near-optimal overlay network for inter-datacenter data replication. In *Proceedings of the Thirteenth EuroSys Conference*, EuroSys '18, New York, NY, USA, 2018. Association for Computing Machinery.
- [154] Hucheng Zhou, Jian-Guang Lou, Hongyu Zhang, Haibo Lin, Haoxiang Lin, and Tingting Qin. An empirical study on quality issues of production big data platform. In *Proceedings of the 37th International Conference on Software Engineering-Volume 2*, pages 17–26. IEEE Press, 2015.

APPENDICES

Appendix A: The Probability of VoltDB Cluster Shutdown

We consider a VoltDB cluster with N nodes. The cluster stores S shards with a replication factor of R . When a partial network partition happens, VoltDB identifies the surviving clique and all the nodes that are not part of the clique shutdown. We denote the number of nodes that shutdown due to a partial partition as F (Since F is not in the surviving clique then $F < \frac{N}{2}$), leaving the system with $(N - F)$ surviving nodes.

Assumptions. We assume that:

1. The system selects R nodes to hold the replicas of a given shard using a uniform random distribution.
2. Shard placement is independent of other shards locations.
3. Each node has enough capacity to store all the shards.

VoltDB will shut down if the surviving clique does not have all the shards. This means that if the F failed nodes contain all the R replicas of any of the S shards, then VoltDB shuts down. In other terms, the VoltDB cluster will survive a partial partition if every shard has at least one replica in the surviving clique.

Step I. Single Shard Probability. Consider the case of a system with a single shard. The system will survive in all cases in which the surviving clique has at least a single replica of the shard. To compute the probability a system will survive a partial partition, we will compute the number of possible replica placements in the cluster, then compute how many of those placements would fail when losing F nodes. Finally, we will use these two numbers to compute the probability a system survives a partial partitioning fault.

Number of possible combinations to place a shard. The system selects R nodes to hold the replicas for a shard. The selection is without replacement since no two copies of the shard can be placed on the same node, and the order of the selected nodes is not important.

The number of possible combinations for placing a shard is:

$$C(N, R) = \frac{N!}{(N - R)! \times R!} \quad (1)$$

If F nodes fail, the number of combinations in which all the replicas of the shard are on the F failed nodes is (this is again without replacement and ordering is not important)

$$C(F, R) = \frac{F!}{(F - R)! \times R!} \quad (2)$$

The probability the system shuts down when F nodes fail is

$$\begin{aligned} P(\text{single_shard_shutdown}) &= \frac{C(N, R)}{C(F, R)} \\ &= \frac{\frac{N!}{(N-R)! \times R!}}{\frac{F!}{(F-R)! \times R!}} = \frac{F!(N - R)!}{N!(F - R)!} \end{aligned} \quad (3)$$

And the probability of a system surviving a partial partition that shuts down F nodes is

$$P(\text{single_shard_surviving}) = 1 - \frac{F!(N - R)!}{N!(F - R)!} \quad (4)$$

Step II. Multi-shard probability. Assuming that shards are placed independently, the probability of a system with S shard surviving a partial partition with F nodes shutting down is:

$$P(\text{system_surviving}) = \left(1 - \frac{F!(N - R)!}{N!(F - R)!}\right)^S \quad (5)$$

The probability the system shuts down is

$$P(\text{system_shutdown}) = 1 - \left(1 - \frac{F!(N - R)!}{N!(F - R)!}\right)^S \quad (6)$$

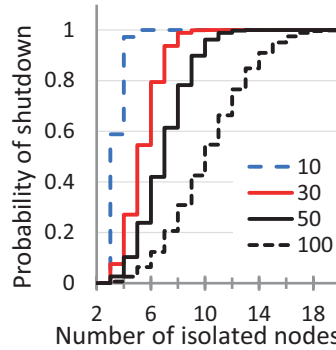


Figure 1: Probability of a VoltDB system shut down.

Example. We used this formula to compute the probability of failure of VoltDB on different cluster sizes. The number of shards VoltDB allocates to nodes is equal to the number of cores. Figure 1 shows the probability of VoltDB shutting down for different cluster sizes, with replication level of 3, and assuming nodes with 32 CPU cores. The figure shows that isolating only 10% of the nodes leads to over 50% probability of shutting down the entire cluster, and isolating only 20% of the nodes leads to a staggering 90% chance for a complete cluster shutdown.

Appendix B: The Impact of Partial Partitions on RabbitMQ

RabbitMQ’s has two main policies for handling partial partitions. The first policy changes a partial partition to a complete partition which may lead to multiple inconsistent copies of the data. The second is the *pause* policy which preserves data consistency but may lead to pausing the entire system or the majority of its nodes.

To determine how many nodes pause when using the pause policy, we conducted an experiment in which we deployed a 15-node RabbitMQ cluster, introduced a partial partition, and observed how many nodes paused. In all experiments, we inject a partition such that one node remained unaffected and able to reach all nodes. Figures 4.4, 2, and 3 show the median number of paused nodes under various partition configurations. We ran each configuration 30 times. Note that the maximum number of nodes that can pause is 14, 12, and 10 in Figures 4.4, 2, and 3, respectively, because the rest of the nodes are bridge nodes that can reach all the nodes in the cluster. Surprisingly, under all partial partition scenario

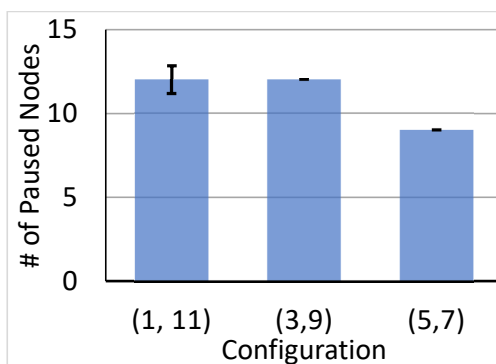


Figure 2: The median number of paused nodes in a cluster of 15 nodes. In all runs, 3 node are unaffected by the partition. The notation (i, j) shows the number of nodes on each side of the partition.

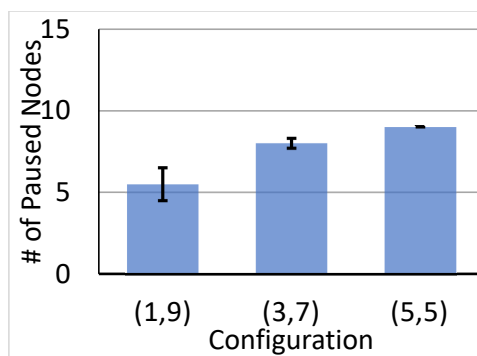


Figure 3: The median number of paused nodes in a cluster of 15 nodes. In all runs, 5 node are unaffected by the partition. The notation (i, j) shows the number of nodes on each side of the partition.

a significant number of the affected nodes are paused. Our investigation of this failure scenario revealed that nodes declare another node unreachable after missing its heartbeats for a *timeout* period. In RabbitMQ, the timeout period is 1 minute by default, which gives enough time for many nodes to detect the partition and pause. We experimented with significantly shorter timeout periods, but that caused some nodes to prematurely declare that other nodes had failed, even without a partial partition.