

1 Minha: Large-Scale Distributed Systems Testing 2 Made Practical

3 **Nuno Machado** 

4 Teradata & INESC TEC
5 nuno.machado@teradata.com

6 **Francisco Maia** 

7 INESC TEC
8 franciso.a.maia@inesctec.pt

9 **Francisco Neves** 

10 INESC TEC & U.Minho
11 franciso.t.neves@inesctec.pt

12 **Fábio Coelho** 

13 INESC TEC & U.Minho
14 fabio.a.coelho@inesctec.pt

15 **José Pereira** 

16 INESC TEC & U.Minho
17 jop@di.uminho.pt

18 — Abstract —

19 Testing large-scale distributed system software is still far from practical as the sheer scale needed
20 and the inherent non-determinism make it very expensive to deploy and use realistically large
21 environments, even with cloud computing and state-of-the-art automation. Moreover, observing
22 global states without disturbing the system under test is itself difficult. This is particularly troubling
23 as the gap between distributed algorithms and their implementations can easily introduce subtle
24 bugs that are disclosed only with suitably large scale tests.

25 We address this challenge with MINHA, a framework that virtualizes multiple JVM instances
26 in a single JVM, thus simulating a distributed environment where each host runs on a separate
27 machine, accessing dedicated network and CPU resources. The key contributions are the ability
28 to run off-the-shelf concurrent and distributed JVM bytecode programs while at the same time
29 scaling up to thousands of virtual nodes; and enabling global observation within standard software
30 testing frameworks. Our experiments with two distributed systems show the usefulness of MINHA
31 in disclosing errors, evaluating global properties, and in scaling tests orders of magnitude with the
32 same hardware resources.

33 **2012 ACM Subject Classification** Computing methodologies → Distributed computing methodolo-
34 gies; Software and its engineering → Software testing and debugging

35 **Keywords and phrases** Distributed software testing; Large scale distributed systems; Simulation

36 **Digital Object Identifier** 10.4230/LIPIcs.OPODIS.2019.32

37 **Funding** This work is financed by National Funds through the Portuguese funding agency, FCT –
38 Fundação para a Ciência e a Tecnologia within project: UID/EEA/50014/2019.

39 **1** Introduction

40 Formal validation and verification tools are increasingly practical and of paramount import-
41 ance in developing distributed algorithms. However, they work over models rather than actual
42 runnable code [21, 8] and with assumptions that do not hold in practice [11]. In addition to
43 outright bugs introduced by the translation from algorithm to runnable code, such as race
44 conditions, a major problem lies in aspects that are abstracted by models and are revealed



© Nuno Machado and Francisco Maia and Francisco Neves and Fábio Coelho and José Pereira;
licensed under Creative Commons License CC-BY

23rd International Conference on Principles of Distributed Systems (OPODIS 2019).

Editors: Pascal Felber, Roy Friedman, Seth Gilbert, and Avery Miller; Article No. 32; pp. 32:1–32:16

Leibniz International Proceedings in Informatics



LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

45 only in large scale tests [19]. A common problem is a situation in which the complexity of an
46 operation (e.g., searching) is not apparent in testing as some data structure (e.g., a buffer) is
47 mostly empty unless there is a large number of participants or with congestion, that never
48 happens with small scale testing.

49 For example, the correctness proof of the accrual failure detector employed by Cas-
50 sandra [15] assumed a negligible processing time of heartbeat messages. However, in practice,
51 Cassandra uses variable length messages during membership changes in the cluster that may
52 consume significant transmission and computation time. This mismatch between protocol
53 design and implementation caused constant *flapping* problems in 500+ node deployments,
54 preventing the cluster from stabilizing and scaling [3]. Flapping is a cluster instability problem
55 where the status of the nodes is continuously switching between up and down.

56 Unfortunately, testing and debugging such systems is extremely challenging, mainly due
57 to the following reasons:

58 *Non-determinism and huge state space.* Distributed executions often entail thousands of
59 non-deterministic, concurrent events (e.g. message arrivals, node crashes, and timeouts)
60 that, depending on the order in which they occur, cause the system to behave differently.
61 Although most event sequences are correct, some incorrect timings of events result in severe
62 damage, such as data corruption or system downtime [23]. In particular, previous work has
63 shown that real-world distributed applications are especially prone to bugs stemming from
64 *message races* [23, 26]. For instance, bug 2212 in ZooKeeper’s issue repository reports a
65 (rare) scenario in which a new node joining the cluster cannot become a leader due to a race
66 condition between a message from the atomic broadcast protocol and one from the leader
67 election protocol [4]. This bug causes a 3-node ZooKeeper cluster to stop working in the
68 presence of a single failure, when in fact the existence of a majority of two nodes alive is
69 sufficient for the system to operate.

70 *Lack of resources for testing at scale.* Distributed systems are typically developed to be
71 deployed on a massive number of independent nodes (e.g. Cassandra [1], Hadoop [2], etc).
72 Alas, as testing with close-to-production conditions can be prohibitively costly and time
73 consuming, these applications are often debugged on small/medium-size deployments that
74 prevent certain faulty behavior from manifesting [24]. Attempts to address scalability in
75 testing environments include cloud computing and virtualization technologies [14, 38], but
76 software validation is still the limiting factor in distributed software development [12].

77 *Difficulty in checking global properties at runtime.* In general, large-scale distributed protocols
78 are designed in such a way that nodes make decisions based only on local information (or
79 a partial view of the system). On the other hand, the correctness of these protocols often
80 implies certain global properties or invariants to hold, which can be hard to verify without a
81 globally-consistent snapshot of the system. For example, Pastry [33] is a distributed hash
82 table whose routing algorithm requires each node to maintain a list with its physically closest
83 peers. Thus, to assess the correct execution of Pastry, one needs to first obtain the complete
84 network overlay, then collect the neighbor lists from all nodes and, finally, compute the
85 distances in both cases to check whether the references in the lists actually correspond to
86 the closest peers.

87 **Contribution.** In this paper we propose MINHA,¹ a framework that combines virtualization
88 and simulation techniques for making large-scale distributed systems testing practical by

¹ MINHA is available as open source at <http://www.minha.pt>.

89 providing a unique trade-off between scale and observation completeness. In particular,
90 MINHA addresses the problems introduced in the translation from algorithm to runnable
91 code and message races, such as experienced in Cassandra [3] and Zookeeper [4], with two
92 contributions:

93 *Scaling-up centralized simulation.* The technique proposed in CESIUM [6] is scaled up to
94 thousands of distributed nodes by reducing the processing and memory requirements for
95 node isolation. Moreover, it is scaled to execute off-the-shelf distributed applications written
96 in modern Java while simulating key environment components, reproducing the concurrency,
97 distribution, and performance characteristics of real-world deployments.

98 *Providing meta-interfaces for observation and automation.* MINHA provides a programming
99 interface to orchestrate large-scale virtual deployments of complete programs or standalone
100 middleware layers with application stubs, aimed at being used within standardized testing
101 frameworks. The same interface eases the collection of consistent snapshots and traces from
102 distributed executions, suitable for visualization or evaluating global properties.

103 We evaluated MINHA on a peer sampling service protocol and a large-scale key-value
104 store. The results show that MINHA is not only able to assess properties over a coherent,
105 global snapshot of the system without imposing runtime overhead while at the same time
106 allowing tests with a large number of nodes to run in cost effective way.

107 The rest of this paper is structured as follows. Section 2 outlines the state-of-the-art
108 in simulation and emulation for testing distributed systems. Section 3 describes how the
109 design and implementation of MINHA provide a new trade-off between scale and observation
110 completeness and scale. Sections 4 and 5 evaluate MINHA. Section 6 concludes the paper.

111 2 Related Work

112 The ideal approach to test distributed systems software for race conditions is the use of
113 implementation-level model checkers. Model checkers, such as MaceMC [18], Demeter [13],
114 MoDist [39], dBug [35] and SAMC [22], intercept non-deterministic events of local and dis-
115 tributed programs (e.g. message arrivals, node crashes, and timeouts), and permute their
116 ordering in systematic runs. This approach is very effective in discovering concurrency bugs,
117 as it virtually explores the whole system state space. However, for large-scale applications,
118 distributed model checking becomes unpractical and starts suffering from scalability issues
119 due to state space explosion [22].

120 The next best approach is to run tests and then evaluate system-wide properties by
121 logging the local state of each node independently along the execution and, periodically,
122 send those logs to a centralized machine to be combined into a globally-consistent snapshot
123 of the system that allows checking the desired predicates. For deployed systems, D³S [27]
124 proposes a simple language for writing distributed predicates that are checked on-the-fly
125 at runtime. DCatch [26], in turn, aims to detect distributed concurrency bugs by resorting
126 to a *happens-before* (HB) model. This work encodes the event causality into HB rules and
127 builds a graph representing the timing relationships of several distributed concurrency and
128 communication mechanisms. This approach has some drawbacks though. First, the execution
129 details to be recorded at runtime must be defined *a priori*. Second, monitoring the nodes'
130 local state is intrusive and induces both performance and space overhead, and finding the
131 sweet spot between overhead and the necessary amount of information to be traced is far
132 from trivial [27]. Third, setting up and running large scale tests requires a corresponding
133 large scale distributed infrastructure to be available.

134 An important class of bugs introduced in the translation from algorithms to implement-
135 ations, for instance, caused by the complexity of library operations and data structures
136 used, can often be disclosed by running the system at large scale. Distributed systems
137 researchers have been building platforms to address this challenge. EmuLab [16] provided
138 a set of dedicated computer nodes and networking hardware that could be reconfigured to
139 mimic different large scale systems. PlanetLab [32] uses a decentralized approach, therefore
140 enabling a much larger and realistic platform to run off-of-the-shelf code. Unfortunately,
141 experiments in PlanetLab are cumbersome to configure, deploy, and run. Splay [25] aims
142 at easing the task of configuring the environment and running large-scale experiments, but
143 limits the development to a specific framework and the Lua language.

144 The use of virtual machines and containers in public clouds, together with state-of-the-art
145 orchestration software [12] makes it much easier for any developer to set up and run large tests.
146 However, virtual nodes compete over the same physical resources, hence impacting negatively
147 the performance of the system being evaluated and the accuracy of the measurements [34],
148 possibly hiding the problems. Moreover, even if possible, it is still expensive to conduct
149 extensive testing in public clouds.

150 Given the cost of running large scale tests and the difficulties in obtaining reliable,
151 reproducible results, a number of proposals have focused on exploiting simulation to test
152 actual implementations. Simulation is used extensively for distributed systems research,
153 allowing simplified models to be tested in a very large scale [31], but they don't capture
154 timeliness properties of implementation decisions and require code to be written in event-
155 driven style, hence, with inversion of control. As an example, Neko [36] offers the ability to
156 use simulation code as actual code, as long as its event-driven API is used in place of the
157 standard Java classes.

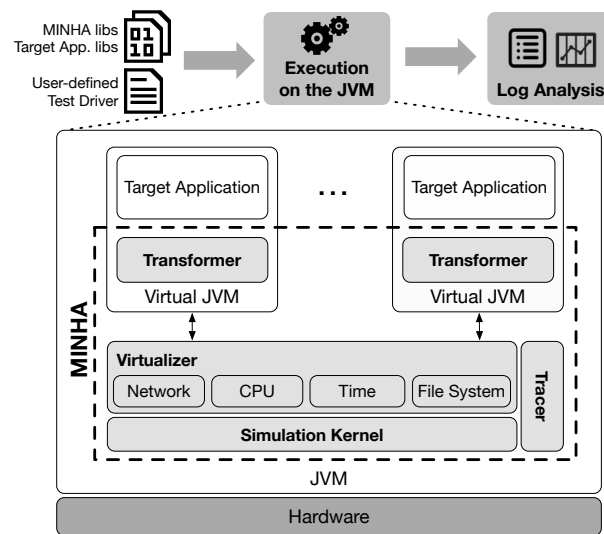
158 An interesting trade-off is achieved by JiST [7] (Java in Simulation Time), a simulation
159 kernel that allows event-driven simulation code to be written as Java threaded code, but
160 avoids the overhead of a native thread for each simulated thread by using continuations.
161 Unfortunately, this simulation kernel does not virtualize Java APIs and thus cannot be used
162 to run most of the existing Java code. Moreover, it does not reflect the actual overhead of
163 Java code in simulation time.

164 CESIUM [6] proposes the centralized simulation approach, in which the time effectively
165 used to execute implementation code for each event is measured and reflected into simulation
166 time. This is useful to detect issues such as the Cassandra failure detector bug [3], as it
167 captures the timeliness properties of implementation code. This does not however address
168 the general Java platform API and thus does not allow general code to be run. Moreover, by
169 using Java class loaders for virtualization it imposes a large memory overhead and restricts
170 simulations to a small number of nodes. UMLsim [5] is a similar proposal at the operating
171 system level, that virtualizes Linux while providing a simulated timeline and network. By
172 requiring a full Linux installation for each node, it restricts attainable scale even more than
173 CESIUM.

174 The ideal approach would thus have the ability to run unmodified implementation code,
175 such as possible by using cloud computing, with the frugality and scalability of simulation, the
176 ability to reproduce timeliness properties of CESIUM, and the ability to capture distributed
177 snapshots of distributed debuggers. This the challenge addressed in this paper with MINHA.

178 **3** Minha Framework

179 MINHA is a practical framework for testing large-scale distributed systems. MINHA virtualizes



■ **Figure 1** Execution flow and architecture of MINHA.

180 multiple JVM instances within a single JVM and simulates key environment components,
 181 thus allowing reproducing the concurrency, distribution, and performance characteristics of
 182 an actual distributed system at a much lower resource cost.

183 The usage of MINHA is shown at the top of Figure 1. From left to right, a *test driver*
 184 defines the execution scenario, such as the number of instances of the target application to
 185 be created and the global properties to be checked. The test driver is then executed, along
 186 with the target application and MINHA libraries, on an off-the-shelf JVM. Properties can be
 187 checked in runtime and logs stored for further off-line checking and visualization.

188 At runtime, MINHA acts as an interposition layer that takes control of the execution and
 189 steers it according to the testing scenario defined in the test driver. As show also in Figure 1,
 190 the main components of MINHA are: the *simulation kernel*, the *virtualizer*, the *transformer*,
 191 and the *tracer*. Briefly, the transformer converts application and middleware code into an
 192 event-driven simulation that interacts with models provided by the virtualizer. Both run
 193 on an event-driven simulation kernel. The tracer collects information for off-line use. The
 194 remainder of this section describes how the design and implementation of MINHA address
 195 the twin challenge of scale and observability of distributed systems software.

196 3.1 Achieving scale

197 Simulating multiple processes in the JVM requires isolating their code and preventing their
 198 executions from interfering with each other. Prior work typically addresses this issue by using
 199 a separate Java class loader for each process [6], that provides private copies of code and data
 200 for each virtual process. However, this approach severely limits attainable numerical scale,
 201 as each virtual process has to load, transform, compile, and store its own copy of each class.

202 A second aspect of scale is the size and complexity of the application and middleware
 203 that can be loaded. Current systems make use of large portions of the Java platform API in
 204 addition to the basic networking and time interfaces that have been intercepted in previous
 205 proposals [6]. Simply replacing all Java API with simulation models would lead to a very
 206 large development effort while impairing compatibility.

207 MINHA's transformer addresses these challenges by using single class loader for all virtual

208 processes and converts the original platform libraries to use simulation models of external
 209 resources, as happens for user provided middleware. As process isolation still needs to
 210 be enforced and the JVM forbids class loaders from rewriting native classes, MINHA uses
 211 the ASM Java bytecode manipulation and analysis framework [9] to perform the following
 212 bytecode modifications:

213 *Redirecting references to virtualized classes.* Direct references to classes that are replaced
 214 by simulation models (e.g., `java.net.Socket`) are rewritten. Simulation models need then
 215 to be written, with same same interfaces, and containing simulation logic to reproduce
 216 their behaviors. This can however be done incrementally: As new applications demand
 217 new platform interfaces that haven't yet been modeled, they fail and report the problem.
 218 Experience shows that having implemented models for a moderate subset of the platform
 219 API supports many interesting distributed middleware components.

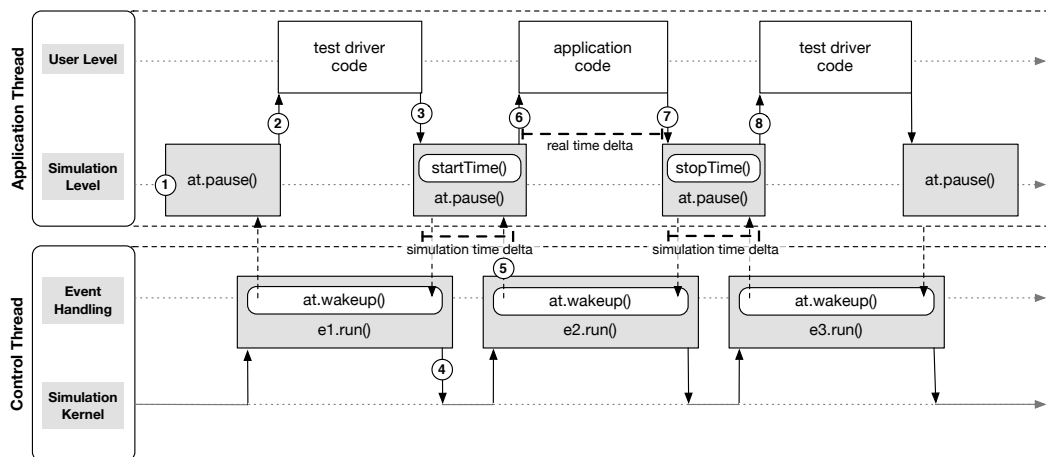
220 *Redirecting static instance variables.* This transformation moves static instance variables to
 221 regular instance variables in an auxiliary class. It then creates and uses static setter and
 222 getter methods for each of them. These methods use a map in the simulation model of each
 223 process to store and retrieve the correct instance, thus enforcing isolation.

224 *Redirecting references to renamed classes and methods.* Transformed platform classes are
 225 renamed to a separate package, thus circumventing the restriction to modifying them.
 226 Therefore, references to these classes and to classes that are replaced with simulated ver-
 227 sions are re-directed to the new package by transforming their respective callers. Indi-
 228 vidual static methods in platform classes that do not need to be replaced as a whole (e.g.,
 229 `System.currentTimeMillis()`) are simply redirected to simulated versions.

230 Moreover, a subset of classes, containing the simulation kernel and environment models,
 231 are kept *global* by delegating their load to the system's class loader. In contrast to isolation
 232 provided by using multiple class loaders and the Java security manager, the isolation between
 233 virtualized nodes in MINHA is not designed for containing any malicious code or attack. This
 234 is however not needed, as MINHA is used only for testing and all nodes are inherently trusted.
 235 This has the advantage of providing a controlled channel for virtual JVMs to interact with
 236 each other, that can easily be exploited by the user-defined test driver.

237 Finally, the last challenge to scale is in the discrete event simulation kernel. It keeps
 238 a list of future events per simulated timeline, scheduled to execute at target simulation
 239 instants. To scale up to large simulations, MINHA's simulation kernel supports multiple
 240 timelines for parallel execution in multi-core systems. Timelines in a single simulation are
 241 conservatively synchronized with a time window [28]: current time for events executing in
 242 parallel in different timelines differs by at most a constant W . This exploits the fact that
 243 MINHA maps one or more simulated hosts to each timeline and that the network latency is at
 244 least W . Therefore, message delivery events scheduled on a different timeline are guaranteed
 245 to be properly ordered.

246 Parallel simulation in MINHA is optimized to use one timeline for each available processor
 247 core. This is achieved by using a concurrent non-blocking data structure to store the event
 248 list and a non-blocking algorithm to keep track of the synchronization window. In detail,
 249 each thread tries to update the global lowest time of an event executing across all timelines,
 250 spinning until its next event handler can be safely executed. This approach works well in cases
 251 where events are evenly spread across all timelines, as typically happens when simulating
 252 large distributed systems in a few processor cores.



■ **Figure 2** Example of how simulation time and execution time are coordinated in MINHA.

253 3.2 Achieving observability

254 MINHA’s programming interface is targeted at automated tests and scripting and combines
 255 reflection, to represent entities in the simulation domain, with an interface to inject events
 256 and interact with those entities.

257 To set up and control the simulation, the test driver API provides the following entities
 258 that describe and manipulate the simulated system: `World` represents the system as a whole,
 259 keeping track of global simulated time and allowing hosts to be created and enumerated.
 260 A `Host` describes a simulated host with processing and storage resources and attached to
 261 the network with an address. It allows creating and enumerating simulated processes. A
 262 `Process` keeps a private address space with its own copy of `static` variables and allows
 263 invocation of methods. Both hosts and processes can be terminated, to simulate crash faults.

264 Execution in the context of a simulated process is started by creating an `Entry` proxy for
 265 some interface and scheduling an invocation. Each entry point corresponds to a thread in the
 266 simulated process. Using the interface, it can specify arbitrary delays or an absolute schedule
 267 and if it is executed synchronously, implicitly running the simulation until the invocation
 268 returns, or asynchronously, providing a future to wait for and retrieve the result whenever
 269 the simulation has advanced enough. The simulation can also be run for predetermined
 270 periods of time, easing periodic observation. The API also provides the ability to exit the
 271 simulation to execute code that performs global observation or logging. This is achieved with
 272 an `Exit` proxy that ensures that the simulation is stopped on invocation and restarted on
 273 return. An example using the API is shown in Figure 3 and further discussed in Section 4.

274 The main challenge addressed is ensuring that only consistent global states can be
 275 inspected, regardless of the simulation containing multiple threads in a number of virtual
 276 processes. First, all blocking operations, such as synchronization and calls to the platform,
 277 are replaced with calls to simulated synchronization primitives by modifying the compiled
 278 bytecode at load time. Second, the transformed code executes consistently with simulation
 279 time. This ensures that: *i*) the application thread advances only in the context of an
 280 simulation event; *ii*) the execution time observed is reflected in the usage of a simulated CPU
 281 core; and *iii*) the waiting time (e.g., when reading from disk) is computed by the simulation
 282 and not by actual contention.

283 Figure 2 shows in detail how this is achieved. From the bottom to the top, the Control

284 Thread (CT) originates in the Simulation Kernel and executes discrete Event Handling
285 procedures. In contrast to common discrete event simulation practice, these do not directly
286 modify model variables. Instead, they allow an Application Thread (AT) that executes test
287 driver and application code at User Level to advance. An AT thus progresses as follows:
288 (1) When started, the AT stopped using the `pause()` method to wait for its turn; when the
289 corresponding event is scheduled by the Simulation Kernel, the CT signals the AT, that can
290 then execute test driver code (2). When the test driver is done, (3) it starts accounting real
291 time with `startTime()`. This pauses the AT and returns control to the Simulation Kernel,
292 to wait for its turn (4). When simulation time has advanced, the AT is finally signaled to
293 start executing application code (6). The time elapsed while executing application code is
294 measured with the CPU cycle counter in `stopTime()` (7) and used to advance simulation
295 time accordingly. This means that further events, either in application or test driver code (8)
296 will be scheduled appropriately in simulation time.

297 This allows the passing of real time while executing code to influence simulation time,
298 thus reproducing the performance characteristics as needed to disclose scaling bugs. However,
299 since it is an event driven simulation, global state is consistent and can be observed by direct
300 inspection while the simulation is stopped. Moreover, as logging is done in user-defined
301 test driver code, outside the periods that measure real time, it has no impact in measured
302 performance and does not disturb the system under test.

303 Note that the test driver thread could get blocked in synchronization primitives within
304 the target application and middleware, leading to a deadlock. However, as synchronization
305 primitives have been replaced by simulated counterparts, these will recognize that the
306 invoking thread is the test driver and avoid blocking it. The developer writing the test driver
307 has however to be careful to make sure that the code used for observation is not susceptible
308 to inconsistent internal state.

309 In addition to directly checking properties, the tracer component allows logging events of
310 interest at runtime for off-line processing. Currently, MINHA is configured to trace events
311 regarding: thread synchronization (i.e. *fork/join/start/end* events), inter-node communic-
312 ation (i.e. socket *send/receive* events), and *read/write* accesses to variables indicated by
313 the programmer, if any. Recall that capturing memory accesses requires instructing the
314 transformer to dynamically instrument the target application's bytecode with calls to the
315 tracer, which may incur additional overhead during the execution. In contrast, all the
316 remaining events are captured by MINHA at the simulator side (i.e., outside the application),
317 hence they do not impose any slowdown, as opposed to what happens in a real deployment.

318 An additional feature of the traces produced by MINHA is that events are logged in a
319 coherent global order, due to the framework's centralized and virtualized nature. This is
320 particularly useful for debugging. In fact, MINHA comes with a built-in diagram generator
321 that provides a visual representation of the execution according to the information stored in
322 the trace file. Section 4.2 illustrates the benefits of this feature using the Cyclon example of
323 Section 4.

324 The event trace is produced in JSON format, which can later be consumed by external
325 tools. Events are stored as JSON objects, containing fields regarding: the thread identifier,
326 the type of event (as indicated in tracer's description in Section 3), the timestamp, and, for
327 inter-node communication events, a message unique identifier, and the source and destination
328 node addresses. MINHA's built-in visualizer consists of a Javascript module that uses the
329 `SVG.js` library to generate a graphical representation of the event trace as a space-time
330 diagram [20]. An example is shown in Figure 4 and discussed in the next section.

Algorithm 1 – Cyclon Active Thread at Node p

```

Init:
 $V_{size} \leftarrow$  size of  $p$ 's partial view
 $view \leftarrow$   $p$ 's view containing  $V_{size}$  references to other peers
 $S_{size} \leftarrow$  number of peers exchanged during a shuffling
operation ( $S_{size} \leq V_{size}$ )

for every  $T$  time units do
  // increase by one the age of all neighbors
  AGEGLOBAL( $view$ )
  // pick node  $q$  with highest age among all neighbors
   $q =$  GETOLDESTPEER( $view$ )
  // select a random subset of  $S_{size}$  neighbors
   $peers =$  SELECTPEERS( $view, S_{size}$ )
  // send list of peers to  $q$ , indicating that this message is a request for a shuffle
  SEND( $q, \{REQ, p, peers\}$ )
end for

```

4 Use Case: Peer Sampling Service

In this section, we show how MINHA can be used to test the properties of the overlays generated by the peer sampling service (PSS) Cyclon [37, 29]. A PSS is a mechanism, widely used by gossip-based peer-to-peer systems, that provides each node of the system with a partial *view* of the network. The overlay network used by the gossip protocol to spread information is thus defined by the logical graph that emerges from the union of all nodes' views at a given instant. Since the effectiveness and efficiency of the dissemination depends heavily on the properties of the overlay, the PSS refreshes each node's view from time to time to account for peer joins and leaves.

4.1 Test application

For the purpose of this example, we will focus on a particular PSS named Cyclon [37]. In a nutshell, the Cyclon protocol consists in a series of *shuffle* rounds, where pairs of neighbor nodes exchange a subset of randomly sampled peers from their views. The shuffle procedure works as follows. Each node assigns an *age* value to the node references in its view. This value is incremented by one at the beginning of a new exchange round, every T units of time. Upon starting a new shuffle, nodes pick the neighbor with highest age and send to it a random subset of other peers in their view, replacing the oldest node's reference with a self-reference of age 0. In turn, when a shuffle message is received, nodes first reply with a subset of peers randomly selected from their own view, and then update their views with the references received, replacing the entries previously sent. As a result of shuffle operations, nodes alive have their views refreshed and nodes that left the system are eventually removed from every view. This way, Cyclon is able to cope with dynamism.

The experimental results in Cyclon's original paper [37] shown that this protocol is able to generate network overlays with properties similar to random graphs, even starting from non-random topologies. This ability is particularly relevant for peer-to-peer systems that have to handle high churn, as random overlays exhibit low diameter and are able to maintain connectivity even in the presence of massive node failures.

The results presented in the original paper were obtained solely from simulations though, so it remains unclear how would the protocol behave in an actual distributed system. In such a real-world setting, the Cyclon protocol can be implemented by means of two execution threads: an *active thread* that is responsible for initiating a new shuffle with a neighbor, and a *passive thread* that operates as a message handler, which receives and processes the

Algorithm 2 – Cyclon Passive Thread at Node p

```

while true do
  // receive a list of peers sent by node q
  {t, q, peersRcv} = RECEIVE()
  if t = REQ then
    // select a random subset of Ssize neighbors
    peersSnd = SELECTPEERS(view, Ssize)
    // send list of peers to q, indicating that this message is a reply to the shuffle
    SEND(q, {REP, p, peersSnd})
  end if
  // incorporate list of peers received into its own view
  UPDATEVIEW(view, peersRcv)
end while

```

```

1 //create a new simulation with 500 Cyclon peers
2 World world = new Simulation();
3 Entry<ICyclon>[] peers = world.createEntries(500, ICyclon.class, CyclonImpl.class.
   getName());
4
5 //start Cyclon peers
6 for (int i = 0; i < 500; i++){
7   peers[i].queue().run();
8 }
9
10 //allow nodes to bootstrap
11 world.run(1, TimeUnit.SECONDS);
12
13 //let the application code execute for 5s
14 //and then observe the overlay
15 for(int i=0; i <= 100; i++) {
16   world.run(5 , TimeUnit.SECONDS);
17   for(ICyclon c : peers){
18     //inspect node's view and store it into log
19     logView(c.getView());
20   }
21 }

```

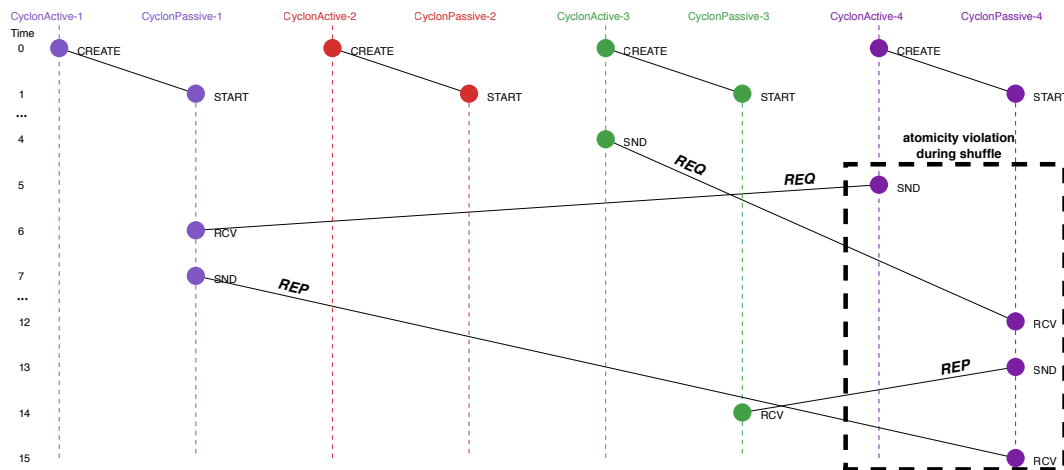
■ **Figure 3** Test driver for running Cyclon on MINHA

363 messages sent by other nodes. The portions of Cyclon executed by the two threads are
 364 detailed in Algorithms 1 and 2, respectively.

365 4.2 Test driver

366 To conduct the experiment in MINHA, one can use a test driver like the one presented in
 367 Figure 3. The test driver starts by creating a new simulation scenario with 500 hosts, each
 368 corresponding to a Cyclon node, using MINHA's API (lines 2-3). Nodes are then prepared
 369 for execution by enqueueing the method `run()` in the simulator (lines 6-8). This method is
 370 responsible for spawning Cyclon's passive and active threads, as well as initiating the node's
 371 view with references to its 11 subsequent neighbors.

372 The test driver proceeds with an instruction to let the simulation run for one second,
 373 thus allowing the Cyclon nodes to bootstrap (line 11). The core of the test run consists of
 374 100 simulation cycles, in which MINHA lets the application code of each instance execute for
 375 5 seconds, before giving back the control to the test driver code (lines 15-21). At this point,
 376 the test driver consults the local state of each node, namely the composition of its view, and
 377 logs it into a trace file (lines 17-20). We highlight that these observations are performed



■ **Figure 4** Space-time diagram generated by MINHA’s visualizer from the event trace captured during the simulation of Cyclon. The values on the left represent logical time ticks, whereas events *CREATE*, *START*, *SND*, *RCV*, denote, respectively, the creation and beginning of a thread, and the send and receive of a message. Messages labeled as *REQ* indicate shuffle requests and as *REP* indicate shuffle replies. The dashed box highlights the fact that atomicity is not guaranteed during a view exchange procedure, as new shuffle requests may arrive in the meantime.

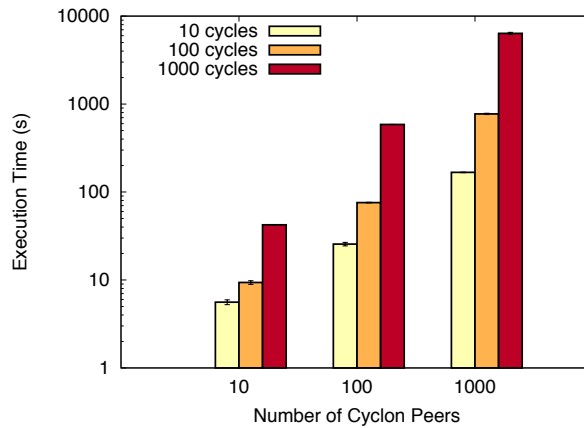
378 transparently to the application and without incurring runtime overhead. Since all nodes
 379 are *paused*, the local views are obtained from a global snapshot of the system, thus allowing
 380 reenacting the actual overlay network at that exact moment.

381 4.3 Log analysis

382 Off-line log analysis is particularly well suited to discover message race conditions. As an
 383 example, we search logs obtained with MINHA for cases in which a node is involved in more
 384 than one concurrent shuffle operations. This fact is not contemplated in the original Cyclon
 385 algorithm, nor in the simulations presented in the original paper [37].

386 On the other hand, these atomicity violations can easily occur in real settings, and may
 387 hamper the properties of the overlay generated, especially in highly dynamic environments.
 388 Figure 4 illustrates a detail of a time-space diagram of the execution plotted from the event
 389 trace captured during the simulation for our experiment with Cyclon. To improve readability,
 390 we only depict the timelines of the two Cyclon threads (namely, the active thread and the
 391 passive thread) for the first four nodes. From the diagram, it can be observed that the active
 392 thread starts by spawning the passive thread and then proceeds to sending shuffle requests.
 393 In turn, the passive thread is responsible for receiving incoming messages and reply back to
 394 complete the view exchange, as defined by the protocol.

395 This time-space diagram confirms the interesting scenario: The atomicity of a shuffle
 396 operation is not guaranteed in practice (see the dashed box in Figure 4). In fact, a node can
 397 receive a new shuffle request from a third node in-between swapping views with a neighbor.
 398 Since the view updates performed by the passive thread upon receiving shuffle requests and
 399 replies are not commutative, it results in a view that is not anticipated in the algorithm.
 400 A possible solution to address this issue is to store incoming requests in a queue until the
 401 awaited reply arrives [17].



■ **Figure 5** MINHA’s execution time (in log scale) for simulations considering different configurations of Cyclon.

402 4.4 Scalability

403 This section evaluates how MINHA’s performance scales with the number of nodes simulated.
 404 In particular, we measured the execution time of MINHA when varying both the number of
 405 nodes and the number of cycles considered in the simulation within the range {10, 100, 1000}.
 406 The experiments were performed on a machine with a 3.4 GHz Dual-Core Intel i3 Processor,
 407 8 GB of RAM and a 7200 RPMs SATA disk. The results, averaged for three runs, are
 408 depicted in Figure 5.

409 As expected, the figure shows that MINHA’s execution time grows proportionally to
 410 the size of the simulation both in terms of number of Cyclon peers and number of cycles,
 411 making it efficient and practical for in-house testing of large-scale distributed systems as a
 412 simulation of the Cyclon protocol with 1000 nodes and 100 cycles (which suffice to assess
 413 most properties of the overlay, as shown in Section 4.2) takes only around 13 minutes with a
 414 minimal hardware configuration.

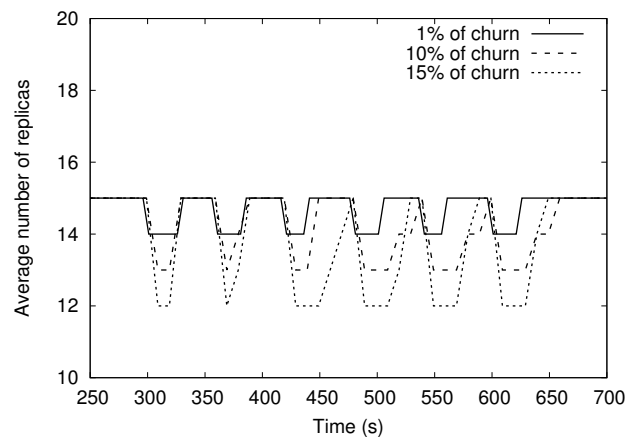
415 5 Use Case: Data Store

416 In this section, we provide a second example of MINHA using the DATAFLASKS peer-to-
 417 peer database system.² This is used to show that MINHA copes with an unmodified larger
 418 application that makes use of more features of the Java platform and external libraries.
 419 Moreover, we compare the scalability and cost of tests with the common alternatives of
 420 setting up a real distributed system or using virtual machines.

421 5.1 Global Property Checking

422 The DATAFLASKS distributed data store was designed to ensure availability of data in the
 423 presence of varying levels of node churn [30]. To assess this property, we ran DATAFLASKS
 424 on MINHA with different levels of churn and wrote a test driver to compute the number
 425 of data replicas stored in the system at a given instant Churn is easily injected in using

² We used the version of DATAFLASKS publicly available at github.com/fmaia/dataflasks



■ **Figure 6** MINHA’s dynamic property checking applied to the evaluation of DATAFLASKS’ replica maintenance properties.

426 the test driver API by invoking `close()` on `Host` objects, to remove them from the system,
 427 or by creating new `Host` instances to bring them back. Computing the number of keys in
 428 the system requires inspecting the storage of each node to check the data it is holding and
 429 combining such information with that of the other nodes.

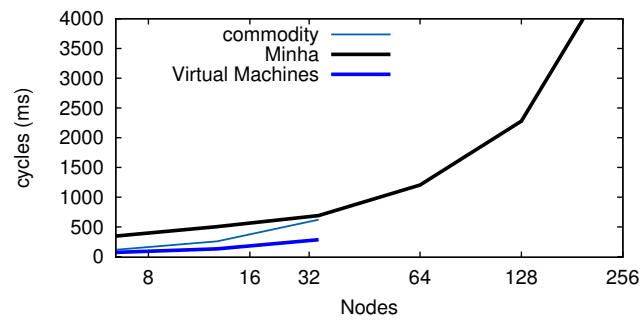
430 In the experiments, we instructed MINHA to run the verification code and compute the
 431 average number of replicas that the system holds at execution intervals of 1 second. We
 432 plot the results of the experiments in Figure 6 for a simulation of 700 seconds. The results
 433 show that, as expected, the average number of data replicas varies with churn (the higher
 434 the churn, the fewer replicas exist in the system), although DATAFLASKS is always able
 435 to eventually recover the keys lost. In fact, the system ends the experiment maintaining
 436 the expected mean number of data replicas. Checking the same property in a real-world
 437 deployment would require extensive logging and cumbersome synchronization mechanisms.

438 5.2 Performance and Resource Usage

439 For this experiment, we considered three different three different deployment configurations
 440 for DATAFLASKS:

441 *Configuration 1 (Commodity).* We considered a deployment built from a set of several
 442 commodity hosts. Each commodity host is randomly selected at startup time from a pool of
 443 resources equipped with either *i*) a 3.1 GHz Dual-Core Intel i3 Processor, 8 GB of RAM and
 444 a 7200 RPMs SATA disk, *ii*) a 3.4 GHz Dual-Core Intel i3 Processor, 8 GB of RAM and a
 445 7200 RPMs SATA disk, or *iii*) a 3.7 GHz Dual-Core Intel i3 Processor, 8 GB of RAM and
 446 a SSD disk. All hosts are interconnected through a switched Gigabit Ethernet. The total
 447 number of commodity hosts available at deployment time was limited to 36, which aims to
 448 represent a scenario where the available resources for testing are not comparable with the
 449 equivalent to those expected in a production deployment.

450 *Configuration 2 (Virtual Machines).* We considered a deployment built from a single server
 451 grade machine equipped with an 2.3 GHz Intel Xeon E5-2670 v3 Processor, 94GB of RAM
 452 and a 7200 RPMs SATA disk. This machine was configured to deploy a set of 32 virtual
 453 machines interconnected through a virtualized switched Ethernet.



■ **Figure 7** Time to execute the workload in the DATAFLASKS key-value store for the three deployment configurations.

454 *Configuration 3 (Minha)*. We considered a single commodity host, extracted from the pool
 455 of resources described in the first scenario, on which we deployed MINHA.

456 For each one of the previous configurations, we deployed DATAFLASKS with an increasing
 457 number of nodes within the range $\{8, 16, 32, 64, 128, 256\}$. We then performed the experiments
 458 by running a simple write workload on top of DATAFLASKS using YCSB [10], in a total
 459 of 5 independent runs for each configuration. For each experiment, we measured the time
 460 required to replicate data across all the necessary replicas and used it to compare the different
 461 configurations. Figure 7 depicts the experimental results.

462 Figure 7 shows that besides scaling to larger system sizes than a configuration with
 463 multiple virtual machines, MINHA is able to do it on a single commodity host that represents
 464 $1/36$ of the cost of a real deployment depicted in either the first or second configurations.
 465 These results support the claim that MINHA is able to accurately simulate large-scale
 466 distributed systems with much less resources than traditional approaches.

467 6 Conclusions and Future Work

468 Distributed systems are notoriously hard to test, mainly due to their large state space and
 469 the difficulty in deploying large infrastructures. This paper addresses the challenges of scale
 470 and observability in MINHA, a simulation framework aimed at easing the burden of testing
 471 large-scale distributed systems. MINHA virtualizes multiple JVM instances within a single
 472 JVM, while simulating key environment components, thus reproducing the concurrency,
 473 distribution, and performance characteristics of an actual system. MINHA also helps assessing
 474 the application’s correctness by allowing checking distributed properties on globally-consistent
 475 snapshots of the system. Moreover, due to time virtualization, these system-wide assertions
 476 can be performed transparently to the target application and without affecting its runtime
 477 performance.

478 Our experiments with a large-scale key-value store and a peer sampling service show that
 479 MINHA can accurately reproduce the characteristics of real-world deployments with fewer
 480 resources than traditional approaches, and is effective in assessing system-wide properties.

481 We believe that MINHA opens a number of interesting research opportunities in the field of
 482 distributed systems testing and debugging. For example, MINHA can be extended with model
 483 checking capabilities to systematically explore the execution space of distributed systems
 484 and discover new bugs. Furthermore, MINHA’s event traces, which are totally ordered and
 485 logged without any runtime overhead for the target application, can be combined with bug

486 detection techniques to also detect problems automatically.

487 ——— **References** ———

- 488 1 Apache Cassandra. <http://cassandra.apache.org>.
- 489 2 Apache Hadoop. <http://hadoop.apache.org>.
- 490 3 Bug CASSANDRA-6127: vnodes don't scale to hundreds of nodes. <https://issues.apache.org/jira/browse/CASSANDRA-6127>.
- 491 4 Bug ZOOKEEPER-2212: distributed race condition related to qv version. <https://issues.apache.org/jira/browse/ZOOKEEPER-2212>.
- 492 5 Werner Almesberger. umlsim-a uml-based simulator. In *10th International Linux System*
- 493 *Technology Conference (Linux-Kongress 2003)*, pages 202–213, 2003.
- 494 6 G. A. Alvarez and F. Cristian. Applying simulation to the design and performance evaluation
- 495 of fault-tolerant systems. In *SRDS '97*, pages 35–42, Oct 1997.
- 496 7 Rimon Barr, Zygmunt J Haas, and Robbert van Renesse. Jist: An efficient approach to
- 497 simulation using virtual machines. *Software: Practice and Experience*, 35(6):539–576, 2005.
- 498 8 Y Bertot and P Castéran. Interactive theorem proving and program development—coq'art:
- 499 The calculus of inductive constructions (2004).
- 500 9 Eric Bruneton, Romain Lenglet, and Thierry Coupaye. Asm: A code manipulation tool to
- 501 implement adaptable systems. In *In Adaptable and extensible component systems*, 2002.
- 502 10 Brian F. Cooper, Adam Silberstein, Erwin Tam, Raghu Ramakrishnan, and Russell Sears.
- 503 Benchmarking cloud serving systems with ycsb. In *SoCC '10*, pages 143–154, New York, NY,
- 504 USA, 2010. ACM.
- 505 11 Pedro Fonseca, Kaiyuan Zhang, Xi Wang, and Arvind Krishnamurthy. An empirical study on
- 506 the correctness of formally verified distributed systems. In *EuroSys '17*, New York, NY, USA,
- 507 2017. ACM.
- 508 12 F. Gortázar, M. Gallego, M. Donato, E. Pages, A. Edmonds, G. Tuñón, A. Bertolino, G. De An-
- 509 gelis, A. Cervantes, T. Bohnert, A. Willner, and V. Gowtham. The ElasTest platform: Support-
- 510 ing automation of end-to-end testing of large complex applications. ElasTest project whitepaper,
- 511 November 2018. URL: https://elastest.io/resources/ElasTest_white_paper.pdf.
- 512 13 Huayang Guo, Ming Wu, Lidong Zhou, Gang Hu, Junfeng Yang, and Lintao Zhang. Practical
- 513 software model checking via dynamic interface reduction. In *SOSP '11*, pages 265–278. ACM,
- 514 2011.
- 515 14 Diwaker Gupta, Kashi Venkatesh Vishwanath, Marvin McNett, Amin Vahdat, Ken Yocum,
- 516 Alex Snoeren, and Geoffrey M. Voelker. Diecast: Testing distributed systems with an accurate
- 517 scale model. *ACM Trans. Comput. Syst.*, 29(2):4:1–4:48, May 2011.
- 518 15 Naohiro Hayashibara, Xavier Défago, Rami Yared, and Takuya Katayama. The phi accrual
- 519 failure detector. In *SRDS*. IEEE Computer Society, 2004.
- 520 16 Mike Hibler, Robert Ricci, Leigh Stoller, Jonathon Duerig, Shashi Guruprasad, Tim Stack,
- 521 Kirk Webb, and Jay Lepreau. Large-scale virtualization in the emulab network testbed.
- 522 17 Márk Jelasity, Spyros Voulgaris, Rachid Guerraoui, Anne-Marie Kermarrec, and Maarten
- 523 Van Steen. Gossip-based peer sampling. *ACM Transactions on Computer Systems (TOCS)*,
- 524 25(3):8, 2007.
- 525 18 Charles Killian, James W Anderson, Ranjit Jhala, and Amin Vahdat. Life, death, and the
- 526 critical transition: Finding liveness bugs in systems code. In *NSDI '07*, 2007.
- 527 19 Ignacio Laguna, Dong H. Ahn, Bronis R. de Supinski, Todd Gamblin, Gregory L. Lee, Martin
- 528 Schulz, Saurabh Bagchi, Milind Kulkarni, Bowen Zhou, Zhezhe Chen, and Feng Qin. Debugging
- 529 high-performance computing applications at massive scales. *Commun. ACM*, 58(9):72–81,
- 530 August 2015.
- 531 20 Leslie Lamport. Time clocks, and the ordering of events in a distributed system. *Commun.*
- 532 *ACM*, 21:558–565, July 1978.
- 533 21 Leslie Lamport. *Specifying systems: the TLA+ language and tools for hardware and software*
- 534 *engineers*. Addison-Wesley Longman Publishing Co., Inc., 2002.
- 535
- 536

- 537 22 Tanakorn Leesatapornwongsa, Mingzhe Hao, Pallavi Joshi, Jeffrey F. Lukman, and Haryadi S.
538 Gunawi. SAMC: Semantic-aware model checking for fast discovery of deep bugs in cloud
539 systems. In *OSDI '14*. USENIX Association, 2014.
- 540 23 Tanakorn Leesatapornwongsa, Jeffrey F. Lukman, Shan Lu, and Haryadi S. Gunawi. TaxDC:
541 A taxonomy of non-deterministic concurrency bugs in datacenter distributed systems. In
542 *ASPLOS '16*. ACM, 2016.
- 543 24 Tanakorn Leesatapornwongsa, Cesar A. Stuardo, Riza O. Suminto, Huan Ke, Jeffrey F.
544 Lukman, and Haryadi S. Gunawi. Scalability bugs: When 100-node testing is not enough. In
545 *HotOS '17*, pages 24–29. ACM, 2017.
- 546 25 Lorenzo Leonini, Étienne Rivière, and Pascal Felber. Splay: Distributed systems evaluation
547 made simple (or how to turn ideas into live systems in a breeze). In *NSDI*, volume 9, pages
548 185–198, 2009.
- 549 26 Haopeng Liu, Guangpu Li, Jeffrey F. Lukman, Jiabin Li, Shan Lu, Haryadi S. Gunawi, and
550 Chen Tian. DCatch: Automatically detecting distributed concurrency bugs in cloud systems.
551 In *ASPLOS '17*. ACM, 2017.
- 552 27 Xuezheng Liu, Zhenyu Guo, Xi Wang, Feibo Chen, Xiaochen Lian, Jian Tang, Ming Wu,
553 M. Frans Kaashoek, and Zheng Zhang. D3s: Debugging deployed distributed systems. In
554 *NSDI '08*. USENIX Association, 2008.
- 555 28 B. D. Lubachevsky. Efficient distributed event-driven simulations of multiple-loop networks.
556 *Commun. ACM*, 32(1):111–123, January 1989.
- 557 29 Nuno Machado, Francisco Maia, Miguel Matos, and Rui Oliveira. BuzzPSS: A dependable
558 and adaptive peer sampling service. In *LADC '16*. IEEE Computer Society, 2016.
- 559 30 Francisco Maia, Miguel Matos, Ricardo Vilaça, José Pereira, Rui Oliveira, and Etienne Rivière.
560 Dataflasks: Epidemic store for massive scale systems. In *SRDS '14*. IEEE Computer Society,
561 2014.
- 562 31 A. Montresor and M. Jelasity. Peersim: A scalable p2p simulator. In *International Conference*
563 *on Peer-to-Peer Computing*, 2009.
- 564 32 Larry Peterson and Timothy Roscoe. The design principles of planetlab. *ACM SIGOPS*
565 *operating systems review*, 40(1):11–16, 2006.
- 566 33 Antony I. T. Rowstron and Peter Druschel. Pastry: Scalable, decentralized object location,
567 and routing for large-scale peer-to-peer systems. In *Middleware '01*, London, UK, UK.
568 Springer-Verlag.
- 569 34 Jörg Schad, Jens Dittrich, and Jorge-Arnulfo Quiané-Ruiz. Runtime measurements in
570 the cloud: Observing, analyzing, and reducing variance. *Proc. VLDB Endow.*, 3(1-
571 2):460–471, September 2010. URL: <http://dx.doi.org/10.14778/1920841.1920902>, doi:
572 10.14778/1920841.1920902.
- 573 35 Jiri Simsa, Randy Bryant, and Garth A Gibson. dbug: Systematic evaluation of distributed
574 systems. In *SSV '10*, 2010.
- 575 36 Peter Urban, Xavier Défago, and André Schiper. Neko: A single environment to simulate
576 and prototype distributed algorithms. In *Information Networking, 2001. Proceedings. 15th*
577 *International Conference on*, pages 503–511. IEEE, 2001.
- 578 37 Spyros Voulgaris, Daniela Gavidia, and Maarten van Steen. CYCLON: Inexpensive membership
579 management for unstructured p2p overlays. *Journal of Network and Systems Management*,
580 13(2):197–217, 2005.
- 581 38 Yang Wang, Manos Kapritsos, Lara Schmidt, Lorenzo Alvisi, and Mike Dahlin. Exalt:
582 Empowering researchers to evaluate large-scale storage systems. In *NSDI'14*, pages 129–141,
583 Berkeley, CA, USA, 2014. USENIX Association.
- 584 39 Junfeng Yang, Tisheng Chen, Ming Wu, Zhilei Xu, Xuezheng Liu, Haoxiang Lin, Mao
585 Yang, Fan Long, Lintao Zhang, and Lidong Zhou. MODIST: Transparent model checking of
586 unmodified distributed systems. In *NSDI '09*. USENIX Association, 2009.