

Journal Pre-proof

Update with care: Testing candidate bug fixes and integrating selective updates through binary rewriting

Anthony Saieva, Gail Kaiser



PII: S0164-1212(22)00105-4
DOI: <https://doi.org/10.1016/j.jss.2022.111381>
Reference: JSS 111381

To appear in: *The Journal of Systems & Software*

Received date: 5 March 2021
Revised date: 16 March 2022
Accepted date: 20 May 2022

Please cite this article as: A. Saieva and G. Kaiser, Update with care: Testing candidate bug fixes and integrating selective updates through binary rewriting. *The Journal of Systems & Software* (2022), doi: <https://doi.org/10.1016/j.jss.2022.111381>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2022 Elsevier Inc. All rights reserved.

Title page with author details

Update with Care: Testing Candidate Patches and
Integrating Selective Updates

Anthony Saieva^a, Gail Kaiser^b

*Columbia University
Department of Computer Science
530 W 120th Street, New York, NY 10027*

^a*ant@cs.columbia.edu*

^b*kaiser@cs.columbia.edu (Corresponding author)*

1 Update with Care: Testing Candidate Bug Fixes and
2 Integrating Selective Updates through Binary Rewriting

3 Anthony Saieva^a, Gail Kaiser^b

4 *Columbia University*
5 *Department of Computer Science*
6 *New York, NY 10027*

7 ^a*ant@cs.columbia.edu*

8 ^b*kaiser@cs.columbia.edu*

9 **Abstract**

10 Enterprise software updates depend on the interaction between user and devel-
11 oper organizations. This interaction becomes especially complex when a single
12 developer organization writes software that services hundreds of different user
13 organizations. Miscommunication during patching and deployment efforts lead
14 to insecure or malfunctioning software installations. While developers oversee
15 the code, the update process starts and ends outside their control. Since devel-
16 oper test suites may fail to capture buggy behavior finding and fixing these bugs
17 starts with user generated bug reports and 3rd party disclosures. The process
18 ends when the fixed code is deployed in production. Any friction between user,
19 and developer results in a delay patching critical bugs.

20 Two common causes for friction are a failure to replicate user specific cir-
21 cumstances that cause buggy behavior and incompatible software releases that
22 break critical functionality. Existing test generation techniques are insufficient.
23 They fail to test candidate patches for post-deployment bugs and to test whether
24 the new release adversely effects customer workloads. With existing test gener-
25 ation and deployment techniques, users can't choose (nor validate) compatible
26 portions of new versions and retain their previous version's functionality.

27 We present two new technologies to alleviate this friction. First, Test Gen-
28 eration for Ad Hoc Circumstances transforms buggy executions into test cases.
29 Second, Binary Patch Decomposition allows users to select the compatible pieces

30 of update releases. By sharing specific context around buggy behavior and devel-
31 opers can create specific test cases that demonstrate if their fixes are appropri-
32 ate. When fixes are distributed by including extra context users can incorporate
33 only updates that guarantee compatibility between buggy and fixed versions.

34 We use change analysis in combination with binary rewriting to transform
35 the old executable and buggy execution into a test case including the developer's
36 prospective changes that let us generate and run targeted tests for the candidate
37 patch. We also provide analogous support to users, to selectively validate and
38 patch their production environments with only the desired bug-fixes from new
39 version releases.

This paper presents a new patching workflow that allows developers to val-
idate prospective patches and users to select which updates they would like to
apply, along with two new technologies that make it possible. We demonstrate
our technique constructs tests cases more effectively and more efficiently than
traditional test case generation on a collection of real world bugs compared
to traditional test generation techniques, and provides the ability for flexible
updates in real world scenarios.

40 *Keywords:* test generation, change analysis, binary analysis, binary rewriting

41 1. Introduction

42 Developer testing may not be representative of how software is used in the
43 field [1] and many test suites are insufficient [2]. User bug reports [3, 4] and vul-
44 nerability disclosures [5, 6] are populated primarily with bugs discovered in the
45 field when users or third-party security analysts use the software in ways that
46 the developers had not tested before deployment. User bug reports typically
47 include some evidence of the bug, such as memory dumps, stack traces, system
48 logs, error messages, screenshots, and so on, but are often insufficient for the
49 developers to reproduce the bug [7]. Even vulnerability disclosures are some-
50 times incomplete, making it difficult for developers to reproduce the reported
51 exploits [8]. Thus when a security vulnerability or other critical bug is not

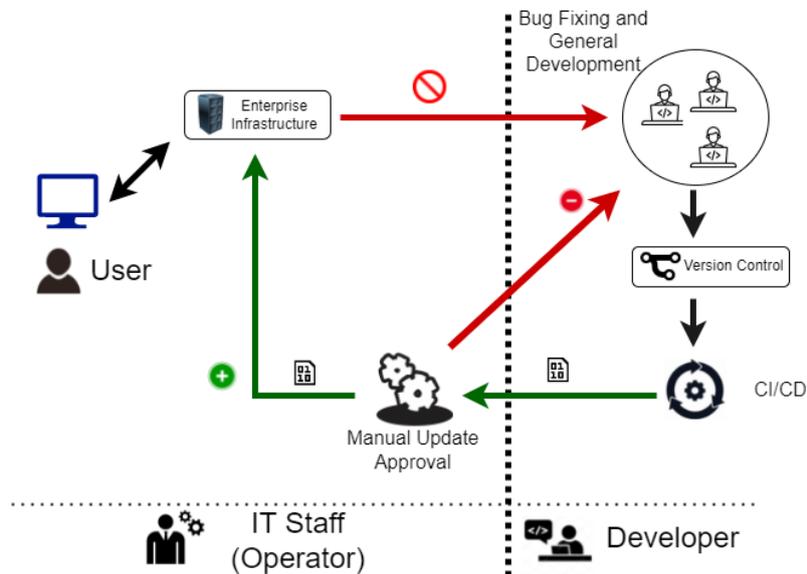


Figure 1: Ad hoc Test Generation Context

52 detected by developer testing prior to deployment, but reported by users, developer
 53 need to construct a new test that both reproduces the bug in the original
 54 version of the code and verifies absence of the bug in the patched code. Aside
 55 from patching, deployment presents another issue since software updates may
 56 not be compatible with existing infrastructure in some user environments. As
 57 a result, customer organizations may avoid updating their installations, leaving
 58 buggy code in production.

59 Enterprise software update procedures revolve around an interaction between
 60 user and development organizations shown in Figure 1. This interaction
 61 gets increasingly complex as a single software provider services hundreds of
 62 different user organizations. The developer organization writes code that gets
 63 checked into a version control system (VCS) and built with continuous deployment/
 64 continuous integration (CI/CD) which ships the resulting executables to
 65 user organizations that deploy the software. Software operators (IT staff) within

66 the user or customer organization approve the update and install the new exe-
67 cutable on machines.

68 Awkward interactions between the user organization and the development
69 organization often cause the traditional software update model to fail leading to
70 insecure or nonfunctional installations. When users report bugs, developers need
71 to reproduce the buggy behavior in the developer environment, update their
72 test suites, and develop a prospective patch. The current update paradigm may
73 fail to incorporate information from the specific user instance that caused the
74 buggy behavior relying solely on bug reports to assist developers. This means
75 a developer must manually build a representative test case that reproduces the
76 bug in the original code and verifies the absence of the bug in the patched code.

77 Every change has the potential to include unwanted side effects and while
78 CI and CD provide some protection it only considers the perspective of the
79 developer organization. In the event of a problematic update, the operators
80 responsible for deployment have no recourse other than to submit new bug
81 reports. This interaction gets further complicated by the fact that most bug
82 fixes are incorporated as part of more general releases which include changes
83 other than the bug fix. These additional changes may in fact break existing
84 functionality on any given installation even if they pass tests during CI/CD.

85 Interaction during reporting and distribution fail for the same reason: re-
86 stricted context. The user organization will not have access to the source code
87 producing the software and cannot make tailored modifications, and the devel-
88 oper organization must support many different user organizations without access
89 to any specific installation. This handshake between parties demands operators
90 and developers have intimate knowledge of what the other organization needs
91 while also inherently preventing them from sharing information.

92 We propose a new update paradigm that exposes precisely the relevant in-
93 formation needed by both organizations that supports the interaction between
94 these groups while still keeping their roles distinctly separate without imposing
95 prohibitive overhead. We give an overview in the next section.

96

111 production. This lightweight recording captures all sources of non-determinism
112 required to recreate the buggy execution. Then the log is augmented with addi-
113 tional information with an offline heavyweight recording process that includes
114 required additional information as outlined in 3.1. The verbose log shares ap-
115 propriate context between customer organization software operators (IT staff)
116 and developers maintaining the code. With this augmented log the developers
117 can re-create the bug from the production environment. Using our novel test
118 generation technique for ad hoc circumstances (Ad Hoc Test Generation), out-
119 lined in 3.2 and 3.3, developers turn this augmented log into a repeatable test
120 case for prospective patches. That test case becomes part of the test suite and
121 the developer approved patch gets added to the version control system (VCS)
122 along with any other changes as part of standard development. In the event mul-
123 tiple users report the same bug, there may be redundant test cases. The VCS
124 still integrates with continuous integration and continuous deployment systems
125 (CI/CD), but in order to expose additional context to the operator it interacts
126 with our BPD changelog datastore. Should an update received from developers
127 break critical functionality and fail manual approval, an operator has the op-
128 portunity to leverage the context in the BPD datastore to craft a custom partial
129 update for their customer organization.

130 By sharing specific context in the verbose trace and the BPD datastore
131 between parties the developers automatically have a test case to fix bugs and
132 operators deploying software have the flexibility to build updates that meet their
133 needs. Developer practices limit the level of granularity in the current prototype
134 of the BPD datstore so tangled commits (single commits with multiple unrelated
135 or weakly related changes [11]) create some confusion. With additional effort
136 from the developing party or minor changes to the BPD datastore prototype of
137 course these commits can be untangled.

138 The first of the two technologies that make this possible is test generation
139 for ad hoc circumstances which we call *ad hoc test generation* because the gen-
140 erated test emulates the *ad hoc* user context that manifested the bug. The key
141 observation that makes ad hoc test generation plausible is analogous to Tucek

142 et al.'s "delta execution" [12], whose large-scale study of patch size found that
143 security and other patches solely to fix bugs tend to be modest in size and
144 scope, rarely change core program semantics, shared memory layout or pro-
145 cess/thread layout. Nonetheless, bug reproduction is difficult. The premise of
146 record-replay technology highlights the difficulty in capturing all of the con-
147 ditions that led to erroneous behavior and recreating those conditions in the
148 developer environment. Standard bug reports often fail to capture all the rele-
149 vant state information, and this paper addresses the feasibility of using ad hoc
150 test [generation](#) in such scenarios. We have developed ATTUNE as a solution
151 that combines the buggy executable with the modified version to emulate what
152 would have happened had the modified version been deployed during the buggy
153 execution.

154 Instead of requiring developers to build doubles, mocks or other test scaf-
155 folding to fake the user environment for its tests, ATTUNE builds on existing
156 record-replay tools. It takes all non-deterministic inputs and replays them as
157 they happened when the bug manifested. This eliminates the common com-
158 munication failures when a developer tries to recreate the execution from user
159 reports. Furthermore, since all non-deterministic inputs are replayed Ad Hoc
160 Testing eliminates the possibility of confusing flaky tests [13].

161 There are two main challenges to technically implementing Ad Hoc testing.
162 1) How do you accurately identify changed portions of the executable once the
163 source code level abstractions have been stripped away? and 2) How do you
164 replay events after executions have diverged?

165 ATTUNE solves these challenges with two key insights: 1) The identifiers
166 used in the source code rarely change, and are still represented in the executable.
167 Software patches rarely modify function names and global variable names. In the
168 event they are changed, a mapping exists between the old and new identifiers so
169 these points still provide consistency between the old version and new versions.
170 These locations provide landmarks amidst the unstructured binary data to guide
171 ATTUNE's manipulation. 2) The recorded log does not need to be replayed
172 verbatim in order. Events in the log can be skipped or swapped, and new events

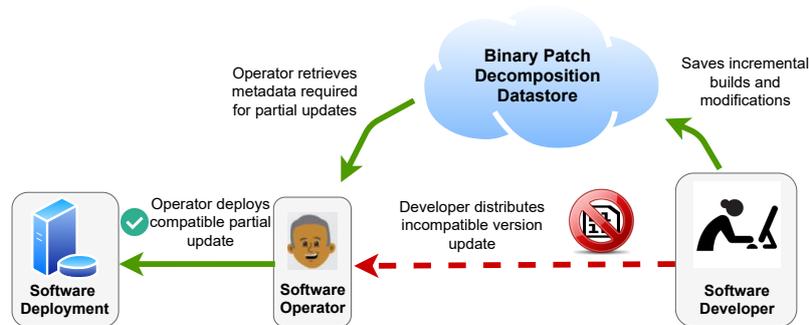


Figure 3: Ad hoc Test Generation and Partial Updates for Customer Organizations

173 can be derived from those in the log, to match the patch. Our runtime emulation
 174 algorithm selects which events to replay and when to support execution after
 175 divergence.

176 In addition to Ad Hoc Test Generation we offer binary patch decomposition
 177 to support deployment of software that would otherwise be impossible. Customer
 178 organizations tend to update software only when necessary for fear of
 179 updates introducing side effects that disrupt service. Software releases accord-
 180 ing to a common schedule, often contain many modifications most of which
 181 singular user would deem unnecessary. These pose an unnecessary risk and may
 182 not be able to integrate new versions if relevant interfaces have been replaced
 183 or wiped away. This leaves many users in an awkward position: they have code
 184 with known deficiencies and the corresponding updates, but they also can't
 185 apply those updates. As an example Equifax's well known breach in 2017 [14]
 186 exposed in 145.5 million people even though a bug fix was available

187 Figure 3 outlines our solution to this problem. The customer organization
 188 software operators (IT staff) monitor the deployed software and submits the
 189 bug report to the developer. While the developer works, the binary patch
 190 decomposition datastore (detailed in Figure 13) runs incremental builds and
 191 tracks changes at the binary level. When the developer distributes the new
 192 release, operators can instead apply partial updates from the BPD datastore

193 if the new release is incompatible with existing infrastructure. Ad hoc test
194 generation allows the operators to test the partial updates on recorded workloads
195 to verify the partially updated software functions correctly.

196 An earlier version of this work [15] introduced ad hoc test generation, and
197 briefly discussed our technique for adding developer environment metadata to
198 patch releases, enabling operators to validate patched versions with their own
199 workloads. In this paper, we build on our previous ad hoc test generation
200 workflow to enable a more complete solution. Furthermore, we added new
201 functionality to allow for partial updates, e.g. when a full update would break
202 mission critical functionality, based on ideas we previously sketched in [16].

203 The new contributions of this expanded paper are:

- 204 • A method for decomposing full version updates, with multiple bug-fixes
205 (and possibly new features), into its component pieces to enable partial
206 updates.
- 207 • A testing framework for determining if a partial update is compatible with
208 existing user environment infrastructure.
- 209 • A patching technique that allows users to apply partial updates despite
210 not having access to the source code.

211 We explain our requirements for verbose execution traces and the techni-
212 cal details of our binary rewriting techniques in Section 3. Our evaluation in
213 Section 4 describes how a developer would use ATTUNE to test candidate
214 patches for a variety of CVE security vulnerabilities and other bugs from well-
215 known open-source projects. Section 4 also gives an example where the user
216 records their own workload with the original program and replays with the
217 modified program to convince themselves that the bug has been fixed and the
218 patch does not break other behavior. We also test our partial patching in-
219 frastructure applied in the user environment to show partial updates can in
220 fact fix the bug. We analyze threats to validity in section 4.4 and compare
221 to related work in Section 5, and then summarize this work. Our open-source

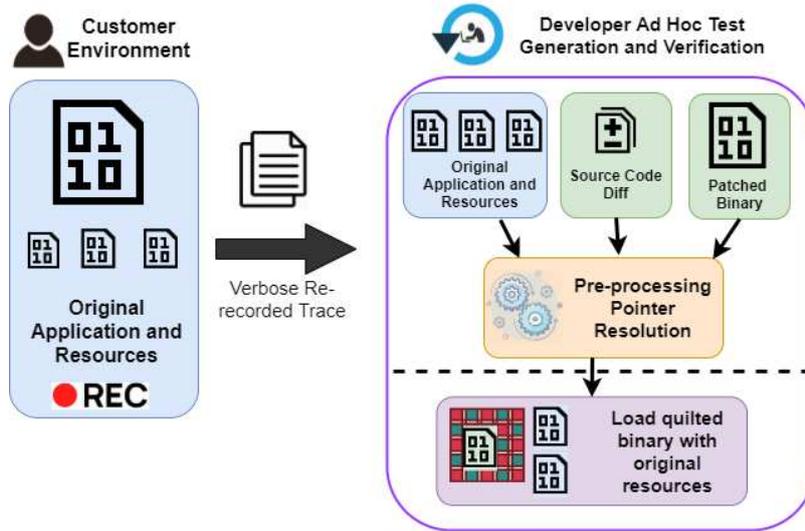


Figure 4: Recording and Preparation for Ad Hoc Test Generation

222 prototype implementation, portable across Linux distributions, is available at
 223 <https://github.com/Programming-Systems-Lab/ATTUNE>.

224 3. Implementation

225 Our ad hoc test generation workflow has four main components: *recording*,
 226 *static preprocessing*, *load-time quilting* and the *runtime replay decisions*. We
 227 detail support for Ad Hoc Test generation in the *customer* environment when
 228 source code is unavailable in section 3.5. Recording and the two preparation
 229 stages are shown in Figure 4, with runtime depicted in Figure 12. Both prepara-
 230 tion stages leverage the open-source Egalito recompilation framework [17].

231 3.1. Recording

232 We assume production recording with the user’s choice of lightweight tool
 233 and, when warranted by some external mechanism that detects an error or ex-
 234 ploit, offline replaying that tool’s recording while re-recording with rr’s recorder
 235 [18] as in Figure 2. Instead of rr, any other recording engine that constructs

236 sufficiently verbose traces would suffice, but we do not know of any actively-
237 supported open-source alternatives. Specifically, the trace must provide the
238 details needed for ATTUNE to recreate the successive register contents and
239 memory layouts leading up to when the bug manifested. Thus the recorded
240 sequence of events must include register values before and after system calls,
241 files that are `mmap`d into memory, and points at which thread interleaving and
242 signal delivery occur during execution.

243 ATTUNE imposes almost no restrictions on the user after constructing the
244 verbose trace. ATTUNE's technical design decisions enable ATTUNE to run
245 without privileges in user-space, with conventional hardware, operating system,
246 compiler, libraries, build processes etc. and no changes to the application or the
247 accompanying libraries. This represents a stark contrast to other test generation
248 techniques like symbolic or concolic execution. While the technical details of
249 our binary rewriting mechanisms are specific to our implementation ad hoc test
250 generation is not, and in principle ATTUNE prototypes could be built with any
251 record-replay technology that supports sufficiently detailed execution traces on
252 any architecture.

253 There are some changes that ATTUNE does not support. While our tech-
254 nique will capture concurrency related bugs, it cannot verify patches since it is
255 impossible to verify what nondeterministic thread interleavings might do after
256 a change. Any significant changes to data structures e.g., changing the size of a
257 struct on the stack or in the heap, that would require changes to memory allo-
258 cation would not be tolerated. Any changes to preprocessor macros don't have
259 symbols associated with them and so are not supported. Of course any major
260 feature addition that fundamentally changes software behavior is not supported.
261 We support all other changes that can be associated with symbols in the binary.

262 3.2. Static Preprocessing

263 **Source Code and Binary Preprocessing.** Figure 5 shows an abbreviated
264 example patch file from a libpng bug-fix [19]. Patch files document which files
265 changed, which function in the file changed, and which lines within that function

```

--- a/pngutil.c // file info
+++ b/pngutil.c
@@ -3167,10 +3167,13 @@ png_check_chunk_length(...)
{ ...
- (png_ptr->width * png_ptr->channels
...
+ (size_t)png_ptr->width
+ * (size_t)png_ptr->channels

```

Figure 5: libpng-1 Abbreviated Example Patch file

266 were inserted and deleted. Patch files are created with a standard format so we
 267 are not limited to a single diff implementation.

268 **Dwarf Information & Symbol Table.** Patch files don't provide any
 269 information about the resulting binary. Since the recorded trace relies on bi-
 270 nary/OS level information (register values, pointers, file descriptors, thread ids,
 271 etc.), we need to translate from changes in the source to changes in the binary.

```

182: 00000000000003fe0 56 FUNC
      GLOBAL DEFAULT 1 png_check_chunk_name
183: 00000000000004020 221 FUNC
      GLOBAL DEFAULT 1 png_check_chunk_length
184: 00000000000004100 172 FUNC
      GLOBAL DEFAULT 1 png_read_chunk_header

```

Figure 6: libpng-1 Symbol Table Entries

```

...
<c> DW_AT_producer: (indirect string, offset:
      0x1d90): GNU C11 7.4.0 ...
<10> DW_AT_language 12 (ANSI C99)
<11> DW_AT_name: (indirect string, offset: 0x1c8e):
      pngutil.c
...
0x0000402b [3156, 0] NS
0x0000403a [3166, 0] NS
0x00004046 [3182, 0] NS

```

Figure 7: libpng-1 Relevant DWARF Line Entries

272 Two mechanisms enable this translation: The first is the symbol table stan-

273 dard in all ELF files and the 2nd is DWARF information. The key insight is that
 274 the **symbols act as a point of reference between the old and the mod-**
 275 **ified binaries.** They remain unchanged even if their addresses and references
 276 change. After processing the [patch file](#) we use the symbol tables to find the
 277 locations of functions and global variables, and we use DWARF information for
 278 finding changed lines and identifying source files. These two sources combined
 279 contain all the information in the source level diff at the binary level. Refer to
 280 Figures 6 and 7 for concrete examples.

281 Most real-world builds create multiple binaries and associated libraries, so
 282 it may be unclear which binary contains the associated change. In order to
 283 generalize to sophisticated build processes ATTUNE uses DWARF information
 284 to search through all re-compiled binaries to find the modified file.

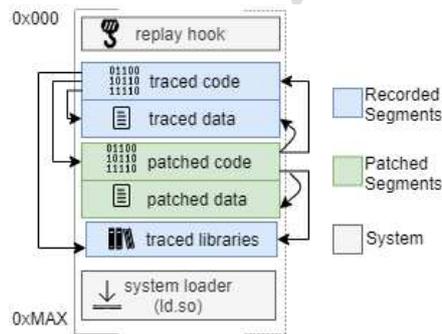


Figure 8: Address Space Detail

285 3.3. Load Time Quilting

286 **Pre-Load Steps for Quilting.** Once the function and line addresses have
 287 been resolved, and a prospective patched binary has been compiled, we can gener-
 288 ate our test code. In order for the newly compiled patched code to remain a
 289 viable test case, it must maintain the binary context of the original code. While
 290 most of the binary context remains unchanged, code pointers and data pointers
 291 that point somewhere inside the modified functions or that point from the mod-
 292 ified functions to any location outside of the modified binary must be updated

293 accordingly. To create the most accurate test we point to the original binary
294 context wherever possible. In order to fully integrate the patched code with the
295 recording, references to shared libraries must point to where the shared libraries
296 were loaded in the recording, references to places in the modified section of the
297 code must point to the appropriate place in the patched code, and references to
298 unmodified contents of the patched binary must point to the appropriate place
299 in the original binary as in Figure 8.

300 In order to prepare for load time quilting resolution (explained shortly),
301 static reference identification needs to occur for bookkeeping purposes. The
302 patched function is scanned for all symbol references that need to be resolved
303 to integrate with the recorded context. Some references like references to lo-
304 cations within the modified function (such as jump and conditional jump in-
305 structions) can remain unaltered in position independent code. So after all
306 references are accounted for, they are trimmed to the subset of references that
307 need to be changed during the quilting procedure. This includes references to
308 strings, shared library functions, functions that only exist in either the original
309 or the modified binary, functions that exist in both, procedure linkage table
310 (PLT) entries, and global variables. Since symbols are the points of reference
311 between original and patched binaries, because recompilation renders addresses
312 meaningless, references to be resolved are defined as a symbol and an offset from
313 that symbol.

314 **Loading Replication & Custom Loading.** In modern Linux systems
315 the system loader is responsible for parsing the executable's header, loading it
316 into memory, and dynamic linking. Since shared libraries are not always loaded
317 at the same positions, references related to the global offset table (GOT), and
318 procedure linkage table (PLT) are resolved after loading completes. Even though
319 ATTUNE knows pre-load which references need resolution, it can't actually
320 resolve those references until load time. To preserve the integrity of the replay,
321 all required shared libraries, executables, and system libraries must be loaded
322 into the recorded memory locations. The trace includes shared libraries and
323 executables required for replay, and non-recorded libraries loaded during replay

324 are limited to the system loader, which is required at the start of any process.

325 In order to replicate the recorded loading activity, ATTUNE begins by load-
 326 ing a small entry point program (replay hook) [that](#) hijacks execution from the
 327 system loader and begins the replay process. Since some references in the
 328 patched code can't be resolved until the original code is loaded into memory,
 329 so initially loading replicates exactly what was recorded. Once the original
 330 segments are loaded into memory and GOT/PLT relocations are completed,
 331 ATTUNE resolves remaining references in the patched code (described below).

332 Finally, ATTUNE's loader loads the quilted code after finding an appropriate
 333 place to put it. Note quilting has to be repeated on every replay, and the files
 334 containing the original and patched executables are not modified. The loader
 335 searches the address space for the lowest slot large enough to accommodate all of
 336 the patched code, then loads the patch following the Linux loading conventions.
 337 Figure 8 depicts the address space when loading has completed, and Algorithm
 338 1 outlines the loading procedure.

Algorithm 1: Custom Loading Algorithm

Result: Load patched code into the address space

```

code_seg_size = 0; char* code_buf;
for func in mod_funcs do
  | code_seg_size += func.size
end
for segment in addr_space do
  | space = next_segment.start - segment.end;
  | if space > code_seg_size then
  | | start = segment.end;
  | | for func in mod_funcs do
  | | | patched_code = func.gen_code;
  | | | code_buf += patched_code;
  | | end
  | end
end
end

```

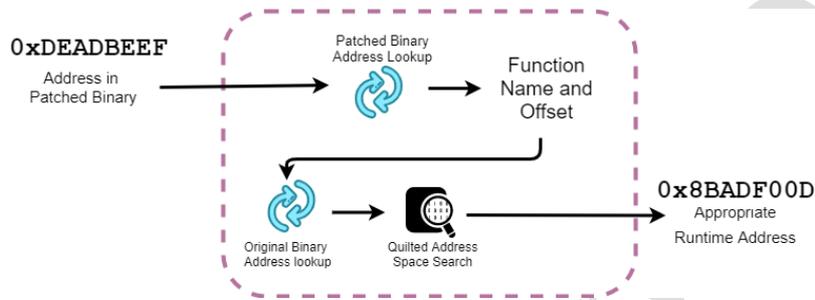


Figure 9: Pointer Translation Procedure

339 **Address Translation Procedure.** A summary of the procedure to trans-
 340 late pointers from the context of the modified binary to the context of the
 341 original binary is given in Figure 9, and consists of both pre-load and load-
 342 time actions. The process starts from the address of the modified function as
 343 determined from the [patch file](#) and DWARF processing. ATTUNE scans the
 344 modified function for references. If a reference is affected by the quilting process,
 345 then ATTUNE’s translation procedure corrects the pointer.

346 The log messages in Figure 10 explain the process in detail: An instruction
 347 in the patched binary at `0x1b214` points to `0xaa60`. In order to update the
 348 instruction to point to the same position in the original binary we need to
 349 identify the correct symbol and offset in the original. First we convert the
 350 target address `0xaa60` into a symbol and offset in the patched binary. Since this
 351 instruction is just calling a function, the target symbol is the function name
 352 and the target offset is 0. Then ATTUNE searches the original binary for the
 353 same symbol and offset, and in this case the function was generated at the same
 354 address in original binary. Resolving string references, global variable references,
 355 and PLT references require slightly different procedures and are described below.
 356 Finally the patched code is generated with instructions pointing to the correct
 357 locations at runtime.

358 **PIC Code, PLT Entries & Trampolines.** Position independent code
 359 compilation has become the standard for security and efficiency reasons, so

```
Linking function: png_check_chunk_length
    in module pngutil
Updating Instruction Reference
    from [0x1b214] to [0xaa60]

//identifying reference point
Target Symbol: png_chunk_error
Offset From Symbol: 0
Symbol Location in original binary:
    0xaa60

//target address in the original binary
Target Address: 0xaa60
...
//patch references string
Resolving string reference at: 0x1b2cd
Resolving offset ...
    for "chunk data is too large"
//identified string in original binary
Found string: "chunk data is too large"
    at 0x320e
... module pngutil code found at 0x000000
... module pngutil data found at 0x200000
... generating quilted code
```

Figure 10: libpng-1 abbreviated linking example



Figure 11: PLT Transformation

360 modern binaries can be loaded anywhere in the address space. As a result the
 361 locations of external functions and symbols are not known until those symbols
 362 actually exist in the address space. Since most library functions aren't called,
 363 they aren't all resolved at load time and instead are resolved only after they
 364 are called. The procedure linkage table (PLT) acts as a table of tiny functions
 365 that perform a function lookup and trampoline to where the code for external
 366 functions are defined.

367 Unfortunately, we can't rely on a PLT because the system loader that per-
 368 forms the runtime function resolution doesn't know about ATTUNE's special
 369 memory configuration. Two key differences let us implement static trampolines
 370 instead of relying on the traditional PLT mechanism. 1) We only need to re-
 371 solve the PLT entries that are referenced by the modified code, which comprise
 372 a small fraction of the overall PLT, and 2) we can resolve these beforehand
 373 without relying on the PLT's lazy loading mechanism because the shared li-
 374 braries have already been loaded by the time this code is injected. The x86_64
 375 architecture only allows call instructions with a 32-bit offset, but we need to call
 376 functions across the 64-bit address space to reference shared library functions.
 377 To accomplish this we transform calls to PLT entries into a move instruction
 378 that loads an address into a register, and then a call instruction to the address
 379 in the register, as shown in Figure 11.

380 **Resolving String & Data Sections.** The patched code may also reference
 381 data section variables like global data and strings. The patched code must
 382 reference the old code where possible and the patched code where required.
 383 Identical symbols and strings function act as points of reference between the
 384 modified and the original binary.

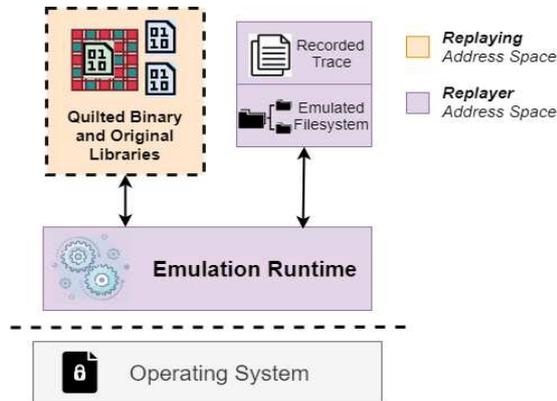


Figure 12: Runtime Architecture

385 These translations are similar to Figure 9, with a few minor differences:
 386 String tables don't have an associated symbol table. The modified code refer-
 387 ences the string directly, but to lookup the location of a specific string in
 388 the original, we have to iterate through all of the read-only data. If the string
 389 exists in the original binary, then we point at it, otherwise ATTUNE points
 390 to the appropriate location in the new data section. Note the binary normally
 391 accesses data through a global offset table entry, but cannot use it here because
 392 the global offset table was compiled for the modified code. Instead, ATTUNE
 393 transforms the binary to point to the data directly, since it knows where the
 394 data has been loaded.

395 3.4. Runtime Replay Decisions

396 The runtime architecture is shown in Figure 12. At runtime we continue to
 397 leverage developer environment information to aid ATTUNE's decision making,
 398 e.g., we know exactly which functions have been modified and perform a strict
 399 replay until a modified function is called. We break at that point and move
 400 to the patched code, where we use information about added or deleted lines to
 401 inform decision making.

402 For any non-deterministic event that takes place during replay, we must
 403 decide whether to use a corresponding event recorded in the log or to actually

404 submit the event for operation by the kernel, i.e., execute live as would be
405 required if the inserted code makes a new system call. We emulate kernel state
406 and kernel events whenever possible, and only ask the kernel to perform the
407 replaying action when necessary, following the greedy approach shown by the
408 pseudocode in Algorithm 2. It should be noted that system calls which depend
409 on process state, like `malloc`, and `mmap`, don't require emulation since this
410 state is actually recreated during replay. All file operations performed during
411 replay are based on information available from the recorded trace, essentially
412 recreating how the program would have acted at the time of the bug except
413 now (for successful patches) without the bug. If there is no further information
414 available, the emulation ends.

415 **System Calls.** The simplest event types to replay are system calls that
416 don't involve file IO. We can reuse results from the log if the parameters for
417 the syscall match what is in the log. It won't match the log exactly since the
418 log contains checks for all registers including the instruction pointer `that` is
419 obviously different, but we relax these checks once replay has diverged to only
420 check registers containing syscall parameters.

421 **File IO.** System calls involving file, network or device IO are harder to
422 replay since they require a specific kernel state. We have to recreate the file
423 state so we track `open`, `close`, `stat`, `read`, `write`, and `seek` operations for all
424 file descriptors during replay. At the point the replay diverges we have a partial
425 view of the file system. Of course we can't recreate any data that doesn't exist,
426 but if a file operation can't be satisfied during replay we can look forward in the
427 recorded trace to see if we have enough information to satisfy the operation. If
428 we do then we emulate it, and unfortunately if we don't we have to die. Another
429 approach would be to supply random bytes, but we feel this wouldn't accurately
430 reflect a realistic state if the full file system were available.

431 **Signal Delivery.** If a signal is intercepted by the emulation engine, we
432 need to decide if that signal should be delivered to the replaying process. Our
433 normal replay mechanism based on `rr`'s replay mechanism determines signal de-
434 livery based on the value of the *retired conditional branches* (RCB) performance

435 counter standard in Intel chips. For signals that have been recorded, we check if
 436 we are in an inserted line. If we are then we deliver the signal and assume it is
 437 created by the patch (such as a segfault from an incorrect memory reference in
 438 the patch). However, if a recorded signal is delivered and we are not currently
 439 in the inserted section of the code we can do our best to estimate at what RCB
 440 count it should be delivered by taking the target RCB count and adding the
 441 number of RCB's caused by inserted lines. While this isn't perfect it does allow
 442 for a rough idea as to when the signal should be delivered. In the event an
 443 unrecorded signal fires we allow that signal to be delivered without interference
 444 since there is no recorded timing information to guide delivery.

Algorithm 2: Runtime Replay Algorithm

Input : e : an event that stops replay

Output: The next event to replay

Function `getResult(e):`

```

  if !diverged then
    | return next_recorded_result;
  if is_syscall_without_file_io  $\&\&$  exists_unused_in_log then
    | return recorded_result;
  if is_syscall_with_file_io  $\&\&$  supported_operation
    exists_unused_in_log then
    | return recorded_result;
  if is_signal  $\&\&$  signal_is_recorded then
    | if current_pos == inserted_code then
    | | return nullptr; // execute live
    | | return DELAY; // delay until RCB count
    | return nullptr; //execute live
  
```

445 3.5. Binary Patch Decomposition and Partial Updates

446 Since all the information for ad hoc testing is at the binary level at runtime,
 447 ATTUNE supports software operators (customer organization IT staff) who
 448 don't have source code but still want test potential updates with their own spe-
 449 cific workloads before deployment. The only requirement is that the developers

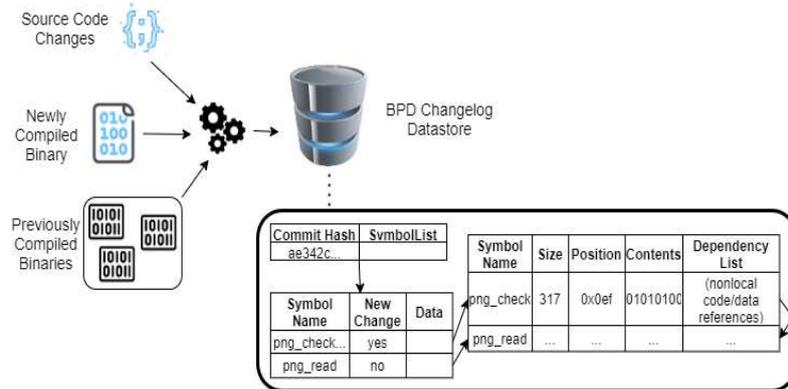


Figure 13: Binary Patch Decomposition Datastore

450 are willing to include metadata describing where the changes in the executable
 451 took place and what they consisted of. To support the developer distributing
 452 this information we developed a novel technique we call binary patch decompo-
 453 sition (BPD), which integrates with the existing build process. In simple terms,
 454 BPD breaks a full update down into *its* component pieces and their contents. *Its*
 455 complexity lies in tracking dependencies between updates such that the version
 456 of the software in the test has the proper contents.

457 Figure 13 outlines the data structure that makes this possible. We integrate
 458 with the version control system as the software is developed and track which
 459 modifications are associated with each commit. In detail, since our ad hoc
 460 testing generation technique depends on symbols in the binary we construct the
 461 metadata that allows the operator to apply the patch based on the symbols
 462 that change. We also track the symbols each piece of code depends on and
 463 the associated versions. Along with the symbols and versions we also store the
 464 contents of those symbols to apply when the ad hoc test is run.

465 Based on the contents of the datastore, the developers can release the binary
 466 specific metadata which details the changes that comprise the update. The
 467 algorithm for constructing the partial update metadata is described in Algorithm

468 3. The algorithm takes the original and new binaries as inputs, and then for
469 every changed symbol it has to do two things. First, if the symbol exists in the
470 old binary, it must add the old size and position data so any references to this
471 piece can be removed. Second, BPD must search through all the dependencies
472 of each changed symbol (a dependency is any nonlocal reference); if the depen-
473 dency existed in the old binary (as per the symbol name), then the metadata
474 can simply add the new code piece and the location of its own dependencies.
475 If the dependency doesn't exist, it must be added to the metadata as a newly
476 changed symbol and its dependencies searched as well.

477 The software operator can check the update for compatibility, and in the
478 event a full update is incompatible with existing infrastructure they can apply
479 a partial update that may still support the old infrastructure. [That partial
480 update may consist of as many individual patches as they would like to include.
481 In the event that selected patches are incompatible \(including multiple versions
482 of the same symbol\) the newest version of the symbol is used.](#)

483 Unlike ad hoc test generation in the developer environment the operator
484 needs to export the test to a modified executable which can be deployed. This
485 distinctly differs from run time quilting as the requirements to keep memory
486 layout the same no longer make sense. Leveraging Egalito's binary rewriting
487 capabilities we completely remove the modified symbols, and replace them with
488 the correct versions. Since Egalito provides arbitrary rewriting, we do so with-
489 out leaving behind any software bloat or extra instructions that would impose
490 performance problems. Effectively we are recompiling the binary to some in-
491 termediate build between version releases despite not having the source code.
492 [ATTUNE's binary rewriting technique avoids the need for recompilation mak-
493 ing it more efficient by both saving compute cycles, and eliminating the need to
494 store source code changes in the BPD datastore. Furthermore by not recompil-
495 ing from selected source code snippets developer practices remain uninterrupted
496 without exposing any source code to the operators \(customer organization IT
497 staff\) deploying the software.](#)

498 As currently constructed ATTUNE requires the developers to ship the entire

499 BPD datastore to [customer organizations](#), but in a commercial setting this
500 datastore would be made available through a shared resource as depicted in
501 Figure 3.

502 Since our BPD technique integrates with git to perform intermediate builds,
503 BPD's ability to create the minimal update is somewhat limited by developer
504 practices. The following circumstances merit further discussion: 1) In the simple
505 case when one bug-fix (involving an arbitrary number of functions) corresponds
506 with one commit, BPD handles this easily. 2) When multiple bug fixes are
507 intertwined in the same commit, git does not provide any way for the developer
508 to distinguish which functions/symbols are associated with each bug-fix so BPD
509 requires that developers update the datastore appropriately. 3) If the same
510 function is updated as part of multiple bug-fixes in different commits, the BPD
511 datastore provides operators access to each version of that function. 4) If the
512 same function is changed as part of multiple bug-fixes in the same commit then
513 the BPD datastore does not support this because there are no intermediate
514 builds per bug-fix to extract different versions of the same function. 5) [If a
515 bug fix involves thread interleavings while ATTUNE's ad hoc test generation
516 provides no guarantees, binary patch decomposition supports apply selected
517 patches of this nature.](#)

518 4. Evaluation

519 We evaluated ATTUNE on a Dell OptiPlex 7040 with Intel core i7-6700 CPU
520 at 3.4GHz with 32GB memory, running Ubuntu 18.04 64bit, using gcc/g++
521 version 7.4.0 and python 3.4.7. ATTUNE is built using CMake version 3.10.2
522 and Make version 4.1.

523 Since we want to evaluate ATTUNE on an unbiased selection of patches
524 for both security vulnerabilities (CVEs) and other kinds of bugs, and know
525 of no benchmark that provides user environment execution traces or scripts
526 to set up the user context for recording traces, we recruited (for one semester
527 of academic credit) an independent challenge team of three graduate students

Algorithm 3: Pseudocode: Metadata Construction

Result: Metadata(old.info, new.info)**Input** : parsed original binary OV; parsed new binary NV; BPD
datastore DB**Output:** metadata information required to construct patch MD**getMetadata** (*OV, NV, DB*)

```

original_code = DB.get_code_pieces(OV);
new_code = DB.get_code_pieces(NV);
res.new_info ← ∅, res.old_info ← ∅;
changed_symbols = DB.getChangedSymbols(NV);
foreach symbol ∈ changed_symbols do
  res.new_info.add(symbol);
  if symbol ∈ original_code then
    res.old_info.add(symbol);
    foreach cp ∈ symbol.dependency_list do
      if cp ∉ original_code then
        sym = Symbol(cp, newChange=True);
        res.new_code.add(sym);
      end
    end
  end
end
return Metadata(res.new_info, res.old_info)

```

528 who were not involved in developing ATTUNE nor versed in how it works.
529 They were tasked to identify a diverse collection of around 20 bugs in widely
530 used C/Linux programs. The bugs had to have been patched during 2016–2019
531 and the students had to construct user contexts that demonstrated the buggy
532 behavior. For example, in order to recreate the circumstances leading up to
533 the redis-1 bug, first one needs to run the server with a specific configuration,
534 connect to the server in MONITOR mode, and then send a specific byte stream
535 to the server. Note the team could script creation of such contexts given the
536 bug and its root cause is already known; record/replay is for capturing and
537 reproducing the contexts of previously unknown bugs. The team identified the
538 21 bugs listed in Table 4.

539 *4.1. ATTUNE successfully validates a wide range of patches provided that cor-*
540 *responding metadata is available*

541 ATTUNE successfully validated the real developer patches in both the [devel-](#)
542 [oper](#) and operator environments for 19 and failed for 2 of the bugs the challenge
543 team collected, marked with ✓ and ✗ in Table 4 resp. We organize the 19 bugs
544 successfully handled into several different types and describe how the developer
545 employs ATTUNE in each case, then explain the 2 failures.

546 **String Parsing** bugs are fairly common as there are often many corner
547 cases, which can have significant security implications since input strings may
548 act as attack vectors. Figure 14 [20] adjusts Curl’s treatment of URLs that end
549 in a single colon. In the buggy version, Curl incorrectly throws an error and
550 never initiates a valid http request. The patch modifies one file. Since ATTUNE
551 replaces the entire modified function instead of individual lines of code, it needs
552 to resolve all references in the new version.

553 ATTUNE uses the recorded execution to recreate the context that triggered
554 the bug, and then jumps to the patched code upon entering the modified func-
555 tion. Since the only change was adding an if statement that doesn’t trigger a
556 recorded event, the ad hoc test continues past the point where the bug occurred,
557 without divergence other than instruction pointer and base pointer. The devel-

```

...
+ if(!portptr[1]) {
+   *portptr = '\0';
+   return CURLUE_OK;
+ }
-   if(rest != &portptr[1]) { ...
-   ...
+ *portptr++ = '\0'; /* cut off the name there */
+ *rest = 0;
+ msnprintf(portbuf, sizeof(portbuf), "%ld", port);
+ u->portnum = port;
...

```

Figure 14: Curl-1 URL Parsing

oper can set a breakpoint at the patched section, watch the if statement process the input correctly and verify the string in **portptr*. The test then ends since the log has no information regarding how the network would have responded to the http request had it been sent.

Figure 15 [21] deals with mishandling URL strings crafted with special characters, e.g., the "#@" in *http://example.com#@evil.com* caused Curl to erroneously send a request to a malicious URL. The patch calls *sscanf* with a different filter string. Since the surrounding function handles all the URL parsing for the application, it is rather large with lots of references. Unlike the above bug, which only requires resolving pointers to old strings, the new filter string needs to be loaded into a new data section and referenced appropriately. ATTUNE recreates the state that caused the initial behavior and then jumps to the modified code. There the developer can verify the patch by checking the values in *protobuf* and *slashbuf*.

Mathematical Errors can have security implications when related to pointer errors or integer overflows. For example, a malicious PNG image triggers a bad calculation of *row_factor* in Figure 16 [19], causing a divide-by-zero error and Denial-of-Service (DoS). With traditional bug reports, the user would need to send the image as an attachment, but a legitimate user affected by the DoS is

```

static CURLcode parseurlandfillconn(...) {
    path[0]=0;
    rc = sscanf(data->change.url,
-   "%15[^\n]:%3[/]%^[\n/?]%^[\n]",
+   "%15[^\n]:%3[/]%^[\n/?#]%^[\n]", /*new data*/
        protobuf, slashbuf, conn->host.name, path);
    if(2 == rc) {
        ....
    }
}

```

Figure 15: Curl-12 String Parsing

```

png_check_chunk_length(...) {
    ...
    size_t row_factor =
-   (png_ptr->width * png_ptr->channels
-   * (png_ptr->bit_depth > 8? 2: 1)
-   + 1 + (png_ptr->interlaced? 6: 0));
+   (size_t)png_ptr->width
+   * (size_t)png_ptr->channels
+   * (png_ptr->bit_depth > 8? 2: 1)
+   + 1
+   + (png_ptr->interlaced? 6: 0);
}

```

Figure 16: libpng-1 Mathematical Error

```

+/* Return non zero if a non breaking space. */
+ static int iswnbospace (wint_t wc) {
+ return ! possibly_correct && (wc == 0x00A0 ...
+ static int isnbospace (int c) {
+ return iswnbospace (btowc (c));
+}
+
wc (args) {
- if (isspace (wide_char))
+ if (isspace (wide_char)
    || isnbospace(wide_char))
    goto mb_word_separator;
    ...
- if (isspace (to_uchar (p[-1])))
+ if (isspace (to_uchar (p[-1]))
    || isnbospace (to_uchar (p[-1])))
    goto word_separator;
}
...

```

Figure 17: wc-1 New Function and Refactoring

577 unlikely to be aware of the carefully crafted malicious image uploaded by an
578 attacker. ATTUNE does not require attachments besides the execution trace,
579 since the re-recorded trace includes the image. After the developer writes the
580 patch, they use ATTUNE to verify that *row_factor* is no longer 0. The patch
581 doesn't trigger any new events so the function returns gracefully.

582 **New Functions & Function Parameter Refactoring.** Many fixes, es-
583 pecially those that pertain to size miscalculations, involve refactoring the buggy
584 function to require a new parameter or writing an entirely new function (with
585 new DWARF and ELF metadata). While not particularly strenuous for the
586 developer, these types of fixes create a challenge for ATTUNE. Since both the
587 function that has been refactored or inserted and the functions that call the
588 new/refactored function need to be modified, ATTUNE must replace all these
589 functions in the executable and properly link them.

590 A patch for the *wc* file processing utility adds special character parsing

```

void addReplyErrorLength
    (client *c, const char *s ...)
{
- if (c->flags & (CLIENT_MASTER|CLIENT_SLAVE)) {
+ if (c->flags & (CLIENT_MASTER|CLIENT_SLAVE)
+   && !(c->flags & CLIENT_MONITOR)) {
+   char* to = c->flags &
+   CLIENT_MASTER? "master": "replica";
...

```

Figure 18: redis-1 Erroneous Conditional

```

url_parse (const char *url ...) {
...
+ /* check for invalid control characters in host
name */
+ for (p = u->host; *p; p++) {
+   if (c_iscntrl(*p)) {
+     url_free(u);
+     error_code = PE_INVALID_HOST_NAME;
+     goto error;
+   }
+ }

```

Figure 19: wget-2 New Loop

591 functions as shown in Figure 17 [22]. ATTUNE loads patched versions of the
592 new function and those functions that call the new function into the address
593 space. The new function is loaded to point towards the original libraries and
594 executables where appropriate, and the modified calling functions point to the
595 new function. There is no need to send a file with the problematic non-standard
596 characters in the bug report to the developer, since it is included in the recorded
597 log. These types of bugs can be difficult for conventional bug reports as files in
598 transit may arrive with modified encoding types and changed contents.

599 ATTUNE provides the input from the recorded file and successfully returns
600 from the modified functions displaying the patched output. Testing the modi-
601 fied *wc* code doesn't diverge drastically from the original execution trace. The
602 developer can verify the patch by letting the program run to termination and
603 inspecting the calculated value.

604 **Adding Conditionals.** Perhaps the most common patch we saw involved
605 adding conditionals. Many security-critical patches make one-line changes to
606 correct conditional checks. We examined one such example in *redis*. Such
607 services are particularly hard to test and debug using conventional mocks, as
608 complex network inputs can be difficult to recreate in mocking frameworks. *Re-*
609 *dis* allows monitor connections to send logging and status checking commands.
610 The buggy version in Figure 18 [23] didn't check the client flags for the monitor,
611 which resulted in a kernel panic. While this was one of the smaller patches, the
612 validation process varied substantially from the log. ATTUNE enables the de-
613 veloper to step through the program and watch progress through the modified
614 control flow past the point of the crash.

615 **New or Changing Loop Conditions.** Bad loop conditionals are also
616 common. Reference resolution is performed as before, but these patches vary
617 greatly from an ad hoc testing perspective because loop conditionals do not
618 necessarily exhibit the bug on the loop's first iteration. One such example from
619 the *wget* utility demonstrates how ATTUNE handles this sort of change in a
620 security-critical situation. The bug allowed attackers to inject arbitrary HTTP
621 headers via CRLF sequences into the URL's host subcomponent. Attackers

622 could insert arbitrary cookies and other header info, perhaps granting access
623 to unauthorized resources. The developer modified the `url_parse` functions in
624 Figure 19 [24] to check each character in the host name and throw an appropriate
625 error. During ad hoc testing the developer verifies the patch works by watching
626 the program check each character, and upon entering the if statement freeing
627 the URL pointer and proceeding correctly to the error handling code.

628 **Swapped Code:** ATTUNE successfully constructed test cases in scenarios
629 that swapped library function calls `yes-1` [25] and swapped control flow blocks
630 `df-1` [26]. The `yes-1` patch makes far-reaching changes across the code base to
631 address the same bug in multiple places (15 files). Assuming the recorded log
632 only manifests one instance of the bug, then the generated ad hoc test case can
633 only check for that instance, not changes elsewhere in the code base.

634 **Failures:** ATTUNE successfully generated ad hoc test cases for those chal-
635 lenging patches where the compiled binaries included complete metadata. How-
636 ever, it failed on **functions with no ELF symbol table entry:** We were
637 initially surprised that a removed break statement in `shred-1` [27] caused an
638 error, since the change is so small. Upon investigation, we found this behavior
639 should be expected, since the function (used only in one place) was inlined by
640 the compiler – thus no symbol table entry for cross-referencing the function.
641 ATTUNE also failed due to **DWARF omissions:** Applying ATTUNE to pa-
642 rameter changes in `curl-8` [28] was unsuccessful. We expected to be able to locate
643 the modified function in the loaded binaries to link the patch, but the DWARF
644 metadata generated by the compiler did not include the filename for the file
645 containing that function. ATTUNE depends on the compiler’s compliance with
646 the DWARF specification.

647 4.2. ATTUNE’s wait time and memory overhead is small

648 To get some perspective of ATTUNE’s overhead, we compared ATTUNE
649 with KLEE [29, 30], a state of the art test suite generation tool. We used
650 KLEE version 2.1 [31] compiled with LLVM 9.0.1. We limited this comparison
651 to those bugs in Table 2 from CoreUtils, since KLEE supports CoreUtils eas-

652 ily. The other bugs we studied have more external libraries, aside from libc, so
653 would require additional engineering effort for KLEE to accommodate. KLEE’s
654 test generation time was budgeted to timeout after 60 minutes, as in [29]. We
655 omitted a comparison to other testing tools that only detect crashing and per-
656 formance bugs, like typical AFL-based fuzzers, since most of the bugs we studied
657 were not crashing bugs and we did not consider performance bugs. We consider
658 this type of testing to be a completely different testing methodology that is not
659 comparable.

660 **Ad Hoc Test Construction Time** ATTUNE’s quilting occurs at load time
661 so runs when each candidate patch is tested. However, since recording allows
662 for targeted test construction, almost all the overhead introduced by symbolic
663 execution searching the program space is removed. As shown in Table 2, our
664 worst case was just under 4 seconds.

665 **Memory Footprint:** ATTUNE inserts patched code prior to each test ex-
666 ecution, so it incurs some memory overhead at test time, as shown in Table 2.
667 *need more here* Symbolic execution, on the other hand, requires significant re-
668 sources to maintain the intermediate program states required to develop test
669 cases. We found on the studied bugs that ATTUNE reduced memory over-
670 head over 90% in all cases and could reduce memory usage by as much as 97%
671 compared to KLEE.

672 4.3. Operators validate released patches with their own workloads and apply par- 673 tial updates if necessary

674 In the last (optional) stage of the patching workflow, the operator validates
675 the patch in their own environment to verify no needed functionality has bro-
676 ken. We integrated our binary patch decomposition datastore so ATTUNE
677 produces correctly formatted metadata enables operators to select individual
678 bug-fix patches from new releases containing other unrelated changes. Since
679 ATTUNE operates entirely in user-space, without the support of hardware,
680 operating system, and so on, it can run in both developer and operator envi-
681 ronments. ATTUNE summarizes the “diffs” in source and binary code, and

Bug	ATTUNE	KLEE	ATTUNE Time	KLEE Time	Speedup	ATTUNE Mem	KLEE Mem	Overhead Reduction
wc-1 [22]	✓	✓	1.37s	300.046s (5m)	99.5%	5.9 KB	108.388 KB	94.5%
wc-2 [32]	✓	✗	1.277s	na	na	2.8 KB	107.7 KB	97.4%
yes-1 [25]	✓	✓	3.4s	8.569s	60.3%	10.6 KB	107.09 KB	90.1%
shred-1 [27]	✗	✗	na	na	na	na	na	na
ls-1 [33]	✓	✓	1.6s	19.57s	91.82%	7.4 KB	132.9 KB	94.4%
mv-1 [34]	✓	✓	3.6s	58.4s	93.84%	4.3 KB	208.2 KB	97%
df-1 [26]	✓	✓	1.48s	18.869s	92.15%	5.97 KB	151 KB	96.05%
bs-1 [35]	✓	✗	1.2s	na	na	5.6 KB	113.37 KB	95.06%

Table 2: Comparison to KLEE Test Generation

```

inserted line addresses:
    0x6b
    0x6e
deleted line addresses:
    0x495AD
    0x495B7
patched code:
...
69:  jne    0xb9
6b:  and   0x2,%eax
6e:  lea   -0x58090939(%rip),%rdx
75:  mov   0x58(%rbx),%rax
...

```

Figure 20: redis-bug-1 Metadata for User Validation

682 exports metadata allowing for operator validation and partial updates.

683 For sample user environment workloads, we used the redis benchmark [36],
684 which simulates thousands of different requests to the server, and the *httperf*
685 benchmarking tool [37] making thousands of connections. Similar to the re-
686 dis discussion above, ATTUNE’s validation procedure for the redis patch [36]
687 utilizes only the metadata it added to the released patch, shown in Figure 20.

688 ATTUNE needs inserted and deleted line addresses for its runtime deci-
689 sion algorithm. The metadata’s ”inserted line addresses” and ”deleted line
690 addresses” are offsets into the relevant files while deleted lines from the original
691 binary are offsets into the original executable. Inserted lines only appear in the
692 patch release so their addresses are offsets into the patched codefile that gets
693 [mmap](#)ed into memory.

694 4.4. Threats to Validity

695 To test our ability to both export and apply partial updates, for each bug
696 we inspected we exported metadata describing each individual change in the
697 complete version update. Then we quilted the singular change that fixed the
698 patch as an operator would apply a single change at a time. Unlike ad hoc test
699 generation in the developer environment, when the modified executable exists
700 statically, in the operator environment we provide the ability to export the in

701 memory ad hoc test to a static file. For every bug in the table the partially
702 updated patched version successfully fixed the bug.

703 It's important to note that the size and scope of the change is not accurately
704 measured only by the lines of code changed, but also how many references need
705 to be resolved in the quilting procedure. Table 4 describes the extent of the
706 changes at the binary level, by tracking how many data reference and code
707 reference resolutions need to be performed to successfully quilt the patch in.
708 Table 4 also shows how many individual changes are in each version update.
709 For each partial update that was applied, the exported version of the binary
710 successfully fixed the buggy behavior.

711 **Internal.** As far as we know, no execution traces were recorded when any
712 of the studied bugs were discovered, so we needed bug-triggering user contexts
713 that could be recorded. We recorded directly with rr, rather than first using a
714 lightweight recorder and then re-recording the lightweight replay using rr. Ar-
715 guably, these user contexts could have been designed to facilitate ATTUNE's
716 test generation. This threat is partially mitigated since the carefully crafted
717 scenarios were developed by three grad students who were not ATTUNE de-
718 velopers and did not know how ATTUNE operates. We did, however, instruct
719 them on how to use rr. Further, we describe how we imagine a developer would
720 verify their candidate patches using ATTUNE, but we are not developers on
721 these projects and lack the developers' knowledge. This is mitigated to some
722 extent since ATTUNE generated ad hoc tests for the real developer patches.
723 Ideally, we would also use ATTUNE to generate ad hoc tests for candidate
724 patches discarded by the developers, to illustrate how we envision a developer
725 would leverage ATTUNE to determine that their attempted patch fails to fix
726 the bug, but we could not find any such commits in the version repositories.
727 Lastly, since do not have access to production users for any of the programs in
728 our dataset, we simulated a production workload using a standard redis bench-
729 mark, which may not be representative of the workloads that production users
730 would construct to validate the redis patch in their own environment.

731 **External.** We demonstrate that ATTUNE supports a wide variety of single-

Bug	Data Resolutions	Code Resolutions	Buggy Version Tag	Patch Version Tag	Distinct Changes Between Versions	Partial Update Success
curl-1 [20]	4	31	curl.7.63.0	curl.7.64.0	128	✓
curl-2 [38]	69	318	curl.7.63.0	curl.7.64.0	128	✓
curl-5 [39]	6	53	curl.7.33.0	curl.7.34.0	246	✓
curl-6 [40]	6	71	curl.7.63.0	curl.7.64.0	128	✓
curl-8 [28]	n/a	n/a	curl.7.61.0	curl.7.60.0	223	✗
curl-9 [41]	8	26	curl.7.62.0	curl.7.63.0	122	✓
curl-10 [42]	3	21	curl.7.62.0	curl.7.63.0	122	✓
curl-11 [43]	273	1012	curl.7.62.0	curl.7.63.0	122	✓
curl-12 [21]	37	103	curl.7.50.0	curl.7.51.0	333	✓
libpng-1 [19]	1	6	v1.6.34	v1.6.35	53	✓
libpng-2 [44]	1	6	v1.6.32beta02	v1.6.33beta02	97	✓
wc-1 [22]	109	298	v8.30	v8.31	90	✓
wc-2 [32]	79	155	v8.26	v8.27	69	✓
yes-1 [25]	234	399	v8.30	v8.31	90	✓
shred-1 [27]	n/a	n/a	v8.27	v8.28	72	✗
ls-1 [33]	380	387	v8.29	v8.30	68	✓
mv-1 [34]	89	204	v8.29	v8.30	68	✓
df-1 [26]	164	348	v8.28	v8.29	65	✓
bs-1 [35]	140	296	v8.28	v8.29	65	✓
wget-1 [45]	8	16	5.0.6	5.0.7	30	✓
redis-1 [23]	3	10	v1.19.5	v1.20	51	✓

Table 4: Partial Update Tests – Partial updates applying a single commit that fixes a patch but each individual change from a version update is available.

732 line and multi-line patches for security vulnerabilities and other bugs in real
733 programs. ATTUNE resolved references between modified and original exe-
734 cutables and program state with binary transformations, but we cannot claim
735 that ATTUNE’s set of transformations will resolve all types of references sup-
736 ported by the expansive x86-64 instruction set. We have not yet studied C++
737 or other non-C programs and we have not yet investigated ARM or other ar-
738 chitectures. The bugs we studied may not be representative of real-world bugs;
739 notably we have not yet studied GUI bugs.

740 **Construct.** The overhead measurements comparing ATTUNE to Klee are
741 arguably unfair, since the symbolic execution explores “from scratch” even
742 though, in principle, Klee’s symbolic execution engine could be modified to
743 leverage rr’s verbose execution traces. We considered integrating Klee with
744 record/replay to be a major research effort, outside the scope of this work. Zuo
745 et al. [10] recently completed such an effort, going even further by skipping
746 ATTUNE’s verbose re-recording entirely, and integrating Klee with lightweight
747 hardware-assisted control and data tracing. Zuo et al. present what they
748 call *shepherded symbolic execution*, where a new production release cycle is in-
749 curred whenever constraint solving bogs down while trying to match the lightly
750 recorded trace. In each new production build, instrumentation is added to cap-
751 ture key data values involved in complex constraint dependencies (long chains
752 of symbolic writes and accesses to large symbolic memory objects). Assum-
753 ing the bug reoccurs sufficiently often in production, after several release cycles
754 the shepherded symbolic execution will eventually find inputs that reproduce
755 the bug (not necessarily the same inputs that triggered the bug when it was
756 originally discovered). Of the thirteen bugs in Zuo et al.’s dataset, two were re-
757 produced from the initial lightweight recording, while the other eleven required
758 from 2 to 10 re-occurrences in production. Their paper did not specify the real-
759 world calendar time involved, but we think it is safe to assume it was longer
760 than the 60 minutes we allowed for Klee timeout.

761 *4.5. Limitations*

762 Our ATTUNE prototype extending rr inherits rr’s design decision to replay
763 multi-threaded recordings on a single thread and simulate thread interleaving
764 by interrupting that single thread’s execution [46, 47]. Although ATTUNE ac-
765 commodates thread synchronization and faithfully emulates the error state, rr’s
766 approach makes it impossible for ATTUNE to accurately verify patches for con-
767 currency bugs that manifest due to the true parallelism of multi-core execution.
768 There is nothing in ATTUNE itself that inherently prevents it from addressing
769 concurrency bugs, but we would need to find a faithfully multi-threading re-
770 placement for rr, ATTUNE also relies on rr to re-record the execution trace in
771 the user environment and to replay that recording in the developer environment
772 with the original version of the program [46, 47, 18]. Since rr was designed to
773 be used during developer testing, with too high overhead for production [46],
774 we adopt the re-recording model shown in Figure 4. In theory, lightweight pro-
775 duction recorders could fail to capture sufficient detail to faithfully replay some
776 behaviors even in the same user environment, in which case the re-recording
777 might not manifest the bug, but Mashtizadeh et al. [9] explain this limitation
778 is generally unimportant in practice.

779 A few ATTUNE limitations are orthogonal to the rr recorder. ATTUNE
780 does not currently verify patches to preprocessor macros, since it compares the
781 source file versions rather than the results of preprocessing the source files.
782 ATTUNE also does not currently support generating tests for patches that
783 change the size of a data structure on the stack or in the heap. We allow new
784 values to be put on the stack and heap, but don’t adjust memory allocation
785 when replaying logged values.

786 Ideally, ATTUNE would address the privacy concerns inherent in all bug
787 report systems that send information gathered in the user environment to the
788 developer. This might be achieved by adding an anonymization phase during
789 or after re-recording with rr, prior to sending to the developer. For example, we
790 could use path conditions and a constraint solver to generate new anonymous
791 data forcing the same execution paths, as was done in [48, 49]. Something

792 like trace wringing [50] might also be an option. We see anonymizing user
793 environment recordings as a major engineering effort that is outside the scope
794 of our research. However, we note that unlike third-party website session script
795 recordings [51], ATTUNE does not run surreptitiously: the user has to select
796 lightweight recordings to re-record and submit to the developer.

797 5. Related Work

798 iFixR [3] automatically generates candidate patches from bug reports, but
799 relies on conventional regression testing even though those tests initially failed
800 to detect the bug. In future work, we plan to investigate integrating ATTUNE
801 with automatic program repair (APR) technology. Differential unit tests [52]
802 construct unit tests using in-memory program state immediately prior to invoking
803 the target method, but cannot reproduce bugs not detected by the original
804 developer tests. [53] similarly extracts unit tests from developer execution
805 traces. In future work, we will investigate constructing unit tests from the ad
806 hoc tests generated by ATTUNE.

807 KATCH [54] combines symbolic execution with heuristics to generate test
808 cases that cover the patched part of the code, while shadow symbol execu-
809 tion [55] symbolically explores divergences between original and patched ver-
810 sions. Neither leverages execution traces recorded in the user environment, nor
811 fully models system calls, so the generated test cases may not reflect the bug-
812 triggering circumstances. However, symbolic execution enables reaching parts
813 of the program not exercised by the recording, complementing ATTUNE.

814 Parallel retro-logging [56] allows developers to change their logging instru-
815 mentation so previous executions produce augmented logs, but the program is
816 not modified. Network-level traffic cloning tools can relay or replay the net-
817 work inputs for service-oriented and microservices architectures. For example,
818 in Parikshan [57] the traffic is fed to a forked copy of an architectural compo-
819 nent in a sandbox, for debugging, or to a modified version of a component, for
820 testing patches. But the replay is not necessarily faithful when there are other

821 sources of non-determinism besides network traffic

822 Kuchta et al. [55] generates tests for software patches using "shadow sym-
823 bolic execution". The old program shadows the new version as the two are
824 symbolically executed in tandem. Whenever new and old diverge, their Shadow
825 tool generates a test exercising the divergence, to comprehensively test new
826 behaviors. Shadow's symbolic execution time budget might permit reaching
827 parts of the program not exercised by available user execution, complement-
828 ing ATTUNE. Shadow does not leverage user execution traces and may not
829 model all system calls, so its tests may not reflect known bug-triggering user
830 environments. It only considers control-flow divergences, not data-only diver-
831 gences, whereas ATTUNE relies on the program and environment state (data)
832 from the verbose re-recording. Further, as explained by Kuchta et al., Shadow
833 suffers from the incompleteness of symbolic execution, the impact of the initial
834 set of inputs, multi-hunk patches (several of our studied patches cross multi-
835 ple files), and the technical limitations of building on top of KLEE and LLVM
836 bitcode — external calls to native code, such as library and system calls, are
837 challenging. ATTUNE assumes the faithful recording of these calls.

838 Elbaum et al. [52] introduced "differential unit tests" generated from the
839 execution traces of developer system tests. Their CR (Carving and Replaying)
840 tool extracts and combines the trace segments that construct in-memory pro-
841 gram state as it was just prior to invoking the target Java method, which then
842 serves as a unit test. CR also complements ATTUNE, since its system tests
843 would likely exercise the program more broadly than available user execution
844 traces. Since CR does not leverage user execution traces and its system traces
845 support only in-memory events, its tests may not reflect known bug-triggering
846 user environments. Other work similarly extracts unit tests from developer exe-
847 cution traces, e.g., [53], with analogous advantages and disadvantages. CR does
848 not attempt to continue replay through the execution of the method under test,
849 the method's return to its caller in the full system execution, and beyond. In
850 contrast, ATTUNE's ad hoc tests are generated from system recordings made
851 in user rather than developer environments, and the primary goal is indeed to

852 continue replay of the full system through the execution of every changed func-
853 tion until its clear the bug no longer manifests — a more challenging problem.
854 ATTUNE requires faithful replay as a baseline, including emulation of interac-
855 tions with files, databases and other resources in the user environment, whereas
856 CR replays only in-memory program state. Tiwari et al. [2] take a similar ap-
857 proach to Elbaum et al., but cull their unit tests from production executions
858 rather than from developer system tests. Their focus is on devising test ora-
859 cles for pseudo-tested methods, where the test suite exercises the methods but
860 no test oracle specifies their expected behavior. As in our user environment
861 validation, Tiwari et al. assume that previous executions in the production
862 environment produced the desired results.

863 A problem posed by Kravets and Tsafrir [58] is more similar to ad hoc test
864 generation. They proposed "mutable replay", where a record-replay engine tries
865 to execute a modified program by closely matching a recorded execution trace
866 from a previous program version. They sketched a hypothetical design based
867 on a then-recent record-replay system, Scribe [59].

868 The Kravets and Tsafrir paper motivated the Scribe developers to implement
869 "mutable replay" themselves [60]. They leveraged checkpoint/restart [61] in
870 a backtracking search algorithm that sought to minimize adds/deletes to the
871 recorded event log. Although successful on many bug-fix examples in the sense
872 that the "mutable replay" continued through the modified portion of the code,
873 the constructed execution was not necessarily the same execution that would
874 have occurred had the modified code been in place in the user environment
875 at the time the original code encountered the bug, which is what ATTUNE
876 aims. Scribe's implementation centered on a special Linux kernel module that
877 intercepted system calls and other non-deterministic events within the operating
878 system kernel, granting complete control over how the kernel responded to the
879 events, whereas ATTUNE runs without privileges in user-space with no changes
880 to the operating system. Scribe required a shared file system (copy on write)
881 between the recording and replaying environments, so was impractical for the
882 post-deployment scenarios we envision, where no files, databases, etc. are shared

883 between user and developer environment other than the contents accessed during
884 verbose re-recording and thus included in the log sent to developers.

885 There are numerous other record-replay tools in the literature, e.g., [62, 63,
886 64, 65]. Some versions of gdb build-in recording and replaying debugging ses-
887 sions [66], as does Microsoft’s IntelliTrace [67]. Some work trades off faithful
888 replay guarantees to better address long-lived latent bugs for time-travel debug-
889 ging [68, 69]. In some record-replay papers the recorded log is referred to as a
890 test case, but most replays only reproduce the buggy execution of the recorded
891 version of the program, and cannot be used to test patched versions.

892 Many record-replay tools focus on reproducing concurrency bugs, e.g., [70,
893 71, 72]. `tsan11rec` [73] combines a custom scheduler for detecting data races
894 with a sparse approach to record/replay: it records only those sources of non-
895 determinism configured per application. `tsan11rec` can record and replay I/O-
896 intensive software like video games, but cannot faithfully replay applications
897 where memory layout non-determinism significantly affects application behav-
898 ior. `rr` faithfully reproduces memory layout, but is not sufficiently perfor-
899 mant for I/O-bound applications – thus our re-recording architecture. While
900 ATTUNE supports ad hoc test generation for multi-threaded programs, our
901 prototype built on `rr` cannot generate tests for patches aimed specifically at
902 concurrency bugs due to how the `rr` implements multi-threading (it simulates
903 multiple threads within a single thread).

904 Much research focuses on reducing recording overhead, trading off lower
905 production overhead (thus better production performance) for faithful replay
906 guarantees, e.g., [74, 75, 76]. REPT [77] combines hardware tracing and binary
907 analysis to try to reconstruct execution traces, which can then be replayed with
908 the same program version. Castor [9] records multi-core applications by leverag-
909 ing hardware-optimized logging, transactional memory, and a custom compiler.
910 Its successful replays allow for slightly modified binaries that do not impact
911 program state. Zuo et al.’s approach outlined in Section 4.4, called Execution
912 Reconstruction (ER) [10], begins with hardware tracing of control and data flow.
913 After a failure, ER uses symbolic execution to find an input that is consistent

914 with the trace, i.e., reproduces the bug. When constraint solving bogs down, ER
915 releases a production patch that records selected data values chosen to shepherd
916 the symbolic execution further. This process iterates. If the failure occurs often
917 enough, eventually sufficient production data will be accumulated to allow the
918 symbolic execution to complete. Thus ER trades off lower production overhead
919 for potentially quite long calendar-time delays in bug reproduction, and there-
920 fore bug repair. We assume some lightweight record/replay system for pervasive
921 recording during production, even though none of them are guaranteed to repro-
922 duce every bug the first time it is encountered. ATTUNE's re-recording with rr
923 in the user environment kicks in only when the lightweight recorder succeeds in
924 reproducing the bug, so could be quite prompt with REPT or Caster but would
925 inherit ER's wait time.

926 Although some papers about record-replay systems refer to capture-replay,
927 e.g., [78], record-replay as discussed in this paper is different from most capture-
928 playback tools. These record or script user actions to repeat for GUI com-
929 patibility testing across multiple operating system versions, browsers, or de-
930 vices [79, 80, 81, 82, 83]. Capture-playback is conceptually similar to ad hoc
931 test generation, but these tools focus on externally visible behavior and are not
932 intended for faithful bug reproduction or testing patches.

933 Multi-Version Execution (MVE) provides an alternative approach to user
934 validation. In MVE, the patched and original versions run *simultaneously* on
935 production user workloads, adding runtime overhead but enabling immediate
936 detection of undesirable divergences [84, 85, 86, 87]. In contrast, we envision
937 that the user records production workloads with the old version and re-records
938 offline as in Figure 2, but skips the developer stage and uses ATTUNE locally
939 to generate ad hoc tests that replay the workloads with the patch. If all is
940 satisfactory, production switches to the new version via some mechanism outside
941 ATTUNE, e.g., live-update. Live-update tools deploy software updates without
942 restarting running programs, e.g., by enabling the new version to resume a
943 checkpoint from the old version similar to a fresh initialization [88, 89]. Dynamic
944 software updating [90] combines multi-version execution with live-update, where

945 the update is applied to a forked copy while the original system continues to
946 operate. The new version shadows the original for a warmup period, and if
947 there are no problems production execution switches over. Unlike MVE systems
948 running different code versions, as in ATTUNE, LDX [91] runs two instances of
949 the same code to infer causality between events. The slave changes one event
950 from the master execution to find divergent impacts, orthogonal to our work.

951 Fuzzing seeks inputs that induce crashes and other problems [92, 93, 94, 95].
952 Other approaches also strive to induce bad behaviors, e.g., [96, 97]. Soltani
953 et al. [98] builds on EvoSuite’s search-based testing [99] to reproduce crashes.
954 Symbolic execution [54], fuzzing [100] and other approaches generate test suites
955 to achieve coverage goals. There is a rich literature concerned with generat-
956 ing inputs intended to trigger or reproduce bugs. Generally, the same generat-
957 ed tests could be applied to multiple program versions — unless those tests
958 are “flaky”. There has also been much work towards making tests repeatable,
959 which is sometimes difficult even in the developer environment on the exact
960 same system build [101]. These tools, as well as regression testing, complement
961 ATTUNE by providing generic testing methodologies, but ATTUNE’s targeted
962 approach based on the specific buggy execution provides a more efficient alter-
963 native. Compared to fuzzing, symbolic execution, and coverage based testing
964 ATTUNE targets the specific buggy execution without relying on approximate
965 heuristics like symbolic conditions and code coverage that cannot guarantee bug
966 reproduction.

967 Research in continuous integration and deployment techniques like those out-
968 lined by Shahin et al. in [102] provide a different functionality than ATTUNE.
969 Some deployment production environments even have test monitoring tools
970 built-in [103]. While they do incremental builds that test software during de-
971 velopment, they do not provide the deploying users a chance to make changes
972 outside of what the developer has distributed.

973 Test case prioritization techniques like the one described by Srivastava et
974 al. [104] and those described by Bajaj et al. [105] look to improve efficiency
975 in regression testing by looking to reduce the total number of tests when re-

976 gression bugs are introduced. While customer organizations could employ such
977 techniques to determine which commits they would like to introduce to their
978 distribution, we leave such discussions to future work. Other code analysis tech-
979 niques like software slicing [106, 107] and branch coverage [108, 109] that are
980 used to augment testing techniques provide a different function than ATTUNE.
981 Slicing identifies groups of statements that are most effected by a given change
982 to inform testing strategies. ATTUNE’s specialized tests remove the need for
983 software slicing as the log guarantees only the critical program paths are un-
984 der test. The popular branch coverage testing techniques are orthogonal to
985 ATTUNE as branch coverage is not a metric that is used by ATTUNE. Of
986 course ATTUNE’s test would cover specific branches of code, but maintaining
987 adequate coverage across the entire suite is still left to developer practices.

988 It should be noted that this technology is significantly different than auto-
989 mated program repair outlined by Le Goues et al. in [110] since we still rely on
990 the human developer to actually write the repair. It also differs from regression
991 testing techniques like those reviewed in [111] by Khatibsyarbini et al.

992 Binary rewriting has been used for many reasons including implementing de-
993 fenses, automatic program repair, hot patching, and optimization. Hot patching
994 is an interesting example since it requires conserving dynamic program state at
995 the time the repair is applied, similar to binary quilting. Katana [112] has highly
996 sophisticated mechanisms for handling this problem, many of which would aug-
997 ment our current quilting procedure, but relies on trampolines to apply the
998 patches that could incur significant overhead the same as [113]. Other binary
999 rewriting mechanisms like Zipr [114, 115] raise the binary to a higher level IR
1000 that allows for increased efficiency in the reassembly process similar to Egal-
1001 ito [17], but have demonstrated generic binary level defense transformations
1002 instead of semantically complex bug specific patching.

1003 6. Conclusion

1004 ATTUNE (Ad hoc Test generation ThroUgh biNary rEwriting) leverages
1005 record-replay and binary rewriting technologies to automate test generation
1006 for security vulnerabilities and other critical bugs discovered post-deployment,
1007 when there are no existing tests for testing candidate patches, and little time for
1008 constructing and vetting new tests. ATTUNE first quilts the modified functions
1009 (the patch) into the original binary and then interprets the recorded execution
1010 trace from the original binary, as it executed in the user environment, to “replay”
1011 on the patched binary in the developer environment. The developer monitors
1012 the progress of the ad hoc test to check that the bug no longer manifests, but
1013 does not intervene in test generation and does not need to build test scaffolding.
1014 We have augmented our original implementation of ATTUNE with binary patch
1015 decomposition, which integrates with the build process to give software oper-
1016 ators (user organization IT staff) the ability to test and apply partial updates
1017 in the event the developer’s full release breaks functionality, e.g., because of
1018 incompatibilities with user environment infrastructure. Our BPD datastore lets
1019 operators selectively apply patches leaving most of the production version un-
1020 touched. We showed that ATTUNE successfully generates tests for a wide range
1021 of known security vulnerabilities and other bugs in recent versions of open-source
1022 software, with minimal developer effort, both quickly and efficiently. We also
1023 demonstrated that BPD can successfully construct updated binaries to address
1024 installation-specific problematic behavior. Our open-source implementation is
1025 available at <https://github.com/Programming-Systems-Lab/ATTUNE>.

1026 7. Acknowledgment

1027 Funding: This material is based upon work supported in part by the National
1028 Science Foundation under Grants No. CNS-1563555 and CCF-1815494.

1029 We thank Yangruibo Ding, Victor Xu, and Ziao Wang for compiling the patch-
1030 testing dataset. We thank Robert O’Callahan and David Williams-King for
1031 their suggestions.

1032 **References**

- 1033 [1] Q. Wang, Y. Brun, A. Orso, Behavioral Execution Comparison: Are Tests
1034 Representative of Field Behavior?, in: IEEE International Conference on
1035 Software Testing, Verification and Validation (ICST), 2017, pp. 321–332.
1036 URL <https://doi.org/10.1109/ICST.2017.36>
- 1037 [2] D. Tiwari, L. Zhang, M. Monperrus, B. Baudry, Production Monitoring
1038 to Improve Test Suites (December 2020). [arXiv:2012.01198](https://arxiv.org/abs/2012.01198).
1039 URL <https://arxiv.org/abs/2012.01198>
- 1040 [3] A. Koyuncu, K. Liu, T. F. Bissyandé, D. Kim, M. Monperrus, J. Klein,
1041 Y. Le Traon, iFixR: Bug Report Driven Program Repair, in: 27th ACM
1042 Joint Meeting on European Software Engineering Conference and Sym-
1043 posium on the Foundations of Software Engineering (ESEC/FSE), 2019,
1044 pp. 314–325.
1045 URL <http://doi.acm.org/10.1145/3338906.3338935>
- 1046 [4] G. Catolino, F. Palomba, A. Zaidman, F. Ferrucci, Not All Bugs Are the
1047 Same: Understanding, Characterizing, and Classifying the Root Cause of
1048 Bugs, *Journal of Systems and Software* 152 (2019) 165–181.
1049 URL [http://www.sciencedirect.com/science/article/pii/](http://www.sciencedirect.com/science/article/pii/S0164121219300536)
1050 [S0164121219300536](http://www.sciencedirect.com/science/article/pii/S0164121219300536)
- 1051 [5] Cybersecurity & Infrastructure Security Agency, CISA COORDINATED
1052 VULNERABILITY DISCLOSURE (CVD) PROCESS, [https://www.](https://www.cisa.gov/coordinated-vulnerability-disclosure-process)
1053 [cisa.gov/coordinated-vulnerability-disclosure-process](https://www.cisa.gov/coordinated-vulnerability-disclosure-process) (Decem-
1054 ber 2019).
- 1055 [6] hackerone, Vulnerability Disclosure Guidelines, [https://www.](https://www.hackerone.com/disclosure-guidelines)
1056 [hackerone.com/disclosure-guidelines](https://www.hackerone.com/disclosure-guidelines) (July 2019).
- 1057 [7] O. Chaparro, C. Bernal-Cárdenas, J. Lu, K. Moran, A. Marcus,
1058 M. Di Penta, D. Poshyanyk, V. Ng, Assessing the Quality of the Steps
1059 to Reproduce in Bug Reports, in: 27th ACM Joint Meeting on European

- 1060 Software Engineering Conference and Symposium on the Foundations of
1061 Software Engineering (ESEC/FSE), 2019, pp. 86–96.
1062 URL <http://doi.acm.org/10.1145/3338906.3338947>
- 1063 [8] D. Mu, A. Cuevas, L. Yang, H. Hu, X. Xing, B. Mao, G. Wang, Under-
1064 standing the Reproducibility of Crowd-reported Security Vulnerabilities,
1065 in: 27th USENIX Security Symposium (USENIX Security 18), USENIX
1066 Association, Baltimore, MD, 2018, pp. 919–936.
1067 URL [https://www.usenix.org/conference/usenixsecurity18/
1068 presentation/mu](https://www.usenix.org/conference/usenixsecurity18/presentation/mu)
- 1069 [9] A. J. Mashtizadeh, T. Garfinkel, D. Terei, D. Mazieres, M. Rosenblum,
1070 Towards Practical Default-On Multi-Core Record/Replay, in: 22nd In-
1071 ternational Conference on Architectural Support for Programming Lan-
1072 guages and Operating Systems (ASPLOS), 2017, pp. 693–708.
1073 URL <http://doi.acm.org/10.1145/3037697.3037751>
- 1074 [10] G. Zuo, J. Ma, A. Quinn, P. Bhatotia, P. Fonseca, B. Kasikci, Execution
1075 Reconstruction: Harnessing Failure Reoccurrences for Failure Reproduc-
1076 tion, in: Proceedings of the 42nd ACM SIGPLAN International Confer-
1077 ence on Programming Language Design and Implementation, PLDI 2021,
1078 Association for Computing Machinery, New York, NY, USA, 2021, p.
1079 1155–1170. doi:10.1145/3453483.3454101.
1080 URL <https://doi.org/10.1145/3453483.3454101>
- 1081 [11] M. Dias, A. Bacchelli, G. Gousios, D. Cassou, S. Ducasse, Untangling
1082 Fine-Grained Code Changes, in: 2015 IEEE 22nd International Confer-
1083 ence on Software Analysis, Evolution, and Reengineering (SANER), 2015,
1084 pp. 341–350.
1085 URL <https://doi.org/10.1109/SANER.2015.7081844>
- 1086 [12] J. Tucek, W. Xiong, Y. Zhou, Efficient Online Validation with Delta Ex-
1087 ecution, in: 14th International Conference on Architectural Support for
1088 Programming Languages and Operating Systems (ASPLOS), 2009, pp.

- 1089 193–204.
1090 URL <http://doi.acm.org/10.1145/1508244.1508267>
- 1091 [13] W. Lam, S. Winter, A. Wei, T. Xie, D. Marinov, J. Bell, A Large-Scale
1092 Longitudinal Study of Flaky Tests, Proceedings of the ACM on Program-
1093 ming Languages 4 (OOPSLA) (2020) 1–29.
1094 URL <https://doi.org/10.1145/3428270>
- 1095 [14] A. Ng, How the Equifax hack happened, and what still needs to be done
1096 (September 7 2018).
1097 URL <https://tinyurl.com/27633662>
- 1098 [15] A. Saieva, S. Singh, G. Kaiser, Ad hoc Test Generation Through Binary
1099 Rewriting, in: IEEE 20th International Working Conference on Source
1100 Code Analysis and Manipulation (SCAM), 2020, pp. 115–126.
1101 URL <https://doi.org/10.1109/SCAM51674.2020.00018>
- 1102 [16] A. Saieva, G. Kaiser, Binary Quilting to Generate Patched Executables
1103 without Compilation, in: Workshop on Forming an Ecosystem Around
1104 Software Transformation (FEAST), 2020, pp. 3–8.
1105 URL <https://doi.org/10.1145/3411502.3418424>
- 1106 [17] D. Williams-King, H. Kobayashi, K. Williams-King, G. Patterson,
1107 F. Spano, Y. J. Wu, J. Yang, V. P. Kemerlis, Egalito: Layout-Agnostic
1108 Binary Recompilation, in: 25th International Conference on Architectural
1109 Support for Programming Languages and Operating Systems (ASPLOS),
1110 2020, pp. 133–147.
1111 URL <https://doi.org/10.1145/3373376.3378470>
- 1112 [18] Mozilla, what rr does (2021).
1113 URL <https://rr-project.org/>
- 1114 [19] Red Hat Bugzilla – Bug 1599943, libpng: Integer overflow and resul-
1115 tant divide-by-zero, https://bugzilla.redhat.com/show_bug.cgi?id=

- 1116 1599943, CVE: <https://nvd.nist.gov/vuln/detail/CVE-2018-13785>
1117 (2019).
- 1118 [20] Curl string parsing bug, <https://github.com/curl/curl/pull/3365>
1119 (2019).
- 1120 [21] Curl string parsing vulnerability, [https://github.com/curl/](https://github.com/curl/curl/commit/3bb273db7)
1121 [curl/commit/3bb273db7](https://github.com/curl/curl/commit/3bb273db7), CVE: [https://curl.haxx.se/docs/](https://curl.haxx.se/docs/CVE-2016-8624.html)
1122 [CVE-2016-8624.html](https://curl.haxx.se/docs/CVE-2016-8624.html) (2019).
- 1123 [22] wc special character bug, [https://github.com/coreutils/coreutils/](https://github.com/coreutils/coreutils/commit/a5202bd58531923e)
1124 [commit/a5202bd58531923e](https://github.com/coreutils/coreutils/commit/a5202bd58531923e) (2019).
- 1125 [23] Redis monitor request causes crash, [https://github.com/antirez/](https://github.com/antirez/redis/commit/e2c1f80b)
1126 [redis/commit/e2c1f80b](https://github.com/antirez/redis/commit/e2c1f80b) (2019).
- 1127 [24] wget insert new loop to parse URL's, [https://github.com/mirror/](https://github.com/mirror/wget/commit/4d729e322fae)
1128 [wget/commit/4d729e322fae](https://github.com/mirror/wget/commit/4d729e322fae), CVE: [https://nvd.nist.gov/vuln/](https://nvd.nist.gov/vuln/detail/CVE-2017-6508)
1129 [detail/CVE-2017-6508](https://nvd.nist.gov/vuln/detail/CVE-2017-6508). (2019).
- 1130 [25] yes coreutils library function, [https://github.com/coreutils/](https://github.com/coreutils/coreutils/commit/44af84263e)
1131 [coreutils/commit/44af84263e](https://github.com/coreutils/coreutils/commit/44af84263e) (2019).
- 1132 [26] df coreutils library function, [https://github.com/coreutils/](https://github.com/coreutils/coreutils/commit/b04ce61958c)
1133 [coreutils/commit/b04ce61958c](https://github.com/coreutils/coreutils/commit/b04ce61958c) (2019).
- 1134 [27] shred coreutils library function, [https://github.com/coreutils/](https://github.com/coreutils/coreutils/commit/c34f8d5c787e6)
1135 [coreutils/commit/c34f8d5c787e6](https://github.com/coreutils/coreutils/commit/c34f8d5c787e6) (2019).
- 1136 [28] Curl change parameter fix, [https://github.com/curl/curl/commit/](https://github.com/curl/curl/commit/e50a2002)
1137 [e50a2002](https://github.com/curl/curl/commit/e50a2002) (2019).
- 1138 [29] C. Cadar, D. Dunbar, D. Engler, KLEE: Unassisted and Automatic Gen-
1139 eration of High-Coverage Tests for Complex Systems Programs, in: 8th
1140 USENIX Conference on Operating Systems Design and Implementation
1141 (OSDI), 2008, pp. 209—224.

- 1142 URL [https://www.usenix.org/legacy/events/osdi08/tech/full_](https://www.usenix.org/legacy/events/osdi08/tech/full_papers/cadar/cadar_html/index.html)
1143 [papers/cadar/cadar_html/index.html](https://www.usenix.org/legacy/events/osdi08/tech/full_papers/cadar/cadar_html/index.html)
- 1144 [30] KLEE Team, KLEE LLVM execution engine, <http://klee.github.io/>
1145 (2021).
- 1146 [31] KLEE Team, Stable releases of KLEE (2021).
1147 URL <http://klee.github.io/releases/>
- 1148 [32] wc reports wrong byte counts when using '-from-files0=-', [https://](https://debbugs.gnu.org/cgi/bugreport.cgi?bug=23073)
1149 debbugs.gnu.org/cgi/bugreport.cgi?bug=23073 (2016).
- 1150 [33] ls -aA shows . and .. in an empty directory, [https://debbugs.gnu.org/](https://debbugs.gnu.org/cgi/bugreport.cgi?bug=30963)
1151 [cgi/bugreport.cgi?bug=30963](https://debbugs.gnu.org/cgi/bugreport.cgi?bug=30963) (2018).
- 1152 [34] 'cp -n -u' and 'mv -n -u' now consistently ignore the -u
1153 option, [https://github.com/coreutils/coreutils/commit/](https://github.com/coreutils/coreutils/commit/7e244891b0c41bbf9f5b5917d1a71c183a8367ac)
1154 [7e244891b0c41bbf9f5b5917d1a71c183a8367ac](https://github.com/coreutils/coreutils/commit/7e244891b0c41bbf9f5b5917d1a71c183a8367ac) (2018).
- 1155 [35] Running b2sum with -check option, and simply providing a
1156 string "BLAKE2", [https://debbugs.gnu.org/cgi/bugreport.cgi?](https://debbugs.gnu.org/cgi/bugreport.cgi?bug=28860)
1157 [bug=28860](https://debbugs.gnu.org/cgi/bugreport.cgi?bug=28860) (2017).
- 1158 [36] Redis benchmark tests server functionality, [https://github.com/](https://github.com/antirez/redis)
1159 [antirez/redis](https://github.com/antirez/redis) (2019).
- 1160 [37] The httperf HTTP load generator, <https://github.com/httplib> (2019).
- 1161 [38] Curl string parsing bug, <https://github.com/curl/curl/pull/3381>
1162 (2019).
- 1163 [39] Curl info leak, <https://github.com/curl/curl/pull/3381>, CVE:
1164 <https://curl.haxx.se/docs/CVE-2017-1000101.html> (2019).
- 1165 [40] Curl security vulnerability, <https://github.com/curl/curl/pull/3433>
1166 (2019).

- 1167 [41] Curl_follow: accept non-supported schemes for "fake"
1168 redirects, [https://github.com/curl/curl/commit/
1169 2c5ec339ea67f43ac370ae77636a0f915cc5fbeb](https://github.com/curl/curl/commit/2c5ec339ea67f43ac370ae77636a0f915cc5fbeb) (2018).
- 1170 [42] URL: fix IPv6 numeral address parser, [https://github.com/curl/curl/
1171 pull/3219](https://github.com/curl/curl/pull/3219) (2018).
- 1172 [43] Curl globbing error, [https://github.com/curl/curl/issues/3251
1173](https://github.com/curl/curl/issues/3251) (2019).
- 1174 [44] libpng IDAT miscalculation, [https://sourceforge.net/p/libpng/
1175 bugs/270/](https://sourceforge.net/p/libpng/bugs/270/) (2019).
- 1176 [45] Simple fix stops creating the log when using -o and -q in
1177 the background, [https://github.com/mirror/wget/commit/
1178 7ddcebd61e170fb03d361f82bf8f5550ee62a1ae](https://github.com/mirror/wget/commit/7ddcebd61e170fb03d361f82bf8f5550ee62a1ae) (2018).
- 1179 [46] R. O'Callahan, C. Jones, N. Froyd, K. Huey, A. Noll, N. Partush,
1180 Engineering Record and Replay for Deployability, in: USENIX Annual
1181 Technical Conference (USENIX ATC), 2017, pp. 377–389.
1182 URL [https://www.usenix.org/conference/atc17/
1183 technical-sessions/presentation/ocallahan](https://www.usenix.org/conference/atc17/technical-sessions/presentation/ocallahan)
- 1184 [47] R. O'Callahan, C. Jones, N. Froyd, K. Huey, A. Noll, N. Partush, Engi-
1185 neering Record And Replay For Deployability: Extended Technical Report
1186 (May 2017).
1187 URL <http://arxiv.org/abs/1705.05937>
- 1188 [48] M. Castro, M. Costa, J.-P. Martin, Better Bug Reporting with Better
1189 Privacy, in: 13th International Conference on Architectural Support for
1190 Programming Languages and Operating Systems (ASPLOS), 2008, pp.
1191 319—328.
1192 URL <https://doi.org/10.1145/1346281.1346322>
- 1193 [49] J. Clause, A. Orso, Camouflage: Automated Anonymization of Field Data,
1194 in: 33rd International Conference on Software Engineering (ICSE), 2011,

- 1195 pp. 21—30.
1196 URL <https://doi.org/10.1145/1985793.1985797>
- 1197 [50] D. Dangwal, W. Cui, J. McMahan, T. Sherwood, Safer Program Behavior
1198 Sharing Through Trace Wrangling, in: 24th International Conference on
1199 Architectural Support for Programming Languages and Operating Sys-
1200 tems (ASPLOS), 2019, pp. 1059—1072.
1201 URL <https://doi.org/10.1145/3297858.3304074>
- 1202 [51] S. Englehardt, G. Acar, A. Narayanan, No boundaries: Exfiltration of
1203 personal data by session-replay scripts, Freedom to Tinker.
1204 URL [https://freedom-to-tinker.com/2017/11/15/
1205 no-boundaries-exfiltration-of-personal-data-by-session-replay-scripts/](https://freedom-to-tinker.com/2017/11/15/no-boundaries-exfiltration-of-personal-data-by-session-replay-scripts/)
- 1206 [52] S. Elbaum, H. N. Chin, M. B. Dwyer, M. Jorde, Carving and Replaying
1207 Differential Unit Test Cases from System Test Cases, IEEE Transactions
1208 on Software Engineering (TSE) 35 (1) (2009) 29–45.
1209 URL <https://doi.org/10.1109/TSE.2008.103>
- 1210 [53] F. Křikava, J. Vitek, Tests from Traces: Automated Unit Test Extraction
1211 for R, in: 27th ACM SIGSOFT International Symposium on Software
1212 Testing and Analysis (ISSTA), 2018, pp. 232—241.
1213 URL <https://doi.org/10.1145/3213846.3213863>
- 1214 [54] P. D. Marinescu, C. Cadar, KATCH: High-Coverage Testing of Software
1215 Patches, in: 9th Joint Meeting of the European Software Engineering
1216 Conference / ACM SIGSOFT Symposium on the Foundations of Software
1217 Engineering (ESEC/FSE), 2013, pp. 235–245.
1218 URL <https://doi.org/10.1145/2491411.2491438>
- 1219 [55] T. Kuchta, H. Palikareva, C. Cadar, Shadow Symbolic Execution for test-
1220 ing software patches, ACM Transactions on Software Engineering and
1221 Methodology (TOSEM) 27 (3) (2018) 10:1–10:32.
1222 URL <http://doi.acm.org/10.1145/3208952>

- 1223 [56] A. Quinn, J. Flinn, M. Cafarella, Sledgehammer: Cluster-fueled debug-
1224 ging, in: 12th USENIX Conference on Operating Systems Design and
1225 Implementation (OSDI), 2018, pp. 545–560.
1226 URL <https://dl.acm.org/doi/10.5555/3291168.3291208>
- 1227 [57] N. Arora, J. Bell, F. Ivančić, G. Kaiser, B. Ray, Replay Without Record-
1228 ing of Production Bugs for Service Oriented Applications, in: 33rd
1229 ACM/IEEE International Conference on Automated Software Engineer-
1230 ing (ASE), 2018, pp. 452–463.
1231 URL <http://doi.acm.org/10.1145/3238147.3238186>
- 1232 [58] I. Kravets, D. Tsafir, Feasibility of Mutable Replay for Automated
1233 Regression Testing of Security Updates, in: 2nd Workshop on Runtime
1234 Environments, Systems, Layering and Virtualized Environments (RE-
1235 SoLVE), 2012, pp. 1–6.
1236 URL [http://www.dcs.gla.ac.uk/conferences/resolve12/papers/
1237 session4_paper2.pdf](http://www.dcs.gla.ac.uk/conferences/resolve12/papers/session4_paper2.pdf)
- 1238 [59] O. Laadan, N. Viennot, J. Nieh, Transparent, Lightweight Application
1239 Execution Replay on Commodity Multiprocessor Operating Systems, in:
1240 ACM SIGMETRICS International Conference on Measurement and Mod-
1241 eling of Computer Systems, 2010, pp. 155–166.
1242 URL <http://doi.acm.org/10.1145/1811039.1811057>
- 1243 [60] N. Viennot, S. Nair, J. Nieh, Transparent Mutable Replay for Multicore
1244 Debugging and Patch Validation, in: 18th International Conference on
1245 Architectural Support for Programming Languages and Operating Sys-
1246 tems (ASPLOS), 2013, pp. 127–138.
1247 URL <http://doi.acm.org/10.1145/2451116.2451130>
- 1248 [61] O. Laadan, J. Nieh, Transparent Checkpoint-Restart of Multiple Processes
1249 on Commodity Operating Systems, in: USENIX Annual Technical Con-
1250 ference (ATC), 2007, pp. 25:1–25:14.
1251 URL <https://dl.acm.org/doi/10.5555/1364385.1364410>

- 1252 [62] Y. Zhao, T. Yu, T. Su, Y. Liu, W. Zheng, J. Zhang, W. G. J. Halfond,
1253 ReCDroid: Automatically Reproducing Android Application Crashes
1254 from Bug Reports, in: 41st International Conference on Software En-
1255 gineering (ICSE), 2019, pp. 128–139.
1256 URL <https://doi.org/10.1109/ICSE.2019.00030>
- 1257 [63] E. Pobe, W. K. Chan, AggrePlay: Efficient Record and Replay of Multi-
1258 threaded Programs, in: 27th ACM Joint Meeting on European Software
1259 Engineering Conference and Symposium on the Foundations of Software
1260 Engineering (ESEC/FSE), 2019, pp. 567–577.
1261 URL <http://doi.acm.org/10.1145/3338906.3338959>
- 1262 [64] H. Liu, S. Silvestro, W. Wang, C. Tian, T. Liu, iReplayer: In-situ
1263 and Identical Record-and-replay for Multithreaded Applications, in: 39th
1264 ACM SIGPLAN Conference on Programming Language Design and Im-
1265 plementation (PLDI), 2018, pp. 344–358.
1266 URL <http://doi.acm.org/10.1145/3192366.3192380>
- 1267 [65] Y. Shalabi, M. Yan, N. Honarmand, R. B. Lee, J. Torrellas, Record-Replay
1268 Architecture as a General Security Framework, in: IEEE International
1269 Symposium on High Performance Computer Architecture (HPCA), 2018,
1270 pp. 180–193.
1271 URL <https://doi.org/10.1109/HPCA.2018.00025>
- 1272 [66] GDB Wiki, Process Record and Replay (2013).
1273 URL <https://sourceware.org/gdb/wiki/ProcessRecord>
- 1274 [67] Microsoft, IntelliTrace for Visual Studio Enterprise (C#, Visual Basic,
1275 C++) (2018).
1276 URL [https://docs.microsoft.com/en-us/visualstudio/debugger/
1277 intellitrace?view=vs-2019](https://docs.microsoft.com/en-us/visualstudio/debugger/intellitrace?view=vs-2019)
- 1278 [68] A. Miraglia, D. Vogt, H. Bos, A. Tanenbaum, C. Giuffrida, Peeking into
1279 the Past: Efficient Checkpoint-Assisted Time-Traveling Debugging, in:

- 1280 2016 IEEE 27th International Symposium on Software Reliability Engi-
1281 neering (ISSRE), 2016, pp. 455–466.
1282 URL <https://doi.org/10.1109/ISSRE.2016.9>
- 1283 [69] D. Vogt, Efficient High Frequency Checkpointing for Recovery and
1284 Debugging, Ph.D. thesis, Vrije Universiteit Amsterdam (2019).
1285 URL [https://research.vu.nl/ws/portalfiles/portal/77028965/
1286 complete+dissertation.pdf](https://research.vu.nl/ws/portalfiles/portal/77028965/complete+dissertation.pdf)
- 1287 [70] Y. Cai, B. Zhu, R. Meng, H. Yun, L. He, P. Su, B. Liang, Detecting
1288 Concurrency Memory Corruption Vulnerabilities, in: 27th ACM Joint
1289 Meeting on European Software Engineering Conference and Symposium
1290 on the Foundations of Software Engineering (ESEC/FSE), 2019, pp. 706—
1291 -717.
1292 URL <https://doi.org/10.1145/3338906.3338927>
- 1293 [71] Y. Hu, I. Neamtiu, A. Alavi, Automatically Verifying and Reproducing
1294 Event-Based Races in Android Apps, in: 25th International Symposium
1295 on Software Testing and Analysis (ISSTA), 2016, pp. 377–388.
1296 URL <http://doi.acm.org/10.1145/2931037.2931069>
- 1297 [72] S. Rattanasuksun, T. Yu, W. Srisa-An, G. Rothermel, RRF: A Race Re-
1298 production Framework for Use in Debugging Process-Level Races, in: 2016
1299 IEEE 27th International Symposium on Software Reliability Engineering
1300 (ISSRE), 2016, pp. 162–172.
1301 URL <https://doi.org/10.1109/ISSRE.2016.35>
- 1302 [73] C. Lidbury, A. F. Donaldson, Sparse Record and Replay with Controlled
1303 Scheduling, in: 40th ACM SIGPLAN Conference on Programming Lan-
1304 guage Design and Implementation (PLDI), 2019, pp. 576–593.
1305 URL <http://doi.acm.org/10.1145/3314221.3314635>
- 1306 [74] A. Orso, B. Kennedy, Selective Capture and Replay of Program Execu-
1307 tions, in: 3rd International Workshop on Dynamic Analysis (WODA),

- 1308 2005, pp. 1—7.
1309 URL <https://doi.org/10.1145/1083246.1083251>
- 1310 [75] Y. Hu, T. Azim, I. Neamtiu, Versatile yet Lightweight Record-and-Replay
1311 for Android, in: ACM SIGPLAN International Conference on Object-
1312 Oriented Programming, Systems, Languages, and Applications (OOP-
1313 SLA), 2015, pp. 349–366.
1314 URL <http://doi.acm.org/10.1145/2814270.2814320>
- 1315 [76] S. Joshi, A. Orso, SCARPE: A Technique and Tool for Selective Capture
1316 and Replay of Program Executions, in: 23rd IEEE International Confer-
1317 ence on Software Maintenance (ICSM), 2007, pp. 234–243.
1318 URL <https://doi.org/10.1109/ICSM.2007.4362636>
- 1319 [77] W. Cui, X. Ge, B. Kasikci, B. Niu, U. Sharma, R. Wang, I. Yun, REPT:
1320 Reverse Debugging of Failures in Deployed Software, in: 12th USENIX
1321 Conference on Operating Systems Design and Implementation (OSDI),
1322 2018, pp. 17–32.
1323 URL <https://dl.acm.org/doi/10.5555/3291168.3291171>
- 1324 [78] J. Steven, P. Chandra, B. Fleck, A. Podgurski, JRapture: A Cap-
1325 ture/Replay Tool for Observation-Based Testing, in: ACM SIGSOFT In-
1326 ternational Symposium on Software Testing and Analysis (ISSTA), 2000,
1327 pp. 158—167.
1328 URL <https://doi.org/10.1145/347324.348993>
- 1329 [79] Microsoft, WinAppDriver (2021).
1330 URL <https://github.com/Microsoft/WinAppDriver>
- 1331 [80] SeleniumHQ, Browser automation (2021).
1332 URL <https://www.seleniumhq.org>
- 1333 [81] appium, Automation for apps (2021).
1334 URL <http://appium.io/>

- 1335 [82] S. Negara, N. Esfahani, R. P. L. Buse, Practical Android Test Recording
1336 with Espresso Test Recorder, in: 41st International Conference on Soft-
1337 ware Engineering: Software Engineering in Practice (ICSE-SEIP), 2019,
1338 pp. 193–202.
1339 URL <https://doi.org/10.1109/ICSE-SEIP.2019.00029>
- 1340 [83] T. Ki, C. M. Park, K. Dantu, S. Y. Ko, L. Ziarek, Mimic: Ui compatibility
1341 testing system for Android apps, in: 41st International Conference on
1342 Software Engineering (ICSE), 2019, pp. 246–256.
1343 URL <https://doi.org/10.1109/ICSE.2019.00040>
- 1344 [84] P. Hosek, C. Cadar, Safe Software Updates via Multi-version Execution,
1345 in: International Conference on Software Engineering (ICSE), 2013, pp.
1346 612–621.
1347 URL <https://dl.acm.org/doi/10.5555/2486788.2486869>
- 1348 [85] P. Hosek, C. Cadar, VARAN the Unbelievable: An Efficient N-version
1349 Execution Framework, in: 20th International Conference on Architectural
1350 Support for Programming Languages and Operating Systems (ASPLOS),
1351 2015, pp. 339–353.
1352 URL <http://doi.acm.org/10.1145/2694344.2694390>
- 1353 [86] D. Kim, Y. Kwon, W. N. Sumner, X. Zhang, D. Xu, Dual Execution for
1354 On the Fly Fine Grained Execution Comparison, in: 20th International
1355 Conference on Architectural Support for Programming Languages and
1356 Operating Systems (ASPLOS), 2015, pp. 325–338.
1357 URL <http://doi.acm.org/10.1145/2694344.2694394>
- 1358 [87] S. Österlund, K. Koning, P. Olivier, A. Barbalace, H. Bos, C. Giuffrida,
1359 kMVX: Detecting Kernel Information Leaks with Multi-variant Execu-
1360 tion, in: 24th International Conference on Architectural Support for Pro-
1361 gramming Languages and Operating System (ASPLOS), 2019, pp. 559–
1362 572.
1363 URL <http://doi.acm.org/10.1145/3297858.3304054>

- 1364 [88] C. Giuffrida, C. Iorgulescu, G. Tamburrelli, A. S. Tanenbaum, Automat-
1365 ing Live Update for Generic Server Programs, *IEEE Transactions on Soft-*
1366 *ware Engineering* 43 (3) (2017) 207–225.
1367 URL <https://doi.org/10.1109/TSE.2016.2584066>
- 1368 [89] S. Kashyap, C. Min, B. Lee, T. Kim, P. Emelyanov, Instant OS Updates
1369 via Userspace Checkpoint-and-Restart, in: *USENIX Annual Technical*
1370 *Conference (USENIX ATC)*, 2016, pp. 605–619.
1371 URL [https://www.usenix.org/conference/atc16/](https://www.usenix.org/conference/atc16/technical-sessions/presentation/kashyap)
1372 [technical-sessions/presentation/kashyap](https://www.usenix.org/conference/atc16/technical-sessions/presentation/kashyap)
- 1373 [90] L. Pina, A. Andronidis, M. Hicks, C. Cadar, MVEDSUA: Higher avail-
1374 ability dynamic software updates via multi-version execution, in: *24th*
1375 *International Conference on Architectural Support for Programming Lan-*
1376 *guages and Operating Systems (ASPLOS)*, 2019, pp. 573–585.
1377 URL <http://doi.acm.org/10.1145/3297858.3304063>
- 1378 [91] Y. Kwon, D. Kim, W. N. Sumner, K. Kim, B. Saltaformaggio, X. Zhang,
1379 D. Xu, LDX: Causality Inference by Lightweight Dual Execution, in: *21st*
1380 *ACM International Conference on Architectural Support for Programming*
1381 *Languages and Operating Systems (ASPLOS)*, 2016, pp. 503–515.
1382 URL <https://doi.org/10.1145/2872362.2872395>
- 1383 [92] R. Padhye, C. Lemieux, K. Sen, L. Simon, H. Vijayakumar, FuzzFactory:
1384 Domain-Specific Fuzzing with Waypoints, *Proceedings of the ACM on*
1385 *Programming Languages (PACMPL)* (2019) 3.
1386 URL <https://doi.org/10.1145/3360600>
- 1387 [93] L. Lampropoulos, M. Hicks, B. C. Pierce, Coverage Guided, Property
1388 Based Testing, *Proceedings of the ACM on Programming Languages*
1389 *(PACMPL)* (2019) 3.
1390 URL <https://doi.org/10.1145/3360607>
- 1391 [94] C. Lemieux, K. Sen, FairFuzz: A Targeted Mutation Strategy for Increas-
1392 ing Greybox Fuzz Testing Coverage, in: *33rd ACM/IEEE International*

- 1393 Conference on Automated Software Engineering (ASE), 2018, pp. 475–
1394 485.
1395 URL <http://doi.acm.org/10.1145/3238147.3238176>
- 1396 [95] J. Choi, J. Jang, C. Han, S. K. Cha, Grey-box Concolic Testing on Binary
1397 Code, in: 41st International Conference on Software Engineering (ICSE),
1398 2019, pp. 736–747.
1399 URL <https://doi.org/10.1109/ICSE.2019.00082>
- 1400 [96] M. Biagiola, A. Stocco, F. Ricca, P. Tonella, Diversity-based Web Test
1401 Generation, in: 27th ACM Joint Meeting on European Software Engi-
1402 neering Conference and Symposium on the Foundations of Software En-
1403 gineering (ESEC/FSE), 2019, pp. 142–153.
1404 URL <http://doi.acm.org/10.1145/3338906.3338970>
- 1405 [97] H. Wu, N. Changhai, J. Petke, Y. Jia, M. Harman, An Empirical Com-
1406 parison of Combinatorial Testing, Random Testing and Adaptive Random
1407 Testing, IEEE Transactions on Software Engineering (TSE) 46 (3) (2018)
1408 302–320.
1409 URL <https://doi.org/10.1109/TSE.2018.2852744>
- 1410 [98] M. Soltani, A. Panichella, A. van Deursen, A Guided Genetic Algorithm
1411 for Automated Crash Reproduction, in: 39th International Conference on
1412 Software Engineering (ICSE), 2017, pp. 209–220.
1413 URL <https://doi.org/10.1109/ICSE.2017.27>
- 1414 [99] G. Fraser, A. Arcuri, EvoSuite: Automatic Test Suite Generation for
1415 Object-Oriented Software, in: 19th ACM SIGSOFT Symposium and
1416 the 13th European Conference on Foundations of Software Engineering
1417 (ESEC/FSE), 2011, pp. 416–419.
1418 URL <https://doi.org/10.1145/2025113.2025179>
- 1419 [100] C. Lemieux, K. Sen, FairFuzz: A Targeted Mutation Strategy for In-
1420 creasing Greybox Fuzz Testing Coverage, in: Proceedings of the 33rd

- 1421 ACM/IEEE International Conference on Automated Software Engineer-
1422 ing, Association for Computing Machinery, New York, NY, USA, 2018, p.
1423 475–485.
1424 URL <https://doi.org/10.1145/3238147.3238176>
- 1425 [101] W. Lam, K. Muşlu, H. Sajjani, S. Thummalapenta, A Study on the Life-
1426 cycle of Flaky Tests, in: 42nd International Conference on Software En-
1427 gineering (ICSE), 2020, pp. 1471–1482.
1428 URL <https://doi.org/10.1145/3377811.3381749>
- 1429 [102] M. Shahin, M. Ali Babar, L. Zhu, Continuous Integration, Delivery and
1430 Deployment: A Systematic Review on Approaches, Tools, Challenges and
1431 Practices, IEEE Access 5 (2017) 3909–3943.
1432 URL <https://doi.org/10.1109/ACCESS.2017.2685629>
- 1433 [103] D. Tiwari, L. Zhang, M. Monperrus, B. Baudry, Production Monitoring
1434 to Improve Test Suites, IEEE Transactions on Reliability (2021) 1–17.
1435 URL <https://ieeexplore.ieee.org/document/9526340>
- 1436 [104] P. R. Srivastava, TEST CASE PRIORITIZATION, Journal of Theoretical
1437 & Applied Information Technology 4 (3).
1438 URL [https://www.researchgate.net/profile/
1439 Dr-Praveen-Srivastava/publication/235799411_Test_
1440 case_prioritization/links/0c960531bf241d5bee000000/
1441 Test-case-prioritization.pdf](https://www.researchgate.net/profile/Dr-Praveen-Srivastava/publication/235799411_Test_case_prioritization/links/0c960531bf241d5bee000000/Test-case-prioritization.pdf)
- 1442 [105] A. Bajaj, O. P. Sangwan, A Systematic Literature Review of Test Case
1443 Prioritization Using Genetic Algorithms, IEEE Access 7 (2019) 126355–
1444 126375.
1445 URL <https://doi.org/10.1109/ACCESS.2019.2938260>
- 1446 [106] X. Li, A. Orso, More Accurate Dynamic Slicing for Better Supporting
1447 Software Debugging, in: 2020 IEEE 13th International Conference on
1448 Software Testing, Validation and Verification (ICST), 2020, pp. 28–38.
1449 URL <https://doi.org/10.1109/ICST46399.2020.00014>

- 1450 [107] B. Stoica, S. K. Sahoo, J. R. Larus, V. S. Adve, Statistical Program Slic-
1451 ing: a Hybrid Slicing Technique for Analyzing Deployed Software, CoRR
1452 abs/2201.00060. arXiv:2201.00060.
1453 URL <https://arxiv.org/abs/2201.00060>
- 1454 [108] S. Chatterjee, A. Shukla, A unified approach of testing coverage-based
1455 software reliability growth modelling with fault detection probability, im-
1456 perfect debugging, and change point, Journal of Software: Evolution and
1457 Process 31 (3) (2019) e2150.
1458 URL <https://doi.org/10.1002/smr.2150>
- 1459 [109] Q. Yang, J. J. Li, D. M. Weiss, A Survey of Coverage Based Testing Tools,
1460 The Computer Journal 52 (5) (2009) 589–597.
1461 URL <https://doi.org/10.1093/comjnl/bxm021>
- 1462 [110] C. Le Goues, M. Pradel, A. Roychoudhury, Automated Program Repair,
1463 Communications of the ACM 62 (12) (2019) 56–65.
1464 URL <https://doi.org/10.1145/3318162>
- 1465 [111] M. Khatibsyarbini, M. A. Isa, D. N. Jawawi, R. Tumeng, Test case prior-
1466 itization approaches in regression testing: A systematic literature review,
1467 Information and Software Technology 93 (2018) 74–93.
1468 URL <https://doi.org/10.1016/j.infsof.2017.08.014>
- 1469 [112] A. Ramaswamy, S. Bratus, S. W. Smith, M. E. Locasto, Katana: A Hot
1470 Patching Framework for ELF Executables, in: International Conference
1471 on Availability, Reliability and Security (ARES), 2010, pp. 507–512.
1472 URL <https://doi.org/10.1109/ARES.2010.112>
- 1473 [113] H. Jeong, J. Baik, K. Kang, Functional Level Hot-patching Platform for
1474 Executable and Linkable Format Binaries, in: IEEE International Con-
1475 ference on Systems, Man, and Cybernetics (SMC), 2017, pp. 489–494.
1476 URL <https://doi.org/10.1109/SMC.2017.8122653>

- 1477 [114] W. H. Hawkins, J. D. Hiser, M. Co, A. Nguyen-Tuong, J. W. David-
1478 son, Zipr: Efficient Static Binary Rewriting for Security, in: 47th Annual
1479 IEEE/IFIP International Conference on Dependable Systems and Net-
1480 works (DSN), 2017, pp. 559–566.
1481 URL <https://doi.org/10.1109/DSN.2017.27>
- 1482 [115] J. Hiser, A. Nguyen-Tuong, W. Hawkins, M. McGill, M. Co, J. Davidson,
1483 Zipr++: Exceptional Binary Rewriting, in: Workshop on Forming an
1484 Ecosystem Around Software Transformation (FEAST), 2017, pp. 9–15.
1485 URL <https://doi.org/10.1145/3141235.3141240>

Generates test cases from record-replay execution trace

Tests candidate fixes in addition to reproducing bug

Decomposes version update binaries into partial patches

Customer selectively tests and applies partial patches

Leverages binary rewriting without requiring source code

Journal Pre-proof

Biography

Anthony Saieva received his BA in Computer Science from Stonehill College, his BS in Computer Engineering from University of Notre Dame, and his MS in Computer Science from Columbia University. Anthony was a PhD student in Computer Science at Columbia University at the time this work was conducted.

Gail Kaiser is a Professor of Computer Science at Columbia University. She received her ScB in Computer Science and Engineering from Massachusetts Institute of Technology, her MS in Computer Science from Carnegie Mellon University, and her PhD in Computer Science from Carnegie Mellon University.

Declaration of interests

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests:

Journal Pre-proof