

Automatic Test-Case Reduction in Proof Assistants: A Case Study in Coq

Jason Gross   

CSAIL, Massachusetts Institute of Technology, 77 Massachusetts Ave., Cambridge, MA 02139, USA
MIRI, USA

Théo Zimmermann   

Inria, Université de Paris, CNRS, IRIF, F-75013, Paris, France

Miraya Poddar-Agrawal  

Reed College, 3203 SE Woodstock Blvd, Portland, OR 97202, USA

Adam Chlipala   

CSAIL, Massachusetts Institute of Technology, 77 Massachusetts Ave., Cambridge, MA 02139, USA

Abstract

As the adoption of proof assistants increases, there is a need for efficiency in identifying, documenting, and fixing compatibility issues that arise from proof assistant evolution. We present the Coq Bug Minimizer, a tool for *reproducing buggy behavior* with *minimal* and *standalone* files, integrated with *coqbot* to trigger *automatically* on Coq reverse CI failures. Our tool eliminates the overhead of having to download, set up, compile, and then explore and understand large developments: enabling Coq developers to easily obtain modular test-case files for fast experimentation. In this paper, we describe insights about how test-case reduction is different in Coq than in traditional compilers. We expect that our insights will generalize to other proof assistants. We evaluate the Coq Bug Minimizer on over 150 CI failures. Our tool succeeds in reducing failures to smaller test cases in roughly 75% of the time. The minimizer produces a fully standalone test case 89% of the time, and it is on average about one-third the size of the original test. The average reduced test case compiles in 1.25 seconds, with 75% taking under half a second.

2012 ACM Subject Classification Software and its engineering → Software evolution; Software and its engineering → Maintaining software; Software and its engineering → Compilers; Software and its engineering → Formal software verification

Keywords and phrases debugging, automatic test-case reduction, Coq, bug minimizer

Digital Object Identifier 10.4230/LIPIcs.CVIT.2016.23

Related Version *Earlier*: <https://jasongross.github.io/papers/2015-coq-bug-minimizer.pdf> [6]

Supplementary Material <https://figshare.com/s/60037cb0c9bbf464e686>

1 Introduction

In the world of machine verification, the dream is to prove the correctness of every program. Projects such as Coq Coq Correct! [13] make significant progress towards this dream for even our most foundational tools: proof assistants themselves. However, large swathes of proof-assistant software—such as tactic languages, elaboration hints, and document managers—remain unproven, lacking even adequate test-suite coverage!

As a solution to expanding the test-suite coverage for the proof assistant Coq, developers adopted “reverse” continuous integration (CI) [18, 10] wherein changes in Coq are tested against a crowdsourced suite of external Coq projects maintained by different teams in different repositories. In this manner, user-centric concerns are well addressed. To prevent the crowdsourced test suite from shrinking, when Coq evolves in a desired direction but breaks some external project in the process, developers of Coq will fix the compatibility



© Jason Gross and Théo Zimmermann and Miraya Poddar-Agrawal and Adam Chlipala; licensed under Creative Commons License CC-BY 4.0

42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:17

Leibniz International Proceedings in Informatics



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

44 issue in the external project. *We believe that to facilitate the use of proof assistants in*
 45 *industry-scale projects, it is essential to make it easy to find, understand, and fix compatibility*
 46 *issues as the proof assistant continues to evolve.*

47 Since the external projects in the Coq test suite are large and intricate, debugging and
 48 fixing failures reported by the reverse CI is a time- and effort-intensive process for developers.
 49 They must perform many steps before beginning to understand and work on the bug. First,
 50 developers will tediously sit through the process of downloading, setting up, and compiling
 51 the external project. Then, they may have to take on the daunting task of figuring out the
 52 larger project context which is not even directly relevant to the bug!

53 The current debugging process can be significantly optimized for developer experience.
 54 Additionally, the current process does not easily yield test cases to add to Coq’s internal test
 55 suite. Instead the test cases remain buried in external developments whereas we would like
 56 to bring bugs to the center! In order to improve the debugging process, we built the Coq
 57 Bug Minimizer¹ which **reproduces buggy behavior in minimal and standalone** files.
 58 Typically, minimized files reduce the total number of lines of code involved in exhibiting
 59 buggy behavior by about a factor of three, making it significantly easier for developers to
 60 observe, play with, understand, and fix bugs. Furthermore, we have integrated the Coq Bug
 61 Minimizer with coqbot [19] to trigger **automatically** on reverse-CI failures, reducing the
 62 friction of building minimized files.

63 Test-case reduction has already a rich literature [3]. However, it is focused mostly on
 64 traditional languages such as C, and even generic reduction techniques may not apply so
 65 well to proof assistants. In this paper, we share what we have learnt about where test-case
 66 reduction is harder and where it is easier in Coq than in traditional compilers, and describe
 67 how we got around the difficulties. Drawing on empirical results from nearly a year of use
 68 in Coq’s production CI system, we reflect on how effective our style of test-case reduction
 69 has been and where the biggest opportunities for improvement remain. We believe that our
 70 methods may be of interest for developers of other proof assistants who are also facing a
 71 tradeoff between enabling evolution and preserving stability, in a context of industrial use.

72 For the mobile reader, Section 2 introduces a constructed example of test-case reduction in
 73 Coq, and articulates desiderata for test-case reduction in the proof-assistant setting. Section 3
 74 details aspects of traditional-setting test-case reduction that are simpler or irrelevant in
 75 Coq. Then Sections 4, 5, 6, and 7 explore the four desiderata and describe the details of our
 76 solution to the more important challenges of each. Section 8 forays into the applicability of
 77 the Coq Bug Minimizer for bug reporter workflow as a secondary use case. Finally, Section 9
 78 presents our deployment in Coq’s production CI, with analysis of how effectively different
 79 test cases were minimized; Section 10 describes connections to related work; and Section 11
 80 discusses our thoughts on the most worthwhile improvements to make to our tooling.

81 **2** Desiderata

82 To add color to our picture, let’s begin with a constructed example of minimizing a reverse-CI
 83 failure. Our objective is to explore the space of file modifications that will aid human
 84 understanding of the bug. Consider the following Coq source file:

```
85 Require Import UsefulTactics.  
86 Definition zero := 0. Definition one := 1.
```

¹ Available on GitHub in JasonGross/coq-tools

```

87 Definition two := 2. Definition three := 3.
88 Lemma foo : forall x, x = zero -> S x = one.
89 Proof. crush. Qed.

```

90 Suppose the `crush` tactic triggered a new bug in Coq. The most obvious move is to find
 91 deletable sentences and delete them, producing a smaller file:

```

92 Require Import UsefulTactics.
93 Definition zero := 0. Definition one := 1.
94 Lemma foo : forall x, x = zero -> S x = one.
95 Proof. crush.

```

96 The file still depends on an imported module not native to the Coq standard library. The
 97 next move is to inline this dependency, producing a standalone file:

```

98 Module UsefulTactics.
99 Ltac head expr := match expr with | ?f _ => head f | _ => expr end.
100 Ltac head_hnf expr := let expr' := eval hnf in expr in head expr'.
101 Ltac crush := intros; subst; try reflexivity.
102 End UsefulTactics.
103 Import UsefulTactics.
104 Definition zero := 0. Definition one := 1.
105 Lemma foo : forall x, x = zero -> S x = one.
106 Proof. crush.

```

107 Now we may look for any more opportunities to delete lines, producing a standalone,
 108 reduced file:

```

109 Ltac crush := intros; subst; try reflexivity.
110 Definition zero := 0. Definition one := 1.
111 Lemma foo : forall x, x = zero -> S x = one.
112 Proof. crush.

```

113 From the above process we can extrapolate desiderata for the Coq Bug Minimizer.

- 114 **1. Reproducing buggy behavior:** Deciding when two source files indicate the same bug.
 115 Many reasonable file simplifications lead to incidental changes in error messages. The
 116 Coq Bug Minimizer must tradeoff between preserving specific details of error messages
 117 and aiding human understanding of the underlying bug.
- 118 **2. Minimal files:** Exploring the space of program simplifications in a smart way with
 119 respect to constraints of the proof-assistant setting. Many research papers in the software-
 120 engineering community have been written on just this topic [5, 17, 2, 14, 16], but
 121 constraints in a proof-assistant setting are uncommon in conventional programming. For
 122 instance, highly automated Coq developments often have long compile times even for
 123 single files, so we may need to be more frugal in how many program variants we test.
- 124 **3. Standalone files:** Creating standalone files that illuminate new test cases and can be
 125 added to Coq's internal test suite. This is difficult in dependently typed languages with
 126 metaprogramming facilities such as Coq. For instance, eliminating needless dependencies
 127 in simply typed languages may be trivial, but dependently typed languages eliminate the
 128 distinction between runtime and compile time resulting in tight coupling between files.

129 **4. Smooth developer experience:** Automatically finding which file triggered a bug, with
 130 which compilation settings, including path information to find dependencies. The Coq
 131 Bug Minimizer must work with the wide variety of build systems used in different Coq
 132 libraries.

133 Achieving each desideratum posed interesting challenges, and required making several
 134 design choices. Before proceeding to share solutions to these challenges, we note the ways in
 135 which test-case reduction is *simpler* in the proof-assistant setting than in other settings.

136 **3 Simplifications of the Proof-Assistant Setting**

137 Classic delta debugging [16] is a technique in test-case reduction for traditional compilers. It
 138 employs binary search through program structure to find subprograms that can be removed
 139 while preserving properties relevant to triggering specific bugs for the chosen compiler. Coq’s
 140 lack of forward references permits a simpler method: first remove everything after the
 141 error-message-generating line, and then try removing the syntactic units beforehand in-order,
 142 one-at-a-time. Unlike in languages from Java to Haskell, where all functions in a file are
 143 considered mutually recursive, in Coq there should be no way for one error-message-generating
 144 line of a file to change behavior based on modifications to later lines. In this manner we
 145 reduce the number of “experiments” on program variants, which is especially useful when
 146 each program variant requires significant processing time as is often the case in Coq.

147 Our empirical evaluation (Section 9) demonstrates that this strategy is adequately
 148 performant. We conjecture that the reason for this adequate performance is that dependency
 149 trees of Coq theorems and proofs tend to be relatively deep compared to the number of
 150 definitions and theorems in any single file. This hypothesis is borne out by the fact that our
 151 typical “minimal” test case tends to be only about a third the size of the total amount of
 152 code in all files in the dependency tree of the initial test case. If instead there were orders
 153 of magnitude more useless lines than true dependent lines, we expect that a binary-search
 154 strategy would be required for adequate performance.

155 **4 Reproducing Buggy Behavior**

156 How do we know modifications to source files are genuine simplifications that have not
 157 masked bugs? What does it mean to reproduce the “same” bug? We generate a file that
 158 succeeds on the previous version of Coq and continues to fail on the modified version of Coq,
 159 with the same error message that showed up in the reverse CI. However, the error message of
 160 the generated file does not need to be *exactly* the same as in the original file, so long as the
 161 reason for the error message is the same. Thus, we modify our goal to reproducing buggy
 162 behavior in place of reproducing the “same” bug.

163 We apply the following relaxations in comparing error messages.

- 164 1. Universe inconsistencies are how Coq prevents users from proving absurdity by assuming
 165 a “set of all sets.” The explanations of universe inconsistencies in error messages are
 166 sensitive to how many universes are floating around and in what order constraints were
 167 added. Rather than requiring output files to mimic the error messages exactly, we only
 168 require that they result in *some* universe inconsistency.
- 169 2. Any two error messages about “forgotten universes” are considered matching, since these
 170 tend to arise only from very specific Coq internal errors.
- 171 3. Usually differences in numbering, e.g. in universes or autogenerated identifiers, are
 172 incidental and are not treated as implying different error messages. One special case

173 is lengths of universe instances, so we look for the text “Universe instance should have
174 length” in the error message and only use number-insensitive comparison if this text is
175 not found.

176 4. We consider any error messages containing “Unsatisfied constraints: ... (maybe a bugged
177 tactic)” as equivalent, since related bugs are localized to one relatively small part of
178 the Coq implementation, and small changes to a source file can modify constraints
179 significantly.

180 5. We also ignore filenames, line references, and word wrapping in comparing error messages.

181 **5 Minimal Files**

182 Test-case reduction is powerful in making long source files more comprehensible to developers.
183 In addition to this, external projects in Coq can take minutes or hours to compile, so the
184 edit-compile-test-debug loop is long. We have two additional goals to improve this workflow.

185 1. Finding minimal test cases as fast as possible, given that experimenting with each program
186 variant has long compilation time.

187 2. Compilation of the test case in seconds or fractions of a second so that developers can
188 fluidly try hypotheses for solutions.

189 **5.1 Making the Minimization Process Itself Fast**

190 In our goal to get the shortest reproducing test case as quickly as possible, it helps to first
191 make any changes that might significantly speed up the execution time, and only after we’re
192 done with all of the changes that might improve running time should we try to further
193 minimize the file with changes that are unlikely to impact compile time.

194 The slowest part of almost all Coq developments is proof scripts. Hence we attempt to
195 remove proof scripts as early as possible. Since proof assistants check that proofs are valid,
196 we cannot simply remove a proof, like we might remove a function body in a traditional
197 programming language. However, most proof assistants have some mechanism for “giving
198 up” on a proof or “trusting” the user, and Coq is no exception. Its mechanism involves any
199 of `Admitted`, `Admit Obligations`, or the `admit` tactic. Replacing proof blocks with these
200 commands, rather than just removing proof scripts, allows us to make much smaller and
201 faster examples than might otherwise be possible.

202 **5.2 Finding Textually Smaller Test Cases**

203 The simplest function of the bug minimizer is to remove unneeded lines. As noted in the
204 prior section, we try removing one syntactic unit at a time, moving backwards from the unit
205 that triggered the error message.

206 However, we can easily enough get stuck in local minima, when we remove single commands
207 and check that bug behavior is unchanged. For instance, there may be an irrelevant lemma
208 that we want to remove.

209 `Lemma irrelevant : two = 2.`

210 `Proof. reflexivity. Qed.`

211 Since Coq forbids nested lemmas, removing statements one-at-a-time will not work, as the
212 state

213 `Lemma irrelevant : two = 2.`

214 `Proof. reflexivity.`

215 results in an error about nested proofs, if there is a theorem afterward.

216 We instead group statements into *definition* blocks to be removed all at once. We get
 217 information about definitions by parsing the output of `coqtop -emacs -time`. This way, we
 218 can remove the lemma block all at once.

219 We could in theory deal with more complicated nesting structure, for example trying to
 220 remove an entire section or module at a time. The delta tool [14] is in fact built around
 221 preprocessing the file into one that exposes nested structure clearly, then removing well-
 222 parenthesized blocks. However, removing statements, grouped into definitions as necessary,
 223 suffices for removing time-consuming code.

224 5.2.1 The Program Construct

225 One Coq construct that does not fit neatly into this approach is `Program`, where a function
 226 definition is associated with following proofs of obligations related to dependent typing. We
 227 cannot just look for `Program` statements followed by `Obligation` blocks to remove all together,
 228 because `Obligation` blocks can be interleaved with other definitions. Luckily, we can replace
 229 any obligation block with a use of the `Admit Obligations` command, which admits all
 230 remaining obligations—and it happily handles any case with *no* remaining obligations, so we
 231 need not worry about introducing duplicate invocations.

232 5.2.2 Empty Sections and Modules

233 Removing statements one-at-a-time will not always be able to remove empty sections (nor
 234 empty `Modules` or `Module Types`). That is why we have a pass dedicated to removing empty
 235 sections, modules, and module types.

236 5.2.3 Exporting Modules

237 Coq’s features to import and export modules (e.g., including all definitions of one module
 238 inside another) can create some particularly thorny situations for statement-at-a-time shrink-
 239 ing. If we remove just an `Import` commands, then later commands fail because important
 240 identifiers are out-of-scope. If we remove just the definition of the imported module, then the
 241 `Import` fails. The solution is to merge these two commands, so that they become a candidate
 242 for removal together. We change `Module` commands into `Module Export` commands to this
 243 end. Often that change renders later `Import` commands redundant, so they are removed by
 244 later passes.

245 5.2.4 Splitting Definitions

246 One pass in the minimizer tries to replace traditional definitions with uses of the interactive
 247 proof mode, which is a first step toward admitting those proof bodies (i.e., postulating
 248 existence of identifiers rather than giving their definitions) in later steps.

249 5.2.5 Early Removal of Unused Constants

250 There are some likely-to-succeed steps that we try early on, which are superseded by removing
 251 each and every structured block one at a time but may result in faster minimization. The
 252 primary example of this sort of step is removing tactics, `Variable` and `Context` statements,
 253 and definitions which are not referred to at all after their definitions.

254 5.2.6 Splitting Import and Exports

255 It may be the case that users import modules that they never use, such as in `Import`
 256 `unused1 used unused2`. To allow eventual removal of `unused1` and `unused2` even when
 257 the `Import used` statement cannot be removed, we have a pass that attempts to split such
 258 statements into separate `Import` statements, resulting in `Import unused1. Import used.`
 259 `Import unused2`.

260 5.3 Finding Test Cases That Coq Processes More Quickly

261 We mentioned how admitting proofs is a very handy step to shrink files and get them
 262 processed more quickly. There are, however, a few gotchas to keep in mind.

263 The first quirk is around transparency vs. opacity of lemma definitions; that is, whether
 264 the generated proof term is accessible to later definitions. Either choice (transparent vs.
 265 opaque) can break some developments. Marking a proof-mode definition *opaque* will break
 266 later definitions that unfold the definition and then perform further tactic-based surgery on it,
 267 while marking a proof-mode definition *transparent* could cause previously failing unfoldings
 268 to succeed. Therefore, we always try both styles of marking a lemma admitted.

269 Some lemma proofs are declared as transparent rather than opaque, where later steps
 270 really do depend on their details. If those dependencies are too specific, then our shrinking
 271 heuristics are not going to work well. However, one common-enough case is where a later
 272 definition uses tactics to *unfold* an earlier definition, going on to use other tactics that may
 273 very well be able to adapt to changes in that definition. There are at least two different ways
 274 to mark a proof as admitted (`Admitted` vs. using a preexisting `Axiom`), which can switch up
 275 whether the associated definition is considered transparent or opaque.

276 Additionally, we may want to admit some parts of the proof script without replacing all
 277 of it. Currently, we use a rather conservative heuristic: Coq has a tactical `abstract` that
 278 executes the tactic it is passed as an argument, making the resulting proof term opaque.
 279 Such subproofs should be able to be replaced with `admit` without changing the behavior of
 280 the proof script. The details are a little subtle, e.g. to avoid changing which section variables
 281 a proof depends on and thus changing its type outside the section.

282 6 Standalone Files

283 While the complex structure of external developments is a boon to stress-testing Coq, there
 284 are three reasons for wanting to reproduce bugs in standalone files.

- 285 1. It is challenging for developers to understand the intricacies of external developments
 286 well enough to diagnose root causes.
- 287 2. Build systems are necessary to handle multiple files, but using them adds unnecessary
 288 overhead in the debugging workflow.
- 289 3. Intricate file-dependency structure complicates test-suite infrastructure, whereas having
 290 self-contained files results in a simpler test suite.

291 Naïvely, the way to produce a standalone file is to *linearize* the dependency tree and
 292 combine the contents of all files. We saw an example of roughly this strategy in Section 2,
 293 and e.g. C compilers follow this strategy in preprocessing `#include` statements.

294 Two difficulties arise when following this strategy in Coq:

- 295 1. As in all languages that allow shadowing of global symbols, inlining files changes what
 296 names are available and hence may result in unintended changes of behavior. The de-
 297 pendent typing and metaprogramming facilities of Coq largely eliminate the distinction

298 between runtime and compile time. As a result, we have to inline not just function
 299 declarations but also function bodies, and thus the problem of name resolution is com-
 300 paratively harder in Coq and similar languages than in those with simple types and
 301 without metaprogramming facilities. Furthermore, Coq has additional quirks around
 302 name resolution and (lack of) namespacing that have to be managed and worked around.

303 2. Coq has a great deal of global state (e.g., notations, universe polymorphism, the default
 304 tactic mode) that changes the way sentences are interpreted. Because there is no way to
 305 isolate changes on this global state fully, there may not even be *any* linearization that
 306 reproduces the same behavior.

307 6.1 Addressing Shadowing and Name Resolution

308 Coq assigns names based on three components: the name and location of the file in which
 309 the identifier is defined, the module structure surrounding the identifier, and the final
 310 name. For example, the constant `Coq.MSets.MSetPositive.PositiveSet.t` is defined in
 311 the file `MSets/MSetPositive.v`, which is bound to `Coq.MSets.MSetPositive`, in the module
 312 `PositiveSet`, with the name `t`.

313 If we were to inline this file into some other file `bug.v`, then the constant becomes
 314 `bug.PositiveSet.t`. We now have two choices: we can attempt to adjust the name of the
 315 constant on inlining, or we can adjust references to the constant.

316 We combine these strategies to maximize the chance of successfully inlining dependencies.

317 First, as shown in the example in Section 2, we wrap the contents in a module whose name
 318 matches that of the file (in this case, we wrap the contents in `Module MSetPositive`). Fur-
 319 thermore, since users can refer to this constant as `Coq.MSets.MSetPositive.PositiveSet.t`,
 320 `MSets.MSetPositive.PositiveSet.t`, or `MSetPositive.PositiveSet.t`, we can wrap this
 321 module in further modules (`Coq` and `MSets`) and `Export` them to make this naming scheme
 322 available. Finally, because Coq forbids multiple modules with the same absolute kernel
 323 name, we must wrap the top-level module in yet another module, with a uniquely generated
 324 identifier. While this strategy is not perfect, running afoul of `coq/coq#14587` for example,
 325 we try a couple of variations on this strategy, and very often one of them is adequate for
 326 reproducing buggy behavior.

327 Second, we want to adjust references so that they still point at the same underlying object
 328 after inlining. Coq helpfully emits *globalization* files, which contain information about how
 329 Coq resolves almost all names in the file. Since Coq generates and installs these `.glob` files,
 330 we can use this information to transform both the names in the files we inline and the names
 331 that refer to constants in that file.

332 However, we cannot just blindly update all names, because these `.glob` files are not
 333 perfectly accurate² and are not complete³. Instead, we have found in practice that the most
 334 important names to resolve are those used in `Require`, `Import`, and `Export` statements.
 335 `Require` statements are sensitive to the searchpath flags (`-Q` and `-R`) passed to Coq. If we
 336 are inlining a file from `Flocq` into a file from `VST`, for example, the `Requires` in the `Flocq` file
 337 may not resolve to the same files on disk when compiling with the compiler flags that `VST`
 338 uses. `Import` and `Export` statements, while not dependent on searchpath flags to the same
 339 extent as `Require`, still seem empirically more likely to refer to potentially ambiguous names

² See `coq/coq#15497` and `coq/coq#14537`.

³ They are missing information, for example, on tactic-name resolution and notation interpretation.

340 than most other statements. Hence we choose to resolve the names used in `Require`, `Import`,
341 and `Export` statements when inlining, letting Coq determine all other name resolution.

342 6.2 Addressing Nonlinearizability of Global State

343 While shadowing and name resolution are mechanically resolvable at least in theory, the
344 global state of Coq is sufficiently disorganized that we are not aware of any fully general
345 technical means of linearizing Coq files.⁴ Hence our approach here consists of several partial
346 workarounds.

347 The most basic technique to attempt to isolate global state is to wrap the inlined file in a
348 module. Most state not explicitly marked as `Global` does not escape the boundaries of the
349 module it is defined inside. As we already use module wrapping to handle name resolution
350 as discussed in Subsection 6.1, we already reap the benefits of this technique.

351 Our only other technique is to try multiple linearizations and hope that one of them
352 is adequate. We try inserting the file being inlined at the top of the file, as well as at the
353 location where it is `Required`. In the future, we might also want to try moving `Requires` up
354 higher in the file, to try to handle more situations.

355 In Section 11, we discuss a few potential future avenues to better handling of global state.
356 For example, we may want to more explicitly manage the state before and after inlining a
357 file by taking advantage of Coq's ability to print the current settings of flags with `Print`
358 `Options`.

359 6.3 Getting to Standalone Files Quickly

360 We have a flag that allows inlining dependencies all at once, much like `gcc` inlines all
361 `#included` files at once. While originally all files were minimized in that way, having to
362 process such a large file slowed down minimization drastically, often resulting in minimization
363 times of multiple weeks. As a result, the current default behavior is to minimize the current
364 file before inlining other files.

365 Furthermore, we want to ensure that we only inline files that are actually used. Much
366 like we want to split `Import` and `Export` statements in Subsubsection 5.2.6, we also want to
367 split `Require` statements, for example from `Require unused1 used unused2.` to `Require`
368 `unused1. Require used. Require unused2.`

369 Additionally, if the buggy behavior depends on a file only for its own dependencies, we
370 prefer to inline the transitive dependency directly rather than needing to inline the entire
371 intermediate file. To that end, we have a pass that performs the transitive closure of the
372 dependency relation, inserting `Require` statements at the top of the file for all transitive
373 dependencies of the file being minimized. Because we insert the `Requires` in dependency
374 order, removing one statement at a time in reverse order will give us the minimal `Requires`
375 needed to reproduce the error message. This strategy ensures that we only inline dependencies
376 that are actually necessary.

⁴ The `Require` command results in many side effects, including global setting of flags, opacity, and argument status; behavior of `auto with *`; hint databases; global overwriting of Ltac definitions; presence or absence of constants that change the behavior of built-in tactics such as `tauto`; and even the presence of constants with certain kernel names can change shadowing behavior. Some of these effects can even be set on the command-line, and at present there is no way to determine what flags were used to compile a given installed file.

377 **7** Smooth Developer Experience

378 In order to analyze a specific source file, we need to take a few steps.

- 379 1. Unpack and install both the succeeding and failing versions of Coq and corresponding
380 developments.
- 381 2. Replace the Coq binaries with wrappers that print out the arguments that Coq was called
382 with, as well as `COQPATH` (an environment variable listing directories to be searched for
383 imported modules) and the current directory.
- 384 3. Run Coq on the succeeding and failing developments, ensuring that the version that
385 should pass does in fact pass, and the version that should fail has a recognizable error
386 message.
- 387 4. Parse the build log to determine the buggy file name and the arguments to pass to Coq,
388 using the extra logging introduced by our wrappers. This workflow means that we need
389 not interface directly with varied build systems of different contributions on the CI.
- 390 5. Run Coq on the buggy file.
- 391 6. Parse the error message, ensuring that it matches with the error message from the build
392 log. (See Section 4 for subtleties in that comparison)

393 Again, the goal of the minimizer is to take a CI development that succeeds on the tip of
394 the master branch and fails on a given pull request (PR), emitting a small, standalone file
395 that succeeds on master and fails in the same way on the PR. In order to do so efficiently, we
396 reuse the CI artifacts from Coq. We download the prebuilt versions of Coq from master and
397 from the tip of the PR. From just these artifacts and the name of the failing CI development,
398 we must assemble enough information to run the bug minimizer. We replicate Coq's generic
399 CI workflow to install Coq as well as any dependencies of this CI development, into different
400 directories: one for the version of Coq expected to pass and another for the version of Coq
401 expected to fail. We also reuse Coq's generic CI workflow to figure out the error message
402 and the failing file we want to minimize.

403 Let us justify the extra information that our Coq wrapper programs log. We need `COQPATH`
404 to ensure that we have the right search path for the dependencies of `coqc`, the command-line
405 Coq compiler. We need the command-line arguments so that we know what flags to tell the
406 bug minimizer to pass to `coqc`. Note that we *must not* change relative paths to absolute
407 ones when passing arguments along to `coqc`, because the output of `coqc` is sensitive to the
408 difference between relative and absolute paths, so changes can muddle tests that are meant
409 to produce output files (and did in the past, for example with `ci-elpi`). We can locate the
410 error message by looking for the last instance of `File "f", line ℓ , characters $n-m$:`
411 followed immediately by a line beginning with `Error`. (Note that warning messages also emit
412 the `File ...` line, but we do not want to catch warnings.) We look for the last instance of
413 the wrapper debug printout information that points at the same file, though, so long as we
414 were careful always to build single-threadedly, we could instead just look for the most recent
415 debug printout before the error message.

416 Given this information, we adjust the arguments so that we can tell the bug minimizer
417 where the dependencies live both for the passing and failing versions of Coq. We then pass
418 this information to the bug minimizer:

- 419 ■ the location of the file to be minimized;
- 420 ■ the log file containing the error message, which must match the error message that the
421 minimizer believes the file produces;
- 422 ■ the locations of the `coqc`, `coqtop`, and `coq_makefile` programs for the tip of the PR;
- 423 ■ the location of the `coqc` program for the master branch;

- 424 ■ the locations of the dependencies for both the passing and failing versions of Coq, parsed
- 425 from the command-line arguments and from walking the directories in `COQPATH`;
- 426 ■ any arguments to `coqc` that are neither naming dependency locations nor known to be
- 427 both irrelevant to the processing of the file and counterproductive to the minimizer's
- 428 operation (such arguments are `-batch`, which applies only to `coqtop`; `-time`, which will
- 429 only make logs of the minimizer much longer; and `-noglob`, `-dump-glob`, and `-o`, which
- 430 interfere with the generation of outputs used by the minimizer).

431 **8 An Alternative Usage Mode**

432 Up to this point, we have talked about using the Coq Bug Minimizer exclusively to minimize
 433 reverse-CI failures for debugging faulty changes in Coq. Our tool can also be used to minimize
 434 test cases for newly found bugs in Coq. In this mode, a bug reporter can write a shell script
 435 that invokes a *single* version of Coq to produce buggy behavior on some Coq file, asking
 436 `coqbot` to produce a minimal example from this script. When running in this mode, we place
 437 an additional constraint on the minimizer that the proof script generating the error message
 438 should be left untouched, which allows bug reporters to write proof scripts such as

```
439 some_tactic; lazymatch goal with
440 | buggy_goal => fail 0 "bug remains"
441 | [ |- ?G ] => fail 0 "bug disappeared!" G end.
```

442 to customize the desired reproducing case, trusting that the entire file will not be minimized
 443 to something silly like `Goal False. fail 0 "bug remains"`.

444 **9 Integration in Coq's CI and Evaluation of Results**

445 **9.1 Triggering the Minimizer**

446 The Coq project uses a custom, multi-task bot to automate everyday tasks, including
 447 triggering CI and reporting its results to the GitHub repository [19]. We have extended
 448 this bot to automatically propose and manage the minimization of failing test cases. The
 449 bot posts a comment to propose to run minimization when a PR has passed Coq's internal
 450 test suite but has failures with external projects, and these external projects have built
 451 successfully on the base commit (on the master branch).

452 If someone answers with a comment to trigger minimization, then the bot prepares a
 453 branch with all the information needed by the minimizer and pushes this branch to an
 454 external repository dedicated to running the minimizer. This triggers a GitHub Action
 455 workflow which will proceed with the minimization process. GitHub Action jobs have a
 456 6-hour timeout, so by the limit, the bot answers back with the results of the minimization
 457 process. If the minimization was stopped because of the timeout, then the bot automatically
 458 restarts it by reusing the file obtained at the previous step.

459 **9.2 Research Questions**

460 To evaluate the usefulness of our bug minimizer, we investigate several research questions:

461 **RQ1:** How often does the minimizer successfully produce a reduced test case from the CI
 462 failures it was triggered on?

463 **RQ2:** How often is this reduced test case fully standalone (no dependencies other than Coq's
 464 standard library)?

23:12 Automatic Test-Case Reduction in Coq

465 **RQ3:** How long does it take to produce such reduced test cases?

466 **RQ4:** What is the size of the reduced cases?

467 **RQ5:** How long do the reduced cases take to run?

468 **RQ6:** What is the amount of code reduction?

469 9.3 Data Collection and Analysis

470 To support reproducing the results, we provide our data collection and analysis code (as a
471 Jupyter notebook) and our dataset (as a CSV file) in the supplementary materials.

472 We retrieve the runs of the bug minimizer by looking for PRs in the Coq GitHub repository
473 with the words “coqbot ci minimize”, and we fetch all comments from the bot (timestamp
474 and body text) from these PRs using GitHub’s GraphQL API. We exclude PRs opened
475 by the first author, as most of these PRs were for testing the minimizer integration and
476 debugging issues. When the minimizer is triggered, the bot answers with a comment “I have
477 initiated minimization . . .” or “I am now running minimization . . .”, providing the list of
478 projects on which it is being run. Then, when it finishes minimizing a project, it produces a
479 comment with the minimized file. This file starts with header comments containing useful
480 information about the minimization process. The comment may also contain “interrupted by
481 timeout, being automatically continued” if the minimization process timed out and has to
482 be restarted to go further, which the bot automatically does. We ignore these comments,
483 only looking for final reduction outputs. Finally, the bot posts a comment starting with
484 “Error: Could not minimize file” when it was not able to minimize the requested failure, for
485 instance, because it could not reproduce it or could not reproduce the successful run on the
486 base branch.

487 We match comments indicating the start of the minimization with comments indicating
488 the end of it, using these two comments to determine if the minimizer was able to produce a
489 reduced test case, find how long it took, and answer our other research questions. To avoid
490 double-counting multiple runs on the same CI failure, we only look at the first bug-minimizer
491 trigger on a given PR and a given project.

492 9.4 Results

493 9.4.1 RQ1: How often does the minimizer produce a reduced test 494 case?

495 Looking only at the first minimization runs for a given PR and project, we have identified
496 191 runs on 51 PRs (very often, several minimization runs are started in the same PR on
497 different projects). On these 191 runs, 75% succeeded in producing reduced test cases. We
498 count as failed runs the ones where the bot reported “Error: Could not minimize file”, the
499 ones where we could not find a comment marking the end of minimization, and the ones
500 where the bot answered with a minimized file but this file was not actually reduced from the
501 initial test case (which we can detect from the header comments).

502 There were 5 runs for which we found no comment marking the end of minimization.
503 By manually looking at them, we have determined that 4 out of 5 were caught in infinite
504 loops and had to be canceled manually. Loops can arise when the 6-hour timeout of the
505 minimization process is not enough to make any new progress and thus the minimization
506 gets stuck without ever reaching its end. Typical circumstances are when testing out a single
507 change takes over 20 seconds, since we only have enough time to compile a 20-second-long

508 file about a thousand times in six hours. The last case of our 5 seems to be coqbot having
509 failed to post the comment marking the end of the minimization process.

510 There were 19 runs that concluded with an explicit “Error: could not minimize file”
511 comment. These errors are often due to issues downloading CI artifacts (9 runs), for instance
512 because the corresponding base CI jobs have been skipped or the CI artifacts have expired.
513 Runs concluding with errors can also happen because of bugs in Coq or in the tested projects’
514 build infrastructure that prevent minimization. Virtually all these issues were reported, and
515 most of them are already fixed. For instance, the MetaCoq project alone was responsible for
516 5 failures because of issues in its build system.

517 Finally, there were 23 runs ending with comments reporting on supposedly minimized
518 files but where (from the header comments or their absence thereof) we can conclude that
519 the minimization process failed to start properly (e.g., because it could not reproduce the
520 error message). Most of these problems were related to error-message parsing, namespace
521 management, or similar issues that have been fixed by making the bug minimizer more robust
522 to them (see Section 4 to Section 6). A few of these issues have been noted but not yet fixed.
523 Finally, a few of these failed runs were due to the minimizer being misused or called on a
524 project that had failed for a reason that was unrelated to the PR.

525 The accompanying notebook contains specific comments for each of the failed runs.

526 **9.4.2 RQ2: How often is this reduced test case fully standalone?**

527 We consider that a reduced test case will be most useful if any dependency beyond Coq’s
528 standard library was successfully inlined, leaving it possible to run the reduced test case
529 without needing to import any additional dependency. As a result, it is more likely that the
530 test case can be added to Coq’s test suite.

531 To measure how often the reduced test case is standalone, we rely on the minimizer
532 recording when it failed to inline a dependency in the header comments of the minimized
533 file. This feature was only added recently, so we only perform this measurement on the
534 47 successful runs of the minimizer that had this information available. On these 47 runs,
535 there were only 5 failures to inline dependencies fully, i.e., the minimizer produced a fully
536 standalone file in 89% of the cases.

537 Looking at the 5 failures to inline dependencies, we observe several types of reasons.
538 One case was related to robustness to changing error messages, one case was related to a
539 build-system issue in the project being minimized, and 3 cases were due to a common issue
540 blocking attempts at all inlining methods. All of these issues have been fixed since then.

541 **9.4.3 RQ3: How long does it take to produce such reduced test cases?**

542 We compute the duration of minimization as the time delta between the start and the end
543 comments. This method overapproximates the actual time spent in the minimization process,
544 since it also includes time setting up a VM and possibly waiting in the queue for an available
545 runner. We can look at this duration for both successful and failed runs.

546 For failed runs, we observe that the average duration for the minimization to conclude is
547 5 minutes (306 seconds) and that the maximum duration is 15 minutes (890 seconds).

548 For successful runs, we observe more variety. The minimum duration is 4 minutes (232
549 seconds), the maximum duration is 20 hours (73072 seconds), and the average is 104 minutes
550 (6238 seconds). 50% of the successful runs finish in under 20 minutes (1218 seconds), and
551 80% finish in under 140 minutes (8396 seconds). This number is still reasonable compared to
552 the time that contributors routinely spend waiting for the results of Coq’s CI [18].

553 9.4.4 RQ4: What is the size of the reduced cases?

554 For the last three questions, we focus mainly on the 42 recent minimization runs that are
555 known to have produced standalone files.

556 The shorter the reduced test case, the more useful it is: it can help developers understand
557 the problem more quickly, and it makes it more likely that it will be added to Coq’s test
558 suite. Here again, there is some variety in the size of the reduced cases (counted in number
559 of lines). The average size is 270 lines, and the maximum size is 2648 lines. However, 25% of
560 the reduced cases are under 39 lines, 50% are under 114 lines, and 75% are under 262 lines.

561 Results on the full set of 144 successful minimization runs are of the same order of
562 magnitude, with an average at 367 lines and a maximum size of 3804 lines.

563 Developers have the option to perform additional minimization manually and restart the
564 automatic minimization process on their manually reduced cases, which can help obtain even
565 more reduced cases, but we have not evaluated this feature quantitatively.

566 9.4.5 RQ5: How long do the reduced cases take to run?

567 Following a recent addition, the minimizer has reported the expected `coqc` compile time as
568 part of the header comments in the minimized file. Our recent 42 standalone cases all had
569 this field available. We observe that the reduced cases take on average 1.25 seconds to run,
570 although 75% of them take under half-a-second, while the maximum time is 26.5 seconds.

571 9.4.6 RQ6: What is the amount of code reduction?

572 To compute how much code reduction there was, we use data that the minimizer records
573 about each minimization step (how many lines it started from and how many lines it ended up
574 with). These numbers go up at times because of the process of inlining external dependencies.
575 On the other hand, dependencies are only inlined if they could not simply be removed, so
576 these numbers do not include the size of the files that were previously imported but did not
577 need to be inlined during the minimization process.

578 We aggregate these numbers by simply taking the sum of the differences in line count
579 at the beginning and the end of each minimization step. We compute the amount of code
580 reduction by taking the ratio of the final size over the total test-case size, defined as being the
581 sum of the final size and the total number of removed lines. We obtain an average figure of
582 31%, which means that the final test-case size is on average one-third the size of the original
583 test (including the dependencies that actually matter for the test case).

584 If we compute the size difference only looking at the initial file we started from and
585 the final file we obtained, without accounting for the inlined dependencies, then we get an
586 average ratio of 50%, which means that the final file is on average half the size of the file we
587 started from. Note that because of dependency inlining, nothing prevents the reduced test
588 case from being longer than the file we started from, which does happen in 6 out of 42 cases.
589 If we look only at the 36 cases for which there was some code reduction, we get that the
590 average reduction is by a factor of 4 to 5. If we look only at the 6 cases for which there was
591 code expansion, we get that the average expansion is by a factor of 2.

592 9.5 Limitations of our Evaluation

593 Evaluating a bug minimizer for a proof assistant such as Coq is difficult because there is
594 no preexisting benchmark that it could be run on. In this paper, we have decided to take

595 advantage of the integration of our minimizer in the CI infrastructure of Coq to evaluate it
596 on real use cases where Coq developers have felt the need for it.

597 While we have taken steps to ensure that the evaluation is as unbiased as possible (such
598 as not using reruns of the minimization on the same project in the same PR), our evaluation
599 is still limited by our choice to use real use cases. In particular, it should be noted that our
600 evaluation results are not obtained on a fixed version of the minimizer. On the contrary,
601 the minimizer has evolved (and is still evolving) in reaction to the very same cases on
602 which we have evaluated it. Since we always account only for first runs, many cases where
603 the minimizer has been counted as failing have been eventually fixed and would result in
604 successful runs today. Subsequent runs on other projects or other PRs may have succeeded
605 thanks to earlier fixes.

606 Other limitations are that our computation of the minimization duration is an overap-
607 proximation that also includes the time for things such as setting up a VM to run the process,
608 and that our evaluation of several research questions is based only on a subset of recent
609 minimizer runs.

610 Due to all these limitations, our evaluation should only be understood as demonstrating
611 the feasibility of our approach and the usefulness of its application to the development of Coq.
612 However, it should not be understood as a basis that future versions of the minimizer, or
613 alternative minimizers, can compare to, since today's version would already obtain different
614 results if it were rerun on all these cases.

615 **10** Related Work

616 Our work is at the crossing of two research areas: research on debugging techniques, which
617 is a subdomain of software-engineering research, and research on proof assistants.

618 Debugging is a largely explored topic, but mostly with a focus on more mainstream and
619 less formal languages than Coq. In this research domain, test-case-reduction techniques
620 have been studied for standard programming languages and compilers [3]. There are two
621 types of approaches that have been proposed. First, there are generic approaches that are
622 supposed to work for any programming language, by using structure information on the
623 program being reduced. Examples include delta debugging [16] but also the generalized
624 tree-reduction algorithm [7] and the syntax-guided Perses tool [7]. These generic techniques
625 would not be likely to work well for Coq programs without careful adaptation, because many
626 Coq programs can be considered syntactically valid even if completely nonsensical. For
627 instance, we have already mentioned the issue with removing a `Qed` statement at the end of a
628 tactic-based proof. Despite breaking a semantic block of code, this change does not actually
629 produce a syntactically invalid Coq program.

630 Second, there are programming-language-specific approaches, which take advantage of
631 specific knowledge to make the test reduction more performant. Our own work is related to
632 this second category, where most tools focus on mainstream languages like C. Some are even
633 dedicated to reducing the output of specific test-generation frameworks such as Csmith [12].

634 However, work on generating many diverse test cases from nothing has complementary
635 value. Csmith [15] has an effective algorithm based on knowledge of C semantics, to provoke
636 undefined behavior. Techniques like equivalence modulo inputs [9] find compiler bugs via
637 differential testing, where a compiler is run on programs that are known to have the same
638 semantics. Perhaps this generative approach would also be useful for proof assistants,
639 composed fruitfully with test-case reduction as we have presented.

640 Finally, the literature has identified the issue of *slippage* in test reduction [4, 8], which is

641 when the initial and reduced cases produce different compiler bugs. This challenge was one
 642 of the main ones we had to account for in designing our bug minimizer (see Section 4).

643 Proof-assistant ecosystems were already no stranger to testing techniques. For instance,
 644 Isabelle/HOL’s Nitpick [1] uses Boolean satisfiability to find theorem counterexamples.
 645 QuickChick [11] does random test generation to try to falsify Coq theorems. These tools are
 646 handy to save users from investing time in trying to prove false theorems. Testing-based
 647 approaches to debugging *proof assistants themselves* are a complementary topic.

648 **11** Future Work

649 We were pleasantly surprised to find that several “shortcuts” in the logistics behind the
 650 minimizer led to good results empirically, but some of these may be worth revisiting to
 651 improve results even more. In various places, we use workarounds (like `.glob` files) to avoid
 652 integrating a proper Coq parser, but there would be advantages like being able to remove
 653 specific fields from record types. We remove single commands at a time, rather than removing
 654 entire well-balanced command blocks, which probably costs us in minimization time.

655 A broader opportunity is finding related groups of commands that need to be removed
 656 together, to avoid changing the error message. For instance, we might want to move a lemma
 657 out of a module, to the top level of a file. Removing the commands that open and close
 658 the module might suffice, even if removing either one alone disturbs the error message. A
 659 general-enough version of this process could replace many specific passes.

660 One remaining aggravation is proper handling of lemma proofs within sections, where the
 661 details of the lemma proof influence which section variables are kept in the lemma’s type.
 662 We could use the `Set Suggest Proof Using` command to insert `Proof using` clauses.

663 As mentioned in Subsection 6.2, we would like to improve the ability of the minimizer
 664 to linearize dependency trees and handle Coq’s global state. We could, for example, print
 665 out the full table of flag settings at a particular point, reset them to the initial values before
 666 inlining a file, and then restore them after inlining. To fully handle global state, we would
 667 need some way to reconstruct the command-line flags used to compile installed files.

668 There are further-out ideas that could speed minimization significantly but might require
 669 significant modifications to Coq itself. Incremental compilation would be helpful, to save us
 670 from rerunning long proof scripts every time we change single lines below them. Minimizing
 671 multiple files in parallel, rather than only inlining files, would allow us to take advantage of
 672 multicore execution within single minimization jobs.

673 References

- 674 1 Jasmin Christian Blanchette and Tobias Nipkow. Nitpick: A counterexample generator for
 675 higher-order logic based on a relational model finder. In Matt Kaufmann and Lawrence C.
 676 Paulson, editors, *Interactive Theorem Proving*, pages 131–146, Berlin, Heidelberg, 2010.
 677 Springer Berlin Heidelberg. doi:10.1007/978-3-642-14052-5_11.
- 678 2 Martin Burger, Karsten Lehmann, and Andreas Zeller. Automated debugging in Eclipse. In
 679 *Companion to the 20th Annual ACM SIGPLAN Conference on Object-Oriented Programming,*
 680 *Systems, Languages, and Applications*, OOPSLA ’05, pages 184–185, New York, NY, USA,
 681 2005. Association for Computing Machinery. doi:10.1145/1094855.1094926.
- 682 3 Junjie Chen, Jibesh Patra, Michael Pradel, Yingfei Xiong, Hongyu Zhang, Dan Hao, and
 683 Lu Zhang. A survey of compiler testing. *ACM Comput. Surv.*, 53(1), February 2020. doi:
 684 10.1145/3363562.
- 685 4 Yang Chen, Alex Groce, Chaoqiang Zhang, Weng-Keen Wong, Xiaoli Fern, Eric Eide, and
 686 John Regehr. Taming compiler fuzzers. In *Proceedings of the 34th ACM SIGPLAN Conference*

- 687 on *Programming Language Design and Implementation*, PLDI '13, pages 197–208, New York,
688 NY, USA, 2013. Association for Computing Machinery. doi:10.1145/2491956.2462173.
- 689 5 Holger Cleve and Andreas Zeller. Finding failure causes through automated testing. In Mireille
690 Ducassé, editor, *Proceedings of the Fourth International Workshop on Automated Debugging*,
691 *AADEBUG 2000, Munich, Germany, August 28-30th, 2000*, 2000. arXiv:cs/0012009.
- 692 6 Jason Gross. Coq bug minimizer, January 2015. Presented at The First International
693 Workshop on Coq for PL (CoqPL'15). URL: [https://jasongross.github.io/papers/
694 2015-coq-bug-minimizer.pdf](https://jasongross.github.io/papers/2015-coq-bug-minimizer.pdf).
- 695 7 Satia Herfert, Jibesh Patra, and Michael Pradel. Automatically reducing tree-structured
696 test inputs. In *Proceedings of the 32nd IEEE/ACM International Conference on Automated
697 Software Engineering*, ASE 2017, pages 861–871, Urbana-Champaign, IL, USA, 2017. IEEE
698 Press. doi:10.1109/ase.2017.8115697.
- 699 8 Josie Holmes, Alex Groce, and Mohammad Amin Alipour. Mitigating (and exploiting) test
700 reduction slippage. In *Proceedings of the 7th International Workshop on Automating Test
701 Case Design, Selection, and Evaluation*, A-TEST 2016, pages 66–69, New York, NY, USA,
702 2016. Association for Computing Machinery. doi:10.1145/2994291.2994301.
- 703 9 Vu Le, Mehrdad Afshari, and Zhendong Su. Compiler validation via equivalence modulo
704 inputs. In *Proceedings of the 35th ACM SIGPLAN Conference on Programming Language
705 Design and Implementation*, PLDI '14, pages 216–226, New York, NY, USA, 2014. Association
706 for Computing Machinery. doi:10.1145/2594291.2594334.
- 707 10 Lina Ochoa, Thomas Degueule, and Jean-Rémy Falleri. BreakBot: Analyzing the impact of
708 breaking changes to assist library evolution. In *44th IEEE/ACM International Conference on
709 Software Engineering: New Ideas and Emerging Results, ICSE (NIER) 2022*. IEEE, 2022.
- 710 11 Zoe Paraskevopoulou, Cătălin Hrițcu, Maxime Dénès, Leonidas Lampropoulos, and Ben-
711 jamin C. Pierce. Foundational property-based testing. In *ITP 2015 - 6th conference on
712 Interactive Theorem Proving*, volume 9236 of *Lecture Notes in Computer Science*, Nan-
713 jing, China, August 2015. Springer. URL: <https://hal.inria.fr/hal-01162898>, doi:
714 10.1007/978-3-319-22102-1_22.
- 715 12 John Regehr, Yang Chen, Pascal Cuoq, Eric Eide, Chucky Ellison, and Xuejun Yang. Test-case
716 reduction for C compiler bugs. In *Proceedings of the 33rd ACM SIGPLAN Conference on
717 Programming Language Design and Implementation*, PLDI '12, pages 335–346, New York, NY,
718 USA, 2012. Association for Computing Machinery. doi:10.1145/2254064.2254104.
- 719 13 Matthieu Sozeau, Simon Boulrier, Yannick Forster, Nicolas Tabareau, and Théo Winterhalter.
720 Coq Coq Correct! verification of type checking and erasure for Coq, in Coq. *Proc. ACM
721 Program. Lang.*, 4(POPL), December 2019. doi:10.1145/3371076.
- 722 14 Daniel S. Wilkerson and Scott McPeak. delta - delta assists you in minimizing “interesting”
723 files subject to a test of their interestingness, February 2006. Presented at CodeCon 2006.
724 URL: <https://github.com/dsw/delta>.
- 725 15 Xuejun Yang, Yang Chen, Eric Eide, and John Regehr. Finding and understanding bugs in C
726 compilers. In *Proceedings of the 32nd ACM SIGPLAN Conference on Programming Language
727 Design and Implementation*, PLDI '11, pages 283–294, New York, NY, USA, 2011. Association
728 for Computing Machinery. doi:10.1145/1993498.1993532.
- 729 16 Andreas Zeller. Isolating cause-effect chains from computer programs. In *Proceedings of
730 the 10th ACM SIGSOFT Symposium on Foundations of Software Engineering*, SIGSOFT
731 '02/FSE-10, pages 1–10, New York, NY, USA, 2002. Association for Computing Machinery.
732 doi:10.1145/587051.587053.
- 733 17 Andreas Zeller. *Why Programs Fail: A Guide to Systematic Debugging*. Elsevier, 2009.
- 734 18 Théo Zimmermann. *Challenges in the collaborative evolution of a proof language and its
735 ecosystem*. PhD thesis, Université de Paris, 2019. URL: <https://hal.inria.fr/tel-02451322>.
- 736 19 Théo Zimmermann, Julien Coolen, Jason Gross, Pierre-Marie Pédro, and Gaëtan Gilbert.
737 Extending the team with a project-specific bot. working paper, December 2021. URL:
738 <https://hal.inria.fr/hal-03479327>.