

1 Proof Repair Infrastructure for Supervised Models: 2 Building a Large Proof Repair Dataset

3 Tom Reichel ✉

4 University of Illinois Urbana-Champaign, USA

5 R. Wesley Henderson ✉

6 Radiance Technologies, Inc., Huntsville, AL, USA

7 Andrew Touchet ✉

8 Radiance Technologies, Inc., Huntsville, AL, USA

9 Andrew Gardner* ✉

10 Radiance Technologies, Inc., Huntsville, AL, USA

11 Talia Ringer* ✉

12 University of Illinois Urbana-Champaign, USA

13 — Abstract —

14 We introduce a new, large proof-repair dataset and benchmark suite for the Coq proof assistant. The
15 dataset is made up of Git commits from dozens of open-source projects with old and new versions of
16 definitions and proofs aligned across commits. Building this dataset was a significant undertaking,
17 highlighting a number of challenges and gaps in existing infrastructure. We discuss these challenges
18 and gaps, and we provide recommendations for how the proof assistant community can address them.
19 Our hope is to make it easier to build datasets and benchmark suites so that machine-learning tools
20 for proofs will move to target the tasks that matter most and do so equitably across proof assistants.

21 **2012 ACM Subject Classification** Computing methodologies → Machine learning; Software and its
22 engineering → Software maintenance tools; Security and privacy → Logic and verification

23 **Keywords and phrases** proof repair, datasets, benchmarks, machine learning, formal proof

24 **Digital Object Identifier** 10.4230/LIPIcs.CVIT.2016.23

25 **Supplementary Material** Further dataset details and a sample can be found as follows:

26 *Appendix:* https://tomreichel.web.illinois.edu/Repair_Data_Appendix.pdf

27 *Dataset (Sample):* https://tomreichel.web.illinois.edu/Repair_Data_Sample.zip

28 **Funding** This research was developed with funding from the Defense Advanced Research Projects
29 Agency. The views, opinions, and/or findings expressed are those of the author(s) and should not be
30 interpreted as representing the official views or policies of the Department of Defense or the U.S.
31 Government.

32 **1** Introduction

33 Machine learning (ML) is coming for proofs. Recent years have seen a surge in interest in ML
34 for proofs—one reflected by the many recent research venues [7, 4, 12], papers [26, 45, 37],
35 tools [2, 13, 35], industrial research groups [47, 5], and funding opportunities [3, 6] centering
36 on or prominently featuring ML for proofs. The surge in interest blurs the line between
37 proofs and data so that any proof development, once released, may itself become data to
38 improve proof automation for future proof developments.

* Co-senior authors

Distribution Statement A (Approved for Public Release, Distribution Unlimited).



© Tom Reichel, R. Wesley Henderson, Andrew Touchet, Andrew Gardner, and Talia Ringer;
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. 23:1–23:20



Leibniz International Proceedings in Informatics

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

23:2 Proof Repair Infrastructure for Supervised Models

```
Lemma proc_rspec_crash_refines_op T (p : proc C_0p T)
  (rec : proc C_0p unit) spec (op : A_0p T) :
  (forall sA sC,
-   absr sA sC tt -> proc_rspec c_sem p rec (refine_spec spec sA)) ->
-   (forall sA sC, absr sA sC tt -> (spec sA).(pre)) ->
+   absr sA (Val sC tt) -> proc_rspec c_sem p rec (refine_spec spec sA)) ->
+   (forall sA sC, absr sA (Val sC tt) -> (spec sA).(pre)) ->
  (forall sA sC sA' v,
-   absr sA' sC tt ->
+   absr sA' (Val sC tt) ->
    (spec sA).(post) sA' v -> (op_spec a_sem op sA).(post) sA' v ->
  (forall sA sC sA' v,
-   absr sA sC tt ->
+   absr sA (Val sC tt) ->
    (spec sA).(alternate) sA' v -> (op_spec a_sem op sA).(alternate) sA' v ->
  crash_refines absr c_sem p rec (a_sem.(step) op)
  (a_sem.(crash_step) + (a_sem.(step) op;; a_sem.(crash_step))).
```

■ **Figure 1** Changes made to a lemma by a participant in a recent user study of proof engineers, from the REPLICA user study paper [52].

39 We in the proof engineering community have agency in how this surge of interest plays out.
40 We can develop datasets and benchmark suites that steer the ML community toward the tasks
41 that matter most. We can build infrastructure that makes it easy to develop those datasets
42 and benchmark suites, or to work on those tasks. And we can build evaluation methodologies
43 that measure success on those tasks in ways that truly matter, so that state-of-the-art results
44 on benchmarks will transfer smoothly to real-world improvements in proof automation.

45 This paper takes a step in that direction. In particular, it presents Proof Repair Infras-
46 tructure for Supervised Models (PRISM)—a dataset and benchmark suite for an important
47 proof automation task in Coq: *proof repair* [49], which comprises the automatic correction of
48 proofs in response to breaking changes in programs or specifications. Proof repair is crucially
49 important for reducing costs in large proof developments and for enabling the application of
50 formal methods to broader and more diverse contexts. Unfortunately, data for proof repair
51 is scarce and challenging to collect [52]. This paper highlights the challenges involved in
52 collecting repair data with an emphasis on how the proof engineering community can adapt
53 to those challenges as ML becomes increasingly relevant. Its contributions are the following:

- 54 1. a Coq proof repair dataset and benchmark suite accessible to ML experts (Section 3),
- 55 2. reusable tools for building and extracting information from Coq projects (Section 4), and
- 56 3. a discussion of the challenges we encountered and how to overcome them (Section 5).

57 Our overarching goal is to build the infrastructure and proof assistant community support
58 we need to steer the ML community toward the tasks that matter most (Section 2).

59 2 Overview

60 The necessity and difficulty of proof maintenance has been borne out empirically. A recent
61 user study of eight intermediate through expert proof engineers showed that maintenance
62 happened constantly for participating proof engineers during proof development [52] and
63 that even experts sometimes gave up in the face of change [49].

64 Consider, for example, the change in Figure 1, in which a user study participant updated a
65 lemma statement in response to a change in a dependency. As noted in the user study paper,

66 this was part of a larger change, with 10 other definitions or lemma statements changing in
67 analogous ways. Furthermore, this change broke at least five proofs, four of which the user
68 study participant—an expert proof engineer—admitted or aborted rather than repair.

69 The ubiquity of maintenance and the challenges of repair have been largely neglected
70 in ML tools for proofs. ML tools for proofs have instead historically fixated on tasks like
71 predicting tactics or automatically formalizing natural language [13, 37, 45]. The lack of a
72 good dataset and benchmark suite obstructs progress; what currently exists is not sufficient
73 for training and evaluating proof repair models (see Section 6). If the datasets and benchmark
74 suites are not fit for maintenance tasks, the ML community may neglect those tasks entirely,
75 instead chasing state-of-the-art results on tasks for which existing benchmark suites suffice.

76 Our experiences interacting with ML experts and building datasets ourselves suggest that
77 the choice of datasets and benchmark suites for a domain is not driven solely by what is
78 likely to be useful—it is also driven by barriers imposed by infrastructure, lack of domain
79 expertise, or social factors. Things we may take for granted, like parsing and checking proofs,
80 can become prohibitive infrastructure challenges for ML experts.

81 In this paper, we overcome a number of these challenges and build a large proof repair
82 dataset for Coq. We also discuss the barriers we do encounter, and we describe both how we
83 overcome those barriers, and what kind of work the proof assistant community would need
84 to put in to make it so that they cease to be barriers at all.

85 We find this work especially prudent given that the danger of chasing benchmarks that
86 may not transfer to real life workflows has been realized quite dramatically in other domains,
87 from incorrect patches to programs [46] to incorrect clinical interpretation of x-ray results [67].
88 Furthermore, these challenges can influence not just the tasks that the ML community chooses
89 to tackle but even the very proof assistants for which the ML community chooses to build
90 supporting tools. It is in the community’s best interest to drive strong, practical results for
91 useful tasks in a way that is equitable across proof assistants.

92 Our hopes are twofold. First, we hope that our dataset will be immediately useful for
93 proof repair. Second, and perhaps more prudently, we hope our discussion of the challenges
94 involved in building it will serve as a call for action to improve infrastructure. The proof
95 engineering community can then ensure that ML experts focus on the automation tasks that
96 matter most to the community, that they measure success on those tasks in ways that transfer
97 smoothly to real-world usefulness, and that they do so equitably across proof assistants.

98 **3 A Proof Repair Dataset**

99 The PRISM dataset and benchmark suite that we have assembled will be publicly available.¹
100 The task that PRISM focuses on is proof repair (Section 3.1). The data comprise aligned
101 Git commits that correspond to existing changes in proof developments found on GitHub
102 (Section 3.2). Success on the resulting benchmark is evaluated in terms of successful proof
103 checking for repaired proofs (Section 3.3).

104 **3.1 The Task: Proof Repair**

105 In ML, a *task* refers to a high-level input/output specification of what is being learned.
106 A dataset and benchmark suite typically organizes itself around a particular task while

¹ We plan to release the dataset before the camera ready. The team includes researchers subject to US government approval processes, and cannot release datasets without a month’s notice. We have included a sample as supplementary material, which we were able to approve a month in advance.

23:4 Proof Repair Infrastructure for Supervised Models

107 remaining agnostic to the details of the model implementation.

108 We define proof repair as an ML task with **inputs** comprising an old theorem, proof of
109 the old theorem, and a new theorem; and **outputs** comprising a proof of the new theorem.
110 Note that this particular task assumes that we already know how to repair the theorem
111 statement and its dependencies. We could also consider a second task that allows the model
112 to repair the specification itself. This second repair task would be harder to evaluate, so we
113 do not focus our benchmark suite on it at this time, even though PRISM supports it.

114 **Inputs** We aim to provide sufficient context in the data to support a wide range of ML
115 approaches. At a high level, the input to the ML repair model is the entire state of a
116 project where the approach dictates how much of this state (and in what form) actually
117 reaches the model. More precisely, the input comprises the statement of the theorem whose
118 proof should be repaired, any contextual definitions on which it depends, the step-by-step
119 goals and hypotheses for each sentence in the old proof, and known changes to the project
120 up to and including changes in the theorem statement and its dependencies. For each of
121 these components, we supply raw text representations, abstract syntax trees (ASTs), and
122 identifiers. In the case of goals and hypotheses, serialized Coq kernel representations supply
123 detailed internal proof states. Environmental dependencies such as Coq compiler versions are
124 captured for errors induced by external application programming interface (API) changes.

125 **Outputs** The ultimate output is the repaired proof, which takes the form of text that
126 may be generated one sentence at a time, all at once, or through targeted modifications of
127 existing sentences. In the case of supervised repair learning, we supply ground truth targets
128 in the same form as the inputs: raw text, ASTs, etc. for the entire repaired project’s state
129 (compactly represented as a ‘diff’ relative to the input). Sufficient context is provided in the
130 data to programmatically execute up to the error in an interactive REPL such as `coqtop` or
131 `sertop`, where one may apply reinforcement learning akin to CoqGym [65].

132 3.2 The Data: Aligned Git Commits

133 The training and testing data comprise aligned Git commits for a selection of realistic Coq
134 projects. The initial release of PRISM will span 60 Coq projects and roughly 20K commits.
135 A full summary of all projects in the dataset can be found in supplementary material.

136 Projects were originally selected by querying GitHub’s API for projects that contained
137 Coq source code, had a file called “Makefile” in the project’s root directory, and had at least
138 100 commits from which to mine repair data. Eventually, we also included projects from
139 CoqGym [65] and filtered to projects that were listed in OCaml Package Manager (opam)
140 repositories (opam acts as the primary distributor of Coq projects). We excluded projects
141 that did not contain any proofs, or that had ulterior motives in their builds (e.g., projects
142 that intended to test the performance of the Coq compiler `coqc`). We hope in the future to
143 include additional projects, though this will require us to support more build environments
144 and expand upon the work detailed in Section 4.2.

145 **Repair Examples** Within each Coq project, the data comprises a number of repair examples—
146 that is, changes to definitions or proofs. A repair example is constructed by comparing
147 a definition or proof before and after a change. Since sentences and files may be moved,
148 renamed, added, deleted, or otherwise altered between commits, they must first be *aligned*
149 to ensure the right changes are compared. This means that Vernacular commands in one
150 commit are assigned one-by-one to commands in another commit, where these assignments

151 may cross file boundaries. Note that each command may not get a partner, indicating that
152 it was either added or deleted. We describe this in more detail in Section 4.5.

153 After alignment, proof repair examples are constructed by partially applying changes,
154 e.g., by omitting the changes to a proof that accompanied a change to the proposition.
155 Thus, one pair of commits may give rise to multiple examples. The examples are compactly
156 represented by commit hashes and diffs that indicate the state before and after a repair.
157 This compact representation enables dissemination of the dataset without the accompanying
158 projects, although we note supplementary tools for efficiently extracting project data will
159 still be vitally important for eliminating redundant computations and effort in practice.

160 **Data Split** ML datasets and benchmark suites often include a data split between training,
161 validation, and testing data. We do not commit to a split ahead of time, but we consider two
162 different ways of splitting data: across projects and chronologically within projects. These
163 two splits test two different kinds of generalization beyond the training data. We plan to
164 include defaults for both splits in the final release of PRISM.

165 The first split—across projects—chooses distinct sets of projects to use for training,
166 validation, and test data. This approach measures generalization of the learned model to new
167 projects not seen at training time. The second split—chronologically within projects—uses
168 the same set of projects for training and test data, but splits them by time of commit, so
169 that training data includes earlier commits, and testing data includes later commits for the
170 same projects. This approach measures generalization of the learned model to new changes
171 within a given project, when the model was trained on older data for that project.

172 3.3 The Metrics: Proof Checking

173 Changes in proof developments that break proofs can be fixed in two ways: by repairing the
174 proofs themselves or by repairing some other definition such as a program or specification [49].
175 PRISM includes both kinds of changes. We focus our benchmark suite on the former
176 (repairing proofs), as the metric for success is immediately clear. We hope our benchmark
177 suite will also be useful for the latter (repairing definitions), but we believe the problem of
178 choosing a good metric for success for repairing definitions to be an open research problem.

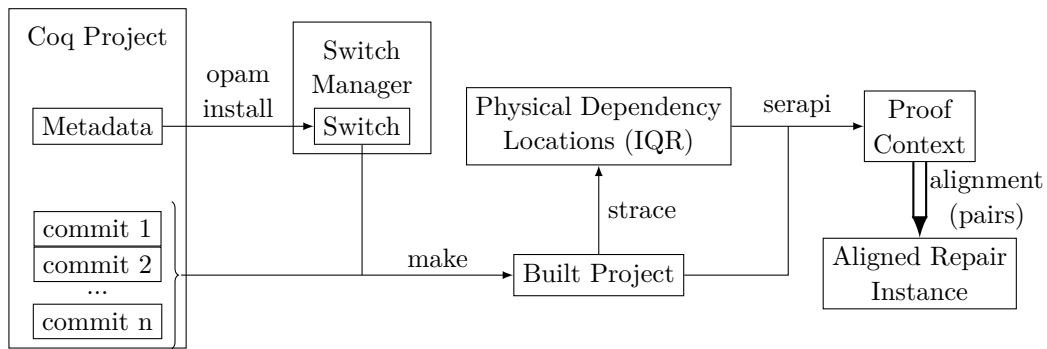
179 **Repairing Proofs** We focus our benchmarks on the problem of repairing a proof script
180 assuming that the statement of the repaired theorem is already known. In this case, checking
181 the correctness of the repaired proof amounts to using Coq’s kernel to proof check the type of
182 the repaired proof against the type that represents the desired repaired theorem statement.

183 The proof checking metric is the same as that used for the standard CoqGym [65] proof
184 generation benchmark suite for Coq. This metric is sound and complete (up to the correctness
185 of Coq’s kernel with nonterminating proof scripts designated as incorrect): any proof that
186 checks with the desired type is a proof of the theorem the type encodes (**soundness**), and
187 all proofs that prove that theorem will check with the desired type (**completeness**).

188 For this flavor of proof repair, we are able to take advantage of the fact that proof checking
189 is a perfect oracle when the theorem statement is known. Perfect oracles have been hugely
190 beneficial for existing ML work for proof generation [26] and for early symbolic work on proof
191 repair when specifications do not change [53]. They continue to benefit ML for proof repair.

192 **Repairing Definitions** Helping users fix the definitions that the specification depends on—or
193 the theorem statement itself—is also desirable. The REPLICIA user study, for example, found
194 that 75% of the time proof engineers fixed a broken proof, they did so by fixing something

23:6 Proof Repair Infrastructure for Supervised Models



■ **Figure 2** Process of extraction for a Coq project commit.

195 else, like a program or specification [52]. Supporting this use case may actually be *more*
196 helpful to proof engineers than supporting the original flavor of proof repair. Unfortunately,
197 existing metrics are insufficient for measuring success on this task:

- 198 ■ The **proof checking** metric is insufficient when the repaired specification is unknown;
199 showing that a proof type checks is not meaningful unless we know its intended type.
- 200 ■ The metric of **exact equality** with an expected repaired definition is too conservative, as
201 there are many equivalent ways to state the same theorems or write the same definitions.
- 202 ■ Common notions of **definitional** or **propositional equality** in Coq are less conservative,
203 but are still too far from complete.
- 204 ■ PUMPKIN PI [51] repairs definitions to be equivalent up to **univalent transport**, but
205 checking this automatically is undecidable.
- 206 ■ Common natural language **distance metrics** like BLEU [42] are poor measures of success
207 in code tasks [24], so we expect them to be inadequate for proofs.

208 If you choose to evaluate a model for repairing definitions, we recommend a conservative
209 metric like exact or definitional equality to avoid the danger of chasing misleading benchmarks.
210 We hope to eventually develop a suitable less conservative metric, particularly one that
211 captures what makes a change “close to correct” for some suitable notion of correctness.

212 4 Building the Proof Repair Dataset

213 We now take a step back and describe the processes behind our data collection efforts.
214 As stated, the foundation of the dataset comprises open-source Coq projects. Mining the
215 commits of these projects eventually yields examples of refactors or repairs. Each project is
216 accompanied by per-commit metadata containing project dependencies, source URL, and
217 build commands that is in parts manually curated and programmatically inferred. The
218 process of generating repair data from a project comprises the following steps (see Figure 2):

- 219 ■ We obtain a *switch* (opam virtual environment) that satisfies as many of the project’s
220 dependencies as possible using the **Switch Manager (SwiM)** (Section 4.1).
- 221 ■ Once we have a switch, we run the build command in the generated switch to produce a
222 **Built Project** (Section 4.2).
- 223 ■ We **strace** the build process to scrape the **Physical Dependency Locations (IQR**
224 **flags)** of each document in the project (Section 4.3).
- 225 ■ Using the IQR flags for each document in a built project along with the Coq serializer
226 SerAPI [8], we extract a **Proof Context** corresponding to each intermediate proof state
227 for the project by querying Coq’s state during the execution of proofs (Section 4.4).

228 ■ Finally, we align changed proofs across commits and save those along with their interme-
229 mediate proof contexts to arrive at an **Aligned Repair Instance** (Section 4.5).
230 Building this infrastructure was a significant undertaking with many challenges encountered
231 along the way; we discuss these challenges in Section 5. Our hope is that the infrastructure
232 we have built will make it easier to collect similar datasets in the future.

233 4.1 The Switch Manager

234 To extract information about projects like intermediate proof states, we must build them.
235 This requirement is nontrivial because different projects can depend on different versions of
236 Coq, the Ocaml compiler, or other dependencies.

237 To resolve dependencies and make it possible to build many different commits of one or
238 more projects, we introduce a novel SwiM capability that works in tandem with opam to
239 model the build environment for a given commit of a given project as a Python object. In
240 particular, the object models an opam switch, which is opam’s representation of an isolated
241 collection of installed packages.

242 This capability subverts the typical manual opam workflow to create and activate a
243 switch. This manual workflow would be intractable at the scale of hundreds of commits for
244 each of dozens of projects. With the SwiM, we can automate this functionality and extract a
245 dataset at scale.

246 Several benefits arise from the SwiM’s design. The SwiM enables build sandboxing by
247 providing switch clones that last for just the duration of the commit’s extraction, and it also
248 minimizes the time needed to obtain a clone by maintaining a pool of switches across all
249 threads with pre-installed packages, and choosing the one upon request that is closest to
250 satisfying a commit’s requirements. Implementation of this capability required reflection of
251 opam’s dependency formula parsing and evaluation logic from OCaml to Python.

252 As new commits are built, switches containing their dependencies are added to the
253 managed pool of switches. Since switches range in size from hundreds of megabytes to a few
254 gigabytes, a least-recently-used cache maintains the total disk consumption below an implicit
255 limit by deleting stale, infrequently used switches.

256 4.2 Built Projects

257 Using the switch provided by the SwiM and a build command from the metadata, we may
258 be able to build the project. Confounding issues that may prevent building include undefined
259 opam variables within dependency formulas. In practice, we have so far seen a build failure
260 rate of about 68%. We attempt to build each commit with seven different major versions of
261 Coq ranging from 8.9 to 8.15 corresponding to versions of SerAPI that support capabilities
262 we deemed necessary. Since Coq releases are rarely backwards compatible, many of the build
263 failures can be explained by the fact that each commit can only be expected to build for the
264 single Coq version for which it was written. Furthermore, since we pin one of the seven Coq
265 versions in the switch supplied by the SwiM, conflicting version requirements may yield an
266 opam command that has no solution. Consequently, some build errors are inevitable.

267 However, other errors are due to mistakes or missing information in the human-sourced
268 metadata. We plan to address this latter class of build errors over time by fixing problems
269 in the metadata through automated inference mechanisms. If the project build fails, we
270 hope in the future to be able to recover proofs from the documents that built before the
271 failure as well as subsequent independent proofs. We are also exploring ways to automatically

272 recover from simple build errors such as dependency mismatches between the switch and the
 273 project’s requirements by using the date the commit was made as a version hint.

274 4.3 Physical Dependency Locations (IQR Flags)

275 In order to run any of the Coq or SerAPI tools (e.g., `coqc`, `coqtop`, `sertop`) on a given Coq
 276 source file, one or more flags regularly need to be passed to these commands to specify the
 277 physical location of dependencies. These flags are described below:

- 278 ■ The `-I` flag allows a directory to be added to the OCaml loadpath.
- 279 ■ The `-Q` flag adds a physical directory to the loadpath and binds it to a given logical path.
- 280 ■ The `-R` flag acts like the `-Q` flag, but also makes subdirectories available recursively.

281 In publicly available Coq projects, these flags (referred to as “IQR” flags from here on) are
 282 specified in one or more build or configuration files. No single standardized approach for
 283 specifying IQR flags exists, making it difficult to automatically infer them from configuration
 284 and build files alone. While projects will be able to build successfully without our knowledge
 285 of these flags, we must infer them to use SerAPI tools in other stages of our framework.

286 Our solution to this problem builds off an approach developed in IBM’s PyCoq [20].
 287 Following PyCoq, we use `strace` to inspect the actual commands run during the build
 288 process for a Coq project. Each build command is captured and any present IQR flags are
 289 extracted using regular expressions. In some projects, build files may be nested, and IQR
 290 flags may specify physical paths that are relative to the nested directories. We need to ensure
 291 that the inferred IQR flags are relative to the project root directory, so before we store the
 292 inferred IQR flags, their paths are resolved to the project root directory.

293 4.4 Proof Contexts

294 Once the project has been built, the individual Coq source files are parsed into sentences
 295 and then interactively executed with `sertop` to capture intermediate proof states.

296 A parser (`sercomp`) is available through SerAPI, but it only works on Coq source files
 297 whose dependencies are already compiled, which prohibits its use in recoveries from partial
 298 builds. Furthermore, `sercomp` introduces significant redundant computation with respect
 299 to `sertop`. As a more efficient alternative, we developed a simple regular-expression-based
 300 “heuristic parser” to perform sentence extraction and approximate proof identification.

301 From `sertop`, we can collect thorough context from the document, like whitespace-
 302 normalized text, ASTs, command types, and intermediate proof steps with goals and
 303 hypotheses. Each command is accompanied by inferred identifiers of the command itself (e.g.,
 304 an inductive type’s name and constructors) and a list of fully-qualified identifiers referenced
 305 within the command, which enables models to more easily incorporate local context or apply
 306 graph-based approaches. Accompanying source code locations allow for accurate provenance
 307 of data and application of proposed repairs to appropriate destinations for testing.

308 4.5 Aligned Repair Instances

309 The last step in our data collection process is extracting proof repair examples from different
 310 versions of projects, accounting for the fact that definitions and proofs may be changed,
 311 moved, renamed, or deleted between commits. We must establish a robust mapping between
 312 Vernacular commands in a pair of commits that preserves some notion of command identity.
 313 Our objective is similar to that of the ‘diff’ utility, which describes changes made between
 314 files. Where our objective differs is that we seek to match Vernacular commands rather than
 315 lines, and we seek to do so within the entire project directory structure rather than a file.

316 A traditional order-preserving alignment between two sequences, e.g., the Smith-Waterman
 317 algorithm [58], is not quite an appropriate approach to resolve this issue as it cannot correctly
 318 align two independent definitions whose order has been reversed during a refactor (perhaps
 319 due to an introduced dependency). Therefore, we approach the problem as a bipartite
 320 matching or *assignment* between the unordered elements of two sets such that the overall
 321 similarity of matched elements is maximized. We can formally specify the desired assignment
 322 between two commits X and Y considered as respective sets of commands across one or
 323 more files as the solution to the following optimization problem:

$$\begin{aligned} & \underset{U, W}{\text{minimize:}} \sum_{u \in U} C(u, f(u)) \\ & \text{subject to: } f : U \leftrightarrow W \wedge U \subseteq X \wedge W \subseteq Y \wedge \|U\| = \min(\|X\|, \|Y\|) \end{aligned} \quad (1)$$

325 Here, $C(x, y)$ is a non-negative cost function that measures the *distance* between the
 326 commands x and y , and f is a bijection between subsets of X and Y —a partial alignment
 327 between commands of X and Y . To align as many commands as possible, the domain of f
 328 must have at least as many members of the smaller of X and Y , which is our final constraint
 329 above. This optimization is an instance of the well-known *assignment problem*, which one
 330 can solve exactly in polynomial time, e.g., with the Hungarian algorithm [41].

331 The optimization is parameterized by the choice of C , for which we choose a normalized
 332 variant of the Levenshtein edit distance E :

$$333 \quad C(x, y) = \frac{2E(x, y)}{\|x\| + \|y\| + E(x, y)}, \quad (2)$$

334 where $\|x\|$ and $\|y\|$ give the character lengths of x and y considered as text (not including proof
 335 bodies). This normalization is an instance of the biotope transform [22], which preserves the
 336 metric properties (such as the triangle inequality) of the edit distance. We further threshold
 337 the distance by a constant t such that $C_t(x, t) = \min\{C(x, y), t\}$, which also preserves metric
 338 properties [43]. After solving for f , commands x and y assigned to one another ($f(x) = y$)
 339 with a cost of t are considered to be unassigned (i.e., we determine that x was dropped
 340 between commits and y was added). We choose $t = 0.4$, which roughly corresponds to 50%
 341 of a command’s text being changed before it is considered to have been dropped.

342 Solving this assignment problem for two entire commits can be costly: solving exactly is
 343 cubic complexity, and calculating the edit distance between all pairs of commands from both
 344 commits is necessarily quadratic complexity. Furthermore, the assignments produced may
 345 be somewhat spurious, especially in the event of multiple global optima. We mitigate these
 346 issues by applying the assignment problem only to those commands known to have changed
 347 in some manner between the commits according to their intersection with a (Git) ‘diff’. The
 348 final resolution of the problem is thus somewhere in between alignment and assignment.

349 Once we determine an alignment, we create examples of proof errors by leaving out
 350 changes to individual proofs one at a time, thus providing the context for each change to a
 351 proof that required repair but not the repair itself. The left-out change to the proof then
 352 accompanies the error as a ground truth target for supervised learning.

353 **5 Challenges**

354 Collecting and building datasets and benchmark suites for many tasks is still extremely
 355 challenging, and it is challenging in a way that is not at all equitable across proof assistants.
 356 Here, we discuss our experiences dealing with challenges encountered during the creation

23:10 Proof Repair Infrastructure for Supervised Models

357 of PRISM (Section 5.1), what we believe the Coq community can learn from other proof
358 assistant communities (Section 5.2), and how the proof assistant community at large could
359 address them more sustainably going forward (Section 5.3).

360 5.1 Our Experiences

361 The major challenges we faced in building this dataset and benchmark suite chiefly fall into two
362 categories: Project Management (Section 5.1.1) and Parsing & Serialization (Section 5.1.2).
363 For each of these categories, we discuss our experiences dealing with each stated challenge.

364 5.1.1 Project Management

365 One of the greatest barriers to building this dataset was the lack of a centralized archive for
366 Coq proof data. In the absence of this centralized archive, we resorted to looser collections
367 of projects organized by package management. The package manager `opam` gives us a
368 programmatic interface to build compatible environments for the dataset’s constituent
369 projects. However, it was designed to service individual developers using a few switches,
370 whereas we must spin up *dozens* of switches efficiently. We thus had to reimplement and
371 expand upon some of `opam`’s capabilities. We faced three challenges in so doing:

- 372 1. Significant **build system variation** across different proof developments;
- 373 2. **Expressive dependencies** in `opam` packages that complicate efficient installs;
- 374 3. Insufficient caching of `opam` build artifacts that necessitated **copying switches** to avoid
375 rebuilding the same packages.

376 **Build System Variation** Over the years, the recommended build system for Coq proof
377 developments has been in flux. In 2019, for example, the Coq development team urged proof
378 engineers to move their proof developments to Dune [19]. This effort did not fully succeed,
379 and the documentation for the latest Coq version includes instructions for both Dune and the
380 native Coq build system [19]. The native build system itself has also changed over time, losing
381 compatibility with its previous versions. Because of this fragmented build infrastructure, we
382 had to employ extremely abstract methods to extract arguments for SerAPI tools, namely
383 by using `strace` to grab IQR flags passed to Coq’s compiler `coqc` (described in Section 4.3)
384 while making almost no assumptions about the process invoking `coqc`.

385 **Expressive Dependencies** The `opam` package manager provides a powerful and expressive
386 syntax (package formulae) for packages to specify dependencies. Package formulae allow
387 developers to restrict the versions of dependencies that can be installed, to conjunct and
388 disjunct formulae into more complicated expressions, and to refer to variables declared
389 elsewhere in the environment. This feature benefits the library developer that can precisely
390 specify the environment for running code, but for our purposes it poses a challenge: packages
391 can be picky about their environments and force `opam` to rebuild existing libraries. Since we
392 need to install many versions of many packages, we need efficient ways to create or select
393 compatible switches, which means interpreting these formulae. As a result, we reimplemented
394 a majority of `opam`’s package formula features, including parsing the custom grammar for
395 package formulae and implementing package version comparison, to reason about which
396 existing switches would require the least time to install a given package with `opam`.

■ **Listing 1** A notation that breaks CoqIDE’s parser. This example was found in the Coq Discourse [17].

```
Notation "( a . b )" := (a, b).
Check (1 . 2).
```

397 **Copying Switches** To work with conflicting packages or different versions of the same
 398 packages, we must use different environments. The opam “switch” abstraction allows us to
 399 sandbox environments, but creating many switches incurs exorbitant overhead as each new
 400 switch rebuilds packages from source. Building a package *once* and deploying it in multiple
 401 switches is preferable, but many executables built by opam contain their absolute path as a
 402 hardcoded variable, which means they stop working if the name or location of the containing
 403 switch changes. That is, a built opam package only necessarily works in one switch. Our
 404 workaround is to copy switches and use the `bwrap` utility (which is also used internally in
 405 opam) to bind-mount the copied switch over the original such that the clone is in the original
 406 hardcoded location from the perspective of the running process. This solution allows copies of
 407 switches to act as if they are the original. Of course, handling these cloned switches requires
 408 extra bookkeeping and infrastructure, which the SwiM (Section 4.1) ultimately handles.

409 5.1.2 Parsing & Serialization

410 No matter how sophisticated the build system, we cannot get detailed data about individual
 411 proofs without parsing Coq files and serializing proof state to text. SerAPI [8] is the de facto
 412 standard for serializing Coq, providing a query protocol for exposing internal Coq data like
 413 definitions in the global environment, syntax trees, goals, types, and more. We used the
 414 CoqGym [65] Python wrapper as a starting point for our implementation, taking care to
 415 decouple it from CoqGym’s custom versions of Coq and SerAPI since we need to support
 416 multiple versions of each coinciding with chosen projects’ Git histories. This need to support
 417 multiple versions of Coq exacerbated challenges arising from gaps in SerAPI’s query protocol,
 418 requiring us to implement workarounds using the most public and arguably stable interface
 419 Coq possesses: its Vernacular query commands. We faced four challenges related to parsing
 420 & serialization:

- 421 1. Executing a file one Coq sentence at a time requires accurately **parsing sentence**
 422 **boundaries**, but parsing requires execution: a catch-22.
- 423 2. **Identifying dependencies between commands** (e.g., which lemmas a theorem uses)
 424 is critical to providing locally relevant repairs but is not a capability provided by SerAPI.
- 425 3. **Determining the scope of a conjecture** is complicated by the potential presence of
 426 nested proofs/definitions and arbitrary grammar extensions.
- 427 4. SerAPI is experimental software, which leads to breaking **changes between versions**.

428 **Parsing Sentence Boundaries** A Coq statement or ‘sentence’ ends with a period (`.`), but
 429 Coq also uses the symbol for import paths and module members so that one cannot identify
 430 sentences in a file merely by splitting on periods. To further complicate matters, Coq boasts
 431 an extensible syntax that enables users to define syntax that allows periods to show up in
 432 even more situations. For example, Listing 1 defines syntax using a period that complicates
 433 sentence splitting to the point where the latest version of CoqIDE—the official editor for
 434 Coq—cannot correctly parse and run this code even though Coq can. We did not discover any
 435 public or officially supported mechanism to extract the sentences of a Coq document, which

23:12 Proof Repair Infrastructure for Supervised Models

■ **Listing 2** A simple example showing that proofs may be interleaved and that multiple proofs (obligations) may be associated with one term.

```
Require Coq.Program.Tactics.
Set Nested Proofs Allowed.
Program Definition foo := let x := _ : unit in _ : x = tt.
Next Obligation. (* Start first obligation of foo *)
  Definition foobar : unit. (* Interject with new conjecture. *)
    exact tt.
  Next Obligation. (* Switch back to first obligation of foo *)
    exact tt.
  Qed. (* Finish proof of foo's first obligation *)
Defined. (* Finish proof of foobar *)
Next Obligation. (* Start next obligation of foo *)
  simpl; match goal with | ?a = _ => now destruct a end.
Qed. (* foo is defined *)
```

436 led us to develop the Python-based heuristic parser mentioned in Section 4.3 for simplicity
437 and maximal portability between build environments.

438 **Identifying Command Dependencies** Identifying dependencies decomposes into two sub-
439 problems: detecting the definitions (if any) introduced by a given command and resolving
440 referenced names unambiguously.

441 No SerAPI query resolves the first subproblem, nor is there any reliable syntactic clue
442 in the text that generalizes across unforeseen grammar extensions. Instead, we rely upon
443 parsing user-level feedback that notes the introduction of new identifiers (e.g., “*X* is defined”)
444 and Vernacular queries. Since feedback is not guaranteed for all definition types (particularly
445 propositions, depending on the Coq version), we also monitor for changes in the set of all
446 locally defined names yielded from Vernacular `Print All` command. One can thus reliably
447 identify a command with names introduced immediately after its execution.

448 The second subproblem arises from the fact that identifiers within ASTs yielded from
449 SerAPI are not necessarily fully qualified. Correcting this deficiency requires locating the
450 identifiers within the AST and issuing a Vernacular `Locate` query for each one. Care must
451 be taken to ensure that variables within local binders, patterns, or other sub-expressions
452 do not get mistaken as any top-level definition that they may shadow. Given the lack of
453 insight available into Coq’s internal name resolution, the accuracy is ultimately limited by
454 handcrafted scope rules. We also note one restriction on resolving globally bound identifiers:
455 if a definition shadows an existing one, then it cannot also use the shadowed one. Violation
456 of this assumption is possible (consider a recursive function `nat` that expects arguments of
457 type `nat`) but not expected to pose a significant risk as it is unlikely in the first place and
458 would generally be considered poor practice. If the restriction is violated, then the shadowed
459 definition will simply be mistaken for its shadower within the shadower’s definition.

460 **Determining Conjecture Scope** Determining conjecture scope decomposes into two subprob-
461 lems: attribution of proof steps to the correct conjecture and detection of proof (conjecture)
462 completion. Each is complicated by potentially intermingled or nested proof steps as shown
463 in Listing 2 and by the lack of a SerAPI query of the active conjecture’s identity.

464 A Vernacular command—`Show Conjectures`—again provides the solution. This com-
465 mand lists the names of currently stated but unproved conjectures and by all observations is

466 guaranteed to list the conjecture actively being proved first. We rely upon this presumed order
467 to identify the current conjecture, accumulating proof steps in stacks per open conjecture.
468 The method’s accuracy depends upon the assumption that each conjecture enters proof
469 mode once its first sentence is executed. The only known exceptions to this rule comprise
470 Programs, which do not enter proof mode until their first Obligation’s proof is begun.

471 Special handling is required to associate each Obligation with the correct Program since
472 `Show Conjectures` reveals a unique name for each Obligation. However, the special handling
473 means any grammar extension that defines its own Obligation or Program equivalents (e.g.,
474 multi-block proofs) cannot be serialized to the same level of accuracy. If any extension does
475 so, then each Obligation-equivalent is expected to be serialized as an unrelated theorem.

476 We rely upon detection of definitions to determine when and if a conjecture was proved,
477 assuming that no conjecture emits an identifier before it is defined (i.e., before it is proved).
478 Only subproofs (generally delimited by bullets and braces) are allowed to violate this rule.
479 However, one cannot assume that the first detected definition in the midst of a proof
480 corresponds to the conjecture, nor can one assume that the name of the conjecture once
481 defined will actually match its name as returned by `Show Conjectures`.

482 We ultimately detect the completion of a proof by requiring two conditions: a change in
483 the currently detected conjecture and the detection of a new definition. This rule necessarily
484 invokes an additional assumption: a change in the current conjecture implies that either a
485 new proof has begun or the current proof has ended (but not both). Since we assume that
486 conjectures cannot emit identifiers before they are done, we deduce that the emission of an
487 identifier upon the change of the current conjecture implies the completion of the prior one.

488 Finally, if the conjecture is aborted, then it will never be detected as a definition at all
489 even though its proof has ended. We detect aborted proofs simply by checking the type of
490 the command, assuming that no grammar extension defines `Abort` or `Abort All` equivalents.

491 **Serialization and Version Changes** SerAPI was in theory supposed to help with proof
492 assistant versioning problems. In practice, though, SerAPI itself depends on the version of
493 Coq, and we found we had to break the SerAPI abstraction barrier often as the Coq version
494 changed. In other words, while SerAPI provides a convenient interface to expose certain Coq
495 internals, those internals are not necessarily stable. For example, SerAPI had “can’t-fix” bugs
496 involving nested proofs because the serialization errors occur in the Coq codebase itself [28].
497 SerAPI itself has as of a few days ago been deprecated in favor of a new serializer [29, 18].

498 5.2 Other Proof Assistants

499 Here, we discuss features in other proof assistants in the context of our experiences above.

500 **Project Management** In summary, we are not aware of an elegant and effective solution to
501 package management for other proof assistants, but we believe Isabelle/HOL’s rich archival
502 culture sets a good example to follow. In Isabelle/HOL, the Archive of Formal Proofs (AFP)
503 provides a highly centralized, standard host for proof developments and eases their association
504 with metadata that may be useful for ML. The AFP also neatly versions proof developments
505 for every minor and major version of Isabelle/HOL and semantically groups them in different
506 folders. At the time of writing, the AFP includes 725 proof developments [1], and it already
507 forms the basis of a static dataset for Isabelle/HOL [32]. We suspect the AFP would also
508 make a very strong basis for a proof repair dataset due to its neat versioning.

509 Agda possesses its own library management system, which it uses in combination with
510 Hackage, the Haskell package repository. Anecdotally, researchers we have spoken to cite

23:14 Proof Repair Infrastructure for Supervised Models

511 installation difficulties as a barrier to learning Agda. Lean also has its own package manager
512 but lacks advanced features to address the problems we faced. Isabelle/HOL in general takes
513 an IDE-centric approach to builds and other tooling [50, 60], but does include a notion of
514 sessions that can inherit from other sessions, thereby allowing for some cached dependency
515 sharing. However, one of the authors has found that students learning Isabelle/HOL in a
516 proof automation course struggle to understand how to build dependencies.

517 **Parsing & Serialization** A particularly successful example of an interoperable proof system
518 is MetaMath [39], whose syntax and semantics are so simple that its verifier has been
519 reimplemented in under 1000 lines of Python [61]. This and its centralized proof database
520 has made it a popular choice for ML experts as a benchmark for ML applications in theorem
521 proving [37, 45] as all barriers to serialization can be avoided by modifying a very small
522 parser/verifier. As a trade off, MetaMath does not have a comparable feature set to ITPs
523 like Coq, Lean or Isabelle.

524 More complicated and featureful ITPs have more varied methods: PISA [33] is a bleeding
525 edge Isabelle interaction and proof serialization tool written to support an ML experiment [32].
526 This complements `scala-isabelle` [57], an earlier, less ML-oriented tool which is also actively
527 maintained. As for Lean, PACT [31] presented a dataset (LeanStep) aimed at ML applications
528 that uses Lean’s meta-programming facilities to serialize Lean. LeanStep’s tools weigh under
529 1500 lines of code, which is light compared to line-counts for other serialization efforts.

530 5.3 Recommendations

531 **Project Management** A strong archival effort is the best way forward, though even the
532 best archival infrastructures for proof assistants fall short in multiple ways. For example,
533 while Isabelle/HOL’s AFP makes a natural data source for ML tools for proofs, it does not
534 include any processes for informed consent—the datasets that build on it assume that all
535 publicly available data is fair game. While this assumption is standard in ML for code and
536 proofs, it is not ideal; archival is a natural place to consider it.

537 In addition, though archival makes it possible to associate metadata with proof develop-
538 ments, little consideration is given to metadata that associates definitions and proofs *across*
539 *versions* of a proof development. Presently, we rely on package management tools such as
540 `opam`, which also posed challenges. Though a legitimate argument can be made that package
541 management targets a very different use case from ours and that existing tools are sufficient
542 for that use case, shared high-level libraries and tools on top of existing package managers
543 and in support of bulk efforts like our own would be especially advantageous since such
544 efforts are common when building ML datasets and tools. Nonetheless, package managers
545 themselves warrant some improvement. For example, the problem we encountered of copying
546 switches was due to poor caching of build dependencies, which itself was due to some degree
547 to hard-coding of paths.

548 **Parsing & Serialization** The great flexibility afforded Coq by its extensible grammar allows
549 documents to be more human-friendly and readable, but the lack of syntactic assurances
550 introduces major headaches for automated systems. Ideally, future languages will be struc-
551 tured to be machine-readable human/compiler-out-of-the-loop or to at least provide a public
552 parsing API. For Coq, exposing the classification of a Vernacular command² in SerAPI

² See the `vernac_classification` type.

553 would help substantially and obviate the need for many of the workarounds detailed above.

554 We also recommend a greater emphasis on backwards compatibility and backporting as
555 several useful and even critical features that exist in newer versions of Coq or SerAPI were
556 not suitable for our use. To this end, SerAPI is “still a research, experimental project, and
557 it is expected to evolve considerably” [27] For instance, future plans rebase SerAPI on the
558 language server protocol standard [29], which exposes features like document overviews that
559 appear to list all the definitions in the file and the ability to fold proofs, implying that it has
560 the capability to list theorems and gather the associated lines—one of our current challenges.

561 **Overall** Based on our observations, we make the following broad recommendations for the
562 proof assistant community going forward:

- 563 1. Work with the Isabelle/HOL community to learn how to build **centralized archives** as
564 successful as the AFP for other proof assistants.
- 565 2. Include an **informed consent form** in any centralized archive that allows proof engineers
566 to opt in or out of their developments being used for ML tools.
- 567 3. Determine what **kinds of metadata** within and across proof developments would ease
568 the creation of ML tools for the tasks that matter most, make it easy to **track that**
569 **metadata** inside of any centralized archive, and create standard ways of **associating**
570 **that metadata** with proof developments even outside of centralized archives.
- 571 4. Establish or adopt **open standards** for managing proof developments and interfacing
572 with external tools like modern IDEs.
- 573 5. Develop tools based on these standards to enable **extraction of metadata** relevant to
574 ML tools from non-archival proof developments.
- 575 6. Help proof engineers **port legacy proof developments** to meet those standards, and
576 continue to work on tools for **proof repair and reuse** that ease this burden.
- 577 7. Consider building **shared libraries and tools** optimized for the problems of bulk builds.
- 578 8. Consider limiting the scope of possible **dependency complexity**.
- 579 9. Improve **build caching** across multiple or bulk builds, for example by avoiding hard-coded
580 paths seen in opam.
- 581 10. Consider opening conversations with **language developers and companies** inside and
582 outside of verification about their solutions for package management, distribution, and
583 release management, as these problems are pervasive across all software.

584 **6** Related Work

585 **Datasets & Benchmark Suites** The REPLICA [52] user study collected incremental edit
586 data from eight proof engineers over the course of a month. Due to difficulties recruiting
587 participants, the dataset is too small for data-hungry ML tools. PRISM is less incremental,
588 but much larger. The REPLICA data may make a useful supplement to PRISM.

589 A number of datasets and benchmark suites target *autoformalization*: the automatic
590 translation of natural language mathematics to formal mathematics. Autoformalization
591 datasets consisting of aligned natural and formal language include ProofNet [9] and the
592 Isabelle Parallel Corpus [15]. MiniF2F [68] includes math Olympiad problems formalized in
593 different proof assistants and is used as a benchmark for autoformalization and synthesis.

594 A few datasets and benchmark suites exist for proof synthesis, including CoqGym [65] for
595 Coq and HOList [11] for HOL Light. These datasets include static data from fixed project
596 versions. The distinguishing feature of PRISM is that it describes the project’s history, which
597 is necessary to produce repair examples.

598 We expect there is much that we can learn from ML for code, given the similarities
 599 between code and proofs. A summary of recent work in this space can be found in a survey
 600 paper on neurosymbolic programming [16]. Of particular interest for our work is the question
 601 of whether code distance metrics like CodeBLEU [48] will work well for formal proof.

602 In the field of software engineering, accessible datasets facilitate new research. For
 603 example, Defects4 [34] is a collection of bugs and patches in Java that is frequently used as a
 604 benchmark for program repair [23, 59, 38]. We hope that PRISM will spur new research in
 605 proof repair. We also hope that, by focusing on good benchmarks and metrics for success
 606 early on, we can avoid some of the methodology challenges faced in program repair [46].

607 **Proof Repair** The ML task that our dataset focuses on is proof repair, which is summarized
 608 in the namesake thesis [49]. There is not yet published work we are aware of for ML for
 609 proof repair, though we are aware of ongoing work by other teams in proof assistants other
 610 than Coq. We plan to train and evaluate at least two distinct proof repair models in Coq
 611 using PRISM, and we hope that PRISM makes it easy for others to do the same.

612 Proof repair is closely related to work in proof reuse [25, 54, 14], proof refactoring [63, 62,
 613 55], and proof transformation [44]. These and other related topics in proof engineering have
 614 a long history, described in detail in the proof engineering survey paper QED at Large [50],
 615 as well as in the proof repair namesake thesis [49].

616 Proof repair can be viewed as program repair [40, 30] for proofs. There is a large amount
 617 of work on learning to repair programs, both symbolically (for example, in Getafix [10]) and
 618 neurally (for example, in Break-It-Fix-It [66]). This work may provide useful insights when
 619 building ML datasets, benchmark suites, and models for proof repair, though care must be
 620 taken to consider the differences between typical programs and formal proof developments [49].

621 **Machine Learning for Proofs** Advances in ML have had a transformative effect on many
 622 fields, and theorem provers are not excluded. Examples of recent work on ML for synthesizing
 623 formal proofs include GPT-f [45] and HTPS [36] for Metamath and Lean; Proverbot9001 [56],
 624 ASTactic [65], Tactician [13], and DIVA [26] for Coq; and DeepHOL [11] for HOL Light.
 625 Also of note is recent work on autoformalization in Isabelle/HOL [64], Lean [9], and Coq [21].
 626 More ML work for proofs can be found in QED at Large [50]. Our main goal is to expand the
 627 scope of tasks covered in ML for proofs, reaching important tasks not previously explored.

628 **7 Conclusions & Future Work**

629 We have described a novel dataset and benchmark suite for the ML task of proof repair
 630 centered around the Coq Proof Assistant. The proposed dataset is significantly larger than
 631 any existing alternative, spanning years-long developments collected from a corpus of open-
 632 source Github repositories. We discussed challenges that we encountered during the creation
 633 of the dataset and ramifications for the proof engineering community going forward. The
 634 tools developed to overcome these challenges enable subsequent expansion of the dataset
 635 with supplementary Coq projects and are likely to be useful for creating and interfacing with
 636 datasets for other proof-related ML tasks including, for example, proof synthesis.

637 Moving forward, our immediate plan is to build ML models for proof repair in Coq.
 638 We would also like to develop better metrics for measuring success at repairing definitions.
 639 Finally, we hope to work with the rest of the proof assistant community to address the many
 640 challenges we have highlighted, so that we may steer ML for proofs in the right direction.

641 — References —

- 642 1 Statistics - archive of formal proofs. URL: <https://www.isa-afp.org/statistics/>.
- 643 2 Arpan Agrawal, Emily First, Zhanna Kaufman, Tom Reichel, Shizhuo Zhang, Timothy Zhou,
644 Alex Sanchez-Stern, and Talia Ringer. Proofster. URL: <https://www.alexsanchezstern.com/papers/proofster.pdf>.
- 645 3 Europroofnet. URL: <https://europroofnet.github.io/>.
- 646 4 2nd MATH-AI Workshop at NeurIPS'22, 2021-2022. URL: <https://mathai2022.github.io/>.
- 647 5 Openai. URL: <https://openai.com/>.
- 648 6 Proof engineering, adaptation, repair, and learning for software (pearls). URL: <https://sam.gov/opp/da84366306554cc981f37f703a78c698/view>.
- 649 7 AI for Theorem Proving, 2016-2022. URL: <http://aitp-conference.org/>.
- 650 8 Emilio Jesús Gallego Arias. SerAPI: Machine-Friendly, Data-Centric Serialization for
651 COQ. Technical Report hal-01384408, HAL, 2016. URL: <http://dml.mathdoc.fr/item/hal-01384408/>.
- 652 9 Zhangir Azerbayev, Bartosz Piotrowski, and Jeremy Avigad. Proofnet: A benchmark for
653 autoformalizing and formally proving undergraduate-level mathematics problems. In *Second*
654 *MATH-AI Workshop*, 2022. URL: <https://mathai2022.github.io/>.
- 655 10 Johannes Bader, Andrew Scott, Michael Pradel, and Satish Chandra. Getafix: Learning to fix
656 bugs automatically. *Proc. ACM Program. Lang.*, 3(OOPSLA), oct 2019. doi:10.1145/3360585.
- 657 11 Kshitij Bansal, Sarah M Loos, Markus N Rabe, Christian Szegedy, and Stewart Wilcox. Holist:
658 An environment for machine learning of higher order logic theorem proving. In *ICML*, 2019.
- 659 12 Beyond Bayes: Paths Towards Universal Reasoning Systems, 2022. URL: <https://beyond-bayes.github.io/>.
- 660 13 Lasse Blaauwbroek, Josef Urban, and Herman Geuvers. The tactician: A seamless, interactive
661 tactic learner and prover for coq. In *Intelligent Computer Mathematics: 13th International*
662 *Conference, CICM 2020, Bertinoro, Italy, July 26–31, 2020, Proceedings*, page 271–277, Berlin,
663 Heidelberg, 2020. Springer-Verlag. doi:10.1007/978-3-030-53518-6_17.
- 664 14 Olivier Boite. Proof reuse with extended inductive types. In *Theorem Proving in Higher Order*
665 *Logics: 17th International Conference, TPHOLs 2004, Park City, Utah, USA, September*
666 *14-17, 2004. Proceedings*, pages 50–65, Berlin, Heidelberg, 2004. Springer. doi:10.1007/
667 978-3-540-30142-4_4.
- 668 15 Anthony Bordg, Yiannos A Stathopoulos, and Lawrence C Paulson. A parallel corpus of
669 natural language and isabelle artefacts. In *7th Conference on Artificial Intelligence and*
670 *Theorem Proving (AITP)*, 2022. URL: http://aitp-conference.org/2022/abstract/AITP_2022_paper_8.pdf.
- 671 16 Swarat Chaudhuri, Kevin Ellis, Oleksandr Polozov, Rishabh Singh, Armando Solar-Lezama,
672 Yisong Yue, et al. Neurosymbolic programming. *Foundations and Trends® in Programming*
673 *Languages*, 7(3):158–243, 2021.
- 674 17 Is there a full documentation of coq's grammar? URL: <https://coq.discourse.group/t/is-there-a-full-documentation-of-coqs-grammar/647/10>.
- 675 18 coq_lsp. URL: <https://github.com/ejgallego/coq-lsp>.
- 676 19 Proposal: a custom build tool for coq projects. URL: <https://coq.discourse.group/t/proposal-a-custom-build-tool-for-coq-projects/239/2>.
- 677 20 pycoq. URL: <https://github.com/IBM/pycoq>.
- 678 21 Garrett Cunningham, Razvan C. Bunescu, and David Juedes. Towards autoformalization
679 of mathematics and code correctness: Experiments with elementary proofs, 2023. URL:
680 <https://arxiv.org/abs/2301.02195>, doi:10.48550/ARXIV.2301.02195.
- 681 22 Michel Marie Deza and Elena Deza. *Encyclopedia of Distances*. Springer Berlin Heidel-
682 berg, Berlin, 3 edition, October 2014. URL: <https://link.springer.com/book/10.1007/978-3-642-30958-8>.
- 683
- 684
- 685
- 686
- 687
- 688
- 689
- 690

- 691 23 Thomas Durieux, Matias Martinez, Martin Monperrus, Romain Sommerard, and Jifeng
692 Xuan. Automatic repair of real bugs: An experience report on the defects4j dataset. *CoRR*,
693 abs/1505.07002, 2015. URL: <http://arxiv.org/abs/1505.07002>, arXiv:1505.07002.
- 694 24 Mikhail Evtikhiev, Egor Bogomolov, Yaroslav Sokolov, and Timofey Bryksin. Out of the bleu:
695 how should we assess quality of the code generation models? *arXiv preprint arXiv:2208.03133*,
696 2022.
- 697 25 Amy Felty and Douglas Howe. Generalization and reuse of tactic proofs. In *Logic Programming*
698 *and Automated Reasoning: 5th International Conference, LPAR '94*, pages 1–15, Berlin,
699 Heidelberg, 1994. Springer. doi:10.1007/3-540-58216-9_25.
- 700 26 Emily First and Yuriy Brun. Diversity-driven automated formal verification. In *Proceedings*
701 *of the 44th International Conference on Software Engineering (ICSE)(22–27)*. Pittsburgh, PA,
702 USA. <https://doi.org/10.1145/3510003.3510138>, 2022.
- 703 27 General roadmap. URL: <https://github.com/ejgallego/coq-serapi/issues/252>.
- 704 28 Query ast returns empty result. URL: [https://github.com/ejgallego/coq-serapi/issues/](https://github.com/ejgallego/coq-serapi/issues/117)
705 117.
- 706 29 Serapi 'classic mode' final release notice. URL: [https://github.com/ejgallego/coq-serapi/](https://github.com/ejgallego/coq-serapi/issues/252#issuecomment-1365510329)
707 [issues/252#issuecomment-1365510329](https://github.com/ejgallego/coq-serapi/issues/252#issuecomment-1365510329).
- 708 30 Luca Gazzola, Daniela Micucci, and Leonardo Mariani. Automatic software repair: A survey. In
709 *Proceedings of the 40th International Conference on Software Engineering*, ICSE '18, page 1219,
710 New York, NY, USA, 2018. Association for Computing Machinery. doi:10.1145/3180155.
711 3182526.
- 712 31 Jesse Michael Han, Jason Rute, Yuhuai Wu, Edward W. Ayers, and Stanislas Polu. Proof
713 artifact co-training for theorem proving with language models. *CoRR*, abs/2102.06203, 2021.
714 URL: <https://arxiv.org/abs/2102.06203>, arXiv:2102.06203.
- 715 32 Albert Jiang, Wenda Li, Jesse Han, and Wu Yuhuai. Lisa: Language models of isabelle proofs.
716 In *6th Conference on Artificial Intelligence and Theorem Proving (AITP)*, 2021.
- 717 33 Portal-to-isabelle. URL: <https://github.com/albertqjiang/Portal-to-ISAbelle>.
- 718 34 René Just, Darioush Jalali, and Michael D. Ernst. Defects4J: A Database of existing faults to
719 enable controlled testing studies for Java programs. In *ISSTA 2014, Proceedings of the 2014*
720 *International Symposium on Software Testing and Analysis*, pages 437–440, San Jose, CA,
721 USA, July 2014. Tool demo.
- 722 35 Daniel Kühlwein, Jasmin Christian Blanchette, Cezary Kaliszyk, and Josef Urban. Mash:
723 Machine learning for sledgehammer. In Sandrine Blazy, Christine Paulin-Mohring, and David
724 Pichardie, editors, *Interactive Theorem Proving*, pages 35–50, Berlin, Heidelberg, 2013. Springer
725 Berlin Heidelberg.
- 726 36 Guillaume Lample, Marie-Anne Lachaux, Thibaut Lavril, Xavier Martinet, Amaury Hayat,
727 Gabriel Ebner, Aurélien Rodriguez, and Timothée Lacroix. Hypertree proof search for neural
728 theorem proving. *arXiv preprint arXiv:2205.11491*, 2022.
- 729 37 Guillaume Lample, Marie-Anne Lachaux, Thibaut Lavril, Xavier Martinet, Amaury Hayat,
730 Gabriel Ebner, Aurélien Rodriguez, and Timothée Lacroix. Hypertree proof search for neural
731 theorem proving, 2022. URL: <https://arxiv.org/abs/2205.11491>, doi:10.48550/ARXIV.
732 2205.11491.
- 733 38 Thibaud Lutellier, Hung Viet Pham, Lawrence Pang, Yitong Li, Moshi Wei, and Lin Tan.
734 Coconut: Combining context-aware neural translation models using ensemble for program
735 repair. In *Proceedings of the 29th ACM SIGSOFT International Symposium on Software*
736 *Testing and Analysis*, ISSTA 2020, page 101–114, New York, NY, USA, 2020. Association for
737 Computing Machinery. doi:10.1145/3395363.3397369.
- 738 39 Norman D. Megill and David A. Wheeler. *Metamath: A Computer Lan-*
739 *guage for Mathematical Proofs*. Lulu Press, Morrisville, North Carolina, 2019.
740 <http://us.metamath.org/downloads/metamath.pdf>.
- 741 40 Martin Monperrus. Automatic software repair: A bibliography. *ACM Comput. Surv.*, 51(1),
742 jan 2018. doi:10.1145/3105906.

- 743 41 James Munkres. Algorithms for the Assignment and Transportation Problems. *Journal of*
744 *the Society for Industrial and Applied Mathematics*, 5(1):32–38, 1957. Publisher: Society for
745 Industrial and Applied Mathematics. URL: <https://www.jstor.org/stable/2098689>.
- 746 42 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: A method for automatic
747 evaluation of machine translation. In *Proceedings of the 40th Annual Meeting on Association for*
748 *Computational Linguistics*, ACL '02, page 311–318, USA, 2002. Association for Computational
749 Linguistics. doi:10.3115/1073083.1073135.
- 750 43 Ofir Pele and Michael Werman. Fast and robust earth mover’s distances. In *2009 IEEE 12th*
751 *International Conference on Computer Vision*, pages 460–467, 2009. doi:10.1109/ICCV.2009.
752 5459199.
- 753 44 Frank Pfenning. *Proof Transformations in Higher-Order Logic*. PhD thesis, Carnegie Mellon
754 University Pittsburgh, 1987.
- 755 45 Stanislas Polu and Ilya Sutskever. Generative language modeling for automated theorem
756 proving. *CoRR*, abs/2009.03393, 2020. URL: <https://arxiv.org/abs/2009.03393>, arXiv:
757 2009.03393.
- 758 46 Zichao Qi, Fan Long, Sara Achour, and Martin Rinard. An analysis of patch plausibility and
759 correctness for generate-and-validate patch generation systems. In *Proceedings of the 2015*
760 *International Symposium on Software Testing and Analysis*, ISSTA 2015, page 24–36, New
761 York, NY, USA, 2015. Association for Computing Machinery. doi:10.1145/2771783.2771791.
- 762 47 Markus N. Rabe and Christian Szegedy. Towards the automatic mathematician. In André
763 Platzer and Geoff Sutcliffe, editors, *Automated Deduction – CADE 28*, pages 25–37, Cham,
764 2021. Springer International Publishing.
- 765 48 Shuo Ren, Daya Guo, Shuai Lu, Long Zhou, Shujie Liu, Duyu Tang, Neel Sundaresan, Ming
766 Zhou, Ambrosio Blanco, and Shuai Ma. Codebleu: a method for automatic evaluation of
767 code synthesis. *CoRR*, abs/2009.10297, 2020. URL: <https://arxiv.org/abs/2009.10297>,
768 arXiv:2009.10297.
- 769 49 Talia Ringer. *Proof Repair*. PhD thesis, University of Washington, 2021.
- 770 50 Talia Ringer, Karl Palmkog, Ilya Sergey, Milos Gligoric, and Zachary Tatlock. QED at large:
771 A survey of engineering of formally verified software. *CoRR*, abs/2003.06458, 2020. URL:
772 <https://arxiv.org/abs/2003.06458>, arXiv:2003.06458.
- 773 51 Talia Ringer, RanDair Porter, Nathaniel Yazdani, John Leo, and Dan Grossman. Proof repair
774 across type equivalences. In *Proceedings of the 42nd ACM SIGPLAN International Conference*
775 *on Programming Language Design and Implementation*, PLDI 2021, page 112–127, New York,
776 NY, USA, 2021. Association for Computing Machinery. doi:10.1145/3453483.3454033.
- 777 52 Talia Ringer, Alex Sanchez-Stern, Dan Grossman, and Sorin Lerner. REPLica: REPL
778 instrumentation for Coq analysis. In *Proceedings of the 9th ACM SIGPLAN International*
779 *Conference on Certified Programs and Proofs*, CPP 2020, page 99–113, New York, NY, USA,
780 2020. Association for Computing Machinery. doi:10.1145/3372885.3373823.
- 781 53 Talia Ringer, Nathaniel Yazdani, John Leo, and Dan Grossman. Adapting proof automation to
782 adapt proofs. In *Proceedings of the 7th ACM SIGPLAN International Conference on Certified*
783 *Programs and Proofs*, CPP 2018, page 115–129, New York, NY, USA, 2018. Association for
784 Computing Machinery. doi:10.1145/3167094.
- 785 54 Talia Ringer, Nathaniel Yazdani, John Leo, and Dan Grossman. Ornaments for proof reuse in
786 coq. In *Interactive Theorem Proving*, 2019.
- 787 55 Valentin Robert. *Front-end tooling for building and maintaining dependently-typed functional*
788 *programs*. PhD thesis, UC San Diego, 2018.
- 789 56 Alex Sanchez-Stern, Yousef Alhessi, Lawrence Saul, and Sorin Lerner. Generating cor-
790 rectness proofs with neural networks, 2019. URL: <https://arxiv.org/abs/1907.07794>,
791 doi:10.48550/ARXIV.1907.07794.
- 792 57 scala-isabelle. URL: <https://dominique-unruh.github.io/scala-isabelle/>.

- 793 58 T.F. Smith and M.S. Waterman. Identification of common molecular subsequences. *Journal of*
794 *Molecular Biology*, 147(1):195–197, 1981. URL: [https://www.sciencedirect.com/science/](https://www.sciencedirect.com/science/article/pii/0022283681900875)
795 [article/pii/0022283681900875](https://www.sciencedirect.com/science/article/pii/0022283681900875), doi:[https://doi.org/10.1016/0022-2836\(81\)90087-5](https://doi.org/10.1016/0022-2836(81)90087-5).
- 796 59 Ming Wen, Junjie Chen, Rongxin Wu, Dan Hao, and Shing-Chi Cheung. Context-aware patch
797 generation for better automated program repair. In *Proceedings of the 40th International*
798 *Conference on Software Engineering, ICSE '18*, page 1–11, New York, NY, USA, 2018.
799 Association for Computing Machinery. doi:10.1145/3180155.3180233.
- 800 60 Makarius Wenzel. Isabelle/jedit — a prover IDE within the PIDE framework. *CoRR*,
801 [abs/1207.3441](https://arxiv.org/abs/1207.3441), 2012. URL: <http://arxiv.org/abs/1207.3441>, arXiv:1207.3441.
- 802 61 mmverify.py. URL: <https://github.com/david-a-wheeler/mmverify.py>.
- 803 62 Iain Johnston Whiteside. *Refactoring proofs*. PhD thesis, University of Edinburgh, November
804 2013. URL: <http://hdl.handle.net/1842/7970>.
- 805 63 Karin Wibergh. Automatic refactoring for agda. Master’s thesis, Chalmers University of
806 Technology and University of Gothenburg, 2019.
- 807 64 Yuhuai Wu, Albert Q Jiang, Wenda Li, Markus N Rabe, Charles Staats, Mateja Jamnik,
808 and Christian Szegedy. Autoformalization with large language models. *arXiv preprint*
809 *arXiv:2205.12615*, 2022.
- 810 65 Kaiyu Yang and Jia Deng. Learning to prove theorems via interacting with proof assistants.
811 In *International Conference on Machine Learning (ICML)*, Long Beach, CA, USA, 2019. URL:
812 <http://proceedings.mlr.press/v97/yang19a/yang19a.pdf>.
- 813 66 Michihiro Yasunaga and Percy Liang. Break-it-fix-it: Unsupervised learning for program
814 repair. In *International Conference on Machine Learning*, pages 11941–11952. PMLR, 2021.
- 815 67 John R. Zech, Marcus A. Badgeley, Manway Liu, Anthony B. Costa, Joseph J. Titano, and
816 Eric Karl Oermann. Variable generalization performance of a deep learning model to detect
817 pneumonia in chest radiographs: A cross-sectional study. *PLOS Medicine*, 15(11):e1002683,
818 November 2018. doi:10.1371/journal.pmed.1002683.
- 819 68 Kunhao Zheng, Jesse Michael Han, and Stanislas Polu. minif2f: a cross-system benchmark for
820 formal olympiad-level mathematics. In *International Conference on Learning Representations*,
821 2022. URL: <https://openreview.net/forum?id=9ZPegFuFTFv>.